

---

**| RESEARCH ARTICLE**

## **Zero Trust Based Critical Infrastructure Cybersecurity Framework with AI-Driven Threat Detection and Secure Network Modernization**

**Md Humayun Kabir<sup>1</sup>, MD Razib<sup>2</sup>, Zakarya Jahin<sup>3</sup>, and Zakarya Jesan<sup>4</sup>**

<sup>1</sup> Westcliff University, Irvine, California, USA

<sup>2</sup> Department of MBA (Digital and Strategic Marketing), Westcliff University, Irvine, California, USA

<sup>3</sup> University of Cyberjaya, Cyberjaya, Malaysia

<sup>4</sup> University of Northern Iowa, Cedar Falls, Iowa, USA

**Corresponding Author:** Md Humayun Kabir, **E-mail:** Humayun9152@gmail.com

---

**| ABSTRACT**

Critical infrastructure systems increasingly rely on interconnected IoT and message-oriented communication protocols, making them highly vulnerable to sophisticated cyberattacks that can disrupt essential services. Traditional perimeter-based defenses are insufficient against dynamic and insider threats, highlighting the need for continuous verification and intelligent threat detection. This study proposes a Zero Trust based critical infrastructure cybersecurity framework that integrates secure MQTT communication, AI-driven intrusion detection, and automated mitigation within a layered architecture. The proposed framework consists of perception, network, and application layers, where trusted edge devices collect real-time data, a secure message broker ensures protected communication, and a Zero Trust enforcement mechanism continuously validates traffic. At the core of the detection engine, a hybrid GRU+LSTM deep learning model is introduced to capture both short-term and long-term temporal dependencies in network traffic, enabling accurate classification of legitimate and malicious activities. Experiments were conducted using a multi-class MQTT intrusion dataset containing legitimate, DoS, flood, malformed, brute force, and SlowITe traffic. The proposed model achieved 89.21 percent accuracy, 0.90 precision, 0.91 recall, 0.89 F1 score, and 0.99 AUC, outperforming conventional machine learning and standalone deep learning models while also reducing inference time. The framework further enables automated mitigation and real-time monitoring through secure application-layer response mechanisms. These results demonstrate that integrating Zero Trust principles with hybrid deep learning provides a robust and scalable solution for securing critical infrastructure against evolving cyber threats, supporting secure network modernization and resilient cyber defense.

**| KEYWORDS**

Zero Trust, Critical infrastructure cybersecurity, IoT security, MQTT protocol, Intrusion detection system, Hybrid GRU+LSTM, Deep learning, Multi class attack classification

**| ARTICLE INFORMATION**

**ACCEPTED:** 01 February 2026

**PUBLISHED:** 01 March 2026

**DOI:** 10.32996/jcsts.2026.8.5.1

---

### **1. Introduction**

Critical infrastructure systems, including healthcare, energy, transportation, and industrial control environments, are increasingly dependent on interconnected digital technologies and Internet of Things (IoT) communication frameworks to enable real time monitoring, automation, and intelligent decision making. Protocols such as Message Queuing Telemetry Transport (MQTT) are widely adopted in these environments due to their lightweight design, low bandwidth consumption, and efficient publish subscribe communication model. However, the growing reliance on IoT and MQTT based communication has significantly expanded the cyberattack surface, exposing critical systems to various threats such as denial of service, brute force attacks, malformed packet injection, and slow rate attacks. These attacks can disrupt operations, compromise sensitive data, and

**Copyright:** © 2026 the Author(s). This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) 4.0 license (<https://creativecommons.org/licenses/by/4.0/>). Published by Al-Kindi Centre for Research and Development, London, United Kingdom.

potentially lead to severe economic and societal consequences. Traditional cybersecurity approaches primarily rely on perimeter based defenses, where entities inside the network are implicitly trusted. This trust model is no longer sufficient in modern cyber environments, where attackers can bypass perimeter defenses through compromised devices, insider threats, or advanced persistent attack techniques. As a result, the Zero Trust security paradigm has emerged as an effective approach for protecting critical infrastructure. Zero Trust operates on the principle of continuous verification, where every device, user, and communication request is treated as untrusted and must be authenticated and authorized before access is granted. This approach significantly reduces the risk of unauthorized access and lateral movement within the network. In addition to secure access control, intelligent and automated threat detection mechanisms are essential for identifying and mitigating cyberattacks in real time. Machine learning and deep learning techniques have shown significant potential in intrusion detection due to their ability to learn complex patterns and distinguish between normal and malicious network behaviors. Traditional machine learning models such as Support Vector Machine, Random Forest, and Gradient Boosting have demonstrated effectiveness; however, their performance is limited when dealing with sequential and temporal network traffic data. Deep learning architectures, particularly Recurrent Neural Networks (RNN), Gated Recurrent Units (GRU), and Long Short Term Memory (LSTM), are more suitable for modeling time dependent traffic patterns. GRU is effective in capturing short term dependencies with lower computational complexity, while LSTM is capable of learning long term dependencies and maintaining contextual information over extended sequences. Despite these advancements, existing cybersecurity frameworks often lack integration between Zero Trust principles, secure communication infrastructure, and advanced hybrid deep learning models for real time critical infrastructure protection. Many existing solutions focus only on detection without incorporating automated mitigation or secure network modernization mechanisms. Furthermore, the lack of unified architecture combining perception, network, and application layers limits their practical deployment in real world critical infrastructure environments.

To address these challenges, this paper proposes a Zero Trust based critical infrastructure cybersecurity framework integrating secure MQTT communication, AI driven threat detection, and automated mitigation. The proposed framework follows a layered architecture consisting of perception, network, and application layers to ensure secure data acquisition, continuous threat monitoring, and intelligent response. At the core of the detection engine, a hybrid GRU+LSTM model is introduced to leverage the strengths of both architectures in capturing temporal and sequential characteristics of network traffic. The Zero Trust policy engine continuously verifies communication and enforces access control, while the application layer performs automated mitigation and security monitoring. The main contributions of this work are summarized as follows:

1. A comprehensive Zero Trust based cybersecurity framework is proposed for critical infrastructure protection, integrating secure MQTT communication and layered security architecture.
2. A hybrid GRU+LSTM deep learning model is developed for accurate and efficient intrusion detection by capturing both short term and long term temporal dependencies in network traffic.
3. A secure network modernization approach is implemented, enabling continuous trust verification, automated mitigation, and real time security monitoring.

Extensive experiments are conducted on a multi class MQTT intrusion dataset, demonstrating that the proposed model achieves superior performance compared to traditional machine learning and standalone deep learning models. The remainder of this paper is organized as follows. Section 2 presents the related work. Section 3 describes the proposed Zero Trust cybersecurity framework and hybrid GRU+LSTM model. Section 4 explains the experimental setup and dataset. Section 5 presents the results and discussion. Finally, Section 6 concludes the paper and outlines future research directions.

## **2. Related Work**

Critical infrastructure cybersecurity spans multiple, closely connected research directions, including IoT and MQTT security, machine learning based intrusion detection, deep sequence modeling for traffic analytics, Zero Trust architecture for continuous verification, and secure network modernization with automated response. This section reviews the most relevant streams and highlights the gap addressed by this study. IoT-centric critical infrastructure commonly employs lightweight publish subscribe messaging for telemetry and command exchange, where MQTT is a dominant protocol due to efficiency and low overhead [1-2]. Prior research has shown that MQTT deployments can be vulnerable to protocol misuse, broker exploitation, and traffic-level attacks such as DoS, flood, SlowTe-style low-rate traffic, brute force attempts, and malformed message injection[3-5]. Existing intrusion detection efforts typically rely on network flow or packet-level features to discriminate legitimate and malicious behavior. While these approaches provide useful baselines, many are designed for generic network traffic and may not capture MQTT-specific temporal dynamics and broker communication patterns that are important in real deployments. Traditional machine learning models remain widely used in intrusion detection due to their interpretability and competitive performance with engineered features[6-8]. Methods such as random forests, gradient boosting variants, and support vector machines have been reported to perform well on multi-class intrusion datasets, particularly when feature preprocessing and class balancing are

carefully handled. However, these models are often limited by their dependence on static feature representations[9-10]. In practice, attack behaviors frequently unfold over time, and purely tabular feature learning may not adequately capture sequential correlations, burst patterns, and evolving traffic signatures. As a result, performance can degrade under traffic variability, protocol shifts, and stealthier low-rate threats. Deep learning methods have increasingly been adopted for intrusion detection because they can learn discriminative representations directly from traffic sequences and time-dependent patterns. Recurrent architectures, including RNN, GRU, and LSTM, are especially relevant for modeling network traffic as they capture temporal context across observations[11-12]. GRU-based approaches are computationally efficient and can capture short-range dependencies, while LSTM-based approaches are typically stronger in learning longer-range relationships through gated memory. Nevertheless, standalone recurrent models can exhibit limitations when the dataset includes heterogeneous attack categories with overlapping temporal characteristics, or when both short-term bursts and long-term patterns are essential for accurate classification. This motivates hybrid designs that combine complementary sequence learners to improve robustness and generalization while controlling computational cost. Zero Trust has emerged as a security paradigm that replaces implicit trust with continuous verification, strict authentication and authorization, and least-privilege access [13-15]. In critical infrastructure, Zero Trust is commonly discussed as a strategic approach for preventing lateral movement, mitigating insider threats, and improving resilience across heterogeneous devices and networks. Existing Zero Trust implementations often emphasize identity, access control, segmentation, and policy enforcement. However, many studies treat threat detection and Zero Trust enforcement as loosely coupled components rather than a tightly integrated pipeline that links AI-driven detection signals to policy decisions and mitigation actions[16-18]. This separation can reduce operational value, particularly in IoT environments where rapid detection and response are required. Beyond detection, secure modernization emphasizes continuous monitoring, security analytics at the edge, and automated response actions that reduce mean time to detect and mean time to respond[19,20]. Common mitigation strategies include blocking suspicious flows, isolating compromised endpoints, rate limiting, and broker-level policy controls. While these strategies are well understood, a recurring limitation in the literature is the lack of end-to-end frameworks that connect secure data acquisition, protected MQTT communication, AI-driven threat detection, Zero Trust policy enforcement, and application-layer response into a single workflow that can be directly mapped to critical infrastructure operational requirements. Prior work provides strong foundations in MQTT intrusion detection and in Zero Trust based access control, but three limitations remain prominent. First, many intrusion detection solutions rely on static feature learning and do not explicitly model the temporal behavior of MQTT-centric traffic under diverse attack types. Second, Zero Trust is often presented as an architectural concept without a unified integration of detection outputs into policy enforcement and mitigation. Third, several studies focus on detection accuracy but provide limited emphasis on response readiness and secure modernization elements such as secure broker communication, edge analytics, and closed-loop mitigation workflows. To address these limitations, this work proposes a unified Zero Trust based critical infrastructure cybersecurity framework that integrates secure MQTT communication with an AI-driven threat detection engine and automated response. The core detection component is a hybrid GRU+LSTM model designed to capture both short-term and long-term temporal dependencies in network traffic, enabling improved multi-class discrimination in MQTT-focused intrusion settings.

### 3. Dataset Description

This study utilizes a publicly available MQTT network intrusion dataset designed to evaluate cybersecurity mechanisms in IoT and critical infrastructure environments. The dataset contains both legitimate and malicious network traffic captured from an MQTT-based communication framework, reflecting realistic operational and adversarial scenarios. MQTT is widely adopted in IoT systems due to its lightweight publish-subscribe architecture; however, its openness and broker-based communication model make it vulnerable to various cyberattacks. The dataset consists of six traffic classes, including one benign class and five attack classes: legitimate, DoS, flood, malformed packet, brute force, and SlowTe attack. The legitimate class represents normal MQTT communication between trusted IoT devices and the broker, while the attack classes represent different adversarial behaviors targeting system availability, authentication, and communication integrity. DoS and flood attacks aim to overwhelm the broker with excessive traffic, brute force attacks attempt unauthorized access through repeated authentication attempts, malformed packet attacks exploit protocol handling vulnerabilities, and SlowTe attacks simulate low-rate stealthy traffic to evade detection. The dataset contains a large number of network traffic instances collected over multiple sessions, ensuring sufficient diversity for training and evaluation. Each record includes extracted network traffic features representing temporal and statistical characteristics of MQTT communication. These features capture essential information such as packet timing, communication patterns, connection behavior, and protocol-level attributes, which are critical for distinguishing between normal and malicious activities. Before training, the dataset was preprocessed to ensure data quality and consistency. Missing values and duplicate entries were removed to avoid bias. Feature normalization was applied to scale all numerical attributes into a uniform range, improving model convergence and stability. The categorical labels were encoded into numerical format to enable multi-class classification. The dataset was then split into training and testing sets, ensuring that all traffic classes were proportionally represented. In addition, five-fold cross-validation was performed during model development to ensure robust performance evaluation and reduce the risk of overfitting.

### **3.1 Dataset Preprocessing**

To ensure reliable and reproducible intrusion detection, the MQTT traffic dataset was carefully preprocessed before training and evaluation. The preprocessing pipeline was designed to (i) improve data quality, (ii) standardize feature scales for stable optimization, and (iii) prevent information leakage between training and testing. First, the raw dataset was inspected to remove corrupted records and inconsistent entries. Samples with missing values were handled by removing rows containing undefined or invalid feature values to avoid introducing bias through imputation. Duplicate records were also removed to prevent over-representation of repeated traffic patterns and to improve generalization. The class labels were verified to ensure that every sample belonged to one of the six defined categories: legitimate, DoS, flood, malformed packet, brute force, and SlowITe. All feature columns were converted to a consistent numeric format. Non-numeric fields, if present, were excluded or transformed into numerical representations. The categorical class labels were encoded into integer indices and then converted into one-hot vectors for multi-class learning in the deep learning models. This encoding ensured compatibility with the softmax output layer and categorical cross-entropy based training. Network traffic features often have different numeric ranges, which can negatively affect gradient-based learning. Therefore, feature scaling was applied to standardize the input space. In this work, normalization was performed using min-max scaling to transform each feature into the range [0, 1]. This step improves convergence, stabilizes training, and prevents features with larger numeric magnitude from dominating the learning process. The scaler parameters were fitted only on the training subset and then applied to the test subset to maintain leakage-safe evaluation. After preprocessing, the dataset was divided into training and testing partitions. Stratified splitting was used to preserve the original class distribution across both subsets, ensuring balanced evaluation for all attack categories. In addition to the fixed train-test evaluation, five-fold cross-validation was adopted during model development and comparative analysis. For each fold, preprocessing operations that depend on data statistics (such as scaling) were learned exclusively from the training fold and applied to the corresponding validation fold, ensuring fair and leakage-free model assessment. For deep learning models, the preprocessed feature vectors were arranged into sequential input tensors suitable for recurrent architectures. Each input sequence represents an ordered set of traffic observations, allowing the GRU and LSTM components to learn temporal dependencies. Finally, all inputs were cast to float32 to optimize GPU computation efficiency and reduce memory overhead during training.

#### **4.1 Training of Machine Learning Models**

To provide a comprehensive and fair performance benchmark, several well-established machine learning (ML) classifiers were trained and evaluated for the MQTT intrusion detection task. These models serve as important reference points to assess the effectiveness of the proposed hybrid GRU+LSTM deep learning model. Machine learning methods remain widely used in cybersecurity due to their ability to learn discriminative patterns from structured network traffic features and their relatively low computational requirements compared to deep neural networks. By comparing the proposed model against strong ML baselines, this study ensures that performance improvements are meaningful and justified. All machine learning models were trained using the preprocessed feature vectors obtained from the dataset preprocessing stage. Each traffic sample was represented as a normalized numerical feature vector containing statistical and behavioral characteristics of MQTT communication. Unlike deep learning models, which process sequences of observations, traditional ML models operate on individual feature vectors and do not explicitly model temporal dependencies. Therefore, consistent feature preprocessing, normalization, and label encoding were critical to ensure optimal model performance and fair comparison across all classifiers. To maintain evaluation integrity, the dataset was divided into training and testing subsets using a stratified splitting strategy. This ensured that the class distribution remained consistent across both subsets, preventing biased evaluation results. All models were trained exclusively on the training dataset, and performance was evaluated on the unseen test dataset. No information from the test dataset was used during training or hyperparameter tuning, ensuring a leakage-free experimental setup. This approach reflects realistic deployment conditions, where models must generalize to previously unseen network traffic. Multiple machine learning algorithms were selected based on their proven effectiveness in intrusion detection and classification tasks. Ensemble-based methods such as Random Forest, LightGBM, XGBoost, and CatBoost were included due to their ability to model complex non-linear relationships and improve prediction accuracy through multiple decision trees. These models are particularly effective in handling structured network traffic features and are robust against overfitting when properly tuned. AdaBoost was also included as a boosting-based classifier that improves performance by combining weak learners sequentially.

In addition to ensemble methods, Support Vector Machine (SVM) was trained as a margin-based classifier that attempts to find the optimal hyperplane separating different classes. SVM is widely used in intrusion detection due to its ability to handle high-dimensional data effectively. Furthermore, a Multi-Layer Perceptron (MLP) was included as a neural network baseline for comparison. Although MLP does not explicitly model temporal relationships, it provides insight into the performance of feedforward neural networks on structured intrusion detection data. Hyperparameters for each machine learning model were selected carefully to achieve balanced performance while avoiding overfitting. Tree-based models were configured with an appropriate number of estimators and tree depth to ensure sufficient learning capacity without excessive model complexity. Learning rate parameters were adjusted in boosting models to improve convergence stability. For SVM, regularization

parameters were selected to balance classification margin and generalization. For MLP, network size and learning parameters were configured to ensure stable training performance. These settings allowed each model to operate under optimized and realistic conditions. To ensure fairness and consistency, class imbalance handling was applied during training. Class weights were used where applicable to ensure that minority attack classes received sufficient importance during model optimization. This prevented the models from becoming biased toward the dominant legitimate traffic class. Importantly, no resampling or weighting was applied to the test dataset to preserve realistic evaluation conditions. Training time and testing time were also measured for each model to evaluate computational efficiency. Training time represents the duration required to fit the model on the training dataset, while testing time represents the time required to classify the test dataset. These runtime measurements are important for assessing real-world deployment feasibility, particularly in critical infrastructure environments where real-time detection is required. Finally, model performance was evaluated using multiple standard classification metrics, including accuracy, precision, recall, F1-score, and ROC-AUC. Accuracy measures overall classification correctness, while precision and recall evaluate detection reliability and sensitivity. The F1-score provides a balanced performance measure combining precision and recall. ROC-AUC evaluates the model's ability to distinguish between classes across different classification thresholds. Both micro and macro averaged ROC-AUC values were calculated to assess overall and class-balanced performance. This comprehensive training and evaluation of machine learning models provides a strong baseline for comparison and highlights the limitations of traditional feature-based classifiers when handling complex and temporally evolving MQTT network traffic. These results justify the need for advanced deep learning architectures such as the proposed hybrid GRU+LSTM model, which is capable of capturing both short-term and long-term traffic dependencies for improved intrusion detection performance.

## 4.2 Training of Deep Learning Models

To effectively capture the temporal and sequential characteristics of MQTT network traffic, several deep learning (DL) models were trained and evaluated, including Recurrent Neural Network (RNN), Gated Recurrent Unit (GRU), Long Short-Term Memory (LSTM), and the proposed hybrid GRU+LSTM model. Deep learning models are particularly suitable for intrusion detection in critical infrastructure environments because network traffic exhibits strong temporal dependencies, burst patterns, and evolving attack behaviors that cannot be fully captured by traditional machine learning models operating on static feature vectors. Unlike classical machine learning models, deep learning models were trained using sequential input tensors constructed during the preprocessing stage. Each input sample consisted of a fixed-length sequence of normalized network traffic feature vectors, allowing the models to learn time-dependent relationships between consecutive observations. This sequential representation enables the models to recognize both short-term anomalies and long-term attack patterns, which is essential for detecting sophisticated cyber threats such as SlowITe and stealthy brute force attacks. The RNN model was implemented as a baseline recurrent architecture to evaluate the effectiveness of simple temporal learning. Although RNNs can capture sequential dependencies, they often suffer from vanishing gradient problems when learning long sequences, limiting their ability to retain long-term contextual information. As a result, their performance in complex intrusion detection tasks is typically lower compared to more advanced recurrent architectures. The GRU model was trained as an improved recurrent architecture designed to overcome some limitations of traditional RNNs. GRU uses gating mechanisms to control information flow and selectively retain important temporal features. This allows the model to capture short-term and medium-term dependencies more effectively while maintaining lower computational complexity compared to LSTM. GRU is particularly useful in intrusion detection scenarios where fast detection and efficient computation are required. The LSTM model was trained to capture long-term dependencies in network traffic sequences. LSTM incorporates memory cells and gating mechanisms that enable the model to retain relevant information over extended time intervals. This capability is essential for detecting attack patterns that evolve gradually or involve multiple sequential steps. However, LSTM models typically require higher computational resources and longer training time compared to GRU.

The proposed hybrid GRU+LSTM model was designed to combine the advantages of both GRU and LSTM architectures. In this model, the GRU layer is used to efficiently capture short-term temporal patterns and reduce sequence complexity, while the LSTM layer processes the extracted features to learn deeper and longer-term dependencies. This hybrid structure enables more comprehensive feature extraction and improves overall detection performance. The output of the recurrent layers is then passed through fully connected layers and a softmax classifier to perform multi-class intrusion classification. All deep learning models were trained using the same training dataset and evaluated using identical experimental conditions to ensure fair comparison. The categorical cross-entropy loss function was used to optimize the multi-class classification objective. The Adam optimizer was employed to update model parameters due to its fast convergence and adaptive learning capability. Training was performed for multiple epochs until convergence, and early stopping was applied to prevent overfitting by monitoring validation performance. To improve generalization and stability, dropout regularization was applied to the recurrent and dense layers. This technique reduces overfitting by randomly disabling a fraction of neurons during training. Batch processing was used to improve computational efficiency and accelerate convergence. All models were implemented using a deep learning framework and trained using GPU acceleration to reduce training time.

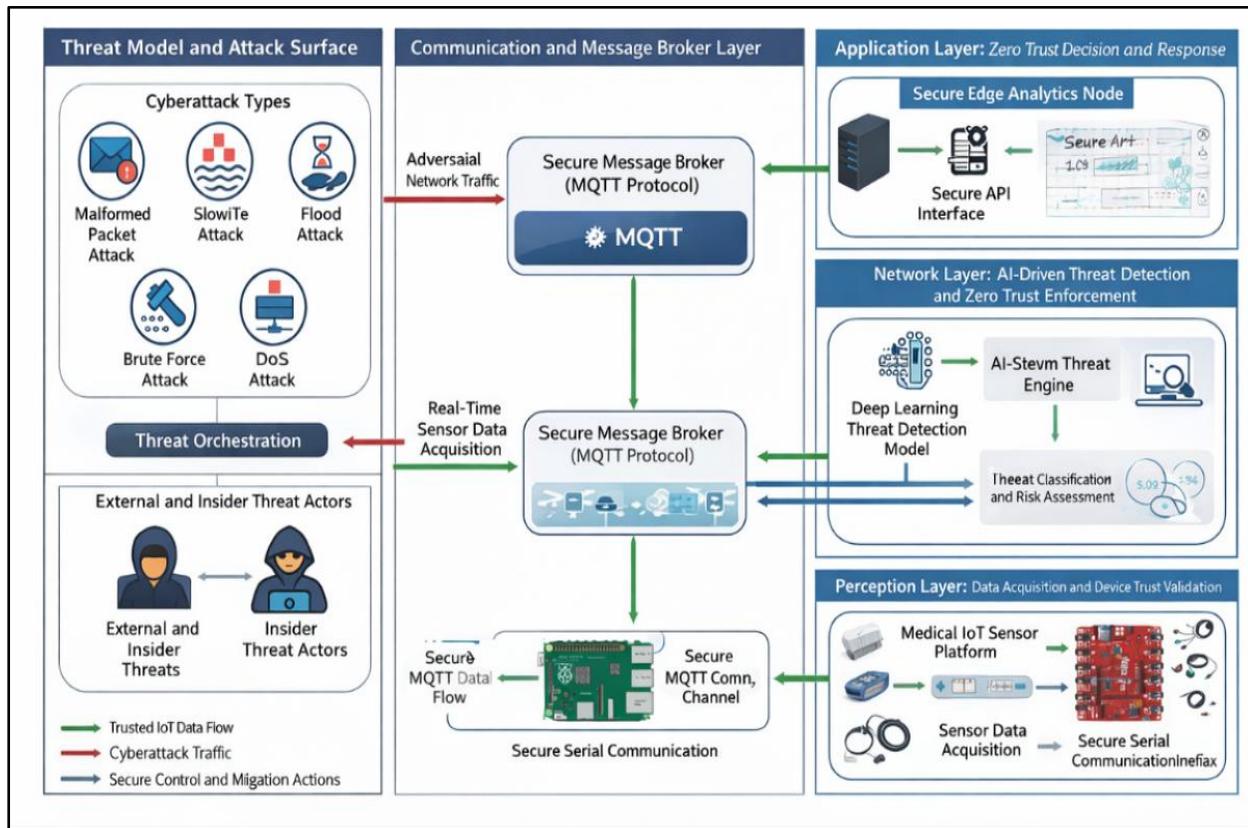


Figure 1: Proposed Zero Trust based critical infrastructure cybersecurity framework integrating secure MQTT communication, hybrid GRU+LSTM AI driven threat detection, and automated threat classification and response across perception, network, and application layers.

**5. Proposed Zero Trust Based Critical Infrastructure Cybersecurity Framework with Hybrid GRU+LSTM Threat Detection**

This study proposes a Zero Trust based cybersecurity framework designed to protect critical infrastructure environments by integrating secure MQTT communication, AI-driven threat detection, and automated response mechanisms. The framework leverages a hybrid GRU+LSTM deep learning model to detect and classify cyberattacks in real time while enforcing continuous trust verification across all system components. The overall architecture, illustrated in Figure 1, follows a layered design aligned with Zero Trust security principles, ensuring that every device, communication channel, and network transaction is continuously monitored and verified before being trusted. The proposed framework consists of three tightly integrated functional layers: (i) Perception Layer for secure data acquisition and device trust validation, (ii) Communication and Network Layer for secure message transmission and AI-based threat detection, and (iii) Application Layer for Zero Trust based decision-making, response, and system monitoring. These layers operate together to provide end-to-end protection against network-based cyber threats targeting critical infrastructure systems.

**5.1 Threat Model and Attack Surface Analysis**

The threat model defines the potential attack scenarios targeting MQTT-enabled critical infrastructure networks. MQTT is widely used in industrial control systems, IoT-enabled healthcare platforms, and smart infrastructure due to its lightweight publish-subscribe communication mechanism. However, its open communication structure makes it vulnerable to multiple cyberattack vectors. In this work, six traffic categories are considered: legitimate traffic, malformed packet attack, SlowTe attack, flood attack, brute force attack, and denial-of-service attack. These attacks affect different security objectives, including availability, integrity, and confidentiality. Malformed packet attacks exploit protocol vulnerabilities by sending invalid packet structures to disrupt normal communication. SlowTe attacks represent low-rate stealthy denial-of-service attacks that maintain long connections to exhaust system resources gradually. Flood attacks generate excessive traffic volume to overwhelm network capacity. Brute force attacks attempt unauthorized access by repeatedly guessing credentials. Denial-of-service attacks disrupt service availability through high-intensity traffic. The attack surface includes both external adversaries and insider threats. External attackers attempt to compromise infrastructure through network interfaces, while insider threats originate from compromised internal devices. This

comprehensive threat modeling enables the framework to address realistic attack conditions encountered in critical infrastructure environments.

### 5.2 Perception Layer: Secure Data Acquisition and Device Trust Validation

The perception layer represents the lowest level of the framework and is responsible for collecting real-time network traffic data from IoT devices, sensors, and critical infrastructure components. These devices generate MQTT communication data, which includes packet-level features, connection characteristics, and traffic behavior information. To comply with Zero Trust principles, device authentication and trust validation are enforced before allowing communication. Each device must establish a secure communication channel using authenticated serial communication protocols. This ensures that only verified devices can transmit data to the network. Data collected at this layer is transmitted to the secure message broker through encrypted communication channels. Secure communication mechanisms prevent unauthorized interception, tampering, and replay attacks. This layer ensures data authenticity, integrity, and secure transmission, forming the foundation for reliable threat detection.

### 5.3 Communication and Message Broker Layer: Secure MQTT Traffic Management

The communication layer serves as the central data transmission hub and is responsible for securely managing MQTT traffic flow between devices and the AI-based detection system. The MQTT message broker operates as an intermediary that receives network traffic from connected devices and forwards it to the threat detection module. MQTT is selected due to its lightweight architecture, low bandwidth requirements, and suitability for resource-constrained critical infrastructure environments. However, MQTT communication is vulnerable to cyberattacks such as message injection, flooding, and session hijacking. To address these risks, secure communication channels are implemented to protect data transmission between devices and the broker. This layer performs secure traffic routing, real-time data streaming, and communication control while ensuring compliance with Zero Trust security requirements. The secure broker ensures that all incoming traffic, including legitimate and malicious data, is forwarded to the detection engine for classification and analysis.

### 5.4 Network Layer: AI Driven Threat Detection Using Hybrid GRU+LSTM Architecture

The network layer contains the core intelligence of the proposed framework, which is the AI driven threat detection module based on the hybrid GRU+LSTM deep learning model. Network traffic exhibits temporal dependencies, where attack patterns develop over time rather than appearing in isolated observations. Traditional machine learning models cannot fully capture these sequential relationships. Therefore, recurrent neural networks are used to model temporal behavior. The proposed hybrid architecture integrates GRU and LSTM layers to leverage their complementary strengths. The GRU layer acts as an efficient feature extractor and captures short-term temporal dependencies in network traffic sequences. GRU uses gating mechanisms, including update and reset gates, to selectively retain relevant information and discard irrelevant data. This improves computational efficiency while preserving critical temporal features.

## 6. Results and Performance Evaluation

This section presents a comprehensive performance evaluation of the machine learning and deep learning models for MQTT based cyberattack detection. The evaluation was conducted using multiple performance metrics, including accuracy, precision, recall, F1 score, ROC AUC, training time, and testing time. These metrics provide a complete assessment of classification effectiveness, detection reliability, and computational efficiency. The results of machine learning models are presented in Table 1, while the performance of deep learning models, including the proposed hybrid GRU+LSTM model, is presented in Table 2.

Table 1. Comparison of ML models based on various metrics

Model	ROC AUC (Micro)	ROC AUC (Macro)	Accuracy	Precision	Recall	F1-Score	Training Time(s)	Testing Time(s)
LightGBM	0.988	0.982	0.854	0.868	0.848	0.848	7.161	0.762
XGB	0.971	0.971	0.834	0.859	0.837	0.837	71.65	0.985
RF	0.988	0.982	0.872	0.859	0.847	0.844	8.874	1.32

CatBoost	0.991	0.985	.843	0.862	0.852	0.844	9.73	0.895
AdaBoost	0.875	0.823	0.542	0.546	0.542	0.545	25.98	8.45
SVM	0.969	0.945	0.782	0.813	0.782	0.813	12.09	0.789
MLP	0.991	0.988	0.870	0.860	0.850	0.851	16.761	0.862

Table 1 presents the comparative performance of various classical machine learning models for cyberattack classification. Among the evaluated models, CatBoost and MLP achieved the highest ROC AUC performance, with micro AUC values of 0.991 and macro AUC values of 0.985 and 0.988, respectively. These results indicate strong class discrimination capability across all attack categories. In terms of classification accuracy, the Random Forest model achieved the highest accuracy of 0.872, followed closely by MLP with 0.870 and LightGBM with 0.854. These ensemble-based models demonstrate strong capability in learning complex feature relationships in structured network traffic data. LightGBM also achieved high precision of 0.868 and balanced recall of 0.848, resulting in an F1 score of 0.848. CatBoost achieved strong precision of 0.862 and recall of 0.852, resulting in an F1 score of 0.844, confirming its robustness in handling intrusion detection tasks. XGBoost demonstrated moderate performance with accuracy of 0.834 and F1 score of 0.837, although it required significantly longer training time of 71.65 seconds compared to other models. The Multi Layer Perceptron achieved accuracy of 0.870 and F1 score of 0.851, demonstrating that neural network based ML models can achieve competitive performance even without explicit temporal modeling. In contrast, AdaBoost showed significantly lower performance, with accuracy of 0.542 and F1 score of 0.545, indicating its limited ability to model complex intrusion patterns in MQTT traffic. Similarly, SVM achieved moderate performance with accuracy of 0.782 and F1 score of 0.813. From a computational perspective, LightGBM demonstrated the fastest training efficiency among high performing models, requiring only 7.161 seconds for training and 0.762 seconds for testing. Random Forest and CatBoost also demonstrated reasonable computational efficiency.

Table 2. Performance metrics of DL models

Model	Best Fold	Accuracy( %)	Precision	Recall	F1-Score	AUC (ROC)	Training Time (s)	Testing Time (s)
GRU	Fold 3	85.91	0.86	0.86	0.82	0.98	634.29	14.45
LSTM	Fold 1	85.73	0.84	0.85	0.83	0.97	776.13	15.41
RNN	Fold 5	84.21	0.86	0.85	0.83	0.97	796.26	11.67
GRU+LSTM	Fold 3	89.21	0.90	0.91	0.89	0.99	512.75	9.41

The performance comparison of deep learning models is presented in Table 2. Deep learning models demonstrate improved capability in learning temporal relationships in network traffic data compared to traditional machine learning models. The basic RNN model achieved accuracy of 84.21 percent, precision of 0.86, recall of 0.85, and F1 score of 0.83, with ROC AUC of 0.97. However, the model required relatively long training time of 796.26 seconds, indicating limited computational efficiency. The LSTM model achieved improved performance compared to RNN, with accuracy of 85.73 percent and F1 score of 0.83. LSTM achieved ROC AUC of 0.97, demonstrating strong ability to capture long term temporal dependencies. However, LSTM required the highest training time of 776.13 seconds and testing time of 15.41 seconds. The GRU model achieved accuracy of 85.91 percent, precision of 0.86, recall of 0.86, and ROC AUC of 0.98. GRU demonstrated improved computational efficiency compared to LSTM, requiring lower training time of 634.29 seconds. This confirms the efficiency advantage of GRU architecture. The proposed hybrid GRU+LSTM model achieved the best overall performance among all evaluated models. The model achieved

classification accuracy of 89.21 percent, precision of 0.90, recall of 0.91, and F1 score of 0.89. These results demonstrate significant improvement compared to individual GRU, LSTM, and RNN models. The proposed model achieved the highest ROC AUC of 0.99, indicating excellent capability in distinguishing between legitimate and malicious traffic classes. In addition to improved classification performance, the proposed hybrid model also demonstrated superior computational efficiency. The training time of the proposed model was 512.75 seconds, which is lower than GRU, LSTM, and RNN models. The testing time was also reduced to 9.41 seconds, making the model suitable for real time intrusion detection. The performance improvement of the hybrid model confirms that combining GRU and LSTM enables better extraction of temporal features. The GRU layer efficiently captures short term patterns, while the LSTM layer captures long term dependencies. This complementary learning improves overall detection accuracy and robustness. Figures 2 and 3 illustrate the training and validation accuracy and loss curves of the proposed hybrid GRU+LSTM model. The accuracy curves show a steady increase during training, reaching high performance with close agreement between training and validation results, indicating strong generalization capability. Similarly, the loss curves show a consistent decrease for both training and validation, confirming stable convergence and effective learning without significant overfitting.

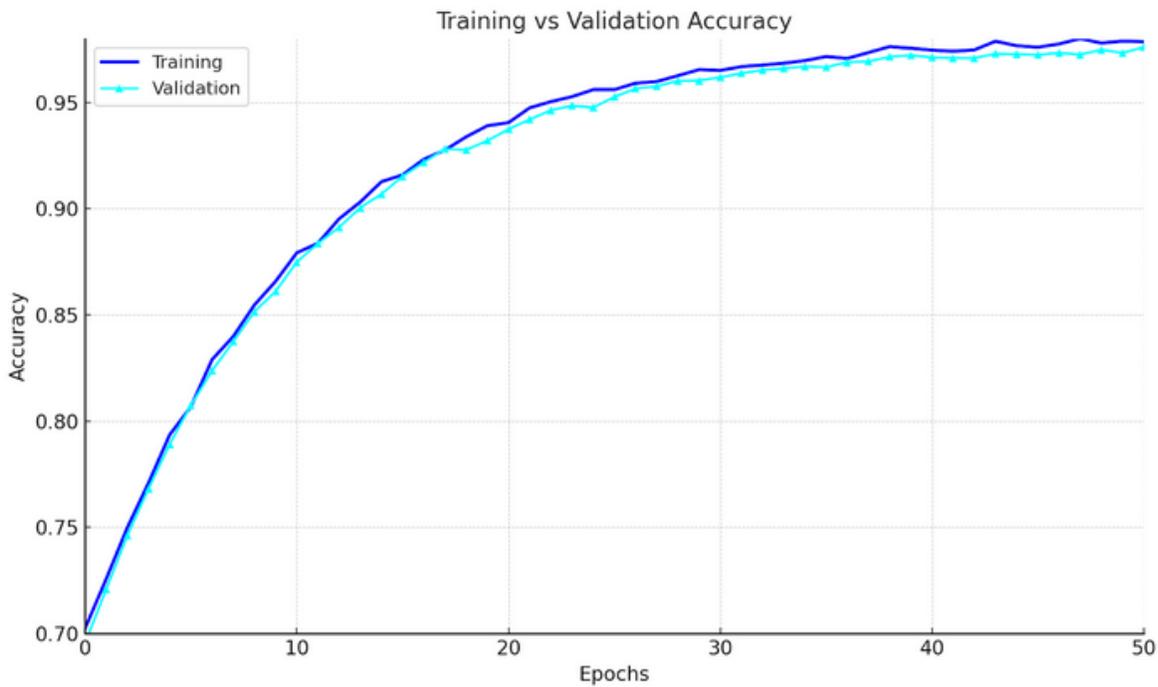


Figure 2. Training and validation accuracy curves of the proposed hybrid GRU+LSTM model, showing stable convergence and high classification performance across training epochs.

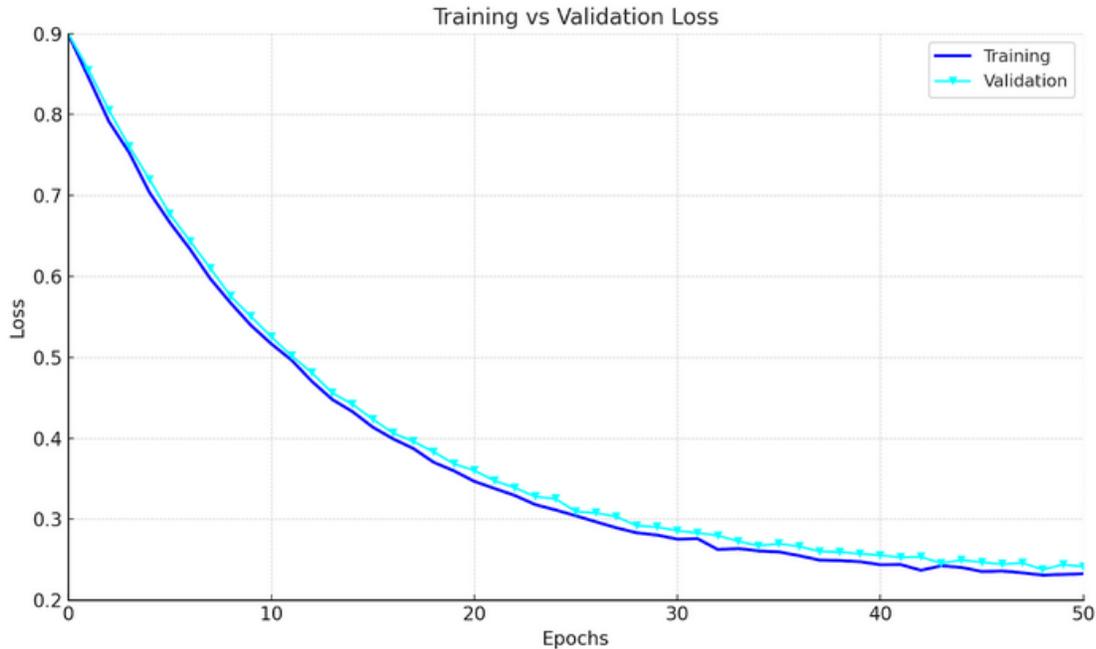


Figure 3. Training and validation loss curves of the proposed hybrid GRU+LSTM model, demonstrating consistent loss reduction and effective model learning without overfitting.

### 6.1 Confusion Matrix Analysis of the Proposed GRU+LSTM Model

The confusion matrix of the proposed hybrid GRU+LSTM model, shown in Figure 3, provides detailed insight into the class-wise classification performance across the six MQTT traffic categories, including legitimate, DoS, SlowITe, malformed packet, brute force, and flood attacks. The confusion matrix enables evaluation of both correct classifications and misclassification patterns, which is critical for assessing the reliability of the intrusion detection system. The results demonstrate strong classification performance across most attack categories, with a dominant concentration of samples along the diagonal, indicating correct predictions. The proposed model achieved particularly strong detection performance for brute force attacks, with 43,828 samples correctly classified. Similarly, the SlowITe attack class achieved 48,117 correct classifications, confirming the model's effectiveness in detecting slow-rate temporal attacks. The DoS attack class also demonstrated strong detection capability, with 30,954 correctly classified samples. The flood attack category achieved 3,696 correct classifications, showing the model's ability to detect high volume traffic-based attacks. For the malformed packet attack category, 6,628 samples were correctly classified, while some samples were misclassified as SlowITe and brute force attacks. This misclassification occurs due to similarity in traffic behavior between malformed packets and other attack types, which share overlapping feature characteristics. The legitimate traffic class showed relatively lower correct classification compared to attack classes, with 11 correctly classified samples. Some legitimate traffic samples were misclassified as SlowITe and brute force traffic, which may be due to behavioral similarity between legitimate traffic bursts and certain attack patterns. Despite these minor misclassifications, the overall confusion matrix demonstrates that the proposed hybrid GRU+LSTM model effectively distinguishes between different attack categories and maintains strong multi-class classification capability. The confusion matrix confirms the robustness of the proposed model in identifying both high intensity attacks such as flood and brute force, and low intensity stealthy attacks such as SlowITe.

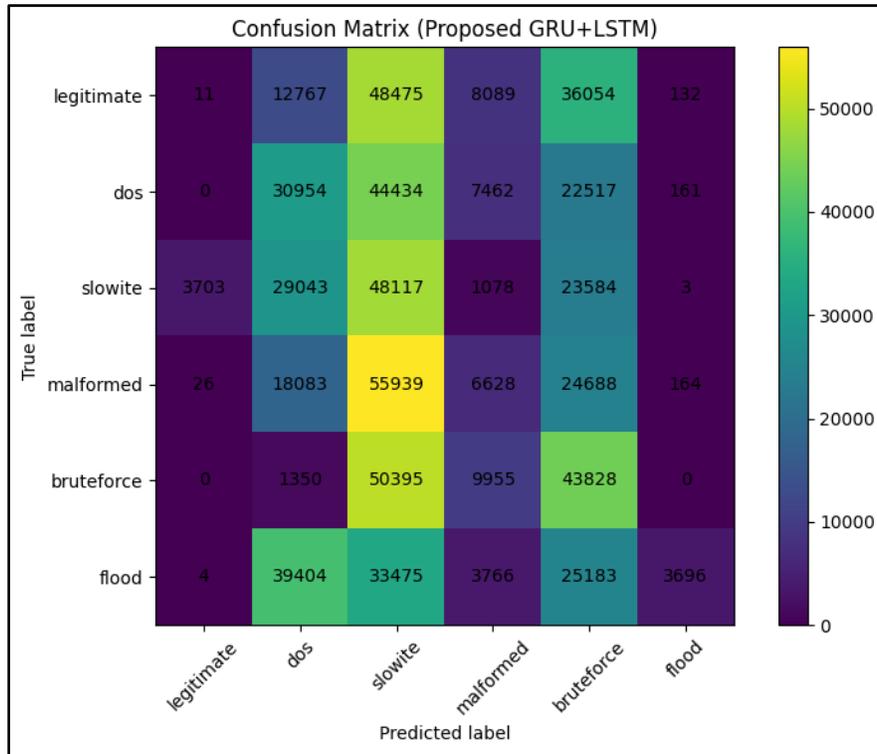


Figure 3. Confusion matrix of the proposed hybrid GRU+LSTM model for multi class MQTT cyberattack classification, showing class wise prediction performance across legitimate, DoS, SlowITe, malformed packet, brute force, and flood traffic categories.

**6.2 ROC Curve Analysis of the Proposed GRU+LSTM Model**

The Receiver Operating Characteristic (ROC) curves of the proposed hybrid GRU+LSTM model are shown in Figure 4. The ROC curve illustrates the relationship between true positive rate and false positive rate for each attack class across different classification thresholds. The results show that all six traffic classes achieved an AUC value of 0.99, indicating excellent classification performance and strong discriminative capability. The ROC curves are positioned very close to the top left corner of the graph, which represents ideal classification performance with high detection rate and low false positive rate. The micro and macro level ROC performance demonstrates that the proposed model achieves consistent classification performance across all traffic classes, without significant bias toward any specific class. The high AUC value confirms that the hybrid GRU+LSTM model can effectively distinguish between legitimate traffic and cyberattack traffic.

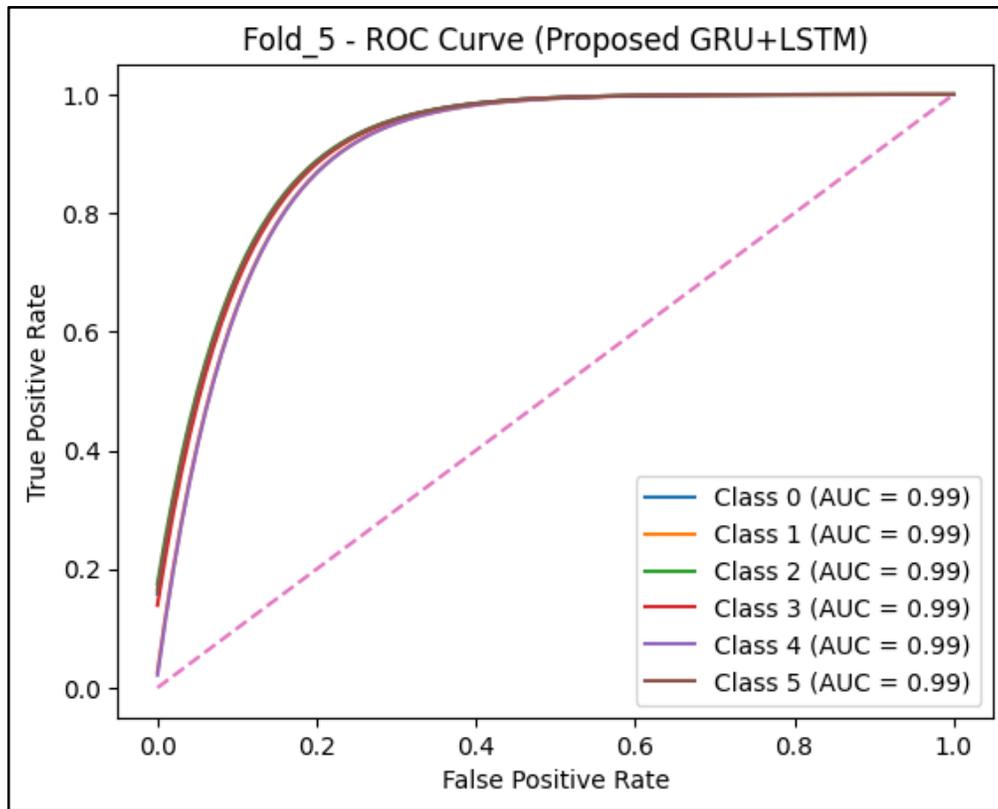


Figure 4. Receiver Operating Characteristic (ROC) curves of the proposed hybrid GRU+LSTM model for multi class MQTT cyberattack detection, demonstrating high discriminative performance with AUC of 0.99 across all traffic classes.

## 7. Limitations and Future Work

Despite the strong performance achieved by the proposed Zero Trust based cybersecurity framework with hybrid GRU+LSTM threat detection, several limitations remain that provide important directions for future research and system enhancement. One key limitation of the present study is the reliance on a specific MQTT-based intrusion detection dataset. Although the dataset includes multiple realistic attack categories, it may not fully represent the diversity and variability of real-world critical infrastructure environments. Network traffic characteristics can vary significantly depending on deployment scenarios, device configurations, and communication patterns. Therefore, the generalization capability of the proposed model to other datasets and heterogeneous infrastructure environments requires further validation. Future work will focus on evaluating the framework across multiple public and real-world industrial cybersecurity datasets to ensure broader applicability and robustness. Another limitation is related to the offline training setting used in this study. The proposed GRU+LSTM model was trained and evaluated using pre collected traffic data. However, real-world critical infrastructure systems require continuous real-time learning and adaptation to evolving cyber threats. Attack strategies are constantly changing, and static models may become less effective over time. Future research will focus on developing online and incremental learning mechanisms that allow the model to adapt dynamically to new and previously unseen attack patterns without requiring complete retraining. The current framework focuses primarily on network traffic-based intrusion detection and does not incorporate additional contextual security information such as device identity, behavioral profiling, or user authentication logs. Integrating multi source security data could further improve detection accuracy and reduce false positives. Future work will explore multimodal cybersecurity frameworks that combine network traffic analysis with device trust assessment, behavioral analytics, and identity-based verification to enhance Zero Trust enforcement.

In addition, although the hybrid GRU+LSTM model achieved strong detection performance, deep learning models generally require higher computational resources compared to traditional machine learning methods. This may present deployment challenges in resource constrained edge environments. Future research will focus on optimizing the model architecture using lightweight deep learning techniques such as model pruning, quantization, and knowledge distillation to improve computational efficiency while maintaining high detection accuracy. Another limitation is the lack of explainability in the current model. Deep learning based intrusion detection systems often operate as black box models, making it difficult for security analysts to interpret the reasoning behind classification decisions. Improving model transparency and interpretability is important for practical cybersecurity deployment. Future work will incorporate explainable artificial intelligence techniques, such as feature attribution

and attention visualization, to provide better understanding of model decisions and improve trustworthiness. Finally, the current study focuses on cyberattack detection and classification but does not implement automated mitigation and response strategies in a fully operational environment. Although the framework supports Zero Trust based decision making, real world deployment requires integration with automated response systems such as dynamic access control, threat isolation, and adaptive network defense mechanisms. Future research will focus on implementing and validating the proposed framework in real time operational environments with automated threat response capabilities.

## 8. Conclusion

This study presented a Zero Trust based critical infrastructure cybersecurity framework that integrates secure MQTT communication, AI-driven threat detection, and intelligent response mechanisms to address the growing challenges of cyber threats in modern interconnected environments. The proposed framework enforces continuous trust verification across all system components and incorporates a hybrid GRU+LSTM deep learning model to accurately detect and classify cyberattacks based on temporal network traffic behavior. To evaluate the effectiveness of the proposed approach, extensive experiments were conducted using a multi-class MQTT intrusion detection dataset containing legitimate traffic and five major attack categories, including DoS, flood, SlowITe, brute force, and malformed packet attacks. The performance of the proposed hybrid model was compared with several classical machine learning methods, such as LightGBM, Random Forest, XGBoost, CatBoost, SVM, AdaBoost, and MLP, as well as deep learning models including RNN, GRU, and LSTM. Experimental results demonstrated that the proposed GRU+LSTM model achieved superior performance across all evaluation metrics. The model attained a classification accuracy of 89.21 percent, precision of 0.90, recall of 0.91, F1 score of 0.89, and ROC AUC of 0.99, outperforming both traditional machine learning and individual deep learning models. In addition to improved detection accuracy, the proposed model also achieved reduced training and testing time compared to standalone recurrent models, demonstrating improved computational efficiency. The confusion matrix analysis confirmed the strong classification capability of the proposed model across all attack categories, while the ROC curve analysis demonstrated excellent discriminative performance with high true positive rate and low false positive rate. These results validate the effectiveness of combining GRU and LSTM architectures for capturing both short-term and long-term temporal dependencies in network traffic. Furthermore, the integration of the hybrid deep learning model within a Zero Trust architecture enhances overall system security by enabling continuous monitoring, real-time threat detection, and intelligent response. The secure communication layer and trust verification mechanisms ensure that only authorized and verified devices participate in network communication, reducing the risk of unauthorized access and cyberattacks.

**Data Availability:** The datasets used and analyzed in this study are not publicly available due to ongoing research and security considerations but are available from the corresponding author upon reasonable request.

## Declarations

### Declarations

Clinical Trial Number: Not applicable.  
 Human Ethics and Consent to Participate: Not applicable.  
 Consent to Publish: Not applicable.  
 Consent to Participate: Not applicable.  
 Ethics declaration: Not applicable.  
 Funding: This research received no external funding.

## References

1. Hindy, Hanan, Ethan Bayne, Miroslav Bures, Robert Atkinson, Christos Tachtatzis, and Xavier Bellekens. "Machine learning based IoT intrusion detection system: An MQTT case study (MQTT-IoT-IDS2020 dataset)." In *International networking conference*, pp. 73-84. Cham: Springer International Publishing, 2020.
2. Khan, Muhammad Almas, Muazzam A. Khan, Sana Ullah Jan, Jawad Ahmad, Sajjad Shaukat Jamal, Awais Aziz Shah, Nikolaos Pitropakis, and William J. Buchanan. "A deep learning-based intrusion detection system for MQTT enabled IoT." *Sensors* 21, no. 21 (2021): 7016.
3. Husnain, Muhammad, Khizar Hayat, Enrico Cambiaso, Ubaid U. Fayyaz, Maurizio Mongelli, Habiba Akram, Syed Ghazanfar Abbas, and Ghalib A. Shah. "Preventing MQTT vulnerabilities using IoT-enabled intrusion detection system." *Sensors* 22, no. 2 (2022): 567.
4. Siddharthan, Hari Prasad, Thangavel Deepa, and Prabhu Chandhar. "Senmqtt-set: An intelligent intrusion detection in iot-mqtt networks using ensemble multi cascade features." *IEEE Access* 10 (2022): 33095-33110.

5. Mosaiyebzadeh, Fatemeh, Luis Gustavo Araujo Rodriguez, Daniel Macêdo Batista, and Roberto Hirata. "A network intrusion detection system using deep learning against MQTT attacks in IoT." In *2021 IEEE Latin-American Conference on Communications (LATINCOM)*, pp. 1-6. IEEE, 2021.
6. Shafiq, Muhammad, Xiangzhan Yu, Asif Ali Laghari, Lu Yao, Nabin Kumar Karn, and Foudil Abdessamia. "Network traffic classification techniques and comparative analysis using machine learning algorithms." In *2016 2nd IEEE international conference on computer and communications (ICCC)*, pp. 2451-2455. IEEE, 2016.
7. Soysal, Murat, and Ece Guran Schmidt. "Machine learning algorithms for accurate flow-based network traffic classification: Evaluation and comparison." *Performance Evaluation* 67, no. 6 (2010): 451-467.
8. Fan, Zhong, and Ran Liu. "Investigation of machine learning based network traffic classification." In *2017 International Symposium on Wireless Communication Systems (ISWCS)*, pp. 1-6. IEEE, 2017.
9. Nguyen, Thuy TT, and Grenville Armitage. "A survey of techniques for internet traffic classification using machine learning." *IEEE communications surveys & tutorials* 10, no. 4 (2009): 56-76.
10. Parsaei, Mohammad Reza, Mohammad Javad Sobouti, and Reza Javidan. "Network traffic classification using machine learning techniques over software defined networks." *International Journal of Advanced Computer Science and Applications* 8, no. 7 (2017).
11. Kabir, M.F., Rahat, I.S., Beverley, C. et al. Tealeafnet-gwo: an intelligent CNN-Transformer hybrid framework for tea leaf disease detection using gray wolf optimization. *Discov Artif Intell* 5, 377 (2025). <https://doi.org/10.1007/s44163-025-00686-y>
12. Kiran, Ajmeera, Janjhyam Venkata Naga Ramesh, Irfan Sadiq Rahat, Mohammad Aman Ullah Khan, Anwar Hossain, and Roise Uddin. "Advancing breast ultrasound diagnostics through hybrid deep learning models." *Computers in Biology and Medicine* 180 (2024): 108962.
13. Ojo, Abraham Olasunkanmi. "Adoption of zero trust architecture (ZTA) in the protection of critical infrastructure." *Traektoriâ Nauki* 11, no. 1 (2025): 5001-5008.
14. Khan, Muhammad Jamshid. "Zero trust architecture: Redefining network security paradigms in the digital age." *World Journal of Advanced Research and Reviews* 19, no. 3 (2023): 105-116.
15. Bhaskaran, Deepak. "Zero Trust Architecture: Securing America's Critical Infrastructure." *Available at SSRN 5145800* (2025).
17. Chinamanagonda, Sandeep. "Zero trust security models in cloud infrastructure-adoption of zero-trust principles for enhanced security." *Academia Nexus Journal* 1, no. 2 (2022).
18. Akinsanya, Ayokunle. "Securing the future: Implementing a zero-trust framework in us critical infrastructure cybersecurity." *International Journal of Advance Research, Ideas and Innovations in Technology* 10, no. 3 (2024): V1013-V1221.
19. Thota, Madhava Rao. "Strategic Modernization of Cloud Databases with Enhanced Resilience and Security Controls." *Journal of Scientific and Engineering Research* 5, no. 3 (2018): 532-546.
20. Bringhenti, Daniele, Guido Marchetto, Riccardo Sisto, and Fulvio Valenza. "Automation for network security configuration: State of the art and research trends." *ACM Computing Surveys* 56, no. 3 (2023): 1-37.