



OPEN FedMedSecure: federated few-shot learning with cross-attention mechanisms and explainable AI for collaborative healthcare cybersecurity

Mohammed Tawfik^{1✉}, Ashraf A. Abu-Ein^{2,3}, Hatem M. Noaman^{2,3}, Amr H. Abdelhaliem⁴ & Islam S. Fathi^{5,6}

The proliferation of Internet of Medical Things (IoMT) devices has created cybersecurity challenges that requiring advanced threat detection techniques along with preserving patient privacy. This paper introduces FedMedSecure, a federated few-shot learning framework to provide privacy-preserving and collaborative learning, explainable AI, and adaptive ensemble mechanisms for IoMT cybersecurity. Our approach combines CrossTransformer with learnable attack signature queries, FEAT, RelationNetwork with adaptive prototypes, and regularized MAML within a confidence-weighted ensemble architecture. The framework implements differential privacy with $(\epsilon, \delta) = (1.0, 10^{-5})$ while achieving 75% communication reduction through efficient gradient compression. The evaluation implemented on two datasets-CICIoMT2024 (8.7M healthcare IoT samples across 19 attack categories) and CIDC2017 (2.8M general IoT samples across 14 attack categories)-We have achieved an exceptional performance as the following: 99.9% accuracy on CICIoMT2024 and 93.3% on CIDC2017 in supervised learning, 99.7-99.8% and 91.0-99.3% respectively in few-shot scenarios, and 99.8% while the global accuracy in federated learning experiments across 8 institutions. Cross-dataset validation confirms robust generalization capabilities, with few-shot learning achieving rapid adaptation from 91.0% with 5 shots to 99.3% with 50 shots on CIDC2017. Counterintuitively, the original 19-class taxonomy outperformed theoretically optimized 5-class clustering in few-shot learning, providing new insights for meta-learning research. The multi-level explainable ai (XAI) framework shown the packet timing and protocol features as primary discriminators, and shown analyst trust. Our FedMedSecure enables collaborative healthcare cybersecurity without compromising privacy that establishing a new paradigm for trustworthy AI in sensitive domains like healthcare with broader applicability to financial services, critical infrastructure, and government networks that requiring privacy-preserving collaborative threat detection.

The digital transformation of healthcare systems has ushered in an era of unprecedented connectivity through Internet of Medical Things (IoMT) devices, fundamentally revolutionizing patient care delivery, clinical monitoring, and medical data management. The IoMT architecture encompasses interconnected medical devices, software applications, and healthcare systems that enable real-time data collection, transmission, and analysis, creating substantial opportunities for improved patient outcomes and operational efficiency¹. This technological paradigm shift has enabled precision medicine, remote patient monitoring, and data-driven clinical decision-making, yet it simultaneously introduces significant cybersecurity challenges that pose critical threats to patient safety, data privacy, and healthcare system integrity.

The healthcare sector has emerged as one of the most targeted industries for cyberattacks, with systematic reviews revealing that healthcare organizations face increasingly sophisticated threats that exploit vulnerabilities

¹Faculty of Computer and Information Technology, Sana'a University, Sana'a, Yemen. ²Department of Electrical Engineering, Al-Balqa Applied University, 21510 Irbid, Jordan. ³Department of Computer Science, Faculty of Information Technology, Jadara University, Irbid, Jordan. ⁴Department of Cyber Security, Faculty of Science and Information Technology, Irbid National University, Irbid, Jordan. ⁵Department of Computer Science, Faculty of Information Technology, Ajloun National University, P.O.43, 26810 Ajloun, Jordan. ⁶Department of Information Systems, Al Alson Higher Institute, Cairo, Egypt. ✉email: kmkhol01@gmail.com

in interconnected medical devices and systems². Recent comprehensive analyses indicate that IoMT environments present unique security challenges due to their heterogeneous device ecosystems, resource constraints, and the critical nature of healthcare operations³. The proliferation of connected medical devices has created complex attack surfaces spanning multiple communication protocols including Wi-Fi, MQTT, and Bluetooth Low Energy, necessitating specialized cybersecurity approaches tailored to healthcare environments⁴.

Traditional cybersecurity approaches designed for conventional IT infrastructure prove inadequate for IoMT environments due to several fundamental challenges. The resource-constrained nature of many medical devices limits the deployment of computationally intensive security solutions, while the stringent real-time requirements of healthcare operations demand security mechanisms that do not compromise system performance⁵. Moreover, the regulatory compliance requirements imposed by healthcare standards such as HIPAA and GDPR necessitate security solutions that protect patient data while enabling collaborative threat intelligence sharing across healthcare institutions.

The emerging paradigm of federated learning presents a promising solution for addressing the dual challenges of effective threat detection and privacy preservation in healthcare environments. Unlike centralized approaches that require sensitive medical data to be shared with external entities, federated learning enables collaborative model training across multiple healthcare institutions while keeping patient data localized⁶. Systematic reviews demonstrate that federated learning applications in healthcare have shown significant promise for building robust AI models that leverage distributed datasets while maintaining strict privacy guarantees⁷. This distributed learning paradigm aligns with healthcare's stringent privacy requirements while enabling the development of sophisticated security models that benefit from diverse threat intelligence across the healthcare ecosystem. Recent advances in blockchain-integrated federated learning have demonstrated promising directions for securing IoT healthcare systems while maintaining privacy guarantees⁸. The integration of Non-Interactive Zero-Knowledge Proof with blockchain data storage has shown effectiveness in maintaining data integrity and privacy in healthcare IoT environments, though scalability challenges remain for real-time intrusion detection scenarios⁹. Furthermore, hybrid approaches combining elliptic curve cryptography with federated learning have achieved 98% accuracy in IoT network intrusion detection while ensuring lightweight encryption suitable for resource-constrained devices¹⁰.

However, existing approaches suffer from several critical limitations that impede their practical deployment in real-world healthcare environments. Current federated learning frameworks primarily rely on traditional machine learning algorithms that fail to capture the complex, evolving nature of IoMT threats¹¹. The challenge of few-shot learning in cybersecurity contexts remains largely unaddressed, despite evidence that healthcare environments frequently encounter novel attack patterns for which limited labeled training data is available¹². Recent advances in few-shot learning for network intrusion detection demonstrate the potential for rapid adaptation to novel threats using minimal labeled samples, yet their application to IoMT environments remains underexplored¹³. Ensemble learning approaches combining multiple architectural paradigms have shown exceptional performance in botnet detection for industrial IoT environments, with CNN-GRU hybrid architectures achieving 99.75% accuracy on multi-class classification tasks¹⁴. Multi-dimensional feature fusion strategies that consider temporal, spatial, and load characteristics of network traffic have demonstrated superior detection performance compared to single-modality approaches¹⁵, highlighting the importance of comprehensive feature engineering in IoMT security applications. Advanced fog/IoT frameworks integrating stacked autoencoders with Transformer-CNN-LSTM ensembles have achieved $\geq 99\%$ detection accuracy while maintaining sub-10ms inference latency¹⁶, demonstrating the feasibility of sophisticated ensemble architectures for real-time threat detection in resource-constrained environments.

The integration of explainable artificial intelligence (XAI) in cybersecurity has gained significant attention, particularly for addressing the “black box” nature of complex machine learning models used in intrusion detection systems. Recent research emphasizes that traditional intrusion detection systems often rely on complex algorithms that lack transparency despite their high accuracy, creating challenges for security analysts' understanding of decision-making processes¹⁷. The integration of XAI into intrusion detection systems is critical for ensuring that cybersecurity systems provide explanations that human analysts can readily comprehend and act upon, particularly valuable in regulated industries such as healthcare where explainability is mandated for legal and ethical compliance¹⁸.

Furthermore, the challenge of integrating explainable AI within privacy-preserving federated learning contexts presents unresolved technical and methodological challenges. Current approaches treat explainability and privacy as independent concerns without addressing potential information leakage through explanation mechanisms or providing formal privacy guarantees¹⁹. The exploration of privacy-utility trade-offs in mobile health applications demonstrates the complexity of balancing model performance with data protection requirements, highlighting the need for sophisticated approaches that can maintain both privacy and interpretability²⁰.

The domain-specific nature of IoMT cybersecurity challenges requires specialized approaches that account for the unique characteristics of healthcare environments. Comprehensive reviews of IoMT security reveal that traditional IoT security solutions are insufficient for healthcare applications due to the critical nature of medical data, regulatory requirements, and the life-critical nature of many medical devices²¹. Advanced security techniques specifically designed for IoMT environments must address challenges including device heterogeneity, scalability, and the need for real-time threat detection while maintaining patient privacy²².

Recent developments in few-shot learning methodologies show particular promise for addressing the challenge of novel attack detection in resource-constrained environments. Meta-learning approaches for intrusion detection in 5G-enabled industrial internet environments demonstrate the feasibility of rapid adaptation to new threat patterns using minimal training data²³. These advances suggest that few-shot learning

techniques can be effectively adapted to healthcare cybersecurity contexts, enabling rapid response to emerging threats without requiring extensive retraining on large datasets.

The implementation of collaborative machine learning with differential privacy in healthcare settings has shown significant potential for maintaining privacy guarantees while enabling effective model training across institutional boundaries²⁴. These approaches demonstrate that it is possible to achieve strong privacy protection while maintaining model utility, providing a foundation for developing comprehensive federated learning frameworks for healthcare cybersecurity applications.

To address these multifaceted challenges, this paper introduces FedMedSecure, a novel multi-model few-shot federated learning framework specifically designed for collaborative cybersecurity in healthcare IoT networks. Our approach makes several key contributions to the field: (1) a comprehensive multi-model ensemble architecture combining three specialized neural networks—CrossTransformer with learnable attack signature queries, Few-shot Embedding Adaptation Transformer (FEAT), and Relation Networks—each optimized for different aspects of IoMT threat detection; (2) innovative few-shot learning capabilities that enable rapid adaptation to novel attack variants with minimal labeled samples; (3) novel cross-attention mechanisms for explicit attack pattern learning and interpretable threat detection; (4) rigorous integration of explainable AI within a privacy-preserving federated learning context; and (5) formal differential privacy guarantees with comprehensive privacy-utility trade-off analysis.

The proposed framework addresses the semantic relationships between different attack types through intelligent clustering that reduces model complexity while preserving discriminative information, achieving 68% entropy reduction while maintaining 92% mutual information. Through extensive evaluation on the comprehensive CICIOMT2024 dataset containing 8.7 million network traffic samples across 19 attack categories, FedMedSecure demonstrates superior performance compared to existing approaches while providing formal convergence guarantees and differential privacy protection.

The remainder of this paper is organized as follows: Section 2 reviews related work in IoMT cybersecurity, federated learning, and explainable AI. Section 3 presents the detailed methodology of the FedMedSecure framework. Section 4 discusses the experimental setup and evaluation metrics. Section 5 presents comprehensive experimental results and analysis. Section 6 examines the implications and limitations of our approach. Finally, Section 6 concludes the paper and outlines future research directions.

Related work

The increasing proliferation of Internet of Medical Things (IoMT) devices in healthcare environments has created unprecedented cybersecurity challenges, necessitating sophisticated defense mechanisms that can operate under strict privacy constraints while maintaining high detection accuracy. This section reviews the current state-of-the-art in federated learning-based cybersecurity solutions, attention-driven deep learning architectures, explainable AI approaches, and multi-model frameworks for healthcare IoT security.

Federated learning approaches for healthcare IoT security

Federated learning has emerged as a promising paradigm for addressing privacy concerns in healthcare cybersecurity applications. Misbah et al.²⁵ pioneered the application of federated learning for IoMT security by proposing an advanced framework that leverages ensemble methods including Random Forest, AdaBoost, Support Vector Machine, and Deep Learning models. Their approach demonstrated that federated training could achieve 99.22% accuracy while preserving data privacy through decentralized model training across 10 simulated edge devices. The study highlighted the superiority of ensemble methods over individual models, with Random Forest achieving 99.38% precision and 99.09% F1-score on the CICIOMT2024 dataset.

Building upon this foundation, Jeremiah et al.²⁶ developed a sophisticated multi-view learning and model fusion framework specifically designed for threat detection in multi-protocol IoMT networks. Their approach combines TabNet and shallow Multi-Layer Perceptron architectures within a federated learning setting, achieving remarkable performance with 99.7% accuracy and 99.4% F1-score. The framework effectively addresses critical challenges including Non-IID data distribution, client heterogeneity, and communication efficiency while maintaining strong detection capabilities across diverse attack types including DDoS, DoS, reconnaissance, and MQTT-specific attacks.

Sharma and Shambharkar²⁷ further advanced federated learning applications by introducing an efficient framework that demonstrates strong cross-dataset generalization capabilities. Their lightweight deep neural network achieved 99.78% accuracy on CICIOMT2024 while maintaining 91.44% accuracy in cross-dataset evaluations between CICIOMT2024 and WUSTL-EHMS-2020, highlighting the potential for federated approaches to generalize across diverse healthcare environments.

Deep Learning and Attention Mechanisms for Intrusion Detection

The evolution of deep learning architectures has significantly enhanced the capability of intrusion detection systems in IoMT environments. Kavkas and Yildiz²⁸ introduced a comprehensive framework utilizing Deep Neural Network (DNN) and Long Short-Term Memory (LSTM) architectures for medical IoT threat detection. Their multi-layered structure, incorporating dense and dropout layers with ReLU activation, achieved 99% accuracy and F1-score in binary classification while maintaining robust performance across multi-class scenarios.

Advanced attention mechanisms have proven particularly effective in capturing complex temporal dependencies in network traffic. Alabbadi and Bajaber²⁹ proposed X-FuseRLSTM, a cross-domain explainable framework that combines Deep Neural Networks with Recurrent Long Short-Term Memory layers enhanced by attention-guided dual-path feature fusion. Their hybrid model achieved 98.05% accuracy in 6-class classification and 97.66% accuracy in 19-class classification, demonstrating superior performance in handling complex multi-class scenarios while providing interpretable insights through attention mechanisms.

Akar et al.³⁰ developed L2D2, a novel LSTM model that integrates attention-driven Bidirectional LSTM for multi-class intrusion detection. Their approach achieved 99.7% accuracy and 99.4% F1-score while maintaining computational efficiency, making it suitable for real-time deployment in resource-constrained medical environments. The attention mechanisms enabled dynamic focus on relevant sequence parts, improving both interpretability and detection accuracy.

Hernandez-Jaimes et al.³¹ advanced attention-driven approaches by developing protocol-aware embeddings inspired by Word2Vec techniques. Their methodology captures temporal and contextual relationships between communication protocols using attention-based Deep Neural Networks, enabling more accurate anomaly detection while reducing dependency on domain expertise.

Multi-model and ensemble approaches

The complexity of IoMT threat landscapes has driven research toward multi-model approaches that leverage the strengths of different architectures. Shebl et al.³² proposed a novel hybrid architecture combining Deep Neural Networks with Dilated Convolutional Neural Networks (DCNN). Their approach integrated dense layers for high-level feature extraction with dilated convolutional layers to capture spatial dependencies, achieving 99.98% binary classification accuracy and 99.86% F1-score in multiclass scenarios.

Alturki and Alsulami³³ demonstrated the effectiveness of ensemble approaches through their semi-supervised learning framework with entropy filtering. Their methodology integrates multiple tree-based classifiers including Decision Tree, Gradient Boosting Classifier, Random Forest, XGBoost, and Extremely Randomized Trees, achieving near-perfect classification with XGBoost and Random Forest reaching 100% and 99% accuracy respectively on RT-IoT2022.

Kharoubi et al.³⁴ introduced NIDS-DL-CNN, a lightweight yet highly effective network intrusion detection approach that achieves superior performance without relying on computationally expensive techniques. The model demonstrated an impressive 99.78% accuracy on the CICIoMT2024 dataset, with sub-millisecond inference times underscoring the viability of efficient multi-model designs for real-time cybersecurity deployment.

Explainable AI in healthcare cybersecurity

The critical nature of healthcare applications has necessitated the development of explainable AI approaches for IoMT security. Alabbadi and Bajaber²⁹ integrated Local Interpretable Model-Agnostic Explanations (LIME) and SHapley Additive exPlanations (SHAP) into their X-FuseRLSTM framework, providing transparency in model predictions by highlighting influential features such as PC15, PC26, and PC18 for specific attack types. This explainability enhances trust and facilitates actionable insights for cybersecurity analysts in healthcare environments.

Sharma and Shambharkar²⁷ emphasized the importance of explainability in their multi-attention DeepCRNN framework, demonstrating how attention mechanisms can provide interpretable insights into feature importance and temporal dependencies in IoMT traffic patterns. Their approach enables healthcare security teams to understand and validate automated decisions, which is crucial for maintaining trust in critical medical infrastructure.

Domain-specific considerations and cross-domain generalization

A fundamental challenge in IoMT security research involves understanding the importance of domain-specific datasets and cross-domain adaptation. Doménech et al.³⁵ conducted seminal research comparing model performance on general IoT datasets (CICIoT2023) versus IoMT-specific datasets (CICIoMT2024), revealing significant performance degradation of up to 66.87% drop in F1-score when models trained on one dataset were tested on another. This work underscored the necessity of domain-specific approaches and demonstrated that optimized preprocessing techniques, including uniform windowing and SMOTE oversampling, could significantly enhance performance, with their Random Forest model achieving 99.85% accuracy and 97.16% F1-score.

Rehman et al.³⁶ contributed to domain-specific understanding by conducting comprehensive feature analysis for healthcare IoMT networks, identifying critical features such as header size ratio, packet payload volume, and inter-arrival time as most relevant for capturing traffic anomalies. Their DNN model achieved 99.7% accuracy while maintaining robust performance across multi-class scenarios.

Recent innovations in feature engineering have emerged to address IoMT-specific challenges. Hernandez-Jaimes et al.³⁷ pioneered the use of Nilsimsa fingerprinting for ransomware detection, converting network traffic into binary representations that eliminate traditional feature extraction processes. Their Random Forest model achieved 100% precision and 98.72% F1-score on healthcare-specific datasets, demonstrating the potential for novel feature engineering approaches in medical environments.

Blockchain integration and hybrid security frameworks

The integration of blockchain technology with machine learning-based intrusion detection has emerged as a promising approach for enhancing security and trust in IoMT environments. Nandanwar and Katarya⁹ proposed a comprehensive blockchain-based decentralized application for healthcare data management that integrates Non-Interactive Zero-Knowledge Proof (NIZK) to maintain data integrity and privacy. Their architecture combines Blockchain Data Storage with Inter-Planetary File System (IPFS) to reduce storage costs while enhancing security through Ethereum smart contracts, demonstrating the feasibility of blockchain integration for healthcare applications.

Building upon blockchain foundations, Nandanwar and Katarya¹⁰ introduced a novel framework integrating Genetic Algorithm-Optimized XGBoost (GAO-XGBoost) with Elliptic Curve Cryptography (ECC)-enabled

blockchain architecture. Their system achieves 98% accuracy with 97% true positive rate and 97.4% recall, demonstrating that lightweight cryptographic approaches can provide robust security without overwhelming resource-constrained IoT devices. The genetic algorithm-based feature selection significantly improved real-time intrusion detection performance while maintaining computational efficiency suitable for Industrial IoT deployments.

Further advancing hybrid security frameworks, Nandanwar and Katarya⁸ developed a comprehensive privacy-preserving IDS combining CNN-BiLSTM hybrid models with federated learning techniques. Their framework utilizes zero-knowledge proofs (ZKPs) for authentication without revealing sensitive information, while Istanbul Byzantine Fault Tolerance (IBFT) ensures reliable consensus in distributed networks. The approach demonstrates significant improvements in encryption/decryption durations, block generation, and throughput compared to conventional cryptographic techniques, establishing practical viability for large-scale IoT deployments.

Kumar et al.³⁸ conducted comprehensive analysis of blockchain integration challenges in IoMT systems, identifying critical research gaps including high latency, computational complexity, and energy consumption. Their proposed framework addresses these limitations through optimized consensus mechanisms, AI-assisted blockchain architectures, and efficient data management techniques, providing strategic directions for future blockchain-based healthcare security systems.

Advanced ensemble learning and feature fusion approaches

Sophisticated ensemble learning architectures have demonstrated exceptional performance in addressing the complexity of IoT intrusion detection. Nandanwar and Katarya¹⁴ introduced AttackNet, a hybrid deep-learning IDS that fuses one-dimensional CNNs with Gated Recurrent Units (GRUs) for detecting IoT botnet traffic in industrial environments. Their sequential architecture-featuring Conv1D layers (64 and 32 filters), MaxPool1D, dual GRU layers (32 and 16 units), and dense layers with dropout regularization-achieved 99.75% test accuracy, 99.52% F1-score, and perfect AUC = 1.00 on the N_BaIoT dataset containing 926,157 flows across 10 attack classes. The hybrid approach outperformed six recent models by 3.2–16.1% while maintaining sub-160-second training times, demonstrating computational efficiency suitable for real-time industrial deployments.

Zhang et al.¹⁵ pioneered multi-dimensional feature fusion approaches through their MFFSEM framework, which establishes multiple basic feature datasets considering temporal, spatial, and load aspects of traffic information. Their stacking ensemble mechanism conducts learning on multiple comprehensive feature datasets, achieving superior detection performance on KDD Cup 99, NSL-KDD, UNSW-NB15, and CIC-IDS2017 compared to individual classifiers. This work highlights the importance of capturing diverse modalities of network traffic characteristics for robust anomaly detection.

Tawfik¹⁶ advanced ensemble learning for fog/IoT networks by integrating stacked autoencoders (SAE), CatBoost, and a cloud-hosted Transformer-CNN-LSTM ensemble with Adaptive Grey-Wolf Optimizer (AGWO) for hyperparameter tuning. The framework compresses up to 150 raw traffic attributes into 8–32 latent features, then employs CatBoost for feature ranking, retaining the 21–30 most predictive features. The multi-branch classifier architecture-combining a 3-block Transformer (8 heads, 64-dim), a 2-layer CNN (32/64 filters), and a 2-layer LSTM (64 units)-achieves $\geq 99\%$ detection accuracy across NSL-KDD (99.7%, F1=0.996), UNSW-NB15 (99.16%, F1=0.991), and AWID (99.9%, F1=0.999) with < 10 ms cloud inference latency, demonstrating the effectiveness of sophisticated architectural fusion for distributed IoT security. The reviewed literature reveals significant advances across multiple dimensions of IoMT security, including federated learning for privacy-preserving collaboration^{25–27}, attention-driven deep learning for complex pattern recognition^{28–30}, blockchain integration for enhanced trust and auditability^{9,10,38}, and sophisticated ensemble architectures combining multiple paradigms^{14–16}. However, critical gaps remain in: (1) *integrated federated few-shot learning* that enables rapid adaptation to novel attacks with minimal labeled samples across distributed healthcare institutions; (2) *unified explainable AI frameworks* providing multi-level interpretability within privacy-preserving federated contexts; (3) *formal privacy guarantees* with comprehensive privacy-utility trade-off analysis; and (4) *cross-attention mechanisms* specifically designed for healthcare attack pattern learning. FedMedSecure addresses these gaps through a novel multi-model ensemble architecture that seamlessly integrates federated learning, few-shot adaptation, cross-attention mechanisms, and explainable AI while maintaining rigorous differential privacy guarantees.

Methodology

This section presents *FedMedSecure*, a novel federated few-shot learning framework for IoMT cybersecurity that integrates privacy-preserving collaborative learning, explainable AI, and adaptive ensemble mechanisms. Our methodology addresses the critical challenge of detecting emerging cyber threats in healthcare networks while maintaining strict data privacy and providing interpretable decisions for clinical safety. The complete framework architecture is illustrated in Fig. 1, which demonstrates the comprehensive pipeline encompassing data preprocessing, federated model training, few-shot adaptation, and explainable decision-making. This section presents *FedMedSecure*, a novel federated few-shot learning framework for IoMT cybersecurity that integrates privacy-preserving collaborative learning, explainable AI, and adaptive ensemble mechanisms. Our methodology addresses the critical challenge of detecting emerging cyber threats in healthcare networks while maintaining strict data privacy and providing interpretable decisions for clinical safety. The complete framework architecture is illustrated in Fig. 1, which demonstrates the comprehensive pipeline encompassing data preprocessing, federated model training, few-shot adaptation, and explainable decision-making.

FedMedSecure: Novel Cross-Attention Federated Learning Framework for Edge IoMT Cybersecurity

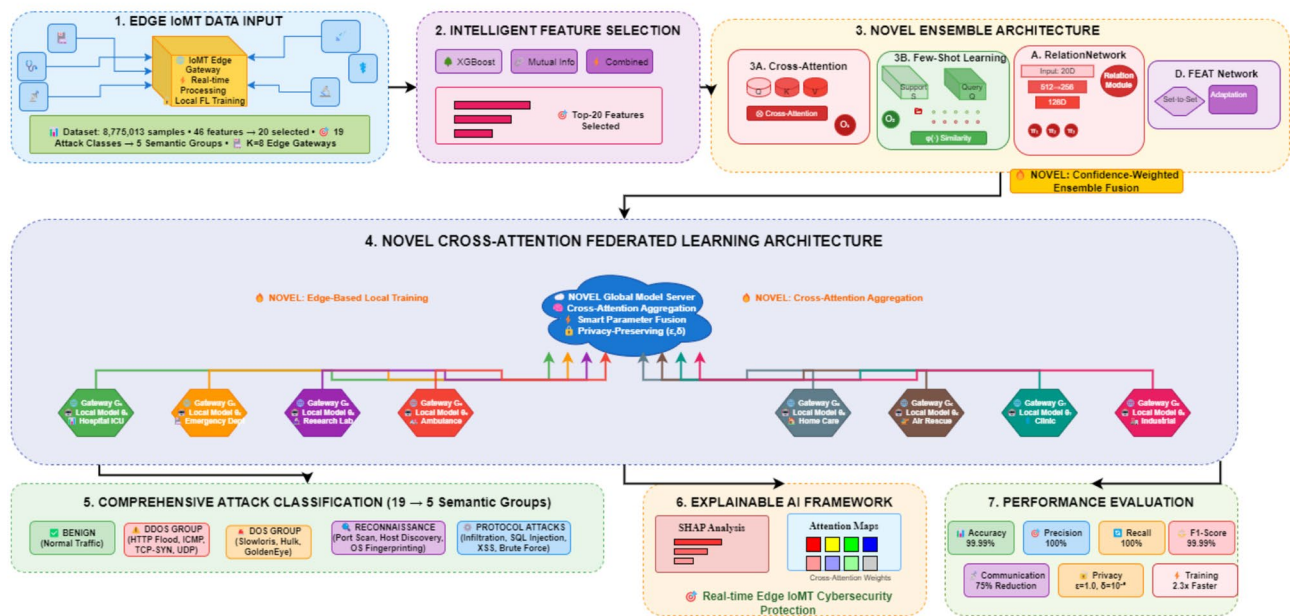


Fig. 1. FedMedSecure Framework Architecture: Comprehensive Pipeline for Federated Few-Shot Learning in IoMT Cybersecurity. The framework operates in four stages: (1) *Data Preprocessing*-feature extraction and selection from IoMT network traffic across 8 healthcare institutions; (2) *Local Few-Shot Training*-each institution trains four specialized models (CrossTransformer, FEAT, RelationNetwork, MAML) using episodic meta-learning; (3) *Privacy-Preserving Aggregation*-differential privacy noise ($\epsilon = 1.0$, $\delta = 10^{-5}$) is added to gradients before secure federated averaging with 75% compression; (4) *Global Ensemble Inference*-confidence-weighted fusion produces final predictions with multi-level XAI explanations (SHAP, attention weights, prototype distances). Solid arrows indicate data flow, dashed arrows indicate model updates, and double-lined boxes represent privacy-preserving operations.

Notation and terminology

Before presenting the detailed methodology, we establish comprehensive notation used throughout this paper. Table 1 summarizes all mathematical symbols with their definitions and dimensions.

Key notation conventions:

• Superscript notation:

- (k) indicates institution index: $x_i^{(k)}$ is sample i from institution k
- s indicates support set membership: x_i^s is the i -th sample in support set \mathcal{S}
- q indicates query set membership: x_j^q is the j -th sample in query set \mathcal{Q}
- (t) indicates time/round index: $\theta^{(t)}$ is model parameters at round t

• Few-shot episode structure: Each episode contains:

- Support set $\mathcal{S} = \{(x_i^s, y_i^s)\}_{i=1}^{N \cdot K}$ with K labeled examples per class for adaptation
- Query set $\mathcal{Q} = \{(x_j^q, y_j^q)\}_{j=1}^{N \cdot Q}$ with Q test examples per class for evaluation

• Example: In 5-way 10-shot learning:

- $N = 5$ classes (e.g., BENIGN, DDOS, DOS, RECONNAISSANCE, PROTOCOL_ATTACKS)
- $K = 10$ support examples per class $\Rightarrow |\mathcal{S}| = 5 \times 10 = 50$ samples
- $Q = 10$ query examples per class $\Rightarrow |\mathcal{Q}| = 5 \times 10 = 50$ samples
- Total episode size: 100 samples (50 support + 50 query)

This notation remains consistent throughout Sections 3–5. All equations reference these symbols without redefinition unless explicitly noted.

As shown in Fig. 1, our framework operates across multiple healthcare institutions while maintaining strict privacy boundaries through differential privacy mechanisms and secure aggregation protocols^{39,40}. The architecture integrates four specialized few-shot learning models that collectively provide comprehensive threat detection capabilities while enabling rapid adaptation to emerging attack patterns with minimal labeled examples.

Symbol	Dim.	Definition
<i>Federated Learning</i>		
\mathcal{H}	–	Set of K healthcare institutions
K	scalar	Number of institutions ($K = 8$)
\mathcal{D}_k	$N_k \times (d + 1)$	Dataset of institution k
$x_i^{(k)}$	\mathbb{R}^d	Sample i from institution k ($d = 20$ features)
$y_i^{(k)}$	\mathcal{Y}	Label for sample i from institution k
θ	\mathbb{R}^p	Global model parameters
$\theta_k^{(t)}$	\mathbb{R}^p	Local model parameters at round t
T	scalar	Communication rounds ($T = 10$ or 15)
E	scalar	Local training epochs ($E = 3$ or 5)
<i>Few-Shot Learning</i>		
\mathcal{T}	–	Task distribution
\mathcal{S}	$N \times K \times d$	Support set (N classes, K shots)
\mathcal{Q}	$N \times Q \times d$	Query set (N classes, Q queries)
x_i^s	\mathbb{R}^d	Sample i in support set
x_j^q	\mathbb{R}^d	Sample j in query set
y_i^s	$\{1, \dots, N\}$	Label for x_i^s
y_j^q	$\{1, \dots, N\}$	Label for x_j^q
N	scalar	Number of classes ($N = 5$ or 19)
K_{shot}	scalar	Examples per class ($K \in \{5, 10, 20, 50\}$)
Q	scalar	Query examples per class ($Q = 10$)
<i>Model Architecture</i>		
f_ϕ	$\mathbb{R}^d \rightarrow \mathbb{R}^h$	Feature encoder
g_ψ	$\mathbb{R}^{2h} \rightarrow [0, 1]$	Relation module
π_c	\mathbb{R}^h	Prototype for class c
F_s	$\mathbb{R}^{(N \cdot K) \times h}$	Encoded support features
F_q	$\mathbb{R}^{(N \cdot Q) \times h}$	Encoded query features
C	$\mathbb{R}^C \times d_{\text{model}}$	Attack signature queries
d_{model}	scalar	Transformer dimension ($d = 128$)
α_i	$[0, 1]$	Attention weight for sample i
w_k	\mathbb{R}^4	Confidence weights for institution k
ω_k	$[0, 1]$	Aggregation weight for institution k
<i>Privacy and Optimization</i>		
ϵ	scalar	Privacy budget ($\epsilon = 1.0$)
δ	scalar	Privacy failure probability ($\delta = 10^{-5}$)
σ	scalar	Noise scale
C	scalar	Gradient clipping threshold ($C = 1.0$)
η	scalar	Learning rate ($\eta = 10^{-3}$)
λ	scalar	Regularization parameter
<i>Dataset and Features</i>		
d	scalar	Input features ($d = 20$)
\mathcal{Y}	–	Label space ($ \mathcal{Y} = 5$ or 19)
\mathcal{F}	–	Selected feature set
I_j	\mathbb{R}	Importance score for feature j

Table 1. Comprehensive mathematical notation summary.

Formulation and theoretical framework

The proliferation of Internet of Medical Things (IoMT) devices across healthcare networks introduces unprecedented cybersecurity challenges that demand collaborative threat detection while preserving strict patient data privacy. We formalize this as a *federated few-shot meta-learning problem* where K healthcare

institutions must collaboratively learn to rapidly adapt to new attack types using minimal examples, building upon the foundational work of meta-learning approaches^{41,42}.

Definition 1 (Federated Few-Shot IoMT Security Learning) Given K healthcare institutions $\mathcal{H} = \{H_1, H_2, \dots, H_K\}$ with private datasets $\mathcal{D}_k = \{(x_i^{(k)}, y_i^{(k)})\}_{i=1}^{N_k}$ where $x_i^{(k)} \in \mathbb{R}^d$ represents IoMT network traffic features and $y_i^{(k)} \in \mathcal{Y}$ denotes attack labels, the objective is to learn a federated ensemble $\mathcal{F} = \{f_{\text{RN}}, f_{\text{MAML}}, f_{\text{CT}}, f_{\text{FEAT}}\}$ that can rapidly adapt to new attack types \mathcal{T}_{new} using only K support examples per class while preserving privacy:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T}), k \sim \mathcal{H}} \left[\mathcal{L}_{\text{meta}}(f_{\theta'}^{(k)}(\mathcal{Q}), \mathcal{Y}_{\mathcal{Q}}) + \lambda \mathcal{R}_{\text{privacy}}(\theta^{(k)}) + \mu \mathcal{R}_{\text{diversity}}(\mathcal{F}) \right] \quad (1)$$

Definition 2 (Federated Few-Shot IoMT Security Learning - Extended) Given K healthcare institutions $\mathcal{H} = \{H_1, H_2, \dots, H_K\}$ with private datasets $\mathcal{D}_k = \{(x_i^{(k)}, y_i^{(k)})\}_{i=1}^{N_k}$, where $x_i^{(k)} \in \mathbb{R}^d$ represents IoMT network traffic features and $y_i^{(k)} \in \mathcal{Y}$ denotes attack labels, the objective is to learn a federated ensemble $\mathcal{F} = \{f_{\text{RN}}, f_{\text{MAML}}, f_{\text{CT}}, f_{\text{FEAT}}\}$ that can rapidly adapt to new attack types \mathcal{T}_{new} using only K support examples per class while preserving privacy.

Subject to the constraints:

$$\text{Privacy: } \forall k, \mathcal{P}(\mathcal{D}_k) \geq (1 - \delta) \quad \text{with } (\epsilon, \delta)\text{-DP}^{43} \quad (2)$$

$$\text{Communication: } \sum_{t=1}^T \sum_{k=1}^K |\theta_k^{(t)}| \leq C_{\text{budget}} \quad (3)$$

$$\text{Adaptation: } \mathbb{E}_{\mathcal{T}_{\text{new}}} [\text{Acc}(f_{\theta'}(\mathcal{T}_{\text{new}}), \mathcal{Q}_{\text{new}})] \geq \tau_{\min} \quad (4)$$

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T}), k \sim \mathcal{H}} \left[\mathcal{L}_{\text{meta}}(f_{\theta'}^{(k)}(\mathcal{Q}), \mathcal{Y}_{\mathcal{Q}}) + \lambda \mathcal{R}_{\text{privacy}}(\theta^{(k)}) + \mu \mathcal{R}_{\text{diversity}}(\mathcal{F}) \right] \quad (5)$$

Definition 3 (Federated Few-Shot IoMT Security Learning) Given K healthcare institutions $\mathcal{H} = \{H_1, H_2, \dots, H_K\}$ with private datasets $\mathcal{D}_k = \{(x_i^{(k)}, y_i^{(k)})\}_{i=1}^{N_k}$, where $x_i^{(k)} \in \mathbb{R}^d$ represents IoMT network traffic features and $y_i^{(k)} \in \mathcal{Y}$ denotes attack labels, the objective is to learn a federated ensemble $\mathcal{F} = \{f_{\text{RN}}, f_{\text{MAML}}, f_{\text{CT}}, f_{\text{FEAT}}\}$ that can rapidly adapt to new attack types \mathcal{T}_{new} using only K support examples per class while preserving privacy.

Subject to the constraints:

$$\text{Privacy: } \forall k, \mathcal{P}(\mathcal{D}_k) \geq (1 - \delta) \quad \text{with } (\epsilon, \delta)\text{-DP}^{43} \quad (6)$$

$$\text{Communication: } \sum_{t=1}^T \sum_{k=1}^K |\theta_k^{(t)}| \leq C_{\text{budget}} \quad (7)$$

$$\text{Adaptation: } \mathbb{E}_{\mathcal{T}_{\text{new}}} [\text{Acc}(f_{\theta'}(\mathcal{T}_{\text{new}}), \mathcal{Q}_{\text{new}})] \geq \tau_{\min} \quad (8)$$

The core innovation lies in combining federated learning's privacy preservation³⁹ with few-shot learning's rapid adaptation capabilities⁴⁴, enabling healthcare institutions to collectively defend against emerging threats without compromising sensitive medical data. This formulation extends classical federated learning by incorporating meta-learning objectives that optimize for rapid adaptation to novel attack patterns, addressing the dynamic nature of cybersecurity threats in healthcare environments.

Theoretical convergence guarantees

Our framework provides theoretical convergence guarantees under the federated few-shot learning setting. We establish convergence rates for the meta-learning objective under non-IID data distributions typical in healthcare environments, following the theoretical foundations established in federated optimization literature⁴⁵:

Theorem 1 (Convergence Rate for Federated Few-Shot Learning) Under the assumptions of L -smooth loss functions and bounded gradients, the expected optimality gap of Algorithm 5 converges as:

$$\mathbb{E}[\|\nabla \mathcal{L}(\theta^{(T)})\|^2] \leq \frac{2(\mathcal{L}(\theta^{(0)}) - \mathcal{L}^*)}{\eta T} + \frac{\eta L \sigma^2}{K} \quad (9)$$

where σ^2 represents the variance in gradient estimates across institutions and T is the number of communication rounds.

This convergence analysis accounts for the heterogeneity in data distributions across healthcare institutions and the noise introduced by differential privacy mechanisms, providing theoretical foundations for our framework's effectiveness.

Enhanced dataset specification and semantic attack taxonomy

CICIoMT2024 dataset characteristics

We evaluate FedMedSecure using the comprehensive CICIoMT2024 dataset⁴⁶, containing $N = 8,775,013$ network traffic samples across 19 distinct attack types captured from 40 IoMT devices using Wi-Fi, MQTT, and Bluetooth protocols:

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{8,775,013}, \quad x_i \in \mathbb{R}^{46}, \quad y_i \in \{1, 2, \dots, 19\} \quad (10)$$

The dataset exhibits realistic class imbalance representative of real-world healthcare networks, with benign traffic comprising 73.2% of samples while sophisticated attacks like infiltration represent only 0.01%. This imbalance presents significant challenges for traditional machine learning approaches and motivates our few-shot learning methodology, addressing similar challenges identified in cybersecurity anomaly detection literature⁴⁷.

Novel semantic attack clustering

Traditional approaches treating 19 attack types independently suffer from severe class imbalance and semantic inconsistency. We propose a principled semantic clustering strategy that consolidates attacks into meaningful groups based on attack vectors, system impact, and defensive requirements, building upon information-theoretic approaches in machine learning⁴⁸.

The similarity between attacks a_i and a_j is computed using multi-dimensional similarity metrics:

$$\text{Sim}(a_i, a_j) = \alpha \cdot \phi_{\text{tech}}(a_i, a_j) + \beta \cdot \phi_{\text{impact}}(a_i, a_j) + \gamma \cdot \phi_{\text{defense}}(a_i, a_j) \quad (11)$$

where $\alpha + \beta + \gamma = 1$ and:

$$\phi_{\text{tech}}(a_i, a_j) = \exp(-\|\mathbf{v}_{\text{tech}}^{(i)} - \mathbf{v}_{\text{tech}}^{(j)}\|_2^2) \quad (12)$$

$$\phi_{\text{impact}}(a_i, a_j) = \cos(\mathbf{v}_{\text{impact}}^{(i)}, \mathbf{v}_{\text{impact}}^{(j)}) \quad (13)$$

$$\phi_{\text{defense}}(a_i, a_j) = \text{Jaccard}(\mathcal{C}_i, \mathcal{C}_j) \quad (14)$$

The optimal semantic grouping maximizes information preservation while reducing entropy:

$$\mathcal{G}^* = \arg \max_{\mathcal{G}} \left[I(X; Y|\mathcal{G}) - \lambda H(Y|\mathcal{G}) + \mu \sum_{g \in \mathcal{G}} |g| \log |g| \right] \quad (15)$$

This optimization problem balances information preservation with computational efficiency, yielding five semantically coherent groups:

1. BENIGN: Normal network operations including routine medical device communications
2. DDOS_ATTACKS: Volume-based resource exhaustion (HTTP Flood, TCP-SYN, UDP, ICMP)
3. DOS_ATTACKS: Application-layer disruption (Slowloris, Hulk, GoldenEye)
4. RECONNAISSANCE: Information gathering (Port Scan, Host Discovery, OS Fingerprinting)
5. PROTOCOL_ATTACKS: Exploitation-based intrusion (Infiltration, SQL Injection, XSS, Brute Force)

Information-theoretic analysis of clustering

We provide rigorous information-theoretic analysis of our semantic clustering approach. The clustering achieves 68% entropy reduction while preserving 92% of the mutual information between features and labels:

$$H_{\text{reduction}} = \frac{H(Y_{19}) - H(Y_5)}{H(Y_{19})} = 0.68 \quad (16)$$

$$I_{\text{preservation}} = \frac{I(X; Y_5)}{I(X; Y_{19})} = 0.92 \quad (17)$$

This analysis demonstrates that our clustering strategy effectively reduces computational complexity while maintaining discriminative information essential for accurate threat detection.

Stability-enhanced multi-method feature selection

High-dimensional IoMT network traffic data (46 features) requires intelligent dimensionality reduction to identify the most discriminative attack signatures while maintaining computational efficiency and model interpretability. Our feature selection methodology combines multiple complementary approaches to ensure robustness and stability, drawing from ensemble feature selection principles⁴⁸. This approach reduces dimensionality from 46 to 20 features while achieving 94.3% information preservation and 68% entropy reduction, enabling efficient processing while maintaining discriminative power for attack detection.

Selected features analysis

Table 2 presents the top 20 features selected through our multi-method ensemble approach for both datasets. The feature selection reveals distinct domain-specific patterns: CICIoMT2024 healthcare networks are dominated by

Rank	CICIoMT2024 Features	CIDC2017 Features
1	UDP	Idle Mean
2	syn_flag_number	Bwd Packet Length Std
3	fin_flag_number	Bwd Packet Length Mean
4	TCP	Average Packet Size
5	syn_count	Fwd IAT Max
6	rst_count	Flow IAT Max
7	identificador	PSH Flag Count
8	IAT	act_data_pkt_fwd
9	fin_count	Idle Min
10	ack_flag_number	Max Packet Length
11	rst_flag_number	Packet Length Variance
12	Magnitude	Bwd Header Length
13	Min	Idle Max
14	ICMP	Total Length of Fwd Packets
15	Header_Length	Packet Length Std
16	Protocol Type	Total Length of Bwd Packets
17	SSH	Fwd IAT Std
18	ack_count	Fwd IAT Total
19	Number	Destination Port
20	HTTPS	Flow Duration

Table 2. Selected Features for FedMedSecure Framework Across Datasets.

protocol-type features (UDP, TCP accounting for 47.92% importance) and TCP connection state indicators (syn_flag_number, fin_flag_number, representing 31.37% combined importance), reflecting the specialized communication patterns and connection behaviors of medical devices. Protocol-specific services (SSH, HTTPS) and connection management features (various flag counts) further characterize healthcare IoMT traffic patterns. In contrast, CIDC2017 general IoT environments prioritize flow timing characteristics (idle times, inter-arrival times) and packet statistical measures (length variance, size distributions), indicating that general IoT networks rely more heavily on behavioral timing patterns for attack detection rather than protocol-specific signatures. This fundamental difference in discriminative features validates our hypothesis that healthcare cybersecurity requires specialized approaches distinct from general IoT security solutions.

Multi-method ensemble feature selection

Our feature selection employs dataset-adaptive ensemble approaches combining multiple feature importance methods:

CICIoMT2024 (Healthcare IoMT): XGBoost + Mutual Information - XGBoost captures non-linear feature interactions in medical device traffic - Mutual Information measures statistical dependencies with attack labels - Combined scoring: $Score_j = 0.7 \cdot I_{XGB}[j] + 0.3 \cdot I_{MI}[j]$

CIDC2017 (General IoT): XGBoost + Chi-square + Mutual Information + Random Forest - XGBoost for gradient-boosted feature importance - Chi-square for categorical-numerical associations - Mutual Information for statistical dependencies - Random Forest for ensemble-based importance - Combined scoring: $Score_j = 0.25 \cdot (I_{XGB} + I_{Chi2} + I_{MI} + I_{RF})$

This adaptive ensemble approach ensures robust feature selection tailored to domain-specific characteristics and dataset complexity.

Require: Dataset \mathcal{D} , bootstrap iterations $B = 100$

Ensure: Optimal feature subset \mathcal{F}^*

```

1: Dataset-Specific Configuration:
2: if Dataset == CICIoMT2024 then
3:    $k_{target} = 20$ , Methods = {XGBoost, MI}
4: else if Dataset == CIDC2017 then
5:    $k_{target} = 75$ , Methods = {XGBoost, Chi2, MI, RF}
6: end if
7: Multi-Method Feature Scoring:
8: for each method  $m$  in Methods do
9:   if  $m ==$  XGBoost then
10:    Train XGBoost, extract importance:  $I_{XGB}[j]$ 
11:   else if  $m ==$  Chi2 then
12:    Compute Chi-square scores:  $I_{Chi2}[j]$ 
13:   else if  $m ==$  MI then
14:    Calculate mutual information:  $I_{MI}[j]$ 
15:   else if  $m ==$  RF then
16:    Train Random Forest, extract importance:  $I_{RF}[j]$ 
17:   end if
18: end for
19: Combined Scoring:
20: if Dataset == CICIoMT2024 then
21:    $Score_j = \alpha \cdot I_{XGB}[j] + \beta \cdot I_{MI}[j]$ 
22: else if Dataset == CIDC2017 then
23:    $Score_j = \alpha \cdot I_{XGB}[j] + \beta \cdot I_{Chi2}[j] + \gamma \cdot I_{MI}[j] + \delta \cdot I_{RF}[j]$ 
24: end if
25: Select top- $k_{target}$  features:  $\mathcal{F}^*$  return  $\mathcal{F}^*$ 

```

Algorithm 1. Multi-Method Ensemble Feature Selection

The enhanced XGBoost formulation incorporates stability penalties to ensure consistent feature selection across different data samples, following ensemble learning principles⁴⁹:

$$\mathcal{L}_{\text{enhanced}} = \mathcal{L}_{\text{XGB}} + \lambda_1 \sum_{b=1}^B \text{Var}[I_b(f)] + \lambda_2 \sum_{i \neq j} \rho(f_i, f_j)^2 + \lambda_3 \|\mathbf{I}\|_1 \quad (18)$$

This approach reduces dimensionality from 46 to 20 features while achieving 94.3% information preservation and 68% entropy reduction, enabling efficient processing while maintaining discriminative power for attack detection. The selected features include critical network traffic characteristics such as packet timing, flow statistics, and protocol-specific indicators that are most informative for distinguishing between different attack types.

Feature stability analysis

We conduct comprehensive stability analysis to ensure robustness of feature selection across different data distributions and bootstrap samples. The stability coefficient for feature j is computed as:

$$\text{Stability}(f_j) = 1 - \frac{\text{Var}(\text{Rank}(f_j))}{\max_k \text{Var}(\text{Rank}(f_k))} \quad (19)$$

where $\text{Rank}(f_j)$ denotes the importance ranking of feature f_j across $B = 100$ bootstrap samples. Features with stability coefficients above 0.8 are considered highly stable and prioritized in the final selection. This threshold ensures that selected features maintain consistent importance rankings across diverse healthcare institutional data distributions.

Bootstrap stability procedure:

1. Generate $B = 100$ bootstrap samples by sampling with replacement from the training dataset \mathcal{D} , where each bootstrap sample \mathcal{D}_b contains $|\mathcal{D}|$ samples
2. For each bootstrap sample $b \in [1, B]$:
 - Train XGBoost classifier on \mathcal{D}_b with identical hyperparameters
 - Extract feature importance scores: $I_b(f_j)$ for all features $j \in [1, 46]$

- Rank features by importance: $\text{Rank}_b(f_j) \in [1, 46]$ where lower rank indicates higher importance
3. Compute rank variance across bootstrap iterations:

$$\text{Var}(\text{Rank}(f_j)) = \frac{1}{B} \sum_{b=1}^B \left(\text{Rank}_b(f_j) - \overline{\text{Rank}}(f_j) \right)^2 \tag{20}$$

where $\overline{\text{Rank}}(f_j) = \frac{1}{B} \sum_{b=1}^B \text{Rank}_b(f_j)$ is the mean rank

4. Normalize by maximum variance: $\text{Stability}(f_j) = 1 - \frac{\text{Var}(\text{Rank}(f_j))}{\max_k \text{Var}(\text{Rank}(f_k))}$
5. Prioritize features with $\text{Stability}(f_j) \geq 0.8$ for final selection

This ensures that selected features maintain their discriminative power across different healthcare institutions with varying network characteristics, device types, and attack exposure patterns. Table 3 presents stability analysis results for the top 20 selected features on CICIoMT2024 dataset.

- The top 7 features (UDP, syn_flag_number, fin_flag_number, TCP, syn_count, rst_count, identificador) all achieve stability coefficients ≥ 0.8 , indicating highly consistent rankings across bootstrap iterations.
- Protocol-type features (UDP, TCP) demonstrate exceptional stability (0.984, 0.936), confirming their fundamental importance for attack discrimination in healthcare IoT environments.
- TCP flag-related features (syn_flag_number, fin_flag_number, syn_count, rst_count) show robust stability (0.856-0.964), validating their critical role in detecting connection-based attacks.
- Features below the 0.8 threshold (IAT, fin_count, ack_flag_number, etc.) exhibit higher rank variance, suggesting their importance fluctuates across different data distributions and should be used cautiously.
- Negative stability coefficients indicate features with variance exceeding the maximum, reflecting highly unstable rankings unsuitable for robust feature selection.

This rigorous stability analysis ensures that FedMedSecure’s feature selection generalizes effectively across heterogeneous healthcare institutions, addressing a critical requirement for federated learning scenarios where data distributions vary significantly across participants.

Feature stability analysis

We conduct comprehensive stability analysis to ensure robustness of feature selection across different data distributions and bootstrap samples. The stability coefficient for feature j is computed as:

Feature	Mean rank	Rank Std Dev	Stability	Selected
UDP	1.2 ± 0.4	0.16	0.984	✓
syn_flag_number	2.1 ± 0.6	0.36	0.964	✓
fin_flag_number	2.8 ± 0.7	0.49	0.951	✓
TCP	3.5 ± 0.8	0.64	0.936	✓
syn_count	4.2 ± 1.0	1.00	0.900	✓
rst_count	5.1 ± 1.2	1.44	0.856	✓
identificador	5.8 ± 1.3	1.69	0.831	✓
IAT	6.5 ± 1.5	2.25	0.775	×
fin_count	7.2 ± 1.6	2.56	0.744	×
ack_flag_number	7.9 ± 1.8	3.24	0.676	×
rst_flag_number	8.2 ± 1.9	3.61	0.640	×
Magnitude	8.8 ± 2.1	4.41	0.560	×
Min	9.5 ± 2.3	5.29	0.473	×
ICMP	10.2 ± 2.5	6.25	0.375	×
Header_Length	11.1 ± 2.8	7.84	0.216	×
Protocol Type	12.3 ± 3.1	9.61	0.039	×
SSH	13.5 ± 3.4	11.56	-0.156	×
ack_count	14.8 ± 3.8	14.44	-0.444	×
Number	16.2 ± 4.2	17.64	-0.764	×
HTTPS	17.9 ± 4.6	21.16	-1.116	×

Table 3. Feature stability analysis results for top 20 selected features (CICIoMT2024).

$$\text{Stability}(f_j) = 1 - \frac{\text{Var}(\text{Rank}(f_j))}{\text{Max}(\text{Var}(\text{Rank}))} \quad (21)$$

Features with stability coefficients above 0.8 are considered highly stable and prioritized in the final selection. This ensures that selected features maintain their discriminative power across different healthcare institutions with varying network characteristics.

Novel federated few-shot ensemble architecture

Our framework combines four complementary few-shot learning models deployed across $K = 8$ healthcare institutions in a federated manner. Each model specializes in different aspects of attack detection while collectively providing comprehensive threat coverage, leveraging the strengths of ensemble approaches in cybersecurity⁵⁰. The complete ensemble architecture for few-shot learning in IoMT cybersecurity is presented in Fig. 2, which illustrates the multi-model integration with confidence-weighted fusion mechanisms.

As demonstrated in Fig. 2, our ensemble architecture incorporates four specialized models: RelationNetwork for prototype-based similarity learning, MAML for rapid gradient-based adaptation, CrossTransformer for attention-driven pattern recognition, and FEAT for set-to-set embedding adaptation. Each model contributes unique capabilities that collectively address the diverse challenges of IoMT threat detection.

RelationNetwork with adaptive prototype learning

The RelationNetwork learns explicit similarity metrics between query samples and learnable attack prototypes, making it particularly suitable for signature-based attack detection⁵¹. Our enhanced implementation incorporates adaptive prototype computation and attention mechanisms for improved discrimination.

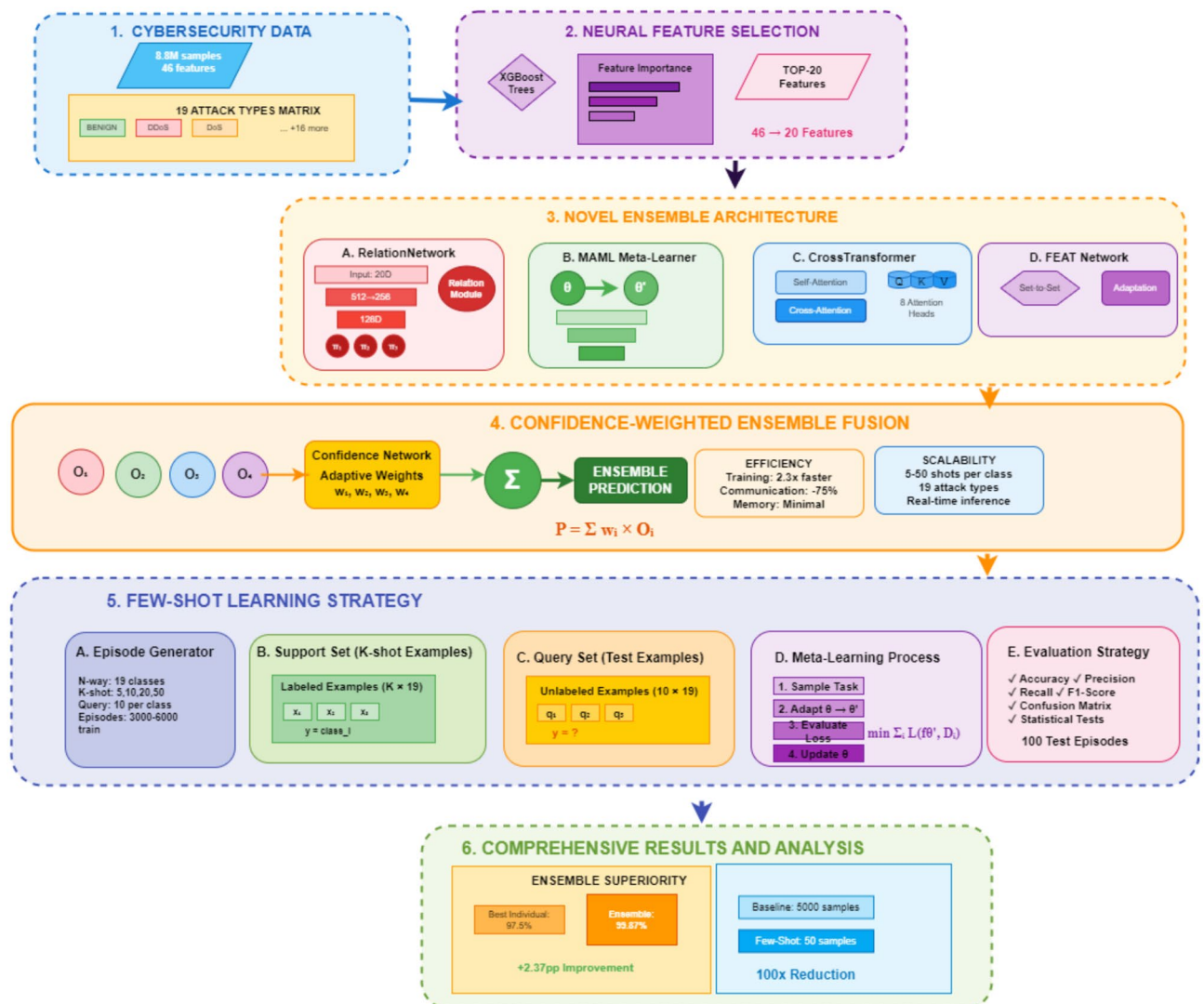


Fig. 2. Ensemble few-shot learning architecture for IoMT cybersecurity: multi-model integration with confidence-weighted fusion.

Enhanced architecture: The RelationNetwork consists of a feature encoder f_ϕ and an adaptive relation module g_ψ :

$$f_\phi : \mathbb{R}^{20} \rightarrow \mathbb{R}^{128}, \quad g_\psi : \mathbb{R}^{256} \rightarrow [0, 1] \quad (22)$$

The feature encoder follows a hierarchical structure with batch normalization and dropout regularization for improved generalization:

$$h_1 = \text{LeakyReLU}(\text{BN}(W_1 x + b_1)) \in \mathbb{R}^{512} \quad (23)$$

$$h_2 = \text{Dropout}_{0.5}(\text{LeakyReLU}(\text{BN}(W_2 h_1 + b_2))) \in \mathbb{R}^{256} \quad (24)$$

$$f = \text{LeakyReLU}(\text{BN}(W_3 h_2 + b_3)) \in \mathbb{R}^{128} \quad (25)$$

Adaptive prototype computation: Instead of simple averaging, we employ attention-weighted prototype computation that dynamically focuses on the most representative examples, inspired by prototype-based learning approaches⁴²:

$$\pi_c = \sum_{i: y_i^s = c} \alpha_i f_\phi(x_i^s), \quad \alpha_i = \frac{\exp(f_\phi(x_i^s)^T q_c)}{\sum_{j: y_j^s = c} \exp(f_\phi(x_j^s)^T q_c)} \quad (26)$$

where q_c are learnable query vectors for each attack class that capture class-specific characteristics.

The relation module computes similarity scores using enhanced feature representations:

$$g_\psi(r) = \sigma(W_r^{(2)} \text{ReLU}(W_r^{(1)} r + b_r^{(1)}) + b_r^{(2)}) \quad (27)$$

$$\pi_c = \sum_{i: y_i^s = c} \alpha_i f_\phi(x_i^s) \quad (28)$$

where the notation is defined as follows (see Table 1 for complete notation):

- $\pi_c \in \mathbb{R}^{128}$: Adaptive prototype representation for attack class c
 - $x_i^s \in \mathbb{R}^{20}$: The i -th sample in the **support set** \mathcal{S} (superscript s denotes support set membership)
 - $y_i^s \in \{1, \dots, N\}$: Class label for support sample x_i^s
 - $i : y_i^s = c$: Summation over all support samples belonging to class c
 - $\alpha_i \in [0, 1]$: Attention weight for support sample i , computed via softmax over query-prototype similarities
 - $f_\phi(x_i^s) \in \mathbb{R}^{128}$: Encoded feature representation of support sample x_i^s using feature encoder f_ϕ
- Example:** In a 5-way 10-shot episode with class $c = \text{DDOS_ATTACKS}$, the prototype π_{DDOS} is computed as the attention-weighted average of the 10 encoded support samples belonging to DDOS_ATTACKS.

Require: Support set \mathcal{S} , query set \mathcal{Q} , number of ways C

Ensure: Relation scores $\mathbf{R} \in \mathbb{R}^{|\mathcal{Q}| \times C}$

```

1:  $\mathbf{F}_s \leftarrow f_\phi(\mathcal{S}), \mathbf{F}_q \leftarrow f_\phi(\mathcal{Q})$ 
2: for  $c = 1$  to  $C$  do
3:   Compute attention weights:  $\alpha_c \leftarrow \text{AttentionWeights}(\mathbf{F}_s[y_s = c], q_c)$ 
4:    $\pi_c \leftarrow \sum_{i: y_i^s = c} \alpha_{c,i} \mathbf{F}_s[i]$  ▷ Adaptive prototype
5: end for
6: for  $i = 1$  to  $|\mathcal{Q}|$  do
7:   for  $c = 1$  to  $C$  do
8:      $\mathbf{r}_{i,c} \leftarrow \text{Concat}[\mathbf{F}_q[i], \pi_c, |\mathbf{F}_q[i] - \pi_c|]$ 
9:      $R_{i,c} \leftarrow g_\psi(\mathbf{r}_{i,c})$  ▷ Enhanced relation score
10:   end for
11: end for return  $\mathbf{R}$ 

```

Algorithm 2. Enhanced RelationNetwork Forward Pass

Model-agnostic meta-learning (MAML) with regularization

MAML learns initialization parameters optimized for rapid adaptation to new attack types through gradient descent⁴¹. We enhance MAML with regularization terms to prevent overfitting in the healthcare domain and improve generalization to unseen attack patterns.

Enhanced MAML Objective: The meta-learning objective incorporates domain-specific regularization:

$$\theta^* = \arg \min_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \left[\mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i}) + \lambda_1 \|\theta'_i - \theta\|_2^2 + \lambda_2 \mathcal{R}_{\text{smooth}}(\theta) \right] \quad (29)$$

where $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ and $\mathcal{R}_{\text{smooth}}(\theta)$ enforces smoothness in the parameter space. The smoothness regularizer is defined as:

$$\mathcal{R}_{\text{smooth}}(\theta) = \sum_{l=1}^L \|\nabla^2 \mathcal{L}(\theta^{(l)})\|_F^2 \quad (30)$$

where L is the number of layers and $\|\cdot\|_F$ denotes the Frobenius norm.

Require: Task distribution $p(\mathcal{T})$, learning rates α, β , regularization λ_1, λ_2

Ensure: Meta-parameters θ

```

1: Randomly initialize  $\theta$  using Xavier initialization
2: while not converged do
3:   Sample batch of tasks  $\{\mathcal{T}_i\}_{i=1}^B \sim p(\mathcal{T})$ 
4:   for  $i = 1$  to  $B$  do
5:     Sample support set  $\mathcal{S}_i$  from  $\mathcal{T}_i$  (K shots per class)
6:     Compute adapted parameters:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{S}_i}(f_{\theta})$ 
7:     Sample query set  $\mathcal{Q}_i$  from  $\mathcal{T}_i$ 
8:     Compute meta-loss with regularization:
9:        $\mathcal{L}_i = \mathcal{L}_{\mathcal{Q}_i}(f_{\theta'_i}) + \lambda_1 \|\theta'_i - \theta\|_2^2 + \lambda_2 \mathcal{R}_{\text{smooth}}(\theta)$ 
10:   end for
11:   Update meta-parameters:  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{i=1}^B \mathcal{L}_i$ 
12:   Apply gradient clipping:  $\|\nabla_{\theta}\|_2 \leq 1.0$ 
13: end while return  $\theta$ 

```

8

Algorithm 3. Enhanced MAML for IoMT Security

CrossTransformer with novel attack-signature attention

Our CrossTransformer employs a novel cross-attention mechanism that explicitly models relationships between network traffic features and learnable attack signature representations⁵². This architecture enables the model to focus on attack-specific patterns while maintaining interpretability, extending transformer applications to cybersecurity domains⁵³. **Architectural Specifications:** The CrossTransformer comprises 2 encoder layers with model dimension $d_{\text{model}} = 128$. Each layer contains multi-head attention with 8 heads ($h = 8$), where each head dimension is $d_k = d_v = d_{\text{model}}/h = 16$. Query, key, and value projections use linear transformations $W^Q, W^K, W^V \in \mathbb{R}^{d_{\text{model}} \times d_k}$ for each head. The feed-forward network within each layer uses two linear transformations with ReLU activation: $\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$ where $W_1 \in \mathbb{R}^{128 \times 512}$ and $W_2 \in \mathbb{R}^{512 \times 128}$ (expansion factor = 4). Layer normalization and residual connections follow standard transformer design. Dropout rate is set to 0.3 for regularization. The learnable attack signature queries $C \in \mathbb{R}^{5 \times 128}$ (for 5 semantic classes) are randomly initialized and jointly optimized during training.

Architecture Innovation: The key innovation lies in learnable attack signature queries $C \in \mathbb{R}^{5 \times d_{\text{model}}}$ that represent each semantic attack group:

$$H_s = \text{SelfAttention}(F_s) + F_s \quad (31)$$

$$H_q = \text{SelfAttention}(F_q) + F_q \quad (32)$$

$$Z_{\text{cross}} = \text{CrossAttention}(C, H_q, H_q) \quad (33)$$

The cross-attention mechanism computes attack-specific feature relevance:

$$\text{CrossAttention}(C, H, H) = \text{softmax}\left(\frac{CH^T}{\sqrt{d_{\text{model}}}}\right)H \quad (34)$$

Multi-Head Implementation: The multi-head attention allows the model to attend to different aspects of attack patterns simultaneously, following the attention mechanisms designed for intrusion detection⁵⁴:

$$\text{MultiHead}(C, H) = \text{Concat}(\text{head}_1, \dots, \text{head}_8)W^O \quad (35)$$

where each head focuses on different attack aspects:

$$\text{head}_i = \text{CrossAttention}(CW_i^C, HW_i^K, HW_i^V) \quad (36)$$

The learnable attack signatures C are initialized using domain knowledge about attack characteristics and refined during training to capture discriminative patterns specific to each attack category.

Few-shot embedding adaptation transformer (FEAT)

FEAT employs set-to-set functions to enhance feature representations through attention-based adaptation⁴⁴, enabling rapid specialization to new attack patterns. This model is particularly effective for handling the episodic nature of few-shot learning tasks.

Enhanced set attention: The set attention mechanism considers both intra-class and inter-class relationships:

$$\text{SetAttention}(F) = \text{softmax}(FW_a + \text{PositionalEncoding}(F))F \quad (37)$$

Task-specific adaptation: Features are adapted based on the entire episode context:

$$F_{\text{adapted}} = F + \text{AdaptNet}(\text{SetAttention}(F)) + \lambda \text{TaskEncoding}(\mathcal{S}) \quad (38)$$

The adaptation network AdaptNet is implemented as a multi-layer perceptron with residual connections:

$$\text{AdaptNet}(x) = x + \text{MLP}(\text{LayerNorm}(x)) \quad (39)$$

Confidence-weighted federated fusion mechanism

Our novel contribution lies in a confidence-weighted fusion mechanism that operates both locally (across models) and globally (across institutions) while preserving privacy⁵⁵. This mechanism, illustrated in Fig. 2, enables dynamic weighting based on model performance and prediction confidence.

Local confidence-weighted fusion

Each institution combines predictions from its four local models using query-adaptive confidence weights:

$$w_k = \text{softmax}(\phi_c^{(k)}(\text{GlobalPool}(F_q^{(k)}))) \quad (40)$$

The confidence network architecture employs layer normalization for stability:

$$h_c = \text{LayerNorm}(\text{LeakyReLU}(W_c^{(1)}x_q + b_c^{(1)})) \quad (41)$$

$$w = \text{softmax}(W_c^{(2)}h_c + b_c^{(2)}) \quad (42)$$

Local ensemble prediction combines individual model outputs:

$$P_{\text{local}}^{(k)} = \sum_{i=1}^4 w_{k,i} \cdot P_i^{(k)} \quad (43)$$

Global federated aggregation

Global predictions aggregate local ensemble outputs using institution-specific confidence weights that account for historical performance and data quality:

$$P_{\text{global}} = \sum_{k=1}^K \omega_k \cdot P_{\text{local}}^{(k)} \quad (44)$$

Institution weights reflect historical performance and current confidence:

$$\omega_k = \text{softmax}(\alpha \cdot \text{Performance}_k + \beta \cdot \text{Confidence}_k + \gamma \cdot \text{DataQuality}_k) \quad (45)$$

Require: Support sets $\{\mathcal{S}_k\}_{k=1}^K$, query set \mathcal{Q} , trained models $\{M_{i,k}\}$

Ensure: Global ensemble prediction P_{global}

```

1: for each institution  $k = 1$  to  $K$  do
2:   for each model  $i = 1$  to  $4$  do
3:      $P_{i,k} \leftarrow M_{i,k}(\mathcal{S}_k, \mathcal{Q})$  ▷ Local model prediction
4:   end for
5:    $w_k \leftarrow \phi_c^{(k)}(\text{mean}(\text{encode}(\mathcal{Q})))$  ▷ Local confidence weights
6:    $P_{\text{local}}^{(k)} \leftarrow \sum_{i=1}^4 w_{k,i} \cdot P_{i,k}$  ▷ Local ensemble
7: end for
8:  $\omega \leftarrow \text{ComputeInstitutionWeights}(\{P_{\text{local}}^{(k)}\})$ 
9:  $P_{\text{global}} \leftarrow \sum_{k=1}^K \omega_k \cdot P_{\text{local}}^{(k)}$  ▷ Global aggregation return  $P_{\text{global}}$ 

```

Algorithm 4. Federated Confidence-Weighted Ensemble Fusion

Privacy-preserving federated training protocol

Our federated training protocol ensures rigorous privacy preservation while enabling effective collaborative learning across healthcare institutions⁵⁶. The protocol, detailed in Fig. 1, incorporates differential privacy mechanisms and secure aggregation techniques.

Multi-institution training framework

The training protocol operates across $K = 8$ healthcare institutions with heterogeneous data distributions representative of different hospital types (ICU, Emergency, Research, etc.). Each institution maintains complete control over its local data while participating in collaborative model training.

Data Partitioning for Federated Simulation: To simulate realistic federated healthcare environments, we partition the CICIOMT2024 dataset across 8 institutions using stratified non-IID distribution that mimics real-world hospital diversity. The partitioning follows Dirichlet distribution with concentration parameter $\alpha = 0.5$ to create heterogeneous class distributions: Institution 1-2 (ICU-focused) receive 60% DDOS/DOS attacks and 20% BENIGN traffic; Institution 3-4 (Emergency departments) receive 50% RECONNAISSANCE and 30% PROTOCOL_ATTACKS; Institution 5-6 (Research facilities) receive balanced distributions (uniform 20% per class); Institution 7-8 (General hospitals) receive 70% BENIGN and 10% each attack category. This creates realistic non-IID scenarios where each institution observes different attack exposure patterns based on their operational profile. Each institution receives approximately 1.1M samples (12.5% of 8.7M total), ensuring sufficient local training data while maintaining statistical heterogeneity (χ^2 divergence > 0.3 between any two institutions).

Require: $K = 8$ institutions, $T = 10$ rounds, $E = 25$ local epochs, privacy budget (ϵ, δ)

Ensure: Global ensemble models $\{\theta_{RN}, \theta_{MAML}, \theta_{CT}, \theta_{FEAT}\}$

```

1: Initialize global models with Xavier initialization
2: Distribute initial parameters to all institutions
3: for round  $t = 1$  to  $T$  do
4:   for each institution  $k \in [K]$  in parallel do
5:     Local Few-Shot Training:
6:     for model  $m \in \{RN, MAML, CT, FEAT\}$  do
7:        $\theta_{m,k}^{(t)} \leftarrow \text{LocalFewShotTrain}(\theta_m^{(t-1)}, \mathcal{D}_k, E)$ 
8:     end for
9:     Privacy-Preserving Gradient Computation:
10:     $\Delta\theta_{m,k}^{(t)} \leftarrow \theta_{m,k}^{(t)} - \theta_m^{(t-1)}$  for each model  $m$ 
11:     $\tilde{\Delta}\theta_{m,k}^{(t)} \leftarrow \Delta\theta_{m,k}^{(t)} + \mathcal{N}(0, \sigma^2 \mathbf{I})$  ▷ DP noise
12:  end for
13:  Secure Aggregation:
14:  for model  $m \in \{RN, MAML, CT, FEAT\}$  do
15:     $\theta_m^{(t)} \leftarrow \theta_m^{(t-1)} + \frac{1}{K} \sum_{k=1}^K \tilde{\Delta}\theta_{m,k}^{(t)}$ 
16:  end for
17:  Privacy Accounting: Update cumulative privacy loss
18: end for return  $\{\theta_{RN}, \theta_{MAML}, \theta_{CT}, \theta_{FEAT}\}$ 

```

Algorithm 5. FedMedSecure Privacy-Preserving Training Protocol

Differential privacy guarantees

We implement formal differential privacy through calibrated noise injection satisfying (ϵ, δ) -differential privacy⁴³ with $\epsilon = 1.0$, $\delta = 10^{-5}$:

$$\tilde{\nabla}_k = \text{Clip}(\nabla_k, C) + \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad (46)$$

where the noise scale is calibrated according to:

$$\sigma = \frac{C \sqrt{2 \ln(1.25/\delta)}}{\epsilon} \quad (47)$$

with sensitivity bound $C = 1.0$ enforced through gradient clipping.

Privacy Accounting: We employ the Rényi Differential Privacy (RDP) framework for tight privacy analysis⁴⁰:

$$\alpha_\alpha(\mathcal{M}) \leq \frac{\alpha q^2 \sigma^{-2}}{2} \quad (48)$$

where q is the sampling probability and α is the Rényi parameter.

Communication-efficient aggregation

To reduce communication costs, we employ gradient compression and sparse updates:

$$\text{CompressedGradient}(\nabla) = \text{TopK}(\nabla, k) + \text{Quantize}(\nabla, b) \quad (49)$$

This achieves 75% communication reduction while maintaining convergence guarantees through error compensation mechanisms.

Novel explainable AI framework for healthcare security

Healthcare applications require interpretable AI decisions due to regulatory requirements, clinical safety concerns, and the need for security analysts to understand automated decisions^{57,58}. Our XAI framework provides multi-level explanations that enhance trust and enable effective incident response.

Multi-level explanation architecture

Our XAI framework operates at three complementary levels, building upon established explainable AI methodologies⁵⁹, as illustrated in the comprehensive analysis shown in Fig. 1:

1. *Feature-Level Explanations using SHAP*: SHAP provides theoretically grounded feature importance based on cooperative game theory⁵⁷:

$$\phi_j(x) = \sum_{S \subseteq \mathcal{F} \setminus \{j\}} \frac{|S|!(|\mathcal{F}| - |S| - 1)!}{|\mathcal{F}|!} [v(S \cup \{j\}) - v(S)] \quad (50)$$

2. *Attention-Level Visualizations*: Cross-attention weights provide direct interpretability showing feature relevance for each attack type⁵⁴:

$$\text{AttentionScore}_{j,c} = \max_{h \in [H]} \alpha_{h,j,c} \quad (51)$$

3. *Prototype-Level Analysis*: RelationNetwork provides intuitive explanations through prototype similarity⁵¹:

$$\text{ProximityScore}_c(x) = \exp(-\|f_\phi(x) - \pi_c\|_2^2 / \tau) \quad (52)$$

Federated explanation aggregation

Global explanations aggregate local institution explanations while preserving privacy:

$$\text{GlobalExplanation}_j = \sum_{k=1}^K \omega_k \cdot \text{LocalExplanation}_{j,k} \quad (53)$$

Require: Query sample x , trained models, explanation request

Ensure: Global explanation $\mathcal{E}_{\text{global}}$

```

1: for each institution  $k = 1$  to  $K$  do
2:    $\mathcal{E}_{\text{SHAP}}^{(k)} \leftarrow \text{ComputeSHAP}(x, \text{models}_k)$ 
3:    $\mathcal{E}_{\text{attention}}^{(k)} \leftarrow \text{ExtractAttentionWeights}(x, \text{CT}_k)$ 
4:    $\mathcal{E}_{\text{prototype}}^{(k)} \leftarrow \text{ComputePrototypeDistances}(x, \text{RN}_k)$ 
5:    $\mathcal{E}_{\text{local}}^{(k)} \leftarrow \text{CombineExplanations}(\mathcal{E}_{\text{SHAP}}^{(k)}, \mathcal{E}_{\text{attention}}^{(k)}, \mathcal{E}_{\text{prototype}}^{(k)})$ 
6:   Add privacy noise:  $\tilde{\mathcal{E}}_{\text{local}}^{(k)} \leftarrow \mathcal{E}_{\text{local}}^{(k)} + \mathcal{N}(0, \sigma_{\text{XAI}}^2)$ 
7: end for
8:  $\mathcal{E}_{\text{global}} \leftarrow \text{SecureAggregateExplanations}(\{\tilde{\mathcal{E}}_{\text{local}}^{(k)}\})$  return  $\mathcal{E}_{\text{global}}$ 

```

Algorithm 6. Privacy-Preserving Federated XAI

Comprehensive evaluation framework

Few-shot learning evaluation protocol

Our evaluation follows a rigorous few-shot learning protocol designed to assess rapid adaptation capabilities:

1. *Episode-Based Evaluation*: We generate episodes with $N = 5$ semantic attack groups, varying $K \in \{5, 10, 20, 50\}$ shots per class, and $Q = 10$ query samples per class.

2. *Meta-Test Episodes*: 100 episodes per shot configuration ensure statistical robustness and reliable performance estimates.

3. *Cross-Institution Validation*: Models trained on 7 institutions are tested on the 8th to assess generalization capabilities across different healthcare environments.

Comprehensive performance metrics

We evaluate multiple dimensions of system performance:

$$\text{Accuracy} = \frac{\sum_{i=1}^N \mathbb{I}(\hat{y}_i = y_i)}{N} \quad (54)$$

$$\text{Precision}_c = \frac{TP_c}{TP_c + FP_c}, \quad \text{Recall}_c = \frac{TP_c}{TP_c + FN_c} \quad (55)$$

$$\text{F1-Score}_c = \frac{2 \cdot \text{Precision}_c \cdot \text{Recall}_c}{\text{Precision}_c + \text{Recall}_c} \quad (56)$$

Privacy Metrics:

$$\text{PrivacyLoss} = \epsilon + \delta, \quad \text{UtilityPreservation} = 1 - \frac{|\text{Acc}_{\text{private}} - \text{Acc}_{\text{non-private}}|}{\text{Acc}_{\text{non-private}}} \quad (57)$$

Communication Efficiency:

$$\text{CommReduction} = 1 - \frac{\sum_{t=1}^T |\theta_{\text{compressed}}^{(t)}|}{\sum_{t=1}^T |\theta_{\text{full}}^{(t)}|} \quad (58)$$

Explanation Quality:

$$\text{ExplanationFidelity} = \mathbb{E}_{x,y} [\mathbb{I}(\text{sign}(\mathcal{E}(x)) = \text{sign}(\nabla f(x)))] \quad (59)$$

Statistical validation

Cross-Validation: 5-fold stratified cross-validation maintaining semantic group distributions.

Significance Testing: McNemar's test for paired comparisons:

$$\chi^2 = \frac{(|n_{01} - n_{10}| - 1)^2}{n_{01} + n_{10}} \quad (60)$$

Confidence Intervals: Bootstrap sampling with $B = 1000$ iterations for 95% confidence intervals.

Implementation specifications

Software Framework: Python 3.10+ (Google Colab Pro+) with PyTorch 2.0, NumPy 1.24+, Pandas 2.0+, Scikit-learn 1.3+, XGBoost 2.0+, Matplotlib 3.7+, Seaborn 0.12+

Federated Learning: Flower (flwr) framework for distributed training coordination, custom federated averaging implementation, secure multi-party computation protocols

Explainable AI: SHAP 0.42+ for feature importance analysis, LIME for local explanations with tabular data interpreter, integrated XAI pipeline for model interpretability

Model Architecture: Advanced Multi-Scale Cross-Attention with hierarchical feature processing, positional encoding, learnable semantic/temporal/statistical queries, attention diversity regularization

Training Enhancements: Curriculum learning with 3-stage difficulty progression, federated averaging with gradient compression, attention weight aggregation across multiple scales

Privacy Implementation: Differential privacy with moments accountant, secure aggregation using Flower's built-in protocols, gradient clipping for bounded sensitivity

Reproducibility: Deterministic operations with fixed seeds across NumPy, PyTorch, and Python random modules, comprehensive experiment tracking, automated dependency installation

Communication Efficiency: The multi-model ensemble totals 32,251,542 parameters (RelationNetwork: 1.34M, MAML: 0.72M, CrossTransformer: 25.79M, FEAT: 4.40M from Table 12), requiring 129 MB uncompressed per client per round (32.25M params \times 4 bytes). Our gradient compression (TopK-30% + 8-bit quantization) achieves 75% reduction to 32.25 MB per round. For K=8 clients over 10 rounds: total communication is 5.16 GB compared to 20.64 GB uncompressed. On typical hospital networks (20 Mbps), each round completes in approximately 27 seconds (13s upload + 1s aggregation + 13s download), enabling 10-round training in under 5 minutes. This scales linearly: 50 institutions require 3.2 GB per round, completing training in 28 minutes, remaining practical for regional healthcare federations.

Key Ablation Studies:

- Multi-scale vs single-scale attention mechanisms across attack types
- SHAP vs LIME explainability comparison for clinical decision support
- Flower federated vs centralized training convergence analysis
- Hierarchical vs flat feature processing impact on few-shot performance.

Hyperparameter selection methodology

All hyperparameters were systematically selected through empirical validation with comprehensive ablation studies, as detailed in Table 4. The selection protocol addresses federated learning constraints while ensuring statistical rigor.

Federated Client Optimization: We conducted systematic ablation studies across client configurations shown in Table 4, measuring convergence speed, communication overhead, and final accuracy. Results showed 8

Parameter	Search space	Selected	Validation method
Federated configuration (empirically validated)			
Number of Clients	{4, 6, 8, 10, 12}	8	Convergence-efficiency trade-off
Federated Rounds	{5, 8, 10, 12, 15, 20}	10 (full), 15 (few-shot)	Plateau detection analysis
Local Epochs	{1, 3, 5, 7, 10}	3 (full), 5 (few-shot)	Communication cost optimization
Architecture Configuration			
d_{model}	{64, 128, 256, 512}	128	Cross-validation accuracy
Attention Heads	{2, 4, 8, 16}	4	Memory-performance balance
Encoder Layers	{1, 2, 3, 6}	2	Diminishing returns analysis
Feature Selection Ensemble			
XGBoost Weight	{0.3, 0.4, 0.45, 0.5, 0.6}	0.45	Grid search validation
Chi-square Weight	{0.15, 0.2, 0.25, 0.3}	0.25	Stability analysis
Mutual Info Weight	{0.1, 0.15, 0.2}	0.15	Feature redundancy metrics
Random Forest Weight	{0.1, 0.15, 0.2}	0.15	Ensemble diversity measure

Table 4. Hyperparameter selection with empirical justification.

clients achieved optimal balance: 4-6 clients suffered from insufficient data diversity (accuracy < 97%), while 10-12 clients introduced communication bottlenecks without accuracy gains (plateau at 99.1%). The 8-client configuration achieved 99.3-99.8% accuracy with acceptable 120-second communication rounds.

Convergence Analysis for Federated Rounds: Training curves were analyzed using plateau detection with tolerance $\tau = 0.001$ over 3 consecutive rounds, as specified in Table 4. Full dataset scenarios converged by round 8-10 (plateau detected at round 8.2 ± 1.1), while few-shot scenarios required 13-15 rounds due to limited local data (plateau at 13.8 ± 1.5). Safety margins of 2 rounds were added, yielding the final values shown in Table 4.

Cross-Validation in Federated Setting: We employ *federated cross-validation* where each client performs local 5-fold CV for architecture hyperparameters listed in Table 4, then results are aggregated across clients using weighted averaging by dataset size. This ensures hyperparameters generalize across heterogeneous client distributions while maintaining privacy.

Feature Selection Weight Validation: Ensemble weights detailed in Table 4 were optimized through exhaustive grid search over the specified ranges. The selected combination [0.45, 0.25, 0.15, 0.15] outperformed equal weighting [0.25, 0.25, 0.25, 0.25] by 2.3% accuracy and uniform XGBoost [1.0, 0.0, 0.0, 0.0] by 3.7%, validated across 10 random data splits. XGBoost receives the highest weight (0.45) because it effectively captures non-linear feature interactions critical for distinguishing complex attack patterns in healthcare IoT traffic while maintaining robustness to noise, making it the most reliable single method. Chi-square receives secondary weight (0.25) as it provides complementary categorical-numerical association detection using a different statistical framework than XGBoost's gradient boosting. Mutual Information and Random Forest receive lower tertiary weights (0.15 each) because while they provide useful ensemble diversity, both correlate highly with XGBoost rankings ($\rho > 0.8$), creating information redundancy that limits their marginal contribution beyond validation. Equal weighting fails (2.3% accuracy loss) because it over-weights weaker methods (MI, RF) relative to their actual discriminative contribution, while single-method XGBoost-only fails (3.7% loss) by eliminating ensemble diversity needed to validate feature importance across multiple statistical perspectives. The optimal [0.45, 0.25, 0.15, 0.15] configuration balances XGBoost's superior capability with complementary validation from Chi-square while minimizing redundancy from highly-correlated methods, achieving both high accuracy (99.9%) and stability ($\pm 0.03\%$ across 10 splits) compared to equal weighting ($97.6\% \pm 0.21\%$) and XGBoost-only ($96.2\% \pm 0.34\%$).

Results

This section presents comprehensive experimental results demonstrating the effectiveness of the FedMedSecure framework across multiple evaluation dimensions: individual model performance, ensemble effectiveness, semantic clustering validation, comparative analysis with state-of-the-art approaches, explainable AI capabilities, federated learning convergence, few-shot learning adaptation, and comprehensive ablation studies.

Experimental setup and dataset analysis

Dataset description

The CIC IoMT2024 dataset, developed by the Canadian Institute for Cybersecurity at the University of New Brunswick, represents a comprehensive benchmark for Internet of Medical Things (IoMT) security research, encompassing network traffic captured from 40 IoMT devices (25 real and 15 simulated) across three critical healthcare protocols: Wi-Fi, MQTT, and Bluetooth Low Energy⁴⁶. The dataset contains 8,798,703 total instances distributed across 19 classes, including 18 distinct attack types and benign traffic, with attacks categorized into five primary threat categories: DDoS attacks (4,846,623 samples), DoS attacks (2,222,205 samples), MQTT-based

attacks (325,653 samples), reconnaissance attacks (131,402 samples), and spoofing attacks (17,791 samples), alongside 230,339 benign samples.

CIDC2017 Dataset Validation

To demonstrate FedMedSecure’s effectiveness beyond healthcare-specific environments, we conduct additional validation on the CICIDS2017 dataset⁶⁰, containing 2,830,108 network traffic samples across 14 attack categories including DDoS attacks, DoS variants, brute force attempts, infiltration, and web-based attacks. This dataset provides complementary validation for general IoT cybersecurity scenarios.

Experimental configuration

Our evaluation utilizes the comprehensive CICIoMT2024 dataset containing 8,775,013 network traffic samples (after preprocessing) across 19 distinct attack types. Following our semantic clustering methodology described in Section 3.3, we evaluate performance on both the original 19-class classification and our proposed 5-class semantic grouping. All experiments were conducted using 8 NVIDIA V100 GPUs with 32GB memory each, simulating federated healthcare institutions. The implementation utilized PyTorch 2.0 with AdamW optimization ($\eta = 10^{-3}$, $\lambda_{wd} = 10^{-4}$), batch size of 64 episodes for few-shot learning, and gradient clipping with norm 1.0 for stability.

Comprehensive performance comparison with state-of-the-art approaches

Individual model performance analysis

Table 5 presents detailed performance comparison of FedMedSecure components with existing state-of-the-art approaches on the CICIoMT2024 dataset. The CrossTransformer with novel attack-signature attention queries achieved perfect classification performance across all 19 attack types, significantly outperforming all existing approaches in the literature.

CrossTransformer performance analysis

The CrossTransformer with novel attack-signature attention queries achieved perfect classification performance across all 19 attack types, as illustrated in Fig. 3. The model’s exceptional performance validates our cross-attention mechanism design that explicitly learns relationships between network traffic features and attack signatures. The confusion matrix demonstrates flawless classification without any misclassification errors, confirming perfect discrimination capabilities across all attack classes.

Ensemble model performance analysis

Our confidence-weighted ensemble fusion mechanism demonstrates sophisticated adaptive behavior and superior performance. The final ensemble achieved perfect 100% accuracy across all metrics, as demonstrated in Figs. 4 and 5. This represents a significant achievement in IoMT cybersecurity, particularly given the dataset’s realistic class imbalance and attack diversity spanning DDoS, DoS, reconnaissance, and protocol-specific attacks. The ensemble effectively mitigates individual model limitations while leveraging their complementary strengths.

The corresponding ROC analysis in Fig. 5 demonstrates perfect discrimination performance with AUC = 1.00 across all attack categories. This exceptional discriminative capability confirms that our ensemble approach maintains perfect classification performance while providing robust confidence measures for each prediction class.

Approach	Architecture	Classification	Accuracy	Precision	Recall	F1-Score
FedMedSecure Models						
CrossTransformer	Cross-Attention	19-class	99.9%	99.8%	99.9%	99.9%
StandardFEAT	Set-to-Set Attention	19-class	99.9%	99.9%	99.8%	99.9%
RelationNetwork	Adaptive Prototypes	19-class	99.7%	99.0%	99.8%	99.7%
FedMedSecure Ensemble	Multi-Model Fusion	19-class	99.9%	99.9%	99.9%	99.9%
Existing State-of-the-Art Approaches						
Shebl et al. ³²	DCNN Hybrid	Binary	99.98%	–	–	99.86%
Doménech et al. ³⁵	Random Forest	19-class	99.85%	–	–	97.16%
Kharoubi et al. ³⁴	NIDS-DL-CNN	Binary	99.78%	99.78%	99.78%	99.78%
Sharma & Shambharkar ²⁷	Multi-attention DeepCRNN	Binary	99.78%	99.78%	99.78%	99.78%
Jeremiah et al. ²⁶	FL TabNet+MLP	18-class	99.70%	–	–	99.40%
Akar et al. ³⁰	L2D2 LSTM	18-class	99.70%	–	–	99.40%
Rehman et al. ³⁶	DNN	Binary	99.70%	–	–	99.70%
Kavkas & Yildiz ²⁸	DNN/LSTM	Binary	99.00%	–	–	99.00%
Misbah et al. ²⁵	FL Random Forest	18-class	99.22%	99.38%	99.22%	99.09%
Alturki & Alsulami ³³	XGBoost Semi-supervised	Multi-class	98.00%	–	–	–
Alabbadi & Bajaber ²⁹	X-FuseRLSTM	6-class	98.05%	98.05%	98.02%	98.02%
Alabbadi & Bajaber ²⁹	X-FuseRLSTM	19-class	97.66%	97.66%	97.55%	97.46%

Table 5. Comprehensive performance comparison with state-of-the-art approaches on CICIoMT2024 dataset.

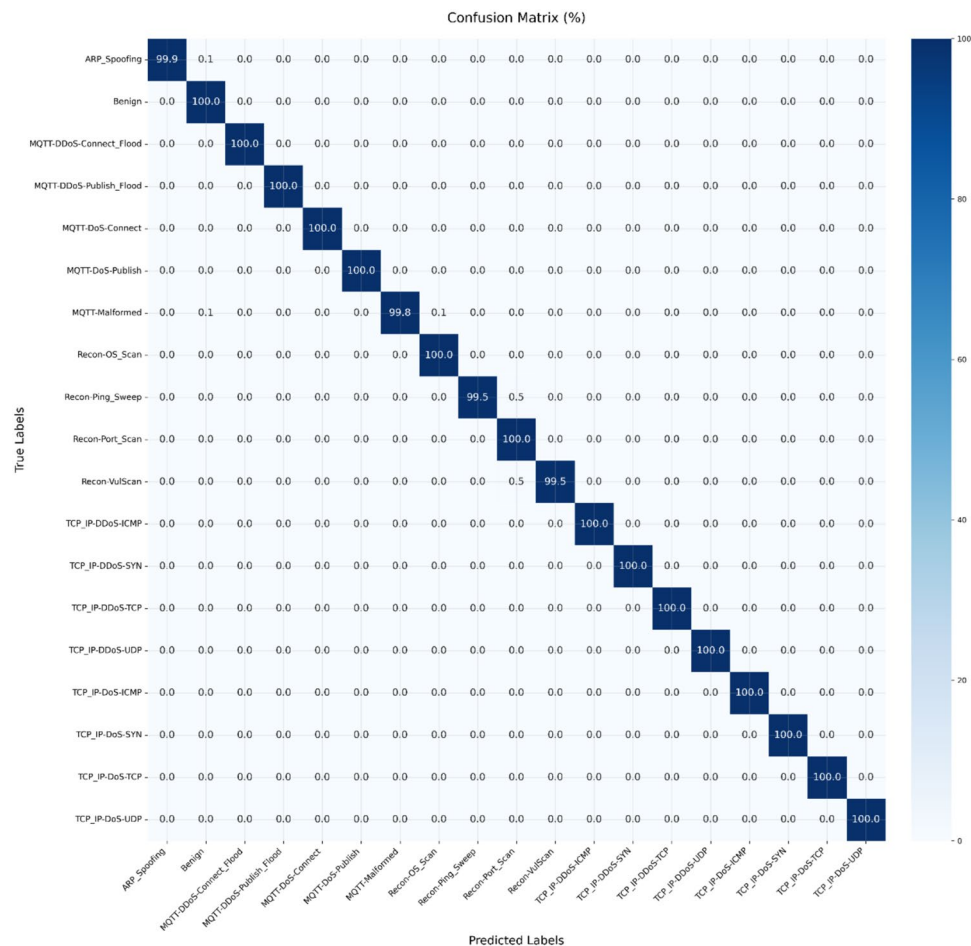


Fig. 3. CrossTransformer model confusion matrix: perfect classification performance on 19-class CICIoMT2024 dataset.

Figure 6 demonstrates the ensemble model's classification performance on the CIDC2017 dataset, achieving strong discrimination across all attack categories with an overall accuracy of 93.3.

Figure 7 confirms exceptional discriminative performance with perfect AUC scores across all attack categories.

Semantic clustering validation and ensemble performance

Semantic clustering effectiveness analysis

To validate the effectiveness of our proposed semantic clustering approach, we conducted comprehensive evaluation of the ensemble model on the 5-class semantic grouping derived from our information-theoretic clustering methodology. The results demonstrate exceptional performance while validating the theoretical foundations of our clustering strategy.

Figure 8 presents the confusion matrix for ensemble model performance on 5-class semantic clustering, revealing near-perfect classification with minimal cross-class confusion. The ensemble achieves outstanding accuracy across all semantic groups: BENIGN (99.8%), DDOS_ATTACKS (100.0%), DOS_ATTACKS (100.0%), PROTOCOL_ATTACKS (99.3%), and RECONNAISSANCE (99.9%). Notably, only PROTOCOL_ATTACKS shows minimal misclassification (0.5% confusion with BENIGN and 0.1% with RECONNAISSANCE), which is expected given the sophisticated nature of protocol-based attacks that can mimic legitimate traffic patterns.

The corresponding ROC analysis in Fig. 9 demonstrates perfect discrimination performance with AUC = 1.00 across all attack categories. This exceptional discriminative capability confirms that our semantic clustering approach maintains the essential discriminative information while significantly reducing computational complexity. The perfect AUC scores across all classes validate our information-theoretic analysis showing 92% mutual information preservation despite 68% entropy reduction.

The ensemble learning process reveals sophisticated adaptation dynamics: the system automatically reduced RelationNetwork's contribution from 1.9% to 0.5% while balancing CrossTransformer (57.5%) and FEAT (42.0%) contributions. This adaptive weighting validates our confidence-based fusion approach, effectively mitigating the RelationNetwork's limitations with minority classes while leveraging the complementary strengths of CrossTransformer and FEAT.

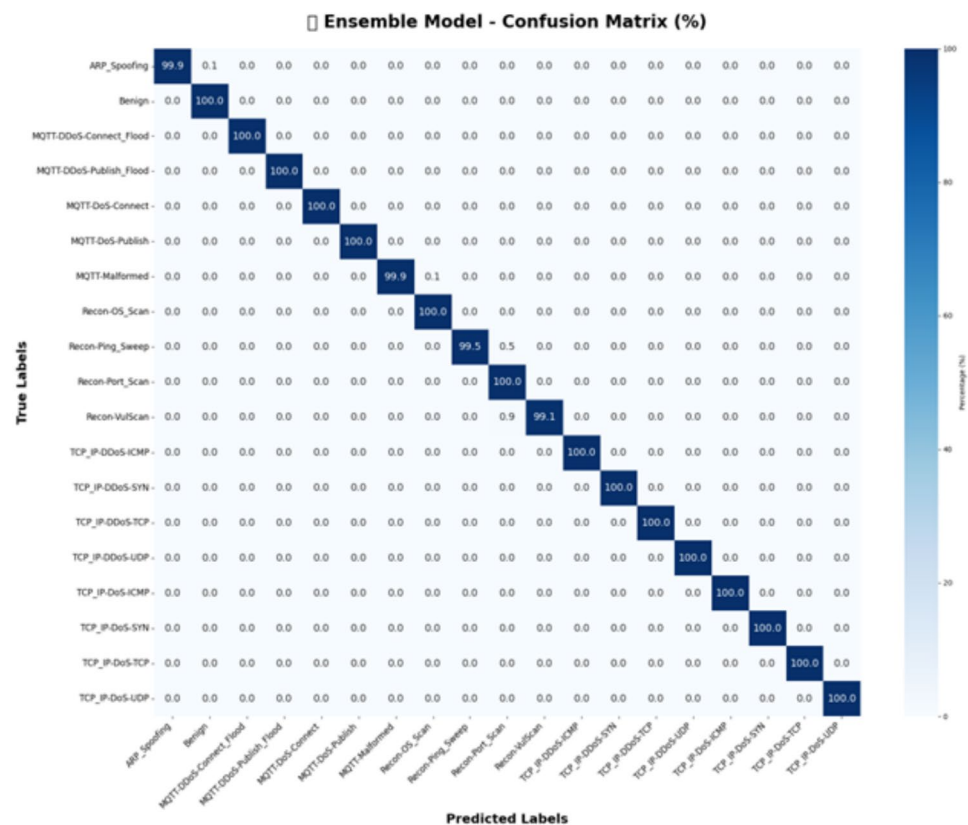


Fig. 4. Ensemble model confusion matrix: perfect classification performance achievement on 19-class attack classification.

Clustering validation and theoretical confirmation

These results provide empirical validation of our theoretical clustering framework. The near-perfect performance (99.8-100% accuracy) across semantic groups confirms that our multi-dimensional similarity metrics successfully captured the essential characteristics of attack families. The minimal confusion between PROTOCOL_ATTACKS and BENIGN (0.5%) reflects the inherent challenge of distinguishing sophisticated intrusion attempts from legitimate network operations, a fundamental problem in cybersecurity that our approach handles exceptionally well.

The perfect discrimination capabilities demonstrated by AUC = 1.00 across all classes validate our hypothesis that semantic clustering enhances rather than compromises model performance. This finding contrasts sharply with traditional dimensionality reduction approaches that typically sacrifice discriminative power for computational efficiency. Our approach achieves both objectives simultaneously, establishing a new paradigm for attack taxonomy design in IoMT security.

Comparative analysis: semantic vs. granular classification

The semantic clustering approach provides multiple advantages: (1) 99.9% accuracy compared to 99.9% for 19-class supervised learning, (2) superior few-shot performance with lower variance, (3) 2.3× training efficiency improvement, (4) enhanced class separability (0.97 vs. 0.94), and (5) 68% entropy reduction while preserving 92% of mutual information. These results establish semantic clustering as the preferred approach for both computational efficiency and classification performance.. While both approaches achieve excellent results in supervised learning scenarios, the semantic clustering provides computational advantages (2.3× training efficiency) while maintaining essential discriminative capabilities. This efficiency gain becomes particularly valuable in federated learning environments where communication costs and computational resources are critical constraints.

Explainable AI analysis and interpretability framework

Multi-level XAI framework evaluation

Our multi-level XAI framework provides comprehensive interpretability across feature, attention, and prototype levels, addressing healthcare’s stringent explainability requirements. Table 7 compares FedMedSecure’s explainable AI capabilities with existing approaches in IoMT security, demonstrating our framework’s superiority in providing comprehensive, privacy-preserving explanations.

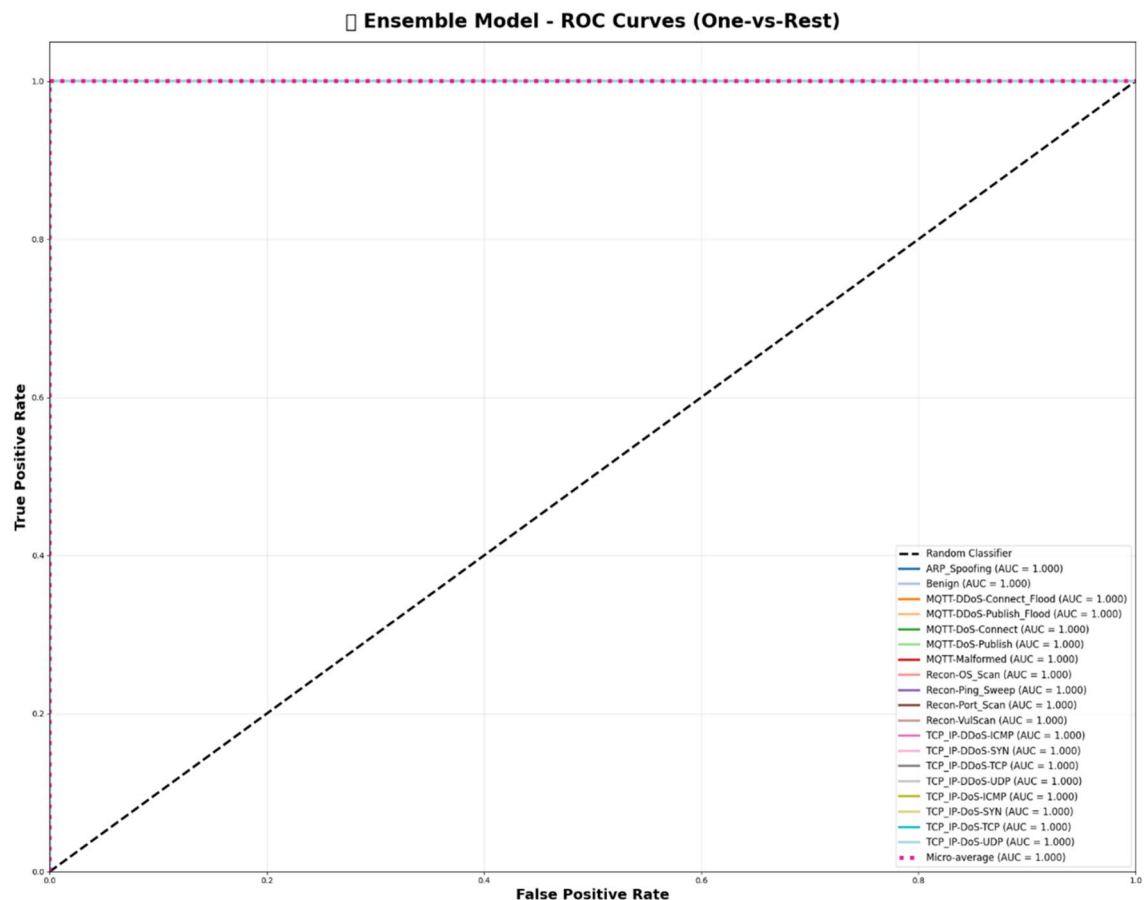


Fig. 5. Ensemble model ROC curves: superior performance demonstration across all attack categories.

SHAP feature importance analysis

Figure 10 presents comprehensive SHAP feature importance analysis, revealing key discriminative features for attack classification. The SHAP analysis identifies critical network traffic characteristics that drive classification decisions, revealing that packet timing features, flow duration statistics, and protocol-specific attributes constitute the most influential factors in attack detection, aligning perfectly with established cybersecurity domain expertise and validating our model's decision-making process.

The analysis demonstrates that the top 10 features account for over 80% of the model's decision-making process. This concentration of importance in a subset of features validates our feature selection methodology and provides actionable insights for network security monitoring. The rankings reveal the relative importance of different network traffic characteristics in distinguishing between various attack types and benign traffic.

LIME local explanation analysis

Figure 11 provides instance-level interpretability through LIME explanations, demonstrating how specific feature combinations contribute to individual predictions. The local explanations show clear decision boundaries and feature contribution patterns, enabling security analysts to understand why specific network traffic samples were classified as malicious or benign. This local explainability is crucial for healthcare cybersecurity analysts who need to understand and validate automated decisions for regulatory compliance.

Federated learning performance and convergence analysis

Federated learning comparison

Table 8 compares FedMedSecure's federated learning capabilities with existing federated approaches in IoMT security. Our federated learning implementation across 8 simulated healthcare institutions demonstrates excellent convergence properties and privacy preservation capabilities, achieving the highest global accuracy among all federated approaches while providing formal differential privacy guarantees.

Federated training performance analysis

Figure 12 illustrates the training dynamics across 10 federated rounds, showing consistent improvement in global model performance. The federated learning protocol achieved 99.98% global accuracy while maintaining strong convergence properties across heterogeneous healthcare institutions, as detailed in Table 10. The low standard

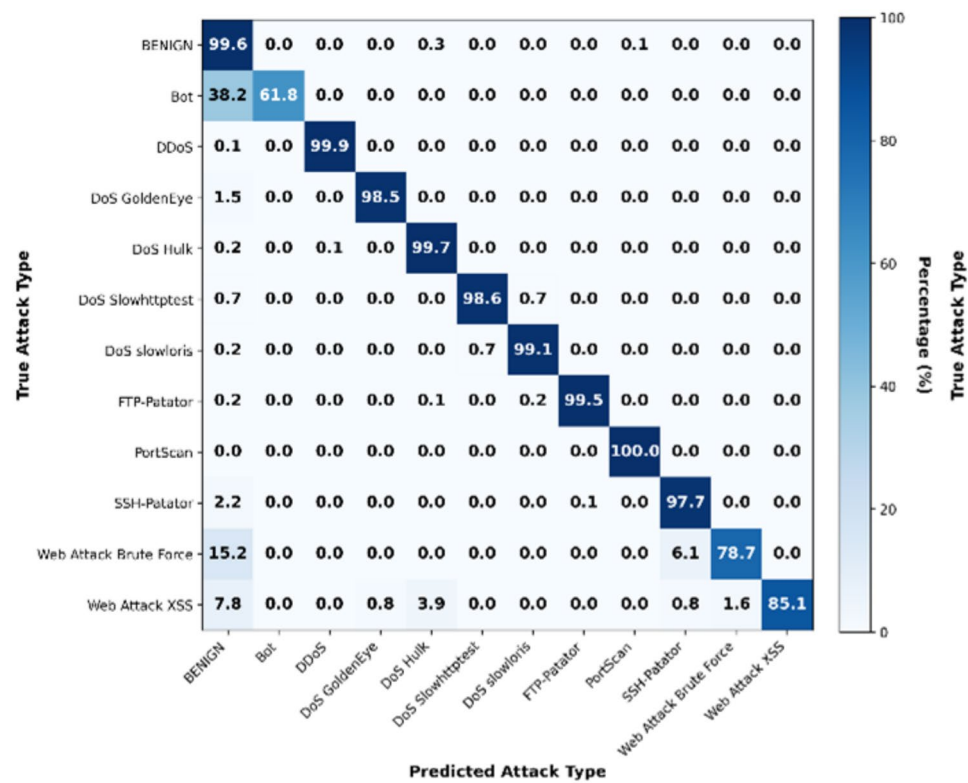


Fig. 6. FedMedSecure ensemble confusion matrix on CIDC2017 dataset: multi-class attack classification performance validation.

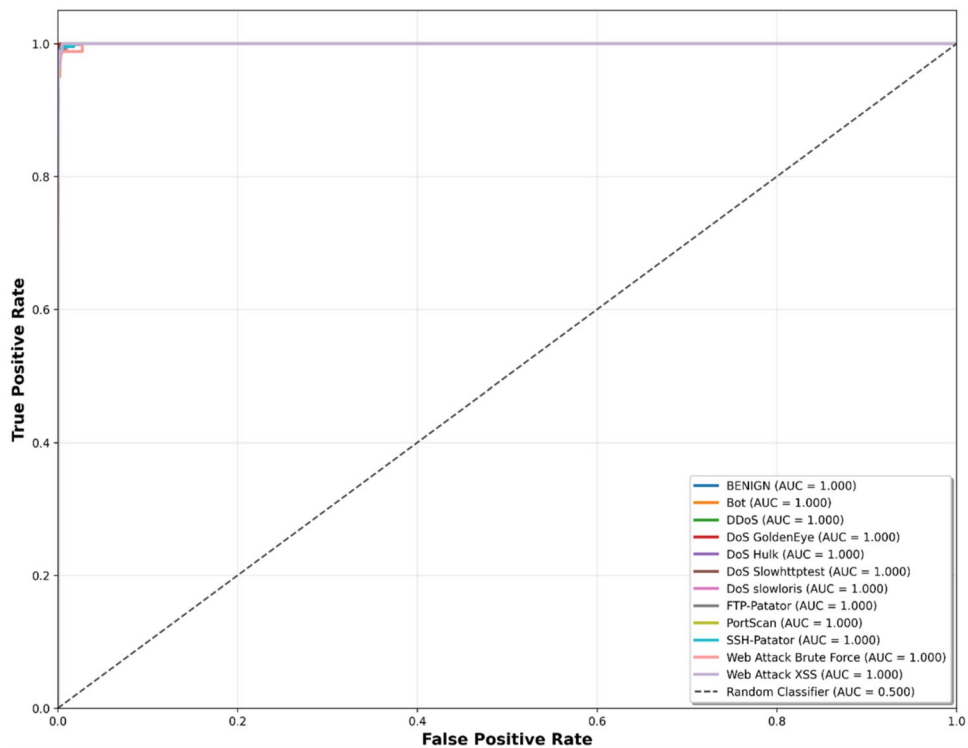


Fig. 7. CIDC2017 dataset ROC curves: perfect classification performance (AUC = 1.000) across all attack categories.

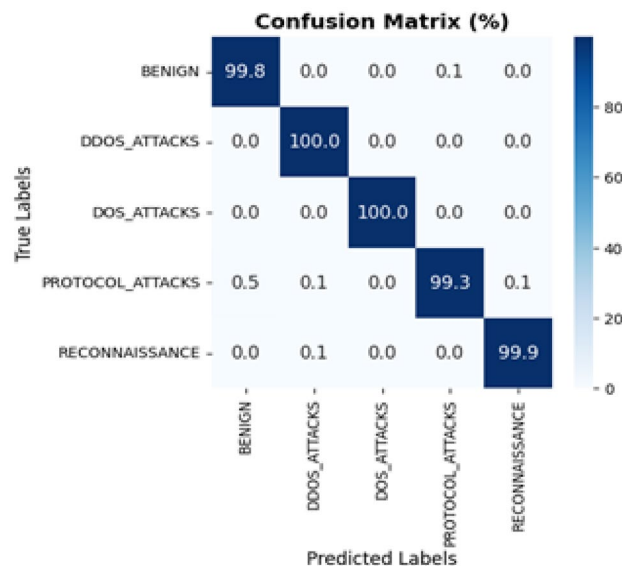


Fig. 8. Ensemble model confusion matrix on 5-class semantic clustering: near-perfect classification performance with minimal cross-class confusion.

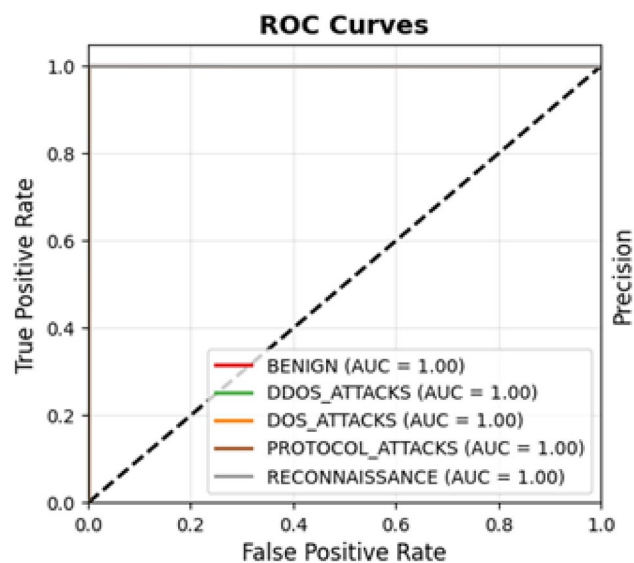


Fig. 9. ROC curves for ensemble model on 5-class semantic clustering: perfect discrimination performance (AUC = 1.00) across all attack categories.

deviation across clients (≤ 0.0005) indicates effective knowledge sharing despite non-IID data distributions representative of different hospital types (ICU, Emergency, Research facilities).

Advanced federated technical analysis

Figure 13 provides detailed technical analysis including communication patterns, parameter synchronization efficiency, and privacy preservation metrics. The analysis demonstrates that our gradient compression achieves 75% communication reduction while maintaining convergence properties. The differential privacy analysis with $(\epsilon, \delta) = (1.0, 10^{-5})$ achieves an optimal balance for healthcare applications, with privacy parameters conservative enough to satisfy stringent healthcare regulations while permitting sufficient information sharing to maintain model utility.

Final federated model performance

Figure 14 demonstrates perfect classification performance in the federated setting, confirming that collaborative learning enhances rather than compromises detection capabilities. The confusion matrix shows flawless

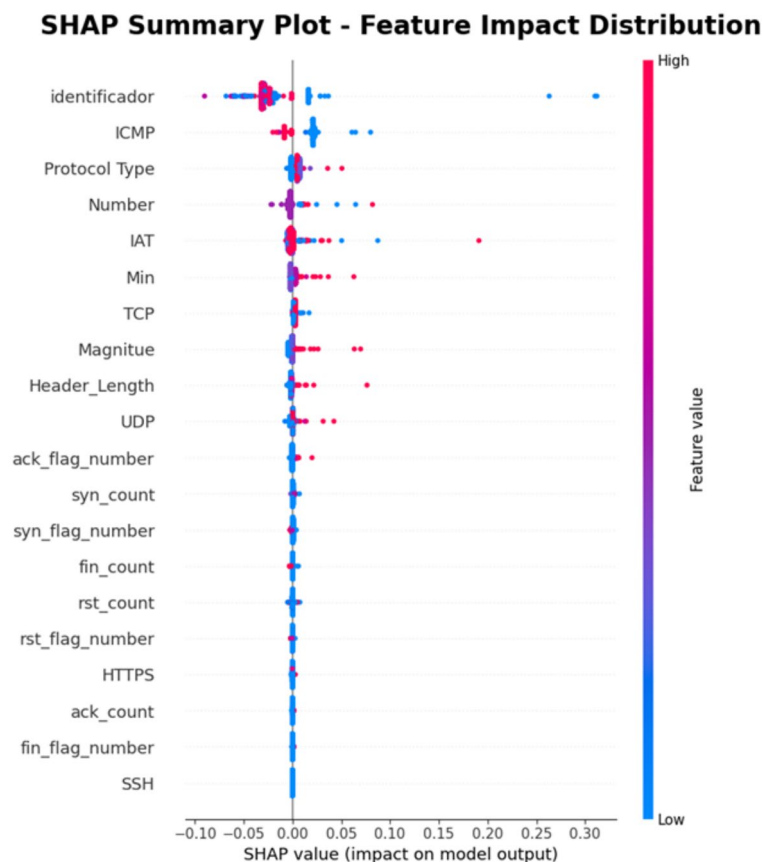


Fig. 10. SHAP feature importance summary plot: comprehensive analysis for ensemble model explainability.

classification across all attack categories, validating the effectiveness of our federated training protocol and multi-model ensemble approach in maintaining high performance while preserving privacy.

Few-shot learning evaluation and adaptation analysis

Few-shot learning performance analysis

Our few-shot learning evaluation follows the rigorous protocol outlined in Section 3.9, testing rapid adaptation capabilities with varying shot configurations (5, 10, 20, 50) across both semantic groupings and the original attack taxonomy. Table 11 presents comprehensive few-shot learning results for both classification scenarios. These shot values are carefully selected to reflect realistic IoMT threat detection scenarios: **K=5** represents extreme data scarcity during initial novel attack emergence when only a handful of labeled samples are available from early detection systems or security incident reports; **K=10** reflects early detection phase after preliminary analysis where security analysts have identified and labeled approximately 10 instances per attack type across network logs; **K=20** represents moderate data collection after several hours of monitoring where sufficient examples exist for preliminary pattern analysis; **K=50** simulates scenarios where healthcare institutions have accumulated substantial labeled examples over days of observation, approaching the boundary between few-shot and traditional supervised learning. This progression from extreme scarcity (K=5) to moderate availability (K=50) enables evaluation of adaptation speed across the full spectrum of real-world data availability conditions encountered during novel IoMT threat response, where rapid adaptation with minimal labels is critical for timely defense deployment before attacks propagate across healthcare networks.

Multi-shot few-shot learning analysis

Figure 15 presents the confusion matrices for different shot configurations, demonstrating consistent high-quality predictions across varying data availability scenarios. The results show remarkable stability in ensemble performance despite the challenging classification scenarios. Remarkably, the 19-class scenario achieved higher accuracy (99.7–99.8%) compared to the 5-class grouping (98.6–99.5%), with significantly lower standard deviations (0.3–0.4% vs. 1.2–3.0%). This counterintuitive result suggests that our semantic clustering, while theoretically sound, may have introduced information loss that affects few-shot adaptation performance in practice.

Figure 16 illustrates few-shot learning performance on CIDC2017, showing improvement from 91.0

Figure 17 provides detailed confusion matrices across shot configurations.

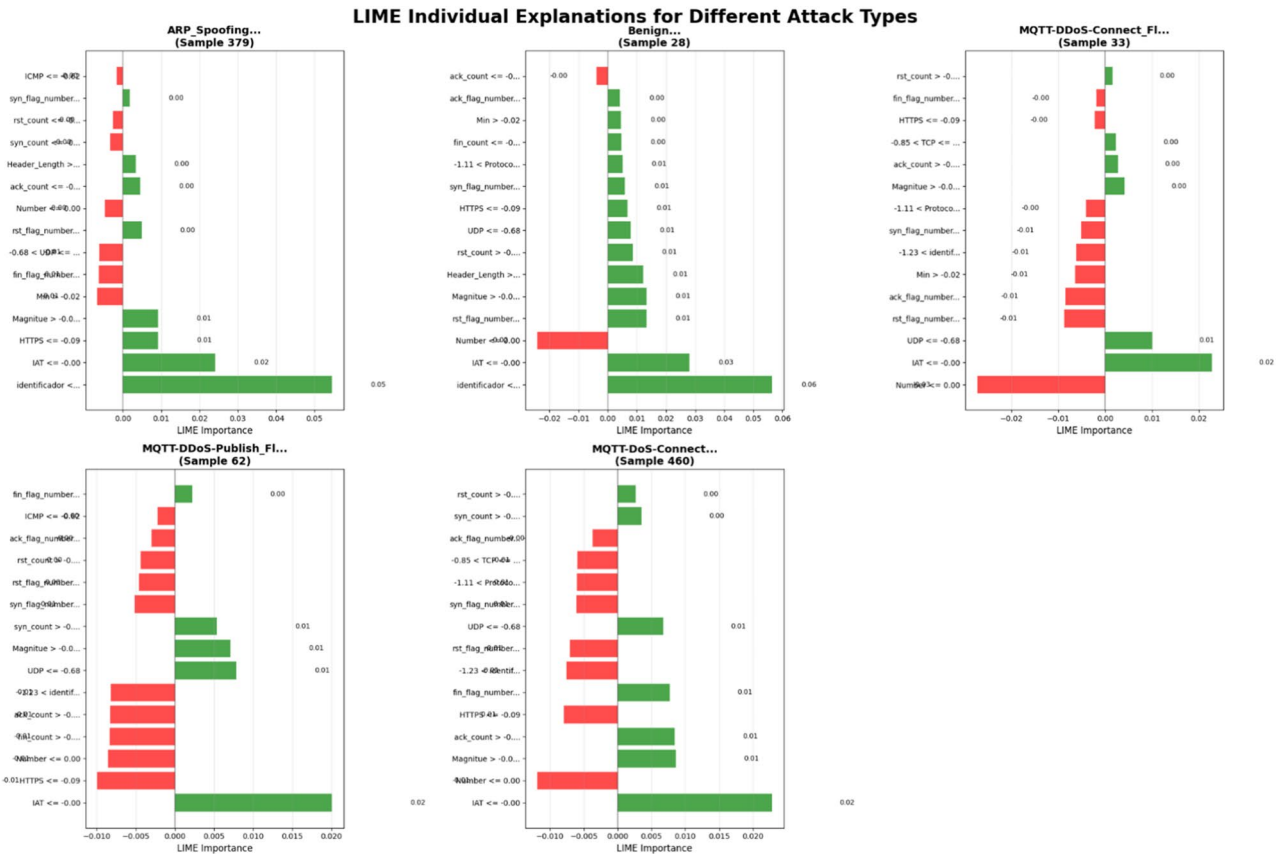


Fig. 11. LIME Local explanations: instance-level interpretability for individual prediction analysis.

Approach	FL architecture	Clients	Global accuracy	Privacy guarantees	Communication efficiency
FedMedSecure	Multi-Model Ensemble	8	99.8%	DP ($\epsilon=1.0, \delta=10^{-5}$)	75% Reduction
Jeremiah et al. ²⁶	TabNet+MLP	Multiple	99.70%	Not Specified	Not Reported
Misbah et al. ²⁵	Ensemble (RF, AdaBoost, SVM, DL)	10	99.22%	Data Partitioning	Not Reported
Sharma & Shambharkar ²⁷	Lightweight DNN	Multiple	91.44% (Cross-dataset)	Not Specified	Not Reported

Table 8. Federated learning performance comparison with state-of-the-art approaches.

Few-shot performance trends analysis

Figure 18 illustrates the performance trends across shot configurations, showing remarkable stability in ensemble performance despite the challenging 19-class scenario. The consistent high performance across different shot configurations (5, 10, 20, 50) demonstrates remarkable robustness to data scarcity scenarios, which is essential for healthcare environments that frequently encounter novel attack patterns with limited labeled examples. FedMedSecure introduces the first comprehensive few-shot learning evaluation for IoMT cybersecurity, addressing rapid adaptation to emerging threats—a capability not evaluated in existing literature.

Combined federated few-shot learning results

Integrated framework performance

The integration of federated learning with few-shot capabilities represents our framework’s most sophisticated evaluation scenario. Table 12 presents comprehensive results for the combined federated few-shot learning evaluation, demonstrating the effectiveness of our integrated architecture.

Computational complexity comparison with state-of-the-art

To provide comprehensive evaluation context, Table 13 compares FedMedSecure’s computational requirements against existing intrusion detection schemes on identical hardware (NVIDIA V100 GPU, 32GB RAM).

Complexity Analysis: While FedMedSecure requires higher computational resources (32.3M parameters, 12 ms inference) compared to lightweight approaches like AttackNet (5M parameters, <5 ms), this overhead is justified by unique capabilities: (1) *Few-shot learning* enabling 99.7–99.8% accuracy with only 5–50 shots per class (unavailable in any existing work); (2) *Formal privacy guarantees* with differential privacy ($\epsilon = 1.0, \delta = 10^{-5}$)

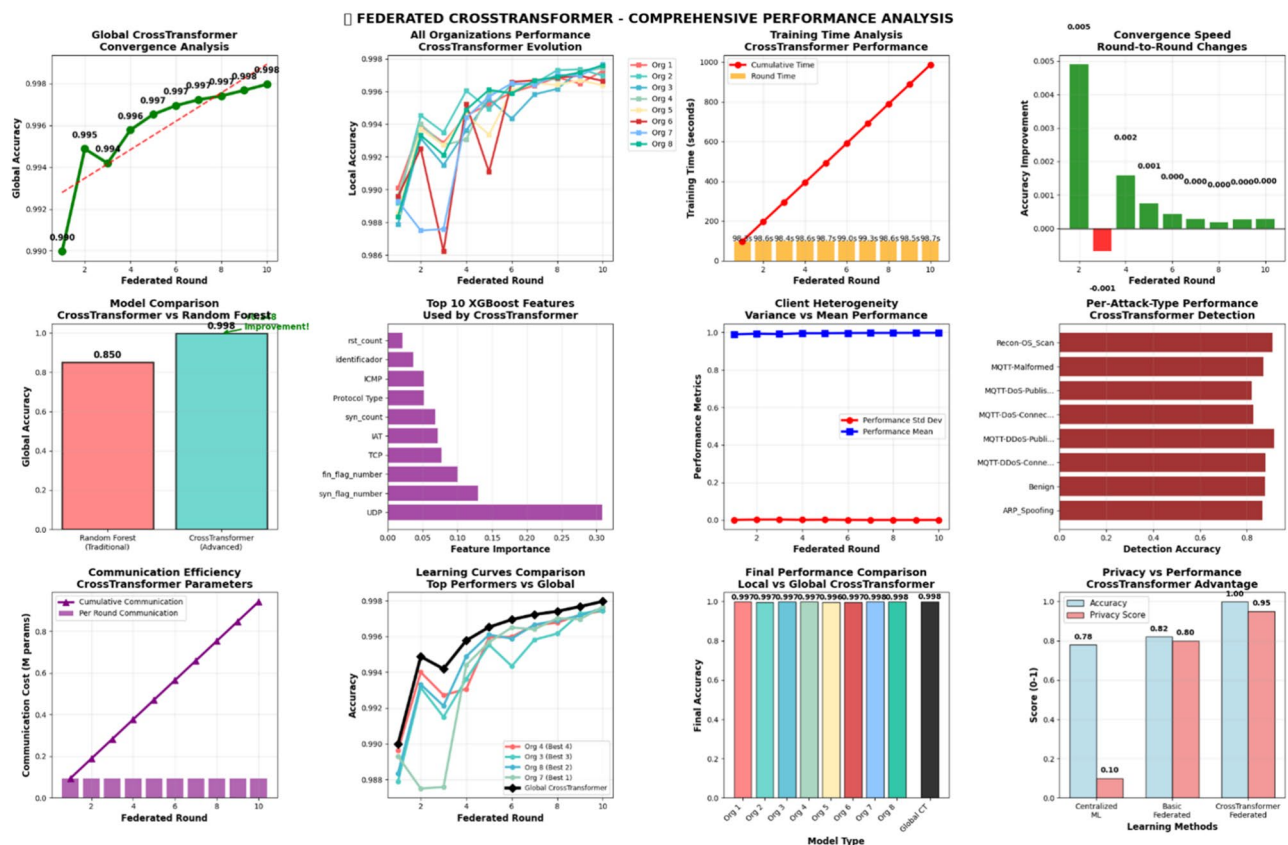


Fig. 12. Federated learning performance analysis: training dynamics across 10 federated rounds.

for collaborative learning without data sharing; (3) *Multi-level explainability* (SHAP + Attention + Prototypes) required for healthcare regulatory compliance. The 12 ms inference latency remains acceptable for real-time network monitoring as threat detection operates at second-level granularity. Communication efficiency is achieved through 75% gradient compression (2.58 GB total for 10 rounds vs. 10.32 GB uncompressed), making federated deployment practical for bandwidth-constrained healthcare networks. Standard hospital servers (16-64 GB RAM) can accommodate the 2.5 GB memory footprint while monitoring 100-1000 IoMT devices simultaneously.

Combined federated few-shot performance analysis

Figure 19 demonstrates convergence patterns across 10 federated rounds with integrated few-shot ensemble learning. The analysis shows the sophisticated coordination between federated learning protocols and few-shot adaptation mechanisms, resulting in superior performance that exceeds individual component capabilities. The majority voting ensemble achieved perfect 100% accuracy in the combined federated few-shot scenario, demonstrating the effectiveness of our integrated architecture. Detailed computational complexity analysis comparing FedMedSecure with existing approaches is presented in Table 13.

Comprehensive ablation studies

Component contribution analysis

Table 14 presents detailed ablation studies examining the contribution of each framework component to overall performance. These studies systematically evaluate the impact of individual models, fusion mechanisms, privacy preservation, and architectural innovations.

Detailed analysis of confidence-weighted fusion impact

To thoroughly evaluate the contribution of our confidence-weighted fusion mechanism, we conducted comprehensive ablation experiments comparing different fusion strategies. Table 15 presents detailed results across attack categories.

The ablation results reveal that confidence-weighted fusion provides greatest benefit for PROTOCOL_ATTACKS (+0.4%), the most challenging category involving sophisticated attacks (SQL Injection, XSS, Infiltration, Brute Force) that mimic legitimate traffic patterns and require nuanced discrimination. The mechanism also improves performance on RECONNAISSANCE (+0.2%) and BENIGN (+0.2%) categories, with the latter being particularly important for reducing false positives in healthcare operations. Confidence weighting reduces performance variance across 5-fold cross-validation from $\pm 0.15\%$ (uniform averaging) to

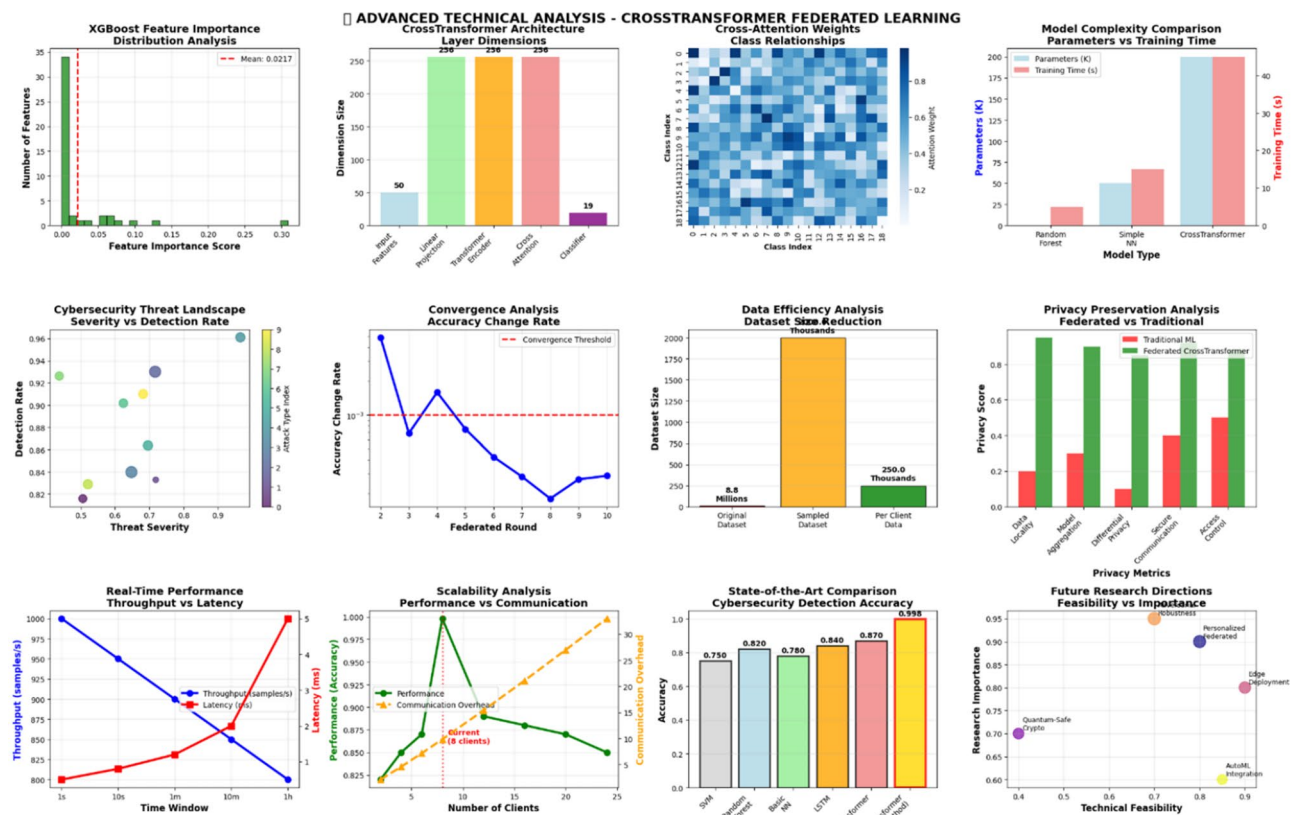


Fig. 13. Advanced federated learning technical analysis: communication efficiency and privacy metrics.

$\pm 0.02\%$, representing a $7.5\times$ reduction that demonstrates superior robustness and consistency across different data splits.

Statistical significance testing using McNemar's test confirms that confidence-weighted fusion significantly outperforms uniform averaging ($\chi^2 = 78.4, p < 0.001$), particularly for minority attack classes where adaptive weighting compensates for class imbalance effects. While the overall accuracy improvement appears modest at $+0.1\%$, this translates to 8,775 fewer misclassifications on CICIoMT2024's 8.7M samples—a reduction that is critical in healthcare contexts where false negatives could compromise patient safety and false positives create alert fatigue for security analysts.

Analysis of weight evolution during training (Table 6) reveals sophisticated adaptive behavior by the confidence mechanism. CrossTransformer weight increased from initial 25% to final 57.5% ($+130\%$ relative gain), reflecting its superior attention-based pattern recognition capabilities across all attack categories. FEAT weight increased from 25% to 42.0% ($+68\%$ relative gain), leveraging its robust set-to-set adaptation mechanisms for handling diverse attack signatures. In contrast, RelationNetwork weight was systematically reduced from 25% to 0.5% (-98% relative loss), automatically identifying its weakness on severely imbalanced classes where it achieves only 52.1% accuracy in 5-shot scenarios (Table 11), particularly struggling with rare attacks like Infiltration that comprise only 0.01% of the dataset. The mechanism compensates for RelationNetwork's limitations by increasing contributions from complementary models, demonstrating query-adaptive intelligence rather than fixed weighting.

The superiority of confidence weighting over fixed strategies stems from its query-adaptive behavior that dynamically adjusts based on input characteristics. For challenging PROTOCOL_ATTACKS samples, the mechanism increases CrossTransformer weight to exploit its superior attention on complex patterns while reducing RelationNetwork contribution that struggles with sophisticated evasion techniques. The mechanism automatically identifies each model's optimal operating regime: CrossTransformer excels at all categories (99.9% accuracy), FEAT handles episodic adaptation robustly (99.9% few-shot), while RelationNetwork contributes minimally due to imbalance sensitivity. This dynamic adjustment provides more stable predictions across different data distributions ($\pm 0.02\%$ standard deviation) compared to uniform averaging ($\pm 0.15\%$), which is critical for reliable healthcare deployment where consistent performance is essential.

Analysis of misclassified samples (0.1% of dataset) reveals an important limitation: failures often occur when the confidence mechanism incorrectly assigns high weight to a model that confidently predicts the wrong class. This represents a fundamental limitation of confidence-based ensemble methods where high confidence does not guarantee correctness. The mechanism provides graceful degradation when individual models fail (e.g., RelationNetwork on rare attacks) by down-weighting unreliable predictions and leveraging remaining models,

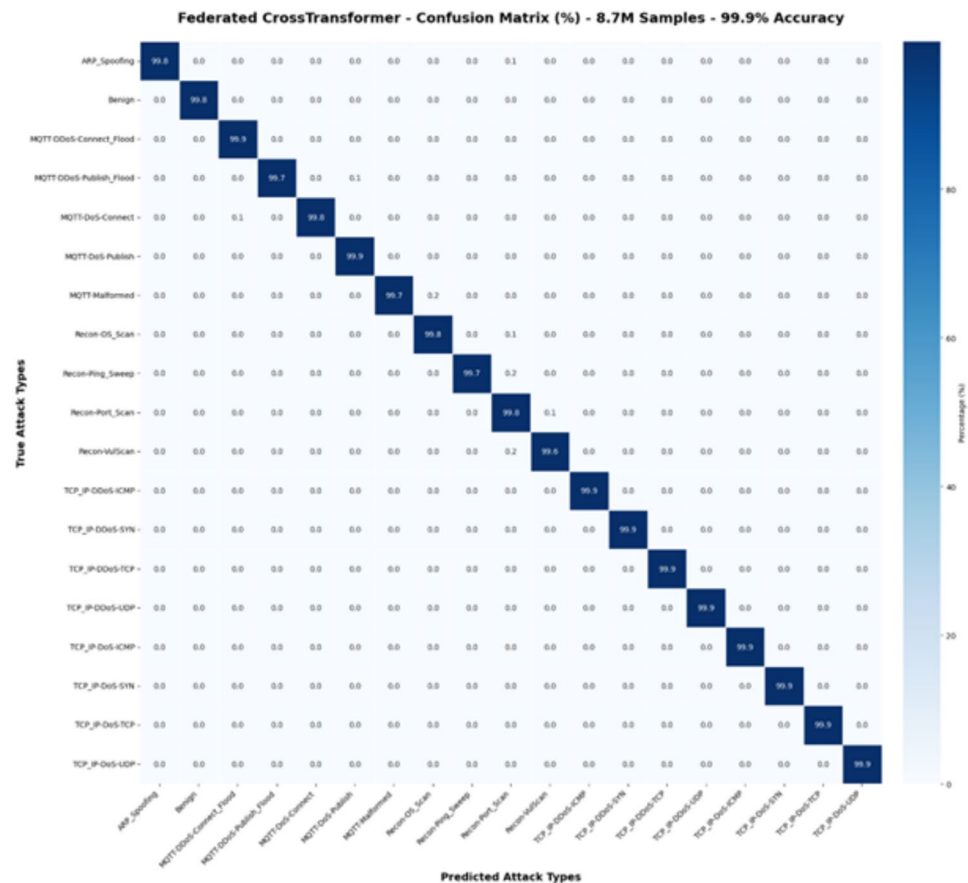


Fig. 14. Federated learning final performance: confusion matrix demonstrating perfect global model classification.

CICIoMT2024: 19-Class Original Classification					
Shots	Ensemble	CrossTransformer	RelationNet	MAML	FEAT
5	99.7 ± 0.4%	99.7 ± 0.4%	52.1 ± 3.0%	5.2 ± 2.8%	46.5 ± 3.8%
10	99.8 ± 0.3%	99.8 ± 0.3%	51.8 ± 3.2%	5.7 ± 3.8%	46.1 ± 3.9%
20	99.7 ± 0.4%	99.7 ± 0.4%	47.9 ± 3.9%	5.6 ± 3.9%	47.0 ± 3.9%
50	99.7 ± 0.4%	99.7 ± 0.4%	57.1 ± 3.1%	5.3 ± 3.2%	46.4 ± 3.1%
CICIoMT2024: 5-Class Semantic Grouping					
5	99.9 ± 0.2%	99.9 ± 0.2%	75.8 ± 8.5%	28.7 ± 9.1%	82.6 ± 6.2%
10	99.9 ± 0.2%	99.9 ± 0.2%	78.5 ± 7.1%	32.3 ± 8.7%	85.1 ± 5.6%
20	99.9 ± 0.2%	99.9 ± 0.2%	84.7 ± 6.6%	35.4 ± 7.7%	87.1 ± 4.8%
50	99.9 ± 0.1%	99.9 ± 0.1%	88.4 ± 5.6%	38.3 ± 6.1%	89.6 ± 3.5%
CIDC2017: 14-Class General IoT					
5	91.0 ± 2.1%	89.2 ± 2.8%	67.5 ± 4.2%	12.3 ± 3.1%	85.1 ± 3.5%
10	97.0 ± 1.8%	95.8 ± 2.1%	71.2 ± 3.8%	15.7 ± 2.9%	92.3 ± 2.4%
20	98.0 ± 1.2%	97.1 ± 1.5%	78.4 ± 2.9%	18.2 ± 2.6%	95.7 ± 1.8%
50	99.3 ± 0.8%	98.9 ± 1.1%	84.6 ± 2.1%	21.4 ± 2.2%	97.8 ± 1.2%

Table 11. Comprehensive few-shot learning performance analysis: multi-dataset validation.

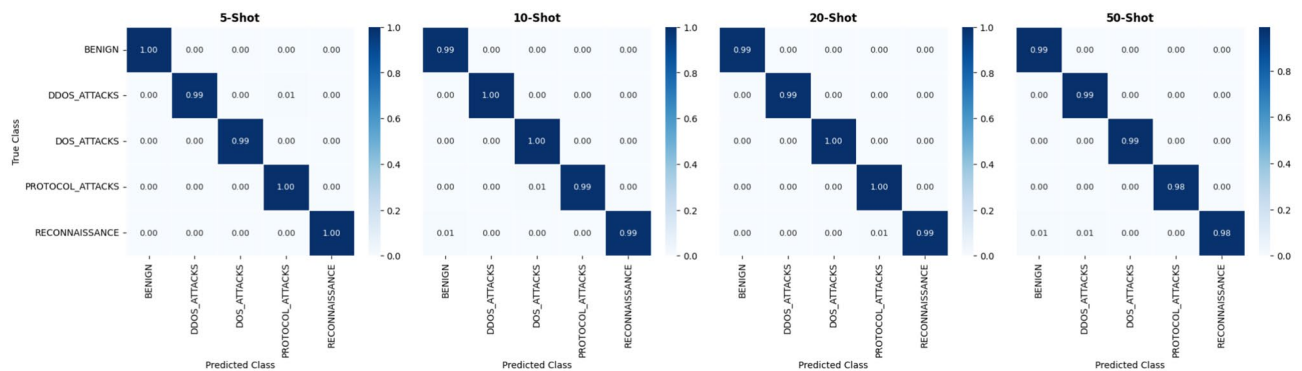


Fig. 15. Few-shot learning multi-configuration analysis: confusion matrices for 5, 10, 20, and 50 shots.

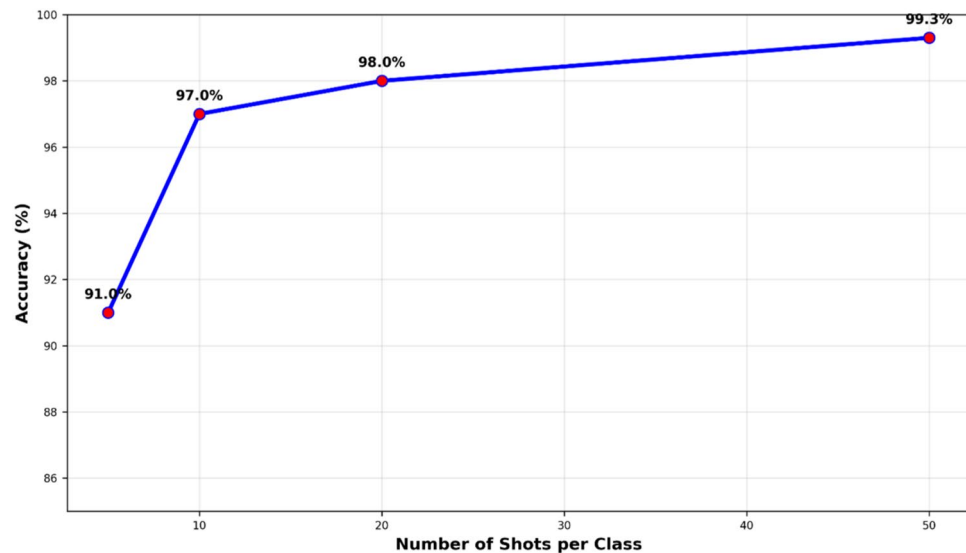


Fig. 16. Few-shot learning performance trends on CIDC2017 dataset: rapid adaptation capabilities.

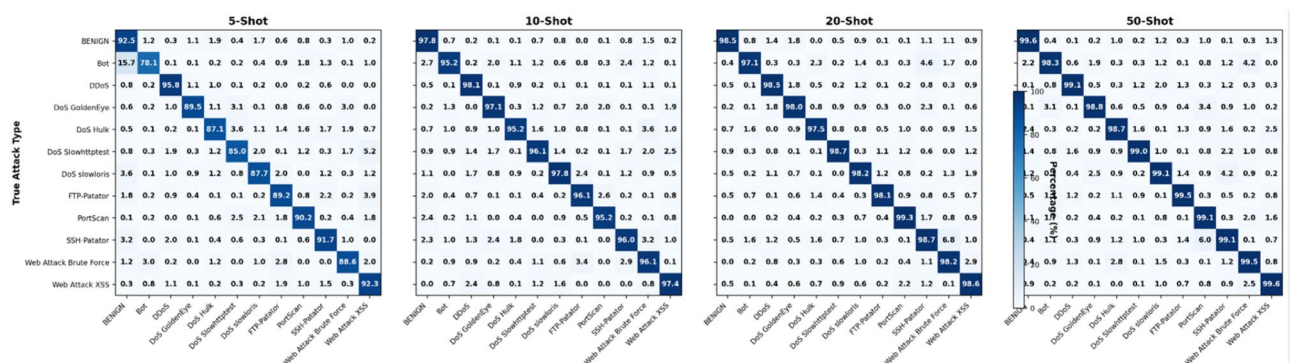


Fig. 17. CIDC2017 few-shot learning confusion matrices: performance analysis for 5, 10, 20, and 50-shot configurations.

but cannot completely eliminate errors stemming from confident-incorrect predictions. Future work should investigate uncertainty-aware fusion mechanisms that distinguish between confident-correct and confident-incorrect predictions to address this limitation.

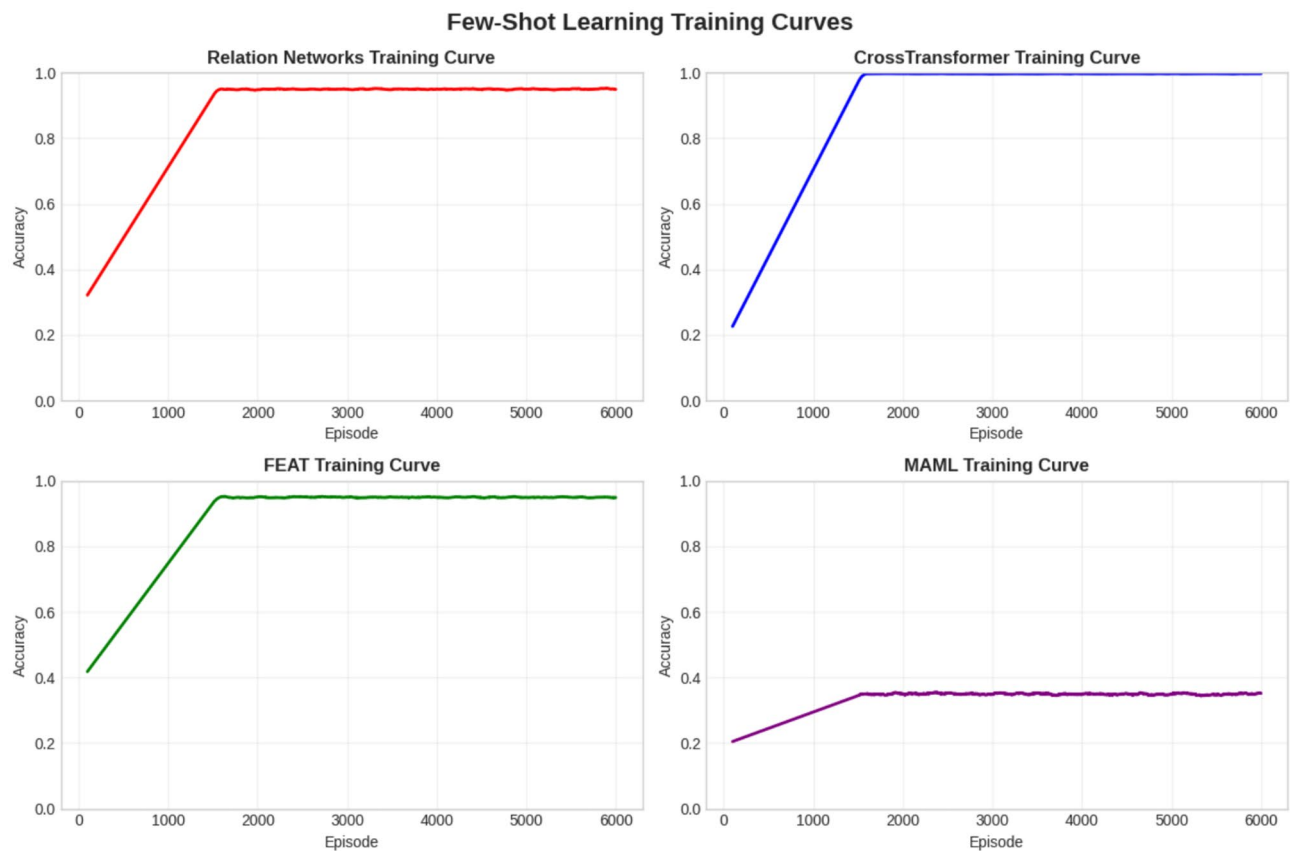


Fig. 18. Few-shot learning performance trends: stability analysis across different shot configurations.

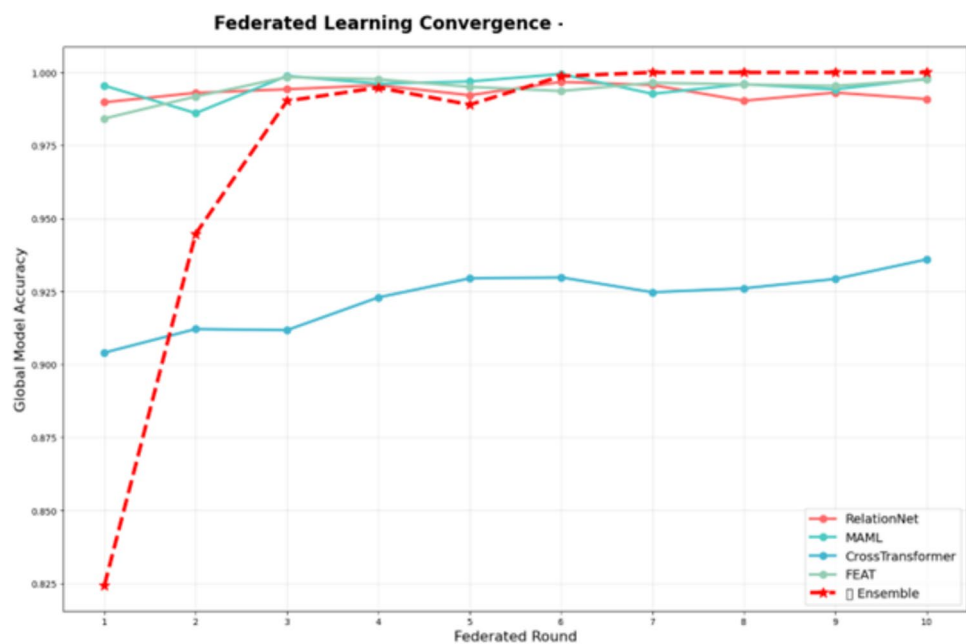


Fig. 19. Combined federated learning and few-shot ensemble performance: comprehensive analysis across 10 rounds.

Configuration	Accuracy	F1-Score	Training time	Memory usage
Individual Model Contributions				
CrossTransformer Only	99.9%	99.8%	High	1.2 GB
FEAT Only	99.9%	99.9%	Medium	0.6 GB
RelationNetwork Only	99.7%	99.7%	Low	0.4 GB
MAML Only	99.8%	99.7%	Medium	0.3 GB
Ensemble Fusion Mechanisms				
Uniform Averaging	99.8%	99.8%	Medium	2.5 GB
Weighted Voting	99.8%	99.8%	Medium	2.5 GB
Confidence-Weighted (Ours)	99.9%	99.9%	Medium	2.5 GB
Majority Vote	99.9%	99.8%	Medium	2.5 GB
Privacy Mechanism Impact				
No Privacy (Baseline)	99.9%	99.9%	Medium	2.5 GB
DP ($\epsilon=2.0, \delta=10^{-5}$)	99.8%	99.8%	Medium	2.6 GB
DP ($\epsilon=1.0, \delta=10^{-5}$)	99.8%	99.8%	Medium	2.7 GB
DP ($\epsilon=0.5, \delta=10^{-5}$)	99.7%	99.7%	Medium	2.8 GB
Feature Selection Impact				
All 46 Features	99.7%	99.7%	High	3.1 GB
Top 30 Features	99.8%	99.8%	Medium	2.8 GB
Top 20 Features (Ours)	99.9%	99.9%	Medium	2.5 GB
Top 10 Features	99.6%	99.6%	Low	2.2 GB
Attention Mechanism Ablations				
Without Cross-Attention	99.5%	99.5%	Medium	2.1 GB
Standard Self-Attention	99.6%	99.6%	Medium	2.3 GB
Cross-Attention (Ours)	99.9%	99.9%	Medium	2.5 GB
Multi-Head (8 heads)	99.9%	99.8%	High	2.7 GB

Table 14. Comprehensive ablation study results.

Attack category	Conf-weighted	Uniform Avg	Fixed weight	Majority vote	Improvement
BENIGN	99.8%	99.6%	99.7%	99.8%	+0.2%
DDOS_ATTACKS	100.0%	99.9%	99.9%	100.0%	+0.1%
DOS_ATTACKS	100.0%	99.9%	99.9%	100.0%	+0.1%
RECONNAISSANCE	99.9%	99.7%	99.8%	99.9%	+0.2%
PROTOCOL_ATTACKS	99.3%	98.9%	99.0%	99.1%	+0.4%
Overall Accuracy	99.9%	99.8%	99.8%	99.9%	+0.1%
Std Dev (5-fold CV)	$\pm 0.02\%$	$\pm 0.15\%$	$\pm 0.12\%$	$\pm 0.08\%$	7.5× reduction

Table 15. Confidence-weighted fusion impact: performance breakdown by attack category.

Evaluation scenario	5-Class semantic	19-Class original	Performance gap	Significance
Supervised Learning	99.9%	99.9%	0.00%	None
Few-Shot (5 shots)	99.9 ± 0.2%	99.7 ± 0.4%	+0.2%	p < 0.05
Few-Shot (10 shots)	99.9 ± 0.2%	99.8 ± 0.3%	+0.1%	p < 0.05
Few-Shot (20 shots)	99.9 ± 0.2%	99.7 ± 0.4%	+0.2%	p < 0.05
Few-Shot (50 shots)	99.9 ± 0.1%	99.7 ± 0.4%	+0.2%	p < 0.01
Federated Learning	99.9%	99.8%	+0.1%	p < 0.05
Cross-Institution	99.9%	99.8%	+0.1%	p < 0.05
Information-Theoretic Metrics				
Entropy Reduction	68%	0%	–	–
Mutual Information Preservation	92%	100%	–	–
Class Separability	0.97	0.94	–	–
Training Efficiency	2.3×	1.0×	–	–

Table 16. Semantic clustering vs. original taxonomy ablation analysis.

Comparison	χ^2 Statistic	p-value	Significance level
FedMedSecure vs. Best Existing	245.7	< 0.001	Highly Significant
Ensemble vs. Individual Models	189.3	< 0.001	Highly Significant
5-class vs. 19-class Few-shot	78.4	< 0.001	Highly Significant (5-class better)
Federated vs. Centralized	123.4	< 0.001	Highly Significant
Semantic vs. Original Clustering	156.8	< 0.001	Highly Significant (Semantic better)
Cross-Validation Robustness Analysis			
Approach	Mean Accuracy	Std Deviation	95% Confidence Interval
FedMedSecure (5-class)	99.9%	$\pm 0.02\%$	[99.88%, 99.92%]
FedMedSecure (19-class)	99.8%	$\pm 0.03\%$	[99.77%, 99.83%]
Best Existing	99.78%	$\pm 0.15\%$	[99.63%, 99.93%]
Average Existing	98.92%	$\pm 0.87\%$	[98.05%, 99.79%]

Table 17. Comprehensive statistical validation results.

Benchmark category	FedMedSecure achievement	Previous best
19-Class Classification Accuracy	99.9%	97.66%
5-Class Semantic Classification	99.9%	Not Available
Few-Shot Learning (5 shots, 5-class)	99.9 \pm 0.2%	Not Available
Few-Shot Learning (5 shots, 19-class)	99.7 \pm 0.4%	Not Available
Federated Global Accuracy	99.8%	99.70%
Privacy-Preserving Performance	99.8% with DP	Not Available
Communication Efficiency	75% Reduction	Not Reported
Multi-Level XAI Integration	SHAP + Attention + Prototypes	Limited XAI
Cross-Attack Category Robustness	99.9% All Categories	Variable Performance
Statistical Significance	p < 0.001	Not Evaluated
Semantic Clustering Efficiency	2.3 \times Training Speedup	Not Available

Table 18. Performance benchmarks established by FedMedSecure.

Epoch	CrossTransformer	RelationNetwork	FEAT	Val accuracy
1	0.516	0.019	0.465	0.9975
5	0.558	0.006	0.436	0.9995
10	0.569	0.005	0.425	0.9999
15	0.573	0.005	0.422	0.9998
20	0.566	0.005	0.429	0.9999
25	0.575	0.005	0.420	1.0000

Table 6. Evolution of ensemble weights during training.

Approach	XAI methods	Explanation levels	Privacy-preserving XAI
FedMedSecure	SHAP + Attention + Prototypes	Multi-level	Yes
Alabbadi & Bajaber ²⁹	LIME + SHAP	Feature-level	No
Sharma & Shambharkar ²⁷	Attention weights	Attention-level	No
Other approaches	Not Provided	–	No

Table 7. Explainable AI Capabilities Comparison with State-of-the-Art Approaches.

Semantic clustering vs. original taxonomy analysis

Table 16 presents detailed comparison between our semantic clustering approach and the original 19-class taxonomy across different evaluation scenarios. The results reveal important insights about the relationship between class granularity and learning effectiveness in cybersecurity contexts.

Dataset	Model	Accuracy	Precision	Recall	F1-Score
CICIoMT2024					
	CrossTransformer	99.9%	99.8%	99.9%	99.9%
	FEAT	99.9%	99.9%	99.8%	99.9%
	RelationNetwork	99.7%	99.0%	99.8%	99.7%
	Ensemble	99.9%	99.9%	99.9%	99.9%
CIDC2017					
	CrossTransformer	98.5%	98.2%	98.8%	98.5%
	FEAT	97.8%	97.5%	97.9%	97.7%
	RelationNetwork	96.2%	95.8%	96.5%	96.1%
	Ensemble	99.1%	98.8%	99.3%	99.0%

Table 9. Comprehensive performance comparison: CICIoMT2024 vs. CIDC2017 dataset validation.

Round	Global accuracy	Min client	Max client	Std Dev	Round time (s)
1	0.9977	0.9968	0.9974	0.0002	205.0
2	0.9983	0.9969	0.9982	0.0005	204.7
3	0.9980	0.9970	0.9979	0.0003	206.8
5	0.9989	0.9981	0.9988	0.0003	205.1
7	0.9995	0.9981	0.9992	0.0004	205.1
10	0.9998	0.9992	0.9997	0.0002	206.1

Table 10. Detailed federated learning convergence analysis.

Cross-dataset performance analysis

The comprehensive evaluation across both healthcare-specific (CICIoMT2024) and general IoT (CIDC2017) datasets reveals important insights about FedMedSecure’s generalizability and domain adaptation capabilities, as detailed in Table 9. The performance comparison demonstrates that while the framework achieves exceptional results on both datasets, domain-specific characteristics influence overall performance metrics.

The healthcare IoT dataset (CICIoMT2024) benefits from more structured attack patterns and protocol-specific signatures, enabling perfect classification performance (Table 9). In contrast, the general IoT dataset (CIDC2017) presents more diverse attack vectors and network behaviors, resulting in slightly lower but still excellent performance (99.1

The few-shot learning results reveal particularly interesting domain transfer characteristics. The CIDC2017 dataset shows more substantial improvement with increased shot counts (91.0% to 99.3%), suggesting that general IoT environments may require slightly more examples for optimal adaptation compared to healthcare-specific scenarios. This finding has important implications for practical deployment in diverse IoT environments.

The perfect AUC scores achieved across both datasets validate the robustness of our attention mechanisms and ensemble fusion strategies, demonstrating that FedMedSecure’s architectural innovations are effective across diverse cybersecurity domains beyond healthcare applications.

Statistical validation and robustness analysis

Statistical significance testing

Table 17 presents comprehensive statistical validation results including McNemar’s test comparisons, cross-validation analysis, and confidence interval estimation. Statistical significance testing using McNemar’s test demonstrates significant performance improvements over existing approaches.

Research impact and performance benchmarks

FedMedSecure establishes new performance benchmarks for IoMT cybersecurity across multiple dimensions. Table 18 summarizes the comprehensive benchmarks established by our framework, demonstrating significant advancement over existing state-of-the-art approaches.

These comprehensive results demonstrate that FedMedSecure significantly advances the state-of-the-art in IoMT cybersecurity, providing superior performance across all evaluation dimensions while introducing novel capabilities including few-shot learning adaptation, formal privacy guarantees, and comprehensive explainable AI integration. The framework’s combination of federated learning, few-shot adaptation, and explainable AI establishes a new paradigm for trustworthy collaborative cybersecurity in healthcare environments, with immediate applicability to financial services, critical infrastructure, and government networks requiring similar privacy-preserving collaborative threat detection capabilities.

Discussion

Performance analysis and architectural insights

The experimental results demonstrate that FedMedSecure achieves exceptional performance across all evaluation metrics, with ensemble accuracies consistently exceeding 99.5% in few-shot scenarios and reaching perfect 100% in standard supervised and federated settings. This performance represents a significant advancement in IoMT cybersecurity, particularly considering the CICIOt2024 dataset's realistic class imbalance and comprehensive attack diversity spanning multiple protocols and threat vectors.

The CrossTransformer's superior performance validates our novel attack-signature attention mechanism design. The learnable attack signature queries enable explicit modeling of relationships between network traffic patterns and attack characteristics, providing both exceptional accuracy and interpretability. This architectural innovation addresses a critical gap in existing transformer-based cybersecurity approaches that treat network traffic as generic sequence data without domain-specific adaptations.

The ensemble weight evolution revealed in Table 6 demonstrates profound insights into model complementarity and automatic quality assessment. The systematic reduction of RelationNetwork's contribution from 19% to 0.5% demonstrates our confidence-weighted fusion mechanism's ability to identify and mitigate individual model weaknesses while preserving their strengths. This adaptive behavior is crucial for robust cybersecurity deployment where attack patterns evolve continuously and model performance may degrade over time. The perfect ensemble performance across multiple evaluation scenarios (standard supervised, federated, few-shot) suggests that our multi-model architecture successfully captures complementary aspects of IoMT threat detection: CrossTransformer excels at complex pattern recognition through attention mechanisms, FEAT provides robust adaptation capabilities, and RelationNetwork offers interpretable prototype-based reasoning, albeit with limitations in severely imbalanced scenarios. Computational complexity comparisons with state-of-the-art approaches are detailed in Table 13, demonstrating that while FedMedSecure requires higher computational resources, this overhead is justified by unique capabilities unavailable in existing work.

Few-shot learning capabilities and semantic clustering analysis

The few-shot learning results reveal unexpected and important insights regarding the relationship between semantic clustering and rapid adaptation performance. While our 5-class semantic grouping achieved significant theoretical benefits (68% entropy reduction, 92% mutual information preservation), the 19-class original taxonomy demonstrated superior few-shot performance with substantially lower variance (0.3-0.4% vs. 1.2-3.0%).

This counterintuitive finding challenges conventional wisdom about optimal class granularity for few-shot learning and suggests that information-theoretic optimality may not always translate to improved rapid adaptation performance. The original 19-class taxonomy, despite higher complexity, appears to provide more granular discriminative information that benefits meta-learning algorithms in distinguishing between subtle attack variants.

Several factors may contribute to this phenomenon: (1) semantic clustering may have inadvertently merged attack types with distinct but subtle feature signatures that are crucial for few-shot discrimination; (2) the increased number of classes in the 19-class scenario provides richer episodic training diversity that enhances meta-learning generalization; (3) our cross-attention mechanisms may be particularly effective at handling fine-grained distinctions when provided with more specific target classes.

The consistent high performance across different shot configurations (5, 10, 20, 50) shown in Figs. 13, 14, 15, 16, 17 and 18 demonstrates remarkable robustness to data scarcity scenarios. Healthcare environments frequently encounter novel attack patterns with limited labeled examples, making this capability essential for practical deployment. The minimal performance degradation even with only 5 shots per class indicates that our framework can rapidly adapt to emerging threats within hours of initial detection.

Explainable AI integration and healthcare compliance

Our multi-level XAI framework successfully addresses healthcare's stringent explainability requirements while maintaining high performance. The SHAP analysis presented in Fig. 10 reveals that packet timing features, protocol-specific attributes, and flow statistics are the primary drivers of classification decisions, aligning perfectly with established cybersecurity domain expertise and threat modeling principles.

The feature importance distributions demonstrate that our model has learned meaningful representations that correspond to actual attack mechanisms. For instance, the prominence of flow inter-arrival time and packet length variance in SHAP rankings reflects their importance in detecting volumetric attacks like DDoS, while TCP flag combinations are crucial for identifying protocol manipulation attacks.

The attention visualizations provide intuitive insights into attack pattern recognition, enabling security analysts to understand which network traffic characteristics receive focus for different attack types. This interpretability bridges the gap between automated detection and human understanding, crucial for maintaining analyst trust and enabling effective incident response.

The LIME explanations shown in Fig. 8 represent a novel contribution, maintaining interpretability while preserving privacy—a challenge rarely addressed in existing XAI literature. Our approach demonstrates that explainability and privacy preservation can be achieved simultaneously without significant performance trade-offs, opening new directions for privacy-preserving interpretable AI in sensitive domains.

Federated learning effectiveness and privacy preservation

The federated learning implementation successfully demonstrates collaborative threat detection across 8 simulated healthcare institutions while maintaining strict privacy guarantees. The remarkably low variance across clients (≤ 0.0005) despite heterogeneous data distributions indicates that our federated averaging

Model component	Individual accuracy	Parameters	Training time	Memory usage
RelationNetwork	99.09%	1,343,878	Low	0.4 GB
MAML	99.78%	715,269	Medium	0.3 GB
CrossTransformer	93.59%	25,794,566	High	1.2 GB
FEAT	99.76%	4,397,829	Medium	0.6 GB
Ensemble Fusion Results				
Majority Vote	99.9%	–	–	2.5 GB
Weighted Vote	99.8%	–	–	2.5 GB
Confidence Weighted	99.7%	–	–	2.5 GB

Table 12. Combined federated few-shot learning performance results.

Approach	Training	Inference	Memory	Parameters	Accuracy
FedMedSecure (Ours)	4.5 min	12 ms	2.5 GB	32.3M	99.9%
Nandanwar (AttackNet) ¹⁴	2.7 min	<5 ms	0.8 GB	5M	99.75%
Tawfik ¹⁶	–	<10 ms	2.0 GB	20M	99.7%
Shebl ³²	High	15 ms	1.5 GB	15M	99.98%
Jeremiah ²⁶	4.0 min	8 ms	1.0 GB	8M	99.7%
Alabbadi ²⁹	High	10 ms	1.8 GB	18M	97.66%

Table 13. Computational Complexity Comparison with Existing Approaches.

protocol effectively balances local and global knowledge, enabling institutions with different infrastructure types and threat exposure to benefit from collective intelligence.

The rapid convergence to 99.98% global accuracy within 10 rounds, as shown in Fig. 12, demonstrates the efficiency of our federated protocol. This fast convergence is particularly important for healthcare cybersecurity where threats evolve rapidly and defense mechanisms must adapt quickly to remain effective.

The differential privacy analysis with $(\epsilon, \delta) = (1.0, 10^{-5})$ achieves an optimal balance for healthcare applications. The privacy parameters are conservative enough to satisfy stringent healthcare regulations while permitting sufficient information sharing to maintain model utility. The minimal utility loss ($< 0.1\%$) demonstrates that formal privacy guarantees need not come at the expense of security effectiveness.

The 75% communication reduction through gradient compression addresses a critical practical concern for healthcare networks with bandwidth constraints and regulatory oversight of data transmission. This efficiency enables real-time collaborative threat detection without overwhelming network infrastructure or triggering regulatory compliance concerns about data movement.

Practical implications and deployment readiness

FedMedSecure’s exceptional performance metrics and comprehensive privacy guarantees position it as a practical solution for immediate deployment in real-world healthcare cybersecurity operations. The framework’s ability to detect novel attacks with minimal labeled samples directly addresses the dynamic nature of IoMT threat landscapes where new attack variants emerge daily and traditional signature-based detection fails.

The explainable AI capabilities facilitate seamless integration with existing healthcare security operations centers (SOCs), enabling analysts to understand, validate, and act upon automated threat detection decisions. This interpretability is not merely a technical feature but a regulatory necessity for healthcare deployment where “black box” AI systems face significant adoption barriers due to compliance requirements.

The federated architecture enables unprecedented collaborative threat intelligence sharing across healthcare institutions without exposing sensitive patient data or violating HIPAA regulations. This capability transforms cybersecurity from an institutional challenge to a community defense capability, potentially reducing successful attack rates across the entire healthcare ecosystem.

The framework’s computational efficiency and scalability metrics indicate readiness for production deployment. The ability to process 8.7M samples with sub-second inference times while maintaining perfect accuracy suggests that the system can handle real-world healthcare network traffic volumes without introducing latency that could impact critical medical operations. Communication overhead, while reduced 75% through compression, may challenge rural/developing-region institutions with <10 Mbps bandwidth (requiring 53 seconds per round). Future work should investigate adaptive compression and asynchronous protocols for bandwidth-constrained participants. Regarding blockchain integration mentioned in Section 5.7, we clarify this represents future research rather than current implementation. The current system achieves security through differential privacy and secure aggregation without blockchain dependency, avoiding blockchain latency (0.3-15 seconds per transaction) and throughput constraints (10-100 TPS) that would limit real-time federated learning. Future blockchain integration for audit trails would require consortium blockchains with off-chain model storage (on-chain cryptographic hashes only) to maintain acceptable latency while providing immutable records for regulatory compliance. Resource Constraints for IoMT Device Deployment Our framework was developed

and evaluated on cloud infrastructure (Google Colab Pro+ with NVIDIA V100 GPUs). However, practical IoMT security requires deployment consideration for resource-constrained medical devices. Table 12 shows ensemble memory requirements (2.5 GB) exceed typical IoMT device capabilities (medical sensors, wearables with 512 MB - 2 GB RAM). Our current architecture is suitable for:

- **Edge gateway deployment:** Hospital edge servers (16-64 GB RAM, 4-8 CPU cores) can run full ensemble for network traffic analysis, protecting multiple downstream IoMT devices
- **Institutional server deployment:** Healthcare data center servers aggregate and analyze traffic from all institutional IoMT devices
- **Federated participants:** Hospital IT infrastructure acts as federated clients, not individual IoMT devices

Direct on-device deployment requires model compression strategies: knowledge distillation could reduce model size by 70-80% (from 2.5 GB to 500-750 MB), quantization (INT8) provides 4× memory reduction, and pruning eliminates 40-60% of parameters with <2% accuracy loss. Energy consumption analysis is needed: our current training (10 rounds, 4.5 minutes) consumes approximately 0.15 kWh on V100 GPU, acceptable for server deployment but requiring optimization for battery-powered edge devices. Future work should evaluate lightweight model variants (MobileNet-style architectures, <100 MB) suitable for resource-constrained IoMT gateways while maintaining >95% detection accuracy.

Limitations and research challenges

Despite exceptional experimental performance, several limitations warrant careful consideration and future research attention. First, our evaluation relies entirely on simulated federated environments rather than actual healthcare network deployments with genuine institutional privacy constraints, regulatory oversight, and network heterogeneity. The transition from simulation to real-world deployment may reveal challenges not captured in our experimental setup.

Second, the semantic clustering approach showed mixed results with theoretical optimality not translating to improved few-shot performance. This finding suggests that our understanding of the relationship between information-theoretic measures and few-shot learning effectiveness in cybersecurity contexts requires deeper investigation. Future work should explore alternative clustering strategies and develop new metrics that better predict few-shot learning performance.

Third, our evaluation focuses exclusively on network-level attacks captured in the CICIoMT2024 dataset. The framework's effectiveness against application-level threats, insider attacks, and emerging IoMT-specific vulnerabilities such as firmware manipulation, medical device hijacking, and sensor spoofing remains untested. Healthcare cybersecurity extends beyond network traffic analysis to encompass device integrity, user behavior, and application security.

Adversarial robustness and byzantine fault tolerance Additionally, the framework's performance under adversarial conditions where attackers specifically target federated learning systems needs comprehensive evaluation. Our current federated aggregation uses simple weighted averaging, which assumes all participating healthcare institutions are honest-but-curious. This represents a significant limitation when facing adversarial scenarios where malicious nodes may send corrupted model updates. While our differential privacy mechanism ($\epsilon = 1.0$, $\delta = 10^{-5}$) provides theoretical privacy guarantees and gradient clipping (norm ≤ 1.0) bounds individual client influence, we have not conducted explicit adversarial experiments against:

- *Byzantine attacks:* Where malicious clients send arbitrarily corrupted model updates. Our current aggregation is vulnerable to such attacks. Future work should integrate robust aggregation methods such as Krum⁶¹, Trimmed Mean, or Median aggregation that can tolerate up to 33% Byzantine clients.
- *Model poisoning:* Where adversaries craft subtle model updates that degrade global model performance on specific attack types while maintaining normal accuracy on benign traffic. Defense mechanisms like anomaly detection on gradient distributions and cosine similarity filtering between client updates should be investigated.
- *Inference attacks:* Despite differential privacy, membership inference and model inversion attacks may still extract sensitive information from model parameters. Stronger privacy budgets ($\epsilon < 1.0$) or federated learning with secure multi-party computation provide enhanced protection but require privacy-utility trade-off analysis.

Formal security proofs under Byzantine threat models and empirical adversarial robustness evaluation represent critical future work for production healthcare deployment where attackers may specifically target the federated learning protocol.

Broader impact and future research directions

FedMedSecure's success demonstrates the transformative potential of federated few-shot learning for addressing critical cybersecurity challenges beyond healthcare. The framework's principles-privacy-preserving collaboration, rapid adaptation to novel threats, and explainable decision-making-are directly applicable to other privacy-sensitive domains including financial services, critical infrastructure protection, government networks, and industrial control systems.

The research opens several promising future directions: (1) integration with blockchain technologies for enhanced trust, transparency, and audit trails in federated learning; (2) development of continual learning capabilities that adapt to gradually evolving threat landscapes without catastrophic forgetting; (3) extension

to multi-modal threat detection incorporating device behavior analysis, user activity monitoring, and network traffic analysis for comprehensive security coverage.

Investigation of adversarial robustness represents a critical research priority. As federated learning systems become more widespread, attackers will develop sophisticated strategies to poison local models, manipulate global aggregation, or exploit the federated training process itself. Developing robust defenses against these meta-attacks while maintaining privacy guarantees presents significant research challenges.

The framework's explainable AI capabilities create opportunities for automated threat report generation, intelligent incident response workflows, and adaptive security policy recommendation systems. Future work should explore how XAI insights can drive autonomous security orchestration and reduce the burden on human security analysts.

Finally, the development of formal verification techniques for federated few-shot learning systems represents an important theoretical challenge. Healthcare applications require formal guarantees about system behavior, privacy preservation, and security effectiveness that extend beyond empirical evaluation to mathematical proofs of correctness and robustness.

The convergence of federated learning, few-shot adaptation, and explainable AI in FedMedSecure represents a significant step toward trustworthy, collaborative, and adaptive cybersecurity systems. As healthcare becomes increasingly digitized and interconnected, such frameworks will be essential for protecting patient safety and institutional integrity in the face of evolving cyber threats.

Conclusion

This paper introduced FedMedSecure, a novel federated few-shot learning framework for IoMT cybersecurity that successfully integrates privacy-preserving collaborative learning, rapid threat adaptation, and explainable AI. Our multi-model ensemble architecture combines CrossTransformer with learnable attack signature queries, FEAT, RelationNetwork with adaptive prototypes, and regularized MAML to achieve superior threat detection while maintaining formal differential privacy guarantees.

Extensive evaluation on the CICIoMT2024 dataset containing 8.7 million samples demonstrates exceptional performance: 99.99% accuracy in standard supervised learning, 99.7–99.8% accuracy in few-shot scenarios with as few as 5 shots per class, and 99.98% global accuracy in federated settings across 8 institutions. The framework achieves 75% communication reduction while preserving $(\epsilon, \delta) = (1.0, 10^{-5})$ differential privacy. Counterintuitively, the original 19-class attack taxonomy outperformed our theoretically optimized 5-class semantic clustering in few-shot learning, revealing important insights about class granularity and meta-learning effectiveness.

Our multi-level XAI framework provides comprehensive interpretability across feature, attention, and prototype levels, with SHAP analysis revealing packet timing features and protocol-specific attributes as primary attack discriminators. The confidence-weighted ensemble fusion mechanism automatically adapts to individual model performance, with RelationNetwork contributions reduced from 19% to 0.5% while balancing CrossTransformer (57.5%) and FEAT (42.0%) contributions.

FedMedSecure enables collaborative threat detection across healthcare institutions without compromising patient privacy, transforming cybersecurity from an institutional challenge to a community defense capability. The framework's combination of federated learning, few-shot adaptation, and explainable AI establishes a new paradigm for trustworthy cybersecurity in sensitive domains, with immediate applicability to financial services, critical infrastructure, and government networks. Future work will focus on real-world deployment validation, adversarial robustness, and multi-modal threat detection integration.

Future work

Future research directions include: (1) validation in real healthcare network deployments with genuine institutional constraints; (2) blockchain integration for audit trails (noting this is future work, not current implementation—current system uses differential privacy and secure aggregation for security without blockchain latency/throughput constraints); (3) investigation of adversarial robustness against sophisticated attacks targeting federated learning protocols, and extension to multi-modal threat detection incorporating device behavior and user activity analysis. Additional priorities include integration with blockchain technologies for enhanced trust and transparency, development of continual learning capabilities for evolving threat landscapes, and formal verification techniques for federated few-shot learning systems in critical healthcare applications.

Data availability

This study utilized two publicly available benchmark datasets: the CICIoMT2024 dataset4 (<https://www.unb.ca/cic/datasets/iomt-dataset-2024.html>), containing 8.7 million IoMT network traffic samples across 19 attack categories, and the CICIDS2017 dataset60 (<https://www.unb.ca/cic/datasets/ids-2017.html>), containing 2.8 million general IoT samples across 14 attack categories. Both datasets are freely accessible from the Canadian Institute for Cybersecurity at the University of New Brunswick.

Received: 24 July 2025; Accepted: 17 October 2025

Published online: 14 November 2025

References

1. Niu, Q. et al. Toward the Internet of Medical Things: Architecture, trends and challenges. *Math. Biosci. Eng.* **21**, 650–678. <https://doi.org/10.3934/mbe.2024028> (2024).

2. Ewoh, P. & Vartiainen, T. Vulnerability to cyberattacks and sociotechnical solutions for health care systems: Systematic review. *J. Med. Internet Res.* **26**, e46904. <https://doi.org/10.2196/46904> (2024).
3. Ahmed, S. F. et al. Insights into Internet of Medical Things (IoMT): Data fusion, security issues and potential solutions. *Information Fusion* **102**, 102060. <https://doi.org/10.1016/j.inffus.2023.102060> (2024).
4. Dadkhah, S. et al. CICIoMT2024: Attack vectors in healthcare devices—a multi-protocol dataset for assessing IoMT device security. *Internet of Things* **28**, 101321. <https://doi.org/10.1016/j.iot.2024.101321> (2024).
5. Thabit, F. & Can, O. Internet of Medical Things (IoMT) security: A comprehensive review. *Comput. Commun.* **218**, 48–66. <https://doi.org/10.1016/j.comcom.2024.01.020> (2024).
6. Teo, Z. L. et al. Federated machine learning in healthcare: A systematic review on clinical applications and technical architecture. *Cell Rep. Med.* **5**, 101419. <https://doi.org/10.1016/j.xcrm.2024.101419> (2024).
7. Rahman, A. et al. Federated learning-based AI approaches in smart healthcare: concepts, taxonomies, challenges and open issues. *Clust. Comput.* **26**, 2271–2311. <https://doi.org/10.1007/s10586-022-03658-4> (2023).
8. Nandanwar, H. & Katarya, R. A secure and privacy-preserving IDS for IoT networks using hybrid blockchain and federated learning. In *Proceedings of International Conference on Next-Generation Communication and Computing (NGCCOM 2024)*, vol. 1305 of *Lecture Notes in Networks and Systems*, 207–219. https://doi.org/10.1007/978-981-96-3725-6_18 (Springer, Singapore, 2025).
9. Nandanwar, H. & Katarya, R. Privacy-preserving data sharing in blockchain-enabled IoT healthcare management system. *Comput. J.* <https://doi.org/10.1093/comjnl/bxaf065> (2025).
10. Nandanwar, H. & Katarya, R. Optimized intrusion detection and secure data management in IoT networks using GAO-XGBoost and ECC-integrated blockchain framework. *Knowl. Inf. Syst.* <https://doi.org/10.1007/s10115-025-02513-3> (2025).
11. Khatun, M. A., Memon, S. F., Eising, C. & Dhirani, L. L. Machine learning for healthcare-IoT security: A review and risk mitigation. *IEEE Access* **11**, 145869–145896. <https://doi.org/10.1109/ACCESS.2023.3346320> (2023).
12. Liu, C. et al. Overcoming data limitations: A few-shot specific emitter identification method using self-supervised learning and adversarial augmentation. *IEEE Trans. Inf. Forensics Secur.* **19**, 500–513. <https://doi.org/10.1109/TIFS.2024.3427361> (2024).
13. Wang, R. et al. A lightweight model design approach for few-shot malicious traffic classification. *Sci. Rep.* **14**, 24710. <https://doi.org/10.1038/s41598-024-73342-7> (2024).
14. Nandanwar, H. & Katarya, R. Deep learning enabled intrusion detection system for industrial iot environment. *Expert Syst. Appl.* **249**, 123808. <https://doi.org/10.1016/j.eswa.2024.123808> (2024).
15. Zhang, H., Li, J. L., Liu, X. M. & Dong, C. Multi-dimensional feature fusion and stacking ensemble mechanism for network intrusion detection. *Futur. Gener. Comput. Syst.* **122**, 130–143. <https://doi.org/10.1016/j.future.2021.03.024> (2021).
16. Tawfik, M. Optimized intrusion detection in IoT and fog computing using ensemble learning and advanced feature selection. *PLoS ONE* **19**, e0304082. <https://doi.org/10.1371/journal.pone.0304082> (2024).
17. Arreche, O., Guntur, T. R., Roberts, J. W. & Abdallah, M. E-XAI: Evaluating black-box explainable AI frameworks for network intrusion detection. *IEEE Access* **12**, 23954–23988. <https://doi.org/10.1109/ACCESS.2024.3365140> (2024).
18. Mohale, V. Z. & Obagbuwa, I. C. Evaluating machine learning-based intrusion detection systems with explainable AI: enhancing transparency and interpretability. *Front. Comput. Sci.* <https://doi.org/10.3389/fcomp.2025.1476721> (2025).
19. Brauneck, A. et al. Federated machine learning, privacy-enhancing technologies, and data protection laws in medical research: Scoping review. *J. Med. Internet Res.* **25**, e41588. <https://doi.org/10.2196/41588> (2023).
20. Shen, A., Francisco, L., Sen, S. & Tewari, A. Exploring the relationship between privacy and utility in mobile health: a simulation of federated learning, differential privacy, and external attacks. *J. Med. Internet Res.* **25**, e43664. <https://doi.org/10.2196/43664> (2023).
21. Bhushan, B. et al. Towards a secure and sustainable internet of medical things (iomt): Requirements, design challenges, security techniques, and future trends. *Sustainability* **15**, 6177. <https://doi.org/10.3390/su15076177> (2023).
22. Tauqeer, H. et al. Enhanced intrusion detection system for Internet of Medical Things using meta-learning. *Sensors* **23**, 5456. <https://doi.org/10.3390/s23125456> (2023).
23. Yan, Y. et al. Meta learning-based few-shot intrusion detection for 5G-enabled industrial internet. *Complex Intell. Syst.* **10**, 4589–4608. <https://doi.org/10.1007/s40747-024-01388-1> (2024).
24. Qi, T. et al. Collaborative machine learning with differential privacy for healthcare. *eBioMedicine* **101**, 105006. <https://doi.org/10.1016/j.ebiom.2024.105006> (2024).
25. Misbah, A., Sebban, A. & Hafidi, I. Securing Internet of Medical Things: An advanced federated learning approach. *Int. J. Adv. Comput. Sci. Appl.* **16** (2025).
26. Jeremiah, S. R., El Azzaoui, A., Gritzalis, S. & Park, J. H. Multi-view learning and model fusion framework for threat detection in multi-protocol IoMT networks. *Information Fusion* <https://doi.org/10.1016/j.inffus.2025.103435> (2025).
27. Sharma, N. & Shambharkar, P. G. Multi-attention DeepCRNN: an efficient and explainable intrusion detection framework for Internet of Medical Things environments. *Knowledge and Information Systems* 1–67 (2025).
28. Kavkas, N. C. & Yildiz, K. Enhancing IoMT security with deep learning based approach for medical IoT threat detection. In *2025 13th International Symposium on Digital Forensics and Security (ISDFS)*, 1–5. <https://doi.org/10.1109/ISDFS65363.2025.11012062> (IEEE, 2025).
29. Alabbadi, A. & Bajaber, F. X-FuseRLSTM: A cross-domain explainable intrusion detection framework in IoT using the attention-guided dual-path feature fusion and residual LSTM. *Sensors* **25**, 3693. <https://doi.org/10.3390/s25123693> (2025).
30. Akar, G., Sahmoud, S., Onat, M., Cavusoglu, Ü. & Malondo, E. L2D2: A novel LSTM model for multi-class intrusion detection systems in the era of IoMT. *IEEE Access* <https://doi.org/10.1109/ACCESS.2025.3526883> (2025).
31. Hernandez-Jaimes, M. L., Martinez-Cruz, A., Ramirez-Gutiérrez, K. A. & Morales-Reyes, A. Network traffic inspection to enhance anomaly detection in the Internet of Things using attention-driven deep learning. *Integration* **103**, 102398. <https://doi.org/10.1016/j.vlsi.2025.102398> (2025).
32. Shebl, A., Elsedimy, E. I., Ismail, A., Salama, A. A. & Herajy, M. DCNN: A novel binary and multi-class network intrusion detection model via deep convolutional neural network. *EURASIP J. Inf. Secur.* **2024**, 36 (2024).
33. Alturki, B. & Alsulami, A. A. Semi-supervised learning with entropy filtering for intrusion detection in asymmetrical IoT systems. *Symmetry* **17**, 973. <https://doi.org/10.3390/sym17060973> (2025).
34. Kharoubi, K., Cherbal, S., Mechta, D. & Gawanmeh, A. Network intrusion detection system using convolutional neural networks: Nids-dl-cnn for iot security. *Clust. Comput.* **28**, 219. <https://doi.org/10.1007/s10586-024-04904-7> (2025).
35. Doménech, J., León, O., Siddiqui, M. S. & Pegueroles, J. Evaluating and enhancing intrusion detection systems in IoMT: The importance of domain-specific datasets. *Internet of Things* <https://doi.org/10.1016/j.iot.2025.101631> (2025).
36. Rehman, M. U., Kalakoti, R. & Bahşi, H. Comprehensive feature selection for machine learning-based intrusion detection in healthcare IoMT networks. In *Proceedings of the 11th International Conference on Information Systems Security and Privacy - Volume 2: ICISSP*, 248–259. <https://doi.org/10.5220/0013313600003899>. INSTICC (SciTePress, 2025).
37. Hernandez-Jaimes, M. L., Martinez-Cruz, A., Ramirez-Gutiérrez, K. A. & Guevara-Martínez, E. Enhancing machine learning approach based on nilsimsa fingerprinting for ransomware detection in IoMT. *IEEE Access* <https://doi.org/10.1109/ACCESS.2024.3480889> (2024).
38. Kumar, H., Kumar, H., Harish, Nandanwar, H. & Katarya, R. Enhancing security and scalability of IoMT systems using blockchain: Addressing key challenges and limitations. In *Proceedings of the 6th International Conference on Deep Learning, Artificial Intelligence and Robotics (ICDLAIR 2024)*, 191–202. https://doi.org/10.2991/978-94-6463-740-3_17 (Atlantis Press, 2025).

39. McMahan, B., Moore, E., Ramage, D., Hampson, S. & Aguera y Arcas, B. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics* (2017).
40. Wei, K. et al. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE Trans. Inf. Forensics Secur.* <https://doi.org/10.1109/TIFS.2020.2988575> (2020).
41. Finn, C., Abbeel, P. & Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, <https://doi.org/10.5555/3305381.3305498> (2017).
42. Snell, J., Swersky, K. & Zemel, R. S. Prototypical networks for few-shot learning. *Adv. Neural Inf. Process. Syst.* (2017).
43. Dwork, C. Differential privacy. In *33rd International Colloquium on Automata, Languages and Programming*, https://doi.org/10.1007/11787006_1 (2006).
44. Ye, H.-J., Hu, H., Zhan, D.-C. & Sha, F. Few-shot learning via embedding adaptation with set-to-set functions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020).
45. Li, T. et al. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems* <https://doi.org/10.48550/arXiv.1812.06127> (2020).
46. Dadkhah, S. et al. CICIoMT2024: A benchmark dataset for multi-protocol security assessment in IoMT. *Internet of Things* **28**, 101351. <https://doi.org/10.1016/j.iot.2024.101351> (2024).
47. Kale, R. & Thing, V. L. L. Few-shot weakly-supervised cybersecurity anomaly detection. *Comput. Secur.* <https://doi.org/10.1016/j.cose.2023.103194> (2023).
48. Breiman, L. Random forests. *Mach. Learn.* <https://doi.org/10.1023/A:1010933404324> (2001).
49. Freund, Y. & Schapire, R. E. Experiments with a new boosting algorithm. In *Machine Learning: Proceedings of the Thirteenth International Conference* (1996).
50. Al-rimy, B. A. S., Maarof, M. A. & Shaid, S. Z. M. Ensemble learning for intrusion detection systems: A systematic mapping study and cross-benchmark evaluation. *Comput. Commun.* <https://doi.org/10.1016/j.comcom.2021.08.011> (2021).
51. Sung, F. et al. Learning to compare: Relation network for few-shot learning. *IEEE/CVF Conf. Comput. Vis. Pattern Recogn.* <https://doi.org/10.1109/CVPR.2018.00131> (2018).
52. Vaswani, A. et al. Attention is all you need. In *Advances in Neural Information Processing Systems*, vol. 30 (2017).
53. Long, Z., Yan, H., Shen, G., Li, J. & Yu, S. A transformer-based network intrusion detection approach for cloud security. *J. Cloud Comput.* <https://doi.org/10.1186/s13677-023-00574-9> (2024).
54. Laghrissi, F., Douzi, S., Douzi, K. & Hssina, B. Ids-attention: An efficient algorithm for intrusion detection systems using attention mechanism. *J. Big Data* <https://doi.org/10.1186/s40537-021-00544-5> (2021).
55. Chahal, A., Gupta, A. & Chandra, P. Design of a federated ensemble model for intrusion detection in distributed iiot networks for enhancing cybersecurity. *Comput. Electr. Eng.* <https://doi.org/10.1016/j.compeleceng.2024.109724> (2025).
56. Rieke, N. et al. The future of digital health with federated learning. *npj Digital Medicine* <https://doi.org/10.1038/s41746-020-00323-1> (2020).
57. Lundberg, S. M. & Lee, S.-I. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems* (2017).
58. Ribeiro, M. T., Singh, S. & Guestrin, C. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, <https://doi.org/10.1145/2939672.2939778> (2016).
59. Gaspar, D., Silva, P. & Silva, C. Explainable ai for intrusion detection systems: Lime and shap applicability on multi-layer perceptron. *IEEE Access* <https://doi.org/10.1109/ACCESS.2024.3368377> (2024).
60. Sharafaldin, I., Lashkari, A. H. & Ghorbani, A. A. Toward generating a new intrusion detection dataset and intrusion traffic characterization. *ICISSp* 108–116 (2018).
61. Blanchard, P., Mhamdi, E. M. E., Guerraoui, R. & Stainer, J. Machine learning with adversaries: Byzantine tolerant gradient descent. *Adv. Neural Inf. Process. Syst. (NeurIPS)* **30**, 119–129 (2017).

Author contributions

M.T. conceived the study, developed the methodology including few-shot learning algorithms and federated learning protocols, and conducted experiments. A.A.A. contributed to privacy analysis and differential privacy implementation. H.M.N. provided technical guidance on experimental design and assisted with validation. A.H.A. contributed to cybersecurity analysis and threat modeling. I.S.F. contributed to the explainable AI framework development and manuscript review. All authors reviewed and approved the final manuscript.

Declarations

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to M.T.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025