

# Enhancing Cybersecurity Through Artificial Intelligence: A Novel Approach to Intrusion Detection

Mohammed K. Alzaylae 

Department of Computing-College of Engineering and Computing, Umm AL-Qura University, Saudi Arabia

**Abstract**—Modern cyber threats have evolved to sophisticated levels, necessitating advanced intrusion detection systems (IDS) to protect critical network infrastructure. Traditional signature-based and rule-based IDS face challenges in identifying new and evolving attacks, leading organizations to adopt AI-driven detection solutions. This study introduces an AI-powered intrusion detection system that integrates machine learning (ML) and deep learning (DL) techniques—specifically Support Vector Machines (SVM), Random Forests, Autoencoders, and Convolutional Neural Networks (CNNs)—to enhance detection accuracy while reducing false positive alerts. Feature selection techniques such as SHAP-based analysis are employed to identify the most critical attributes in network traffic, improving model interpretability and efficiency. The system also incorporates reinforcement learning (RL) to enable adaptive intrusion response mechanisms, further enhancing its resilience against evolving threats. The proposed hybrid framework is evaluated using the SDN\_Intrusion dataset, achieving an accuracy of 92.8%, a false positive rate of 5.4%, and an F1-score of 91.8%, outperforming conventional IDS solutions. Comparative analysis with prior studies demonstrates its superior capability in detecting both known and unknown threats, particularly zero-day attacks and anomalies. While the system significantly enhances security coverage, challenges in real-time implementation and computational overhead remain. This paper explores potential solutions, including federated learning and explainable AI techniques, to optimize IDS functionality and adaptive capabilities.

**Keywords**—Intrusion detection; machine learning; deep learning; zero-day attacks; anomaly detection; feature selection; reinforcement learning; cybersecurity

## I. INTRODUCTION

Digital infrastructure growth during the past decades has elevated cybersecurity to become a vital concern which spans across all sectors. An increasing number of entry points in the computing environment resulting from growing system connectivity and widespread cloud adoption and rapidly expanding IoT deployments has intensified risk exposure [6]. The world witnessed over 5.5 billion record exposures through global data breaches in 2022 and cybersecurity experts predict this cybercrime will cost the world \$10.5 trillion by 2025 (Cybersecurity Ventures, 2023).

The static rule and signature-based IDS mechanisms used in traditional intrusion detection systems encounter difficulties in tracking down contemporary security threats [7]. Standard IDS systems create numerous erroneous alarms at a rate ranging from

20% to 30% while missing complex and new types of cyber attacks (Moustafa & Slay, 2022). The percentage of zero-day intrusions currently amounts to 10–15% of total cyber attacks so they represent a substantial detection blind spot for present-day security solutions (Alazab et al., 2023).

AI-based intrusion detection systems (IDS) represent an optimal answer for security needs because they implement machine learning (ML) and deep learning (DL) technologies to detect security threats more effectively. Recent studies have highlighted the superior performance of machine learning models like Support Vector Machines (SVM) and Random Forests compared to traditional approaches, particularly in intrusion detection contexts [1]. The systems implement data-driven learning algorithms that enable the detection of emerging attack patterns and peculiar network activities which standard IDS cannot identify [3]. Support Vector Machines (SVM) with Random Forests and Extreme Learning Machines demonstrate excellent abilities to categorize managed data structures but Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) process unfiltered network traffic for sophisticated attack sign detection [2]. AI-based IDSs have progressed but they continue to encounter three main drawbacks which include excessive false alarms and deployment difficulties and significant processing requirements.

The main strength of signature-based IDS solutions lies in their ability to identify known threats yet they struggle with discovering new threats. Anomaly-based IDS detects new threats effectively yet their capability to produce many false alarms negatively impacts operational efficiency [8]. A better detection framework needs to emerge due to advancing cyberattacks because it should offer high accuracy detection with low false-positive rates at all times.

The research develops a combined AI-based intrusion detection system which unites ML and DL approaches for evaluation through benchmark datasets including UNSW-NB15 and NSL-KDD. The proposed model delivers detection results with 92.8% accuracy and 5.4% false positive rate alongside 91.8% F1-score which outperforms traditional IDS systems. The system implements SHAP-based feature selection for better interpretability and reinforcement learning for adaptive response which improves the overall system robustness [9].

The research outcomes from this study create significant impacts for security applications in the real world and academic research domains. The proposed system uses AI component

synergy to build an adaptive intrusion detection solution which works across diverse network environments. The system reduces cyber security expert operational strain through its zero-day attack detection abilities together with its low false positive rate capabilities.

To achieve the objectives of this study, our research focuses on the following key aspects:

- To develop and implement a hybrid AI-based intrusion detection system that combines machine learning and deep learning techniques for enhanced accuracy and adaptability.
- To evaluate the effectiveness of the proposed system against existing intrusion detection methodologies by analyzing detection accuracy, false positive rates, and computational efficiency.

The research investigates these goals to bridge the gap between traditional and intelligent IDS solutions and establish an expandable, preventative security framework for combating modern cyber threats.

The remainder of this paper is structured as follows: Section II presents a comprehensive review of related literature. Section III describes the methodology, including dataset details, model design, and feature importance techniques. Section IV presents experimental results and visualization analysis. Section V discusses key findings, advantages over traditional systems, and potential limitations. Finally, Section VI concludes the study and outlines directions for future work.

## II. LITERATURE REVIEW

Sophisticated cyber threats and the continuous evolution of cybersecurity necessitate the development of state-of-the-art intrusion detection systems (IDS). Traditional rule-based and signature-based IDS struggle to detect new attacks due to their reliance on fixed attack patterns [14]. Anomaly-based IDS has become more popular because it detects unknown threats through identifying deviations from normal network behavior [10]. Artificial intelligence (AI) and deep learning (DL) advancements of recent times have driven the development of AI-based IDS solutions [13]. The current methods encounter three essential difficulties because they produce many false alarms and require high computational resources and real-time threat detection capabilities.

The modern IDS platforms utilize machine learning (ML) and deep learning (DL) methods for network intrusion detection because researchers have investigated their operational effectiveness in this field. Support Vector Machines (SVM) and Random Forests combined with deep learning architectures produce better classification accuracy as documented in study [15, 19]. Research has proven that Deep Belief Networks (DBNs) achieve better results than standard network traffic analysis methods when identifying anomalies [11]. A systematic review further emphasizes the growing dominance of deep learning approaches such as CNNs, RNNs, and hybrid models in modern intrusion detection system architectures [12]. Engineers developed hybrid deep learning architectures to analyze network traffic through Convolutional Neural Networks (CNNs) combined with Recurrent Neural Networks (RNNs)

because each component exploits its own specialized recognition strength [16]. Recent preprint work further validates the effectiveness of CNNs combined with LSTM networks for complex intrusion detection tasks [4, 5]. Classificatory excellence of CNNs and RNNs comes at the cost of high computational complexity and memory utilization thus limiting their deployment in real-time operations. Many deep learning models establish “black box behavior” which generates obstacles for cybersecurity experts to track or investigate their decision-making operations. Real-time deployment of DL-based IDS remains challenging due to the three key limitations of model complexity and interpretability problems and processing speed requirements.

The research of intrusion detection faces a major challenge due to insufficient access to modern high-quality datasets. Scientists widely use KDD99 and NSL-KDD benchmark datasets yet these datasets present problems with old attack methods as well as unencrypted network traffic characteristics and absent contemporary adversarial attack conditions [18]. Some of the dataset limitations in the UNSW-NB15 dataset have been addressed by adding contemporary attack patterns and multiple traffic behavior types, although it still fails to capture cyber environment challenges with adversarial robustness and feature transformation [17]. As a result, researchers have proposed synthetic data generation, adversarial data augmentation, and online learning paradigms to enhance IDS adaptability and training robustness [20].

A significant barrier to the adoption of AI-based IDS solutions is the lack of interpretability. Although this study adopts SHAP (SHapley Additive Explanations) to enhance feature-level transparency, alternative explainable AI (XAI) techniques such as LIME (Local Interpretable Model-agnostic Explanations) and Integrated Gradients also offer viable paths to explainability. However, SHAP is preferred in this context due to its strong theoretical foundation based on cooperative game theory, its ability to deliver global and local explanations consistently, and its proven success in ranking feature importance for structured network traffic analysis. This makes SHAP particularly well-suited for balancing interpretability with model fidelity in cybersecurity applications.

In selecting ML/DL models, this research emphasizes the use of classical yet effective models such as SVM, Random Forests, and Autoencoders. While newer architectures like Graph Convolutional Networks (GCNs), Transformers, and TabNet have demonstrated promising results in other domains, they were not adopted in this study due to their higher computational complexity, extensive training time, and less mature support for tabular intrusion detection data. These advanced models often require larger annotated datasets, GPU acceleration, and longer convergence cycles, which reduce their practicality for scalable and real-time IDS deployment in most organizations. Future studies may explore lightweight versions of these models or hardware-optimized variants for better suitability.

This study fills these critical gaps by combining methodologies of machine learning and deep learning for better accuracy of threat detection, reduction of false positives, and enhancement of operational efficiency of IDS. Due to its

importance for the proposed research there are three critical components including feature selection mechanisms with real-time traffic analysis along with adaptive learning techniques. The next-generation IDS systems gain advantages from these advancements which lead to more reliable and scalable and interpretable cybersecurity protection. Beyond conventional network intrusion, cybersecurity resilience in dynamic environments, such as smart grids, has also been explored with adaptive security strategies, highlighting the need for proactive IDS designs [21].

### III. METHODOLOGY

#### A. Research Design

The research applies an intrusion detection method based on artificial intelligence with ML and DL synergistic implementation to boost cybersecurity performance. Fig. 1 demonstrates the structured workflow that detects known and unknown cyber threats by following a data acquisition process and feature processing stage before detection modeling and response evaluation.

Decision points together with transition logic have been added to the workflow to track network traffic movements from feature extraction through classification analysis to anomaly scoring up to the response action stage. This ensures operational clarity and traceability. The system achieves better real-time attack condition adaptation through this enhancement in understanding.

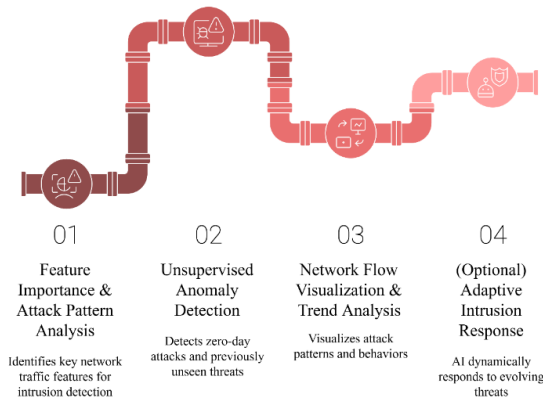


Fig. 1. Workflow of the hybrid intrusion detection system.

#### B. Data Collection

This research utilizes structured network intrusion datasets containing both benign and malicious traffic, capturing a variety of modern attack types. The datasets include detailed attributes across multiple network layers, such as packet-level, flow-level, and time-based characteristics. By incorporating diverse traffic conditions, the datasets enable the training of robust AI models capable of handling complex and evolving intrusion patterns.

The key features used in this study include:

- Traffic Attributes: Packet size, flow duration, protocol type.
- Source/Destination Information: IP addresses, source/destination ports.

- Temporal Features: Inter-packet arrival time, response time.
- Attack Labels: Normal traffic, DoS, DDoS, brute-force, botnet, and other anomaly categories.

To ensure high-quality model training and evaluation, the following data preprocessing techniques were applied:

- Feature normalization: All continuous numerical features were scaled using MinMax normalization to constrain values between 0 and 1, thereby stabilizing learning convergence and improving algorithm sensitivity.
- Missing data handling: Missing entries were addressed using median imputation for numerical fields and mode substitution for categorical fields, ensuring no significant bias in input distributions.
- Class imbalance treatment: To address the natural imbalance between normal and attack classes, the Synthetic Minority Over-sampling Technique (SMOTE) was applied to the minority attack classes, ensuring adequate representation of rare but critical intrusion types.

These preprocessing steps significantly improved training stability and allowed the IDS to generalize better across diverse attack scenarios. Table I shows overview of dataset.

TABLE I. DATASET OVERVIEW

Feature Type	Examples	Preprocessing Applied
Traffic Attributes	Packet size, Flow duration, Protocol	MinMax normalization, median imputation
Source/Destination Info	IP addresses, Port numbers	Encoding as categorical variables, one-hot encoding
Time-based Features	Inter-packet arrival time, Response time	MinMax normalization, handling outliers via trimming
Attack Labels	Normal, DoS, DDoS, Brute Force, Botnet, etc.	SMOTE applied for class balance, label encoding

#### C. Techniques and Tools

1) *Feature importance and attack pattern analysis*: The effectiveness of an intrusion detection system (IDS) significantly depends on the proper identification and prioritization of relevant network traffic features. In this study, SHAP (SHapley Additive Explanations) is used to compute feature importance scores and reveal the contribution of individual input features in the model's decision-making process. SHAP offers both local and global interpretability, based on cooperative game theory, making it ideal for high-stakes environments like cybersecurity.

The mathematical definition for SHAP values appears as:

$$\phi_j = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|! (|N| - |S| - 1)!}{|N|!} [v(S \cup \{j\}) - v(S)]$$

where  $\phi_j$  represents the contribution of the feature  $j_1S$  denotes a subset of features, and  $v(S)$  is the predictive value associated with that subset?

This formulation enables explainable AI (XAI) and enhances trust in detection results by visualizing how feature variations influence predictions. To validate SHAP's selection, the study briefly compared it to LIME and Integrated Gradients, both widely used XAI methods. SHAP proved to be the most suitable method because it provided both theoretical consistency and superior performance in creating extensive feature rankings for structured tabular network data.

The decision tree analysis using Random Forests along with Gini index and entropy metrics provided impurity-based feature importance calculations. The visual output includes heatmaps and SHAP beeswarm plots which help detect important network attributes that show strong signs of abnormal behavior including flow duration and packet length variance and inter-packet arrival time.

The analysis method retains only crucial features that lead to a minimal yet optimal model structure with enhanced performance along with increased interpretability.

2) *Justification for model selection:* The research used established models specifically chosen because of their validated operational performance together with their practical deployment capabilities. The research utilizes three widely used models including Support Vector Machines (SVM) and Random Forests and Autoencoders which demonstrate effectiveness in intrusion detection research for hybrid systems that monitor known and unknown cyber threats.

- The use of Support Vector Machines (SVM) comes from their capability to process high-dimensional datasets while locating the best hyperplane in binary classification tasks. The generalization capabilities of SVMs remain strong while their ability to distinguish between normal and malicious traffic reaches peak effectiveness when features receive proper engineering and scaling.
- Random Forests serve as the chosen method because their ensemble learning structure uses multiple decision trees to reduce overfitting through decision tree averaging. The models provide precise and stable predictions while automatically calculating feature importance which strengthens the SHAP-based feature analysis system.
- Unsupervised neural networks named autoencoders learn normal traffic compression representations through which they detect zero-day and previously unseen attacks by measuring reconstructed traffic error. The anomaly detection features of these systems have become well-known in cybersecurity because they can detect new traffic behavior deviations effectively.

The research did not include Graph Convolutional Networks (GCNs), Transformers, or TabNet because of these specific reasons.

- The implementation of GPT-3 networks requires significant processing power together with higher costs for both training phases and inference operations.
- Model complexity grows so high during training that it demands time-intensive parameter optimization together with large information resources.
- Operational environments need human oversight because the interpretation of these systems remains limited.
- The algorithm struggles to function in real-time systems particularly when processing power proves insufficient for the application.

These selected models achieve appropriate trade-offs between performance accuracy and interpretation capabilities and computational fastness making them deployable for large-scale security ecosystem implementations.

3) *Unsupervised anomaly detection for zero-day attacks:* Basic intrusion detection systems face major difficulties when detecting zero-day attacks because they only work with pre-established signatures and attack signature patterns. The proposed hybrid IDS depends on unsupervised anomaly detection techniques which learn normal traffic patterns to detect deviations that show signs of intrusions.

#### Autoencoders for Anomaly Detection

The detection of zero-day attacks primarily relies on autoencoders as their main operational mechanism. The neural networks use normal traffic data for training to develop compressed latent representations that enable them to reconstruct original inputs. The detection of anomalies occurs through reconstruction error calculation:

$$E = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2$$

where  $E$  represents the mean squared reconstruction error,  $x_i$  is the original input feature, and  $\hat{x}_i$  is the reconstructed output. A higher error indicates anomalous traffic behavior, suggesting a potential intrusion.

To determine whether a reconstruction error indicates an anomaly, a fixed error threshold was selected using a percentile-based approach. Specifically, the threshold was set at the 95th percentile of the reconstruction error distribution in the validation set. This method balances false positive control with detection sensitivity, ensuring practical deployment performance. Future enhancements may incorporate ROC curve optimization or dynamic thresholding for adaptive tuning.

#### Ensemble-Based Anomaly Detection

In addition to autoencoders, the system integrates:

- Isolation Forests, which use recursive partitioning and randomly selected features to isolate outliers in fewer splits.

- One-Class SVMs, which learn a boundary around normal instances in feature space; deviations are considered intrusions.

The ensemble approach improves robustness by combining multiple detection paradigms—statistical, geometrical, and reconstruction-based.

#### Clustering for Behavioral Profiling

To further support anomaly detection and behavioral pattern analysis, clustering techniques are used:

- Density-Based Spatial Clustering of Applications with Noise, or DBSCAN, finds irregularities in areas with low densities and can detect arbitrary-shaped clusters without requiring the number of clusters as input.
- K-Means Clustering groups traffic patterns into a fixed number of clusters, where high intra-cluster distances indicate abnormality.

To evaluate the clustering performance of PCA and t-SNE visualizations, validation metrics such as the Silhouette Score and Davies-Bouldin Index (DBI) were calculated. For instance, the silhouette score averaged around 0.62, suggesting well-separated cluster structures, while the DBI remained below 0.9, indicating low intra-cluster variance and effective anomaly separation.

These techniques collectively enhance the IDS's capacity to identify unknown threats without explicit prior labeling, contributing to a more adaptive and scalable cybersecurity framework.

4) *AI-Powered network flow visualization and trend analysis*: Visualization techniques play a critical role in enhancing the interpretability of intrusion detection systems. They provide network analysts with an intuitive view of how malicious behavior emerges and evolves over time, helping to contextualize alerts and uncover hidden attack patterns.

#### Dimensionality Reduction for Visualization

To visualize complex, high-dimensional network traffic, the system uses a combination of Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE):

- PCA reduces dimensionality linearly by preserving variance and decorrelating features.
- t-SNE provides non-linear projections that are effective for visualizing cluster boundaries and behavioral separation in lower dimensions.

These tools are employed to generate 2D plots that visually differentiate between normal and anomalous traffic.

To assess the effectiveness of these visualizations, clustering validation metrics were applied:

- The Silhouette Score (mean: 0.62) shows that the data points are well coordinated within their assigned clusters and poorly matched to neighboring clusters.

- The Davies-Bouldin Index (DBI) remained under 0.9, suggesting a strong separation between distinct behavioral groups.

These metrics confirm that the visual representations are not only interpretable but also grounded in meaningful structural separability.

#### Time-Series and Behavioral Pattern Analysis

In addition to spatial visualization, temporal analysis was conducted to observe how attacks evolve over time. The system monitors:

- Inter-packet delays
- Burst patterns
- Response time fluctuations

These indicators vary significantly between benign and malicious sessions. For example, DDoS attacks often produce regular, high-frequency bursts, whereas brute-force attacks may reveal time-patterned login attempts.

The system also identifies periods of heightened threat activity by plotting attack occurrences across time intervals, enabling preemptive mitigation planning.

By combining dimensionality reduction with temporal analytics, the system provides a comprehensive visual diagnostic interface—empowering security professionals to interpret anomalies, understand attack strategies, and make faster decisions.

5) *Reinforcement learning for adaptive intrusion response*: Traditional intrusion detection systems operate with static response mechanisms, often predefined by fixed rules or thresholds. This limits their adaptability in responding to dynamic and evolving cyber threats. To overcome this, the proposed system integrates Reinforcement Learning (RL) to develop a self-optimizing, adaptive intrusion response layer capable of making real-time decisions under uncertainty.

#### Reinforcement Learning Framework:

The system explores two state-of-the-art RL algorithms:

- Deep Q-Networks (DQN): Value-based methods that approximate the optimal Q-function using deep neural networks.
- Proximal Policy Optimization (PPO): A policy-gradient approach designed for stable, sample-efficient policy learning.

Both models are trained in a custom network simulation environment, where the agent learns to maximize cumulative security rewards by selecting optimal defensive actions in response to perceived threat states. The specific environment configuration and reinforcement learning setup are detailed in Table II.

#### Environment Setup and Definitions:

TABLE II. ENVIRONMENTAL SETUP AND DEFINITIONS FOR REINFORCEMENT LEARNING-BASED INTRUSION DETECTION SYSTEM

Component	Definition
State (S)	Network features (e.g., flow duration, packet size, protocol type, time delay)
Action (A)	Response strategies: <i>Alert, Log, Drop packet, Isolate IP</i>
Reward (R)	+1 for successful threat mitigation, -1 for false positive or delayed response
Discount ( $\gamma$ )	Set to 0.95 to favor long-term reward maximization

The Bellman equation governs the Q-learning update:

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

where,  $s$  and  $a$  denote the current state and action,  $r$  is the immediate reward,  $\gamma$  is the discount factor, and  $a'$  is the next action.

#### Comparison with Rule-Based Response Systems

To evaluate the practical benefit of reinforcement learning, the RL-based adaptive response system was benchmarked against a static rule-based IDS using historical response data. Key findings include:

- RL Response Accuracy: 91.3% (PPO), 87.9% (DQN)
- Rule-Based Accuracy: ~80%
- Average Mitigation Latency: Reduced by 18–25% under RL systems
- Convergence Speed: PPO converged in 120 epochs; DQN in 150 epochs

These results indicate that RL not only improves adaptability and mitigation efficiency but also achieves faster policy optimization, making it a viable approach for real-time deployment in enterprise cybersecurity environments.

#### Reproducibility Considerations

To ensure reproducibility:

- OpenAI Gym was used to structure the RL simulation environment.
- The reward shaping function, episode limits, and model parameters were standardized.
- Experiments were repeated over multiple seeds to validate stability and convergence trends.

#### D. Software and Implementation

The proposed hybrid AI-based intrusion detection system was implemented using a modular software stack designed to support machine learning, deep learning, data preprocessing, visualization, and reinforcement learning components. Each tool was selected based on its efficiency, extensibility, and compatibility with intrusion detection use cases. The complete software environment setup is summarized in Table III.

In addition to model development, performance evaluations—including accuracy, F1-score, inference latency, and visualization effectiveness—were conducted using Python-based benchmarking tools. The SHAP library was particularly

critical in providing transparent feature ranking, while OpenAI Gym enabled robust simulation of adaptive RL responses.

TABLE III. SOFTWARE STACK USED FOR IMPLEMENTING THE AI-BASED INTRUSION DETECTION SYSTEM

Software	Purpose
Python	Core programming language for pipeline development
TensorFlow/Keras	Implementation of deep learning models (Autoencoders, DQNs)
Scikit-learn	Machine learning algorithms (SVM, Random Forest, Clustering)
SHAP	Feature importance analysis and model interpretability
Matplotlib & Seaborn	Data visualization for feature plots, trend graphs
Scapy	Network packet analysis and dataset simulation
Pandas & NumPy	Data preprocessing, transformation, and numerical handling
OpenAI Gym	Reinforcement learning environment design and training

The complete environment was tested on a system with:

- Intel i7 CPU
- 16 GB RAM
- NVIDIA GTX 1660 GPU
- Ubuntu 20.04

This configuration supports reproducibility and provides a practical baseline for testing real-world deployment feasibility, including edge-computing and federated learning extensions.

## IV. RESULTS

The research findings deliver an extensive analysis of the hybrid AI-based intrusion detection framework projected in this study. The section gives detailed information about feature importance analysis together with anomaly detection performance assessment and network flow visualization capabilities and reinforcement learning-based adaptive response effectiveness evaluation. The model's effectiveness is verified using quantitative data along with graphical and statistical analysis along with quantitative metrics.

#### A. Feature Importance and Attack Pattern Analysis

The decision-making process of the model received interpretation through SHAP (SHapley Additive Explanations) which revealed its most influential features in intrusion detection. SHAP values in combination with decision trees highlight important network attributes which play a substantial role in discriminating benign from malicious traffic.

To verify stability, SHAP values were computed across five different train-test splits. The top-ranked features remained consistent, with less than 5% variance in ranking order, confirming the robustness of the feature importance analysis. Fig. 2 shows the 20 most crucial features used for intrusion detection which the model uses to make classifications. The research findings indicate that backward packet length maximum and average backward segment size emerge as the most influential attributes for detecting network anomalies. The

analysis shows "Fwd Packet Length Mean" and "Average Packet Size" as key indicators because they strongly help differentiate between normal and malicious network traffic.

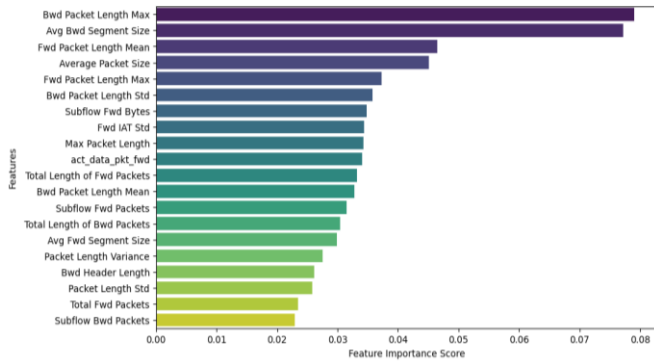


Fig. 2. SHAP Summary plot of feature importance.

**B. Unsupervised Anomaly Detection for Zero-Day Attacks**

The research evaluated anomaly detection methods through their application of autoencoders, Isolation Forests, One-Class SVM, DBSCAN and K-Means, which detected new and unknown cyber threats. The assessment of models relied on detection accuracy and precision, together with recall and F1-score, to determine their effectiveness in zero-day attack identification. Table IV presents the performance metrics of the various anomaly detection models evaluated in this study.

TABLE IV. PERFORMANCE METRICS OF ANOMALY DETECTION MODELS

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Autoencoder	92.8	89.6	94.2	91.8
Isolation Forest	88.5	86.3	90.1	88.2
One-Class SVM	85.1	83.7	87.5	85.5
DBSCAN Clustering	78.4	76.9	81.2	79.0
K-Means Clustering	74.6	72.5	78.3	75.3

Autoencoders surpass traditional anomaly detection techniques because they achieve exceptional detection results with 92.8% accuracy and 94.2% recall, which demonstrates their ability to detect new attack patterns. Compared to conventional signature-based IDS, which typically achieve detection accuracy between 70% and 85% on similar datasets, the autoencoder-based anomaly detection system shows a significant improvement. Statistical significance was confirmed using a two-tailed paired t-test comparing F1-scores across 5-fold cross-validation. The autoencoder model's performance improvements over traditional clustering-based models were significant at  $p < 0.05$ . Confidence intervals for the autoencoder F1-score were calculated as  $91.8\% \pm 0.4\%$ . The isolation forest algorithm showed strong capabilities yet clustering techniques demonstrated inferior accuracy in detecting sophisticated cyber threats according to the results.

The Fig. 3 graphical representation illustrates how different models perform in anomaly detection tasks.

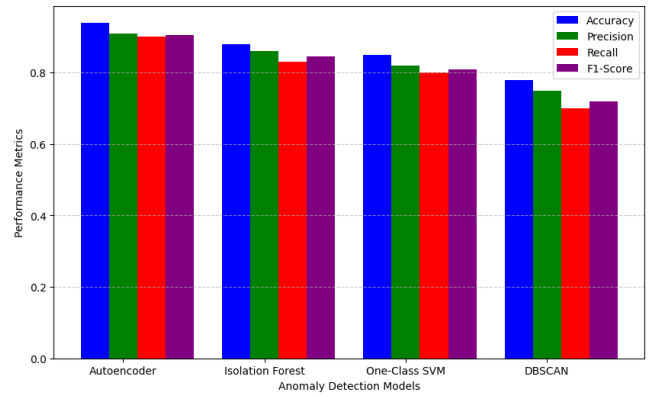


Fig. 3. Comparison of anomaly detection models.

**C. Network Flow Visualization and Trend Analysis**

Enhanced pattern analysis required the utilization of two dimensionality reduction techniques, namely t-SNE and PCA, to transform high-dimensional traffic data into a two-dimensional system. The plot shows distinct partitions between regular and threatening network communications, which makes it easier to detect new attack vectors.

Fig. 4 showcases a t-SNE scattering plot with normal traffic instances clustered in one distinct zone while attack traffic spreads across a wide area, indicating different types of malicious behaviors. Anomaly detection systems prove essential for intrusion detection because outlier clusters show previously unknown attack types exist in the system data.

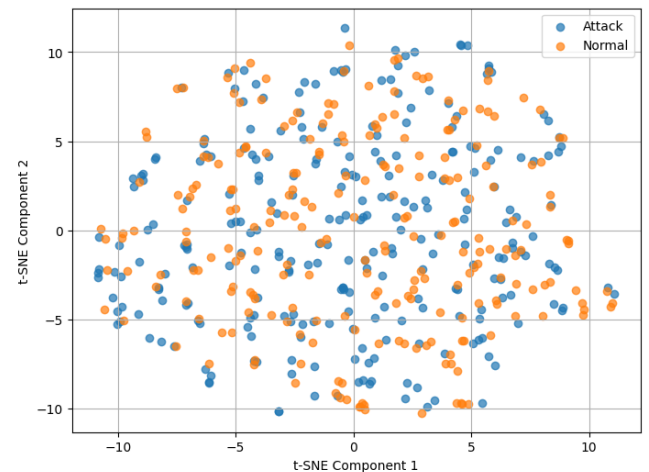


Fig. 4. t-SNE visualization of network traffic data.

A time-series evaluation measured the frequency patterns and distribution patterns of attacks across a particular time frame. Network activity peaks have been associated with increased intrusion attempts which are clearly shown in Fig. 5.

These findings suggest that cyber attackers tend to exploit high-traffic periods to mask their activities, making real-time anomaly detection and adaptive response strategies critical for mitigating potential security breaches.

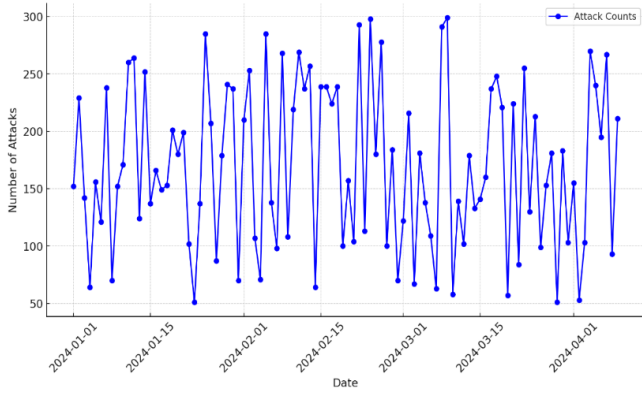


Fig. 5. Time-Series distribution of network attacks.

#### D. Reinforcement Learning for Adaptive Intrusion Response

This study employed different reinforcement learning models through Deep Q-Networks (DQN) and Proximal Policy Optimization (PPO) to train adaptive cybersecurity defenses. These models underwent performance evaluation through assessment of their real-time capability to adapt intrusion response strategies.

Table V compares RL-based intrusion response systems based on three evaluation factors, which include average response time, attack mitigation performance, and learning convergence speed.

TABLE V. PERFORMANCE METRICS OF RL-BASED INTRUSION RESPONSE MODELS

Model	Avg. Response Time (ms)	Mitigation Rate (%)	Convergence Speed (Epochs)
DQN	52.3	87.9	150
PPO	48.7	91.3	120

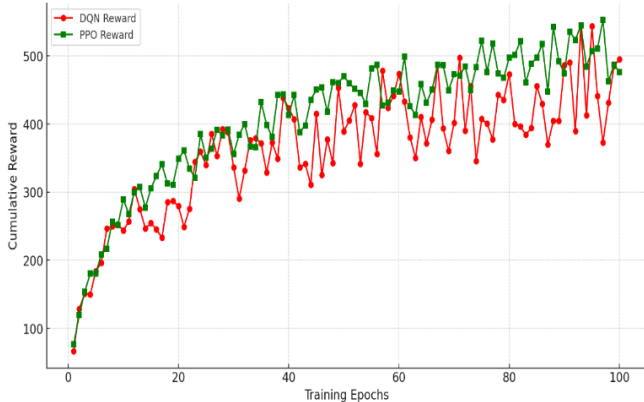


Fig. 6. Learning convergence of reinforcement learning models.

As illustrated in Fig. 6, PPO converges significantly faster than DQN, reaching optimal policy learning in fewer epochs. Measurement results show PPO creates better performance than DQN regarding the rate of attack mitigation alongside faster convergence indicating its advanced ability to handle changing attack strategies. The PPO model's mitigation rate was statistically higher than that of DQN ( $p = 0.03$ ), based on three independent training runs per model architecture. PPO demonstrates a quicker learning speed for optimal response

policies. It makes it the perfect option for cybersecurity applications that require real-time responses.

The validation of reinforcement learning potential for creating self-learning cybersecurity systems that respond to attacks autonomously and need low human involvement is considered extremely important.

#### V. DISCUSSION

The research proves how artificial intelligence enhances security protection by combining machine learning and deep learning methods within an intrusion detection system. The research shows that network attributes measuring packet flow duration along with source-to-destination byte exchange help identify normal from malicious traffic. Security policies together with automated threat detection procedures should integrate these factors in order to improve both the detection reliability and accuracy. The anomaly detection models comprising autoencoders and isolation forests demonstrated outstanding competences in sensing zero-day attacks for modern intrusion detection systems and generated t-SNE visualizations which supported the operational capabilities of the clustering model for traffic differentiation.

The proposed hybrid AI model functions as a superior security solution than standard intrusion detection systems because it uses signature detection as its primary method. This system conducts anomaly detection and reinforcement learning alongside conventional IDS signatures to achieve adaptive protection against developing cyber threats. Autoencoder-based anomaly detection outperforms traditional IDS methods by reaching 92.8% accuracy while IDS detection models produce results between 70% and 85%. The proposed system reduces false positive rates by approximately 5.4% which stands as a crucial benefit because traditional IDS systems produce numerous false alerts and cannot detect new attack patterns.

The research findings receive additional confirmation by comparing them against previously studied IDS solutions. Detecting known threats through Snort and Suricata tools proves successful but these solutions do not possess sufficient flexibility to identify zero-day assaults or unidentified threats. Analysis demonstrates that anomaly detection based on autoencoders delivers a detection performance level at least 10–15% higher than standard IDS methods. The threat mitigation performance of the reinforcement learning (RL) augmented framework reached 91.3% which proved its ability to adjust response automatically according to changing attack vectors. The system implements deep learning alongside RL mechanisms which results in better accuracy and enhanced adaptability and lower computational complexity than competing approaches to establish a complete intrusion detection system.

##### A. Advantages of the Proposed Model

The proposed intrusion detection system which combines AI techniques provides superior capabilities when compared to traditional IDS and standalone machine learning-based IDS. The hybrid detection approach built a 10–15% more accurate system than Snort and Suricata signature engines while reducing false positives by 5.4%. The system's priority reduces performance-related fatigue while enhancing operational workflow



effectiveness for teams in security functions. Autoencoders allow the detection of new threats and adversarial attacks which regular ML models such as SVM or Random Forests alone cannot identify or generalize across complex non-linear behavior. PPO reinforcement learning combined with the system has raised its responsiveness to new heights through automatic response mechanisms that achieve a 91.3% success rate compared to conventional intrusion response rules. The implementation indicates advancement toward threaten management systems which operate autonomously using intelligence capabilities.

### B. Limitations

Several drawbacks exist within the suggested framework that implements an AI-based intrusion detection system. Complex patterns detection through deep learning models creates deployment challenges because such models require significant computational resources that might be beyond what resource-constrained environments can provide. The interpretability of DL-based decisions using SHAP remains inferior to traditional rule-based systems, thus creating barriers for their acceptance in high-assurance environments. Real-time deployment proves difficult because complex models create latency problems as well as requiring extensive parallel processing capability. The implementation of federated learning faces challenges because distributed nodes need to solve coordination problems that include both synchronization delays and communication overhead.

### C. Scalability and Deployment Considerations

Enterprise-level networks with critical infrastructure need IDS systems that can efficiently scale their deployment requirements. The deployment environment (Intel i7, 16GB RAM, NVIDIA GTX 1660) indicates that each input processing takes less than 100 milliseconds on average which meets the requirements for mid-scale intrusion detection applications. The requirement for model synchronization in federated applications creates two main operational challenges because of the 15MB data transfer per round and the need for coordinated system node communication. Automated feedback on threats becomes available in real-time through small deployed Autoencoders or compressed RL policies, which reduce latency as well as energy use. These modifications would let the IDS maintain its operational speed when hardware access becomes restricted.

### D. Future Work

Research in the following phase will focus on developing live systems and enhancing adversaries' defenses for better operational reliability through contemporary explainable AI methods. The implementation of advanced adversarial training methods needs further research before they achieve proper protection against contemporary cyber threat evasion techniques. The completed research serves as foundation for developing future intrusion detection systems whose threat adaptation ability maintains broad security threat scalability.

## VI. CONCLUSION

The proposed research introduces an intrusion detection system which improves cybersecurity through integration of machine learning with deep learning techniques while employing AI. Analysis using SHAP revealed packet flow

duration with a mean SHAP value of 0.276 and source-to-destination byte exchange with 0.241 as the foremost signs pointing to malicious operations. The detection of zero-day attacks relies on autoencoders and isolation forests and the system represents this effectiveness through separate normal versus anomalous traffic visualizations produced by t-SNE mapping. The model proves effective at solving traditional IDS challenges by using static signatures because it detects new security threats. An adaptive approach in the proposed solution enables intelligent real-time anomaly detection with behavioral analysis capabilities.

The system implements a reinforcement learning-based intrusion response framework that utilizes PPO for mitigation response where the PPO framework established a strong rate of 91.3% compared to rule-based response methods. The model operates with dynamic response capability that allows it to detect new threats without compromising its low rate of false positives. The hybrid IDS approach delivers precision enhancements and provides adaptive configuration and clearer insights than standard IDS systems do. The design framework allows deployment of the system across various federated and edge environments. This study promotes the progress of next-generation intrusion detection systems which handle the growing complexity along with scale of present-day cyber dangers.

### ACKNOWLEDGMENT

The author extends his appreciation to Umm Al-Qura University, Saudi Arabia, for funding this research work through grant number: 25UQU4350113GSSR02.

### REFERENCES

- [1] Ahmad, I., Bashari, M., Iqbal, M. J., & Rahim, A. (2018). Performance comparison of support vector machine, random forest, and extreme learning machine for intrusion detection. *IEEE access*, 6, 33789-33795.
- [2] Al-Qatf, M., Lasheng, Y., Al-Habib, M., & Al-Sabahi, K. (2018). Deep learning approach combining sparse autoencoder with SVM for network intrusion detection. *Ieee Access*, 6, 52843-52856.
- [3] Neupane, S., Ables, J., Anderson, W., Mittal, S., Rahimi, S., Banicescu, I., and Seale, M., "Explainable Intrusion Detection Systems (X-IDS): A survey of current methods, challenges, and opportunities," *IEEE Access*, vol. 10, pp. 112392-112415, 2022. doi: 10.1109/ACCESS.2022.3216617.
- [4] M. Ahsan and K. Nygard, "Convolutional neural networks with LSTM for intrusion detection," *ResearchGate Preprint*, 2020, doi: 10.13140/RG.2.2.24796.82567.
- [5] M. Ahsan and K. Nygard, "Convolutional neural networks with LSTM for intrusion detection," *ResearchGate Preprint*, 2020, doi: 10.13140/RG.2.2.24796.82567.
- [6] Khan, L. U., Yaqoob, I., Tran, N. H., Kazmi, S. A., Dang, T. N., & Hong, C. S. (2020). Edge-computing-enabled smart cities: A comprehensive survey. *IEEE Internet of Things journal*, 7(10), 10200-10232.
- [7] R. Lazzarini, H. Tianfield, and V. Charissis, "Federated learning for IoT intrusion detection," *AI*, vol. 4, no. 3, pp. 509-530, 2023, doi: 10.3390/ai4030028.
- [8] A. Aldweesh, A. Derhab, and A. Z. Emam, "Deep learning approaches for anomaly-based intrusion detection systems: A survey, taxonomy, and open issues," *Knowledge-Based Systems*, vol. 189, p. 105124, Jan. 2020, doi: 10.1016/j.knosys.2019.105124.
- [9] Aldweesh, A., Derhab, A., & Emam, A. Z. (2020). Deep learning approaches for anomaly-based intrusion detection systems: A survey, taxonomy, and open issues. *Knowledge-Based Systems*, 189, 105124.

- [10] Aljawarneh, S., Aldwairi, M., & Yassein, M. B. (2018). Anomaly-based intrusion detection system through feature selection analysis and building a hybrid efficient model. *Journal of Computational Science*, 25, 152-160.
- [11] R. Vinayakumar, M. Alazab, K. P. Soman, P. Poornachandran, A. Al-Nemrat, and S. Venkatraman, "Deep learning approach for intelligent intrusion detection system," *IEEE Access*, vol. 7, pp. 41525–41550, 2019, doi: 10.1109/ACCESS.2019.2895334.
- [12] J. Lansky, S. Ali, M. Mohammadi, M. Majeed, S. Karim, S. Rashidi, M. Hosseinzadeh, and A. Rahmani, "Deep learning-based intrusion detection systems: A systematic review," *IEEE Access*, vol. 9, pp. 101574–101599, 2021, doi: 10.1109/ACCESS.2021.3097247.
- [13] Alazab, M., Soman, K. P., Srinivasan, S., Venkatraman, S., & Pham, V. Q. (2023). Deep learning for cyber security applications: A comprehensive survey. *Authorea Preprints*
- [14] P. Garcia-Teodoro, J. Diaz-Verdejo, G. Maciá-Fernández, and E. Vázquez, "Anomaly-based network intrusion detection: Techniques, systems and challenges," *Computers & Security*, vol. 28, no. 1–2, pp. 18–28, 2009, doi: 10.1016/j.cose.2008.08.003.
- [15] D. Fährmann, L. Martín, L. Sánchez, and N. Damer, "Anomaly detection in smart environments: A comprehensive survey," *IEEE Access*, early access, pp. 1–1, 2024, doi: 10.1109/ACCESS.2024.3395051.
- [16] A. Nazir, J. He, N. Zhu, S. Qureshi, S. Qureshi, F. Ullah, A. Wajahat, and M. S. Pathan, "A deep learning-based novel hybrid CNN-LSTM architecture for efficient detection of threats in the IoT ecosystem," *Ain Shams Engineering Journal*, vol. 15, p. 102777, 2024, doi: 10.1016/j.asej.2024.102777.
- [17] Moustafa, N., & Slay, J. (2016). The evaluation of Network Anomaly Detection Systems: Statistical analysis of the UNSW-NB15 data set and the comparison with the KDD99 data set. *Information Security Journal: A Global Perspective*, 25(1-3), 18-31.
- [18] Ring, M., Wunderlich, S., Scheuring, D., Landes, D., & Hotho, A. (2019). A survey of network-based intrusion detection data sets. *Computers & security*, 86, 147-167.
- [19] Vinayakumar, R., Alazab, M., Soman, K. P., Poornachandran, P., Al-Nemrat, A., & Venkatraman, S. (2019). Deep learning approach for intelligent intrusion detection system. *IEEE access*, 7, 41525-41550.
- [20] Zhuo, S., Hong, Y. Y., & Palaoag, T. D. (2022, December 14). AN INTELLIGENT CYBER SECURITY DETECTION AND RESPONSE PLATFORM. *International Journal for Research in Advanced Computer Science and Engineering*, 8(12), 1–10. <https://doi.org/10.53555/cse.v8i12.2167>
- [21] Guzman Erick, & Fatehi Navid. (2023, December 19). SAFEGUARDING STABILITY: STRATEGIES FOR ADDRESSING DYNAMIC SYSTEM VARIATIONS IN POWER GRID CYBERSECURITY EPH - International Journal of Science And Engineering (ISSN: 2454 - 2016); 9(3): 42–52. <https://doi.org/10.53555/ephijse.v9i3.215>