Enhancing Cybersecurity in Critical Infrastructure with LLM-Assisted Explainable IoT Systems

†Ashutosh Ghimire, †Ghazal Ghajari, †Karma Gurung, §Love K. Sah, †*Fathi Amsaad † Department of Computer Science and Engineering, Wright State University § Department of Electrical and Computer Engineering, Western New England University Email: †{ashutosh.ghimire, ghajari.2, gurung.14, fathi.amsaad}@wright.edu, § love.sah@wne.edu

Abstract—Ensuring the security of critical infrastructure has become increasingly vital with the proliferation of Internet of Things (IoT) systems. However, the heterogeneous nature of IoT data and the lack of human-comprehensible insights from anomaly detection models remain significant challenges. This paper presents a hybrid framework that combines numerical anomaly detection using Autoencoders with Large Language Models (LLMs) for enhanced preprocessing and interpretability. Two preprocessing approaches are implemented: a traditional method utilizing Principal Component Analysis (PCA) to reduce dimensionality and an LLM-assisted method where GPT-4 dynamically recommends feature selection, transformation, and encoding strategies.

Experimental results on the KDDCup99 10% corrected dataset demonstrate that the LLM-assisted preprocessing pipeline significantly improves anomaly detection performance. The macroaverage F1 score increased from 0.49 in the traditional PCA-based approach to 0.98 with LLM-driven insights. Additionally, the LLM generates natural language explanations for detected anomalies, providing contextual insights into their causes and implications. This framework highlights the synergy between numerical AI models and LLMs, delivering an accurate, interpretable, and efficient solution for IoT cybersecurity in critical infrastructure.

Index Terms-IOT, LLM, Autoencoder, Explainable AI

I. INTRODUCTION

Critical infrastructure, including power grids, healthcare systems, and transportation networks, relies heavily on Internet of Things (IoT) devices for automation, monitoring, and operational efficiency. However, this integration introduces substantial cybersecurity risks, as IoT systems generate heterogeneous, high-volume data and operate on resource-constrained devices [1]. Traditional security mechanisms often fail to manage these complexities, leaving critical systems vulnerable to threats such as Denial-of-Service (DoS) attacks, network intrusions, and data breaches [2].

Anomaly detection methods for IoT systems predominantly utilize machine learning techniques, including clustering, Support Vector Machines (SVMs), and deep learning models like Autoencoders. While effective at identifying anomalies, these approaches face three significant challenges: (1) operating as opaque "black box" models, which lack interpretability and fail to provide explanations for their decisions, (2) neglecting robust preprocessing despite the noisy and redundant nature of IoT data, and (3) underutilizing valuable unstructured textual

data, such as logs or event descriptions, which hold critical contextual insights [3].

To address these challenges, a hybrid framework is introduced that integrates Large Language Models (LLMs), such as GPT-4, with Autoencoder-based anomaly detection for IoT cybersecurity. The framework enhances anomaly detection performance through optimized, LLM-assisted preprocessing and improves interpretability by generating human-readable explanations for detected anomalies [4]. Specifically, two preprocessing pipelines are employed: a traditional PCA-based approach for dimensionality reduction and an LLM-driven approach, where GPT-4 recommends feature selection, transformation, and encoding strategies.

The contributions of this work are twofold. First, the proposed LLM-assisted preprocessing pipeline intelligently handles noisy, redundant, and categorical features in IoT data, leading to a significant improvement in anomaly detection performance, as demonstrated by an increase in the macroaverage F1 score from 0.49 (PCA-based preprocessing) to 0.98 (LLM-assisted preprocessing). Second, GPT-4 enhances the explainability of the framework by producing natural language descriptions of anomalies, enabling cybersecurity analysts to understand and trust the model's outputs. By addressing these limitations, the proposed framework demonstrates the synergy between numerical AI models and LLMs, achieving robust, interpretable, and efficient anomaly detection. This work not only improves detection accuracy but also fosters trust in AIdriven IoT cybersecurity systems through intelligent preprocessing and explainable results.

II. RELATED WORKS

The integration of artificial intelligence (AI) and machine learning (ML) techniques into IoT cybersecurity has been a topic of extensive research. Given the increasing complexity of IoT networks and their critical role in infrastructure systems, various approaches have been explored to detect, mitigate, and explain cyber threats [5]. This section reviews existing methods, highlights their limitations, and identifies the gap that this research aims to address.

A. AI and ML-Based Anomaly Detection in IoT Cybersecurity

In recent years, traditional and deep learning techniques have been widely applied to anomaly detection in IoT systems. Supervised learning methods such as Support Vector Machines

^{*} Corresponding author

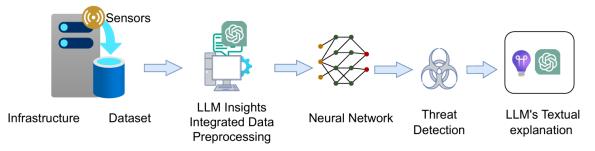


Fig. 1: Proposed framework for LLM-assisted anomaly detection in IoT systems.

(SVMs), Decision Trees, and Random Forests have been employed to classify IoT network traffic into normal and attack categories [6]. For instance, datasets like KDDCup99 and NSL-KDD have been used to train classification models for intrusion detection. While these models are effective when labeled data is available, their reliance on labeled training samples limits their applicability to real-world IoT systems, where data labeling is time-consuming and expensive[7].

To overcome the limitations of supervised learning, unsupervised learning methods such as clustering (e.g., k-Means, DBSCAN) and anomaly detection techniques have been explored. These methods identify patterns or deviations in unlabeled data that signify anomalies. Among these approaches, Autoencoders have shown significant promise. Autoencoders are deep neural networks that learn to reconstruct input data and flag instances with high reconstruction error as anomalies. Researchers have applied Autoencoders to IoT traffic analysis, time-series data, and sensor readings to detect abnormal behaviors indicative of cyberattacks [8]. However, despite their effectiveness in identifying anomalies, Autoencoders suffer from a lack of interpretability. They function as black-box models, providing no meaningful explanation for why a particular instance was flagged as an anomaly.

B. The Role of Natural Language Processing and Large Language Models in Cybersecurity

Natural Language Processing (NLP) has emerged as a promising tool for addressing the limitations of traditional numerical AI models in cybersecurity. While traditional anomaly detection focuses primarily on structured, numerical IoT data, real-world IoT systems generate substantial amounts of unstructured textual data, including logs, alerts, and error messages [9]. NLP techniques can extract meaningful insights from this textual data, enhancing the context and interpretability of anomaly detection systems.

Recently, Large Language Models (LLMs), such as OpenAI's GPT-4, have demonstrated unprecedented capabilities in understanding and generating human-like text. In the context of cybersecurity, LLMs have been applied to various tasks, including log analysis, where they summarize system logs, detect anomalies in textual alerts, and identify patterns in error messages [10]. Additionally, LLMs are used for incident reporting, generating natural language explanations of detected incidents to assist human operators, and for threat intelligence,

which involves extracting knowledge from security reports and describing potential attack vectors [11].

Despite significant advancements, Large Language Models (LLMs) have not been effectively integrated with numerical anomaly detection models. Most existing approaches treat numerical and textual analyses as separate tasks, failing to combine their strengths to create a holistic cybersecurity solution. IoT data is inherently noisy, redundant, and heterogeneous [12]. Existing numerical AI models often depend on manual or simplistic preprocessing methods, which limit their performance. While LLMs offer enhanced interpretability, they remain underutilized in preprocessing IoT data—a critical bottleneck in anomaly detection pipelines.

III. PROPOSED FRAMEWORK

This section outlines our proposed framework, which integrates numerical anomaly detection models with Large Language Models (LLMs) to address limitations in traditional IoT cybersecurity solutions. The framework consists of three main components: data preparation, anomaly detection, and natural language-based anomaly explanation. Two distinct approaches are implemented for data preprocessing: (1) PCA-based numerical preprocessing and (2) LLM-assisted preprocessing.

A. Dataset

The experiments in this study are conducted on the KDD-Cup99 10% corrected dataset, a standard benchmark for intrusion detection systems [13]. The dataset comprises numerical features (e.g., src_bytes, dst_bytes) and categorical features (e.g., protocol_type, flag). Despite its utility, the dataset poses challenges such as redundancy, noise, and feature heterogeneity, which require careful preprocessing for effective anomaly detection.

B. Data Preprocessing

To evaluate the effectiveness of LLMs in enhancing anomaly detection, two distinct preprocessing strategies are employed:

1) Traditional Preprocessing Using PCA: In the traditional preprocessing approach using Principal Component Analysis (PCA), the steps included standardizing numerical features through Min-Max scaling and one-hot encoding of categorical features such as protocol_type and flag. PCA is then applied to reduce the input feature space to 25 principal components, preserving approximately 95% of the variance. The resulting

PCA-transformed data served as input to the anomaly detection model.

2) LLM-Assisted Preprocessing: The second approach leveraged GPT-4 to automate and optimize the preprocessing pipeline. The LLM analyzed the dataset by reviewing feature statistics, such as variance and correlation, to identify redundant and low-variance features. It recommended specific feature transformations, such as binarization for sparse columns and averaging for highly correlated features, and suggested encoding techniques for categorical variables to ensure compatibility with machine learning models. By generating automated preprocessing scripts, the LLM-assisted method reduced manual effort and produced cleaner, more meaningful input data for anomaly detection.

C. Anomaly Detection Using Autoencoder

The preprocessed datasets from both methods (PCA-based and LLM-assisted) are used to train separate Autoencoder models. An Autoencoder is a type of unsupervised neural network that learns to compress and reconstruct input data. During inference, instances with high reconstruction errors are flagged as anomalies.

The Autoencoder model consists of two main components:

$$z = \sigma(W_e x + b_e), \quad z \in \mathbb{R}^m, \ m < n$$
 (1)

$$\hat{x} = \sigma(W_d z + b_d), \quad \hat{x} \in \mathbb{R}^n$$
 (2)

where $x \in \mathbb{R}^n$ represents the input data, z is the latent representation, and \hat{x} is the reconstructed input. The parameters W_e, W_d and biases b_e, b_d are learned during training. The reconstruction loss is computed as:

$$L_{\text{recon}} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{x}_i)^2$$
 (3)

Anomalies are detected by thresholding the reconstruction error.

D. LLM-Driven Anomaly Explanation

To bridge the interpretability gap in anomaly detection, GPT-4 is integrated into the framework to generate humanreadable explanations for detected anomalies. The process involves:

- 1) Extracting anomaly-related features, such as reconstruction error, *src_bytes*, and *protocol_type*.
- Structuring the data into natural language prompts and feeding them to GPT-4.
- 3) Generating textual explanations that describe the anomaly's potential cause and implications.

For example, an anomaly with unusually low traffic volume (*src_bytes*) and a rejected TCP connection may be explained as a failed network scanning attempt.

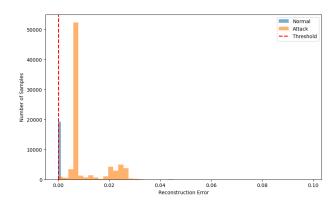


Fig. 2: Reconstruction error distribution of the LLM-integrated Autoencoder.

E. Comparison and Evaluation

The comparison and evaluation of the two preprocessing approaches were conducted using multiple metrics to assess the overall performance of the anomaly detection framework. The evaluation included the reconstruction error distributions of the Autoencoder, as well as training and validation loss curves, to analyze the model's convergence and generalization capabilities. Additionally, standard anomaly detection performance metrics, such as accuracy, false positive rate, and macroaverage precision, recall, and F1-score for both the normal and attack classes, were used to provide a comprehensive assessment. The macro-average F1-score, in particular, highlighted the balance between precision and recall for the two classes, ensuring the model's ability to detect rare anomalies without bias toward the dominant class.

Furthermore, the quality of GPT-4 explanations was evaluated in terms of their ability to generate actionable insights for detected anomalies, enhancing the interpretability of the system. The comparative results clearly demonstrated that the LLM-assisted preprocessing approach not only improved anomaly detection accuracy but also provided more meaningful and interpretable outputs, making the system more reliable for real-world applications.

IV. RESULTS AND DISCUSSION

This section presents the results of our proposed framework, comparing the performance of traditional numerical anomaly detection methods (PCA-based preprocessing) with LLM-assisted preprocessing. Additionally, we demonstrate the ability of LLMs to provide interpretable explanations for detected anomalies. The results are analyzed in terms of training/validation loss, reconstruction error distribution, and performance metrics.

A. Autoencoder Performance Analysis

To evaluate the impact of LLM-assisted preprocessing, two Autoencoder models are trained:

- **Traditional Autoencoder**: Using PCA-based preprocessing with 25 principal components.
- LLM insights Integrated Autoencoder: Using LLMdriven preprocessing, where GPT optimized data preparation by suggesting feature transformations, removal of redundancy, and encoding strategies.

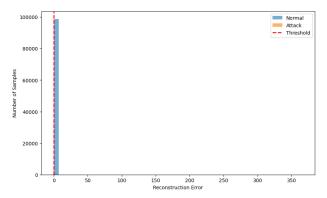


Fig. 3: Reconstruction error distribution of the traditional Autoencoder.

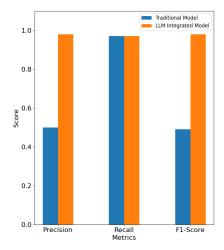


Fig. 4: Comparison of macro-average performance metrics between the traditional Autoencoder and the LLM-integrated Autoencoder.

- 1) Reconstruction Error Distribution: The reconstruction error distributions for both models are compared in Fig. 3 and Fig. 2. The LLM-integrated Autoencoder produces a tighter error distribution for normal samples and more distinct outliers for anomalies. This indicates that the LLM-assisted preprocessing enhanced the model's ability to differentiate between normal and anomalous patterns.
- 2) Performance Comparison: Fig. 4 provides a performance comparison between the traditional Autoencoder and the LLM-integrated Autoencoder. The LLM-integrated model achieves a higher accuracy and lower false positive rate due to its improved preprocessing pipeline.

B. LLM-Driven Explanation of Anomalies

The integration of GPT-4 allows for human-readable explanations of detected anomalies. This section demonstrates the ability of the LLM to analyze anomalies flagged by the Autoencoder and provide actionable insights for cybersecurity analysts.

Figs. 5 and 6 illustrate example anomalies detected by the Autoencoder, along with their GPT-generated explanations. These explanations provide contextual insights, such as potential attack vectors, failed connection attempts, or unusual traffic patterns.

The ability of GPT-4 to describe anomalies in natural language enhances the interpretability of the system. For instance,

Prompt Response: Anomaly 1

Reconstruction Error: 0.0001

Features:

- src_bytes=1.4883707192251517e-06, dst_bytes=0.0
- protocol_type_tcp=0.0, flag_REJ=0.0

Prediction: Attack

Generated Insight: The anomaly suggests minimal source bytes (src_bytes) and zero destination bytes (dst_bytes). The TCP protocol is inactive (protocol_type_tcp=0.0), and there are no rejection flags (flag_REJ=0.0). This behavior might indicate low-volume traffic anomalies or stealth reconnaissance attempts where connections do not exchange meaningful data.

Steps for Further Analysis:

- Network Logs: Verify the source IP and check for repeated connection attempts.
- Port Analysis: Determine if specific ports are targeted.
- Threshold Adjustment: Reassess detection thresholds for anomalies with near-zero payloads.

Fig. 5: Example 1: GPT-generated explanation of a detected anomaly.

anomalies with high reconstruction errors may be explained as unusual TCP connection refusals or potential port scanning activities. Such insights are crucial for cybersecurity analysts to take informed actions.

C. Discussion

The experimental results clearly demonstrate the advantages of the proposed framework. First, improved detection performance was achieved through LLM-assisted preprocessing, which refined the data input and led to more accurate anomaly detection. Second, the faster convergence of the LLM-integrated Autoencoder was evident, as it required fewer epochs to converge, as shown in the loss curves. Third, the framework provided enhanced interpretability by leveraging GPT-4 to generate actionable insights for anomalies, effectively bridging the gap between numerical outputs and human understanding. By combining numerical AI models with advanced language capabilities, the proposed framework addresses the limitations of existing anomaly detection systems and fosters trust in AI-driven cybersecurity for critical infrastructure.

V. CONCLUSION

This paper proposed a hybrid framework for anomaly detection in IoT systems, integrating Autoencoder-based models with GPT-4 for enhanced preprocessing and interpretability. Two approaches were explored: a traditional PCA-based method and an LLM-assisted pipeline where GPT-4 dynamically optimized feature selection and transformation.

Experimental results on the KDDCup99 dataset showed that the LLM-assisted preprocessing improved anomaly detection performance, leading to faster convergence, better

Prompt Response: Anomaly 2

Reconstruction Error: 0.0008 **Features:**

- src_bytes=3.9661041452220616e-07, dst_bytes=6.013033152373364e-05
- protocol_type_tcp=1.0, flag_REJ=0.0

Prediction: Attack

Generated Insight: The anomaly indicates a very small src_bytes value and slightly higher dst_bytes, while the TCP protocol is active (protocol_type_tcp=1.0). There are no rejections (flag_REJ=0.0). This could indicate port scanning or probing behavior where minimal traffic is sent to test open ports or system responses.

Steps for Further Analysis:

- Traffic Context: Analyze surrounding traffic for sequential port scans.
- Protocol Analysis: Verify if the TCP handshake completes or drops midway.
- Source IP Tracking: Look for repeated connections from the same IP.

Fig. 6: Example 2: GPT-generated explanation of a detected anomaly.

reconstruction error separation, and enhanced accuracy. Additionally, GPT-4 provided natural language explanations for detected anomalies, improving interpretability and aiding human decision-making.

In summary, the proposed framework demonstrates the potential of combining numerical AI models with LLMs to achieve accurate and explainable anomaly detection in IoT systems. Future work will focus on real-time IoT log integration and validation across additional critical infrastructure datasets.

REFERENCES

- A. Djenna, S. Harous, and D. E. Saidouni, "Internet of things meet internet of threats: New concern cyber security issues of critical cyber infrastructure," *Applied Sciences*, vol. 11, no. 10, p. 4580, 2021.
- [2] Ö. Aslan, S. S. Aktuğ, M. Ozkan-Okay, A. A. Yilmaz, and E. Akin, "A comprehensive review of cyber security vulnerabilities, threats, attacks, and solutions," *Electronics*, vol. 12, no. 6, p. 1333, 2023.
- [3] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," ACM computing surveys (CSUR), vol. 41, no. 3, pp. 1–58, 2009.
- [4] J. Su, C. Jiang, X. Jin, Y. Qiao, T. Xiao, H. Ma, R. Wei, Z. Jing, J. Xu, and J. Lin, "Large language models for forecasting and anomaly detection: A systematic literature review," arXiv preprint arXiv:2402.10350, 2024.
- [5] X. Liu, C. Qian, W. G. Hatcher, H. Xu, W. Liao, and W. Yu, "Secure internet of things (iot)-based smart-world critical infrastructures: Survey, case study and research opportunities," *IEEE Access*, vol. 7, pp. 79 523– 79 544 2019
- [6] A. Ghimire, A. N. Asiri, B. Hildebrand, and F. Amsaad, "Implementation of secure and privacy-aware ai hardware using distributed federated learning," in 2023 IEEE 16th Dallas Circuits and Systems Conference (DCAS). IEEE, 2023, pp. 1–6.
- [7] R. D. Ravipati and M. Abualkibash, "Intrusion detection system classification using different machine learning algorithms on kdd-99 and nsl-kdd datasets-a review paper," *International Journal of Computer Science & Information Technology (IJCSIT) Vol.*, vol. 11, 2019.
- [8] H. Gao, B. Qiu, R. J. D. Barroso, W. Hussain, Y. Xu, and X. Wang, "Tsmae: a novel anomaly detection approach for internet of things time series data using memory-augmented autoencoder," *IEEE Transactions* on network science and engineering, vol. 10, no. 5, pp. 2978–2990, 2022.

- [9] C. Bertero, M. Roy, C. Sauvanaud, and G. Trédan, "Experience report: Log mining using natural language processing and application to anomaly detection," in 2017 IEEE 28th International Symposium on Software Reliability Engineering (ISSRE). IEEE, 2017, pp. 351–360.
- [10] F. Hadadi, Q. Xu, D. Bianculli, and L. Briand, "Anomaly detection on unstable logs with gpt models," arXiv preprint arXiv:2406.07467, 2024.
- [11] J. O'Brien, S. Ee, and Z. Williams, "Deployment corrections: An incident response framework for frontier ai models," arXiv preprint arXiv:2310.00328, 2023.
- [12] L. Wang, "Heterogeneous data and big data analytics," Automatic Control and Information Sciences, vol. 3, no. 1, pp. 8–15, 2017.
- [13] M. Tavallaee, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the kdd cup 99 data set," in 2009 IEEE symposium on computational intelligence for security and defense applications. Ieee, 2009, pp. 1–6.