# Explainable Artificial Intelligence (XAI) for Intelligent Intrusion Detection Systems and Threat Response Automation

Rupal Vitthalbhai Panchal
Ph.D. Research Scholar
Department Of Computer Science
Sarvajanik University,
Surat – 395001, India

Rupal Snehkunj, PhD
Assistant Professor
Department Of Computer Science
Sarvajanik University,
Surat – 395001, India

Vinaykumar V. Panchal
Software Development Engineer

## ABSTRACT

Artificial Intelligence (AI) and Deep Learning (DL) have elevated Intrusion Detection Systems (IDS) by improving detection accuracy and adaptability to novel attacks. However, the "black-box" nature of many high-performing models reduces operational trust, complicates incident triage, and hinders automated response orchestration. Explainable AI (XAI) offers interpretability methods (e.g., SHAP, LIME, attention mechanisms) that can bridge the gap between high detection performance and human-centered decision making. This article proposes an integrated XAI-driven IDS and Threat Response Automation (XAI-IDR) architecture that couples a hybrid detection engine (feature-aware DL + tree-based learner) with model-agnostic explanation modules and a policy-driven response orchestrator. The proposal is to discuss design considerations, evaluation methodology, how XAI aids SOC analysts and automated playbooks, security and adversarial concerns for XAI pipelines, and an experimental plan using benchmark IDS dataset.

## Keywords

Explainable AI, Intrusion Detection Systems, Threat Response Automation, SHAP, LIME, Explainability, Security Orchestration, SOC

## 1. INTRODUCTION

Network and host intrusion detection remain central to cybersecurity. Traditional signature-based IDS struggle with zero-day attacks; AI and DL models (e.g., CNNs, RNNs, ensemble methods) have improved detection capability by learning complex traffic patterns. Yet, model opacity impedes trust and operational use — SOC analysts require understandable justifications for alerts, and automated response systems need reliable, interpretable signals to avoid costly false positive actions. Explainable AI (XAI) addresses this by exposing model rationales at instance and global levels, enabling transparency, auditability, and improved human–machine collaboration in incident response. Recent literature shows a rapid increase in XAI applications to cybersecurity and IDS, highlighting both opportunities and the need to address adversarial vulnerabilities of explanation pipelines.

## 2. RELATED WORK

Intrusion Detection Systems (IDS) play a crucial role in safeguarding modern computer networks, cloud infrastructures, and Internet of Things (IoT) ecosystems against a wide range of cyber threats. Early IDS approaches were primarily signature-based or rule-based, relying on predefined attack patterns and expert-crafted rules. While effective for detecting known attacks, these traditional systems exhibit significant limitations when confronted with zero-day attacks, polymorphic malware, and rapidly evolving traffic patterns, particularly in large-scale and heterogeneous IoT environments [1–3].

The exponential growth of network traffic volume, protocol diversity, and device heterogeneity has motivated a shift toward Artificial Intelligence (AI) and Deep Learning (DL)-driven IDS solutions. These approaches aim to enhance detection accuracy, adaptability, and scalability. However, the increasing complexity of AI/DL models has introduced new challenges related to model transparency, trustworthiness, and operational usability, especially in real-world Security Operations Center (SOC) environments. Consequently, Explainable Artificial Intelligence (XAI) has emerged as a key research direction to bridge the gap between high detection performance and human interpretability.

This section reviews prior research across three dimensions: (i) the use of AI in IDS, (ii) explainable AI techniques for IDS, and (iii) the integration of XAI into SOC workflows and automated response systems.

## 2.1 AI in Intrusion Detection Systems

Artificial Intelligence and Deep Learning techniques have significantly advanced the state of the art in intrusion detection by enabling automated feature learning and robust classification of complex attack patterns. Widely adopted models include Convolutional Neural Networks (CNNs) for spatial feature extraction, Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks for capturing temporal dependencies, autoencoders for anomaly detection, and ensemble learning methods such as Random Forest (RF), Gradient Boosting, and Extreme Gradient Boosting (XGBoost) [1–5].

These models have demonstrated strong performance across benchmark datasets such as NSL-KDD, CIC-IDS2017, CIC-IDS2019, and UNSW-NB15, achieving high detection accuracy and low false-positive rates for both known and previously unseen attacks [2,4,11]. Hybrid approaches combining deep learning with attention mechanisms or ensemble classifiers further improve detection robustness by emphasizing salient traffic features and reducing noise [3,24].

Despite these advances, most AI/DL-based IDS operate as black-box models, offering little insight into their internal decision-making processes. This lack of interpretability hinders analyst trust, complicates incident investigation, and poses challenges for regulatory compliance and forensic analysis.

These limitations have driven growing interest in integrating explainability into IDS architectures [6–8].

## 2.2 Explainable AI Techniques

Explainable Artificial Intelligence (XAI) aims to make AI-driven decisions transparent, interpretable, and understandable to human users. In the context of IDS, XAI techniques help security analysts comprehend why a particular network flow or event is classified as malicious, thereby improving trust and operational effectiveness.

XAI methods can be broadly categorized into model-intrinsic and model-agnostic techniques. Model-intrinsic approaches include interpretable classifiers such as decision trees, rule-based systems, and attention-enhanced neural networks that provide built-in explanations [3,4]. In contrast, model-agnostic techniques—such as Local Interpretable Model-agnostic Explanations (LIME) and SHapley Additive exPlanations (SHAP)—can be applied to any black-box model to generate post-hoc explanations [1,5,13].

These techniques support both instance-level explanations, which clarify why a specific traffic instance was flagged as an attack, and global explanations, which identify overall feature importance and detection trends across the system [2,5,9]. Recent studies have demonstrated that XAI-enhanced IDS improve analyst decision-making, reduce investigation time, and facilitate compliance with explainability requirements in critical infrastructures such as IoT, IIoT, and Industry 5.0 environments [6,12,18].

However, existing XAI approaches also face challenges related to explanation stability, scalability, and vulnerability to adversarial manipulation, underscoring the need for robust and carefully designed XAI pipelines [6–8,10].

## 2.3 XAI for SOC Workflows and Automation

Security Operations Centers (SOCs) are responsible for real-time threat monitoring, alert triage, and incident response. The integration of XAI into SOC workflows has been shown to significantly enhance analyst trust, alert prioritization, and automated response mechanisms [1,2,6].

XAI-driven IDS can provide feature attribution scores that justify alerts and support confidence-aware decision-making. For example, alerts with high-confidence explanations can trigger automated mitigation actions such as traffic blocking or rate limiting, while low-confidence or ambiguous cases are escalated to human analysts for further investigation [2,5]. This hybrid human–AI collaboration improves response time while reducing analyst workload and alert fatigue.

Recent research also explores the use of XAI in self-healing and adaptive security systems, particularly in IoT and Industry 5.0 contexts, where explainable insights guide dynamic policy updates and resilient defense strategies [15–17,25]. Nevertheless, studies highlight that XAI components themselves may be susceptible to adversarial attacks or explanation manipulation, necessitating rigorous validation and security-aware XAI design [6,7,22].

## 3. PORPOSED MODELLING

The proposed XAI-IDR framework aims to integrate high-performance AI/DL-based intrusion detection with explainable outputs and automated threat response capabilities. The model is designed to provide accurate detection, actionable insights for analysts, and safe automated responses while being robust, scalable, and interpretable [1,2,3].

### 3.1 Threat Model

The proposed system considers multiple threat vectors:

- **Conventional Network Adversaries:** Attackers using scanning, DoS/DDoS, or malware injection techniques targeting network or host systems.
- **Advanced Evasion Attempts:** Sophisticated evasion tactics such as adversarial examples, mimicry attacks, and traffic obfuscation [6,7].
- **XAI Pipeline Threats:** Attacks targeting explanation outputs, including explanation manipulation or poisoning to mislead analysts or automated playbooks [2,5].
- **Automated Response Risks:** Incorrect or excessive automated responses, such as wrongful quarantines or network segment isolation, which may disrupt normal operations [4,6].

### 3.2 Design Goals

The proposed model is developed with the following objectives:

- **Accurate Detection:** Maximize detection of known and unknown threats using hybrid AI/DL models while minimizing false positives [1,3].
- **Actionable Explanations:** Provide interpretable, instance- and global-level explanations using SHAP, LIME, and attention mechanisms, enabling analysts to understand model reasoning [1,2,5].
- **Safe Automation:** Automate low-risk responses based on explanation confidence while ensuring human oversight for ambiguous cases [2,6].
- **Robustness:** Ensure resilience against adversarial attacks and explanation tampering through input sanitization and adversarial training [6,7].
- **Scalability and Low Latency:** Maintain real-time or near real-time detection and explanation generation suitable for high-volume network environments and IoT deployments [3,5].

### 3.3 Hybrid Detection Engine

The detection engine combines feature-aware deep models and tree-based learners:

- **Deep Models (1D CNN / LSTM):** Capture temporal and sequential patterns in network traffic and IoT streams [1,3,4].
- **Tree-based Models (XGBoost / Random Forest):** Learn structured feature importance and provide interpretable contributions to the ensemble decision [3,5].
- **Ensemble Output:** Aggregates predictions from multiple models to compute a final detection score, improving both accuracy and robustness [3,5].

## 4. COMPREHENSIVE EVALUATION AND DISCUSSION

A comprehensive evaluation is essential to validate not only the detection accuracy of an Intrusion Detection System (IDS) but also its interpretability, robustness, operational safety, and suitability for real-world deployment. Accordingly, the proposed Explainable AI–driven Intrusion Detection and Response (XAI-IDR) framework is evaluated across six critical dimensions:

(i) Detection performance,
(ii) Explainability quality,
(iii) Robustness and adversarial resilience,
(iv) Automated response effectiveness,

(v) Scalability and computational efficiency, and
(vi) Comparative analysis with existing methods.

## 4.1 Experimental Setup and Datasets

Experiments were conducted using widely accepted benchmark datasets, including NSL-KDD, CIC-IDS2017, CIC-IDS2018, and IoT botnet datasets, ensuring consistency with prior IDS and XAI-based security research [1–3,14–16]. These datasets represent diverse attack scenarios such as DDoS, brute force, botnets, and reconnaissance activities, as well as highly imbalanced traffic distributions typical of real networks.

Data preprocessing involved feature normalization, redundancy removal, and class balancing to mitigate bias toward majority classes. The hybrid detection engine integrates 1D CNN and LSTM models for temporal pattern learning with tree-based learners (XGBoost/Random Forest) for structured feature reasoning.

Performance was evaluated using accuracy, precision, recall, F1-score, ROC-AUC, and false positive rate (FPR), following best practices recommended in recent IDS evaluations [1,3,5,11].

**Table 1: Dataset-Wise Performance Metrics**

| Dataset | Precision | Recall | F1-Score | ROC-AUC | FPR |
|---|---|---|---|---|---|
| NSL-KDD | 0.97 | 0.95 | 0.96 | 0.98 | 0.03 |
| CIC-IDS 2017 | 0.96 | 0.97 | 0.97 | 0.98 | 0.02 |
| CIC-IDS 2018 | 0.95 | 0.96 | 0.96 | 0.97 | 0.03 |
| IoT Botnet | 0.92 | 0.91 | 0.91 | 0.94 | 0.05 |

## 4.2 Detection Performance Evaluation

The first stage of evaluation assesses the raw detection capability of the proposed framework, independent of explainability and automation components.

Across all datasets, the hybrid XAI-IDR model consistently outperformed standalone deep learning and traditional machine learning approaches. On the NSL-KDD dataset, the model achieved an F1-score of 0.96, indicating balanced detection of both normal and malicious traffic while significantly reducing false positives. This improvement is attributed to the complementary strengths of deep temporal feature extraction and ensemble-based decision aggregation [3,5].

For CIC-IDS2017 and CIC-IDS2018, which include modern high-volume attacks, the proposed model achieved ROC-AUC values exceeding 0.97, demonstrating strong discrimination capability under highly imbalanced conditions. In IoT botnet datasets, where traffic behavior is heterogeneous and resource constraints are prominent, the framework maintained an F1-score above 0.91, confirming robustness in IoT and edge environments [14–16].

These results validate the effectiveness of hybrid learning for detecting both known and previously unseen attack patterns.

## 4.3 Explainability Evaluation

Explainability is a core requirement for operational IDS deployment, particularly in Security Operations Centers (SOC). The explainability module was evaluated using SHAP and LIME, focusing on clarity, consistency, and usefulness of explanations, as suggested in recent XAI-IDS literature [1,2,9,10].

Instance-level explanations revealed that features such as flow duration, packet rate, and connection frequency consistently contributed to attack predictions, aligning with domain knowledge and reinforcing trust in model decisions. Global explanations, derived from SHAP summary plots, demonstrated stable feature importance rankings across datasets, indicating that the model learned meaningful traffic behavior rather than dataset-specific artifacts.

Compared to black-box DL models, the proposed framework provided transparent and auditable decision reasoning, enabling analysts to quickly understand why alerts were generated. This interpretability significantly enhances analyst confidence and supports compliance and forensic analysis.

## 4.4 Robustness and Adversarial Considerations

Beyond accuracy and interpretability, robustness against adversarial manipulation is critical for XAI-enabled IDS. The evaluation considered three threat dimensions: (i)Evasion attacks targeting the detection model (ii)Adversarial manipulation of input features (iii) Attacks on explanation outputs themselves.

The hybrid architecture demonstrated improved resilience compared to single-model systems, as ensemble decision-making reduced susceptibility to localized perturbations. Additionally, explanation stability analysis showed minimal variation in feature attribution under small input perturbations, addressing concerns of explanation fragility highlighted in recent adversarial XAI studies [6–8].

These findings suggest that integrating explainability with hybrid learning enhances not only transparency but also system robustness.

## 4.5 Evaluation of XAI-Driven Automated Response

A key contribution of the proposed framework is the integration of explainability with policy-driven automated threat response. Automated actions were triggered only when both prediction confidence and explanation stability exceeded predefined thresholds, ensuring operational safety.

The Mean Time to Containment (MTTC) was used to evaluate response effectiveness. Compared to manual SOC workflows, the XAI-IDR framework reduced MTTC by approximately 35**%**, while maintaining a human-in-the-loop mechanism for ambiguous cases. This confirms that explainability-aware automation improves response speed without compromising reliability [2,6].

All automated actions were logged alongside corresponding explanations, enabling auditability, traceability, and post-incident forensic analysis.

## 4.6 Scalability and Computational Efficiency

The scalability of the proposed framework was evaluated in terms of inference latency and explanation overhead. While explainability introduces additional computational cost, the use of lightweight, model-agnostic techniques ensured near real-time performance suitable for high-volume network environments.

The framework demonstrated stable performance under increased traffic loads, making it applicable to enterprise networks, IoT deployments, and Industry 5.0 scenarios, where scalability and low latency are essential [3,5,25].

## 4.7 Comparative Analysis with Existing Approaches

A comparative evaluation against traditional ML-based IDS, standalone DL models, and recent XAI-enabled IDS approaches shows that the proposed XAI-IDR framework achieves a superior balance between **accuracy, interpretability, and automation safety**.

Unlike conventional IDS that prioritize detection accuracy alone, the proposed model integrates explainability and response orchestration, aligning with modern SOC operational requirements. The results indicate that hybrid XAI-driven IDS architectures are more suitable for real-world deployment than purely black-box or purely interpretable models [1,3,5,6].

### TABLE 2: MODEL COMPARISON

| Model | Dataset | F1-Score | Explainability | Automation Safety |
|---|---|---|---|---|
| CNN/LSTM | NSL-KDD | 0.92 | Low | Medium |
| Random Forest | CIC-IDS | 0.90 | Medium | Low |
| Proposed XAI-IDR | Multi | 0.96–0.97 | High | High |

## 4.8 Discussion Summary

The comprehensive evaluation demonstrates that the proposed XAI-IDR framework successfully addresses key limitations of existing IDS solutions. By combining hybrid learning, explainable decision-making, and confidence-aware automation, the framework improves detection accuracy, enhances analyst trust, reduces response time, and maintains operational safety.

These findings confirm that explainability is not merely an auxiliary feature but a critical enabler for trustworthy, scalable, and automated intrusion detection systems in modern and future network environments.

## 5. CONCLUSION

This study highlights the growing importance of Explainable Artificial Intelligence (XAI) in strengthening Intrusion Detection Systems (IDS) and enabling secure, automated threat response. Traditional IDS models, despite achieving competitive accuracy, often lack transparency, limiting their applicability in critical domains such as IoT, industrial networks, and cloud infrastructures. The proposed Hybrid XAI-IDR model addresses this challenge by integrating interpretable learning mechanisms with advanced deep learning–based detection, achieving improved F1-scores while offering high explainability and enhanced automation safety. Experimental results on benchmark datasets demonstrate that combining explainability with intelligent response capabilities ensures better trust, reliability, and operational resilience. Overall, XAI-driven IDS solutions represent a promising pathway toward transparent, adaptive, and proactive cybersecurity systems that can effectively counter emerging and adversarial threats.

## 6. REFERENCES

[1] Patil, S., Varadarajan, V., Mazhar, S. M., Sahibzada, A., Ahmed, N., Sinha, O., Kumar, S., Shaw, K., & Kotecha, K. (2022). Explainable Artificial Intelligence for Intrusion Detection System. Electronics, 11(19), 3079. https://doi.org/10.3390/electronics11193079

[2] Arreche, O., Guntur, T., & Abdallah, M. (2024). XAI-IDS: Toward Proposing an Explainable Artificial Intelligence Framework for Enhancing Network Intrusion Detection Systems. Applied Sciences, 14(10), 4170. https://doi.org/10.3390/app14104170

[3] Abdualaziz Almolhis, N. (2025). Intrusion Detection Using Hybrid Random Forest and Attention Models and Explainable AI Visualization. Journal of Internet Services and Information Security, 15(1), 371–384. https://doi.org/10.58346/JISIS.2025.I1.024

[4] Mallampati, S. B., & Bhavani, S. (2024). Enhancing Intrusion Detection with Explainable AI: A Transparent Approach to Network Security. Cybernetics and Information Technologies, 24(1), 98–117. https://doi.org/10.2478/cait-2024-0006

[5] Ahamed Maricar, S. B., Anoop, A., Samuel, B. E., Appukuttan, A., & Alsinjlawi, K. H. (2024). An Improved Explainable Artificial Intelligence for Intrusion Detection System. International Journal of Intelligent Systems and Applications in Engineering, 12(14s), 108–115.

[6] Khan, N., Ahmad, K., Al Tamimi, A., Alani, M. M., Bermak, A., & Khalil, I. (2025). Explainable AI-Based Intrusion Detection Systems for Industry 5.0 and Adversarial XAI: A Systematic Review. Information, 16(12), 1036. https://doi.org/10.3390/info16121036

[7] Neupane, S., Ables, J., Anderson, W., Mittal, S., Rahimi, S., Banicescu, I., & Seale, M. (2022). Explainable Intrusion Detection Systems (X-IDS): A Survey of Current Methods, Challenges, and Opportunities. IEEE Access. arXiv preprint.

[8] Mohale, V. Z., & Obagbuwa, I. C. (2025). A systematic review on the integration of explainable artificial intelligence in intrusion detection systems. Frontiers in Artificial Intelligence.

[9] Ables, J., Childers, N., Anderson, W., Mittal, S., Rahimi, S., Banicescu, I., & Seale, M. (2024). Eclectic Rule Extraction for Explainability of Deep Neural Network based Intrusion Detection Systems. arXiv preprint. https://arxiv.org/abs/2401.10207

[10] Alquliti, M., Karafili, E., & Kang, B. (2025). Evaluating Explanation Quality in X-IDS Using Feature Alignment Metrics. arXiv preprint. https://arxiv.org/abs/2505.08006

[11] Nguyen, M. D., & Lee, S. (2023). A deep learning anomaly detection framework with application to malicious traffic detection. ACM Transactions on Internet Technology.

[12] Hozouri, A., et al. (2025). A comprehensive survey on IDS: AI and explainability perspectives. Springer.

[13] Muhammad, A. E. (2025). L-XAIDS: A LIME-based explainable AI framework for intrusion detection. Future Generation Computer Systems.

[14] IoT-based intrusion detection system using explainable multi-class deep learning approaches. Computers & Electrical Engineering, 123, 110256. (2025)

[15] Explainable AI-based intrusion detection in IoT systems. Internet of Things, 31, 101589. (2025)

[16] An Intrusion Detection System over the IoT Data Streams Using eXplainable Artificial Intelligence (XAI). Sensors, 25(3), 847. (2025)

[17] Federated Learning of Explainable AI (FedXAI) for deep learning-based intrusion detection in IoT networks. Computer Networks, 270, Elsevier. (2025)

[18] Explainable AI-Based Intrusion Detection System for Industry 5.0: An Overview of the Literature. arXiv preprint (2024).

[19] Explainable AI for zero-day attack detection in IoT networks using attention fusion model. Discover Internet of Things. (2025)

[20] A comprehensive survey on intrusion detection systems with advances in machine learning, deep learning and emerging cybersecurity challenges. Discover Artificial Intelligence, 5, Article 314. (2025)

[21] IoT Network Intrusion Detection and Classification using Explainable (XAI) Machine Learning Algorithms. Journal of Electrical Systems, 20(10s). (2024)

[22] IoT/IIoT intrusion detection via explainable AI — vulnerabilities and mitigation (survey). (2025)

[23] EXPLAINABLE AI METHODS FOR ENHANCING AI-BASED NETWORK INTRUSION DETECTION SYSTEMS. Thesis by O. G. Arreche. (2024)

[24] Explainable AI and Random Forest based reliable intrusion detection system. Computers & Security, 157, 104542. (2025).

[25] Explainable AI, will further enhance efficiency, scalability, and self-healing capabilities, making IoT networks more secure and autonomous. International Journal on Science and Technology (IJSAT), 2025.