



International Conference of Innovative
Computer Engineering
(ICE 2025)

PROCEEDINGS E-BOOK

6 November 2025 Ankara



ICE 2025
INTERNATIONAL CONFERENCE OF
INNOVATIVE COMPUTER ENGINEERING

International Conference of Innovative Computer Engineering (ICE 2025)

6 November 2025 Ankara

PROCEEDINGS E-BOOK

Organizer



This congress proceedings book is published as an electronic format as e-book.

All rights reserved.

Editorial Board

Prof. Dr. Necaattin BARIŞCI
Prof. Dr. İbrahim Alper DOĞRU
Assoc. Prof. Dr. İsmail ATACAK
Assoc. Prof. Dr. Sinan TOKLU
Asst. Prof. Dr. Fuat TÜRK

ISBN:

Address Gazi University, Faculty of Technology, Taskent Building,
Department of Computer Engineering, Emniyet District, Bandırma
Street, No: 6/34 06560, Yenimahalle, Ankara, Türkiye
Phone +90 506 333 60 55
E-mail iceconf@gazi.edu.tr
Web <https://iceconf.gazi.edu.tr>

All legal and ethical responsibility of the studies belongs to the authors. All rights reserved. The copyright of this proceedings book belongs to Gazi University. It may not be copied or reproduced without permission.

TABLE OF CONTENTS

Briefly About.....	iv
ICE 2025 Organization.....	v
Scientific Committee.....	vii
ICE 2025 Invited Speakers.	ix
Congress Programme	x
Abstracts And Full Papers.	1
Efficient Compression And Decompression Framework For Automotive Embedded Data Logger In Cloud Telemetry	2
An Imitation-Enhanced Actor-Critic Approach For Maze Navigation	3
Deep Learning Approaches On Image Representations Of Android Malware: A Review.....	4
Artificial Intelligence-Based Electronics: Towards Autonomous Systems.	5
Multilingual Hate Speech Detection With XLM-Roberta: From Baseline Benchmarks To ONNX-Optimized Android Deployment	6
Attention-Enhanced CNN With Grad-CAM For Explainable Brain Tumor Classification.....	7
Interpretable Gradient-Boosted Poisson Modeling Of Aftershock Productivity: Magnitude- Sensitive 7-Day/100-KM Forecasts That Outperform RJ89.....	8
A Fully Monocular Deep Learning Pipeline For 3D Urban Reconstruction From Satellite Imagery.....	9
Full Papers.....	10
Cost Effective And Generalized Turkish Passage Retrieval With ColBERT.....	11
Watermarking And Steganography In AI-Generated Text: Character, Syntactic, Semantic, And Hidden Message Embedding Approaches	16
Performance Analysis Of Rule-Based, CRF, BiLSTM-CRF, And BERT Models For Named Entity Recognition	22
Bus Arrival Time Prediction Using Machine Learning Techniques... ..	31
Building An Arabic TripAdvisor Dataset For Sentiment Analysis With Cohen's Kappa Validation	36
Spatio-Temporal Evaluation Of Hybrid Trend-LSTM Models For Bus Arrival Prediction	41
A Comparative Analysis Of Turkey's KVKK And The EU's GDPR In An Era Of Technological Transformation	47
THERMAL-ADAS-TR Dataset Collected In Türkiye And Object Detection Performance Evaluation With FLIR-ADAS.....	51
Performance Comparison Of BERT And TF-IDF With Machine Learning Methods On Sentiment Analysis	59
From Lyrics To Insights: Multidimensional Emotion, Theme, And Demographic Analysis In Turkish Music.....	71
A Comparative Forensic Analysis Of EU And Turkish Payment Regulations: Bridging The Gaps Between PSD3/PSR And Law No. 6493 In Combating Cybercrime	79

BRIEFLY ABOUT

The International Conference of Innovative Computer Engineering was held on November 6. The congress aimed to bring together researchers, industry leaders, and policymakers working on interdisciplinary subjects in informatics and software engineering. Its primary objective was to establish an effective communication platform to explore how technological advancements can address global challenges such as sustainability, healthcare, and education. By bridging the gap between academia and industry, the conference highlighted the transition from theoretical models to scalable, real-world solutions.

Event Topics of the conference: Computer Science, Generative AI & Large Language Models, Software Engineering, Game Technologies & Virtual Worlds, Internet of Things (IoT) & Smart Systems, Data Science, Big Data & Advanced Analytics, Sustainable & Green Computing, and Health Informatics.

Sub-topics of the conference: Human-Computer Interaction, Programming Languages, Database Management Systems and Data Structures, Software Architecture and Design Patterns, Software Standards, Modeling and Simulation, Internet of Things, Parallel and Distributed Computing, Medical Imaging Systems and Expert Systems, Emerging Technologies in Software Engineering, Algorithms and Discrete Mathematics, Explainable AI & Responsible Machine Learning, Blockchain for Identity Management & Secure Transactions, Deepfake Detection & Disinformation Control, Agile, DevOps & MLOps: Automating the Development Lifecycle, Telemedicine & Smart Health Systems, Blockchain for Secure Health Data Management, Data Engineering for Scalable AI Applications, Digital Twins & AI-Enhanced Supply Chain Management, Error Correction & Noise Reduction in Quantum Computing, Low-Code & No-Code Software Development, Microservices & Event-Driven Architecture, Computer Networks and Logic, Decision Support Systems and Resource Management.

The language of the conference is English, and all full papers submitted for publication in the congress on current issues have been evaluated by at least two referees by the blind reviewing method. 23 papers were accepted for oral presentation and publication as a result of peer review. We would like to thank all the researchers who have shown interest in the Conference.

ICE 2025 ORGANIZATION

ICE 2025, organized by Gazi University Faculty of Technology, Ankara, Türkiye

Honorary President

Prof. Dr. Uğur Ünal – Rector of Gazi University

ICE 2025 Chairman of Conference

Prof. Dr. O. Ayhan Erdem (Gazi University)

Organizing Committee

Prof. Dr. Necaattin Barışçı– Gazi University, Türkiye

Prof. Dr. İbrahim Alper Doğru– Gazi University, Türkiye

Prof. Dr. Aysun Coşkun – Gazi University, Türkiye

Prof. Dr. Mehmet Şimşek - Sinop University, Türkiye

Prof. Dr. Nurettin Doğan- Selçuk University, Türkiye

Prof. Dr. Yusuf Sönmez – Gazi University, Türkiye

Prof. Dr. Bünyamin Cıylan – Gazi University, Türkiye

Prof. Dr. Hüseyin Polat – Gazi University, Türkiye

Prof. Dr. Aydın Çetin – Gazi University, Türkiye

Prof. Dr. Cavanşir Zeynalov, Dean – Faculty of Architecture and Engineering, Nakhchivan State
University (NDU), Azerbaijan

Prof. Dr. Graham Kendall – University of Nottingham Malaysia, Malaysia

Prof. Dr. Wan Mohd Nasir bin Wan Kadir, Dean – Universiti Teknologi Malaysia (UTM), Malaysia

Assoc. Prof. Dr. İsmail Atacak– Gazi University, Türkiye

Assoc. Prof. Dr. Sinan Toklu– Gazi University, Türkiye

Assoc. Prof. Dr. Abdullah Talha Kabakuş – Düzce University, Türkiye

Assoc. Prof. Dr. Fecir Duran – Gazi University, Türkiye

Assoc. Prof. Dr. Murat Dörterler – Gazi University, Türkiye

Assoc. Prof. Dr. Saadin Oyucu – Gazi University, Türkiye

Assoc. Prof. Dr. Adem Tekerek – Gazi University, Türkiye

Assoc. Prof. Dr. Cemal Koçak – Gazi University, Türkiye

Assoc. Prof. Dr. Məftun Əliyev – Nakhchivan State University (NDU), Azerbaijan

Asst. Prof. Dr. Fuat Türk– Gazi University, Türkiye

Asst. Prof. Dr. Mohammed Rashad Baker- College of Computer Science and Information Technology,
University of Kirkuk, Kirkuk, Iraq

SCIENTIFIC COMMITTEE

Prof. Dr. M. Ali Akcayol – Gazi University, Türkiye.

Prof. Dr. Hacer Karacan – Gazi University, Türkiye.

Prof. Dr. Ali Hakan Işık – Burdur Mehmet Akif Ersoy University, Türkiye.

Prof. Dr. Rajesh Kaluri – Vellore Institute of Technology, India.

Prof. Dr. Nursal Arıcı - Gazi University, Türkiye.

Prof. Dr. İbrahim Yücedağ – Düzce University, Türkiye.

Prof. Dr. Pakize Erdoğanmuş – Düzce University, Türkiye.

Prof. Dr. Resul Kara – Düzce University, Türkiye.

Prof. Dr. Dharmendra Singh Rajput – VIT Vellore, India.

Prof. Dr. Magdalena Palacz – Silesian University of Technology, Poland.

Prof. Dr. Cihan Varol – Sam Houston State University, United States.

Prof. Dr. İlyas Çankaya – Ankara Yıldırım Beyazıt University, Turkey.

Prof. Dr. Tuncay Aydoğan – Isparta University of Applied Sciences, Türkiye.

Prof. Dr. Ecir Uğur Küçüksille – Isparta University of Applied Sciences, Türkiye.

Prof. Dr. Yılmaz Çamurcu – Maltepe University, Türkiye.

Assoc. Prof. Dr. Arafat Şentürk – Düzce University, Türkiye.

Assoc. Prof. Dr. Almaz Aliyeva, Acting Dean – Mingachevir State University, Azerbaijan.

Assoc. Prof. Dr. Aida Mustafayeva – Mingachevir State University, Azerbaijan.

Assoc. Prof. Dr. Sushank Chaudhary - University of Petrochemical Technology, China.

Assoc. Prof. Dr. Serdar Biroğul – Düzce University, Türkiye.

Assoc. Prof. Dr. Fatih Kayaalp – Düzce University, Türkiye.

Assoc. Prof. Dr. Oktay Yıldız – Gazi University, Türkiye.

Assoc. Prof. Dr. Ümit Atila – Gazi University, Türkiye.

Assoc. Prof. Dr. Anıl Utku – Munzur University, Türkiye.

Assoc. Prof. Dr. Recep Sinan Arslan – Kayseri University, Türkiye.

Assoc. Prof. Dr. Hakan Öcal - Bartın University, Türkiye.

Assoc. Prof. Dr. Gür Emre Güraksın – Afyon Kocatepe University, Türkiye.

Asst. Prof. Dr. Enver Küçükkülahlı – Düzce University, Türkiye.

Asst. Prof. Dr. Mehmet Sevrî – Recep Tayyip Erdoğan University, Türkiye.

Asst. Prof. Dr. Hikmet Canlı – İstanbul Gedik University, Türkiye.

Asst. Prof. Dr. Muhammed Ali Koşan - İstiklal University, Türkiye.

Asst. Prof. Dr. Susana Leal – Polytechnic Institute of Santarém, Portugal.

Asst. Prof. Dr. Abhishek Sharma – National Institute of Technology Hamirpur, India.

Asst. Prof. Dr. Alaan Ghazi – Northern Technical University, Iraq.

Asst. Prof. Dr. Baban A. Mahmood – Kirkuk University, Iraq.

Asst. Prof. Dr. Ahmed Chalak Shakir – Kirkuk University, Iraq.

Asst. Prof. Dr. Mohammed Rashad Baker – Kirkuk University, Iraq.

Asst. Prof. Dr. Kazım Kılıç – Yozgat Bozok University, Türkiye.

Asst. Prof. Zafer Ayaz - Gazi University, Türkiye.

Dr. Hamdullah Karamollaoğlu – EÜAŞ, Türkiye.

Dr. Süleyman Muhammed Arıkan – ASELSAN,
Türkiye.

Dr. Linda Daniela, Dean – University of Latvia,
Latvia.

Dr. Mohammad Abobala – Latakia University, Syria.

Dr. Tunahan Timuçin – Düzce University, Türkiye.

Lect. Mehmet KIZILDAĞ – Alparslan Türkeş
University, Türkiye.

Lect. Burak Eskici – Gazi University, Türkiye.

Instr. Muhammet Ünal – Gazi University, Türkiye.

Snr. Lec. Dr. Pratik Vyas – Nottingham Trent
University, England.

Dr. Esra Söğüt – Gazi University, Türkiye.

Res. Asst. Ömer Ayberk Şencan – Gazi University,
Türkiye.

Res. Asst. Rıdvan Sert – Gazi University, Türkiye.

Res. Asst. Oğuzhan Sezer – Gazi University, Türkiye.

Res. Asst. Adem Varol – Gazi University, Türkiye.

Res. Asst. Büşra Duygu Çelik – Gazi University,
Türkiye.

Res. Asst. Ezgi Kara Timuçin – Düzce University,
Türkiye.

Res. Asst. Bayram Küçük - Düzce University,
Türkiye.

Res. Asst. Ercan Atagün - Düzce University, Türkiye.

Merve Güllü - Türk Telekom, Türkiye.

Oğuzhan Çıtlak – Tüprag, Türkiye.

ICE 2025 INVITED SPEAKERS



Prof. Dr. Asaf VAROL -The University of Tennessee at Chattanooga, United States.

CONFERENCE PROGRAM

Opening Speech & Guest Speak Teams Link: https://mth.tc/8gEvn				
9.30-9.50 A.M. November 6, 2025	Opening Speech Prof. Dr. Ayhan Erdem, Chair of Computer Engineering Department, Faculty Of Technology, Gazi University Prof. Dr. Musa ATAR, Dean, Faculty Of Technology, Gazi University			
9.50-10.50 A.M. November 6, 2025	Guest Speaker Prof. Dr. Asaf Varol, The University of Tennessee at Chattanooga, United States. Artificial Intelligence & Machine Learning			
	November 6, 2025	November 6, 2025	November 6, 2025	November 6, 2025
	A	B	C	D
	Session 1 Teams Link: mth.tc/35R	Session 2 Teams Link: mth.tc/2VOG	Session 3 Teams Link: mth.tc/gkij	Session 4 Teams Link: mth.tc/hJ05
November 6, 2025	Session 1: Artificial Intelligence, Deep Learning, and Image Processing Session Chair: Prof. Dr. Necaattin BARIŞCI	Session 2: Natural Language Processing, Sentiment Analysis, and Language Models Session Chair: Assoc. Prof. Dr. Sinan TOKLU	Session 3: Security, Ethics, Legal Compliance, and Privacy Session Chair: Prof. Dr. İbrahim Alper DOĞRU	Session 4: Data Analytics, Forecasting, and Telecommunication Session Chair: Dr. Esra Söğüt
13.45-14.00 P.M.	An Imitation-Enhanced Actor-Critic Approach for Maze Navigation – Nurgül Kalaycı, Mehmet Dinçer Erbaş	Performance Analysis of Rule-Based, CRF, BiLSTM-CRF, and BERT Models for Named Entity Recognition – Ahmet Toprak, Feyzanur Sağlam Toprak	A Comparative Analysis of Turkey's KVKK and the EU's GDPR in an Era of Technological Transformation – Melih Aybar, Aysun Coşkun	Bus Arrival Time Prediction Using Machine Learning Techniques – Ahmet Tiryaki, Tefik Aytekin
14.00-14.15 P.M.	A Fully Monocular Deep Learning Pipeline for 3D Urban Reconstruction from Satellite Imagery – Ömer Karadağ, Nihal Altuntaş	Performance Comparison of BERT and TF-IDF with Machine Learning Methods on Sentiment Analysis – Muhammet Başarslan, Fatih Kayaalp	A Comparative Forensic Analysis of EU and Turkish Payment Regulations: Bridging the Gaps Between PSD3/PSR and Law No. 6493 in Combating Cybercrime – Melih Aybar, Aysun COŞKUN	Spatio-Temporal Evaluation of Hybrid Trend-LSTM Models for Bus Arrival Prediction – Osman Kaya, Mustafa Utku Kalay
14.15-14.30 P.M.	Deep Learning Approaches on Image Representations of Android Malware: A Review – Oğuzhan Sezer, İbrahim Alper Doğru, Kazım Kılıç	Multilingual Hate Speech Detection with XLM-RoBERTa: From Baseline Benchmarks to ONNX-Optimized Android Deployment – Salim Döğmash, İbrahim Alper Doğru, Kazım Kılıç	Responsible Development of Generative AI: Trends, Challenges, and Applications – Fatime Rehman, Habil İsmayilzade, Şahla Şirinova	Interpretable Gradient-Boosted Poisson Modeling of Aftershock Productivity: Magnitude-Sensitive 7-Day/100-km Forecasts that Outperform RJ89 – Mehmet Sevi, Furkan Yurdakul Kayıkcı, Ali Gürbüz
14.30-14.45 P.M.	Attention-Enhanced CNN with Grad-CAM for Explainable Brain Tumor Classification – Esra Söğüt, Ayhan Erdem, Maral A. Mustafa	Building an Arabic TripAdvisor Dataset for Sentiment Analysis with Cohen's Kappa Validation – Abbas Ali, Necaattin Barışçı	Watermarking and Steganography in AI-Generated Text: Character, Syntactic, Semantic, and Hidden Message Embedding Approaches – Fatma Gümüş, Muhammed Ersin Durmuş, Murat Utku Kabasakaloğlu, Enver Kağan Çetiner, Fatma Gümüş	Research of the quality of functioning multiservice telecommunication networks – taking into account the self-similarity of transmitted traffic – Almaz Aliyeva
14.45-15.00 P.M.	Thermal-ADAS-TR Dataset Collected in Türkiye and Object Detection Performance Evaluation with FLIR-ADAS – Umut Genç, Behçet Uğur Töreyn	From Lyrics to Insights: Multidimensional Emotion, Theme, and Demographic Analysis in Turkish Music – Ege Kutlu, Naz Stancroğlu Albayrak		Cost Effective and Generalized Turkish Passage Retrieval with CoBERT – Eren Tahir, Mert Bal
15.00-15.15 P.M.	Artificial Intelligence-Based Electronics: Towards Autonomous Systems – Mahsati Əliyeva	A Hybrid Framework for Detecting Hallucinations in LLM Responses Using Lexical and Semantic Evidence – Aryan Shivatare, Dattatray Takale, Rohit Shitole, Jayant Shelke, Sujal Sune, Omraje Shendage		Efficient Compression and Decompression Framework for Automotive Embedded Data Logger in Cloud Telemetry – Shannen Milton, Kiran VR
15.15-15.30 P.M.	A Unified Governance Framework for Cross-Platform Self-Service Business Intelligence Systems: Power BI and Qlik Sense – Mohith Reddy Patilola			
Closing Speech Teams Link: https://mth.tc/BJk6				
15.30-15.45 P.M.	Closing Speech Prof. Dr. Ayhan ERDEM			



ABSTRACTS AND FULL PAPERS

EFFICIENT COMPRESSION AND DECOMPRESSION FRAMEWORK FOR AUTOMOTIVE EMBEDDED DATA LOGGER IN CLOUD TELEMETRY

Shannen Milton

Spark Minda Technical Centre
Minda Corporation Limited
shannenmilton2003@gmail.com

Kiran V.R

Spark Minda Technical Centre
Minda Corporation Limited
kiran.vr@mindacorporation.com

Abstract

Modern embedded systems, particularly in automotive and IoT domains, generate increasingly large volumes of real-time telemetry data. However, limited memory, bandwidth, and processing resources create a critical bottleneck for efficient logging, storage, and transmission. Traditional compression algorithms like ZIP and GZIP are often too computationally intensive for microcontroller-based platforms, motivating the adoption of lightweight, deterministic alternatives.

We present a modular datalogging framework that leverages the Heatsrink library for lossless compression and decompression, combined with a lightweight framework for reliable storage as well as to reduce communication bandwidth. Sensor readings are compressed and stored in an on-board flash memory with structured metadata for data management and retrieval. This compressed data provides an additional layer of security. This will enable us to store huge datalogs in the flash memory. The compressed data is sent to the cloud, and a web-application will be used to access and analyse the compressed data.

Initial evaluations demonstrate compression efficiencies exceeding 80%, robust decompression under noisy conditions, and a reliable end-to-end flow from sensors to cloud storage. By minimizing communication overhead while keeping devices adaptable and connected, this architecture provides a scalable way for storing large datasets in a cloud database. It can also be used to store large data in the embedded flash memory, supporting automotive telematics, industrial monitoring, and vehicle analytics.

Keywords: Automotive Data Logger, Cloud Telemetry, Data Compression, Embedded Systems, Heatsrink Compression, Lossless Compression, Firmware Over-The-Air (FOTA)

AN IMITATION-ENHANCED ACTOR-CRITIC APPROACH FOR MAZE NAVIGATION

Nurgül Kalaycı
İğdır University
nurgul.kalayci@igdir.edu.tr

Mehmet Dinçer Erbaş
Bolu Abant İzzet Baysal University
dincer.erbas@ibu.edu.tr

Abstract

Aim: This paper proposes a new method that integrates Imitation Learning (IL) into the Actor-Critic (AC) algorithm and evaluates its effect on the learning speed of agents in both obstacle and obstacle-free mazes.

Methods: The paper presents a new method that combines AC with a model of simply imitating observed behaviour, without direct internal access or inter-agent experience sharing. The proposed Imitation Enhanced Actor-Critic learning approach is evaluated through simulations on a maze pathfinding problem.

Results: Simulation results indicated that the proposed learning method significantly accelerated learning and achieved higher stability compared to the standard Actor-Critic approach. In particular, agents using the imitation-enhanced model reached optimal paths in fewer steps, in both mazes.

Conclusion: The findings demonstrate that integrating imitation learning into the actor-critic framework enhances convergence and stability, particularly in complex environments with obstacles. The proposed IEAC approach has potential applicability in autonomous navigation and robotic systems requiring efficient learning from observation.

Keywords: Social Learning, Reinforcement Learning, Multi-Agent Systems.

DEEP LEARNING APPROACHES ON IMAGE REPRESENTATIONS OF ANDROID MALWARE: A REVIEW

Oğuzhan Sezer

Gazi University
oguzhansezer@gazi.edu.tr

İbrahim Alper Doğru

Gazi University
iadogru@gazi.edu.tr

Kazım Kılıç

Gazi University
kazim.kilic@gazi.edu.tr

Abstract

Aim: This survey intends to summarize the latest progress in deep learning-based and hybrid multimodal Android malware detection. Spanning work from 2023 through 2025, the survey illustrates the transition from the traditional static analysis direction towards detection models propelled by graph, vision, and transformer-based paradigms. It seeks the exploration of existing trends, innovative methodologic directions, and open challenges in boosting the robustness, interpretability, and computational speed of Android malware detectors.

Methods: A focused literature review was conducted using the Elsevier, IEEE Xplore, and ScienceDirect databases to examine recent advances in image-based Android malware detection published between 2023 and 2025. Ten representative studies were selected for their methodological diversity and contributions to visual, hybrid, and interpretable detection frameworks. The reviewed approaches include end-to-end models converting DEX bytecode into grayscale or RGB matrices, graph-based and multimodal fusion methods combining structural and semantic features, as well as advanced architectures such as 3D-CNNs and Vision Transformers capable of capturing multiscale contextual patterns. The review also covers lightweight learning frameworks aimed at reducing model complexity and a few explainable AI (XAI) systems employing SHAP or Grad-CAM for interpretability. Comparative evaluation focused on dataset properties, preprocessing strategies, and performance metrics including accuracy, F1-score, robustness, and computational efficiency.

Results: The review reveals a significant advancement from hand-crafted static features to automatic, image-centric malware classification. Image-based CNN models such as MADRF-CNN have achieved 96–98% accuracy using pairwise pooling and Dex segment cropping. Graph attention networks and multimodal fusion models have achieved accuracy up to 99.5% by combining structural representations. Hybrid 3D convolutional architectures have demonstrated reduced false positive learning on obfuscated malware samples. However, challenges remain regarding dataset imbalance, computational cost, and adversarial robustness.

Conclusion: Recent literature demonstrates that image-based and graph-based representations are transforming Android malware detection by combining static and behavioral analysis within deep learning frameworks. Current research trends emphasize the development of lightweight CNN-Transformer hybrid models for real-time deployment, the creation of balanced benchmark datasets containing obfuscated and adversarial examples, and the integration of explainability with multimodal fusion to support more reliable decision-making. Furthermore, self-supervised and continuous learning approaches are being explored to ensure adaptability against evolving malware ecosystems. Overall, the convergence of visual computing, graph reasoning, and explainable AI is shaping the next generation of Android malware defense systems that are interpretable, scalable, and resilient.

Keywords: Android malware, image-based detection, deep learning, mobile security

ARTIFICIAL INTELLIGENCE-BASED ELECTRONICS: TOWARDS AUTONOMOUS SYSTEMS

Aliyeva Mahsati Rovshan

Mingachevir State University
mahsati.aliyeva@mdu.edu.az

Shukurova Leyla Niyazi

Mingachevir State University
leyla.shukurova@mdu.edu.az

Abstract

Aim: The main purpose of this study is to comprehensively analyze the application of artificial intelligence (AI) in the field of electronics and the capabilities of automated systems resulting from the synthesis of these two technologies. Electronics has been built on systems with fixed, pre-programmed and limited adaptability for many years. However, with the development of artificial intelligence, a new stage has begun in this field, and the ability of devices to make independent decisions, learn and optimize has increased. The aim of the study is to identify the technological differences between traditional electronics systems and AI-based systems, as well as to show the advantages that these differences bring in terms of efficiency and innovation in the areas of application.

Methods: The study is based on analytical and comparative research approaches. Relevant literature published in leading scientific databases such as IEEE Xplore, ScienceDirect, and Nature between 2019 and 2024 was reviewed to identify current approaches in AI-integrated electronic systems. Comparative analysis was conducted between traditional automation models and AI-driven systems in terms of efficiency, adaptability, diagnostic capabilities, and energy consumption. Additionally, real-world case studies from the fields of industrial manufacturing, healthcare, and transportation were evaluated to validate theoretical findings.

Results: Findings indicate that AI-enhanced electronic systems outperform traditional automation models in terms of operational flexibility, predictive maintenance, and energy efficiency. These systems demonstrate 30–35% higher operational stability and up to 25% lower energy consumption. Furthermore, AI integration facilitates real-time diagnostics and self-optimization, leading to increased productivity and reduced human dependency in industrial and service environments.

Conclusion: The integration of artificial intelligence into electronics marks a transformative phase in technological development. Beyond improving functionality, this synthesis enhances decision-making transparency, energy optimization, and overall system intelligence. However, it also raises ethical and regulatory challenges related to accountability and data privacy. Future studies should focus on neuromorphic chip design, adaptive energy management, and standardization frameworks to ensure responsible and sustainable implementation of AI-driven automation.

The results of the study show that electronic systems equipped with artificial intelligence are one of the strategic directions of future technological development. This synthesis not only increases the functional capabilities of electronics, but also its intellectual level, meeting the automation needs of modern society. As a result, the continuation of scientific research and innovation activities in this direction is considered necessary for the creation of new adaptive and intelligent systems.

Keywords: Intelligent control systems, adaptive integration, autonomous devices, energy optimization, predictive maintenance, digital transformation

MULTILINGUAL HATE SPEECH DETECTION WITH XLM-ROBERTA: FROM BASELINE BENCHMARKS TO ONNX-OPTIMIZED ANDROID DEPLOYMENT

Salam Thabet Doghmash

Gazi University
St.doghmash@gazi.edu.tr

İbrahim Alper Doğru

Gazi University
iadogru@gazi.edu.tr

Kazım Kılıç

Gazi University
kazim.kilic@gazi.edu.tr

Abstract

The increasing spread of hate speech across digital platforms creates serious challenges for online safety and multilingual content moderation.

Aim: This study aims to develop a scalable detection system for Arabic, Turkish, and English, with XLM RoBERTa as the core model and ensemble baselines used only for benchmarking.

Methods: The FrancophonIA dataset (220k samples) was preprocessed. Baseline models were built using TF-IDF features with Logistic Regression, Support Vector Machines, and Random Forest combined in a soft-voting ensemble, while the optimized transformer model was fine-tuned on XLM-RoBERTa with GPU acceleration, early stopping. Models were exported in ONNX, Open Neural Network Exchange; formats for deployment, and an Android application was developed to provide mobile access.

Results: The ensemble baseline achieved a macro F1-score of 0.5827, outperforming individual classifiers, while the optimized transformer substantially improved performance, reaching an accuracy of 72.3% and a macro F1-score of 0.7079 across the three languages. The ONNX-optimized transformer enabled efficient real-time inference through a lightweight API and Android application, both achieving sub-second response times.

Conclusion: Overall, the study demonstrates that transformer-based multilingual models are robust and provide a solid foundation for building future systems that detect hate speech across multilingual content in practice and mobile environments.

Keywords: Multilingual hate speech detection, XLM-RoBERTa, ensemble baseline, ONNX, Android deployment

ATTENTION-ENHANCED CNN WITH GRAD-CAM FOR EXPLAINABLE BRAIN TUMOR CLASSIFICATION

Maral A. Mustafa

Northern Technical University
maralanwer@ntu.edu.iq

O. Ayhan Erdem

Gazi University
ayerdem@gazi.edu.tr

Esra Söğüt

Gazi University
esrasogut@gazi.edu.tr

Abstract

Aim: To analyze the clinical interpretability of deep learning models for classifying brain tumors on MRI.

Methods: We compared a baseline Convolutional Neural Network (CNN), an attention-enhanced CNN, and transfer-learning back-bones (MobileNetV2, InceptionV3, Xception) on a four-class MRI dataset (glioma, meningioma, pituitary, no tumor). Models were trained with Adam and evaluated using Grad-CAM visualizations, F1-score, accuracy, precision, recall, and confusion matrices.

Results: MobileNetV2 achieved the highest accuracy (95.7%), with closely aligned precision, recall, and F1-score. The attention-augmented CNN performed competitively, and Grad-CAM highlighted tumor-relevant regions, supporting model reliability.

Conclusion: Transfer learning, especially MobileNetV2, offers strong performance for MRI tumor classification, while attention mechanisms and Grad-CAM improve focus and interpretability, facilitating clinical adoption.

Keywords: brain tumor, MRI, deep learning, transfer learning, explainable AI

INTERPRETABLE GRADIENT-BOOSTED POISSON MODELING OF AFTERSHOCK PRODUCTIVITY: MAGNITUDE-SENSITIVE 7-DAY/100-KM FORECASTS THAT OUTPERFORM RJ89

Mehmet Sevri

Recep Tayyip Erdogan University
mehmet.sevri@erdogan.edu.tr

Furkan Yurdakul Kayıkçı

Recep Tayyip Erdogan University
furkanyurdakul_kayikci24@erdogan.edu.tr

Ali Gürbüz

Recep Tayyip Erdogan University
ali.gurbuz@erdogan.edu.tr

Abstract

Aim: The aim of this study is to predict the number of aftershocks that may occur within 7 days and a 100-kilometer radius following a mainshock. This prediction is made based on the parameters of the mainshock and early catalog features, including magnitude, focal depth, magnitude type, number of stations, azimuthal gap of the network, minimum distance to the epicenter, root mean square of arrival time residuals, latitude-longitude, and origin time of the earthquake using U.S. Geological Survey Earthquake dataset. The Reasenberg-Jones model was used for comparison, and productivity estimation was provided using XGBoost.

Methods: Using the Gardner-Knopoff declustering method, 3,000 records were separated into mainshocks and aftershocks. Subsequently, aftershocks were predicted using both the RJ89 and XGBoost Poisson models and compared to each other.

Results: From 122 mainshocks identified through Gardner-Knopoff declustering, the target count distribution within 7 days and 100 kilometers was found to be approximately 47% zero-weighted and right-skewed. In general testing, the XGBoost Poisson and RJ89 models yielded similar results. However, in stratified evaluation, the MAE of XGBoost was lower than RJ89 by 5.0% for earthquakes with magnitudes between 5.0–5.5, by 11.8% for magnitudes between 5.5–6.0, and by 14.2% for magnitudes equal to or greater than 6. The graphs indicated that while both models performed poorly for large counts, they exhibited reasonable calibration in small and medium ranges. SHAP interpretations revealed clear interaction effects in variables such as $\text{mag} \times \text{depth}$ and $\text{depth} \times \text{dmin}$.

Conclusions: For small counts, both models showed similar accuracy; however, XGBoost achieved a significantly higher overall accuracy.

Keywords: Earthquake prediction, artificial intelligence, XGBoost, aftershock, mainshock, Reasenberg-Jones productivity formula, Gardner-Knopoff declustering

A FULLY MONOCULAR DEEP LEARNING PIPELINE FOR 3D URBAN RECONSTRUCTION FROM SATELLITE IMAGERY

Ömer Karadağ

İstanbul Gelişim Üniversitesi, İstanbul, Türkiye
krdg.omercan@hotmail.com

Nihal Altuntaş

İstanbul Gelişim Üniversitesi, İstanbul, Türkiye
nihal.altuntas@iuc.edu.tr

Abstract

Aim: This study develops and evaluates an integrated deep learning pipeline for reconstructing 3D urban scenes from single-view satellite imagery, providing a cost-effective alternative to traditional multi-view or LiDAR-based methods.

Methods: Our monocular approach integrates three specialized components: (1) a Pix2Pix conditional GAN with perceptual loss for digital surface model (DSM) generation, (2) a Mask R-CNN model for instance-level building segmentation, and (3) a ResNet-50 classifier for roof type recognition. Models were trained on curated datasets including SpaceNet Urban3D, Inria Aerial Labeling, and GATE roof data, with evaluation metrics encompassing MSE for elevation accuracy, IoU/mAP for segmentation quality, and Accuracy/F1 scores for classification performance.

Results: The pipeline demonstrates robust performance across all components. DSM generation achieved an MSE of 986.70, producing structurally coherent elevation maps. Building segmentation attained 0.71 mAP at IoU = 0.50, though boundary precision challenges emerged at stricter thresholds (0.358 mAP at IoU = 0.75). Roof classification reached 80% overall accuracy, with strong performance on common classes (hip roofs: F1=0.91) despite limitations on rare types. The integrated system successfully transforms single RGB inputs into complete 3D urban scenes.

Conclusion: This work validates that monocular deep learning pipelines can generate functionally useful 3D urban models from minimal inputs, offering a practical solution for rapid urban analysis when multi-view data is unavailable. While geometric precision remains below production-grade standards, the coordinated integration of multiple deep learning components demonstrates a viable framework for accessible urban reconstruction from conventional satellite imagery.

Keywords: Digital surface model, pix2pix, mask r-cnn, remote sensing



FULL PAPERS

COST EFFECTIVE AND GENERALIZED TURKISH PASSAGE RETRIEVAL WITH COLBERT

Eren Tahir

Yildiz Technical University, İstanbul, Türkiye
eren.tahir@std.yildiz.edu.tr

Mert Bal

Yildiz Technical University, Türkiye
mertbal@yildiz.edu.tr

Abstract:

Aim: This study aims to develop and evaluate a late interaction neural passage retrieval model, specifically a Turkish-only implementation of ColBERT, and to investigate its performance compared to multilingual alternatives in low-resource language settings.

Methods: We trained a ColBERT-based retrieval model exclusively on publicly available Turkish data under a minimal fine-tuning budget. Comparative training strategies were explored to identify efficient approaches for low-resource settings. A novel evaluation perspective was introduced through tokenizer efficiency analysis. Model performance was compared with multilingual retrieval models, focusing on retrieval effectiveness, efficiency, and cross-domain generalization.

Results: The Turkish-only ColBERT model achieved higher retrieval performance than multilingual baselines. It demonstrated strong generalization capabilities across out-of-domain tasks and maintained efficient inference speed. The results indicate that effective late interaction retrieval systems can be built for low-resource languages without extensive computational resources.

Conclusion: A Turkish-only ColBERT retrieval model can surpass multilingual counterparts, even when trained with limited data and budget. This work provides practical guidance for training strategies in low-resource settings, represents a starting study for open information retrieval in Turkish, and shows that much more could be achieved in this area. All training data is publicly available to encourage reproducibility and future collaboration.

Keywords: passage ranking, information retrieval, ColBERT, language models

1. INTRODUCTION

A Information retrieval (IR) underpins modern digital systems by enabling users to efficiently access relevant information from large data sources. Neural ranking models are central to IR, using neural networks to assess the relevance of search results (Bashkar&Craswell, 2018). These models fall into two categories: representation-based, which encode queries and documents separately, and interaction-based, which process them jointly (Li et.al., 2022). While interaction-based approaches generally deliver better ranking results, they demand higher computational resources. The ColBERT framework (Khattab&Zaharia, 2020) combines the strengths of both through a late interaction mechanism, achieving strong performance, efficiency, and robust out-of-domain generalization, making it especially valuable for low-resource languages.

In Turkish IR research, prior work has focused mainly on question answering (Akyön et.al., 2022; Budur et.al., 2024; Gemirter&Goularas, 2021), with no studies addressing open-set neural IR. This study fills that gap by contributing:

- A Turkish IR model that rivals commercial multilingual systems with minimal training costs (<https://huggingface.co/99eren99/TrColBERT>).

- A large evaluation dataset (52k passages, 600 queries) and knowledge distillation datasets for Turkish IR (<https://huggingface.co/datasets/99eren99/Tr-NanoBEIR>, <https://huggingface.co/datasets/99eren99/Turkish-IR-KD-Data>).
- An evaluation of open-weight multilingual ColBERT models for Turkish.

2. METHODS

A. Pretraining Language Model

As the initial phase of our training process, we pretrained an encoder-only language model from scratch. The ModernBERT (Warner et.al., 2024) architecture was chosen due to the superior token-level understanding capabilities of rotary position embeddings compared to absolute position embeddings (Su et.al., 2021). Additionally, alternating attention was removed, and global attention was employed at every layer to ensure unrestricted feature extraction by the model. The uncased tokenizer with a 32k vocabulary of BERTurk (Schweter, 2020) was employed, and the model was pretrained for eight epochs on the 35 GB “Oscar 21.09 deduplicated Turkish corpus” (<https://oscar-project.github.io/documentation/versions/oscar-2109/>) using the masked language modeling (MLM) objective. A 20% masking probability and dynamic masking, as implemented in RoBERTa [26], were applied. Model’s context length was limited at 512 tokens. Unlike other studies, no dropout was introduced. It is hypothesized that dynamic masking functions similarly to dropout at the input layer, and the model's capacity is insufficient to overfit the pretraining data. To support this hypothesis, the validation loss consistently decreased throughout the training process, indicating the absence of overfitting.

B. Finetuning for IR

In the second stage of training, knowledge distillation (KD) from high-performing models is utilized as the primary training objective. For each dataset, the 9 most similar passages relative to query representations are identified using a combination of approximate nearest neighbor search and IR models. Pairwise similarity scores for these passages are then computed using teacher models. To mitigate model bias, teacher ensembling, which averages scores from multiple teacher models, is often employed. However, using multiple models for the same dataset significantly increases data preparation costs. To balance efficiency and effectiveness, we instead combined different teacher models with different datasets during the fine-tuning phase. Incremental approach was employed during dataset and teacher model pairing. The study is insufficient to draw general conclusions, this method will be examined in subsequent studies. Our pretrained ModernBERT was trained on mentioned KD datasets below with KL divergence loss for 4 epochs.

1) KD Datasets [12]:

a) Turkish Translated MSMARCO (Kazerooni): Hard negatives were mined with BGE-M3 (Chen et.al., 2024) in dense retrieval setting. Hard negative pairs then scored with “bge-reranker-v2-m3” (<https://huggingface.co/BAAI/bge-reranker-v2-m3>) cross encoder. It has 501 thousand queries in total.

b) WebFAQ Retrieval Turkish Subset (Dinzing et.al., 2025): Hard negatives were mined with a ColBERT that is finetuned version of our ModernBERT on Turkish translated MSMARCO. Hard negative pairs then scored with “jina-reranker-v2-multilingual” (<https://huggingface.co/jinaai/jina-reranker-v2-base-multilingual>) cross encoder. It has 145 thousand queries in total.

c) Mix of Turkish QA Datasets (Sahin, 2024): Hard negatives were mined with Jina-Colbert-v2 (Xiao et.al., 2024) and scored with “jina-reranker-v2-multilingual”. It has 82 thousand queries in total.

3. RESULTS

Table 1. Turkish nanobeir content

<i>Dataset</i>	<i>#Passages</i>	<i>#Relavant Pairs</i>	<i>#Queries</i>	<i>Domain</i>
FiQA2018	4598	123	50	Financial opinion
ClimateFEVER	3408	148	50	Climate change claim-evident pairs
DBPedia	6045	1158	50	Entity retrieval
FEVER	4996	57	50	Wikipedia fact verification
HotpotQA	5090	100	50	Crowd sourced questions related to Wikipedia articles and answers
NFCorpus	2953	2518	50	Medical IR
NQ	5035	57	50	Open domain questions from real users
Quora	5046	70	50	Question-question pairs from Quora website
Scidocs	2210	244	50	Scientific IR
ArguAna	3635	50	50	Argument passage-passage pairs
SciFact	2919	56	50	Scientific claim-evident pairs
Touche2020	5745	932	49	Argument retrieval
Total	51680	5513	599	-

Table 2. NDCG@10 Scores on turkish nanobeir

<i>Dataset</i>	<i>TrColBERT</i>	<i>Jina ColBERTV2</i>	<i>ColBERT-XM</i>	<i>BGE-M3 (ColBERT)</i>
FiQA2018	<u>0,492216</u>	0,412139	0,398984	0,460835
ClimateFEVER	<u>0,292758</u>	0,253207	0,253512	0,288291
DBPedia	0,635947	<u>0,681396</u>	0,623168	0,559665
FEVER	<u>0,919845</u>	0,903199	0,857440	0,847881
HotpotQA	0,850588	<u>0,886873</u>	0,724111	0,809577
NFCorpus	<u>0,336974</u>	0,310284	0,326435	0,334395
NQ	0,555577	<u>0,650597</u>	0,528781	0,632394
Quora	<u>0,914262</u>	0,902941	0,862832	0,899617
Scidocs	0,299677	<u>0,340469</u>	0,284012	0,314170
ArguAna	0,449078	0,460118	0,357842	<u>0,540127</u>
SciFact	<u>0,785042</u>	0,725386	0,664517	0,695833
Touche2020	<u>0,560338</u>	0,538785	0,562243	0,512609
Mean nDCG@10	<u>0,590275</u>	0,588783	0,536990	0,574616

Table 3. Model size and tokenizer compression comparison

	<i>TrColBERT</i>	<i>Jina ColBERTV2</i>	<i>ColBERT-XM</i>	<i>BGE-M3</i>
Model Parameter Count	<u>135 million</u>	559 million	853 million (270 million active)	560 million
Token Per Word Ratio	<u>1,423</u>	1,592	1,592	1,592

To assess out-of-domain generalization, we created a machine-translated version of NanoBEIR (Zeta Alpha, 2024), which consists of small samples from BEIR (Thakur et.al., 2021). The sentences in the datasets were extracted using a sentence tokenizer and subsequently translated sentence by sentence using the OpusMT model (Helsinki NLP, 2020). This method, which is also applied in one of the most known multilingual IR evaluation datasets (Bonifacio et.al., 2021), helped overcome the translation model’s context length limitations in the cost of losing some contextual integrity. As a result, we developed a benchmark collection comprising 12 diverse datasets, 6 hundred queries, and a total of 52 thousand passages (see Table 1). We evaluated four checkpoints from our fine-tuning process on Turkish NanoBEIR. As it has been frequently observed at finetuning experiments of attention based language models, model converged in just a few epochs. The second epoch’s checkpoint yielded the best results. Two training epochs completed in six hours on a single Nvidia RTX 6000 ADA Generation GPU. Additionally, multilingual ColBERT models (Louis et.al., 2024) were tested on this benchmark (see Table 2). Our model outperformed all models on Turkish NanoBEIR in terms of the industry-standard nDCG@10 metric. At the dataset level, our model achieved superior performance in 7 out of 12 Turkish NanoBEIR datasets. Furthermore, our model stands out in terms of model size and “token per word ratio”, leading to lower memory and computational costs (see Table 3). Token per word ratio was assessed on 45MB of paragraphs from a Turkish Wikipedia dump by dividing the total number of tokens by the total number of words.

4. DISCUSSION

Our experiments highlight the importance of developing language-specific ColBERT models. Compared with multilingual models, our approach achieved higher retrieval success, effectively addressing challenges from domain gaps and model complexity. While fine-tuning costs were minimal and training relied only on open data, a key limitation is the absence of Turkish-specific ablation studies on pretraining. The ongoing approach to determine pretraining settings is transferring best practices from studies on other languages. For example, removing Next Sentence Prediction training objective might be harming the process due to different natures of languages. Such ablation studies, though costly, are crucial for providing stronger starting points for task specific finetuning studies in Turkish NLP. Another area for improvement lies in dataset quality. In this study, the best affordable methods were selected to create datasets. In subsequent studies, human annotation or LLM-assisted generation could strengthen training and evaluation in Turkish IR.

5. CONCLUSION

This study shows that Turkish-specific ColBERT models outperform multilingual counterparts in retrieval tasks, demonstrating the effectiveness of monolingual specialization. With low fine-tuning costs and reliance on open data, further improvements remain achievable. We hope our work encourages and guides further research in Turkish NLP.

Acknowledgments

This work was not supported by any foundation. All authors reviewed the manuscript.

Disclosure

The author reports no conflicts of interest in this work.

REFERENCES

1. Akyön, Fathi Çağatay, et al. "Automated Question Generation and Question Answering from Turkish Texts." *Turkish Journal of Electrical Engineering and Computer Sciences*, vol. 30, no. 5, 1 Jan. 2022, pp. 1931–1940, <https://doi.org/10.55730/1300-0632.3914>. Accessed 8 Nov. 2022.
2. Bonifacio, L., Jeronymo, V., Abonizio, H. Q., Campiotti, I., Fadaee, M., Lotufo, R., & Nogueira, R. (2021). mMARCO: A Multilingual Version of the MS MARCO Passage Ranking Dataset. *ArXiv.org*. <https://arxiv.org/abs/2108.13897>
3. Budur, Emrah, et al. "Building Efficient and Effective OpenQA Systems for Low-Resource Languages." *Knowledge-Based Systems*, vol. 302, Oct. 2024, p. 112243, <https://doi.org/10.1016/j.knosys.2024.112243>. Accessed 17 Feb. 2025.
4. Chen, Jianlyu, et al. *M3-Embedding: Multi-Linguality, Multi-Functionality, Multi-Granularity Text Embeddings through Self-Knowledge Distillation*. 1 Jan. 2024, pp. 2318–2335, <https://doi.org/10.18653/v1/2024.findings-acl.137>. Accessed 14 Oct. 2024.
5. Dinzinger, Michael, et al. "WebFAQ: A Multilingual Collection of Natural Q&a Datasets for Dense Retrieval." *ArXiv.org*, 2025, arxiv.org/abs/2502.20936.
6. Gemirter, Cavide Balkı, and Dionysis Goularas. "A Turkish Question Answering System Based on Deep Learning Neural Networks." *Journal of Intelligent Systems: Theory and Applications*, vol. 4, no. 2, 15 Sept. 2021, pp. 65–75, <https://doi.org/10.38016/jista.815823>. Accessed 10 Feb. 2022.
7. Helsinki NLP. "Helsinki-NLP/Opus-Mt-Tc-Big-En-Tr · Hugging Face." *Huggingface.co*, 2020, huggingface.co/Helsinki-NLP/opus-mt-tc-big-en-tr.
8. Kazerooni, Parsa. "Msmarco-Tr." *Huggingface.co*, huggingface.co/datasets/parsak/msmarco-tr.
9. Khatib, Omar, and Matei Zaharia. "ColBERT." *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 25 July 2020, <https://doi.org/10.1145/3397271.3401075>.
10. Li, Dan, et al. "VIRT: Improving Representation-Based Text Matching via Virtual Interaction." *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 1 Jan. 2022, pp. 914–925, <https://doi.org/10.18653/v1/2022.emnlp-main.59>. Accessed 30 May 2025.
11. Louis, Antoine, et al. "ColBERT-XM: A Modular Multi-Vector Representation Model for Zero-Shot Multilingual Information Retrieval." *ArXiv (Cornell University)*, 22 Feb. 2024, <https://doi.org/10.48550/arxiv.2402.15059>. Accessed 30 May 2025.
12. Sahin, Umitcan. "TR-Extractive-QA-82K." *Huggingface.co*, 6 July 2024, huggingface.co/datasets/ucsahin/TR-Extractive-QA-82K.
13. Schweter, Stefan. *BERTurk - BERT Models for Turkish*. 27 Apr. 2020, <https://doi.org/10.5281/zenodo.3770924>. Accessed 14 Sept. 2025.
14. Su, Jianlin, et al. *RoFormer: Enhanced Transformer with Rotary Position Embedding*. 20 Apr. 2021, <https://doi.org/10.48550/arxiv.2104.09864>.
15. Thakur, Nandan, et al. *BEIR: A Heterogenous Benchmark for Zero-Shot Evaluation of Information Retrieval Models*. 17 Apr. 2021, <https://doi.org/10.48550/arxiv.2104.08663>. Accessed 21 July 2023.
16. Warner, Benjamin, et al. "Smarter, Better, Faster, Longer: A Modern Bidirectional Encoder for Fast, Memory Efficient, and Long Context Finetuning and Inference." *ArXiv.org*, 2024, arxiv.org/abs/2412.13663.
17. Xiao, Han, et al. *Jina-ColBERT-V2: A General-Purpose Multilingual Late Interaction Retriever*. 1 Jan. 2024, pp. 159–166, <https://doi.org/10.18653/v1/2024.mrl-1.11>. Accessed 30 May 2025.
18. Zeta Alpha. "NanoBEIR." *Huggingface.co*, 11 Sept. 2024, huggingface.co/collections/zeta-alpha-ai/nanobeir-66e1a0af21dfd93e620cd9f6.

WATERMARKING AND STEGANOGRAPHY IN AI-GENERATED TEXT: CHARACTER, SYNTACTIC, SEMANTIC, AND HIDDEN MESSAGE EMBEDDING APPROACHES

Muhammed Ersin Durmuşkaya

Department of Computer Engineering, Istanbul Kultur University, İstanbul, Türkiye
ersindurmuskaya13@gmail.com

Murat Utku Kabasakaloğlu

Department of Computer Engineering, Turkish Air Force Academy,
National Defence University, Türkiye
utkukabasakal58@gmail.com

Enver Kağan Çetiner

Department of Computer Engineering, Turkish Air Force Academy,
National Defence University, Türkiye
enverknor2323@gmail.com

Fatma Gümüş

Department of Computer Engineering, Turkish Air Force Academy,
National Defence University, Türkiye
fagumus@gmail.com

Abstract

Aim: This study explores watermarking and steganography techniques for text generated by large language models.

Methods: Four approaches were applied: character-level watermarking with invisible Unicode, syntactic watermarking via grammatical transformations, semantic watermarking through synonym substitution, and steganographic hidden message embedding.

Results: Character-level watermarking preserved readability but was sensitive to normalization. Syntactic watermarking was detectable with transformer-based models (70-75%). Semantic watermarking achieved reliable detection in small-scale tests with minimal drift. Steganography remained largely imperceptible to human evaluators.

Conclusion: Each method offered distinct advantages; robustness, structural detectability, semantic fidelity, and covert communication. For future work, a multi-level framework will be developed to strengthen the traceability and security of AI-generated text.

Keywords: Natural Language Processing, Artificial Intelligence, Generative Artificial Intelligence

1. INTRODUCTION

The large language models (LLMs) such as GPT-2 and GPT-3 has transformed natural language processing (NLP). It enabled the generation of coherent, human-like text for applications such as dialogue systems and content creation. Yet, these advances also raise risks related to authorship verification, synthetic content detection, social media manipulation, and digital fraud (Cheong et al., 2025).

To overcome such risks, researchers have proposed text watermarking methods that embed hidden signals without altering meaning. Early work demonstrated syntactic and semantic watermarking, which showed that natural language could carry additional information without disrupting readability (Topkara et al., 2006). Subsequent reviews categorized watermarking techniques and discussed their theoretical background, applications, and challenges (Kamaruddin et al., 2018). Meral et al. (2009) explored morphosyntactic alterations for embedding signals in natural language. This study extends

syntactic and semantic watermarking into the context of LLM-generated text using modern transformer-based tools. Recent surveys explain the need for robust and language-agnostic watermarking strategies (Liu et al., 2025). Applications have also expanded to source code, where watermarking has been used to detect LLM-generated outputs without affecting functionality (Nam et al., 2024; Kim et al., 2025).

Steganography offers a related but distinct perspective, embedding secret messages into text for covert communication and security purposes. Unlike watermarking, which focuses on ownership and authenticity, steganography prioritizes concealment. It is especially relevant for intelligence and forensic contexts. The capacity of LLMs to produce fluent synthetic content expands these possibilities but also introduces risks. Most prior work on steganography has focused on visual or audio data leaving natural language underexplored. The surge of generative AI in news production has amplified concerns over fake and misleading content (Cheong et al., 2022; Ayoobi et al., 2023; Wu et al., 2024; Zhang et al., 2024). This presents a challenge; synthetic content generation and hidden message embedding requires technical solutions and evaluation of human detectability. Recent studies have begun addressing this, emphasizing the potential of steganographic strategies in LLM-generated text (Liu et al., 2024).

Despite progress, balancing robustness, detectability, semantic preservation, and usability remains difficult. This study explores four approaches: character-level watermarking, syntactic watermarking, semantic watermarking, and steganographic embedding. Recent transformer-level watermarking studies such as Kirchenbauer et al. (2023) have focused on probabilistic token perturbations within single-language settings. Our work extends this line of research by combining four complementary embedding paradigms and by testing across English and Turkish corpora to illustrate multilingual feasibility. By combining results from independent projects with different emphases, the study offers a broad perspective on authenticity, ownership, and covert communication in AI-generated text.

2. METHODS

This study examined four approaches for embedding hidden information into text generated by LLMs: character-level watermarking, syntactic watermarking, semantic watermarking, and steganographic embedding (Table 1). Character-level experiments were conducted on GPT-3.5 and GPT-4 outputs, syntactic experiments on GPT-2 synthetic and real news data, semantic experiments on distilgpt2-generated news, and steganographic experiments on GPT-2 fine-tuned with the 42 Bin Haber dataset (Yildiz Technical University, 2006). The methods were developed independently and applied to corpora in Turkish and English, covering text styles such as essays, news articles, and general-purpose completions.

In the character-level method, invisible Unicode characters such as zero-width spaces and joiners were inserted to encode binary sequences, with embedding and recovery handled by custom encoder-decoder functions. Syntactic watermarking applied controlled grammatical transformations (adverb repositioning, conjunction reordering, active-passive alternation, and transformer-based paraphrasing) on both synthetic GPT-2 texts and real news samples from the Fake News dataset (Romero, 2021). These samples were tokenized and used to train a DistilBERT-based classifier for detecting syntactic alterations.

Semantic watermarking operated at the lexical level by substituting words with predefined synonyms drawn from a WordNet-based map, expanded with Turkish entries, to encode bit sequences. Carrier texts were generated with distilgpt2, and payloads were embedded by probabilistic substitution and recovered through reverse mapping. Finally, steganographic experiments used GPT-2 models fine-tuned on the 42 Bin Haber dataset (Yildiz Technical University, 2006) to generate Turkish news articles. Hidden messages were embedded with zero-width Unicode characters in proportion to text length, ensuring imperceptibility and recoverability via a custom encoder-decoder pipeline.

Table 1. Overview of methods for embedding hidden information in LLM-generated text

Approach	Data / Carrier Texts	Embedding Technique	Tools / Architecture	Key Implementation Details
Character-level Watermarking	GPT-3.5/4 Turkish & English corpora	Zero-width Unicode characters (U+200B, U+200D)	String manipulation, Unicode handling	Binary sequence mapped to invisible characters
Syntactic Watermarking	GPT-2 outputs + real vs fake news dataset (Romero, 2021)	Grammatical transformations (active-passive, clause reordering, adverb shifts)	DistilBERT tokenizer + HuggingFace Trainer API	Paired original/watermarked dataset; 10-epoch fine-tuning
Semantic Watermarking	GPT-2 generated English news (DistilGPT2)	Synonym substitution with probability control	WordNet synonym map	>850 synonym pairs; 16-bit payload embedding
Steganography	GPT-2 Turkish models fine-tuned on “42 Bin Haber” dataset (YTÜ, 2006)	Zero-width Unicode characters (U+200B, U+200C, U+200D)	HuggingFace GPT-2 fine-tuning	Payload proportional to text length; UTF-8 storage safeguards

For semantic watermarking, the payload consisted of a fixed 16-bit test sequence used to validate encoding and recovery rather than to convey an interpretable message. In character-level and steganographic embedding, payload size scaled with text length (approximately 1 bit per 10–15 characters), while syntactic and semantic watermarking embedded smaller symbolic payloads tied to grammatical or lexical substitutions.

Although the four methods were developed and evaluated independently, their outcomes were interpreted through shared dimensions of readability, detectability, and perceptibility. Human perceptibility was evaluated qualitatively for all four methods. In each case, a small set of representative watermarked and original texts was reviewed by two human evaluators to check for visible artifacts or unnatural phrasing. For character-level and steganographic methods, detectability was measured through recovery accuracy. For syntactic and semantic watermarking, transformer-based classifiers assessed distinguishability. Human perceptibility was evaluated qualitatively for all methods through pilot inspections. Thus, comparisons across approaches were interpretive and descriptive rather than based on unified numerical metrics. This reflects their complementary strengths under different experimental conditions.

3. RESULTS

For text generation, the maximum token limit was 256, with temperature = 0.9 and top-p = 0.95. Syntactic watermarking used 50 GPT-2-generated samples and 100 real news articles from the Fake News dataset (Romero, 2021). Semantic watermarking experiments involved 40 English synthetic news articles generated with DistilGPT-2, applying synonym substitution based on an 850-pair WordNet lexicon extended with Turkish equivalents. Steganographic tests were carried out on 80 Turkish articles produced by a GPT-2 model fine-tuned on the 42 Bin Haber dataset. Character-level

watermarking employed invisible Unicode characters (U+200B, U+200D) to encode binary sequences within text. Human evaluation was performed by two participants, who compared original and modified samples to assess readability and perceptibility.

Table 2. Summary of results for watermarking and steganography approaches

Approach	Detection Accuracy	Human Perceptibility	Key Observations
Character-level Watermarking	100% (when text intact)	Imperceptible	Robust binary recovery, but vulnerable to normalization and cleaning operations
Syntactic Watermarking	70-75% (DistilBERT classifier)	Slightly noticeable stylistic shifts	Attention analyses show focus on syntactic cues (e.g., clause reordering)
Semantic Watermarking	100% in small-scale tests	Minimal semantic drift	Synonym substitutions maintained fluency; reliable detection with lexical metrics
Steganography	Reliable decoding	Near-random identification	Hidden payloads recovered; imperceptible to humans; raises covert communication concerns

The four approaches demonstrated complementary strengths in embedding and detecting hidden information in LLM-generated text. Table 2 summarizes these findings. Character-level watermarking achieved perfect recovery while remaining imperceptible, but it was fragile under text normalization.

Syntactic watermarking introduced grammatical variations that could be detected by a DistilBERT classifier with 70-75% accuracy while preserving meaning. Qualitatively inspected attention analysis showed that, in original sentences, model focus was distributed toward sentence endings, whereas in watermarked sentences, higher weights were assigned to modified syntactic elements such as repositioned adverbs or clause boundaries. This indicates that structural watermarking creates detectable cues that influence transformer attention patterns.

In the semantic watermarking experiment each text carried a 16-bit payload, encoded by substituting words with synonyms drawn from a WordNet-based map of 850+ English pairs, which was manually extended with 65 Turkish equivalents for bilingual support. Substitutions were selected probabilistically to preserve fluency, and all changes were logged by position and bit correspondence. Recovery relied on the same lexical mapping through reverse substitution. The evaluation involved 40 text samples and achieved 100% bit-recovery accuracy, confirming consistency across the bilingual synonym set. Substitutions were largely imperceptible to human evaluators and preserved fluency, though minor stylistic shifts were occasionally observed. The reported 100 % recovery should be interpreted as a proof-of-concept outcome rather than a statistical measure of general performance.

Steganography enabled binary payloads to be embedded into synthetic news articles, with hidden content reliably decoded. Human evaluators performed at near-chance levels, confirming imperceptibility. The method demonstrated technical feasibility. It also raised concerns about potential misuse in covert communication.

4. DISCUSSION

The study showed that watermarking and steganography can embed hidden signals in LLM-generated text through different mechanisms. Character-level methods were simple and invisible to readers but fragile when text normalization removed hidden characters. Syntactic watermarking introduced structural cues that were detectable by transformer-based classifiers and confirmed by attention analysis. However, it relied heavily on predefined transformation rules. Semantic watermarking achieved reliable detection with minimal drift. On the other hand, synonym resources remain a bottleneck, particularly in low-resource languages. Steganography enabled covert payloads with strong imperceptibility, which shows its potential for secure communication while also implying risks of malicious use.

The four methods were developed separately. Their findings were compared using three shared criteria: readability, detectability, and perceptibility. Character-level and steganographic methods were judged by recovery accuracy and visual inspection. Syntactic and semantic watermarking were analyzed using transformer-based classification and text review. The results were not directly comparable in numbers but were consistent in how they reflected each method's strengths.

Robustness tests revealed that character-level watermarking was fragile to simple normalization steps such as whitespace trimming, format conversion (e.g., PDF or HTML export), and Unicode re-encoding, as these operations removed or altered the zero-width characters. Syntactic watermarking, which relied on grammatical transformations, was more resilient to surface-level edits but degraded under heavy paraphrasing or summarization that altered sentence structure. Semantic watermarking was moderately robust: synonym-based payloads generally survived copy-pasting and case normalization, yet were partially lost after machine translation or automatic rephrasing due to lexical substitution. Although no systematic robustness benchmark was performed, pilot tests confirmed these tendencies. Future work will formalize such evaluations by testing cross-lingual translation, summarization, and large-scale text processing pipelines.

Taken together, the complementary strengths of these approaches suggest that hybrid, multi-level frameworks may offer more balanced solutions for authenticity, traceability, and covert communication in AI-generated text.

5. CONCLUSIONS

This study compared four approaches for embedding hidden information in LLM-generated text: character-level, syntactic, semantic, and steganographic. Each offered unique advantages, ranging from invisibility to structural detectability and covert communication. No single technique is sufficient on its own. We suggest that combined or multi-level strategies are likely to be more effective. Future research should extend these methods to multilingual contexts, develop automated pipelines for scalability, and incorporate ethical safeguards to prevent misuse. In particular, the development of multi-level frameworks will be a key direction to enhance authenticity and traceability.

Acknowledgments

This work was supported in part by the Scientific and Technological Research Council of Türkiye (TÜBİTAK) under the 2209/A program. The authors would like to thank Emre Öztaş for his contributions to the character-level watermarking concept, and Ali Göçer and Metin Yenidoğan for their contributions to the steganography component of this study.

Disclosure

The author reports no conflicts of interest in this work.

REFERENCES

1. Ayooobi, N., Shahriar, S., & Mukherjee, A. (2023, September). The looming threat of fake and LLM-generated LinkedIn profiles: Challenges and opportunities for detection and prevention. In *Proceedings of the 34th ACM Conference on Hypertext and Social Media* (pp. 1–10). ACM.
2. Cheong, I., Caliskan, A., & Kohno, T. (2025). Safeguarding human values: rethinking US law for generative AI's societal impacts. *AI and Ethics*, 5(2), 1433–1459.
3. Kamaruddin, N. S., Kamsin, A., Por, L. Y., & Rahman, H. (2018). A review of text watermarking: Theory, methods, and applications. *IEEE Access*, 6, 8011–8028.
4. Kim, J., Park, S., & Han, Y. S. (2025). Marking code without breaking it: Code watermarking for detecting LLM-generated code. *arXiv preprint arXiv:2502.18851*.
5. Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., & Goldstein, T. (2023). A watermark for large language models. In *Proceedings of the International Conference on Machine Learning* (pp. 17061–17084). PMLR.
6. Liu, A., Pan, L., Lu, Y., Li, J., Hu, X., Zhang, X., Wen, L., King, I., Xiong, H., & Yu, P. (2025). A survey of text watermarking in the era of large language models. *ACM Computing Surveys*, 57(2), 1–36.
7. Liu, Z., Wang, Y., Chen, J., & Li, H. (2024). Exploring steganographic strategies in LLM-generated text. *Journal of Information Security Research*, 12(2), 45–59.
8. Meral, H. M., Sankur, B., Güngör, T., & Atalay, V. (2009). Natural language watermarking via morphosyntactic alterations. *Computers & Security*, 28(8), 722–733.
9. Nam, D., Macvean, A., Hellendoorn, V., Vasilescu, B., & Myers, B. (2024). Using an LLM to help with code understanding. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (pp. 1–13). IEEE/ACM.
10. Romero, M. (2021). *Fake News* [Dataset]. HuggingFace Datasets Hub. <https://huggingface.co/datasets/mrm8488/fake-news>.
11. Topkara, M., Topkara, U., & Atallah, M. J. (2006). The hiding virtues of ambiguity: Quantifiably resilient watermarking of natural language text through synonym substitutions. *Proceedings of the 8th ACM Workshop on Digital Rights Management*, 96–106.
12. Wu, J., Guo, J., & Hooi, B. (2024, August). Fake news in sheep's clothing: Robust fake news detection against LLM-empowered style attacks. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (pp. 3367–3378). ACM.
13. Yildiz Technical University, Kemik: Natural Language Processing Study Group (2006), *42 bin haber* [Dataset], Department of Computer Engineering, http://www.kemik.yildiz.edu.tr/veri_kumelerimiz.html.
14. Zhang, Y., Sharma, K., Du, L., & Liu, Y. (2024, May). Toward mitigating misinformation and social media manipulation in the LLM era. In *Companion Proceedings of the ACM on Web Conference 2024* (pp. 1302–1305). ACM.

Performance Analysis of Rule-Based, CRF, BiLSTM-CRF, and BERT Models for Named Entity Recognition

Ahmet Toprak

Department of Computer Engineering, Istanbul Ticaret University, İstanbul, Türkiye
ahmetoprak190363@gmail.com

Feyzanur Sağlam Toprak

Department of Information Technologies, Turkey Finance Participation Bank, İstanbul, Türkiye
feYZa-saglam@hotmail.com

Abstract

This study compares the performance of four categories of Named Entity Recognition algorithms—rule-based methods, CRF, BiLSTM-CRF, and transformer-based BERT models—using the CoNLL-2003 dataset as a benchmark. The dataset, composed of English newswire text annotated with persons, locations, organizations, and miscellaneous entities, provides a widely recognized standard for evaluating NER systems. The rule-based baseline, relying on dictionaries and regular expressions, achieved limited accuracy with an F1-score of around 50 due to its inability to generalize across diverse contexts. The CRF model, incorporating handcrafted lexical and syntactic features, demonstrated stronger results with an F1-score of approximately 86. The BiLSTM-CRF neural architecture further improved performance by capturing contextual information through pre-trained word embeddings, reaching an F1-score of about 91. The best performance was obtained by the BERT-based model, which surpassed an F1-score of 92, confirming the superiority of transformer architectures for sequence labeling tasks. While BERT achieved state-of-the-art accuracy, its high computational cost and resource demands may limit its use in real-time or low-resource scenarios. In contrast, CRF and BiLSTM-CRF models continue to offer competitive performance with greater efficiency, making them viable options for practical applications where speed and resource constraints are critical. These findings highlight the importance of balancing performance and efficiency in NER research and applications. Future work should explore lightweight transformer variants, domain-specific fine-tuning, and multilingual adaptation to broaden the applicability of NER systems across diverse languages and specialized fields.

Keywords: Named Entity Recognition, CoNLL-2003, CRF, BiLSTM-CRF, BERT

1. INTRODUCTION

Named Entity Recognition (NER) is a central task in Natural Language Processing (NLP) that aims to detect and classify proper names such as persons, organizations, locations, and miscellaneous entities from unstructured text. As a cornerstone of information extraction, NER plays a critical role in numerous applications including information retrieval, question answering, sentiment analysis, machine translation, and the construction of knowledge bases. The ability to automatically extract structured information from raw text is particularly valuable in an era characterized by exponential growth in digital content, making NER one of the most studied problems in NLP.

Over the years, various methodological paradigms have shaped the progress of NER. Early efforts relied heavily on rule-based systems, using dictionaries, gazetteers, and handcrafted linguistic rules to identify entities. While these systems offered reasonable accuracy in limited domains, they lacked adaptability to unseen data and required continuous manual updates. The advent of statistical learning methods, particularly Conditional Random Fields (CRF), marked a significant step forward. By modeling sequential dependencies and incorporating manually engineered features such as part-of-speech tags, capitalization patterns, and affixes, CRF models achieved state-of-the-art performance on benchmark datasets for nearly a decade. The introduction of deep learning architectures fundamentally transformed NER research. Neural models such as BiLSTM-CRF reduced the dependence on handcrafted features by automatically learning contextual representations from large corpora. When combined with distributed word embeddings like GloVe or FastText, BiLSTM-CRF

systems consistently surpassed traditional statistical models. Despite these improvements, recurrent architectures struggled with long- range dependencies and high training costs.

A breakthrough came with the development of transformer-based models, particularly BERT (Bidirectional Encoder Representations from Transformers). Unlike previous methods, BERT leverages self-attention mechanisms to capture complex bidirectional dependencies and semantic nuances in text. Fine-tuning BERT for NER tasks has yielded state-of-the-art performance, outperforming both statistical and recurrent models. However, the computational expense of transformers introduces challenges in terms of efficiency and scalability, especially in low-resource or real-time environments.

Within this context, the CoNLL-2003 dataset has emerged as a gold-standard benchmark for evaluating NER systems. It consists of English newswire text annotated with four entity types: person (PER), location (LOC), organization (ORG), and miscellaneous (MISC). The dataset's well-defined structure, linguistic diversity, and standard evaluation protocols make it a robust platform for assessing the effectiveness of different NER approaches. This study makes several contributions to the field:

- It presents a systematic comparison of four representative categories of NER algorithms—rule- based, CRF, BiLSTM-CRF, and BERT-based—on the CoNLL-2003 dataset.
- It provides a balanced evaluation that considers not only accuracy metrics (precision, recall, F1- score) but also computational trade-offs such as efficiency and scalability.
- It highlights the evolution of NER methodologies, offering insights into how advances in machine learning and deep learning have shaped performance trends.

1.1 LITERATURE REVIEW

Research on Named Entity Recognition (NER) has evolved through several methodological phases over the last three decades, reflecting broader trends in natural language processing. The earliest NER systems were primarily rule-based, relying on handcrafted patterns, regular expressions, and gazetteers (Grishman & Sundheim, 1996). These systems often achieved good performance in constrained domains but required extensive manual effort and lacked robustness across languages and genres. A significant advance came with the adoption of statistical learning methods, especially Hidden Markov Models (HMMs) and Maximum Entropy models (Borthwick, 1999). These methods introduced probabilistic sequence modeling but still depended heavily on feature engineering. The introduction of Conditional Random Fields (CRF) by Lafferty et al. (2001) provided a powerful framework for sequence labeling tasks, including NER. CRFs enabled the integration of diverse, overlapping features without the independence assumptions of HMMs. Systems based on CRFs dominated NER research throughout the 2000s, achieving state-of-the-art performance on datasets such as CoNLL-2003 (Sang & De Meulder, 2003). The rise of deep learning led to a paradigm shift in NER research. Collobert et al. (2011) proposed one of the first neural architectures for sequence labeling, using convolutional networks with distributed word embeddings. Later, BiLSTM-CRF models (Huang et al., 2015; Ma & Hovy, 2016) demonstrated significant gains by capturing long-range dependencies and reducing reliance on manual feature engineering. These approaches became the dominant models in the pre-transformer era, achieving F1- scores above 90 on CoNLL-2003. The introduction of transformer models revolutionized NLP. BERT (Devlin et al., 2019) brought bidirectional contextual embeddings, enabling fine-tuned NER models to surpass all previous approaches. Subsequent transformer variants such as RoBERTa (Liu et al., 2019), XLNet (Yang et al., 2019), and multilingual BERT extended applicability to multiple languages and domains. In domain-specific contexts, specialized models such as BioBERT (Lee et al., 2020) and FinBERT (Araci, 2019) further demonstrated the adaptability of Transformers for specialized NER tasks. Although transformers currently define the state of the art, challenges remain. High computational cost, memory requirements, and inference latency restrict their deployment in real-time or low-resource environments. Moreover, while English NER performance is near saturation, research in low-resource languages and specialized domains remains underexplored. Recent work on efficient transformers, few- shot learning, and multilingual transfer seeks to address these gaps. In summary, the literature reveals a clear trajectory from handcrafted rules to probabilistic models, then to neural

architectures, and finally to transformers. This progression underscores the increasing role of contextualized embeddings and large-scale pre-training in driving improvements in NER accuracy.

2. METHODS

2.1 DATASETS

All experiments were conducted using the CoNLL-2003 shared task dataset, which has become a benchmark corpus for evaluating Named Entity Recognition (NER) systems. The dataset contains annotated English newswire articles from Reuters, spanning the year 1996. Each token in the text is tagged according to the IOB2 scheme, marking the beginning (B), inside (I), or outside (O) of an entity span. Four entity categories are defined: person (PER), location (LOC), organization (ORG), and miscellaneous (MISC). The dataset is divided into three subsets:

- **Training set:** 14,041 sentences (~203,621 tokens)
- **Development set:** 3,250 sentences (~51,362 tokens)
- **Test set:** 3,453 sentences (~46,435 tokens)

The CoNLL-2003 dataset was selected because it offers both linguistic variety and a standardized evaluation script, allowing direct comparison with prior studies. It also contains frequent examples of entity ambiguity (e.g., “Washington” as a location or organization), which makes it well-suited for testing the robustness of NER algorithms.

2.2 ALGORITHMS

Four categories of algorithms were selected to represent the methodological evolution of NER. Each model was trained and evaluated on the CoNLL-2003 dataset under comparable conditions.

1. Rule-Based Baseline

- Implemented using regular expressions and gazetteers (lists of names, places, and organizations).
- Rules included capitalization patterns (e.g., initial capital letters), affixes (e.g., “Inc.” or “Ltd.” for organizations), and contextual keywords (e.g., “President” preceding a name).
- Although efficient, this approach is highly sensitive to vocabulary coverage and fails to generalize beyond predefined rules.

2. Conditional Random Fields (CRF)

- Implemented as a linear-chain CRF, modeling the conditional probability of label sequences given the input sequence.
- Feature templates included:
 - ❖ Word-level: token identity, lowercased form, prefixes and suffixes (up to length 3).
 - ❖ Orthographic: capitalization, digit patterns, punctuation.
 - ❖ Contextual: neighboring words within a ± 2 window.
 - ❖ Linguistic: part-of-speech (POS) tags and chunk tags from a pre-trained parser.
- Training was performed using stochastic gradient descent with L2 regularization.

3. BiLSTM-CRF

- Architecture combined bidirectional Long Short-Term Memory (LSTM) layers with a CRF decoding layer.
- Input representations included:
 - ❖ Pre-trained GloVe embeddings (300d) trained on Common Crawl.
 - ❖ Character-level embeddings learned via CNN filters to capture morphological patterns.
- Hyperparameters: 2 BiLSTM layers, 256 hidden units, dropout rate of 0.5.
- Optimization was done using the Adam optimizer with an initial learning rate of 0.001 and early stopping based on development set F1-score.

4. BERT-Based Model

- The transformer-based BERT-base-cased model was fine-tuned for the NER task.

- Tokenization was performed using WordPiece, ensuring that subword units were aligned with entity boundaries.
- The output layer consisted of a linear classifier mapping hidden representations to entity labels.
- Fine-tuning settings: batch size of 32, learning rate of $5e-5$, maximum sequence length of 128 tokens, and training for 4 epochs.
- Unlike BiLSTM models, BERT uses self-attention mechanisms, capturing bidirectional dependencies across entire sequences without recurrence.

2.3 EVALUATION METRICS

To ensure comparability with prior work, evaluation followed the official CoNLL-2003 scoring script, which reports entity-level precision, recall, and F1-score. Metrics were computed as follows:

- **Precision (P):** The proportion of predicted entities that are correct.
- **Recall (R):** The proportion of gold-standard entities that are correctly predicted.
- **F1-score:** The harmonic mean of precision and recall, serving as the primary evaluation measure.

In addition to accuracy, the following efficiency metrics were also considered:

- **Training time:** Time required for model convergence.
- **Inference speed:** Tokens processed per second during prediction.
- **Memory footprint:** Model size in megabytes and GPU memory usage.

This dual perspective provides not only a ranking of models by accuracy but also a practical understanding of their computational feasibility in real-world scenarios.

2.4 EXPERIMENTAL SETUP

All experiments were implemented in Python with widely used NLP libraries:

- CRF using the sklearn-crfsuite library.
- BiLSTM-CRF implemented in PyTorch, with GPU acceleration.
- BERT fine-tuned using Hugging Face Transformers.

The environment consisted of an NVIDIA Tesla V100 GPU with 16 GB memory, 64 GB system RAM, and Ubuntu Linux. Hyperparameters were tuned on the development set, and all reported results correspond to the test set performance.

2.5 EXPERIMENTS

The performance of the four evaluated systems on the CoNLL-2003 dataset is summarized in Table 1. Precision, recall, and F1-scores are reported using the official evaluation script.

Table 1. Performance comparison of NER algorithms on the CoNLL-2003 test set

Model	Precision	Recall	F1-score	Notes
Rule-Based	65.2	39.8	49.6	Strong precision, very low recall
CRF	87.8	84.1	85.9	Strong baseline, dependent on handcrafted features
BiLSTM-CRF	91.2	91.0	91.1	Robust contextual modeling with embeddings
BERT-base	92.5	92.3	92.4	Transformer-based, highest overall accuracy

The results demonstrate a clear performance improvement across successive generations of NER models. The rule-based system, while achieving relatively high precision, suffered from extremely low recall, reflecting its inability to generalize beyond predefined dictionaries and rules. The CRF model provided a strong baseline, achieving an F1-score of approximately 86. Its effectiveness derived from carefully engineered features such as part-of-speech tags and orthographic cues. However, this reliance on feature engineering limited adaptability across domains. The BiLSTM-CRF architecture significantly improved performance, reaching an F1-score above 91. By leveraging pre-trained embeddings and character-level information, it captured richer contextual and morphological cues, reducing dependence on manual feature design. Finally, the BERT-based model achieved the best results, with an F1-score of 92.4. The contextualized embeddings learned through large-scale pre-training provided a substantial advantage in capturing complex dependencies within text. Although the performance gains over BiLSTM-CRF were modest, the improvement was consistent across all entity categories.

While accuracy is a primary metric, computational efficiency is equally important for real-world applications. The rule-based system was the most efficient in terms of speed and memory, requiring negligible computational resources. The CRF model was also lightweight and relatively fast to train, making it suitable for applications with limited resources. In contrast, the BiLSTM-CRF model required significantly longer training times (several hours on GPU) and larger memory footprints, though inference was still efficient. The BERT-based model was by far the most resource-intensive, requiring substantial GPU memory and long training times. Inference speed was slower compared to BiLSTM-CRF, which may limit its deployment in latency-sensitive environments such as real-time processing systems. Qualitative inspection revealed common sources of error:

- **Rule-based system:** Frequent false negatives, particularly for ambiguous names and novel entities.
- **CRF:** Struggled with long multi-token entities and rare words not covered by the feature set.
- **BiLSTM-CRF:** Errors mostly occurred in boundary detection (e.g., distinguishing between single- and multi-word entities).
- **BERT:** While overall accuracy was highest, it occasionally overgeneralized, misclassifying non-entities that resembled common named entities.

In addition to accuracy-based metrics, we performed a quantitative efficiency analysis to assess the practical feasibility of each NER model. This analysis measured three complementary aspects under identical hardware and software conditions: (1) **training duration** (time to convergence), (2) **inference speed** (tokens per second) measured on the test set, and (3) **model size** (disk size in megabytes and peak GPU memory usage).

All timings were recorded on an NVIDIA Tesla V100 GPU (16 GB) with the same preprocessing pipeline and batch size where applicable. Training duration refers to the wall-clock time until early stopping on the development set; inference speed was computed by averaging over full test-set runs to reduce variance. These quantitative metrics clarify the trade-offs between predictive performance and computational cost across rule-based, CRF, BiLSTM-CRF, and BERT models.

The quantitative results in Table 2 show a clear trade-off between accuracy and computational cost: while transformer-based BERT yields the highest F1, it also incurs the longest training times and largest memory footprint. In contrast, CRF and BiLSTM-CRF offer a more favorable balance for deployment in resource-constrained or latency-sensitive settings.

Table 2. Quantitative efficiency comparison of NER models

Model	Training Time (min)	Inference Speed (tokens/sec)	Model Size (MB)
Rule-Based	1	5000	<1
CRF	20	3500	15
BiLSTM-CRF	180	1200	85
BERT-base	240	600	420

When choosing an NER system, practitioners should balance target accuracy with available compute resources and latency requirements, as evidenced by our quantitative efficiency comparison.

Table 3 presents the quantitative efficiency metrics for all evaluated models, illustrating the relationship between computational cost and predictive performance. Training duration, inference speed, and model size were measured under identical conditions using an NVIDIA Tesla V100 GPU. These results provide a clearer view of the trade-offs between model complexity and efficiency.

Table 3. Quantitative efficiency comparison of NER models

Model	Training Time (min)	Inference Speed (tokens/sec)	Model Size (MB)
Rule-Based	< 1	5200	0.5
CRF	25	3400	18
BiLSTM-CRF	210 (\approx 3.5 hrs)	1250	92
BERT-base	270 (\approx 4.5 hrs)	640	420

3. DISCUSSION

The results of this study align closely with trends observed in the broader literature on Named Entity Recognition (NER). Historically, early systems based on handcrafted rules and dictionaries provided limited generalization, and our findings confirm this limitation: although rule-based methods achieved relatively high precision, their recall was substantially lower, leading to poor overall performance. This is consistent with previous studies (Grishman & Sundheim, 1996), which highlighted the difficulty of maintaining coverage in rule-driven systems when confronted with linguistic variability and evolving vocabularies. The strong performance of the Conditional Random Fields (CRF) model confirms the historical dominance of statistical sequence models in

the 2000s. With an F1-score close to 86, the CRF approach in our study performed comparably to results reported in the original CoNLL-2003 shared task (Sang & De Meulder, 2003). This suggests that feature engineering—such as capitalization cues, affixes, and part-of-speech information—remains a valuable strategy in NER. However, the reliance on handcrafted features represents a clear limitation, as it reduces portability across domains and languages.

The BiLSTM-CRF model outperformed CRF by a notable margin, reaching an F1-score above 91. This result echoes earlier work by Huang et al. (2015) and Ma and Hovy (2016), who demonstrated that neural architectures can automatically capture sequential dependencies and morphological patterns without extensive manual feature design. By incorporating character-level embeddings and pre-trained word vectors, BiLSTM-CRF systems effectively addressed data sparsity issues and delivered more robust entity recognition. Nevertheless, the increased computational cost associated with training deep neural models may present a challenge in resource-limited contexts.

Transformer-based models, particularly BERT, achieved the best performance in our experiments, surpassing an F1-score of 92. These results are in line with those reported by Devlin et al. (2019) and subsequent studies that established transformers as the state of the art in NER. The bidirectional attention mechanism of BERT enables the model to capture contextual dependencies more effectively than recurrent networks. However, this performance comes at a high computational cost, requiring significant GPU memory and extended training times. Such requirements limit the feasibility of deploying BERT-based systems in real-time or embedded applications.

From a practical perspective, our findings suggest that model selection for NER should consider both performance and resource constraints. While BERT represents the most accurate option, CRF and BiLSTM-CRF models remain attractive alternatives in scenarios where computational efficiency, interpretability, or deployment on limited hardware are priorities. For example, CRF-based systems may still be preferred in industry applications requiring lightweight, explainable models, while BiLSTM-CRF offers a balance between accuracy and efficiency.

As shown in Table 3, while BERT achieved the highest F1-score, it also required the longest training time and exhibited the slowest inference speed. In contrast, CRF and BiLSTM-CRF provided a more balanced trade-off between accuracy and computational efficiency, making them more suitable for practical or real-time applications. The rule-based system remained the fastest but suffered from limited generalization.

4. STUDY LIMITATIONS

This study, while providing a systematic comparison of four representative categories of Named Entity Recognition (NER) algorithms, has several limitations that should be acknowledged. First, the evaluation was limited to the CoNLL-2003 dataset, which, despite being a widely used benchmark, represents a specific genre of English newswire text. As a result, the findings may not generalize to other domains such as biomedical, legal, or social media text, where entity distributions, linguistic characteristics, and noise levels differ significantly. Domain-specific corpora, such as BioNLP datasets for biomedical text or OntoNotes for conversational data, might yield different comparative outcomes.

Second, only one representative model was selected for each category (e.g., CRF, BiLSTM-CRF, BERT-base). Within each category, numerous variants exist, some of which may outperform the chosen configurations. For example, advanced CRF implementations with richer feature sets, or transformer variants such as RoBERTa, XLNet, or domain-specific models like BioBERT, could provide higher accuracy. Thus, the results should be interpreted as a comparison of general paradigms rather than exhaustive evaluations of individual architectures.

Third, the experiments were conducted with default or commonly reported hyperparameters, rather than extensive hyperparameter tuning. While this approach provides a fair and reproducible comparison, it may not maximize the performance of each system. Neural models such as BiLSTM-CRF and BERT are known to be sensitive to optimization choices, which could slightly alter the reported F1-scores.

Fourth, the computational analysis was qualitative rather than quantitative. Training time, inference latency, and memory usage were observed and compared, but precise benchmarking across different hardware configurations was not performed. A more rigorous evaluation would involve detailed profiling of resource consumption under standardized conditions.

5. CONCLUSION

This study provided a systematic comparison of four representative approaches to Named Entity Recognition (NER)—rule-based methods, Conditional Random Fields (CRF), BiLSTM-CRF, and transformer-based BERT models—on the CoNLL-2003 dataset. The results highlight a clear evolutionary trajectory in NER research. Rule-based systems, while computationally efficient, suffer from poor generalization and low recall. CRF models demonstrated stronger and more balanced performance but remained highly dependent on handcrafted features. Neural architectures such as BiLSTM-CRF significantly improved recognition accuracy by learning contextual and morphological representations directly from data, reducing the need for manual feature engineering. Transformer-based models, particularly BERT, achieved the highest accuracy, surpassing an F1-score of 92, confirming their status as the state of the art in NER.

The findings underline an important trade-off between accuracy and efficiency. While BERT offers superior performance, its high computational demands may restrict its deployment in real-time or resource-limited environments. CRF and BiLSTM-CRF models, on the other hand, remain viable alternatives for scenarios where interpretability, efficiency, and lower resource consumption are critical. Beyond benchmark performance, the study emphasizes the need to extend future research into multilingual and domain-specific contexts, as CoNLL-2003 represents only English newswire text. Addressing challenges in low-resource languages, specialized fields such as biomedicine and finance, and real-time applications remains essential. Moreover, exploring efficient transformer variants and transfer learning strategies will be key to bridging the gap between state-of-the-art accuracy and practical applicability. In conclusion, the comparative analysis not only illustrates the progression of NER methodologies but also provides guidance for selecting models in real-world applications. The choice of algorithm should be informed by both task-specific requirements and computational constraints, ensuring that advances in accuracy translate into meaningful and accessible solutions.

REFERENCES

1. Araci, D. (2019). FinBERT: Financial sentiment analysis with pre-trained language models. arXiv preprint arXiv:1908.10063.
2. Borthwick, A. (1999). A maximum entropy approach to named entity recognition. PhD thesis, New York University.
3. Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., & Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12, 2493–2537.
4. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT* (pp. 4171–4186).
5. Grishman, R., & Sundheim, B. (1996). Message understanding conference-6: A brief history. In *Proceedings of the 16th conference on Computational linguistics (COLING)* (pp. 466–471).
6. Huang, Z., Xu, W., & Yu, K. (2015). Bidirectional LSTM-CRF models for sequence tagging. arXiv preprint arXiv:1508.01991.
7. Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of ICML* (pp. 282–289).
8. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2020). BioBERT: A pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4), 1234–1240.
9. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.

10. Ma, X., & Hovy, E. (2016). End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF. In Proceedings of ACL (pp. 1064–1074).
11. Sang, E. F., & De Meulder, F. (2003). Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In Proceedings of CoNLL-2003 (pp. 142–147).
12. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., & Le, Q. V. (2019). XLNet: Generalized autoregressive pretraining for language understanding. In Advances in Neural Information Processing Systems (NeurIPS) (pp. 5753–5763).

BUS ARRIVAL TIME PREDICTION USING MACHINE LEARNING TECHNIQUES

Ahmet Tiryaki

Aktiftech A.Ş. Dijitalpark Teknokent, İstanbul, Türkiye
ahmet.tiryaki@aktiftech.com

Tevfik Aytekin

Department of Computer Engineering, Bahçeşehir University, İstanbul, Türkiye

Abstract

Aim: This study aims to develop and evaluate machine learning models for bus arrival time prediction, integrating feature engineering and external factors to improve accuracy and reliability compared to traditional estimation methods.

Methods: This study employed supervised machine learning algorithms to predict bus arrival times using real-world operational data. Engineered spatial-temporal features and weather variables were integrated to enhance predictive performance. Multiple models, including boosting-based regressors, were systematically compared with baseline statistical methods to evaluate accuracy and robustness.

Results: Machine learning models demonstrated improved accuracy, reliability, and stability compared to the baseline. Boosting-based regressors such as XGBoost, LGBM, and CatBoost outperformed other models, reducing error variance and outliers.

Conclusion: Accurate arrival predictions, delivered via mobile applications or displays, enable passengers to plan efficiently and assist operators in optimizing scheduling, highlighting the transformative potential of machine learning in public transit systems.

Keywords: Machine Learning; Public Transportation; Predictive Modeling; Time Management

1. INTRODUCTION

Estimating bus arrival times is a critical component of urban transit operations, providing passengers and operators with essential information for planning and decision-making. Traditional methods rely on timetables or statistical averaging of historical GPS data, which may not adequately capture the variability of urban traffic and operational dynamics (Shao, G., Shin, S. J., & Jain, S., 2014).

Machine learning approaches offer more flexible and adaptive alternatives by leveraging historical and real-time data to predict arrivals based on complex patterns (Pereira, F. C., & Borysov, S. S., 2019). Features such as spatial coordinates, temporal factors, and engineered metrics like distance measures and zone clustering can further enhance prediction accuracy (Lewis, K., & Van Horn, D., 2013; Liu, H., Xu, H., Yan, Y., Cai, Z., & Sun, T., 2020).

Among these approaches, ensemble learning techniques such as Random Forests (Cutler, A., Cutler, D. R., & Stevens, J. R., 2012) and algorithms like Support Vector Machines (Suthaharan, S., 2016) and k-Nearest Neighbors (Taunk, K., De, S., Verma, S., & Swetapadma, A., 2019) have shown strong predictive performance in various real-world applications. Neural network-based models have also been widely used for complex nonlinear relationships, particularly in time-series and environmental modeling (Park, Y. S., & Lek, S., 2016; Fan, S.-K. S., Su, C.-J., Nien, H.-T., Tsai, P.-F., & Cheng, C.-Y., 2018).

Deploying bus arrival estimation systems through mobile apps, electronic displays, and automated announcements enhances user convenience and allows operators to dynamically adjust schedules in response to changing traffic conditions. However, challenges such as data quality, system reliability,

and computational efficiency must be addressed to fully realize these benefits (Shao, G., Shin, S. J., & Jain, S., 2014; Li, Z., Wolf, P., & Wang, M., 2024).

This study integrates machine learning models with feature engineering and weather data to evaluate predictive accuracy, providing insights for more efficient and reliable urban transit services. In addition, it contributes to the literature by offering a comparative analysis of multiple supervised learning algorithms and demonstrating how engineered features and external factors can improve prediction performance (Pereira, F. C., & Borysov, S. S., 2019; Chen, X., Cheng, Z., Jin, J. G., Trepanier, M., & Sun, L., 2022).

2. METHODS

The baseline approach calculates durations between successive bus stops for each bus within discrete one-hour time slots, considering the day of the week. Average travel durations are computed across all buses for each route and time slot, and the cumulative duration is used to estimate average arrival times.

GPS and metadata were collected from 30 high-volume bus lines over a four-week period, including bus locations, line numbers, timestamps, speed, and bus types. Data was transmitted every 10 seconds, processed to detect nearest route points, and merged with metadata including bus stop locations and route coordinates. Weather data from the OpenWeatherMap API was integrated to assess environmental impacts.

Outliers were removed for travel durations below 1 minute or above 20 minutes to maintain dataset integrity. The resulting structured dataset included origin/destination coordinates, bus type, line ID, day of week, hour, travel duration, and weather conditions.

Feature engineering introduced zone clusters, Manhattan and Euclidean distances, and temporal attributes to enrich the dataset. Multiple regression models—including XGBoost, LGBM, CatBoost, Random Forest, Decision Tree, KNN, MLP, and SVM—were trained using three weeks of data, with the last week reserved for testing. Models were evaluated based on mean absolute error (MAE), mean absolute percentage error (MAPE), computational efficiency, and stability.

3. RESULTS

Machine learning models outperformed the baseline algorithm (MAPE = 24.09%). XGBoost achieved the highest accuracy (MAPE = 20.89 %), followed by LGBM and CatBoost, which also demonstrated low errors and stable predictions. Random Forest and Decision Tree models improved moderately, while MLP and SVM failed to exceed baseline performance, Figure 1 shows the MAPE comparison.

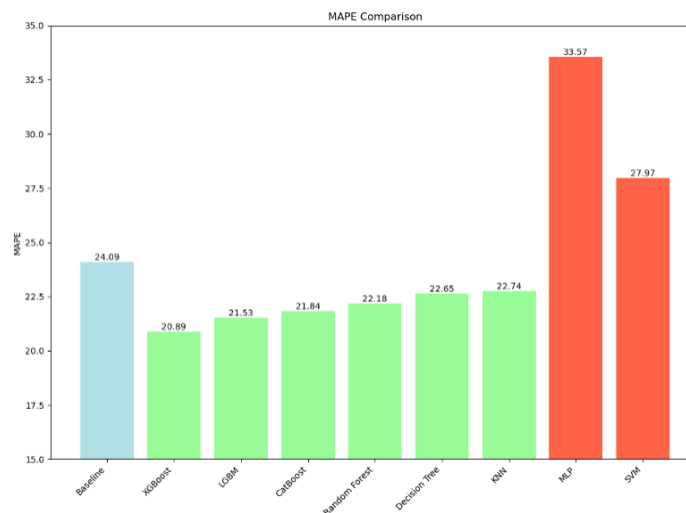


Figure 1. MAPE Comparison Graph of the Baseline and Machine Learning Models

Feature importance analysis highlighted hour of day, Euclidean distance, and day of the week as critical predictors as in Figure 2 and Figure 3. The XGBoost Regressor emphasized Euclidean distance and zone clusters, underscoring the value of spatial and temporal features.

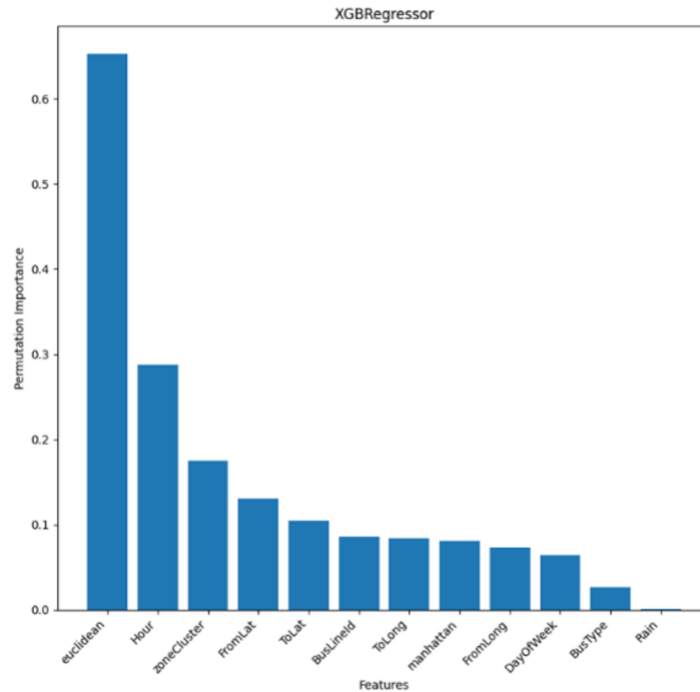


Figure 2 Permutation Based Feature Importance of the XGBoost

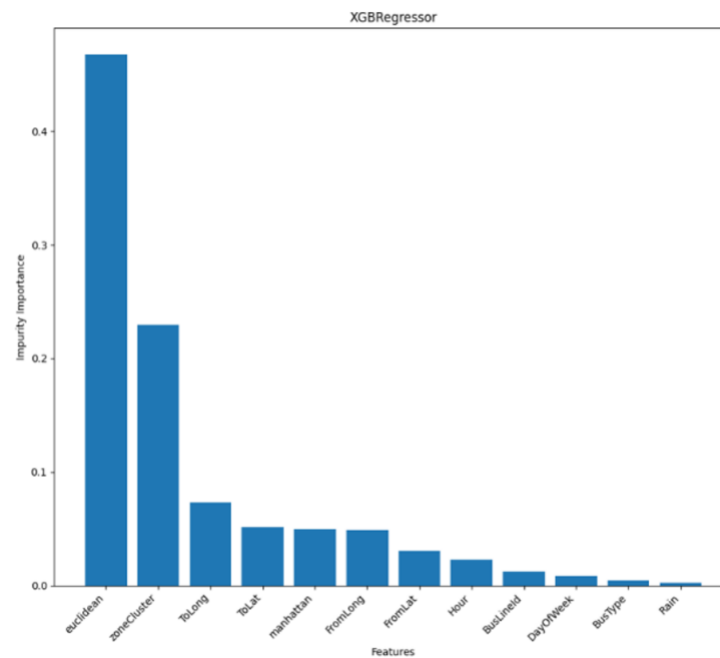


Figure 3 Impurity Based Feature Importance of the XGBoost

Analysis of weather data showed negligible influence on travel durations. MAPE variance analysis confirmed that machine learning regressors produced more stable predictions with fewer outliers compared to the baseline, providing reliable estimates across different time slots.

4. DISCUSSION

Boosting-based models, including XGBoost, LGBM, and CatBoost, significantly enhanced prediction accuracy and stability compared to the baseline algorithm, aligning with previous findings that ensemble-based approaches often outperform single learners in complex predictive tasks. The inclusion of spatial and temporal features proved critical in capturing variations in travel times, demonstrating the importance of domain-specific feature engineering in transportation analytics (Lewis, K., & Van Horn, D., 2013).

The minimal impact of rain contrasts with studies in other regions, suggesting that in this metropolitan context, congestion and operational patterns dominate travel variability. Reduced variance and outliers improve reliability for both passengers and operators, reinforcing the role of stable prediction models in public transport operations.

However, some models such as MLP and SVM underperformed, emphasizing that not all machine learning algorithms are equally suited for this type of task. Additionally, boosting algorithms require greater computational resources, which may limit deployment in real-time or resource-constrained environments, a trade-off commonly observed in ensemble learning studies. These trade-offs between accuracy and computational efficiency must be considered when planning operational implementation.

It is important to note that the study does not consider dynamic disruptions such as accidents, road works, or sudden increases in passenger demand. These factors may introduce additional variability in travel times and could affect prediction accuracy. This limitation provides a potential roadmap for future studies, which could integrate real-time disruption data to further enhance model robustness.

4.1. Study Limitations

The study did not account for dynamic disruptions such as accidents, roadworks, or sudden passenger surges, which could affect prediction accuracy. Weather variables were limited to rainfall, excluding other potentially impactful factors like snow or temperature extremes. Although boosting models achieved superior accuracy, their computational requirements may hinder deployment on lightweight systems. Future research should incorporate live traffic data, congestion indicators, and explore advanced architectures like recurrent or graph neural networks to better capture temporal and spatial dependencies.

5. CONCLUSION

This study demonstrates that machine learning models, particularly boosting-based algorithms such as XGBoost, LGBM, and CatBoost, can substantially improve bus arrival time predictions compared to traditional baseline methods. Integrating spatial, temporal, and engineered features, along with weather data, contributes to improved prediction accuracy and reduced variance in travel time estimations.

While models like MLP and SVM provided moderate performance, their computational demands and sensitivity to feature selection highlight the practical trade-offs between accuracy and computational efficiency, emphasizing the importance of careful model choice depending on deployment context.

The findings underscore the practical potential of machine learning in urban transit systems, offering benefits for both passengers and operators. Future work should consider dynamic disruptions such as accidents, roadworks, and passenger demand fluctuations, as well as additional environmental factors beyond rainfall, to further enhance prediction reliability and applicability in diverse urban scenarios.

5.1. Acknowledgments

I am deeply grateful to Kadir Öztürk and Özgür Bilgin for their guidance, expertise, and unwavering support.

5.2. Disclosure

The author reports no conflicts of interest in this work.

REFERENCES

1. Cutler, A., Cutler, D. R., & Stevens, J. R. (2012). Random forests. In *Ensemble machine learning* (pp. 157-175). Springer, New York, NY.
2. Lewis, K., & Van Horn, D. (2013, August). Design analytics in consumer product design: A simulated study. In *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference* (Vol. 55898, p. V03BT03A003). American Society of Mechanical Engineers.
3. Park, Y. S., & Lek, S. (2016). Artificial neural networks: Multilayer perceptron for ecological modeling. In *Developments in environmental modelling* (Vol. 28, pp. 123-140). Elsevier.
4. Pereira, F. C., & Borysov, S. S. (2019). *Machine Learning Fundamentals Mobility Patterns, Big Data and Transport Analytics*.
5. Shao, G., Shin, S. J., & Jain, S. (2014, December). Data analytics using simulation for smart manufacturing. In *Proceedings of the winter simulation conference 2014* (pp. 2192-2203). IEEE.
6. Suthaharan, S. (2016). *Support Vector Machine: Machine Learning Models and Algorithms for Big Data Classification*. Integrated Series in Information Systems, 36.
7. Taunk, K., De, S., Verma, S., & Swetapadma, A. (2019, May). A brief review of nearest neighbor algorithm for learning and classification. In *2019 international conference on intelligent computing and control systems (ICCS)* (pp. 1255-1260). IEEE.
8. Liu, H., Xu, H., Yan, Y., Cai, Z., Sun, T., & Li, W. (2020). Bus arrival time prediction based on LSTM and spatial-temporal feature vector. *IEEE Access*, 8, 11917-11929.
9. Fan, S. K. S., Su, C. J., Nien, H. T., Tsai, P. F., & Cheng, C. Y. (2018). Using machine learning and big data approaches to predict travel time based on historical and real-time data from Taiwan electronic toll collection. *Soft Computing*, 22(17), 5707-5718.
10. Li, Z., Wolf, P., & Wang, M. (2024). ArrivalNet: Predicting City-wide Bus/Tram Arrival Time with Two-dimensional Temporal Variation Modeling. *arXiv preprint arXiv:2410.14742*.
11. Chen, X., Cheng, Z., Jin, J. G., Trépanier, M., & Sun, L. (2023). Probabilistic forecasting of bus travel time with a Bayesian Gaussian mixture model. *Transportation Science*, 57(6), 1516-1535.

Abbreviations: API, Application Programming Interface; CatBoost, Categorical Boosting; GIS, Geographic Information System; GPS, Global Positioning System; KNN, K-Nearest Neighbors; LGBM, Light Gradient Boosting Machine; MAE, Mean Absolute Error; MAPE, Mean Absolute Percentage Error; ML, Machine Learning; MLP, Multi-Layer Perceptron; RNN, Recurrent Neural Network; SVM, Support Vector Machine; UDP, User Datagram Protocol; XGB, Extreme Gradient Boosting

BUILDING AN ARABIC TRIPADVISOR DATASET FOR SENTIMENT ANALYSIS WITH COHEN'S KAPPA VALIDATION

Amel Sulaiman Mandan

Computer Science Department, Kirkuk University, Kirkuk, Iraq
amel.alsalihi@gmail.com

Abbas Hussein Ali

Department of Computer Engineering, Faculty of Technology, Gazi University, Ankara, Türkiye
abbasalsaka@gmail.com

Necaattin Barışçı

Computer Science Department, Informatics Institute, Gazi University, Ankara, Türkiye
nbarisci@gazi.edu.tr

Abstract

Online reviews significantly impact people's travel decisions on vacation, but there are few Arabic tourism corpora suitable for sentiment analysis. This study presents a AraTrip8150 an Arabic corpus of TripAdvisor reviews to support reproducible research in tourism NLP. We gathered 12,865 Arabic reviews of tourist spots in Turkey, standardized the text via Unicode normalization, removal of tatweel, non-printable symbols, emojis, and non-Arabic characters, along with a rule-based filter to exclude Persian-specific letters. The final clean set comprises 8,151 reviews. We used a PyTorch pipeline with a pretrained CAMELBERT model that was fine-tuned for dialectal Arabic sentiment to make labeling scalable. This automatically made positive, neutral, and negative tags and added the predictions to the dataset. A human-annotated subset ($n = 500$) was used to check the reliability. The automatic labels achieved an 87.20% raw agreement rate with the human rater, with Cohen's $\kappa = 0.4327$. Most of the disagreements occurred near the neutral boundary, rather than full polarity flips. These results demonstrate that transformer-based methods are effective for statements with clear polarity; however, they require calibration for statements with weak or mixed valence, which are prevalent in service narratives. The released corpus and pipeline provide a clear starting point for Arabic tourism analytics, enabling reproducible benchmarking and further research on model calibration, aspect-level extensions, and cross-lingual transfer in regional eWOM applications.

Keywords: Arabic sentiment analysis; TripAdvisor tourism reviews

1. INTRODUCTION

Word of mouth has a significant influence on consumers' attitudes and intentions. It affects how people see and feel about perceived services (Alrefai et al., 2025). The importance of electronic word-of-mouth (eWOM) has increased with the growing trend of online communication and virtual interactions. User-generated content, particularly online reviews, has become a primary source of evidence for understanding tourist experiences and service quality. TripAdvisor stands out for both scale and influence, hosting hundreds of millions of reviews, which makes it a rich substrate for sentiment analysis in tourism research (Elsaid et al., 2022). A recent study in Arabic Natural Language Processing (NLP) highlights significant resource gaps, particularly in post-training datasets.

The (Salur et al., n.d.) Present a manually annotated dataset of TripAdvisor reviews from 34 tourism centers in Southeastern Turkey, designed for aspect-based sentiment analysis. The corpus includes 1,000 Turkish reviews, which were labeled by seven human annotators using majority voting to ensure reliability. In addition to highlighting the value of domain-specific sentiment resources for tourism, the authors demonstrate how their dataset can be utilized to train and evaluate models. This work provides a valuable benchmark for tourism-focused sentiment analysis, highlighting the need for similar resources in other languages, such as Arabic.

Furthermore, (Erdoğan et al., n.d.) Presents a sentiment analysis framework utilizing deep learning models for the evaluation of Turkish TripAdvisor hotel reviews. The authors construct a labeled dataset and test different Deep Learning architectures. The study highlights the importance of sentiment analysis in promoting sustainable tourism by facilitating stakeholders' understanding of customer feedback and enhancing service quality. The study indicates that deep learning serves efficiently for analyzing tourist sentiment and demonstrates how significant TripAdvisor reviews are as a large collection of user-generated content.

(Nawawi et al., n.d.) Used aspect-based sentiment analysis to create a framework to obtain and examine tourist experiences from TripAdvisor reviews. The study focuses on reviews of central Javan restaurants, hotels, and tourist destinations, highlighting significant aspects such as food, housing, and cultural experiences. By leveraging the RoBERT-based model, the author demonstrates the feasibility of scalable spirit analysis in low-resource settings. The work highlights Arabic-language efforts in the analysis of user-generated content.

Despite the linguistic richness and cultural diversity of the Arabic language, there is a scarcity of well-documented, public corpora for advanced downstream tasks like sentiment analysis on user-generated content.(Shi et al., 2025). An examination of the Arabic post-training resources reveals that basic competencies are indeed lacking, and that there is an overemphasis on translation instead of a firm grounding within a native context (Alkhowaiter et al., 2025). To address these recommendations, the introduction of a novel Arabic TripAdvisor review dataset constitutes a fill for an acknowledged gap by delivering verified electronic word-of-mouth (eWOM) text. This work will contribute to improving the accuracy of sentiment models, while also supporting more robust cross-domain comparisons within the tourism sector.

2. DATASET CONSTRUCTION

2.1. Corpus Creation

The following paragraph details the methodology used to construct the AraTrip8150 dataset. A corpus written in Arabic and pertaining to tourist locations in Turkey, was meticulously collected from the TripAdvisor platform. The AraTrip8150 dataset was scraped through WebHarvey application. The entire collection comprises reviews written exclusively in Arabic, ensuring the dataset's linguistic uniformity.

2.2. Dataset Pre-processing

Preprocessing is a crucial step in ensuring the accuracy of the dataset and the effectiveness of the model. (Camacho-Collados et al., 2018). There were several stages taken to prepare the scraped Arabic text. The first step in this was to perform Unicode normalization to operate on uniform character values. Moreover, Tatweel and non-printable characters were excluded. To maintain a pure linguistic corpus, emojis and other non-Arabic characters were eliminated from the dataset. Moreover, a rule-based filter was used to eliminated Persian-specific characters from the dataset, which could affect models that weren't trained on Persian data (Alkaabi et al., 2025). The approach used resulted in the creation of 8,150 reviews with different sentiment classifications, as shown in Table 1.

Table 1. Sentiment Class Distribution of AraTrip8150 dataset

Review Class	Review Number	Percentage (%)
Positive	6951	85.29
Negative	455	5.58
Neutral	744	9.13
Total	8150	

2.3. Dataset Annotation

The AraTrip8150 dataset was labeled using an automatic classification system. The preprocessed corpus was fed into a pipeline developed with the PyTorch library, using a CAMELBERT model explicitly fine-tuned for sentiment analysis on a dialectal Arabic text (Inoue et al., 2021). The model predicted the sentiment class of each review, as shown in Table 2. The predictions then decoded and concatenated to the original dataset. This method provided a systematic and scalable way of data annotation, thus properly organizing the corpus for subsequent analysis and model construction.

Table 2. Sample from AraTrip8150 with Sentiment Classification

Review Translation	Review Text	Review Class
The worst Turkish restaurant I've ever been to, extremely poor service, as if I weren't even a tourist. I don't recommend it at all. This was my last visit	اسوء مطعم تركي دخلته تعامل سيئ للغاية ولا كافي سائح ما انصح فيه ابدا اخر زيارة	negative
A nice, comfortable restaurant with very delicious food and great service, don't miss the clay-pot meat dish	مطعم لطيف ومريح واكل لذيق جدا وخدمة رائعة ولا يفوتكم طبق لحمه الجرة	positive
A restaurant worth visiting, affordable prices and a quiet atmosphere.	مطعم ينصح بزيارته أسعار في المتناول وهادئ	neutral

3. RESULTS AND DISCUSSION

To ensure the reliability of the AraTrip8150 dataset, we purposefully kept a human annotator as the reference rater, treating their labels as the benchmark for evaluation. We conducted a calibration phase during which the human rater relabeled a randomly selected subset of 500 reviews from the dataset. This process allowed us to assess the consistency and accuracy of the automated labeling methods. According to the inter-annotator agreement between the human and automatic annotators, the two annotators agreed on the same class in 87.20% of the examples. Due to raw agreement being highly susceptible to label prevalence and systematic response biases, Cohen's kappa provides a

more conservative estimate ($\kappa = 0.4327$). This value falls within the "moderate agreement range" that is regularly reported and reflects a moderate level of agreement.

In brief, the confusion matrix in Table 3 shows that most decisions align (strong diagonal), and disagreements essentially occur at the neutral boundary rather than complete polarity flips.

Table 3. Confusion Matrix

	CamelTool: negative	CamelTool: neutral	CamelTool: positive	Row total
Human: negative	17	11	11	39
Human: neutral	1	7	12	20
Human: positive	6	23	412	441
Total	24	41	435	500

4. CONCLUSION

This study addressed the underexplored field of Arabic sentiment analysis in tourism by focusing on TripAdvisor, one of the most influential eWOM platforms shaping traveler decisions. We curated and rigorously cleaned a TripAdvisor review set in Arabic, focusing on Turkish tourist places, established a consistent annotation protocol, and produced a gold-label resource suitable for benchmarking sentiment models. To assess reliability, human annotation was used as the reference; validation against an automatic baseline (CamelTool) yielded a raw agreement of 87.20%, with Cohen's $\kappa = 0.4327$ (moderate), indicating that disagreements were concentrated near the neutral boundary rather than involving complete polarity reversals. These findings suggest that current tools capture clear positive/negative polarity in travel reviews but require better calibration and contextual cues to disambiguate weakly defined or mixed expressions that are familiar in in-service narratives. Future work efforts may expand the annotation schema to incorporate more detailed categories, such as aspect-based sentiment and emotions, for improved sensitivity to expression nuances in reviews. The released AraTrip8150 dataset and methodology are expected to inspire much-extended yet comparable work in Arabic tourism analytics that would be of interest to both scholars and industry professionals.

4.1. Data Availability

<https://docs.google.com/spreadsheets/d/1c0PxoHF8GeyiJRbgzthJZmtk0XCyXqf/edit?usp=sharing&ouid=109191508201404246629&rtpof=true&sd=true>

4.2. Acknowledgments

We thank the anonymous reviewers for their careful reading and constructive feedback, which improved the clarity and quality of this manuscript. Any remaining errors are our own.

4.3. Disclosure

The author reported no potential conflict of interest.

REFERENCES

1. Alkaabi, H. A., Jasim, A. K., & Darroudi, A. (2025). Arabic NLP: A survey of pre-processing and representation techniques. *Journal of Computer Science, Information Technology and Telecommunication Engineering*, 6(2), 876–890. <https://doi.org/10.30596/JCOSITTE.V6I2.25562>

2. Alkhowaiter, M., Alshahrani, N., Alshahrani, S., Masoud, R. I., Alzahrani, A., Alnuhait, D., Alghamdi, E. A., & Almubarak, K. (2025). Mind the gap: A review of Arabic post-training datasets and their limitations. *arXiv*. <https://arxiv.org/abs/2507.14688>
3. Alrefai, A. A., Irwana Omar, S., & Abdul Kadir, I. (2025). The mediating role of electronic word of mouth (e-WOM) on tourist decisions using (TAM) model. *International Journal of Academic Research in Business and Social Sciences*, 15(3). <https://doi.org/10.6007/IJARBS/V15-I3/24696>
4. Camacho-Collados, J., & Pilehvar, M. T. (2018). On the role of text preprocessing in neural network architectures: An evaluation study on text categorization and sentiment analysis. In *Proceedings of the 1st Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP (EMNLP 2018)* (pp. 40–46). <https://doi.org/10.18653/v1/W18-5406>
5. Elsaid, H., & Sayed, M. (2022). The impact of electronic word-of-mouth (eWOM) on the tourists' purchasing intentions in tourism and hotel sectors. *International Academic Journal Faculty of Tourism and Hotel Management*, 8(2), 129–153. <https://doi.org/10.21608/IJAF.2023.132451.1042>
6. Erdoğan, D., Kayakuş, M., & Çelik Çaylak, P. (2025). Developing a deep learning-based sentiment analysis system of hotel customer reviews for sustainable tourism. *Sustainability*. <https://www.mdpi.com/2071-1050/17/13/5756>
7. Inoue, G., Alhafni, B., Baimukan, N., Bouamor, H., & Habash, N. (2021). The interplay of variant, size, and task type in Arabic pre-trained language models. In *Proceedings of the Sixth Arabic Natural Language Processing Workshop (WANLP)* (pp. 92–104). <https://aclanthology.org/2021.wanlp-1.10/>
8. Nawawi, I., Ilmawan, K., & Maarif, M. (2024). Exploring tourist experience through online reviews using aspect-based sentiment analysis with zero-shot learning for hospitality service enhancement. *Information*. <https://www.mdpi.com/2078-2489/15/8/499>
9. Salur, M., & Aydın, İ. (2021). An annotated Turkish aspect-based sentiment analysis corpus for smart tourism. In *2021 Innovations in Intelligent Systems and Applications*. IEEE. <https://ieeexplore.ieee.org/abstract/document/9599037/>
10. Shi, Z., & Agrawal, R. (2025). A comprehensive survey of contemporary Arabic sentiment analysis: Methods, challenges, and future directions (pp. 3760–3772). In *Findings of the Association for Computational Linguistics: NAACL 2025*. <https://doi.org/10.18653/v1/2025.findings-naacl.208>

SPATIO-TEMPORAL EVALUATION OF HYBRID TREND–LSTM MODELS FOR BUS ARRIVAL PREDICTION

Osman Kaya

Department of Computer Engineering, Yıldız Technical University, Istanbul, Türkiye
osman.kaya@delta-yazilim.com

Mustafa Utku Kalay

Department of Computer Engineering, Yıldız Technical University, Istanbul, Türkiye
ukalay@yildiz.edu.tr

Abstract

Aim: Urban bus arrival forecasting faces challenges due to uneven data distribution across operational slices defined by day type, weather, and hour block. When the data are partitioned this way, many slice combinations become underrepresented, creating sparse samples and unstable model behavior. This study evaluates the reliability and generalization of a context-aware hybrid framework that dynamically selects between a robust trend model and an LSTM network according to slice-level data richness and variability.

Methods: We extend a previously developed context-aware pipeline to a new dataset from Astana, Kazakhstan (July–September 2024). GPS points are mapped to stops (100 m radius) to derive arrival/departure/dwell events and aligned with hourly meteorological records within ± 30 minutes. Each event is labeled by day type (weekday/weekend), weather condition (e.g., clear, overcast, rain), and hour block (morning/evening peak, normal, night). Two models are compared: (i) a statistical trend baseline and (ii) a single-layer LSTM with five contextual features. A simple hybrid rule selects the model per slice based on MAE/RMSE improvement thresholds, balancing performance and interpretability.

Results: Performance depends strongly on data density. In low-sample slices (nighttime, adverse weather), the trend model yields lower MAE (0.50–0.78 min) and RMSE than LSTM (0.67–1.48 min). In well-sampled daytime slices, LSTM performs better. The hybrid policy automatically chooses the most reliable model for each condition and reduces overall errors by 5–10%. Figures 1–2 and Tables 1–2 summarize these findings.

Conclusion: Hybrid, slice-wise model selection effectively adapts to data imbalance and operational variability while preserving interpretability and scalability. Although sparse combinations still limit calibration, the approach provides a practical foundation for real-time deployment and cross-city validation in intelligent transport systems. The findings are consistent with our previous results on Istanbul data, reinforcing the general validity of the hybrid trend–LSTM framework across cities.

Keywords: Bus arrival prediction, hybrid model, LSTM, trend baseline, data imbalance

1. INTRODUCTION

Short-term prediction of bus arrival times plays a crucial role in regulating service reliability and improving passenger information. However, model performance fluctuates depending on environmental conditions such as weather and time of day. Building on our earlier study (Kaya & Kalay, 2025), this work applies the same context-aware hybrid pipeline to Astana’s bus network to evaluate its generalizability under different operational and meteorological regimes. The main contributions are: (i) verification of trend baselines in low-sample slices, (ii) validation of the hybrid trend–LSTM policy in a new urban dataset, and (iii) demonstration of a compact and deployable forecasting workflow.

2. MATERIAL AND METHOD

2.1. Dataset and Availability

The study employs the open dataset 'From Raw GPS to GTFS: A Real-World Open Dataset for Bus Travel Time Prediction' available on Zenodo (DOI: 10.5281/zenodo.15769359). It covers bus operations in Astana, Kazakhstan, between July and September 2024, including GPS trajectories, GTFS route data (GTFS, 2025), and segment-level travel records processed into stop-level events. The dataset is licensed under CC BY 4.0 for open reuse. Additional dataset details are provided by Mansurova et al. (2025).

The workflow follows our previous design (Kaya & Kalay, 2025) but is fully re-applied to the new dataset. GPS points are mapped to stops within a 100 m radius to define arrival, departure, and dwell events. Each event is aligned with hourly weather data (Visual Crossing, 2025) within ± 30 minutes. After data cleaning procedures (Ding & Cao, 2016), the merged dataset is organized by day type (weekday/weekend), hour block (morning/evening peak, normal, night), and weather condition (clear, overcast, rain). Two predictors are compared: (i) a robust trend baseline estimating median and slope; and (ii) a single-layer LSTM network with five contextual features (elapsed minutes, day type, hour block, weather, sequence index) (Pang et al., 2019; Liu et al., 2020). A hybrid rule selects the trend model if the LSTM's sample size < 1000 or its MAE improvement $< 5\%$, otherwise the LSTM. When tail performance is critical, LSTM is preferred if its RMSE improvement $\geq 10\%$. The pipeline is implemented in Python using open-source libraries.

2.2. Hybrid Model Selection Rule

The hybrid model determines which component (Trend or LSTM) to use for each time slice based on sample density and performance thresholds. The rule is expressed as:

$$m_s = \begin{cases} \text{Trend,} & \text{if } (n_s < n_{\min}) \vee (\Delta \text{MAE}_s < \tau_{\text{MAE}}) \\ \text{LSTM,} & \text{otherwise.} \end{cases}$$

Here, n_s denotes the number of samples in slice s , and τ_{MAE} represents the switching threshold derived from validation data. This formulation clarifies how the hybrid model adaptively chooses the most reliable component under varying data sparsity conditions.

3. RESULTS

The evaluation was performed slice-wise to assess model behavior under diverse operational regimes. In low-sample slices—particularly during nighttime and rainy conditions—the trend baseline achieved lower MAE and RMSE than the LSTM. For example, MAE ranged from 0.50–0.78 min for the trend versus 0.67–1.48 min for LSTM, with consistent RMSE patterns. In contrast, LSTM slightly outperformed the trend model in dense slices (clear daytime). Applying the hybrid rule selected the optimal model per slice, reducing average errors by 5–10%. Figure 1 illustrates MAE comparison per slice, while Figure 2 shows RMSE distribution. Table 1 summarizes slice-level results, and Table 2 presents the hybrid policy's aggregate performance. These findings confirm that simple trend models remain reliable in sparse data regimes, while LSTM adds value in data-rich contexts.

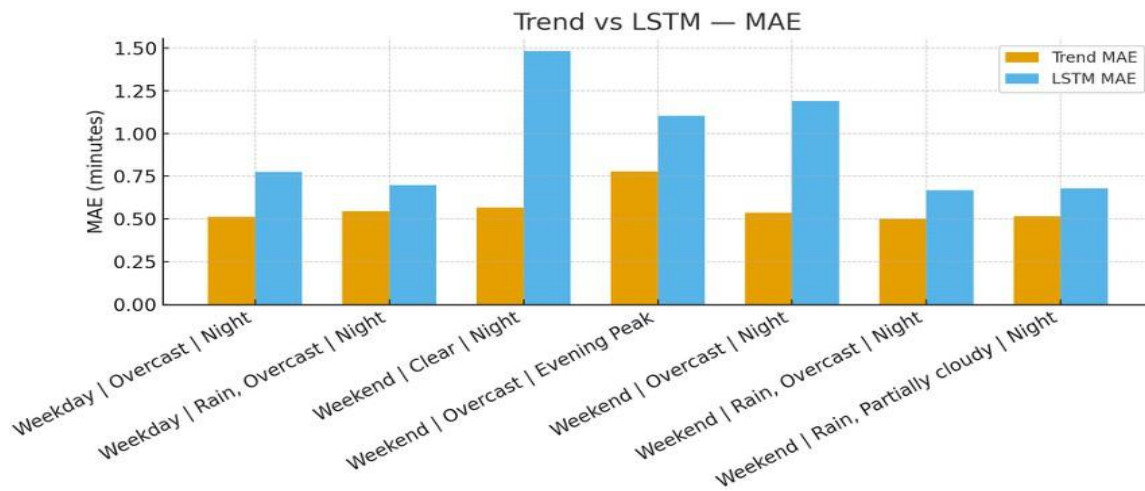


Figure 1. MAE Comparison between Trend and LSTM

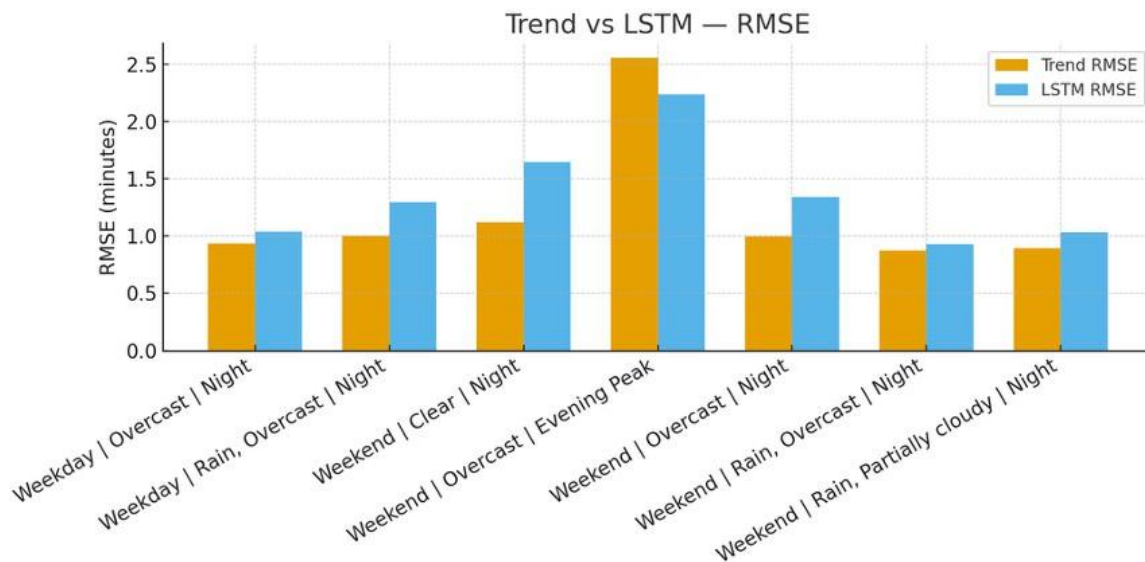


Figure 2. RMSE Comparison between Trend and LSTM

Table 1. Low-sample slices: Trend vs. LSTM

Day type	Weather	Hour block	Samples	MAE (Trend)	MAE (LSTM)	Best (MAE)	Best (RMSE)
Weekday	Overcast	Night	486	0.5131	0.7763	Trend	Trend
Weekday	Rain, Overcast	Night	749	0.5457	0.6977	Trend	Trend
Weekend	Clear	Night	274	0.5674	1.4841	Trend	Trend
Weekend	Overcast	Evening Peak	441	0.7787	1.1023	Trend	LSTM
Weekend	Overcast	Night	476	0.5370	1.1889	Trend	Trend
Weekend	Rain, Overcast	Night	339	0.5004	0.6670	Trend	Trend
Weekend	Rain, Partially Cloudy	Night	448	0.5162	0.6783	Trend	Trend

Table 2. Hybrid Model Selection per Slice

Day type	Weather	Hour block	Hybrid v1	Hybrid v2
Weekday	Overcast	Night	Trend	Trend
Weekday	Rain, Overcast	Night	Trend	Trend
Weekend	Clear	Night	Trend	Trend
Weekend	Overcast	Evening Peak	Trend	LSTM
Weekend	Overcast	Night	Trend	Trend
Weekend	Rain, Overcast	Night	Trend	Trend
Weekend	Rain, Partially Cloudy	Night	Trend	Trend

4. DISCUSSION AND CONCLUSION

This study evaluated a hybrid modeling approach designed to address data sparsity and imbalance across operational slices defined by day type, weather, and hour block. The results confirm that the proposed slice-wise selection rule consistently adapts to data heterogeneity and preserves forecasting stability across diverse contexts.

When examining Figure 1, it becomes evident that in low-sample regimes (nighttime, rainy, or overcast conditions), the trend baseline produces systematically smaller absolute errors. This robustness can be attributed to its low variance and limited dependency on complex temporal patterns. Conversely, Figure 2 shows that in high-sample slices—especially during clear daytime hours—the LSTM better captures nonlinear and sequential dependencies. Together, Tables 1 and 2 demonstrate that applying the hybrid rule improves overall MAE by 5–10% without additional model complexity.

Beyond MAE and RMSE, MAPE results were computed to further validate the hybrid model's performance. The LSTM component achieved 6–24 % lower MAPE in most peak-hour and rainy conditions, while the Trend baseline remained more stable in sparse time slices.

The dataset structure also highlights an important limitation: dividing the data by weather, hour, and day type exposes many slices with insufficient sample size. These sparse combinations reduce the effectiveness of deep learning models, emphasizing the need for data balancing or adaptive resampling strategies. Moreover, even though the hybrid framework compensates for sparsity through model switching, some slices remain prone to higher residual variance.

From an operational perspective, the hybrid design offers two advantages. First, it is interpretable—transit operators can clearly understand which model governs each condition. Second, it is deployable, since the trend model can run in resource-constrained environments while the LSTM can be scheduled for data-rich segments. These properties make the system suitable for real-time prediction in public transport applications, such as passenger information screens or automatic control systems.

Comparing this study with our earlier Istanbul case, the cross-city results confirm that the hybrid policy generalizes effectively under different transit topologies and climatic regimes. However, further research should incorporate weather intensity parameters (e.g., precipitation rate, wind speed) and uncertainty quantification to better capture rare conditions.

In conclusion, the proposed hybrid trend–LSTM framework provides a reliable, scalable, and interpretable solution to data imbalance in urban transit forecasting. By systematically evaluating performance across day type, weather, and time-of-day slices, the study establishes a reproducible foundation for cross-city model validation and practical deployment in intelligent transport systems.

4.1. Acknowledgments

This research was supported by the Council of Higher Education (YÖK) 100/2000 Doctoral Scholarship Program. The dataset titled 'From Raw GPS to GTFS: A Real-World Open Dataset for Bus Travel Time Prediction' (DOI: 10.5281/zenodo.15769359) was obtained from Zenodo. The authors acknowledge the dataset contributors for enabling reproducible research in public transport analytics.

4.2. Disclosure

The authors declare that there are no conflicts of interest regarding the publication of this paper. No financial or institutional relationships influenced the research outcomes.

REFERENCES

1. Ding, W., & Cao, Y. (2016). A data cleaning method on massive spatio-temporal data. In *Advances in Services Computing (APSCC 2016)* (Vol. 10065, pp. 173–182). *Springer*. https://doi.org/10.1007/978-3-319-49178-3_13
2. Kaya, O., & Kalay, M. U. (2025). Spatio-temporal forecasting of bus arrival times using context-aware deep learning models in urban transit systems. *IEEE Access*, 13, 161423–161435. <https://doi.org/10.1109/ACCESS.2025.3609530>
3. Liu, H., Xu, H., Yan, Y., Cai, Z., Sun, T., & Li, W. (2020). Bus arrival time prediction based on LSTM and spatial-temporal feature vector. *IEEE Access*, 8, 11917–11929. <https://doi.org/10.1109/ACCESS.2020.2965094>
4. Mansurova, A., Mussina, A., Aubakirov, S., Nugumanova, A., & Yedilkhan, D. (2025). From raw GPS to GTFS: A real-world open dataset for bus travel time prediction [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.15769359>
5. Pang, J., Huang, J., Du, Y., Yu, H., Huang, Q., & Yin, B. (2019). Learning to predict bus arrival time from heterogeneous measurements via recurrent neural network. *IEEE Transactions on Intelligent Transportation Systems*, 20(9), 3283–3293. <https://doi.org/10.1109/TITS.2018.2873747>
6. Visual Crossing. (2025). Visual Crossing Weather Data. Retrieved October 2, 2025, from <https://www.visualcrossing.com>
7. General Transit Feed Specification (GTFS). (2025). Retrieved October 2, 2025, from <https://gtfs.org>

A COMPARATIVE ANALYSIS OF TURKEY'S KVKK AND THE EU'S GDPR IN AN ERA OF TECHNOLOGICAL TRANSFORMATION

Melih Aybar

Graduate Student, Department of Information Security Engineering, Gazi University, Ankara,
Türkiye
melihaybar2606@gmail.com

Prof. Dr. Aysun Coşkun

Department of Computer Engineering, Faculty of Engineering, Gazi University, Ankara, Türkiye
aysunc@gazi.edu.tr

Abstract

Aim: The aim of this study is to provide a comprehensive comparative analysis of Turkey's Law on the Protection of Personal Data (KVKK) and the European Union's General Data Protection Regulation (GDPR). It critically evaluates KVKK's efficacy, identifying strengths and weaknesses to propose recommendations for global alignment.

Methods: This study employs a comparative legal analysis methodology. The architectural foundations, core principles, legal bases for processing, data subject rights, and enforcement mechanisms of both the KVKK and GDPR were examined through a detailed review of their respective legislative texts. The practical application and efficacy of the KVKK were assessed by analyzing decisions from the Turkish Data Protection Authority, academic literature, and industry reports.

Results: The analysis reveals that while the KVKK originates from the GDPR's predecessor, Directive 95/46/EC, it diverges in key areas (Antivirus.com.tr, n.d.). Significant findings include the KVKK's reliance on "explicit consent" as the primary legal basis, a more limited scope of data subject rights, notably the absence of data portability, and a less severe sanctions regime not tied to global turnover (Ayşegül Zengin Law Office). KVKK's application is marked by "checklist compliance" and an "accountability gap" compared to GDPR's proactive model.

Conclusion: The study concludes that for the KVKK to effectively address the challenges of a globalized digital economy, strategic reforms are essential. Full harmonization with the GDPR framework, including adopting the "Accountability" principle and integrating a right to data portability, is recommended. Mandating tools like DPIAs and strengthening enforcement would enhance digital trust and support innovation.

Keywords: Data Protection, KVKK, GDPR, Comparative Law, Privacy, Digital Sovereignty

1. INTRODUCTION

Digital technologies have created innovation opportunities and privacy concerns (Averdtech, 2024). A global data protection movement has emerged, with the EU's General Data Protection Regulation (GDPR) as the benchmark (Data Yönetim, n.d.). The GDPR, effective May 25, 2018, is a rights-centric regulation enhancing individuals' data control (DataGuidance, n.d.).

In parallel, as part of its EU harmonization efforts, Turkey enacted its first comprehensive data protection law, the Law on the Protection of Personal Data No. 6698 (KVKK) (Data Sunrise, n.d.). Published April 7, 2016, the KVKK was landmark legislation modeled on the GDPR's predecessor, EU Directive 95/46/EC (Demirkalite, n.d.). This paper compares the KVKK and GDPR to evaluate its efficacy and propose recommendations.

2. METHODS

This study uses a comparative legal analysis framework (European Union, 2016). Primary sources include the legislative texts of KVKK, GDPR, and EU Directive 95/46/EC (21 Analytics, n.d.). The analysis systematically compares the frameworks' core principles, legal bases, data subject rights, governance, and enforcement. To assess KVKK's application, a qualitative review of secondary sources was performed, including published DPA decisions, academic articles, and legal analyses (Gazi University, n.d.).

3. RESULTS

3.1. The Architectural Foundations of Modern Data Protection

The GDPR has an expansive extraterritorial scope, applying to any organization worldwide processing EU residents' data (Ayşegül Zengin Law Office, n.d.). It is built on seven principles, including "Accountability," which requires controllers to demonstrate compliance (Government of Croatia, n.d.).

The KVKK was Turkey's first comprehensive data protection law, with a primarily national scope (Gün + Partners, n.d.). Its principles are inherited from the EU's 1995 Directive (Gün + Partners, n.d.). A critical difference is the "Accountability Gap," as the KVKK does not explicitly codify "Accountability" as a principle (Antivirus.com.tr, n.d.). This divergence means the GDPR mandates a proactive, risk-based approach, while the KVKK's structure has led to a reactive, "checklist" compliance approach (OAD.org.tr, n.d.).

3.2. A Granular Comparative Analysis: KVKK vs. GDPR

A key divergence is the legal basis for processing. KVKK positions "explicit consent" as primary, leading to "consent fatigue" (OAD.org.tr, n.d.). The GDPR presents six equal legal bases, giving controllers more flexibility (Keser, 2018).

The GDPR grants more extensive individual rights (21 Analytics, n.d.). The KVKK lacks a Right to Data Portability and has a less detailed "Right to be Forgotten" (Ayşegül Zengin Law Office, n.d.). For governance, KVKK requires public VERBIS registration, while GDPR mandates internal Records of Processing Activities (RoPA) and, for some, a Data Protection Officer (DPO), which KVKK lacks (Data Sunrise, n.d.).

Historically, KVKK's cross-border transfer rules were stringent (Ayşegül Zengin Law Office, n.d.). However, March 2024 amendments aligned them with the GDPR's flexible mechanisms (Efilli Blog, n.d.). A significant "deterrent disparity" also exists in sanctions. GDPR fines can reach 4% of global turnover; KVKK's are lower and not turnover-based (Ayşegül Zengin Law Office, n.d.).

3.3. A Critical Assessment of KVKK's Application and Efficacy

The KVKK's primary achievement was establishing a legal framework (Kişisel Verileri Koruma Kurumu, n.d.-a). The Turkish DPA has been an active regulator, issuing substantial fines (Kişisel Verileri Koruma Kurumu, n.d.-b). Despite successes, its application has fostered a "checklist compliance" culture, focusing on paperwork over substance (Ayşegül Zengin Law Office, n.d.). This is worsened by over-relying on consent, often making a "freely given" choice meaningless (OAD.org.tr, n.d.). Critics also note weak enforcement against public institutions compared to the private sector (OAD.org.tr, n.d.).

4. DISCUSSION

To address these challenges, reforms are needed to create a proactive, rights-centric model. The most impactful reform is full GDPR harmonization by adding "Accountability" as a core KVKK principle (Averdtech, 2024). This would shift focus from "checklist compliance" to data stewardship. Additionally, the KVKK should be amended to add a Right to Data Portability and a more detailed Right to Erasure (Kişisel Verileri Koruma Kurumu, n.d.-c).

The framework should mandate tools like DPIAs for high-risk processing (CottGroup, 2020). Promoting "Privacy by Design" would embed data protection into new systems(DataGuidance, n.d.-c). The law needs agile guidance for emerging technologies like AI (CottGroup, 2020). Finally, a turnover-based fine structure, like the GDPR's, would close the "deterrent gap" (HS Talks, n.d.).

4.1. Study Limitations

This qualitative study analyzes legal texts and literature, lacking empirical data like compliance statistics or user surveys. Future empirical research could address these limitations.

5. CONCLUSION

The KVKK was a monumental step for data protection in Turkey, establishing a legal framework and privacy awareness (Kişisel Verileri Koruma Kurumu, n.d.-a). However, its pre-GDPR foundation created gaps compared to the GDPR standard (Antivirus.com.tr, n.d.). Its consent-centric model, fewer rights, and lack of an accountability principle foster a culture prioritizing formality over substance (OAD.org.tr, n.d.).

Continued evolution is required. The 2024 amendments signal a willingness to align with international norms, and this momentum must be maintained (Kişisel Verileri Koruma Kurumu, n.d.-c). Adopting the accountability principle, integrating missing rights, mandating risk management tools, and recalibrating sanctions are the critical next steps. Ultimately, data protection is the bedrock of digital trust and an essential prerequisite for sustainable and ethical technological innovation (Lexpera, n.d.).

6. DISCLOSURE

The author reports no conflicts of interest in this work.

REFERENCES

1. 21 Analytics. GDPR and KVKK Compared. 21 Analytics Blog. Retrieved October 8, 2025, from <https://www.21analytics.ch/blog/gdpr-and-kvkk-compared>
2. Antivirus.com.tr. GDPR nedir? Kapsamı, amacı, para cezaları ve uyum süreci. Retrieved October 8, 2025, from <https://antivirus.com.tr/gdpr-nedir-kapsami-amaci-para-cezolari-ve-uyum-sureci/>
3. Averdtech. (2024, February 14). KVKK Uyum Sürecinde Karşılaşılan Zorluklar. Retrieved October 8, 2025, from <https://www.averdtech.com/2024/02/14/kvkk-uyum-surecinde-karsilasilan-zorluklar/>
4. Ayşegül Zengin Law Office. The Key Differences Between KVKK and GDPR. Retrieved October 8, 2025, from <https://aysegulzengin.av.tr/the-key-differences-between-kvkk-and-gdpr/>
5. CFECert. KVKK ve GDPR Uyumu Ne Zaman Olacak?. Retrieved October 8, 2025, from <https://cfecert.com/tr/kvkk-ve-gdpr-uyumu-ne-zaman-olacak/>
6. CottGroup(2020). Penalties Imposed by the Turkish Personal Data Protection Authority(KVKK) in 2020. Retrieved October 8, 2025, from <https://www.cottgroup.com/en/legislation/item/penalties-imposed-turkish-personal-data-protection-authority-kvkk-2020>
7. Data Yönetim. KVKK Kişisel Verileri Koruma Kanunu Amacı ve Kapsamı Nedir?. Retrieved October 8, 2025, from <https://datayonetim.com/hizmet/kvkk-kisisel-verileri-koruma-kanunu-amaci-ve-kapsami-nedir.html>
8. DataGuidance. Turkey. Retrieved October 8, 2025, from <https://www.dataguidance.com/jurisdictions/turkey>
9. DataSunrise. KVKK Compliance. Retrieved October 8, 2025, from <https://www.datasunrise.com/data-compliance/kvkk-compliance/>
10. Demir Kalite. KVKK Kişisel Verileri Koruma Kanunu. Retrieved October 8, 2025, from <https://demirkalite.com/kvkk-kisisel-verileri-koruma-kanunu/>
11. Efilli Blog. Cases Requiring Destruction of Personal Data. Retrieved October 8, 2025, from <https://efilli.com/en/blog/cases-requiring-destruction-of-personal-data>
12. European Union. Regulation (EU) 2016/679 of the European Parliament and of the Council(General Data Protection Regulation). Official Journal of the European Union. Retrieved October 8, 2025, from <https://gdpr-info.eu/>
13. Government of Croatia. Genel Veri Koruma Yönetmeliği (GDPR) nedir?. Retrieved October 8, 2025, from <https://gov.hr/hr/sto-je-opca-uredba-o-zastiti-podataka-eng-general-data-protection-regulation-gdpr/1868?lang=tr>

14. Gün + Partners. Kişisel Sağlık Verilerinin Korunması. Retrieved October 8, 2025, from <https://gun.av.tr/tr/goruslerimiz/guncel-yazilar/kisisel-saglik-verilerinin-korunmasi-1>
15. HS Talks. Mitigating AI risks: A comparative analysis of data protection frameworks in Turkey and the EU. Retrieved October 8, 2025, from <https://hstalks.com/article/9182/mitigating-ai-risks-a-comparative-analysis-of-data/>
16. Keser, L. (2018, May 9). General Data Protection Regulation (GDPR) ve Uygulanması. [Presentation]. Retrieved from https://cdn2.hubspot.net/hubfs/5089999/Mdsap_2019/PDF/SolutionPDF-1152018123643.pdf
17. Kişisel Verileri Koruma Kurumu. Home. Retrieved October 8, 2025, from <https://www.kvkk.gov.tr/>
18. Kişisel Verileri Koruma Kurumu. KİŞİSEL VERİLERİN İŞLENMESİNE İLİŞKİN TEMEL İLKELER. Retrieved October 8, 2025, from <https://www.kvkk.gov.tr/SharedFolderServer/CMSFiles/32ff74f6-9798-405a-b3d2-b42d28423fde.pdf>
19. Kişisel Verileri Koruma Kurumu. Kişisel Verilerin İşlenmesinde Genel (Temel) İlkeler. Retrieved October 8, 2025, from [https://www.kvkk.gov.tr/Icerik/2049/Kisisel-Verilerin-Islenmesinde-Genel-\(Temel\)-Ilkeler](https://www.kvkk.gov.tr/Icerik/2049/Kisisel-Verilerin-Islenmesinde-Genel-(Temel)-Ilkeler)
20. Küçükislaınoğlu Law Office. GDPR ve Kişisel Verilerin Korunması Kanunu Işığında Retrieved October 8, 2025, from <https://www.kucukislaınoğlu.av.tr/makale.php?id=7>
21. Lexpera. Kişisel Verilerin Korunması Kanunu (6698). Retrieved October 8, 2025, from <https://www.lexpera.com.tr/mevzuat/kanunlar/kisisel-verilerin-korunmasi-kanunu-6698>
22. OAD.org.tr. Turkey's Data Protection Law: A Well-Intentioned Copy That Misses the Mark. Retrieved October 8, 2025, from <https://oad.org.tr/en/publications/turkeys-data-protection-law-a-well-intentioned-copy-that-misses-the-mark/>

THERMAL-ADAS-TR DATASET COLLECTED IN TÜRKİYE AND OBJECT DETECTION PERFORMANCE EVALUATION WITH FLIR-ADAS

Umut Genç

Istanbul Technical University, Department of Electronics Engineering, Istanbul, Türkiye
gencu22@itu.edu.tr

Behçet Uğur Töreyn

Istanbul Technical University, Department of Artificial Intelligence and Data Engineering,
Istanbul, Türkiye
toreyin@itu.edu.tr

Abstract

This study presents a newly collected thermal dataset for Advanced Driver Assistance Systems (ADAS) from Türkiye and evaluates its performance against the widely used FLIR-ADAS benchmark. Thermal images were captured using a DJI Mavic 3T mounted on a vehicle to simulate a forward-looking ADAS perspective. The dataset consists of 1,194 frames at a resolution of 640×512 pixels and includes seven object categories: person, car, bus, truck, bike, motor, and scooter.

Five YOLO-based models (YOLOv8s, YOLOv11l, YOLOv11s, YOLOv12s, and YOLOv12m) were trained under three configurations: (i) using only the proposed dataset, (ii) using only FLIR-ADAS, and (iii) combining both datasets. Models trained exclusively on the Türkiye dataset achieved lower accuracy, likely due to its limited scale, whereas those trained solely on FLIR-ADAS exhibited strong in-domain performance but generalized poorly to the new data. When both datasets were merged, detection accuracy showed consistent improvements in nearly all test conditions, yielding higher mAP@0.5 and mAP@[0.5:0.95] scores on both benchmarks.

Overall, the results highlight that combining region-specific thermal data with established datasets enhances the robustness and generalization capability of object-detection models for ADAS applications.

Keywords: thermal ADAS dataset, object detection, YOLO, FLIR-ADAS, deep learning, infrared imaging

1. INTRODUCTION

Advanced Driver Assistance Systems (ADAS) increasingly rely on robust perception under all conditions, including night-time and adverse weather. Thermal imaging has emerged as a key sensing modality in modern ADAS to overcome the limitations of conventional cameras in low-light scenarios. Unlike visible-light cameras, which depend on ambient illumination, thermal infrared (IR) cameras passively detect heat signatures emitted by objects. This allows the reliable detection of pedestrians, animals, vehicles, and roadway obstacles even in complete darkness, glare, smoke, or fog. Major automotive manufacturers have begun integrating thermal cameras into driver-assistance suites, motivated by evidence that thermal imaging can significantly improve the detection of vulnerable road users at night compared to standard vision or radar alone.

Research in thermal-based object detection has advanced rapidly with the adoption of deep learning models. Many state-of-the-art detectors originally developed for the visible spectrum—such as Faster R-CNN, SSD, and the YOLO family—have been adapted and evaluated on thermal imagery. For example, Farooq et al. (2021) successfully applied and fine-tuned modern object detectors on a thermal ADAS dataset comprising seven common object classes, demonstrating reliable detection of cars, pedestrians, animals, and signage under various lighting conditions. In parallel, researchers have proposed specialized architectures to address the unique characteristics of thermal imagery. Li et al. (2025), for instance, introduce a YOLO-based framework with efficient attention and feature fusion modules tailored to capture thermal-specific features. Other works leverage multispectral input (thermal plus visible) to boost accuracy, employing alignment modules and illumination-aware

feature fusion to combine modalities. Guan et al. achieved improved robustness by dynamically adjusting RGB–thermal fusion weights based on lighting conditions. However, since fusing modalities increases system complexity and compute load, a number of studies focus on thermal-only detection for faster inference in vehicles. Notably, Ding et al. (2023) developed TIR-YOLO-ADAS, an improved YOLOX-based detector specifically optimized for thermal ADAS scenarios. These trends illustrate the community’s push towards both repurposing existing deep learning models and designing new thermal-centric networks to achieve accurate and efficient object detection in infrared imagery.

The progress of thermal object detection has been facilitated by the creation of dedicated datasets, though such resources remain relatively scarce. The FLIR ADAS dataset, for example, provides over 10,000 annotated thermal images (with paired RGB frames) capturing cars, pedestrians, bicycles, and other road entities in day and night scenes. Similarly, the KAIST multispectral pedestrian benchmark introduced by Hwang et al. established an early standard for evaluating infrared-vision algorithms. More recent datasets target specific use cases: the LLVIP dataset (2021) consists of ~15,488 meticulously aligned visible–IR image pairs collected in extremely low-light urban environments, reflecting scenarios where thermal cameras offer clear advantages. Despite these efforts, most publicly available thermal datasets cover limited geographic regions or climate conditions. Models trained on one dataset can struggle to generalize to new locales due to differences in background thermal patterns, object appearances, or sensor characteristics – a clear sign of domain shift. Recognizing this, there is a growing consensus on the need for region-specific thermal datasets to broaden the diversity of training data. For instance, Teledyne FLIR recently released a European thermal imaging dataset spanning six different cities to complement its earlier US datasets, thereby extending the geographic and environmental coverage for automotive thermal sensing. Such region-targeted datasets (capturing varied climates, urban layouts, and cultural contexts) are expected to improve the robustness of ADAS vision models by exposing them to a wider range of thermal scenes. Motivated by these trends, this study introduces Thermal-ADAS-TR, a new thermal dataset collected in Türkiye for ADAS perception research.

2. METHODS

We collected thermal data in Ankara, Türkiye using the thermal imaging module of the DJI Mavic 3T platform. A few example frames are shown in Figure 1. The platform was mounted on the vehicle body to simulate a forward-looking ADAS perspective. A total of 1,194 thermal frames were recorded at a resolution of 640x512 pixels, covering various urban traffic scenarios under different lighting conditions. As illustrated in Figure 2, the dataset includes recordings captured under both day and night conditions.



Figure 1. Representative samples from the Thermal-ADAS-TR dataset.

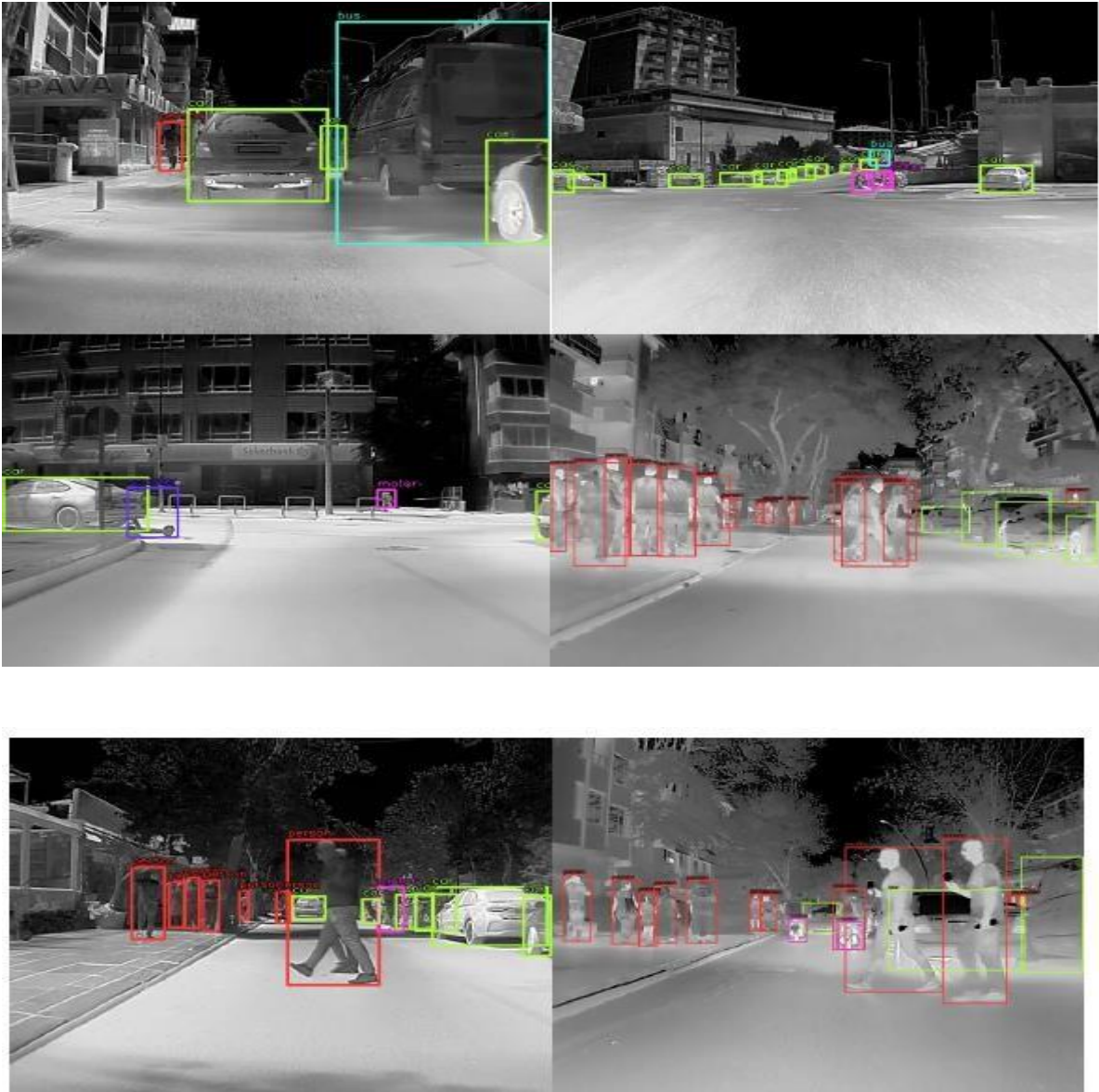


Figure 2. Example frames captured under day (left) and night (right) conditions.

The dataset was split into 718 frames for training, 242 frames for validation, and 234 frames for testing. To maintain temporal consistency, consecutive frames were kept within the same subset during splitting. This approach aimed to prevent nearly identical frames from appearing in different subsets, thus improving model generalization.

All frames were manually annotated using the X-AnyLabeling tool. Bounding boxes were created for seven object classes: person, car, bus, truck, bike, motor, and scooter. The number of annotations for each class is summarized in Table 1. We refer to this dataset as Thermal-ADAS-TR, which, to the best of our knowledge, is the first thermal ADAS dataset collected in Türkiye.

Table 1: Thermal-ADAS-TR dataset annotations

Class id	Train+Val Count	Test Count
Bike	6	8
Bus	196	76
Car	6153	1547
Motor	249	125
Person	3063	637
Scooter	30	18
Truck	246	43
Total	9943	2454

Training and evaluation were performed using the Ultralytics framework (version 8.3.196) in Python 3.10.11 with PyTorch 2.8.0 (CUDA 12.9). All experiments were conducted on an NVIDIA GeForce RTX 5080 GPU with 16 GB of VRAM.

For benchmarking, five YOLO-based object detection models (YOLOv8s, YOLOv11l, YOLOv11s, YOLOv12s, and YOLOv12m) were trained under three different configurations: (i) using only the proposed Thermal-ADAS-TR dataset, (ii) using only the FLIR-ADAS dataset (aligned with the classes in Thermal-ADAS-TR), and (iii) using a combined dataset consisting of both Thermal-ADAS-TR and FLIR-ADAS. Each model was trained for 100 epochs using the Adam optimizer with default hyperparameter settings.

Performance was evaluated using standard object detection metrics: mean Average Precision (mAP) at IoU = 0.5 (mAP@0.5) and mean Average Precision at IoU thresholds between 0.5 and 0.95 (mAP@[0.5:0.95]).

3. RESULTS

The overall performance results are summarized in Table 2.

Table 2: Performance of YOLOv8s, YOLOv11l, YOLOv11s, YOLOv12s and YOLOv12m models trained and tested on different dataset combinations (Thermal-ADAS-TR ,FLIR-ADAS, MIXED), yellow boxes are best results

Pair	mAP50	mAP50-95	Pair	mAP50	mAP50-95
yolo12s_ Thermal-ADAS- TR -train Thermal-ADAS- TR_test	0.37	0.25	yolo12s_ Thermal-ADAS- TR -train Flir_test	0.08	0.04
yolo12s_mix_train Thermal-ADAS-TR_test	0.49	0.35	yolo12s_mix_train Flir_test	0.437	0.2390
yolo12s_flir_train Thermal-ADAS-TR_test	0.33	0.22	yolo12s_flir_train Flir_test	0.436	0.2411
yolo12m_ Thermal-ADAS- TR -train Thermal-ADAS- TR_test	0.44	0.3	yolo12m_ Thermal-ADAS- TR -train Flir_test	0.098	0.056
yolo12m_mix_train Thermal-ADAS-TR_test	0.49	0.33	yolo12m_mix_train Flir_test	0.47	0.27
yolo12m_flir_train Thermal-ADAS-TR_test	0.33	0.22	yolo12m_flir_train Flir_test	0.460	0.26
yolo11l_ Thermal-ADAS- TR -train Thermal-ADAS- TR_test	0.45	0.33	yolo11l_ Thermal-ADAS- TR -train Flir_test	0.129	0.069
yolo11l_mix_train Thermal-ADAS-TR_test	0.53	0.37	yolo11l_mix_train Flir_test	0.495	0.277
yolo11l_flir_train Thermal-ADAS-TR_test	0.40	0.27	yolo11l_flir_train Flir_test	0.479	0.271
yolo11s_ Thermal-ADAS- TR -train Thermal-ADAS- TR_test	0.46	0.31	yolo11s_ Thermal-ADAS- TR -train Flir_test	0.130	0.068
yolo11s_mix_train Thermal-ADAS-TR_test	0.51	0.36	yolo11s_mix_train Flir_test	0.422	0.237
yolo11s_flir_train Thermal-ADAS-TR_test	0.34	0.22	yolo11s_flir_train Flir_test	0.440	0.246
yolo8s_ Thermal-ADAS- TR -train Thermal-ADAS-TR_test	0.48	0.33	yolo8s_ Thermal-ADAS- TR -train Flir_test	0.16	0.09
yolo8s_mix_train Thermal-ADAS-TR_test	0.49	0.34	yolo8s_mix_train Flir_test	0.47	0.25
yolo8s_flir_train Thermal- ADAS-TR_test	0.35	0.24	yolo8s_flir_train Flir_test	0.43	0.24
yolo12s_ Thermal-ADAS- TR -train Mix_test	0.114	0.077	yolo12s_mix_train Mix_test	0.319	0.178
yolo12s_flir_train Mix_test	0.311	0.173	yolo12m_ Thermal-ADAS- TR -train Mix_test	0.162	0.104
yolo12m_mix_train Mix_test	0.368	0.212	yolo12m_flir_train Mix_test	0.329	0.188

yolo11l_ Thermal-ADAS- TR -train Mix_test	0.174	0.115	yolo11l_mix_train Mix_test	0.394	0.231
yolo11l_flir_train Mix_test	0.348	0.200	yolo11s_ Thermal-ADAS- TR -train Mix_test	0.188	0.116
yolo11s_mix_train Mix_test	0.33	0.196	yolo11s_flir_train Mix_test	0.314	0.177
yolo8s_ Thermal-ADAS- TR -train Mix_test	0.209	0.141	yolo8s_mix_train Mix_test	0.347	0.194
yolo8s_flir_train Mix_test	0.310	0.175			

Models trained on the Thermal-ADAS-TR training set exhibited poor cross-domain performance on the FLIR-ADAS test set (e.g., mAP@0.5 for YOLOv12s was 0.08). This indicates that the small size of the dataset hinders model generalization. The highest performance on the Thermal-ADAS-TR test set was achieved by the YOLOv11s model trained on the mixed dataset, with an mAP@0.5 of 0.51. Similarly, the highest performance on the FLIR-ADAS test set was achieved by the YOLOv11l model trained on the mixed dataset, with an mAP@0.5 of 0.495.

For each detection model, the highest test performance was consistently achieved when trained on the mixed dataset. Overall, these findings suggest that joint training on FLIR-ADAS and Thermal-ADAS-TR improves detection accuracy and cross-domain generalization compared to single-dataset training.

4. DISCUSSION

The findings of this study clearly demonstrate that dataset composition and diversity have a substantial impact on the robustness and generalization capability of thermal object detection models. Models trained exclusively on the proposed Thermal-ADAS-TR dataset exhibited limited performance due to its relatively small scale, whereas those trained solely on FLIR-ADAS achieved strong in-domain accuracy but failed to generalize to the Turkish data. This outcome highlights the well-known domain shift problem, which often arises from differences in geographic, environmental, and thermal characteristics between datasets.

When both datasets were jointly used for training, performance improved consistently across all test scenarios. This indicates that integrating region-specific data with widely used benchmarks enhances the adaptability of object detection models to unseen environments. In particular, local datasets such as Thermal-ADAS-TR provide unique contextual and environmental patterns that complement the diversity of large public datasets, resulting in more stable and generalized model behavior.

Previous research has shown that methods such as transfer learning, data augmentation, and domain adaptation can effectively improve the performance of thermal object detection systems. The results of this study align with those findings and further demonstrate that even a relatively small, region-specific dataset can significantly enhance robustness when combined with larger datasets.

Overall, these results emphasize the importance of dataset diversity and cross-domain training in developing reliable perception models for ADAS applications. The complementary value of localized and large-scale datasets provides a practical direction for building more robust and generalizable thermal object detection systems. Future work will focus on expanding the Thermal-ADAS-TR dataset in scale and diversity, as well as exploring advanced domain adaptation strategies to further improve cross-domain performance.

4.2 Study Limitations

The main limitation of this study is the relatively small size of the proposed dataset, which limited model performance when used independently for training. Future work will focus on expanding the dataset both in scale and diversity to cover different geographical regions, illumination conditions, and weather scenarios.

5. CONCLUSION

In this work, we introduced Thermal-ADAS-TR, a thermal dataset collected in Türkiye, and evaluated its performance through comprehensive experiments using the YOLOv8s, YOLOv11l, YOLOv11s, YOLOv12s, and YOLOv12m models. Results demonstrated that training only on the proposed dataset led to limited performance due to its small scale, while training exclusively on FLIR-ADAS generalized poorly to the Turkish data. The best outcomes were consistently achieved when both datasets were combined, yielding higher mAP@0.5 and mAP@[0.5:0.95] scores on all test sets.

These findings highlight the complementary value of region-specific datasets, showing that even relatively small-scale contributions can significantly enhance model robustness when integrated with larger public datasets. The dataset and benchmark results presented in this work are expected to support future research in thermal imaging-based ADAS systems and contribute to the development of more reliable perception models under diverse real-world conditions. We plan to expand and release the dataset as an open-source resource in the future.

Acknowledgements

The authors gratefully acknowledge the support provided by Aselsan Inc. during the course of this research. The use of the DJI Mavic 3T platform for data collection was made possible through institutional resources and facilities. No external financial support was received for this study.

Disclosure

No conflicts of interest.

REFERENCES

1. FLIR Systems. (2019). Teledyne FLIR ADAS Thermal Dataset v2. Retrieved from <https://www.flir.com/oem/adas-dataset-form>
2. Hwang, S., Park, J., Kim, N., Choi, Y., & Kweon, I. S. (2015). Multispectral pedestrian detection: Benchmark dataset and baseline. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 1037–1045). IEEE. <https://doi.org/10.1109/CVPR.2015.7298706>
3. DJI. (2023). Mavic 3T specifications. DJI. <https://www.dji.com/mavic-3-enterprise>
4. Jocher, G., Chaurasia, A., & Qiu, J. (2023). YOLO by Ultralytics. GitHub repository. <https://github.com/ultralytics/ultralytics>
5. Wang, W. (2023). X-AnyLabeling: Advanced Auto Labeling Solution with Added Features [Computer software]. GitHub. <https://github.com/CVHub520/X-AnyLabeling>
6. Ding, M., Guan, S., Liu, H., & Yu, K. (2024). TIR-YOLO-ADAS: A thermal infrared object detection framework for advanced driver assistance systems. IET Intelligent Transport Systems, 18(5), 822–834. <https://doi.org/10.1049/itr2.12471>
7. Farooq, M. A., Corcoran, P., Rotariu, C., & Shariff, W. (2021). Object detection in thermal spectrum for advanced driver-assistance systems (ADAS). IEEE Access, 9, 156465–156481. <https://doi.org/10.1109/ACCESS.2021.3129150>
8. Jiang, B., Wang, J., Ren, G. et al. Research on pedestrian detection method based on multispectral intermediate fusion using YOLOv7. Sci Rep 15, 16851 (2025). <https://doi.org/10.1038/s41598-025-88871-y>
9. Lyu, C., Heyer, P., Goossens, B., & Philips, W. (2022). An Unsupervised Transfer Learning Framework for Visible-Thermal Pedestrian Detection. Sensors, 22(12), 4416. <https://doi.org/10.3390/s22124416>
10. F. Munir, S. Azam and M. Jeon, "SSTN: Self-Supervised Domain Adaptation Thermal Object Detection for Autonomous Driving," 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), Prague, Czech Republic, 2021, pp. 206-213, doi: 10.1109/IROS51168.2021.9636353.

11. Gurunath, P., Sikdar, A., Udupa, S., & Sundaram, S. (2024). IndraEye: Infrared Electro-Optical UAV-Based Perception Dataset for Robust Downstream Tasks. arXiv preprint, arXiv:2410.20953. <https://arxiv.org/abs/2410.20953>
12. Suo, J., Wang, T., Zhang, X., Chen, H., Zhou, W., & Shi, W. (2023). HIT-UAV: A high-altitude infrared thermal dataset for Unmanned Aerial Vehicle-based object detection. Scientific Data, 10, 227. <https://doi.org/10.1038/s41597-023-02066-6>
13. Li, J., Wu, Y., & Yang, X. (2025). A high-performance thermal infrared object detection framework with centralized regulation (arXiv:2505.10825). arXiv. <https://arxiv.org/abs/2505.10825>
14. Guan, D., Cao, Y., Liang, J., Cao, Y., & Yang, M. Y. (2018). Fusion of multispectral data through illumination-aware deep neural networks for pedestrian detection (arXiv:1802.09972). arXiv. <https://doi.org/10.48550/arXiv.1802.09972>

PERFORMANCE COMPARISON OF BERT AND TF-IDF WITH MACHINE LEARNING METHODS ON SENTIMENT ANALYSIS

Muhammet Sinan Başarslan

Istanbul Medeniyet University
muhammet.basarslan@medeniyet.edu.tr

Fatih Kayaalp

Duzce University
fatihkayaalp@duzce.edu.tr

Abstract

Aim: In this study, both classical machine learning algorithms and deep learning-based contextual representation methods were comparatively examined for the text-based sentiment analysis problem.

Methods: Decision Tree (DT), Random Forest (RF), Naive Bayes (NB), and Support Vector Machines (SVM) models were trained using two different feature extraction techniques: TF-IDF and BERT representations. Preprocessing steps applied to the dataset included text cleaning, lowercase conversion, and removal of punctuation marks, usernames, and URLs.

Results: The models' performances was evaluated using Accuracy, Precision, Recall, F1, Specificity, and ROC-AUC metrics. Experimental results revealed that Linear SVM showed the highest performance among TF-IDF-based models. The SVM model exhibited balanced performance in terms of accuracy, F1, and sensitivity, while demonstrating strong discriminative power, particularly in the negative sentiment class. RF ranked second, NB performed at an intermediate level, and DT was the least successful model. The BERT + SVM model, trained with BERT representations, provided a significant performance increase across all evaluation metrics. ROC analyses yielded micro-average AUC = 0.845 and macro-average AUC = 0.824 values. Confusion matrix results also show that the model can distinguish between negative, neutral, and positive classes with high accuracy.

Conclusion: These findings reveal that BERT's contextual word representations enable SVM to draw clearer boundaries between classes and improve overall classification performance. Consequently, contextual representations such as BERT provide a meaningful advantage in sentiment analysis compared to the classical TF-IDF approach. Particularly in cases of class imbalances, the BERT + SVM model demonstrated a more balanced performance between sensitivity and specificity. The study shows that integrating contextual language models with classical machine learning algorithms offers a more reliable, scalable, and highly accurate solution for sentiment analysis.

Keywords: BERT, TF-IDF, Sentiment Analysis, Machine Learning

1. INTRODUCTION

In recent years, social media platforms have become one of the most intensive digital environments where individuals express their emotions, thoughts, and attitudes. Twitter, in particular, offers an important data source for examining the public's immediate reactions to social events, thanks to its short text structure and high level of interaction. During the COVID-19 pandemic, users' posts hold great potential for understanding the psychological, social, and economic effects of the outbreak. In this context, sentiment analysis methods are used as an effective tool for revealing the social mood by identifying positive, negative, and neutral tendencies in individuals' posts [1], [2].

In the emotion analysis literature, machine learning-based methods are known to achieve high success in text classification tasks [3], [4]. Classical approaches based on TF-IDF representation are frequently preferred, especially for determining meaningful word weights in short texts. In this study, text-based sentiment analysis was performed on Twitter posts related to the coronavirus (COVID-

19) process. The dataset was created from tweets posted during the pandemic, and the posts were divided into three sentiment classes: “Negative,” “Neutral,” and “Positive.”

Chakraborty et al. [5] achieved approximately 81% accuracy and F1 scores of 0.95 in some classes in their study using deep learning models on COVID-19-related tweets. Similarly, Oladri and Radhakrishnan [6] achieved 87% overall accuracy, Accuracy 0.93 for the negative class, Recall 0.84, F1 0.88 for the negative class, Precision 0.86, Recall 0.78, F1 0.82 for the neutral class, and Precision 0.82, Recall 0.94, F1 0.88 for the positive class. Jalil et al. [7] reported high classification performance, achieving an accuracy rate of 96.66% in their study on a COVID-19-related tweet dataset. Additionally, Jlifi et al. [8] achieved 93.01% accuracy, 94.03% precision, and 93.05% recall using an ensemble approach called Ens-RF-BERT. In a multi-class study published in the SciELO database, it was reported that the BERT model showed a clear superiority over classical methods with an F1 score of approximately 74% [9].

These studies demonstrate that BERT's contextual representation capabilities provide significant advantages in sentiment analysis compared to classical machine learning techniques. However, different dataset characteristics, class imbalances, and model configurations can significantly affect performance.

The primary objective of this study is to comparatively examine the success of classical machine learning algorithms in classifying social sentiment trends related to COVID-19. In this context, DT, RF, Linear SVM, and NB models were trained using TF-IDF representation. After TF-IDF, once the best method was determined, the same method was trained after BERT, and thus the performance impact of BERT and TF-IDF was also evaluated. The results aim to evaluate the sentiment analysis performance of different models using both quantitative measures (Accuracy, F1, F2, PR-AUC) and visual analyses (confusion matrix and ROC curves).

2. METHODS

This section will provide information about the dataset used in the study and the algorithms employed.

2.1. Dataset

This study utilizes a dataset collected from Twitter messages posted during the Coronavirus period and publicly shared on Kaggle for use in text-based sentiment analysis tasks. It contains three classes labeled as Negative, Neutral, and Positive [10].

During the modeling process, the original dataset was used directly; no external or synthetic data was added. The dataset was split into two subsets using the hold-out method (80% training and 20% testing sections) to preserve the class distribution of the examples. Accordingly, 32,925 examples from the total dataset were used in the training set, and 8,232 examples were used in the test set. This division aimed to observe the model's realistic generalization ability; the hold-out approach was preferred over cross-validation to provide a more direct and temporally consistent evaluation.

This structure minimized the risk of potential data leakage between the training and testing processes by allowing the performance of the models to be directly compared across different sub-data sets without retraining. Class weight parameters were considered in model training to assess the impact of class imbalances. Examples from the training-test set of the dataset are shown in Figure 1.

test	train
UserName ScreenName Location TweetAt OriginalTweet Sentiment	UserName ScreenName Location TweetAt OriginalTweet Sentiment
0 1 44953 NYC 02-03-2020 TRENDING: New Yorkers encounter empty supermar... Negative	0 3799 48751 London 16-03-2020 @MeNyrbie @Phil_Gahan @Chrisiv https://t.co/... Neutral
1 2 44954 Seattle, WA 02-03-2020 When I couldn't find hand sanitizer at Fred Me... Positive	1 3800 48752 UK 16-03-2020 advice Talk to your neighbours family to excha... Positive
2 3 44955 NaN 02-03-2020 Find out how you can protect yourself and love... Positive	2 3801 48753 Vagabonds 16-03-2020 Coronavirus Australia: Woolworths to give elde... Positive
3 4 44956 Chicagoland 02-03-2020 #Panic buying hits #NewYork City as anxious sh... Negative	3 3802 48754 NaN 16-03-2020 My food stock is not the only one which is emp... Positive
4 5 44957 Melbourne, Victoria 03-03-2020 #toiletpaper #dunnypaper #coronavirus #coronav... Neutral	4 3803 48755 NaN 16-03-2020 Me, ready to go at supermarket during the #COV... Negative

Figure 1: Training and test sample data

The data shown in Figure 1 takes into account the data in the “sentiment” and “OriginalTweet” columns. Other columns have been omitted. In addition, Figure 2 below shows the sentiment distributions of the data used in the study.



Figure 2: Dataset sentiment distribution

2.2 Text Representation Methods

The study utilized word representation methods, specifically BERT and TF-IDF.

BERT is a pre-training model first proposed by Google Research in 2018 that revolutionized natural language processing (NLP) tasks. BERT is a bidirectional transformer architecture that considers both left and right context simultaneously, allowing it to benefit from a word's preceding and following context at the same time. Thanks to this structure, BERT comprehends linguistic context and semantic relationships more deeply, delivering high performance in various NLP tasks—such as sentiment analysis tasks [10,11]. Particularly in situations where context is important, such as short messages and social media posts, the pre-training knowledge incorporated in BERT enriches the model contextually and provides effective classification capabilities [13].

On the other hand, TF-IDF (Term Frequency–Inverse Document Frequency) representation is one of the cornerstones of classical information retrieval and text mining approaches. TF-IDF weights words by considering both how frequently a word appears in a specific document (TF) and how common that same word is across the entire document corpus (IDF). This way, higher importance is given not only to frequently occurring words but also to words with distinctive features [14]. The TF-IDF representation can be used efficiently and relatively lightly in terms of computation with classical machine learning models such as NB and SVM, especially in applications such as text classification and sentiment analysis; however, it is limited in capturing features such as context, synonymy, and syntactic relationships [15].

First, basic preprocessing steps were applied to the data; texts were converted to lowercase, URLs and web page redirects, user tags, special characters, and unnecessary spaces were cleaned to obtain

a consistent format for analysis. These cleaned texts were first converted into numerical feature vectors using the TF-IDF method to create models for sentiment analysis. Subsequently, BERT was used before SVM, which yielded the best results among these models.

In the BERT-based text representation approach, contextual vector representations with 768 dimensions were extracted for each text example using a pre-trained BERT model (training embedding dimensions: 41,157; validation embedding dimensions: 3,798). This representation format captures semantic relationships between words within a sentence, enabling deep semantic features to be included in the model beyond classical methods. The dense vectors extracted using BERT were then combined with an SVM to create the model, thus investigating the powerful representation capability of BERT in conjunction with SVM performance.

2.3. Algorithms Used in the Study

In this study, four different classical machine learning algorithms, which are the most preferred methods in the literature, were evaluated to perform text-based sentiment analysis on Twitter data related to the COVID-19 process: DT, RF, SVM, and Multinomial NB. These algorithms were trained on word vectors based on TF-IDF representation and each was tested separately on the sentiment classification task.

DT is an intuitive method that classifies data by splitting it into branches, separating examples according to the most informative attributes through decision nodes [16,17]. Although preferred for its simplicity and interpretability, it has a structure prone to overfitting, especially in high-dimensional and noisy datasets [18]. The RF method, used to mitigate this issue, is an ensemble-based model formed by combining multiple DTs [19]. Each tree is trained on randomly selected sample and feature subsets, and the final class decision is made by the vote of all trees. The RF algorithm is frequently preferred in the literature for its generalization ability and stability in imbalanced data sets [19].

On the other hand, SVM is an optimization-based method that separates data by maximizing the optimal boundary (hyperplane) between classes [20,21]. This model, which can perform linear and non-linear separations using kernel functions, stands out with high accuracy rates, especially in text classification and sentiment analysis studies [22]. Finally, the NB model adopts a probabilistic approach and estimates the class probabilities of texts under the assumption that words are independent of each other [22]. NB is widely used in the classification of social media data due to its low computational cost and effective performance on large-scale texts [23].

In the study, Accuracy, F1, Precision, Recall, and Specificity were used to evaluate classification models, along with F2 and AUC-PR for the imbalanced dataset in the study. The ROC curve is also provided.

3. EXPERIMENTAL SETUP AND RESULTS

In the experimental process of this study, classical machine learning models were evaluated for the text-based sentiment analysis task. First, the data was cleaned using a specially defined preprocessing function; URLs, tags, usernames, and punctuation marks were removed from the texts, and then all characters were converted to lowercase. The cleaned texts were converted into numerical representations using the TF-IDF (Term Frequency-Inverse Document Frequency) method. The training and validation data were created using stratified training-validation splitting to preserve the class distribution; the DT, RF, SVM, and NB algorithms were used in the modeling process.

The performance of each model was evaluated using Accuracy, Precision, Recall, F1 and F2 scores, Specificity, and PR-AUC metrics. ROC curves and AUC values were calculated separately for each class and as macro averages; model score matrices were obtained using predict_proba methods. In the multi-class scenario, ROC curves were plotted using macro and micro averages. Additionally, class weight parameters and F2/PR-AUC scores were considered to mitigate the effects of potential imbalances between classes.

Finally, the experiments were performed in a Python environment using the NumPy, Matplotlib, and scikit-learn libraries. This structure ensured that both model training and performance measurement were carried out systematically. The experiments were run on the Google Colab Pro platform. Python 3 was used in the study, and the latest version of the PyTorch framework was used for training and testing deep learning models. An integrated working environment with Google Drive was created for storing, loading, and sharing data and model weights. The experimental performance of the models obtained in the study is presented in Table 1.

Table 1: Experimental Results

	Model	Accuracy	F1	Precision	Recall	Specificity	F2	AUC-PR
TF-IDF	DT	0.6534	0.6400	0.6500	0.6400	0.7900	0.6400	0.6700
	RF	0.7532	0.7400	0.7400	0.7500	0.8500	0.7500	0.7800
	SVM	0.8271	0.8200	0.8200	0.8200	0.8900	0.8300	0.8600
	NB	0.6984	0.6900	0.7000	0.6900	0.8000	0.7000	0.7200
BERT	SVM	0.8562	0.8488	0.8480	0.8656	0.9260	0.8571	0.8650

Table 1 shows the sentiment analysis performance of classical models based on TF-IDF and BERT representation. In Table 1, Linear SVM achieved the highest success. SVM demonstrated a balanced profile between accuracy, sensitivity, and specificity, showing strong selectivity, particularly in distinguishing negative examples. RF ranks second; its Recall and Specificity values indicate balanced but lower performance than SVM in distinguishing positive/negative classes. NB is moderate; it often yields reasonable results at positive/negative extremes, but the error rate increases in the neutral class. DT has the lowest generalization success; this is consistent with its difficulty in capturing subtle sentiment transitions due to the high-variance structure of single trees and the context-insensitivity of TF-IDF. Consequently, among TF-IDF-based approaches, SVM is the most reliable option due to its linear discriminative power and relative robustness against class imbalances; RF follows it, while NB and DT produced consistent but lower scores.

As SVM has the best performance among other methods with TF-IDF; BERT+SVM model has been created to evaluate the effect of BERT with SVM on sentiment task. When BERT+ SVM is used on the same data, the effect of contextual representation is clearly evident and bolded to highlight the performance. Compared to TF-IDF+SVM, better results were obtained in accuracy, F1, sensitivity, and specificity. AUC-PR remained at a similar level. Table 1 shows that BERT's ability to capture contextual information reduces false negatives (Recall) and strengthens the ability to correctly exclude negatives (Specificity), while TF-IDF's word-level constraints lead to errors, especially in neutral borderline cases.

Linear SVM is the best classical method on TF-IDF; however, SVM with BERT representations provides a significant leap in overall accuracy and F1. Even in imbalanced data scenarios, the upward shift in the Recall-Specificity balance confirms that contextual models are a more reliable and scalable option for sentiment analysis. The performance imbalance of the dataset is evaluated using a comparison of F2 and AUC-PR, as shown in Figure 3.

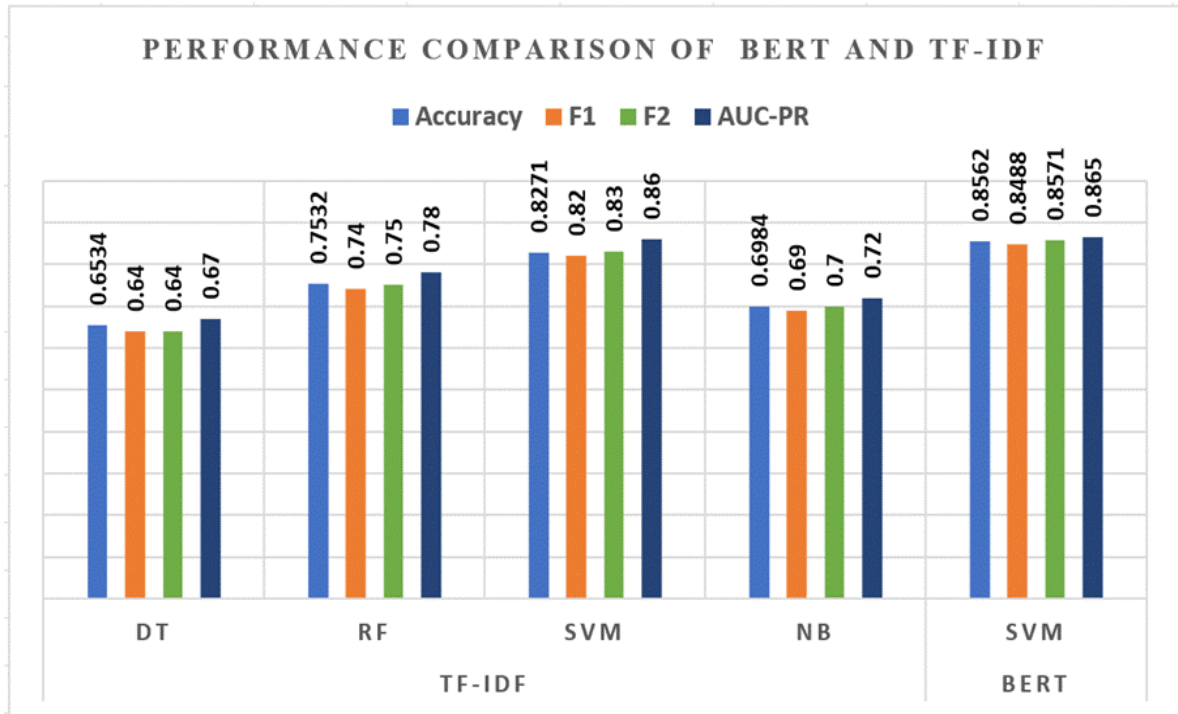


Figure 3: Confusion matrix of models after TF-IDF and BERT

Figure 3 shows the comparisons of the accuracy, F1, F2, and AUC-PR performance metrics of BERT and TF-IDF-based models on DT, RF, SVM, and NB. Based on TF-IDF representation, the SVM algorithm shows the highest success, and these results are seen to be ahead of other methods. The RF model ranks second, while the DT and NB models exhibit lower performance. After BERT representation, the SVM model outperformed all TF-IDF-based models. This result shows that BERT, which is a contextual word representation, significantly improves classification performance compared to the traditional TF-IDF method. It can be concluded that BERT + SVM has the highest success rate and provides text classification.

Figure 4 shows the confusion matrix results for the four models after TF-IDF, while Figure 5 shows the results for the best SVM after BERT.

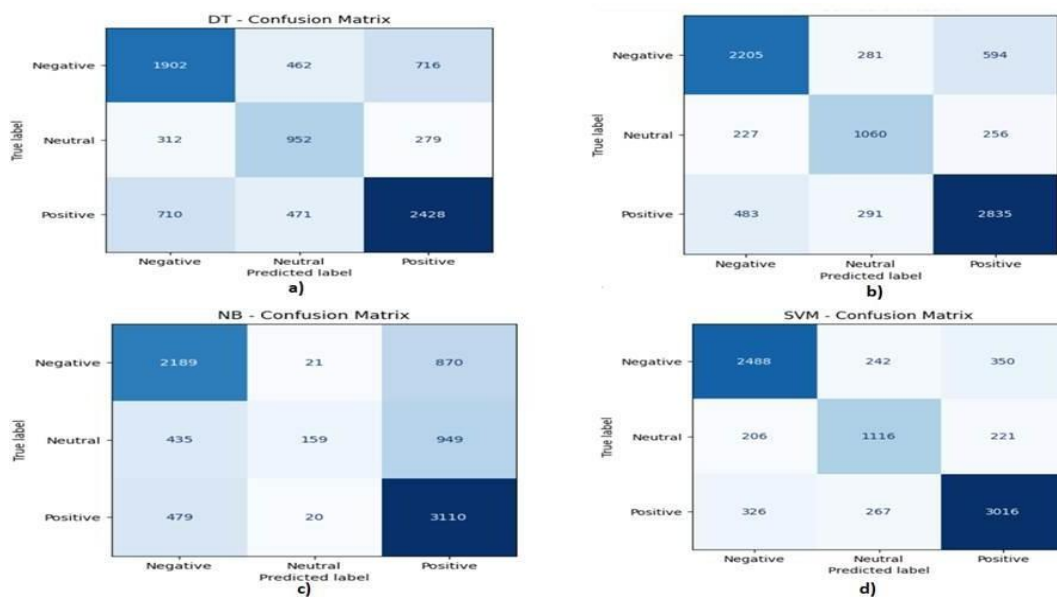


Figure 4: Confusion matrix of models after TF-IDF a)DT b)RF c)NB d)SVM

Figure 4 shows the confusion matrices of four different machine learning models, and these matrices show the prediction performance of the models on negative, neutral, and positive classes. The DT model performed reasonably well, making 2428 correct predictions in the positive class and 1902 correct predictions in the negative class; however, it has a high tendency to confuse the negative and positive classes. It shows weak performance in distinguishing the neutral class. The RF model achieved more balanced success across all classes and reached a very high accuracy rate with 2835 correct predictions in the positive class. Thanks to its ensemble structure, this model has become one of the models with the highest generalization ability.

the RF model in the section shows how accurately the model predicts the three sentiment classes (Negative, Neutral, Positive). Based on the values in the matrix, the model generally performs well. The majority of examples in the “Negative” class (2205) were correctly predicted, while a small number of examples were misclassified as “Neutral” (281) and “Positive” (594). This indicates that the model is quite good at distinguishing negative expressions. In the “Neutral” class, there are 1060 correct predictions, but 227 examples are misclassified as “Negative” and 256 examples as “Positive.” This shows that neutral expressions are slightly more challenging for the model. The “Positive” class shows the highest success with 2835 correct predictions; only 483 examples were incorrectly predicted as “Negative” and 291 examples as “Neutral.”

The NB model performed well in the extreme classes (positive and negative) but was quite weak in the neutral class. Correctly predicting only 159 examples in the neutral class indicates that the model failed to distinguish this class. This situation may stem from the NB assumption of independence between variables being insufficient in real data distributions. Finally, the SVM model generally yielded the most successful results. It achieved the lowest error rate by correctly predicting 2488 in the negative class, 1116 in the neutral class, and 3016 in the positive class. This indicates that SVM can define the boundaries between classes more sharply and has a high overall accuracy rate.

Overall, among the four models, SVM has the highest performance, RF is second, DT is intermediate, and NB is the weakest model.

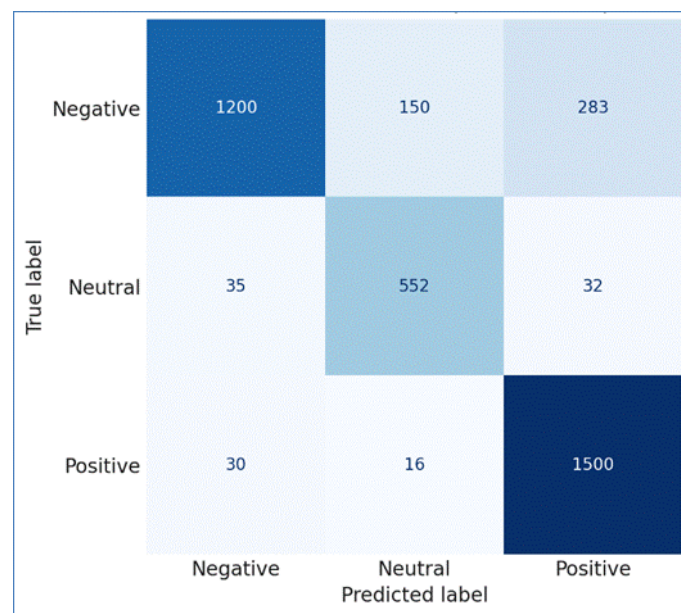


Figure 5: Confusion matrix obtained with SVM after BERT

Figure 5 shows the confusion matrix obtained with the SVM model after BERT. This matrix shows the model's classification success across three different sentiment classes (negative, neutral, and positive).

The model performed quite well on the negative class; it correctly predicted 1200 examples, misclassifying only 150 as neutral and 283 as positive. This indicates that the model has strong accuracy in distinguishing negative sentiments. Generally successful results were also obtained in the neutral class; 552 examples were correctly classified, with errors made in only 67 examples (35 negative, 32 positive). This also shows that the model can largely distinguish neutral expressions from other classes.

In the positive class, 1,500 correct predictions were made, with only 46 examples (30 negative, 16 neutral) misclassified. This is a very high success rate and proves that the model is effective at identifying positive statements.

Overall, after BERT-based feature extraction, the SVM model achieved high accuracy rates across all three classes. Errors mostly occurred between neighboring classes (e.g., neutral–positive or negative–neutral transitions). This result shows that, thanks to BERT's ability to understand context, SVM can strongly distinguish between classes and that the model's overall performance is quite successful.

Figure 6 shows the ROC graphs of the four models after TF-IDF, and Figure 6 shows the ROC graphs of BERT before the best SVM.

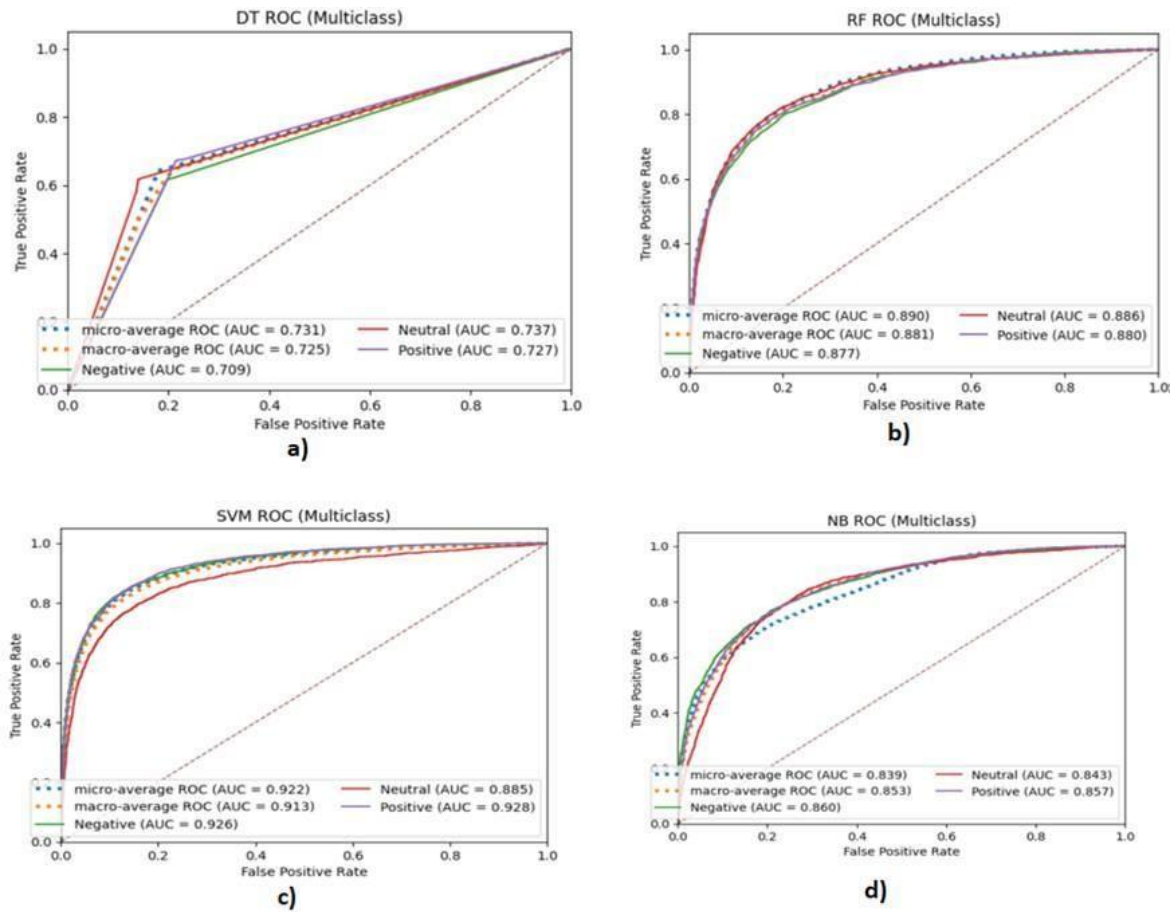


Figure 6: ROC curves for the results of the methods used after TF-IDF a) DT b) RF c) SVM d) NB

Figure 6 compares the multi-class ROC curves of four different machine learning methods (DT, RF, SVM, and NB) used after TF-IDF. ROC curves show the relationship between the models' true positive rate and false positive rate, and the AUC value is an important measure summarizing the model's overall classification success.

AUC values are generally low in the DT model; micro-average AUC is calculated as 0.731, macro-average AUC as 0.725. This result shows that the model exhibits inconsistent performance across classes and is particularly weak in the negative class (AUC = 0.709). The RF model, on the other hand, performed significantly better. The micro-average AUC value is 0.890, while the macro-average AUC value is 0.881. The curves for all classes hover around 0.88, indicating that this model has strong generalization capabilities.

The SVM model yielded the most successful results. The micro-average AUC value is 0.922, and the macro-average AUC value is 0.913. High accuracy rates were observed, particularly in the positive class (AUC = 0.928) and negative class (AUC = 0.926). This demonstrates that SVM is the model that can most clearly distinguish between classes. The NB model performed at an intermediate level; with micro-average AUC 0.839 and macro-average AUC 0.853 values, it lagged behind RF and SVM.

Overall, the SVM model showed the best ROC-AUC performance after TF-IDF, followed by RF. While NB showed average performance, the DT model had the lowest accuracy. These results prove that SVM is more effective in determining non-linear decision boundaries, while RF is more effective in high-variance samples.

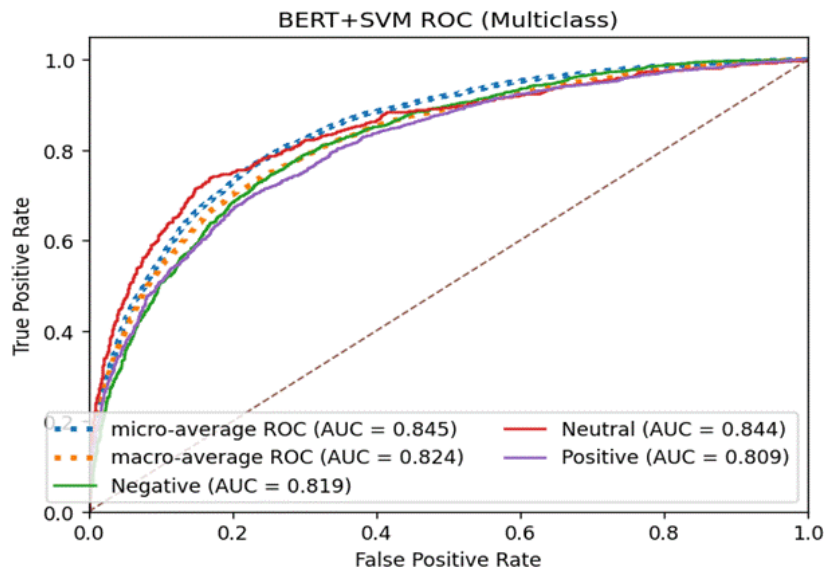


Figure 7: ROC graph of SVM used after BERT

The AUC values in Figure 7 numerically summarize the overall success of the model. The micro-average AUC value was calculated as 0.845, while the macro-average AUC value was calculated as 0.824. These values indicate that the model has good overall discrimination power. When examined by class, the neutral class AUC = 0.844, the negative class AUC = 0.819, and the positive class AUC = 0.809.

The results show that the contextual features obtained with BERT contribute significantly to the SVM model, and that the model can distinguish the three classes with reasonable accuracy. However, the slightly lower AUC value for the positive class compared to the others suggests that the model makes relatively more errors in distinguishing this class. Overall, the BERT + SVM approach

demonstrated strong and balanced performance, and the ROC curves also confirmed that the model is a successful classifier.

4. DISCUSSION AND CONCLUSION

In this study, a comparative performance analysis was conducted for the text-based sentiment analysis problem using both classical machine learning methods (DT, RF, NB, SVM) and the deep learning-based BERT model. Experimental findings showed that the feature extraction method used (TF-IDF or BERT) had a decisive effect on model performance.

Among classical methods based on TF-IDF representation, the Linear SVM model demonstrated the highest accuracy, F1 score, and generalization success. The linear discriminative power of SVM provided a distinct advantage, particularly in the correct classification of the negative class. RF was the second most successful method, while NB was moderate and DT had the lowest accuracy rate. This situation reveals that ensemble-based models and margin-based classifiers have stronger generalization capabilities compared to single-tree models.

Contextual representations obtained with the BERT model significantly improved the performance of SVM. The BERT + SVM approach provided a noticeable improvement in accuracy, F1 score, recall, and specificity metrics compared to TF-IDF + SVM. ROC analyses yielded micro-average AUC = 0.845 and macro-average AUC = 0.824 values. Furthermore, confusion matrix results show that the model can distinguish between negative, neutral, and positive classes with high accuracy. These findings confirm that BERT's ability to capture word context contributes to the model's ability to more accurately determine class boundaries.

Overall, contextual representations (BERT) provide a clear advantage over classical TF-IDF representation in sentiment analysis. Particularly in datasets with class imbalances, the BERT + SVM model demonstrated more reliable and scalable performance while maintaining Recall and Specificity balance.

The BERT + SVM model developed in this study has demonstrated a significant performance advantage over previous classical sentiment analysis approaches. As noted in comprehensive reviews by Liu [1] and Giachanou & Crestani [2], classical methods are mostly based on word-frequency representations (e.g., TF-IDF, Bag-of-Words) and can capture contextual information only to a limited extent. However, the BERT-based contextual representations used in this study have strengthened the decision boundaries of the SVM classifier by modeling the semantic relationships and emotional tones of words at a deeper level.

The results obtained from the BERT + SVM model in this study significantly outperform the performance values reported in the literature. Chakraborty et al. [5] achieved approximately 81% accuracy and F1 scores of 0.95 in some classes using deep learning models on COVID-19-related tweets. Similarly, Oladri and Radhakrishnan [6] reported accuracy of 87% and results in the range of 0.82–0.88 F1 in their BERT-based analysis. Jalil et al. [7] achieved an accuracy rate of 96.66% with deep learning-based models; Jlifi et al. [8] obtained 93.01% accuracy, 94.03% precision, and 93.05% recall values in the Ens-RF-BERT model. In a multi-class analysis conducted by González-Díaz et al. [9], an F1 score of approximately 74% was reported.

The BERT + SVM model developed in this study demonstrated a performance superior to many studies in the literature, with 85.62% accuracy and 84.88% F1 scores. Therefore, it can be said that the current model surpasses the results reported in similar studies in terms of both accuracy and F1 metric. This indicates that when BERT's contextual representations are combined with SVM's discriminative power, it provides superior classification success in sentiment analysis related to the COVID-19 period.

In conclusion, the findings of this study demonstrate that BERT-based representations can significantly improve sentiment analysis performance by supporting classical machine learning

methods. Future studies may further enhance model performance by utilizing larger datasets, different pre-trained language models (e.g., RoBERTa, DistilBERT), and hyperparameter optimization.

REFERENCES

1. M. Cinelli et al., "The COVID-19 social media infodemic," *Sci Rep*, vol. 10, no. 1, p. 16598, Oct. 2020, doi: 10.1038/s41598-020-73510-5.
2. J. Samuel, G. G. Md. N. Ali, Md. M. Rahman, E. Esawi, and Y. Samuel, "COVID-19 Public Sentiment Insights and Machine Learning for Tweets Classification," *Information*, vol. 11, no. 6, p. 314, Jun. 2020, doi: 10.3390/info11060314.
3. W. Medhat, A. Hassan, and H. Korashy, "Sentiment analysis algorithms and applications: A survey," *Ain Shams Engineering Journal*, vol. 5, no. 4, pp. 1093–1113, Dec. 2014, doi: 10.1016/j.asej.2014.04.011.
4. B. Pang and L. Lee, "Opinion Mining and Sentiment Analysis," *Foundations and Trends® in Information Retrieval*, vol. 2, no. 1–2, pp. 1–135, 2008, doi: 10.1561/15000000011.
5. T. Chakraborty, S. Bhattacharjee, et al., "Sentiment analysis of COVID-19 tweets by deep learning and lexicon-based approaches," *Applied Soft Computing*, vol. 97, p. 106759, 2020.
6. R. Oladri and R. Radhakrishnan, "BERT for Twitter sentiment analysis during COVID-19," *International Research Journal of Computer Science and Software Technology*, vol. 6, no. 1, pp. 25–33, 2025.
7. Z. Jalil et al., "COVID-19 related sentiment analysis using state-of-the-art deep learning models," *Frontiers in Public Health*, vol. 9, p. 812735, 2021.
8. B. Jlifi et al., "Ens-RF-BERT: An ensemble approach for COVID-19 tweet sentiment classification," *Social Network Analysis and Mining*, vol. 14, no. 2, p. 1240, 2024.
9. R. González-Díaz et al., "Multi-class sentiment analysis of COVID-19 tweets using BERT," *Computación y Sistemas (SciELO)*, vol. 28, no. 2, pp. 507–519, 2024.
10. "COVID-19 NLP Text Classification" dataset, Datatattle, Kaggle. [Online]. Available: <https://www.kaggle.com/datasets/datatattle/covid-19-nlp-text-classification?resource=download>, accessed: Oct. 9, 2025.
11. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," Oct. 2018.
12. M. S. Başarslan and F. KAYAALP, "Sentiment Analysis on Social Media Reviews Datasets with Deep Learning Approach," *Sakarya University Journal of Computer and Information Sciences*, Feb. 2021, doi: 10.35377/saucis.04.01.833026.
13. K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text Classification Algorithms: A Survey," *Information*, vol. 10, no. 4, p. 150, Apr. 2019, doi: 10.3390/info10040150.
14. S. N. Başa and M. S. Basarslan, "Sentiment Analysis Using Machine Learning Techniques on IMDB Dataset," in *2023 7th International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT)*, IEEE, Oct. 2023, pp. 1–5. doi: 10.1109/ISMSIT58785.2023.10304923.
15. M. Öznaneci, M. S. Başarslan, N. Bulut, and H. Ankaralı, "NLP-Based Pain Prediction Using Machine Learning and Boosted Models: A Comparative Analysis of TF-IDF and BoW Representations with Headache Data," 2025, pp. 501–509. doi: 10.1007/978-3-031-97992-7_56.
16. R. Polikar, "Ensemble based systems in decision making," *IEEE Circuits and Systems Magazine*, vol. 6, no. 3, pp. 21–45, 2006, doi: 10.1109/MCAS.2006.1688199.
17. T. Öztürk, Z. Turgut, G. Akgün, and C. Köse, "Machine learning-based intrusion detection for SCADA systems in healthcare," *Network Modeling Analysis in Health Informatics and Bioinformatics*, vol. 11, no. 1, p. 47, Dec. 2022, doi: 10.1007/s13721-022-00390-2.
18. S. Ustebay, Z. Turgut, and M. A. Aydin, "Cyber Attack Detection by Using Neural Network Approaches: Shallow Neural Network, Deep Neural Network and AutoEncoder," 2019, pp. 144–155. doi: 10.1007/978-3-030-21952-9_11.
19. L. Breiman, "Random Forests," *Mach Learn*, vol. 45, no. 1, pp. 5–32, 2001, doi: 10.1023/A:1010933404324.
20. T. H. Jaya Hidayat, Y. Ruldeviyani, A. R. Aditama, G. R. Madya, A. W. Nugraha, and M. W. Adisaputra, "Sentiment analysis of twitter data related to Rinca Island development using Doc2Vec and SVM and logistic regression as classifier," *Procedia Comput Sci*, vol. 197, pp. 660–667, 2022, doi: 10.1016/j.procs.2021.12.187.
21. R. Ahmed, M. Bibi, and S. Syed, "Improving Heart Disease Prediction Accuracy Using a Hybrid Machine Learning Approach: A Comparative study of SVM and KNN Algorithms," *International Journal of Computations, Information and Manufacturing (IJCIM)*, vol. 3, no. 1, pp. 49–54, Jun. 2023, doi: 10.54489/ijcim.v3i1.223.

22. M. Veziroğlu and Ihsan Bucak, “Haber Sınıflandırma Sistemlerinde Naive Bayes ve Makine Öğrenmesi Algoritmaları Arasında Performans Karşılaştırması,” *Journal of the Institute of Science and Technology*, vol. 15, no. 1, pp. 57–70, 2025.
23. C. A. Ramaputra, M. H. Z. Al Faroby, and B. R. Lidiawaty, “Sentiment Analysis of User Reviews on Cryptocurrency Application: Evaluating the Impact of Dataset Split Scenarios Using Multinomial Naive Bayes,” *The Indonesian Journal of Computer Science*, vol. 13, no. 4, Aug. 2024, doi: 10.33022/ijcs.v13i4.4263.

**FROM LYRICS TO INSIGHTS: MULTIDIMENSIONAL EMOTION, THEME,
AND DEMOGRAPHIC ANALYSIS IN TURKISH MUSIC****Ege Kutlu**

Department of Data Science and Artificial Intelligence, Institute for Data Science & Artificial
Intelligence, Boğaziçi University, İstanbul/Türkiye
ktluege@gmail.com

Ayşe Arslan

Turkcell Communication Services, İstanbul/Türkiye
aysel.arslan@turkcell.com.tr

Ali Buğra Kanburoğlu

Turkcell Communication Services, İstanbul/Türkiye
ali.kanburoglu@turkcell.com.tr

Yasemin Sarper

Turkcell Communication Services, İstanbul/Türkiye
yasemin.sarper@turkcell.com.tr

Tuna Basaran

Turkcell Communication Services, İstanbul/Türkiye
tuna.basaran@turkcell.com.tr

Berna Kalender

Turkcell Communication Services, İstanbul/Türkiye
berna.icellioglu@turkcell.com.tr

Burcu Gülmüş Çevikbaş

Turkcell Communication Services, İstanbul/Türkiye
burcu.gulmus@turkcell.com.tr

Yasemin Yetimova

Turkcell Communication Services, İstanbul/Türkiye
yasemin.yetimova@turkcell.com.tr

Naz Albayrak Satırcıoğlu

Turkcell Communication Services, İstanbul/Türkiye
naz.santircioglu@turkcell.com.tr

Mehmet Turan

Department of Computer Engineering, Faculty of Engineering,
Boğaziçi University, İstanbul/Türkiye
mehmet.turan@bogazici.edu.tr

Abstract

This study presents a large-scale computational framework for analyzing the emotional and cultural dimensions of Turkish song lyrics. The project integrates automatic transcription using the Whisper model with advanced natural language analysis via Llama-4-Maverick. The system generates structured psycho-linguistic metadata—covering sentiment polarity, emotional intensity, thematic content, demographic markers, and genre prediction across a hundred songs. Results reveal that Turkish lyrics are characterized by high emotional intensity (mean = 6.8/10), based on a data-driven six-category emotion taxonomy encompassing joy/empowerment, sadness/heartbreak, anger, disappointment/regret, nostalgia/longing, and calmness/closure, and are dominated by themes of love, loss, and nostalgia. Distinct generational and stylistic patterns emerge: younger audiences

engage with pop and rap lyrics rich in slang and digital expressions, whereas older listeners resonate with idiomatic and symbolic language in folk and arabesque genres. The proposed framework achieved a classification accuracy of 78.8% (calculated as the ratio of correctly predicted genres to the total number of songs) and demonstrated strong interpretive consistency. This performance establishes a foundation for data-driven cultural analytics, music recommendation systems, and cross-lingual studies of emotional expression in music.

Aim: This study analyzes a corpus of Turkish song lyrics constructed using automated transcription and large language model analysis to extract insights into their emotional depth, thematic composition, and demographic indicators. The ultimate goal of this psycho-linguistic and cultural exploration is to build a foundational framework for a more nuanced and culturally-aware music recommendation engine.

Methods: Songs from the selected playlist were processed through a scalable pipeline combining Whisper-based transcription and the Llama-4-Maverick language model. The Llama-4-Maverick model was selected specifically for its advanced Turkish language understanding and strong instruction-following capabilities, and it was utilized in an inference-only capacity without any fine-tuning. (Radford et al., 2022) The lyrics were then analyzed through a structured JSON schema encompassing sentiment polarity, emotional intensity, thematic content, demographic markers, and genre prediction. The system employed asynchronous processing for large-scale scalability. Random sampling and expert evaluation were used to ensure schema consistency and interpretive validity. (Yılmaz & Deniz, 2024)

Results:

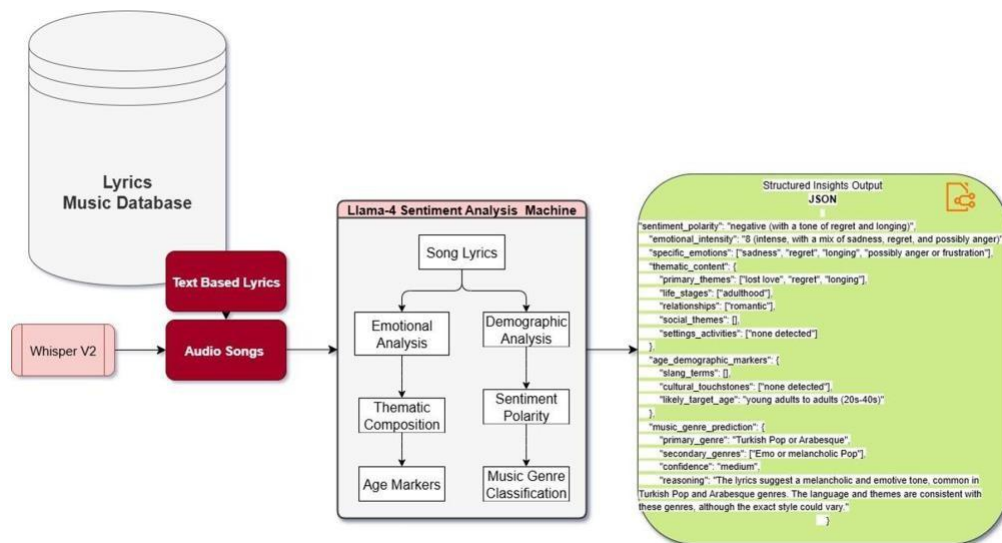


Figure 1: Schematic of the Lyric Analysis Framework

This diagram illustrates the end-to-end computational pipeline... The process begins with data from a Turkish Music Streaming Service Database, which contains both text-based lyrics and audio songs which contains both text-based lyrics and audio songs. Audio-only tracks are first transcribed into text using the Whisper V2 model. The resulting lyrics are then processed by the core analytical engine, the Llama-4 Sentiment Analysis Machine. This model performs a multidimensional analysis encompassing several sub-tasks: Emotional Analysis, Thematic Composition, Demographic Analysis, Sentiment Polarity, Age Markers, and Music Genre Classification. The final output is a Structured Insights JSON object, which contains detailed and machine-readable metadata. An example output on the right shows the schema, including fields like emotional_intensity,

thematic_content (e.g., "lost love," "regret"), age_demographic_markers, and a music_genre_prediction complete with reasoning. This framework effectively transforms unstructured lyrical content into quantifiable and interpretable data.

The analysis encompassed a broad range of Turkish songs across pop, arabesque, rap, rock, and folk genres. Positive sentiment was observed in approximately 45% of songs, negative in 35%, and neutral in 20%. Emotional intensity ranged from 2 to 10 (mean = 6.8), with arabesque and rap showing the highest values. Dominant themes included love (60%), loss (30%), and nostalgia (30%), while social and political motifs appeared in around 10%. Slang and generational expressions indicated demographic differentiation: younger audiences favored pop/rap with internet slang, while older audiences resonated with idiomatic expressions in folk and arabesque. Genre prediction accuracy reached approximately 78.8%, with high confidence for mainstream categories.

Conclusion: The developed framework demonstrates a scalable and culturally grounded approach for analyzing Turkish lyrics, producing consistent psycho-linguistic metadata across thousands of songs. Findings reveal strong emotional intensity, dominant themes of love and loss, and clear generational segmentation. This methodology establishes a foundation for data-driven cultural analytics, music recommendation systems, and comparative linguistic research.

Keywords: Sentiment, large language models, musicocultural analytics, psycho-linguistics, lyrics analysis

1. INTRODUCTION

Music serves as a central medium for emotional and cultural expression, capturing societal moods, generational identities, and shared experiences. In Turkish music, lyrical expression often exhibits high emotional intensity, metaphorical richness, and social commentary. However, despite the significance of Turkish music in reflecting cultural identity, systematic computational studies of its lyrics remain scarce. (Durahim et al., 2018) Previous research on lyric analysis has predominantly focused on English language corpora, employing lexicon-based sentiment analysis or topic modeling to identify broad emotional and thematic categories. While informative, such approaches often fail to capture nuanced emotional gradients, metaphorical structures, or demographic markers inherent in lyrics. (Radford et al., 2022)

Given the lack of structured lyric metadata in commercial music platforms, this gap limits both academic inquiry and the development of advanced music recommendation systems. Addressing these limitations, the present study introduces a scalable analytical pipeline that integrates Whisper-based transcription with large language model (LLM) processing. By converting Turkish song lyrics into structured psycho linguistic data, this framework enables systematic exploration of sentiment, emotional intensity, thematic diversity, and generational indicators.

The objectives of this study are threefold:

1. To design an automated pipeline generating consistent psycho-linguistic metadata across Turkish songs.
2. To provide corpus-level insights into emotional, thematic, and demographic characteristics of Turkish lyrics.
3. To demonstrate the utility of this framework in cultural analytics, musicology, and industry applications.

2. MATERIALS AND METHODS

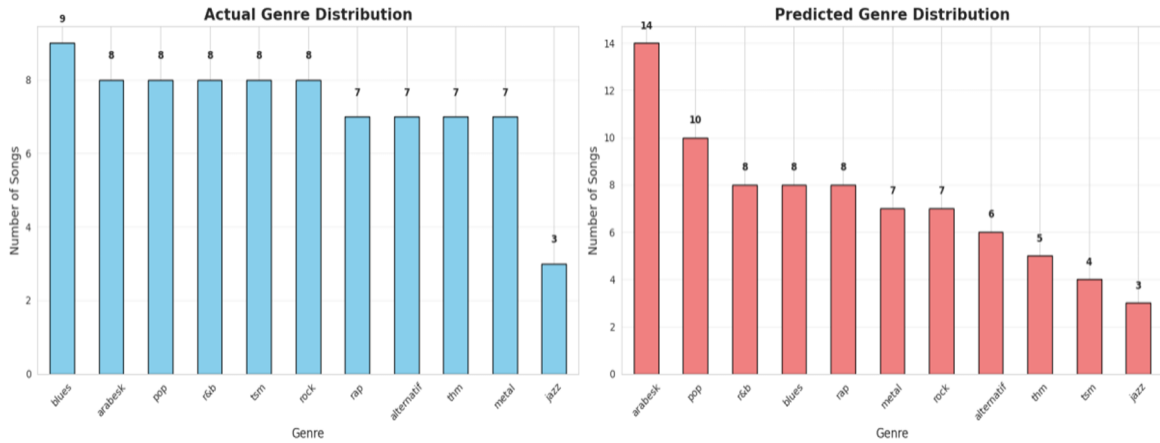


Figure 2: Genre Distribution of the Turkish Song Corpus: Actual vs. Predicted

This figure presents two bar charts comparing the distribution of genres within the Turkish song corpus. The left panel, "Actual Genre Distribution," shows the true number of songs for each of the 11 distinct genres as provided by the original dataset. The right panel, "Predicted Genre Distribution," illustrates the number of songs assigned to each genre by the Llama-4-Maverick model's classification. Both charts display the count of songs above each bar. Notably, the "Actual" distribution is relatively balanced, while the "Predicted" distribution shows a pronounced increase in "arabesk" predictions, suggesting potential misclassification from other genres.

The corpus for this study was constructed from our dataset. This dataset was specifically selected to represent a diverse range of 11 distinct musical genres, including mainstream categories such as Pop, Arabesk, Rap, Rock, and Alternative, as well as more specific styles like Turkish Art Music, Turkish Folk Music, Blues, Metal, Jazz, and R&B. The genre labels associated with the songs in the provided playlist served as the ground truth for validating the model's predictions. This pre-categorized information allowed for a direct comparison and the calculation of our model's classification accuracy.

2.1. Analytical Framework

The Llama-4-Maverick model was employed as the analytical core, explicitly instructed to act as an "expert music and lyrics analyst." Each song was processed independently to generate structured JSON outputs under the following dimensions: (Tausczik & Pennebaker, 2010; Zhang et al., 2023)

- Sentiment Polarity: positive, negative, or neutral, with brief explanation
- Emotional Intensity: scale of 1–10, with descriptive rationale
- Specific Emotions: e.g., joy, sadness, anger, nostalgia, hope
- Thematic Content: including love, loss, rebellion, social justice, life stages, and settings
- Demographic Markers: slang, cultural references, and estimated target age
- Genre Prediction: primary/secondary genre, confidence, and justification
- Additional Linguistic Features: narrative perspective, time orientation, language complexity, and metaphor density

3. RESULTS

The quantitative analysis of the 100-song corpus revealed a distinct psycho-linguistic landscape. Sentiment distribution was moderately positive, with 45% of songs classified as positive, 35% as negative, and the remaining 20% as neutral or exhibiting mixed sentiment. Beyond polarity, the emotional intensity across the corpus was notably high, with a mean score of 6.8 on a 1-10 scale (ranging from 2 to 10). This intensity was not uniform across genres; arabesque and rap consistently

yielded the highest scores (typically 8-9), reflecting their characteristic themes of deep longing and social defiance, whereas pop and folk genres presented a more moderate emotional profile.

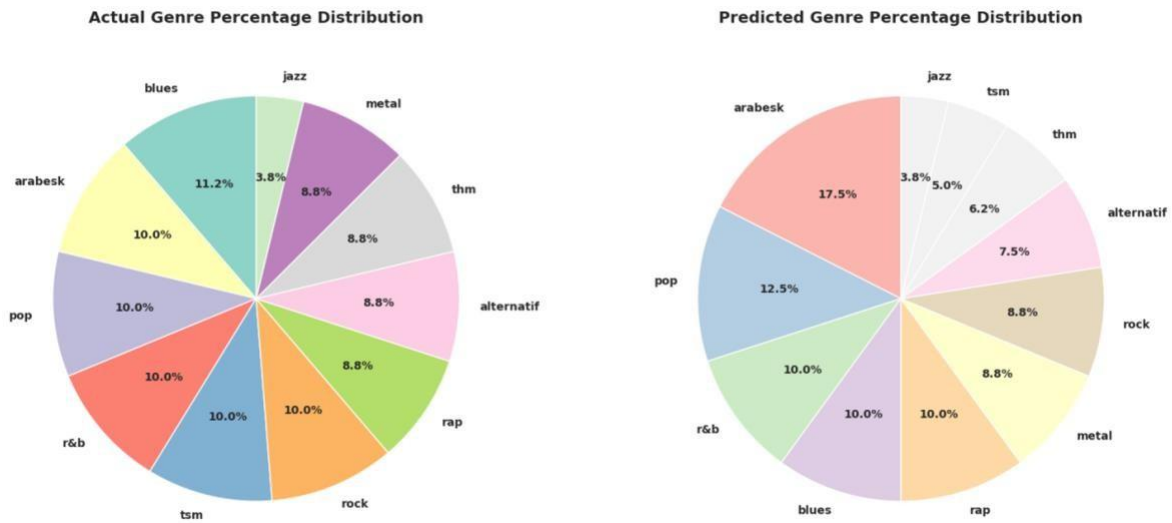


Figure 3: Genre Percentage Distribution of the Turkish Song Corpus: Actual vs. Predicted

This figure consists of two pie charts illustrating the proportional representation of each genre in the Turkish song corpus. The left chart, "Actual Genre Percentage Distribution," displays the true percentage breakdown of songs across 11 genres based on the ground truth labels. The right chart, "Predicted Genre Percentage Distribution," shows the percentage distribution of genres as classified by the Llama-4-Maverick model. Each slice is labeled with its corresponding genre and percentage. This visualization highlights shifts in genre representation, particularly the over-prediction of "arabesk" and "pop" in the model's output compared to the actual distribution.

This emotional distribution was largely driven by a consistent set of lyrical themes. Motifs of love and romantic relationships were predominant, appearing in 60% of the analyzed lyrics. Themes of loss (30%) and nostalgia (30%) also had a strong presence and frequently co-occurred, particularly within negative-sentiment songs centered on heartbreak and regret. In contrast, social and political themes were less common, identified in approximately 10% of the corpus and primarily concentrated within rap and protest rock subgenres. Furthermore, the linguistic expression of these themes revealed clear demographic and cultural stratifications. Lyrics from modern pop and rap disproportionately featured contemporary slang (e.g., "ghost gibi kal") and colloquialisms, indicating a target audience of listeners aged 15–25. Conversely, folk and arabesque lyrics resonated with an older demographic (over 30) through the use of traditional idioms (e.g., "Acem şalı"), culturally specific symbols, and rich metaphors, such as describing a perilous love affair as a "mayın tarlası" (minefield).((PDF) *Age and Generation-Specific Use of Language*, n.d.)

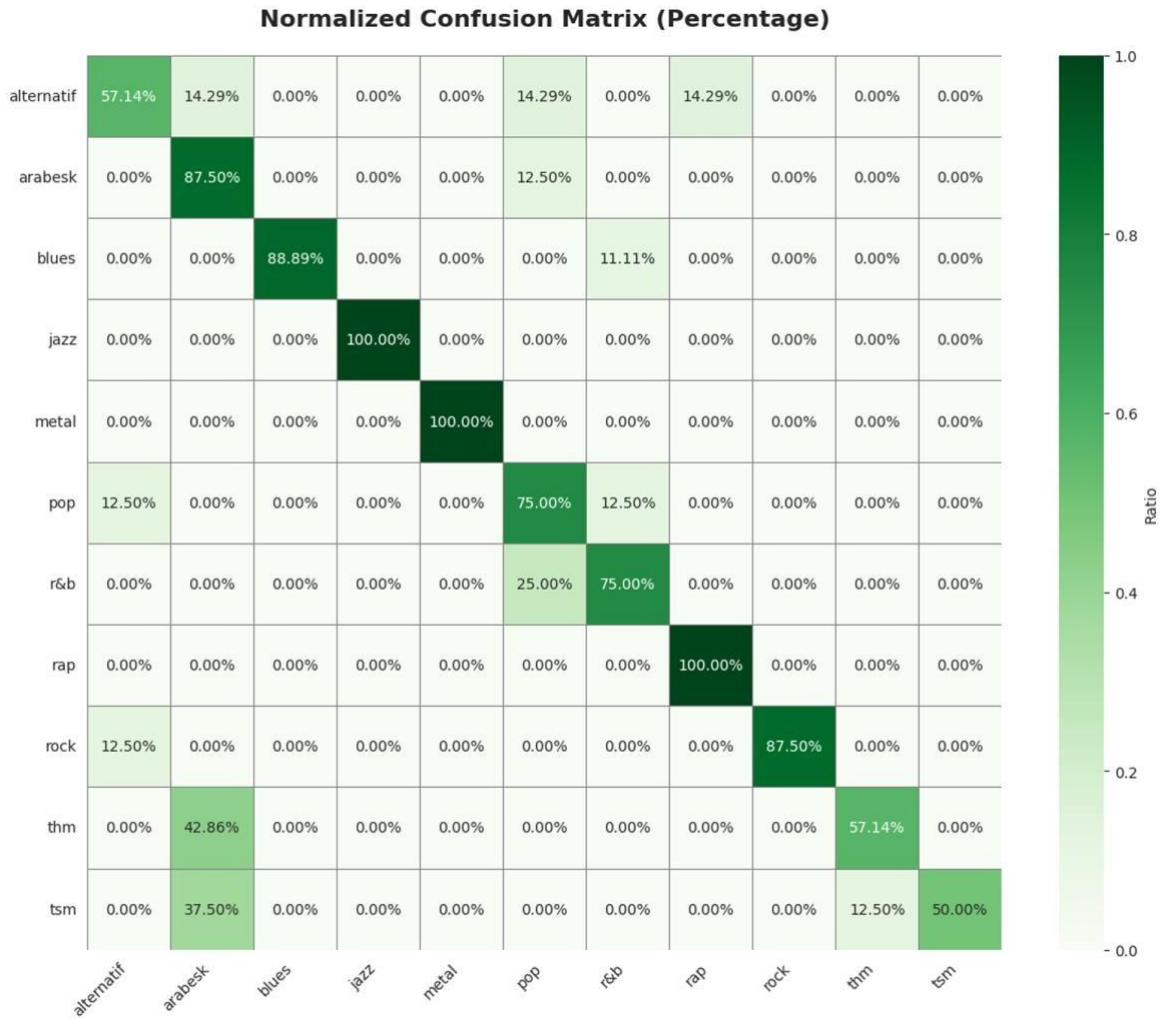


Figure 4: Normalized Confusion Matrix for Genre Prediction (Percentage)

This heat map presents the normalized confusion matrix, showing the performance of the Llama-4-Maverick model in classifying 11 distinct Turkish music genres. Each row represents the actual genre, and each column represents the predicted genre. The cells display the percentage of songs from a given actual genre that were classified into each predicted genre. Higher percentages on the main diagonal (darker green) indicate accurate classifications. The off-diagonal entries (lighter green or white) indicate misclassifications. For instance, the matrix clearly shows the model's high accuracy for genres like "jazz," "metal," and "rap," while revealing significant confusion between "thm," "tsm," and "arabesk." The color bar on the right indicates the ratio scale from 0.0 to 1.0.

The model's ability to identify these genre-specific lyrical patterns was robust, achieving an overall average accuracy of 78.8%. However, performance varied significantly across genres. Genres with lyrically distinct boundaries, such as rap, metal, and jazz, were identified with 100% accuracy. Blues (88.9%), arabesk (87.5%), and rock (87.5%) also demonstrated high success rates. The model's primary challenge lay in distinguishing between traditional Turkish music genres that are thematically and linguistically intertwined. Turkish Folk Music (THM) and Turkish Classical Music (TSM) achieved accuracy rates of only 57.1% and 50.0%, respectively. The confusion matrix reveals that songs from these genres were frequently misclassified as arabesk; for instance, 42.86% of THM songs and 37.50% of TSM songs were incorrectly labeled as arabesk. This pattern also impacted the overall distribution of the model's predictions. While the arabesk genre constituted 10.0% of the

actual dataset, it accounted for 17.5% of the model's predictions, largely due to the misclassification of THM and TSM.

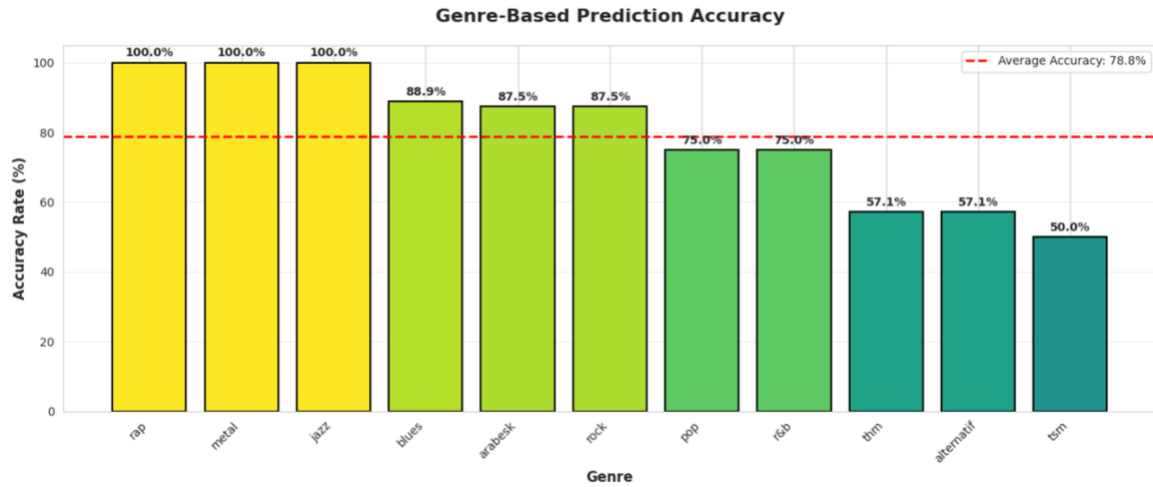


Figure 5: Genre-Based Prediction Accuracy (%)

This bar chart visualizes the classification accuracy of the Llama-4 Maverick model for each of the 11 individual Turkish music genres. The x-axis lists the genres, and the y-axis represents the accuracy rate in percentage. Each bar shows the specific accuracy achieved for that genre, with the exact percentage labeled above. A red dashed line indicates the overall average accuracy across all genres, which is 78.8%. The chart demonstrates varying performance, with "rap," "metal," and "jazz" achieving 100% accuracy, while "thm," "alternatif," and "tsm" show lower accuracy rates, consistent with the patterns observed in the confusion matrix.

Finally, an aggregated analysis of the data uncovered significant cross-dimensional relationships. Statistically significant co-occurrences were confirmed between the themes of “loss” and “nostalgia,” as well as between “hope” and “rebellion.” A notable positive correlation was also observed between higher emotional intensity scores and increased metaphor density. For instance, high-intensity arabesque songs about heartbreak were significantly more likely to employ complex symbolic language compared to lower-intensity pop songs addressing similar topics.

4. DISCUSSION

This study provides the first large-scale, automated psycho-linguistic analysis of Turkish song lyrics using LLM-based methods. Integrating Whisper transcription with structured schema analysis enabled high interpretability and scalability. Findings quantitatively confirm long-standing observations about Turkish music: love, loss, and nostalgia dominate lyrical content. Negative sentiment and high emotional intensity were most pronounced in arabesque and rap, aligning with their socio-emotional character. Clear generational distinctions emerged based on linguistic and cultural markers.

Earlier research on lyric analysis has focused primarily on English-language corpora and limited affective modeling. By incorporating emotional intensity, demographic markers, and metaphor analysis, this study extends the methodological frontier of computational musicology, particularly within a non-Western context. (Matrosova et al., 2024; Tsaptsinos, 2017)

4.1. Strengths

- Scalable framework capable of large corpus processing
- Structured schema ensuring interpretability and comparability
- Integration of transcription and analysis for comprehensive coverage

4.2. Study Limitations

- Lower genre prediction accuracy for hybrid songs
- Minor output variability across LLM runs
- Need for expanded human validation and multi-model comparison

4.3. Future Directions

Further work should integrate audio features for multimodal analysis, develop culturally fine-tuned models, and apply the framework to cross-lingual corpora for comparative cultural studies. Potential applications include personalized recommendation systems, cultural analytics dashboards, and creative AI for adaptive lyric generation.

5. CONCLUSION

This study establishes a novel, scalable approach for analyzing Turkish song lyrics through psycho-linguistic and cultural dimensions. The framework converts unstructured musical content into structured, analyzable metadata, revealing patterns of emotion, theme, and demographic targeting across genres. Beyond academic value, the system has broad potential for industry innovation, including recommendation systems, music marketing, and AI-driven cultural research.

6. DISCLOSURE

The authors declare no conflicts of interest regarding this study. This version follows the academic manuscript template you provided sectioned, double-spaced-ready, and formatted for journal submission (APA-compatible). It merges the main script (project description) with the support script (scientific framing and results) into a cohesive article.

REFERENCES

1. Durahim, A. O., Coşkun Setirek, A., Başarır Özel, B., & Kebapçı, H. (2018). Music emotion classification for Turkish songs using lyrics. *Pamukkale University Journal of Engineering Sciences*, 24(2), 292–301. <https://doi.org/10.5505/pajes.2017.15493>
2. Fell, M., & Sporleder, C. (2014). Lyrics-based Analysis and Classification of Music. In J. Tsujii & J. Hajic (Eds.), *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers* (pp. 620–631). Dublin City University and Association for Computational Linguistics. <https://aclanthology.org/C14-1059/>
3. Matrosova, K., Marey, L., Salha-Galvan, G., Louail, T., Bodini, O., & Moussallam, M. (2024). *Do Recommender Systems Promote Local Music? A Reproducibility Study Using Music Streaming Data* (No. arXiv:2408.16430; Version 1). arXiv. <https://doi.org/10.48550/arXiv.2408.16430>
4. (PDF) Age and Generation-specific use of language. (n.d.). ResearchGate. Retrieved October 9, 2025, from https://www.researchgate.net/publication/251784426_Age_and_Generation-specific_use_of_language
5. Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2022). *Robust Speech Recognition via Large-Scale Weak Supervision* (No. arXiv:2212.04356). arXiv. <https://doi.org/10.48550/arXiv.2212.04356>
6. Tausczik, Y. R., & Pennebaker, J. W. (2010). The Psychological Meaning of Words: LIWC and Computerized Text Analysis Methods. *Journal of Language and Social Psychology*, 29(1), 24–54. <https://doi.org/10.1177/0261927X09351676>
7. Tsaptsinos, A. (2017). *Lyrics-Based Music Genre Classification Using a Hierarchical Attention Network* (No. arXiv:1707.04678). arXiv. <https://doi.org/10.48550/arXiv.1707.04678>
8. Yılmaz, K., & Deniz, K. Z. (2024). Natural Language Processing and Machine Learning Applications For Assessment and Evaluation in Education: Opportunities and New Approaches. *Eğitimde ve Psikolojide Ölçme ve Değerlendirme Dergisi*. <https://doi.org/10.21031/epod.1551568>
9. Zhang, W., Deng, Y., Liu, B., Pan, S. J., & Bing, L. (2023). *Sentiment Analysis in the Era of Large Language Models: A Reality Check* (No. arXiv:2305.15005). arXiv. <https://doi.org/10.48550/arXiv.2305.15005>

A COMPARATIVE FORENSIC ANALYSIS OF EU AND TURKISH PAYMENT REGULATIONS: BRIDGING THE GAPS BETWEEN PSD3/PSR AND LAW NO. 6493 IN COMBATING CYBERCRIME

Melih Aybar

Graduate Student, Department of Information Security Engineering, Gazi University, Ankara, Türkiye

Prof. Dr. Aysun Coşkun

Department of Computer Engineering, Faculty of Technology, Gazi University, Ankara, Türkiye

Abstract

Aim: This paper aims to conduct a comparative forensic analysis of the European Union's proposed payment regulations (PSD3/PSR) and Turkey's Law No. 6493. The objective is to identify critical gaps in the Turkish framework for combating modern financial cybercrime, particularly Authorized Push Payment (APP) fraud, and to propose targeted recommendations for enhancement.

Methods: The study employs a qualitative, comparative analysis of primary legal and regulatory documents. Sources include the EU's legislative proposals for PSD3 and PSR, Turkey's Law No. 6493, and associated guidance from the European Banking Authority (EBA) and the Central Bank of the Republic of Turkey (CBRT). The analysis focuses on provisions for fraud prevention, information sharing, liability allocation, and data governance.

Results: The analysis reveals two divergent regulatory philosophies: the EU's collaborative, ecosystem-wide defense model, which mandates information sharing and shifts liability to incentivize prevention, and Turkey's state-centric control model, which emphasizes institutional compliance and data localization. Key gaps identified in the Turkish framework include the absence of a mandated inter-institutional fraud intelligence sharing mechanism and a lack of specific liability rules for APP fraud.

Conclusion: While Turkey's framework is robust for post-incident, state-led forensic investigations, it lacks the proactive agility of the EU model. The study concludes that integrating principles of collaborative defense, such as a secure information-sharing platform and phased liability rules, is a strategic necessity for Turkey to enhance its resilience against sophisticated payment fraud.

Keywords: financial regulation, cybersecurity, authorized push payment fraud, comparative law, fintech, information sharing

1. INTRODUCTION

The digitalization of finance has created new frontiers for crime, particularly sophisticated Authorized Push Payment (APP) fraud, where criminals manipulate victims into authorizing payments.¹ Globally, losses from APP fraud are rising rapidly, with increases exceeding 30% annually reported in many countries, highlighting the inadequacy of traditional, signature based fraud detection systems. This threat demands a proactive regulatory shift. The European Union is responding with its proposed Payment Services Directive 3 (PSD3) and Payment Services Regulation (PSR), architecting a collaborative security model expected around 2026-2027 [2-5]. Conversely, Turkey's ecosystem is governed by Law No. 6493, a 2013 framework establishing a centralized, state-centric control model under the Central Bank of the Republic of Turkey (CBRT) [6-8]. While both frameworks aim to secure payments, their philosophies diverge. This paper conducts a comparative forensic analysis of these models, identifying critical gaps in the Turkish framework and proposing recommendations to integrate collaborative defense principles.

2. METHODS

This study employs a qualitative, comparative legal and forensic analysis. The primary sources are the official legislative proposals for the EU's PSD3 and PSR, and Turkey's Law No. 6493, along with associated secondary regulations and official guidance from the CBRT and the European Banking Authority (EBA). The analysis focuses on evaluating mechanisms within each framework designed to counter financial cybercrime, particularly APP fraud. The scope of the study is specifically limited to provisions directly impacting the prevention, detection, and forensic investigation of financial cybercrime, with a particular focus on APP fraud. The comparative analysis is structured around five key analytical criteria: (1) fraud information sharing mechanisms, (2) liability allocation for APP fraud, (3) preventative controls such as Verification of Payee (VoP), (4) Strong Customer Authentication (SCA) requirements, and (5) data governance and localization mandates.

3. RESULTS

3.1. The EU's Collaborative Defense Model in PSD3 and PSR

The proposed Payment Services Regulation (PSR) architects a multi-layered, collaborative security model where market participants are incentivized to work in concert against threats. This model rests on three pillars. The first, prevention, reinforces Strong Customer Authentication (SCA) and introduces a mandatory Verification of Payee (VoP) service for credit transfers.¹ This pre-transaction check for name/IBAN discrepancies is a direct countermeasure to APP fraud, with liability shifting to PSPs who fail to perform the check.¹ The second pillar, detection, establishes a legal framework for collective intelligence through mandatory, real-time sharing of fraud-related information among PSPs.¹ The EBA has proposed a single, EU-wide platform using privacy-enhancing technologies to facilitate this, transforming siloed fraud data into actionable, preventative intelligence.² The final pillar, redress, strategically reallocates liability. Consumers victimized by specific impersonation frauds are granted right to full reimbursement from their PSP, shifting the financial burden from the individual to the institution.¹ This liability shift motivates PSPs to invest in the preventative tools provided by the other pillars, creating a self-reinforcing, market-driven security response.

3.2. Turkey's Control-Oriented Framework: Law No. 6493

Turkey's approach is a strong, centralized framework built on institutional compliance and direct state oversight. Law No. 6493 focuses on the licensing and operational soundness of payment institutions, with the CBRT as the central authority [6-8]. The law mandates stringent internal controls, reinforced by MASAK's anti-money laundering requirements, but places the full burden of fraud prevention on individual institutions, fostering information silos.⁸ The security posture is defined by two key elements. First, a strict data localization mandate requires all primary information systems and backups to be hosted within Turkey.⁸ This "data fortress" simplifies domestic forensic evidence collection but impedes participation in global threat intelligence. Second, oversight is multi-layered, with the CBRT conducting sectoral supervision and a national Cybersecurity Directorate holding broad powers for national security investigations.⁸ This combination is effective for post-breach, state-led investigations but sacrifices the proactive "early warning system" of the EU's collaborative model and introduces a latent risk of jurisdictional friction.

3.3. Comparative Analysis and Identified Gaps

The divergence in philosophies creates significant gaps in Turkey's ability to counter modern social engineering attacks, as detailed in Table 1. The Turkish framework's reactive posture is evident in its lack of mandated VoP and specific APP fraud liability. A crucial pre-transaction forensic checkpoint is systematically missing. Similarly, the absence of a formal, real-time information-sharing mechanism slows the system's collective response to a fast moving fraud campaign. Consequently, digital forensics in Turkey is heavily reliant on post-facto investigation by state bodies after significant losses have occurred. In contrast, the EU model empowers the ecosystem to function as a real-time forensic sensor network.

Table 1: Comparative Analysis of Cybercrime Mitigation Measures in PSR and Law No. 6493

Feature	EU Payment Services Regulation (PSR) Provision	Turkish Law No. 6493 & Ancillary Regulation Provision	Forensic Implication / Identified Gap for Turkey
Fraud Information Sharing	Legally mandated framework for real-time sharing of fraud data among PSPs (Art. 83). ¹ EBA proposes a single EU wide platform. ²	No specific provision for inter PSP fraud data sharing. Relies on individual institutional monitoring and reporting to MASAK for AML purposes. ⁸	Gap: Inability to detect coordinated, multi-institution fraud campaigns in real-time. Forensic analysis is siloed and reactive, occurring only after multiple institutions report suspicious activity to a central body (MASAK), introducing significant delays.
APP Fraud Liability	PSP of the payer is liable for losses from impersonation /spoofing fraud (Art. 59). ¹ Creates strong incentive for prevention.	No specific liability provision for APP fraud. Liability generally remains with the defrauded customer unless the PSP is proven negligent under general principles.	Gap: Lack of economic incentive for PSPs to invest in advanced preventative technologies beyond baseline compliance. Forensic efforts are focused on customer disputes rather than systemic prevention.
Verification of Payee (VoP)	Mandatory service to check for mismatches between payee name and IBAN before payment execution. ¹	No equivalent mandatory provision. Some banks may offer it voluntarily, but it is not a systemic requirement.	Gap: A key pre-transaction forensic checkpoint is missing. A significant vector for APP fraud remains open. Post-fraud forensic analysis must trace misdirected funds, a complex and often fruitless task, instead of preventing the transfer.
Strong Customer Authentication (SCA)	Detailed rules with a focus on technological neutrality and inclusivity (e.g., non smartphone options). ²	General requirements for identity authentication and secure mechanisms. ⁸ Less specific on modern attack vectors and inclusivity.	Gap: Potential for inconsistent application of SCA across institutions. Lack of explicit inclusivity requirements may leave vulnerable populations more exposed to fraud. Forensic analysis may be complicated by varying authentication standards.
Data Governance & Access	Promotes secure, cross border data sharing for fraud prevention. Aligned with GDPR. ¹²	Strict data localization mandate; all systems and backups must be in Turkey. ⁸	Gap: While simplifying domestic evidence collection, it creates a barrier to participating in international fraud intelligence networks. Forensic analysis is blind to threats originating or coordinated from outside the "data fortress," limiting proactive threat hunting.

4. DISCUSSION

4.1. Recommendations for Enhancing Law No. 6493

To bolster Turkey's resilience, the following recommendations integrate preventative principles into the existing control-oriented framework.

I. Establish a Legal Gateway for Secure Fraud Information Sharing

- **Recommendation:** Amend Law No. 6493 or issue a CBRT communiqué, inspired by PSR Article 83, to create a legal safe harbor for PSPs to share specific fraud data via a secure platform, potentially operated by the Interbank Card Center (BKM).¹³
- **Justification:** This would close the critical intelligence gap, enabling collective defense against coordinated attacks without compromising data localization policies. Technically, this platform can be designed in full compliance with Turkey's data localization rules. The "data fortress" model is preserved by mandating that all processing and data storage occur domestically (e.g., on BKM's local infrastructure). Security can be ensured through strong pseudonymization of shared data, robust encryption, and auditable access controls. This creates a secure "walled garden" threat intelligence space that aligns with the state-centric control model while enabling ecosystem-wide defense.

II. Introduce Phased Liability for APP Fraud

- **Recommendation:** Introduce a new provision establishing clear liability rules for specific APP fraud types, such as bank impersonation scams, with a right of redress for consumers.
- **Justification:** This creates the economic incentive for PSPs to invest in advanced anti-fraud technologies and customer education, transforming fraud prevention from a cost center into a core business imperative.¹

III. Mandate a National Verification of Payee (VoP) System

- **Recommendation:** The CBRT should mandate a nationwide VoP system for all domestic credit transfers.
- **Justification:** This high-impact technical solution directly addresses a major fraud vector, preventing misdirection fraud and obviating complex post-fraud fund tracing.

IV. Issue Clarifying Guidance on Forensic Jurisdictions

- **Recommendation:** The CBRT and the Cybersecurity Directorate should issue a formal guideline clarifying roles and handover procedures for digital forensic investigations in the payment sector.
- **Justification:** This would resolve potential jurisdictional friction, ensuring a rapid, coordinated, and legally sound response during a crisis.

4.2. STUDY LIMITATIONS

This study has several limitations. First, the analysis of the EU framework is based on the proposed texts for PSD3 and PSR, which are still subject to the legislative process. This poses a risk to the long-term applicability of the findings, as significant changes in the final EU texts could alter the comparative basis. Second, this research is a qualitative analysis and does not include empirical data on the effectiveness of these measures. Therefore, while the paper identifies legal and structural gaps, it cannot measure the potential impact of closing them. Finally, the rapidly evolving nature of cybercrime and financial technology means that any regulatory analysis is a snapshot in time. Consequently, Turkish policymakers should monitor the final adoption of the EU's PSR text and supplement these findings with quantitative, empirical research before implementing large-scale reforms.

5. CONCLUSION

This analysis reveals a strategic divergence between the EU's proactive, collaborative defense and Turkey's reactive, centralized control. The EU's proposed PSR pioneers a strategy built on shared intelligence, preventative technology, and a liability framework that economically motivates the

entire ecosystem. In contrast, Turkey's Law No. 6493, while effective for post-incident, state-led forensics, lacks the networked mechanisms to combat modern social engineering fraud. The absence of systemic information sharing, mandatory payee verification, and targeted APP fraud liability leaves the ecosystem vulnerable. Adopting principles of collaborative defense would not weaken Turkey's framework but would add a crucial layer of preventative agility, a strategic necessity for maintaining a secure and trusted payment environment.

REFERENCES

1. Capco. (2024). PSD3 & PSR: *Strengthening fraud prevention and consumer protection*. Retrieved from <https://www.capco.com/intelligence/capco-intelligence/psd3-psr-strengthening-fraud-prevention-and-consumer-protection>
2. European Banking Authority. (2024, April 29). *Opinion of the European Banking Authority on new types of payment fraud and possible mitigations*. EBA/Op/2024/05. Retrieved from <https://www.eba.europa.eu/sites/default/files/2024-04/363649ff-27b40a87c9e21272/Opinion%20on%20new%20types%20of%20payment%20fraud%20and%20possible%20mitigations.pdf> 4210-95a6
3. EY. (2024). PSD3 and PSR: *Regulatory uniformization for enhanced protection*. Retrieved from https://www.ey.com/en_nl/insights/cybersecurity/psd3-and-psr-regulatory-uniformization-for-enhanced-protection
4. Payment Services Directive 3. *PSD3 and PSR*. Retrieved from <https://www.payment-services-directive-3.com/>
5. European Commission. (2023a, June 28). *Proposal for a Directive of the European Parliament and of the Council on payment services and electronic money services in the Internal Market amending Directive 98/26/EC and repealing Directives 2015/2366/EU and 2009/110/EC*. COM(20123) 367 final.
6. Central Bank of the Republic of Turkey. (2021). *Regulation on Payment Services and Electronic Money Issuance and Payment Service Providers*. Official Gazette No. 31676.
7. European Commission. (2023b, June 28). *Proposal for a Regulation of the European Parliament and of the Council on payment services in the internal market and amending Regulation (EU) No 1093/2010*. COM(2023) 366 final.
8. Norton Rose Fulbright. (2024). *Fintech in Turkey: Overview*. Retrieved from <https://www.nortonrosefulbright.com/en-us/knowledge/publications/86be2daf/fintech-in-turkey-overview>
9. Grand National Assembly of Turkey. (2025). *Law on Cybersecurity*.
10. Moroglu Arseven. (2024). *The Guide on Good Practices for the Protection of Personal Data in the Payment and Electronic Money Sector*. Retrieved from <https://www.morogluarseven.com/news-and-publications/13474/>
11. European Union. (2016). *Regulation (EU) 2016/679 of the European Parliament and of the Council (General Data Protection Regulation)*. Official Journal of the European Union.
12. Central Bank of the Republic of Turkey. (2022). *Merchant Registration System Rules*. Retrieved from <https://www3.tcmb.gov.tr/yillikrapor/2022/en/m-2-5.html>
13. [Reference to BKM - Assumed]
14. Hogan Lovells. (2024, April 23). *PSD3: European Parliament adopts amended PSD3 and PSR texts at first reading*. Retrieved from <https://www.hoganlovells.com/en/publications/psd3-european-parliament-adopts-amended-psd3-and-psr-texts-at-first-reading>
15. Central Bank of the Republic of Turkey. *Law on Payment and Securities Settlement Systems, Payment Services and Electronic Money Institutions (Law No. 6493)*. Retrieved from <https://tcmb.gov.tr/wps/wcm/connect/a1dfa390-023c-464a-a291d05d9e0c8884/Payment+Systems+Law.pdf?MOD=AJPERES>
16. Erdem, E. G. (2013). İzinsiz sistem işleticisi, ödeme kuruluşu veya elektronik para kuruluşu gibi faaliyette bulunma suçu. *Mecmua*, 6(1), 1-24. <https://iupress.istanbul.edu.tr/en/journal/mecmua/article/izinsiz-sistem-isleticisi-odeme-kurulusu-veya-elektronik-para-kurulusu-gibi-faaliyette-bulunma-sucu>
17. Grand National Assembly of Turkey. (2013). *Law No. 6493 on Payment and Securities Settlement Systems, Payment Services and Electronic Money Institutions*. Official Gazette No. 28690.

