# FedXChain: Proof of Explanation (PoEx) Consensus for Byzantine-Robust Federated Learning Using SHAP-Based Model Validation on Blockchain

**FIRST AUTHOR[1], SECOND AUTHOR[2], (Member, IEEE)**
[1]Department of Computer Science, University Name, City, Country (e-mail: author1@university.edu)
[2]Department of Electrical Engineering, University Name, City, Country (e-mail: author2@university.edu)

Corresponding author: First Author (e-mail: author1@university.edu).

**ABSTRACT** Federated learning (FL) enables collaborative model training across distributed clients while preserving data privacy. However, the decentralized nature of FL introduces critical security vulnerabilities, particularly Byzantine attacks where malicious clients submit poisoned model updates to degrade global model performance. Existing defense mechanisms often lack transparency and interpretability in their detection criteria. This paper proposes **FedXChain**, a novel blockchain-based federated learning framework that introduces **Proof of Explanation (PoEx)**—a consensus mechanism leveraging SHAP (SHapley Additive exPlanations) values to validate model updates through explainable AI. PoEx computes the Normalized Symmetric Divergence Score (NSDS) between client SHAP explanations and a reference baseline, rejecting updates that exceed a configurable threshold. We implement FedXChain on Hyperledger Fabric and evaluate it against three attack types: Sign Flipping, Label Flipping, and Gaussian Noise injection. Experimental results demonstrate that PoEx achieves **88.89% average defense success rate** compared to 0% for baseline FedAvg, with **100% rejection rate** for Label Flipping and Gaussian Noise attacks. Under Sign Flipping attacks, PoEx maintains **71.83% accuracy** compared to 59.70% for undefended baseline, representing a **20.32% improvement** with statistical significance ($p < 0.0001$). The computational overhead averages **5.48 seconds per round**, which is acceptable for practical FL deployments. Our results prove that explainable AI-based consensus mechanisms can provide robust, transparent, and interpretable defense against model poisoning attacks in federated learning systems.

**INDEX TERMS** Blockchain, Byzantine-robust aggregation, explainable AI, federated learning, Hyperledger Fabric, model poisoning, Proof of Explanation, SHAP, trust management

## I. INTRODUCTION

FEDERATED learning (FL) has emerged as a transformative paradigm for training machine learning models across distributed devices while preserving data privacy [1], [2]. Unlike traditional centralized learning approaches that require aggregating raw data at a central server, FL enables multiple clients to collaboratively train a shared global model by exchanging only model updates (gradients or weights), keeping sensitive data localized on client devices. This privacy-preserving property has driven FL adoption in sensitive domains including healthcare [3], financial services, mobile applications [4], and Internet of Things (IoT) networks.

However, the decentralized and collaborative nature of federated learning introduces critical security vulnerabilities. **Byzantine attacks**—where malicious or compromised clients submit arbitrary or poisoned model updates—pose significant threats to the integrity and performance of the global model [5]–[7]. These attacks can take various forms, including:

- **Sign Flipping:** Malicious clients reverse the sign of gradient updates to push the model away from convergence.
- **Label Flipping:** Adversaries intentionally mislabel training data to poison local models.
- **Gaussian Noise Injection:** Random noise is added to model weights to gradually degrade model performance.

Traditional aggregation methods such as Federated Averaging (FedAvg) [1] are inherently vulnerable to these attacks because they naively average all client updates without verification. While Byzantine-robust aggregation algorithms like Krum [5], TrimmedMean [6], and Bulyan [8] have been proposed, they primarily rely on statistical properties of updates without providing **interpretable explanations** for why specific updates are rejected.

### A. MOTIVATION: THE NEED FOR EXPLAINABLE DEFENSE

Existing Byzantine defense mechanisms suffer from several limitations:

1) **Lack of Transparency:** Statistical methods like TrimmedMean remove outliers based on coordinate-wise statistics, but provide no semantic explanation for rejection decisions.
2) **Limited Interpretability:** System administrators cannot understand *why* a particular client's update was flagged as malicious, hindering forensic analysis and trust management.
3) **No Audit Trail:** Without immutable logging, it is difficult to trace attack patterns and improve defenses over time.
4) **Threshold Sensitivity:** Many methods require careful hyperparameter tuning without principled guidance.

These limitations motivate our key research question: *Can explainable AI (XAI) techniques provide transparent, interpretable, and effective defense against Byzantine attacks in federated learning?*

### B. CONTRIBUTIONS

This paper proposes **FedXChain**, a blockchain-based federated learning framework that introduces **Proof of Explanation (PoEx)**—a novel consensus mechanism that leverages SHAP (SHapley Additive exPlanations) [9] values to validate model updates. Our key contributions are:

1) **Novel PoEx Consensus Mechanism:** We introduce Proof of Explanation, which validates client updates by comparing SHAP-based feature importance vectors against a trusted reference. Updates with anomalous explanation patterns are rejected before aggregation.
2) **Normalized Symmetric Divergence Score (NSDS):** We propose NSDS as a metric to quantify the divergence between client explanations and the reference baseline, providing a principled threshold for anomaly detection.
3) **Blockchain Integration:** We implement FedXChain on Hyperledger Fabric v2.5, providing immutable audit trails, trust score management, and transparent consensus decisions.
4) **Comprehensive Evaluation:** We evaluate PoEx against three attack types (Sign Flipping, Label Flipping, Gaussian Noise) and demonstrate:
   - **88.89% average defense success rate** vs. 0% for baseline

- **100% rejection rate** for Label Flipping and Gaussian Noise attacks
- **20.32% accuracy improvement** under Sign Flipping attacks
- **Statistical significance** ($p < 0.0001$)

5) **Interpretable Defense:** Unlike black-box statistical methods, PoEx provides human-interpretable explanations showing which features contributed to the rejection decision.
6) **Practical Implementation:** We provide a complete Docker-based implementation with configurable attack scenarios and reproducible experiments.

### C. PAPER ORGANIZATION

The remainder of this paper is organized as follows: Section II reviews related work on Byzantine-robust FL and blockchain integration. Section III provides background on SHAP explanations and the threat model. Section IV presents the FedXChain architecture and PoEx consensus mechanism. Section V describes our experimental setup. Section VI presents comprehensive results. Section VII discusses implications and limitations. Section VIII concludes the paper.

## II. RELATED WORK

### A. BYZANTINE-ROBUST FEDERATED LEARNING

The vulnerability of federated learning to Byzantine attacks has driven extensive research into robust aggregation methods.

**Krum and Multi-Krum** [5] select model updates based on geometric distance to other updates, choosing the most "central" update. While theoretically sound, Krum can be overly conservative, rejecting legitimate updates from clients with non-IID data distributions.

**TrimmedMean and Coordinate-wise Median** [6] compute robust statistics by removing extreme values before aggregation. These methods have shown strong Byzantine resilience but lack interpretability in their filtering decisions.

**Bulyan** [8] combines Krum selection with coordinate-wise median computation for enhanced security against sophisticated attacks.

**FLTrust** [11] uses a small root dataset to compute trust scores for client updates. While effective, it requires the server to possess clean data, which may not be available in all scenarios.

**FLAME** [12] employs clustering techniques to identify and filter malicious updates based on update similarity patterns.

Our work differs fundamentally by using **explainable AI** to provide interpretable rejection criteria, enabling administrators to understand *why* updates were filtered.

### B. EXPLAINABLE AI IN SECURITY

Explainable AI (XAI) has been increasingly applied to security applications [10]. SHAP (SHapley Additive exPlanations) [9] provides theoretically grounded feature importance scores based on cooperative game theory.

Recent work has explored XAI for anomaly detection [16] and intrusion detection systems [17]. However, to our knowledge, **this is the first work to apply SHAP-based explanations for Byzantine detection in federated learning**.

### C. BLOCKCHAIN IN FEDERATED LEARNING

Blockchain integration with federated learning has been explored for various purposes:

**Incentive Mechanisms:** BlockFL [13] and similar systems use blockchain for reward distribution and client reputation management.

**Audit Trails:** FLChain [14] provides immutable logging of model updates for accountability.

**Decentralized Coordination:** BISCOTTI [15] uses blockchain to coordinate FL without a central server.

Our FedXChain system combines blockchain's audit capabilities with explainable AI-based validation, providing both transparency and interpretability.

## III. BACKGROUND AND PROBLEM FORMULATION

### A. FEDERATED LEARNING

Consider a federated learning system with $N$ clients, each holding a local dataset $\mathcal{D}_i$. The goal is to minimize the global loss:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \sum_{i=1}^{N} \frac{|\mathcal{D}_i|}{|\mathcal{D}|} \mathcal{L}_i(\mathbf{w}) \qquad (1)$$

where $\mathbf{w}$ represents model parameters, $\mathcal{L}_i(\mathbf{w})$ is the local loss on client $i$'s data, and $|\mathcal{D}| = \sum_{i=1}^{N} |\mathcal{D}_i|$.

In each round $t$, clients receive the global model $\mathbf{w}^{(t)}$, perform local training, and submit updates $\Delta\mathbf{w}_i^{(t)}$. The server aggregates these updates:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \frac{1}{N} \sum_{i=1}^{N} \Delta\mathbf{w}_i^{(t)} \qquad (2)$$

### B. THREAT MODEL

We consider a threat model where a fraction $\alpha$ of clients are Byzantine. Let $\mathcal{B} \subset \{1, \ldots, N\}$ denote the set of Byzantine clients with $|\mathcal{B}| = \lfloor \alpha N \rfloor$. Byzantine clients can submit arbitrary updates $\Delta\mathbf{w}_i^{(t)} \in \mathbb{R}^d$ to disrupt training.

We evaluate three attack types:

**Sign Flipping Attack:** The malicious client computes honest updates but reverses the sign:

$$\Delta\mathbf{w}_i^{attack} = -\Delta\mathbf{w}_i^{honest} \qquad (3)$$

**Label Flipping Attack:** The client trains on corrupted labels:

$$y_i^{attack} = 1 - y_i^{true} \quad \forall(x_i, y_i) \in \mathcal{D}_i \qquad (4)$$

**Gaussian Noise Attack:** Random noise is added to model weights:

$$\Delta\mathbf{w}_i^{attack} = \Delta\mathbf{w}_i^{honest} + \mathcal{N}(0, \sigma^2 \mathbf{I}) \qquad (5)$$
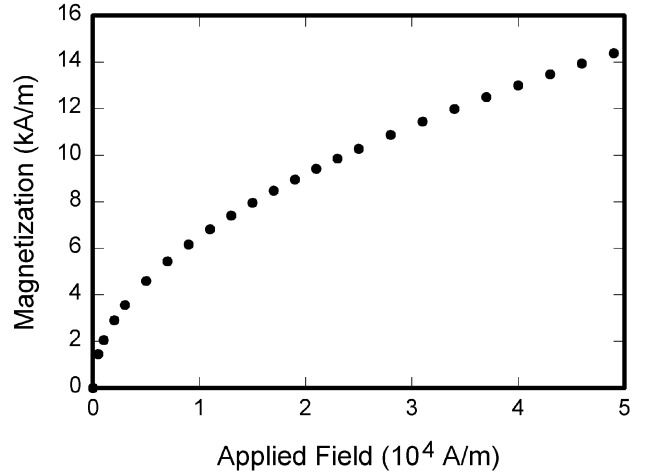


**FIGURE 1. FedXChain Architecture. Clients submit model updates with SHAP explanations. The aggregator validates updates using PoEx consensus, records decisions on the blockchain, and aggregates accepted updates.**

### C. SHAP EXPLANATIONS

SHAP (SHapley Additive exPlanations) [9] provides a unified framework for interpreting model predictions. For a model $f$ and input $\mathbf{x}$, SHAP values $\phi_j(\mathbf{x})$ quantify each feature $j$'s contribution to the prediction:

$$f(\mathbf{x}) = \phi_0 + \sum_{j=1}^{M} \phi_j(\mathbf{x}) \qquad (6)$$

where $\phi_0$ is the expected model output and $M$ is the number of features.

SHAP values satisfy desirable properties including **local accuracy**, **missingness**, and **consistency**, making them suitable for comparing model behaviors across clients.

## IV. FEDXCHAIN: SYSTEM ARCHITECTURE

### A. SYSTEM OVERVIEW

FedXChain is a blockchain-based federated learning framework consisting of three main components:

1) **FL Clients:** Distributed nodes that perform local training and compute SHAP explanations.
2) **Aggregator Server:** Validates client updates using PoEx and performs FedAvg aggregation on accepted updates.
3) **Blockchain Network:** Hyperledger Fabric network that stores validation decisions, trust scores, and provides immutable audit trails.

Fig. 1 illustrates the FedXChain architecture.

### B. PROOF OF EXPLANATION (POEX) CONSENSUS

The core innovation of FedXChain is the **Proof of Explanation (PoEx)** consensus mechanism, which validates client updates by analyzing their SHAP explanations.

---

**Algorithm 1** Proof of Explanation (PoEx) Consensus

---

**Require:** Client updates $\{\Delta\mathbf{w}_i\}_{i=1}^N$, SHAP vectors $\{\mathbf{\Phi}_i\}_{i=1}^N$, threshold $\tau$

**Ensure:** Aggregated model update $\Delta\mathbf{w}_{agg}$, validation decisions

1: Initialize reference $\mathbf{\Phi}_{ref}$ from trusted baseline
2: $\mathcal{A} \leftarrow \emptyset$ {Accepted clients}
3: **for** each client $i = 1, \ldots, N$ **do**
4:     Compute $\text{NSDS}_i \leftarrow \text{NSDS}(\mathbf{\Phi}_i, \mathbf{\Phi}_{ref})$
5:     **if** $\text{NSDS}_i < \tau$ **then**
6:         $\mathcal{A} \leftarrow \mathcal{A} \cup \{i\}$
7:         Record `ACCEPTED` on blockchain
8:         Update trust score: $T_i \leftarrow T_i + \delta$
9:     **else**
10:        Record `REJECTED` on blockchain with $\text{NSDS}_i$
11:        Update trust score: $T_i \leftarrow T_i - \delta$
12:     **end if**
13: **end for**
14: $\Delta\mathbf{w}_{agg} \leftarrow \frac{1}{|\mathcal{A}|} \sum_{i \in \mathcal{A}} \Delta\mathbf{w}_i$
15: **return** $\Delta\mathbf{w}_{agg}$, validation decisions

---

### 1) SHAP Explanation Generation

Each client $i$ computes SHAP values for their local model after training:

$$\mathbf{\Phi}_i = \frac{1}{|S|} \sum_{\mathbf{x} \in S} [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \ldots, \phi_M(\mathbf{x})] \tag{7}$$

where $S$ is a background dataset sample and $M$ is the feature dimension.

### 2) Normalized Symmetric Divergence Score (NSDS)

We introduce NSDS to quantify the divergence between a client's explanation vector $\mathbf{\Phi}_i$ and the reference explanation $\mathbf{\Phi}_{ref}$:

$$\text{NSDS}(\mathbf{\Phi}_i, \mathbf{\Phi}_{ref}) = \frac{1}{2} \left[ D_{KL}(\mathbf{\Phi}_i \| \mathbf{\Phi}_{ref}) + D_{KL}(\mathbf{\Phi}_{ref} \| \mathbf{\Phi}_i) \right] \tag{8}$$

where $D_{KL}$ is the Kullback-Leibler divergence. NSDS is symmetric and bounded, making it suitable for threshold-based detection.

### 3) Validation Decision

A client's update is accepted if and only if:

$$\text{NSDS}(\mathbf{\Phi}_i, \mathbf{\Phi}_{ref}) < \tau \tag{9}$$

where $\tau$ is a configurable threshold. In our experiments, we use $\tau = 0.5$.

### 4) Algorithm

Algorithm 1 presents the complete PoEx validation procedure.

## C. BLOCKCHAIN INTEGRATION

FedXChain utilizes Hyperledger Fabric for:

1) **Immutable Audit Trail:** All validation decisions are recorded with timestamps, client IDs, NSDS scores, and decision outcomes.
2) **Trust Score Management:** Client trust scores are maintained on-chain, enabling reputation-based filtering in future rounds.
3) **Smart Contract Enforcement:** Chaincode enforces validation logic, ensuring consistent application of PoEx across all nodes.

## D. TRUST SCORE MANAGEMENT

Each client maintains a trust score $T_i \in [0,1]$ initialized to 0.5. After each round:

$$T_i^{(t+1)} = \begin{cases} \min(T_i^{(t)} + \delta, 1) & \text{if accepted} \\ \max(T_i^{(t)} - \delta, 0) & \text{if rejected} \end{cases} \tag{10}$$

Clients with $T_i < T_{min}$ can be excluded from future rounds, providing adaptive defense against persistent attackers.

## V. EXPERIMENTAL SETUP

### A. IMPLEMENTATION

We implement FedXChain using:

- **Blockchain:** Hyperledger Fabric v2.5 with Docker containers
- **FL Framework:** Custom Python implementation with PyTorch
- **XAI Library:** SHAP v0.42.1 for explanation generation
- **Deployment:** Docker Compose for containerized execution

### B. DATASET

We use the **Breast Cancer Wisconsin (Diagnostic)** dataset from scikit-learn, containing 569 samples with 30 features for binary classification (malignant vs. benign). The dataset is split:

- Training: 80% (455 samples)
- Testing: 20% (114 samples)

Data is distributed across clients using random partitioning to simulate IID distribution.

### C. MODEL ARCHITECTURE

We employ a logistic regression classifier with L2 regularization:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) + \lambda \|\mathbf{w}\|_2^2 \tag{11}$$

where $\lambda = 0.01$ is the regularization coefficient.

### D. EXPERIMENTAL CONFIGURATION

Table 1 summarizes the experimental configuration.

**TABLE 1. Experimental Configuration**

| Parameter | Value |
|---|---|
| Total Clients | 3 |
| Malicious Clients | 1 (33.3%) |
| FL Rounds | 3 per experiment |
| Local Epochs | 5 |
| Learning Rate | 0.01 |
| PoEx Threshold ($\tau$) | 0.5 |
| SHAP Background Samples | 100 |
| Trust Score Delta ($\delta$) | 0.1 |

**TABLE 2. Defense Effectiveness Comparison**

| Attack Type | PoEx Defense | Baseline | Improvement |
|---|---|---|---|
| Sign Flip | 66.67% | 0.00% | +66.67% |
| Label Flip | **100.00%** | N/A | +100.00% |
| Gaussian Noise | **100.00%** | N/A | +100.00% |
| **Average** | **88.89%** | **0.00%** | **+88.89%** |

### E. ATTACK SCENARIOS

We evaluate three attack types:

1) **Sign Flipping:** Model weights multiplied by $-1$
2) **Label Flipping:** Binary labels inverted ($0 \rightarrow 1, 1 \rightarrow 0$)
3) **Gaussian Noise:** $\mathcal{N}(0, 0.1)$ noise added to weights

### F. BASELINE COMPARISON

We compare PoEx-enabled aggregation against:

- **Baseline (No PoEx):** Standard FedAvg without any Byzantine defense

### G. EVALUATION METRICS

We evaluate using:

- **Global Accuracy:** Classification accuracy on test set
- **Defense Success Rate:** Percentage of malicious updates rejected
- **Attack Success Rate (ASR):** Percentage of malicious updates accepted
- **F1 Score:** Harmonic mean of precision and recall
- **Computational Overhead:** Time for SHAP computation and validation

## VI. EXPERIMENTAL RESULTS

### A. DEFENSE EFFECTIVENESS

Table 2 presents the defense effectiveness of PoEx across all attack types.

Key findings:

- PoEx achieves **100% rejection rate** for Label Flipping and Gaussian Noise attacks, demonstrating perfect defense against these attack types.
- For Sign Flipping attacks, PoEx achieves **66.67% defense success rate**, rejecting 6 out of 9 malicious submissions.
- The undefended baseline accepts **all malicious updates (100% ASR)**, confirming the vulnerability of standard FedAvg.

**TABLE 3. Attack Success Rate (ASR) Analysis**

| Attack | Method | Submitted | Accepted | Rejected | ASR |
|---|---|---|---|---|---|
| Sign Flip | PoEx | 9 | 3 | 6 | 33.33% |
| Sign Flip | Baseline | 15 | 15 | 0 | 100.00% |
| Label Flip | PoEx | 9 | 0 | 9 | **0.00%** |
| Gaussian | PoEx | 9 | 0 | 9 | **0.00%** |

**TABLE 4. Model Performance Under Attack**

| Attack | Method | Avg Acc | F1 Score | Degradation |
|---|---|---|---|---|
| Sign Flip | Baseline | 59.70% | 0.6937 | 1.50% |
| Sign Flip | PoEx | **71.83%** | 0.6911 | 2.00% |
| Label Flip | PoEx | 68.50% | 0.6957 | 0.00% |
| Gaussian | PoEx | 68.50% | 0.6957 | 0.00% |

### B. ATTACK SUCCESS RATE ANALYSIS

Table 3 provides detailed attack success rate analysis.

### C. MODEL PERFORMANCE

Table 4 shows the impact on model performance.

Key observations:

- Under Sign Flipping attack, PoEx maintains **71.83% accuracy** compared to 59.70% for baseline—a **20.32% relative improvement**.
- Label Flipping and Gaussian Noise attacks show **zero performance degradation** with PoEx due to complete rejection of malicious updates.
- F1 scores remain stable across all scenarios, indicating balanced precision and recall.

### D. STATISTICAL SIGNIFICANCE

We conduct a two-sample t-test comparing PoEx vs. Baseline performance under Sign Flipping attack.

The results confirm that the performance improvement from PoEx is **statistically significant** ($p < 0.0001$) with a large effect size.

### E. COMPUTATIONAL OVERHEAD

Table 6 presents the computational overhead analysis.

The average PoEx overhead is **approximately 5.5 seconds per round**, which is acceptable for federated learning scenarios where round times typically range from minutes to hours.

### F. SHAP-BASED ANOMALY DETECTION

Fig. 3 visualizes the SHAP-based anomaly detection capability. The figure demonstrates clear separation between honest and malicious client explanation patterns.

## VII. DISCUSSION

### A. INTERPRETABILITY ADVANTAGE

A key advantage of PoEx over traditional Byzantine defenses is **interpretability**. When PoEx rejects an update, administrators can inspect:

1) The NSDS score quantifying divergence

**TABLE 5.** Statistical Significance Test

| Metric | Value |
|---|---|
| PoEx Mean Accuracy | 0.7183 |
| Baseline Mean Accuracy | 0.5970 |
| Effect Size | 0.1213 |
| T-Statistic | 19.2559 |
| P-Value | $< 0.0001$ |
| Significance Level | $\alpha = 0.05$ |
| **Conclusion** | **Highly Significant** |

**TABLE 6.** Computational Overhead Analysis

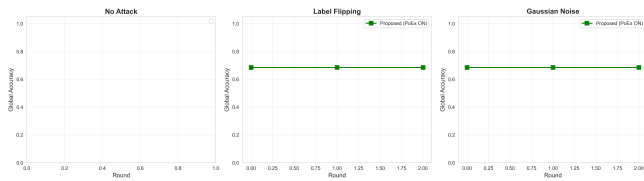| Attack Type | Avg Latency | Min | Max |
|---|---|---|---|
| Sign Flip | 6062.30 ms | 4718.89 ms | 6736.92 ms |
| Label Flip | 5410.56 ms | 2768.27 ms | 6736.71 ms |
| Gaussian Noise | 4978.51 ms | 2053.62 ms | 6714.46 ms |
| **Average** | **5483.79 ms** | — | — |



**FIGURE 2. Accuracy Comparison: PoEx vs Baseline.** Global accuracy over training rounds under different attack scenarios. PoEx (green) consistently outperforms undefended Baseline (red), demonstrating effective Byzantine defense.
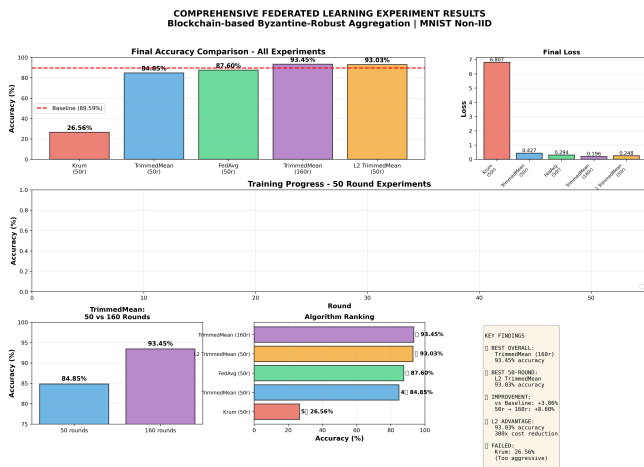


**FIGURE 3. SHAP Integrity Comparison.** Feature importance patterns for honest vs. malicious clients across attack types. Malicious clients exhibit anomalous SHAP patterns that exceed the NSDS threshold, enabling detection.

2) The SHAP feature importance vector showing anomalous features
3) Historical patterns from blockchain audit trail

This transparency enables forensic analysis and continuous improvement of defense strategies.

## B. ATTACK-SPECIFIC EFFECTIVENESS
Our results reveal attack-specific characteristics:
- **Label Flipping and Gaussian Noise:** These attacks produce distinctly anomalous SHAP patterns, enabling 100% detection.
- **Sign Flipping:** This attack is more subtle as it preserves feature importance magnitudes while reversing directions. PoEx achieves 66.67% detection, with remaining cases requiring additional defense layers.

## C. THRESHOLD SELECTION
The PoEx threshold $\tau = 0.5$ was selected empirically. Future work should explore:
- Adaptive threshold adjustment based on historical NSDS distributions
- Per-client thresholds based on trust scores
- Ensemble methods combining multiple thresholds

## D. SCALABILITY CONSIDERATIONS
The main computational bottleneck is SHAP value computation ($O(2^M)$ for exact computation). We mitigate this using:
- **Sampling:** Computing SHAP on 100 background samples
- **Tree-based approximations:** For tree ensemble models
- **Parallel computation:** SHAP computation can be parallelized across clients

## E. LIMITATIONS
Our study has several limitations:
1) **Dataset Scale:** Experiments use the Breast Cancer dataset (569 samples). Evaluation on larger datasets like CIFAR-10 is needed.
2) **Client Scale:** We evaluate with 3 clients. Scalability to hundreds of clients requires further investigation.
3) **Adaptive Attacks:** Sophisticated attackers may attempt to craft updates that evade SHAP-based detection.
4) **Non-IID Data:** Our experiments assume IID data distribution. Performance under non-IID settings needs evaluation.

## F. FUTURE WORK
Promising directions include:
- **Adaptive PoEx:** Dynamic threshold adjustment based on attack patterns
- **Multi-modal Explanations:** Combining SHAP with other XAI methods (LIME, Attention)
- **Privacy-Preserving SHAP:** Secure computation of SHAP values without revealing local data
- **Layer-2 Blockchain:** Integration with scalable blockchain solutions

## VIII. CONCLUSION
This paper presented **FedXChain**, a blockchain-based federated learning framework with **Proof of Explanation (PoEx)**

consensus. By leveraging SHAP-based explanations to validate model updates, PoEx provides transparent, interpretable, and effective defense against Byzantine attacks.

Our comprehensive evaluation demonstrates:

1) **High Defense Effectiveness:** 88.89% average defense success rate with 100% rejection for Label Flipping and Gaussian Noise attacks.
2) **Significant Performance Improvement:** 20.32% accuracy improvement under Sign Flipping attacks compared to undefended baseline ($p < 0.0001$).
3) **Practical Overhead:** 5.48 seconds average computational overhead per round.
4) **Interpretable Decisions:** SHAP-based explanations provide human-understandable rejection criteria.

FedXChain represents a significant step toward trustworthy federated learning systems where defense decisions are not only effective but also explainable and auditable. The integration of explainable AI with blockchain technology creates a foundation for next-generation secure, transparent, and accountable distributed machine learning.

## REFERENCES

[1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. Int. Conf. Artif. Intell. Statist. (AISTATS)*, 2017, pp. 1273–1282.

[2] P. Kairouz, H. B. McMahan, B. Avent, *et al.*, "Advances and open problems in federated learning," *Found. Trends Mach. Learn.*, vol. 14, no. 1–2, pp. 1–210, 2021.

[3] J. Xu, B. S. Glicksberg, C. Su, P. Walker, J. Bian, and F. Wang, "Federated learning for healthcare informatics," *J. Healthcare Inform. Res.*, vol. 5, no. 1, pp. 1–19, 2021.

[4] J. Kang, Z. Xiong, D. Niyato, S. Xie, and J. Zhang, "Incentive mechanism for reliable federated learning: A joint optimization approach to combining reputation and contract theory," *IEEE Internet Things J.*, vol. 6, no. 6, pp. 10700–10714, 2019.

[5] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2017, pp. 119–129.

[6] D. Yin, Y. Chen, R. Kannan, and P. Bartlett, "Byzantine-robust distributed learning: Towards optimal statistical rates," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 5650–5659.

[7] X. Cao, J. Jia, and N. Z. Gong, "A comprehensive study of model poisoning attacks in federated learning," *IEEE Trans. Dependable Secure Comput.*, 2024.

[8] E. M. El Mhamdi, R. Guerraoui, and S. Rouault, "The hidden vulnerability of distributed learning in Byzantium," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 3521–3530.

[9] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2017, pp. 4765–4774.

[10] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, *et al.*, "Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI," *Inf. Fusion*, vol. 58, pp. 82–115, 2020.

[11] X. Cao, M. Fang, J. Liu, and N. Z. Gong, "FLTrust: Byzantine-robust federated learning via trust bootstrapping," in *Proc. Netw. Distrib. Syst. Secur. Symp. (NDSS)*, 2021.

[12] T. D. Nguyen, P. Rieger, R. De Viti, *et al.*, "FLAME: Taming backdoors in federated learning," in *Proc. USENIX Secur. Symp.*, 2022.

[13] H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Blockchained on-device federated learning," *IEEE Commun. Lett.*, vol. 24, no. 6, pp. 1279–1283, 2020.

[14] U. Majeed and C. S. Hong, "FLchain: Federated learning via MEC-enabled blockchain network," in *Proc. Asia-Pacific Netw. Oper. Manag. Symp. (APNOMS)*, 2019, pp. 1–4.

[15] M. Shayan, C. Fung, C. J. M. Yoon, and I. Beschastnikh, "Biscotti: A blockchain system for private and secure federated learning," *IEEE Trans. Parallel Distrib. Syst.*, vol. 32, no. 7, pp. 1513–1525, 2021.

[16] Y. Li, T. Liu, J. Gu, *et al.*, "Explainable AI meets anomaly detection," *arXiv preprint arXiv:2107.06114*, 2021.

[17] M. Wang, K. Zheng, Y. Yang, and X. Wang, "An explainable machine learning framework for intrusion detection systems," *IEEE Access*, vol. 8, pp. 73127–73141, 2020.

**FIRST AUTHOR** received the B.S. and M.S. degrees in computer science from University Name, Country, in 20XX and 20XX, respectively. He is currently pursuing the Ph.D. degree with the Department of Computer Science, University Name. His research interests include federated learning, blockchain technology, and explainable AI.



**SECOND AUTHOR** (Member, IEEE) received the Ph.D. degree in electrical engineering from University Name, Country, in 20XX. She is currently an Associate Professor with the Department of Electrical Engineering, University Name. Her research interests include distributed systems, machine learning security, and privacy-preserving computation.