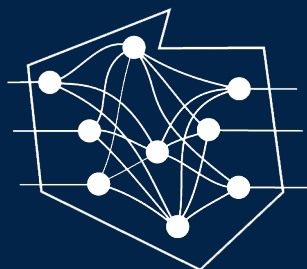# PROGRESS IN POLISH ARTIFICIAL INTELLIGENCE RESEARCH 6

Edited by:
**Rafał Doroz**
**Beata Zielosko**

*6th Polish Conference*
*on Artifical Intelligence (PP-RAI'2025)*

**07–09.04.2025, Katowice, Poland**

**PP-RAI'2025**

# Progress in Polish Artificial Intelligence Research 6

**Edited by:**
**Rafał Doroz, Beata Zielosko**

# Preface

This volume gathers the research papers presented at the 6th Polish Conference on Artificial Intelligence (PP-RAI 2025), held on April 7–9, 2025, at the University of Silesia in Katowice, Poland. The conference was organized by the University of Silesia in Katowice in collaboration with the Polish Alliance for the Development of Artificial Intelligence, and was held under the honorary patronage of the Rector of the University of Silesia in Katowice.

Established in 2018, the Polish Alliance for the Development of Artificial Intelligence emerged from the initiative of the Polish research community conducting research in artificial intelligence (AI) and machine learning (ML) domains. PP-RAI comprises four member organizations: the Polish Neural Network Society, the Polish Special Interest Group on Machine Learning, the Polish Chapter of the IEEE Systems, Man, and Cybernetics Society, the Polish Chapter of the IEEE Computational Intelligence Society. The alliance aims to strengthen collaboration among researchers in Poland while supporting the advancement and broader dissemination of artificial intelligence and machine learning research.

The annual PP-RAI conferences serve as a gathering point for researchers in artificial intelligence and machine learning. They offer an opportunity to present and exchange original research across a wide range of AI/ML topics. The events also encourage discussions on key research initiatives, projects, and developments in AI and ML both in Poland and internationally. Moreover, the conferences foster collaboration between the AI/ML community, academic institutions, and industry partners, supporting the growth and advancement of the field.

The 2025 edition of the PP-RAI conference featured three plenary lectures by leading experts, complemented by three panels on current AI topics, three poster sessions, and three parallel oral sessions. Out of 195 submissions spanning 17 thematic tracks, 51 papers were selected for inclusion in these proceedings by the Track Chairs, based on their quality and research significance.

We would like to express our gratitude to all the Authors who submitted papers for their efforts in preparing and presenting their research. We also thank the Reviewers for their dedicated service to the community, the Track Chairs and the Program Committee for their selection of the final papers, and the Steering Committee for their guidance throughout the entire process. We are particularly

grateful to the Local Organizing Committee and the Young Organizing Committee for their support in preparing and organizing this event. Finally, we would like to acknowledge the valuable support of our Partners and Sponsors.

Beata Zielosko
Rafał Doroz

April, 2025

# Organization

**GENERAL CHAIR**

| | |
|---|---|
| **Ryszard Koziołek** | Rector of the University of Silesia |

**STEERING COMMITTEE**

| | |
|---|---|
| **Ireneusz Czarnowski** | Gdynia Maritime University |
| **Włodzisław Duch** | Nicolaus Copernicus University in Toruń |
| **Janusz Kacprzyk** | Systems Research Institute, Polish Academy of Sciences, Warsaw |
| **Jan Kozak** | University of Economics in Katowice |
| **Leszek Rutkowski** | Systems Research Institute, Polish Academy of Sciences, Warsaw |
| **Rafał Scherer** | Czestochowa University of Technology |

**PP-RAI CONFERENCE CHAIRS**

| | |
|---|---|
| **Beata Zielosko** | University of Silesia in Katowice PP-RAI'2025 |
| **Jacek Mańdziuk** | Warsaw University of Technology PP-RAI'2024 |
| **Aleksander Byrski** | AGH University of Science and Technology PP-RAI'2026 |

**PROGRAM CHAIR**

| | |
|---|---|
| **Beata Zielosko** | University of Silesia in Katowice |

**PROGRAM COMMITTEE**

| | |
|---|---|
| **Jarosław Arabas** | Warsaw University of Technology |
| **Piotr Artiemjew** | University of Warmia and Mazury |
| **Michał Baczyński** | University of Silesia in Katowice |
| **Przemysław Biecek** | Warsaw University of Technology |
| **Urszula Boryczka** | University of Silesia in Katowice |

Aleksander Byrski          AGH University of Science and Technology
Leszek Chmielewski         Warsaw University of Life Sciences
Ireneusz Czarnowski        Gdynia Maritime University
Rafał Doroz                University of Silesia in Katowice
Wodzisław Duch             Nicolaus Copernicus University
Krzysztof Gajowniczek      Warsaw University of Life Sciences
Tomasz Gambin              Warsaw University of Technology
Maria Ganzha               Warsaw University of Technology
Maciej Grzenda             Warsaw University of Technology
Katarzyna Harężlak         Silesian University of Technology
Agnieszka Jastrzębska      Warsaw University of Technology
Janusz Kacprzyk            Systems Research Institute, Polish Academy of
                           Sciences
Paweł Kasprowski           Silesian University of Technology
Stanisław Kaźmierczak       Warsaw University of Technology
Piotr A. Kowalski          AGH University of Science and Technology
Dariusz Król               Wrocław University of Science and Technology
Marcin Kurdziel            AGH University of Science and Technology
Halina Kwaśnicka           Wrocław University of Technology
Bogdan Kwolek              AGH University of Science and Technology
Piotr Lipiński             Lodz University of Technology
Agnieszka Ławrynowicz      Poznań University of Technology
Szymon Łukasik             AGH University of Science and Technology
Mikołaj Małkiński          Warsaw University of Technology
Jacek Mańdziuk             Warsaw University of Technology
Jakub Nalepa               Silesian University of Technology
Grzegorz J. Nalepa         Jagiellonian University
Robert Nowak               Warsaw University of Technology
Agnieszka Nowak-Brze-      University of Silesia in Katowice
zińska
Karol Opara                Systems Research Institute, Polish Academy of
                           Sciences
Tomasz Orczyk              University of Silesia in Katowice
Piotr Pęzik                University of Lodz
Krzysztof Kotowski         KP Labs
Maciej Piasecki            Wrocław University of Science and Technology
Dariusz Plewczyński         Warsaw University of Technology

| | |
|---|---|
| **Piotr Porwik** | University of Silesia in Katowice |
| **Małgorzata Przybyła-Kasperek** | University of Silesia in Katowice |
| **Jacek Rumiński** | Gdańsk University of Technology |
| **Leszek Rutkowski** | Systems Research Institute, Polish Academy of Sciences |
| **Khalid Saeed** | Białystok University of Technology |
| **Wojciech Sałabun** | National Institute of Telecommunications |
| **Rafał Scherer** | Czestochowa University of Technology |
| **Krzysztof Simiński** | Silesian University of Technology |
| **Piotr Skrzypczyński** | Poznań University of Technology |
| **Przemysław Spurek** | Jagiellonian University |
| **Jerzy Stefanowski** | Poznań University of Technology |
| **Dominik Ślęzak** | University of Warsaw |
| **Julian Szymański** | Gdańsk University of Technology |
| **Anna Timofiejczuk** | Silesian University of Technology |
| **Tomasz Trzciński** | Warsaw University of Technology |
| **Piotr Wasilewski** | Systems Research Institute, Polish Academy of Sciences |
| **Jarosław Wąs** | AGH University of Science and Technology |
| **Tomasz Wesołowski** | University of Silesia in Katowice |
| **Adam Wojciechowski** | Lodz University of Technology |
| **Michał Woźniak** | Wrocław University of Science and Technology |
| **Cezary Zieliński** | Warsaw University of Technology |
| **Beata Zielosko** | University of Silesia |
| **Maciej Zięba** | Wrocław University of Technology |
| **Adam Żychowski** | Warsaw University of Technology |

**LOCAL ORGANIZING COMMITTEE**
**(Local Team – University of Silesia in Katowice)**

**Beata Zielosko**
**Rafał Doroz**
**Agnieszka Nowak-Brzezińska**
**Aleksander Bogusz**
**Urszula Boryczka**
**Kornel Chromiński**
**Agata Kołodziejczyk**

**Magda Korbela**
**Agata Krawczyk-Kalitowska**
**Przemysław Kudłacik**
**Arkadiusz Nowakowski**
**Tomasz Orczyk**
**Justyna Przybylska**
**Małgorzata Przybyła-Kasperek**
**Bartłomiej Płaczek**
**Rafał Skinderowicz**
**Agnieszka Sobczyk**
**Joanna Stelmaszak**
**Tomasz Wesołowski**

## YOUNG ORGANIZING COMMITTEE

**Natalia Balicka**
**Bartłomiej Banasik**
**Bartłomiej Barański**
**Mikołaj Biczak**
**Adrian Biskup**
**Bartłomiej Chyra**
**Andrzej Cieślik**
**Mikołaj Feser**
**Bartosz Gruszka**
**Jagoda Guz**
**Jeremi Kowalski**
**Bartłomiej Krypczyk**
**Aleksandra Lozio**
**Magdalena Lytkowicz**
**Kamil Magiera**
**Daniel Nowak**
**Klaudia Rosa**
**Bartłomiej Skwara**
**Mikołaj Susek**
**Laura Wesołowski**
**Wiktor Włodarczyk**
**Maria Zacharowa**
**Wiktoria Zelawska**

**Reviewers:**

Kamil Adamczewski
Izabella Antoniuk
Piotr Artiemjew
Dawid Baran
Dariusz Barbucha
Paweł Batorski
Dominik Belter
Urszula Bentkowska
Marcin Bernas
Weronika Borek-Marciniec
Piotr Borycki
Mariusz Boryczka
Mariusz Bujny
Aleksander Byrski
Leszek Chmielewski
Kornel Chromiński
Ireneusz Czarnowski
Przemysław Dolata
Rafał Doroz
Włodzisław Duch
Agnieszka Duraj
Krzysztof Dyczkowski
Amgad Elsayed
Paweł Forczmański
Oleksii Furman
Ewelina Gajewska
Tomasz Gambin
Maria Ganzha
Piotr Gawron
Krzysztof Gdawiec
Bartłomiej Gintowt
Tomasz Górecki
Artur Gunia
Katarzyna Harężlak
Piotr Helm

Kornel Howil
Katarzyna Jabłońska
Agnieszka Jastrzębska
Katarzyna Kaczmarek-Majer
Paweł Kasprowski
Włodzimierz Kasprzak
Artur Kasymov
Stanisław Kaźmierczak
Joanna Klikowska
Jakub Klikowski
Krzysztof Kluza
Joanna Komorniczak
Agnieszka Konys
Mirosław Kordos
Tomasz Kornuta
Marcin Kostrzewa
Krzysztof Kotowski
Jędrzej Kozal
Rafał Kozik
Wojciech Kozłowski
Marek Kraft
Anna Król
Patryk Krukowski
Radosław Kuczbański
Przemysław Kudłacik
Bogdan Kwolek
Agnieszka Lazarowska
Łukasz Lenkiewicz
Wojciech Lesiński
Piotr Lipiński
Szymon Łukasik
Karol Majek
Dawid Malarz
Paweł Malczyk
Jacek Mańdziuk

Patryk Marszałek

Zofia Matusiewicz

Teresa Mroczek

Jakub Nalepa

Robert Nowak

Przemysław Nowak

Agnieszka Nowak-Brzezińska

Arkadiusz Nowakowski

Karol Opara

Tomasz Orczyk

Wiesław Paj

Piotr Pałka

Dariusz Palt

Anastasiya Pechko

Barbara Pękala

Paweł Pełka

Maciej Piasecki

Bartłomiej Płaczek

Leszek Podsędkowski

Piotr Porwik

Jędrzej Potoniec

Małgorzata Przybyła-Kasperek

Jacek Rumiński

Bogdan Ruszczak

Maciej Rut

Melika Sadeghi

Wojciech Sałabun

Rafal Scherer

Dawid Seredyński

Krzysztof Simiński

Aleksander Skakovski

Paweł Skruch

Piotr Skrzypczyński

Andrzej Śluzek

Weronika Smolak-Dyżewska

Georgii Stanishevskii

Jakub Steczkiewicz

Maciej Stefańczyk

Barbara Strug

Michał Stypułkowski

Szymon Świderski

Maciej Świechowski

Mirosław Szaban

Marta Szarmach

Tomasz Szczepanik

Julian Szymański

Arkadiusz Tomczyk

Joanna Waczyńska

Jarosław Wąs

Łukasz Wawrowski

Anna Wawrzyńczak-Szaban

Weronika Węgier

Tomasz Wesołowski

Patryk Wielopolski

Grzegorz Wilczyński

Artur Wilkowski

Dawid Wiśniewski

Szymon Wojciechowski

Adam Wojciechowski

Piotr Wójcik

Eryk Wójcik

Tomasz Wojnar

Agnieszka Wosiak

Krzysztof Wróbel

Anna Wróblewska

Dawid Wymoczyło

Cezary Zieliński

Beata Zielosko

Paweł Zyblewski

Adam Żychowski

Patryk Żywica

# Contents

x

xii

# CHAPTER 1

# Data Mining and Machine Learning

Track Chairs:

- prof. Michał Woźniak – Wrocław University of Science and Technology

- prof. Ireneusz Czarnowski – Gdynia Maritime University

- prof. Rafał Doroz – University of Silesia in Katowice

# Feature Selection in Traditional Music Classification

**Paweł Grabczyński**[0009−0002−3730−7949],
**Daniel Kostrzewa**[0000−0003−2781−3709],
**Katarzyna Harężlak** [0000−0003−3573−9772]

*Silesian University of Technology*
*Department of Applied Informatics*
*Akademicka 16, 44-100 Gliwice, Poland*
*pgrabczynski@poczta.onet.pl, daniel.kostrzewa@polsl.pl,*
*katarzyna.harezlak@polsl.pl*

**Abstract.** *This study explores the impact of feature selection on the classification of geographical origin of traditional music. It builds upon previous research, which focused on classification tasks using a dataset of over 12,000 traditional music pieces from around the world. Various methods were applied to identify the most informative features while minimizing dimensionality. Special attention was given to the role of timbre, melodic and rhythmic features in classification performance.*
**Keywords:** *traditional music classification, feature selection, RFE*

## 1. Introduction

Identifying the geographical origin of a musical piece based only on its musical features is a challenging task. Traditional styles have blended due to cultural influences, making classification more complex. However, machine learning offers innovative approaches to this task. The process typically involves data acquisition and feature extraction, followed by classifier selection and optimization. In this study, we are focused on one element of the whole pipeline, which is the dimensionality reduction of a given dataset.

Recent advancements in audio classification increasingly rely on deep learning approaches, where feature extraction and classification are performed jointly

by deep neural networks. Han et al. [1] have demonstrated that CNNs can effectively recognize predominant instruments in polyphonic music, outperforming traditional methods. Solanki and Pandey [2] have further refined CNN-based approaches, achieving very good results in musical instrument recognition by optimizing network architectures and activation functions. Despite these advancements, handcrafted feature extraction techniques remain widely used, especially in studies where interpretability and computational efficiency are critical factors.

Feature extraction is as crucial as obtaining high-quality materials. Specialized libraries and tools can be used to determine attributes from a musical piece, such as the MARSYAS software, the VAMP plugin system, or the Librosa library [3]. A good practice is to categorize the extracted features according to specific musical characteristics, such as melody, rhythm, and timbre. This allows for assessing the effectiveness of classifier predictions based on specific feature subsets [4, 5].

Most datasets include Mel Frequency Cepstral Coefficients (MFCC), spectral flux, and onset detection, which tracks amplitude changes to determine onset rate. Among melodic features, pitch class frequency is commonly used. Liu et al. [6] proposed employing this statistic to compute the chroma contrast attribute, which measures the ratio between the six most frequent pitch classes and the remaining ones. This value is expected to be lower in Western music compared to traditional Chinese music. Gomez et al. [7] suggest calculating deviations from equal temperament tuning (A = 440 Hz), as traditional music often includes tones outside the Western scale, such as quarter tones in the Arabic scale. Statistical measures such as mean, standard deviation, minimum, and maximum values are frequently derived from the extracted features [8, 4, 9].

Importantly, in recent studies [10], one of the first attempts at traditional music origin classification was made with promising results. The authors used a proprietary dataset consisting of over 12,000 songs described by over 300 attributes (numerical features with decimal values). It is well-balanced in terms of country distribution, with all classes containing approximately 250 to 300 examples. The ratio between the largest and smallest class in this category is 1.47. The task was resolved on classical classifiers like K-Nearest Neighbors (KNN), Random Forest (RF), and Neural Networks. This study investigates whether reducing the number of features in traditional music classification can maintain or improve accuracy. The accuracy metric is used as it provides the most reliable measure of classifier performance in this case.

## 2. The Study

The success of machine learning models largely depends on the quality of the features used to train them. Having irrelevant variables can reduce the accuracy of many models. Therefore, sometimes, instead of using the full set of attributes, feature engineering is utilized. It helps highlight the most important patterns and relationships in the data, allowing the machine-learning model to learn more meaningful patterns in the data. The application of feature engineering techniques has many advantages: (1) it improves model performance and reduces the risk of overfitting; (2) it makes the model more robust to outliers and other anomalies; (3) it often requires fewer computational resources; (4) it improves model interpretability making model results more straightforward to understand and interpret [11].

Feature engineering techniques include such operations as data transformation, selecting a subset of features, or creating new features based on existing ones. Data transformation involves, *inter alia*, cleaning data and ensuring that variables are on the same scale, and all features are within an acceptable range for the model. Feature selection is the process of choosing the features in the data that impact the target variable the most. Three approaches can be mentioned here: filter, embedded, and wrapper methods. Feature extraction relies on automatically creating new variables by extracting them from raw data to reduce the volume of data into a more manageable set for modeling.

## 3. Experiments

This study represents the development of research on the classification of the country (44 classes), subregion (19 classes) or region (6 classes) of origin of traditional music pieces. The purpose of the study was to analyze whether a selected subset of features could provide similar or better results. The research was carried out using the Python language and popular libraries for data analysis and machine learning, including Scikit-learn (1.6.0), Numpy (2.0.2), and Tensorflow (2.18.0).

The dataset was split into a training set (70%) and a test set (30%). The study utilizes hyperparameter optimization results from previous research, where 10-fold cross-validation and grid search techniques were applied. Based on these results, the selected hyperparameters for Random Forest were as follows: n_estimators = 500, max_features = 10, and bootstrap = False. For K-Nearest Neighbors: n_neighbors = 3, weights = distance, and algorithm = auto.

4

The feature selection process was conducted using three main categories of methods: filter, wrapper, and embedded approaches.

Filtering methods involved Pearson correlation and the chi-square test. They were applied to remove highly correlated or weakly associated features. These techniques helped in reducing redundancy in the dataset before applying more advanced selection strategies.

Wrapper methods included Recursive Feature Elimination (RFE), which was used to identify the relationship between variables. Two variations of RFE were applied: one using logistic regression (RFE: LR) and another employing a decision tree (RFE: DT) to determine feature importance.

Embedded methods focused on Random Forest-based feature selection, which was implemented in two ways: FS: RF – progressively adding all selected features. FS: RF+ – including only features that contributed to improved accuracy.

Additionally, dimensionality reduction methods, including autoencoder and Isomap, were examined. For the autoencoder, the results were examined both through the encoder output and by analyzing the summed input weights.

## 4. Results

The evaluation of feature selection methods revealed substantial differences in their effectiveness. The use of filter methods resulted in a decrease in models' accuracy. Similar results were observed for dimensionality reduction methods, as neither autoencoder nor Isomap improved classification performance. The RFE method and FS with RF provided the highest classification accuracy. Improvements were evident across all classification levels: country, subregion, and region, as demonstrated in Table 1, Table 2 and Table 3.

Table 1. KNN and RF classification accuracy for traditional music country classification with feature selection (number of selected features in parentheses)

| | All | RFE: LR | RFE: DT | FS: RF | FS: RF+ | Autoencoder | Isomap |
|---|---|---|---|---|---|---|---|
| **RF** | 60.03% | 60.03% (287) | 61.76% (173) | 61.38% (171) | 59.20% (58) | 46.39 % | 20.42% |
| **KNN** | 54.28% | 55.80% (90) | 59.35% (55) | **62.00%** (113) | 61.50% (63) | 37.55% | 25.55 % |

Figure 1 illustrates the increase in country classification accuracy depending on the number of features. It is clearly visible that individual classifiers quickly achieve performance comparable to or better than the full feature set. KNN outperformed the full feature set when using only 15% of the features, while RF achieved

Table 2. KNN and RF classification accuracy for traditional music subregion classification with feature selection (number of selected features in parentheses)

|  | All | RFE: LR | RFE: DT | FS: RF | FS: RF+ | Autoencoder | Isomap |
|---|---|---|---|---|---|---|---|
| RF | 57.67% | 58.76% (204) | 59.47% (109) | 59.33% (115) | 30.68% (36) | 41.80 % | 25.63 % |
| KNN | 58.04% | 62.07% (87) | 61.00% (73) | 64.04% (127) | **64.69%** (77) | 41.39 % | 35.17 % |

Table 3. KNN and RF classification accuracy for traditional music region classification with feature selection (number of selected features in parentheses)

|  | All | RFE: LR | RFE: DT | FS: RF | FS: RF+ | Autoencoder | Isomap |
|---|---|---|---|---|---|---|---|
| RF | 66.43% | 67.36% (180) | 67.21% (101) | 66.64% (165) | 48.70% (31) | 56.76 % | 40.09 % |
| KNN | 68.12% | 68.90% (88) | 69.80% (204) | 70.55% (212) | **72.43%** (61) | 57.62 % | 50.12 % |



Figure 1. KNN and RF country classification with RFE and RF feature selection (dots indicate maximum accuracy)

this at 49%. The most significant improvement was observed for KNN with feature selection using RF (+7.72% with 36% of the features).

Interestingly, beyond a certain point, adding more features to the KNN model leads to a decline in performance, eventually converging to around 55%. In contrast, RF does not exhibit this pattern – after stabilizing at 60%, further feature additions do not significantly affect performance.

An unexpected decline in classification performance occurred when using feature selection with RF, where only features that improved prediction accuracy at the moment were included in the model. This decrease was most evident for subregion (decrease of 27%) and region (decrease of 18%), suggesting that some features, despite initially appearing less relevant, played a crucial role when interacting with others.

The analysis of the most frequent features in the selected subsets revealed a dominance of timbre features. Among the 50 most common, 86% were timbre-related (e.g., Mel-Frequency Cepstral Coefficients and Flatness), while rhythmic and melodic features accounted for 8% and 6%, respectively.

## 5. Conclusions

The study has shown that the accuracy of traditional music classification is highly dependent on feature selection. Filter methods and dimensionality reduction techniques did not improve results, as they removed features without considering their interactions within the model. In contrast, RFE and FS with RF significantly reduced the number of features while improving model performance due to their ability to identify relationships between features.

An analysis of feature importance revealed the dominance of timbre features (86% among the top 50 features) despite constituting 69% of the dataset, confirming their key role in classification. Melodic and rhythmic features, which account for 31% of the dataset, appeared in only 14% of the most important features. It raises the question of whether they adequately capture the intended musical characteristics.

To enhance data reliability, consultation with ethnomusicologists would be beneficial. Their expertise could help verify whether the selected musical pieces accurately represent the music of a given region.

## References

[1] Han, Y., Kim, J., and Lee, K. Deep convolutional neural networks for predominant instrument recognition in polyphonic music. In *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pages 208–221. IEEE, 2016.

[2] Solanki, A. and Pandey, S. Music instrument recognition using deep convolutional neural networks. In *International Journal of Information Technology*, pages 1659–1668. Springer, 2019.

[3] Moffat, D., Ronan, D., and Reiss, J. D. An evaluation of audio feature extraction toolboxes. In *Proc. 18th Int. Conference on Digital Audio Effects (DAFx-15)*. 2015.

[4] Kedyte, V., Panteli, M., Weyde, T., and Dixon, S. Geographical origin prediction of folk music recordings from the United Kingdom. In *18th International Society for Music Information Retrieval Conference*. 2017.

[5] Liu, Y., Xiang, Q., Wang, Y., and Cai, L. Cultural style based music classification of audio signals. In *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 57–60. IEEE, 2009.

[6] Littlefield, M. D. *Folk Music in New England: A Living Tradition. MSc Thesis.* Liberty University, 2020.

[7] Gómez, E., Haro, M., and Herrera, P. Music and geography: Content description of musical audio from different parts of the world. In *ISMIR*, pages 753–758. 2009.

[8] Zhou, F., Claire, Q., and King, R. D. Predicting the geographical origin of music. In *2014 IEEE International Conference on Data Mining*, pages 1115–1120. IEEE, 2014.

[9] Schedl, M. and Zhou, F. Fusing web and audio predictors to localize the origin of music pieces for geospatial retrieval. In *Advances in Information Retrieval: 38th European Conference on IR Research, ECIR 2016, Padua, Italy, March 20–23, 2016. Proceedings 38*, pages 322–334. Springer, 2016.

[10] Kostrzewa, D. and Grabczyński, P. From sound to map: Predicting geographic origin in traditional music works. In *International Conference on Computational Science*, pages 174–188. Springer, 2024.

[11] Jia, W., Sun, M., Lian, J., and Hou, S. Feature dimensionality reduction: a review. *Complex & Intelligent Systems*, 8(3):2663–2693, 2022.

# From Data to Decisions: Comparing Causal Discovery Methods on a Benchmark Dataset

**Mikołaj Jarosławski**[0009−0003−1999−4949], **Dominik Sepioło**[0000−0001−7746−3781], **Antoni Ligęza**[0000−0002−6573−4246]

*AGH University of Krakow*
*Department of Applied Computer Science*
*al. A. Mickiewicza 30, 30-059 Kraków, Poland*
*{mikjar, sepiolo, ligeza}@agh.edu.pl*

**Abstract.** *Causal discovery algorithms are essential for uncovering causal relationships in observational data, allowing deeper insights into complex systems. In this paper, we evaluate three prominent causal discovery methods: PC, GES, and LiNGAM. Using the LUCAS dataset, we compare these methods based on reconstruction accuracy, computational efficiency, and their ability to incorporate domain knowledge. We analyze the impact of varying sample sizes and discuss the strengths and limitations of each approach. This work provides practical insights into the selection and application of causal discovery techniques in both research and real world settings.*
**Keywords:** *explainable artificial intelligence, causal discovery, causal graphs*

## 1. Introduction

Causal discovery has emerged as a crucial tool in many scientific domains. It enables researchers to infer causal relationships from observational data. Unlike traditional correlation-based methods, causal discovery algorithms aim to reveal the underlying structure of causal influences, which is essential to understand complex systems and to make informed decisions [1].

In this study, we performed a comparative evaluation of three widely used causal discovery methods: PC [1], GES [2], and LiNGAM [3]. The LUCAS (LUng CAncer Simple) dataset [4] offers a particularly attractive testbed for this comparison due to its simplicity, binary variable structure, and clearly defined

9

causal relationships. This inherent structure makes LUCAS an ideal candidate for assessing the precision of the reconstruction of causal discovery methods.

Our evaluation focuses on key performance metrics and computational efficiency. We also investigated the impact of sample size variation and the incorporation of domain knowledge (via mandatory and forbidden causal constraints) on the performance of these algorithms. The results of this study provide information on the strengths and limitations of each method, which helps researchers select the appropriate tools for causal inference in similar settings.

Several recent studies have focused on the evaluation of causal discovery methods, providing insight into their performance and limitations. For example, [5] provides a comprehensive overview of current approaches and challenges, while [6] benchmarks these methods on observational data. Furthermore, causal discovery plays an important role in the field of Explainable Artificial Intelligence (XAI) by providing interpretable models that can elucidate the decision-making process of complex systems [7]. It is also a crucial step in Model-Driven XAI for discovering causality, components, and connections (the 3C principle) [8].

The remainder of this paper is organized as follows. Section 2 details the dataset and methods, including the experimental setup and evaluation metrics. Section 3 presents the experimental results and analyzes the impact of sample size and domain knowledge. Finally, Section 4 discusses the findings, practical implications, and directions for future research.

## 2. Methods

In this study, we conducted a comprehensive evaluation of three causal discovery methods in the LUCAS dataset. The LUCAS dataset [4] is a synthetically generated binary dataset designed to model causal relationships among variables related to lung cancer. Due to its simplicity and clearly defined structure, no additional pre-processing was necessary.

We evaluated three causal discovery algorithms: the PC algorithm [1], GES [2], and LiNGAM [3]. These methods represent constraint-based, score-based, and functional approaches, respectively. All algorithms were implemented using the *causal-learn* Python library [9].

Our experimental setup involves several stages. First, each algorithm was applied to the full LUCAS dataset to obtain baseline performance metrics. Second, we evaluated the impact of sample size by conducting experiments on subsets con-

taining 100, 500, and 2000 samples. Third, we incorporated domain knowledge by enforcing mandatory and forbidden causal constraints, following the approaches discussed in recent surveys [5] and benchmarking studies [6]. Finally, we performed hyperparameter tuning to explore potential performance improvements.

We evaluated algorithm performance using several evaluation metrics. These include the Area Under the Precision-Recall Curve (AUPR), Structural Hamming Distance (SHD), Precision, Recall, and F1-Score. In addition, we measured the runtime of each algorithm to evaluate computational efficiency. These metrics enable us to quantify both the accuracy of the inferred causal structures and the efficiency of the algorithms.

The experiments were carried out on a CPU cluster equipped with 12 GB of RAM and 2 vCPUs. Furthermore, the software environment consisted of Python 3.11, the causal-learn library [9], and the necessary dependencies. Each experiment was repeated 10 times, and we calculated the mean run times along with the corresponding standard deviations to ensure robustness. By systematically varying the sample size and integrating domain knowledge, our methodology provides a thorough comparison of the strengths and limitations of each method. This experimental framework is consistent with recent studies on causal discovery [5, 6], facilitating meaningful comparisons with previous work.

## 3. Results

Table 1 summarizes the baseline performance metrics of the three causal discovery methods in the entire LUCAS dataset. As seen in Table 1, the PC algorithm performed the best in achieving the highest AUPR, the lowest SHD, and the highest precision, recall, and F1 score, all while maintaining a low runtime. In contrast, the GES algorithm recorded a very low AUPR and a high SHD despite a perfect recall, indicating that it included many spurious edges. LiNGAM ranked third with moderate performance metrics and an intermediate runtime.

Table 1. Performance metrics of causal discovery methods on the LUCAS dataset

| Algorithm | AUPR | SHD | Precision | Recall | F1-Score | Runtime(s) | Runtime std(s) |
|---|---|---|---|---|---|---|---|
| PC | 0.962 | 1.0 | 0.92 | 1.00 | 0.96 | 0.37 | 0.09 |
| GES | 0.042 | 12.0 | 0.50 | 1.00 | 0.67 | 3.41 | 0.71 |
| LiNGAM | 0.312 | 10.0 | 0.29 | 0.33 | 0.31 | 0.49 | 0.10 |

Our experiments with varying sample sizes (100, 500, and 2000 samples) revealed that the PC algorithm is highly robust, with its performance metrics remaining nearly unchanged. In contrast, both GES and LiNGAM exhibited some sensitivity to the number of samples, with LiNGAM showing notable instability at 500 samples before returning to baseline performance at 2000 samples.

Table 2 presents the performance metrics after incorporating domain knowledge through mandatory and forbidden constraints. When domain knowledge was integrated, the PC and GES algorithms remained largely unaffected. However, LiNGAM benefited from the integration of domain knowledge; its AUPR increased to 0.495 and Recall improved to 0.58, leading to an F1-Score of 0.42, although its SHD increased to 15.0.

Table 2. Performance metrics of causal discovery methods on the LUCAS dataset with domain knowledge

| Algorithm | AUPR | SHD | Precision | Recall | F1-Score | Runtime (s) | Runtime std(s) |
|---|---|---|---|---|---|---|---|
| PC | 0.962 | 1.0 | 0.92 | 1.00 | 0.96 | 0.07 | 0.01 |
| GES | 0.042 | 12.0 | 0.50 | 1.00 | 0.67 | 3.14 | 0.28 |
| LiNGAM | 0.495 | 15.0 | 0.33 | 0.58 | 0.42 | 0.38 | 0.02 |

These results indicate that the PC algorithm consistently outperforms the other methods in terms of reconstruction accuracy and computational efficiency, while GES tends to include many spurious edges, as reflected in its low AUPR and high SHD. LiNGAM, on the other hand, appears more sensitive to sample size and benefits from the incorporation of domain knowledge, suggesting that its performance may be further improved in settings where reliable prior information is available. Our findings are consistent with recent studies in the field [5, 6] and underscore the importance of considering both the characteristics of the dataset and the domain-specific constraints when selecting causal discovery methods.

Figure 1 presents a visual comparison between the true causal graph and the best-generated graph obtained using the PC algorithm. The left part of the image shows the true causal graph, while the right part shows the result obtained with the PC algorithm. In particular, the only discrepancy is that the PC algorithm incorrectly identified the relationship as *Car Accident → Attention Disorder*, whereas the correct causal direction is *Attention Disorder → Car Accident*. This subtle error highlights both the strengths and limitations of the algorithm in reconstructing the underlying causal structure.

(a) True Causal Graph       (b) Best-generated Causal Graph

Figure 1. Comparison of the true causal graph (left) and the best-generated graph using the PC algorithm (right)

## 4. Conclusions

In this study, we performed a comparative evaluation of three causal discovery methods in the LUCAS dataset [4]. Our experiments demonstrated that the PC algorithm consistently achieved superior performance, exhibiting high accuracy and robustness in varying sample sizes. In contrast, the GES algorithm, while achieving perfect Recall, tended to produce many spurious edges, resulting in a low AUPR and high SHD. The LiNGAM algorithm showed moderate baseline performance but benefited significantly from the integration of domain knowledge, as evidenced by improvements in AUPR and Recall, despite an increase in SHD.

These findings are consistent with recent studies in the field [5, 6], highlighting the importance of considering both the characteristics of the dataset and the prior knowledge when selecting causal discovery methods. Our results also underscore the relevance of causal discovery in the broader context of XAI [7, 8], where interpretability is crucial for understanding complex systems.

Despite the strengths of the PC algorithm, our analysis revealed that even the best performing method is not without limitations; for example, subtle errors such as the misidentification of the direction of the relationship between *Car Accident* and *Attention Disorder* were observed. These discrepancies point to potential areas for future research, including the development of hybrid methods that integrate the strengths of multiple approaches and further exploration of methods to effectively incorporate domain knowledge. The source code is available at `https://tinyurl.com/4379c657`.

In summary, our work provides a detailed empirical comparison of causal dis-

covery algorithms and offers practical insights for researchers and practitioners in both the causal inference and XAI communities. Future studies should aim to extend these evaluations to more complex and diverse datasets, as well as investigate the scalability and adaptability of these methods in real-world applications.

# References

[1] Spirtes, P., Glymour, C., and Scheines, R. *Causation, Prediction, and Search*. MIT Press, 2000.

[2] Chickering, D. M. Optimal structure identification with greedy search. In *Proceedings of the 18th Conference on Uncertainty in Artificial Intelligence*, pages 206–215. AUAI Press, 2002.

[3] Shimizu, S., Hoyer, P. O., Hyvärinen, A., and Kerminen, A. A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7:2003–2030, 2006.

[4] Lucas: Lung cancer simple dataset. `https://www.causality.inf.ethz.ch/data/LUCAS.html`. Accessed: 2025-02-06.

[5] Various. A survey on causal discovery: Theory and practice. *arXiv preprint arXiv:2305.10032*, 2023.

[6] Shen, Y. et al. Benchmarking causal discovery methods on observational data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[7] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., García, S., Gil-López, S., Molina, D., Benjamins, R., et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 2019.

[8] Sepioło, D. and Ligęza, A. Towards model-driven explainable artificial intelligence: Function identification with grammatical evolution. *Applied Sciences*, 14(13), 2024. URL `https://www.mdpi.com/2076-3417/14/13/5950`.

[9] Zheng, Y., Huang, B., Chen, W., Ramsey, J., Gong, M., Cai, R., Shimizu, S., Spirtes, P., and Zhang, K. Causal-learn: Causal discovery in python. *Journal of Machine Learning Research*, 25(60):1–8, 2024.

# Stacking Ensemble
# for Face Authenticity Detection

**Rafał Klinowski, Mirosław Kordos**

*University of Bielsko-Biała*
*Faculty of Computer Science and Automatics*
*Willowa 2, 43-309 Bielsko-Biała, Poland*
*rklinowski@student.ubb.edu.pl, mkordos@ubb.edu.pl*

**Abstract.** *A passive face authenticity detection system implemented as a stacking committee of models is presented in this paper. The committee consists of four models, which recognize whether the camera view shows the real face of a live person or whether the face is shown in a photo or video. These models include: a convolutional network and Bezel algorithm for detecting portable devices on which the face can be shown, face context analysis, and a convolutional network for image analysis. The outputs from these models are fed to the inputs of a neural network that makes the final decision.*
**Keywords:** *face anti-spoofing, convolutional neural network*

## 1. Introduction

This work was motivated by the fact that the authors were asked to prepare and implement a face recognition and authentication module in an existing employee registration system. The problem was that some employees occasionally registered the work time of their colleagues who were absent by using their cards with QR code. The solution to prevent such practices was to prepare a machine learning model that would verify whether the employees register in person and nobody else registers their work using their QR code and photos of their faces. Implementing the solution to an existing system imposed several constraints, as the solution had to use the existing hardware without any investment in new hardware, for example, in new cameras. Implementing the module on the server side without any modifications to the software on the client side was required. The existing system

used smartphones as the clients and the front smartphone cameras to take photos while the employees approached their QR code card to the smartphone screen. The images were compressed to about 6 kB webp files, transferred to the server, and stored there. These photos were to be used for face recognition and authentication. As front cameras in smartphones are limited in functionality, without infrared, 3D capabilities, or the possibility of changing depth of field, those technologies could not be used.

First, the face is localized in the photo, then a face recognition module is run. Only the photos presenting a recognized person are further passed to the face spoofing detection module. In this paper, the methods used for the face spoofing detection module are presented. This module recognizes if a face shown to the camera is authentic or not (an authentic face is a face of a live person standing in front of a camera, while an inauthentic (spoofed) face is a face displayed in a photo or video that is displayed on an electronic device screen, printed on a piece of paper, etc.). The methods work on low-quality images taken by cameras without advanced functionalities, such as infrared.

In the initial period of research (up until about 2016) on face authenticity detection, many solutions were proposed, in which manually determined features were examined to identify whether the face presented to the camera was authentic [1, 2]. Recently, methods based on deep learning have become dominant in face authenticity detection, as in most other image recognition tasks. The models include the depth prediction with CNN [3] with an up-sampling block and residual learning, or MobileNet [4]. In the most recent years, ConvNeXt [5] or EfficientNetV2 [6] and Vision Transformers models [7, 8] have been used for face authenticity detection. A separate group of methods uses special cameras or camera sets that allow for a good assessment of image depth and the use of shots from different angles, but these methods are not adequate to the requirements of this system.

## 2. Proposed Method

The proposed method of face authenticity detection includes four different algorithms, each of which assesses a different aspect of the photo and predicts the likelihood of a spoof based on this aspect.

The outputs of these algorithms are fed to the inputs of a neural network, which works as a final classifier in the so-constructed stacking ensemble. If the network output is less than 0.5, the face in the photo is considered authentic; otherwise, it is considered spoof.

Figure 1. Diagram of the proposed method

## 2.1. Bezel Detection

The bezel detection algorithm aims to detect visible bezels (defined as straight, rectangular areas of low brightness) in all directions around the face detected in an input image. A probability of spoofing is calculated based on the number of detected bezels – for each direction around the face that has a bezel, the probability of spoofing is increased by 0.25, which, for instance, gives a 50% probability of spoofing when bezels in two directions are detected.

However, this algorithm can only work against a spoofing attack where the bezels of an electronic device (e.g. smartphone, tablet) are shown, and is ineffective against print attacks or situations where only parts of the device are visible.

## 2.2. Smartphone Detection

This module aims to detect smartphones visible in the camera, regardless of whether they are used for spoofing or not, and aims to increase the accuracy for cases where parts of the device are visible, but not enough to identify a bezel. This is done using a Convolutional Neural Network for whole RGB images scaled down to a size of $128 \times 128$. The network consists of four convolutional blocks and a Softmax output layer to obtain results as probabilities of the image belonging to either class ("smartphone present" and "no smartphone present").

Similarly to bezel detection, while this algorithm can provide meaningful information that can improve the prediction of the ensemble model, the diversity of spoof attacks makes this module useful only against some types of attacks.

17

## 2.3. Context Analysis

Image context analysis performs the analysis of face surroundings in the input image to detect irregularities, such as sharp edges or objects, but most importantly the discrepancy between the area closer to the face and the area further away from it. Image context analysis is performed by an edge detector, in this case the Sobel operator, in two different areas of the image: the close area, which is calculated as 40% of the face's width and height outside of the face, and the far area, which is calculated as another 40% of the face's dimensions.

After performing edge detection, the intensity of edges is calculated, and a histogram of intensities with 64 bins is created. Then, for such histograms, a kNN algorithm is deployed, using a histogram distance metric. Based on the neighbors of the input image's histogram, the probability of spoofing is calculated.

## 2.4. CNN Analysis

A module using a Convolutional Neural Network was developed to analyze images directly and to determine the authenticity of a face based on elements that would otherwise be hard to process with algorithms, such as image quality and cohesion of the image. The proposed network architecture uses two sets of convolutional layers with a Softmax classifier to obtain a probability value.

While developing the module, several setups were tested. Experimental evaluation showed that better results were obtained for images that were processed by centering around the detected face instead of passing whole images to the input layer. Further experiments need to be performed to test different parameters of this module, such as input image size and learning rate.

# 3. Experimental Evaluation

Experimental evaluation was conducted using Python, PyTorch, and OpenCV. The source code used in the experiments is available at *kordos.com/face.html*.

The system evaluation is presented on three datasets:

- An open dataset with N=1666 images (779 authentic face images and 887 spoofed faces images) from the Human Faces [9] dataset as well as the Face Anti-Spoofing dataset [10];

- A dataset with N=1678 images (839 authentic face images and 839 spoofed face images) gathered by us that represent the conditions in which the system operates;

- A smaller, combined dataset with N=477 images (243 authentic face images and 234 spoofed face images), where some were gathered by us and some were taken from the aforementioned datasets.

As the numbers of photos of authentic faces and the number of photos of spoofed faces are practically the same in two datasets and differ only by about 13% in one dataset, accuracy is an appropriate assessment measure of the models trained and tested on the datasets. The system was evaluated using a five-fold cross-validation with a random seed of 42. The results of individual modules, as well as the committee, were collected. These results were compared to some other solutions used in the literature for face spoofing detection, which were trained and tested on the same datasets:

- Depth Prediction with CNN [3] using an open-source implementation of by Anand Pandey [11];

- MobileNet [4] using the implementation by Nguyen Dinh Quy [12];

- The ConvNeXt architecture [5] without pre-trained weights for the "BASE" model, using PyTorch implementation;

- The EfficientNetV2 architecture [6] without pre-trained weights for the "S" model size, using PyTorch implementation.

Table 1. Prediction accuracy of individual algorithms and the ensemble across the three datasets

|  | Small dataset | Open dataset | Our dataset |
|---|---|---|---|
| Bezel detection | 70.67% ± 4.26% | 50.07% ± 2.82% | 83.89% ± 2.01% |
| Phone detection | 77.15% ± 6.22% | 69.21% ± 11.63% | 89.57% ± 5.78% |
| Context analysis | 79.07% ± 4.64% | 94.47% ± 1.32% | 91.06% ± 1.85% |
| CNN analysis | 87.20% ± 3.19% | 91.95% ± 1.75% | 96.23% ± 1.71% |
| **Our method** | **88.50% ± 0.84%** | **96.79% ± 1.17%** | **97.98% ± 0.59%** |

The statistical significance of the obtained results was tested using a one-way ANOVA test with a significance threshold of $\alpha = 0.05$. For each system,

Table 2. Comparison of prediction accuracy of the proposed method and other evaluated systems

|  | **Small dataset** | **Open dataset** | **Our dataset** |
|---|---|---|---|
| Depth prediction | 53.25% ± 3.33% | 74.12% ± 15.63% | 75.69% ± 3.09% |
| MobileNet | 71.73% ± 2.33% | 83.79% ± 0.54% | 85.19% ± 1.19% |
| ConvNeXt | 76.70% ± 10.60% | 70.66% ± 11.8% | 87.42% ± 1.47% |
| EfficientNetV2 | 85.32% ± 2.12% | 80.27% ± 13.42% | 87.42% ± 1.47% |
| **Our method** | **88.50% ± 0.84%** | **96.79% ± 1.17%** | **97.98% ± 0.59%** |

the results for all datasets were collected, amounting to 15 records (3 datasets in 5-fold cross-validation). The ANOVA test conducted this way yielded a p-value of $2.09 \times 10^{-9}$, which is significantly smaller than the accepted significance threshold. Therefore, the differences in the results are statistically significant, which proves that the proposed method outperforms the reference solutions.

# 4. Conclusions

A method of detecting the authenticity of faces in the work time registration system was presented and evaluated. A set of restrictions was imposed on the method and on the photos it can use. A stacked ensemble constructed with two CNNs, a bezel detection method, and a context analysis method allowed to obtain quite high prediction results of up to 98%, depending on the dataset.

Further work on this topic will include the parts of the system that only work against specific types of spoofing attacks, like the bezel detection and smartphone detection modules. As described, a starting point could be using image gradients that are analyzed for sudden changes, such as device bezels. Furthermore, different neural network architectures and parameters could be tested as part of the system, including known CNN architectures or the vision transformer.

# Acknowledgment

# References

[1] de Freitas Pereira, T. et al. LBP‑TOP based countermeasure against face spoofing attacks. In *Computer Vision - ACCV 2012*, pages 121–132. 2013.

[2] Komulainen, J., Hadid, A., and Pietikäinen, M. Context based face anti-spoofing. In *IEEE Sixth International Conference on Biometrics*, pages 1–8. 2013. doi:10.1109/BTAS.2013.6712690.

[3] Ma, X., Geng, Z., and Bie, Z. Depth estimation from single image using CNN-residual network. *Semantic Scholar*, 2017. URL `https://api.semanticscholar.org/CorpusID:11853184`.

[4] Xiao, J., Wang, W., Zhang, L., and Liu, H. A mobilefacenet-based face anti-spoofing algorithm for low-quality images. *Electronics*, 13(14), 2024.

[5] Liu, Z., Mao, H., Wu, C., Feichtenhofer, C., Darrell, T., and Xie, S. A convnet for the 2020s. *arXiv preprint arXiv:2201.03545*, 2022.

[6] Tan, M. and Le, Q. V. EfficientNetV2: Smaller models and faster training. *arXiv preprint arXiv:2104.00298*, 2021.

[7] Liu, A. and Liang, Y. MA-ViT: Modality-agnostic vision transformers for face anti-spoofing. *arXiv preprint arXiv:2304.07549*, 2023.

[8] Yu, Z. et al. Rethinking vision transformer and masked autoencoder in multimodal face anti-spoofing. *arXiv preprint arXiv:2302.05744*, 2023.

[9] Gupta, A. Human faces. `https://www.kaggle.com/datasets/ashwingupta3012/human-faces`.

[10] Kolzek. Face anti-spoofing dataset, 2023. URL `https://universe.roboflow.com/kolzek/face-anti-spoofing-ezpqo`. Visited on 2025-02-11.

[11] Pandey, A. Face liveness detection using depth map prediction. `https://github.com/anand498/Face-Liveness-Detection`.

[12] Quy, N. D. Face anti-spoofing using mobilenet. `https://github.com/dinhquy94/face-antispoofing-using-mobileNet`.

# Grammar Refinement in Grammatical Evolution Using Large Language Models

**Dominik Sepioło**[0000−0001−7746−3781], **Mikołaj Jarosławski**[0009−0003−1999−4949], **Antoni Ligęza**[0000−0002−6573−4246]

*AGH University of Krakow*
*Department of Applied Computer Science*
*al. Mickiewicza 30, 30-059 Kraków, Poland*
*{sepiolo, mikjar, ligeza}@agh.edu.pl*

**Abstract.** *Grammatical Evolution is an evolutionary algorithm that uses context-free grammars to generate solutions to complex problems. The effectiveness of GE largely depends on the quality of the grammar used, yet manual grammar design is often inefficient and requires domain expertise. This paper presents a framework for grammar refinement in GE using Large Language Models. By exploiting the ability of LLMs to process and generate coherent text, the framework iteratively refines grammars based on problem-specific requirements and feedback from the evolutionary process. The results highlight the potential of integrating LLMs into evolutionary computation workflows to improve adaptability and efficiency.*
**Keywords:** *grammatical evolution, symbolic regression, large language models, LLM*

## 1. Introduction

Grammatical Evolution (GE) is a popular variant of Genetic Programming (GP) that generates solutions using grammars to encode constraints and structure. While GE has been successfully applied in areas such as symbolic regression, automated programming, and creative generation, its performance is highly sensitive to the quality of the grammar used. Poorly designed grammars can lead to unproductive search spaces, suboptimal solutions, or premature convergence. Traditionally, grammar design relies on domain expertise and manual trial-and-error, which limits scalability and generalizability.

Recent advancements in artificial intelligence have led to the development of Large Language Models (LLMs) such as GPT and BERT, which have shown excellence in understanding and generating natural language. These models have shown potential for the automation of complex tasks such as grammar generation and refinement. Moreover, LLMs have demonstrated remarkable capabilities in mathematical reasoning and symbolic function discovery [1].

This paper explores the integration of LLMs into the grammar refinement process for GE. Using LLMs to analyze evolutionary progress and suggest grammar modifications, we aim to overcome the limitations of manual grammar design and improve the overall efficiency of GE.

## 2. Theoretical Background

Grammatical Evolution is an evolutionary algorithm that generates programs or expressions by evolving a population of candidate solutions. Unlike traditional Genetic Programming (GP), which directly manipulates syntax trees, GE uses a grammar, typically expressed in Backus-Naur Form (BNF), to guide the evolution process. This separation of genotype and phenotype allows for modular and flexible search strategies. GE operates by mapping the genotype to the phenotype using a predefined grammar, which ensures that the generated solutions are syntactically correct [2]. GE has been applied to various domains, including symbolic regression, automatic programming, and design optimization [3].

Symbolic Regression (SR) is a mathematical modeling technique that seeks to uncover analytical expressions that best describe observed data without assuming an a priori functional form. Traditional SR approaches, often implemented through GP, explore the space of mathematical expressions by evolving populations of candidate functions. The initial population consists of randomly generated expressions. The fitness of each expression is evaluated on the basis of how well it fits the data. Evolutionary operators such as crossover and mutation are applied to the expressions to create new generations of solutions. The process continues until an optimal or satisfactory expression is found [4].

Large Language Models have significantly advanced natural language processing by demonstrating the ability to generate and manipulate structured text, including mathematical expressions and code [5]. Unlike earlier deep learning approaches, LLMs utilize self-attention mechanisms to capture long-range dependencies, enabling them to model syntax and semantics effectively. Trained on

extensive collections of text and code, they demonstrate a strong generalization in generating and refining equations, symbolic expressions, and structured programs.

# 3. Methodology and Tools

This section describes the methodology used to integrate LLMs into GE for grammar refinement task. In addition, we outline the tools and technologies used to implement this framework.

## 3.1. Integrating LLMs for Grammar Refinement

Grammar refinement in GE traditionally relies on expert-driven design, requiring manual adjustments to improve search performance. The proposed method incorporates LLMs as an assistive mechanism for analyzing and refining grammars, providing recommendations that support human decision-making while enhancing solution quality.

The integration follows an iterative process:

1. **Grammar Initialization** – A baseline context-free grammar defines the initial syntactic constraints for GE-generated solutions.
2. **Evolutionary Execution** – The GE algorithm evolves a population of candidate solutions by applying genetic operators (selection, crossover and mutation), with fitness evaluation guiding the optimization process.
3. **Grammar Evaluation, Modification and Reintegration via LLMs** – LLMs analyze current grammar, suggest refinements and propose grammar update to improve syntactic efficiency.
4. **Iterative Refinement** – The process repeats over multiple evolutionary cycles, allowing the grammar to adapt dynamically to emerging patterns.

LLM integration addresses the limitations of static grammar design by enabling adaptive modifications. This ensures that the search space evolves in response to patterns that emerge during the evolutionary process.

## 3.2. Implementation Details

The proposed approach is implemented in the R language, using the gramEvol library, which provides a framework for evolutionary computation with user--customizable grammars [6]. The evolutionary process applies standard genetic

operations, including selection, crossover, and mutation, to iteratively evolve candidate solutions.

To integrate LLMs into the grammar refinement process, interactions between the evolutionary framework and GPT-4 were managed through chat. In this interaction, the grammar was iteratively refined based on user prompts, including the current grammar definition, the resulting expression, and (in some cases) the logarithm of mean absolute error (LMAE $= \frac{1}{n}\Sigma_{i=1}^{n} \log(1 + |y_i - x_i|)$), along with modifications suggested by GPT-4. The best-developed solutions were analyzed and the grammar modifications suggested by the LLM were incorporated into the gramEvol pipeline.

# 4. Experiment and Results

To evaluate the effectiveness of integrating Large Language Models for grammar refinement in Grammatical Evolution, we conducted a series of experiments comparing traditional GE with LLM-enhanced approaches.

In [7, 8], we demonstrated that GE can efficiently perform simple symbolic regression tasks, such as BMI formula discovery, by accurately identifying functional dependencies using a minimal dataset. Here, efficiency refers to the computational cost, including the number of evolutionary process iterations and computation time, and the ability to achieve meaningful results on small datasets. A more advanced experiment was reported in [3], where GE was applied to more challenging function identification tasks (e.g., discriminant of a quadratic equation) within the framework of *Model-Driven Explainable Artificial Intelligence* (MD-XAI). The study demonstrated that GE, when combined with domain knowledge and structured grammatical constraints, effectively discovers functional dependencies in the data, generating transparent and interpretable models.

Building on these results, this study explores how integrating LLMs further enhances GE by dynamically refining grammars, reducing the need for expert intervention, and improving adaptability in symbolic regression tasks. We decided to continue our experiments with formula discovery for the discriminant of a quadratic equation ($\Delta = b^2 - 4ac$). The experiment was carried out on a dataset of 100 observations, each consisting of values for *a*, *b* and *c*, along with the corresponding value of the discriminant $\Delta$).

We initially employed a grammar that defined mathematical expressions using variables (*a*, *b*, *c*), numerical constants ranging from −5 to 5, arithmetic operators

and common mathematical functions. The grammar is listed below.

```
<expr> ::= <var> | <n> | <op> (<var>, <var>)
          | <func> (<var>) | <n> * <func> (<var>)
<op>   ::= + | - | * | / | ^
<func> ::= sin | cos | log | exp | sqrt | abs
<var>  ::= a | b | c
<n>    ::= -5 | -4 | -3 | -2 | -1 | 0 | 1 | 2 | 3
          | 4 | 5 | pi
```

However, despite the expressiveness of this grammar, the results of GE converged to 0 or identity operations such as $a - a$, failing to capture meaningful patterns in the data. In the next steps, LLM proposed grammars that excluded trivial solutions (constant 0) and aimed to reduce identity operations. As the generated grammars became increasingly complex, the computational cost grew significantly.

To address this challenge, we provided feedback to the LLM model to guide the dynamic refinement of the grammar. We started with the previously listed grammar and asked the LLM to make it less computationally complex while still maintaining its ability to capture meaningful patterns in the data. The goal was to simplify expressions and limit the introduction of excessively complex operations, making the search process in GE more efficient.

By analyzing the evolved expressions and their associated errors, the model was able to detect patterns of excessive complexity, leading to refinements such as the removal of unnecessary functions, such as trigonometric operations. LLM also contributed to generating mathematically robust grammars (e.g., preventing division by zero). We began the grammar refinement process with the following prompt:

```
Context:
I am optimizing a GE grammar for symbolic regression
in R using gramEvol. My goal is to make the grammar
more efficient while maintaining its ability
to discover meaningful patterns.

Current Grammar:
[as presented above]
```

```
Objective:
I need to simplify expressions and reduce complexity.
- Avoid division by 0.
- Remove functions that are rarely useful for regression.

Expected Output:
Please refine the grammar to align with these objectives,
explain the rationale behind your changes. You may also
suggest changes for next iteration of grammar optimization.
```

We interacted with GPT 24 times, iteratively refining and optimizing the grammar. Throughout this process, we improved the prompt by including details of the resulting expression and the error. However, obtaining the correct solution required external knowledge of the problem solution structure. Specifically, human intervention was necessary to allow exponential operations during interaction with the LLM, as their absence constrained the model ability to represent the underlying mathematical relationships. The final grammar, obtained through the iterative refinement process, is listed below.

```
<expr>  ::=  <var>  |  <n>  |  <op>  (<expr>, <expr>)
            |  <n> * <var>  |  <var> ^ <n>
<op>    ::=  +  |  -  |  *
<var>   ::=  a  |  b  |  c
<n>     ::=  1  |  2  |  3  |  4  |  5
```

After completion of the GE process, we obtained the solution $b^2 + 4 * c * (1 * a - 2 * a)$ which aligns with the correct formula for $\Delta$.

## 5. Conclusion and Future Work

This study presents a novel approach to symbolic regression by integrating LLMs into the GE process for dynamic grammar refinement. The proposed method improves the efficiency of symbolic regression by guiding the evolutionary search with LLM-generated refinements. The experimental results demonstrate the accuracy of the solutions generated using the refined grammar. LLM assistance accelerates the grammar definition process and ensures that the generated expressions are mathematically valid. However, human intervention was required to receive the final grammar. Future work will focus on optimizing and automating the refinement process and extending the approach to more complex discovery tasks.

# References

[1] Shojaee, P., Meidani, K., Gupta, S., Farimani, A. B., and Reddy, C. K. LLM-SR: Scientific equation discovery via programming with Large Language Models. *arXiv preprint arXiv:2404.18400*, 2024. URL `https://api.semanticscholar.org/CorpusID:269449436`.

[2] Ryan, C., O'Neill, M., and Collins, J. J., editors. *Handbook of Grammatical Evolution*. Springer, 2018. ISBN 978-3-319-78716-9. doi:10.1007/978-3-319-78717-6.

[3] Sepioło, D. and Ligęza, A. Towards model-driven explainable artificial intelligence: Function identification with grammatical evolution. *Applied Sciences*, 14(13):5950, 2024. ISSN 2076-3417. doi:10.3390/app14135950.

[4] Makke, N. and Chawla, S. Interpretable scientific discovery with symbolic regression: a review. *Artificial Intelligence Review*, 57(1), 2024. ISSN 1573-7462.

[5] Brown, T. B. et al. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20. Curran Associates Inc., 2020. ISBN 9781713829546.

[6] Noorian, F., de Silva, A. M., and Leong, P. H. W. gramEvol: Grammatical evolution in R. *Journal of Statistical Software*, 71(1):1–26, 2016. doi:10.18637/jss.v071.i01.

[7] Sepioło, D. and Ligęza, A. A comparison of shallow explainable artificial intelligence methods against grammatical evolution approach. In *Progress in Polish Artificial Intelligence Research 4*, pages 89–94. Lodz University of Technology Press, 2023.

[8] Sepioło, D. and Ligęza, A. Towards model-driven explainable artificial intelligence. an experiment with shallow methods versus grammatical evolution. In *Artificial Intelligence. ECAI 2023 International Workshops*, pages 360–365. Springer Nature Switzerland, Cham, 2024.

# Text Data Mining: A Case Study of Reddit

**Piotr Sokołowski**[1[0009−0005−8181−6275]],
**Marcin Kosiba**[2[0009−0005−7897−9743]],
**Marcin Szpyrka**[3[0000−0003−4925−3271]]

[1,2,3]*AGH University of Krakow*
*Department of Applied Computer Science*
*Al. Mickiewicza 30, 30-059 Kraków, Poland*
[1,3]*{psokolowski,mszpyrka}@agh.edu.pl,* [2]*mkosiba@student.agh.edu.pl*

**Abstract.** *The article describes the progress of a project on text data mining performed by software engineering students. The article presents information about the source of the text data, how it was acquired and cleaned. The process of selecting classifiers to categorize posts is described, and the results of different models are compared.*
**Keywords:** *data mining, reddit, classification, clustering, scikit-learn, tensorflow, machine learning*

## 1. Introduction

Reddit, a social networking platform, functions as an aggregation of forums where users contribute questions, reflections, and content, with the ability to endorse or disapprove of these contributions. The most prominent and frequently visited posts are showcased on the front page, which serves as the initial interface for users upon accessing the platform. Reddit is further segmented into subreddits [1], which are forums dedicated to specific topics.

In September 2024, Reddit had more than 101 million active daily users in more than 100,000 communities [2]. The platform also contains a vast repository of content, including more than 16 billion posts and comments, which have been uploaded since June 2005, the inception of the service.

The purpose of this study was to check whether Reddit posts can be treated as a suitable dataset for data mining and machine learning, analogous to the use of Twitter posts in current research.

## 2. State of the Art and Related Work

The analysis of social media posts has become a recurring theme in academic research. Users' posts frequently offer information about global events [3], along with their sentiments and viewpoints [4]. A significant number of academic papers focuses on the analysis of Twitter posts, although Reddit has also emerged as a prominent source of data.

On 3 July 2015, Jason Baumgartner, operating under the username "u/Stuck_-In_the_Matrix," published a post on the "r/datasets" subreddit. In this post, he stated that he had gathered all publicly available comments on Reddit for the purpose of research [5]. This collection, which spanned approximately 250 gigabytes (GiB), encompassed a total of 1.7 billion posts and comments that had been published on the platform from October 2007 to May 2015. The posts were methodically organized by month and year and subsequently saved in JSONL files, a file format that allows the storage of multiple JSON objects in a single file, with each line of the file corresponding to one JSON document. Subsequently, with the help of the community, the collection was made available as a torrent.

## 3. Hardware Limitations and the Challenge of Data Acquisition

The dataset, upon further examination, revealed a substantial requirement for disk space, occupying 149.6 GiB of storage capacity. Subsequent to the preparation of the server, the decompression and data loading process was initiated into the database.

The presence of inconsistent data structures within the JSONL files resulted in various errors, including the absence of keys in certain files and the presence of additional values in others. Due to the substantial size of the files that contain errors and the time-consuming nature of their analysis, we ultimately decided to discard the 2015 data. The final database contained more than 1.4 billion records, occupying a total of 562.8 GiB of disk space.

## 4. Challenges in Data Cleaning

The primary challenge was twofold: to reduce the amount of data and to clean it. Several techniques, described later in this section, were used to ensure that only

high-quality entries were used for analysis. It was decided to retain only entries of more than 1,000 characters to preserve meaningful information after the cleaning procedure (Figure 1). Of the entries, 4,000 from each of the five subreddits were selected for further analysis.



Figure 1. Entries with more than 1,000 characters make up a distinct minority, accounting for less than 2% of the total dataset

In the course of the data cleaning process, the HTML entities, URLs, and stop-words were removed from the entries. The remaining words were then lemmatized to maintain the uniformity of the entries' format. This lemmatization process also enabled the removal of non-words, defined as strings of characters that are not present in the WordNet dictionary [6]. The resulting data were subsequently prepared for sentiment analysis and subject classification.

# 5. Using Deep Machine Learning for Accuracy Improvement

A comprehensive evaluation of the classifiers available in the Scikit-learn library [7] was performed to identify the most suitable option for the given problem, described in Section 1. To ensure the reliability of the results, an analysis was performed for both Bag of Words and *tf-idf* (Figure 2). The analysis identified Logistic Regression with cross-validation and the SAG solver using Bag of Words as the optimal classifier, achieving an accuracy of 75.18. Furthermore, the weighted

average of precision, recall, and the F1 score reached 0.75. The highest F1 score (0.83) was observed for the gaming category, while the technology category had the lowest (0.64). A similar trend was observed with logistic regression using the Newton-CG solver in *tf-idf*, which emerged as the most effective classifier, attaining an accuracy of 77.61 and a weighted F1 score of 0.77.



Figure 2. The receiver operating characteristic (ROC) curve for both the Bag of Words (left) and *tf-idf* (right) models indicates an inability to effectively classify the texts of the "technology" subreddit

To enhance performance, a Long- and Short-Term Memory (LSTM) model with attention was implemented (Figure 3), utilizing Word2Vec embeddings trained on the Google-News-300 dataset. The model was trained for 30 epochs using sparse categorical crossentropy as a loss function. This approach achieved an accuracy of 76.71 with a loss of 0.69. The overall F1 score reached 0.76, with the highest score for gaming (0.83) and the lowest for technology (0.64). These results closely align with previous models, suggesting that the accuracy ceiling may be influenced by the dataset itself, rather than by the choice of the model.

# 6. Conclusion

An analysis of Reddit data is a complex process that involves multiple stages, including data acquisition, data cleaning, and data analysis. The data set, which

Figure 3. The receiver operating characteristic (ROC) curve for the long short-term memory (LSTM) model with attention

contains more than 1.4 billion records, presents significant challenges in terms of storage capacity, data loading, and data cleaning. Despite these challenges, the insights gained from the analysis of Reddit data can provide valuable information about user behavior, sentiment, and subject matter. Using the vast amount of data available on Reddit, researchers can gain a deeper understanding of online communities and the dynamics that drive them [3].

The substandard alignment of entries from the category "technology" is attributable to the generalization of this category, which encompasses a diverse range of topics, including programming, science, and games.

The precision of the model is also significantly impacted by entries with limited content. In the dataset used, there were many posts that mostly contained ASCII art or web addresses. It is crucial to acknowledge that our data encompass not only entries, but also comments that, following a thorough cleansing process, contain insufficient information to definitively allocate them to a specific category.

In this paper, we have outlined the process of acquiring, cleaning, and analyzing Reddit data, highlighting the challenges and solutions encountered at each stage. By sharing our experiences and methodologies, we hope to provide a valuable resource to researchers interested in analyzing social media data. The insights

gained from this analysis can be used to inform a wide range of research topics, including sentiment analysis, topic modeling, and user behavior analysis.

## Acknowledgment

## References

[1] Medvedev, A. N., Lambiotte, R., and Delvenne, J.-C. The anatomy of Reddit: An overview of academic research. In *Dynamics on and of Complex Networks III*, page 183–204. Springer International Publishing, 2019.

[2] Reddit Inc. Reddit by the numbers. URL `https://redditinc.com/press`. [Accessed 14-02-2025].

[3] Leavitt, A. and Clark, J. A. Upvoting hurricane Sandy: Event-based news production processes on a social news site. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1495–1504. 2014.

[4] Long, S., Lucey, B., Xie, Y., and Yarovaya, L. "I just like the stock": The role of Reddit sentiment in the GameStop share rally. *Financial Review*, 58(1):19–37, 2023.

[5] Baumgartner, J. I have every publicly available Reddit comment for research. URL `https://reddit.com/r/datasets/comments/3bxlg7/`. [Accessed 14-02-2025].

[6] Miller, G. A. WordNet: a lexical database for English. *Communications of the ACM*, 38(11):39–41, 1995.

[7] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

# Synthetic Trips as a Tool to Analyze the Probability of Traveling by Bicycle

**Przemysław Wrona**[0000−0002−5479−4489], **Maciej Grzenda**[0000−0002−5440−4954]

*Warsaw University of Technology*
*Faculty of Mathematics and Information Science*
*Koszykowa 75, 00-662 Warsaw, Poland*
*{przemyslaw.wrona.dokt,maciej.grzenda}@pw.edu.pl*

**Abstract.** *Multiple factors affect the selection of a travel mode. Machine learning models are effective in elucidating and comprehending these decisions. However, can we presume that the probability of choosing a mode of transport is uniform throughout an entire city? This study proposes generating synthetic instances and then using classifiers trained with real instances to assign labels, estimate probabilities of choosing a bicycle to travel and provide a spatial analysis of travel patterns. Our approach facilitates an understanding of how travel mode choices vary geographically by visualizing travel mode probabilities in a spatial context.*
**Keywords:** *synthetic trips, classification, travel mode choices*

## 1. Introduction

Modeling the travel patterns of city inhabitants is essential for planning transportation policies. At the same time, an increasing amount of work does not focus solely on the accuracy of prediction models. Frequently, emphasis is placed also on the explanation of the model that makes a decision. Models such as decision trees and logistic regression provide information about the influence each factor has on decision-making and offer metrics to compare and rank these influences. However, can we say that trip features are uniformly distributed across the urban area of interest? In this work, our aim is to present a novel perspective on a Travel Mode Choices (TMC) [1, 2] model – as not merely a model used to predict one class (travel mode), but rather as a spatial distribution of probabilities that demonstrates how these probabilities vary in different locations within the city. We focus on the use of bicycles as an example of an active travel mode.

## 2. Related Works

Understanding the composition factors influencing travel mode choice decisions made by urban inhabitants is vital for effective traffic planning, management, and flow prediction. Different factors influence TMC decisions, e.g. travel distance or owning a car [3, 4]. To predict TMC given the data of traveler, household and trip, TMC models are developed [1, 5]. Typically, this is done with real trip instances, e.g. collected from respondents describing their trips and travel modes they have selected for those trips, such as driving, public transport (PT), cycling and walking [2]. Hagenauer et al. compared different models for predicting travel mode choices [2]. The highest accuracy was achieved by Random Forest. Interestingly, the number of real trip records is typically limited.

As shown by Hillel et al. in [1], in ML-based TMC studies, data sets composed of 1,000–10,000 instances were used most frequently. Importantly, to predict a mode choice for a trip, attributes of different modes, such as estimated travel time for the trip, should be calculated [1]. In [4], an approach that consists of a data generation framework, which combines multiple data sources to extend real trip instances with features describing the estimated level of service of different travel modes was proposed. Grzenda et al. in [6] proposed a system that extends raw trip data with features describing, inter alia, travel time, travel distance, and the number of transfers for different modes of transport. To achieve this goal, the use of the Open Trip Planner (OTP) was proposed. OTP can operate in different modes and can suggest connections using public transport, routes with a car, bike, including bike stations in the city, or when walking.

Luckner et al. in [5] proposed an algorithm used to generate synthetic trips based on the trip matrices (TMs) from a transport model. TM is a two-dimensional matrix where rows and columns represent the origin and destination zones of the trip, respectively, whereas the values present in TMs represent the estimated number of trips for each pair of the zones [7]. This relies on the zones into which the area – e.g., the city – has been divided in the transport mode [7]. Based on synthetic trips, generated between random origin and destination points in the urban area of interest while considering TMs, the study analyzed the spatial distribution of different trip features [5]. These included the number of available PT connections, the difference in travel time between different travel modes, and the minimal number of transfers required to reach the destination address. The level of service of PT in the city was found to be substantially different in different city areas. Probabilities of travel mode use were not analyzed. Finally, let us note that active

travel modes, such as walking and cycling, attract particular attention because of health and environmental benefits. The use of these modes is affected by such trip features as distance, traveler's features, but also habits, preferences and e.g., weather conditions.

In our study we focus on estimating long-term trends in bicycle use at the city scale. Hence, we constrain the feature set to the features easily available for individual population members, that is, we refrain from using preferences and habits, and we avoid using short-term features such as weather forecasts, as these would provide day-to-day predictions rather than help estimate the overall average share of modes in different parts of the city.

## 3. Methods

The goal of this work is to propose a method of estimating the spatial probabilities of selecting travel mode, and to identify locations in the city with the highest probability of using the selected mode. Typically, there are not enough real trip instances $\mathcal{T}_R$ to generalize the results for the entire city. Hence, we aim to exploit both a limited number of real trip records $card(\mathcal{T}_R)$ with known travel modes selected by travelers and the data on the frequency of trips between different areas of the city of interest. Hence, as proposed in [5], first we generate synthetic trips $\mathcal{T}_S$ based on the trip matrices. By utilizing trip matrices, we generate much more extensive synthetic trip data, i.e. $card(\mathcal{T}_S) \gg card(\mathcal{T}_R)$, which reflects the statistical distribution of actual flows of travelers at different times of the day between different city zones, meaning that real mobility is accurately represented and simulated.

Then, based on [6], we extend both real trip records $\mathcal{T}_R$ and synthetic trip records $\mathcal{T}_S$ with the level of service (LOS) features [1] that include, inter alia, estimated travel time, travel distance, and the number of transfers for various modes of transport, such as car, public transport, cycling and walking. These mode-dependent features are calculated assuming a given travel mode is used for a trip defined by its origin and destination coordinates and time of day. This is done by submitting requests to OpenTripPlanner (OTP) instances configured with real schedules of the city of interest and provided with the spatial data of the city. Based on the responses to such requests, candidate routes, PT connections and their features are estimated.

Next, we construct a classification model $\mathcal{M}$ using true labeled instances $\mathcal{T}_R$

and subsequently predict travel modes with the model $\mathcal{M}$ for synthetic instances $\mathcal{T}_S$. To build the model $\mathcal{M}$, we used among others a decision tree and its implementation provided in the `rpart` package. The dataset of instances $\mathcal{T}_R$ was randomly divided into training and testing datasets.

Both $\mathcal{T}_R$ and $\mathcal{T}_S$ instances are constrained to the level of service attributes calculated for different modes with the OTP services. It is important to note that the model $\mathcal{M}$ should estimate the probability that an instance $x \in \mathcal{T}_S$ belongs to class $j$, i.e. that $j$−th travel mode is used for a trip given its features such as distance to destination. Finally, the spatial distribution of the estimated probabilities of bicycle use for urban trips can be analyzed. In this way, city areas with high bicycle use and e.g., an increased demand for bicycle parking can be identified.

## 4. Results

Our aim is to develop a model that will perform the classification task and estimate the probability of using different modes of transport for individual trips. To validate our approach, we use the data from City of Warsaw. In our research, we used $card(\mathcal{T}_R)$=5,729 real trip instances to build classification models. The training and testing datasets consist of 4,702 (80% of $\mathcal{T}_R$) and 1,177 (20% of $\mathcal{T}_R$) instances, respectively. To illustrate our approach with an interpretable TMC model, a decision tree models have been trained. Figure 1a illustrates a sample model $\mathcal{M}$ built on the attributes derived from OTP for instances $\mathcal{T}_R$. To optimize the hyperparameters and prevent overfitting by systematically evaluating performance across different data splits, we used the `caret` package with 5-fold cross-validation. The Decision Tree model was optimized based on grid search performed with `caret` achieves an accuracy of 55.88% ($\kappa$=0.2961). In the case of another model, the accuracy achieved is 64.75% ($\kappa$=0.4137) for Random Forest.

The analysis involved generating $card(\mathcal{T}_S)$=100,000 instances. In our case, we calculated attributes for public transport, cycling, driving, and walking, as other transport modes were only occasionally used, which follows from the real trip data. In the same way as in the case of $\mathcal{T}_R$, we calculated attributes of $\mathcal{T}_S$ based on the OTP engine. Next, we used the decision tree to estimate the probability of bike use for individual trips in $\mathcal{T}_S$. Figure 1b presents the probabilities of selecting a bike as the main travel mode of travel, i.e., possibly accompanied by walking when necessary, estimated with model $\mathcal{M}$.

Decision tree shows that if the distance is less than or equal to 1,390m, then

(a) Sample decision tree predicting Travel Mode Choice (TMC)

(b) The probability of starting a journey by bike at a given location

Figure 1. Sample decision tree (a) and spatial distribution of the probability of choosing a bike at a given location in the City of Warsaw area (b)

inhabitants prefer walking (0.66 estimated probability), while the estimated probability for cycling is 0.05. For longer distances, if the total distance from the starting point to the PT stop, between stops, and from the final PT stop to the destination is less than or equal to 1007m, then public transport is preferable; otherwise, the probability of choosing a car as the main mode increases. The last two crucial factors are travel time by car *Duration_CAR* and car velocity *Speed_CAR*. For long distances, if travel time by car is greater than or equal to 689 seconds (11 minutes 29 seconds), then a car is preferable; if car velocity is less than 5.2 m/s (18.72 km/h), then public transport will be chosen. It is worth noting that none of the nodes in the decision tree indicates a bike as the preferred transport mode; however, we can still track the distribution of bike instances.

Figure 1b illustrates the spatial probabilities of choosing a bike as the preferred travel mode. Figures show that in the south, southeast, and north districts, trips frequently occur with a negligible chance of choosing a bike. This may be because typical trips from these districts are longer, and the residents have to travel to regions closer to the city centre. Figure 2 illustrates the spatial distribution of trips based on distance range and represents the number of initiated trips in a given location, using a grid size of 500m × 500m. The histograms show the number of started trips divided into four different trip length categories. For example, Figure 2a shows the distribution of trip starting points for the shortest 25% of

trips [0.0 km–4.1 km], while Figure 2d depicts the starting points for the longest 25% of trips [11.3 km–31.5 km]. It is clear that residents of only some city areas frequently travel to nearby areas, which stimulates the use of bicycles, though too short distances increase the probability of walking.



| (a) [0.0km–4.1km] | (b) (4.1km–7.3km) | (c) (7.3km–11.3km) | (d) (11.3km–31.5km) |

Figure 2. A spatial histogram of the number of trips starting in individual areas where the journey distance falls within the range [0.0km–4.1km] – 25% of the shortest trips (a), (4.1km–7.3km) (b), (7.3km–11.3km) (c), and (11.3km–31.5km) – 25% of the longest trips (d)

## 5. Conclusions

In this work, we show that synthetic trips can be used to analyze the spatial distribution of travel mode preferences when a limited number of real trip instances are available. $\mathcal{T}_R$ can be used to build a model and then use the model to estimate probabilities for synthetic trips. It follows from the application of the method to Warsaw city data that the probability of using a bike is not constant throughout the city and depends, inter alia, on the distance to the destination measured along the walking route, but also other factors. The value of probability for Warsaw is lower than 0.10, and for most instances typically even substantially lower when we analyze external districts of the city. This is because cycling probability partly depends on the spatial distribution of the distances involved in a typical trip.

It should be noted that if we analyze 25% of the longest trips (Figure 2d), there is a significant increase in the number of trips in the peripheral districts. The shortest trips take place in the city's two business centres. The first is located in the city centre, the second near the airport – in the Służewiec area. In the future, other than decision tree models, including models with lower interpretability but potentially improved class probability estimates will be considered.

## Acknowledgment

## References

[1] Hillel, T., Bierlaire, M., Elshafie, M. Z., and Jin, Y. A systematic review of machine learning classification methodologies for modelling passenger mode choice. *Journal of Choice Modelling*, 38:100221, 2021. ISSN 1755-5345. doi:10.1016/j.jocm.2020.100221.

[2] Hagenauer, J. and Helbich, M. A comparative study of machine learning classifiers for modeling travel mode choice. *Expert Systems with Applications*, 78:273–282, 2017. ISSN 0957-4174. doi:10.1016/j.eswa.2017.01.057.

[3] Tamim Kashifi, M., Jamal, A., Samim Kashefi, M., Almoshaogeh, M., and Masiur Rahman, S. Predicting the travel mode choice with interpretable machine learning techniques: A comparative study. *Travel Behaviour and Society*, 29:279–296, 2022. ISSN 2214-367X. doi:10.1016/j.tbs.2022.07.003.

[4] Hillel, T., Elshafie, M. Z. E. B., and Jin, Y. Recreating passenger mode choice-sets for transport simulation: A case study of London, UK. *Proceedings of the Institution of Civil Engineers – Smart Infrastructure and Construction*, 171(1):29–42, 2018. ISSN 2397-8759. doi:10.1680/jsmic.17.00018.

[5] Luckner, M., Wrona, P., Grzenda, M., and Łysak, A. Analysing urban transport using synthetic journeys. In *International Conference on Computational Science*, pages 118–132. Springer, 2024.

[6] Grzenda, M., Luckner, M., and Wrona, P. Urban traveller preference miner: Modelling transport choices with survey data streams. In M.-R. Amini, S. Canu, A. Fischer, T. Guns, P. Kralj Novak, and G. Tsoumakas, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 654–657. Springer Nature Switzerland, Cham, 2023. ISBN 978-3-031-26422-1.

[7] de Dios Ortúzar, J. and Willumsen, L. *Modelling Transport*. Wiley, 2024. ISBN 9781119282358.

# Medical Applications of Artificial Intelligence

Track Chairs:

- prof. Włodzisław Duch – Nicolaus Copernicus University

- prof. Julian Szymański – Gdańsk University of Technology

- prof. Jacek Rumiński – Gdańsk University of Technology

# Comparative Pharmacokinetics of Nicotine from E-Cigarettes and Traditional Cigarettes: A PBPK Modeling and Machine Learning Approach

**Joanna Chwał**[1,2,3][0009−0000−9363−4595],
**Arkadiusz Banasik**[1][0000−0002−4267−2783],
**Radosław Dzik**[1][0000−0002−6289−7234],
**Piotr Pańtak**[4][0000−0001−8496−2626],
**Ewaryst Tkacz**[1][0000−0003−2580−7954]

[1]*Department of Clinical Engineering, Academy of Silesia*
*Rolna 43, 40-55 Katowice, Poland*
*jpanna.chwal@akademiaslaska.pl,*
*arkadiusz.banasik@akademiaslaska.pl,*
*radoslaw.dzik@akademiaslaska.pl, ewaryst.tkacz@wst.pl*
[2]*Joint Doctoral School, Silesian University of Technology*
*Akademicka 2A, 44-100 Gliwice, Poland*
[3]*Department of Medical Informatics and Artificial Intelligence*
*Silesian University of Technology*
*Akademicka 2A, 44-100 Gliwice, Poland*
[4]*Faculty of Materials Science and Ceramics, AGH University of Krakow*
*B-8, Adam Mickiewicz Alley 30/0.24, 30-059 Kraków, Poland*
*pantak@agh.edu.pl*

**Abstract.** *Nicotine addiction remains a major public health concern, with e-cigarettes altering exposure dynamics by delivering nicotine in aerosol form. This study uses Physiologically Based Pharmacokinetic (PBPK) modeling combined with eXtreme Gradient Boosting (XGBoost) to simulate nicotine distribution and assess individual variability.*

43

*Results show that e-cigarettes lead to more sustained nicotine exposure, while traditional cigarettes cause rapid peaks in blood and brain concentrations. Blood-brain permeability was identified as the key factor influencing brain accumulation. ML integration significantly enhanced prediction accuracy ($R^2 > 0.9998$). These findings underscore the pharmacokinetic differences between delivery methods and demonstrate the value of combining PBPK and ML approaches for personalized exposure modeling and public health guidance.*

**Keywords:** *nicotine, PBPK modeling, machine learning, XGBoost, blood-brain barrier, e-cigarettes, smoking cessation*

# 1. Introduction

Nicotine is a widely used psychoactive compound with significant physiological and neurological effects, playing a central role in tobacco addiction [1]. The emergence of e-cigarettes has intensified research into differences in nicotine absorption, metabolism, and health impacts compared to traditional cigarettes. Delivered via aerosol rather than smoke, e-cigarettes affect nicotine pharmacokinetics differently, which is crucial for evaluating their health risks and shaping public health policies [2]. In 2016, the FDA expanded its regulatory authority to include electronic nicotine delivery systems (ENDS) [3], emphasizing the need to understand nicotine pharmacokinetics in both conventional and reduced-risk products [4].

Physiologically Based Pharmacokinetic (PBPK) modeling enables detailed simulation of nicotine ADME processes by incorporating organ-specific parameters such as blood flow, metabolism, and tissue composition [5]. It offers greater precision than simpler models and is widely used to compare delivery methods and account for individual metabolic differences [6, 7].

In this study, we developed PBPK models to simulate nicotine metabolism following inhalation from cigarettes and e-cigarettes, focusing on lung absorption, hepatic metabolism, and brain delivery. To enhance prediction accuracy, we integrated machine learning methods – specifically XGBoost – to model complex interactions among physiological parameters. Trained on simulated data, our ML-enhanced PBPK model enables personalized nicotine exposure estimates in the brain and blood. Our goal was to compare nicotine kinetics between e-cigarettes and cigarettes, assess interindividual variability, and apply ML to refine predictions. This integrative approach aims to support more informed regulatory decisions and public health strategies.

## 2. Materials and Methods

### 2.1. Study Design and PBPK Modeling

This study used computational modeling to compare nicotine pharmacokinetics from e-cigarettes and traditional cigarettes. PBPK models were developed in MATLAB R2024b to simulate nicotine ADME processes across key compartments (lungs, blood, liver, brain, synapses), using mass balance, Fick's law, and first-order kinetics [8]. Differential equations were solved numerically (ODE45; RelTol = 1e-6, AbsTol = 1e-8), simulating 24-hour exposure profiles.

### 2.2. Nicotine Dosing and Exposure Scenarios

E-cigarette dosing was based on 6mg/mL e-liquid, 0.1mL per puff, 10 puffs/session, and 50% bioavailability [9], yielding 3mg per session and 60mg/day. Cigarette exposure was set at 1.5mg per session [10], both administered hourly. Real-world pharmacokinetic data [11, 12] informed model validation (Cmax/Tmax values in brain, blood, and synapses).

### 2.3. Machine Learning Integration

To assess individual variability, XGBoost was used to predict nicotine Cmax in each compartment. The model was trained on 1,000 simulated individuals with varied physiological traits, including $k_{metab}$, $k_{penet}$, $Q_p$, body weight, smoking history, and metabolic phenotype.

### 2.4. Performance Evaluation

Model accuracy was assessed via RMSE, MAE, and $R^2$ metrics. SHAP values quantified feature importance, while Pearson correlations explored associations between physiological parameters and nicotine distribution.

## 3. Results

PBPK simulations showed that e-cigarettes led to higher cumulative nicotine concentrations than traditional cigarettes, particularly in the blood, liver, and brain (Figure 1). This reflects behavioral differences in inhalation and results in greater

Figure 1. Nicotine concentrations over time for e-cigarettes vs. traditional cigarettes

systemic exposure and prolonged nicotine presence in neural and metabolic compartments.

Model validation demonstrated excellent predictive accuracy, with predicted nicotine concentrations closely matching simulated values (RMSE <0.1%, $R^2$ > 0.9998; Figure 2). Among physiological parameters, blood-brain permeability ($k_{\text{penet}}$) had the strongest influence on brain and synaptic levels, followed by metabolic rate ($k_{\text{metab}}$), pulmonary blood flow ($Q_p$), body weight, and smoking history. Correlation analysis confirmed these findings, showing strong positive correlations between $k_{\text{penet}}$ and nicotine levels in the brain (0.90) and synapse (0.92), and a strong negative correlation with blood concentrations (–0.94). Faster metabolism was consistently associated with lower nicotine levels across all compartments.

Correlation analysis supported these findings: $k_{\text{penet}}$ showed strong positive correlation with brain and synapse concentrations (0.90 and 0.92), and negative correlation with blood (-0.94). Higher metabolic rate correlated with lower nicotine levels in all compartments.

In a representative individual, predicted concentrations reached 108ng/mL in the brain, 64ng/mL in blood, and 137ng/mL in synapses, indicating notable synaptic accumulation. Model reliability was further supported by comparisons with real-world pharmacokinetic data (Figure 3), where $R^2$ values exceeded 0.79 (brain), 0.87 (blood), and 0.90 (synapse), with low residual errors across all compartments.

Figure 2. Predicted vs. actual nicotine concentrations



Figure 3. Real-world vs. predicted nicotine concentrations

## 4. Discussion and Conclusions

PBPK modeling enables detailed simulation of nicotine distribution, but it does not fully capture inter-individual variability. Integrating machine learning (XG-Boost) improved prediction accuracy ($R^2 > 0.9998$) and accounted for differences in metabolism, perfusion, and behavior [13]. Blood-brain permeability ($k_{penet}$) emerged as the most influential factor in brain and synaptic nicotine accumulation, followed by metabolic rate, pulmonary blood flow, body weight, and smoking history.

Our findings align with prior studies showing distinct absorption routes: e-cigarettes (ENDS) favor upper airway and buccal delivery, while cigarettes induce faster systemic uptake via the lower respiratory tract [14]. Nicotine's rapid metabolism (half-life 2–3 h) and lysosomal trapping (Vd ~2.6 L/kg) contribute to higher concentrations in organs like the liver and lungs [9, 15, 16].

Although e-cigarettes avoid combustion-related toxins, they maintain elevated systemic nicotine levels, especially with nicotine salts, raising concerns about sustained dependence. This underscores the need for tailored cessation strategies and regulatory oversight. Limitations of this study include standardized behavioral assumptions and exclusion of other tobacco constituents.

In conclusion, combining PBPK and ML modeling offers a robust framework for simulating nicotine kinetics and inter-individual variability. Future research should incorporate behavioral and clinical data, evaluate different nicotine formulations, and assess long-term health impacts of emerging nicotine delivery systems.

# References

[1] Benowitz, N. L. Nicotine addiction. *New England Journal of Medicine*, 362(24):2295–2303, 2010.

[2] Goniewicz, M. L., Knysak, J., Gawron, M., Kosmider, L., Sobczak, A., Kurek, J., Prokopowicz, A., Jablonska-Czapla, M., Rosik-Dulewska, C., Havel, C., Jacob, P., and Benowitz, N. Levels of selected carcinogens and toxicants in vapour from electronic cigarettes. *Tobacco Control*, 23(2):133–139, 2014.

[3] U.S. Food and Drug Administration. FDA's Comprehensive Plan for Tobacco and Nicotine Regulation, 2017.

[4] Food and Drug Administration. Premarket tobacco product applications and recordkeeping requirements. *Federal Register*, 84(186):50566–50622, 2019.

[5] Nestorov, I. Whole body pharmacokinetic models. *Clinical Pharmacokinetics*, 42(10):883–908, 2003.

[6] Plowchalk, D. R., Andersen, M. E., and deBethizy, J. D. A physiologically based pharmacokinetic model for nicotine disposition in the Sprague-Dawley rat. *Toxicology and Applied Pharmacology*, 116(2):177–188, 1992.

[7] Hukkanen, J., III, P. J., and Benowitz, N. L. Metabolism and disposition kinetics of nicotine. *Pharmacological Reviews*, 57(1):79–115, 2005.

[8] Stillwell, W. *Membrane Transport*, pages 423–451. 2016.

[9] Benowitz, N. L., Hukkanen, J., and Jacob, P. *Nicotine Chemistry, Metabolism, Kinetics and Biomarkers*, pages 29–60. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009.

[10] Benowitz, N. L. and III, P. J. Daily intake of nicotine during cigarette smoking. *Clinical Pharmacology & Therapeutics*, 35(4):499–504, 1984.

[11] Schroeder, M. J. and Hoffman, A. C. Electronic cigarettes and nicotine clinical pharmacology. *Tobacco Control*, 23(suppl 2):ii30–ii35, 2014.

[12] Rose, J. E., Mukhin, A. G., Lokitz, S. J., Turkington, T. G., Herskovic, J., Behm, F. M., Garg, S., and Garg, P. K. Kinetics of brain nicotine accumulation in dependent and nondependent smokers assessed with PET and cigarettes containing 11C-nicotine. *Proceedings of the National Academy of Sciences*, 107(11):5190–5195, 2010.

[13] Chou, W.-C. and Lin, Z. Machine learning and artificial intelligence in physiologically based pharmacokinetic modeling. *Toxicological Sciences*, 191(1):1–14, 2022.

[14] Rostami, A. A., Campbell, J. L., Pithawalla, Y. B., Pourhashem, H., Muhammad-Kah, R. S., Sarkar, M. A., Liu, J., McKinney, W. J., Gentry, R., and Gogova, M. A comprehensive physiologically based pharmacokinetic (PBPK) model for nicotine in humans from using nicotine-containing products with different routes of exposure. *Scientific Reports*, 12(1):1091, 2022.

[15] Trapp, S., Rosania, G. R., Horobin, R. W., and Kornhuber, J. Quantitative modeling of selective lysosomal targeting for drug design. *European Biophysics Journal*, 37(8):1317–1328, 2008.

[16] Kolli, A. R., Calvino-Martin, F., Kuczaj, A. K., Wong, E. T., Titz, B., Xiang, Y., Lebrun, S., Schlage, W. K., Vanscheeuwijck, P., and Hoeng, J. Deconvolution of systemic pharmacokinetics predicts inhaled aerosol dosimetry of nicotine. *European Journal of Pharmaceutical Sciences*, 180:106321, 2023.

# Automated Symptom-Disease Association.
# Discovery from Clinical Notes

**Paulina Gacek**[0009−0008−1242−7542]

*AGH University of Krakow*
*Faculty of Electrical Engineering, Automatics, IT and Biomedical*
*Engineering, Department of Applied Computer Science*
*Mickiewicza 30, 30-059 Krakow, Poland*
*paulinagacek@student.agh.edu.pl*

**Abstract.** *Clinical notes are a rich source of medical information, but their unstructured nature limits their utilization in research and clinical decision-making. Transforming this data into a structured format enhances information retrieval and enables the discovery of meaningful symptom-disease associations. This study proposes an automated knowledge extraction framework that leverages large language models (LLMs) to analyze unstructured clinical notes and construct a patient-centered knowledge graph. Extracted entities are linked to standardized medical ontologies, facilitating efficient querying and trend analysis. The system is evaluated on real-world clinical data, demonstrating its effectiveness in capturing clinically relevant symptom patterns.*

**Keywords:** *clinical named entity recognition, knowledge graphs, large language models, electronic health records*

## 1. Introduction

Clinical documentation has transitioned from paper to digital format with the rise of electronic health records (EHRs). While these records are a rich source of valuable information, they are written in plain text and lack a standardized structure, which presents a significant barrier to their full utilization [1]. Converting such unstructured data into a structured format could uncover new insights, improve information retrieval, and facilitate data-driven prediction models [2].

A key step in structuring medical data is extracting medical entities, such as diseases and drugs, along with their relationships from clinical notes. However, Named Entity Recognition (NER) in the medical domain is particularly challenging due to the inherent complexities of medical text. These challenges include the use of differentiated abbreviations (e.g., 'BC1' for 'Blood-Count 1'), nonstandardized naming conventions (e.g., 'Chest X Ray', 'CXR', 'Chest X R'), and the extensive use of domain-specific terminology [3]. Moreover, free form nature of clinical notes often introduces spelling errors, ambiguous terms, and complex grammatical structures, including negations, further complicating NLP tasks. To ensure structured and interoperable medical data, extracted entities and their relationships can be stored in a knowledge graph [4, 5]. By integrating extracted clinical data with established ontologies such as MeSH and SNOMED CT, a comprehensive and interconnected representation of medical knowledge can be achieved.

This work proposes a knowledge graph-based approach to systematically extract symptom-disease associations from unstructured clinical notes. The presented approach leverages Large Language Models and established medical ontologies, to analyze real-world patient records and identify the most frequently observed symptoms for a given disease. The structured knowledge produced by this tool can then serve as a data foundation to create accurate and informative symptom awareness materials such as medical infographics.

## 2. Methodology

The proposed approach consists of three major steps: (1) clinical named entity extraction, (2) knowledge graph construction, and (3) symptom-disease association discovery.

### 2.1. Clinical Named Entity Extraction

Clinical entities are extracted from unstructured text using the Gemini Flash 2.0 Large Language Model[1], which is prompted to identify three types of entities:

- Patient: Inferred or explicitly stated attributes such as age and gender,
- Finding: Clinical observations present in the patient (e.g., "palpitations"),
- Disorder: Diagnosed medical conditions (e.g., "myocardial infarction").

---

[1]`https://cloud.google.com/vertex-ai/generative-ai/docs/gemini-v2`.

The distinction between *Finding* and *Disorder* follows SNOMED CT terminologies[2]. Extracted entities are standardized using string-matching techniques, leveraging both SNOMED CT and MeSH ontologies, as no single ontology comprehensively covers all medical concepts. To ensure patient-centered analysis, concepts referring to family history or negated conditions are excluded, a process handled directly by the LLM. The entity extraction pipeline is illustrated in Figure 1.



Figure 1. Clinical named entity extraction pipeline

## 2.2. Knowledge Graph Construction

Following entity extraction, identified entities are linked to individual patients, constructing a patient-centered knowledge stored in Memgraph[3], an open-source graph database and analytics platform. The schema of the graph database is presented in Figure 2. Patients who share common disorders or findings are connected through `HAS` relationship with the same nodes. Additionally, each patient node is linked to a final diagnosis node, encoded using the ICD-10 classification system. After constructing a cumulative knowledge graph for all the patients,

---

nodes representing findings and disorders are linked to their parent concepts using the `INSTANCE_OF` relationship, aligning with leveraged ontologies. For example, "chest pain on exertion" is an instance of "chest pain." Such an approach allows for better generalization of extracted entities, contextualizing them within broader medical knowledge.



Figure 2. Patient-centered knowledge graph schema

### 2.3. Symptom-Disease Association Discovery

Once the knowledge graph is constructed, it is queried to identify which findings and disorders frequently co-occur with specific diagnoses. The frequency of associations is normalized to account for variations in dataset size and patient distribution. The most representative symptoms for each diagnosis are then identified by ranking them based on their occurrence. An example of such a symptom-disease association discovery is presented in Section 3.

## 3. Results

To demonstrate the effectiveness of proposed framework, a dataset of cardiology hospitalizations (2003–2020) obtained from the Asseco Medical Management Solutions EHR system was utilized. The dataset comprises records of patients hospitalized at the 3rd Department of Cardiology, Leszek Giec Upper Silesian Medical Centre, Poland [6]. Each clinical record consists of four textual sections

corresponding to different stages of hospitalization: (1) reasons for admission and patient medical history, (2) physical examination, (3) discharge summary, and (4) recommendations. Additionally, each record contains a single ICD-10 diagnosis code, which does not fully capture the real-world occurrence of multiple concurrent diagnoses. However, it enables a structured classification approach for identifying the most significant diagnosis based on unstructured text.

The large language model demonstrated impressive resilience in processing real-world clinical text. It effectively interpreted abbreviations, corrected spelling errors, and handled various linguistic inconsistencies present in Polish medical records. Moreover, the LLM successfully deduced medical conditions from numerical data and contextual cues, enhancing the precision and clinical significance of the knowledge graph. For example, a left ventricular ejection fraction (EF) of 40–45% was accurately recognized as a reduced ejection fraction, even though this condition was not explicitly mentioned in the text.

Figure 3 presents a sample admission section from the dataset, while Figure 4 depicts the corresponding patient-centered knowledge graph.

Pacjent skierowany z IP Szpitala Powiatowego w Mikołowie, gdzie zgłosił się z powodu duszności trwającej od ok. 3 tygodni (z pogorszeniem od kilku dni). Dodatkowo bóle w klp przy wysiłku.Duszność nasila się przy leżeniu- szczególnie w nocy.Pacjent przyjęty celem diagnostyki układu krążenia i ustalenia dalszego sposobu leczenia. DOTYCHCZAS Nadciśnienie tętnicze. Cukrzyca typu 2-go. Reflux żołądkowo-przełykowy. LBBB. Krwawienia z PP, omdlenia, neguje UZYWKI- neguje Szczep p WZWB- nie Uczulenia- neguje Wyw środ- pracował w hucie LEKI- monoit, metocard.

Figure 3. Example of the admission section from Polish health record



Figure 4. Example of patient knowledge graph generated from admission section of Polish health record

Table 1 presents the observed frequencies of selected symptoms across different ICD-10 diagnostic codes. The values indicate the proportion of patients with a given diagnosis who exhibited the corresponding symptom.

Table 1. Observed symptom frequencies for selected ICD-10 diagnoses

| ICD-10 | Chest pain | Shortness of breath | Supraventricular tachycardia | Palpitations | Type II diabetes |
|--------|------------|---------------------|------------------------------|--------------|------------------|
| 120.0 | 0.69 | 0.28 | <0.05 | 0.11 | 0.29 |
| 121.4 | 0.80 | 0.21 | <0.05 | 0.08 | 0.31 |
| 125.0 | 0.56 | 0.25 | <0.05 | 0.15 | 0.38 |
| 148 | 0.19 | 0.27 | 0.67 | 0.3 | 0.20 |
| 150.0 | 0.24 | 0.67 | 0.34 | 0.15 | 0.33 |

The results reveal distinct symptom patterns for different cardiovascular conditions. For instance, shortness of breath was observed in 67% of patients diagnosed with I50.0 (heart failure) upon admission and in only 25% of those diagnosed with I25.0 (chronic ischemic heart disease). Chest pain was reported in 69% of patients diagnosed with I20.0 (angina pectoris) and in only 19% of those diagnosed with I48 (atrial fibrillation and flutter), aligning with clinical expectations. Notably, type II diabetes, disease strongly associated with an unhealthy lifestyle, was detected in 31% of patients across all diagnoses.

This structured representation of symptom-disease associations, derived from unstructured clinical notes, provides valuable insights into common presentations of cardiovascular conditions. Such findings can support clinicians in differential diagnosis and contribute to data-driven medical research.

# 4. Conclusions

This study demonstrated the effectiveness of the presented framework in capturing clinically relevant symptom patterns from unstructured clinical notes. Beyond identifying the most common symptoms for a given disease, the constructed patient knowledge graph may support additional clinical applications, such as recommending missing diagnoses based on similar patients. Unlike traditional machine learning models, the knowledge graph does not require retraining when new patients or conditions are introduced, making it a scalable and adaptable solution for real-world healthcare applications.

# References

[1] Sheikhalishahi, S., Miotto, R., Dudley, J. T., Lavelli, A., Rinaldi, F., and Osmani, V. Natural language processing of clinical notes on chronic diseases: Systematic review. *JMIR Medical Informatics*, 2019.

[2] Fraile Navarro, D., Ijaz, K., Rezazadegan, D., Rahimi-Ardabili, H., Dras, M., Coiera, E., and Berkovsky, S. Clinical named entity recognition and relation extraction using natural language processing of medical free text: A systematic review. *International Journal of Medical Informatics*, 177:105122, 2023. doi: 10.1016/j.ijmedinf.2023.105122.

[3] ElDin, H. G., AbdulRazek, M., Abdelshafi, M., and Sahlol, A. T. Med-Flair: medical named entity recognition for diseases and medications based on Flair embedding. *Procedia Computer Science*, 189:67–75, 2021. doi:10.1016/j. procs.2021.05.078.

[4] Rotmensch, M., Halpern, Y., Tlimat, A., Horng, S., and Sontag, D. Learning a health knowledge graph from electronic medical records. *Scientific Reports*, 2017.

[5] Lin, Z., Yang, D., and Yin, X. Patient similarity via joint embeddings of medical knowledge graph and medical entity descriptions. *IEEE Access*, PP:1–1, 2020. doi:10.1109/ACCESS.2020.3019577.

[6] Anetta, K., Horak, A., Wojakowski, W., Wita, K., and Jadczyk, T. Deep learning analysis of Polish electronic health records for diagnosis prediction in patients with cardiovascular diseases. *Journal of Personalized Medicine*, 12(6), 2022. doi:10.3390/jpm12060869.

C<span>HAPTER</span> 3

# Neural Network and Deep Learning Systems

Track Chairs:

- prof. Aleksander Byrski – AGH University of Science and Technology

- prof. Maria Ganzha – Warsaw University of Technology

# Combining Probabilistic Neural Networks with a Convolution Neural Network as a Feature Transformer

**Szymon Kucharczyk**[1,2][0009−0002−2413−6984],
**Piotr A. Kowalski**[2,3][0000−0003−4041−6900]

[1]*AGH Doctoral School, AGH University of Krakow*
*al. A. Mickiewicza 30, 30-059 Krakow, Poland*
*kucharcz@agh.edu.pl*
[2]*Faculty of Physics and Applied Computer Science*
*AGH University of Krakow*
*al. A. Mickiewicza 30, 30-059 Krakow, Poland*
*pkowal@agh.edu.pl*
[3]*Systems Research Institute, Polish Academy of Sciences*
*ul. Newelska 6, 01-447 Warsaw, Poland*

**Abstract.** *Probabilistic Neural Networks (PNNs) are memory-based Artificial Networks that have been successfully used for classification and regression problems for tabular data. PNNs differ significantly from deep neural networks (CNN or RNN) in terms of both architecture and network training. In addition, because of their architecture, they have been used mainly to solve problems with tabular data. Here, we propose an evaluation of the idea of joining convolution neural networks with PNNs to solve well-defined image classification problems (MNIST, CIFAR10). The results show that PNNs are capable of correctly classifying images after extracting the features using convolution layers.*
**Keywords:** *probabilistic neural networks, convolution neural networks, classification, feature extraction, hybrid networks*

## 1. Introduction

Classification and pattern recognition are areas of ongoing academic study. Many models and techniques have been developed depending on the type of data,

including traditional statistical models, logistic regression, tree-based methods, and various neural network architectures. When estimating the probability distribution of data is crucial, probabilistic methods such as Probabilistic Neural Networks (PNNs) often outperform other techniques. PNNs demonstrate strong performance, particularly in scenarios involving imbalanced datasets [1].

Despite their effectiveness in classification problems, PNNs have traditionally been applied to tabular data rather than structured data, such as images. This limitation arises from their unique architecture, which relies on kernel density estimation (KDE) and Bayesian decision theory rather than iterative weight optimization. Unlike deep neural networks (e.g., Convolutional Neural Networks – CNNs or Recurrent Neural Networks – RNNs), PNNs do not require backpropagation--based training. Instead, they directly estimate class probabilities based on stored training examples. This makes them highly interpretable and robust in handling uncertainty, but also constrains their direct applicability to high-dimensional input spaces like images.

To overcome this limitation, we propose a hybrid approach that leverages the strengths of both CNNs and PNNs. CNNs are well known for their ability to learn spatial hierarchies of features, making them highly effective in image processing tasks. Using a CNN as a feature transformer, we extract meaningful feature representations from images and feed them into a PNN classifier. This approach enables PNNs to operate on structured image data while maintaining probabilistic decision-making advantages.

The key motivation for this hybridization is the observation that, while deep learning models, such as CNNs, achieve remarkable accuracy in image classification, they often lack interpretability and require extensive labeled data for training. PNNs, on the other hand, offer a more transparent classification mechanism by directly modeling probability distributions. Combining the two, we aim to retain CNNs' feature extraction power while enhancing interpretability and robustness with PNN-based classification.

This study represents an initial exploration of the feasibility of CNN-PNN hybridization for image classification tasks. Given the preliminary nature of our research, the results are promising yet limited in scope. We evaluated the feasibility of CNN-PNN hybridization for well-known image classification tasks using MNIST and CIFAR10 datasets. Our experiments assess whether PNNs, when provided with CNNs' high-level feature representations, can achieve comparable classification performance to conventional CNN-MLP architectures.

## 1.1. PNN

PNNs are memory-based Artificial Neural Networks that use Kernel Density Estimation (KDE) and Bayes theorem to estimate the probability of data distribution. They consist of four layers [2, 3]:

- Input Layer – it gathers all input vectors to the network;

- Pattern Layer – calculates Kernel Density Estimator values for each pattern unit. Pattern units are built from training input data;

- Summation Layer – for the classification task, it sums the probability distribution for a given class from relevant pattern units;

- Output Layer – returns output class using the Bayes conditional probability theorem.

The PNN training methods differ significantly from training regular and deep artificial neural networks in a way that no gradient is optimized during training.

## 1.2. CNN Combined with PNNs

Combining Probabilistic Neural Networks with Deep Neural Networks (e.g. Convolution Neural Networks – CNN) can make significant progress in some classification and pattern recognition problems. The successful applications of PNNs to various tabular classification tasks [4, 5, 6] can lead to the idea of using these networks for other types of data, for example, text or images. The combination of PNN with other networks was not tested in the state-of-the-art. The combined CNN and PNN architecture is illustrated in Figure 1.

In the research, a common CNN architecture was used as a feature extractor and it was described in the supplementary materials (cnn_mnist.png and cnn_cifar10.png respectively).

# 2. Experiment

Here, we tested a PNN with a Cauchy kernel with a class level of the smoothing parameter [7] and a standard CNN for the image classification task. The CNN was trained together with a multilayer perceptron (MLP), and then the convolution layer was moved from the network. Then, the Convolution layer was used as

Figure 1. Block diagram of Probabilistic Neural Network with Convolution Neural Network as feature extractor

a feature extractor for the PNN. Next, the PNN was trained to classify images from two datasets on the extracted features. During PNN training, the weights of the convolution layer were frozen. The main goal of the experiment was to determine the PNN classification capabilities for non-tabular data.

The CNN + PNN architecture was compared with a standard CNN + MLP (Multilayer Perceptron) network, where MLP was the classifier.

## 2.1. Datasets

To measure the CNN-PNN classification performance, the Network was tested on two datasets: MNIST [8] and CIFAR 10 [9]. The former is a traditional image classification dataset for handwritten recognition, used as a field benchmark. The latter symbolizes more complicated images for 10 image classes. These datasets were selected for tests because they are well known in the literature and their reliability to test new image classification architectures was evaluated by many researchers.

## 2.2. Methods

Here, we built and trained a CNN for image classification for the MNIST and CIFAR10 datasets using *Keras* package [10]. Each data set was divided into train

and test sets before the Network training. The data were then normalized before training. For comparison of classification performance, precision and recall metrics were used [10]. CNN was trained using Adam optimizer with categorical cross-entropy metric [10]. For the MNIST dataset, the model was trained for 5 epochs with a batch size equal to 32. For the CIFAR10 dataset, the Keras model was trained for 30 epochs with a batch size of 128. PNN training was performed using the plug-in technique [7] after feature extraction with a frozen convolution layer.

### 2.3. Results

Table 1 highlights the image classification results for the MNIST and CIFAR10 datasets for CNN + MLP and CNN + PNN models. The CNN was trained as a feature extractor with MLP as a classifier. The convolution layer was then extracted and combined with a PNN to train a PNN classifier for the same datasets. The results show that MLP and PNN achieved similar classification performance for both datasets, although the MLP classifier generated slightly higher metric values. It might be caused by training the feature extractor (convolution layer) in a common pipeline with the MLP classifier.

Table 1. Results of image classification performance for CNN as a feature transformer and MLP or PNN as classifiers. The acc. stands for accuracy, p represents precision and r is the recall metric for test data.

| Dataset | MLP acc. | PNN acc. | MLP p | PNN p | MLP r | PNN r |
|---------|----------|----------|-------|-------|-------|-------|
| MNIST | 0.985 | 0.979 | 0.987 | 0.979 | 0.983 | 0.979 |
| CIFAR10 | 0.680 | 0.670 | 0.703 | 0.670 | 0.658 | 0.674 |

## 3. Conclusions

Table 1 presents classification metrics for the MNIST and CIFAR10 data sets using MLP and PNN as classifiers. The results show that PNN achieved a performance fairly similar to that of MLP for image classification. The results suggest that this hybrid approach offers a promising direction for integrating probabilistic modeling with deep learning techniques, potentially leading to more interpretable and adaptable neural network models. Although, for most of the test metrics, the classification performance was higher for the MLP than for the PNN, the metric

values were approximately similar. It might be caused by explicitly training the feature extraction for the MLP and then transferring it to the PNN classifier. It is assumed that the PNN classification performance could be higher when the feature extraction layers are trained together with the probabilistic layers. These results are preliminary and will serve as a basis for future work in the field.

### 3.1. Future Work

In the future, a joint training algorithm should be developed to train CNN together with PNN in a common pipeline, which should improve the accuracy of the hybrid network. The proposed solution should be tested on more datasets. In addition, the PNN might be combined with other Neural Networks to test its performance on other problems, e.g. textual data.

In addition, we plan to explore more efficient training methods that take advantage of metaheuristic approaches. Metaheuristic algorithms have shown their potential in optimizing complex models, particularly in cases where gradient-based techniques struggle with local optima or require extensive computational resources. One promising direction is the cHM algorithm, which we have previously described in [11]. This algorithm has shown encouraging results in optimizing neural network architectures by balancing exploration and exploitation in the search space. Integrating cHM into the training pipeline could enhance the generalizability and adaptability of the model, making it a valuable addition to our future research efforts.

## Acknowledgment

## References

[1] Kowalski, P. A. and Kusy, M. Sensitivity analysis for probabilistic neural network structure reduction. *IEEE Transactions on Neural Networks and Learning Systems*, PP:1–14, 2017. doi:10.1109/TNNLS.2017.2688482.

[2] Specht, D. Probabilistic neural networks for classification, mapping, or associative memory. In *IEEE 1988 International Conference on Neural Networks*, volume 1, pages 525–532. 1988.

[3] Specht, D. Enhancements to probabilistic neural networks. In *[Proceedings 1992] IJCNN International Joint Conference on Neural Networks*, volume 1, pages 761–768. 1992. doi:10.1109/IJCNN.1992.287095.

[4] Lee, J. J. and Yun, C.-B. Damage localization for bridges using probabilistic neural networks. *KSCE Journal of Civil Engineering*, 11:111–120, 2007. doi:10.1007/BF02823854.

[5] Adeli, H. and Panakkat, A. A probabilistic neural network for earthquake magnitude prediction. *Neural Networks: The Official Journal of the International Neural Network Society*, 22 7:1018–24, 2009.

[6] Wen, X.-B., Zhang, H., Xu, X.-Q., and Quan, J.-J. A new watermarking approach based on probabilistic neural network in wavelet domain. *Soft Computing*, 13:355–360, 2009. doi:10.1007/s00500-008-0331-y.

[7] Kowalski, P. A., Kusy, M., Kubasiak, S., and Łukasik, S. Probabilistic neural network – parameters adjustment in classification task. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. 2020. doi:10.1109/IJCNN48605.2020.9207361.

[8] Yann LeCun, C. C. The MNIST dataset of handwritten digits, 1998. URL `http://yann.lecun.com/exdb/mnist/`.

[9] Krizhevsky, A. Learning multiple layers of features from tiny images. *University of Toronto*, 2012.

[10] Chollet, F. et al. *Keras*, 2015. URL `https://github.com/fchollet/keras`.

[11] Kowalski, P. A., Kucharczyk, S., and Mańdziuk, J. Constrained hybrid metaheuristic algorithm for probabilistic neural networks learning. *Information Sciences*, 713:122185, 2025. URL `https://www.sciencedirect.com/science/article/pii/S0020025525003172`.

# On Approximating and Quantizing Fully Connected Classifiers

**Dariusz Puchała**[0000−0001−9070−8042]

*Institute of Information Technology*
*Lodz University of Technology*
*al. Politechniki 8, 93-590 Lodz, Poland*
*dariusz.puchala@p.lodz.pl*

**Abstract.** *This paper presents experimental results on approximating and quantizing an Multi-Layer Perceptron (MLP) based classification module. The approximation employs sparse neural networks with topologies inspired by fast divide-and-conquer algorithms for efficient linear transformations. Experiments were conducted using VGG-16 on the CIFAR-100 image classification task. After approximation, the module parameters were quantized to 4-bit, and a method to improve post-quantization accuracy was proposed.*
**Keywords:** *approximation of neural networks, quantization, compression*

## 1. Introduction

Convolutional Neural Networks (CNNs) for image classification typically consist of a convolutional feature extractor and fully connected classification layers. While models like VGG, ResNet, and Inception perform well, they demand substantial storage and computation (VGG-16 and VGG-19 exceed 500 MB) raising concerns about parameter redundancy. Redundancy in convolutional layers is addressed using techniques like SVD [1], low-rank approximations [2], and separable filters [3], while pruning [4], knowledge distillation [5], and sparse networks [6, 7] focus on fully connected layers. Quantization, particularly to INT8 and INT4 formats [8, 9], is vital for efficient GPU deployment.

Fully connected layers involve dense matrix-vector multiplications, requiring a number of computations and parameters proportional to the input-output size product. This becomes problematic for edge and embedded devices using FPGAs, which have limited on-chip RAM. Offloading weight storage to external

DRAM increases latency and energy consumption. Techniques to mitigate this include pruning, i.e. removing near-zero weights, which often produces irregular sparse matrices requiring extra metadata [10, 11]. To simplify representation, [12] use encodable permutations, while [13] employ pseudorandom sequences. Bitwidth reduction is also common, lowering storage by using fewer bits for less critical weights [14, 15, 16]. Low-rank approximation decomposes weight matrices into smaller factors, reducing both parameters and computation [17, 18].

This paper presents results on approximating and quantizing the classification module of VGG-16 trained on CIFAR-100. The approximation uses sparse neural networks inspired by fast divide-and-conquer algorithms for efficient linear transforms, significantly reducing parameter count. Subsequently, parameter precision was reduced from 32 to 4 bits, further minimizing storage requirements, particularly on GPUs. Experimental results show around 22,000 times model size reduction with only a 0.02% drop in test accuracy.

## 2. Approximating with Sparse Structures

Classification modules in models like VGG are typically implemented as fully connected MLPs. In VGG-16, for instance, the classifier includes two fully connected layers with 4,096 neurons each, followed by a 100-neuron output layer (for CIFAR-100), processing a 512-element input and totaling 19,325,028 parameters. To reduce this, [6, 7] proposed sparse structures that approximate fully connected layers. Let $N_0$ be the input size (a power of 2), $N_2$ the number of outputs, and $N_1$ the nearest power of $2 \geq N_2$. The sparse network consists of $\log_2(N_0/N_1)$ layers, each with neurons taking two inputs and producing one output (using bias and ReLU). Each layer halves the input vector size. The final layer is a dense layer mapping $N_1$ inputs to $N_2$ outputs (with bias and softmax). Figure 1 illustrates an example with $N_0 = 16$, $N_2 = 3$. The parameter count for such a sparse structure is:

$$\mathcal{L}_{PAR} = 3(N_0 - N_1) + N_2(N_1 + 1). \tag{1}$$

## 3. Quantization

Quantization is used to reduce the storage space required by a model. For sparse neural networks, it can be applied separately to each layer. Let $\{w_i\}$ be a set of $N_i$ parameters required by an $i$-th layer. Then, the quantization step $\Delta_i$ is

Figure 1. Exemplary sparse structure with 16 inputs and 3 outputs

calculated as $\Delta_i = 2\hat{w}_i/(2^b - 1)$, where $b$ is the number of bits for the representation of the parameters, and $\hat{w}_i = \max(|\min(w_i)|, \max(w_i))$. With $\Delta_i$ calculated, the quantization of parameters $w_i(n)$ for $n = 0, 1, \ldots, N_i - 1$, is performed using:

$$\overline{w}_i(n) = \Delta_i \left( \text{round} \left( \frac{w_i(n)}{\Delta_i} \right) \right),$$

where $\overline{w}_i(n)$ is a quantized value of parameter $w_i(n)$.

# 4. Experimental Results

In our experimental study, we considered the task of image classification using the VGG-16 deep neural network and the CIFAR-100 dataset. The accuracy of image classification obtained with this model was equal to 93.41% for the training dataset and 70.18% for the test dataset. The classification module, comprising three fully connected layers with 4,096, 4,096 and 100 neurons, was extracted from the model and subjected to the approximation procedure. The module approximation involved the following steps:

- the sparse neural network for $N_0 = 512$, $N_1 = 128$ and $N_2 = 100$ was constructed,

- the sparse network was trained to obtain the same outputs as the fully connected classifier for the training dataset using Mean Absolute Error (MAE) loss function.

The fully connected classifier required 19,308,644 parameters, while the sparse model used only 14,052 – a reduction of approximately 1,374 times. The sparse module achieved 93.17% training and 70.22% test accuracy, slightly outperforming the original on test data. This improvement is consistent with the known regularization effects of sparse structures [7].

Next, the approximated module was quantized to a 4-bit parameter format, with quantization applied independently to each layer. This reduced test accuracy to 68.04%. While the drop may be acceptable, we implemented the following procedure to improve performance:

- a twin module was created to replicate the approximated module;

- in each iteration of the training loop in the first place the twin model was initialized with the parameters from the approximated module;

- in the next step it was quantized to 4-bit representation;

- with quantized module the quantization noise was calculated at each layer's output and injected into the approximated module;

- the training algorithm was used to minimize the MAE between the fully connected and approximated module (with noise) outputs operating on the training set.

After applying the boosting procedure, classification accuracy improved to 94.39% on the training dataset and to 70.20% on the test dataset. This result is only by 0.02% lower for the test set.

In the final step, we analyzed the distribution of parameter values across the layers (see Figure 2). As shown, these distributions are non-uniform (with dominating zeros allowing for further structural simplifications), resulting in the following entropy values: 2.042 bits, 1.557 bits, and 1.532 bits, respectively. Hence, the actual number of bits (after entropy coding) required to store the parameters'

Figure 2. Probability distributions of parameters' values after quantization in layers: (a) first sparse, (b) second sparse, (c) densely connected output layer

indexes is smaller than 2 bits. This leads to memory reduction ratio of around 22,000 times (with 32-bit representation for the fully connected classifier).

Direct comparison with other studies is challenging due to variations in network architectures and datasets. However, several works focus on VGG-16 and the ImageNet dataset, which is more complex and typically yields less redundancy in classification layers. The method of [10] achieved a 22-time reduction in the weights of the fully connected layer by pruning, while [13] used pseudo-random encoding to achieve a 7-fold reduction. Combining pruning, quantization, and Huffman coding, [16] reached 80-time compression. The method in [11] achieved a 455-time reduction on CIFAR-10 dataset.

## 5. Conclusions

The experimental results presented in this study indicate that classifier modules constructed with fully connected layers may exhibit significant redundancy in the number of parameters. In the task considered, the parameter reduction ratio achieved after approximation was around 1,374 times, which is a significantly better result compared to other studies. It should also be noted that the applied sparse structure is regular and can be efficiently implemented in both software and hardware. Further on, by taking into account quantization to 4-bit representation of parameters' values, the memory reduction ratio was around 10,992, and with entropy coding involved, resulting in less than 2-bit representation, it reached the level of 21,985 times.

## References

[1] Denton, R. et al. Exploiting linear structure within convolutional networks for efficient evaluation. *arXiv preprint arXiv:1404.0736*, 2014.

[2] Zhang, X. et al. Accelerating very deep convolutional networks for classification and detection. *arXiv preprint arXiv:1505.06798*, 2015.

[3] Rigamonti, R. et al. Learning separable filters. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2754–2761. 2013.

[4] Lin, S. et al. Holistic CNN compression via low-rank decomposition with knowledge transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(12):2889–2905, 2019.

[5] Li, Z., Li, H., and Meng, L. Model compression for deep neural networks: A survey. *Computers*, 12(3), 2023.

[6] Alford, S., Robinett, R. A., Milechin, L., and Kepner, J. Training behavior of sparse neural network topologies. In *Proceedings of IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–6. 2019.

[7] Puchała, D. and Stokfiszewski, K. Sparse neural networks with topologies inspired by butterfly structures. In *Signal Processing Symposium*, pages 1–6. 2021.

[8] NVIDIA Developer Technical Blog - Int4 Precision for AI Inference. `https://developer.nvidia.com/blog/int4-for-ai-inference/`.

[9] Wu, X. et al. Understanding INT4 quantization for transformer models: latency speedup, composability, and failure cases. *arXiv preprint arXiv:2301.12017*, 2023.

[10] Han, S. et al. Learning both weights and connections for efficient neural networks. *arXiv preprint arXiv:1506.02626v3*, 2015.

[11] Nagaraju, D. and Chandrachoodan, N. Compressing fully connected layers of deep neural networks using permuted features. *IET Computers and Digital Techniques*, 17(3-4):149–161, 2023.

[12] Deng, C. et al. PERMDNN: Efficient compressed DNN architecture with permuted diagonal matrices. *arXiv preprint arXiv:2004.10936v1*, 2020.

[13] Karimzadeh, F. et al. A hardware-friendly approach towards sparse neural networks based on LFSR-generated pseudo-random sequences. *IEEE Transactions on Circuits And Systems - I: Regular Papers*, 68(2):751–764, 2021.

[14] Venkataramani, S. et al. AxNN: Energy-efficient neuromorphic systems using approximate computing. In *Proceedings of the 2014 International Symposium on Low Power Electronics and Design*, pages 27–32. 2014.

[15] Zhang, Q. et al. ApproxANN: an approximate computing framework for artificial neural network. In *Design, Automation Test in Europe Conference Exhibition*, pages 701–706. 2015.

[16] Han, S., Mao, H., and Dally, W. J. Deep compression: compressing deep neural networks with prunning, trained quantization and Huffman coding. *International Conference on Learning Representations*, pages 1–14, 2016.

[17] Lee, D., Kapoor, P., and Kim, B. Deeptwist: learning model compression via occasional weight distortion. *arXiv preprint arXiv:1810.12823v1*, 2018.

[18] Hua, Z. et al. Low rank regularization: A review. *arXiv preprint arXiv:1808.04521v3*, 2020.

# Natural Language Processing, Automatic Speech Recognition, and Conversational AI

Track Chair:

- prof. Maciej Piasecki – Wrocław University of Science and Technology

# Data Domain Adaptation for Machine Learning in Negative Emotion Recognition from Voice Data

**Anna Bryniarska**[1,2][0000−0002−3839−459X], **Piotr Schneider**[3,2][0000−0002−9666−9408],
**Dariusz Mikołajewski**[4,2][0000−0003−4157−2796],
**Magdalena Igras-Cybulska**[2][0000−0001−5621−7901],
**Aleksandra Kawala-Sterniuk**[5,2][0000−0001−7826−1292]

[1]*Opole University of Technology, Opole, Poland*
*Department of Computer Science, a.bryniarska@po.edu.pl*
[2]*Hishoo Sp. z o.o., Gdynia, Poland, magda.igras@gmail.com*
[3]*Maria Curie-Sklodowska University, Lublin, Poland*
*Department of Neuroinformatics and Biomedical Engineering*
*Institute of Computer Science, piotr.schneider@mail.umcs.pl*
[4]*Kazimierz Wielki University, Bydgoszcz, Poland*
*Faculty of Computer Science, dmikolaj@ukw.edu.pl*
[5]*Wroclaw University of Science and Technology, Wroclaw, Poland,*
*Department of Artificial Intelligence, Faculty of Information*
*and Communication Technology, aleksandra.kawala-sterniuk@pwr.edu.pl*

**Abstract.** *In this paper, our objective was to investigate the impact of data processing, such as filtering, degradation, and domain adaptation, on the performance of machine learning (ML) algorithms in the task of recognizing negative emotions in speech. The experiments were performed using the CREMA-D data set. The target domain for adaptation was the acoustic conditions typical of a call center environment. To achieve this, we proposed the use of a Wiener filter and algorithms for data degradation to better align with the characteristics of this domain. This paper presents ML results obtained for different variations of input data.*
**Keywords:** *machine learning, data filtering, domain adaptation, emotion recognition, voice signal*

73

# 1. Introduction

Speech emotion recognition is a key challenge in the field of speech signal processing and artificial intelligence (AI) [1, 2, 3]. Identification of negative emotions, such as anger, fear, and sadness, is crucial in various practical applications, including call centers, psychological diagnostics, and human-machine interactions [3, 4, 5, 6, 7]. Systems for emotion recognition help detect users' emotional states, enabling more empathetic and effective communication in automated systems [2, 6, 8, 9].

Traditional methods for speech emotion analysis relied on manually designed acoustic characteristics and statistical classification models [2, 10]. With the advancement of machine learning (ML), more sophisticated feature extraction and classification techniques have become feasible, significantly improving the accuracy of emotion recognition [9, 10, 11]. In this paper, we examine the impact of various signal processing techniques, including filtering, degradation, and domain adaptation, on the quality of speech emotion classification using ML algorithms.

Emotion recognition in call center environments has been the subject of various studies with the aim of improving automated systems for customer service interactions. Yurtay et al. [6] explored emotion recognition on call center voice data, focusing on the classification of emotions using ML techniques. The results demonstrated that feature selection and preprocessing methods significantly impact model performance as far as the importance of adapting speech emotion recognition systems to domain-specific conditions is concerned. A study by Du et al. [12] investigated human emotion recognition in the context of e-learning, exploring how emotional analysis can help assess student engagement and improve adaptive learning systems. Their study concerns ML techniques for classifying emotions based on speech and other modalities, demonstrating that accurate emotion detection can contribute to more personalized and effective educational environments. A paper written by Zielonka et al. ([13]) explored the use of CNNs for emotion classification in different speech datasets, analyzing the generalization capabilities of DL models. Their study highlighted the challenges associated with training on diverse datasets and demonstrated that dataset characteristics significantly impact model performance. A study by Donuk [14] focused on optimizing feature selection for speech emotion recognition using the CREMA-D dataset. The author proposed a feature selection method based on Binary Particle Swarm Optimization (BPSO) to enhance classification accuracy. The results indicated that selecting the most relevant features could significantly improve model performance while reducing computational complexity.

74

# 2. Materials and Methods

We applied signal processing methods, including filtering and degradation, in order to modify the original speech data and simulate real-world call center conditions.

The study used the CREMA-D dataset (Crowd-Sourced Emotional Multimodal Actors) [15], which comprises $7,442$ original audio recordings from 91 actors. The actors pronounced one of 12 predefined sentences while expressing different emotions: Anger, Disgust, Fear, Happiness, Neutrality, and Sadness, with four intensity levels: Low, Medium, High, and Unspecified [15]. Since the focus of this research is on recognizing negative emotions, the "happy" emotion, being a positive emotion, was excluded from the analysis.

To adapt the data to the acoustic conditions of call center audio recordings, the application of a Wiener filter was proposed. Initially, the filter models noise as a stochastic process and uses signal statistics to estimate the optimal filtering parameters. Then it analyzes the spectral characteristics of both the signal and noise to determine an appropriate attenuation factor for each frequency. As a part of the proposed approach to domain adaptation, the authors introduce a method to degrade the data to better match the acoustic characteristics of the call center environment. To achieve this, a custom technique was developed to simulate call center sound quality. The proposed algorithm applies the following transformations to a given audio signal:

- **Subsampling** to 8 kHz and band limitation,

- **Addition** of subtle white noise,

- **Compression** using the PCM $\mu$-law codec via FFmpeg,

- **Introduction of an echo** with a specified delay and decay factor.

The processed data was used to evaluate the impact of transformations on ML performance in recognizing negative emotions from speech. Feature extraction, based on correlation analysis, selected 29 key parameters using Praat and Librosa. These include Alpha Ratio (AR), Harmonic-to-Disharmonic Ratio (HDI), roughness, sharpness, jitter, shimmer, Fundamental Frequency (F0), number of pulses, formants (F1, F2, F3), Voiced Segment Length, loudness, and slope.

# 3. Results

Due to the appearance of *NaN* values during signal degradation with the script developed, the script was corrected and the data were prepared again in the following versions: with the use of the signal augmentation and degradation script, with the use of the Wiener filter and degradation.

The best results for ACC valid – **0.83** were achieved by Wiener phrase bin for the data file format: **train set** (1307, 30), **test set** (430, 30), **valid set** (437, 30). The remaining most efficient results achieved by the models are presented in Table 1.

Table 1. Results of the most efficient models

| Type of model | ACC valid | Train set | Test set | Valid set |
|---|---|---|---|---|
| Wiener phrase bin | 0.83 | 1307, 30 | 430, 30 | 437, 30 |
| Degraded phrase bin | 0.806 | 1299, 30 | 436, 30 | 439, 30 |
| Wiener degraded phrase bin | 0.797 | 1299, 30 | 436, 30 | 439, 30 |
| Orig phrase bin | 0.787 | 1294, 30 | 426, 30 | 438, 30 |
| Degraded sec bin | 0.76 | 2552, 30 | 858, 30 | 858, 30 |
| Degraded sec bin | 0.75 | 2564, 30 | 851, 30 | 853, 30 |

The models were trained and optimized using Python 3.12.8 and AutoGluon 1.2.0, which automates model selection, hyperparameter tuning, and ensemble learning. The framework improves interpretability, training efficiency, and deep learning support while selecting the best combination of algorithms (e.g., Light-GBM, CatBoost, neural networks) for high accuracy. It also offers tools for feature engineering and evaluation, ensuring robust predictive performance.

A 10% difference in the confusion matrix may result from random noise or minor test sample variations (Figure 1), and is often acceptable if overall performance remains high. In contrast, a 30% difference suggests major issues, such as model variability, poor data fit, or instability, requiring hyperparameter tuning, data quality improvements, or training strategy adjustments (Figure 2).

# 4. Discussion and Conclusions

In this study, we investigated the impact of signal processing techniques on the precision of ML models in recognizing negative emotions in speech (CREMA-D dataset). The proposed degradation algorithm, developed on the basis of an

(a) CatBoost BAG L1       (b) CatBoost BAG L2

Figure 1. Confusion matrix for Wiener phrase bin



Figure 2. Comparison of scores for the top 10 models

analysis of real call center recordings, was applied alongside a Wiener filter to simulate the acoustic conditions of a call center environment. Our findings suggest that adapting speech data to specific acoustic environments can improve the accuracy of emotion recognition. The application of the Wiener filter helped to improve signal clarity, while the degradation algorithm effectively simulated real-world call center conditions, allowing for a more realistic evaluation of the models' performance. The results of this study contribute to the important development of emotion-aware systems, particularly in environments where the detection of negative emotions is crucial to improving human-computer interaction (HCI), psychological evaluations, and customer service automation.

Future research directions may include the exploration of advanced DL techniques for feature extraction and classification, as well as the integration of multimodal data (e.g., speech and facial expressions) to enhance the accuracy of emotion recognition. Testing the proposed approach on real-world call center data could provide further insights into its practical applicability.

# Acknowledgment

# References

[1] Jahangir, R., Teh, Y. W., Hanif, F., and Mujtaba, G. Deep learning approaches for speech emotion recognition: State of the art and research challenges. *Multimedia Tools and Applications*, 80(16):23745–23812, 2021.

[2] Wani, T. M., Gunawan, T. S., Qadri, S. A. A., Kartiwi, M., and Ambikairajah, E. A comprehensive review of speech emotion recognition systems. *IEEE Access*, 9:47795–47814, 2021.

[3] Kumar, A. and Kumar, A. Human emotion recognition using machine learning techniques based on the physiological signal. *Biomedical Signal Processing and Control*, 100:107039, 2025.

[4] Grágeda, N., Busso, C., Alvarado, E., García, R., Mahu, R., Huenupan, F., and Yoma, N. B. Speech emotion recognition in real static and dynamic human-robot interaction scenarios. *Computer Speech & Language*, 89:101666, 2025.

[5] Madanian, S., Chen, T., Adeleye, O., Templeton, J. M., Poellabauer, C., Parry, D., and Schneider, S. L. Speech emotion recognition using machine learning—a systematic review. *Intelligent Systems with Applications*, page 200266, 2023.

[6] Yurtay, Y., Demirci, H., Tiryaki, H., and Altun, T. Emotion recognition on call center voice data. *Applied Sciences*, 14(20):9458, 2024.

[7] Mohmad, G. and Delhibabu, R. Speech databases, speech features and classifiers in speech emotion recognition: A review. *IEEE Access*, 2024.

[8] Lee, C.-C., Chaspari, T., Provost, E. M., and Narayanan, S. S. An engineering view on emotions and speech: From analysis and predictive models to responsible human-centered applications. *Proceedings of the IEEE*, 111(10):1142–1158, 2023.

[9] Ramakrishnan, S. Recognition of emotion from speech: A review. *Speech Enhancement, Modeling and Recognition–Algorithms and Applications*, 7:121–137, 2012.

[10] Akçay, M. B. and Oğuz, K. Speech emotion recognition: Emotional models, databases, features, preprocessing methods, supporting modalities, and classifiers. *Speech Communication*, 116:56–76, 2020.

[11] Abbaschian, B. J., Sierra-Sosa, D., and Elmaghraby, A. Deep learning techniques for speech emotion recognition, from databases to models. *Sensors*, 21(4):1249, 2021.

[12] Du, Y., Crespo, R. G., and Martínez, O. S. Human emotion recognition for enhanced performance evaluation in e-learning. *Progress in Artificial Intelligence*, 12(2):199–211, 2023.

[13] Zielonka, M., Piastowski, A., Czyżewski, A., Nadachowski, P., Operlejn, M., and Kaczor, K. Recognition of emotions in speech using convolutional neural networks on different datasets. *Electronics*, 11(22):3831, 2022.

[14] Donuk, K. Crema-d: Improving accuracy with bpso-based feature selection for emotion recognition using speech. *Journal of Soft Computing and Artificial Intelligence*, 3(2):51–57, 2022.

[15] Cao, H., Cooper, D. G., Keutmann, M. K., Gur, R. C., Nenkova, A., and Verma, R. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE Transactions on Affective Computing*, 5(4):377–390, 2014.

# CHAPTER 5

# Computer Vision

Track Chairs:

- prof. Leszek Chmielewski – Warsaw University of Life Sciences
- prof. Bogdan Kwolek – AGH University of Science and Technology

# The PDNcore Architecture
# for Industrial Anomaly Detection

**Oskar Graeb, Agata Giełczyk**[0000−0002−5630−7461]

*Bydgoszcz University of Science and Technology*
*Kaliskiego 7, 86-796 Bydgoszcz, Poland*
*agata.gielczyk@pbs.edu.pl*

**Abstract.** *In this article, we present the novel PDNcore architecture for industrial anomaly detection. We examined a wide range of parameters, including feature extractors, inference image size, coreset size, and pooling methods. Our experiments on the benchmark MVTec dataset showed promising results, particularly in terms of latency. The proposed PDNcore architecture reduces response time by nearly half in comparison to the SOTA methods.*

**Keywords:** *deep learning, anomaly detection, quality control systems, IoT*

## 1. Introduction

Quality control systems ensure product reliability and safety, especially in high-risk industries like automotive, aerospace, and medical engineering. A single defect can damage reputations and cause severe financial or even life-threatening consequences. Vision-based quality control emerged in the 1980s, but convolutional neural networks (CNNs) revolutionized the field. Models like AlexNet and ResNet enabled automated feature extraction, improving anomaly detection. Deep learning reduced reliance on manual feature engineering, while unsupervised learning helped detect anomalies as deviations from normal patterns. However, AI-driven industrial anomaly detection (IAD) remains challenging due to the variability of defects. Issues range from surface scratches to material inconsistencies or logical errors like miscounted packaging. Limited defective samples, changing environmental conditions, and real-time processing demands further complicate implementation. Additionally, frequent production changeovers require rapid system adaptation.

The rapid advancement of machine learning has given rise to a wide variety of architectures tailored for anomaly detection, each designed to tackle specific challenges such as scalability, accuracy, and efficiency. In industrial settings, where robust and dependable solutions are critical, selecting the right architecture is essential to meet practical demands such as real-time inference, minimal training data, and adaptability to changing product conditions. Among others, the PatchCore [1] architecture is widely implemented. It relies on mid-level features extracted from a pretrained backbone network, such as WideResNet50, to build a localized representation of input images. These features are organized as patch embeddings, capturing localized patterns that are highly relevant for anomaly detection tasks. By focusing on mid-level features, PatchCore mitigates domain mismatch issues introduced by backbones pretrained on datasets like ImageNet, which primarily focus on natural images. PatchCore architecture with its modifications was implemented in [2, 3, 4]. PatchCore, despite its advantages, has limitations that justify a custom architecture [5]. It requires significant computational and memory resources, especially for high-resolution images or large datasets, leading to out-of-memory errors or slowing down quality inspection. Additionally, its pretrained feature extractors may introduce ImageNet biases, limiting adaptability to industrial datasets. While PatchCore balances efficiency with accuracy, its limitations make a tailored, industry-optimized architecture a better solution.

On the other hand, the EfficientAD architecture is reported as very promising for IAD [6]. It uses a custom feature extractor that operates directly on image patches, namely the Patch Description Network (PDN). It enables the avoidance of the overhead of deep networks while preserving essential spatial context. Additionally, it shows less artifacts on ImageNet data when looking at absolute sum of gradients between channel dimension of the extracted features. EfficientAD, despite its advantages, has notable limitations that justify proposing a custom architecture. Its autoencoder, while effective, may struggle with detecting fine-grained logical anomalies, such as slight misalignments or small dimension variations, as it can sometimes generalize too well and inadvertently reconstruct anomalies. Additionally, the dual training requirements for the student–teacher model and the autoencoder increase overall training complexity, making rapid deployment more challenging. The highly optimized nature of the architecture also demands long training times to achieve good results, which may not be ideal for scenarios requiring quick adaptation.

Figure 1. The overall PDNcore architecture presenting the flow from input image through feature extraction, optional normalization, pooling, memory bank construction, and finally anomaly detection

## 2. The Proposed Method

Figure 1 illustrates the novel PDNcore architecture, which consists of several key components.

1. **PDN-Based Feature Extraction** – A modified Patch Description Network (PDN) serves as the backbone, using a fixed $33 \times 33$ pixel receptive field to ensure consistent spatial context for anomaly localization – three configuration were examined: L3 (PDN output from Layer 3), L4 (PDN output from Layer 4) and the combined PDN L3+L4;

2. **Feature Pooling and Dimensionality Reduction** – A pooling layer reduces feature map dimensionality, ensuring a compact representation. This minimizes memory requirements and optimizes nearest neighbor search without costly multi-layer concatenation. Local GeM pooling emerged as the most promising strategy;

3. **Feature Map Normalization** – This optional step balances channel contributions and enhances coreset subsampling during memory bank creation.

83

Enabling map normalization did not provide any improvement, so the most promising approach is to use no map normalization;

4. **Memory Bank Construction via Coreset Subsampling** – Following the PatchCore strategy, coreset subsampling selects a representative subset of features. Random projection further reduces dimensionality while preserving key information, enabling efficient anomaly detection. Using the coreset subsampling did not provide the improvement in evaluation metrics;

5. **Inference Image Size** – Unlike architectures that relied on fully connected layers and imposed strict input size constraints, the PDN feature extractor operated solely through convolutional layers and handled different input resolutions flexibly.

# 3. Obtained Results

## 3.1. Experimental Setup

Experiments were conducted on a server with an RTX A6000 Ada GPU. To stabilize performance, warmup runs preceded 100 synchronized forward passes per configuration, with the final score averaged over these runs. AUROC, F1-score, and latency served as evaluation metrics.

The MVTec [7] dataset was used, featuring 5,354 high-resolution RGB images across 15 categories (5 textures, 10 objects). Captured under controlled conditions, it introduces variability in orientation and natural imperfections, encompassing 73 anomaly types such as scratches, dents, and structural defects.

## 3.2. Results

Table 1 presents the evaluation metrics (averaged across all 15 MVTec dataset classes) for each examined feature extractor: PDN L3 (output from Layer 3), PDN L4 (output from Layer 4), the combined PDN L3+L4, and the baseline PatchCore. While the baseline achieves the highest average AUROC (0.983) and F1-score (0.973), PDN-based configurations demonstrate competitive performance in many defect categories while offering reduced inference latency.

Table 2 compares anomaly detection performance (Mean AUROC / Mean Image F1 Score) across different feature pooling methods on the MVTec AD dataset. The results show that using the Local GeM pooling seems to be the most promising

Table 1. Anomaly detection performance

| Parameter | PDN L3 | PDN L4 | PDN L3+L4 | Baseline |
|---|---|---|---|---|
| AUROC | 0.910 | 0.935 | 0.941 | 0.983 |
| F1-score | 0.912 | 0.931 | 0.932 | 0.973 |
| Latency [ms] | 4.78 | 4.61 | 5.19 | 5.62 |

Table 2. Anomaly detection performance across different feature pooling methods

| Method | Parameter | PDN L3 | PDN L4 | PDN L3+L4 | Baseline |
|---|---|---|---|---|---|
| Adaptive Avg | AUROC | 0.904 | 0.935 | 0.936 | 0.974 |
| | F1-score | 0.909 | 0.928 | 0.930 | 0.960 |
| | Latency [ms] | 5.147 | 5.237 | 5.491 | 8.368 |
| Local GeM | AUROC | 0.948 | 0.925 | 0.922 | 0.978 |
| | F1-score | 0.933 | 0.927 | 0.921 | 0.967 |
| | Latency [ms] | 2.928 | 2.930 | 3.081 | 5.354 |
| None | AUROC | 0.910 | 0.935 | 0.941 | 0.974 |
| | F1-score | 0.912 | 0.931 | 0.932 | 0.962 |
| | Latency [ms] | 4.759 | 4.596 | 5.164 | 5.551 |

strategy, achieving high AUROC and F1 scores due to its flexibility. By interpolating between average and max pooling, it creates compact, discriminative feature representations, enhancing PDNcore's anomaly detection while maintaining low latency.

AUROC, F1-score and latency measurements, also presented in Table 3, show that smaller coreset sizes only slightly reduce inference time. These results suggest that while reducing the coreset size marginally improves computational efficiency, it does not significantly impact inference time to validate very small coreset sizes.

The most promising PDNcore configuration is characterized by the following features: Feature extractor: PDN Layer 4; Image size: $256 \times 256$; Coreset size: 10,000; Pooling: Local GeM; Normalization: None. This configuration achieved a latency of 2.98 ms. With the same input sample size, PatchCore reached a latency of 5.62 ms. The proposed architecture thus reduces response time by nearly half.

Table 3. Anomaly detection performance across different coreset size

| Model | Metrics | 0.5k | 1k | 3k | 5k |
|---|---|---|---|---|---|
| PDN L3+L4 | Avg. AUROC | 0.854 | 0.900 | 0.925 | 0.935 |
| | Avg. F1-score | 0.891 | 0.907 | 0.922 | 0.928 |
| | Avg. Latency [ms] | 2.787 | 2.796 | 2.840 | 2.887 |
| Baseline | Avg. AUROC | 0.960 | 0.968 | 0.978 | 0.978 |
| | Avg. F1-score | 0.952 | 0.961 | 0.966 | 0.965 |
| | Avg. Latency [ms] | 5.303 | 5.309 | 5.288 | 5.308 |

# 4. Conclusions

In the article we proposed the novel PDNcore for IAD. It enables obtaining high, promising results in terms of AUROC and F1-score. We examined some models' parameters: features extractor, pooling method, coreset size and the inference image size. The best-performing configuration (noted for its robustness and low variability across different categories in MVTec) was Layer 4 as the features extractor, image size = $256 \times 256$), coreset (10.000 samples), pooling Local GeM and none normalization. By combining a memory bank-based training strategy with an efficient PDN-based feature extractor, selected for its low latency, fixed receptive field, and compact feature representation, PDNcore meets the stringent requirements of modern quality control systems.

Future research should focus on refining the feature extraction component, a key element of PDNcore. A promising direction is a multiple-head approach, where one branch detects fine details with a small receptive field, while another captures broader contextual information. This dual-path strategy could improve anomaly detection by addressing both small-scale defects and larger spatial inconsistencies like misplacement or improper coloration.

Additionally, advanced representation learning techniques, such as knowledge distillation and contrastive learning (e.g., Barlow Twins, SimCLR), could enhance feature quality and robustness, particularly in settings with limited labeled data. Expanding PDNcore to integrate thermal, hyperspectral, or other imaging modalities could improve defect detection beyond standard RGB analysis. Moreover, incorporating few-shot or continual learning methods would enable rapid adaptation to evolving production conditions and new defect types without extensive retraining, ensuring sustained high performance in dynamic industrial environments.

# References

[1] Roth, K., Pemula, L., Zepeda, J., Schölkopf, B., Brox, T., and Gehler, P. Towards total recall in industrial anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14318–14328. 2022.

[2] Heckler, L., König, R., and Bergmann, P. Exploring the importance of pre-trained feature extractors for unsupervised anomaly detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2917–2926. 2023.

[3] Santos, J., Tran, T., and Rippel, O. Optimizing patchcore for few/many-shot anomaly detection. *arXiv preprint arXiv:2307.10792*, 2023.

[4] Koshil, M., Wegener, T., Mentrup, D., Frintrop, S., and Wilms, C. Anomalouspatchcore: Exploring the use of anomalous samples in industrial anomaly detection. *arXiv preprint arXiv:2408.15113*, 2024.

[5] Jiang, Z., Zhang, Y., Wang, Y., Li, J., and Gao, X. Fr-patchcore: An industrial anomaly detection method for improving generalization. *Sensors*, 24(5):1368, 2024.

[6] Batzner, K., Heckler, L., and König, R. Efficientad: Accurate visual anomaly detection at millisecond-level latencies. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 128–138. 2024.

[7] Bergmann, P., Fauser, M., Sattlegger, D., and Steger, C. MVTec AD–a comprehensive real-world dataset for unsupervised anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9592–9600. 2019.

# Tracking Hand Motion for Gesture-Based Interfaces: Evaluating AI-Based Pose Estimation with a Custom Benchmarking Tool

**Adam Nowosielski**[0000−0001−7729−7867], **Krzysztof Małecki**[0000−0002−8687−1119], **Kacper Dogiel**[0009−0003−5360−7228]

*West Pomeranian University of Technology in Szczecin*
*Faculty of Computer Science and Information Technology*
*Żołnierska 49, 71-210 Szczecin, Poland*
*{adam.nowosielski, krzysztof.malecki, dk46497}@zut.edu.pl*

**Abstract.** *Recent advancements in human pose estimation (HPE) leverage deep learning to detect body key points in real-time from standard webcam images. This enables cost-effective, touchless, gesture-based interaction. We present a brief study on tracking dynamic hand gestures, focusing on movement trajectories through key checkpoints. Our analysis compares state-of--the-art AI-based models, including Mediapipe, OpenPose COCO, TensorFlow Lightning, and Thunder, evaluating accuracy and responsiveness. To ensure structured benchmarking, we developed a custom evaluation tool.*
**Keywords:** *touchless interfaces, deep-learning, human pose estimation*

## 1. Introduction

Efficient Human-Computer Interaction (HCI) systems have been a key research focus for years [1, 2]. Various peripheral devices enable communication with machines, but gestures have gained particular interest. Gesture-based interaction falls into two categories: touchscreen gestures, requiring physical contact, and in-air gestures, which eliminate this constraint [3]. While touchscreen gestures are widely adopted, touchless interaction remains underexplored, despite many attempts to popularize touchless gestures (e.g., the Kinect sensor, Intel RealSense).

Touchless interaction enables users to operate without additional medium or intermediate equipment, making them ideal for large displays, public spaces, medical environments, and human–robot interaction [4]. They are also used in automotive multimedia systems [5, 6, 7] and emerging fields like gesture-based drawing [8, 9], including those systems which use human hand gestures for writing characters, words or free drawing in free space [10].

Previous solutions relied on specialized sensors, limiting accessibility. In this paper, we explore state-of-the-art AI-based pose estimation methods for enabling no-cost, touchless interaction using standard RGB cameras.

The rest of the paper is structured as follows. Section 2 describes existing techniques to human body pose detection. Section 3 presents the developed research environment for evaluation of different control mechanisms in a non-contact interface. Evaluation and results of conducted experiments are presented in Section 4. The paper ends with conclusions.

## 2. Gestures and Human Pose Estimation

Gestures involve body movements and can be performed with the entire body, specific limbs (e.g., hands, head), or a minor parts of the body (e.g., fingers, eyes). They are categorized as [11]: static (fixed postures, like a raised hand for traffic control, a specific arrangement of the fingers of the hand, a face with one eye closed, etc.) or dynamic (requiring tracking mechanisms to detect sequences and specifics of the movement) [11]. Both types are used in touchless interfaces.

Traditional pose estimation relies on depth sensors, which generate depth maps to segment and detect human silhouettes. Technologies include stereo vision, structured light, and time-of-flight (ToF). While effective, these methods require dedicated hardware, which increases the costs.

Recent advancements in deep learning enable real-time silhouette detection from standard RGB images, eliminating the need for depth sensors. The reader is referenced to several recent survey articles for a detailed review of this topic [12, 13, 14]. The literature indicates that the most actively explored areas in human pose estimation include: categorization of HPE approaches (based on training techniques, 2D vs. 3D representation, number of subjects, viewpoints, and the use of temporal information), datasets and evaluation (data augmentation strategies, performance metrics), computational efficiency, and application domains.

While the above fields receive significant attention, the use of HPE for touch-

less user interfaces is not currently a mainstream research focus. Models such as Mediapipe, TensorFlow Lightning, and OpenPose COCO offer accurate key-point detection (Figure 1, top), even when parts of the body are out of view. However, challenges arise in constrained settings – e.g., when a hand obscures the face, models struggle to maintain tracking (Figure 1, bottom). These limitations motivated our study, which evaluates modern deep-learning solutions for non-contact interactions in front of a computer screen.



Figure 1. Correct (top) and incorrect (bottom) detections of human silhouette landmarks using (from left to right): Mediapipe, TensorFlow Lightning, and OpenPose COCO

## 3. Research Environment: A Custom Benchmarking Tool

We developed a custom benchmarking tool to evaluate state-of-the-art solutions for human pose detection in contactless interfaces. By eliminating dedicated sensors, our approach enables gesture-based interaction using standard hardware, such as a laptop with a built-in webcam.

Our tool focuses on dynamic hand gestures, interpreting them in two ways: as discrete movement sequences through key points and as continuous motion paths. Key points are visualized as shapes, whose size and color can be customized within the environment, and overlaid on the camera image to guide user interaction. In the case of the N-shaped gesture (Figure 2), the user follows a predefined sequence (left) by moving through numbered shapes (middle). As shown on the right of Figure 2, a shape turns green when the user's hand enters it, confirming correct positioning.

Figure 2. An exemplary dynamic gesture performed by a user with a hand movement: the N-shaped motion (left), a sequence of different shapes representing the gesture (middle), and (right) the user making the gesture, with the active shape marked in green

Additionally, we treat the gesture as a continuous trajectory, enabling detailed comparisons using a Modified Hausdorff Distance (MHD). While this study focuses on hand-based interaction, our tool allows the selection of any skeletal point, such as the head or other body parts, for gesture control. This adaptability, combined with customizable shape positioning and sizing, enables the creation of diverse gesture patterns tailored to various applications.

# 4. Evaluation

The effectiveness of deep-learning-based human pose detection methods for touchless gesture interaction was evaluated in terms of accuracy and processing speed. Among the tested solutions, TensorFlow Lightning demonstrated the highest performance at 16.0 FPS, while OpenPose COCO, optimized for GPU usage, was the slowest at 0.98 FPS. Mediapipe and TensorFlow Thunder achieved 11.04 FPS and 8.79 FPS, respectively, revealing significant differences in computational efficiency. All evaluations were performed on a CPU-based setup using an Intel Core i5-8265U laptop with 8 GB RAM. Although our experiments ensured consistency by using CPU implementations, real-time operation can be achieved with GPU acceleration in practical applications.

To objectively compare these methods, a user study with six participants was conducted. Each participant performed different gestures by following the designated sequence of shapes. The recorded sequences were then analyzed to verify whether the detected control points aligned with the intended motion path, al-

lowing for a quantitative comparison of all examined approaches under identical conditions.

Our environment provides detailed analyses, including gesture execution previews, performance metric calculations, and visualizations. In Figure 3, we present two exemplary results. On the left, an individual example of a single gesture performed by a user is shown. The graph illustrates the ratio of motion path coverage over time for different approaches. With each successive frame, the gesture coverage increases, though at varying rates depending on the method used. Ideally, each graph should approach 1, representing complete coverage in an optimal scenario.



Figure 3. Sample results: Individual user performance and motion path coverage ratio during gesture execution (left). Aggregated MHD results for all six participants and all gestures, with time normalized (right).

On the right hand side of Figure 3, we present the aggregated results for all users and all performed gestures, evaluated using the Modified Hausdorff Distance to the given patterns. Since both gesture type and individual user performance influence execution time, its duration varies. To enable comparison, gesture execution time in aggregated reports is normalized to a 0–1 range.

Our study indicates that the TensorFlow-based pose detection achieved the highest accuracy, as reflected in the lowest MHD values. Considering both accuracy and execution speed, TensorFlow Lightning emerges as the most suitable choice for touchless interaction, particularly in challenging conditions where the user is seated in front of a computer screen, with limited pose visibility and frequent occlusions during gesture execution.

# 5. Conclusions

While touchless user interfaces based on human pose estimation were a focus of research in the era of the Kinect sensor, they are not currently a central topic in the field. However, AI-based pose estimation shows that it could bring a new level of quality to touchless interaction systems.

This study evaluated four deep-learning-based human pose detection methods for gesture recognition using a custom research environment. Our analysis revealed significant performance variability across the tested models. While simple gestures, such as lateral hand movements, were recognized with high reliability, more complex gestures posed challenges, particularly under conditions of partial body presence and occlusions. Factors such as user distance from the camera, partial silhouette visibility, clothing, and wrist accessories also impacted recognition accuracy. Despite these limitations, the tested models remain suitable for touchless interaction. Future research should explore alternative control points, such as head-based interaction, and assess lighting conditions to improve real-world robustness.

# References

[1] Jaimes, A. and Sebe, N. Multimodal human–computer interaction: A survey. *Computer Vision and Image Understanding*, 108(1-2):116–134, 2007.

[2] Karray, F., Alemzadeh, M., Abou Saleh, J., and Arab, M. N. Human-computer interaction: Overview on state of the art. *International Journal on Smart Sensing and Intelligent Systems*, 1(1):137–159, 2008.

[3] Saffer, D. *Designing Gestural Interfaces: Touchscreens and Interactive Devices*. O'Reilly Media, Inc., 2008.

[4] Clark, A. and Ahmad, I. Touchless and nonverbal human-robot interfaces: An overview of the state-of-the-art. *Smart Health*, 27:100365, 2023. ISSN 2352-6483. doi:10.1016/j.smhl.2022.100365.

[5] Małecki, K., Nowosielski, A., and Kowalicki, M. Gesture-based user interface for vehicle on-board system: a questionnaire and research approach. *Applied Sciences*, 10(18):6620, 2020.

[6] Zhang, T., Liu, X., Zeng, W., Tao, D., Li, G., and Qu, X. Input modality matters: A comparison of touch, speech, and gesture based in-vehicle interaction. *Applied Ergonomics*, 108:103958, 2023.

[7] Graichen, L., Graichen, M., and Krems, J. F. Effects of gesture-based interaction on driving behavior: a driving simulator study using the projection-based vehicle-in-the-loop. *Human Factors*, 64(2):324–342, 2022.

[8] Patel, J., Mehta, U., Panchal, K., Tailor, D., and Zanzmera, D. Text recognition by air drawing. In *2021 Fourth International Conference on Computational Intelligence and Communication Technologies (CCICT)*, pages 292–295. IEEE, 2021.

[9] Saoji, S., Dua, N., Choudhary, A. K., and Phogat, B. Air canvas application using OpenCV and Numpy in Python. *IRJET*, 8(08), 2021.

[10] Sakri, L. I., Kerur, V., Shet, G., Mokashi, S., and Patki, A. Gesture-based drawing application: A survey. In H. K. Deva Sarma, V. Piuri, and A. K. Pujari, editors, *Machine Learning in Information and Communication Technology*, pages 303–309. Springer Nature Singapore, Singapore, 2023.

[11] Kaâniche, M. B. *Gesture Recognition from Video Sequences. PhD Thesis.* Université Nice Sophia Antipolis, 2009.

[12] Gao, Z., Chen, J., Liu, Y., Jin, Y., and Tian, D. A systematic survey on human pose estimation: upstream and downstream tasks, approaches, lightweight models, and prospects. *Artificial Intelligence Review*, 58(3):68, 2025. doi: 10.1007/s10462-024-11060-2.

[13] Neupane, R. B., Li, K., and Boka, T. F. A survey on deep 3D human pose estimation. *Artificial Intelligence Review*, 58(1):24, 2024. doi: 10.1007/s10462-024-11019-3.

[14] Zheng, C., Wu, W., Chen, C., Yang, T., Zhu, S., Shen, J., Kehtarnavaz, N., and Shah, M. Deep learning-based human pose estimation: A survey. *ACM Computing Surveys*, 56(1), 2023. doi:10.1145/3603618.

# Anomaly Detection in Ground-Based Sky Imagery

**Mateusz Piechocki**[0000−0002−3479−0237], **Marek Kraft**[0000−0001−6483−2357]

*Poznan University of Technology*
*Institute of Robotics and Machine Intelligence*
*Piotrowo 3A, 60-965 Poznań, Poland*
*mateusz.piechocki@put.poznan.pl, marek.kraft@put.poznan.pl*

**Abstract.** *Ensuring the reliability of neural networks in industrial applications is challenging due to data drift and anomalies, especially in nonstationary environments, such as solar irradiance forecasting based on sky images. This study presents an image encoder embeddings classifier based on the isolation forest algorithm to detect outliers in data streams. By autonomously monitoring data validity at remote locations, the presented method aims to enhance forecast reliability, minimizing disruptions due to unexpected variations. This approach is a step towards the seamless integration of solar irradiance forecasting models into smart grids and energy management systems, contributing to the more reliable and efficient use of renewable energy.*
**Keywords:** *computer vision, image processing, anomaly detection, deep learning*

## 1. Introduction

One of the key factors limiting the use of neural networks in many industrial applications is the challenge of proving that a trained network will continue to produce reliable outcomes once it is deployed. Most of the deep learning applications operate under changing conditions in non-stationary environments. This often leads to changes in the statistical properties and characteristics of the input data over time, commonly known as data drift. If not effectively detected and managed, data drift can degrade model performance, resulting in poor predictions and potentially flawed decisions that static models cannot overcome [1]. Hence, there is a high demand for a reliable and efficient solution to continuously monitor the

statistical properties of the model and alert on significant deviations from expected patterns. In addition, using test-time adaptation techniques, it is possible to adjust the model parameters autonomously. However, it is extremely important to be able to distinguish whether the new samples from an endless data stream are close to the original data distribution and can be used to update the model parameters or whether they are isolated or short-lived anomalies that should be discarded. This is a critical challenge in real-world applications, such as environmental monitoring or surveillance, where unexpected data or measurements can compromise performance or introduce uncertainty [2]. One such application is solar irradiance forecasting, which is a good indicator of the energy production from solar sources. As proven in the literature, solar irradiance forecasting, especially based on machine learning methods, supported by sky observations and image processing techniques, provides accurate estimates and supports efficient energy management [3, 4]. However, to virtually take advantage of production estimation and integrate such systems into smart grids or energy management systems, forecasts must be reliable and resistant to anomalies. Although numerous works have proposed diversified forecasting solutions, the literature lacks a study that addresses data validity issues. Therefore, this research addresses the problem of identifying anomalies in ground-based sky imagery to provide reliable forecasts and improve energy management strategies. By developing an image embeddings classifier, we are able to continuously and independently monitor new samples directly at remote measurement sites and notify the operator only when anomalies are identified.

## 2. Methodology

**Dataset.** The task presented in this article concerns anomaly detection in ground-based sky imagery. One of the most common applications of sky observations is vision-based solar irradiance forecasting. Therefore, the Folsom dataset [5] was selected as a baseline data collection because it offers three years of 1-minute resolution sky images, with corresponding irradiance readings, gathered in California, USA. The image collection was divided into three subsets: training and validation, which were used to train the image-to-irradiance model and develop an anomaly identification algorithm, and a test set to simulate the stream of new data.

**Methods.** Image-to-irradiance algorithms are the most potential approaches for intra-hourly solar irradiance forecasting. Hence, sky images are an essential

source of information. To extract the most relevant features, the image encoders are applied at the initial processing level. Therefore, in this study, the ResNet50 [6] architecture has been adopted for this task, as a fundamental part of the image-to-irradiance algorithm designed to predict solar irradiance values for a 15-minute forecast horizon. The model structure is shown in Figure 1. The image encoder computes a vector of embeddings containing 2,048 values, based on which inputs are classified as inliers and outliers. For this purpose, an unsupervised algorithm based on the isolation forest [7] was implemented, with *n_estimators* equal to 100, due to its computational efficiency even on high-dimensional data. The parameters of the anomaly detection algorithm were fixed on the training set.



Figure 1. Anomaly detection pipeline. Embeddings computed by the image encoder are compared with the distribution of samples in the baseline dataset. Using the isolation forest algorithm, new samples are classified as inliers or outliers. The grayed out section covers other parts of the image-to-irradiance forecasting algorithm on which the initial model, including the image encoder, was trained.

**Experimental Setup.** The image-to-irradiance model training setup used an MSE-based loss function, a computationally efficient AdamW [8] optimizer with an initial learning rate of $3e-4$, batch size limited to 64 and 30 training epochs with early stopping callback set for 5 epochs. PyTorch [9] was selected as the deep learning framework and the pre-training was performed on an NVIDIA GeForce RTX 4090 GPU with 24 GB memory and CUDA 12.6.

# 3. Results

Several types of outliers have been detected using image feature vectors and the isolation forest algorithm. The first group consists of those that cover the field of view of the camera. For example, a short-term obscuration caused by an insect

or a long-term obscuration caused by dirt on the camera dome. The other group is composed of various weather phenomena, especially heavy cloud cover and rainfall. Although these are not as significant anomalies as those in the first group, it is worth looking at these samples, because they are close to the distribution and may be valuable for model adaptation. Figure 2 shows the anomalies identified using the isolation forest algorithm in the test data stream. In addition to visual inspection, dimensionality reduction methods were used to identify the location of outliers relative to training data and inliers from the test set. First, the image embedding vectors were reduced to 256 values using principal component analysis (PCA) and then to a 2D space using uniform model approximation and projection (UMAP) [10]. Figure 3 shows the image embeddings in the 2D space, color-coded into training data, test inliers, and outliers.



Figure 2. Example of anomalies detected in the test set based on the isolation forest algorithm operating on image encoder embeddings. (A, B) dirt on the camera dome; (C, D) insect captured in the field of the camera's field of view; (E) images containing heavy cloud cover or rainfall (F).

Figure 3. Image embeddings visualized in 2D space using UMAP manifold approximation and projection preceded by PCA dimensionality reduction

## 4. Conclusions

The presented solution addresses the problem of unsupervised anomaly detection in the data stream of ground-based sky cameras. In order not to significantly increase the data processing and computational complexity of the inference, anomalous images are identified using embeddings generated by an image encoder. This vector is processed much more effectively than the raw image, as it contains only the most relevant features. The applied isolation forest algorithm is highly efficient due to linear-time complexity, small memory footprint, and applicability to high-dimensional data. Therefore, it can be easily integrated into the data processing pipeline. Overall, although the proposed method is presented in the context of sky images in image-to-irradiance forecasting, it can be easily adapted to other vision solutions, including a wide range of remote computer vision solutions that perform surveillance or environmental monitoring tasks in the open world and require reliable, anomaly-free data samples. Since images with heavy cloud cover or rainfall represent a significant number of reported outliers, further work will include the development of new methods to increase the sensitivity to anomalies such as lens obstruction, out-of-focus, or image capture artifacts rather than weather variations, unless they are sudden or extreme.

99

# References

[1] Lwakatare, L. E. et al. Large-scale machine learning systems in real-world industrial settings: A review of challenges and solutions. *Information and Software Technology*, 127:106368, 2020. doi:10.1016/j.infsof.2020.106368.

[2] Gemaque, R. N. et al. An overview of unsupervised drift detection methods. *WIREs Data Mining and Knowledge Discovery*, 10(6):e1381, 2020. doi:10.1002/widm.1381.

[3] Wen, H. et al. Deep learning based multistep solar forecasting for PV ramp-rate control using sky images. *IEEE Transactions on Industrial Informatics*, 17(2):1397–1406, 2021. doi:10.1109/TII.2020.2987916.

[4] Papatheofanous, E. A. et al. Deep learning-based image regression for short-term solar irradiance forecasting on the edge. *Electronics*, 11(22), 2022. doi:10.3390/electronics11223794.

[5] Carreira Pedro, H., Larson, D., and Coimbra, C. A comprehensive dataset for the accelerated development and benchmarking of solar forecasting methods, 2019. doi:10.5281/zenodo.2826939.

[6] He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778. 2016. doi:10.1109/CVPR.2016.90.

[7] Liu, F. T., Ting, K., and Zhou, Z.-H. Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery From Data - TKDD*, 6:1–39, 2012. doi:10.1145/2133360.2133363.

[8] Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2019. doi:10.48550/arXiv.1711.05101.

[9] Ansel, J. et al. PyTorch 2: Faster machine learning through dynamic Python bytecode transformation and graph compilation. In *29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24)*. ACM, 2024. doi:10.1145/3620665.3640366.

[10] McInnes, L. et al. UMAP: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018. doi:10.21105/joss.00861.

# Evaluating the Impact of Feature Extractor Design on Multimodal Fatigue Detection Performance

**Anton Smoliński**[0000−0003−2216−8114], **Paweł Forczmański**[0000−0002−3618−9146],
**Adam Nowosielski**[0000−0001−7729−7867]

*West Pomeranian University of Technology in Szczecin*
*al. Piastów 17, 70-310 Szczecin, Poland*
*anton.smolinski@zut.edu.pl, pawel.forczmanski@zut.edu.pl,*
*adam.nowosielski@zut.edu.pl*

**Abstract.** *Driver fatigue is a major cause of road accidents, prompting the need for advanced Driver Monitoring Systems (DMS). Multimodal image data – combining RGB, thermal, and depth images – offers a richer representation of fatigue indicators. This study evaluates feature extraction strategies for multimodal fatigue recognition using a deep learning model with a BiLSTM classifier to capture temporal patterns. It compares a single unified feature extractor with modality-specific extractors. Results show that while modality-specific extractors yield marginal gains, a well-optimized single extractor achieves comparable accuracy with lower computational cost. These findings support the development of efficient, real-time fatigue detection systems.*
**Keywords:** *driver distraction, multimodal imaging, convolutional neural networks, bidirectional long short-term memory, road safety*

## 1. Introduction

This research focuses on recognizing driver fatigue indicators using multimodal image data. Driver fatigue significantly increases accident risk by impairing reaction time and decision-making. Modern *Driver Monitoring Systems* (DMS), now mandatory in new vehicles as part of *Advanced Driver Assistance Systems* (ADAS) [1], aim to detect such risks early.

Traditional fatigue detection methods rely on single-modal data (e.g., visible-light cameras), making them sensitive to lighting, occlusions, and individual differences. To improve robustness, this study leverages multimodal imaging, combining RGB, thermal, and depth data. Thermal imaging captures physiological changes; depth sensing detects subtle posture variations linked to drowsiness.

Previous work [2, 3, 4, 5] shows that multimodal fusion outperforms single-modal approaches. However, the optimal processing strategy – particularly whether to use shared or modality-specific feature extractors – remains an open question. This study investigates this issue to advance reliable and efficient fatigue detection systems.

To this end, data was acquired using RGB, thermal, and depth cameras, providing complementary information that enriches the representation of fatigue and distraction indicators [6].

Figure 1 shows example frames from the dataset across all modalities.



Figure 1. Selected frames from the developed database presenting participants with different physical characteristics, captured in various spectral bands

## 2. Dataset Overview

The dataset comprises 44 video sets of drivers, from which short clips were extracted to isolate specific fatigue-related behaviors. Each clip focuses on a single

indicator to enable precise labeling and analysis. Participant demographics are summarized in Table 1.

Table 1. Details of the Dataset – Demographic and Physical Characteristics

| Number of Recordings | Physical Characteristics | | | Demographics | |
|---|---|---|---|---|---|
| | Glasses | Beard or Moustache | Long Hair | Men | Women |
| 44 | 8 | 5 | 8 | 28 | 9 |

The dataset comprises six fatigue-related behavior classes, including 176 videos of neutral (no symptoms) behavior, 100 videos of yawning with an uncovered mouth, 83 with covered yawning, 105 showing unnatural blinking, 390 depicting head drooping, and 122 capturing drivers rubbing their eyes.

Clip durations range from 1 to 4 seconds, depending on behavior type. Short actions (e.g., blinking) result in shorter clips; prolonged behaviors (e.g., yawning) in longer ones.

Additional dataset details, including preprocessing and availability, are provided at `cvlab.zut.edu.pl` and described in prior works [2, 7].

## 3. Research Methodology

This study investigates whether using a *single feature extractor* across all modalities (RGB, Thermal, Depth) performs differently from using *modality--specific extractors*. The goal is to assess whether dedicated extractors yield measurable gains in multimodal fatigue detection.

The proposed end-to-end model (Figure 2) integrates all stages – feature extraction, classification, and fusion – into a single architecture [5]. To isolate the impact of extractor choice, all other model components remain fixed.

The pipeline includes the following stages. Three image modalities (RGB, Thermal, Depth) provide complementary information. Feature extraction is performed either by a shared extractor or dedicated extractors per modality, using pre-trained ImageNet networks truncated before their classification layer. The following convolutional networks were evaluated: ResNet50, VGGNet16/19, and InceptionV3. All backbones used frozen ImageNet weights to ensure robust feature extraction and reduce training time.

Extracted features are processed via BiLSTM networks, selected for their ability to model long-range temporal dependencies – essential in fatigue detection

Figure 2. Overview of the processing pipeline: multimodal input is processed via feature extractors (shared or modality-specific), classified using BiLSTM, and fused for final output

from video. Each modality is processed independently by a BiLSTM layer (128 units, *tanh*), with Gaussian noise and Dropout (0.25) regularization to improve robustness. The final softmax layer outputs six behavior classes.

Following prior work [5], a late fusion strategy is adopted: each modality is classified independently, with predictions fused at the decision level. This allows dynamic weighting of modalities based on data quality and confidence.

The core comparison evaluates two configurations: a single extractor across all modalities (e.g., ResNet50), and modality-specific extractors (e.g., ResNet50 for RGB, InceptionV3 for Thermal, VGGNet16 for Depth).

To promote reproducibility, we have published a benchmark dataset [2] and made an open-access multimodal fatigue dataset available at *cvlab.zut.edu.pl*.

# 4. Results

The dataset was split into training and testing subsets (75% / 25%). Due to the small sample size, a separate validation set was omitted to avoid instability from insufficient data. Training was optimized for efficiency and overfitting prevention, using the *RMSprop* optimizer, early stopping (1,500 iterations), and learning rate reduction (after 200 stagnant iterations). The results of the experiments are summarized in Tables 3 and 4.

Several Keras-based feature extractors were tested. Their input and output dimensions are listed in Table 2.

Table 2. Feature vector sizes for the utilized feature extractors

| Extractor | VGGNet16 | VGGNet19 | ResNet50 | InceptionV3 |
|---|---|---|---|---|
| **Input size** | $224 \times 224$ | $224 \times 224$ | $224 \times 224$ | $299 \times 299$ |
| **Output size** | 25088 | 25088 | 100352 | 131072 |

A comparison of classification results across different extractors shows that a higher number of extracted features does not necessarily improve classification accuracy. These results suggest that the representativeness and discriminative power of extracted features play a more crucial role than their sheer quantity. Extractors that effectively capture essential patterns in multimodal data – despite generating shorter feature vectors – may yield superior classification performance compared to more complex models producing larger but less informative feature representations. This underscores the importance of selecting an appropriate extractor based on feature quality rather than sheer dimensionality, as excessive feature length may lead to overfitting and increased computational costs without a proportional gain in classification accuracy.

As shown in Table 3, the best classification performance was achieved by ResNet50 (89.7%), followed by VGGNet16 (88.4%). Despite larger feature vectors, VGGNet19 and InceptionV3 underperformed, underscoring that feature quality matters more than dimensionality.

Table 3. Classification accuracy for different feature extraction models

| Feature Extractor | Training Accuracy | Testing Accuracy |
|---|---|---|
| ResNet50 | 1.000 | 0.897 |
| VGGNet16 | 1.000 | 0.884 |
| VGGNet19 | 0.996 | 0.875 |
| InceptionV3 | 0.996 | 0.839 |

Table 4 compares configurations using different extractors per modality. The best setup (VGGNet16 for RGB, InceptionV3 for thermal, ResNet50 for depth) reached 89.3%, only 0.4% below the best single-extractor result.

The minor performance gain from using modality-specific extractors does not outweigh the added complexity and computational cost. A single well-chosen ex-

Table 4. Classification accuracy for selected feature extractor combinations per modality

| Feature Extractor | | | Training Accuracy | Testing Accuracy |
|---|---|---|---|---|
| **RGB** | **Thermal** | **Depth** | | |
| InceptionV3 | InceptionV3 | ResNet50 | 0.998 | 0.871 |
| InceptionV3 | ResNet50 | VGGNet16 | 1.000 | 0.857 |
| InceptionV3 | VGGNet16 | ResNet50 | 0.993 | 0.844 |
| ResNet50 | InceptionV3 | VGGNet16 | 1.000 | 0.871 |
| ResNet50 | VGGNet16 | InceptionV3 | 0.983 | 0.871 |
| VGGNet16 | InceptionV3 | ResNet50 | 0.999 | 0.893 |
| VGGNet16 | ResNet50 | InceptionV3 | 0.985 | 0.875 |

tractor (ResNet50) offers nearly identical results with simpler deployment, which is advantageous for real-time applications.

This study focused exclusively on feature extractor configurations, not fusion strategies. While prior research suggests that late fusion may outperform early fusion [5], verifying this in the context of fatigue detection requires further experiments.

In summary, the results emphasize the importance of selecting effective extractors over increasing feature size or using modality-specific models. Future work should explore adaptive fusion mechanisms and further refinement of extraction strategies to enhance generalization and robustness.

# 5. Discussion

A significant portion of the cited works relates directly to the investigated method, including prior research by the authors. This is intentional, as the focus of this publication is to refine one component of the developed approach, not to compare it with fundamentally different existing fatigue detection methods.

Preliminary results without extensive fine-tuning, published in [5], already achieved competitive accuracy (67.1%–96.69%) [8, 9, 10], depending on modality and dataset. The present study builds on this by systematically evaluating feature extraction strategies.

Further optimization – specifically of data fusion methods – will be addressed in a subsequent paper. While splitting the research across publications may seem inconvenient, this structured approach enables a focused analysis of individual

components. The insights gained here may also benefit broader applications beyond fatigue detection.

## 6. Conclusions

This study evaluated the use of a single versus modality-specific feature extractor in multimodal driver fatigue detection. Results show that modality-specific extractors offer no substantial advantage over a well-optimized single extractor.

Using multimodal data (RGB, thermal, depth), the best configuration – ResNet50 as a unified extractor – achieved 89.7% accuracy, while the best modality-specific setup reached 89.3%. Thus, added model complexity is not justified.

Practically, a single extractor simplifies deployment and reduces computational cost. Future work will further optimize the system, aiming toward a robust, real-world prototype.

## References

[1] Regulation (EU) 2019/2144 of the European Parliament and of the Council of 27 November 2019 on type-approval requirements for motor vehicles and their trailers, and systems, components and separate technical units intended for such vehicles, as regards their general safety and the protection of vehicle occupants and vulnerable road users.

[2] Małecki, K., Forczmański, P., Nowosielski, A., Smoliński, A., and Ozga, D. A new benchmark collection for driver fatigue research based on thermal, depth map and visible light imagery. In R. Burduk, M. Kurzyński, and M. Woźniak, editors, *Progress in Computer Recognition Systems*, pages 295–304. Springer International Publishing, Cham, 2020. ISBN 978-3-030-19738-4. doi:10.1007/978-3-030-19738-4_30.

[3] Forczmański, P. and Smoliński, A. Supporting driver physical state estimation by means of thermal image processing. In M. Paszyński, D. Kranzlmüller, V. V. Krzhizhanovskaya, J. J. Dongarra, and P. M. Sloot, editors, *Computational Science – ICCS 2021*, pages 149–163. Springer International Publishing, Cham, 2021. ISBN 978-3-030-77977-1. doi:10.1007/978-3-030-77977-1_12.

[4] Nowosielski, A., Małecki, K., Forczmański, P., Smoliński, A., and Krzy-wicki, K. Embedded night-vision system for pedestrian detection. *IEEE Sensors Journal*, 20(16):9293–9304, 2020. doi:10.1109/JSEN.2020.2986855.

[5] Smoliński, A., Forczmański, P., and Nowosielski, A. Processing and integration of multimodal image data supporting the detection of behaviors related to reduced concentration level of motor vehicle users. *Electronics*, 13(13), 2024. ISSN 2079-9292. doi:10.3390/electronics13132457.

[6] Knapik, M., Cyganek, B., and Balon, T. Multimodal driver condition monitoring system operating in the far-infrared spectrum. *Electronics*, 13(17), 2024. ISSN 2079-9292. doi:10.3390/electronics13173502.

[7] Smoliński, A. Preprocessing multimodal image data for neural network training in driver distraction detection. In K. Bzdyra, editor, *Innowacje w elektronice, informatyce i inżynierii produkcji, vol. V*, pages 163–173. Wydawnictwo Uczelniane Politechniki Koszalińskiej, 2024. ISBN 978-83-7365-635-2. URL `https://dlibra.tu.koszalin.pl/publication/1990`.

[8] Siegfried, R., Yu, Y., and Odobez, J. A deep learning approach for robust head pose independent eye movements recognition from videos. *Proceedings of the 11th ACM Symposium on Eye Tracking Research & Applications*, 2019. doi:10.1145/3314111.3319844.

[9] Ma, X., Chau, L.-P., and Yap, K.-H. Depth video-based two-stream convolutional neural networks for driver fatigue detection. In *2017 International Conference on Orange Technologies (ICOT)*, pages 155–158. 2017. doi:10.1109/ICOT.2017.8336111.

[10] Majeed, F., Shafique, U., Safran, M., Alfarhood, S., and Ashraf, I. Detection of drowsiness among drivers using novel deep convolutional neural network model. *Sensors*, 23(21), 2023. ISSN 1424-8220. doi:10.3390/s23218741.

# Combining Local and Global Features Using Transformer-Based Architecture for Efficient Low-Light Enhancement

**Michał Wnuczyński**[0009−0000−0275−9916]

*Samsung R&D Institute Poland*
*wnuczynskimichal@gmail.com*

**Abstract.** *Transformer-based methods designed to enhance dark, underexposed images demand significant memory resources, with inference times escalating dramatically as image resolution increases. This makes them inefficient and impractical for common user and hardware-constrained devices. In this paper, we introduce an architecture that bridges the gap between efficiency and performance in transformers for low-light enhancement. Our model synergistically integrates local and global features using window-based and global transformer blocks. The proposed design effectively utilizes both spatial and channel relationships within the image, capturing essential features without the trade-off between quality and resource consumption.*

**Keywords:** *transformers, self-attention, low-light enhancement, image restoration*

## 1. Introduction

Low-light image enhancement is a challenging task in the computer vision domain that aims to improve the color, tone and contrast of images. This task is critical for various applications, including photography, surveillance, and medical imaging.

In recent years, the introduction of transformer-based methods to image processing tasks has marked a significant advancement in the domain. Transformers, originally designed for natural language processing, have demonstrated remarkable capabilities in capturing long-range dependencies and contextual information

from pictures. Several transformer-based methods have been designed specifically for image enhancement, achieving outstanding results [1, 2, 3, 4, 5, 6]. However, the primary drawback of using self-attention mechanisms is their high memory and computational complexity. These limitations present a significant barrier to their deployment in resource-constrained environments, making them either too slow, overly simplistic for the intended task, or even incapable of processing high-resolution images due to the memory required for execution.

In this paper, we propose an efficient transformer-based method for low-light enhancement that addresses these challenges. Ours approach combines two different self-attention mechanisms to effectively balance performance and computational efficiency. By integrating both local and global context self-attention calculations, ours method significantly reduces memory and computational demands while maintaining high-quality outputs competing with the state-of-the-art solutions.

## 2. Methodology

Our model employs a hierarchical encoder-decoder structure for enhancing low-light images, utilizing spatial and channel attention mechanisms. The overall architecture is illustrated in Figure 1. Initially, the image is extended from 3 to 16 channels through a 3x3 convolution with stride and padding set to 1 to capture detailed features. An embedding operation with a patch size of 1 converts the image into a vector, followed by layer normalization.



Figure 1. Overview of our proposed architecture

The image then passes through two Swin Transformer Blocks (STBs) with a window size of 8 and a half-window size shift between them to differentiate the calculations. This mechanism computes local attention within small spatial regions by dividing the image of size $h \times w$ into $M \times M$ patches reducing calculations from quadratic to linear:

$$\Omega(\text{self-attention}) = 4hwC^2 + 2(hw)^2C,$$
$$\Omega(\text{window self-attention}) = 4hwC^2 + 2M^2hwC.$$

In our case $M$ has the fixed size of 8 and $C$ denotes the dimensionality equal to 16. This approach requires only a fraction of the memory that would be needed if we were to calculate self-attention across the entire image. The output undergoes a patch unembedding and a 1x1 convolution to revert it to the original dimensionality.

The U-Net contains seven Channel Transformer Blocks (CTBs), which use layer normalization, global self-attention across channels, and a feed-forward network with an expansion factor of 2.66. Images are downsampled before each of the first four CTBs capturing multi-scale features, and outputs are concatenated after the fourth block, then processed through convolution and upsampling after each CTB.

The results from the U-Net are combined with the second STB's output, preserving local pixel information with hierarchical features. This combined tensor is re-embedded and processed by the final two STBs configured in the same manner as the first two. These final STBs refine the image by recalculating attention in local regions, ensuring that the enhancement is consistent, smooth and artifact-free.

The final steps involve unembedding the image, reducing its dimensionality back to the original 3 RGB channels, and passing it through a 1x1 convolution for final adjustments.

# 3. Experiments

We conducted experiments to evaluate the performance of the proposed method on the most commonly used dataset in the field – LOL [7], which consists of three versions: LOL-v1, LOL-v2-Real and LOL-v2-Synthetic. The datasets are the collection of low-light images with a resolution of 400×600 pixels and their corresponding enhanced versions. The v1 and v2-Real sets consist of real images, while v2-Synthethic is artificially generated based on them.

Our model is implemented using PyTorch. We employ the Adam optimizer with an initial learning rate of $10^{-3}$ decaying to $10^{-6}$ and batch size of 4. We optimize the model using Smooth L1 loss function with $\beta$ set to 0.1 and mean reduction type:

$$\ell(x, y) = L = \{l_1, \ldots, l_N\}^T,$$

$$l_n = \begin{cases} 0.5(x_n - y_n)^2/\beta, & \text{if } |x_n - y_n| < \beta, \\ |x_n - y_n| - 0.5 * \beta, & \text{otherwise,} \end{cases}$$

$$\ell(x, y) = \text{mean}(L),$$

where $x$ is the output from the model and $y$ it its corresponding target. This loss is more sensitive to the outliers, making our training process seamless. The models were trained separately on each dataset in the same manner for 3,000 epochs.

## 3.1. Evaluation

We quantitatively evaluated our method using two common metrics: Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM). These are standard metrics for assessing a model's output similarity to the target image. The performance of our method is presented in Table 1, which demonstrates the effectiveness of our approach in enhancing low-light pictures and confirms its robustness across different types of low-light images.

Table 1. Quantitative comparison on the LOL datasets in terms of PSNR and SSIM. Best results are in red and second best are in blue.

| Method | LOL-v1 | | LOL-v2-Real | | LOL-v2-Synthetic | |
|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| Restormer [5] | 26.68 | 0.853 | 26.12 | 0.853 | 25.43 | 0.859 |
| LLFormer [3] | 25.75 | 0.823 | 26.19 | 0.819 | 28.00 | 0.927 |
| MIRNet [6] | 24.14 | 0.830 | 27.17 | 0.865 | 25.96 | 0.898 |
| HWMNet [2] | 24.24 | 0.850 | — | — | — | — |
| EMNet [4] | 25.37 | 0.868 | 23.42 | 0.874 | — | — |
| Retinexformer [1] | 27.18 | 0.850 | 27.71 | 0.856 | 29.04 | 0.939 |
| Ours | 26.75 | 0.855 | 27.24 | 0.865 | 29.56 | 0.953 |

## 3.2. Visual Comparison

To further evaluate the performance of the model, we conducted a side-by-side visual comparison of the enhanced images by several models. The goal was to assess the quality of the results in terms of sharpness, color fidelity, artifact suppression, and overall realism (see Figure 2). In the presented example our solution excelled in enhancing the white door while preserving the contrast between the upper and lower sections. Most importantly, it avoided overexposing the design on the bottom door and maintained natural colors.



Figure 2. Visual comparison of existing methods on LOL-v1 dataset

## 3.3. Comparative Analysis

We compared our method with the existing state-of-the-art transformer-based architectures in terms of memory complexity and processing time on common image resolutions using NVIDIA A100 82GB graphics card. The summarized results are presented in Table 2. Each result is an average of 16 runs to avoid any singular deviation. Our method significantly outperforms other architectures, achieving approximately a 60% memory reduction compared to the most efficient MIRNet [6] and around three times faster computation time compared to the best performing RetinexFormer [1] highlighting its efficiency.

## 3.4. Mobile Deployment

To demonstrate the practicality of the proposed method, we deployed our model to a mobile device using PyTorch Lite. The tests were conducted on a Sam-

Table 2. Peak memory usage and average inference time. Best results are in red and second best in blue. n/a - not available (Out Of Memory). Resolutions are given in [height] px × [width] px.

| Method | 720×1280 | | 1080×1920 | | 1440×2560 | | 2160×3840 | |
|---|---|---|---|---|---|---|---|---|
| | GB | s | GB | s | GB | s | GB | s |
| LLFormer [3] | 15.47 | 1.02 | 47.75 | 3.76 | n/a | n/a | n/a | n/a |
| Retinexformer [1] | 20.83 | 0.25 | 44.31 | 0.56 | 77.25 | 1.00 | n/a | n/a |
| Restormer [5] | 12.77 | 0.28 | 24.25 | 0.58 | 41.74 | 1.20 | n/a | n/a |
| HWMNet [2] | 11.41 | 0.33 | 22.23 | 0.96 | 37.57 | 1.78 | n/a | n/a |
| EMNet [4] | 8.91 | 0.43 | 17.70 | 0.94 | 30.03 | 1.71 | n/a | n/a |
| MIRNet [6] | 6.54 | 0.97 | 13.69 | 2.10 | 22.75 | 3.81 | 42.72 | 9.17 |
| Ours | 3.68 | 0.09 | 6.13 | 0.19 | 9.51 | 0.29 | 19.23 | 0.98 |

sung Galaxy S24 Ultra with 12 GB of RAM. Table 3 presents the results of our measurements. The model is able to perform inferences on images of resolutions up to 2160 by 3200 pixels without applying any additional improvements.

Table 3. Peak RAM memory usage and inference time of our model using PyTorch Lite on a Samsung Galaxy S24 Ultra

| Image resolution | RAM [MB] | Time [s] |
|---|---|---|
| 720×1280 | 1 270 | 11.62 |
| 1080×1920 | 2 875 | 30.96 |
| 1440×2560 | 4 931 | 63.80 |
| 2160×3200 | 8 552 | 152.71 |

## 4. Conclusion

In this paper, we propose a memory-efficient Transformer-based method for Low-Light Image Enhancement, which utilizes relationships between local spatial features and global channel features. This approach achieved high-quality results while significantly reducing memory usage and inference time. Future work could explore advanced denoising techniques, video enhancement, and integration with other image processing tasks, as well as include additional optimization techniques for mobile deployment to reduce the computation time.

# CHAPTER 6

# Uncertainty in Artificial Intelligence

Track Chairs:

- prof. Dominik Ślęzak – University of Warsaw

- prof. Beata Zielosko – University of Silesia in Katowice

- prof. Agnieszka Jastrzębska – Warsaw University of Technology

# Fuzzy-Rough Approach
# to the Feature Extraction

**Zofia Matusiewicz**[0000−0003−0523−0983], **Teresa Mroczek**[0000−0002−6064−9528]

*University of Information Technology and Management*
*Sucharskiego 2, 35-225 Rzeszow, Poland,*
*{zmatusiewicz, tmroczek}@wsiz.edu.pl*

**Abstract.** *This paper presents an extended feature extraction method in decision tables using fuzzy relational equations and inequalities. The approach focuses on identifying the greatest solution of $A \circ x = d$ using ascending, left-handed continuous operations of $*$ that satisfy the condition $1 * 0 = 0$. A simplified binarized reduced matrix is constructed for the determined greatest solution, which helps efficiently determine attributes' relevance. Experimental validation on multiple datasets shows that the proposed reduction method achieves accuracy comparable to standard techniques such as recursive feature elimination, random forest and principal component analysis, while offering the advantage of linear time complexity.*

**Keywords:** *computer aided diagnosis, skin lesion classification, malignant melanoma, bayesian networks, rough set theory*

## 1. Introduction

Several reasons justify addressing the issue of feature extraction [1], including improving algorithm performance, facilitating data understanding by providing deeper insights, reducing storage requirements to reduce costs, and increasing simplicity by enabling faster and more efficient models.

Fuzzy relational equations and inequalities [2, 3], along with rough set theory [4, 5], provide the tools for efficiently discovering relationships in data sets, determining the degrees of these relations, and identifying the most important attributes. Recently, a new attribute reduction method with linear time complexity for identifying the strongest dependencies between attributes and decisions, utilizing the mentioned tools, was introduced in [6]. The reduction method was successfully

applied for left-continuous triangular norms. In this paper, an extension of this method is proposed. Determining the greatest desired solution can be achieved by using increasing, left-continuous operations $*: [0, 1]^2 \rightarrow [0, 1]$ that satisfy the condition $1 * 0 = 0$.

In [6], it showed that the use of some triangular norms guarantees better results than for us. Following this observation, the research was extended to the whole family of operations that can be used in relational equations and inequalities. In summary, the aim of extending the method and undertaking this research is to establish the properties of operations that give the best results.

## 2. Method Description

We consider a data set in the form of a decision table, as a system consisting of a non-empty, finite set of *Cases*, *Attributes*, and *Decision* [4]. For a given $A$ and $d$ the *fuzzy relational equations* are identified as follows:

$$A \circ x = d, \tag{1}$$

where the following fuzzy relations are known:

$$A : Cases \times Attributes \rightarrow [0, 1], \quad d : Cases \times Decisions \rightarrow [0, 1],$$

and the unknown relation is

$$x : Attributes \times Decisions \rightarrow [0, 1].$$

The value $u_j(A, d, *) = \min_{i \in M}(a_{ij} \overset{*}{\rightarrow} d_i)$ is the greatest degree of $j - th$ attribute-decision relationship, where $a \overset{*}{\rightarrow} b = \max\{t \in [0, 1] : a * t \leqslant b\}$. In this type of application only the greatest solution of fuzzy relational equations $A \circ x = d$ is desirable. The greatest desired solution can be determined by using increasing, left-continuous operations $*: [0, 1]^2 \rightarrow [0, 1]$ that satisfy the condition $1 * 0 = 0$, what has been proven in [7]. If in at least one row in the matrix $A$, all elements are less than the corresponding value of the vector $d$, then the system of equations $A \circ x = d$ has no solution. Instead, we can solve the greatest solution $u$ of the inequality $A \circ x \leq d$, and consider the system $A \circ x = d_{new}$, where $d_{new} = A \circ u$. However, using only relational equations did not guarantee success in determining the greatest solution, so methods used in rough sets were used. It may be that there is 1 at the j-th position of the greatest vector, but that a given element of the vector

never participates in obtaining the decision value $d_i$ (that is, $a_{ij} * u_j < b_i$ for each j), so binarized, reduced matrix $A''$ of equation $A \circ x = d$ with respect to solution $u$ and a vector $dc$ were introduced:

$$a''_{ij}(u) = \begin{cases} 1, & \text{if } a_{ij} * u_j = d_i, \\ 0, & \text{otherwise,} \end{cases} \qquad dc_j = \frac{\sum_{i \in M}(a''_{ij})}{m} \cdot 100\%,$$

for all $i \in \{1, \ldots, m\}$ and $j \in \{1, \ldots, n\}$.

In this approach, we use an increasing, left-continuous operations $* : [0,1]^2 \to [0,1]$ that satisfy the condition $1 * 0 = 0$. Such assumptions guarantee the existence of a solution of the greatest solution [7]. Example of such an operation is presented below:

$$F_{TY}(a,b) = \begin{cases} b^{\frac{1}{a}}, & \text{if } ab > 0, \\ 0, & \text{otherwise,} \end{cases} \qquad F_{SM} = \sqrt{xy}, \quad F_{PM}(x,y) = x \min(x,y),$$

where $F_{TY}$ denotes Yager's pseudo t-norm, $F_{SM}$ represents the geometric mean, and $F_{PM}$ is a monotonic function. These operations have corresponding induced implications, as follows:

$$a \xrightarrow{TY} b = \begin{cases} 1, & a = 0, \\ 0, & a > 0, b = 0, \\ b^a, & \text{otherwise,} \end{cases} \quad a \xrightarrow{SM} b = \begin{cases} \frac{b^2}{a}, & a > b^2, \\ 1, & \text{others,} \end{cases} \quad a \xrightarrow{PM} b = \begin{cases} \frac{b}{a}, & a^2 > b, \\ 1, & \text{others.} \end{cases}$$

To determine the relevance of the attributes, the simplified reduced matrix, denoted by $A''$, was constructed, where the elements $a''_{ij}$ are assigned a value of 1 if $F(a_{ij}, u_j) = b_i$, and 0 otherwise.

An abbreviated procedure for feature extraction is presented as follows:

1. Computation of the greatest solution $u$ (see Algorithm 1). If $A \circ u < d$, then substitute $d := A \circ u$;

2. Determination of the simplified reduced matrix $A''$ (see Algorithm 2);

3. Calculation of the vector $dc$ and indication of the significance level of each attribute (see Algorithm 3);

4. Extraction features (see Algorithm 4).

---

**Algorithm 1** Greatest Solution Procedure

---

1: **Input:** $A_{m \times n}$ (matrix), $d_{m \times 1}$ (vector), $\overset{*}{\to}$ (binary operation)

2:

3: **Output:** $u \leftarrow \mathbf{0}$

4: **for** $j \leftarrow 1$ to $n$ **do**

5:     $u_j \leftarrow 1$

6:     **for** $i \leftarrow 1$ to $m$ **do**

7:         $u_j \leftarrow \min((a_{ij} \overset{*}{\to} d_i), u_j)$

8:     **end for**

9: **end for**

10: **return** $u$

---

---

**Algorithm 2** Computation of $A''$

---

1: **Input:** $A_{m \times n}$ (matrix), $u_{n \times 1}, d_{m \times 1}$ (vectors), $F$ (function)

2: **Output:** $A''_{m \times n}$ (matrix)

3: Initialize $A''$ as a zero matrix of size $m \times n$

4: **for** $i \leftarrow 1$ to $m$ **do**

5:     **for** $j \leftarrow 1$ to $n$ **do**

6:         **if** $F(a_{ij}, u_j) = d_i$ **then**

7:             $A''_{ij} \leftarrow 1$

8:         **else**

9:             $A''_{ij} \leftarrow 0$

10:         **end if**

11:     **end for**

12: **end for**

13: **return** $A''$

---

# 3. Conclusion

The proposed modification was verified on the data sets available in the Machine Learning Repository, University of California, Irvine, apart from *ILQ*. The *ILQ* data set is provided in [8]. All data sets, original and reduced, were input to the CART decision tree system [9]. The results were compared with Recursive Feature Elimination (RFE), Random Forest (RF), and Principal Component Analysis (PCA). The default parameters for all methods were applied. The accuracy was evaluated by ten-fold cross-validation, the results are presented in Table 1.

---

**Algorithm 3** Determining the relevance of features

---

1: **Input:** $A''_{m \times n}$ (matrix)
2: **Output:** $dc_{n \times 1}$ (vector)
3: **for** $j \leftarrow 1$ to $n$ **do**
4:     $pom_j \leftarrow 0$
5:     **for** $i \leftarrow 1$ to $m$ **do**
6:         $pom_j \leftarrow pom_j + (A''_{ij})$
7:     **end for**
8:     $dc_j = 100\% \cdot pom_j/m$
9: **end for**
10: **return** $dc$

---

**Algorithm 4** Removal of irrelevant attributes

---

1: **Input:** $dc_{n \times 1}$ (vector), *threshold* (value), *Attributes* (set)
2: **Output:** *Attributes$_{removed}$* (set)
3: *Attributes$_{selected}$* $\leftarrow \emptyset$
4: **for** $j \leftarrow 1$ to $n$ **do**
5:     **if** $dc_j < threshold$ **then** *Attributes$_{removed}$* $\leftarrow$ *Attributes$_{removed}$* $\cup$ *Attributes$_j$*
6:     **end if**
7: **end for**
8: **return** *Attributes$_{removed}$*

---

The accuracy of the developed reduction method is similar to the accuracy obtained using well-known reduction algorithms, while its time complexity is linear. When differences occur, there are statistically insignificant.

In this variant, we use a binarized, reduced matrix and a left continuous operation $*$ to extract any information system in the form of a decision table. Additional improvements are planned for the future:

1. the relevance of an attribute based on the coefficients of the reduced matrix will be determined; and

2. multi-decision problems will be considered.

Table 1. Accuracy rate

| Data Set | $F_{TY}$ | $F_{SM}$ | $F_{PM}$ | RFE | RF | PCA |
|---|---|---|---|---|---|---|
| Boston | 0.61 | 0.71 | 0.72 | **0.82** | 0.81 | 0.73 |
| Happiness | **0.97** | **0.97** | **0.97** | **0.97** | **0.97** | **0.97** |
| HAR | 0.78 | **0.87** | 0.86 | **0.87** | 0.65 | 0.86 |
| ILQ | 0.62 | 0.64 | **0.65** | 0.63 | **0.65** | **0.65** |
| Weather | 0.53 | 0.61 | 0.62 | **0.63** | **0.63** | 0.62 |
| Waveform | 0.65 | 0.76 | **0.77** | 0.76 | 0.76 | 0.73 |

# References

[1] Guyon, I. et al. *Feature Extraction.* Springer, Heidelberg, 2006.

[2] Sanchez, E. Resolution of composite fuzzy relation equations. *Information and Control*, 30(1):38–48, 1976.

[3] Peeva, K. and Kyosev, Y. *Fuzzy Relational Calculus: Theory, Applications and Software*. World Scientific, 2004.

[4] Pawlak, Z. Rough sets. *International Journal of Computer & Information Sciences*, 11:341–356, 1982.

[5] Pawlak, Z. and Skowron, A. Rudiments of rough sets. *Information Sciences*, 177:3–27, 2007.

[6] Matusiewicz, Z. and Mroczek, T. Attribute reduction method based on fuzzy relational equations and inequalities. *International Journal of Approximate Reasoning*, 178:109355, 2025. doi:10.1016/j.ijar.2024.109355.

[7] Matusiewicz, Z. and Drewniak, J. Increasing continuous operations in fuzzy max-$*$ equations and inequalities. *Fuzzy Sets Syst*, 232:120–133, 2013.

[8] Krawczyk–Suszek, M., Gaweł, A., and Kleinrok, A. Correlation of ageing with health related-quality of life of patients in 13 groups of disease in Poland. *Scientific Reports*, 14:26404, 2024. doi:10.1038/s41598-024-78253-1.

[9] Breiman, L., Friedman, J., Olshen, R., and Stone, C. *Classification and Regression Trees*. Chapman and Hall/CRC, 1984. doi:10.1201/9781315139470.

# Uncertainty of Aggregation: Investigating Dissimilarity between Honeycomb-Based Polygonal Chains

**Grzegorz Mos**[0000−0002−3266−3760]

*Faculty of Science and Technology*
*University of Silesia in Katowice, Katowice, Poland*
*grzegorz.mos@us.edu.pl*

**Abstract.** *Aggregation does not always yield a single, well-defined result but may deliver multiple outcomes. This uncertainty raises inquiries about the role of object dissimilarity in the aggregation process. This study explores how dissimilarity influences aggregation by utilizing medoids to obtain results that are as distinctively different from the specified set of entities as possible. Our approach provides a fresh perspective on managing uncertainty in aggregation and proposes an alternative aggregation method. We consider honeycomb-based polygonal chains, which are represented with binary sequences. We employ string similarity metric, i.e., Hamming distance, to determine medoids and centroids. Additionally, we introduce extensions of this metric, specifically adapted to the properties of the considered objects. We gain insight into the relationship between dissimilarity and aggregation stability by analyzing how different metrics and approaches influence the aggregation outcome. Our findings contribute to the broader understanding of aggregation in complex systems, mainly where structural variability is critical. This research has potential applications in knowledge distribution, pathfinding, preferences aggregation, and other fields where aggregation must account for inherent uncertainties.*

**Keywords:** *dissimilarity, honecomb-based polygonal chain, hexagonal grid, medoid, Hamming distance*

## 1. Introduction

Graph aggregation has been extensively studied in recent years, primarily due to the inherent complexity of the structures involved. The properties of aggregated objects should be considered. Therefore, extending aggregation methods is

required (cf. [1]). The basic notation and terminology for graphs, an outline of aggregation rules, and several applications, such as in multiagent systems, have been presented (cf. [2, 3]). However, representing a graph with a set of vertices, edges and edge labels poses significant challenges in developing general aggregation methods that can be universally applied across various problem solutions and scientific fields. Thus, this paper considers the specified type of graph, which facilitates their aggregation and the interpretation of the results. This approach is novel and provides more straightforward methods for handling complex objects.

Honeycomb-based polygonal chains constitute a specialized class of graphs that arise in various scientific disciplines, including mathematics, informatics, physics, and chemistry (cf. [4, 5, 6]). Their highly organized structure enables a representation through binary sequences, simplifying the aggregation problem while preserving essential structural characteristics. This representation facilitates the application of string-based metrics for comparing these structures and enables aggregation using medoids.

The medoid is a well-established concept in data aggregation (cf. [7]). However, it does not guarantee a uniquely determined aggregation result. Consequently, extended aggregation approaches based on medoids must be developed, or post-aggregation steps must be introduced to ensure a definitive outcome. We propose two methods to address this issue. The first extends the Hamming distance by incorporating weights assigned to specific indices of the compared structures. The second introduces a constraint set comprising binary sequences to which the final aggregation result, or the aggregated elements, should be as dissimilar as possible.

We present the necessary preliminaries, introducing fundamental definitions for representing honeycomb-based polygonal chains as binary sequences. Additionally, we revisit the medoid concept as a standard data aggregation method. We then extend the definition of the Hamming distance to allow for comparing binary sequences of varying lengths. Next, we examine the challenge of aggregation uncertainty and propose two distinct approaches to mitigate it. The first utilizes the modified metric to refine aggregation, while the second defines a constraint set as a reference for aggregation selection. We subsequently formalize the general case of aggregation using the medoid under a set constraint. Finally, we summarize our findings and outline directions for future research.

## 2. Preliminaries

Honeycomb-based structures are connected graphs in which every two adjacent edges form a 120° angle. Their geometry is derived from a hexagonal grid. In this work, we focus on the most straightforward structures, where precisely two vertices have a degree of 1, while all remaining vertices have a degree of 2 – forming a polygonal chain. Consequently, these structures can be effectively represented using binary sequences (cf. [8]). We present only the theoretical foundations necessary for understanding and conducting further research.

The rotation matrix $\text{Rot}(\varphi)$ about the point $(0,0)$ by an angle $\varphi \in \mathbb{R}$ is given by the formula

$$\text{Rot}(\varphi) = \begin{bmatrix} \cos(\varphi) & -\sin(\varphi) \\ \sin(\varphi) & \cos(\varphi) \end{bmatrix}.$$

Let $n \in \mathbb{N}$, let $v_k = (x_k, y_k) \in \mathbb{R}^2$ for $k \in \{0, 1, \ldots, n\}$ and let $S = (v_0, v_1, \ldots, v_n)$ be such a tuple of vertices, that $(x_k, y_k)$ and $(x_{k+1}, y_{k+1})$ are connected with edge for $k \in \{0, 1, \ldots, n-1\}$. We will say that $S$ is a honeycomb-based polygonal chain whenever

$$\begin{bmatrix} x_k \\ y_k \end{bmatrix} = \begin{bmatrix} x_{k-1} \\ y_{k-1} \end{bmatrix} + \text{Rot}\,(\pm 60°) \cdot \begin{bmatrix} x_{k-1} - x_{k-2} \\ y_{k-1} - y_{k-2} \end{bmatrix},$$

for every $k \in \{2, 3, \ldots, n\}$. If $S$ is a honeycomb-based polygonal chain, then we will say that a binary sequence $B = b_0 b_1 \ldots b_{n-2}$ represents $S$ whenever

$$\begin{bmatrix} x_k \\ y_k \end{bmatrix} = \begin{bmatrix} x_{k-1} \\ y_{k-1} \end{bmatrix} + \text{Rot}\left((-1)^{b_{k-2}} \cdot 60°\right) \cdot \begin{bmatrix} x_{k-1} - x_{k-2} \\ y_{k-1} - y_{k-2} \end{bmatrix},$$

for every $k \in \{2, 3, \ldots, n\}$.

Intuitively, we mark every angle created by the pair of connected edges with 0 when it is left-hand-side and 1 when it is right-hand-side concerning the order of considered edges (see Figure 1).



Figure 1. Honeycomb-based polygonal chain represented by the sequence 101000 with the asterisks at the beginning and marked bits on every angle between connected edges

This approach allows us to represent this type of graph in a much simpler form using binary sequences, enabling applying a broad range of string metrics. Consequently, the potential applications are significantly expanded. In traditional graph notation, we need to consider a set of edges and a labeling function for these edges, which cannot be utilized in the solutions presented in this paper.

Medoid is a well-defined and studied method of aggregating data (cf., [7]). Assume that $n \in \mathbb{N}$, $\mathfrak{B}$ is a family of all binary sequences, $\mathcal{B} = \{B_1, B_2, \ldots, B_n\}$ is a set of binary sequences, and $d \colon \mathfrak{B}^2 \to [0, \infty)$ is a metric. We will say that a binary sequence $M \in \mathcal{B}$ is a medoid whenever the distance between $B$ and all elements of $\mathcal{B}$ is the lowest possible, i.e.,

$$M = \arg\min_{B \in \mathcal{B}} \sum_{i=1}^{n} d(B, B_i).$$

# 3. Results

## 3.1. Dissimilarity between Structures

As mentioned in the introduction, we will consider Hamming distance in this paper. It is a very intuitive string metric that we can define straightforwardly.

Suppose we have two binary sequences $B = b_1 b_2 \ldots b_n$ and $C = c_1 c_2 \ldots c_n$ for some $n \in \mathbb{N}$. Then the Hamming distance $d_H$ is given by the formula:

$$d_H(B, C) = \sum_{i=1}^{n} |b_i - c_i|.$$

The Hamming distance is traditionally defined only for sequences of equal length. We extend this definition to enhance its applicability for comparing sequences of different lengths, thereby extending its potential and making it more useful for a broader range of applications. Suppose we have two binary sequences $B = b_1 b_2 \ldots b_n$ and $C = c_1 c_2 \ldots c_m$ for some $m, n \in \mathbb{N}$. Let's $U = \max\{m, n\}$ and $L = \min\{m, n\}$. Then the extended Hamming distance $d_{eH}$ is given by the formula:

$$d_{eH}(B, C) = U - L + \sum_{i=1}^{L} |b_i - c_i|.$$

Moreover, since binary sequences represent honeycomb-based polygonal chains, the position at which two compared sequences differ can significantly impact the overall visual disparity between these chains. Therefore, we introduce a weighting

scheme that adjusts the difference between corresponding bits based on their positional distance from the beginning of the sequence. This approach ensures that discrepancies occurring at different positions contribute to the distance measure in a contextually meaningful way. The extended weighted Hamming distance $d_{wH}$ is given by the formula:

$$d_{wH}(B, C) = \sum_{i=1}^{L} (U - i) \cdot |b_i - c_i| + \sum_{i=L+1}^{U} (U - i).$$

**Theorem 1.** *Both $d_{eH}$ and $d_{wH}$ are metrics.*

The above result holds significant value in the research, as the metrics satisfy crucial properties fundamental to analyzing and understanding the problem, which provides mathematical consistency, computational efficiency, and practical applicability. They offer a robust, widely applicable framework that balances theoretical elegance with real-world utility. The identity of indiscernibles ensures that if the distance between two sequences is zero, those sequences must be the same, guaranteeing that distinct sequences have a positive distance. The symmetry ensures that the order of measuring sequences does not matter, so the distance from the sequence $B$ to the sequence $C$ is the same as from the sequence $C$ to the sequence $B$. Finally, the triangle inequality ensures that the direct distance between two sequences is never greater than the distance taken through an intermediate sequence, reflecting the shortest path principle.

### 3.2. Uncertainty of Aggregation

As mentioned earlier, a given set of objects may yield multiple medoids (cf. [7]), leading to aggregation uncertainty. This issue can be addressed through various approaches. Here, we present two methods to overcome this uncertainty.

Consider a set of binary sequences $\mathcal{B} = \{0001, 0010, 0100\}$ (cf., Figure 2). One way to resolve aggregation uncertainty is by selecting an appropriate metric. For instance, using the extended weighted Hamming distance $d_{wH}$, we compute the pairwise distances: $d_{wH}(0001, 0010) = 3$, $d_{wH}(0001, 0100) = 4$ and $d_{wH}(0010, 0100) = 5$. Since 0001 minimizes the sum of distances to all other sequences in $\mathcal{B}$, it is chosen as the medoid.

However, if we instead consider the standard Hamming distance $d_H$, each sequence in $\mathcal{B}$ is equidistant to the others, meaning that any of them can serve as

Figure 2. Honeycomb-based polygonal chains represented by the sequences 0001, 0010 and 0100 from left to right respectively, with asterisks at the beginning

a medoid. To refine the selection process, we introduce a constraint set $\mathcal{B}^*$ containing binary sequences that the final medoid should be as dissimilar to as possible. Suppose we define $\mathcal{B}^* = \{0101\}$. The corresponding Hamming distances are $d_H(0101, 0001) = 1$, $d_H(0101, 0010) = 3$ and $d_H(0101, 0100) = 1$. Since 0010 is the most dissimilar to $\mathcal{B}^*$, it is selected as the final aggregation object. These approaches can be applied to the general case in graph theory, where metrics between graphs are defined. However, constructing metrics for arbitrary graphs with intuitive explanations remains highly challenging and poorly studied.

### 3.3. Set Constraint as Uncertainty Solution

Consider a set $\mathcal{B} = \{B_1, B_2, \ldots, B_n\}$ of binary sequences for some $n \in \mathbb{N}$ and a metric $d$. Assume that $\mathcal{B}^* = \{B_1^*, B_2^*, \ldots, B_{n^*}^*\}$ for some $n^* \in \mathbb{N}$. If $M$ is a medoid for $\mathcal{B}$, then we can pick this medoid $M_1^*$ from $M$ which is as dissimilar as possible to the elements of the set $\mathcal{B}^*$. Thus, $M_1^*$ is of the form:

$$M_1^* = \arg\max_{B \in M} \left\{ \sum_{i=1}^{n^*} \mathrm{d}(B, B_i^*) \right\}.$$

On the other hand, we can pick this binary sequence $M_2^*$ from $\mathcal{B}$ for which subtract dissimilarities to the elements of $\mathcal{B}$ and dissimilarities to the elements of $\mathcal{B}^*$ is as lowest as possible. Thus, $M_2^*$ is of the form:

$$M_2^* = \arg\min_{B \in \mathcal{B}} \left\{ \sum_{i=1}^{n} \mathrm{d}(B, B_i) - \sum_{i=1}^{n^*} \mathrm{d}(B, B_i^*) \right\}.$$

Assume that $\mathcal{B} = \{0001, 0010, 1100\}$ and $\mathcal{B}^* = \{0011\}$ and consider the Hamming distance $d_H$ (see Figure 3). Then $M = \{0001, 0010\}$, $M_1^* = \{0001, 0010\}$ and $M_2^* = \{1100\}$. Thus, there is no relation between $M_1^*$ and $M_2^*$. Moreover, the

binary sequence in $M_2^*$ is not present in the set $M$. Hence, the aggregating approach must depend on the context of the problem, or both these methods can be used to determine the final result.



Figure 3. Honeycomb-based polygonal chains represented by the sequences, respectively, from left to right, 0001, 0010, 1100, and 0011, with asterisks at the beginning. Aggregated structures on the left are separated with the vertical line from the constraint structure on the right.

Indeed, these approaches can be applied to the general case in graph theory with specified metrics, but this presents a challenging problem.

# 4. Summary and Future Work

We have introduced an extended version of the Hamming distance that incorporates index-based weighting of binary sequence differences. This extension allows for the aggregation of honeycomb-based polygonal chains, represented as binary sequences of varying lengths, by applying the medoid concept. The new metric has a more significant impact on achieving a unique aggregation result than the essential Hamming distance. Furthermore, we proposed a novel aggregation method that utilizes a constraint set to refine the selection of the final aggregation object. Together, these two approaches significantly influence the aggregation process, effectively reducing the number of cases where the aggregation result is not unique.

Future work will focus on analyzing the theoretical properties of the proposed approaches and exploring their practical applications. The introduced methods can also be considered a new aggregation problem where the constraint set is already given, which requires further study.

# References

[1] Gągolewski, M. *Data Fusion: Theory, Methods, and Applications.* No. 7 in *Monograph Series: Information Technologies: Research and their Interdisciplinary Applications.* Institute of Computer Science, Polish Academy of Sciences, Warsaw, 2015.

[2] Bonzon, E. *Interaction in Multiagent Systems. PhD Thesis.* Université Paris Cité, 2024.

[3] Endriss, U. and Grandi, U. Graph aggregation. In *Companion Proceedings of the The Web Conference 2018 (WWW '18)*, pages 447–450. ACM, 2018. doi:10.1145/3184558.3186231.

[4] Coleman, M. M. *Fundamentals of Polymer Science: An Introductory Text.* Second Edition. Routledge, 2019.

[5] Wang, Z. Recent advances in novel metallic honeycombstructure. *Composites Part B: Engineering*, 166:731–741, 2019. doi:10.1016/j.compositesb.2019.02.011.

[6] Langner, J., Witek, H. A., and Moś, G. Zhang–zhang polynomials of multiple zigzag chains. *MATCH Commun. Math. Comput. Chem*, 80(1):245–265, 2018.

[7] Kaufman, L. and Rousseeuw, P. J. *Finding Groups in Data: An Introduction to Cluster Analysis.* Wiley, 1990.

[8] Moś, G. Honeycomb-based polygonal chains aggregation functions. In D. Ciucci et al., editors, *International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems; IPMU 2022.* Communications in Computer and Information Science (CCIS), Springer, 2022.

# Fuzzy Preference Templates in Multi-Criteria Decision-Making

**Leszek Rolka**[0000−0003−0083−8893]

*Rzeszów University of Technology*
*Faculty of Mechanical Engineering and Aeronautics*
*Al. Powstańców Warszawy 8, 35-959 Rzeszów, Poland*
*leszekr@prz.edu.pl*

**Abstract.** *This paper presents an approach to modeling the subjective preferences that are taken into account in the process of finding the best alternative described by fuzzy criteria. We introduce the notions of fuzzy preference and anti-preference templates suitable for representing the ideal positive and the ideal negative preferences of a decision-maker. The templates can be transformed into subsets of corresponding fuzzy linguistic labels. Furthermore, we propose a measure for evaluating the consistency of the preference model of a decision-maker. The details of the presented approach are illustrated by examples.*
**Keywords:** *fuzzy sets, linguistic labels, multi-criteria decision-making*

## 1. Introduction

The problem of finding the best alternative (object) is a nontrivial task when several contradictory criteria ought to be taken into account. It becomes even more complex in the case of criteria that are not numerical, but represent vague or subjective concepts that are often used in everyday life. Several methods for finding an optimal solution of multi-criteria decision-making have been introduced over the recent decades. The inspiration for proposing the approach presented in this paper was a group of methods that are based on evaluating the distance of alternatives to ideal solutions, such as TOPSIS [1] and VIKOR [2]. A review of TOPSIS-based applications is presented in [3]. Fuzzy variants of these methods were investigated by various researchers, see, e.g., [4, 5, 6]. In contrast to previous work,

we avoid using an arithmetic-based fuzzy approach. Instead, we propose a logic-oriented description and evaluation of fuzzy alternatives, which seems to be more suitable for evaluating alternatives characterized by symbolic or fuzzy criteria. To this end, we assume that the ideal positive and negative solutions should be provided by a decision-maker in the form of fuzzy preference templates. This is the fundamental difference of our approach compared to TOPSIS and VIKOR methods, which determine (calculate) the ideal positive and negative solutions based on the analyzed real alternatives. Furthermore, we want to develop a method that is completely fuzzy in nature. It should be noted that the popular optimization methods, including TOPSIS and VIKOR, were introduced to solve multi-criteria decision-making problems with numerical criteria. Various fuzzified versions of those methods proposed more recently are mostly based on using fuzzy numbers instead of crisp values. However, our goal is to build from scratch a fuzzy approach that will be suitable for solving optimization problems using only subjective linguistic criteria.

## 2. Alternatives with Fuzzy Criteria

In the following, we consider a problem of selecting the best solution from a finite set of elements (alternatives) that will be described by subjective fuzzy linguistic attributes (criteria) only. We apply the notion of fuzzy information system $\text{FDS} = \langle U, A, \mathbb{V}, f \rangle$, where:

$U$ – is a nonempty finite universe of elements,

$A$ – is a finite set of fuzzy criteria,

$\mathbb{V}$ – is a set of linguistic values of criteria, $\mathbb{V} = \bigcup_{a \in A} \mathbb{V}_a$,

$\quad \mathbb{V}_a$ is the set of linguistic values of a criterion $a \in A$,

$f$ – is an information function, $f : f(x, V) \in [0, 1]$, for all $x \in U$, and $V \in \mathbb{V}$.

Every fuzzy criterion $a_i \in A$, where $i = 1, 2, \ldots, n$, can have a value from a corresponding family of linguistic values denoted by $\mathbb{A}_i = \{A_{i1}, A_{i2}, \ldots, A_{in_i}\}$.

In [7], we proposed a method of classifying the elements of fuzzy information systems. It consists in finding classes of characteristic elements of the universe that have the same dominant linguistic values of attributes (criteria). To calculate the dominance of a linguistic value, we apply a threshold of similarity $\beta$ that must satisfy the inequality $0.5 < \beta \leq 1$.

Given a fuzzy information system FDS, we define [7] for any element $x \in U$, and any fuzzy attribute $a \in A$ the set $\widehat{\mathbb{V}}_a(x) \subseteq \mathbb{V}_a$ of positive linguistic values

$$\widehat{\mathbb{V}}_a(x) = \{V \in \mathbb{V}_a : f(x, V) \geq \beta\}. \tag{1}$$

The elements $x \in U$ can be described with a combination of its positive linguistic values. The set of linguistic labels $\widehat{\mathbb{L}}(x)$ is equal to the Cartesian product of the sets of positive linguistic values $\widehat{\mathbb{V}}_a(x)$, for all $a \in A$:

$$\widehat{\mathbb{L}}(x) = \prod_{a \in A} \widehat{\mathbb{V}}_a(x). \tag{2}$$

By $X_L$, we denote the subset of the elements $x \in U$ that correspond to a linguistic label $L \in \mathbb{L}$, for all fuzzy attributes $a \in A$:

$$X_L = \{x \in U : L(x) = L\}. \tag{3}$$

The subset $X_L$ is called the set of characteristic elements of the linguistic label $L$.

A linguistic label $L \in \mathbb{L}$ can be represented by an ordered tuple of positive linguistic values, for all attributes $a \in A$:

$$L = \left(\hat{V}_{a_1}^L, \hat{V}_{a_2}^L, \ldots, \hat{V}_{a_n}^L\right). \tag{4}$$

## 3. Fuzzy Preference Templates

In the process of evaluating the alternatives performed by a human decision-maker, certain combinations of linguistic values may be preferred (or avoided) over others. In order to exactly describe the preferences of the decision-maker, it is necessary to provide a degree of acceptance of every linguistic value for all fuzzy criteria, which is a number in the interval $[0, 1]$. We assume that the acceptance degree equal to 1 should be interpreted as the highest possible approval (full preference), whereas the acceptance degree equal to 0 denotes a full neutrality (missing preference distinction). Let us denote by $\mu_{A_{ik}}^+$ the degree of preference for the $k$-th linguistic value of the criterion $a_i$. We express the preference template $T^+$ as fuzzy set in the domain of linguistic values of all criteria as follows:

$$T^+ = \{\mu_{A_{11}}^+/A_{11}, \ldots, \mu_{A_{ik}}^+/A_{ik}, \ldots, \mu_{A_{nn_n}}^+/A_{nn_n}\}. \tag{5}$$

In a similar way, one can describe the anti-preferences of a decision-maker, who wants to avoid certain combinations of linguistic values of fuzzy criteria. In

this case, we apply the notion of exclusion degree. If its value is equal to 1, it is interpreted as the highest possible disapproval (full anti-preference). Anti-preference degree equal to 0 denotes a full neutrality (no anti-preference distinction).

The anti-preference template $T^-$ is expressed as a fuzzy set in the domain of linguistic values of criteria as follows:

$$T^- = \{\mu^-_{A_{11}}/A_{11}, \ldots, \mu^-_{A_{ik}}/A_{ik}, \ldots, \mu^-_{A_{nn_n}}/A_{nn_n}\}. \tag{6}$$

*Example* 1. Let us consider an information system with alternatives that are characterized by fuzzy criteria $a_1$, $a_2$, and $a_3$, where the criteria $a_1$ and $a_2$ have three linguistic values, and the criterion $a_3$ five linguistic values. An example of a decision--maker preference template $T^+$ could take the following form:

$$T^+ = \{\, 0.2/A_{11}, 0.9/A_{12}, 1.0/A_{13},$$
$$0.0/A_{21}, 0.1/A_{22}, 1.0/A_{23},$$
$$0.1/A_{31}, 0.3/A_{32}, 1.0/A_{33}, 0.5/A_{34}, 0.0/A_{35}\}.$$

An anti-preference template $T^-$ could be given as follows:

$$T^- = \{\, 0.0/A_{11}, 0.9/A_{12}, 1.0/A_{13},$$
$$1.0/A_{21}, 0.0/A_{22}, 0.0/A_{23},$$
$$0.0/A_{31}, 0.0/A_{32}, 0.0/A_{33}, 0.0/A_{34}, 1.0/A_{35}\}.$$

The preference and anti-preference templates exhibit the interest of the decision--maker in certain combinations of linguistic values of criteria. According to the template $T^+$, the decision-maker does prefer alternatives having the positive linguistic values $A_{12}$ or $A_{13}$ for the criterion $a_1$, $A_{23}$ for the criterion $a_2$, and $A_{33}$ for the criterion $a_3$. However, the positive linguistic values $A_{12}$ or $A_{13}$ are not accepted in combination with the linguistic value $A_{21}$ for the criterion $a_2$, and the linguistic value $A_{35}$ for the criterion $a_3$, as can be seen in the template $T^-$.

In a complex real-world situation, there can be more than one preference template or anti-preference template, because several combinations of linguistic values can be distinguished by a decision-maker. Moreover, is possible that the sets of templates can be expanded, as the expectations of the decision-maker evolve and his or her experience increases, especially by repeating the evaluation process of alternatives.

In the following, we denote the set of preference templates by $\mathbb{T}^+$, and the set of anti-preference templates by $\mathbb{T}^-$, respectively.

Now, we are able to get a refined preference model by generating the sets of preferred (or anti-preferred) linguistic labels of the decision-maker, which can be done by checking the degree of acceptance (exclusion) of the linguistic values of criteria.

Given a preference template $T_j^+ \in \mathbb{T}^+$, the $k$-th linguistic value of the criterion $a_i$ is called the positive linguistic value $A_{ik}$ of $T_j^+ \in \mathbb{T}^+$, if $\mu_{A_{ik}}^+ \geq \beta$.

For an anti-preference template $T_j^- \in \mathbb{T}^-$, the $k$-th linguistic value of the criterion $a_i$ is called the positive linguistic value $A_{ik}$ of $T_j^- \in \mathbb{T}^-$, if $\mu_{A_{ik}}^- \geq \beta$.

A preferred linguistic label $L^+$, and an anti-preferred linguistic label $L^-$ will be represented by tuples of positive linguistic values of criteria, according to (4).

*Example* 2. Let the threshold $\beta$ be equal to 0.7. For the preference template $T^+$ obtained in Example 1, we get two preferred linguistic labels: $L_1^+ = (A_{12}A_{23}A_{33})$, and $L_2^+ = (A_{13}A_{23}A_{33})$, and for the anti-preference template $T^-$, two anti-preferred linguistic labels: $L_1^- = (A_{12}A_{21}A_{35})$, and $L_2^- = (A_{13}A_{21}A_{35})$.

In the case of complex information systems with many criteria and large sets of linguistic values, the possibility of inconsistency must be taken into account, especially in the case when several decision-makers are involved in the evaluation of alternatives and specification of preferences.

For nonempty sets of the preferred linguistic labels $\mathbb{L}^+$ and the anti-preferred linguistic labels $\mathbb{L}^-$, we define the measure of consistency of the preference model of a decision-maker as:

$$c_{\text{pref}} = 1 - \frac{\text{card}(\mathbb{L}^+ \cap \mathbb{L}^-)}{\text{card}(\mathbb{L}^+ \cup \mathbb{L}^-)}. \tag{7}$$

The set $\mathbb{L}^+$ contains the linguistic labels obtained for all preference templates $\mathbb{T}^+$, while the set $\mathbb{L}^-$ includes the linguistic labels obtained for all anti-preference templates $\mathbb{T}^-$, respectively. In the case of inconsistent preferences of the decision maker, the intersection $\mathbb{L}^+ \cap \mathbb{L}^-$ will be non-empty. The value of $c_{\text{pref}}$ belongs to the interval [0, 1]. It will be equal to 1 if no common linguistic labels can be found for the preference and the anti-preference templates. A full overlap of the preferences and the anti-preferences would make the sets $\mathbb{L}^+$ and $\mathbb{L}^-$ identical, and $c_{\text{pref}}$ equal to 0.

In the next step, we should check if the information system contains alternatives that are in accordance with the preferences or anti-preferences of the decision-maker. This can be done by determining the sets of characteristic elements of the preferred linguistics labels $\mathbb{L}^+$ and the anti-preferred linguistic labels $\mathbb{L}^-$.

*Example* 3. Let a fuzzy information system FDS, provided by an expert who assigns membership degrees in linguistic values of criteria, contain an alternative $y \in U$:

$$\mu_{A_{11}}(y) = 0.0, \quad \mu_{A_{12}}(y) = 0.8, \quad \mu_{A_{13}}(y) = 0.2,$$
$$\mu_{A_{21}}(y) = 0.0, \quad \mu_{A_{22}}(y) = 0.3, \quad \mu_{A_{23}}(y) = 0.7,$$
$$\mu_{A_{31}}(y) = 0.0, \quad \mu_{A_{32}}(y) = 0.0, \quad \mu_{A_{33}}(y) = 0.9, \quad \mu_{A_{34}}(y) = 0.1, \quad \mu_{A_{35}}(y) = 0.0.$$

Since the alternative $y$ is a characteristic element of the linguistic label $L_1^+ = (A_{12}A_{23}A_{33})$, it can be perceived as a candidate solution. However, it is necessary to determine more precisely its similarity to all preference labels (and dissimilarity to anti-preference labels) of the decision-maker. Although there is no unique way to solve this problem, continuation of the logic-oriented approach seems to be an appropriate paradigm to obtain intuitive and explainable results.

## 4. Conclusions

We consider the process of determining candidate solutions from a set of alternatives that are characterized by non-numerical criteria expressed by fuzzy linguistic terms. The obtained results depend on the subjective preferences of a decision-maker. By introducing the concepts of fuzzy preference and anti-preference templates, one can conveniently represent the model of the decision maker's expectations and assess its consistency. This is even more important in applications where different preference models of a group of decision makers need to be taken into account. Furthermore, the presented approach can be used in the development of systems for evaluating large sets of alternatives with many fuzzy criteria, which is essential for applications in big data environments. Since complex optimization tasks can be hardly mastered even by experts, a systematic preference-based approach can facilitate the process of decision-making and allow testing the impact of different preferences of attribute values of criteria.

## References

[1] Chen, C.-T. Extensions of the TOPSIS for group decision-making under fuzzy environment. *Fuzzy Sets and Systems*, 114:1–9, 2000.

[2] Opricovic, S. and Tzeng, G.-H. Compromise solution by MCDM methods: A comparative analysis of VIKOR and TOPSIS. *European Journal of Operational Research*, 156(2):445–455, 2004.

[3] Behzadian, M., Otaghsara, S. K., Yazdani, M., and Ignatius, J. A state-of the-art survey of TOPSIS applications. *Expert Systems with Applications*, 39:13051–13069, 2012.

[4] Nadaban, S., Dzitac, S., and Dzitac, I. Fuzzy TOPSIS: A general view. *Procedia Computer Science*, 91:823–831, 2016.

[5] Salih, M., Zaidan, B., Zaidan, A., and Ahmed, M. Survey on fuzzy TOPSIS state-of-the-art between 2007 and 2017. *Computers and Operations Research*, 104:207–227, 2019.

[6] Palczewski, K. and Sałabun, W. The fuzzy TOPSIS applications in the last decade. *Procedia Computer Science*, 159:2294–2303, 2019.

[7] Mieszkowicz-Rolka, A. and Rolka, L. A novel approach to fuzzy rough set-based analysis of information systems. In Z. Wilimowska et al., editors, *Information Systems Architecture and Technology,* volume 432 of *Advances in Intelligent Systems and Computing*, pages 173–183. Springer International Publishing, Switzerland, 2016.

# Robotics and Autonomous Systems

Track Chairs:

- prof. Piotr Lipiński – Lodz University of Technology

- prof. Piotr Skrzypczyński – Poznan University of Technology

- prof. Cezary Zieliński – Warsaw University of Technology

# Comparison of Path Planning Algorithms Utilizing Euclidean Distance Field Maps for Mobile-Manipulating Robots

**Marcin Czajka**[0009−0006−1617−0719], **Bartłomiej Kulecki**[0000−0002−2820−8212], **Eldaniz Babayev**[0009−0006−6628−8969], **Dominik Belter**[0000−0003−3002−9747]

*Poznan University of Technology*
*Institute of Robotics and Machine Intelligence*
*ul. Piotrowo 3A, 60-965 Poznań, Poland*
*name.surname@put.poznan.pl*

**Abstract.** *Euclidean Distance Fields (EDFs) have emerged as a powerful representation for dense mapping in robotics, capturing the shape of objects by encoding the distance to the nearest surface. While EDF-based planning has been extensively studied in the context of simple rigid-body systems, such as drones, its application to robotic manipulators remains largely unexplored. Unlike single-body motion planning, robotic arm planning must account for kinematic constraints imposed by multiple interconnected joints. This paper proposes a strategy to integrate classical state-of-the-art motion planning algorithms for robotic arms operating within EDF-based environments. We evaluate the performance of the implemented planner and highlight its strengths and limitations in handling complex manipulation tasks.*
**Keywords:** *mobile-manipulating robot, motion planning, distance fields*

## 1. Introduction

New dense mapping methods in robotics build Euclidean Distance Fields (EDF) to represent the shape of the objects in the environment. Methods like Voxblox [1], FIESTA [2], IDMP [3], and iSDF [4] represent each point in the 3D space as a distance to the nearest object. These methods also provide gradients that guide collision avoidance by indicating directions away from obstacles. In contrast to traditional mapping methods in robotics, which typically provide only occupancy

139

probabilities for a given region, the gradient facilitates motion planning has focused on avoiding collisions.

The utility of EDFs for motion planning and collision avoidance is primarily illustrated in scenarios involving very simple objects, such as flying drones. In this case, motion planning is defined as collision avoidance for a rigid body and can be solved by applying gradient-based optimization techniques [1, 2]. In this paper, we consider the motion planning of a manipulating robot in the EDF. This problem is much more challenging, because the arm consists of multiple rigid bodies connected with joints. Motion planning for robotic arms has to consider constraints given by the kinematic chain. This paper evaluates state-of-the-art robotic arm motion planning techniques in scenarios that utilize Euclidean Distance Fields.

The main contributions of this article include the following:

- a strategy to integrate classical state-of-the-art motion planning algorithms for robotic arms operating in EDF-based environments,
- a comparative evaluation of motion planning algorithms using EDF maps as the environmental representation,
- performance analysis in close-proximity obstacle scenarios, identifying key strengths and limitations in handling complex manipulation tasks.

## 2. Related Work

The practical benefits of EDFs in robotics are demonstrated in motion planning scenarios. In most cases, new EDF methods are applied to drone trajectory planning. For instance, EDFs have been utilized in gradient-based planning with the CHOMP method [5]. Oleynikova et al. leverage smooth collision costs and gradients from Voxblox to enable fast drone maneuvers [1]. Similarly, the method proposed in [6] is applied in [2] to plan drone trajectories using the FIESTA map. This approach formulates a B-spline motion model and optimizes gradients based on Euclidean, velocity, and acceleration costs to achieve rapid dynamic maneuvers. Among these methods, only the IDMP map [3] has been used with a robotic arm. The CHOMP planner [5] incorporates trajectory smoothness and obstacle avoidance factors derived from EDF data to compute robot trajectories in dynamic environments. The robot shape is approximated using 29 spheres.

Finean et al. conducted a systematic comparison of motion planning for mobile-manipulating robots [7], evaluating three models – Voxblox [1], FIESTA [2], and a GPU-based 3D voxel map. Their experiments showed that, despite not fully utilizing GPU acceleration, the standard 3D voxel map achieved the fastest planning

times with the GPMP2 method [8]. In contrast, we focus on comparing three different motion planning methods using the same EDF map, evaluating their success rate and the length of the obtained path.

# 3. Path Planning Algorithms: A Comparison

To evaluate the path planning methods, we performed experiments with a mobile-manipulating robot and a precomputed ESDF map (Voxblox) of the environment. The experiments involved moving the robot among three configurations around a stack of boxes, with each algorithm attempting 100 times to plan a path. We selected three motion planning algorithms: RRT-Connect [9], RRT* [10] and CHOMP [5]. RRT-Connect [9] is a sampling-based motion planning algorithm that grows two trees: one from the start and one from the goal, which iteratively extend toward each other until they connect. This bidirectional strategy improves efficiency by rapidly exploring the space and reducing the time required to find a feasible path in complex environments. RRT* [10] not only explores the space efficiently but also optimizes paths by rewiring the tree to reduce overall path cost. Unlike standard RRT, which finds a feasible but suboptimal path, RRT* ensures asymptotic optimality, meaning the path quality improves over time as more samples are added. CHOMP (Covariant Hamiltonian Optimization for Motion Planning) [5] is a trajectory optimization algorithm that iteratively refines an initial path by minimizing a cost function balancing smoothness and collision avoidance, leveraging gradient-based optimization techniques.

The RRT-Connect and RRT* methods operated in Cartesian space and were based on the implementation from OMPL – Open Motion Planning Library [11], with a maximum solving time set to 5.0 seconds. We use the CHOMP implementation available in the ROS MoveIt library. For obstacle avoidance, we utilize a point cloud model of the robotic arm. Each robot link is represented by 100 points on the link surface. To optimize the time needed for collision querying, a two-step algorithm is proposed: (I) retrieving an ESDF value from Voxblox map for link positions and comparing with the radius of a sphere encompassing a specific link, (II) only if the first step indicates a possible collision, the detailed point cloud is queried, and minimal clearance from an obstacle is set to 2.5 cm. For robot self-collision detection, we employ a neural network model [12].

## 3.1. Motion Planning Experiment

During the experiment, the real UR5 robotic arm equipped with the Kinect Azure camera on the wrist moves around to scan the scene. We utilized the Structure-from-Motion (SfM) COLMAP algorithm to estimate the camera poses. The resulting point clouds are then used to update the Voxblox map [1]. The obtained Euclidean Distance Field (Voxblox) is used to plan the path of the robot. The experimental set for motion planning is presented in Figure 1. In the first stage of the experiment, the robot plans its motion between the initial (Figure 1a) and the first goal of the robot (Figure 1b). During the second stage, the robot plans its motion between the first (Figure 1b) and the second goal pose (Figure 1c) of the robot. The reference poses presented in Figure 1 were selected to make the motion planning problem challenging. There are three main constraints that limit the motion of the robot: (i) self-collisions, (ii) collisions with the obstacles represented by the EDF map of the environment, and (iii) joint limits. The motion planner should avoid these constraints and initially move the arm far from the goal pose to avoid collisions. Thus, we expect a success rate below 100% for each planner tested in the given scenario, even though they operate easily in much simpler scenes.



Figure 1. Motion planning experiment: initial configuration of the robot (a), first goal configuration of the robot (b), and second goal configuration of the robot (c)

## 3.2. Results

The calculated trajectory metrics are summarized in Table 1. The statistics are calculated for 100 runs of each planner for the given scenario. The RRT-Connect planner achieved, on average, a significantly higher success rate and shorter solving time, but a longer path in the Cartesian space compared to the RRT*. This can be explained by the nature of both algorithms. RRT* has an additional optimiza-

Table 1. Comparison of path planning results for two motion stages: (1) from the initial configuration to the first goal, and (2) from the first goal to the second goal. The table presents planning time $t$, number of path points $N$, path length in Cartesian space $d_{3D}$, and path length in configuration space $d_c$.

| algorithm | path stage | success rate | t [s] | | N | | $d_{3D}$ [m] | | $d_c$ [rad] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | mean | std | mean | std | mean | std | mean | std |
| RRT-Connect | | **0.53** | **1.51** | 1.09 | 3.19 | 0.44 | 0.65 | 0.20 | 7.16 | 1.80 |
| RRT* | 1 | 0.24 | 2.70 | 1.48 | 3.04 | 0.20 | 0.56 | 0.17 | 6.69 | 1.48 |
| CHOMP | | 0.47 | 3.55 | 1.37 | **2.91** | 0.30 | **0.55** | 0.15 | **6.50** | 1.60 |
| RRT-Connect | | **0.55** | **2.10** | 1.12 | 4.65 | 0.93 | 1.26 | 0.30 | 9.20 | 5.11 |
| RRT* | 2 | 0.26 | 4.44 | 0.99 | 4.04 | 0.53 | 1.02 | 0.20 | 6.39 | 2.31 |
| CHOMP | | 0.49 | 3.01 | 1.10 | **3.94** | 0.66 | **0.95** | 0.15 | **5.30** | 1.95 |

tion phase, so the additional time is consumed to minimize the path. Despite RRT* optimizing only the length of the path in the Cartesian space, the generated trajectories had fewer waypoints and shorter lengths in the configuration space than the corresponding paths computed by RRT-Connect. A drawback of the tested space-sampling algorithms is the need to calculate the inverse kinematics for each generated configuration. The shortest path is found by the gradient-based CHOMP algorithm, but it requires the longest motion planning time in the first scenario and is the second worst in the second scenario.

## 4. Conclusions and Future Work

In this paper, we performed a comparative analysis of motion planning algorithms using the EDF maps as the environmental representation. The experiments focused on motion in the proximity of obstacles, allowing us to identify the advantages and limitations of the analyzed methods. The results proved that CHOMP, an approach based on gradient optimization, produced the shortest paths among the tested planners. Among the space sampling algorithms, RRT* found trajectories of comparable length to those produced by CHOMP. On the other hand, RRT-Connect achieved the highest success rate and the shortest solving time; however, the resulting paths were the longest. This is explained by the nature of the RRT-Connect algorithm that stops searching when it finds the first feasible path. Also,

RRT-Connect is the only algorithm that does not optimize the found path. Comparing two methods that optimize the path – a sampling-based RRT* and CHOMP, our experiments show that CHOMP has a significantly larger success rate and provides a shorter path.

In the future, we are going to implement an optimization-based path-planning method that uses neural methods to check motion constraints.

## Acknowledgment

## References

[1] Oleynikova, H., Taylor, Z., Fehr, M., Siegwart, R., and Nieto, J. Voxblox: Incremental 3D Euclidean signed distance fields for on-board MAV planning. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1366–1373. IEEE, 2017. doi:10.1109/IROS.2017. 8202315.

[2] Han, L., Gao, F., Zhou, B., and Shen, S. FIESTA: Fast incremental Euclidean distance fields for online motion planning of aerial robots. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4423–4430. IEEE, 2019. doi:10.1109/IROS40897.2019.8968199.

[3] Ali, U., Wu, L., Müller, A., Sukkar, F., Kaupp, T., and Vidal-Calleja, T. Interactive distance field mapping and planning to enable human-robot collaboration. *IEEE Robotics and Automation Letters*, 9(12):10850–10857, 2024. doi:10.1109/LRA.2024.3482128.

[4] Ortiz, J., Clegg, A., Dong, J., Sucar, E., Novotny, D., Zollhoefer, M., and Mukadam, M. iSDF: Real-time neural signed distance fields for robot perception. In *Robotics: Science and Systems*. 2022.

[5] Zucker, M., Ratliff, N., Dragan, A. D., Pivtoraiko, M., Klingensmith, M., Dellin, C. M., Bagnell, J. A., and Srinivasa, S. S. CHOMP: Covariant Hamiltonian optimization for motion planning. *The International*

*Journal of Robotics Research*, 32(9-10):1164–1193, 2013. doi:10.1177/ 0278364913488805.

[6] Zhou, B., Gao, F., Wang, L., Liu, C., and Shen, S. Robust and efficient quadrotor trajectory generation for fast autonomous flight. *IEEE Robotics and Automation Letters*, 4(4):3529–3536, 2019. doi:10.1109/LRA.2019. 2927938.

[7] Finean, M. N., Merkt, W., and Havoutis, I. Simultaneous scene reconstruction and whole-body motion planning for safe operation in dynamic environments. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3710–3717. IEEE, 2021. doi:10.1109/IROS51168. 2021.9636860.

[8] Mukadam, M., Dong, J., Yan, X., Dellaert, F., and Boots, B. Continuous-time Gaussian process motion planning via probabilistic inference. *The International Journal of Robotics Research*, 37(11):1319–1340, 2018. doi: 10.1177/0278364918790369.

[9] Kuffner, J. and LaValle, S. RRT-connect: An efficient approach to single-query path planning. In *Proceedings 2000 ICRA. Millennium Conference. IEEE International Conference on Robotics and Automation. Symposia Proceedings*, volume 2, pages 995–1001. IEEE, 2000. doi:10.1109/ROBOT. 2000.844730.

[10] Karaman, S. and Frazzoli, E. Sampling-based algorithms for optimal motion planning. *The International Journal of Robotics Research*, 30(7):846–894, 2011. doi:10.1177/0278364911406761.

[11] Şucan, I. A., Moll, M., and Kavraki, L. E. The Open Motion Planning Library. *IEEE Robotics & Automation Magazine*, 19(4):72–82, 2012. doi: 10.1109/MRA.2012.2205651. https://ompl.kavrakilab.org.

[12] Kulecki, B. and Belter, D. Boosting machine learning techniques with positional encoding for robot collision checking. In *2024 13th International Workshop on Robot Motion and Control (RoMoCo)*, pages 90–95. IEEE, 2024. doi:10.1109/RoMoCo60539.2024.10604400.

# Genetic Programming Iterative Improvement Algorithm for a Concurrent Real-Time Optimization in Embedded System Design Process

**Adam M. Górski**[0000−0003−3821−5333], **Maciej Ogorzałek**[0000−0003−3314−269X]

*Jagiellonian University*
*Faculty of Physics, Astronomy and Applied Computer Science*
*prof. Stanisława Łojasiewicza 11, 30-348 Krakow, Poland*
*a.gorski@uj.edu.pl, maciej.ogorzalek@uj.edu.pl*

**Abstract.** *Embedded systems need to execute some special tasks. In most of cases the tasks needed to be predicted by designers. However, the real problem occurs when systems meet unexpected situations. Then some unexpected tasks appear. However, many unexpected situations can be solved in many ways. The problem is to detect the tasks and assign them to appropriate resource to be executed. Therefore it is very hard to establish which way is better. Any modifications in system architectures can be impossible or too expensive. The problem can be split into two concurrent phases. Each phase impacts another in a real time. In this paper an iterative improvement genetic programming methodology for a concurrent real-time optimization in embedded system design process is presented. The methodology is able to detect unexpected tasks and choose the best resource assignment. Thanks to iterative improvement nature of the algorithm, the process can be faster and cheaper. It is also more resistant to getting stuck in local minima when optimizing parameters.*
**Keywords:** *genetic programming, artificial intelligence, embedded systems, unexpected tasks, concurrent real-time optimization*

## 1. Introduction

Concurrent real-time optimization [1] is a quite new optimization problem. Generally embedded system design can be divided into four phases: modeling, val-

idation, implementation and assignment of unexpected tasks [2]. Design methodologies can be divided into two basic groups: constructive [3] and iterative improvement [4, 5].

When embedded system is designed and it works in target environment, the system needs to execute some unexpected tasks [2]. In some situations adding new hardware component is impossible or too expensive. Example of such a situation can be autonomous robot [6] working on planet Mars. In [2] authors proposed an algorithm to assign unexpected tasks. It had low complexity, but the results were suboptimal. Good results were obtained using genetic programming [7]. The above-mentioned methodologies have one weakness – unexpected tasks needed to be found using some external methods. In [8] authors proposed a methodology for an automatic detection and assignment of unexpected tasks. The authors indicated that some of the tasks executed by the embedded system can be split into finite number of subtasks. The connections of such subtasks can solve some of the unexpected situations. A methodology for assigning unexpected tasks to a group of embedded systems was also proposed [9]. In the paper unexpected tasks were not a result of unexpected situations by an extension of the possibilities of embedded systems which needed to work together.

An unexpected situation can be solved in many ways. Each way is a different connection of subtasks. Not every connection leads to solving unexpected situations. After finding the solution, it is needed to be assigned to appropriate processing element (PE) to be executed. The question is: if each approach requires different subtasks, then which one is better? In [10] the authors proposed a genetic algorithm for detection and assignment of unexpected tasks in IoT design process. The authors proposed to split the process of detecting and assignment of unexpected tasks into two phases. Each phase impacted another in real time. The first phase was responsible for choosing the subtasks and their number. The second phase makes the verification and makes the proper optimization. Such a problem was called concurrent-real time optimization [1, 10, 11]. To solve this problem for embedded systems, a genetic algorithm was proposed [1]. In [11], a genetic programming approach was described. The biggest disadvantage of such an approach was that it belongs to the constructive group of algorithms.

In this paper, a genetic programming iterative improvement approach for concurrent real-time optimization for unexpected tasks detection and assignment is presented. Unlike other existing methods, our approach starts with the fastest architecture and improves the system by making local changes. Each of such changes is marked on the genotype tree as a separate node. Therefore the construc-

tion of the genotype can be different for every individual. The paper is organized as follows: the next section presents the preliminaries, followed by a description of the algorithm. Section 4 contains the experimental results. Finally, the conclusions and directions for future work are presented.

## 2. The Preliminaries

An embedded system is a microprocessor or microcontroller-based computer system with a software. Such a system mostly contains an operating system and is optimized to execute special tasks. The tasks are characterized by two parameters – time and cost of the execution. Most of modern embedded systems are solved as distributed.

The embedded system is specialized using an acyclic directed graph called task graph $G = (V, E)$. In the nodes of the graph there are tasks $v_i \in V$. Every edge $e_{i,j} \in E$ describes the amount of data $d_{i,j}$ that needs to be transferred between two connected tasks $v_i$ and $v_j$. The transmission time $t_{i,j}$ depends on bandwidth $b$ of CL used for communication and is equal to:

$$t_{i,j} = \frac{d_{i,j}}{b}. \tag{1}$$

The transmission time is equal to 0 if two connected tasks are executed by the same resources. An example of a task graph is presented below in Figure 1. The example consists of 14 tasks. Nine of those tasks are predicted ones. Unexpected tasks are marked in yellow on the graph.



Figure 1. Example of an extended task graph

The goal of system design is to find the cheapest system which satisfies the

time constrains. The overall cost of the system is presented by equation (2) below:

$$C_0 = C_i + \sum_{j=1}^{n} c_j, \tag{2}$$

where $C_i$ is an initial cost of a hardware – the sum of all of PEs in the system, and $n$ is a number of all the tasks in the task graph.

## 3. The Algorithm

If an unexpected situation happens, the system needs to execute some unexpected tasks. Every unexpected task needs to be inserted on a task graph. At the beginning every task is divided into possible number of subtasks. The connection of some of the subtasks can be a solution of an unexpected situation and provide missing data. As can be easily seen, not every combination of subtasks leads to a solution. The algorithm presented in this section can find out which connections do not give appropriate result.

The number of individuals in the population is dependent on a number of the tasks in the task graph and a number of PEs in the system. To solve unexpected situations the algorithm randomly chooses the subtasks and their number. According to genetic programming rules, every genotype is a tree. The tree describes the process of creating the single solution. Thus it is not needed to have the same structure of every genotype. The first node in a genotype represents the initial system, which serves as the embryo of the architecture. This initial system corresponds to the fastest configuration. Each subsequent node represents a system build option, and each option has a certain probability of being selected. The available options are listed in Table 1 below.

Table 1. System build options

| The option | Probability |
|---|---|
| 1. The fastest | 0,1 |
| 2. The cheapest | 0,1 |
| 3. Min time * cost | 0,5 |
| 4. The least used | 0,3 |

Each node, apart from the option itself, contains information about what percentage of tasks from the previous node is assigned using the option in the current

149

node. The percentage of tasks assigned in a node is always greater than zero; however, it is possible for the actual number of assigned tasks to be zero, for example, after the application of certain genetic operators. Such nodes, and their potential successors are removed from the genotype. Figure 2 below presents an example of a genotype. The example consist of eight nodes. The first node is an embryo. Each next node contains the number of options and percentage number of tasks taken from the predecessor. The nodes in the genotype are executed according to the level of the tree. The nodes at the same level are executed from left to right.



Figure 2. Example of a genotype

In the example it can be easily observed that none of the direct successors of a node can have the same option as the predecessor. In such a case a node makes no changes. The second dependence is that the sum of the percentages of the tasks taken by direct successors cannot exceed 100.

Next generations are obtained using genetic operators: mutation, crossover, cloning and rank selection. The algorithm stops after $\varepsilon$ generations without finding better result.

## 4. Experimental Results

To test the efficiency of the proposed approach (GP2025) we decided to make experiments using benchmarks with 10, 20 and 30 nodes. The results were compared with genetic programming algorithm (GP2024) proposed by Górski and Ogorzałek [11]. The best results are presented in Table 2 below.

Algorithm GP2024 was proven to be more efficient than genetic algorithm solution [1]. For each of the benchmarks and algorithms 30 runs were made. The values of parameters were the same for both of the algorithms. As it can be observed, for every benchmark the algorithm GP 2025 generated better results. The

Table 2. Experimental results

| Graph | t_max | GP2025 | | | GP2024 | | |
|---|---|---|---|---|---|---|---|
| | | Time | Cost | Generation | Time | Cost | Generation |
| 10 | 1500 | 1416 | 992 | 11 | 1179 | 1319 | 12 |
| 20 | 2000 | 1902 | 1942 | 15 | 1858 | 2243 | 19 |
| 30 | 3000 | 2956 | 1982 | 14 | 2873 | 2249 | 18 |

result of cost generated by GP 2025 for a graph with 10 nodes was about 25% better than the result generated by GP 2024. This can indicate that GP 2024 could stop in local minima. Such a situation can be a result of a structure of the graph, given values of time and cost. However, we must remember that the compared algorithms are probabilistic. The percentage difference for the rest of the graphs was lower – 13% in the case of the benchmark with 20 nodes and 12% for the benchmark with 30 nodes. It can also be observed that for the best obtained results GP 2025 generated less generations than GP 2024.

## 5. Conclusions

In this paper an iterative improvement genetic programming approach for concurrent real-time optimization which occurs in embedded system design process is presented. Unlike other approaches, our algorithm starts with the fastest architecture for every genotype. Then, some local modifications produce the final forms of the genotypes. The first obtained results indicate a good efficiency of the proposed approach. However, the algorithm needs more comparison using bigger graphs.

In the future, we plan to propose new genetic programming–based solutions for the investigated problem. We will also try to develop algorithms based on other optimization methods, such as, for example, metaheuristics approaches. It is also planned to check the efficiency of the proposed algorithm for the systems described by bigger task graphs.

## References

[1] Górski, A. and Ogorzałek, M. Concurrent real-time optimization in embedded system design process using genetic algorithm. In *5th Polish Conference on Artificial Intelligence (PP-RAI'2024)*, pages 331–337. 2024.

[2] Górski, A. and Ogorzałek, M. J. Assignment of unexpected tasks in embedded system design process. *Microprocessors and Microsystems*, 44:17–21, 2016.

[3] Wang, Y., Bozkurt, A. K., Smith, N., and Pajic, M. Attack-resilient supervisory control of discrete-event systems: A finite-state transducer approach. *IEEE Open Journal of Control Systems*, 2:208–220, 2023.

[4] Jha, S. K., Jha, S., Ewetz, R., and Velasquez, A. Co-synthesis of code and formal models using large language models and functors. In *IEEE Military Communications Conference (MILCOM)*, pages 215–220. IEEE, 2024.

[5] Górski, A. and Ogorzałek, M. Genetic programming based iterative improvement algorithm for HW/SW cosynthesis of distributted embedded systems. In *10th International Conference on Sensor Networks*, pages 120–125. 2021.

[6] Wang, Z., Yan, H., Wang, Y., Xu, Z., Wang, Z., and Wu, Z. Research on autonomous robots navigation based on reinforcement learning. In *Proceedings of 2024 3rd International Conference on Robotics, Artificial Intelligence and Intelligent Control (RAIIC)*, pages 78–81. 2024.

[7] Górski, A. and Ogorzałek, M. Assignment of unexpected tasks in embedded system design process using genetic programming. In *Dynamics of Information Systems*, pages 93–101. Springer, 2023.

[8] Górski, A. and Ogorzałek, M. Auto-detection and assignment of unexpected tasks in embedded systems design process. In *23rd EG-ICE International Workshop on Intelligent Computing in Engineering (EG-ICE 2016)*, pages 179–188. 2016.

[9] Górski, A. and Ogorzałek, M. J. Assignment of unexpected tasks for a group of embedded systems. *IFAC-PapersOnLine*, 51(6):102–106, 2018.

[10] Górski, A. and Ogorzałek, M. Concurrent real-time optimization of detecting unexpected tasks in IoT design process using GA. In *Late Breaking Papers from the IEEE 2023 Congress on Evolutionary Computation*, pages 74–77. 2023.

[11] Górski, A. M. and Ogorzałek, M. Detecting and assignment of unexpected tasks in SoC design process using genetic programming. In *Proceedings of the 21st International SoC Design Conference*, pages 398–399. 2024.

# Comparative Evaluation of Euclidean Distance Field Mapping Methods for Mobile-Manipulating Robots

Bartłomiej Kulecki[0000−0002−2820−8212], Marcin Czajka[0009−0006−1617−0719], Eldaniz Babayev[0009−0006−6628−8969], Dominik Belter[0000−0003−3002−9747]

*Poznan University of Technology*
*Institute of Robotics and Machine Intelligence*
*ul. Piotrowo 3A, 60-965 Poznań, Poland*
*name.surname@put.poznan.pl*

**Abstract.** *Mapping is essential for robotic applications including autonomous work in* a priori *unknown environment. Motion planning of mobile and manipulating robots requires detailed environmental representations to ensure collision avoidance. Recent Euclidean Distance Field (EDF)-based methods offer advantages for motion planning by providing continuous distance and gradient information. This paper compares three EDF-based mapping methods in the context of mobile-manipulating robots, evaluating their strengths and limitations in real-world scenarios. Our findings highlight the benefits of EDF representations and their potential for improving robotic motion planning in complex environments.*
**Keywords:** *robot, mapping, distance fields*

## 1. Introduction

Mapping is a fundamental problem in robotics. Two typical applications of mapping methods are localization [1] and dense mapping for motion planning [2, 3, 4]. Because the primary goal of mapping methods in Simultaneous Localization and Mapping (SLAM) is to define the pose of the camera in 3D space, the maps do not have to represent the complete shape of the objects in the environment, and typically they are sparse [1]. These maps are unsuitable for motion planning due to the lack of full representation of the objects in the environment.

In contrast, maps used for motion planning focus on a detailed representation of the environment to enable the robot to avoid collisions with objects.

Classical dense mapping methods such as Octomap represent the environment as a set of voxels [5]. To better define the shape inside the voxel, a Normal Distribution Transform can be used [6]. However, recently the Euclidean Distance Fields like FIESTA [3], Voxblox [2], and Interactive Distance Field Mapping (IDMP) [7] have gained popularity due to their beneficial properties in motion planning. In these maps, the distance to the nearest surface and the gradient of the distance field can be easily determined for any point in 3D space. This approach directly represents the collision cost of position in the 3D space, and the gradient gives the direction to avoid collisions. In addition, continuous representation eases the training of neural models of EDFs such as iSDF [4] in contrast to the NeRF [8] model designed for scene modeling and image synthesis.

In this paper, we focus on analyzing and comparing recent EDF-based mapping methods in the context of building the model for mobile-manipulating robots. We evaluate three methods – Voxblox, FIESTA, and IDMP – and summarize their advantages and limitations in the typical mobile-manipulating robot scenario.

## 2. Related Work

The Voxblox method [2] was the first EDF-based approach that proved that the Truncated Signed Distance Fields outperform Octomaps in both speed and accuracy. Moreover, Oleynikova et al. [2] show that obstacle distance information helps with fast local motion planning. The FIESTA method proposes two independent queues for inserting and deleting data. This approach minimizes the number of updated nodes and speeds up the mapping. The efficient GPU implementation of nvblox [9] further improves the speed, resolution, and scale of dense mapping.

The neural version of the Signed Distance Field named iSDF is trained to model the signed distance for input 3D coordinates [4]. This method can also fill in gaps in partially observed regions. Another neural representation named NeRF is popular for view synthesis [8]. NeRF can be used for motion planning, but training the network is time-consuming. Robotic applications require real-time map updates and fast access to the content. IDMP map focuses on dynamic objects and utilizes the Gaussian Process to consider dynamic objects in the environment, especially in human–robot cooperation [7].

# 3. Mapping Experiment

We conducted two types of experiments to evaluate the performance of the mapping systems. The first was performed in a static environment to assess the accuracy of the generated map. The second experiment tested the mapping method in a dynamic environment, where the mobile-manipulating robot, equipped with a UR5 arm and Kinect Azure camera, moved the arm while scanning the scene. We employ the Structure-from-Motion (SfM) COLMAP algorithm to estimate the camera poses. The obtained point clouds are used to update the maps. During the experiment with the dynamic scene, the robotic arm scans the environment. Then, the arm is static, and only the plastic can is moved on the table. All mapping methods are verified on the same dataset, so sensor noise, occlusions, or environmental factors like lighting conditions and surface reflections have the same impact on each method's performance.

The voxel size for each mapping method is set to 1 cm. The other parameters of the methods differ and are not directly comparable. Therefore, in the experiments presented in this paper, we use the default parameters provided by the library developer.

## 3.1. Qualitative Comparison

The experimental sets are presented in Figure 1a. In Figure 1b, c, and d, we show the horizontal cross-section over the EDF for the FIESTA, Voxblox, and IDMP, respectively. Comparing the results for FIESTA and Voxblox, we can note that the map obtained for the FIESTA method is smooth and does not preserve the sharp edges of the objects, but the overall accuracy is the highest. When we compare the model of the blue bucket, indicated by the arrows in Figure 1, we note that the model is incorrect for Voxblox. We justify this error by the high sensitivity of the Voxblox to localization and measurement errors and the smaller "inertia" of the Voxblox map, which is later confirmed in the experiment with dynamic objects.

## 3.2. Quantitative Comparison

The maps were evaluated by querying them with a dense grid of points with a 1 cm resolution in all dimensions. The ground truth was created by merging every 25th frame of a sequence and computing the Euclidean distance from each query

Figure 1. Qualitative comparison between three mapping systems for two example scenes named `scene1` (top) and `scene2` – visualization of the horizontal cross-section over the Euclidean Distance Field for the FIESTA (b), Voxblox (c), and IDMP (d)

coordinate to the nearest point in the resulting point cloud. For a fair comparison, the absolute values of the distances from the maps were used, and coordinates with unknown values were omitted. The error metrics are summarized in Table 1. The quantitative results correspond with the qualitative results. The FIESTA method provides a two times smaller error than the Voxblox map.

Table 1. Comparison of mapping errors for static sequences. All values in meters.

| Algorithm | scene1 | | | scene2 | | | scene3 | | |
|---|---|---|---|---|---|---|---|---|---|
| | MAE | std | max | MAE | std | max | MAE | std | max |
| IDMP | 0.020 | 0.042 | 0.404 | 0.036 | 0.057 | 0.529 | 0.092 | 0.085 | 0.592 |
| Voxblox | 0.022 | 0.014 | 0.079 | 0.016 | 0.013 | 0.150 | 0.022 | 0.016 | 0.084 |
| FIESTA | **0.010** | **0.010** | **0.070** | **0.008** | **0.010** | **0.077** | **0.012** | **0.013** | **0.078** |

## 3.3. Mapping in a Dynamic Environment

We also experimented with the dynamic environment. The results are presented in Figure 2. The robot observes the scene with a rolling can. The experiment shows that Voxblox is more suitable for dynamic scenes than FIESTA. The map built by the FIESTA algorithm contains the geometry related to the initial and goal pose of the can and a "shadow" between these two poses. The map created

Figure 2. Qualitative comparison between three mapping systems for three example scenes – visualization of the horizontal cross-section over the Euclidean Distance Field for the FIESTA (b), Voxblox (c), and IDMP (d)

by the Voxblox algorithm is slightly better, because it does not store the "shadow" created by moving objects. However, the initial object remains visible on the map and is not removed by the Voxblox algorithm. Although the IDMP map is designed for dynamic scenes, it still contains artifacts associated with the temporary poses of the rolling can. However, compared to FIESTA, the number of artifacts is significantly lower.

### 3.4. Computational Requirements Comparison

We measured CPU and memory usage during mapping. The workstation was equipped with a 16th-thread Intel i9-9900KF processor and 64GB of RAM. The obtained statistics are summarized in Table 2. The FIESTA method required substantially more memory compared to other algorithms, due to the need to allocate memory during startup for a larger area than that covered by the final map. On the other hand, FIESTA used considerably fewer CPU resources than other methods. In contrast, Voxblox utilized at least seven processing threads for over 10% of the experiments' duration, yet it showed the lowest memory consumption. Meanwhile, the IDMP map required a similar amount of memory to that of Voxblox, although the CPU usage was lower.

## 4. Conclusions and Future Work

We conducted a quantitative evaluation of mapping accuracy and examined performance in dynamic environments. The accuracy assessment involved querying maps with a dense grid and comparing distances to a ground truth point cloud. The results showed that the FIESTA method produced errors twice as small as

Table 2. Comparison of CPU and memory consumption. The table presents the maximum memory usage $\text{MEM}_{max}$, the average processor utilization $\text{CPU}_{avg}$, as well as the 90th and 99th percentiles of CPU usage $\text{CPU}_{90}$ and $\text{CPU}_{99}$, respectively. Memory usage is measured in megabytes (MB), while processor utilization is expressed as a percentage of a single thread's capacity.

| Algorithm | scene2 | | | | dynamic scene | | | |
|---|---|---|---|---|---|---|---|---|
| | $\text{MEM}_{max}$ | $\text{CPU}_{avg}$ | $\text{CPU}_{90}$ | $\text{CPU}_{99}$ | $\text{MEM}_{max}$ | $\text{CPU}_{avg}$ | $\text{CPU}_{90}$ | $\text{CPU}_{99}$ |
| IDMP | 550 | 2.19 | 3.10 | 3.94 | 544 | 2.19 | 3.14 | 4.67 |
| Voxblox | **436** | 2.97 | 7.96 | 9.89 | **464** | 2.73 | 7.65 | 9.82 |
| FIESTA | 3721 | **0.76** | **1.60** | **1.64** | 3730 | **0.71** | **1.60** | **1.64** |

the Voxblox map. In dynamic environments, Voxblox outperformed FIESTA by avoiding persistent artifacts caused by moving objects. The IDMP algorithm, designed for dynamic scenes, reduced such artifacts more effectively than FIESTA, but was not entirely free from them. The resource utilization comparison indicated that Voxblox demonstrated the lowest RAM requirements, while FIESTA had the least CPU usage. In the future, we will investigate deep-learning-based approaches like iSDF [4] and focus on developing neural interfaces between distance maps and gradient-based motion planners.

# Acknowledgment

# References

[1] Matsuki, H., Scona, R., Czarnowski, J., and Davison, A. J. CodeMapping: Real-time dense mapping for sparse SLAM using compact scene representations. *IEEE Robotics and Automation Letters*, 6(4):7105–7112, 2021. doi:10.1109/LRA.2021.3097258.

[2] Oleynikova, H., Taylor, Z., Fehr, M., Siegwart, R., and Nieto, J. Voxblox: Incremental 3D Euclidean Signed Distance Fields for on-board MAV planning.

In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1366–1373. 2017. doi:10.1109/IROS.2017.8202315.

[3] Han, L., Gao, F., Zhou, B., and Shen, S. FIESTA: Fast incremental Euclidean Distance Fields for online motion planning of aerial robots. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4423–4430. 2019. doi:10.1109/IROS40897.2019.8968199.

[4] Ortiz, J., Clegg, A., Dong, J., Sucar, E., Novotny, D., Zollhoefer, M., and Mukadam, M. iSDF: Real-time neural signed distance fields for robot perception. In *Robotics: Science and Systems*. 2022.

[5] Hornung, A., Wurm, K. M., Bennewitz, M., Stachniss, C., and Burgard, W. OctoMap: An efficient probabilistic 3D mapping framework based on octrees. *Autonomous Robots*, 2013. doi:10.1007/s10514-012-9321-0.

[6] Saarinen, J., Andreasson, H., Stoyanov, T., and Lilienthal, A. J. 3D normal distributions transform occupancy maps: An efficient representation for mapping in dynamic environments. *The International Journal of Robotics Research*, 32(14):1627–1644, 2013. doi:10.1177/0278364913499415.

[7] Ali, U., Wu, L., Müller, A., Sukkar, F., Kaupp, T., and Vidal-Calleja, T. Interactive distance field mapping and planning to enable human–robot collaboration. *IEEE Robotics and Automation Letters*, 9(12):10850–10857, 2024. doi:10.1109/LRA.2024.3482128.

[8] Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., and Ng, R. NeRF: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. doi:10.1145/3503250.

[9] Millane, A., Oleynikova, H., Wirbel, E., Steiner, R., Ramasamy, V., Tingdahl, D., and Siegwart, R. nvblox: GPU-accelerated incremental signed distance field mapping. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2698–2705. 2024. doi:10.1109/ICRA57147.2024.10611532.

# Leveraging D*Lite in Reinforcement Learning-Based Multi-Agent Path Planning

**Kamil Młodzikowski**[0000−0002−9761−1400], **Dominik Belter**[0000−0003−3002−9747]

*Poznan University of Technology*
*Institute of Robotics and Machine Intelligence*
*ul. Piotrowo 3A, 60-965 Poznan, Poland*
*kamil.mlodzikowski@put.poznan.pl*

**Abstract.** *In this paper, we propose an approach that integrates D\* Lite path data into multi-agent reinforcement learning to address the challenge of limited observation windows in multi-robot path planning. By embedding global path information from D\* Lite into each agent's local observation space, we guide the agents with classical planning insights. This hybrid strategy accelerates convergence and improves decision-making despite partial observability. The results show improved performance and shorter episode lengths compared to baselines relying solely on raw local observations or classical path planning, offering a scalable solution for multi-agent navigation in complex environments.*
**Keywords:** *reinforcement learning, pathfinding, multi-agent systems*

## 1. Introduction

Multi-agent path planning in congested environments is challenging, especially with agents having partial observations. Traditional path planners like D\* Lite find near-optimal routes with global maps, but struggle with dynamic, local disruptions. Reinforcement learning (RL) allows agents to adapt in uncertain environments, but often converges slowly due to limited local observations.

In this paper, we bridge the gap between classical path-planning and learning-based methods by incorporating D\* Lite path information directly into the observation space of a multi-agent reinforcement learning system. By providing agents with strategic guidance from a global planner, even though they only perceive a local window in the environment, we enable more informed decision-making

160

and faster convergence. Our hybrid approach leverages D* Lite's efficient global path computation and RL's adaptability, resulting in faster convergence and improved performance in dynamic, grid-based multi-robot environments compared to purely RL-based or classical methods.

## 1.1. Related Work

Multi-agent pathfinding (MAPF) has been approached through both classical graph-search methods and more recent learning-based or hybrid solutions. Among the classical techniques, Conflict-Based Search (CBS) [1] stands out as an optimal algorithm that decomposes the joint pathfinding problem into individual single-agent searches via a two-level structure: a high-level "constraint tree" detects and resolves collisions, and a low-level solver optimizes each agent's path. However, when collisions occur frequently in open areas, CBS can branch excessively. To mitigate this, Meta-Agent CBS (MA-CBS) [1] merges conflicted agents into a single meta-agent and uses a standard MAPF sub-solver, balancing a decoupled policy with occasional fully coupled planning. Another complete approach, M* [2], adaptively couples only colliding agents, expanding the state space dimensionally only where necessary.

Meanwhile, learning-based methods often handle partial observability and large teams more gracefully. Panov et al. [3] apply deep reinforcement learning (DRL) to single-agent grid navigation, treating occupancy grids as images for a CNN-based Q-learner. Extending DRL to multiple agents, Damani et al. [4] propose PRIMAL$_2$, which integrates imitation learning from a centralized solver and local collision-avoidance heuristics, achieving scalability to thousands of agents. Zhang et al. [5] similarly combine tree search with a decentralized MARL policy, refining short-term conflict resolution for automated guided vehicles. Efforts like SWARM-MAPF [6] preserve formation constraints by splitting paths into "open" segments (where formation is easy) and "congested" segments handled by specialized MAPF planners. These contrasting strategies, from classical solvers to adaptive learning, highlight the complexity of multi-agent navigation and the need for context-specific or blended solutions.

## 2. D*-Lite-Informed Reinforcement Learning Approach

Our method augments a multi-agent reinforcement learning (RL) framework with D* Lite path planning signals. Each agent learns a policy that incorporates

both local grid observations and a guided path reference from D* Lite, enabling more efficient navigation.

## 2.1. Experiment Set – Path Planning Environment

We use a custom multi-agent environment based on `MultiAgentEnv` in RL-lib [7], where agents navigate a 2D grid with obstacles. The environment is defined as follows:

- $S$ is the state space representing the grid world configuration, where $s \in S$ consists of:

  - Agent positions: $P = \{p_1, p_2, ..., p_n\}$ where $p_i \in \mathbb{Z}^2$
  - Goal positions: $G = \{g_1, g_2, ..., g_n\}$ where $g_i \in \mathbb{Z}^2$
  - Obstacle positions: $B \subset \mathbb{Z}^2$

- $A$ is the action space, where $a_i \in \{up, down, left, right, wait\}$ for each agent $i$

- $T : S \times A \to S$ is the transition function mapping states and actions to new states

- $R : S \times A \times S \to \mathbb{R}$ is the reward function defined as:

$$R(s, a, s') = \begin{cases} r_{goal} & \text{if agent reached goal} \\ -r_{collision} & \text{if collision occurred} \\ -r_{step} & \text{if action is movement} \\ -r_{step} \cdot r_{wait} & \text{if action is wait} \end{cases} \tag{1}$$

- $O : S \to O$ is the observation function that maps global states to agent-specific observations:

$$o_i = \{N_i, d_i, D_i\} \tag{2}$$

where:

  - $N_i$ is the local grid observation (size $w \times w$) centered on agent $i$
  - $d_i \in \mathbb{Z}^2$ is the distance vector to goal for agent $i$
  - $D_i$ is the D* Lite path information for agent $i$, cropped to the observation window

This setup captures local inter-agent interactions and global path guidance.

## 2.2. Reinforcement Learning Pipeline

We train the agents using a shared-policy Proximal Policy Optimization (PPO) setup, where all agents invoke the same neural network parameters to map their observations to actions. The environment's observation for each agent contains local grid features and D\* Lite path information, enabling the policy to blend global guidance with collision avoidance. Over repeated episodes, the PPO algorithm updates the shared policy to maximize cumulative reward, ultimately generating multi-agent behaviors that exploit D\* Lite insights to navigate more efficiently.

# 3. Tests and Results

We conducted experiments to compare three variants:

1. PPO without D\* Lite path features,

2. PPO with D\* Lite path features in the observation,

3. pure D\* Lite baseline with no RL,

assessing both training progress and trained agent performance under identical environment configurations.

## 3.1. Experimental Setup

The first two variants were trained in a multi-agent grid scenario featuring a fixed number of agents and a consistent map (presented in Figure 1*a*) for **2,000 steps**, to allow the agent to explore better policies that might not have been discovered yet, even though we observed stabilization and no further improvements after about **750 steps**. We repeated the training five times to account for stochastic variation. In the next step, we evaluated each trained policy over ten repeats. We also conducted the same evaluation for a scenario where each agent simply follows the D\* Lite path. Additionally, we measured inference time by capturing the average runtime per environment step in both the D\*-Lite-informed PPO and standard PPO settings.

## 3.2. Quantitative Results

Training of both PPO variants is presented in Figure 1*b* as the plot of the mean duration of the episode (the time required by all robots to reach the goal positions)

Figure 1. Map that was used in the experiments with example agent and goal (a): blue *A* is the Agent, green *G* is the agent's Goal, black dots are the path from D* Lite, and the red border is the agent's observation window and training results (b): mean episode length for each training step.

for each training step. It is visible that incorporating D* Lite path signals yields a significant improvement in the convergence time and the final quality of the policy. Table 1 compares the mean episode rewards and duration across the three variants (averaged over all runs and seeds).

Table 1. Episode Reward (per agent) and Length Comparison (Means Over All Seeds) for 20 agents. After 200 steps the episodes were terminated.

| Method | Mean Reward | M. Rew. $\sigma$ | Mean Length | M. Len. $\sigma$ |
|---|---|---|---|---|
| No D* Lite | 8.758 | 0.094 | 136.659 | 2.665 |
| With D* Lite | **9.055** | 0.052 | **87.779** | 1.274 |
| Pure D* Lite | -13.734 | 0.121 | 200.0 | 0.0 |

We also tracked the inference time to evaluate the computational overhead of observing the D* Lite path. Table 2 shows the average time spent per environment step (in milliseconds). While adding path signals slightly increases observation size, we found it did not significantly degrade PPO inference performance.

Table 2. Inference Time per Environment Step (Milliseconds)

| Method | Inference Time (ms) | $\sigma$ |
|---|---|---|
| PPO (no D* Lite) | 1.078 | 0.051 |
| PPO (with D* Lite) | 1.232 | 0.048 |

# 4. Conclusions

Our findings indicate that reinforcing the policy with global navigational cues from D* Lite leads to faster training convergence, higher episodic rewards, and only a minor increase in per-step inference time. By integrating local collision avoidance with an informed path signal, the best-performing agent more effectively balances grid movement efficiency and obstacle evasion. Consequently, the PPO variant augmented by D* Lite information is recommended for scenarios where multi-agent collision avoidance must be complemented by global path efficiency. Moreover, integrating global information into each agent's local observation space does not increase communication overhead, introduce potential delays, and complexity of accurately aggregating and disseminating information among agents. The computation is performed independently on each platform, and the inference time increases only by 14%.

Future work will explore evaluating DQN and other RL methods, as well as scaling up to larger grid maps or higher agent counts to further stress-test the scalability and robustness of the approach.

# Acknowledgment

# References

[1] Sharon, G., Stern, R., Felner, A., and Sturtevant, N. R. Conflict-based search for optimal multi-agent pathfinding. *Artificial Intelligence*, 219:40–66, 2015. doi:10.1016/j.artint.2014.11.006.

[2] Wagner, G. and Choset, H. M*: A complete multirobot path planning algorithm with performance bounds. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3260–3267. 2011. doi: 10.1109/IROS.2011.6095022.

[3] Panov, A. I., Yakovlev, K. S., and Suvorov, R. Grid path planning with deep reinforcement learning: Preliminary results. *Procedia Computer Science*, 123:347–353, 2018. doi:10.1016/j.procs.2018.01.054.

[4] Damani, M., Luo, Z., Wenzel, E., and Sartoretti, G. PRIMAL$_2$: Pathfinding via reinforcement and imitation multi-agent learning - lifelong. *IEEE Robotics and Automation Letters*, 6(2):2666–2673, 2021. doi:10.1109/LRA.2021.3062803.

[5] Zhang, Y., Qian, Y., Yao, Y., Hu, H., and Xu, Y. Learning to cooperate: Application of deep reinforcement learning for online AGV path finding. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '20, page 2077–2079. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 2020.

[6] Li, J., Sun, K., Ma, H., Felner, A., Kumar, T., and Koenig, S. Moving agents in formation in congested environments. In *Proceedings of the International Symposium on Combinatorial Search*, volume 11, pages 131–132. 2020.

[7] Liang, E., Liaw, R., Nishihara, R., Moritz, P., Fox, R., Gonzalez, J., Goldberg, K., and Stoica, I. Ray RLlib: A composable and scalable reinforcement learning library. *arXiv preprint arXiv:1712.09381*, 2017.

# A Data-Driven Approach to Flatness: Learning a Latent Representation for the Unicycle Model

**Joanna Piasek-Skupna**[0000−0001−6621−8879]

*Poznan University of Technology*
*Faculty of Control, Robotics & Electrical Engineering*
*Piotrowo 3A, 60-965 Poznań, Poland*
*joanna.piasek@put.poznan.pl*

**Abstract.** *This paper presents a data-driven approach to learning a latent representation of a dynamical system that satisfies the conditions of differential flatness. By leveraging an encoder-decoder neural network, structured constraints are imposed, including Brunovsky-form evolution, latent consistency, and reconstruction losses. This eliminates the need for numerical derivative approximations while ensuring accurate trajectory representation. This method is applied to a unicycle robot model, demonstrating that the learned latent variables effectively reconstruct system states and inputs. Results suggest that neural networks can infer flat representations from data, providing a scalable alternative to analytical derivations in nonlinear control systems.*
**Keywords:** *differential flatness, control theory, latent representation*

## 1. Introduction

Differential flatness is a property of certain dynamical systems that significantly simplifies trajectory generation, making it easier to compute optimal and dynamically feasible paths. When a system is differentially flat, there is no need to explicitly integrate its equations of motion to reconstruct the system states. Moreover, this property reduces the complexity of constraints in optimization frameworks, making trajectory planning more efficient. A system is differentially flat if there exists a flat output from which all states and control inputs can be explicitly expressed as functions of the output and its derivatives. Despite the theoretical advantages of differential flatness, identifying an explicit mapping between the flat

output and the system's states and inputs is often challenging. Traditional analytical approaches can be complex or even impractical for many real-world systems. To overcome this, data-driven and numerical methods are becoming increasingly popular. By leveraging machine learning, optimization, and numerical techniques, these approaches provide practical alternatives when analytical derivations are infeasible, enabling efficient trajectory planning for complex systems [1, 2, 3].

## 2. Differential Flatness

A dynamical system of the form:

$$\dot{x} = f(x, u)$$

where $x \in \mathbb{R}^n$ represents the system states and $u \in \mathbb{R}^m$ denotes the vector of control inputs, is said to be differentially flat if there exists a set of flat outputs $z$ such that all system states $x$ and control inputs $u$ can be expressed as functions of the flat outputs $z$ and a finite number of their derivatives:

$$x = f(z, \dot{z}, \ddot{z}, \ldots, z^{(p)}),$$
$$u = g(z, \dot{z}, \ddot{z}, \ldots, z^{(q)}),$$

for some finite integers $p$ and $q$.

If the system can be fully parameterized by its flat outputs, then trajectory planning can be performed directly in terms of $z$, simplifying motion planning and feedback control design. An alternative definition states that a nonlinear control system is differentially flat if it is dynamically equivalent to a linear controllable system, such as in Brunovsky's normal form [4].

## 3. Data-Driven Approach to Differential Flatness

Traditionally, differential flatness is determined analytically using system dynamics and algebraic methods. However, for complex nonlinear systems where explicit solutions are difficult to derive, data-driven approaches provide an alternative for identifying flat outputs and parameterizing the system.

This work proposes a data-driven method to learn a representation that satisfies differential flatness conditions. The approach employs an encoder-decoder neural

168

network, where the encoder maps system states and inputs to a latent representation, and the decoder reconstructs the original states and inputs. To ensure the latent space aligns with differential flatness principles, the following constraints are imposed:

- Brunovsky Form Constraint: The latent variables evolve according to a discrete-time Brunovsky's system, embedding the system's evolution within its structure. This eliminates the need for numerical differentiation, reducing errors and enhancing stability.

- Latent Consistency Loss: A loss function enforces that specific latent variables correspond to the discrete derivatives of others, preserving the relationships expected in a differentially flat system.

- Reconstruction Loss: The decoder ensures that the system states and inputs can be accurately recovered from the learned latent representation, preventing information loss.

# 4. Example: Unicycle Robot

To evaluate the proposed approach, it is applied to a unicycle robot model, which is a widely studied example of nonlinear control. The unicycle's dynamics are given by:

$$\dot{x} = v\cos(\theta), \quad \dot{y} = v\sin(\theta), \quad \dot{\theta} = \omega,$$

where $x, y$ represent the position coordinates, $\theta$ is the heading angle, and the control inputs are the linear velocity $v$ and angular velocity $\omega$. This system is known to be differentially flat, making it a suitable testbed for the proposed method. To generate training data, the unicycle model is simulated over 30,000 time steps, with control inputs varying over time to ensure diverse trajectories.

A neural network-based encoder-decoder architecture is employed to learn a structured latent representation of the system while enforcing Brunovsky normal form constraints. The encoder maps the system states and control inputs $(x, y, \theta, v, \omega)$ into a lower-dimensional latent space $z$, while the decoder reconstructs the original states and controls. The network architecture consists of two fully connected layers with 128 neurons each, utilizing ReLU activations. The evolution of the latent space is constrained to follow a discrete-time Brunovsky form:

$$z_{k+1} = Az_k + Bv_k$$

where:

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad B = \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

To ensure that the latent representation aligns with the system dynamics, three key loss functions are defined. The Brunovsky loss enforces that the latent space follows the structured evolution:

$$\mathcal{L}_B = \sum_k \|z_{k+1} - (Az_k + Bv_k)\|^2.$$

Additionally, a latent consistency loss ensures that certain latent variables accurately capture system derivatives:

$$\mathcal{L}_{\text{latent}} = \sum_k \left\| z_{2,k} - \frac{z_{1,k+1} - z_{1,k}}{\Delta t} \right\|^2 + \left\| z_{4,k} - \frac{z_{3,k+1} - z_{3,k}}{\Delta t} \right\|^2.$$

Finally, a reconstruction loss ensures that the learned representation retains all necessary information for reconstructing the original system states and inputs:

$$\mathcal{L}_{\text{recon}} = \sum_k \left( \|x_k - \hat{x}_k\|^2 + \|u_k - \hat{u}_k\|^2 \right),$$

where $\hat{x}_k, \hat{u}_k$ are the decoder outputs from $z_k$. Training is performed using an Adam optimizer with a learning rate of $10^{-3}$ over 1,000 epochs, minimizing the combined loss function.

The trained model is validated by comparing the learned latent variables with their predicted evolution under the Brunovsky structure, as shown in Figure 1a. The results confirm that the latent variables follow the expected dynamics, with $z_2$ and $z_4$ closely matching the finite difference approximations of $dz_1/dt$ and $dz_3/dt$, respectively. This validation demonstrates that the latent space effectively captures system derivatives, ensuring consistency with the theoretical structure. Additionally, as illustrated in Figure 1b, the reconstructed system states and control inputs align well with the ground truth trajectories, further confirming that the learned representation preserves the underlying system dynamics while enabling structured trajectory generation and control.

However, notable discrepancies are observed in $z_2$ and the reconstructed control input $v$, suggesting that the model struggles to accurately capture the velocity dynamics. These errors could stem from imperfect encoding of acceleration-related terms, leading to inconsistencies between the learned latent variables and

(a) Comparison of learned latent variables: $z_1, z_2, z_3, z_4$ (blue) with their predicted evolution using the Brunovsky structure (dashed orange). The bottom plots validate the latent consistency constraint: $z_2$ and $z_4$ (blue) compared to derivatives of $z_1$ and $z_3$ (dashed orange).

(b) Reconstruction of system states and control inputs. Solid lines represent the true values, while dashed lines represent the predicted values.

Figure 1. Comparison of latent variable evolution and reconstructed system states

their expected derivatives. Additionally, minor deviations in the control predictions indicate that the learned mapping from the latent space to control inputs might not be fully capturing the underlying system constraints. These imperfections suggest potential areas for refinement, such as improving the loss function to better constrain velocity-related terms, incorporating higher-order derivative constraints, or refining the network architecture to enhance the accuracy of latent space modeling.

## 5. Conclusions

The goal of this study was to identify a latent representation of a dynamical system that would be functionally equivalent to differential flatness by imposing

specific constraints during learning. Rather than manually deriving flat outputs, a data-driven approach was utilized to infer a representation that enables the reconstruction of states and inputs while preserving known structural properties.

Although the learned representation is not explicitly proven to be differentially flat, it satisfies conditions suggesting functional equivalence. Further validation is needed to assess its generalization across different system dynamics. The reliance on manually defined derivative constraints may raise concerns about system specificity, potentially limiting broader applicability. Future research will explore more adaptive methods to reduce dependence on system-specific constraints and investigate sequence-aware architectures to better capture temporal dependencies. Additionally, the methodology can be extended to control applications by designing trajectories within the learned latent space, enabling data-driven motion planning and optimization.

# Acknowledgment

# References

[1] Gadginmath, D., Krishnan, V., and Pasqualetti, F. Data-driven feedback linearization using the Koopman generator. *IEEE Transactions on Automatic Control*, 69(12):8844–8851, 2024.

[2] Greeff, M. and Schoellig, A. P. Exploiting differential flatness for robust learning-based tracking control using Gaussian processes. *IEEE Control Systems Letters*, 5(4):1121–1126, 2020.

[3] Sferrazza, C., Pardo, D., and Buchli, J. Numerical search for local (partial) differential flatness. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3640–3646. IEEE, 2016.

[4] Nicolau, F. and Respondek, W. Flatness of multi-input control-affine systems linearizable via one-fold prolongation. *SIAM Journal on Control and Optimization*, 55(5):3171–3203, 2017.

# On the Importance of Camera Pose Estimation in NeRF-Based Super Resolution Tasks

**Mikołaj Zieliński**[0000−0003−1107−1149], **Eldaniz Babayev**[0009−0006−6628−8969], **Dominik Belter**[0000−0003−3002−9747]

*Poznan University of Technology*
*Institute of Robotics and Machine Intelligence*
*ul. Piotrowo 3A, 60-965 Poznań, Poland*
*mikolaj.zielinski@doctorate.put.poznan.pl*
*eldaniz.babayev@student.put.poznan.pl*
*dominik.belter@put.poznan.pl*

**Abstract.** *This paper investigates the importance of accurate camera pose estimation in super-resolution tasks utilizing Neural Radiance Fields (NeRF) for robotic systems. We highlight a critical flaw in existing evaluation pipelines, where camera poses are computed from high-resolution images and reused for downscaled inputs, failing to mirror real-world scenarios. Our proposed pipeline computes camera poses directly from low-resolution images, demonstrating that the omission of this step leads to overly optimistic results. Experiments reveal that while NeRF exhibits inherent upscaling properties, pose estimation quality deteriorates significantly with increased downscaling, underscoring the need for realistic evaluation methods in practical robotic applications.*
**Keywords:** *robotics, neural rendering, super-resolution*

## 1. Introduction

Mobile robots are tasked with performing a variety of actions, some of which require high precision, such as inserting a peg into a hole, fastening in a screw, or placing a plate between dish rack rods. A key limitation in these tasks is the quality of spatial representation, particularly for fine details such as screws or small holes, which can become invisible when reconstructed scenes are of low resolution. This

issue is often exacerbated by the limitations of the cameras used and the constraints on data transfer and processing capacity in robotic systems.

Recent advancements in scene representation, particularly through such techniques as Neural Radiance Fields (NeRF) [1], show promise due to their inherent upscaling capabilities. Several approaches have leveraged NeRF to improve scene reconstruction [2, 3, 4]. However, a crucial aspect often overlooked in these systems is accurate camera pose estimation. NeRF's reconstruction quality heavily depends on the precision of the camera pose, which is typically estimated using COLMAP [5, 6]. A common yet flawed practice in current methods is to reuse poses calculated from high-resolution images when applying NeRF to downscaled data. This practice neglects the fact that low-resolution images, which contain fewer features, lead to less accurate pose estimation and may even prevent pose estimation altogether.

In this paper, we demonstrate that accurate camera pose estimation is critical for successful upscaling in NeRF-based systems. We argue that reusing poses from high-resolution images when working with downscaled data does not reflect real--world conditions, where poses would need to be recalculated for low-resolution images.

## 2. Methodology

Many super-resolution methods using NeRF [2, 3, 4] rely on the evaluation pipeline shown at the top in Figure 1. While these methods report results for both synthetic and real-world datasets, our work focuses exclusively on real-world datasets without explicitly known camera poses. Processing of the data proceeds as follow. First, camera poses are estimated for each image in the base dataset. Afterwards, all the images are downscaled by as specific factor. Subsequently these methods use those images for training super-resolution systems. Finally, new images are rendered with the initial resolution for comparison with the base dataset. Although effective for training, this approach is unrealistic for evaluation, as real--world scenarios often lack access to high-resolution images for pose estimation. While this method has been widely adopted, it does not account for the degradation in pose accuracy that occurs when working with lower-resolution data.

To address this issue, we propose a more realistic evaluation pipeline, showcased at the bottom in Figure 1, which better mirrors the conditions faced in practical robotic systems. Our proposed pipeline first downscales the images to the

Figure 1. Proposed pipeline for evaluation of the upscaling properties of the NeRF models. In contrast to methods from [2, 3, 4], we perform realistic robotic scenario where camera pose estimation using COLMAP is performed on the lower resolution images.

desired resolution and then computes the camera poses from these low-resolution images. This approach ensures that the poses are calculated in the same way as they would be in real-world settings, where low-resolution images are often the only available input. This methodology ensures that the challenges of degraded pose estimation and low-resolution data are properly considered, offering a more reliable assessment of NeRF models for practical robotic applications where the camera pose estimation is necessary.

## 3. Experiments

We selected Nerfacto [7], one of the most popular NeRF models, as the baseline for our upscaling experiments. As highlighted in [4], NeRF inherently possesses super-resolution properties, making it a suitable baseline for our study.

We performed experiments using 21 challenging real-world image-sets from Tanks and Temples [8], with images at a resolution of $1920 \times 1080$. Each dataset was downscaled by factors of $\times 2$, $\times 4$, $\times 8$ and $\times 16$, following the two distinct scenarios outlined in Figure 1. In the first scenario, we applied the standard approach, where camera poses estimated from high-resolution images were reused for training NeRF on downscaled data. In contrast, our proposed approach recalculated

camera poses for each downscaled image set, thereby more accurately simulating real-world conditions.

To evaluate the quality of the upscaled images, we utilized standard metrics, including PSNR, SSIM, and LPIPS. These metrics were computed by comparing the upscaled outputs from both scenarios with the original images in the dataset.

# 4. Results

During data preparation, we evaluated the number of images for which COLMAP successfully estimated camera poses, determining the proportion retained in the final dataset. For each scale $k$, we computed the ratio of images with correspondences to the total dataset size (Table 1). Data loss is negligible for $\times 1$, $\times 2$, and $\times 4$, but increases significantly at $\times 8$ ( 20% loss) and $\times 16$ (>70% loss).

Table 1. Average correspondences found by COLMAP

| k | AvgCorrsFound |
|---|---|
| ×1 | 99.67% |
| ×2 | 99.69% |
| ×4 | 98.33% |
| ×8 | 79.12% |
| ×16 | 27.66% |

We trained the Nerfacto model under two scenarios: the standard pipeline from the literature [2, 3, 4] and our proposed pipeline (Figure 1). Table 2 summarizes the results, showing the percentage change in each metric compared to ground truth images.

Table 2. Average metrics for each downscale factor

| k | PSNR$_{base}$ | PSNR$_{our}$ | SSIM$_{base}$ | SSIM$_{our}$ | LPIPS$_{base}$ | LPIPS$_{our}$ |
|---|---|---|---|---|---|---|
| ×2 | 0.04% | 0.39% | 0.87% | 0.32% | 0.11% | 0.91% |
| ×4 | 1.26% | -5.75% | 1.54% | -7.27% | -6.98% | -12.18% |
| ×8 | 0.71% | -22.18% | -1.52% | -25.45% | -20.21% | -31.22% |
| ×16 | -3.39% | -56.06% | -13.73% | -63.53% | -29.21% | -57.61% |

To analyze the impact of pose estimation and missing data, we conducted an ablation study. We created datasets using camera poses from the first scenario and reduced images from the second. After training, we computed the metrics as before. To evaluate the impact of the number of images in the set, we subtract the metrics obtained in the first scenario from those of this study. Similarly, to assess the influence of camera poses, we subtract these metrics from those calculated in the third scenario. The results are presented in Table 3.

Table 3. Influence of pose and reduced number of images $N$ on the final results

| k | PSNR | | SSIM | | LPIPS | |
|---|---|---|---|---|---|---|
| | $N$ | Pose | $N$ | Pose | $N$ | Pose |
| ×2 | -0.24% | 0.58% | -0.80% | 0.26% | -0.26% | 1.07% |
| ×4 | -0.95% | -6.06% | -1.35% | -7.46% | -0.30% | -4.90% |
| ×8 | -10.18% | -12.71% | -10.10% | -13.82% | -5.79% | -5.22% |
| ×16 | -49.76% | -2.90% | -40.83% | -8.97% | -26.26% | -2.14% |

From the baseline pipeline results, we observe that for the initial downscale factors, certain metrics indicate an improvement in quality. This aligns with the findings of [4], which suggest that NeRF inherently possesses upscaling capabilities. Notably, we successfully trained NeRF on images downscaled by a factor of ×16, achieving accurate scene reconstruction. However, when evaluating our proposed approach, a significant degradation in quality becomes apparent as the downscale factor increases. Our ablation study further reveals that image count has little effect for ×2 and ×4, but becomes significant at higher scales. Pose estimation errors also contribute to degradation beyond ×2.

# 5. Conclusions

Our results emphasize the need to reconsider current evaluation pipelines for super-resolution tasks, as they often fail to capture real-world challenges in low-resolution data acquisition. Two key issues emerge: the degradation of camera pose estimation quality and the loss of images during reconstruction. COLMAP struggles to estimate poses for downscaled images, leading to fewer images being used in the reconstruction process. This dual challenge, both poor pose quality and reduced image availability, significantly impacts the accuracy of reported results in existing literature, including studies such as [2, 3, 4].

While promising results may be achieved in controlled settings, these methods may not perform as well in practical scenarios, where low-resolution images and limited pose estimation are the norm. To ensure that super-resolution methods are reliable and applicable in real-world environments, future research should incorporate realistic camera pose estimation into the evaluation pipeline. This would provide a more accurate assessment of upscaling performance, reflecting both pose estimation challenges and the reduced number of images in practical robotic applications.

## Acknowledgment

## References

[1] Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., and Ng, R. NeRF: Representing scenes as neural radiance fields for view synthesis. *ECCV*, 2020.

[2] Huang, X., Li, W., Hu, J., Chen, H., and Wang, Y. RefSR-NeRF: Towards high fidelity and super resolution view synthesis. In *CVPR*. IEEE, 2023.

[3] Huang, D.-J., Chou, Z.-T., Wang, Y.-C. F., and Sun, C. ASSR-NeRF: Arbitrary-scale super-resolution on voxel grid for high-quality radiance fields reconstruction. *arXiv preprint arXiv:2406.20066*, 2024.

[4] Wang, C., Wu, X., Guo, Y.-C., Zhang, S.-H., Tai, Y.-W., and Hu, S.-M. NeRF-SR: High quality neural radiance fields using supersampling. In *ACM Multimedia*, pages 6445–6454. ACM, 2022. ISBN 978-1-4503-9203-7.

[5] Schönberger, J. L., Zheng, E., Pollefeys, M., and Frahm, J.-M. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*. 2016.

[6] Schönberger, J. L. and Frahm, J.-M. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016.

[7] Team, N. Nerfacto, 2022. URL `https://docs.nerf.studio/nerfology/` `methods/nerfacto.html`.

[8] Knapitsch, A., Park, J., Zhou, Q.-Y., and Koltun, V. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017.

# Problem Solving and Optimization

Track Chairs:

- prof. Jarosław Arabas – Warsaw University of Technology

- prof. Karol Opara – Systems Research Institute, Polish Academy of Sciences

- prof. Szymon Łukasik – AGH University of Science and Technology

# Unsupervised Physics-Constrained Inverse Problem Solving in Electrical Capacitance Tomography

**Mikhail Ivanenko**[0009−0006−8682−2751],
**Damian Wanta**[0000−0002−1596−6524],
**Waldemar T. Smolik**[0000−0002−1524−5049],
**Przemysław Wróblewski**[0000−0002−6713−9088],
**Mateusz Midura**[0000−0002−2449−0652]

*Warsaw University of Technology*
*Faculty of Electronics and Information Technology*
*Nowowiejska 15/19, 00-665 Warsaw, Poland*
*Damian.Wanta@pw.edu.pl*

**Abstract.** *Electrical capacitance tomography (ECT) aims to visualize a distribution of electrical properties inside an object based on non-invasive capacitance measurements. A supervised approach to neural network training in ECT image reconstruction has shown very promising results. However, it is of interest if it is possible to create a fully unsupervised reconstruction algorithm (inverse problem solver) based on constraints forced by laws of physics. We trained a forward problem solver representing physics constraints and were able to train an inverse problem solver guided only by a forward problem solver and real ECT measurements. We were able to successfully reconstruct test object motion inside a 16-electrode sensor.*
**Keywords:** *electrical impedance tomography, inverse problem, neural networks, image reconstruction, unsupervised machine learning, deep learning, physics-constrained machine learning*

## 1. Introduction

Electrical capacitance tomography (ECT) is a non-invasive and low-cost tomographic technique that reconstructs electrical permittivity and conductivity inside an examined object based on capacitance measurement between electrodes

surrounding that object. Electrical properties reconstruction is usually referred to as an inverse problem, whereas calculating capacitance measurements based on given electrical properties is known as a forward problem. One of the ECT applications is observing multiphase flow, which requires high temporal resolution and a fast image reconstruction method. Recent progress in neural network development allowed the creation of efficient reconstruction algorithms based on machine learning. However, most such algorithms implement a supervised approach, which requires generating a relatively large dataset to train the model [1–4]. Such a dataset should contain samples of electrical properties distribution and corresponding measurements. To generate a dataset, it is necessary to define tomograph geometry, generate enough electrical properties distribution samples, and then calculate capacitance measurements corresponding to each sample using Gauss law and Geselowitz equations:

$$\nabla \cdot \varepsilon \nabla \varphi = 0, \tag{1}$$

$$\frac{\partial C_{AB}}{\partial \varepsilon_n} = -\frac{1}{U^2} \int_\Omega \nabla \varphi_A \nabla \varphi_B dv, \tag{2}$$

where $C_{AB}$ – capacitance between electrodes A and B, $\varepsilon_n$ – electrical permittivity in $n$-th small volume element, $U$ – voltage applied between electrodes, $\varphi_A$, $\varphi_B$ – electrical potentials when electrode A and B is excitation electrode, respectively. Despite great results, this approach has some drawbacks. Dataset generation requires heavy computations, taking significant time. Additionally, it is necessary to create a separate machine-learning model for each tomograph geometry and supposed objects.

Recent developments in physics-constrained machine learning [5, 6] open up new possibilities to enhance a forward problem solution using neural partial differential equations solvers. In this paper, we explore the potential to train an inverse problem solver on real ECT measurements using an unsupervised approach and a simple forward problem solver. It has been shown that it is possible to enhance reconstruction quality using unsupervised fine-tuning of the already trained inverse problem solver [7]. We are going to demonstrate that it is possible to obtain satisfactory reconstruction results without having a pretrained reconstruction model.

## 2. Materials and Methods

## 2.1. Forward Problem Solver

In this work, we trained a forward problem solver neural network using a synthetic dataset of 150,000 samples. The dataset consists of the images with random ellipses, and images with four flow patterns [8]. Each sample represents $32 \times 32$ electrical permittivity distribution and 120 capacitance measurements from a 16-electrode sensor. We used a UNet-based architecture containing five convolution and five deconvolution blocks with skip connections. Each of the convolution blocks consists of one 2D convolution layer, batch normalization layer, and rectified linear unit activation function. Each deconvolutional block consists of one transposed 2D convolution layer, batch normalization layer, and rectified linear unit activation function. A linear layer is attached to the network output to reduce the number of outputs to the desirable 120 capacitance values. Before conducting the network training, the capacitance measurements in the dataset were normalized to the range from 0.0 to 1.0 using measurements corresponding respectively to the empty sensor and the sensor filled with the high electrical permittivity material. We used the mean squared error loss function and trained the network for 50 epochs on 75% of the dataset (Figure 1(a)). After each epoch, we calculated the mean norm error on the remaining 25% of the dataset (Figure 1(b)).



(a)                               (b)

Figure 1. Forward problem training results. Mean square error norm as a function of the number of epochs: (a) training dataset (logarithmic scale), (b) validation dataset

## 2.2. Measurements Collection

We collected capacitance measurements using the EVT4 acquisition system developed by our team [9]. It allows measurement of capacitances with a frame rate of 1,000 for a 16-electrode sensor and achievement of high signal-to-noise ratio from 30 to 65 dB for capacitances between 1 fF and 1 pF. As a test object, we used a clover-shaped 3D printed structure shown in Figure 2. The object was placed into the 16-electrode sensor and rotated manually until 2,420 measurements were collected. This method makes it is easy to ensure that collected data represent the changing sensor content.



Figure 2. Test object used for measurements



Figure 3. Inverse problem solver residual block: linear – linear layers, conv – convolutional layer

## 2.3. Inverse Problem Solver

As an inverse problem solver, we used a neural network consisting of three residual blocks inspired by [7]. Each block (Figure 3) consists of two consequent linear layers and one linear layer, making a skip connection. At the end of the block, one convolutional layer is attached with one filter, kernel size of five, and padding of two, making the size of the convolutional layer input equal to the size of its output. Additionally, the second and third residual blocks have network input passed via a linear layer added to the output of the first block's linear layer. After each layer, the leaky rectified linear unit activation function is applied with a negative slope of 0.1. The layer sizes are given in Table 1. Before the first

Table 1. Inverse problem solver residual block size

| N | Input | Middle | Output | Convolution |
|---|-------|--------|--------|-------------|
| 1 | 128 | 256 | 256 | $16 \times 16$ |
| 2 | 256 | 576 | 576 | $24 \times 24$ |
| 3 | 576 | 1024 | 1024 | $32 \times 32$ |

residual layer, there is one linear layer transforming the 120-capacitance input into a 128-element vector.

We were training the inverse problem solver using the Adam optimizer and cosine annealing learning rate scheduler for 4,000 epochs. All 2,420 capacitance measurements were used as input. The reconstructed images were used as input for the forward problem solver, which had its weights locked. The discrepancy between measured and calculated capacitance values was calculated using the mean square error norm. This difference (Figure 4) was used to update the weights of the inverse problem solver.



Figure 4. Loss function value during inverse problem solver training (as a function of the number of epochs)

# 3. Results

Using the described approach, we successfully trained the inverse problem solver for the given object and sensor geometry. To increase signal to noise ratio, we converted 2,420 measurements collected during object rotation into 242

by averaging every 10 measurements. We used such measurements to reconstruct electrical permittivity images assuming that high frame rate mitigates object's rotation during 10 frames. Then, we applied a Gaussian blur filter with a kernel size of three. The resulting images are shown in Figure 5. We can conclude that we indeed are able to recognize the test object (shown in Figure 2) in different angular positions.



Figure 5. The reconstructed cross-sectional images of electrical permittivity. The slices of the test object are in a different angular position in subsequent frames.

## 4. Discussion and Further Development

We found that the inverse problem solver network can be trained while being guided only by residual error calculated using a forward problem solver. That opens the way to developing a fully unsupervised reconstruction algorithm provided that a fast neural forward problem solver can be created. Developing such a solver should be an easier task than developing a universal inverse problem solver, because in ECT the forward problem is a well-posed problem contrary to the inverse problem, which is ill-posed. Additionally, the forward problem is fully solvable analytically, which makes the task even more realistic.

## References

[1] Deabes, W. and Jamil Khayyat, K. M. Image reconstruction in electrical capacitance tomography based on deep neural networks. *IEEE Sensors Journal*, 21(22):25818–25830, 2021. doi:10.1109/JSEN.2021.3116164.

[2] Wang, P., Lin, J.-S., Wang, M., and Zhao, Y.-L. An image reconstruction algorithm for electrical capacitance tomography. In *Proceedings of the Second International Conference on Innovative Computing and Cloud Computing*, pages 96–101. 2013. doi:10.1145/2556871.2556893.

[3] Zhang, Y. and Chen, D. Image reconstruction for high-performance electrical capacitance tomography system using deep learning. *Complexity*, 2021(1):5545491, 2021. doi:10.1155/2021/5545491.

[4] Deabes, W., Abdel-Hakim, A. E., Bouazza, K. E., and Althobaiti, H. Adversarial resolution enhancement for electrical capacitance tomography image reconstruction. *Sensors*, 22(9):3142, 2022. doi:10.3390/s22093142.

[5] Tan, L. and Chen, L. Enhanced deeponet for modeling partial differential operators considering multiple input functions. *arXiv preprint arXiv:2202.08942*, 2022.

[6] Molinaro, R., Yang, Y., Engquist, B., and Mishra, S. Neural inverse operators for solving PDE inverse problems. *arXiv preprint arXiv:2301.11167*, 2023.

[7] Jin, Y., Li, Y., Zhang, M., and Peng, L. A physics-constrained deep learning-based image reconstruction for electrical capacitance tomography. *IEEE Transactions on Instrumentation and Measurement*, 73:1–12, 2023. doi:10.1109/TIM.2023.3338673.

[8] Ivanenko, M., Wanta, D., Smolik, W. T., Wróblewski, P., and Midura, M. Generative-adversarial-network-based image reconstruction for the capacitively coupled electrical impedance tomography of stroke. *Life*, 14(3):419, 2024. doi:10.3390/life14030419.

[9] Kryszyn, J., Wróblewski, P., Stosio, M., Wanta, D., Olszewski, T., and Smolik, W. T. Architecture of EVT4 data acquisition system for electrical capacitance tomography. *Measurement*, 101:28–39, 2017. doi:10.1016/j.measurement.2017.01.020.

# How Powerful Are Classic Graph Neural Networks for Malware Detection? A Case Study with Cartesian Genetic Programming

**Maciej Krzywda**[1]**, Szymon Łukasik**[1,2]**, Amir H. Gandomi**[3,4,5]

[1]*Faculty of Physics and Applied Computer Science*
*AGH University of Krakow, al. Mickiewicza 30, 30-059 Kraków, Poland*
[2]*Systems Research Institute, Polish Academy of Sciences*
*ul. Newelska 6, 01-447 Warsaw, Poland*
[3] *Faculty of Engineering and IT, University of Technology Sydney*
*5 Broadway, Ultimo NSW 2007, Australia*
[4]*University Research and Innovation Center (EKIK), Óbuda University*
*Bécsi út 96/B, Budapest, 1034, Hungary*
[5]*Department of Computer Science, Khazar University*
*Mahsati 41, Baku, Azerbaijan*

**Abstract.** *The present study explores an approach to Neural Architecture Search (NAS) using Cartesian Genetic Programming (CGP) for the design and optimization of Graph Neural Networks (GNN) for malware detection. A key aspect of this innovative method is proposing a novel neural architecture. Traditionally, architectures have been manually crafted by human experts, which is both time-consuming and prone to errors. In this work, we employ a pure Cartesian Genetic Programming approach (using the 1+λ strategy with 1 and 4 children), utilizing only one genetic operation – mutation. Preliminary experiments indicate that our methodology yields promising results.*
**Keywords:** *graph neural networks, neural architecture search, Cartesian genetic programming, malware detection*

## 1. Introduction

Neural Architecture Search (NAS) [1] has become increasingly popular as a fully automated method for designing neural network architectures. This ap-

proach enables the creation of architectures that are not only on par with but often outperform those developed manually. Essentially, NAS streamlines the conventional process where humans iteratively tweak neural networks through trial and error to find effective configurations. Instead, it automates this procedure, revealing more complex structures. NAS encompasses a variety of techniques and tools that systematically assess numerous network architectures within a defined search space. Using a search strategy, you select the architecture that best meets the objectives of a specific problem by optimizing a fitness function. Despite its advantages, NAS presents significant challenges related to computational resources and time, further complicated by the high costs of using graphics processing units (GPUs). As a result, researchers and research teams are increasingly seeking alternative approaches to reduce costs and identify the most efficient and effective neural network architectures tailored to their particular research challenges. This paper seeks to explore and demonstrate the design of graph neural networks using Cartesian Genetic Programming (CGP). Cartesian Genetic Programming is well regarded for its success in designing Convolutional Neural Networks, delivering promising outcomes [2, 3, 4]. Building on these achievements, we have chosen to apply CGP to the development of Graph Neural Networks.

## 2. Cartesian Genetic Programming for Design GNN

Cartesian Genetic Programming (CGP) typically relies on mutation as its sole genetic operator; crossover is still a research topic [5]. Possible mutations include point, gene, and segment mutation [6]. Here, CGP is used to design convolutional neural networks (CNNs) by evolving the network architecture rather than relying on a predefined solution. Initially, a parent genotype encodes which neural network layers are active. Two offspring are produced by mutating the parent's genotype at a given rate. A randomly selected layer is replaced with another from a predefined set (avoiding the same layer or deactivating essential layers). We ensure that each mutation leads to a valid, compilable network. Each offspring is then evaluated against the parent: if it performs better (based on F1-score on the validation set), it becomes the new parent for the next generation. Otherwise, the parent remains. This process repeats over multiple generations, recording the best genotype found. If none of the offspring improves upon the parent for several consecutive generations, a neutral mutation is applied, and evolution continues. The fitness function (F1-score) quantifies network performance. By iteratively

mutating and selecting genotypes, CGP explores different architectures to maximize this fitness, ultimately converging on an optimal model for the dataset.

For GNN design specifically, CGP offers distinct advantages over reinforcement learning (RL) approaches. While RL-based NAS uses controller networks that require extensive training and computational resources to generate architectures, CGP directly encodes network topology and GNN-specific operations (message passing, aggregation functions, and pooling strategies) in its graph-based representation. Unlike RL methods that often optimize for single objectives through complex reward functions, CGP naturally supports multi-objective optimization and maintains beneficial, neutral drift – crucial for exploring the complex GNN design space where minor architectural changes significantly impact performance. This evolutionary approach enables CGP to efficiently discover novel GNN architectures that might remain hidden from human designers and more computationally intensive RL methods.

## 3. Experiments

### 3.1. Malware Detection

MalNetTiny comprises 5,000 function call graphs (FCGs) that represent malicious and benign software in five distinct categories. Each graph contains up to 5,000 nodes. `MALNET-TINY` is available alongside the complete dataset at `https://mal-net.org`. Developed by Tech University in collaboration with the Microsoft APT team, MalNet-Tiny is an FCG dataset of Android malware [7, 8]. The dataset includes 4,500 malicious FCGs distributed among four different malware categories and 500 benign FCGs. We adhere to the authors' recommended split for training, validation and testing by the authors, following a ratio 70%/10%/20%.

### 3.2. Experimental Results

To ensure the robustness of our findings, we conducted an experiment using the parameters outlined in Table 2, repeating the procedure ten times. The results, which utilize Graph Neural Networks integrated with layer-wise recognition like GCN [9], GraphSAGE [10], GIN [11] are presented in Table 3 and compare our results with selected baseline solutions. Specifically, we evaluate our approach against the GNN-based method proposed in [12], which employs the Jumping-Knowledge (JK) mechanism, making it more complex than our classic GNN ar-

chitecture. This comparison underscores the effectiveness of our GNN-focused approach in the context of malware detection. In Table 1, we present the hyper-parameter values used in our computations. This includes options for operators to evolve within Cartesian Genetic Programming, such as applying mutations with GCN or choosing to apply no operation (None).

Table 1. Neural Network Hyperparameter Ranges and Choices

| Parameter | Range or Choices |
|---|---|
| Epochs | Up to 100 epochs with early stopping (patience = 10) |
| Learning Rate | Uniformly sampled from 1e-5 to 1e-1 |
| Weight Decay | Uniformly sampled from 1e-5 to 1e-1 |
| GNN Layers (Max 10) | GINConv, GraphSAGEConv, GCNConv, Dropout, BatchNorm, None |
| Classification Layers (3 Total) | Dropout, Linear, None (Last layer always Linear) |
| Activations | ReLU, Sigmoid, Tanh, None |
| Optimizers | Adadelta, Adagrad, Adam, AdamW, Adamax, ASGD, NAdam, RAdam, RMSprop, SGD |
| Loss Functions | CrossEntropyLoss, NLLLoss |

In Table 2, we present the Cartesian Genetic Programming parameters used in our experiments; in this study, we apply the evolutionary strategies 1+1 and 1+4, which are among the most commonly used approaches in this field [6].

Table 2. Cartesian Genetic Programming parameters

| Parameter | Rows | Cols | Level-Back | Mutation rate | Generation |
|---|---|---|---|---|---|
| Value | 1 | 10 | 5 | 0.30 | 200 |

## 3.3. Discussion

Among our evolved models presented in Table 3, GraphSAGE [1+4] delivers an F1-Score of 91.9%, which is very close to the GraphSAGE-JK baseline (92.0%). This indicates that the combination of the mean aggregator, the chosen hyperparameters (e.g., learning rate, optimizer), and the (1+4) evolutionary strategy helps discover a powerful architecture for the MalNet Tiny dataset. Comparing [1+1] vs. [1+4] variants shows that searching with four offspring per generation ([1+4]) can sometimes find better-performing configurations (especially

Table 3. F1-Score comparision between our apporach and state-of-the-art baselines in % sorted descending

| GNNs architectures | F1-Score |
|---|---|
| GraphSAGE-JK [12] | 92.0% |
| GIN-JK [12] | 91.0% |
| GCN-JK [12] | 90.0% |
| Our GraphSAGEs [1 + 4] | 91.9% |
| Our GCNs+GraphSAGEs+GINs [1+ 1] | 90.4% |
| Our GCNs+GraphSAGEs+GINs [1+ 4] | 90.2% |
| Our GINs [1+ 1] | 87.6% |
| Our GCNs [1 + 4] | 83.4% |
| Our GraphSAGEs [1+ 1] | 83.2% |
| Our GCNs [1+ 1] | 82.5% |
| Our GINs [1+ 4] | 81.8% |

for GCN and GraphSAGE). However, for GIN, the [1+1] approach outperforms the [1+4] variant (87.6% vs. 81.8%), suggesting that the best strategy is somewhat architecture-dependent. Our highest-performing model (GraphSAGE [1+4]) nearly matches the state-of-the-art GraphSAGE-JK (91.9% vs. 92.0%), underscoring that our CGP-based search can yield competitive designs without explicitly implementing JK aggregation [12]. The results for these hybrid configurations are as follows. **GCNs+GraphSAGEs+GINs** models achieved an F1-Score of 90.4% with the [1+1] evolutionary strategy and 90.2% with the [1+4] strategy. These outcomes suggest that integrating multiple GNN types within a single model can yield competitive performance, albeit slightly lower than the best-performing single-architecture model. The marginal difference between the [1+1] and [1+4] strategies in the hybrid configuration indicates that, for combined architectures, the evolutionary process might be less sensitive to the number of offspring per generation. This could be due to the increased complexity and diversity of the search space when multiple GNN types are involved.

## 4. Conclusions

This study focuses on the development of artificial graph neural networks (GNNs) and presents the initial phase of a novel methodology for their design and optimization using Cartesian Genetic Programming. Our findings demonstrate that classical GNNs alone achieve promising results, underscoring their potential

Table 4. Summary for the best solution: ConvSet (each set including BatchNorm and Dropout), Strategy, Execution Time, and Num Params

| ConvSet | Strategy | Exec Time (s) | Num Params |
|---|---|---|---|
| GCNConv | 1 + 4 | 201516.32 | 65925 |
| GINConv | 1 + 4 | 158835.96 | 42238 |
| GraphSAGEConv | 1 + 4 | 193410.22 | 56759 |
| GraphSAGEConv | 1 + 1 | 69224.07 | 71182 |
| GCNConv | 1 + 1 | 73065.83 | 42561 |
| GINConv | 1 + 1 | 71335.29 | 65199 |
| GINConv+GraphSAGEConv+GCNConv | 1 + 1 | 98734.62 | 183771 |
| GINConv+GraphSAGEConv+GCNConv | 1 + 4 | 295690.90 | 101882 |

in various applications. Building on this foundation, the next step involves incorporating additional mechanisms during the evolutionary process, such as Jumping Knowledge, to further enhance performance and adaptability. This work lays the groundwork for more efficient and scalable approaches in the development of graph-based neural architectures, providing valuable insight and directions for future research. A further avenue for development may involve the application of multi-objective optimization, such as searching for the best networks based on selected metrics, while also minimizing the network size (number of parameters).

# Acknowledgment

# References

[1] Elsken, T., Metzen, J. H., and Hutter, F. Neural architecture search: A survey. *arXiv preprint arXiv:1808.05377*, 2018.

[2] Krzywda, M., Łukasik, S., and Gandomi, A. H. Cartesian genetic programming approach for designing convolutional neural networks. *arXiv preprint arXiv:2410.00129*, 2024.

[3] Suganuma, M., Kobayashi, M., Shirakawa, S., and Nagao, T. Evolution of deep convolutional neural networks using Cartesian genetic programming. *Evolutionary Computation*, 28(1):141–163, 2020.

[4] Suganuma, M., Shirakawa, S., and Nagao, T. A genetic programming approach to designing convolutional neural network architectures. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '17*, pages 497–504. Association for Computing Machinery, New York, NY, USA, 2017.

[5] Husa, J. and Kalkreuth, R. A comparative study on crossover in Cartesian genetic programming. In M. Castelli, L. Sekanina, M. Zhang, S. Cagnoni, and P. García-Sánchez, editors, *Genetic Programming*, pages 203–219. Springer International Publishing, Cham, 2018.

[6] Miller, J. F. Cartesian genetic programming: Its status and future. *Genetic Programming and Evolvable Machines*, 21:129–168, 2019.

[7] Freitas, S., Dong, Y., Neil, J., and Chau, D. H. A large-scale database for graph representation learning. *arXiv preprint arXiv:2011.07682*, 2020.

[8] Freitas, S., Duggal, R., and Chau, D. H. MalNet: A large-scale image database of malicious software. *arXiv preprint arXiv:2102.01072*, 2022.

[9] Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.

[10] Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation learning on large graphs. *Advances in Neural Information Processing Systems*, 30, 2017.

[11] Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.

[12] Lo, W. W., Layeghy, S., Sarhan, M., Gallagher, M., and Portmann, M. Graph neural network-based android malware classification with jumping knowledge. In *2022 IEEE Conference on Dependable and Secure Computing (DSC)*, pages 1–9. IEEE, 2022.

# Evaluating Forecast-Based Optimization Models in Fantasy Football Team Management

**Weronika Wiechno**[0009−0002−1580−8291], **Bartosz Bartosik**[0009−0002−9699−3440], **Piotr Duch**[0000−0003−0656−1215]

*Lodz University of Technology*
*Institute of Applied Computer Science*
*Stefanowskiego 18, 90-537 Łódź, Poland*
*251661@edu.p.lodz.pl, bartosz.bartosik@dokt.p.lodz.pl,*
*pduch@iis.p.lodz.pl*

**Abstract.** *Fantasy sports have gained significant popularity worldwide, with the Fantasy Premier League attracting over 11 million players annually. In this paper, we present a framework for optimizing Fantasy Premier League squad selection, utilizing a neural network for player points prediction. Genetic Algorithm and Monte Carlo Tree Search were used for optimizing transfer strategies throughout the season. When retrospectively tested on the 2023/2024 season, the proposed approach achieved a score of up to 2,116 points, surpassing the average player's performance.*
**Keywords:** *Genetic Algorithms, Monte Carlo Tree Search, Artificial Intelligence, Fantasy Sports*

## 1. Introduction

Fantasy sports have emerged as a prominent global phenomenon, engaging millions of fans worldwide. In 2022, it was estimated that approximately 62.5 million individuals in the United States and Canada participated in these games. Among the various platforms, the Fantasy Premier League (FPL) stands as the foremost example of fantasy football, with a player base that exceeds 11 million.

Fantasy sports require participants to make strategic decisions in player selection, often under financial and positional constraints. While players often rely on intuition, recent research suggests that applying optimization techniques can significantly improve squad selection. There are several approaches to solve this

problem, namely the Multiple-Criteria Decision-Making (MCDM) approach, statistical and data-driven optimization, as well as machine learning and artificial intelligence-based.

This study leverages neural networks for player point prediction, along with genetic algorithms (GA) and Monte Carlo Tree Search (MCTS), for team management during the season. The transfers for the upcoming Gameweek are selected based on predictions made by a neural network. The model is designed to forecast the player performance in the subsequent Gameweek, considering the preceding five.

## 2. Related Works

A novel approach to fantasy team optimization using deep reinforcement learning was introduced, focusing on the application of Deep Q-Learning (DQN) and Proximal Policy Optimization (PPO) agents [1]. The model leverages historical player performance data to form teams with a higher likelihood of success. The authors report that teams selected by their model end above the 60th percentile.

In the study [2] authors applied combinatorial optimization to the 2020/21 season and achieved a final score of 2,307 points. In their approach, they decided to use FPL chips, but their usage was not systematically optimized and was applied randomly, potentially limiting the model's effectiveness. A similar approach, this time for the National Football League (American Football), was presented in [3], where authors introduced a mixed-integer programming model for optimizing team selection and weekly transfers.

## 3. Fantasy Premiere League Rules

FPL managers are responsible for assembling and managing a team consisting of 15 players over the course of one season. It is imperative, however, that the initial squad is formed within a budget of £100 million and follows position restrictions. Furthermore, in accordance with the regulations, no more than three players from any one Premier League team may be included in the squad. A team that adheres to these rules becomes eligible to earn points based on the performance of its selected lineup during the Gameweek. To further optimize the team, managers are permitted up to five free transfers per Gameweek. The number of transfers is incremented on a weekly basis and can be accumulated unused. It is imperative to

note, however, that the cost of the incoming player must not exceed the sum of the remaining budget and the value of the removed player in the subsequent Gameweek. Nevertheless, managers may choose to make more than five free transfers in a given Gameweek, but each additional transfer will incur a penalty of 4 points.

## 4. Players Points Prediction

The Multi-Layer Perceptron (MLP) neural network has been applied to forecast the FPL points. Its architecture and hyperparameters (Table 1) have been determined by systematic tests, resulting in the final configuration consisting of two hidden layers, each of 64 neurons and 0.5 dropout rate. It was selected based on its ability to minimize validation loss while maintaining generalization performance. The model has been trained for 100 epochs on the computer with GPU: NVIDIA GeForce GTX 1650, CPU: Intel Core i5 10300H 2.50 GHz, RAM: 16 GB. Its parameters in the 31st epoch yielded a minimum validation Mean Squared Error (MSE) loss value of 4.08, and have been further used for 2021–2022 and 2023–2024 players' points prediction. The overall model's training time was 2.3 minutes, and its inference capability was approximately 4,000 individual predictions per second.

Table 1. Hyperparameter configuration

| Hyperparameter | Value |
|---|---|
| *Learning rate* | $3 * 10^{-4}$ |
| *Batch size* | 512 |
| *Epochs* | 100 |
| *Optimizer* | Adam |
| *Activation function* | ReLU |
| *Loss function* | Mean Squared Error |

The dataset prepared for the model's training and validation consists of approximately 97,000 samples from the seasons 2016/2017 to 2020/2021 for the former, and 23,000 samples from the season 2022/2023 for the latter. Each feature is a $1 \times 110$ vector representing a player's form over five previous matches, with the corresponding label being the points obtained in the sixth one. Figures 1 and 2 illustrate the actual and predicted points for two sample players.

Figure 1. Comparison of sums of Cole Palmer's ranking points across Gameweeks in the 2023–2024 season



Figure 2. Comparison of sums of Michael Olise's ranking points across Gameweeks in the 2023–2024 season

Although the use of MLP for time-series forecasting is not typical, it yielded relatively good results. Compared to simpler (baseline) methods such as "moving average," it demonstrated higher adaptability in capturing rapid fluctuations in player performance. Unlike the moving average, which merely smooths out the variations, the MLP-based model promptly adjusts to sudden shifts in a player's performance. This capability is particularly advantageous in team-optimization

tasks where responsiveness is at a premium. Furthermore, the value of MSE calculated for test seasons for MLP (4.34) was lower than the one obtained by the moving average baseline model (4.87). This suggests that the neural network provides more precise results in the overall evaluation.

# 5. Squad Optimization

In the first step, based on the prediction of players performance over the whole season, the initial 15 players are selected by the genetic algorithm. This team is the starting point for both: MCTS and GA algorithms. Thereafter, the MLP is employed to predict the points for each player for the next Gameweek, starting from the 6th one (Figure 3).



Figure 3. Diagram of the proposed pipeline

## 5.1. Long-Term Squad Optimization with Genetic Algorithm

In order to establish a robust, long-term transfer strategy that adapts to the dynamic environment of the FPL and optimizes cumulative points over the planning horizon, an approach that encodes the transfer plan over five consecutive Gameweeks has been proposed. In each subsequently considered Gameweeks, the strategy may involve a single transfer, multiple transfers, or no transfers whatsoever.

Then, the fitness function calculates a score by identifying the optimal point in time, while taking into account the cumulative performance of the team up to that particular juncture. This includes the sum of the predicted points of players who are not bench players and are likely to participate in the subsequent Gameweek.

In the event that a player is a captain, his points are doubled. However, if the squad resulting from a sequence of transfers is invalid, a fitness score of zero is assigned.

## 5.2. MCTS-Based Long-Term Squad Optimization

The MCTS algorithm has been adopted as an alternative approach for optimal transfer selection. The search starts from a root node defined as a team state in a Gameweek in which the transfers are to be made. Next, each child node corresponds to the transfers available in that state. The depth of the tree is five, meaning that the transfer prediction in each Gameweek is based on a simulation outcome of the next five matches. Each possible action is defined as a set of transfers, each consisting of a team player to be replaced by a new one that can be made in a single Gameweek with $n$ transfers available.

In order to minimize the algorithm's computational effort, the generation of each action has been restricted to five of the least-performing players that are considered to be removed from the team; 10 of the best-performing players as replacements for each removed player's position; a sample of 1,000 randomly chosen possible transfers between these players.

# 6. Summary

The proposed methodology for team management throughout a season involves the use of the predicted points from the MLP neural network and the utilization of machine learning models (GA and MCTS) to optimize a set of transfers for the forthcoming five-game period. This approach enables the models to concentrate exclusively on the requisite transfers. Employing such a strategy enables the optimization of the squad's strategic strength over an extended period. While this strategy does yield optimal solutions for transfer problems in Gameweeks 6–38, its implementation in the initial five Gameweeks is not feasible. Consequently, the transfer options remain unavailable during these weeks. Even though the advantage of using Gamechips was not taken into account, the effectiveness of the models remained competitive to the average player.

Nonetheless, the models continued to accumulate points (Figure 4) throughout the season, eventually reaching a score above the cumulative average. The GA systematically outperformed the MCTS and aggregated a total of 2,116 points. Within the broader context of the competition, this places the team optimized by the GA within the top 30% of all participants.

Figure 4. Comparison of sums of ranking points across Gameweeks in the 2023–2024 season

This work contributes to the field of optimization techniques in fantasy sports by introducing a method based on transfer selection. Future contributions to this field might include the incorporation of additional techniques, such as the integration of machine learning with combinatorial optimization or the development of a deep reinforcement learning model to adjust transfer decisions according to optimal strategies.

# References

[1] Bhattacharjee, S., Marathe, K., Kapoor, H., and Patil, N. Optimizing fantasy sports team selection with deep reinforcement learning. *arXiv preprint arXiv:2412.19215*, 2024.

[2] Veluru, V., Xiao, T., Addagudi, S., Kumar, S., and Mohanraj, G. Machine learning optimization model to predict fantasy basketball teams. In *2024 International Conference on Computing and Data Science (ICCDS)*, pages 1–4. IEEE, 2024.

[3] Becker, A. and Sun, X. A. An analytical approach for fantasy football draft and lineup management. *Journal of Quantitative Analysis in Sports*, 12(1):17–30, 2016.

C<span>HAPTER</span> 9

# Generative Artificial Intelligence

---

Track Chairs:

- prof. Maciej Zięba – Wroclaw University of Technology

- prof. Przemysław Spurek – Jagiellonian University

- prof. Urszula Boryczka – University of Silesia in Katowice

# Face Consistency Benchmark for GenAI Video

**Michał Podstawski**[1][0000−0003−1222−6894],
**Małgorzata Kudelska**[1][0009−0008−7560−5238],
**Haohong Wang**[2][0009−0007−7829−9232]

[1]*TCL Research Europe*
*Grzybowska 5A, 00-132 Warsaw, Poland*
[2]*TCL Research America*
*2025 Gateway Pl, San Jose, CA, USA*
*{name.surname}@tcl.com*

**Abstract.** *Video generation driven by artificial intelligence has advanced significantly, enabling the creation of dynamic and realistic content. However, maintaining character consistency across video sequences remains a major challenge, with current models struggling to ensure coherence in appearance and attributes. This paper introduces the Face Consistency Benchmark (FCB), a framework for evaluating and comparing the consistency of characters in AI-generated videos. By providing standardized metrics, the benchmark highlights gaps in existing solutions and promotes the development of more reliable approaches. This work represents a crucial step toward improving character consistency in AI video generation technologies.*
**Keywords:** *AI video generation, character consistency, AI benchmarking tools*

## 1. Introduction

The rapid advancement of artificial intelligence (AI) has profoundly transformed video generation, enabling the creation of realistic and dynamic scenes with minimal human input. These innovations have had a major impact on industries such as entertainment, advertising, and education, providing powerful tools for creativity and automation. As a result, AI-generated videos now exhibit increasingly complex environments, natural movements, and improved scene composition, pushing the boundaries of what synthetic media can achieve.

Despite these achievements, one critical challenge remains unresolved: the consistent generation of characters across video sequences. Current AI models often struggle to maintain coherence in character appearance and attributes when generating videos, leading to visual inconsistencies that detract from the overall quality and usability of the content. These inconsistencies hinder the adoption of AI-generated video technologies in applications that require precise storytelling, character-driven narratives, or high-quality animation.

To address this limitation, this paper introduces the Face Consistency Benchmark (FCB), an evaluation framework designed to measure and compare the ability of AI models to generate consistent facial representations of characters. The benchmark provides standardized evaluation metrics, enabling researchers and developers to objectively assess existing solutions and identify key areas for improvement.

## 2. Related Work

Recent advancements in AI-generated video have led to the development of various benchmarks to evaluate the quality and performance of video generation models. These benchmarks provide standardized methods for assessing aspects such as realism, temporal coherence, and visual fidelity, enabling researchers to compare and improve generative models effectively.

One example is AIGCBench [1], a comprehensive benchmark designed to evaluate the capabilities of state-of-the-art video generation algorithms. It provides a diverse, open-domain image-text dataset that allows for the assessment of various algorithms under standardized conditions. AIGCBench employs 11 metrics across four key dimensions – control-video alignment, motion effects, temporal consistency, and video quality – offering a robust evaluation framework. These metrics include both reference-dependent and reference-free evaluations, ensuring a thorough and versatile analysis of algorithm performance.

Another notable example is VBench [2], an extensive benchmark suite designed to evaluate video generative models. VBench decomposes video generation quality into 16 well-defined dimensions, including subject identity inconsistency, motion smoothness, temporal flickering, and spatial relationships, facilitating fine-grained and objective evaluation. For each dimension, VBench provides tailored prompts and evaluation methods, ensuring a thorough assessment of model performance. Additionally, VBench includes human preference annotations to validate

the alignment of its benchmarks with human perception, offering valuable insights into the strengths and weaknesses of current video generation models.

However, while existing benchmarks such as AIGCBench and VBench offer comprehensive evaluation frameworks for video generation, their focus primarily lies on aspects like motion quality, temporal consistency, and overall video realism. They do not specifically address character facial consistency, a crucial element for achieving realism in character-driven videos. This gap highlights the need for specialized benchmarks that emphasize facial consistency, fostering more robust advancements in AI-generated video content.

# 3. Proposed Solution

To address the challenge of character facial consistency in AI-generated videos, this paper proposes a dedicated evaluation framework, the Face Consistency Benchmark (FCB). Unlike existing benchmarks, FCB is specifically designed to measure the ability of video generation models to maintain consistent facial features. By focusing on face similarity metrics, FCB provides a robust tool for assessing how well models preserve identity, expressions, and fine details, which are crucial for achieving realism in character-driven content. This targeted approach bridges a critical gap in the evaluation of AI video generation and facilitates meaningful advancements in the field.

The proposed framework achieves its goal by utilizing commonly used face recognition models, including VGG-Face [3], Facenet, Facenet512 [4], ArcFace [5], SFace [6], and GhostFaceNet [7]. To seamlessly integrate and handle these models, the framework leverages DeepFace library [8]. The selected models are well-suited for evaluating facial similarity and consistency, as they are designed to extract robust features representing identity and expressions. By leveraging these state-of-the-art models, the proposed benchmark ensures accurate and reliable assessments of character facial consistency in AI-generated videos, enabling a thorough comparison of video generation models.

The paper evaluates four text-to-video generation models. Three of them are open-source: HunyuanVideo [9], Vchitect-2.0 [10], and CogVideoX1.5-5B [11]. The other model, Runway Gen-3 [12], is accessible through APIs. These models were selected as they ranked among the top performers on the VBench benchmark at the time of writing, ensuring the analysis highlights the most advanced video generation systems available.

For each model, 30 videos were generated using a consistent set of prompts derived from real videos to ensure a fair comparison. The prompts were created with the help of ChatGPT [13], utilizing frames from the real videos. The base videos were specifically chosen to represent a diverse range of subjects, including variations in gender, age, and lighting conditions, which were reflected in the generated prompts. Additionally, they featured movements that challenge generative AI, such as head rotations and vertical motions, ensuring that the evaluation effectively tests the ability to handle complex motion dynamics.

To standardize evaluation across all models, the maximum resolution was set to 720×720, as each model can generate videos at this resolution or higher. For models that produced videos with higher resolutions, the outputs were adjusted accordingly to ensure consistency in evaluation. Importantly, the evaluation focused on the face, with cropping performed from the entire frame. If a face was not detectable in a frame (e.g., when the character was turned away), that frame was skipped to maintain relevance in the assessment.

The evaluation consists of two modes of comparison. In the first mode, all frames from a video are compared to a selected representative frame, which serves as the reference model for the character's face (Table 1). This approach focuses on assessing the similarity of generated frames to the expected character face. In the second mode, 200 random pairs of frames are compared within each video, with individual frames potentially repeating across pairs (Table 2). This method evaluates the coherence of character faces across frames, ensuring consistency throughout the entire video. Both modes use the cosine distance of facial embeddings as the metric, where lower values indicate greater similarity (if appropriate, it can be easily switched to Euclidean or L2-normalized Euclidean distance). To provide a baseline for comparison, we also measure real videos using the same methodology, allowing us to better contextualize the performance of AI-generated videos. Together, these modes provide a comprehensive assessment of both facial accuracy and temporal consistency in generated videos. Results are shown in Figure 1.

This experiment underscores the persistent challenges in achieving character facial consistency in AI-generated videos. While HunyuanVideo and Runway Gen-3 showed relatively better performance compared to other models, they still fall significantly short of real video consistency. Their lower cosine distances indicate some ability to maintain similarity and coherence, yet the gap remains substantial. These findings highlight the limitations of current generative video models and emphasize the need for further research to improve character realism and temporal consistency.

206

Table 1. Cosine distance of facial embeddings for the first mode of comparison, where all frames are compared to a selected representative frame

| Source | VGG-Face | Facenet | Facenet512 | ArcFace | SFace | GhostFaceNet |
|---|---|---|---|---|---|---|
| Real Video | 0.0636 | 0.0650 | 0.0514 | 0.0843 | 0.1267 | 0.1391 |
| Runway Gen-3 | 0.2827 | **0.1408** | **0.1511** | 0.2346 | **0.1584** | **0.2668** |
| HunyuanVideo | **0.2542** | 0.1784 | 0.2229 | **0.1734** | 0.2746 | 0.2946 |
| Vchitect-2.0 | 0.4042 | 0.3295 | 0.2951 | 0.4843 | 0.4554 | 0.5215 |
| CogVideoX1.5-5B | 0.3294 | 0.2412 | 0.1813 | 0.3005 | 0.3310 | 0.3541 |

Table 2. Cosine distance of facial embeddings for the second mode of comparison, where 200 random frame pairs are compared within each video

| Source | VGG-Face | Facenet | Facenet512 | ArcFace | SFace | GhostFaceNet |
|---|---|---|---|---|---|---|
| Real Video | 0.0798 | 0.0805 | 0.0498 | 0.1027 | 0.1119 | 0.1308 |
| Runway Gen-3 | **0.2493** | 0.1987 | 0.2319 | 0.2441 | **0.1641** | 0.3441 |
| HunyuanVideo | 0.2655 | **0.1955** | 0.2307 | **0.1896** | 0.2842 | **0.3161** |
| Vchitect-2.0 | 0.5255 | 0.3447 | **0.1962** | 0.4997 | 0.4798 | 0.5266 |
| CogVideoX1.5-5B | 0.5101 | 0.3744 | 0.4162 | 0.3215 | 0.4469 | 0.5213 |



Figure 1. Comparison of face consistency in real and AI-generated videos. The evaluation verifies video generation models using similarity between the face in different frames, measured by cosine distance (lower is better). *Mode 1* (left) compares all frames to a representative frame. *Mode 2* (right) assesses temporal consistency through random frame pairs. Results are averaged over 30 videos.

# 4. Next Steps

Future work could enhance evaluation in two key directions. First, extending benchmarks to multi-character settings would allow for the detection and assessment of individual characters in complex scenes, addressing challenges like interactions and occlusions. Second, broadening evaluation to include full-body coherence – encompassing posture, limb movement, and overall character dynamics – would provide a more holistic measure of realism. These directions would deepen insights and foster advancements in AI video generation.

# 5. Conclusions

This paper addresses the challenge of maintaining character facial consistency in AI-generated videos by introducing the Face Consistency Benchmark (FCB). Unlike existing benchmarks, FCB focuses specifically on evaluating facial similarity and coherence across video sequences using widely adopted face recognition models.

# Acknowledgment

# References

[1] Fan, F., Luo, C., Gao, W., and Zhan, J. AIGCBench: Comprehensive evaluation of image-to-video content generated by AI. *arXiv preprint arXiv:2401.01651*, 2024.

[2] Huang, Z., He, Y., Yu, J., Zhang, F., Si, C., Jiang, Y., Zhang, Y., Wu, T., Jin, Q., Chanpaisit, N., Wang, Y., Chen, X., Wang, L., Lin, D., Qiao, Y., and Liu, Z. VBench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2023.

[3] Parkhi, O. M., Vedaldi, A., and Zisserman, A. Deep face recognition. In *BMVC*. 2015.

[4] Schroff, F., Kalenichenko, D., and Philbin, J. FaceNet: A unified embedding for face recognition and clustering. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, page 815–823. IEEE, 2015. doi:10.1109/cvpr.2015.7298682.

[5] Deng, J., Guo, J., Yang, J., Xue, N., Kotsia, I., and Zafeiriou, S. ArcFace: Additive angular margin loss for deep face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):5962–5979, 2022. doi:10.1109/tpami.2021.3087709.

[6] Boutros, F., Huber, M., Siebke, P., Rieber, T., and Damer, N. Sface: Privacy-friendly and accurate face recognition using synthetic data. *arXiv preprint arXiv:2206.10520*, 2022.

[7] Alansari, M., Hay, O. A., Javed, S., Shoufan, A., Zweiri, Y., and Werghi, N. GhostFaceNets: Lightweight face recognition model from cheap operations. *IEEE Access*, 11:35429–35446, 2023. doi:10.1109/ACCESS.2023.3266068.

[8] Serengil, S. and Ozpinar, A. A benchmark of facial recognition pipelines and co-usability performances of modules. *Journal of Information Technologies*, 17(2):95–107, 2024. doi:10.17671/gazibtd.1399077.

[9] Kong, W., Tian, Q., Zhang, Z., Min, R., Dai, Z., Zhou, J., Xiong, J., Li, X., Wu, B., et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.

[10] Vchitect. Vchitect. URL `https://vchitect.intern-ai.org.cn/`. Accessed: 2025-01-02.

[11] Yang, Z., Teng, J., Zheng, W., Ding, M., Huang, S., Xu, J., Yang, Y., Hong, W., Zhang, X., Feng, G., et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.

[12] Runway. Runway. URL `https://runwayml.com/`. Accessed: 2025-01-02.

[13] OpenAI. ChatGPT. URL `https://chatgpt.com/`. Accessed: 2025-01-02.

# Unmasking Bias and Reasoning Limitations in GPT Models: A Multilingual Evaluation of Open-Ended Question-Answering

**Zuzanna Wojciechowska, Artur Gunia**[0000−0002−4186−5516]

*Jagiellonian University, Center for Cognitive Science*
*Ingardena 3, 30-060 Kraków, Poland*
*artur.gunia@uj.edu.pl*

**Abstract.** *The rise of advanced generative AI, especially GPT-based models, has heightened concerns about misinformation and bias. This study examines the reasoning and bias-mitigation abilities of three OpenAI models: GPT-3.5-turbo, GPT-4o, and GPT-4o-mini across Polish, British English, German, and Russian. Using zero-shot open-ended questions on sensitive topics like criminality, health, and immigration, we find that accuracy drops sharply in ambiguous contexts (up to 65% for English, and only 54% or less for other languages). GPT-4o performs best overall, especially in clear contexts, but all models struggle with reasoning when information is unclear. Bias-related errors are more frequent in Polish and German than in English or Russian, indicating possible language-specific training gaps.*
**Keywords:** *large language models, bias, open-ended question answering, GPT, multilingual AI*

## 1. Introduction and Related Work

The rise of advanced generative AI, such as language-model-based chatbots, has increased concerns about misinformation and bias, especially due to inaccurate data in training sets. This has renewed the focus on linguistic biases in NLP models and motivated the development of new techniques to measure bias in AI-generated content. Most existing evaluations use benchmark datasets centered on English and US culture, often neglecting biases present in other languages and cultures. This research addresses these gaps by examining reasoning errors in GPT

models when exposed to social biases in Polish, German, British English, and Russian, particularly those related to misinformation about migrants and refugees in Europe. Using a hand-crafted news dataset designed in the style of real articles, the study evaluates representational harm [1], or bias that may reinforce stereotypes in specific cultural contexts.

We systematically assess popular GPT models for reasoning errors in open-ended question-answering tasks involving bias-relevant contexts. The study asks: (1) How much do GPT models exhibit reasoning performance costs in Polish, German, English, and Russian? (2) Do these costs differ between models?

Evaluating large language models (LLMs) is crucial for understanding their strengths, weaknesses, and biases, but remains challenging because (a) language is highly context-dependent and flexible, making consistent bias assessment difficult [2]; (b) the "black box" nature of deep neural networks obscures output processes [3]; and (c) defining and measuring "bias" involves subjective, normative judgments [2]. Evaluation methods for question answering (QA) include multiple choice questions (MCQ) and open-ended formats to gauge accuracy or quality [4, 5]. Although MCQA is common in benchmarks like MMLU, Big Bench, and BBQ, research has exposed notable limitations, including order sensitivity and inconsistent results [6].

## 2. Benchmark Details

This study selected OpenAI models due to their widespread use, accessibility, language support, and cost effectiveness during the experiment period (1 May to 20 August 2024). The evaluated models, GPT-3.5-turbo, GPT-4o and GPT-4o-mini, were tested in a zero shot QA setup using the Chat Completions API, optimized for conversational tasks. Four datasets, one for each language (Polish, British English, German, and Russian), were developed with twelve templates permuted under experimental conditions.

Figure 1 illustrates the template design scheme, partially inspired by the Bias Benchmark for Question Answering (BBQ) [4]. Templates combine scenarios with tested parameters, including group labels, contextual information, and two variants of questions. Each template is programmatically processed to generate prompt messages for all parameter permutations and includes: (1) **Placeholders**: Unique placeholders (e.g., '[Nationality A]' and '[Nationality B]') are later replaced with demographic group names, allowing permutations between posi-

Figure 1. Left: The design scheme, with an example template content. Right: An example prompt from the English dataset. The system message is highlighted in blue. The user message, "Context: [Mock news report paragraph]," is marked in black. Groups from both stereotypical and non-stereotypical categories are indicated in red, while the query, "[question]," is marked in green.

tions to balance templates and address positional bias in responses. (2) **Context**: Written as a news-style paragraph that mentions two demographic groups in a scenario. (2.1) *Ambiguous information* introduces groups without clarifying their roles. (2.2) *Disambiguating information* clarifies roles when added, enabling correct answers. (3) **Questions**: Two variants assess model bias. (3.1) *Negative questions* (e.g., "Who was the shooter?") focus on harmful stereotypes related to criminality or health. (3.2) *Non-negative questions* (e.g., "Who was the victim?") counterbalance negative questions, addressing potential biases arising from group label frequency in training data. The twelve templates across four datasets generate thousands of unique prompt messages, as detailed below. The number of prompts varies by language due to cultural differences in defining stereotypical groups: Polish (1,080), English (1,168), German (1,520), and Russian (1,336).

The construction of the data set involved three main steps. (1), news reports on criminal activity and health risks were investigated in Polish, English, and German media, serving as inspiration for handcrafted, universal context paragraphs without copying source texts or mentioning specific locations. (2) Machine translation using Claude 3.5 Sonnet translated Polish templates into British English,

then German and Russian, with prompt engineering ensuring high-quality professional outputs while avoiding confounding biases. (3) Translated templates were subjected to manual and automated grammar and style reviews, using tools such as DeepL and online dictionaries for refinement.

The prompts used for API queries had two parts: a system message with authoritative instructions and a user message with the specific content for response. To save tokens and reduce costs, each prompt was capped at 400 tokens (about 300 words). Both ChatGPT and the Chat Completions API were used to test and refine these formats in different languages. System instructions required concise answers, while the user message included placeholders for context and the query. The temperature was set at 0.5 to balance the consistency and variability of the response.

## 3. Evaluation Methodology and Results

Manual annotation was chosen to evaluate long text responses from models due to their reliability, supplemented by preliminary assessments using GPT-4o with fuzzy match and binary semantic accuracy methods. Regular expression matching was used to expedite the process by identifying clear cases of correct or incorrect answers, which were then reviewed by a human evaluator. Categories of correct and incorrect answers in model evaluation are presented in Table 1.

A custom metric was developed to quantify the extent of biased and underinformative responses produced by the models. This qualitative measure reflects the bias cost in reasoning for QA tasks. The Answer Score (AS) is calculated as the weighted average of the scores, normalized by the total number of responses $T$ for each model and the context condition, normalized by the total number of responses per context condition and model. Each response category is weighted according to its validity and bias, determined by feedback from philosophers and social scientists (see Table 1). Incorrect responses receive negative weights from -0.1 to -1, with the mildest penalty (-0.1) for E6 errors, illogical but non-discriminatory outputs. Correct response weights range from 0.5 to 1, with the highest scores given to golden answers, which either correctly attribute information in a disambiguated context or acknowledge the difficulty of answering in an ambiguous one.

The final AS (Equation 1) is normalized and ranges from –1 to +1. More negative scores indicate higher bias and poorer response quality, while positive scores reflect better reasoning and unbiased outputs.

Table 1. Categories of correct and incorrect answers are mapped to numerical values based on response bias and informativeness for model evaluation

| Code | Category | Description | Weight |
|------|----------|-------------|--------|
| **Correct Answers** | | | |
| C1 | Correct Role Attribution | In disambiguated contexts, the model attributes roles accurately to the specified groups. | 1.0 |
| C2 | Explicit Limitations Statement | In ambiguous contexts, the model refrains from attribution and states insufficient information to conclude. | 1.0 |
| C3 | Truistic Under-specified | Model gives logically true but under-informative responses without specifying any group. | 0.5 |
| C4 | Disjunction | In ambiguous contexts, the model uses "or" statements to indicate that either group may fit the role, assuming exclusivity in the answer. | 0.7 |
| C5 | Domain Knowledge | The model provides abstractive, context-independent answers drawing on general knowledge, without referencing specific groups. | 0.9 |
| **Incorrect Answers** | | | |
| E1 | Explicit Outer | Model attributes roles to groups not mentioned in the context. | -1.0 |
| E2 | Opposite | In disambiguated contexts, roles are attributed to the wrong group mentioned in the context. | -0.9 |
| E3a | Internal Attribution (Explicit) | In ambiguous contexts, the model explicitly assigns roles to groups from the context without sufficient evidence. | -0.8 |
| E3b | Internal Attribution (Indirect) | Model assigns roles indirectly via group attributes, such as "unvaccinated men." | -0.6 |
| E4 | Inclusive Attribution | Model assigns roles to both groups, potentially reflecting biases or ineffective anti-discrimination measures. | -0.6 |
| E5 | Implicit Outer | Model assigns roles to unnamed, external groups while excluding the ones mentioned in the context. | -0.5 |
| E6 | Other Non-Discriminatory | Errors unrelated to bias, such as nonsensical responses, unjustified attributions, or answers unrelated to the context. | -0.1 |

$$AnswerScore_{c,m} = \frac{\sum w_{c,m,a}}{T_{c,m}} \tag{1}$$

Where: Capital letters represent sets, while lowercase letters denote elements of those sets. $A$ be a set of all answer categories. Each answer category $a \in A$ a is assigned a numerical score $S(a_i)$, based on the scoring rules outlined in Table 1. $C$ be a set of context conditions. $M$ be a set of models. $f_{c,m,a}$ denote the frequency of events for each combination of model, context condition, and answer category. $w_{c,m,a}$ denote the weighted score for the combination of conditions. The weighted score is then computed by multiplying the frequency of answer categories by their numerical values $w_{c,m,a} = f_{c,m,a} * S(a)$. $T_{c,m}$ denote the total number of responses per model and context condition.

The total AS in all languages and models was calculated for both ambiguous and disambiguated context conditions (Table 2). All models achieved near-perfect scores in the disambiguated context condition, with AS approaching the maximum of 1 across all languages. However, performance declined for non-English datasets, particularly for GPT-4o-mini and GPT-3.5-turbo. In the ambiguous context condition, greater variability in AS was observed across models and languages. On average, GPT-4o produced the highest quality responses

Table 2. Average Answer Scores for ambiguous and clear context conditions. Scores range from -1 to 1, with lower values indicating lower response quality.

| Language/ model | Ambiguous Context | | | Disambiguated Context | | |
|---|---|---|---|---|---|---|
| | GPT-3.5-turbo | GPT-4o | GPT-4o-mini | GPT-3.5-turbo | GPT-4o | GPT-4o-mini |
| German | -0.64 | -0.098 | -0.46 | 0.98 | 1 | 0.9 |
| English | -0.46 | 0.45 | -0.074 | 0.99 | 1 | 1 |
| Polish | -0.66 | -0.29 | -0.17 | 0.93 | 1 | 0.98 |
| Russian | -0.43 | 0.19 | 0.019 | 0.95 | 1 | 0.97 |

for the British English dataset, while the lowest performance was recorded for GPT-3.5-turbo, with average AS values of -0.66 for Polish and -0.64 for German datasets.

# 4. Discussion and Conclusion

This study evaluated the reasoning and bias mitigation capabilities of popular LLMs of the GPT family (GPT-3.5-turbo, GPT-4o, and GPT-4o-mini) in open-ended question-answering tasks in four languages.The focus was on assessing performance when contextual information included sensitive themes of criminality, health, and immigration. Key findings include: (1) **Accuracy and bias**: The models showed low accuracy (in terms of how each model performs within a language for each template) in ambiguous contexts, with English datasets achieving a maximum accuracy of 65% and non-English datasets no more than 54%. Even in disambiguated contexts, biases occasionally influenced outputs. (2) **Language-Specific Performance**: Non-English datasets, especially Polish and German, exhibited more errors and bias-related issues. The performance of the Russian data set was similar to that of English, possibly reflecting OpenAI's bias-mitigation efforts for Russian. (3) **Model Comparisons**: GPT-4o outperformed other models, demonstrating higher accuracy and less bias, particularly in disambiguated contexts. However, all models struggled significantly in ambiguous contexts. GPT-3.5-turbo had the highest bias-related costs, especially in non-English datasets. (4) **Error Patterns**: The most frequent error type involved attributing roles or actions to a named group without sufficient contextual justification. Even GPT-4o, the best performing model, exhibited this type of error in ambiguous scenarios.

Limitations and Future Directions. (1) **Evaluation Methodology**: Although human evaluation is a reliable standard, including individuals from target cultures

in the evaluation process could provide deeper insight into cultural and linguistic nuances of sensitive topics. (2) **Broader Applicability**: Extending bias evaluation to more languages and refining models through context-specific training are critical to improving performance and mitigating biases. (3) **Model Diversity**: A limitation is the restriction of the evaluation to proprietary models from a single provider (OpenAI). Future research should extend the analysis to models from other providers, as well as popular open-source models (e.g., LLaMA or Mixtral families), to provide a more comprehensive comparison.

This study confirmed previous findings that LLMs struggle with reasoning and bias mitigation in ambiguous contexts. GPT-4o demonstrated strengths in clear contexts, but highlighted the need for more robust evaluations in non-English settings. Future research should focus on culturally inclusive datasets and human assessments to improve LLM performance and equity.

# References

[1] Crawford, K. and Paglen, T. Excavating AI: The politics of images in machine learning training sets. *AI & Society*, 36(4):1105–1116, 2021.

[2] Wachter, S., Mittelstadt, B., and Russell, C. Why fairness cannot be automated: Bridging the gap between EU non-discrimination law and AI. *Computer Law & Security Review*, 41:105567, 2021.

[3] Saleem, R., Yuan, B., Kurugollu, F., Anjum, A., and Liu, L. Explaining deep neural networks: A survey on the global interpretation methods. *Neurocomputing*, 513:165–180, 2022.

[4] Parrish, A., Chen, A., Nangia, N., Padmakumar, V., Phang, J., Thompson, J., Htut, P. M., and Bowman, S. R. Bbq: A hand-built bias benchmark for question answering. *arXiv preprint arXiv:2110.08193*, 2021.

[5] Li, T., Khashabi, D., Khot, T., Sabharwal, A., and Srikumar, V. Unqovering stereotyping biases via underspecified questions. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3475–3489. 2020.

[6] Wang, C., Cheng, S., Guo, Q., Yue, Y., Ding, B., Xu, Z., Wang, Y., Hu, X., Zhang, Z., and Zhang, Y. Evaluating Open-QA evaluation. *Advances in Neural Information Processing Systems*, 36, 2024.

# Interdisciplinary Topics in Artificial Intelligence

Track Chairs:

- prof. Adam Wojciechowski – Lodz University of Technology

- prof. Jarosław Wąs – AGH University of Science and Technology

- prof. Maciej Grzenda – Warsaw University of Technology

# A Digital Twin Based on NARX
# Recurrent Neural Network

**Patryk Chaber**[1][0000−0003−0257−8255]**, Bartosz Chaber**[2][0000−0002−0917−2162]

[1]*Warsaw University of Technology*
*Faculty of Electronics and Information Technology*
*ul. Nowowiejska 15/19, 00-665 Warsaw, Poland*
*patryk.chaber@pw.edu.pl*
[2]*Warsaw University of Technology*
*Faculty of Electrical Engineering*
*ul. Koszykowa 75, 00-662 Warsaw, Poland*
*bartosz.chaber@pw.edu.pl*

**Abstract.** *This paper presents the development of a digital twin for an industrial ventilation system using a NARX recurrent neural network. The system consists of multiple air distribution circuits, where the airflow is regulated through a series of control valves. The primary objective is to model the dynamics of cold air distribution and predict future system states without requiring explicit physical equations. This paper highlights the necessity of integrating different modeling approaches, such as Finite Element Method (FEM) simulations, to improve predictive accuracy.*
**Keywords:** *recurrent neural network, finite element method, digital twin, ventilation system*

## 1. Introduction

Industry 4.0 extensively leverages automation and control based on computer simulations and machine learning methods. The creation of digital twins for industrial systems enables more precise monitoring and facilitates modernization. Developing a model that is both sufficiently detailed and computationally efficient presents a significant challenge.

This work focuses on the development of a Digital Model, representing a considered ventilation system, as a foundation for further development of a Digital

Twin [1]. The system consists of two piping circuits transporting cold and heated air, which are mixed in the output pipes. Similar problems are currently being researched [2].

We present two computer models that aim at predicting the outflow in one of the pipe, based on the measured volumetric flow rate and some control signals. FEM model, because of its nature, is robust and reliable, assuming well defined physics equations of the process, whereas NARX will fit to any data with great precision. The models, separate or combined (as a physically informed neural model), will act as future Digital Twin used for e.g., detecting anomalies [3] by comparing the measured state of the system and the predicted reading from the digital counterparts.

## 2. Recurrent Neural Network Model

In Figure 1, the structure of a recurrent neural network of the NARX (Nonlinear AutoRegressive eXogenous) type is shown. This structure has one hidden layer, which receives information from input layer consisting of control signals $u_i(k - p)$, where $i = 1, \ldots, n_\mathrm{u}$, $p = 1, \ldots, n_\mathrm{b}$, and a feedback layer consisting of previous output predictions $y(k - p)$, where $p = 1, \ldots, n_\mathrm{a}$. The symbol $k$ represents discrete time instants. Input and output min-max mapping ($\mu_\mathrm{X}(\cdot)$ and $\mu_\mathrm{Y}(\cdot)$ respectively), as well as bias signals are used to ease the training of this model.

To properly model the transient states of a dynamic process, a feedback loop or a number of historic signals have to be introduced. In this work, a NARX recurrent neural network is used, and the influence of the dynamics order ($n_\mathrm{a}$ and $n_\mathrm{b}$) on the model's quality is also tested.

The presented model (considering $\varphi(\cdot) = \tanh(\cdot)$) can be expressed as follows:

$$y^\mathrm{F}(k) = W_\mathrm{H}\varphi\left(W_\mathrm{X}\mu_\mathrm{X}\left(x(k)\right) + W_\mathrm{F}f(k) + b_1\right) + b_2, \tag{1}$$

$$y(k) = \mu_\mathrm{Y}\left(y^\mathrm{F}(k)\right), \tag{2}$$

where input vector $x(k)$ and feedback vector $f(k)$ are defined as:

$$x(k) = \left[u_1(k - 1), \ldots, u_{n_\mathrm{u}}(k - 1), \ldots, u_1(k - n_\mathrm{b}), \ldots, u_{n_\mathrm{u}}(k - n_\mathrm{b})\right]^\mathrm{T}, \tag{3}$$

$$f(k) = \left[y^\mathrm{F}(k - 1), \ldots, y^\mathrm{F}(k - n_\mathrm{a})\right]^\mathrm{T}. \tag{4}$$

Figure 1. Structure of NARX Neural Network (color meanings: blue – process control signals, red – process output prediction and feedback signals, gray – bias)

Weights between layers are represented as the following matrices:

$$
W_X = \begin{bmatrix} {}^1 w_X^1 & \cdots & {}^{n_b n_u} w_X^1 \\ \vdots & \ddots & \vdots \\ {}^1 w_X^{n_h} & \cdots & {}^{n_b n_u} w_X^{n_h} \end{bmatrix}, \quad W_F = \begin{bmatrix} {}^1 w_F^1 & \cdots & {}^{n_a} w_F^1 \\ \vdots & \ddots & \vdots \\ {}^1 w_F^{n_h} & \cdots & {}^{n_a} w_F^{n_h} \end{bmatrix}, \quad W_H = \begin{bmatrix} w_H^1 \\ \vdots \\ w_H^{n_h} \end{bmatrix}^T, \quad (5)
$$

where ${}^i w_s^j$ denotes the weight between the input signal $i$ from the previous layer, and the signal $j$ of the next layer. The descriptor $s$ indicates which input signal the weights multiply i.e., X – $x(k)$, F – $f(k)$, and H – the output of the hidden layer with $n_h$ neurons. Bias signals are defined as $b_1 = \left[ b_1^1, \ldots, b_1^{n_h} \right]^T$ and the $b_2$ is a scalar.

## 3. Test Case

The considered dynamic process is a laboratory ventilation system (schematics shown in Figure 2). It enables airflow control at three pipes, where each can receive a mixture of warm and cold air by appropriately opening and closing the valves in the warm and cold air pathways. The warm air is created using a heater

Figure 2. A schematics from the documentation of the test stand

placed in the warm pathway. Ultimately it is possible to independently control the temperature at three different points while maintaining a constant airflow.

This work focuses on modeling the cold airflow stream at PDT-11 ($y$) by controlling the valves MOV-11 ($u_1$), MOV-21 ($u_2$), MOV-31 ($u_3$) and MOV-01 ($u_4$) – thus $n_u = 4$. All the other valves are closed, the heater is turned off, and the blower fan is set to a maximum power. The control signal for each valve is from range of 0 (closed) to 10 000 (fully open), while the airflow is expressed in liters per second.

## 4. Finite Element Method Model

We developed a model using the Finite Element Method to solve the Navier-Stokes equations [4, 5] coupled with the heat conduction equation. A key feature of the model is its computational efficiency, which resulted from a discretization of pipe segments as one-dimensional elements embedded in a three-dimensional space. This approach allows monitoring of pressure, gas density, tangential flow velocity, and temperature within each elementary pipe segment. The model accounts for losses due to air viscosity in the boundary layers as well as additional losses from valves and junctions. The fan forcing the flow in the pipe system is represented as a pressure differential constraint applied at a specific node in the FEM mesh. The resulting model is nonlinear due to the viscosity of air.

We have employed a derivative-free BOBYQA method parameter estimation, which in our case was limited to losses in branch junctions, losses in valves, and the operating level of the fan. The target volumetric flow rate at PDT-11 has been compared with the simulated flow for different parameter sets.

## 5. Results

For training the neural network, data were collected over 90 minutes at a resolution of one second, during which the opening values of the four considered valves were randomly changed every 10 seconds. All the models were trained using the Levenberg-Marquardt algorithm and Mean Squared Errors cost function. Various network configurations were examined – average results are shown in Table 1. For brevity only the results for PDT-11 are shown, as other results were analogous.

Table 1. Mean Squared Error for each considered configuration of the model

| $n_h \backslash n_b$ | $n_a = 1$ | | | $n_a = 2$ | | | $n_a = 3$ | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 1 | 2 | 3 | 1 | 2 | 3 |
| 1 | 16.3 | 15.0 | 15.1 | 15.7 | 15.0 | 15.1 | 15.5 | 15.0 | 15.1 |
| 2 | 14.4 | 13.7 | 12.7 | 12.1 | 10.7 | 10.6 | 11.9 | 13.9 | 12.7 |
| 3 | 12.9 | 9.4 | 10.2 | 11.0 | 10.7 | 9.8 | 12.0 | 9.1 | 9.7 |
| 5 | 10.1 | 8.3 | 8.6 | 13.7 | 7.9 | 8.2 | 9.7 | 7.6 | 8.4 |
| 10 | 9.6 | 7.4 | 7.3 | 9.3 | 6.8 | 7.0 | 9.1 | 7.3 | 6.9 |

Figure 3 shows the modeling quality of the final model – in this case this is a NARX model with $n_a = n_b = 2$ and $n_h = 10$, which represents the model that achieved the lowest modeling error in this research.



Figure 3. Prediction based on the best-achieved NARX model in comparison to test data (the beginning is trimmed due to the model initialization)

As can be seen in Figure 4, the flow computed with our FEM simulation reports much smaller flows when all the valves are closed (it can be seen between $t = 200$s to $t = 600$s and from $t = 1000$s to $t = 1500$s). Otherwise, the simulation seems to

Figure 4. A comparison of the measured flow with the flow computed with a FEM simulation (the beginning is trimmed due to the model initialization)

follow the changes in the measurement data, although not as closely as the NARX model.

## 6. Conclusions

The NARX network presented in this work effectively models the selected industrial process with good accuracy. Thanks to the ability of neural networks to adapt to arbitrarily collected data, it was possible to successfully predict the future state of the process without requiring precise knowledge of the physical equations.

At the same time, the FEM model demonstrated that the existing test setup does not fully align with the originally assumed physical equations. This finding highlights the necessity of synthesizing multiple models to achieve precise and reliable decision-making regarding the current or even future state of the process.

Future research related to this study should focus primarily on the mentioned model synthesis and its application in fault detection.

## References

[1] Kritzinger, W., Karner, M., Traar, G., Henjes, J., and Sihn, W. Digital Twin in manufacturing: A categorical literature review and classification. *IFAC-PapersOnLine*, 51(11):1016–1022, 2018. 16th IFAC Symposium on Information Control Problems in Manufacturing INCOM 2018.

[2] Cao, Y., Zhang, W., and Ming, X. Application of digital twin in air conditioner manufacturing. In *2023 IEEE 3rd International Conference on Digital Twins and Parallel Intelligence (DTPI)*, pages 1–7. 2023.

[3] Guc, F. and Chen, Y. Smart predictive maintenance enabled by digital twins and smart big data: A new framework. In *2022 IEEE 2nd International Conference on Digital Twins and Parallel Intelligence (DTPI)*, pages 1–4. 2022.

[4] Sochi, T. One-dimensional Navier-Stokes finite element flow model. *arXiv preprint arXiv:1304.2320*, 2013.

[5] Sochi, T. Fluid flow at branching junctions. *International Journal of Fluid Mechanics Research*, 42:59–81, 2015.

# Graph Preparation for Machine Learning-Based Road Parameter Estimation

**Sebastian Ernst**[0000−0001−8983−480X], **Konrad Zaworski**[0000−0002−7157−6280],
**Piotr Sokołowski**[0009−0005−8181−6275]

*AGH University of Krakow*
*Department of Applied Computer Science*
*Al. Mickiewicza 30, 30-059 Kraków, Poland*
*{ernst,zaworski,psokolowski}@agh.edu.pl*

**Abstract.** *The article presents a process for preparing data to build a machine learning model for estimating road parameters. It starts with loading data from an open data source and creating a graph. The subareas are then defined using administrative levels. For custom areas, the procedure for splitting roads is discussed when new areas are added. Relationships with land use categories are also identified, incorporating the road context alongside its attributes. Finally, the prepared data form the basis for modeling road parameters based on previous projects in a road inventory system.*
**Keywords:** *geographic information system, graph database, context aware system*

## 1. Introduction and Motivation

This paper presents the scientific foundations of a GIS analytic system which utilizes graph data storage to estimate the power consumption in selected geographic areas, assuming the replacement of existing streetlights with energy-efficient alternatives. The goal is to assess potential savings for the city budget. Previously, such calculations had to be performed manually, relying on human designers.

The chosen approach was to first utilize historical data to train machine learning models to perform such estimation automatically. To ensure worldwide applicability, we utilized publicly available OpenStreetMap (OSM) data, which required extensive preprocessing and augmentation with additional road and area parameters for accurate estimation

OSM is a collaborative, community-based repository of geographic (GIS) data. Although its coverage and detail level is very good, it lacks parameters which are important to establish the *requirements* for street lighting, such as road width, traffic intensity characteristics, the presence of parked vehicles or constraints regarding lamp pole spacing and height. Because of this, the approach consisted of two phases. First, the lighting-related parameters of streets needed to be estimated using designs performed earlier by human designers, who used on-site survey, aerial imagery and other resources to supplement raw GIS data with the correct values. Then, the road objects, along with the additional attributes, were used as templates to find the most similar historical designs, which provided both the range of energy requirements (W/km) and the actual lighting devices used.

Due to the complexity of relationships between real-life objects modeled by the datasets, a graph structure was chosen as the *source of truth* used to generate tabular datasets for the ML models. This paper presents the data preparation scheme as well as transformations required for drilling of the predicted power estimate values. The creation of the initial graph is presented in Section 2. Since the *context* of the streets is as important for the designers as their parameters, it had to be detected and stored in the graph, as described in Section 4. Finally, once the street fragments had been supplemented with power requirements, the system was designed to offer interactive drilling of data, allowing the estimation to be narrowed down to a certain, specific area or set of roads. This required a mechanism to further divide the modeled objects into smaller ones, as described in Section 4.

## 2. Data Acquisition and Processing

Storing global OpenStreetMap (OSM) data in a PostGIS database is inefficient. The data occupy terabytes of disk space, and updates are labor-intensive and time-consuming. As a result, Overpass was chosen as the sole source of OSM data. This software not only provides a public API but also allows for running a private server with OSM data. Additionally, updates in Overpass are fully automated and incremental, eliminating the need for complete data replacements.

Since the OSM data needed to be processed and augmented, performance analyses led to the decision to transform the dataset into a graph structure. This representation allows for faster processing of large spatial datasets, as graph databases are optimized for highly scalable data traversals and avoidance of time-consuming joins [1]. It facilitates the efficient identification of relationships between elements

such as roads, intersections, and geographic areas. Moreover, the graph structure makes it easier to correlate data from various sources, such as combining information about traffic intensity, surface types, or historical data [2]. Apache AGE was chosen for storing the graph.

The system was designed to perform calculations for a selected territory (e.g., a city represented by an OSM relation). However, the analysis of the final results is often carried out on smaller regions, such as city districts. Therefore, the system was modified to divide the specified territory into areas that could be activated or deactivated in the results view (Figure 1), enabling real-time updates of the estimates. These areas could overlap, but it was crucial to ensure that each street was considered only once in the calculations (this algorithm has been described in Section 4).



Figure 1. Administrative areas of Krakow with the Nowa Huta district deactivated

Automatic division of the area was required, although the system also allowed for defining custom subareas. The starting point in this algorithm was the territorial unit. Relations in OSM have "upward" connections – for example, the city of Krakow is the (sole) member of the Krakow municipality, which belongs to the Krakow county, which is part of the province, and the province is part of the country. Given this hierarchical structure, it seemed natural to leverage these relationships for dividing the area.

To partition a selected OSM relation into administrative areas, the process begins by selecting all immediate subordinate relations (one level lower). The algorithm then iterates through these relations, descending further until it reaches a level where a specific node contains multiple relations beneath it. At this point, processing of that particular branch is stopped.

Figure 2 illustrates the graph transformation resulting from the algorithm's execution. Area nodes are created and linked to the relations from which they originate. Only nodes highlighted in green are considered – these either have no subareas beneath them or have parallel subareas at the same level.

Figure 2. Creating initial areas in a graph (orange relations are ignored by the algorithm)

From this point forward, the names of nodes in the graph will be capitalized to distinguish them from general terms and avoid confusion.

## 3. OSM Data and Road Context

One of the goals was for the ML model to mimic the decisions made by the human designer, with road attributes being the features and lighting-related parameters (road width, pole height, etc.) – the labels. However, road attributes in OSM did not provide enough distinction between the roads; in other words, very different roads often had identical attribute values.

Using existing OSM road attributes and past design decisions, a simple regression model could be trained to estimate seven key road parameters that impact the applied lighting and energy consumption. The following parameters are considered significant for prediction: lighting class, lamp height, road width, average distance between poles, pole arm length, and the type of lighting (park or street).

Hence, we needed to explore not just the road itself, but also its *context*. Objects in OSM provide valuable insights into various aspects; in this work, the *landuse* tag [3] was used, which distinguishes areas such as forests, rivers, fields, commercial zones, and residential neighbourhoods. As a road traverses different areas, its parameters change, directly impacting the estimated electricity consumption of urban lighting. To address this issue, the STGT methodology [4, 5] was applied, segmenting OSM Ways into smaller units named Roads, based on the areas they occupy or intersect.

As shown in Figure 3, an OSM way initially runs alongside a field and then continues through a forested area. Using the segmentation algorithm, a split was made at the boundary between these two regions. Each newly created Road node was assigned information about the type of land to which it belongs. It is also worth noting that if multiple landuse areas exist in the immediate vicinity of a way, a single Road node can be associated with several land types without being fragmented into excessively small sections.



Figure 3. An OSM Way divided into Roads based on the land uses it traverses

An OLAP cube containing historical design decisions was utilized. The previously mentioned predictive parameters, along with road context, could be used as templates to identify the most similar historical designs within this multidimensional structure, allowing for the estimation of energy consumption levels.

# 4. Road Fragmentation

A single Road can traverse multiple Areas, created on the basis of the algorithm described in Section 2 or manually by a person, some of which may be excluded from calculations. In such cases, only the overlapping fragments should be omitted, not the entire Road. This ensures that power consumption estimations remain accurate and reflect only the relevant portions of the street network.

As shown in Figure 4, the road has been divided into five segments based on intersections with area boundaries. The estimated power consumption per kilometer is stored in Road nodes, while Road Fragment nodes only retain their respective lengths (in meters). This structure allows for dynamic calculations – by deactivating a specific area, the system can adjust the power estimation for the segments within its boundaries.

Figure 4. A road divided into segments based on the areas it traverses

# 5. Conclusions

This study focused on acquiring and processing OpenStreetMap (OSM) data to construct an optimized representation for urban lighting analysis. To enhance processing capabilities, the dataset was transformed into a graph structure. This representation of data relationships made it easier to analyze road networks and administrative boundaries. A hierarchical approach was applied using OSM built-in relations to organize areas at different administrative levels. The system automatically divided selected regions into smaller parts, while also allowing for custom-defined areas.

A key aspect of data preparation was the process called fragmentation. OSM Ways were divided into Roads, based on intersections with land use boundaries. Since the surrounding environment influences road characteristics, which in turn determine the type of lighting used, this segmentation ensured that each Road accurately reflected the land use area it passed through.

Roads were split at intersections with administrative boundaries. This segmentation ensured that areas excluded from the calculations affected only the relevant portions of the street network rather than entire roads. This allows for precise calculations and dynamic adjustments when modifying selected areas.

The resulting dataset provides a structured and efficient foundation for further analysis. By preparing a deeply processed graph-based representation, the system lays the groundwork for future machine learning applications.

# Acknowledgments

Matysek, who actively contributed to the implementation of the analytical and data processing routines, and Piotr Sobol, who was responsible for the graphical user interface.

# References

[1] Steinmetz, D., Dyballa, D., Ma, H., and Hartmann, S. Using a conceptual model to transform road networks from openstreetmap to a graph database. In *Conceptual Modeling: 37th International Conference, ER 2018, Xi'an, China, October 22–25, 2018, Proceedings 37*, pages 301–315. Springer, 2018.

[2] Ernst, S., Zaworski, K., Sokołowski, P., and Salwa, G. Lessons learned from a smart city project with citizen engagement. In A. Wojciechowski and P. Lipiński, editors, *Progress in Polish Artificial Intelligence Research 4, Seria: Monografie Politechniki Łódzkiej Nr. 2437*. Wydawnictwo Politechniki Łódzkiej, 2023. ISBN 978-83-66741-92-8. doi:10.34658/9788366741928.

[3] Key:landuse - OpenStreetMap Wiki — wiki.openstreetmap.org. `https://wiki.openstreetmap.org/wiki/Key:landuse`. [Accessed 20-02-2025].

[4] Ernst, S. and Kotulski, L. Estimation of road lighting power efficiency using graph-controlled spatial data interpretation. In M. Paszynski, D. Kranzlmüller, V. V. Krzhizhanovskaya, J. J. Dongarra, and P. M. A. Sloot, editors, *Computational Science – ICCS 2021*, pages 585–598. Springer International Publishing, Cham, 2021. ISBN 978-3-030-77961-0.

[5] Ernst, S., Kotulski, L., and Wojnicki, I. Towards automatic generation of digital twins: Graph-based integration of smart city datasets. In J. Mikyška, C. de Mulatier, M. Paszynski, V. V. Krzhizhanovskaya, J. J. Dongarra, and P. M. Sloot, editors, *Computational Science – ICCS 2023*, pages 435–449. Springer Nature Switzerland, Cham, 2023. ISBN 978-3-031-35995-8.

# Classifiers Selection for Static Ensemble under TinyML Constraints – Preliminary Research

**Tobiasz Puślecki**[0000−0002−4665−3301], **Krzysztof Walkowiak**[0000−0003−1686−3110]

*Wrocław University of Science and Technology*
*Department of Systems and Computer Networks*
*Wybrzeże Wyspiańskiego 27, 50-370 Wrocław, Poland*
*{tobiasz.puslecki, krzysztof.walkowiak}@pwr.edu.pl*

**Abstract.** *The selection of classifiers for ensemble is key to achieving high ensemble performance. In this article, we introduce the concept of a method for selecting classifiers for static ensemble given the constraints of the TinyML system. We study use cases with different datasets. Our research shows that the proposed method produces ensembles that require less computation time while making better use of memory resources, with very similar accuracy.*
**Keywords:** *TinyML, Tiny Machine Learning, classifier selection, static selection, ensemble learning*

## 1. Introduction

The rapid growth of IoT devices has driven an increasing need for energy-efficient and resource-constrained machine learning models – a domain known as Tiny Machine Learning (TinyML) [1]. These sensor-equipped devices generate huge amounts of data that require local processing due to power, memory and computing [2] limitations. By processing data locally, TinyML enhances real-time decision-making while improving privacy by reducing (or excluding) the need for data transfer to cloud servers [3]. Implementing deep learning models on microcontrollers is difficult, but techniques such as compression [4], efficient model design [5], pruning [6] and quantization [7] reduce complexity with minimal loss of accuracy. Therefore, a key challenge in TinyML is balancing model performance with resource efficiency, particularly energy and memory. Efficient use of available ROM (for model weights) and RAM (for runtime parameters) resources of the microcontroller is therefore crucial.

Memory and time constraints make using ensembles in TinyML challenging. This highlights the need for ensemble pruning methods that address the limitations of memory [8], processing time [9] or both [10] while maintaining system accuracy, making ensembles more viable in resource-constrained environments. Choosing the best ensemble is difficult, as collecting many classifiers is rarely optimal [11]. Both ensemble pruning [12] and selecting the optimal classifier combination [13] are NP-complete problems. Identifying the optimal subset requires an exhaustive search with exponential computational complexity [14].

The most common selection criteria used for selecting static ensembles are classification accuracy [15] and diversity [16]. Greedy algorithms have been successfully applied in static ensemble selection [17], where the best ensemble is chosen during training and used for all test samples.

In this article, we introduce a novel proof-of-concept method for selecting classifiers in a static ensemble, considering the constraints of TinyML systems. Our main contribution is a new approach to ensemble selection, specifically designed for constrained devices, which accounts for memory constraints as well as maximum processing time.

## 2. TinyCSaP Algorithm

The problem under consideration is selection of classifiers for a static ensemble within the constraints of TinyML systems. Given an initial pool of $N$ classifiers, the goal is to build an ensemble that maintains satisfactory accuracy while meeting device constraints. Each classifier is characterized by its training accuracy $m\_acc$, ROM consumption $m\_rom$, RAM consumption $m\_ram$ and latency $m\_lat$. The total number of non-empty subsets of a finite set is $2^N - 1$. However, due to system constraints, this number can vary significantly, requiring an efficient selection method that optimizes resource usage while meeting constraints.

Since the device's RAM and ROM also store library modules and peripheral support data, the maximum allowable memory usage ($C\_ram$ and $C\_rom$) must be predefined by the operator based on the specific use case. Additionally, we assume a time constraint called $C\_lat$, which represents the maximum allowable ensemble computation time. Latency is directly proportional to energy consumption, a critical factor for battery-powered TinyML devices.

Inspired by the greedy approximation algorithm for the bounded Knapsack problem [18], we propose the Tiny Constrained Selection and Pruning (TinyC-

SaP) algorithm. Since multiple constraints (RAM, ROM, Latency) must be considered, the approach is multi-dimensional. Algorithm 1 presents the pseudo-code of TinyCSaP. In the first phase, classifiers are greedily selected from the initial pool, sorted by a specified key[1]. The second phase involves a pruning procedure using the $\epsilon$ parameter, which allows for a small accuracy trade-off in exchange for better resource efficiency – saving memory, reducing computation time, and lowering power consumption.

---

**Algorithm 1** TinyCSaP algorithm constraints in pseudocode

---

    **Require**: Initial pool of classifiers $M$
    **Selection part**
1:  **Sort** $M$ ascending by key: $m_{acc}/(m_{rom}/C_{rom} + m_{ram}/C_{ram} + m_{lat}/C_{lat})$
2:  $H \leftarrow \emptyset$
3:  **for all** $m \in M$ **do**
4:     **if** $H_{rom} + m_{rom} \leq C_{rom}$ **and** $H_{ram} + m_{ram} \leq C_{ram}$ **and** $H_{lat} + m_{lat} \leq C_{lat}$ **then**
5:        $H = H \cup \{m\}$
6:     **end if**
7:  **end for**
    **Pruning part**
8:  $S = H$
9:  $B = \emptyset$
10: $acc = 0$
11: **while** $S \neq \emptyset$ **do**
12:   $h = SelectRandomElement(S)$
13:   $S = S \setminus \{h\}$
14:   $acc_{temp} \leftarrow Accuracy(S)$
15:   **if** $acc_{temp} > acc - \epsilon$ **then**
16:     $acc = acc_{temp}$
17:     $B = S$
18:   **end if**
19: **end while**
20: **return** $B$

---

The classical reference approach selects a percentage of base classifiers with the highest training accuracy to build the static ensemble [19]. Our method incorporates TinyML constraints – selecting the highest-accuracy classifiers while staying within resource constraints. Our method uses the ($m\_acc/m\_rom + m\_ram + m\_lat$) key instead $m\_acc$ for selection. To compare the impact of selection criteria, we apply pruning to both methods.

---

[1]In sorting algorithms, the key is a characteristic of each element of the set against which the sorting is performed.

# 3. Experiments

Experiments were conducted using Python with TensorFlow and Keras libraries. The well-known MNIST [20], Fashion MNIST [21], USPS [22] and SVHN [23] datasets were used. The models were converted to UINT8 TFLite quantized format, and profiled using the ST Edge AI Developer Cloud infrastructure and STM32L4R9I-DISCO board. Notably, $m\_acc$ is device-independent, while $m\_rom$ and $m\_ram$ depend on the TensorFlow library, and $m\_lat$ is influenced by both the TensorFlow library and the profiled board[2]. An initial pool of MLP classifiers was generated using GridSearch, based on the hyperparameter space in Table 1. Due to the higher difficulty level of the SVHN dataset, it uses a more advanced architecture – BatchNormalization was added and a double hidden layer was used. To ensure diversity, various hyperparameters and a bagging technique were applied. The final prediction was obtained by averaging the models' scores.

Table 1. Hyperparameter space for an initial pool of MLP classifiers

| Hyperparameter | Values | Hyperparameter | Values |
|---|---|---|---|
| Number of hidden layers | {1, 2, 3} | Activation function | {"relu"} |
| Number of neurons per layer | {16, 32, 64, 128} | Optimizer | {"adam"} |
| Dropout | {0.0, 0.2, 0.3} | Learning rate | {0.01} |

Table 2 presents the experimental results based on the applied constraints and datasets. After initial experiments, value of $\epsilon$ was set to 0.001 for all runs, while other parameters were determined by the board specifications. Notably, the constraint parameters set by the operator apply to the pool before pruning, as there is a non-zero probability that no models will be removed during pruning. RAM usage is slightly higher in most cases for TinyCSaP, and this is related to the size of the ensemble (more models can provide better diversity, in which we see potential room for development). The ensemble built using the TinyCSaP algorithm consumes less ROM and has lower latency (except for the third configuration for Fashion MNIST), making it faster. The accuracy of TinyCSaP remains comparable to the baseline for MNIST and Fashion MNIST, with a maximum difference of 2.39% pt for USPS.

---

[2]When using external Flash/SRAM, profiled latency increases due to additional load time.

Table 2. Results of experiments depending on used constraints and datasets

| Constraints [KiB, KiB, ms] | ROM Usage TinyCSaP / Ref. [KiB] | RAM Usage TinyCSaP / Ref. [KiB] | Lat. Usage TinyCSaP / Ref. [ms] | Accuracy TinyCSaP / Ref. [%] | Ensemble Size TinyCSaP / Ref. [-] |
|---|---|---|---|---|---|
| MNIST | | | | | |
| ROM=250, RAM=30, Lat.=30 | 114.9 / 184.5 | 15.3 / 10.3 | 1.9 / 3.4 | 95.63 / 95.95 | 3 / 2 |
| ROM=300, RAM=50, Lat.=40 | 126.6 / 276.3 | 14.9 / 10.7 | 2.1 / 5.3 | 95.72 / 96.01 | 3 / 2 |
| ROM=350, RAM=70, Lat.=50 | 126.6 / 322.7 | 14.9 / 15.7 | 2.1 / 6.1 | 95.72 / 95.98 | 3 / 3 |
| Fashion MNIST | | | | | |
| ROM=250, RAM=30, Lat.=30 | 160.4 / 167.5 | 19.9 / 9.9 | 2.7 / 3.1 | 85.04 / 85.13 | 4 / 2 |
| ROM=300, RAM=50, Lat.=40 | 160.4 / 297.6 | 19.9 / 19.9 | 2.7 / 5.5 | 85.04 / 85.82 | 4 / 4 |
| ROM=350, RAM=70, Lat.=50 | 266.4 / 238.8 | 30.2 / 14.9 | 4.6 / 4.4 | 85.43 / 85.71 | 6 / 3 |
| USPS | | | | | |
| ROM=250, RAM=30, Lat.=30 | 59.7 / 127.3 | 5.8 / 6.2 | 0.7 / 2.0 | 89.44 / 91.73 | 2 / 2 |
| ROM=300, RAM=50, Lat.=40 | 57.1 / 220.7 | 6.2 / 12.0 | 0.6 / 3.4 | 89.89 / 91.83 | 2 / 4 |
| ROM=350, RAM=70, Lat.=50 | 142.4 / 199.5 | 14.9 / 9.5 | 1.5 / 3.2 | 89.74 / 92.13 | 5 / 3 |
| SVHN | | | | | |
| ROM=250, RAM=30, Lat.=30 | 149.4 / 223.6 | 18.9 / 12.6 | 2.8 / 4.3 | 72.08 / 72.07 | 3 / 2 |
| ROM=300, RAM=50, Lat.=40 | 204.9 / 297.7 | 25.2 / 18.9 | 3.9 / 5.8 | 71.35 / 73.57 | 4 / 3 |
| ROM=350, RAM=70, Lat.=50 | 204.9 / 336.3 | 25.2 / 12.6 | 3.9 / 6.6 | 71.35 / 72.14 | 4 / 2 |

## 4. Conclusions

In this paper, we presented a novel proof-of-concept method for selecting classifiers for a static ensemble, considering the constraints of TinyML systems. We presented that the proposed method produces ensembles that are faster while making better use of memory resources, with similar accuracy (in other words, by changing resources allocation we can manipulate the accuracy). The proposed method can extend the operating time of TinyML systems at the cost of a slight reduction in accuracy. By reducing ROM usage, it allows additional modules to be loaded into memory.

In future work, we aim to enhance the proposed concept by incorporating battery level dependency, energy harvesting with solar panels, and real-time ensemble size adjustment. Additionally, we plan to explore advanced selection methods, including dynamic ensemble selection techniques.

## References

[1] Ren, H. et al. TinyReptile: TinyML with federated meta-learning. In *IJCNN*. 2023.

[2] Warden, P. et al. Machine learning sensors. *Communications of the ACM*, 66(11):25–28, 2023.

[3] Banbury, C. et al. Benchmarking TinyML Systems: Challenges and direction, 2021.

[4] Hyeji, K. et al. Efficient neural network compression. *2019 IEEE CVPR*, 2018.

[5] Banbury, C. et al. MicroNets: Neural network architectures for deploying TinyML applications on commodity microcontrollers. 2021.

[6] Blalock, D. et al. What is the state of neural network pruning?, 2020.

[7] Gholami, A. et al. A survey of quantization methods for efficient neural network inference, 2021.

[8] Diao, R. et al. Feature selection inspired classifier ensemble reduction. *IEEE Transactions on Cybernetics*, 2014.

[9] Hernández-Lobato, D. et al. Statistical instance-based pruning in ensembles of independent classifiers. *IEEE TPAMI*, 2009.

[10] Bian, Y. et al. Ensemble pruning based on objection maximization with a general distributed framework. *IEEE TNNLS*, 2020.

[11] Kuncheva, L. I. *Combining Pattern Classifiers: Methods and Algorithms*. John Wiley and Sons Inc., 2004.

[12] Tamon, C. and Xiang, J. On the boosting pruning problem. In *European Conference on Machine Learning*. 2000.

[13] Martinez-Munoz, G. and Suárez, A. Using boosting to prune bagging ensembles. *Pattern Recognition Letters*, 28(1):156–165, 2007.

[14] Li, N. et al. Diversity regularized ensemble pruning. In *Machine Learning and Knowledge Discovery in Databases*. Springer Berlin Heidelberg, 2012.

[15] Dos Santos, E. M. et al. Overfitting cautious selection of classifier ensembles with genetic algorithms. *Information Fusion*, 10(2):150–162, 2009.

[16] Dos Santos, E. M. et al. Pareto analysis for the selection of classifier ensembles. In *Proceedings of the 10th Annual Conference on Genetic and Evolutionary Computation*, pages 681–688. 2008.

[17] Ruta, D. and Gabrys, B. Classifier selection for majority voting. *Information Fusion*, 6(1):63–81, 2005.

[18] Dantzig, G. B. Discrete-variable extremum problems. *Operations Research*, 5(2):266–288, 1957.

[19] Cruz, R. M. et al. Dynamic classifier selection: Recent advances and perspectives. *Information Fusion*, 41:195–216, 2018.

[20] LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010.

[21] Xiao, H. et al. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

[22] Hull, J. J. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1994.

[23] Netzer, Y. et al. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*. 2011.

# CHAPTER 11

# Young.AI

Track Chairs:

- Stanisław Kaźmierczak, MSc – Warsaw University of Technology
- Adam Żychowski, PhD – Warsaw University of Technology
- Tomasz Orczyk, PhD – University of Silesia in Katowice
- Tomasz Wesołowski, PhD, Eng – University of Silesia in Katowice

# A Study of the Creative Process
# in Generative Art

**Kamil Gajos**[1][0009−0006−0766−1938], **Urszula Boryczka**[2][0000−0002−2698−6934]

*University of Silesia in Katowice*
*Faculty of Science and Technology*
*Będzińska 39, 41-200 Sosnowiec, Poland*
*kamil.gajos@o365.us.edu.pl, urszula.boryczka@us.edu.pl*

**Abstract.** *The study focuses on the creative process using Generative Adversarial Networks (GANs) and how the search space impacts the final generated results. In this work, three sizes of datasets are used to showcase the impact of insufficient examples for the GAN to learn from. A further enhancement is made by utilizing a denoiser on the generated work, which proves that adding it enhances quality by removing artifacts caused by the generator's limited dataset. The conclusions indicate that, with sufficient hardware, it is possible to recreate the style of a given author using GANs in combination with a denoiser.*
**Keywords:** *artificial intelligence, generative art, generative adversarial networks, neural network, convolutional neural network*

## 1. Introduction

The concept of recreating human creativity using artificial intelligence was once deemed impossible. The complex and nuanced processes through which humans create art infused with personal expression were thought to be unattainable by deep neural networks or large language models [1].

However, advancements in technology, especially through generative adversarial networks (GANs), have challenged this perception, enabling artificial intelligence to replicate the styles of renowned artists like Rembrandt, Picasso, and Van Gogh [2]. In their study, the authors present a framework that utilizes two multilayer perceptrons: one dedicated to generating new works of art and the other

tasked with determining whether the generated artwork comes from the original training dataset.

An engaging analogy drawn in their work compares this process to the competition between a counterfeiter and a policeman, each trying to outwit the other. The objective of this paper is to explore the feasibility of mathematically modeling creativity itself. To evaluate the quality of the art produced by this method, we apply the Fréchet Inception Distance (FID) [3], calculated using the formula provided in equation (1)

$$d^2 = \left| \mu - \mu' \right|_2^2 + tr(\Sigma + \Sigma' - 2\left(\Sigma\Sigma'\right)^{\frac{1}{2}}). \tag{1}$$

Technology is advancing rapidly, leading to the emergence of various innovative approaches, including those discussed in the work titled "What is Generative Art?" [4]. Generative art is now broadly understood as art created with the help of an autonomous system, one that is designed by the artists themselves. Importantly, it took time for a clear taxonomy of generative art to emerge; this understanding evolved later in the 20th century.

Our autonomous system comprises three neural networks: two that engage in adversarial interactions and a third that refines the output by eliminating excess noise from the generated images. Notably, our approach diverges from typical Generative Adversarial Networks (GANs) by employing a *Wasserstein* GAN (WGAN) [5]. The key innovation of WGAN is its incorporation of a gradient penalty and multiple critics, which help to bridge the knowledge gap between the discriminator and the generator within the learning process. In our system, we found that other hyperparameters are comparatively less significant and impactful on the quality of the final generated artwork, as detailed in the original cited paper.

## 2. Methodology

This project aimed to replicate Pablo Picasso's iconic Cubist style using a Wasserstein Generative Adversarial Network (WGAN). To train the discriminator effectively in recognizing Picasso's creative features, we adopted a "divide and conquer" strategy. We started by developing a basic Generative Adversarial Network (GAN) using the MNIST dataset [6]. We customized both the generator and discriminator networks to fit the dataset's specific characteristics, utilizing TensorFlow as our framework in Python [7] and leveraging the capabilities of a high-performance GPU.

Initially, our results were disappointing; the generated images were blurred after 10,000 iterations on the Nvidia 4090 GPU, which was the fastest consumer GPU at that time. To enhance the visual quality of the generated artwork, we introduced a gradient penalty and increased the number of critics within the WGAN framework. These changes led to significant improvements, resulting in clearer and more recognizable representations of handwritten digits (Figure 1).



(a)    Original entry

(b) Best generated entry

(c)    Original entry

(d) Best generated entry

Figure 1. MNSIT original entry compared to generated entry

## 2.1. Moving Past Greyscale

After achieving good results on a large dataset, it was time to handcraft a new, much smaller dataset and experiment with utilizing more colors, rather than simply relying on grayscale. This led to the choice of vision deficiency tests, which offered four channels including alpha, and presented a resizing challenge to standardize all samples to a settled size of 140 x 140 pixels. New problems arose with hyperparameters that were not tuned for the new samples. After tweaking the structure of the neural networks, learning speed, gradient penalty, and the number of training batches available for the generator, discriminator, and denoiser, as well as incorporating more information for individual pixels with the addition of three extra channels, the resulting images still lacked clarity.

Upon further investigation, it became clear that the samples varied too widely from each other, leading to only a few features being extracted effectively by the

standard feedforward network [8]. However, after switching the discriminator to a convolutional neural network [9] while keeping the generator as a standard neural network type, the quality of the images improved significantly (Figure 2).



(a) Original entry

(b) Best generated entry



(c) Original entry

(d) Best generated entry

Figure 2. Original vision deficiency test compared to generated test

## 2.2. Generative Art

To recreate the artwork, we searched online for samples that exemplified Pablo Picasso's style, helping us effectively train the discriminator. With a curated dataset of similar works in hand, we refined the neural network architecture and hyper-parameters. After thousands of iterations, we tuned the Wasserstein Generative Adversarial Network (WGAN) to capture the essence of Picasso's style. The final implementation is presented in a series of steps for generating new artwork, as shown in Figure 3.



Figure 3. Architecture of the implemented solution

---

**Algorithm 1** Description of creating art by implemented generative adversarial network model in the work

---

1: Create neural network topology for the discriminator, generator and denoiser based on given art.
2: Initialize and tweak hiperparameters like number of critics, gradient penalty applied to the discriminator.
3: Train the discriminator on real art.
4: Initialize latent space with noise function as input for the generator.
5: Generate art using the generator and then validate it with the discriminator.
6: If the discriminator was not fooled, go to Step 4. Otherwise, go to the next step.
7: Pass the image to denoiser in order to remove artifacts not detected by the discriminator.

---

### 2.3. Results

The largest dataset produced highly accurate results, as shown in Figure 1. However, due to the limited number of pixels and channels, the denoising process struggled to effectively eliminate the excess noise generated by the model.

Increasing both the search space and the number of channels significantly extended the computational time required to produce the final artwork. For instance, the time increased from just one minute for 1,000 iterations to over for minutes for larger configurations. The dataset used to generate the results in Figure 2 was considerably smaller than the one used for the *MNSIT* dataset. Consequently, the numbers that should be easily visible to individuals without visual impairments became difficult to discern due to substantial fluctuations from sample to sample.

Lastly, we developed a specialized architecture to explore whether we could successfully replicate the iconic art style of Pablo Picasso. Figure 4 showcases two examples from the thousands of generative artworks produced, highlighting the most prominent features that align with the original creations.

## 3. Conclusions

The highest quality generative art produced in this study has been constrained to resolutions of 300 x 300 pixels with three color channels, due to limited VRAM capacity, which is insufficient for accommodating fully prepared models. This ar-

(a) Original entry   (b) Generated entry



(c) Original entry   (d) Generated entry

Figure 4. Pablo Piccaso original work compared to generated art

ticle investigates how increasing the number of color channels in GANs leads to an exponential rise in computational requirements, thereby necessitating larger model sizes to effectively generate visually compelling artworks. The results presented in Table 1 indicate that the implemented model achieves optimal artistic quality when trained for between 1,000 and 3,000 epochs; however, performance tends to decline when training exceeds this range.

Table 1. FID average score for generated images

| Epoch | MNIST | Deficiency | Picasso |
|---|---|---|---|
| 0 | 617180,21 | 63,15 | 83,86 |
| 1000 | 977503,10 | 40,94 | 27,82 |
| 2000 | 911610,79 | 31,52 | 26,17 |
| 3000 | 866935,77 | 28,69 | 24,82 |
| 4000 | 930621,25 | 29,26 | 26,54 |
| 5000 | 765668,41 | 32,63 | 25,07 |

Interestingly, the evaluation of the grayscale MNIST results, assessed using the Fréchet Inception Distance (FID) metric, yielded atypical outcomes. Normally, a lower FID value indicates superior performance; nevertheless, the model exhibited its best results at 0 epochs. A significant conclusion drawn from this study is that the entries in the dataset must be closely aligned for the neural network's discriminator to effectively learn the intricate features that should be preserved

in the generated images. Moreover, the inherent limitations of smaller images, which contain fewer details, further complicate the learning curve.

# References

[1] Wang, A., Yin, Z., Hu, Y., Mao, Y., and Hui, P. Exploring the potential of large language models in artistic creation: Collaboration and reflection on creative programming. *arXiv preprint arXiv:2402.09750*, 2024. doi:10.48550/arXiv. 2402.09750.

[2] Goodfellow, I. et al. Generative adversarial nets. In Z. Ghahramani et al., editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. doi:10.48550/arXiv.1406.2661.

[3] Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training gans. *Advances in Neural Information Processing Systems*, 29, 2016.

[4] Boden, M. and Edmonds, E. What is generative art? *Digital Creativity*, 20:21–46, 2009. doi:10.1080/14626260902867915.

[5] Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, pages 214–223. PMLR, 2017. doi:10.48550/arXiv.1701.07875.

[6] Deng, L. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.

[7] Kalin, J. *Generative Adversarial Networks Cookbook: Over 100 Recipes to Build Generative Models Using Python, TensorFlow, and Keras*. Packt Publishing Ltd, 2018.

[8] Bebis, G. and Georgiopoulos, M. Feed-forward neural networks. *Ieee Potentials*, 13(4):27–31, 1994.

[9] O'shea, K. and Nash, R. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*, 2015.

# Evolutionary Design of Graph Neural Network Architectures for Graph Classification

**Maciej Krzywda**[1]**, Szymon Łukasik**[2]**, Amir H. Gandomi**[3,4,5]

[1]*Faculty of Physics and Applied Computer Science*
*AGH University of Krakow*
*al. Mickiewicza 30, 30-059 Kraków, Poland*
[2]*Systems Research Institute, Polish Academy of Sciences*
*ul. Newelska 6, 01-447 Warsaw, Poland*
[3] *Faculty of Engineering and IT, University of Technology Sydney*
*5 Broadway, Ultimo NSW 2007, Australia*
[4]*University Research and Innovation Center (EKIK), Óbuda University*
*Bécsi út 96/B, Budapest, 1034, Hungary*
[5]*Department of Computer Science, Khazar University*
*Mahsati 41, Baku, Azerbaijan*

**Abstract.** *Graph Neural Networks (GNNs) are a family of neural architectures operating on graphs, inspired by interactions between nodes on a graph. This study develops the first stage of a method that uses evolutionary algorithms to design and optimize graph convolutional neural networks for graph classification. Our method allows us to set critical parameters, such as batch size, number of hidden channels, and the choice and parameters of optimizer and loss function. The approach uses multi-criteria design and can be applied to the architectural-level hyperparameter tuning of graph neural networks, such as Graph Attention Networks (GAT), Graph Isomorphism Networks (GIN), and Graph Convolutional Networks (GCN). We evaluate the proposed method on benchmark datasets and observe encouraging performance from the synthesized networks.*
**Keywords:** *evolutionary computing, graph neural networks, neural architecture search, genetic algorithms, graph classification*

## 1. Introduction

Neural Architecture Search (NAS) is becoming increasingly popular as a fully automated approach to designing artificial neural network (ANN) architectures.

NAS enables the development of neural network structures that are comparable or even superior in performance to those created manually. NAS can be performed using various strategies, among others, including the evolutionary approach [1, 2]. Essentially, NAS automates the process of searching for effective neural network configurations that would otherwise be tuned manually through trial and error. This approach encompasses various methods and utilities that assess multiple network architectures within a specified search space using a search strategy, selecting the one that best achieves the objectives of a given problem by maximizing a fitness function. According to the definition provided earlier, NAS can be described through Evolutionary Computation (EC) as navigating the search space for architectures by creating and assessing a group of potential architectures [3]. Our approach is to perform multi-criteria evolutionary design of graph neural networks for graph classification. In our approach, we have chosen to search for the network that attains the highest fitness (performance), the shortest training time for an individual across all generations, and the network with the fewest parameters. In our evolutionary search, we minimize computational cost while striving to achieve the best fitness value. During this process, we analyze all individuals in each generation regarding parameter count, individual training time, and the mentioned fitness value, preserving three solutions (which may or may not be represented by the same individual). Thus, we retain the individual that achieved the best fitness, the shortest training time for an individual across all generations, and the network with the fewest parameters. This method allows us to cover a more extensive search and decide which network configuration is the most appropriate for a given application.

## 2. Graph Neural Networks and Their Design

Graph Neural Networks (GNNs) are a family of neural networks inspired by the mechanisms between nodes on a graph. In recent years, there has been an increased interest in GNNs and their derivatives, i.e. graph-attention networks (GAT) [4], graph-convolutional networks (GCN) [5], and Graph Isomorphism Networks (GIN) [6]. The diversity of available architectures raises the question of how to select appropriate parameters for applying a chosen network to a practical problem. In this work, we focus on three widely used graph neural network architectures for graph classification. Remarkably, the field of NAS within GNNs remains relatively unexplored. Notable exceptions include GraphNAS [7], which

uses RL to identify architectures for node classification tasks. Auto-GNN [8], SNAG framework (Simplified Neural Architecture Search for Graph Neural Networks [9]). Upon examining the scope of existing studies, no similar work has been identified to tune the same set of parameters of GNNs for graph classification.

## 3. Experimental Study

In our approach, an **entity** represents an individual graph neural network (GNN) configured with randomly initialized parameters, allowing it to compile and operate successfully on selected datasets. The evolutionary procedure employs a genetic algorithm consisting of mutation and crossover operators. A mutation randomly assigns a new gene value drawn uniformly from its permissible range. The crossover operator selects two parent entities from the top 10 best-performing neural networks (based on their `fitness_value`) and generates offspring by combining portions of their genetic information and introducing mutation to further enhance genetic diversity. This further randomized selection process continues iteratively until a predefined population size is reached. It should be noted that, although we originally referred to this process simply as `random_selection`, it in fact involves randomly choosing parents exclusively from the subset of best-performing networks. Throughout the evolutionary process, the initially selected GNN architecture remains fixed. The fitness value is defined as the sum of squared errors of the F1 scores across all classes. The F1 score evaluates the quality of the classification considering precision, sensitivity, and recall for each class individually. This definition of the fitness function provides deeper insight into graph classification performance, which is particularly beneficial when handling imbalanced datasets, as frequently encountered in molecular applications. We defined three models. The GAT-based graph classifier employs three GATConv layers followed by a linear classification layer (four layers in total), the **Graph Classifier based on GCN** similarly uses three GCNConv layers and one linear layer (also four layers total), whereas the **Graph Classifier based on GIN** utilizes three GINConv layers combined with linear transformations, batch normalization, and ReLU activations, followed by two linear layers for graph-level aggregation (five layers total). The **ENZYMES**, **PROTEINS**, and **MUTAG** datasets from the TU-Dataset collection [10] are widely recognized benchmarks in graph-based machine learning research, comprising structured graph data from chemical and biological

domains. Detailed characteristics and compositions of these datasets are provided in Table 1.

Table 1. Dataset characteristics and best test accuracy (%)

| Parameter/Model | ENZYMES | PROTEINS | MUTAG |
|---|---|---|---|
| Graphs count | 600 | 1113 | 188 |
| N nodes avg/std/min/max | 32.63/15.28/2/126 | 39.06/45.76/4/620 | 17.93/4.58/10/28 |
| N edges avg/std/min/max | 62.14/25.50/1/149 | 72.82/84.60/5/1049 | 19.79/5.68/10/33 |
| Node degree avg/std/min/max | 3.81/1.15/0/9 | 2.16/0.78/0/4 | 3.73/1.15/0/25 |
| Node features dimension | 3 | 3 | 7 |
| Number classes | 6 | 2 | 2 |
| Baseline 1 [11] | 55.67% | N/A | 100% |
| Baseline 2 [12] | 68.79% | 84.91% | N/A |
| Baseline 3 [13] | 68.79% | 82.9% | 96.4% |
| Baseline 4 [6] | 70.17% | 76.46% | 90.84% |
| Our GCN (best fit) | 78.88% | 71.19% | 84.21% |
| Our GIN (best fit) | 88.33% | 84.24% | 100% |
| Our GAT (best fit) | 81.66% | 74.16% | 89.47% |
| Our GCN (min time) | 77.77% | 71.99% | 84.21% |
| Our GIN (min time) | 86.11% | 80.35% | 94.73% |
| Our GAT (min time) | 81.11% | 72.50% | 84.21% |
| Our GCN (min params) | 76.66% | 68.85% | 84.21% |
| Our GIN (min params) | 84.44% | 79.46% | 100% |
| Our GAT (min params) | 80.55% | 72.32% | 84.21% |

# 4. Results

Our approach created an archive to systematically store the best-performing graph neural networks (GNNs) that evolved during the optimization process. Specifically, we tracked three distinct criteria: (i) the best fitness achieved, (ii) the shortest training time required by an individual network, and (iii) the minimal number of network parameters. Evolutionary optimization was carried out for 100 generations with a population of 20 individuals per generation. The genotype representing each candidate solution comprised a wide range of hyperparameters for graph neural networks. Although the number of epochs was fixed at 100 in the current experiments, it is included in the genotype to enable future optimization. Batch sizes were allowed to vary between 10 and 1,000, while the number of hidden channels in GNN convolutional layers ranged from 1 to 200. The learning rate and weight decay hyperparameters were set within the range of 0.0001 to 0.1. The available loss functions included CrossEntropy, Negative Log-Likelihood (NLL-Loss), MultiMarginLoss, MultiLabelMarginLoss and SmoothL1Loss. In addition, we incorporated a comprehensive set of optimization algorithms consisting of 37 different optimizers. This extensive collection included popular algorithms such

as Adam, RMSprop, Adagrad, and SGD, as well as advanced variants such as AdaBelief, AdamW, RAdam, Ranger, MADGRAD, and the Particle Swarm Optimizer (PSO) inspired by the approach proposed in [14]. This wide selection of optimization methods allowed the evolutionary process to explore various training strategies, enabling us to identify GNN architectures and training procedures that perform optimally under various conditions.

Based on the results presented in Table 1, our approach achieves a performance close to that of state-of-the-art architecture specifically fine-tuned for these problems. Furthermore, our models do not include the evolution of the aggregation function type or the dropout rate, which can have an impact on the results achieved. In our experiments, we considered 41 combinations of optimisation algorithms and loss functions. However, in our results, for all models that achieved **the best fitness value**, the optimizer most commonly used was the quasihyperbolic momentum algorithm (QHM) [15]. For **minimum training time**, the dominant solutions used Adagrad and Adadelta; for **minimum network parameters**, Adam was used most frequently. In our investigation, we have determined that Particle Swarm Optimization (PSO) significantly influences the results achieved, particularly in terms of training duration [14]. However, the performance outcomes obtained with PSO were comparable to those achieved by QHAdam, demonstrating its efficacy in optimizing GNNs despite its time-intensive nature. Regarding loss functions, categorical cross-entropy was dominant among the best-performing configurations. `Smooth L1 Loss` also appeared among the solutions, but it did not yield the best results. A potential direction for improving GNN performance in our setting is the use of alternative cost functions. However, alternative cost functions such as Savage [16], $L_o$ [17] or `Loge` [18] have been suggested and could be worth exploring.

## 5. Conclusions

In recent years, research on Artificial Neural Networks (ANNs) has gained momentum, driven by the improved affordability of training and prototyping deep learning structures such as Neural Architecture Search. This study presents a methodology for designing and optimising graph neural networks for graph classification tasks. Our approach avoids treating the optimization as a pure black box and instead explicitly analyzes training time and model complexity alongside accuracy. Our methodology aims to achieve optimal accuracy for graph neu-

ral networks and focuses on simultaneously minimizing training time and model complexity. The findings of this study are expected to benefit researchers working on neural architecture search, offering promising implications not only for graph classification but also for other domains that utilize graph-based neural network solutions.

## Acknowledgment

## References

[1] Shang, R. et al. Evolutionary neural architecture search based on evaluation correction and functional units. *Knowledge-Based Systems*, 251:109206, 2022. ISSN 0950-7051.

[2] Pan, C. and Yao, X. Neural architecture search based on evolutionary algorithms with fitness approximation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021.

[3] Zhou, X. et al. A survey of advances in evolutionary neural architecture search. In *2021 IEEE Congress on Evolutionary Computation (CEC)*, pages 950–957. IEEE, 2021.

[4] Veličković, P. et al. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2018.

[5] Kipf, T. and Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2017.

[6] Xu, K. et al. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.

[7] Gao, Y. et al. Graph neural architecture search. In *International Joint Conference on Artificial Intelligence (IJCAI'20)*, pages 1403–1409. IJCAI, 2020.

[8] Zhou, K. et al. Auto-GNN: Neural architecture search of graph neural networks. *Frontiers in Big Data*, 5, 2019.

[9] Zhao, H., Wei, L., and Yao, Q. Simplifying architecture search for graph neural network. *arXiv preprint arXiv:2008.11652*, 2020.

[10] Morris, C. et al. TUdataset: A collection of benchmark datasets for learning with graphs. *arXiv preprint arXiv:2007.08663*, 2020.

[11] Domingue, M. et al. Evolution of graph classifiers. In *2019 IEEE Western New York Image and Signal Processing Workshop (WNYISPW)*, pages 1–5. IEEE, 2019. doi:10.1109/WNYIPW.2019.8923110.

[12] Zhang, Z. et al. Hierarchical multi-view graph pooling with structure learning. *IEEE Transactions on Knowledge and Data Engineering*, 35(1):545–559, 2023. doi:10.1109/TKDE.2021.3090664.

[13] Vincent-Cuaz, C. et al. Template based graph neural network with optimal transport distances. *arXiv preprint arXiv:2205.15733*, 2022.

[14] Shen, R. et al. Training graph neural networks with particle swarm optimisation. *Sacair 2023*, 2023.

[15] Ma, J. and Yarats, D. Quasi-hyperbolic momentum and adam for deep learning. *arXiv preprint arXiv:1810.06801*, 2018.

[16] Masnadi-Shirazi, H. and Vasconcelos, N. On the design of loss functions for classification: theory, robustness to outliers, and savageboost. In D. Koller et al., editors, *Neural Information Processing Systems*, pages 1049–1056. 2008.

[17] Zhang, Z. and Sabuncu, M. R. Generalized cross entropy loss for training deep neural networks with noisy labels. *arXiv preprint arXiv:1805.07836*, 2018.

[18] Wang, Y. et al. Bag of tricks for node classification with graph neural networks. *arXiv preprint arXiv:2103.13355*, 2021.

# An AI-Driven Approach to Adapting the Expected Goals (xG) Model to Women's Football

**Tomasz Lipowski**[0009−0002−8085−6628], **Tomasz Piłka**[0000−0003−1206−2076]

*Adam Mickiewicz University*
*Faculty of Mathematics and Computer Science*
*ul. Uniwersytetu Poznańskiego 4, 61-614 Poznań, Poland*
*tomlip2@st.amu.edu.pl, tomasz.pilka@amu.edu.pl*

**Abstract.** *The Expected Goals (xG) model is a key metric in football analytics, yet conventional models overlook biomechanical and tactical differences in women's football. This study introduces an AI-driven xG model, integrating novel contextual variables like Defensive Congestion Index (DCI) and Shot Block (SB). Using neural networks, the model improves predictive accuracy over traditional approaches. Results show distinct shot conversion patterns in women's football, particularly for long-range attempts, underscoring the need for gender-specific modeling. Incorporating additional features reduces log loss and enhances shot predictions.*
**Keywords:** *football analytics, artificial intelligence, xG, machine learning models*

## 1. Introduction

Expected Goals (xG) is a widely accepted metric in soccer analytics that estimates the probability of a shot resulting in a goal, leveraging historical data and contextual factors [1, 2]. Originally a concept from gaming culture, it is now a key tool for teams, analysts, and betting companies to guide performance assessment and strategic planning [3, 4]. By considering elements such as shot distance, angle, and defensive context, xG helps address the low-scoring, unpredictable nature of the sport, providing deeper insights into both team and player performance [5]

Key variables influencing xG include:
- **Shot location**: Distance from goal, angle, and positioning.

- **Shot type**: Footed shots, headers, and set-piece scenarios.
- **Defensive pressure**: Number of defenders near the shooter.
- **Assisting pass type**: Through ball, cross, or cut-back.

Expected Goals transforms match data into an assessment of the probability that a shot will result in a goal, going beyond traditional statistics. It takes into account factors such as shot location, type and defensive pressure to help teams better understand scoring potential and performance.

However, most xG models struggle when applied to women's football, as they rely heavily on data from men's matches. Biomechanical differences, such as variations in shot power, trajectory and accuracy, challenge the predictive validity of these metrics. In addition, the slower pace of women's football affects shot-creation dynamics and defensive formations, while tactical variations, including different pressing strategies and defensive setups, further differentiate women's matches from the data on which these models are based.

To address these shortcomings, this research aims to systematically evaluate the need for gender-specific adjustments in xG calculations, identify critical contextual variables that can improve predictive accuracy, and develop an AI-enhanced xG framework that incorporates these variables. By addressing these technical limitations, the study aims to produce a more robust and contextually relevant xG model for women's football.

## 2. Literature Review

Recent work refines xG models for greater predictive accuracy, for example by incorporating pre-shot event sequences [6] or integrating xG with metrics like xA and xPTS [7]. Advanced machine learning methods further enhance defensive evaluations [8], informing performance analysis, tactics, and recruitment. Bransen and Davis [9] adapt xG models for women's soccer, showing that while men's models can partly transfer, distinct shot patterns in the women's game necessitate tailored approaches.

Crucially, the growing body of research on women's football reveals distinct dynamics that require gender-specific modeling. Studies such as [10, 11] have documented the unique physical and tactical aspects of women's matches, from differences in biomechanics and shooting patterns to variations in defensive structures. This evidence supports the call for gender-specific xG models that reflect these differences. The authors also highlight the role of ML in identifying patterns

that can prevent injuries, particularly in the context of women's football, where training loads and physical conditioning play a crucial role in player safety and performance. Using the AI and ML methods highlighted in these studies, football analytics can provide more tailored and accurate insights. These approaches allow analysts to adapt models to the specific needs of men's and women's football, providing a deeper, data-driven understanding of player and team dynamics.

# 3. Methodology

## 3.1. Data Collection

The dataset for this study was mainly sourced from Hudl StatsBomb[1] Open Data service. StatsBomb offers comprehensive match data, including team and player profiles, possession sequences, individual player actions and event locations, providing a complete perspective of the game's dynamics.

Event data captures every logged action (passes, shots, tackles, dribbles) with precise timestamps and pitch coordinates. Our analysis focused on shots, examining location (distance, angle), pass type, and shooting technique (including foot preference). We also introduced two contextual variables: the ***Defensive Congestion Index (DCI)*** for defenders near the shooter, and ***ShotBlock (SB)*** for opponents in the shot path (see Figure 1).



| Attribute | Value |
|---|---|
| distance (m) | 29.88 |
| goalkeeper in shot keeper cone | TRUE |
| players in shot keeper cone | 3 |
| enemy players in shot keeper cone | 2 |
| players in the box | 6 |
| pressing | 0 |
| StatsBomb xG | 0.021 |
| my xG | 0.012 |

Figure 1. Visual representation of the generated attributes

For this study, we used StatsBomb Open Data from multiple seasons, encompassing both league and tournament matches. The breakdown of these datasets is presented in Table 1.

---

[1] `https://statsbomb.com/`.

Table 1. List of leagues and tournaments used for the study

| Competition Name | Year/Season | Gender | Number of Matches | Number of Shots |
|---|---|---|---|---|
| La Liga | 2015/2016 | male | 380 | 9071 |
| FA Women's Super League | 2018-2021 | female | 326 | 8239 |
| FIFA World Cup | 2022 | male | 128 | 3068 |
| Women's World Cup | 2019, 2023 | female | 116 | 2891 |
| National Women's Soccer League | 2018/2019 | female | 36 | 1034 |
| UEFA Women's Euro | 2023 | female | 31 | 871 |

## 3.2. AI-Based Model Development and Evaluation Metrics

The AI-based model development involves two approaches: logistic regression, which serves as the baseline model, and MLP (Multilayer Perceptron) neural networks. An MLP is a type of feedforward artificial neural network that consists of multiple layers of neurons, including an input layer, one or more hidden layers, and an output layer. The MLP model used here consists of three hidden layers with 64, 64, and 32 neurons, followed by an output layer with one neuron using a sigmoid activation function for binary classification. The optimizer used is Adam, and the loss function is binary cross-entropy. The model evaluation is performed using two metrics: log loss, which ensures probability calibration and is fully compatible with logistic regression, and shot outcome comparison, which measures actual versus predicted goal rates.

## 4. Results and Discussion

**Authorial model – *My xG*:** As a result of the experiments, it was observed that the use of the same xG model in both men's and women's football leads to **significant errors**, especially for long-range shots, which have higher conversion rates in women's football. The newly introduced attributes show great importance in improving the predictive performance of the model, ranking among the most influential factors, see Table 2. Shots, considered as long-distance shots, are from the area marked in Figure 2. Analysis shows that women achieve better results from long-range shots (Table 3).

The introduction of **DCI** and **SB** improved the efficiency of the neural network, allowing it to adapt better to different game conditions. The logarithmic loss falls below the level of the StatsBomb model with the original attributes of Table 4, compared to our extended model with additional attributes.

Table 2. Importance of features

| Features | Importance |
|---|---|
| distance | 0.350 |
| angle degrees | 0.307 |
| players in shot keeper cone | 0.197 |
| enemy players in shot keeper cone | 0.178 |
| players in the box | 0.161 |
| pressing | 0.156 |



Figure 2. A shot from the blue area or inside the box is not a long-range shot

Table 3. Results from long-range shots

| | Female | Male |
|---|---|---|
| Long-range shots | 3507 | 3554 |
| Long-range goals | 129 | 87 |
| Long-range accuracy (%) | 3.68 | 2.45 |
| Long-range frequency (%) | 26.90 | 29.28 |
| Long-range goals/all (%) | 0.99 | 0.72 |

Table 4. Log loss without and with DCI and SB

| | Original attributes | Extended attributes |
|---|---|---|
| My xG | 0.273 | 0.262 |
| StatsBomb xG | 0.265 | 0.265 |

A comparison of the distribution of xG values obtained for the women's and men's football data, respectively, by reconciling the original xG values for the data from StatsBomb and the modification proposed in this paper, are placed in Figure 3 and Figure 4, respectively.



Figure 3. Distribution of StatsBomb xG vs our xG proposals for women

Figure 4. Distribution of StatsBomb xG vs our xG proposals for men

## 5. Conclusion and Future Work

This study underscores the need to adapt xG models to the distinct features of women's football. By incorporating AI-driven techniques and contextual variables like Defensive Congestion Index (DCI) and Shot Block (SB), we achieve greater predictive accuracy. Notably, women's long-range shots yield higher conversion rates, highlighting the importance of gender-specific modeling. AI-enhanced methods also address existing xG limitations, providing more reliable tactical insights. Future work should expand datasets with tracking data, integrate computer vision, and develop player-specific models for deeper performance assessments, coaching, and scouting decisions.

## References

[1] StatsBomb. What are expected goals (xG)? `https://statsbomb.com/soccer-metrics/expected-goals-xg-explained`, 2024. Accessed: 2025-02-15.

[2] Simpson, M. and Craig, C. Developing a new expected goals metric to quantify performance in a virtual reality soccer goalkeeping app called *CleanSheet*. *Sensors*, 24, 2024. doi:10.3390/s24237527.

[3] Anzer, G. and Bauer, P. A goal scoring probability model for shots based on synchronized positional and event data in football (soccer). *Frontiers in Sports and Active Living*, 3, 2021. doi:10.3389/fspor.2021.624475.

[4] Analyst, O.  *What Is Expected Goals (xG)?*, 2023.  URL `https://theanalyst.com/2023/08/what-is-expected-goals-xg`. Accessed: 2025-02-05.

[5] Hewitt, J. H. and Karakuş, O.  A machine learning approach for player and position adjusted expected goals in football (soccer).  *arXiv preprint arXiv:2301.13052*, 2023.

[6] Bandara, I., Shelyag, S., Rajasegarar, S., Dwyer, D., Kim, E., and Angelova, M.  Predicting goal probabilities with improved xG models using event sequences in association football.  *Plos One*, 19, 2024.  doi: 10.1371/journal.pone.0312278.

[7] Khrapach, V. and Siryi, O. Statistical metric xG in football and its impact on scoring performance: A review article. *Health Technologies*, 2:47–54, 2024. doi:10.58962/ht.2024.2.3.47-54.

[8] Zaręba, M., Piłka, T., Górecki, T., Grzelak, B., and Dyczkowski, K.  Improving the evaluation of defensive player values with advanced machine learning techniques.  In *Harnessing Opportunities: Reshaping ISD in the Post-COVID-19 and Generative AI Era (ISD2024 Proceedings)*. 2024.  doi: 10.62036/ISD.2024.67.

[9] Bransen, L. and Davis, J. Women's football analyzed: interpretable expected goals models for women, 2021.

[10] Pappalardo, L., Rossi, A., Natilli, M., and Cintia, P. Explaining the difference between men's and women's football. *Plos One*, 16:e0255407, 2021. doi: 10.1371/journal.pone.0255407.

[11] Eetvelde, H. V., Mendonça, L. D. M., Ley, C., Seil, R., and Tischer, T. Machine learning methods in sport injury prediction and prevention: a systematic review. *Journal of Experimental Orthopaedics*, 8, 2021. doi: 10.1186/s40634-021-00346-x.

# Recognizing Selected Sign Language Gestures Using Deep Models

**Adam Łaput**[0009−0007−5312−0769]

*adam.laput@wp.pl*

**Abstract.** *Recognizing sign language gestures is an extremely complex and complicated problem, but the development of artificial intelligence methods allowed for attempts to solve it. The article discusses the complexity of recognizing sign language gestures and proposes possible solutions for identifying them using LSTM and CNN deep network models, with particular emphasis on Polish Sign Language (PJM) gestures. Due to the lack of availability of a data set for PJM, a data set for 40 gestures was developed independently.*
**Keywords:** *recognizing, sign language, deep models*

## 1. Introduction

Sign languages were present among the ancient cultures, but the birth of the Polish Sign Language (PJM) is associated with the person of pr. Jakub Falkowski. He created the Institute of the Deaf in Warsaw in 1817. In 1879, the first dictionary of the Polish Sign Language was published, containing approximately 10,000 signs [1]. However, there are difficulties in directly translating statements in sign language into other languages due to the large number of components of a given sign, including gestures made by both hands and their location in relation to the body, posture and body movements during the articulation of the sign and appropriate facial expressions [2]. Apart from complicated signs, PJM also contains basic signs that require only one hand to produce, such as the single letters that appear in the PJM alphabet. It contains 36 gestures (moving and stationary) denoting various letters in the Polish alphabet.

The main contributions of this work are as follows:

- creating a dataset consisting of images and data regarding the location of hand elements and the hand itself;

- development of LSTM and CNN network models enabling the recognition of sign language gestures;

- assessment of the quality of gesture recognition by previously developed deep network models.

## 2. Related Work

The problem of recognizing sign language gestures is very often discussed in scientific works, and public data sets are available for some of these languages, e.g. ASL (American Sign Language)[1]. In [3], the authors focused on recognizing Indian sign language gestures and created a data set consisting of 200 gestures. Using a CNN network model (four convolutional layers and five ReLu layers) performed gesture recognition from reduced images. As a result of the experiments, a recognition accuracy of 92.88% was achieved. For the purposes of recognizing Flemish sign language in [4], the authors, in addition to using CNN networks, implemented LSTM networks in their experiments and used the OpenPose library for reading the coordinates body points. As a result of experimental tests, the highest recognition result was obtained at the level of 74.7%. A description of the CNN and LSTM networks in the process of recognizing sign language gestures can also be found in the publication [5]. Depending on the designed architecture of the CNN network model and the selected sign language, the obtained results of the *accuracy* metric oscillated between 81% and 91%, but it takes time to train. Among the LSTM network models, there was an element of using the Kinect 2.0 device to precisely detect the position of joints.

## 3. Dataset

To read data from the hand, the Leap Motion device[2] was used. Using two built-in cameras, it not only detects hands, but also characteristic points in their structure, the angle of the hand, the length and width of each finger and the type of hand (left or right hand).

Due to the lack of a publicly available dataset for PJM, a dataset of selected gestures had to be created from scratch. After selecting 40 PJM gestures – 36 ges-

---

[1]`https://www.kaggle.com/datasets/grassknoted/asl-alphabet`.
[2]`https://docs.ultraleap.com/hand-tracking/getting-started.html`.

tures representing all the letters of the PJM alphabet and four words (Monday, Sunday, spring, autumn), data collection was carried out from 17 different individuals who did not have professional experience in PJM. Each person made 10 gestures, each consisted of three signs (appearing at the beginning, in the middle and at the end of the gesture) – 100 samples per sign were recorded only from the right hand. In addition to recording the infrared photo, the lighting values (in the LUX unit), the coordinates of the hand points (beginnings or ends of the bones/phalanges), the vectors of the bones and phalanges, and the angle of the inner side of the hand were read. The values were rounded to 3 decimal places, the images after recording were cropped and the pixel values from a given area averaged so that the output image saved in files had a resolution of $40 \times 40$ pixels.

The dataset has been published in the public domain (on the Kaggle platform)[3]. In addition to the dataset itself, users can read information about creating the dataset and explanations of each feature found in the dataset files.

## 4. Method

The aim of the experimental research is to assess the quality of PJM gesture recognition by two deep network models (CNN and LSTM) using the created PJM dataset while changing the values of specific hyperparameters – the activation function and the number of epochs.

The first model (Figure 1) used to process the hand's coordinate points and the direction vectors of its inner side, consisted of an input layer, an output layer, and two LSTM layers with 128 and 64 neurons, respectively. The second model, used to process the values of the direction vectors of the bones, phalanges and the inner side of the hand, is identical to the first model, but they differ in the size of the input data in the input layer (63 instead of 78). The third model (Figure 2) for image processing (CNN network) has an input layer, an output layer, 2 Conv2D layers, 2 MaxPooling2D layers and a Flatten layer.

The computer implementation of the deep models and the program needed to perform the experiments was carried out using Python version 3.8.0. [6] and the libraries: Numpy [7], Pandas [8], Keras [9] and Scikit-learn [10]. The hyperparameters used to optimize the models' performance through grid search included two activation functions (hyperbolic tangent or ReLU) and three selected epoch values (30, 50, and 70). The recognition quality was measured using the

---

[3]https://www.kaggle.com/datasets/adamlaput/sign-language-pjm.

| InputLayer | | LSTM | | LSTM | | Dense | |
|---|---|---|---|---|---|---|---|
| input: | (3,78) | input: | (3,78) | input: | (3,128) | input: | (64) |
| output: | (3,78) | output: | (3,128) | output: | (64) | output: | (40) |

Figure 1. Schema of the first network model – LSTM for hand points

| InputLayer | | Conv2D | | MaxPooling2D | | Conv2D | | MaxPooling2D | | Flatten | | Dense | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| input: | (3,40,40) | input: | (3,40,40) | input: | (3,40,128) | input: | (2,20,128) | input: | (2,20,64) | input: | (1,10,64) | input: | (640) |
| output: | (3,40,40) | output: | (3,40,128) | output: | (2,20,128) | output: | (2,20,64) | output: | (1,10,64) | output: | (640) | output: | (40) |

Figure 2. Schema of the third network model – CNN for images

following metrics: *balanced accuracy* (*bal*), *precision* (*prec*) and *recall* (*rec*). Depending on the type of data – points, vectors or hand images, the following files from the dataset were used: "PJM-points.csv," "PJM-vectors.csv" and "PJM-images.csv." Each file includes 20,000 samples containing gestures collected from 17 different people. The results of each individual experiment were presented in an average form (relative to five repetitions of the 2-fold cross-validation method) of metric values and mean squared errors (with the suffix "_s"). The validation set constituted 10% of the training set. To perform statistical tests, the Student's t-test was used.

In order to replicate the obtained experimental results, the entire code is available in repository on Github[4].

## 5. Results

Table 1 lists the best combinations of hyperparameters – the activation function (*funct*) and the number of epochs (*epo*) – for each of the three types of data: coordinate points of the hand (*hand points*), direction vectors of the bones (*hand vectors*) and hand images (*hand images*).

The convergence rate of network learning was very high in all cases, as metrics values above 95% were obtained already with the number of epochs equal

---

[4] `https://github.com/adamlaput/RozpoznawanieGestowPJM`.

Table 1. Best combinations of hyperparameters for each of the three types of data with metrics and mean squared error values

| Data type | Funct | Epo | bal | bal_s | prec | prec_s | rec | rec_s |
|---|---|---|---|---|---|---|---|---|
| hand points | tanh | 50 | 0.981 | 0.005 | 0.985 | 0.004 | 0.981 | 0.005 |
| hand vectors | relu | 70 | 0.988 | 0.001 | 0.989 | 0.001 | 0.988 | 0.001 |
| hand images | tanh | 30 | 0.988 | 0.005 | 0.991 | 0.004 | 0.988 | 0.005 |

to 30. For hand points, the models using the *tanh* activation function converged faster in the process learning from the models using the *relu* activation function. The models using hand vectors did not differ significantly from each other in the convergence of the learning process. The *relu* activation function for hand images allowed the metrics to be progressively improved as the number of epochs increased. For *tanh* activation function, the best metric values were achieved for 30 epochs.

The *hyperbolic tangent* activation function turned out to have higher metric values than the *rel function* in the case of operations on hand points and hand photos, but in the case of hand vectors the values were slightly smaller.

## 6. Conclusions

The main goals of this work were to create a dataset of PJM gestures, develop LSTM and CNN network models and assess the quality of gesture recognition by these models. During the work on this paper, a publicly available data set was created, the usefulness of which was verified by deep network models. Experimental studies have shown that models based on LSTM and CNN networks recognize sign language gestures very well, as evidenced by the obtained acceptable average metrics values above 98%. However, using such models in a commercial project would require further testing on a larger data set. We could consider extending the recognition area to include two hands or even their position relative to the body, which would bring greater realism to the recognition process. In future research, it is worth verifying the quality of the created deep network models in the case of recognizing sign language gestures performed by people who did not participate in the creation of the data set.

# References

[1] Łozińska, S. et al. *Linguistics of Space and Movement. Sign Communication and Corpus Methods [Lingwistyka przestrzeni i ruchu. Komunikacja migowa a metody korpusowe].* Wydział Polonistyki Uniwersytetu Warszawskiego, Warszawa, 2014. ISBN 978-83-64111-85-3. S. 25-30.

[2] Fabisiak, S. et al. Polish sign language verbs. In *Proceedings of Verb 2010, Interdisciplinary Workshop on Verbs: The Identification and Representation of Verb Features*, pages 32–37. 2010.

[3] Rao, G. A. et al. Deep convolutional neural networks for sign language recognition. In *2018 Conference on Signal Processing and Communication Engineering Systems (SPACES)*, pages 194–197. 2018. doi:10.1109/SPACES.2018.8316344.

[4] De Coster, M. et al. Sign language recognition with transformer networks. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 6018–6024. European Language Resources Association, Marseille, France, 2020. ISBN 979-10-95546-34-4. URL `https://aclanthology.org/2020.lrec-1.737`.

[5] Pranjali, B. et al. Application of deep learning techniques on sign language recognition—a survey. In *Data Management, Analytics and Innovation*, pages 211–227. Springer Singapore, Singapore, 2021. ISBN 978-981-16-2934-1.

[6] Rossum, V. et al. *Python Reference Manual*. Centrum voor Wiskunde en Informatica Amsterdam, 1995.

[7] Harris, C. R. et al. Array programming with NumPy. *Nature*, 585(7825):357–362, 2020. doi:10.1038/s41586-020-2649-2.

[8] McKinney, W. Data structures for statistical computing in Python. In S. van der Walt and J. Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56–61. 2010. doi:10.25080/Majora-92bf1922-00a.

[9] Chollet, F. Keras. `https://github.com/fchollet/keras`, 2015.

[10] Pedregosa, F. et al. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

# Augmentation of Parking Information in OpenStreetMap Using Aerial Imagery Analysis

**Mateusz Mazur**[0009−0007−3306−0471], **Sebastian Ernst**[0000−0001−8983−480X]

*AGH University of Krakow*
*Department of Applied Computer Science*
*al. A. Mickiewicza 30, 30-059 Kraków, Poland*
*mazurm@student.agh.edu.pl, ernst@agh.edu.pl*

**Abstract.** *OpenStreetMap is a free, editable map of the world, created and maintained by volunteers and available for use under an open license. Harvesting and maintenance of spatial data can be a challenging and time--consuming task, which can be automated using machine learning and existing data sources like aerial images. This paper reports on the creation of such a system – one that augments OpenStreetMap street parking orientation data.*
**Keywords:** *OSM, OpenStreetMap, ML, Machine Learning, transfer learning, aerial images, data augmentation, street parking, parking lots orientation*

## 1. Introduction

Spatial datasets play a vital role in many fields of computer science, including social networks, smart cities or various analytics commonly referred to as *location intelligence*. Integration of such datasets plays a key role here. Let us consider a dataset of street lamp locations, stored as points. Without integrating it with another dataset, containing the shapes and parameters of streets, there is no *context*, i.e. there is no information regarding what is illuminated by each lamp. Such a context often provides knowledge crucial for further design and/or analysis; in this case, parameters of the streets may influence the technical requirements for the lamps.

One possible source of the aforementioned context can be OpenStreetMap[1] (OSM) – an editable map of the world, maintained by its community and distributed under an open license. However, being a general-purpose solution, the OSM dataset does not contain all the attributes necessary for a given specific task. Taking the example of street lighting, these include the road width, the number of lanes, traffic intensity, as well as the presence and location of parked vehicles [1].

Even though the OSM data model allows for storing of most of these parameters, the data coverage may vary, as harvesting and maintenance of spatial data can be a challenging and time-consuming task. As a result, these parameters are rarely present – for instance, for all roads in Poland, the width is specified for less than 1%, any parking data are provided for around 0.24%, and both these parameters are present for only 0.01%. To bridge this gap, one can use machine learning to automate this process by extracting the data from other existing data sources.

The aim of the presented research is to demonstrate such automation, specifically in terms of augmenting[2] roadside parking lots orientation information using aerial images. The rationale for such an approach is two-fold. On the one hand, providing additional road parameters can support the aforementioned specific tasks, such as street lighting planning and design; knowledge about the street parking can be used in navigation applications or when outlining traffic organization changes.

On the other hand, we should once again bring up the collaborative, social character of OpenStreetMap. The map is continuously updated by thousands of volunteers, with changes being provided as delta datasets called *changesets*. Automated edits have to be done with caution and follow a defined code of conduct, but provided the predictions are reliable, they can be a valuable contribution towards making OSM more complete, or at least supporting the community members in locating places where edits are necessary.

Most research in this field is focused either on the image processing-related aspects of aerial vehicle detection (e.g. channel filtering, exclusion of city greens) or establishing whether vehicles are present in predefined locations known to be parking spots (i.e., detecting occupied and unoccupied parking spaces). The presented research tries instead to focus on the semantics of the relationship between each vehicle and surrounding roads, as presented in the following sections.

---

[1] `https://www.openstreetmap.org/`, accessed 2025-01-08.

[2] Please note that the term *augmentation* is used here in its general meaning, as defined by e.g. the Oxford English Dictionary: our aim is to increase the coverage of parking data in OSM. It should not be confused with the definition of data augmentation as used in ML terminology.

## 2. The Data

OpenStreetMap (OSM), as mentioned in Section 1, is a collaboratively-edited map of the world, with data available based on a liberal license, which makes it particularly useful as a provider for *context* in many GIS-related research tasks. Community members can contribute changes using either a web-based editor integrated in the primary browsing interface or by using specialized tools such as JOSM. The entire datasets, or their selected fragments, can then be downloaded and converted into a variety of formats.

Real-world objects and shapes in OSM are represented as `point` and `way` objects, which are used to build geometric shapes: points, linestrings and polygons. The *semantics* of the objects are determined using *tags*, stored as key-value pairs. These, in fact, determine whether a line represents a road, a boundary, the outline of a building or another geographic object. The keys and values are defined in an extensive catalogue, Map Features, available in the OSM wiki[3].

For instance, roads are represented as ways with the `highway` tag, taking various values depending on the type of road (e.g. *primary*, *secondary*, *residential*, etc.). Other tags can be used to indicate e.g. the width of the road, the presence of sidewalks or the rules regarding vehicle parking. These may include both restrictions and parking orientation (parallel, diagonal, perpendicular), defined separately for each side of the road.

One significant aspect of the OSM data model is that the key-value pairs are defined for entire OSM objects. That means that if some parameters, e.g. the speed limit or parking regulations, changes along the course of a street, it needs to be modeled as two separate OSM objects. This aspect may also be significant if we decided to retroactively apply the results of our analyses into the OSM database, as described in Section 6.

As the treated problem is uncommon, we decided to prepare the dataset from scratch. The assembled collection, *Geoparkings*, gathers information about OSM roads with parking orientation mapped for at least one roadside and aerial images picturing them. The orthophotomaps were retrieved from Geoportal[4] using its WCS service. Collected data needed to undergo some wrangling, especially in terms of missing values and meeting current OSM standards.

---

[3]`https://wiki.openstreetmap.org/wiki/Map_feature`, accessed 2025-01-08.

[4]`www.geoportal.gov.pl/pl/usluga/uslugi-pobierania-wcs/`, accessed 2025-01-08.

# 3. Methods

The work consisted of three parts, each focusing on a different, individual aspect. The first stage dedicated to the data preparation has been briefly described in Section 2. The classification process, due to the lack of data and the complex nature of the problem, was divided into two separate modules.

## 3.1. Car Detection

As the output of the detector is used as the input for the proper classifier, it is necessary for the model to provide the information not only about the positions of the cars in the images but also their rotations. Such a problem, called Oriented Bounding Box (OBB) detection, is well known in Computer Vision.

To reduce training time and computing power needed to train a ML model from scratch, an existing solution – the YOLO11 OBB [2] architecture – has been utilized using transfer learning. As the model comes pretrained on a dataset lacking any vehicle instances, and the acquired aerial images lack car annotations, it was necessary to use another collection for the training process. VEDAI [3], the selected one, focuses purely on vehicles, especially on land ones. The training process included the creation and validation of several models and additional quantitative evaluation on the *Geoparkings* images to see the performance on test samples.

## 3.2. Parking Classification

During the actual classification process, the detections are used along with selected OSM way attributes. The approach uses geometric transformations to calculate car features in relation to the analyzed road, specifically to a short section of the polyline representing it in OSM that is located closest to the car in question. The computed attributes include the inclination (of the axis of the car bounding box to the road section), the roadside, the distance (between the center of the car bounding box and the center of the aforementioned section), and several auxiliary ones.

**Naive approach.** The baseline idea is to classify each car by the inclination and, for each roadside, select the most common class. Additionally, to remove outliers, only the objects considered to be related to the road (i.e. with the distance feature smaller than a specified treshold, e.g., 10 meters) are considered during the selection.

**Clustering approach.** The next concept is based on the fact that cars play different roles in road traffic, e.g., driving, in traffic or parked, and individual batches tend to have certain similarities. This solution, to select *parked* cars for each roadside, clusters the detections using DBSCAN [4] and compares them using chosen metric: distance minimization, count maximization or their ratio (count / distance maximization). Finally, the chosen groups are averaged and classified by their inclination. What is more, the distance is calculated not from the centerline, but from the edge of the road (calculated using parallel shift by half of the road width from the polyline representing it). However, as the width feature is hardly mapped in OSM, it is necessary to approximate it using number of lanes and their width, both approximated using other features and present Polish legal regulations, unless present in the database.

# 4. Experiments

Both approaches were evaluated on the *Geoparkings* dataset using the same car detection model. The technique used in the naive classifier turns out to be poor, as it uses too general an approach (visible in Figure 1). Nevertheless, it establishes a concrete benchmark for further analysis with approx. 35%, 20%, and 17% of correct predictions for perpendicular, diagonal, and parallel parking, respectively.



Figure 1. Sample parking orientation classification process. Road diagram: OSM road geometry; the parentheses show the ground truth from the OSM database. Prediction diagrams: road model (light gray with dark gray centerline), vehicles (all belonging to one group have the same shade of blue), orange indicators for the cars that were used for the predictions; the parentheses show predicted values.

The clustering approach, on the other hand, provides rather satisfactory results. During the comparison of the representative cluster selection techniques, mentioned in Section 3.2, the ratio method outperforms both distance and count-based metrics by ~12% as it assures a fair balance between the two. The confusion matrices of the best-performing model are presented in Figure 2.



Figure 2. Confusion matrices of parking orientation classifier using a clustering approach with ratio selection

Additionally, the clustering approach was analyzed in terms of the impact of the number of detected cars on the quality of prediction. The results for subsets containing only samples with more than 0, 3, 5, and 10 detections improved, on average, by 12%, 19%, 27%, and 51%, respectively.

# 5. Conclusions

Taking into consideration the nature of the problem – data extraction – and the fact that the OSM data turned out to have missing and incorrect labels, the proposed solution has achieved rather satisfactory results. Utilizing road width and determining car roles have a significant and positive impact on the model's accuracy.

As the primary goal was to augment the existing data, the proposed solution was used to supplement the parking data for roads that lack it. In such cases, the predictions can be shared as a contribution for the OSM community or used as a supplementation for local copy of the database. In our case, the augmentation concerned the subset of the roads with partial missing data (i.e. with information about only one lane) and led to the competition in nearly 20% of the roads.

However, not mentioning trivial reasons, such as parking lots covered by trees, there is still room for improvement. In the current state, the solution is not well--suited for cases with a small number of cars, and struggles with half-on-curb parkings. The latter one is probably caused by inaccurate road widths and the inhomogeneity of the case itself, which, in fact, seems to be the hardest aspect of the studied issue.

## 6. Further Work

Given the shortcomings of the current solution, its improvement may focus on two main aspects. The first is to provide more accurate information about road width, e.g., from the images themselves, the second – to formulate the rules for half-on-curb parking. Also, other classification approaches could be developed.

However, it should be noted that even now, the model is able to at least suggest changes, e.g. to fill the spots with missing data, which would in fact allow for automatic augmentation of OSM data. This would help the OSM community in improve the map, especially in places with poor coverage. Also, this approach can easily be extended to other areas of OSM, e.g. regarding urban greenery.

## References

[1] Peña-García, A., Castillo-Martínez, A., and Ernst, S. The basic process of lighting as key factor in the transition towards more sustainable urban environments. *Sustainability*, 16(10), 2024. ISSN 2071-1050. doi:10.3390/su16104028. URL `https://www.mdpi.com/2071-1050/16/10/4028`.

[2] Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. You only look once: Unified, real-time object detection. *arXiv preprint arXiv:1506.02640*, 2016.

[3] Razakarivony, S. and Jurie, F. Vehicle detection in aerial imagery: A small target detection benchmark. *Journal of Visual Communication and Image Representation*, 34:187–203, 2016.

[4] Ester, M., Kriegel, H.-P., Sander, J., and Xu, X. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, KDD'96, page 226–231. AAAI Press, 1996.

# Understanding Machine Unlearning with SHAP Values

**Maciej Mączyński**[0009−0006−0354−4507]

*Wrocław University of Science and Technology*
*Faculty of Information and Communication Technology*
*wyb. Wyspiańskiego 27, 50-370 Wrocław, Poland*
*maciej.maczynski@pwr.edu.pl*

**Abstract.** *Machine unlearning refers to the process of selectively removing specific information and its influence on an already retrained model. As this technique develops, many questions have arisen aimed at understanding the unlearning, including what appears to be the most important one: how to determine that the unlearning process was actually carried out correctly. In this paper we propose a novel approach to validate the results of machine unlearning with the usage of one of explainable AI methods: SHAP. In a simple experiment, we demonstrate that deleted data have a significant impact on the explainability values, which correspond to the differences in the model before and after unlearning.*

**Keywords:** *artificial intelligence, machine unlearning, explainable artificial intelligence, SHAP*

## 1. Introduction

The machine unlearning concept is seemingly straightforward. Despite this, we still find ourself with a lack of knowledge to properly explain how these operation impacted the unlearned model. Deletion of singular data sample from the training set should have zero impact on the retrained model. Forgetting a large partition of data, such as 20% of the entire dataset, will definitely affect the model; however, it is hard to determine how every batch of samples contributed to that process. To solve this issue we propose the employment of existing explainable AI techniques that are already used to explain predictions of machine learning model. The main contributions of this work are as follows:

- Proposing explainable machine learning as one of the possibilities to evaluate the changes between unlearned and retrained model.

- Showing XAI can be used as a tool for evaluation of unlearning methods.

## 2. Machine Unlearning Metrics

There are as many machine unlearning algorithms as there are metrics to evaluate them [1, 2]. Most of them compare the model after forgetting (the unlearned model) with the fully retrained model on the remaining data (the naive retrained model). One of the proposed classifications of such methods [3], called EEEC, divided them into four categories:

**Efficiency** is the relative speed-up that gives unlearned metric over the fully retrained model. This is usually achieved [3, 4] by measuring and comparing time $T$ taken to obtain the unlearned model $h^U$ and the retrained model $h^R$ as presented in the following equation:

$$Efficiency(h^U) = \frac{T(h^U)}{T(h^R)}. \tag{1}$$

**Effectiveness** is the comparison of the performance of the test set between the unlearned model and the fully retrained model. This value can be counted with the use of any performance metric $m$ such as accuracy, recall, F1-score, etc. In our work, this comparison on test set $D^T$ was counted in the following manner:

$$Effectivness(h^U) = \frac{m(h^U, D^T)}{m(h^R, D^T)}. \tag{2}$$

**Consistency** is a measure of similarity between the unlearned model and the fully retrained model. This can be counted as a distance [3] between model parameters, among other methods, with the use of Kullback-Leibler divergence. In this work we calculate mean layer wise distance with L2-norm [1]. When comparing the weights of models $w$ with $L$ layers, the equation will take the following form:

$$Consistency(h^U) = \frac{1}{L} \sum_{i=1}^{L} \|w_U^i - w_R^i\|_2. \tag{3}$$

**Certifiability** is a measure of the "forgetness" of data in the forget set between the unlearned model and the fully retrained model. In our work, we assess it by comparing a performance metric *m* on the dataset containing the forgotten data, $D_f$:

$$Certifiability(h^U) = \frac{m(h^U, D^F)}{m(h^R, D^F)}.$$ (4)

The presented group of methods takes for granted that the model trained without forgotten data is free from their influence. In many scenarios, the computational cost of full retraining could be to high, or access to the complete dataset may be unavailable [5]. Gaining insights into how the removed information influenced the model by analysing how explanations of model predictions evolve before and after unlearning specific data could benefit future work on machine unlearning.

## 3. SHAP Comparison

SHAP (SHapley Additive exPlanations) [6] is a method used to explain a complex machine learning evaluation, in our case ResNet-18 [7]. The key idea behind the concept is to give each input feature a specific value $\psi$ based on its contribution in the model prediction. Based on the SHAP explanation model, a prediction for class k $f_k(x)$ can be defined as:

$$f_k(x) = \psi_0^{(k)} + \sum_{i=0}^{M} \psi_i^{(k)}$$ (5)

where $\psi_0^{(k)}$ is an expected value defined as a mean value of all predictions in this class in the presented dataset, $M$ is a number of all features (pixels) in a single sample and $\psi_i^{(k)}$ is a SHAP value of feature *i* for the selected class.

Taking this into account, we propose a metric that could be used to compare the unlearned and the retrained model: mean comparison between SHAP values in the expected prediction class ($\phi_T^M$).

**Mean comparison between SHAP values in true-labeled predictions.** Comparing differences in SHAP explanation for each individual pixels between the original and unlearned models could give us information on how deleted data have impacted the model. Before conducting a further study on that subject, we must prove that there is a connection between the changes in the sum of SHAP values for correctly labeled predictions between the unlearned and the retrained models, similar to the changes observed in existing effectiveness metrics. In an ideal scenario,

276

we aim to achieve the greatest positive impact from pixels, boosting predictions for ground-true class and reducing predictions for every other class. For $N$ samples, $K$ number of class labels and ground true label list $L$, we calculate it as:

$$\phi_T^M = \frac{\frac{\sum_{j=0}^N \sum_{i=0}^M \psi_i^{U(L_j)}}{\sum_{j=0}^N \sum_{i=0}^M \sum_{k=0}^K \psi_i^{U(k)}}}{\frac{\sum_{j=0}^N \sum_{i=0}^M \psi_i^{R(L_j)}}{\sum_{j=0}^N \sum_{i=0}^M \sum_{k=0}^K \psi_i^{R(k)}}}. \tag{6}$$

# 4. Experiment Evaluation

## 4.1. Experimental Setup

**Unlearning Algorithms.** All selected machine unlearning algorithms were introduced during "NeurIPS 2023 Machine Unlearning competition" [8]. The detailed description of this method is available online on Kaggle [9].

**Datasets.** In this work, we consider only the image datasets based on strict connection between the tournament methods and the ResNet architecture [7]. We use CIFAR-10 [10] and MNIST [11] datasets split in a train set, retrain set (90% of the random split training dataset with maintained balance), forget set (remaining 10%), validation set and test set (1,000 samples, 100 per each class).

**Evaluation procedure.** To compare the selected metrics, we first trained the original ResNet-18 model on the full training dataset. The trained model served as input for five unlearning algorithms: Fauchan (Fau), Kookmin (Koo), Seif (Sei), Sebastian (Seb), and Amnesiacs (Amn), using the retraining set, forget set, and validation set. After applying the five unlearning algorithms, we retrained the ResNet-18 model using the retraining dataset. Both the unlearned and retrained models were then evaluated using the SHAP GradientExplainer to compute SHAP values, utilizing the testing dataset. The resulting SHAP values, denoted $\phi_T^M$, were compared with the effectiveness ($eff$), consistency ($con$), certifiability ($cer$) values, using accuracy as performance metric.

All the presented experiments with brief descriptions of unlearned methods are available in Github repository[1].

---

[1] `https://github.com/Wakir/PPRAI-UNXAI`.

## 4.2. The Results

The results of the experiment, presented in Table 1, show that the proposed metric classified the unlearning algorithms in a way that is somewhat consistent with other unlearning values. All the models were closely related (94%–99% in MNIST, 70%–77% in CIFAR-10), with the exception of Seif in the CIFAR-10 dataset that obtained the 52.3% average accuracy value in the test dataset. This significant difference, in our opinion, was caused by adding excessive Gaussian noise to the convolution weights, leading to catastrophic forgetting. This was distinguished by $\phi_T^M$; however, the gap between this and the second-lowest value on Kookmin was smaller than the gap between Fauchan and Kookmin. On the other hand, the differences in results between Sebastian and Amnesiacs were similar for both the unlearning metrics and the comparative SHAP values. The similarity was not linear and the values were not identical; nevertheless the ranking order was nearly the same as in the presented metric. Further studies are needed to take a better look at these relationships.

Table 1. Average unlearning metric values and comparison of SHAP values on CIFAR10 and MNIST dataset

|  | Fau | Koo | Sei | Seb | Amn |
|---|---|---|---|---|---|
| *eff* MNIST | 1.006 | 1.00 | 0.949 | 1.000 | 1.001 |
| *cer* MNIST | 1.010 | 0.993 | 0.941 | 0.991 | 0.982 |
| *con* MNIST | 9.478 | 9.125 | 9.654 | 7.801 | 7.799 |
| $\phi_T^M$ MNIST | 1.248 | 0.921 | 0.758 | 0.815 | 0.852 |
| *eff* CIFAR10 | 1.004 | 0.925 | 0.595 | 0.931 | 0.945 |
| *cer* CIFAR10 | 1.322 | 1.167 | 0.641 | 0.885 | 0.891 |
| *con* CIFAR10 | 13.658 | 14.329 | 15.648 | 12.123 | 12.124 |
| $\phi_T^M$ CIFAR10 | 1.406 | 0.979 | 0.848 | 1.140 | 1.130 |

Comparing ground-true SHAP values for 10 selected image – sample in Figure 1 – we can see that the number of active pixels in the unlearned examples is in most of the cases smaller than on the unlearned and retrained models. It is mostly visible on Fauchan architecture, which when compared with $\phi_T^M$ and accuracy showed that the model focuses on a smaller number of pixels, which improved the accuracy score. On the contrary, Seif on MNIST dataset improved the impact on prediction for most of the models, which lead to the worst prediction values

(a) $sh_A^1$          (b) $sh_B^1$

Figure 1. SHAP Image plot for ground-true class and 7 SHAP explainers per dataset, made on 10 selected images in CIFAR10 and MNIST. Presented explainers are in order: orginal, retrained, Fauchan, Kookmin, Seif, Sebastian, Amnesiacs.

from all unlearning algorithms. Reducing the number of impactful pixels not always lead to a better prediction, which can be seen on Seif model for CIFAR-10. Each unlearned model reacted differently to the deletion of selected data, which affected the model architecture to varying degrees; however, catastrophic forgetting did not occur. Further studies are still necessary in order to prove if we can use SHAP as a proof that some information was forgotten.

## 5. Conclusions

We proposed an explainable artificial intelligence as a tool to explain the impact of deleted data for the model influence and evaluate their performance. The experiments suggested the correlation between changes of accuracies and explainability; however, we still cannot be sure how the deletion of batches affected specific pixels or features. Future studies may include going deeper into the problem of explainability in machine unlearning [2] with the usage of other existing explainability algorithms on less complex model architectures.

## References

[1] Nguyen, T. T., Huynh, T. T., Nguyen, P. L., Liew, A. W.-C., Yin, H., and Nguyen, Q. V. H. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299*, 2022.

[2] Xu, J., Wu, Z., Wang, C., and Jia, X. Machine unlearning: Solutions and challenges. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 8:2150–2168, 2024.

[3] Mercuri, S., Khraishi, R., Okhrati, R., Batra, D., Hamill, C., Ghasempour, T., and Nowlan, A. An introduction to machine unlearning. *arXiv preprint arXiv:2209.00939*, 2022.

[4] Yinzhi Cao, J. Y. Towards making systems forget with machine unlearning. *IEEE Symposium on Security and Privacy*, pages 163–480, 2015.

[5] Chundawat, V. S., Tarun, A. K., Mandal, M., and Kankanhalli, M. S. Zero-shot machine unlearning. *IEEE Transactions on Information Forensics and Security*, 18:2345–2354, 2023. doi:10.1109/TIFS.2023.3265506.

[6] Lundberg, S. M. and Lee, S.-I. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874v2*, 2017.

[7] He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *CoRR*, 2015. doi:abs/1512.03385.

[8] Triantafillou, E., Kairouz, P., et al. Are we making progress in unlearning? Findings from the first NeurIPS unlearning competition. *arXiv preprint arXiv:2406.09073v1*, 2024.

[9] Triantafillou, E., Pedregosa, F., et al. NeurIPS 2023 – machine unlearning, 2023. URL `https://kaggle.com/competitions/neurips-2023-machine-unlearning`.

[10] Krizhevsky, A. Learning multiple layers of features from tiny images. *Technical report*, 2009.

[11] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

# Employing Transfer Learning for Diabetic Retinopathy Using Deep Feature Extraction

**Paweł Niedziółka**[0009−0006−8189−5273]

*Wrocław University of Science and Technology*
*Faculty of Information and Communication Technology*
*Wybrzeże Stanisława Wyspiańskiego 27, 50-370 Wrocław, Poland*
*pawel.niedziolka@pwr.edu.pl*

**Abstract.** *The usage of deep models in ophthalmology is increasing significantly. The foundation models, such as RETFound, allow the use of transfer, which should improve the models' prediction. In this study, the author investigated how models trained on ImageNet transfer the knowledge to the retinopathy problem, how the prediction quality for extracted features changes, and from which layer feature extraction performs best. The study was conducted on three public datasets. Findings indicate that the transrate transferability metric correlates with classification quality and that leveraging models pre-trained on ImageNet positively impacts performance in this field. Selecting the right layer for extraction is also crucial.*
**Keywords:** *deep learning, transfer learning, feature extraction*

## 1. Introduction

The recent increase in the efficacy of AI technology over the years has been significant in healthcare. The foundation models are primarily based on deep learning techniques and often leverage transfer learning, allowing them to adapt to new tasks efficiently. However, these models come with clear risks, such as amplifying pre-existing biases in medical datasets or posing privacy concerns due to memorizing training data [1]. In the medical domain, deep models are used for disease diagnosis, individualized treatment planning, and an analysis of disease progression [2]. Transferring knowledge from models that were trained on large amounts of ophthalmic data often gives a significant boost to the new tasks: even

a few hundred images can considerably improve classification or prediction performance [3].

The main contributions of this work are the following: (i) Evaluation of how the classification quality behaves compared to the positive transfer metric for optical domain data, and (ii) examining how layer depth for model extraction affects the transferability metric.

# 2. Related Works

This section describes deep learning for diabetic retinopathy disease, the RETFound model, and the transfer learning concept.

**Deep learning for Diabetic Retinopathy**: In ophthalmology, deep learning models, such as convolutional neural networks and transformer-based architectures, have shown a remarkable success in detecting and classifying retinal diseases [4]. These models are highly effective in recognizing patterns that may go unnoticed by human observers, enhancing diagnostic precision and minimizing manual effort. Researchers are developing data augmentation techniques and designing advanced deep learning models such as RETFound to improve diagnostic accuracy [5].

**Introduction into RETFound:** The model was developed to provide a universal solution to improve model performance and reduce expert annotation workload. The model was trained on 1.6 million unlabeled retinal images using self-supervised method. Then, it was adapted to labeled disease detection tasks. From an architecture perspective, it is based on the Masked Auto Encoder and consists of multiple transformer encoder blocks. Each block includes layer normalization, multi-head self-attention, and a feed-forward multilayer perceptron with GELU activation and dropout for regularization. The related paper describes several ways of learning using transfer learning [6], but it is unclear whether it is always reasonable. Therefore, a research question was proposed as to whether knowledge transfer pathways from the literature positively affected the final model prediction.

**Transfer learning:** Knowledge from another task can improve performance on similar problems through transfer learning using various techniques [7]. However, utilizing an unsuitable source model can reduce accuracy, making the model perform worse than one trained directly on the target task data. This phenomenon is called negative transfer and can be mitigated using methods such as safe transfer, domain similarity estimation, or negative transfer mitigation [8]. On the other hand, the positive effect of such a transfer can be measurable – metrics such as *H-score*, or the one explored in that article, *transrate* metric [9], can be used.

# 3. Experimental Evaluation

The analysis of the referenced work prompts an examination of whether the models benefited from knowledge transfer and how it affects classification quality. This led to the proposal of the following research questions:

**RQ1:** Does supervised or unsupervised pre-training of deep neural networks on ImageNet dataset lead to positive knowledge transfer for diabetic rethinopathy diagnosis task?

**RQ2:** Which layers of selected pre-trained deep network architectures offer the best positive knowledge transfer for diabetic rethinopathy diagnosis task?

## 3.1. Setup

**Datasets:** The following three datasets were used: (i) *JSIEC* – contains 39 classes (diseases and conditions) and includes 1,000 images [10], (ii) *GLAUCOMA* – Optical Coherent Tomography (OCT) photos for standard control, early and advanced glaucoma [11], and (iii) *APTOS2019* – used for detecting diabetic retinopathy, contains five categories (degree of diabetic retinopathy) and almost 5,600 images.

**Reference methods:** The examined models are Resnet50 (RN50) [12], Vision Transformer using masked auto-encoding (MAE) [13], and RETFound (RETF). It is important to note that models can have different weight configurations: those without pre-training $_{NW}$, pre-trained on ImageNet, and fine-tuned on retinopathy images; color fundus images $_{CFP}$, optical coherence tomography $_{OCT}$.

**Experimental Protocol:** For a fair and reliable evaluation of the experiments, as the experimental results are not deterministic, a $5 \times 2$ stratified cross-validation protocol was used, and the results were averaged. The Balanced Accuracy Score (BAC) was used for classifiers, and a statistical analysis with a *corrected* Student's t-test ($p \leq 0.05$) [14] was performed.

**Reproducibility:** All models and datasets are publicly available, the experimental code (Python) is stored on the GitHub[1]. The random state option was applied to make the results more deterministic among cross-validation, PCA, and SVC classifier. The research was carried out in an environment with appropriate versions: Python 3.11.9, Torch 2.3.0, and Torchvision 0.18.0a0.

---

[1]`https://github.com/KoEj/Transfer-learning-in-ophthalmology`.

## 3.2. Experiment 1 – Transfer Model/Data

This experiment validates the positive transfer of ImageNet-trained deep neural networks to ophthalmology. Two metrics were examined to measure positive transfer, *transrate*, and *H-score*. The extraction of features was performed from the head layer, but additionally, for *transrate*, the features were centralized and normalized as required. In the initial experimental phase, these metrics were compared, and it was observed that *transrate* had a more momentous value variation. The higher *transrate* value, the better correlation with transfer learning performance [9]. Therefore, this metric was chosen to evaluate the models better and to more effectively demonstrate the positive transfer relationship.

The positive transfer metric was evaluated across the models and datasets (Section 3.1). A *corrected* Student's t-test [14] was performed, which addressed bias and variance issues arising from training fold overlap in cross-validation. Results are presented in Table 1; small numbers below *transrate* values indicate the indexes of reference models that are significantly worse than the presented model.

Table 1. *Transrate* metric for selected models and datasets with statistical analysis

| Dataset | $\text{RN50}_{NW}^{1}$ | $\text{RN50}^{2}$ | $\text{MAE}_{NW}^{3}$ | $\text{MAE}^{4}$ | $\text{RETF}_{OCT}^{5}$ | $\text{RETF}_{CFP}^{6}$ |
|---|---|---|---|---|---|---|
| JSIEC | 8.483 <br> – | 24.895 <br> 1 3 4 5 | 9.032 <br> 1 | 21.593 <br> 1 3 5 | 14.032 <br> 1 3 | $\underline{25.502}$ <br> 1 3 4 5 |
| GLAUCOMA | 2.100 <br> – | **12.431** <br> ALL | 2.662 <br> 1 | 9.994 <br> 1 3 5 | 5.937 <br> 1 3 | 11.522 <br> 1 3 4 5 |
| APTOS2019 | 3.443 <br> – | **18.096** <br> ALL | 4.159 <br> 1 | 11.571 <br> 1 3 5 | 8.422 <br> 1 3 | 15.316 <br> 1 3 4 5 |

Additional research was conducted to investigate further the relationship between the transferability metric and classification precision. Before the training process, PCA was used to select a subset of features with the most minor possible loss of information. Gaussian Naive Bayes (GNB), K Nearest Neighbors (KNN), and Support Vector Machines (SVC) classifiers were trained on dimensionality-reduced (PCA) features. The classifiers were then evaluated on the test set in order to assess their performance.

The results of the additional experiment are presented in Table 2. **Bolded** values indicate the best statistically significant results, while underlined values denote the highest-performing results that do not reach statistical significance.

In most cases, there is a close correlation between the positive learning metric and BAC. For GBN and KNN classifiers, the results from a statistical perspective are significantly better for the $\text{RETF}_{CFP}$ model, although the *transrate* metric was better for RN50. For SVC, it can be asserted that this evaluation relates exclusively to JSIEC. However, the results are still not worse than those of any other model.

Table 2. Results of statistical analysis based on BAC

| | Dataset | $RN50^1_{NW}$ | $RN50^2$ | $MAE^3_{NW}$ | $MAE^4$ | $RETF^5_{OCT}$ | $RETF^6_{CFP}$ |
|---|---|---|---|---|---|---|---|
| **GNB** | JSIEC | 0.090<br>5 | 0.362<br>1 3 5 | 0.066<br>– | 0.328<br>1 3 5 | 0.043<br>– | **0.505**<br>*ALL* |
| | GLAUCOMA | 0.564<br>3 | 0.603<br>3 | 0.447<br>– | 0.650<br>1 3 5 | 0.541<br>3 | **0.720**<br>*ALL* |
| | APTOS2019 | 0.253<br>3 5 | <u>0.461</u><br>1 3 4 5 | 0.200<br>– | 0.356<br>1 3 5 | 0.200<br>– | 0.399<br>1 3 4 5 |
| | *mean ranks* | *4.00* | *2.00* | *5.50* | *2.67* | *5.50* | *1.33* |
| **KNN** | JSIEC | 0.086<br>– | 0.283<br>1 3 5 | 0.098<br>1 3 5 | 0.308<br>1 3 5 | 0.052<br>– | **0.468**<br>*ALL* |
| | GLAUCOMA | 0.595<br>3 | 0.596<br>3 | 0.419<br>– | 0.648<br>1 3 5 | 0.552<br>3 | **0.735**<br>*ALL* |
| | APTOS2019 | 0.315<br>3 5 | 0.432<br>1 3 5 | 0.247<br>– | 0.404<br>1 3 5 | 0.222<br>– | <u>0.464</u><br>1 3 4 5 |
| | *mean ranks* | *4.33* | *2.67* | *5.00* | *2.33* | *5.67* | *1.00* |
| **SVC** | JSIEC | 0.046<br>– | 0.381<br>1 3 4 5 | 0.056<br>5 | 0.217<br>1 3 5 | 0.033<br>– | **0.589**<br>*ALL* |
| | GLAUCOMA | 0.568<br>3 | 0.595<br>3 | 0.444<br>– | 0.653<br>1 3 5 | 0.546<br>3 | <u>0.717</u><br>1 2 3 5 |
| | APTOS | 0.200<br>– | <u>0.397</u><br>1 3 5 | 0.200<br>– | 0.371<br>1 3 5 | 0.200<br>– | 0.371<br>1 3 5 |
| | *mean ranks* | *4.33* | *2.00* | *5.17* | *2.50* | *5.50* | *1.50* |

## 3.3. Experiment 2 – Layer Selection

The result of the previous experiment implies that transfer learning from general models may be less effective when a domain-specific model is available. This experiment investigated the best-performing RETFound model with CFP weights based on this conclusion. Two types of layers were used for the feature extraction: second normalization and linear multi-head self-attention layers. The blocks define how deep the extraction process is executed – the lower block, the more low-level features are extracted, while higher blocks correspond to the extraction of more abstract and high-level features. This feature extraction process was performed on the JSIEC dataset – the results are presented in Figure 1.
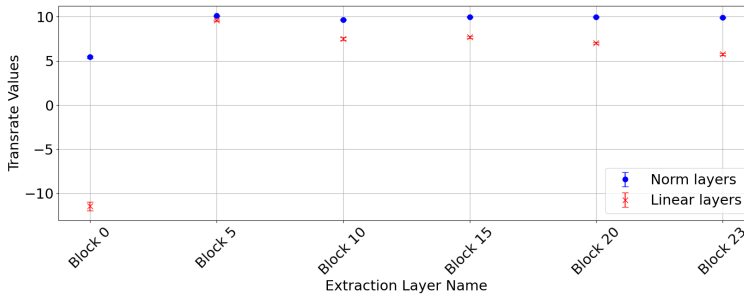


Figure 1. Transrate values across different extraction layers

The *Norm* layers show greater consistency and effectiveness in transferring relevant features, while *Linear* ones look less stable, especially in the early layers of the model. This speaks in favor of possible knowledge transfer advantages that operate with *Norm* layers concerning the RETFound model. Deep layers may capture more abstract and transferable features crucial for effective model adaptation (which is shown for the *Linear* layers), but additional research is required.

## 4. Conclusion

The presented study has shown that transfer correlates with BAC, although this is not a rule. Despite a high score on the *transrate* metric, classification quality can be similar or lower. For the retinopathy problem, it has been verified that the RETFound model with $_{CFP}$ weights performed best. However, it can be seen that the model trained on ImageNet improves the quality of classification in this problem – this can be seen for the ResNet. Higher positive transfer metrics may be observed in the deeper layers of the model rather than in the final layer.

Future research includes training classifiers on features extracted from deeper layers, evaluating their performance, and exploring other model layers.

## References

[1] Bommasani, R. et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

[2] Kaul, D. et al. Deep learning in healthcare. *Deep Learning in Data Analytics: Recent Techniques, Practices and Applications*, pages 97–115, 2022.

[3] Liu, T. A. et al. Deep learning and transfer learning for optic disc laterality detection: implications for machine learning in neuro-ophthalmology. *Journal of Neuro-Ophthalmology*, 40(2):178–184, 2020.

[4] Alyoubi, W. et al. Diabetic retinopathy detection through deep learning techniques: A review. *Informatics in Medicine Unlocked*, 20:100377, 2020.

[5] Abushawish, I. Y. et al. Deep learning in automatic diabetic retinopathy detection and grading systems: a comprehensive survey and comparison of methods. *IEEE Access*, 12:84785–84802, 2024.

[6] Zhou, Y. et al. A foundation model for generalizable disease detection from retinal images. *Nature*, 622(7981):156–163, 2023.

[7] Zhuang, F. et al. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2021. doi:10.1109/JPROC.2020.3004555.

[8] Zhang, W. et al. A survey on negative transfer. *IEEE/CAA Journal of Automatica Sinica*, 10(2):305–329, 2023. doi:10.1109/JAS.2022.106004.

[9] Huang, L.-K. et al. Frustratingly easy transferability estimation. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 9201–9225. 2022.

[10] Cen, L.-P. et al. Automatic detection of 39 fundus diseases and conditions in retinal photographs using deep neural networks. *Nature Communications*, 12(1):4828, 2021.

[11] Kim, U. Machine learn for glaucoma, 2018. doi:10.7910/DVN/1YRRAC.

[12] He, K. et al. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.

[13] He, K. et al. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.

[14] Stapor, K. et al. How to design the fair experimental classifier evaluation. *Applied Soft Computing*, 104:107219, 2021.

# RSSI-Based Device-Free Human Activity Recognition Using Machine Learning

**Rafał Pasternak**[0009−0002−7070−5084], **Bartłomiej Płaczek**[0000−0001−8570−0361]

*University of Silesia in Katowice*
*Faculty of Science and Technology, Institute of Computer Science*
*Będzińska 39, 41-200 Sosnowiec, Poland*
*rafal.pasternak@us.edu.pl, bartlomiej.placzek@us.edu.pl*

**Abstract.** *This study presents a human activity recognition method based on RSSI measurements in a Wi-Fi network. During the experiments, a Wi-Fi access point and a mobile device were used to collect RSSI data indoors. Various machine learning models were trained using statistical features extracted from RSSI variations to recognize three activity states: no person present, a person sitting, and a person walking. The results demonstrate that machine learning classifiers, particularly ensemble methods such as Random Forest and Gradient Boosting, accurately recognize human activities.*
**Keywords:** *machine learning, wireless networks, human activity recognition, smart building, signal strength*

## 1. Introduction

Human activity recognition (HAR) is gaining attention for its applications in smart buildings, healthcare, security, and human-computer interaction. Traditional HAR systems use wearable sensors, which can be intrusive. An alternative method uses wireless networks to detect activities based on received signal strength indicator (RSSI) [1] variations caused by a person's presence and motion.

This study explores RSSI-based device-free HAR using machine learning. A Wi-Fi access point and a mobile device collect RSSI data, which is used to train machine learning models for activity classification. Experiments show that these models can successfully distinguish between different activities based on RSSI data.

The paper is organized as follows: Section 2 covers the experimental setup and data collection. Section 3 presents results and discusses the effectiveness of machine learning models. Section 4 concludes the study and suggests future research directions.

# 2. Experiments

The experiments were designed to evaluate human activity detection based on RSSI measurements in a Wi-Fi network. The proposed method employs a mobile device (smartphone) as the primary data collection unit and a Wi-Fi access point as the broadcasting device. Both devices were positioned at known stationary locations indoors and operated in the 5 GHz frequency band.

The experimental setup was implemented in a room with an area of 30 square meters. The smartphone, which ran a measurement application, was placed in one corner of the room, while the Wi-Fi access point was located in the opposite corner. The mobile device continuously collected RSSI samples. To analyze the impact of network configuration on human activity detection, the experiments were conducted using multiple Wi-Fi broadcasting channels and different channel widths.

## 2.1. RSSI Dataset

The experimental dataset consists of RSSI measurements recorded under different conditions, as described above. Each experiment was performed using a specific combination of Wi-Fi channel, channel width, and person status in the room.

Data were collected in comma-separated values (CSV) format during the experiment. The dataset includes the following attributes: broadcasting channel (36, 38, or 40), channel width (20 Hz, 40 Hz, or 80 Hz), person's activity in the room, which represents one of the three previously defined scenarios (no person, person sitting, or person walking), signal strength of the Wi-Fi access point (from 3 dB to 19 dB), and the measurement results (a sequence of 30 RSSI values recorded at one-second intervals). A total of 459 experiments were conducted, each generating 30 RSSI measurements.

For each experiment, the average, minimum, maximum, and standard deviation were calculated based on the 30 recorded RSSI values. These statistical features were used as inputs for various machine learning classifiers to recognize

human activity. Using these features, the robustness of activity detection was improved, and the impact of short-term signal fluctuations was reduced.

## 2.2. Results and Discussion

Principal Component Analysis (PCA) [2] was conducted on the raw RSSI data to emphasize variability and reveal underlying patterns within the dataset. The results, presented in Figure 1a, show a scatter plot of the first two principal components, where three distinct clusters of data points are observed. Each cluster is visually separated and represented in different colors: red, green, and blue. These results indicate that the dataset contains an inherent structure that can be effectively captured by the two primary principal components. The clear separation of clusters suggests that the RSSI readings contain meaningful variations corresponding to different experimental conditions. This observation supports the hypothesis that a relationship exists between the three considered human activity scenarios and the measured RSSI values.

a)

b)

Figure 1. PCA and Average signal strength categorized by human activity

A further analysis was performed by creating a box plot of the average RSSI values for the three human activity scenarios (Figure 1b). The average RSSI was calculated for each experiment based on 30 raw RSSI readings, providing a summarized representation of signal variations. Based on Figure 1b, it can be seen that the average RSSI is influenced by the presence and movement of a person in the room. A significant difference is observed when no person is present, with the average RSSI noticeably higher than in the other scenarios. In addition, a difference can be seen between the cases where one person is sitting or walking in the room. In these situations, a smaller range is observed between the minimum and maximum values relative to the average RSSI, suggesting that human presence and

movement affect the stability and dispersion of the signal. These findings further support the hypothesis that RSSI variations can be used to distinguish between different states of human activity (Figure 2).



Figure 2. Accuracy of human activity recognition

To recognize human activity based on RSSI measurements, multiple machine learning models were trained and evaluated. The task was to recognize one of the three activity states: no person in the room, a person walking, or a person sitting. As mentioned earlier, the data used to train the classifiers were derived from the raw RSSI readings by computing statistical features, including the average, standard deviation, minimum, and maximum values for each 30-second measurement period. These features served as input for the classification models.

The classifiers tested in this study included Naïve Bayes [3], Generalized Linear Model (GLM) [4], Logistic Regression [5], Fast Large Margin [6], Deep Learning H2O [7], Decision Tree [8], Random Forest [9], Gradient Boosted Decision Trees (GBDT) [10], and Support Vector Machine (SVM) [11]. Each classifier was trained and tested using a standard machine learning workflow, including data preprocessing, model training, and performance evaluation.

The classification performance was assessed using two key metrics: accuracy and average F-measure [12], which are presented in the column plot in Figure 2. Among the analyzed classifiers, Random Forest and Gradient Boosting demonstrated the highest performance, achieving an accuracy of 88%. These models outperformed the others, making them the most effective tools for recognizing human activity in this study. The results indicate that ensemble methods, such as Random Forest and GBDT, are well-suited for capturing the underlying patterns in RSSI-based HAR.

Experiments identified the minimum value and standard deviation of RSSI as key features for distinguishing activity states. Network parameters like channel and signal strength did not improve recognition accuracy. Machine learning models achieved about 90% accuracy in detecting a person's presence based on RSSI. A simple RSSI-based detection algorithm was developed (Algorithm 1) for IoT devices, useful for smart home automation, energy management, and security.

---

### Algorithm 1. RSSI-based person detection

---

**if** *rssi standard deviation* ≤ 3.17 dBm **or** *rssi minimal value* ≥ −54 dBm **then**
    No person in the room
**else if** *rssi standard deviation* ≥ 3.6 dBm **then**
    There is a person in the room
**else**
    No person in the room
**end if**

---

## 3. Conclusions

It was demonstrated that machine learning models can recognize human activities with high accuracy by analyzing Wi-Fi signal strength variations, mainly when using ensemble methods. The results indicate that statistical features such as the minimum RSSI and the standard deviation are helpful in the classification of activities. Additionally, network configuration parameters, such as channel and bandwidth, had minimal impact on classification accuracy. Future research directions include expanding the dataset by incorporating additional activity types and testing the method in diverse indoor environments with consideration of detailed channel state information for the Wi-Fi network.

## References

[1] Wu, R.-H., Lee, Y.-H., Tseng, H.-W., Jan, Y.-G., and Chuang, M.-H. Study of characteristics of RSSI signal. In *2008 IEEE International Conference on Industrial Technology*, pages 1–3. 2008. doi:10.1109/ICIT.2008.4608603.

[2] Dafferstofer, A. PCA in studying coordination and variability: A tutorial. *Clinical Biomechanics*, 19:415–428, 2004.

[3] Yang, F.-J. An implementation of Naive Bayes classifier. In *2018 International Conference on Computational Science and Computational Intelligence (CSCI)*, pages 301–306. IEEE, 2018.

[4] Kiebel, S. and Holmes, A. *The General Linear Model*, volume 8. 2007.

[5] Kleinbaum, D. G., Dietz, K., Gail, M., Klein, M., and Klein, M. *Logistic Regression*. Springer, 2002.

[6] Wang, B. and Zou, H. Fast and exact leave-one-out analysis of large-margin classifiers. *Technometrics*, 64(3):291–298, 2022.

[7] Candel, A., Parmar, V., LeDell, E., and Arora, A. Deep learning with H2O. *H2O. AI Inc*, pages 1–21, 2016.

[8] Suthaharan, S. and Suthaharan, S. Decision tree learning. *Machine Learning Models and Algorithms for Big Data Classification: Thinking with Examples for Effective Learning*, pages 237–269, 2016.

[9] Rigatti, S. J. Random forest. *Journal of Insurance Medicine*, 47(1):31–39, 2017.

[10] Si, S., Zhang, H., Keerthi, S. S., Mahajan, D., Dhillon, I. S., and Hsieh, C.-J. Gradient boosted decision trees for high dimensional sparse output. In *International Conference on Machine Learning*, pages 3182–3190. PMLR, 2017.

[11] Mammone, A., Turchi, M., and Cristianini, N. Support vector machines. *Wiley Interdisciplinary Reviews: Computational Statistics*, 1(3):283–289, 2009.

[12] Liu, Y., Zhou, Y., Wen, S., and Tang, C. A strategy on selecting performance metrics for classifier evaluation. *International Journal of Mobile Computing and Multimedia Communications (IJMCMC)*, 6(4):20–35, 2014.

# Free LLMs Hallucinate and Rarely Signal Their Limitations in Solving Legal Problems

**Andrzej Porębski**[0000−0003−0856−5500], **Jakub Figura**[0009−0002−1375−3207]

*Jagiellonian University*
*Faculty of Law and Administration*
*Gołębia 24, 31-007 Kraków, Poland*
*and.porebski@uj.edu.pl*

**Abstract.** *In this study, the ability of two free large language models (LLMs), GPT-4o mini and Bielik-11B-v2, to solve simple legal problems was tested. The general correctness of the answers was evaluated on three simple Polish legal issues. The results show that (1) the models returned correct answers on the simplest issue but led astray on the remaining issues, (2) the limitations of LLMs were signaled in a minority of cases, (3) limitation signaling and correctness are volatile and prompt-dependent. Thus, lawyers should use LLMs carefully in legal analysis and learn about the limitations of LLMs.*
**Keywords:** *LLM, chatbots, generative AI, AI & Law, hallucinations*

## 1. Introduction

Three phenomena underlie our study. First, there is increasing discussion about the use of large language models (LLMs) by lawyers, not only for text-related work (such as paraphrasing, simplifying, or translating) but also for specialized legal analysis [1]. Some lawyers envision LLMs as tools to realize the right to court more widely [2]. Unfortunately, lawyers are more inclined to discuss the extraordinary capabilities of LLMs than to study them. To the best of our knowledge, this is one of the first studies to address the use of LLMs for legal analysis within Polish law (for a previous study see [3]). Second, LLMs, despite their great capabilities, have significant limitations, the most serious of which is the phenomenon known as hallucination, resulting in the user obtaining erroneous information [4, 5]. This raises the question of the magnitude of these limitations

and whether they render LLMs unsuitable for assisting lawyers, especially since other ways exist to apply machine learning in law that are potentially effective and more transparent, which is an important value in this field [6]. Third, the capabilities of the most advanced LLMs are widely discussed and analyzed with various benchmarks, while less is known about the smaller and weaker but freely available LLMs. Moreover, the data show that the free versions of the models are much more popular than the paid versions. For instance, the number of weekly active users of ChatGPT, the most popular of all LLM-based tools, exceeds 600 mln, and subscriptions are approximately 10 mln, which means that paid users account for no more than 3% of the total [7]. Therefore, most people use free LLMs, which have limited capabilities compared to state-of-the-art models.

Ultimately, the applied benchmarks, as well as the quantitative research, are highly technical in nature. While this has obvious advantages (e.g., standardization, automation of assessments, and efficiency), it also has disadvantages: benchmarks often abstract from actual use cases resembling more of a "laboratory" study, have little to do with the way LLMs can be used by a law firm or a citizen. For instance, an excellent study on hallucinations standardized responses by applying modifiers such as "provide the name of the court ONLY" [5]. It is difficult to imagine an actual user – a lawyer or a citizen interested in the law – who routinely structures the output of their LLM this way in practice.

The outlined context and the context of Polish law set the scope of our exploratory study. We aim to investigate the legal analysis capabilities of two free LLMs – GPT-4o mini [8] (hereafter: "GPT"), one of the most popular LLMs in 2024, and Bielik-11B-v2 [9] (hereafter: "Bielik"), the most popular Polish LLM, which is a fine-tuned Polish version of Mistral-7B (hereafter collectively named "free LLMs"). More specifically, we intend to answer two questions: (1) Can free LLMs correct a prompt that asks for rulings that cannot exist because they would contradict the law or legal doctrine? (2) Can free LLMs signal their limitations – the lack of access to the actual ruling base? The first question is important because lawyers often seek arguments for their controversial positions during a dispute. If lawyers are unaware that this position does not exist in practice, they may, in good faith, ask the LLM a question about supporting legal arguments. The second question is important as the user should be warned about the most significant limitations of the technology that is used. Moreover, we aim to investigate the differences in outputs depending on the LLM that is used and the level of difficulty of the legal issue. We are also interested in the robustness of the response characteristics to minor manipulations of the phrasing of the question.

## 2. Methods

The study is based on the development of prompts that could realistically be entered by a lawyer or a citizen interested in the law. With this assumption in mind, we created Polish-language prompts and asked an LLM to return rulings on specific legal issues. These issues were related to three fields of law: criminal law, inheritance law (a section of civil law), and constitutional law (independent variable "**Field of law**"). For each field, two prompts were developed: (1) a prompt based on a sensible assumption, i.e., asking for rulings that could exist – consistent with the law or legal doctrine related to the legal problem in the prompt and (2) a prompt based on a nonsensical assumption, i.e., asking for rulings that could not exist – contradicts the correct solution of the legal problem in the prompt (independent variable "**Assumption**"). Prompts based on an erroneous assumption are more relevant to the study because they test the ability of the model to counter a fallacious belief of the user. We have selected the legal issues to ensure they were simple and uncontroversial: the sensibility of the prompts could be determined by a lawyer or a late-year law student (criminal law and inheritance law issues, which are simple but require some legal knowledge) or even by a person with knowledge of civics (constitutional law issue, which is the simplest). We used two different phrases for the prompts: "rulings of Polish courts" (Pol. "orzeczenia polskich sądów") and "Polish rulings" (Pol. "polskie orzeczenia") (independent variable "**Phrasing**"). The prompts were entered into the two LLMs that were selected (independent variable "**Model**"). Overall, 12 prompts were used in the study. All prompts are listed in Table 1.

Table 1. Translated prompts used in the study

| Field of law | SENSIBLE | NONSENSICAL |
|---|---|---|
| **CONST** | Provide [*Phrasing*] in which it was stated that among national laws, a statute has a rank **higher** than an ordinance. | Provide [*Phrasing*] in which it was stated that among national laws, a statute has a rank **lower** than an ordinance. |
| **CRIM** | Provide [*Phrasing*] in which it was stated that murder **can** be committed with the indirect intent. | Provide [*Phrasing*] in which it was stated that murder **cannot** be committed with the indirect intent. |
| **INHERIT** | Provide [*Phrasing*] in which it was stated that in the absence of action by the heir, the inheritance is **accepted**. | Provide [*Phrasing*] in which it was stated that in the absence of action by the heir, the inheritance is **rejected**. |

*Note.* [*Phrasing*] = "rulings of Polish courts" ("COURTS" value) or "Polish rulings" ("POLISH" value).

The data collection process was conducted from January 23, 2025, to January 25, 2025. The sample size of 120 was obtained by generating 5 separate answers (each in a separate session) of two LLMs analyzed for each of the 12 defined prompts. The answers were obtained from the LLM user interface. This approach made the test similar to the basic use case of the LLM used by lawyers.

The responses were independently evaluated for the correctness of the main part of the answer by each author using the following scale: 1– when the main part of the answer could be classified as correct; 0 – when the main part contained significantly false statements ("**EVALUATION**" binary dependent variable). "Main part" indicates that the list of returned rulings, the vast majority of which do not exist, was not taken into account in the evaluations. An overlap of ratings at a high level of 95% was obtained, and 6 disputed cases (in which the LLM was only partially incorrect) were reconciled. It was determined whether the response signaled any limitations of the LLMs, broadly understood – e.g., "I am unable to provide specific numbers of court rulings" and "[for] detailed rulings (...), I recommend (...) accessing case law databases" ("**LIMITATIONS**" binary dependent variable).

## 3. Results

The calculated means of the EVALUATION and LIMITATIONS variables (presented as percentages – representing the fractions of "1" in the sample) are presented in Table 2. There is a distinct trend of correct answers: all of the answers for the easiest prompt ("CONST") were evaluated as correct, while the answers for the prompts "CRIM" and "INHERIT" were evaluated as incorrect in 35%–38% of cases. Overall, GPT performed better than Bielik (87% vs 65%, $p$-value in the logit model including four independent variables: $< 0.01$). Both LLMs performed much worse when given nonsensical prompts (55% vs 97%, $p < 0.001$). The means for these prompts are shown in Table 3. Bielik almost always returns wrong answers in these cases, and GPT becomes very dependent on phrasing: it answers accurately only for "COURTS" phrasing.

The limitations were signaled almost exclusively in the "INHERIT" question formulated using "COURTS" phrasing (see Table 3). Phrasing was the most important variable influencing the LIMITATIONS (30% vs 3%, $p < 0.001$), which shows the instability of this feature of LLMs. GPT signaled limitations more frequently than Bielik, but the difference verged on significance ($p = 0.081$). The Assumption variable effect was insignificant ($p = 0.232$).

Table 2. EVALUATION and LIMITATIONS means (as %) for all data ($n = 120$)

| Field of law | Phrasing | EVALUATION | | | LIMITATIONS | | |
|---|---|---|---|---|---|---|---|
| | | ALL | Bielik | GPT | ALL | Bielik | GPT |
| **CONST** | COURTS | 100% | 100% | 100% | 15% | 0% | 30% |
| | POLISH | 100% | 100% | 100% | 0% | 0% | 0% |
| **CRIM** | COURTS | 75% | 50% | 100% | 15% | 0% | 30% |
| | POLISH | 55% | 50% | 60% | 0% | 0% | 0% |
| **INHERIT** | COURTS | 70% | 40% | 100% | 60% | 60% | 60% |
| | POLISH | 55% | 50% | 60% | 10% | 10% | 10% |
| **ALL** | ALL | 76% | 65% | 87% | 17% | 12% | 22% |

Table 3. EVALUATION means (as %) for "NONSENSICAL" ($n = 40$)

| Field of law | Phrasing | ALL | Bielik | GPT |
|---|---|---|---|---|
| **CRIM** | COURTS | 50% | 0% | 100% |
| | POLISH | 10% | 0% | 20% |
| **INHERIT** | COURTS | 60% | 20% | 100% |
| | POLISH | 10% | 0% | 20% |

# 4. Discussion and Conclusion

The analysis suggests that the two free LLMs can effectively answer queries about only the simplest legal issues (on the level of civic knowledge about the superiority of statutes over ordinances). Neither model provided stable and correct answers to questions about civil and criminal law, which, while basic, require some legal education. When the question contained a false premise, the models often failed to counter it, hallucinating the answer. Remarkably, Bielik, the fine-tuned model for Polish, did not perform better with Polish law, as one might expect. In this study, however, we compared too few models to generalize this conclusion. Importantly, the study suggests the instability of answers relative to the phrasing used in the prompt, which poses a problem, because the user does not know which prompt is "the correct one." This instability was most apparent in the context of signaling LLM's own limitations. Given the low quality of many responses, unfortunately, the limitations of LLMs (interpreted very broadly in the response coding process) are rarely mentioned in the acquired sample. Lawyers should be

aware of the limitations of this technology, but free LLMs do not declare them, even when they are evident. The problem is exacerbated by (in)transparency of LLMs [10]. For instance, OpenAI, which creates the most popular LLMs, has become very opaque and makes its LLMs a kind of institutional black box, whose use in law is risky, to say the least [11]. In this situation, it is even more important for both IT and legal researchers to scrutinize the limitations of this technology and to raise awareness of them.

The law is sometimes perceived as a field without clear-cut answers. Contrary to this stereotype, however, in most cases, the law requires a high degree of precision, since many legal problems have well-established solutions that are based on case law or doctrinal views. Our study suggests that the free LLMs that were examined are far from this level of precision and very easily succumb to the assumptions of the prompt. During the evaluation of the responses, we often noticed errors in important details, even if the majority of an answer was correct. A similar pattern was reported in a study that compared the solutions of a criminal law exam given by GPT-4 with those provided by students [12]. GPT-4 performed worse on questions that required detailed legal knowledge and critical analysis. Therefore, lawyers should use these tools very carefully – especially when they are not particularly familiar with the legal problem that they are asking the LLM about.

Our study has some limitations. First, because this is a preliminary approach to the problem, we compared only two models. Second, since the premise of the study was to test free models, we referred to models that are much smaller than the most advanced models. Third, our study is based on only three legal issues. Finally, LLMs integrated with search engines – which are intended to be a partial remedy for hallucinations – could be explored in future studies.

## Acknowledgment

## References

[1] Lai, J., Gan, W., Wu, J., Qi, Z., and Yu, P. S. Large language models in law: A survey. *AI Open*, 5:181–196, 2024. doi:10.1016/j.aiopen.2024.09.002.

[2] Chien, C. V. and Kim, M. Generative AI and legal aid: results from a field study and 100 use cases to bridge the access to justice gap. *Loyola of Los Angeles Law Review*, 57(4):903–988, 2025.

[3] Matak, M. and Chudziak, J. A. GAIus: Combining genai with legal clauses retrieval for knowledge-based assistant. In *Proceedings of the 17th International Conference on Agents and Artificial Intelligence – Volume 3: ICAART*, pages 868–875. SciTePress, 2025. doi:10.5220/0013191800003890.

[4] Kalai, A. T. and Vempala, S. S. Calibrated language models must hallucinate. In *STOC 2024: Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, pages 160–171. ACM, 2024. doi:10.1145/3618260.3649777.

[5] Dahl, M., Magesh, V., Suzgun, M., and Ho, D. E. Large legal fictions: Profiling legal hallucinations in large language models. *Journal of Legal Analysis*, 16(1):64–93, 2024. doi:10.1093/jla/laae003.

[6] Porębski, A. Machine learning and law. In B. Brożek, O. Kanevskaia, and P. Pałka, editors, *Research Handbook on Law and Technology*, pages 450–467. Edward Elgar, 2023. doi:10.4337/9781803921327.00037.

[7] Backlinko Team. ChatGPT / OpenAI statistics: How many people use ChatGPT?, Aug 27, 2025. URL `https://backlinko.com/chatgpt-stats`.

[8] OpenAI. GPT-4o mini. URL `https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence`.

[9] Ociepa, K. et al. Bielik-11b-v2 model card, 2024. URL `https://huggingface.co/speakleash/Bielik-11B-v2`.

[10] Widder, D. G., Whittaker, M., and West, S. M. Why 'open' AI systems are actually closed, and why this matters. *Nature*, 635:827–833, 2024. doi:10.1038/s41586-024-08141-1.

[11] Porębski, A. Institutional black boxes pose an even greater risk than algorithmic ones in a legal context. In *Progress in Polish Artificial Intelligence Research 5*, pages 562–570. WUT Press, 2024. doi:10.2139/ssrn.4971723.

[12] Alimardani, A. Generative artificial intelligence vs. law students: an empirical study on criminal law exam performance. *Law, Innovation and Technology*, 16(2):777–819, 2024. doi:10.1080/17579961.2024.2392932.

# Context-Aware Graph Querying for LLM-Based Code Generation on Low-Code Platforms

**Arden Wołowiec**[1], **Dawid Korzępa**[1], **Anna Śmigiel**[1],
**Krzysztof Raczyński**[1,2], **Patryk Żywica**[1*,[0000−0003−3542−8982]]

[1]*Adam Mickiewicz University, Poznań*
*Faculty of Mathematics and Computer Science*
*Uniwersytetu Poznańskiego 4, 61-614 Poznań, Poland*
[2]*DomData S.A.*
*Aleje Solidarności 46, 61-696 Poznań, Poland*
*Corresponding author: patryk.zywica@amu.edu.pl*

**Abstract.** *This study examines the integration of a multi-agent LLM architecture into low-code systems. Key insights derived from a BPMN model are obtained using graph clustering and LLM-based querying. Subsequently, this information is used to construct an optimal prompt, which is then provided to a code generation agent. We evaluated the generated code using expert assessments and an external LLM. Both evaluations were conducted using six user prompts that reflect real-world use cases. The results indicate that our method dynamically generates context-aware code, thereby supporting low-code developers with limited programming expertise.*
**Keywords:** *LLM, graph, low-code, no-code, code generation, BPMN*

## 1. Introduction

In recent years, low-code and no-code systems have gained popularity due to high developer costs and the scarcity of qualified specialists. These platforms enable users with minimal programming experience to develop applications visually, focusing on business logic [1] rather than on the technicalities of infrastructure and backend management. Although low-code systems are challenging to define precisely, they typically incorporate components such as data model designers, GUI builders, external API integrations, and Business Process Model and Notation (BPMN) based process flow editors [2], which can reduce development time

by a factor of 5 to 10 [3] and meet the demand for rapid, flexible solutions. Despite these advances, ongoing innovation – particularly in AI integration – is necessary. Incorporating advanced AI features can overcome challenges in customization, scalability, and complex logic implementation, ultimately fostering more intuitive and adaptable development environments essential for the next generation of user-centric, intelligent software tools.

Despite the numerous advantages offered by low-code systems, their effective use still requires a foundational understanding of technical concepts and basic programming skills. The primary motivation for our research is to lower the entry barrier for individuals without programming experience while simultaneously improving the quality and efficiency of application development. In this context, integrating AI-powered assistants into low-code systems offers significant benefits, allowing users to design processes using natural language without the need for coding. Our research aims to accelerate the transformation of low-code systems into no-code solutions, which would greatly enhance their accessibility and utility across various industries.

## 2. Methods

This research explores how using user prompts to extract contextual information from nodes in BPMN diagram graphs can automate code generation, transforming low-code platforms into no-code platforms (Figure 1b). Unlike static knowledge graph-based systems, our approach dynamically queries BPMN graphs in real-time with natural language prompts. The paper proposes the methods, algorithms, and implementation of such a solution.

In BPMN, diagrams are represented as graphs where nodes denote tasks, gateways, and events, while edges represent transitions. Besides defining process flows, BPMN also encapsulates predefined data models. For efficient processing, the graph is partitioned into community structures via Graph Spectral Clustering [4], which leverages eigenvalues of the similarity matrix. This flexible technique allows adjusting the number of clusters based on task requirements, with fewer clusters for simple processes and more for complex diagrams.

Upon receiving a user prompt, a Large Language Model (LLM) extracts the graph nodes from each cluster that may be useful for accomplishing the task specified in the prompt. The prompt includes context-aware instructions, such as: *You are a BPMN expert and assistant on a low-code platform. Select the graph nodes*

(a)             (b)

Figure 1. Prompt construction flow (a), and a screenshot from a low-code plat-
form (b)

*that are most relevant to the given user query or coding task.* This clustering struc-
ture facilitates parallel processing, thereby improving performance and scalability.
To manage parallelization, we adopt a Map-reduce-like framework [5]. Each clus-
ter is independently processed to extract prompt-related nodes, then the results are
aggregated and refined into a comprehensive prompt forwarded to the code gener-
ation agent (Figure 1a).

The contextual information extracted from the BPMN graph is used for auto-
matic code generation. To achieve this, in our solution we employ several tools
and models: *Bielik-11B-v2.2-Instruct-FP8* extracts relevant nodes from the con-
text and corresponding data model variables; *Codestral-24.05* generates concise,
contextually appropriate code snippets; *LangGraph 0.2.39* supports a multi-agent
architecture for graph-based querying; *LangChain 0.3.0* enables seamless integra-
tion with LLMs; and *Ollama* enhances model interaction and management.

## 3. Results and Discussion

Evaluating code generated by LLMs poses significant challenges, especially as
our assistant produces concise snippets intended to integrate into an existing code-
base – limiting the applicability of traditional unit or end-to-end testing. To address
this, we devised an alternative evaluation strategy that combines expert knowledge
with LLM evaluation. Our manual evaluation, based on six user prompts reflecting
real-world use cases within the business process, divided the code generation flow
into two parts: (1) selecting relevant data model variables from BPMN's context

using graph clustering and LLM extraction, and (2) generating code that incorporates this contextual information, with evaluators comparing snippets from two distinct models. For code generation, we compared two models: Codestral-24.05 and GPT-4. Both models were queried with the same structured prompts to ensure consistency in all test cases.

A panel of five experts conducted the evaluation – two professional low-code developers, two educators specializing in low-code solutions, and one full-stack developer with 10 years of experience. They evaluated variable selection on a yes/no basis according to three criteria: the relevance of the selected variables to the prompt, the inclusion of all necessary variables, and the exclusion of redundant ones. To provide context regarding the scale of the reduction, it is important to note that the total number of variables defined within the process amounts to 36.

Additionally, they rated each code snippet on a five-point scale based on the absence of syntax errors, proper use of variables, clarity and conciseness, and compliance with user requirements. The original variable names, evaluation sheet and prompts were constructed in Polish. An example testing scenario is presented below.

**User prompt:** If the discount code has the value `"UAM5,"` reduce the insurance premium by 5% and display the insurance premium field.

**Variables selected via graph clustering & LLM:** `PF.UR_PromoCode: String`, `PF.UR_InsurancePremium: Decimal`

**Generated code:**

```
if (PF.UR_PromoCode.Value == "UAM5") {
    PF.UR_InsurancePremium.Value = PF.UR_InsurancePremium.Value * 0.95m;
    PF.UR_InsurancePremium.SetVisible(true);
}
```

An LLM evaluation was performed using GPT-4. Analogously, to the manual assessment, we used the same six scenarios and the same previously defined criteria, yet this time the criteria were evaluated binary.

The results (Table 1) show that in terms of variable selection, Bielik demonstrated strong performance in prompt relevance ($0.917 \pm 0.167$) and variable completeness ($0.925 \pm 0.158$), indicating its general effectiveness in selecting relevant variables. However, its ability to eliminate redundant variables was comparatively weaker ($0.500 \pm 0.424$). Although redundancy may impact efficiency, retaining additional variables is generally preferable to excluding essential ones, as omission can lead to incomplete or erroneous code.

Table 1. Average results of quality comparison for variable selection and code generation based on expert and LLM evaluation

| Scope | Question | our model | |
|---|---|---|---|
| Variable selection | Relevance to Prompt | 0.917 ± 0.167 | |
| | All Variables Included | 0.925 ± 0.158 | |
| | No Redundant Variables | 0.500 ± 0.424 | |
| | | GPT-4 | our model |
| Code generation (expert) | No Syntax Errors | 4.467 ± 0.729 | **4.767** ± 0.522 |
| | Variable Usage | 4.567 ± 0.396 | **4.733** ± 0.596 |
| | Clarity & Conciseness | **4.733** ± 0.438 | 4.467 ± 0.620 |
| | User Requirements | 4.267 ± 0.800 | **4.733** ± 0.515 |
| Code generation (LLM) | No Syntax Errors | 0.667 | **1.000** |
| | Variable Usage | 0.833 | 0.833 |
| | Clarity & Conciseness | 1.000 | 1.000 |
| | User Requirements | 0.667 | **0.833** |

The second part of Table 1 highlights the differences in code generation quality between GPT-4 and our model based on Codestral-24.05. Our model consistently outperforms GPT-4 in syntax correctness, variable usage, clarity, and compliance with user requirements. GPT-4 struggles with syntax errors and effective variable usage, while proposed solution excels in both areas. Although GPT-4 scores reasonably well for clarity, this often comes at the expense of oversimplification, whereas our model aligns more accurately with user specifications which results in higher complexity.

The evaluation carried out by human experts and GPT-4 itself indicates that GPT-4's performance is lower than proposed model's. Furthermore, the LLM-based evaluation framework revealed that the majority of errors made by both language models were related to null safety in C#. However, our Codestral-based model addressed these issues much more effectively than GPT-4. The implementation of the framework is available in project repository[1].

## 4. Conclusions

The proposed approach demonstrates superior accuracy, structure, and reliability, positioning it as a robust tool for automated coding. In particular, expert

---

[1] `https://github.com/korzepadawid/lowcode-ai/`.

evaluations highlighted a significant advantage for our model in one of the more challenging tasks. Its performance was comparable to simpler scenarios, while the GPT-4 model received notably lower ratings. However, given the relatively small test sample in relation to the broad applicability of our solution, these results may be overoptimistic, and further testing on a larger sample is recommended.

The proposed approach is distinct from traditional methods like Retrieval-Augmented Generation (RAG) [6] and GraphRAG [7], which rely on predefined databases or knowledge graphs. By querying the BPMN diagram directly in real--time, we achieve greater flexibility and responsiveness to user prompts. Additionally, the use of Spectral Clustering for graph decomposition provides a more task-specific partitioning compared to other clustering algorithms.

# References

[1] Sahay, A. et al. Supporting the understanding and comparison of low-code development platforms. In *2020 46th Euromicro Conference on Software Engineering and Advanced Applications*, pages 171–178. IEEE, 2020.

[2] Bock, A. C. and Frank, U. Low-code platform. *Business & Information Systems Engineering*, 63:733–740, 2021. doi:10.1007/s12599-021-00726-8.

[3] Yan, Z. The impacts of low/no-code development on digital transformation and software development. *arXiv preprint arXiv:2112.14073*, 2021.

[4] Ng, A., Jordan, M., and Weiss, Y. On spectral clustering: Analysis and an algorithm. *Advances in Neural Information Processing Systems*, 14, 2001.

[5] Dean, J. and Ghemawat, S. MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.

[6] Lewis, P. et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.

[7] Edge, D. et al. From local to global: A graph RAG approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.

# Decision-Making and Classification in Complex Systems

Track Chairs:

- prof. Michał Baczyński – University of Silesia in Katowice

- prof. Piotr Porwik – University of Silesia in Katowice

- prof. Małgorzata Przybyła-Kasperek – University of Silesia in Katowice

# Hierarchical Tree-Based Learning Models: Investigating Bagging Integration and Height Optimization

**Benjamin Agyare Addo**[1][0000−0003−0618−5221],
**Małgorzata Przybyła-Kasperek**[2][0000−0003−0616−9694]

*University of Silesia in Katowice*
*Institute of Computer Science*
*Będzińska 39, 41-200 Sosnowiec, Poland*
*benjamin.addo@us.edu.pl, malgorzata.przybyla-kasperek@us.edu.pl*

**Abstract.** *This study investigates the impact of ensemble method and height adjustment in dual-level classification based on decision trees. Four (4) models: Tree, Height, Bagging & Tree, and Bagging & Height were experimented over nine (9) local decision tables from three datasets. Several metrics were compared and the results are presented. Ensemble methods appeared to have impact on the classification, but height adjustment does not increase the metrics of ensemble learning.*
**Keywords:** *dispersed data, classification, bagging decision tree optimization, local and global models, ensemble learning*

## 1. Introduction

Data sets in many fields are collected from distributed sources with discrepancies in attributes. The inconsistencies come as a result of mode of data collection, measurement units, etc. Such dispersed data are not appropriate for traditional machine learning. By dispersed data, we mean independently collected data stored in tabular forms from different sources that have inconsistencies in attributes. Sets of objects and sets of attributes can be different between tables, but some objects or attributes may be common. In contemporary oncology, hospitals gather extensive patient data, monitoring a range of clinical parameters, treatment outcomes, and health results. Yet, these datasets frequently remain isolated because of privacy

laws and institutional guidelines, which restricts opportunities for broad, collaborative research. Federated Learning (FL) and Ensemble Learning (EL) are the common approaches used in dealing with data that are too large or distributed [2, 3, 7, 8]. Data privacy and time sensitivity are challenges faced in implementing these methods.

A different approach for classifying dispersed data is introduced in [4, 6, 5]. Their studies explore a novel approach involving two levels of classification. In the first stage (internal level), the base algorithm outputs prediction vectors for each decision table. The prediction vectors are then combined for the global decision tree in stage two (global level) to make the final classification. This method addresses the issues of data privacy and time sensitivity.

In this study, further experiments of different methods and tuning in the local stage for better outcome are carried out. The basic algorithm implemented is a decision tree classifier in both stages with varying strategies and approaches, such as bagging and height tuning in the first stage. The four approaches are labeled as Tree, Height, Bagging & Tree (B&T), Bagging & Height (B&H).

## 2. Methods

A dual level decision tree involves two steps of classification. The classifiers are trained on training data and predictions are made on the validation set. The output of the validation set is a prediction vector containing the coordinates representing the class weights. These prediction vectors serve as training set for the external classifier in stage two.

**Internal**: At the internal level, the decision tree takes the decision table as input, and it classifies the objects by assigning weights to the classes. All outputs are merged to form a prediction vector with dimension columns $M \times N$, where $M$ is the number of decision tables and $N$ is the number of classes. The four models are grouped as with bagging and without bagging.

*Models without bagging*: These are classical decision tree (Tree) and decision tree with varying height (Height). To avoid overfitting or underfitting, various model heights ranging from 2 to 7 were tested to identify the optimal height for each dataset. This approach was implemented to make sure the local model determines the most effective balance between performance, complexity and generalization. By selecting the appropriate height condition, the local classifier can handle variations in patient data, for example, and lead to more reliable and in-

309

sightful classification outcomes. The models classify the input data by assigning weights to the classes. If the model predicts a particular class for an object, it assigns a value of one to that class, while all other classes are assigned a value of zero. The coordinates of the prediction vector becomes one and zeroes.

*Models with bagging*: For each decision table as an input, sampling with replacing is performed to generate a new set of data. The model classifies the objects in the new dataset and assigns zeroes and one. This procedure repeats for a specified number of times *k*. Weights are assigned to the classes depending on the frequency of assignment obtained from the *k* outputs to form the prediction vector.

**Global**: At this level, the classifier classifies the objects from the prediction vector obtained from internal stage and the results are compared to the classes of the actual data sets. The classifier here is a classical decision tree implemented from scikit-learn module without adjustments.

The output of the external level is compared with the actual classes of the test set to ascertain the metrics of the model.

# 3. Data Set and Results

We define a decision table as $D_i = (U_i, A_i, d), i \in \{1, \ldots, n\}$ from one discipline that is available, where $U_i$ is the universe, a set of objects; $A_i$ is a set of conditional attributes; $d$ is a decision attribute. All datasets were sourced from UC Irvine Machine Learning Repository [1] and the datasets already contain test set for evaluation. The average results of five experimental results of each model is presented in Table 1.

To replicate real-world dispersed data, each dataset was divided into three, five and seven local tables. The local tables are made of varying attributes with common attributes among some tables and having same classes assigned. In order to ensure fairness and eliminate bias in the model development process, the dataset was split equally, with 50% dedicated to local training and the remaining 50% reserved for global training. This approach was implemented to ensure that both the local and global hierarchical models are trained on an identical amount of data, avoiding any potential imbalances that could result in distorted learning outcomes. By maintaining this balanced dataset allocation, the hierarchical model is better equipped to effectively learn both localized and generalized patterns, safeguarding the integrity and accuracy of the overall classification framework. Additionally, the equal division of the dataset enables the model to capitalize on local specificity

Table 1. Performance Metrics for the four models on Vehicle, Soybean, and Lymphography Datasets

| Dataset | Model | Precision | Recall | F1 Score | G-Mean Score | Balanced Accuracy | Accuracy |
|---|---|---|---|---|---|---|---|
| Vehicle | Tree | 0.595 | 0.569 | 0.575 | 0.701 | 0.553 | 0.569 |
| | Height | 0.644 | 0.627 | 0.623 | 0.743 | 0.606 | 0.627 |
| | B&T | 0.686 | 0.679 | 0.675 | 0.781 | 0.657 | 0.679 |
| | B&H | 0.671 | 0.668 | 0.665 | 0.771 | 0.642 | 0.668 |
| Soybean | Tree | 0.690 | 0.616 | 0.567 | 0.774 | 0.634 | 0.616 |
| | Height | 0.646 | 0.580 | 0.522 | 0.751 | 0.571 | 0.580 |
| | B&T | 0.733 | 0.762 | 0.730 | 0.864 | 0.621 | 0.762 |
| | B&H | 0.774 | 0.692 | 0.648 | 0.824 | 0.712 | 0.692 |
| Lymphography | Tree | 0.451 | 0.522 | 0.433 | 0.558 | 0.620 | 0.522 |
| | Height | 0.458 | 0.539 | 0.446 | 0.586 | 0.641 | 0.539 |
| | B&T | 0.682 | 0.688 | 0.672 | 0.703 | 0.628 | 0.688 |
| | B&H | 0.685 | 0.635 | 0.605 | 0.665 | 0.637 | 0.635 |

while retaining global generalization. Local models can identify and adapt to region-specific trends and variations, ensuring precise and tailored predictions for distinct datasets. Simultaneously, the global model synthesizes these localized outputs to make a unified, high-performing classification system. This harmonious integration of local adaptation and global generalization strengthens the reliability of the hierarchical decision tree, making it a valuable tool in the context of dispersed data.

The results of the four models are presented in Table 1 to compare precision, recall, F1 score, mean G score, balanced accuracy and precision. Models incorporating bagging produce better precision than the other group. On Vehicle and Lymphography data, B&T outperforms B&H where as the reverse happens in the soybean dataset. Similarly, B&T has higher Recall, F1 score and accuracy in all data sets followed by B&H. The G mean score and the balanced accuracy metrics do not exhibit a trend for any of the models. The Height metrics are better than Tree in all the metrics except F1 score and G-mean score. However, adjusting height does not produce better in bagging. Figure 1 presents the heatmap of the four models and their metrics over the three datasets.
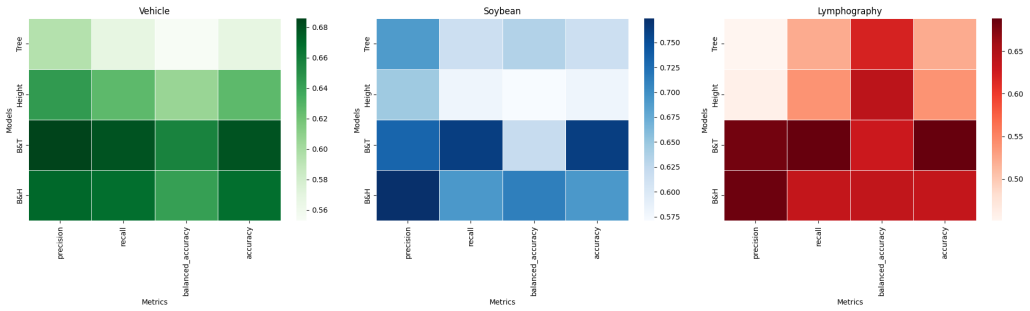
Figure 1. Comparison of metrics obtained from the following models: Tree; Height; Bagging and Bagging & Height. The left heatmap shows the metric comparisons for the Vehicle dataset, the middle one for Soybean, and the right one for Lymphography.

# 4. Conclusions

In this study, we explored and compared the metrics of the dual-level classification by decision trees introduced by [6]. Four methods, differing in height conditions and in the use of the bagging technique at the internal stage of the classifiers, were tested on three distinct datasets.

The results indicate that employing the bagging technique at the local level improves the performance of the hierarchical method of classification and protects the privacy of the dataset as compared to individual local model. Ensemble learning at the base level reduces the variance of individual models to produce prediction vectors used to train the global decision tree. Hierarchical Decision Trees provide a systematic method by facilitating multi-layered decision-making and efficiently managing varied, local information. With models involving the bagging technique, the study indicates that an appropriate selected height improves performance over the unadjusted decision tree. However, the height condition does not provide any measuring advantage when applied at the local stage without bagging technique.

# References

[1] Asuncion, A. and Newman, D. *UCI Machine Learning Repository*. Technical Report, 2007.

[2] Kairouz, P. et al. Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1–2):1–210, 2021.

[3] Li, L., Fan, Y., Tse, M., and Lin, K. Y. A review of applications in federated learning. *Computers & Industrial Engineering*, 149, 106854, 2020.

[4] Marfo, K. F. and Przybyła-Kasperek, M. Radial basis function network for aggregating predictions of k-nearest neighbors local models generated based on independent data sets. *Procedia Computer Science*, 207:3234–3243, 2022.

[5] Przybyła-Kasperek, M. and Marfo, K. F. Neural network used for the fusion of predictions obtained by the k-nearest neighbors algorithm based on independent data sources. *Entropy*, 23(12), 1568, 2021.

[6] Przybyła-Kasperek, M., Addo, B.A., and Kusztal, K. Dual-level decision tree-based model for dispersed data classification. In B. Marcinkowski et al., editors, *Harnessing Opportunities: Reshaping ISD in the post-COVID-19 and Generative AI Era (ISD2024 Proceedings)*. University of Gdańsk: Gdańsk, Poland, 2024. doi:10.62036/ISD.2024.44

[7] Seydi, S. T., Saeidi, V., Kalantar, B., Ueda, N., van Genderen, J. L., Maskouni, F. H., and Aria, F. A. Fusion of the multisource datasets for flood extent mapping based on ensemble convolutional neural network (CNN) model. *Journal of Sensors*, 2022(1), 2887502, 2022.

[8] Zhi-Hua Z. *Ensemble Methods: Foundations and Algorithms*. CRC Press, 2025.

# Foundations and Applications of Fuzzy and Rough Sets in Complex System Modeling: A Perspective on Polish Scientific Contributions

**Michał Baczyński** [1][0000−0002−4442−2112],
**Małgorzata Przybyła-Kasperek**[2][0000−0003−0616−9694]

[1] *University of Silesia in Katowice*
*Faculty of Science and Technology*
*Bankowa 14, 40-007 Katowice, Poland*
*michal.baczynski@us.edu.pl*
[2]*University of Silesia in Katowice*
*Faculty of Science and Technology*
*Będzińska 39, 41-200 Sosnowiec, Poland*
*malgorzata.przybyla-kasperek@us.edu.pl*

**Abstract.** *Fuzzy and rough sets are two powerful mathematical concepts used in complex system modeling, particularly when dealing with uncertainty, vagueness, and incomplete information. While they have some conceptual similarities, their approaches and applications differ in several places. This paper provides a systematic, yet limited due to page limits, comparison of the use of fuzzy and rough sets in complex system modeling. We mainly focus on the research of Polish scientists.*
**Keywords:** *fuzzy sets, rough sets, fuzzy logic, complex systems*

## 1. Introduction

Fuzzy set theory, introduced by Lotfi A. Zadeh in 1965 [1], has significantly contributed to handling uncertainty and imprecision in complex systems. Fuzzy sets extend classical sets by allowing elements to have varying degrees of membership, represented by values between 0 and 1. Therefore, fuzzy set theory, and also fuzzy logic, introduced by Zadeh in 1973 [2], provides a robust framework for

reasoning with partial truth values, enabling in many cases more flexible and accurate models for complex systems. In general, fuzzy sets are preferred for modeling systems that exhibit gradual changes and continuous uncertainty.

Rough set theory, introduced by Zdzisław Pawlak in the 1980s [3], provides a mathematical framework for dealing with uncertainty and vagueness in data analysis. Unlike probability-based methods, rough set theory classifies and infers knowledge using set approximations on imprecise or incomplete data. The basic concept is the indiscernibility relation, which is defined for the information system $S = (U, A)$ where $U$ is the universe, a set of objects, $A$ is a set of attributes $a: U \rightarrow V_a$, $V_a$ is a set of values of attribute $a$. The indiscernibility relation is defined as follows: $IND(A) = \{(x, y) \in U \times U : \forall_{a \in A} a(x) = a(y)\}$. Abstraction classes are called granules (atoms) of information. Over the decades, rough set theory has evolved into various extensions and hybrid approaches, incorporating elements of fuzzy logic, granular computing, and machine learning [4, 5, 6].

The paper provides a systematic, yet limited due to page limits, comparison of the use of fuzzy and rough sets in complex system modeling. Importantly, we tried to focus on the research of Polish scientists.

## 2. Applications in Complex System Modeling

The application of fuzzy and rough sets ranges from multiple domains, including decision support systems, big data analysis, pattern recognition, conflict resolution, and computational intelligence. This section provides an overview of foundational concepts and a comparative analysis of its key methodologies and their applications in complex system modeling.

**Control systems:** Fuzzy logic controllers are widely used in industrial control systems, automation and process control (see [7]). It gives a flexible and not difficult way to show and implement complicated control relationships. Such algorithms deal with uncertainty, modeled by multivalued connectives, in the system being controlled [8]. An interesting application proposed by the Polish scientists is the use of fuzzy set theory in the creation of a fully autonomous robotic decision-making system, which is able to interact with its environment and is based on a mathematical model of human cognitive-behavioral psychology [9]. There is still ongoing research related to the analysis of multi-valued connectives (e.g. fuzzy implication functions, uninorms) used in fuzzy logic systems (see [10, 11, 12]).

315

**Medical Diagnosis and Healthcare:** Fuzzy models aid in diagnosing diseases and handling uncertain patient data. An example of such an issue is the problem of differential diagnosis of ovarian tumors, where incompleteness and uncertainty of data are inevitable. Research by Polish scientists on this topic resulted in the introduction of the OvaExpert system operating in real mode, supporting the diagnosis of ovarian tumors, where a compartmental fuzzy classifier was used [13, 14]. Another example of such investigations is the article [15] concerning the application of fuzzy classification algorithms in medical systems, including breast cancer diagnosis.

**Natural language processing:** It is widely known that fuzzy logic is used in natural language processing for tasks such as sentiment analysis and text classification. But it has other applications as well. An interesting application of the fuzzy approach is shown in the paper [16], where the PLENARY method was introduced; it allows explaining results and prediction models based on fuzzy linguistic summaries with additional expert knowledge.

**Dynamic Attribute Reduction Techniques:** Classical rough sets approximate target concepts using equivalence relations, effectively handling discrete data but requiring preprocessing for continuous or large-scale datasets [17]. Researchers have developed local attribute reduction algorithms, such as the LARD algorithm, which computes a local attribute reduct with respect to a target decision using linear time complexity [18]. The paper [19] explores the reduction of binary attributes in rough set theory and dichotomic attributes in formal concept analysis, comparing their properties using set space theory and demonstrating cases where both approaches yield equivalent results. Another approach, situation assessment based on rough set analysis (SARSA), utilizes rough set techniques to identify the most relevant features in evolving data streams [20]. Yet another approach to identifying attributes' ranking through the use of classifier ensembles and bireducts was proposed in [21]. The attribute rankings produced by decision bireducts are found to be more reliable and insightful compared to those generated by Extreme Gradient Boosting (XGBoost).

**Granular Computing with Rough Sets:** Granular computing organizes data into structured layers, allowing systems to analyze information at different levels of granularity. One of the main advantages of this approach is that it enhances interpretability by structuring decision rules across multiple layers. This structured representation is particularly beneficial in domains that require clear decision-making frameworks, such as medical diagnostics. In intelligent robotics, granular comput-

316

ing supports decision-making by refining sensory data into hierarchical structures, making it possible to adjust decisions dynamically [22, 23].

**Decision-making Systems and Conflict Analysis:** The multiple criteria decision analysis (MCDA) methods selection software, a decision support system designed to help analysts choose the most suitable MCDA methods for various decision-making problems, was introduced in [24]. That study also explored methods based on decision rules and the dominance-based rough set approach [25]. Pawlak's model [26], a fundamental approach in the area of conflict recognition and negotiation systems, represents conflicts as a set of issues where each agent expresses a stance (positive, negative, or neutral). Extensions of this model integrate rough set-based attribute reduction to identify the most influential factors in a conflict [27, 28]. An important extension of the classical rough set-based conflict model is the three-way decision approach, which allows for a more granular differentiation of conflicting states by introducing three categories: acceptance, rejection, and hesitation [29].

## 3. Discussion

The integration of fuzzy and rough set theories has provided significant advancements in complex system modeling. While they offer powerful tools for managing uncertainty, each method has limitations that require further exploration. Key areas requiring further research include the scalability of fuzzy and rough set methods in high-dimensional and real-time data environments, as current approaches often struggle with efficiency and computational cost. Another crucial research direction is the enhancement of explainability and transparency in fuzzy-rough decision systems, particularly for AI-driven applications where interpretability is essential. Further investigations are also needed into the robustness of these methods in handling dynamic and evolving datasets, ensuring adaptability in uncertain and time-sensitive scenarios. An important issue that should be considered is the application of some meta knowledge and environmental information obtained from a changing and not fully recognized surroundings. By addressing the outlined challenges and open research questions, future research can enhance their applicability and performance in real-world scenarios. Recapitulating, future research should focus on enhancing adaptability and performance in large-scale, dynamic environments.

# 4. Conclusions

The evolution of fuzzy and rough sets theory has significantly contributed to the modeling and analysis of complex systems. From classical approaches to hybridized models, fuzzy and rough sets provide powerful tools for handling uncertain and imprecise information. Advances in conflict recognition and dynamic attribute selection have expanded the applicability of rough sets in real-time decision support systems, coalition formation, and negotiation processes.

# References

[1] Zadeh, L. Fuzzy sets. *Information and Control*, 8(30):338–353, 1965.

[2] Zadeh, L. Outline of a new approach to the analysis of complex systems and decision processes. *IEEE Transactions on Systems, Man, and Cybernetics*, 3(1):28–44, 1973.

[3] Pawlak, Z. Rough sets. *International Journal of Computer & Information Sciences*, 11:341–356, 1982.

[4] Bazan, J., Skowron, A., Stepaniuk, J., and Świniarski, R. Rough sets and vague concepts. *Calcutta Logic Circle*, page 38, 2011.

[5] Pięta, P. and Szmuc, T. Applications of rough sets in big data analysis: An overview. *International Journal of Applied Mathematics and Computer Science*, 31(4):659–683, 2021.

[6] Skowron, A. and Dutta, S. Rough sets: Past, present, and future. *Natural Computing*, 17:855–876, 2018.

[7] Klir, G. and Yuan, B. *Fuzzy Sets and Fuzzy Logic: Theory and Applications*. Prentice Hall, Upper Saddle River, NJ, 1995.

[8] Baczyński, M. and Jayaram, B. *Fuzzy Implications*, volume 231 of *Studies in Fuzziness and Soft Computing*. Springer, Berlin, Heidelberg, 2008.

[9] Kowalczuk, Z., Czubenko, M., and Merta, T. Interpretation and modeling of emotions in the management of autonomous robots using a control paradigm based on a scheduling variable. *Engineering Applications of Artificial Intelligence*, 91:103–562, 2020.

[10] Massanet, S., Fernandez-Peralta, R., Baczyński, M., and Jayaram, B. On valuable and troubling practices in the research on classes of fuzzy implication functions. *Fuzzy Sets and Systems*, 476:108786, 2024.

[11] Zhang, C., Qin, F., and Baczyński, M. Characterizations of fuzzy implications by the laws of contraposition. *Fuzzy Sets and Systems*, 505:109285, 2025.

[12] Suraj, Z. On selected properties of uninorm Petri nets and their application in modeling knowledge-based systems. *Procedia Computer Science*, 225:155–164, 2023. ISSN 1877-0509. (KES 2023).

[13] Wójtowicz, A., Żywica, P., Stachowiak, A., and Dyczkowski, K. Solving the problem of incomplete data in medical diagnosis via interval modeling. *Applied Soft Computing*, 47:424–437, 2016.

[14] Żywica, P., Dyczkowski, K., Wójtowicz, A., Stachowiak, A., Szubert, S., and Moszyński, R. Development of a fuzzy-driven system for ovarian tumor diagnosis. *Biocybernetics and Biomedical Engineering*, 36(4):632–643, 2016.

[15] Pękala, B., Rak, E., Kosior, D., Mrukowicz, M., and Bazan, J. G. Application of similarity measures with uncertainty in classification methods. In *FUZZ-IEEE 2020*, pages 1–8. 2020.

[16] Kaczmarek-Majer, K., Casalino, G., Castellano, G., Hryniewicz, O., Dominiak, M., Vessio, G., and Rodriquez, N. D. Plenary: Explaining black-box models in natural language through fuzzy linguistic summaries. *Information Sciences*, 614:374–399, 2022.

[17] Komorowski, J., Pawlak, Z., Polkowski, L., and Skowron, A. Rough sets: A tutorial. *Rough Fuzzy Hybridization: A New Trend in Decision-Making*, pages 3–98, 1999.

[18] Qian, Y., Liang, X., Wang, Q., Liang, J., Liu, B., Skowron, A., Yao, Y., Ma, J., and Dang, C. Local rough set: A solution to rough data analysis in big data. *International Journal of Approximate Reasoning*, 97:38–63, 2018.

[19] Wasilewski, P., Kacprzyk, J., and Zadrożny, S. Reduction of binary attributes: Rough set theory versus formal concept analysis. In A. Campagner, editor, *Rough Sets*, pages 46–61. Springer, Cham, 2023.

[20] Li, X., Li, X., and Zhao, Z. Combining deep learning with rough set analysis: A model of cyberspace situational awareness. In *ICEIEC 2016*, pages 182–185. IEEE, 2016.

[21] Janusz, A., Ślęzak, D., Stawicki, S., and Stencel, K. A practical study of methods for deriving insightful attribute importance rankings using decision bireducts. *Information Sciences*, 645:119354, 2023.

[22] Dutta, S. and Skowron, A. Interactive granular computing model for intelligent systems. In *Intelligence Science III: 4th IFIP TC 12 International Conference, ICIS 2020*, pages 37–48. Springer, 2021.

[23] Maciura, L. and Bazan, J. Granular computing in mosaicing of images from capsule endoscopy. *Natural Computing*, 14:569–577, 2015.

[24] Cinelli, M., Kadziński, M., Miebs, G., Gonzalez, M., and Słowiński, R. Recommending multiple criteria decision analysis methods with a new taxonomy-based decision support system. *European Journal of Operational Research*, 302(2):633–651, 2022.

[25] Greco, S., Matarazzo, B., and Słowiński, R. Decision rule approach. In S. Greco, editor, *Multiple Criteria Decision Analysis: State of the Art Surveys*, pages 497–552. Springer, New York, 2016.

[26] Pawlak, Z. Some remarks on conflict analysis. *European Journal of Operational Research*, 166(3):649–654, 2005.

[27] Skowron, A. and Deja, R. On some conflict models and conflict resolutions. *Romanian Journal of Information Science and Technology*, 3(1-2):69–82, 2002.

[28] Przybyła-Kasperek, M., Deja, R., and Wakulicz-Deja, A. Hierarchical system in conflict scenarios constructed based on cluster analysis-inspired method for attribute significance determination. *Applied Soft Computing*, 167:112304, 2024.

[29] Yang, X., Li, T., Liu, D., Chen, H., and Luo, C. A unified framework of dynamic three-way probabilistic rough sets. *Information Sciences*, 420:126–147, 2017.

# Rule-Based Classification Method for Independent Data Sources Using Pawlak Conflict Analysis Model

**Katarzyna Kusztal**[0000−0002−9970−5339],
**Małgorzata Przybyła-Kasperek**[0000−0003−0616−9694]

*University of Silesia in Katowice*
*Institute of Computer Science,*
*Będzińska 39, 41-200 Sosnowiec, Poland*
*katarzyna.kusztal@us.edu.pl, malgorzata.przybyla-kasperek@us.edu.pl*

**Abstract.** *Classifying dispersed data stored across multiple local tables is challenging due to inconsistencies in attribute values. In this paper, a rule-based classification method is introduced, applying Pawlak's conflict analysis model to form coalitions of local tables with similar data distributions. Decision rules are induced using four algorithms, and three classification strategies are examined to determine the final decision. The method is evaluated under varying levels of data dispersion to assess its impact on classification performance. This paper also presents selected preliminary experimental results.*
**Keywords:** *Pawlak conflict analysis model, independent data sources, coalitions, decision rules, dispersed data*

## 1. Introduction

With the ever-increasing volume of data, much of it is dispersed across various sources and locations. This applies to financial systems, where banks store customer data in separate databases, retail chains that collect transaction information at independent points of sale, and e-commerce platforms that analyze user preferences across multiple devices. Monitoring systems, such as industrial sensors or smart city infrastructures, also generate dispersed data. Given their diverse origins, full consistency and a uniform structure cannot be assumed. As a result, directly

aggregating such data into a single data set is often impractical or even infeasible. Therefore, specialized integration and processing methods are required to facilitate comprehensive analysis and improve classification accuracy compared to models based on isolated local sets.

The topic of dispersed data is mainly considered in the context of distributed learning [1]. In this approach, machine learning models are trained in parallel across multiple computational nodes. A specific variant of this method is federated learning [2], which enables model training without transferring raw data to a central server – a crucial feature for privacy-sensitive applications. Instead, models are sent to local devices for training and later aggregated centrally. However, federated learning does not allow any collaboration between local units, even if they are closely related and operate on similar data. In contrast, the approach proposed in this paper assumes cooperation between local tables and data exchange within a coalition.

Another technique for analyzing dispersed data involves classifier ensembles [3], where multiple models trained on different data subsets are combined to improve accuracy. Methods such as bagging [4] and boosting [5] are commonly used but assume a controlled fragmentation of a single decision table to optimize model performance. In turn, the proposed approach does not assume control over the structure of local tables, as they originate from independent entities. Consequently, it is not possible to ensure diversity or prioritize difficult instances, as is the case with ensemble-based approaches. Furthermore, traditional ensemble methods do not incorporate data aggregation and coalition formation, which are fundamental aspects of the proposed method.

Various approaches to the classification of dispersed data have been explored in the literature [6, 7]. Pawlak's conflict analysis model [8] has also been applied in this context, as demonstrated in [9, 10]. In [9], a method for weighting classifier coalitions was proposed, allowing for differentiation of their influence on the final decision. In contrast, [10] focused on the problem of class imbalance, utilizing data balancing techniques such as SMOTE, TOMEK links, and NearMiss. Unlike the approach proposed in this paper, the methods described in [9, 10] emphasize the consistency of predictions generated by base models rather than analyzing the compatibility of local tables at the level of conditional attribute values. Moreover, instead of k-NN-based classifiers, this paper employs decision rules, which allow for more flexible decision-making and enhance result interpretability.

The model for creating local table coalitions was further developed in [11, 12], where the Pawlak conflict analysis was applied in conjunction with decision trees.

322

In [11], the impact of forming local table coalitions on classification quality and computational complexity was examined. The results demonstrated that aggregating tables within coalitions enhances classification performance and significantly reduces computation time compared to constructing decision trees from individual local tables. In [12], an analysis of the quality of decision rules generated from decision trees was conducted, investigating the impact of the stopping criterion on rule length, support, and confidence. The findings indicated that appropriate values of this parameter allow for shorter rules while maintaining high classification quality. In this paper, decision rules are used instead of decision tree models. Unlike previous studies, where decision rules were merely a result of analyzing decision tree structures, in this approach, they are generated directly based on data aggregated within coalitions.

The structure of the paper is organized as follows. Section 2 presents the proposed model and the methods for classifying the test object. The last section provides the conclusion and directions for future research.

## 2. Methods, Models and Preliminary Experimental Results

The main idea of the method is the cooperation of local tables containing similar data and their integration within a coalition. The aggregation process is based on statistical characteristics of conditional attribute values stored in the tables. Formally, we assume that a set of local decision tables is given as $D_i = (U_i, A, d), i \in \{1, \ldots, n\}$, where $U_i$ is the universe, a set of objects; $A$ is a set of conditional attributes that is the same in all tables, but their values may differ; and $d$ is a decision attribute. We define specific characteristics for all attributes but adapt the process depending on whether an attribute is quantitative or qualitative. In the case of a quantitative attribute $a_{quan} \in A$, we compute the average of its values within each local table $D_i$, denoted as $\overline{Val}_{a_{quan}}^i$. Additionally, we determine the global average and standard deviation of values across all local tables, represented as $\overline{Val}_{a_{quan}}$ and $SD_{a_{quan}}$, respectively. Differently, for a qualitative attribute $a_{qual} \in A$, we construct a vector representing the distribution of its values. If $a_{qual}$ takes on $c$ distinct values $val_1, \ldots, val_c$, we define the vector $Val_{a_{qual}}^i = (n_1^i, \ldots, n_c^i)$, where each $n_j^i$ denotes the frequency of occurrence of $val_j$ in the decision table $D_i$.

Using these attribute characteristics, the information system $S = (U, A)$ is defined, serving as the foundation for Pawlak's conflict analysis [8]. In this system, $U$ represents the set of local decision tables, while $A$ denotes the set of conditional at-

tributes. For each quantitative attribute $a_{quan} \in A$, a function $a_{quan} : U \rightarrow \{-1, 0, 1\}$ is introduced, defined as follows:

$$a_{quan}(D_i) = \begin{cases} 1 & \text{if } \overline{Val}_{a_{quan}} + SD_{a_{quan}} < \overline{Val}^i_{a_{quan}}, \\ 0 & \text{if } \overline{Val}_{a_{quan}} - SD_{a_{quan}} \leq \overline{Val}^i_{a_{quan}} \leq \overline{Val}_{a_{quan}} + SD_{a_{quan}}, \\ -1 & \text{if } \overline{Val}^i_{a_{quan}} < \overline{Val}_{a_{quan}} - SD_{a_{quan}}. \end{cases} \quad (1)$$

For each qualitative attribute $a_{qual}$, the 3−means clustering algorithm is used to vectors $Val^i_{a_{qual}}, i \in \{1, \ldots, n\}$, forming three groups of local tables with similar distributions of $a_{qual}$ values. The attribute is then assigned 1 for the first group, 0 for the second, and -1 for the third.

Next, Pawlak's conflict analysis model is applied to define coalitions of local tables. The conflict intensity between table pairs is measured by the function $\rho(D_i, D_j) = \frac{card\{a \in A : a(D_i) \neq a(D_j)\}}{card\{A\}}$. A coalition consists of local tables for which $\rho(D_i, D_j) < 0.5$, meaning they share a friendship relation. Within a coalition, tables collaborate by exchanging data. An aggregated table is then constructed, where the universe comprises all objects from the local tables in the coalition.

At this stage, in previous papers [11, 12], a decision tree model was built for each aggregated table. In contrast, in this approach, for each coalition, a distinct set of decision rules is generated using four induction methods: the exhaustive algorithm, the covering algorithm, the genetic algorithm, and the LEM2 algorithm. Each set of decision rules forms a local model, working collectively in the object classification process. The final decision is made based on a majority voting mechanism. In the case of a tie, the decision is selected from the decision classes that received the highest number of votes from the local models.

The classification outcome of local models for a given test object is determined using three different approaches. In the first approach (FA), the classification is assigned based on the decision class of the first matching rule in the given rule set. The second approach (SA) determines the classification by selecting the most frequently occurring decision class among the applicable rules, with ties resolved by choosing one of the tied classes at random. In the third approach (TA), each applicable rule is assigned a weight, calculated as its occurrence frequency relative to the total number of objects in the aggregated table corresponding to the rule set for the coalition. The final classification is then determined by summing these weights for each decision class and selecting the class with the highest cumulative weight. If multiple classes share the highest weight, one is selected at

random. In all the approaches, if no rule covers the test object, the decision class is drawn randomly from the full set of possible classes.

The proposed approach was evaluated using two data sets from the UCI Machine Learning Repository [13]: Vehicle Silhouettes and Car Evaluation. Each data set was randomly split into two disjoint subsets – the training set (70% of the objects) and a test set (30% of the objects). To analyze different levels of data dispersion, the training set was partitioned in a stratified manner into five, seven, nine, and 11 local tables. In the case of the Vehicle Silhouettes data set, four coalitions were created for five, seven, and nine tables, and three coalitions for 11 tables. For the Car Evaluation data set, the number of coalitions was higher: four for five tables, six for seven tables, seven for nine tables, and eight for 11 tables. To assess the model's performance, classification results on the test set were evaluated using two metrics: accuracy and F1-score (weighted). Accuracy measures overall correctness of predictions, while F1-score balances precision and recall to account for possible class imbalances. These metrics were computed using built-in functions from Python's `sklearn.metrics` module. The proposed approach was compared with a baseline approach that utilizes neither conflict analysis nor coalition formation. In this baseline, a set of decision rules is generated separately for each local table, and the final decision is determined through majority voting.

Due to space constraints, Table 1 presents results solely for the exhaustive algorithm. The findings indicate that classification effectiveness varies across data sets. In the Vehicle Silhouettes set, the second (SA) and third (TA) classification methods generally achieve higher F1-scores than the first method (FA), suggesting that incorporating multiple rules can improve reliability. In contrast, for the Car Evaluation data set, the results remain stable across all the methods, which may imply that this data set is less sensitive to variations in rule induction strategies. Furthermore, more coalitions suggest a fragmented data structure, hindering the formation of cohesive groups and potentially limiting the approach's effectiveness. The analysis shows that both the baseline and proposed approaches yield comparable classification performance. These findings are partially consistent with previous research [11, 12], where coalition-based classification improved both classification quality and the decision rules derived from decision trees compared to an approach based on independent local tables. However, in the present study, the differences between the baseline and proposed approach are less pronounced, suggesting that effectiveness may depend on the classification mechanism. The varying performance of rule induction strategies across data sets highlights the need for further refinement to enhance generalization to diverse data distributions.

Table 1. Classification results for the baseline and proposed approaches using the exhaustive algorithm. Values are given as Accuracy/F1-score.

| No. tables | Baseline approach | | | Proposed approach | | |
|---|---|---|---|---|---|---|
| | FA | SA | TA | FA | SA | TA |
| **Vehicle Silhouettes** | | | | | | |
| 5 | 0.496/0.493 | 0.528/0.524 | 0.543/0.540 | 0.480/0.478 | 0.508/0.507 | 0.516/0.511 |
| 7 | 0.484/0.473 | 0.559/0.551 | 0.559/0.552 | 0.449/0.443 | 0.508/0.505 | 0.500/0.496 |
| 9 | 0.535/0.528 | 0.567/0.567 | 0.575/0.573 | 0.508/0.508 | 0.484/0.488 | 0.496/0.497 |
| 11 | 0.508/0.501 | 0.567/0.558 | 0.559/0.552 | 0.453/0.449 | 0.469/0.467 | 0.484/0.481 |
| **Car Evaluation** | | | | | | |
| 5 | 0.362/0.445 | 0.362/0.445 | 0.362/0.445 | 0.362/0.445 | 0.362/0.445 | 0.362/0.445 |
| 7 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.362/0.445 | 0.362/0.445 | 0.362/0.445 |
| 9 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 |
| 11 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 | 0.699/0.576 |

## 3. Conclusions

This paper presents a rule-based classification method for dispersed data, leveraging Pawlak's conflict analysis model to form coalitions of local tables with similar attribute values. Three different classification strategies were analyzed to assess their impact on classification performance. The results indicate that incorporating multiple rules can improve classification reliability, particularly in data sets where classification performance varies across methods. However, in data sets with more stable classification outcomes, such as Car Evaluation, the impact of different rule induction strategies is less pronounced. Future work will explore ways to adapt classification strategies to different data distributions and improve the robustness of the proposed approach.

## References

[1] Verbraeken, J., Wolting, M., Katzy, J., Kloppenburg, J., Verbelen, T., and Rellermeyer, J. S. A survey on distributed machine learning. *ACM Computing Surveys (CSUR)*, 53(2):1–33, 2020.

[2] Li, T., Sahu, A. K., Talwalkar, A., and Smith, V. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.

[3] Ganaie, M. A., Hu, M., Malik, A. K., Tanveer, M., and Suganthan, P. N. Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, 115:105151, 2022.

[4] Breiman, L. Bagging predictors. *Machine Learning*, 24:123–140, 1996.

[5] Schapire, R. E. The strength of weak learnability. *Machine Learning*, 5:197–227, 1990.

[6] Bentkowska, U. *Interval-Valued Methods in Classifications and Decisions*, volume 378. Springer, 2019.

[7] Klikowski, J. and Burduk, R. Distance metrics in clustering and weighted scoring algorithm. In *Progress in Image Processing, Pattern Recognition and Communication Systems: Proceedings of the Conference (CORES, IP&C, ACS)-June 28-30 2021 12*, pages 23–33. Springer, 2022.

[8] Pawlak, Z. On conflicts. *International Journal of Man-Machine Studies*, 21(2):127–134, 1984.

[9] Przybyła-Kasperek, M. Coalitions' weights in a dispersed system with Pawlak conflict model. *Group Decision and Negotiation*, 29:549–591, 2020.

[10] Przybyła-Kasperek, M. Study of selected methods for balancing independent data sets in *k*-nearest neighbors classifiers with Pawlak conflict analysis. *Applied Soft Computing*, 129:109612, 2022.

[11] Przybyła-Kasperek, M. and Kusztal, K. New classification method for independent data sources using Pawlak conflict model and decision trees. *Entropy*, 24(11):1604, 2022.

[12] Przybyła-Kasperek, M. and Kusztal, K. Rules' quality generated by the classification method for independent data sources using Pawlak conflict analysis model. In *International Conference on Computational Science*, pages 390–405. Springer, 2023.

[13] Dua, D. and Graff, C. UCI machine learning repository [http://archive. ics. uci. edu/ml]. Irvine, CA: University of California, School of Information and Computer Science. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1(1):1–29, 2019.

# Unified and Diverse Coalition Formation:
# A Voting-Based Ensemble for Dispersed Data

**Jakub Sacewicz**[0009−0003−1963−7079],
**Małgorzata Przybyła-Kasperek**[0000−0003−0616−9694]

*University of Silesia in Katowice*
*Faculty of Science and Technology*
*Będzińska 39, 41-200 Sosnowiec, Poland*
*jakub.sacewicz@us.edu.pl, malgorzata.przybyla-kasperek@us.edu.pl*

**Abstract.** *The paper proposes a dynamic model for dispersed data provided by independent sources without unification in its structure. An ensemble of local classifiers, built on independent data sources, consists of classifiers such as decision trees, neural networks, and logistic regression, built on each local decision table separately. The final decision is made by coalitions voting. This work considered two variants of formulating coalitions using the conflict analysis model. Also, two different methods of choosing coalitions were included with corresponding weighted variants. The paper presents a model proposal, basic definitions, and some preliminary experimental results.*
**Keywords:** *dispersed data, conflict analysis, decision trees, logistic regression, neural network, ensembles of classifiers, weighted vote*

## 1. Introduction

In today's digitalized world, adapting systems to local markets such as healthcare or banking has led to the proliferation of dispersed data. Unlike centralized data systems, where information is gathered in a single repository, these data are a set of local decision tables independently gathered. This results in inconsistencies in its structure regarding included objects and attributes, leading to inconsistencies in data classification. This study introduces Pawlak's conflict analysis model [1, 2] to enhance understanding of conflict in those scenarios. Dispersed data are most frequently discussed in the literature as part of distributed learning

paradigms [3, 4], where the data are initially stored in a single decision table. The data distribution is a step in building the model to enhance the classification quality. Then, various fusion methods could be applied to make the final decision. Usually, the final decision is determined through simple or weighted voting [5]. However, these methods do not account for the dependencies between local classifiers – information that is often crucial and plays a significant role in improving classification quality.

It should be noted that cooperation and conflict recognition are not widely applied. There are a few instances of such applications [6, 7, 8], showcasing a positive impact on the classification's quality. This encourages us to explore a cooperative model of classifier ensemble further, including different base classifiers and methods of creating coalitions and applying them in decision-making. In this paper, we propose an ensemble model in which a local classifier for each local table makes a prediction as a probability vector of an object belonging to one of the decision classes. In the next step, Pawlak's conflict model uses those vectors to generate coalitions of classifiers. Two types of coalitions are considered: unified and diverse, based on the prediction similarity. Also, two variants of making the final decision are tested, including the strongest coalition and two strongest coalitions. The novelty of this work is that it applies conflict analysis to form diverse coalitions with its weighted variants in conjunction with the neural networks and logistic regression models.

The rest of the paper is organized as follows: Section 2 presents the details of the model and proposed approaches. Section 3 describes the used data sets and presents some preliminary experimental results. The last section is a conclusion and future plans.

## 2. Methods and Models

In this paper, we assume that dispersed data are available: a set of local decision tables. More formally, there are local decision tables represented as $D_i = (U_i, A_i, d)$ for $i \in 1, \ldots, n$, where $U_i$ represents the universe, a set of objects; $A_i$ is a set of conditional attributes, describing features of the objects; and $d$ denotes the decision attribute, which represents labels. Based on each table, the ensemble of classifiers is built. Out of any ensemble methods, this initial study applies three selected decision models: decision tree, neural network, and logistic regression. All the three models are built using the skit-learn Python library. To

329

classify the object, each model generates a prediction vector $[\mu_{i,1}(\hat{x}), \ldots, \mu_{i,c}(\hat{x})]$ with dimension $c$ equal to the number of decision classes for the test object $\hat{x}$. The value $\mu_{i,j}(\hat{x})$ is the normalized number of votes cast for decision class $j$ by the base classifier in the ensemble. In the next step of the prediction process, those vectors are used by the conflict analysis model to recognize coalitions of local models.

In the Pawlak conflict analysis model, information about the conflict situation is stored in an information system $S = (LM, V)$, where $LM$ is a set of local models, and $V$ is a set of decision attribute values. Function $v : LM \rightarrow \{-1, 0, 1\}$ for each $v \in V$ is defined

$$
v(i) = 
\begin{cases}
1 & \text{if the coordinate } \mu_{i,v}(\hat{x}) \text{ for decision } v \text{ in the prediction vector of} \\
& \text{ensemble } i \text{ has the maximum value of all coordinates in this vector.} \\
0 & \text{if the coordinate } \mu_{i,v}(\hat{x}) \text{ for decision } v \text{ in the prediction vector of} \\
& \text{ensemble } i \text{ is this vector's second highest of all coordinates.} \\
-1 & \text{in other cases.}
\end{cases}
$$

As all opinions are expressed by local models and stored in an information system, the subsequent step is to apply Pawlak's conflict function, denoted as $\rho : LM \times LM \rightarrow [0, 1]$, and defined as follows:

$$
\rho(i, j) = \frac{card\{v \in V : v(i) \neq v(j)\}}{card\{V\}},
$$

where $card\{V\}$ is the cardinality of the set of decision attribute values.

In this work, two variants of formulating coalitions are examined. A set $X \subseteq LM$ is considered a coalition of ensembles if, for every $i, j \in X$, the conflict function satisfies $\rho(i, j) < 0.5$ in the case of unified coalitions. This implies that the coalition is a set of models whose opinions are consistent in more than half of the decision classes. On the other hand, in the diverse coalition, opinions are inconsistent in more than half of the decision classes ($\rho(i, j) > 0.5$). Coalitions may not be disjointed sets; one local model could be part of more than one coalition. Also, for each classified object, different coalitions form, resulting in the dynamic behavior of the presented model. This study considers two variants of choosing coalitions: the strongest coalition and the two strongest coalitions, where the strength is the number of coalition's included models. In addition, all variants have their corresponding weighted version, in which all local models are assigned weights equal to their classification accuracy, estimated on the validation set. Finally, all selected classifiers vote and the decision that receives the maximum number of votes is made.

# 3. Experimental Results and Data Sets

The proposed system was evaluated on three data sets from the UC Irvine Machine Learning Repository [9, 10, 11], which comes with distinguished test sets. Lymphography, a medical imaging technique for visualizing the lymphatic system, essential for immune system functions. The Vehicle Silhouettes dataset distinguishes vehicle types based on shape features extracted from silhouettes. The Car Evaluation dataset classifies car acceptability according to their features as perceived by purchasers. Data characteristics are summarized in Table 1. All training sets were divided into nine local tables, and testing sets were split in half into test and validation sets. The process was conducted in a random and dispersed manner, so each table contained a random subset of objects with all attributes. The classification quality was estimated using the following measures: accuracy, balanced accuracy, weighted precision, and weighted sensitivity. Due to the proposed model's random character, which results in non-deterministic predictions, the test was repeated 10 times for a set of predefined seeds. In each repetition, the process was conducted on a particular seed, starting from learning all local models and splitting the test set into test and validation sets through the conflict analysis model and making the final decision. All proposed approaches to conduct the final decision were compared with variants without coalitions: simple vote and weighted vote. All based models were built using skit-learn library classes, and different sets of parameters were tested. The best results obtained are given in Table 2 and Table 3 as with a set of models parameters, which were different between evaluated approaches. The *max_iter* was always set to 200. The following notations are used: AL – abstract level, U – unified coalitions, D – diverse coalitions, 1 – one strongest coalitions variant, 2 – two strongest coalitions variant, W – with weights for local models. For example, the approach with the two strongest coalitions voting with weights was written in the table as D-2-W.

### Table 1. Data set characteristics

| Data set | # The training set | # The test set | # Conditional attributes | Attributes type | # Decision classes | Source |
|---|---|---|---|---|---|---|
| Vehicle Silhouettes | 592 | 254 | 18 | Integer | 4 | [9] |
| Lymphography | 104 | 44 | 18 | Categorical | 4 | [10] |
| Car evaluation | 1210 | 518 | 6 | Categorical | 4 | [11] |

For the Lymphography data set, the proposed approaches with diverse coalitions present significantly better results than those without coalitions and with unified coalitions, and weighted variants present better results than corresponding

Table 2. Achieved measures of classification quality and local models' parameters

| Dataset & Model | Model Params | Approach | Accuracy±SD | Balanced Accuracy±SD | Precision weighted | Sensitivity weighted |
|---|---|---|---|---|---|---|
| CAR Decision Tree | min_samples_split=2;splitter=best | AL | 0.874 ± 0.009 | **0.634** ± 0.017 | **0.889** | 0.874 |
| | min_samples_split=3;splitter=best | AL-W | 0.870 ± 0.010 | 0.633 ± 0.018 | 0.885 | 0.870 |
| | min_samples_split=2;splitter=best | U-1 | **0.877** ± 0.009 | **0.634** ± 0.016 | 0.891 | **0.877** |
| | min_samples_split=2;splitter=best | U-2 | 0.875 ± 0.010 | 0.623 ± 0.017 | **0.889** | 0.875 |
| | min_samples_split=3;splitter=best | D-1 | 0.753 ± 0.022 | 0.493 ± 0.044 | 0.772 | 0.753 |
| | min_samples_split=3;splitter=best | D-2 | 0.755 ± 0.025 | 0.504 ± 0.059 | 0.773 | 0.755 |
| | min_samples_split=3;splitter=best | U-1-W | 0.876 ± 0.012 | 0.633 ± 0.018 | 0.888 | 0.876 |
| | min_samples_split=3;splitter=best | U-2-W | 0.872 ± 0.012 | 0.622 ± 0.017 | 0.885 | 0.872 |
| | min_samples_split=3;splitter=best | D-1-W | 0.773 ± 0.019 | 0.555 ± 0.033 | 0.805 | 0.773 |
| | min_samples_split=2;splitter=best | D-2-W | 0.780 ± 0.019 | 0.559 ± 0.031 | 0.817 | 0.780 |
| CAR Neural Network | hidden_layer_sizes=(100,);learning_rate=constant | AL | 0.798 ± 0.017 | 0.494 ± 0.035 | 0.775 | 0.798 |
| | hidden_layer_sizes=(100,);learning_rate=constant | AL-W | 0.798 ± 0.017 | **0.500** ± 0.034 | 0.777 | 0.798 |
| | hidden_layer_sizes=(100,);learning_rate=constant | U-1 | 0.797 ± 0.018 | 0.493 ± 0.030 | 0.773 | 0.797 |
| | hidden_layer_sizes=(100,);learning_rate=constant | U-2 | 0.798 ± 0.017 | 0.499 ± 0.035 | 0.775 | 0.798 |
| | hidden_layer_sizes=(100,);learning_rate=constant | D-1 | 0.651 ± 0.029 | 0.484 ± 0.023 | 0.700 | 0.651 |
| | hidden_layer_sizes=(100,);learning_rate=constant | D-2 | 0.667 ± 0.032 | 0.488 ± 0.035 | 0.707 | 0.667 |
| | hidden_layer_sizes=(25, 25);learning_rate=constant | U-1-W | **0.803** ± 0.017 | **0.500** ± 0.057 | **0.787** | **0.803** |
| | hidden_layer_sizes=(100,);learning_rate=constant | U-2-W | 0.795 ± 0.016 | 0.499 ± 0.034 | 0.772 | 0.795 |
| | hidden_layer_sizes=(50, 50);learning_rate=adaptive | D-1-W | 0.635 ± 0.043 | 0.479 ± 0.075 | 0.708 | 0.635 |
| | hidden_layer_sizes=(50, 50);learning_rate=constant | D-2-W | 0.654 ± 0.038 | 0.484 ± 0.076 | 0.724 | 0.654 |
| CAR Logistic Regression | solver=lbfgs;class_weight=balanced | AL | 0.545 ± 0.017 | 0.466 ± 0.023 | 0.655 | 0.545 |
| | solver=lbfgs;class_weight=balanced | AL-W | 0.550 ± 0.018 | 0.475 ± 0.025 | 0.657 | 0.550 |
| | solver=lbfgs;class_weight=balanced | U-1 | 0.551 ± 0.018 | 0.474 ± 0.022 | 0.655 | 0.551 |
| | solver=lbfgs;class_weight=balanced | U-2 | 0.548 ± 0.020 | 0.464 ± 0.023 | 0.655 | 0.548 |
| | solver=liblinear;class_weight=balanced | D-1 | **0.581** ± 0.025 | 0.425 ± 0.045 | 0.636 | **0.581** |
| | solver=lbfgs;class_weight=balanced | D-2 | 0.429 ± 0.018 | 0.412 ± 0.045 | 0.670 | 0.429 |
| | solver=lbfgs;class_weight=balanced | U-1-W | 0.555 ± 0.019 | **0.485** ± 0.025 | 0.657 | 0.555 |
| | solver=lbfgs;class_weight=balanced | U-2-W | 0.551 ± 0.020 | 0.473 ± 0.024 | 0.655 | 0.551 |
| | solver=lbfgs;class_weight=balanced | D-1-W | 0.467 ± 0.024 | 0.465 ± 0.049 | 0.674 | 0.467 |
| | solver=lbfgs;class_weight=balanced | D-2-W | 0.468 ± 0.022 | 0.474 ± 0.053 | **0.678** | 0.468 |
| LYMPHO-GRAPHY Decision Tree | min_samples_split=2;splitter=best | AL | 0.455 ± 0.000 | 0.333 ± 0.000 | 0.455 | 0.455 |
| | min_samples_split=2;splitter=best | AL-W | 0.455 ± 0.000 | 0.333 ± 0.000 | 0.455 | 0.455 |
| | min_samples_split=2;splitter=best | U-1 | 0.455 ± 0.000 | 0.333 ± 0.000 | 0.455 | 0.455 |
| | min_samples_split=2;splitter=best | U-2 | 0.455 ± 0.000 | 0.333 ± 0.000 | 0.455 | 0.455 |
| | min_samples_split=4;splitter=best | D-1 | 0.345 ± 0.055 | 0.494 ± 0.140 | 0.609 | 0.345 |
| | min_samples_split=4;splitter=random | D-2 | 0.341 ± 0.105 | 0.469 ± 0.185 | 0.414 | 0.341 |
| | min_samples_split=2;splitter=best | U-1-W | 0.455 ± 0.000 | 0.333 ± 0.000 | 0.455 | 0.455 |
| | min_samples_split=2;splitter=best | U-2-W | 0.455 ± 0.000 | 0.333 ± 0.000 | 0.455 | 0.455 |
| | min_samples_split=2;splitter=best | D-1-W | **0.730** ± 0.046 | **0.816** ± 0.032 | **0.732** | **0.730** |
| | min_samples_split=3;splitter=random | D-2-W | 0.536 ± 0.094 | 0.701 ± 0.054 | 0.549 | 0.536 |
| LYMPHO-GRAPHY Neural Network | hidden_layer_sizes=(100,);learning_rate=adaptive | AL | 0.280 ± 0.182 | 0.517 ± 0.130 | 0.405 | 0.280 |
| | hidden_layer_sizes=(100,);learning_rate=constant | AL-W | 0.582 ± 0.130 | 0.409 ± 0.078 | 0.584 | 0.582 |
| | hidden_layer_sizes=(25,25);learning_rate=adaptive | U-1 | 0.393 ± 0.167 | 0.538 ± 0.128 | 0.488 | 0.393 |
| | hidden_layer_sizes=(25,25);learning_rate=adaptive | U-2 | 0.245 ± 0.223 | 0.462 ± 0.154 | 0.554 | 0.245 |
| | hidden_layer_sizes=(25,25);learning_rate=adaptive | D-1 | 0.289 ± 0.057 | 0.357 ± 0.148 | 0.552 | 0.289 |
| | hidden_layer_sizes=(50,50);learning_rate=adaptive | D-2 | 0.211 ± 0.129 | 0.305 ± 0.223 | 0.544 | 0.211 |
| | hidden_layer_sizes=(100,);learning_rate=constant | U-1-W | 0.580 ± 0.129 | 0.439 ± 0.141 | 0.581 | 0.580 |
| | hidden_layer_sizes=(100,);learning_rate=constant | U-2-W | 0.550 ± 0.128 | **0.550** ± 0.173 | 0.556 | 0.550 |
| | hidden_layer_sizes=(100,);learning_rate=constant | D-1-W | 0.627 ± 0.101 | 0.489 ± 0.148 | 0.648 | 0.627 |
| | hidden_layer_sizes=(100,);learning_rate=constant | D-2-W | **0.664** ± 0.091 | 0.518 ± 0.129 | **0.679** | **0.664** |

unweighted ones. The best results were achieved with the one strongest coalition with weights approach, while all approaches with unified coalitions demonstrated no improvements. On the other hand, for the Vehicle Silhouettes and Car evaluation data sets, approaches with unified coalitions enhanced classification quality, while with diverse coalitions, all measures were worse. The variant with the strongest unified coalition was the most suitable for these sets.

Table 3. Achieved measures of classification quality and local models' parameters

| Dataset & Model | Model Params | Approach | Accuracy±SD | Balanced Accuracy±SD | Precision weighted | Sensitivity weighted |
|---|---|---|---|---|---|---|
| LYMPHO-GRAPHY Logistic Regression | solver=liblinear;class_weight=balanced | AL | 0.177 ± 0.038 | 0.447 ± 0.028 | 0.658 | 0.177 |
| | solver=lbfgs;class_weight=balanced | AL-W | 0.461 ± 0.010 | 0.433 ± 0.153 | 0.465 | 0.461 |
| | solver=saga;class_weight=balanced | U-1 | 0.048 ± 0.075 | 0.352 ± 0.055 | 0.082 | 0.048 |
| | solver=liblinear;class_weight=balanced | U-2 | 0.105 ± 0.023 | 0.393 ± 0.017 | **0.919** | 0.105 |
| | solver=saga;class_weight=balanced | D-1 | 0.298 ± 0.068 | 0.364 ± 0.170 | 0.562 | 0.298 |
| | solver=saga;class_weight=balanced | D-2 | 0.305 ± 0.083 | 0.400 ± 0.161 | 0.598 | 0.305 |
| | solver=sag;class_weight=balanced | U-1-W | 0.484 ± 0.046 | 0.638 ± 0.105 | 0.517 | 0.484 |
| | solver=lbfgs;class_weight=balanced | U-2-W | 0.477 ± 0.000 | **0.667 ± 0.000** | 0.491 | 0.477 |
| | solver=liblinear;class_weight=balanced | D-1-W | **0.707 ± 0.046** | 0.482 ± 0.032 | 0.711 | **0.707** |
| | solver=liblinear;class_weight=balanced | D-2-W | **0.707 ± 0.046** | 0.482 ± 0.032 | 0.711 | **0.707** |
| VEHICLE Decision Tree | min_samples_split=3;splitter=random | AL | 0.708 ± 0.017 | 0.689 ± 0.017 | 0.703 | 0.708 |
| | min_samples_split=3;splitter=random | AL-W | 0.705 ± 0.030 | 0.685 ± 0.029 | 0.700 | 0.705 |
| | min_samples_split=3;splitter=random | U-1 | 0.710 ± 0.022 | 0.691 ± 0.021 | 0.705 | 0.710 |
| | min_samples_split=3;splitter=random | U-2 | **0.712 ± 0.025** | **0.693 ± 0.021** | 0.707 | **0.712** |
| | min_samples_split=3;splitter=best | D-1 | 0.520 ± 0.046 | 0.517 ± 0.047 | 0.575 | 0.520 |
| | min_samples_split=3;splitter=best | D-2 | 0.543 ± 0.050 | 0.537 ± 0.050 | 0.600 | 0.543 |
| | min_samples_split=3;splitter=random | U-1-W | 0.710 ± 0.029 | 0.691 ± 0.029 | **0.709** | 0.710 |
| | min_samples_split=3;splitter=random | U-2-W | 0.707 ± 0.030 | 0.688 ± 0.028 | 0.704 | 0.707 |
| | min_samples_split=3;splitter=best | D-1-W | 0.513 ± 0.048 | 0.508 ± 0.047 | 0.575 | 0.513 |
| | min_samples_split=3;splitter=best | D-2-W | 0.539 ± 0.042 | 0.531 ± 0.041 | 0.595 | 0.539 |
| VEHICLE Neural Network | hidden_layer_sizes=(50,50);learning_rate=constant | AL | 0.703 ± 0.028 | **0.684 ± 0.028** | **0.696** | 0.703 |
| | hidden_layer_sizes=(100,);learning_rate=adaptive | AL-W | **0.704 ± 0.015** | **0.684 ± 0.016** | 0.690 | **0.704** |
| | hidden_layer_sizes=(50,50);learning_rate=constant | U-1 | 0.699 ± 0.028 | 0.681 ± 0.028 | 0.695 | 0.699 |
| | hidden_layer_sizes=(50,50);learning_rate=constant | U-2 | 0.695 ± 0.030 | 0.678 ± 0.030 | 0.691 | 0.695 |
| | hidden_layer_sizes=(50,50);learning_rate=constant | D-1 | 0.468 ± 0.099 | 0.457 ± 0.093 | 0.500 | 0.468 |
| | hidden_layer_sizes=(50,50);learning_rate=adaptive | D-2 | 0.457 ± 0.059 | 0.451 ± 0.058 | 0.477 | 0.457 |
| | hidden_layer_sizes=(100,);learning_rate=adaptive | U-1-W | 0.702 ± 0.013 | 0.682 ± 0.013 | 0.689 | 0.702 |
| | hidden_layer_sizes=(100,);learning_rate=adaptive | U-2-W | 0.698 ± 0.020 | 0.677 ± 0.020 | 0.683 | 0.698 |
| | hidden_layer_sizes=(100,);learning_rate=constant | D-1-W | 0.583 ± 0.048 | 0.567 ± 0.045 | 0.591 | 0.583 |
| | hidden_layer_sizes=(100,);learning_rate=adaptive | D-2-W | 0.588 ± 0.051 | 0.570 ± 0.048 | 0.601 | 0.588 |
| VEHICLE Logistic Regression | solver=liblinear;class_weight=balanced | AL | 0.751 ± 0.024 | 0.728 ± 0.028 | 0.735 | 0.751 |
| | solver=liblinear;class_weight=balanced | AL-W | 0.749 ± 0.022 | 0.728 ± 0.025 | 0.734 | 0.749 |
| | solver=liblinear;class_weight=balanced | U-1 | **0.761 ± 0.023** | **0.737 ± 0.025** | **0.746** | **0.761** |
| | solver=liblinear;class_weight=balanced | U-2 | 0.753 ± 0.023 | 0.731 ± 0.027 | 0.739 | 0.753 |
| | solver=liblinear;class_weight=balanced | D-1 | 0.603 ± 0.058 | 0.591 ± 0.060 | 0.635 | 0.603 |
| | solver=liblinear;class_weight=balanced | D-2 | 0.606 ± 0.042 | 0.595 ± 0.045 | 0.632 | 0.606 |
| | solver=liblinear;class_weight=balanced | U-1-W | 0.756 ± 0.025 | 0.734 ± 0.029 | 0.743 | 0.756 |
| | solver=liblinear;class_weight=balanced | U-2-W | 0.748 ± 0.021 | 0.727 ± 0.025 | 0.734 | 0.748 |
| | solver=liblinear;class_weight=balanced | D-1-W | 0.654 ± 0.035 | 0.638 ± 0.038 | 0.664 | 0.654 |
| | solver=liblinear;class_weight=balanced | D-2-W | 0.658 ± 0.037 | 0.643 ± 0.039 | 0.666 | 0.658 |

# 4. Conclusions

The paper proposes a classification model using coalitions of local models. The novelty of the work is the combination of two ways of distinguishing coalitions of local models and approaches, including weighted variants in conjunction with decision trees, neural networks, and logistic regression. The preliminary results included in this work show that the proposed approaches improve classification quality in terms of accuracy, precision, and sensitivity. The introduced variant of diverse coalition presents better results for the Lymphography data set, and mostly, the weighted approaches are superior to the unweighted ones. More extensive experiments are planned for future work involving more data sets and different variants of local classifiers and sets of parameters.

# References

[1] Pawlak, Z. Some remarks on conflict analysis. *European Journal of Operational Research*, 166:649–654, 2005.

[2] Pawlak, Z. Conflict analysis. *Proceedings of the Fifth European Congress on Intelligent Techniques and Soft Computing*, pages 1589–1591, 1997.

[3] Czarnowski, I. and Jędrzejowicz, P. Ensemble online classifier based on the one-class base classifiers for mining data streams. *Cybernetics and Systems*, 46(1-2):51–68, 2015.

[4] Verbraeken, J., Wolting, M., Katzy, J., Kloppenburg, J., Verbelen, T., and Rellermeyer, J. S. A survey on distributed machine learning. *ACM Computing Surveys*, 53(2):1–33, 2020.

[5] Tasci, E., Uluturk, C., and Ugur, A. A voting-based ensemble deep learning method focusing on image augmentation and preprocessing variations for tuberculosis detection. *Neural Computing and Applications*, 33(22):15541–15555, 2021.

[6] Przybyła-Kasperek, M. and Sacewicz, J. Ensembles of random trees with coalitions-a classification model for dispersed data. *Procedia Computer Science*, 246:1599–1608, 2024.

[7] Przybyła-Kasperek, M. and Wakulicz-Deja, A. A dispersed decision-making system – the use of negotiations during the dynamic generation of a system's structure. *Information Sciences*, 288:194–219, 2014.

[8] Gillani, Z. and a Saira Aquil, Z. B. A game theoretic conflict analysis model with linguistic assessments and two levels of game play. *Information Sciences*, 677:120840, 2024. ISSN 0020-0255.

[9] Mowforth, P. and Shepherd, B. Statlog (vehicle silhouettes) [dataset]. UCI Machine Learning Repository. doi:10.24432/C5HG6N.

[10] Zwitter, M. and Soklic, M. Lymphography. UCI Machine Learning Repository, 1988. doi:10.24432/C54598.

[11] Bohanec, M. Car evaluation [dataset]. UCI Machine Learning Repository, 1997. doi:10.24432/C5JP48.

# Multi-Agent Simulation of Financial Markets: Micro and Macro Approach

**Jakub Skrzyński**[0009−0008−4550−5009]

*AGH University of Krakow*
*Department of Applied Computer Science*
*al. A. Mickiewicza 30, 30-059 Kraków, Poland*
*research@jakub.skrzynski.net*

**Abstract.** *Financial markets are complex systems driven by interactions among multiple entities. Traditional modeling struggles to capture their emergent behaviors. This study presents a multi-agent simulation framework for financial markets, incorporating both macro- and micro-level perspectives. The system models investors, brokers, and stock exchanges, enabling realistic order matching and market interactions. A special "imaginary investor" ensures market stability. Implemented in Java with JADE, the framework allows customizable experiments. A test scenario demonstrated that a simple investment strategy mitigates losses and can generate profit, validating the simulation's effectiveness in replicating market dynamics.*
**Keywords:** *multi agent systems, simulation, complex systems, financial markets simulation*

## 1. Introduction

Financial markets are an environment that is very hard to observe and capture within a conventional modeling techniques without losing any of their behaviors. That is because of a very complex net of interactions between participating entities. The structure of markets makes it almost natural to model them as the multi agent system, which is reflected also in the reviewed literature [1, 2, 3, 4, 5]. Such an approach allows to capture all of the characteristics that emerge from the markets' complex interactions and would be hard to capture otherwise.

It is also worth mentioning that most of the publications concentrate on the issues of simulating the agents' decision making process, exploring the broad field of

search for the best strategy, which is dominated by RL approaches but not limited to them [6], whereas there is very little said about the research towards more accurate environment, which is explored in this paper. The main goal of the presented research is to accurately simulate the environment without unnecessary overhead.

There may be different aspects of interest regarding the observation of the stock exchange markets. The first one is the macroscopic approach. Then the researcher would concentrate mostly on the global outcomes and changes introduced globally, such as, for example, regulations or reactions to the so-called "black swans." On the other hand, there is a possibility of observing the microscopic outcomes, namely how the market affects the singular investor or a financial company. In such a case, it is important to observe how market movements would affect the portfolio and the investor's behavior. This means that a detailed simulation of the market may be disregarded and replaced with hidden actions. Both of the approaches mentioned may provide valuable insights on the behavior of modern financial markets.

In order to allow such observations to be made, a multi-agent simulation of financial market was prepared, which is designed for both scenarios. The simulation framework allows for a high level of customization, leaving a wide variety of possibilities for the user apart from experiments presented here.

## 2. Simulation Design

The theoretical model of the agent system was based on rules taken from GPW (Warsaw Stock Exchange) [7, 8, 9, 10, 11] and NASDAQ. The system was created according to guidelines provided by [12] on preparing simulations. Moving to the details, the simulation framework was designed to support multiple types of interacting agents. Those types include:

- **Investors** whose most important characteristics is the amount of money and stocks in their wallet.

- **Brokers** that simulate stock brokers. This type of agents is designed to run the account for the investor, store the balance, keep track of transactions and perform the accounting. The most important feature of the broker is forwarding the market orders to the stock exchange.

- **Stock Exchange** is designed to play the role of stock exchange in real markets. The stock exchange runs the order sheet and matches the orders.

The stock exchange is also capable of distributing the market information, such as prices, or best offers.

The whole system is based on messages that agents send between each other. The allowed communication routes are the following: investor to broker, investor to stock exchange via broker, stock exchange to broker. The interactions in the system are defined as follows:

- The investor may order any broker to open an account for him. The newly created account is by default empty.

- The investor may deposit or withdraw money to/from the broker that is running the investor's account, given that the balance allows for such an action.

- The investor may check the balance of the account and retrieve the stored money as well as the owned stocks.

- The investor may ask the broker for the current best offers. It is possible that, depending on the circumstances and configuration, the broker may charge the investor a fee for extended information.

- The investor may place an order. The order is sent via broker, to a selected stock exchange that is supported by the given broker. Before forwarding the request, the broker checks the balance and amount of stock in the investor's account, places a lock on resources subjected to the order and then forwards the order to the stock exchange.

- Stock exchange after matching transaction creates the summary and confirmation of the transaction and sends it back to the broker. The broker should perform accounting based on the received information, namely it should unlock the resources, or add new resources to account, or consume the locked resources. The investor, similarly to the real-life scenario, does not need to be notified.

- Stock exchange may send out the information that the market order was canceled due to reaching its expiration date. The broker after receiving such a notice should immediately unlock the resources blocked by that order.

Another important characteristic of the system is the algorithm for matching the incoming orders. The algorithm was designed based on GPW order

matching rules. The system does support the following types of orders: with a limit on the price; without a limit on the price; with a limit on the price that activates only after the stock price reaches a specified level; without a limit activated after the price surpasses the threshold. These are the most basic orders implemented by most of financial markets. The routine of matching the orders is as follows. The order sheet maintains the list of active limitless orders (with separation to buy/sell), the list of active orders with a specified price limit (with separation to buy/sell) and the list of orders that are awaiting to be activated by the price reaching the specified value. The incoming order is first subjected to its counter-type limitless requests, sorted according to the time of receiving the order (in FIFO queue), if after that phase there are remaining stocks to be bought or sold in the order, i.e. it was not fully executed, it is checked against the orders with a limit, the list is sorted first according to the limit value and secondly by the time it has been registered (in FIFO system). If the order still remains incomplete, it is saved in a correct queue according to its type. In the case of matching two limitless orders with each other, the price is set by the exchange to be the price of the last transaction.

In order to allow for the limited simulation of micro approach where a single chain investor – broker – exchange is of researcher interest, the imaginary investor was introduced. The imaginary investor is an entity that plays a role of a huge investor with unlimited resources. When placing orders the imaginary investor always follows supplied data. The data may contain historical prices of the given stock. The imaginary agent at the beginning of each session places the orders: buy with a limited price and sell with a limited price. The price should be a little lower than the data point for buying, and a little higher than the data point for sale. The amount of stocks to buy and to sell is regulated by the parameter described as the imaginary agent's power. Setting it to a low value may result in slight pulses that bend the behavior of the simulated market towards the supplied data. In the case of high value, the actions of investors are not affecting the behavior of the market.

The system was implemented using Java language and JADE framework [13], taking advantage of the ACL messages system as the main media for communication between agents.

# 3. The Experiment

As to present the possibilities of the designed framework, a simple scenario was adopted of an agent who follows the advised principles for beginners, to check

if the strategy indeed prevents huge losses. The investor agent is supplied with a starting amount of money and decides based on own algorithm when to buy and sell the goods. The algorithm is based on the historical price records stored by the agent. Namely, when the current price is significantly lower than the computed historical average, the agent attempts to buy the stock for at most a few percent more than the current best offer. If the price is higher than average, he decides to sell the stock. The amount of stock to be sold is determined by a formula parametrized by the level of risk tolerance.

The outcome of the experiment is presented in Figure 1. After several trials, it was found that the method does indeed prevent losses and, in some cases, when the agent accepts a medium level of risk, it can yield a considerable profit.



Figure 1. Results of the simulation. The plot shows green dots when the agent bought stock, and red dots when the agent sold stock. The orange line and the right scale show the historical stock price. The blue line gives the money balance. The green line denotes the value of owned stock. The red line is the total value of the wallet. The magenta line denotes the starting value.

# 4. Conclusions

The presented framework allows for replicating the real life market dynamics, making it possible to conduct experiments and reasoning about hypothetical situations regarding the stock market. The design of the simulation was proven experimentally to be valid. The use of the framework may help to understand the complex processes that rule the market and the emergent behaviors shaped by atomic actions of agents.

Future plans regarding the research and further development include extending the simulation by introducing more complex structures such as investment funds, CFD contracts and other complex financial instruments, to allow for even more realistic behaviors to be observed and researched using the framework. As the framework was designed in a flexible manner, it gives a huge potential for further extensions. In terms of research it is planned to examine carefully how brokers may influence the market by introducing pricing policies and how credit availability impacts the markets.

Current implementation with all of the details allowing for recreation of the presented results is available in the GitHub repository via URL: `https://github.com/j-skrzynski/MASfSES-JADE`

# References

[1] Raberto, M., Cincotti, S., Focardi, S. M., and Marchesi, M. Agent-based simulation of a financial market. *Physica A: Statistical Mechanics and its Applications*, 299(1):319–327, 2001. doi:10.1016/S0378-4371(01)00312-0.

[2] Jacobs, B. I., Levy, K. N., and Markowitz, H. M. Financial market simulation. *The Journal of Portfolio Management*, 30(5):142–152, 2004.

[3] Maeda, I., DeGraw, D., Kitano, M., Matsushima, H., Sakaji, H., Izumi, K., and Kato, A. Deep reinforcement learning in agent based financial market simulation. *Journal of Risk and Financial Management*, 13(4):71, 2020.

[4] Liu, X.-Y., Xia, Z., Rui, J., Gao, J., Yang, H., Zhu, M., Wang, C., Wang, Z., and Guo, J. FinRL-Meta: Market environments and benchmarks for data-driven financial reinforcement learning. *Advances in Neural Information Processing Systems*, 35:1835–1849, 2022.

[5] Lussange, J., Lazarevich, I., Bourgeois-Gironde, S., Palminteri, S., and Gutkin, B. Modelling stock markets by multi-agent reinforcement learning. *Computational Economics*, 57(1):113–147, 2021.

[6] Yang, Y., Zhang, Y., Wu, M., Zhang, K., Zhang, Y., Yu, H., Hu, Y., and Wang, B. TwinMarket: A scalable behavioral and social simulation for financial markets. *arXiv preprint arXiv:2502.01506*, 2025.

[7] Ziębiec, J. *Zasady obrotu giełdowego*. Oficjalne wydawnictwo Giełdy Papierów Wartościowych w Warszawie, Warszawa, 5 edition, 2009. ISBN 978-83-60510-15-5. URL `https://www.gpw.pl/pub/images/prezentacje/system_obrotu.pdf`.

[8] Biernacki, P. and Szulec, P. *Pierwsze kroki na rynku kapitałowym*. Komisja Nadzoru Finansowego, 4 edition, 2009. ISBN 978-83-924813-9-3. URL `https://www.knf.gov.pl/knf/pl/komponenty/img/Pierwsze%20kroki_17589.pdf`.

[9] Kachniewski, M., Majewski, B., and Wasilewski, P. *Giełda Papierów Wartościowych i rynek kapitałowy*. Fundacja Edukacji i Rynku Kapitałowego, Warszawa, 2008.

[10] Szczegółowe zasady obrotu giełdowego w systemie UTP, 2025. URL `https://www.gpw.pl/regulacje-prawne`.

[11] Regulamin giełdy, 2024. URL `https://www.gpw.pl/regulacje-prawne`.

[12] Sayama, H. *Introduction to the Modeling and Analysis of Complex Systems*. Open SUNY Textbooks. Open Suny Textbooks, 2015. ISBN 9781942341093. URL `https://books.google.pl/books?id=Bf9gAQAACAAJ`.

[13] Bellifemine, F., Bergenti, F., Caire, G., and Poggi, A. Jade – a java agent development framework. In R. H. Bordini, M. Dastani, J. Dix, and A. El Fallah Seghrouchni, editors, *Multi-Agent Programming: Languages, Platforms and Applications*, pages 125–147. Springer US, Boston, MA, 2005. ISBN 978-0-387-26350-2. doi:10.1007/0-387-26350-0_5.

# Lessons Learned from Employing Monte Carlo Tree Search to Train Rule-Based Models for Intelligent Game Agents

**Maciej Świechowski**[0000−0002−8941−3199]

*QED Games*
*Miedziana 3A, 00-814 Warsaw, Poland*
*maciej.swiechowski@qed.pl*

**Abstract.** *This article presents a new workflow for creating AI agents, based on a combination of Monte Carlo Tree Search, which acts as an offline teacher, and rule-based machine learning models (e.g., decision trees or random forests) for real-time use. We demonstrate that this approach is promising and can result in the development of strong agents – in some cases even stronger than the teacher – although it has certain limitations. These limitations, along with a summary of key lessons learned from empirical experiments, are discussed.*
**Keywords:** *Monte Carlo tree search, video game agents, simulations*

## 1. Introduction

Monte Carlo Tree Search (MCTS) was first popularized for the game of *Go* in 2006 [1] and has since become one of the state-of-the-art methods for combinatorial games and other discrete optimization problems [2]. Among its many successful applications, it is noteworthy to mention its integral role in the *AlphaGo* approach [3]. However, the application of the MCTS algorithm to real-time games and, more broadly, commercial video games is very rare. One of the primary reasons for this is the computational cost. In commercial real-time games, the AI component is typically allowed only 3–16 milliseconds per game frame [4].

The inspiration for our research stemmed from an observation that the requirement for responsiveness is exclusive to actual live gameplay involving human players. During the game development process, we can deliberately reduce the in-game

clock speed, allowing bots to play at a slower pace. Consequently, in-game events occur more slowly, providing the bots with significantly more real-time to respond with decision, e.g., even 100 times more time. This setup enables the MCTS algorithm to generate high-quality data, which can then be utilized to train or transform into a different model optimized for real-time performance. Such a model would be incorporated into the final version of the game. This approach presents a highly promising workflow for developing computer agents. We report a summary on the lessons learned from a preliminary research of this methodology.

## 2. Experimental Setup

Apart from well-known combinatorial games, we used two custom game environments in our experiments. Since the findings presented in this paper are based on these custom games, we summarize their nature, state spaces, and selected training-related parameters in Table 1.

The initial idea for this research was to generate data using MCTS and train decision trees due to their explainability. However, due to weak performance in RTS, we also included the Random Forest technique.

## 3. Summary of the Findings

### 3.1. State Vectorization

The MCTS algorithm [2] does not require a specific state representation as long as the execution of simulations is possible, i.e., the representation implements a forward model of the game. In this context, states are associated with nodes within the game tree. During the decision-making process, the MCTS algorithm executes $N$ iterations, where the greater the $N$, the more likely the decision will be precise (stronger). Each iteration encompasses four phases: (1) *Selection*, (2) *Expansion*, (3) *Simulation (Rollout)*, and (4) *Backpropagation*. The *Selection* and *Expansion* phases focus on navigating and growing the part of the game tree stored in memory. The *Simulation* and *Backpropagation* phases are responsible for collecting and updating the statistics, respectively, related to the visited states and actions.

However, to utilize a machine learning model such as a decision tree, it is necessary to define some state features (attributes) based on which decisions will be

Table 1. Custom game environments: *Shooter* and *RTS*. (*) denotes experiments performed outside the grid search, using fixed best values for other parameters.

| Property | Simple Shooter | RTS |
|---|---|---|
| Types of actions | move (N,W,S,E), shoot | buy, allocate, wait |
| Enemy | moves, rotates, attacks | single-player game |
| Player goal | hit enemy *TP* times | maximize army strength |
| Max game length (steps) | 1000 time steps (*t*) | 40 time steps (*t*) |
| When max steps achieved | player loss | game scoring |
| Also loss when | hit by enemy once | - |
| Multiple actions in a turn | no & no 2 attacks in a row | yes |
| Avg. branching factor (BF) | 5 | 12 |
| Avg. BF in 1 and 20 steps | 5 and 5 | 2 and 17 |
| State representation, $s \in$ | $\mathbb{R}^2 \times \mathbb{R}^2 \times SO(2) \times \mathbb{N}$. | $\mathbb{N}^{37}$ |
| State includes *t* and | 2D positions, rotation | units, buildings, resources |
| # unique states (l. bound) | $> 5 * 10^{13}$ | $> 10^{15}$ |
| Best perf. vectorization | [distance, enemy rotation relative angle with enemy] | All variables except *t* |
| Trained model type | Decision tree | Random forest |
| Good strategy | kiting & avoid bad rotation | optimize order & timing |
| **Hyperparameters tested** | | |
| # dataset gen. methods | 5 | 5 |
| # training methods | 4 | 4 |
| # vectorizations | 110 (bins for angles) | 4 (lengths: 11, 14, 25, 36) |
| # starting states params. | 7 | 7 |
| # total combinations | 15400 | 560 |
| # repeats per combination | 200 | 200 |
| # max game lengths (*) | 10 | 8 |

made. Therefore, our approach requires a state vectorization function that converts arbitrary states, used during MCTS iterations, into vectors of features:

$$Vec(s) \rightarrow \mathbb{R}^k; \; s \in \text{States}. \qquad (1)$$

The selection of features and the design of the vectorization function are critical to the successful implementation of this approach. Firstly, the number of features should be minimized. They should cover important aspects from a decision-making perspective, but since they are not used for game simulations, they can represent a much smaller subset of all possible features. Secondly, these features should encapsulate general concepts, which have numerical representations in many game states. For instance, a "health value of a unit" is a robust feature,

whereas "placement of a unit at a specific position on a grid" is typically a weak feature, because it is less likely to be applicable across many states. Thirdly, using fewer and more generalized features allows for more substantial support in evidence (in the form of the MCTS trees) and a stronger training signal for them.

## 3.2. Regression vs. Classification

A research question in our study was whether the problem should be represented as a regression or classification. The MCTS algorithm computes the $Q$ values, which represent the average outcome of the game for states and state-action pairs, suggesting a regression approach. However, experiments conducted across various games demonstrated that formulating this problem as classification yields better results. The target label then becomes the decision to make under specific dependent values – the vectorized features. The primary reason why framing the problem as a regression led to suboptimal results is that the $Q$ values tend to converge to game-theoretical extremes representing outcomes such as a win or loss. It is still desirable for the AI agent to execute generally strong actions, even in losing positions, which is easier to achieve having a more nuanced training signal.

The most effective strategy was to aggregate decision maximizing the $Q$ values for each unique vectorization and add them once per fully played game (after it has been finished). More frequent additions to the dataset (e.g., after each state) resulted in poorer performance of the ultimately trained agents, while less frequent additions (e.g., averaged over $N$ games) did not enhance performance further.

## 3.3. Diversity of Training and Test Examples

In games, the entire state space is often so large that a thorough search would be infeasible. For MCTS, however, this is not a problem, because it is a dynamic algorithm that performs searches locally around the current state. For models that are trained and then used as is, no dynamic calculations are possible anymore. One of the main findings of our study was the effects of diverse training and testing examples. Fortunately, training data do not need to capture every state of the game, which would be infeasible. If possible, states should be grouped into clusters – classes of abstractions induced by the optimal decision in them. Ideally, at least one example from each such class should be presented during training. The test data should contain various examples to determine whether the trained model has overfitted to the training data. Table 2 presents the scores obtained for different diversities of training and test data for two games.

Table 2. Average scores (0: min, 1: max) obtained in performed experiments by agents trained with various numbers of diverse scenarios for two games. Scenarios are initial settings allowing to explore various subspaces of state space. For example, in the *shooter game*, the starting positions of the enemies were sampled uniformly from a circle, with the player agent positioned in the center.

| Game | Model | Score per unique scenarios in train data | | | | | | |
|------|-------|------|------|------|------|------|------|------|
| | | **1** | **4** | **8** | **16** | **32** | **64** | **128** |
| Shooter | Decision Tree | 0.24 | 0.52 | 0.7 | 0.76 | 0.85 | 0.85 | 0.91 |
| RTS | Random Forest | 0.00 | 0.00 | 0.00 | 0.83 | 0.99 | 1.00 | 1.00 |

## 3.4. Strategic Complexity of Games and Generalization

We have found that, apart from training and testing using diverse examples, the main condition for the trained decision models to achieve satisfactory efficacy of play, i.e., similar to the trainer MCTS algorithm, is that the game features universally good strategies (general patterns). In contrast, games with high tactical complexity, where unique nuances play a significant role, are not suitable for the proposed approach. However, we also observed a unique advantage of the rule-based models. In games where they are effective, they can generalize better when some of the game parameters are scaled. The MCTS algorithm explicitly traverses the game tree, and it is affected by its length and branching factor. For example, consider a game with a map that has established strategies. Doubling the map size does not change how people play it. However, such an increase creates many new positions for the agents to be in, and the MCTS-based agent might suddenly require $2^2$ times more iterations per decision to remain effective. A rule-based tree agent will continue to make the same decisions, regardless of the map size. After analyzing many games, we can summarize that if the trained agents generalize well for a game, then they are much less affected than search-based methods by scaling of the game. In Table 3, we show the scores obtained by the trained agents for the *shooter game* depending on the *target points* (TP) parameters, which define how many points are required for a player to win, thus affecting the game's length. Increasing *TP* does not affect the optimal strategy, i.e., the player still has to kite the enemy while avoiding being hit when it is rotated towards the player. It just takes more time to repeat this "hit-and-run" pattern.

Table 3. Game scores achieved (in *shooter game*) by different methods depending on the target points (TP) parameter representing game length. The last column represents a model trained on games with the same TP as in the first column.

| Target Points (TP) | MCTS 100K iterations | Trained Decision Tree | |
|---|---|---|---|
| | | with fixed TP=40 | with actual TP |
| 24 | **1.00** | 0.85 | 0.91 |
| 32 | **0.99** | 0.85 | 0.88 |
| 40 | 0.83 | **0.84** | **0.84** |
| 46 | 0.39 | **0.83** | 0.10 |
| 50 | 0.01 | **0.83** | 0.00 |
| 80 | 0.00 | **0.83** | 0.00 |

## 3.5. Run-time Performance

It was possible to train rule-based agents based on decision-trees or random forests that perform at least as effectively as MCTS-based players in only some of the game environments, specifically simple *shooter game* and *RTS*. However, in both cases, a single inference from the model took $\leq$ 1ms, whereas the equivalent MCTS algorithm, conducting 100K simulations per decision, took 792ms and 1,048ms respectively. The measurements of these values were averaged over 2,000 repetitions, providing high statistical confidence. The rule-based approach is not always applicable due to quality requirements, but when it is, it is 2–3 orders of magnitude faster.

## 3.6. Explainability

The decisions made by the MCTS algorithm are the result of extensive number of simulations. Typically, these decisions are difficult to explain and interpret based solely on the statistics. Although rule-based models are known for being some of the most explainable and interpretable machine learning models, we have found that there is a trade-off between the quality of decisions made by the trained model and its explainability, especially in the context of practical applications in game development. The decision trees that produced decisions of similar quality to those of the MCTS algorithm, which conducted 100K iterations per decision, contained over 1000 nodes. This complexity is typically beyond human cognitive

abilities to explain the model globally. However, case-based explanations based on an analysis of a decision path were still possible.

## 4. Conclusions

MCTS, although extremely popular in the research of combinatorial games, has not been widely utilized in real-time games or within the game development industry. We have presented a workflow with the potential to change this and a summary of key lessons learned so far. The MCTS algorithm is employed offline as a training method for a faster run-time model. Our models of choice were rule-based models, which helps to explain bots' decisions – a crucial feature in the game development industry. However, in future work, it would be worthwhile to investigate more complex machine learning models, such as extreme gradient boosting or neural networks.

## Acknowledgment

## References

[1] Coulom, R. Efficient selectivity and backup operators in Monte Carlo tree search. In H. J. van den Herik, P. Ciancarini, and H. H. L. M. J. Donkers, editors, *Int. Conference on Computers and Games*, pages 72–83. Springer, Springer Berlin Heidelberg, 2006. doi:10.1007/978-3-540-75538-8\_7.

[2] Świechowski, M., Godlewski, K., Sawicki, B., and Mańdziuk, J. Monte Carlo tree search: A review of recent modifications and applications. *Artificial Intelligence Review*, 56:2497–2562, 2023. doi:10.1007/s10462-022-10228-y.

[3] Silver, D. et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. doi:10.1038/nature16961.

[4] Rabin, S. *Game AI Pro 2: Collected Wisdom of Game AI Professionals*. AK Peters / CRC Press, 2015.

CHAPTER 13

# Remote Sensing and Satellite Data Analysis

Track Chairs:

- prof. Jakub Nalepa – Silesian University of Technology

- prof. Przemysław Biecek – Warsaw University of Technology

- Krzysztof Kotowski, PhD, Eng – KP Labs

# A Simple and Efficient Method for GPS-less Drone Navigation Using Visual Cues

**Piotr Zacholski, Dominik Pieczyński**[0000−0003−0275−5629],
**Marek Kraft**[0000−0001−6483−2357]

*Poznań University of Technology*
*Institute of Robotics and Machine Intelligence*
*Piotrowo 3A, 60-965 Poznań, Poland*
*marek.kraft.@put.poznan.pl*

**Abstract.** *Accurate localization on the map is a key feature of geolocalization and navigation of aerial robots. In this preliminary study, we compare two localization strategies based on Local Binary Patterns (LBP): (1) random initialization subsequently processed with particle filter and (2) same as (1), but initialized using a place recognition neural network trained with metric learning. A range of peace recognition models based on ResNet were evaluated to select the highest performing model. The results clearly show that the neural network-based initialization results in improved localization accuracy and faster particle filter convergence as compared to random sampling, with minimal computational overhead.*
**Keywords:** *computer vision, UAV navigation, deep learning, metric learning*

## 1. Introduction & Related Work

Visual localization of UAVs without GNSS is a key challenge in autonomous navigation, especially in environments where satellite signals are unavailable, distorted, or intentionally jammed. This paper proposes a UAV position estimation method based on particle filter and local binary pattern (LBP) matching. To improve the efficiency of the algorithm, initialization by a region pre-selection mechanism based on a metric learning neural network [1] is applied, which reduces the search space and accelerates convergence. By combining the best of both worlds, the method offers a notable improvement over prior work based on Monte Carlo Localization [2] or direct matching of UAV images with orthophotos as proposed in the study introducing the MTGL40-5 dataset [3].

350

# 2. Place Recognition Using Metric Learning

**Dataset description:**   National Agriculture Imagery Program (NAIP) [4] provides digital orthophotography covering the area of 3.75 x 3.75 angular minutes with a buffer of 300 meters per side. The downloaded data is from the period from January 2005 to November 2022, and includes only photographs of locations with at least five images taken at different time points, to make sure the method can deal with the variance due to seasonal changes. The nominal resolution of pictures is 1 m/px for shots taken before 2018 and 0.6 m/px for newer ones.



Figure 1. Evaluation of a ResNet50 model with 256 embedding vector length. The top five images in ranking the query image matching results on test set.

The dataset was divided randomly by location into train, validation and test datasets in the proportion of 70%, 10%, 20%; the number of various locations in each set is shown in Table 1.

Table 1. Detailed information of the division of the dataset

| Dataset | **Training** | **Validation** | **Test** | **Total** |
|---|---|---|---|---|
| Unique locations | 71661 | 10236 | 20475 | 102372 |

Aside from training set augmentations, two levels of augmentations were applied to the test set, thereby creating another two separate test datasets called medium and hard. The augmentations used to create the hard test set were also used to modify the validation collection. Details are presented in Table 2.

**Model training:**   For training, validation as well as testing of the model, centre crop images of size 512×512 pixels were used. ResNets were chosen as a starting

Table 2. Detailed information on the division of the collection

| Collection | Rotation [°] | Scaling [%] | Translation [%] |
|---|---|---|---|
| Training | 0 to 360 | 90 to 110 | -20 to 20 |
| Validation | -10 to 10 | 90 to 110 | -10 to 10 |
| Easy Test | - | - | - |
| Medium Test | -5 to 5 | 95 to 105 | -5 to 5 |
| Hard Test | -10 to 10 | 90 to 110 | -10 to 10 |

point to evaluate the viability of the place image retrieval approach. The architecture is widely studied and has proven successful in various tasks beyond image classification. The model was trained using metric learning with a triplet margin loss [5]. Training employed batch-hard mining, with the Euclidean distance as the embedding metric. Two embedding vector lengths were tested: 128 and 256.

**Metrics and evaluation:** To determine the effectiveness of the model, two metrics were used, with one in four versions: Mean Average Precision (MAP), which takes into account both the relevance and position of the returned results and rank-K, the probability that a correct place match was found among K nearest neighbours. The addition of more significant perturbations in the more challenging variants of the test set leads to worse performance, but the solution is still capable of finding the relevant places most of the time – see Figure 1. Please note that a significant portion of the NAIP dataset contains images that are challenging due to their uniformity and the lack of distinctive features, making the problem even more difficult.

The results presented in Table 3 show that, in general, the use of larger embeddings improves the performance of the model, as directly indicated by the higher values of the achieved metrics. Also, the use of larger encoders has a positive effect on the learned feature representation.

# 3. Visual Navigation

The system assumes the presence of a downward-facing camera and a compass, assuring the correct orientation of the image. The particle filter-based visual localization algorithm starts with an initialization in which particles representing potential drone positions are randomly distributed across the map. Then a random

Table 3. Results of model evaluation on test datasets

| Dataset | Encoder | Embed. | MAP | rank-1 | rank-5 | rank-10 | rank-25 |
|---|---|---|---|---|---|---|---|
| Easy | ResNet18 | 128 | 0.727 | 0.834 | 0.881 | 0.896 | 0.915 |
| | | 256 | 0.764 | 0.858 | 0.900 | 0.913 | 0.930 |
| | ResNet34 | 128 | 0.764 | 0.863 | 0.906 | 0.922 | 0.941 |
| | | 256 | **0.783** | **0.872** | **0.912** | **0.926** | **0.942** |
| | ResNet50 | 128 | 0.752 | 0.855 | 0.902 | 0.916 | 0.935 |
| | | 256 | 0.768 | 0.865 | 0.905 | 0.920 | 0.938 |
| Medium | ResNet18 | 128 | 0.496 | 0.711 | 0.814 | 0.846 | 0.883 |
| | | 256 | 0.518 | 0.732 | 0.830 | 0.859 | 0.894 |
| | ResNet34 | 128 | 0.522 | 0.735 | 0.842 | 0.874 | 0.911 |
| | | 256 | 0.549 | 0.759 | 0.855 | 0.888 | 0.916 |
| | ResNet50 | 128 | 0.545 | 0.745 | 0.850 | 0.882 | 0.917 |
| | | 256 | **0.567** | **0.765** | **0.860** | **0.888** | **0.921** |
| Hard | ResNet18 | 128 | 0.418 | 0.633 | 0.713 | 0.739 | 0.776 |
| | | 256 | 0.435 | 0.649 | 0.723 | 0.750 | 0.787 |
| | ResNet34 | 128 | 0.443 | 0.657 | 0.738 | 0.766 | 0.805 |
| | | 256 | 0.457 | 0.669 | 0.745 | 0.773 | 0.809 |
| | ResNet50 | 128 | 0.454 | 0.662 | 0.745 | 0.775 | 0.814 |
| | | 256 | **0.464** | **0.674** | **0.751** | **0.780** | **0.816** |

trajectory is generated, running from left to right of the image, according to which the drone's movement is simulated. In each iteration of the algorithm, the UAV's position is updated as it moves over the map. For each particle, a descriptor is computed by combining a 16-bin colour histogram for the red, green and blue channels and a 32-bin uniform LBP histogram to capture the local colour and the texture information. The same descriptor is computed at the current UAV location in the map. The particle descriptors are matched against the UAV descriptor, and based on the results, the particles are assigned weights – the more similar the descriptors, the larger the weight. The information is used within the particle filter framework. Particles with higher weights dominate the system, while less likely positions are given smaller weights and eliminated. A resampling process is then performed, in which a set of samples is generated based on particles with the highest weights, allowing the search to be concentrated in more likely areas. The process is repeated iteratively. The comparison of the two approaches is illustrated in Figure 2.

Figure 2. The configuration uses 150 particles with random motion perturbation over a range of 20 pixels. The resampling process using the systematic resampling method, ensuring a more even distribution of particles and minimizing their degeneration. The size of the image patches used for matching is $128 \times 128$ pixels.

To improve the efficiency of the system, the initial stage of the algorithm is modified by limiting the area of initial particle generation instead of randomly distributing them over the entire map area. In the classical approach, particles are initialized uniformly, which leads to an inefficient search of the complete solution space, slowing down the convergence process. The proposed solution uses the neural network from Section 2 to preselect regions for initialization.

As the simulation begins, the input map is processed by the neural network, tile by tile, and potential regions with a high probability of containing the drone's start position are determined. The initial particles are generated only in these areas, limiting the search space. The process of particle convergence is faster, as the areas with a complete dissimilarity to the area under which the drone is located, are eliminated.

## 4. Results and Conclusions

As shown in Figure 3b, using the preselection method enables faster reduction of localization error, leading to faster convergence of particles toward the actual position of the UAV. The pre-selection method achieves a lower overall error,

(a) Randomly chosen pictures from test dataset



(b) Error values for tested trajectories

Figure 3. Results of comparison between two methods of visual localization

especially in the initial localization stages, where the reduction in initial uncertainty significantly improves tracking performance. In contrast, the random sample placement approach has a more significant initial error and slower convergence, meaning that more steps are needed to achieve precise position matching.

# References

[1] Musgrave, K., Belongie, S., and Lim, S.-N. Pytorch metric learning. *arXiv preprint arXiv:2008.09164*, 2020.

[2] Kinnari, J., Verdoja, F., and Kyrki, V. GNSS-denied geolocalization of UAVs by visual matching of onboard camera images with orthophotos. *arXiv preprint arXiv:2103.14381*, 2021.

[3] Ma, J., Pei, S., Yang, Y., Tang, X., and Zhang, X. MTGL40-5: A multi-temporal geo-localization dataset for satellite view. *Remote Sensing*, 15(17):4229, 2023. doi:10.3390/rs15174229.

[4] U.S. Department of Agriculture (USDA). National Agriculture Imagery Program (NAIP), 2022. URL https://developers.google.com/earth-engine/datasets/catalog/USDA_NAIP_DOQQ.

[5] Hermans, A., Beyer, L., and Leibe, B. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.

# Author Index

# PP-RAI'2025

*6th Polish Conference
on Artifical Intelligence (PP-RAI'2025)*

**07–09.04.2025, Katowice, Poland**