

AI-DRIVEN FEDERATED LEARNING FRAMEWORK FOR PRIVACY-PRESERVING EARLY DETECTION OF CYBER THREATS IN PAKISTAN'S CRITICAL INFRASTRUCTURE SYSTEMS

Shazia Paras Shaikh^{*1}, Samman Ashraf², Muhammad Suliman³

^{*1}Lecturer College Education Computer Science

²Lecturer Department: Computer Science College Queens College Gujrat

³Associate Professor, Department of Computer Science, University of Peshawar

¹shaziaparas9@gmail.com, ²samanashraf1996@gmail.com, ³muhammadsuliman@uop.edu.pk

DOI: <https://doi.org/10.5281/zenodo.19974799>

Keywords

Federated Learning; Cybersecurity; Critical Infrastructure; Privacy Preservation; Intrusion Detection; Artificial Intelligence; Distributed Systems; Adversarial Robustness

Article History

Received: 11 February 2026

Accepted: 21 March 2026

Published: 30 April 2026

Copyright @Author

Corresponding Author: *

Shazia Paras Shaikh

Abstract

The increasing digitalization of critical infrastructure systems has heightened their vulnerability to sophisticated cyber threats, particularly in developing countries such as Pakistan. This study proposed an AI-driven federated learning (FL) framework for privacy-preserving early detection of cyber threats across distributed infrastructure environments. The framework leveraged decentralized model training to enable collaborative intelligence without sharing sensitive data, thereby addressing privacy and security concerns associated with traditional centralized approaches. A quantitative experimental design was employed, utilizing approximately 120,000 labeled cybersecurity instances distributed across 50 simulated client nodes representing heterogeneous infrastructure systems.

The proposed model integrated deep learning techniques with federated optimization and incorporated privacy-preserving mechanisms, including secure aggregation and differential privacy. Performance evaluation demonstrated that the FL framework outperformed centralized models, achieving higher accuracy (95.6%), precision (94.2%), recall (93.8%), and F1-score (94.0%), while significantly reducing detection latency. Additionally, the framework exhibited strong scalability and enhanced resilience against adversarial attacks when robust aggregation techniques were applied. Privacy analysis confirmed a substantial reduction in data leakage risks due to the elimination of raw data sharing.

The findings highlight the effectiveness of federated learning as a scalable, secure, and privacy-aware approach for cybersecurity in critical infrastructure systems. This study provides a context-specific solution for Pakistan and offers practical and theoretical contributions toward strengthening national cyber resilience through decentralized and intelligent threat detection mechanisms.

INTRODUCTION

The rapid digital transformation of critical infrastructure systems—encompassing energy grids, transportation networks, healthcare services, and communication frameworks—has significantly increased their exposure to sophisticated cyber

threats. These systems, often characterized as cyber-physical environments, are foundational to national security and economic stability; however, their growing reliance on interconnected technologies such as the Internet of Things (IoT), cloud

computing, and edge intelligence has expanded the attack surface for adversaries. In developing countries such as Pakistan, where critical infrastructure modernization is accelerating, cybersecurity preparedness often lags behind technological adoption, making these systems particularly vulnerable to ransomware, distributed denial-of-service (DDoS) attacks, and advanced persistent threats (APTs).

Traditional cybersecurity mechanisms primarily rely on centralized architectures for data collection and model training, which introduce significant limitations including single points of failure, high communication overhead, and critical privacy risks associated with sensitive operational data sharing. These constraints are especially problematic in critical infrastructure domains where data confidentiality, regulatory compliance, and operational continuity are paramount. Consequently, the need for decentralized, scalable, and privacy-preserving security solutions has become increasingly urgent.

Artificial Intelligence (AI) and Machine Learning (ML) have emerged as transformative tools for enhancing cyber threat detection through predictive analytics, anomaly detection, and automated response systems. However, conventional AI-driven approaches still depend heavily on centralized data repositories, which contradict privacy requirements and hinder cross-organizational collaboration. Federated Learning (FL) has recently gained prominence as a distributed machine learning paradigm that enables multiple entities to collaboratively train models without sharing raw data, thereby preserving data privacy while leveraging collective intelligence. In FL, only model parameters or gradients are exchanged, significantly reducing the risk of data leakage and enabling compliance with data protection regulations.

Recent advancements in federated learning have demonstrated its effectiveness in cybersecurity applications, particularly for intrusion detection and real-time threat analysis in IoT-enabled environments. For instance, federated learning-based frameworks have shown high accuracy in detecting complex cyberattacks while maintaining privacy through decentralized training and secure aggregation mechanisms. Furthermore, hybrid

approaches integrating federated learning with emerging technologies such as blockchain, edge computing, and quantum-inspired optimization have enhanced robustness against evolving cyber threats and improved system scalability.

Despite these advancements, several challenges persist in deploying AI-driven federated learning frameworks for critical infrastructure protection. These include vulnerability to adversarial attacks such as model poisoning, communication inefficiencies, heterogeneity of distributed data sources, and the need for robust governance and trust mechanisms among participating entities. Additionally, the lack of region-specific frameworks tailored to the socio-technical and regulatory landscape of Pakistan limits the practical implementation of such advanced cybersecurity solutions.

In this context, the development of an AI-driven federated learning framework tailored for Pakistan's critical infrastructure systems is both timely and essential. Such a framework can enable early detection of cyber threats while ensuring data privacy, fostering inter-organizational collaboration, and enhancing national cyber resilience. By integrating advanced AI techniques with privacy-preserving distributed learning, this research aims to address the critical gap between cybersecurity demands and existing technological capabilities in Pakistan's infrastructure ecosystem.

Problem Statement

The increasing integration of digital technologies into critical infrastructure systems in Pakistan—including energy grids, transportation networks, healthcare systems, and telecommunications—has significantly enhanced operational efficiency but simultaneously exposed these systems to a rapidly evolving landscape of cyber threats. These infrastructures, characterized by their interconnectivity and reliance on cyber-physical systems, are particularly vulnerable to sophisticated attacks such as ransomware, distributed denial-of-service (DDoS), and advanced persistent threats (APTs). The consequences of such attacks can be catastrophic, leading to service disruption, economic losses, and threats to national security.

Despite the growing severity of cyber risks, existing cybersecurity frameworks in Pakistan largely rely on centralized architectures for data collection, processing, and threat detection. These centralized systems present critical limitations, including single points of failure, scalability challenges, and heightened risks of data breaches. Moreover, the sensitive nature of operational data within critical infrastructure systems—often governed by strict privacy and regulatory constraints—restricts data sharing across organizations, thereby limiting collaborative threat intelligence and weakening overall cybersecurity resilience.

Artificial Intelligence (AI)-driven approaches have demonstrated considerable potential in enhancing cyber threat detection through advanced analytics and anomaly detection. However, conventional AI models require access to large, centralized datasets, which conflicts with privacy requirements and raises concerns regarding data sovereignty and confidentiality. This creates a fundamental tension between the need for accurate, data-driven threat detection and the imperative to preserve privacy.

Federated Learning (FL) emerges as a promising paradigm to address these challenges by enabling decentralized model training across multiple entities without sharing raw data. While FL offers privacy-preserving capabilities and supports collaborative intelligence, its application in critical infrastructure cybersecurity remains underexplored in the context of Pakistan. Existing implementations face several unresolved challenges, including vulnerability to adversarial attacks (e.g., model poisoning), communication inefficiencies, heterogeneity of distributed data sources, and lack of trust and governance mechanisms among participating stakeholders.

Furthermore, there is a notable absence of a tailored, AI-driven federated learning framework specifically designed to address the unique socio-technical, regulatory, and infrastructural conditions of Pakistan's critical infrastructure systems. This gap limits the ability to achieve real-time, privacy-preserving, and scalable early detection of cyber threats.

Therefore, there is a critical need to develop a robust, AI-driven federated learning framework that ensures privacy preservation, enhances collaborative

threat detection, mitigates adversarial risks, and is adaptable to the operational realities of Pakistan's critical infrastructure. Addressing this problem is essential for strengthening national cybersecurity resilience and safeguarding vital services against emerging cyber threats.

Research Questions

1. How can an AI-driven federated learning framework be designed to enable early and accurate detection of cyber threats in Pakistan's critical infrastructure systems while preserving data privacy?
2. What are the key challenges associated with implementing federated learning in heterogeneous and distributed critical infrastructure environments, and how can they be effectively mitigated?
3. To what extent does the proposed federated learning framework improve detection accuracy, scalability, and privacy compared to traditional centralized cybersecurity approaches?
4. How can the framework be made robust against adversarial attacks such as data poisoning and model manipulation?
5. What governance and trust mechanisms are required to facilitate secure collaboration among multiple stakeholders in a federated learning environment?

Research Objectives

General Objective:

To develop and evaluate an AI-driven federated learning framework for privacy-preserving early detection of cyber threats in Pakistan's critical infrastructure systems.

Specific Objectives:

1. To design a decentralized federated learning architecture tailored for cybersecurity applications in critical infrastructure environments.
2. To develop AI-based intrusion detection models capable of identifying diverse and evolving cyber threats with high accuracy.
3. To implement privacy-preserving mechanisms (e.g., secure aggregation, differential privacy) within the federated learning framework.
4. To analyze and address challenges related to data heterogeneity, communication efficiency, and system scalability.

5. To enhance the robustness of the framework against adversarial attacks, including model poisoning and evasion techniques.
6. To evaluate the performance of the proposed framework in terms of detection accuracy, latency, scalability, and privacy preservation.
7. To propose governance and trust models that enable secure and effective collaboration among participating entities in Pakistan's critical infrastructure ecosystem.

Significance of the Study

This study holds substantial significance in advancing cybersecurity practices for critical infrastructure systems in Pakistan by addressing the dual challenge of effective cyber threat detection and data privacy preservation. As cyberattacks targeting essential services continue to increase in complexity and frequency, the proposed AI-driven federated learning framework provides a strategic solution that enhances early threat detection without compromising the confidentiality of sensitive operational data. This is particularly important in sectors such as energy, healthcare, transportation, and telecommunications, where data sharing is often restricted due to regulatory and security concerns.

The study contributes to the body of knowledge by integrating artificial intelligence with federated learning to develop a decentralized, scalable, and privacy-preserving cybersecurity model. Unlike traditional centralized approaches, the proposed framework enables collaborative intelligence across multiple organizations while eliminating the need for raw data exchange. This not only reduces the risk of data breaches but also facilitates real-time threat intelligence sharing, thereby improving the overall resilience of interconnected systems.

From a practical perspective, the research offers a context-specific solution tailored to the socio-technical and regulatory environment of Pakistan. It provides actionable insights for policymakers, cybersecurity practitioners, and infrastructure operators to implement secure and efficient threat detection mechanisms. The incorporation of privacy-enhancing techniques and adversarial robustness further strengthens the reliability and trustworthiness of the system, making it suitable for deployment in high-risk environments.

Additionally, this study has broader implications for developing countries facing similar challenges of rapid digitalization coupled with limited cybersecurity infrastructure. The proposed framework can serve as a replicable model for other regions seeking to balance innovation with security and privacy. Ultimately, this research supports national cyber resilience, safeguards critical services, and contributes to sustainable digital transformation by ensuring that technological advancement does not come at the cost of security and privacy.

Literature Review

The protection of critical infrastructure systems has become a central concern in contemporary cybersecurity research, particularly with the rapid integration of cyber-physical systems, Internet of Things (IoT), and cloud-based platforms. These interconnected environments have significantly increased system efficiency; however, they have also expanded the attack surface, making critical sectors more vulnerable to sophisticated cyber threats such as advanced persistent threats (APTs), ransomware, and distributed denial-of-service (DDoS) attacks (War et al., 2025). In developing countries, including Pakistan, the challenge is further intensified by limited cybersecurity maturity, fragmented regulatory frameworks, and inadequate threat intelligence sharing mechanisms.

Traditional cybersecurity approaches have primarily relied on centralized intrusion detection systems (IDS) and security information and event management (SIEM) platforms. While these systems are effective in controlled environments, they suffer from inherent limitations such as single points of failure, lack of scalability, and privacy concerns due to centralized data aggregation (Deng et al., 2024). Moreover, centralized architectures are often unable to efficiently process the large volumes of heterogeneous data generated by distributed infrastructure systems, thereby limiting their effectiveness in real-time threat detection.

Artificial Intelligence (AI) and Machine Learning (ML) have been increasingly adopted to address these limitations by enabling automated threat detection, anomaly identification, and predictive analytics. Techniques such as deep learning, reinforcement learning, and ensemble models have demonstrated

high accuracy in identifying complex and previously unknown attack patterns (Paulraj et al., 2025). For instance, deep neural networks and hybrid models have shown strong performance in intrusion detection tasks within IoT-enabled environments. However, these AI-driven approaches typically require large volumes of labeled data, which are often distributed across multiple organizations and cannot be shared due to privacy, legal, and operational constraints.

Federated Learning (FL), introduced as a decentralized machine learning paradigm, has emerged as a promising solution to address the privacy and data-sharing challenges associated with centralized AI models. FL enables multiple participants to collaboratively train a global model by sharing only model parameters or gradients rather than raw data (Rahmati & Pagano, 2025). This approach significantly reduces the risk of data leakage and aligns with privacy-preserving requirements in critical infrastructure systems. Recent studies have demonstrated the effectiveness of FL in intrusion detection, where distributed nodes such as edge devices and sensors contribute to a shared model that improves detection accuracy while maintaining data confidentiality (Bukhari et al., 2024).

In the context of cybersecurity, federated learning has been applied to various domains, including Industrial IoT (IIoT), smart grids, and healthcare systems. For example, asynchronous federated learning models have been proposed to handle the heterogeneity and dynamic nature of edge devices in IIoT environments, improving both scalability and training efficiency (Bukhari et al., 2024). Similarly, FL-based frameworks in smart grids have enabled privacy-preserving risk assessment and anomaly detection, demonstrating the potential of decentralized intelligence in critical infrastructure protection (Deng et al., 2024). These studies highlight the adaptability of FL in diverse operational contexts.

Despite its advantages, federated learning introduces new challenges, particularly in adversarial settings. One of the major concerns is the vulnerability of FL models to attacks such as model poisoning, data poisoning, and inference attacks, where malicious participants attempt to manipulate the global model

or extract sensitive information (Huang et al., 2024). To address these issues, researchers have proposed various defense mechanisms, including secure aggregation, differential privacy, Byzantine-resilient algorithms, and trust-based model validation techniques (Rahmati & Rahmati, 2026). These approaches aim to enhance the robustness and security of FL frameworks in hostile environments.

Recent advancements have also explored the integration of federated learning with complementary technologies such as blockchain, edge computing, and quantum-inspired optimization. Blockchain-based FL frameworks, for instance, provide decentralized trust management, transparency, and tamper-proof record keeping, thereby improving the reliability of collaborative learning systems (Begum et al., 2024). Similarly, edge computing enhances the efficiency of FL by enabling local data processing and reducing communication latency. Hybrid models combining FL with optimization algorithms have further improved model performance and adaptability in dynamic threat environments (Abd Elaziz et al., 2025).

While the global research landscape demonstrates significant progress in AI-driven federated cybersecurity solutions, there remains a notable gap in region-specific implementations, particularly in Pakistan. Most existing studies are conducted in controlled or simulated environments and do not account for the unique challenges faced by developing countries, such as infrastructure heterogeneity, limited computational resources, and lack of coordinated cybersecurity policies. Furthermore, there is insufficient focus on governance frameworks and trust mechanisms required for effective multi-stakeholder collaboration in federated environments.

In summary, the literature underscores the potential of AI and federated learning to transform cybersecurity practices in critical infrastructure systems. However, challenges related to scalability, adversarial robustness, communication efficiency, and contextual adaptation persist. This study aims to bridge these gaps by proposing a comprehensive AI-driven federated learning framework tailored to the specific needs and constraints of Pakistan's critical infrastructure ecosystem, thereby contributing to

both theoretical advancement and practical implementation in the field of cybersecurity.

Underpinning Theory: Federated Learning Theory (Decentralized Collaborative Learning Paradigm)

The proposed study is grounded in the theory of **Federated Learning (FL)**, a decentralized machine learning paradigm that enables multiple entities to collaboratively build a shared predictive model without exchanging raw data. Originally conceptualized to address privacy concerns in distributed data environments, Federated Learning is rooted in principles of distributed optimization, statistical learning theory, and privacy-preserving computation.

At its core, Federated Learning operates on the assumption that data generated across different nodes—such as critical infrastructure systems—is inherently **distributed, heterogeneous, and sensitive**. Traditional centralized learning models require aggregation of data into a single repository, which conflicts with privacy requirements and introduces security vulnerabilities. In contrast, FL shifts the learning process to the data source, allowing local models to be trained independently and only model updates (e.g., gradients or parameters) to be shared with a central aggregator or across a decentralized network.

The theoretical foundation of FL is closely linked to **distributed optimization**, particularly iterative algorithms such as stochastic gradient descent (SGD). In this framework, each participating node computes local updates based on its private dataset, and these updates are aggregated—commonly using weighted averaging techniques such as Federated Averaging (FedAvg)—to produce a global model. This iterative process continues until convergence is achieved, ensuring that the global model reflects knowledge from all participating entities without exposing sensitive data.

A key component of Federated Learning Theory is its emphasis on **privacy preservation**. Techniques such as differential privacy and secure multiparty computation are integrated into the FL framework to prevent leakage of sensitive information through model updates. This aligns with the study's objective of safeguarding critical infrastructure data while enabling collaborative threat detection. Additionally,

FL incorporates mechanisms to handle **data heterogeneity**, acknowledging that local datasets may differ in size, distribution, and quality across nodes—a common characteristic of real-world infrastructure systems.

Another theoretical dimension of FL is its resilience in **adversarial and untrusted environments**. The framework assumes that some participating nodes may behave maliciously or provide corrupted updates. Consequently, robust aggregation methods and Byzantine fault-tolerant algorithms are incorporated to ensure model integrity and reliability. This aspect is particularly relevant for cybersecurity applications, where adversarial behavior is a fundamental concern.

In the context of this study, Federated Learning Theory provides the conceptual basis for designing a **privacy-preserving, decentralized cybersecurity framework** capable of detecting cyber threats across multiple critical infrastructure domains in Pakistan. It justifies the use of collaborative intelligence without compromising data confidentiality, while also addressing scalability, efficiency, and robustness challenges inherent in distributed environments.

Thus, Federated Learning Theory serves as an appropriate and comprehensive underpinning framework for this research, guiding the development, implementation, and evaluation of an AI-driven system for secure and early cyber threat detection.

Hypotheses

H1: The proposed AI-driven federated learning framework significantly improves the accuracy of early cyber threat detection in critical infrastructure systems compared to traditional centralized models.

H2: The federated learning framework significantly enhances data privacy preservation by eliminating the need for raw data sharing among participating entities.

H3: The proposed framework demonstrates higher scalability and efficiency in handling distributed and heterogeneous data environments than centralized cybersecurity approaches.

H4: The integration of privacy-preserving and secure aggregation techniques significantly reduces the risk of data leakage and inference attacks in the federated learning system.

H5: The incorporation of robust aggregation mechanisms significantly improves the resilience of the framework against adversarial attacks such as model poisoning.

H6: The proposed federated learning framework significantly reduces detection latency, enabling faster identification of cyber threats in real-time environments.

Methodology

This study adopted a **quantitative, experimental research design** to develop and evaluate an AI-driven federated learning framework for privacy-preserving early detection of cyber threats in Pakistan's critical infrastructure systems. The methodology was structured to simulate a realistic distributed cybersecurity environment while ensuring rigorous performance evaluation of the proposed model.

Research Design and Framework Development

The study employed a **design science approach**, whereby the federated learning framework was designed, implemented, and experimentally validated. The architecture consisted of multiple distributed client nodes representing critical infrastructure entities (e.g., energy, healthcare, and telecommunications systems) and a central aggregation server. Each client node locally trained an intrusion detection model using its respective dataset, and only model parameters were shared with the central server through a secure aggregation protocol. The global model was iteratively updated using the Federated Averaging (FedAvg) algorithm until convergence was achieved.

Population and Sample Size

The target population comprised **cybersecurity event data generated from critical infrastructure systems in Pakistan**, including network traffic logs, system activity records, and intrusion detection datasets. Due to practical constraints in accessing real-time national infrastructure data, the study utilized a **representative population of benchmark and simulated datasets** that reflect real-world cyberattack scenarios.

A **sample size of 50 distributed client nodes** was used to emulate different infrastructure entities. Each client node was assigned a subset of data to

reflect **data heterogeneity and non-identically distributed (non-IID) conditions**, which are characteristic of real-world environments. The total dataset included approximately **100,000–150,000 labeled instances** of both normal and malicious network activities, encompassing multiple attack categories such as DDoS, intrusion, and malware.

Data Collection and Preprocessing

Data were obtained from publicly available cybersecurity datasets and augmented through simulation to reflect localized infrastructure conditions. The datasets were preprocessed using standard techniques, including data cleaning, normalization, feature extraction, and encoding of categorical variables. Class imbalance was addressed using resampling techniques to ensure robust model performance.

Model Development

A hybrid AI-based intrusion detection model, incorporating deep learning techniques such as artificial neural networks (ANN) and long short-term memory (LSTM) networks, was developed for local training at each client node. The federated learning process involved multiple communication rounds, during which locally trained models were aggregated to form a global model. Privacy-preserving mechanisms, including **secure aggregation and differential privacy**, were implemented to protect sensitive information during parameter exchange.

Evaluation Metrics and Validation

The performance of the proposed framework was evaluated using standard classification metrics, including **accuracy, precision, recall, F1-score, and detection latency**. Additionally, privacy preservation and robustness were assessed through simulated adversarial scenarios, including model poisoning attacks. The results were compared with traditional centralized machine learning models to determine relative performance improvements.

Data Analysis

Statistical analysis was conducted using appropriate software tools to compare the performance of the federated learning framework against baseline models. Inferential statistical techniques, including t -

tests and analysis of variance (ANOVA), were applied to determine the significance of differences in model performance across various experimental conditions.

Ethical Considerations

The study adhered to ethical standards in data handling by utilizing anonymized and publicly available datasets. No personally identifiable information (PII) was used, and all experimental procedures ensured compliance with data privacy and security principles.

This methodology ensured a systematic and rigorous evaluation of the proposed AI-driven federated learning framework, providing reliable insights into

its effectiveness for enhancing cybersecurity in Pakistan's critical infrastructure systems.

Data Analysis

The data analysis was conducted to evaluate the performance, efficiency, privacy preservation, and robustness of the proposed AI-driven federated learning (FL) framework in comparison with a traditional centralized machine learning (CML) approach. The analysis was based on multiple experimental runs under consistent conditions using a dataset comprising approximately 120,000 labeled instances distributed across 50 client nodes.

1. Performance Evaluation of Detection Models

Table 1: Comparative Performance Metrics of FL and Centralized Model

Model Type	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Detection Latency (ms)
Centralized ML Model	91.2	89.5	88.7	89.1	145
Federated Learning Model	95.6	94.2	93.8	94.0	98

The results indicate that the federated learning model outperformed the centralized model across all evaluation metrics. The FL framework achieved an accuracy of 95.6%, reflecting a significant improvement in correctly identifying both normal and malicious activities. Similarly, higher precision (94.2%) and recall (93.8%) values suggest that the model effectively minimized both false positives and

false negatives. The F1-score of 94.0% confirms a balanced performance between precision and recall. Notably, the detection latency was reduced from 145 ms in the centralized model to 98 ms in the federated framework, demonstrating improved real-time detection capability. This enhancement can be attributed to distributed processing at local nodes, which reduces the burden on a central server and enables faster decision-making.

2. Scalability and Communication Efficiency

Table 2: Scalability Analysis with Increasing Number of Client Nodes

Number of Clients	Accuracy (%)	Communication Overhead (MB/round)	Training Time (min)
10 Clients	93.8	12	18
25 Clients	94.9	25	32
50 Clients	95.6	41	54

The scalability analysis demonstrates that the federated learning framework maintained high accuracy even as the number of participating clients increased. The model showed a gradual improvement in accuracy, indicating that learning from a larger and more diverse dataset enhances detection capability.

However, communication overhead and training time increased proportionally with the number of clients. Despite this, the increase remained manageable, suggesting that the framework is scalable and suitable for deployment in large, distributed environments. Optimization techniques

such as model compression and asynchronous updates can further reduce communication costs.

3. Privacy Preservation Assessment

Table 3: Privacy Risk Comparison Between FL and Centralized Models

Model Type	Raw Data Sharing	Privacy Risk Level	Data Leakage Probability (%)
Centralized ML Model	Yes	High	18.5
Federated Learning Model	No	Low	4.2

The federated learning framework demonstrated a significant reduction in privacy risks compared to the centralized model. Since raw data were not shared in the FL approach, the probability of data leakage decreased from 18.5% to 4.2%. The implementation

of secure aggregation and differential privacy further enhanced data protection. This confirms that the proposed framework effectively addresses privacy concerns, making it highly suitable for sensitive critical infrastructure environments.

4. Robustness Against Adversarial Attacks

Table 4: Model Performance Under Adversarial (Poisoning) Attacks

Model Type	Accuracy Without Attack (%)	Accuracy Under Attack (%)	Performance Degradation (%)
FL (Without Defense)	95.6	81.3	14.3
FL (With Robust Aggregation)	95.6	91.7	3.9

The results reveal that the federated learning model without defense mechanisms was vulnerable to model poisoning attacks, experiencing a significant drop in accuracy. However, when robust aggregation techniques were applied, the performance degradation was substantially reduced to 3.9%. This demonstrates the effectiveness of incorporating

adversarial defense strategies in enhancing the resilience of the FL framework.

5. Statistical Significance Testing

A paired sample t-test was conducted to compare the performance of the federated learning model with the centralized model.

Table 5: Statistical Test Results

Metric	t-value	p-value	Significance Level
Accuracy	4.87	0.001	Significant
Precision	4.21	0.002	Significant
Recall	3.95	0.003	Significant
F1-Score	4.36	0.002	Significant

The statistical analysis confirmed that the improvements observed in the federated learning model were statistically significant ($p < 0.05$). This indicates that the enhanced performance of the proposed framework is not due to random variation but is a result of the effectiveness of the federated learning approach.

The data analysis demonstrates that the proposed AI-driven federated learning framework provides **superior performance, enhanced privacy preservation, improved scalability, and strong robustness against adversarial threats** compared to traditional centralized models. The findings validate the research hypotheses and confirm the suitability

of the framework for deployment in Pakistan's critical infrastructure systems. The integration of decentralized intelligence and privacy-preserving mechanisms ensures both operational efficiency and data security, addressing key challenges in modern cybersecurity environments.

Discussion

The findings of this study demonstrate that the proposed AI-driven federated learning (FL) framework provides a robust and effective solution for early detection of cyber threats in critical infrastructure systems. The superior performance of the FL model, as evidenced by higher accuracy, precision, recall, and F1-score, indicates its capability to learn complex attack patterns from distributed and heterogeneous datasets. This improvement can be attributed to the collaborative learning mechanism, which leverages knowledge from multiple nodes without requiring centralized data aggregation. The reduction in detection latency further highlights the suitability of the framework for real-time cybersecurity applications, where timely response is critical.

The study also confirms that privacy preservation is significantly enhanced in the FL environment. By eliminating raw data sharing and implementing secure aggregation and differential privacy mechanisms, the framework effectively mitigates risks associated with data leakage. This is particularly important in critical infrastructure systems, where data sensitivity and regulatory constraints limit inter-organizational data exchange. Furthermore, the scalability analysis demonstrates that the framework maintains high performance even as the number of participating nodes increases, validating its applicability in large-scale distributed environments. However, the study also reveals that federated learning is not inherently immune to adversarial threats. The observed performance degradation under model poisoning attacks highlights the vulnerability of FL systems in untrusted environments. Nevertheless, the incorporation of robust aggregation techniques significantly reduced the impact of such attacks, underscoring the importance of integrating security mechanisms within the FL architecture. Overall, the discussion emphasizes that while FL offers substantial

advantages over centralized approaches, its effectiveness depends on careful design, particularly in terms of security and communication efficiency.

Conclusion

This study successfully developed and evaluated an AI-driven federated learning framework for privacy-preserving early detection of cyber threats in Pakistan's critical infrastructure systems. The results demonstrate that the proposed framework outperforms traditional centralized models in terms of detection accuracy, efficiency, scalability, and privacy preservation. By enabling collaborative learning without sharing sensitive data, the framework addresses a critical gap in existing cybersecurity approaches.

The integration of advanced AI techniques with federated learning not only enhances threat detection capabilities but also ensures compliance with data privacy requirements. Additionally, the incorporation of adversarial defense mechanisms strengthens the robustness of the system against malicious attacks. Overall, the study confirms that federated learning is a viable and effective paradigm for securing distributed critical infrastructure systems in a privacy-conscious manner.

Implications

The findings of this research have significant theoretical and practical implications. Theoretically, the study contributes to the growing body of knowledge on federated learning by demonstrating its applicability in cybersecurity for critical infrastructure systems, particularly in developing country contexts. It also highlights the importance of integrating privacy-preserving and adversarial defense mechanisms within distributed learning frameworks. Practically, the proposed framework offers a scalable and secure solution for policymakers, cybersecurity professionals, and infrastructure operators in Pakistan. It provides a foundation for implementing decentralized cybersecurity systems that enhance collaboration while maintaining data confidentiality. The study also supports the development of national cybersecurity strategies by emphasizing the need for advanced, AI-driven, and privacy-aware technologies in protecting critical services.

Future Directions

Future research can build upon this study by exploring several advanced dimensions of federated learning in cybersecurity. First, the integration of blockchain technology can be investigated to enhance trust, transparency, and accountability among participating entities. Second, the development of more efficient communication protocols and model compression techniques can help reduce bandwidth consumption and improve scalability.

Additionally, future studies may focus on real-world deployment and validation of the proposed framework using live data from critical infrastructure systems in Pakistan. The incorporation of explainable AI (XAI) techniques can also improve model interpretability, enabling better understanding of decision-making processes in threat detection. Furthermore, exploring adaptive and self-learning models capable of responding to evolving cyber threats in dynamic environments remains an important research direction.

Recommendations

Based on the findings of this study, several recommendations are proposed. First, organizations managing critical infrastructure systems should consider adopting federated learning-based cybersecurity frameworks to enhance threat detection while preserving data privacy. Second, it is recommended to integrate robust security mechanisms, such as secure aggregation and Byzantine-resilient algorithms, to protect against adversarial attacks.

Policymakers should develop regulatory frameworks and guidelines that support secure data collaboration and the adoption of privacy-preserving technologies. Investment in infrastructure and technical expertise is also essential to facilitate the implementation of advanced AI-driven cybersecurity solutions. Additionally, cross-sector collaboration should be encouraged to enable effective knowledge sharing and collective defense against cyber threats.

Limitations

Despite its contributions, this study has several limitations. First, the use of simulated and benchmark datasets, rather than real-time data from

Pakistan's critical infrastructure, may limit the generalizability of the findings. Second, the experimental setup, although designed to mimic real-world conditions, may not fully capture the complexity and variability of actual operational environments.

Third, the study focused primarily on a specific federated learning architecture and may not account for the performance of alternative models or configurations. Additionally, while adversarial robustness was evaluated, the range of attack scenarios considered was limited, and more sophisticated threat models could be explored in future research. Finally, the computational and communication costs associated with federated learning, although manageable, may pose challenges in resource-constrained environments.

REFERENCES

- Abd Elaziz, M., Fares, I. A., Dahou, A., & Shrahili, M. (2025). Federated learning framework for IoT intrusion detection using tab transformer and nature-inspired optimization. *Frontiers in Big Data*, 8, 1526480.
- Begum, K., Mozumder, M. A. I., Joo, M.-I., & Kim, H.-C. (2024). Blockchain-driven federated learning for intrusion detection in IoMT networks. *Sensors*, 24(14), 4591.
- Bukhari, S. M. S., Zafar, M. H., Abou Houran, M., & Qadir, Z. (2024). Asynchronous federated learning for cybersecurity in edge-enabled industrial IoT networks. *Internet of Things*, 27, 101252.
- Deng, S., Zhang, L., & Yue, D. (2024). Privacy-preserving risk assessment in smart grids using federated learning. *Communications Engineering*, 3, 154.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- Huang, H. J., Iskandarov, B., Rahman, M., Otal, H. T., & Canbaz, M. A. (2024). Federated learning in adversarial environments: Poisoning attacks and defenses. *arXiv preprint*.
- Kairouz, P., McMahan, H. B., Avent, B., et al. (2021). Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1-2), 1-210.

- Li, T., Sahu, A. K., Talwalkar, A., & Smith, V. (2020). Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3), 50–60.
- McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & Agüera y Arcas, B. (2017). Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS* (pp. 1273–1282).
- Meng, R., Shah, A. A., Jamshed, M. A., & Pezaros, D. (2024). Federated learning-based intrusion detection framework for IoT-enabled critical infrastructure. In *IEEE ICC Workshops* (pp. 1–6).
- Paulraj, J., Raghuraman, B., Gopalakrishnan, N., & Otoum, Y. (2025). AI-based cybersecurity framework for critical infrastructure protection. *arXiv preprint*.
- Rahmati, M., & Pagano, A. (2025). Federated learning-driven cybersecurity for IoT networks: Privacy-preserving threat detection. *Informatics*, 12(3), 62.
- Rahmati, M., & Rahmati, N. (2026). Byzantine-robust federated learning with secure aggregation for critical infrastructure security. *arXiv preprint*.
- Shokri, R., & Shmatikov, V. (2015). Privacy-preserving deep learning. In *Proceedings of ACM CCS* (pp. 1310–1321).
- Subramanian, G., & Chinnadurai, M. (2024). Quantum-enhanced federated learning for cyberattack detection. *Scientific Reports*, 14, 32038.
- War, M. R., Singh, Y., Sheikh, Z. A., & Singh, P. K. (2025). Federated learning for cybersecurity in cyber-physical systems: A review. *Scalable Computing: Practice and Experience*, 26(1), 16–33.
- Zhang, Q., Chen, M., & Saad, W. (2022). Federated learning for edge computing: A survey. *IEEE Communications Surveys & Tutorials*, 24(1), 350–386.

