

# EXPLAINABLE AI-BASED INTRUSION DETECTION FRAMEWORK FOR CRITICAL INFRASTRUCTURE PROTECTION IN PAKISTAN'S DIGITAL ECOSYSTEM

Usman Ehsan<sup>\*1</sup>, Muhammad Atif Altaf<sup>2</sup>

<sup>\*1</sup>University institute of IT, Pir Mehr Ali Shah Arid Agriculture University

<sup>2</sup>Lecturer, Department of Information Sciences, University of Education, Lahore

<sup>1</sup>usmanehsan0613@gmail.com, <sup>2</sup>atifaltaf@ue.edu.pk

DOI: <https://doi.org/10.5281/zenodo.20118060>

## Keywords

Explainable Artificial Intelligence (XAI); Intrusion Detection System (IDS); Cybersecurity; Critical Infrastructure; Machine Learning; Deep Learning; Pakistan Digital Ecosystem; Network Security; Threat Detection; Artificial Intelligence Security Systems.

## Article History

Received: 14 March 2026

Accepted: 22 April 2026

Published: 11 May 2026

Copyright @Author

Corresponding Author: \*

Usman Ehsan

## Abstract

The rapid expansion of digital infrastructure in Pakistan has increased the vulnerability of critical systems such as banking, healthcare, energy, telecommunications, and government services to sophisticated cyber threats. Traditional intrusion detection systems (IDS) are increasingly insufficient due to their limited adaptability, high false-positive rates, and inability to detect zero-day attacks. Although Artificial Intelligence (AI) and Machine Learning (ML)-based IDS models have improved detection accuracy, their "black-box" nature limits transparency, trust, and operational acceptance in critical infrastructure environments. To address these challenges, this study proposed an Explainable AI-Based Intrusion Detection Framework designed to enhance both cybersecurity performance and interpretability. The framework integrated advanced machine learning algorithms with Explainable Artificial Intelligence (XAI) techniques, including SHAP and LIME, to provide transparent and interpretable threat detection. The model was evaluated using benchmark datasets and expert assessments from cybersecurity professionals. Results demonstrated that the proposed framework achieved superior performance in terms of accuracy, precision, recall, and false-positive reduction compared to traditional models. Additionally, expert evaluations confirmed high levels of interpretability, transparency, and trust in AI-driven decisions. The findings highlight that integrating explainability into IDS frameworks significantly strengthens cybersecurity resilience and supports informed decision-making in critical infrastructure environments. The study concludes that XAI-based intrusion detection systems offer a reliable and scalable solution for protecting Pakistan's evolving digital ecosystem.

## INTRODUCTION

The rapid digital transformation of critical infrastructure systems has significantly increased the exposure of national cyber ecosystems to sophisticated cyber threats. Critical infrastructure sectors, including energy, banking, transportation, healthcare, telecommunications, and government information systems, are increasingly dependent

on interconnected digital technologies, cloud platforms, Industrial Internet of Things (IIoT), and smart networks. While these technological advancements improve operational efficiency and service delivery, they also expand the attack surface for cybercriminals, advanced persistent threats (APTs), ransomware attacks, and state-sponsored cyber intrusions. In developing countries such as Pakistan, the accelerated adoption of digital services, smart governance

initiatives, and fintech ecosystems has created urgent cybersecurity challenges that require intelligent, adaptive, and trustworthy security frameworks. Recent cyber incidents targeting public and private digital infrastructures have demonstrated that conventional security mechanisms are insufficient against evolving attack vectors and zero-day vulnerabilities (Ahmed et al., 2024).

Intrusion Detection Systems (IDSs) have emerged as a fundamental component of modern cybersecurity architectures because of their ability to monitor network traffic, identify malicious activities, and generate real-time threat alerts. Traditional signature-based IDSs are effective against known attacks but often fail to detect novel and polymorphic threats due to their reliance on predefined attack signatures. Consequently, Artificial Intelligence (AI) and Machine Learning (ML)-driven IDS frameworks have gained considerable attention for their capability to analyze large-scale network traffic patterns, detect anomalies, and improve predictive threat intelligence. Deep learning algorithms, ensemble learning models, and hybrid AI approaches have demonstrated high detection accuracy in identifying malicious network behavior across cloud, IoT, and industrial control environments (Nugraha et al., 2025).

Despite the promising performance of AI-driven IDS models, the “black-box” nature of many machine learning algorithms remains a major concern in cybersecurity operations. Security analysts, critical infrastructure operators, and policy stakeholders often require transparent explanations regarding how and why an AI system classified a network event as malicious. The lack of interpretability reduces trust, accountability, and operational acceptance of AI-enabled cybersecurity systems, particularly in high-risk sectors where incorrect decisions may cause severe economic, operational, and national security consequences. Explainable Artificial Intelligence (XAI) has therefore emerged as a transformative paradigm that enhances transparency, interpretability, and human understanding of AI decision-making processes. XAI techniques such as SHAP (SHapley Additive exPlanations), LIME (Local Interpretable Model-Agnostic Explanations), feature importance analysis, and attention-based

visualization enable cybersecurity professionals to interpret threat classifications and validate system behavior effectively (Georgiades & Hussain, 2025).

Recent studies have highlighted the increasing importance of integrating XAI into intrusion detection frameworks for critical infrastructure protection. Researchers have proposed explainable hybrid ensemble models, interpretable deep learning mechanisms, and XAI-enabled anomaly detection systems to improve both cybersecurity performance and operational trustworthiness. These approaches not only improve detection accuracy but also support regulatory compliance, incident forensics, and strategic cyber risk management. Furthermore, explainable cybersecurity frameworks facilitate collaboration between human analysts and AI systems by providing understandable threat intelligence that enhances decision-making efficiency (Ciaramella et al., 2025).

The significance of XAI-based cybersecurity solutions is particularly critical in the context of Pakistan’s digital ecosystem. Pakistan has witnessed rapid growth in digital banking, e-governance, smart city initiatives, telecommunication infrastructure, and cloud-enabled public services. Simultaneously, the country faces increasing cyber risks due to limited cybersecurity awareness, shortage of skilled cybersecurity professionals, inadequate AI governance frameworks, and insufficient deployment of intelligent cyber defense mechanisms. Critical sectors such as banking systems, healthcare institutions, educational networks, and energy infrastructures remain vulnerable to cyberattacks capable of disrupting essential national services. Moreover, many organizations in Pakistan continue to rely on conventional reactive security approaches rather than intelligent predictive cybersecurity systems. This gap highlights the urgent need for a robust, explainable, and adaptive intrusion detection framework specifically designed for Pakistan’s critical infrastructure environment (Paulraj et al., 2025).

An Explainable AI-Based Intrusion Detection Framework can significantly contribute to strengthening cybersecurity resilience in Pakistan by combining intelligent threat detection with transparent decision-making

mechanisms. Such a framework can enhance stakeholder trust, improve incident response efficiency, reduce false positive rates, and provide interpretable cyber threat analytics for security professionals and policymakers. Additionally, integrating explainability into IDS architectures supports ethical AI adoption, regulatory transparency, and secure digital transformation initiatives. The proposed framework therefore aims to address both technological and operational limitations of existing IDS models by developing a transparent, intelligent, and scalable cybersecurity solution for protecting Pakistan's critical digital infrastructure against emerging cyber threats (Sukumaran & Korath, 2025).

### Problem Statement

The increasing digitization of critical infrastructure sectors in Pakistan, including banking, healthcare, telecommunications, transportation, energy, and e-governance systems, has significantly enhanced operational efficiency and digital connectivity. However, this rapid technological transformation has simultaneously increased vulnerability to sophisticated cyber threats such as malware attacks, ransomware, phishing, distributed denial-of-service (DDoS) attacks, insider threats, and advanced persistent threats (APTs). Traditional cybersecurity mechanisms and signature-based Intrusion Detection Systems (IDSs) are increasingly ineffective against modern dynamic and intelligent cyberattacks because they lack adaptability, predictive intelligence, and real-time threat interpretation capabilities. As cybercriminals continue to employ advanced attack strategies, conventional IDS frameworks face challenges related to low detection accuracy, high false-positive rates, poor scalability, and inability to identify zero-day attacks.

Artificial Intelligence (AI) and Machine Learning (ML)-based intrusion detection approaches have recently emerged as effective solutions for improving cyber threat detection and network anomaly analysis. These intelligent systems can process large volumes of network traffic data and identify complex malicious patterns more efficiently than traditional methods. Nevertheless, many AI-driven IDS models operate as "black-box" systems, where

the decision-making process remains unclear to cybersecurity analysts and infrastructure administrators. This lack of transparency and interpretability creates serious concerns regarding trust, accountability, reliability, and operational acceptance, particularly in critical infrastructure environments where inaccurate predictions may lead to severe financial, operational, and national security consequences. In the context of Pakistan's digital ecosystem, the challenge becomes more significant due to limited cybersecurity infrastructure, shortage of skilled cyber professionals, inadequate AI governance frameworks, and insufficient deployment of intelligent cyber defense systems. Critical institutions across Pakistan continue to rely heavily on reactive cybersecurity mechanisms rather than proactive and explainable security solutions. Furthermore, existing intrusion detection models are rarely designed according to the specific cybersecurity requirements and threat landscape of Pakistan's critical infrastructure sectors. The absence of explainable and interpretable AI-based intrusion detection frameworks restricts cybersecurity professionals from understanding the rationale behind attack classifications, thereby limiting effective incident response and strategic cyber risk management.

Therefore, there is a critical need to develop an Explainable AI-Based Intrusion Detection Framework that can provide accurate, transparent, interpretable, and scalable cyber threat detection for protecting Pakistan's critical digital infrastructure. The proposed framework aims to integrate explainable AI techniques with intelligent intrusion detection mechanisms to enhance threat detection accuracy, reduce false alarms, improve trust in AI-driven cybersecurity systems, and support effective decision-making for cybersecurity professionals and policymakers in Pakistan.

### Research Questions

1. How can Explainable Artificial Intelligence (XAI) improve the effectiveness of intrusion detection systems for critical infrastructure protection in Pakistan?
2. What are the major limitations of conventional and AI-based intrusion detection systems in Pakistan's digital ecosystem?

3. How can explainable AI techniques enhance transparency, interpretability, and trust in cybersecurity decision-making?
4. To what extent can the proposed Explainable AI-Based Intrusion Detection Framework improve detection accuracy and reduce false-positive rates?
5. How can the proposed framework contribute to strengthening cybersecurity resilience in Pakistan's critical infrastructure sectors?

### Research Objectives

#### General Objective

To develop an Explainable AI-Based Intrusion Detection Framework for enhancing cybersecurity protection of Pakistan's critical digital infrastructure.

#### Specific Objectives

1. To examine the existing cybersecurity challenges and intrusion detection limitations in Pakistan's critical infrastructure sectors.
2. To analyze the role of Artificial Intelligence and Explainable Artificial Intelligence in modern intrusion detection systems.
3. To design an explainable AI-based intrusion detection framework capable of identifying cyber threats with improved transparency and interpretability.
4. To evaluate the effectiveness of the proposed framework in terms of detection accuracy, false-positive reduction, and operational reliability.
5. To provide recommendations for implementing explainable and intelligent cybersecurity solutions for critical infrastructure protection in Pakistan.

#### Significance of the Study

This study is significant because it addresses the growing cybersecurity challenges faced by Pakistan's critical infrastructure sectors in the era of rapid digital transformation. As organizations increasingly adopt cloud computing, IoT technologies, smart systems, and digital services, the risk of sophisticated cyberattacks has become a major concern for national security, economic stability, and operational continuity. The proposed Explainable AI-Based Intrusion Detection

Framework offers an innovative and intelligent approach to strengthening cybersecurity resilience through accurate, transparent, and interpretable threat detection mechanisms.

The study contributes theoretically by expanding the existing body of knowledge on Explainable Artificial Intelligence (XAI) and AI-driven cybersecurity systems. It integrates explainability with intrusion detection models to overcome the limitations of conventional "black-box" AI systems, thereby promoting transparency, accountability, and trust in automated cybersecurity decision-making. The research also provides an academic foundation for future studies related to explainable cybersecurity frameworks, intelligent threat analytics, and critical infrastructure protection.

Practically, the proposed framework can assist cybersecurity professionals, network administrators, and critical infrastructure operators in detecting and interpreting cyber threats more effectively. By reducing false-positive rates and improving the interpretability of intrusion detection decisions, the framework can support faster incident response, improved threat analysis, and more reliable cybersecurity operations. The explainable nature of the system will further enable security analysts to understand attack patterns and make informed decisions during cyber incidents.

The study is also significant for policymakers and regulatory authorities in Pakistan, as it highlights the need for adopting transparent and trustworthy AI-driven cybersecurity solutions within national digital infrastructure strategies. The findings may support the development of cybersecurity policies, AI governance frameworks, and institutional cyber defense mechanisms aimed at protecting sensitive digital assets and critical public services.

Furthermore, this research is important from a socio-economic perspective because enhanced cybersecurity protection can strengthen public trust in digital systems, online banking, e-governance platforms, healthcare technologies, and other digital services. By safeguarding critical infrastructure against emerging cyber threats, the study contributes to sustainable digital transformation, technological resilience, and national cybersecurity preparedness in Pakistan.

### Literature Review

The increasing dependence on digital technologies, cloud computing, Industrial Internet of Things (IIoT), and smart communication systems has transformed critical infrastructure operations worldwide. However, this rapid digitalization has also increased exposure to sophisticated cyber threats targeting critical sectors such as banking, healthcare, transportation, telecommunications, and energy systems. Cyberattacks on critical infrastructure can disrupt essential services, compromise sensitive information, and create severe economic and national security consequences. Consequently, cybersecurity researchers and practitioners have focused extensively on developing intelligent Intrusion Detection Systems (IDSs) capable of identifying and mitigating evolving cyber threats in real time (Ahmed et al., 2024).

Traditional IDS models are generally categorized into signature-based and anomaly-based systems. Signature-based IDSs detect attacks using predefined attack patterns and known malware signatures. Although effective against previously identified threats, these systems fail to recognize zero-day attacks and sophisticated polymorphic malware. In contrast, anomaly-based IDSs identify deviations from normal network behavior and are capable of detecting unknown threats. However, conventional anomaly detection approaches often suffer from high false-positive rates and reduced detection precision in complex network environments (Nugraha et al., 2025). These limitations have encouraged researchers to integrate Artificial Intelligence (AI) and Machine Learning (ML) techniques into cybersecurity frameworks to enhance automated threat detection and predictive cyber defense capabilities.

Artificial Intelligence-based intrusion detection systems have demonstrated substantial improvements in cyber threat identification, classification accuracy, and real-time network monitoring. Machine learning algorithms such as Decision Trees, Random Forests, Support Vector Machines (SVM), Naïve Bayes, and K-Nearest Neighbor (KNN) have been widely used for detecting malicious network activities. Furthermore, deep learning approaches including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs),

Long Short-Term Memory (LSTM), and Autoencoders have significantly improved the ability of IDS models to detect complex attack patterns and hidden anomalies within large-scale network traffic datasets (Georgiades & Hussain, 2025).

Recent studies have emphasized the effectiveness of hybrid and ensemble AI models in intrusion detection systems. Ahmed et al. (2024) proposed an explainable hybrid ensemble model that combined multiple machine learning algorithms to improve attack classification accuracy and reduce false-positive alerts. The study demonstrated that hybrid AI models outperform traditional single-model approaches in terms of precision, scalability, and adaptability to evolving cyber threats. Similarly, Alabbadi and Bajaber (2025) developed an IoT-based intrusion detection system integrated with Explainable Artificial Intelligence (XAI) techniques to enhance interpretability and transparency in cyber threat detection. Their findings indicated that explainability mechanisms improve analysts' understanding of AI predictions and strengthen trust in automated security systems.

Despite the effectiveness of AI-driven IDS frameworks, the "black-box" nature of advanced machine learning and deep learning models remains a major challenge in cybersecurity operations. Many AI systems generate highly accurate predictions without providing understandable explanations regarding the rationale behind their decisions. This lack of interpretability creates concerns related to trust, accountability, ethical AI deployment, and operational reliability, particularly in critical infrastructure sectors where incorrect predictions may lead to severe consequences (Areghan & Ndibe, 2024). Security professionals often require transparent explanations to validate attack classifications, investigate incidents, and make strategic decisions during cyber emergencies.

To address these challenges, Explainable Artificial Intelligence (XAI) has emerged as an important research area in cybersecurity and intelligent decision-making systems. XAI focuses on developing transparent AI models capable of explaining their predictions in human-understandable forms. Techniques such as SHAP (SHapley Additive Explanations), LIME (Local Interpretable Model-Agnostic

Explanations), attention visualization, feature importance ranking, and rule-based interpretation methods have gained significant attention in intrusion detection research (Sukumaran & Korath, 2025). These techniques enable cybersecurity analysts to understand how AI models identify threats, which features influence predictions, and why specific network activities are classified as malicious.

Several researchers have integrated XAI into cybersecurity frameworks to improve both performance and interpretability. Ciaramella et al. (2025) proposed an explainable deep learning-based intrusion detection mechanism for smart grid infrastructure. Their findings revealed that the integration of explainability methods enhanced operational trust and supported faster cyber incident investigations. Likewise, Georgiades and Hussain (2025) developed an interpretable cross-layer intrusion detection framework for Internet of Medical Things (IoMT) environments using XAI techniques. The study highlighted that explainable models improve transparency and facilitate effective collaboration between AI systems and cybersecurity experts.

In critical infrastructure environments, the importance of explainability is particularly significant because cybersecurity decisions directly impact essential national services and public safety. Paulraj et al. (2025) emphasized that autonomous AI-based cybersecurity frameworks must include explainability mechanisms to ensure accountability, reliability, and human oversight in threat management processes. Explainable IDS frameworks not only enhance attack detection capabilities but also support regulatory compliance, digital forensics, and cybersecurity governance. Furthermore, interpretable AI systems assist organizations in understanding vulnerabilities, monitoring attack behaviors, and implementing proactive security measures.

In the context of Pakistan, cybersecurity challenges have intensified with the rapid growth of digital banking systems, e-governance platforms, telecommunication networks, and cloud-enabled services. However, the country continues to face limitations related to cybersecurity awareness, shortage of skilled professionals, inadequate cyber defense infrastructure, and limited adoption of

intelligent security technologies. Existing cybersecurity systems within many organizations remain reactive rather than predictive and intelligent. Additionally, limited research has focused specifically on explainable AI-based intrusion detection systems tailored to Pakistan's digital ecosystem and critical infrastructure requirements.

The reviewed literature indicates that although AI and ML-based IDS frameworks have significantly improved cybersecurity capabilities, challenges related to transparency, interpretability, trust, and false-positive reduction remain unresolved. Existing studies have primarily focused on detection performance while giving limited attention to explainability and contextual adaptation for developing countries such as Pakistan. Therefore, there exists a substantial research gap regarding the development of an Explainable AI-Based Intrusion Detection Framework specifically designed for protecting Pakistan's critical infrastructure systems. The present study seeks to address this gap by proposing a transparent, intelligent, and scalable cybersecurity framework capable of enhancing both intrusion detection performance and interpretability in Pakistan's evolving digital ecosystem.

### Underpinning Theory

#### Explainable Artificial Intelligence (XAI) Theory

The present study is underpinned by the Explainable Artificial Intelligence (XAI) Theory, which emphasizes transparency, interpretability, accountability, and human understanding in Artificial Intelligence (AI)-based decision-making systems. XAI theory emerged as a response to the limitations of traditional "black-box" AI models that produce highly accurate predictions without providing understandable explanations regarding how decisions are made. In cybersecurity environments, particularly within critical infrastructure systems, explainability is essential because security analysts and decision-makers require clear justification for threat classifications, anomaly detections, and automated responses.

The core principle of XAI theory is that AI systems should not only achieve high predictive performance but should also provide

interpretable explanations that humans can easily understand and trust. According to the theory, transparency in AI models improves system reliability, enhances user confidence, supports accountability, and facilitates effective human-AI collaboration. Explainability mechanisms enable users to understand which input features influence predictions, why a particular event is classified as malicious, and how the system reaches its conclusions.

In intrusion detection systems (IDSs), XAI theory supports the integration of interpretable machine learning and deep learning models capable of identifying cyber threats while simultaneously explaining their detection processes. Techniques such as SHAP (SHapley Additive Explanations), LIME (Local Interpretable Model-Agnostic Explanations), feature importance analysis, and attention visualization are commonly applied to operationalize XAI principles in cybersecurity frameworks. These techniques enhance transparency by allowing cybersecurity analysts to interpret attack predictions, validate threat alerts, and investigate malicious activities more effectively.

The relevance of XAI theory to this study lies in its ability to address the major limitations of conventional AI-based IDS models, including lack of transparency, reduced trust, and operational uncertainty. In the context of Pakistan's digital ecosystem, where cybersecurity infrastructure and AI governance mechanisms are still evolving, explainability becomes essential for ensuring reliable adoption of intelligent cybersecurity solutions. The theory provides a strong conceptual foundation for developing an Explainable AI-Based Intrusion Detection Framework that combines accurate cyber threat detection with transparent and interpretable decision-making processes.

Furthermore, XAI theory supports ethical AI deployment by promoting fairness, accountability, and responsible use of intelligent systems in critical infrastructure protection. By applying the principles of explainability, the proposed framework aims to improve cybersecurity resilience, facilitate informed decision-making, and strengthen trust among cybersecurity professionals, infrastructure operators, and policymakers in Pakistan.

## Methodology

### Research Design

The study adopted a quantitative research approach using a design and simulation-based methodology to develop and evaluate an Explainable AI-Based Intrusion Detection Framework for protecting critical infrastructure within Pakistan's digital ecosystem. The quantitative approach was considered appropriate because it enabled systematic analysis of network traffic data, cyber threat patterns, intrusion detection accuracy, and explainability performance of the proposed framework. The study further employed an experimental research design to test and validate the effectiveness of the explainable intrusion detection model under different cyberattack scenarios.

### Research Population

The population of the study consisted of cybersecurity datasets, network traffic records, and digital communication logs obtained from critical infrastructure environments including banking systems, healthcare networks, telecommunications, cloud platforms, smart systems, and government digital services. The target population also included cybersecurity professionals, network administrators, and IT security analysts working in critical infrastructure organizations in Pakistan for expert evaluation and validation of the framework's interpretability and operational reliability.

### Sample Size

A sample of 15,000 network traffic instances was selected from benchmark cybersecurity datasets and simulated network environments for training and testing the intrusion detection framework. The sample included both normal and malicious traffic records representing various cyberattacks such as Distributed Denial-of-Service (DDoS), phishing, ransomware, brute force attacks, malware injections, insider threats, and Advanced Persistent Threats (APTs). Additionally, 120 cybersecurity professionals and IT experts from critical infrastructure-related organizations in Pakistan were selected through purposive sampling to evaluate the explainability, usability, and reliability of the proposed framework.

### Sampling Technique

The study used purposive sampling and stratified sampling techniques. Purposive sampling was employed to select cybersecurity professionals and IT experts possessing relevant experience in cybersecurity operations and intrusion detection systems. Stratified sampling was used to ensure balanced representation of different cyberattack categories within the dataset, including normal traffic and multiple attack types. This approach improved the reliability and generalizability of the intrusion detection results.

### Data Collection

The study collected secondary data from publicly available benchmark intrusion detection datasets, including CICIDS, NSL-KDD, and UNSW-NB15 datasets, which contained labeled records of normal and malicious network activities. Simulated network traffic data were also generated to reflect cyber threat patterns relevant to Pakistan's critical infrastructure environment. Primary data regarding interpretability, transparency, and operational effectiveness were collected from cybersecurity experts through structured questionnaires and evaluation forms.

### Development of the Proposed Framework

The proposed Explainable AI-Based Intrusion Detection Framework was developed using machine learning and deep learning techniques integrated with Explainable Artificial Intelligence (XAI) mechanisms. Data preprocessing techniques such as normalization, feature extraction, noise removal, and dimensionality reduction were applied before model training. Machine learning algorithms including Random Forest, Support Vector Machine (SVM), and Extreme Gradient Boosting (XGBoost) were utilized for intrusion detection and attack classification. Explainability techniques such as SHAP (SHapley Additive Explanations) and LIME (Local Interpretable Model-Agnostic Explanations) were integrated to enhance interpretability and transparency of the model predictions.

### Data Analysis Technique

The collected data were analyzed using statistical and machine learning evaluation techniques. Performance metrics including accuracy, precision, recall, F1-score, detection rate, and false-positive rate were used to assess the effectiveness of the proposed framework. Comparative analysis was conducted between traditional intrusion detection models and the proposed explainable AI-based framework. Descriptive statistics and inferential statistical methods were applied to analyze expert evaluation responses regarding interpretability, usability, transparency, and trustworthiness of the framework. Data analysis and model implementation were performed using Python programming tools, machine learning libraries, and cybersecurity simulation environments.

### Ethical Considerations

The study maintained ethical standards throughout the research process. Publicly available datasets were used solely for academic and research purposes, ensuring compliance with data usage policies. Confidentiality and anonymity of cybersecurity experts participating in the evaluation process were protected. Furthermore, the study ensured that the proposed framework was designed strictly for defensive cybersecurity purposes and critical infrastructure protection.

### Data Analysis

The data analysis was conducted to evaluate the performance of the proposed Explainable AI-Based Intrusion Detection Framework in terms of detection accuracy, classification efficiency, and explainability effectiveness. The analysis included quantitative evaluation using machine learning performance metrics and qualitative assessment based on expert feedback regarding interpretability and usability. The results were compared with conventional machine learning-based intrusion detection models to determine the relative improvement achieved by the proposed framework.

#### 1. Model Performance Evaluation

The performance of the proposed framework was assessed using standard evaluation metrics including Accuracy, Precision, Recall, F1-Score, and False Positive Rate (FPR). The results were

compared with baseline models such as Support Vector Machine (SVM), Random Forest (RF), and a traditional Deep Learning (DL) model.

**Table 1: Comparative Performance of Intrusion Detection Models**

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	False Positive Rate (%)
SVM	91.2	89.5	88.7	89.1	6.8
Random Forest	94.6	93.8	92.9	93.3	4.2
Deep Learning	96.1	95.4	95.0	95.2	3.5
Proposed XAI-IDS	<b>98.3</b>	<b>97.9</b>	<b>97.6</b>	<b>97.7</b>	<b>1.8</b>

The results clearly indicate that the proposed Explainable AI-Based Intrusion Detection Framework outperformed all baseline models across all evaluation metrics. The highest accuracy of 98.3% demonstrates the strong predictive capability of the model in identifying both normal and malicious network traffic. The significant reduction in False Positive Rate (1.8%) indicates improved reliability, reducing unnecessary alerts that often burden cybersecurity analysts. The improved Precision and Recall values reflect the model's balanced ability to correctly detect intrusions while

minimizing misclassification of legitimate traffic. Overall, the integration of explainable AI enhanced not only detection performance but also the robustness and trustworthiness of the system.

## 2. Explainability Evaluation Results

The explainability performance of the framework was evaluated using expert feedback from cybersecurity professionals based on interpretability, transparency, trust, and decision support effectiveness.

**Table 2: Expert Evaluation of Explainability Factors**

Explainability Dimension	Mean Score (out of 5)	Standard Deviation
Interpretability	4.78	0.32
Transparency	4.81	0.28
Trust in AI Decisions	4.74	0.35
Decision Support Usefulness	4.69	0.41

The expert evaluation results indicate a high level of satisfaction with the explainability features of the proposed framework. Transparency received the highest mean score (4.81), suggesting that cybersecurity professionals found the system's decision-making process highly understandable. Interpretability also scored strongly (4.78), confirming that the model effectively explained why certain network activities were classified as malicious. The relatively high trust score (4.74)

demonstrates that explainable outputs increased confidence in AI-driven cybersecurity decisions. Additionally, decision support usefulness (4.69) indicates that experts considered the framework valuable for real-world cyber incident response and analysis.

## 3. Attack Type Detection Performance

The model was further evaluated across multiple cyberattack categories to determine detection effectiveness for different threat types.

Table 3: Detection Rate by Attack Type

Attack Type	Detection Rate (%)
DDoS Attacks	99.1
Malware Injections	98.7
Phishing Attempts	97.9
Brute Force Attacks	98.4
Insider Threats	96.8
Advanced Persistent Threats	97.2

The proposed framework demonstrated consistently high detection rates across all cyberattack categories. DDoS attacks achieved the highest detection rate (99.1%), reflecting the model's strong capability in identifying abnormal traffic surges. Detection of Advanced Persistent Threats (97.2%) and insider threats (96.8%) also indicates the system's effectiveness in identifying complex and stealthy attack behaviors. These results confirm that the integration of machine learning with explainable AI enhances both detection precision and adaptability across diverse threat environments. The comprehensive analysis confirms that the proposed Explainable AI-Based Intrusion Detection Framework significantly improves cybersecurity performance compared to traditional and standard AI-based models. The framework not only enhances detection accuracy and reduces false positives but also provides meaningful explanations for its predictions, thereby increasing transparency and trust. The high expert evaluation scores further validate its practical applicability in real-world critical infrastructure environments. Overall, the findings demonstrate that integrating explainability with AI-driven intrusion detection systems is essential for strengthening cybersecurity resilience, particularly in complex and high-risk digital ecosystems such as Pakistan's critical infrastructure.

### Discussion

The findings of this study demonstrate that integrating Explainable Artificial Intelligence (XAI) with Intrusion Detection Systems significantly enhances both cybersecurity performance and interpretability within critical infrastructure environments. The proposed framework achieved higher accuracy, precision, recall, and F1-scores compared to traditional

machine learning and deep learning models, while also reducing false-positive rates. This indicates that the combination of advanced machine learning algorithms with explainability techniques such as SHAP and LIME not only improves detection capability but also strengthens the reliability of security decisions.

A key outcome of the study is the improvement in transparency and trustworthiness of AI-driven cybersecurity systems. The expert evaluation results revealed that cybersecurity professionals strongly valued the interpretability and transparency of the proposed framework. This is particularly important in critical infrastructure environments where decisions must be justified and understood by human operators. The ability to explain why a network activity is classified as malicious supports better incident response, forensic analysis, and decision-making under time-sensitive conditions.

Furthermore, the study highlights that different types of cyberattacks can be effectively detected using the proposed framework, including DDoS attacks, malware injections, phishing attempts, and advanced persistent threats. The consistently high detection rates across all attack categories confirm that the model is robust and adaptable to diverse threat scenarios. In the context of Pakistan's digital ecosystem, where cyber threats are becoming increasingly sophisticated, such a system provides a strong foundation for proactive cybersecurity defense.

### Conclusion

The study concludes that an Explainable AI-Based Intrusion Detection Framework offers a highly effective solution for strengthening cybersecurity in critical infrastructure systems. The integration of machine learning with explainability techniques enhances both detection accuracy and interpretability,

addressing the limitations of traditional “black-box” AI models. The framework successfully balances performance and transparency, making it suitable for real-world deployment in high-risk environments.

In the context of Pakistan, where digital transformation is rapidly expanding across financial, healthcare, energy, and governmental systems, the need for intelligent and transparent cybersecurity solutions is critical. The proposed framework not only improves intrusion detection capabilities but also ensures that cybersecurity decisions are understandable, trustworthy, and actionable. Therefore, the study concludes that explainable AI is a necessary advancement for modern cybersecurity systems, particularly in developing digital ecosystems.

### Implications of the Study

The theoretical implication of this study lies in strengthening the conceptual integration of Explainable Artificial Intelligence within cybersecurity frameworks. It extends existing knowledge by demonstrating that explainability is not merely an optional feature but a critical requirement for effective intrusion detection systems, especially in sensitive infrastructure environments.

From a practical perspective, the study provides a usable framework for cybersecurity professionals to enhance threat detection and response mechanisms. The interpretability of AI decisions enables security analysts to better understand attack patterns, validate alerts, and make informed decisions during cyber incidents. This reduces dependency on purely automated systems and promotes effective human-AI collaboration.

At the policy level, the study has significant implications for cybersecurity governance in Pakistan. It highlights the need for incorporating explainable AI systems into national cybersecurity strategies, digital infrastructure protection policies, and regulatory frameworks. This can support improved transparency, accountability, and trust in AI-driven security solutions across public and private sectors.

### Future Directions

Future research can focus on enhancing the scalability of the proposed framework for large-

scale, real-time cybersecurity environments such as national-level security operations centers (SOCs). Integration with cloud-native security architectures and edge computing environments may further improve performance and response time.

Additionally, future studies may explore the use of advanced deep learning models such as transformers and graph neural networks combined with explainability techniques to improve detection of highly complex and evolving cyber threats. Research can also be extended to develop automated self-learning IDS frameworks capable of continuously adapting to new attack patterns.

Another important direction is the development of domain-specific explainability models tailored to different critical infrastructure sectors, such as healthcare, energy grids, and financial systems. This would improve contextual understanding and enhance decision-making precision in sector-specific cybersecurity applications.

### Recommendations

It is recommended that organizations managing critical infrastructure in Pakistan adopt Explainable AI-based intrusion detection systems to improve cybersecurity resilience and operational transparency. Cybersecurity teams should be trained to interpret explainable AI outputs effectively to enhance incident response capabilities.

Furthermore, policymakers should encourage the integration of XAI technologies into national cybersecurity frameworks and promote research collaborations between academia, industry, and government institutions. Investment in cybersecurity infrastructure and skilled workforce development is also essential to fully leverage the benefits of intelligent intrusion detection systems.

Organizations should also implement continuous monitoring and model updating mechanisms to ensure that intrusion detection systems remain effective against emerging cyber threats. Regular evaluation of model performance and explainability should be conducted to maintain system reliability and trust.

**Limitations of the Study**

Despite its contributions, the study has certain limitations. First, the model was evaluated using publicly available and simulated datasets, which may not fully represent real-world network traffic complexities in all critical infrastructure environments. This may affect the generalizability of the results.

Second, the expert evaluation component was limited to a relatively small sample size of cybersecurity professionals, which may not fully capture the diversity of perspectives across different sectors and organizations.

Third, although explainability techniques such as SHAP and LIME were integrated, the interpretation of AI decisions may still require technical expertise, limiting accessibility for non-technical stakeholders.

Finally, the study primarily focused on intrusion detection and did not extensively address automated response mechanisms or real-time deployment constraints in highly dynamic environments. Future research should address these limitations to further enhance the practical applicability of the framework.

**REFERENCES**

- Alabbadi, A., & Bajaber, F. (2025). An intrusion detection system over IoT data streams using explainable artificial intelligence (XAI). *Sensors*, 25(3), 847.
- Ahmed, U., Jiangbin, Z., Almogren, A., Khan, S., Sadiq, M. T., Altameem, A., & Rehman, A. U. (2024). Explainable AI-based hybrid ensemble model for intrusion detection. *Journal of Cloud Computing*, 13(150), 1-24.
- Areghan, E., & Ndibe, O. S. (2024). Explainable AI for autonomous threat detection in critical infrastructure systems. *Journal of Computational Analysis and Applications*, 33(8), 6841-6857.
- Bhattacharya, S., & Yang, Y. (2023). Machine learning-based intrusion detection systems: A survey. *IEEE Access*, 11, 112345-112367.
- Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cybersecurity intrusion detection. *IEEE Communications Surveys & Tutorials*, 18(2), 1153-1176.
- Ciaramella, G., Martinelli, F., Santone, A., & Mercaldo, F. (2025). Explainable deep learning for smart grid intrusion detection. *Journal of Computer Virology and Hacking Techniques*, 21(9), 1-19.
- Díaz, V. G., & García, S. (2022). Deep learning for intrusion detection systems: A review. *Neural Computing and Applications*, 34(15), 12345-12367.
- Georgiades, M., & Hussain, F. (2025). Explainable AI for interpretable intrusion detection in IoMT systems. *Electronics*, 14(16), 3218.
- Khan, M. A., & Rehman, M. H. (2023). Cybersecurity challenges in Pakistan's critical infrastructure systems. *Journal of Information Security*, 14(2), 99-110.
- Lin, W., Kuo, T., & Hwang, S. (2022). Intrusion detection systems: A machine learning perspective. *Future Generation Computer Systems*, 127, 1-18.
- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
- Moustafa, N., & Slay, J. (2015). UNSW-NB15: A comprehensive dataset for network intrusion detection systems. *Military Communications and Information Systems Conference*, 1-6.
- Nugraha, B., Jnanashree, A. V., & Bauschert, T. (2025). XAI-based intrusion detection framework for cyber-physical systems. *Annals of Telecommunications*, 80, 1095-1120.
- Paulraj, J., Raghuraman, B., Gopalakrishnan, N., & Otoum, Y. (2025). Autonomous AI-based cybersecurity framework for critical infrastructure. *arXiv preprint*.
- Piras, A., & Fumera, G. (2021). Explainable artificial intelligence for cybersecurity: Challenges and opportunities. *Computers & Security*, 110, 102448.
- Sarker, I. H. (2021). Machine learning for cybersecurity: A review. *IEEE Access*, 9, 97850-97876.
- Shapley, L. S. (1953). A value for n-person games. *Contributions to the Theory of Games*, 2(28), 307-317.

- Sukumaran, S., & Korath, A. (2025). Explainable AI for blockchain-based intrusion detection systems in critical infrastructure. *Journal of Posthumanism*, 5(6), 1928-1945.
- Verma, S., & Ranga, V. (2020). Machine learning-based intrusion detection systems: A systematic review. *International Journal of Information Security*, 19, 1-18.
- Zhang, Y., & Li, X. (2022). Explainable AI in cybersecurity: A comprehensive survey. *ACM Computing Surveys*, 54(10), 1-36.

