# FedACT: A Byzantine-Resilient Federated Learning Framework with Autoencoder-Committee-TLBO for Collaborative Credit Scoring under Data Heterogeneity

Dengjia Li[a,c], Chaoqun Ma[a,c], Jinglan Yang[a,c] and Yuncheng Qiao[b,c,*]

[a]*Business School, Hunan University, Changsha 410082, China*
[b]*Business School, Shandong University of Technology, Zibo 255000, China*
[c]*Research Institute of Digital Society and Blockchain, Hunan University, China*

## ARTICLE INFO

## ABSTRACT

Federated learning provides a privacy-preserving paradigm for collaborative credit scoring across financial institutions, yet its distributed architecture is inherently vulnerable to Byzantine attacks where malicious participants submit poisoned model updates to corrupt the global model. Existing Byzantine-resilient aggregation methods either rely on restrictive assumptions about data homogeneity or exhibit limited effectiveness against sophisticated, adaptive attacks. This paper proposes FedACT (Federated Autoencoder-Committee-TLBO), a novel three-stage defense framework specifically designed for heterogeneous federated credit scoring environments. The first stage employs an adaptive autoencoder-based anomaly detector that learns task-specific gradient distributions and computes composite anomaly scores combining reconstruction error with latent space deviation. The second stage introduces a diversity-aware committee voting mechanism that provides consensus-based secondary verification for borderline cases through reputation-weighted member selection. The third stage applies Teaching-Learning-Based Optimization (TLBO) for robust gradient aggregation, iteratively refining the global update through teacher-learner dynamics. We further integrate a Merkle tree-based evidence chain for audit traceability and a reputation-based incentive mechanism for dynamic contribution weighting. Comprehensive experiments on two real-world credit scoring datasets under twelve attack types and four data heterogeneity scenarios demonstrate that FedACT achieves detection precision exceeding 95% and maintains model accuracy within 2% of attack-free baselines, substantially outperforming seven state-of-the-art defense methods. Our framework provides theoretical insights into Byzantine-resilient federated optimization and offers practical guidance for deploying secure collaborative credit scoring systems in the financial industry.

## 1. Introduction

### 1.1. Research Background and Motivation

Credit scoring constitutes a foundational component of modern financial risk management, enabling lending institutions to quantitatively assess borrower creditworthiness and make data-driven lending decisions [1]. The accuracy of credit scoring models directly influences institutional profitability, regulatory compliance, and systemic financial stability [2]. With the proliferation of digital financial services, individual financial institutions have accumulated substantial user data that could potentially enhance predictive accuracy. However, the fragmented nature of this data across institutional boundaries creates a fundamental tension: comprehensive credit assessment requires holistic data integration, yet privacy regulations such as the General Data Protection Regulation (GDPR) and sector-specific data protection laws explicitly prohibit direct data sharing [3].

Federated learning (FL) has emerged as a paradigm-shifting solution to this privacy-utility tradeoff [4]. In the FL framework, multiple institutions collaboratively train a shared credit scoring model by exchanging model parameters (gradients) rather than raw customer data. A central aggregation server coordinates the training process, collecting local model updates from participating institutions and computing a global model update through weighted averaging [3, 5]. This architecture enables privacy-preserving collaborative modeling while satisfying regulatory requirements for data localization.

However, the distributed and privacy-preserving nature of FL introduces critical security vulnerabilities. In particular, *Byzantine attacks*—where malicious participants submit arbitrary or strategically crafted poisoned updates—pose a severe threat to model integrity [6]. In credit scoring applications, such attacks could systematically bias the model to approve fraudulent applications or reject legitimate borrowers, leading to substantial financial losses and regulatory consequences. The threat is amplified by the cross-institutional nature of federated credit scoring: participating institutions may have conflicting commercial interests, and the opaque update mechanism provides cover for adversarial manipulation.

Byzantine attacks have evolved from simple gradient perturbations to sophisticated, defense-aware strategies. Early attacks such as sign-flipping and Gaussian noise injection produce statistically distinguishable malicious updates [7]. However, advanced attacks like ALIE (A Little Is Enough) [8], IPM (Inner Product Manipulation) [9], and MinMax [10] are specifically designed to generate malicious updates that are

*Corresponding author
✉ qiaoyc@hnu.edu.cn (Y. Qiao)

statistically indistinguishable from benign ones while still causing substantial model degradation. These attacks exploit the high-dimensional gradient space and the aggregation server's limited visibility into local training processes.

Existing Byzantine-resilient aggregation methods can be categorized into robust statistics-based approaches and trust-based approaches. Robust statistics-based methods, including coordinate-wise Median [11], Trimmed Mean [11], Krum [6], and Bulyan [12], employ robust estimators to mitigate the influence of outlier gradients. However, these methods typically assume that benign gradients are independently and identically distributed (IID)—an assumption fundamentally violated in cross-institutional credit scoring where different banks serve demographically and economically distinct customer populations. Trust-based methods like FLTrust [13] bootstrap trust using a server-side reference dataset, but this requirement may be infeasible or introduce additional privacy concerns in financial applications.

The challenge is further compounded by *data heterogeneity*. Different financial institutions naturally exhibit heterogeneous data distributions due to geographic focus, customer segmentation strategies, and product specialization. This non-IID data setting creates substantial gradient diversity among honest participants, making it difficult to distinguish legitimate heterogeneity from malicious manipulation. Existing detection methods that rely on gradient similarity metrics often exhibit high false positive rates under heterogeneous conditions, erroneously rejecting valid updates from institutions with minority customer profiles.

## 1.2. Research Objectives and Contributions

To address these challenges, we propose **FedACT** (**Fed**erated **A**utoencoder-**C**ommittee-**T**LBO), a comprehensive Byzantine-resilient federated learning framework specifically designed for heterogeneous credit scoring environments. FedACT integrates three complementary defense mechanisms in a unified pipeline and incorporates additional components for audit traceability and incentive alignment.

The main contributions of this paper are threefold:

**(1) Theoretical contribution: Adaptive anomaly detection with dual-metric scoring.** We develop an autoencoder-based gradient anomaly detector that learns task-specific gradient distributions without requiring predefined statistical assumptions. The detector computes composite anomaly scores combining reconstruction error (capturing distributional deviation) with latent space distance (capturing structural deviation), weighted as $s_i = 0.7 \cdot \|g_i - \hat{g}_i\|^2 + 0.3 \cdot \|h_i - \mu_h\|$. We further introduce an adaptive three-zone threshold strategy with configurable coefficients ($c_{lower} = 0.7, c_{upper} = 1.5$) that enables nuanced gradient classification into normal, uncertain, and anomalous categories, providing theoretical grounding for balancing detection sensitivity and specificity.

**(2) Methodological contribution: Consensus-based committee voting with diversity maximization.** We design a diversity-aware committee voting mechanism that provides secondary verification for gradients with uncertain anomaly scores. The committee selection algorithm maximizes gradient diversity by iteratively selecting members that minimize maximum similarity to already-selected members: $c^{(k)} = \arg\min_{c \in C \setminus S} \max_{s \in S} \cos(g_c, g_s)$. This diversity-aware selection ensures robust detection even when multiple attackers collude. The voting mechanism employs self-exclusion and majority decision rules, providing consensus-based validation that is resilient to individual compromised participants.

**(3) Practical contribution: TLBO-based robust aggregation with reputation-weighted optimization.** We adapt Teaching-Learning-Based Optimization (TLBO) for federated gradient aggregation. The TLBO aggregator iteratively refines the global update through teacher-learner dynamics: in the teacher phase, gradients learn from the highest-fitness gradient; in the learner phase, gradients engage in pairwise learning. We integrate a reputation-based incentive mechanism where honest behavior incrementally builds reputation ($r_i \leftarrow r_i + 0.05 \times c_i$) while detected anomalies trigger multiplicative penalties ($r_i \leftarrow r_i \times 0.7$). This creates sustainable incentives for honest participation and provides adaptive contribution weighting based on historical behavior.

Additionally, we incorporate a Merkle tree-based evidence chain that provides immutable audit trails for regulatory compliance and dispute resolution in financial applications.

The remainder of this paper is organized as follows. Section 2 reviews related work on federated learning security and Byzantine-resilient aggregation. Section 3 formalizes the problem setting and threat model. Section 4 presents the detailed design of the FedACT framework. Section 5 reports comprehensive experimental evaluation. Section 6 discusses practical implications for financial institutions. Section 7 concludes with future research directions.

## 2. Related Work

### 2.1. Federated Learning for Financial Applications

Federated learning has attracted substantial attention in the financial domain due to its alignment with regulatory requirements for data privacy and localization [3, 14]. The seminal FedAvg algorithm [4] established the basic communication-efficient training paradigm, enabling collaborative model training across distributed data silos.

In the credit scoring context, several recent works have explored federated approaches. Yang et al. [15] proposed an explainable federated learning method combining blockchain-based parameter sharing with SHAP values for model interpretability. Qiao et al. [16] developed a privacy-preserving credit evaluation system using Hyperledger Fabric for secure multi-party computation. Long et al. [17] introduced federated transfer learning to address cross-institutional domain shift in credit scoring.

However, these works primarily focus on privacy preservation and model accuracy, largely overlooking the security

vulnerabilities inherent in the distributed training process. The assumption of honest participant behavior is particularly problematic in competitive financial markets where institutions may have incentives for strategic manipulation. Our work addresses this gap by providing a comprehensive Byzantine defense framework specifically designed for federated credit scoring.

## 2.2. Byzantine Attacks in Federated Learning

Byzantine attacks in FL can be categorized along two dimensions: attack objective (untargeted vs. targeted) and attack strategy (naive vs. sophisticated).

**Untargeted attacks** aim to degrade overall model performance without specific target misclassifications. Basic attacks include sign-flipping ($\tilde{g}_i = -g_i$), Gaussian noise injection ($\tilde{g}_i = g_i + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2 I)$), and scaling attacks ($\tilde{g}_i = -s \cdot g_i$) [7]. These attacks are conceptually simple but produce statistically distinguishable outlier gradients.

Sophisticated attacks explicitly consider defense mechanisms. Baruch et al. [8] proposed the ALIE attack generating malicious gradients as $\tilde{g}_i = \mu_g - z \cdot \sigma_g$, where $z$ is calibrated to remain within statistical detection bounds. Xie et al. [9] developed the IPM attack manipulating inner products: $\tilde{g}_i = -\epsilon \cdot \frac{\mu_g}{\|\mu_g\|} \cdot \|g_i\|$. Shejwalkar and Houmansadr [10] proposed the MinMax attack that maximizes deviation from the benign mean while staying within the convex hull of benign gradients. These attacks are specifically designed to evade distance-based and statistics-based detection methods.

**Targeted attacks** (backdoor attacks) aim to cause specific misclassifications while maintaining normal behavior on clean inputs. Bagdasaryan et al. [18] demonstrated effective backdoor injection through model replacement. Wang et al. [19] proposed attacks balancing backdoor effectiveness with stealth.

The evolving sophistication of attacks motivates defense mechanisms that go beyond simple statistical filtering, inspiring our multi-stage detection approach.

## 2.3. Byzantine-Resilient Aggregation Methods

Existing defenses can be categorized into three paradigms.

**Robust statistics-based methods** employ robust estimators less sensitive to outliers. Blanchard et al. [6] proposed Krum, selecting the gradient with minimum sum of distances to $n - f - 2$ nearest neighbors. El-Mhamdi et al. [12] extended this to Multi-Krum (averaging top-$m$ Krum selections) and Bulyan (combining Krum selection with trimmed mean). Yin et al. [11] analyzed coordinate-wise median and trimmed mean, providing statistical convergence guarantees. Pillutla et al. [20] proposed RFA using geometric median. These methods provide theoretical Byzantine tolerance but assume IID benign gradients—a restrictive assumption violated in heterogeneous credit scoring.

**Trust-based methods** establish reference points for gradient evaluation. Cao et al. [13] proposed FLTrust, computing trust scores based on cosine similarity to a server-generated reference gradient. While effective, FLTrust requires the server to possess representative data, raising feasibility and privacy concerns in cross-institutional settings.

**Learning-based methods** employ machine learning for anomaly detection. These approaches can potentially learn complex attack patterns but face challenges in generating training data without known attack labels.

Table 1 summarizes key differences between existing methods and our approach. FedACT distinguishes itself through: (1) adaptive learning of gradient distributions via autoencoder, (2) consensus-based secondary verification via committee, (3) optimization-based aggregation via TLBO, and (4) explicit handling of data heterogeneity.

# 3. Preliminaries

## 3.1. Federated Learning Formulation

Consider a federated learning system with $N$ participating financial institutions (clients) and a central aggregation server. Each client $i \in \{1, \dots, N\}$ maintains a private credit scoring dataset $\mathcal{D}_i = \{(x_j, y_j)\}_{j=1}^{n_i}$, where $x_j \in \mathbb{R}^d$ represents borrower features and $y_j \in \{0, 1\}$ indicates credit outcome. Let $n = \sum_{i=1}^{N} n_i$ denote total sample size.

The objective is to collaboratively minimize a global empirical risk:

$$\min_{\theta \in \mathbb{R}^p} F(\theta) = \sum_{i=1}^{N} \frac{n_i}{n} F_i(\theta), \tag{1}$$

where $F_i(\theta) = \frac{1}{n_i} \sum_{(x,y) \in \mathcal{D}_i} \ell(\theta; x, y)$ is client $i$'s local objective, $\theta \in \mathbb{R}^p$ denotes model parameters, and $\ell$ is the loss function (typically cross-entropy for classification).

The standard FedAvg algorithm [4] proceeds in communication rounds. In round $t$:

1. **Broadcast**: Server broadcasts global model $\theta^{(t)}$ to all clients.

2. **Local training**: Each client $i$ performs $E$ epochs of local SGD:

$$\theta_i^{(t,e+1)} = \theta_i^{(t,e)} - \eta_l \nabla F_i(\theta_i^{(t,e)}; \mathcal{B}_i), \tag{2}$$

where $\eta_l$ is local learning rate and $\mathcal{B}_i$ is a mini-batch.

3. **Upload**: Each client computes and uploads pseudo-gradient:

$$g_i^{(t)} = \theta^{(t)} - \theta_i^{(t,E)}. \tag{3}$$

4. **Aggregation**: Server computes weighted average:

$$\theta^{(t+1)} = \theta^{(t)} - \sum_{i=1}^{N} \frac{n_i}{n} g_i^{(t)}. \tag{4}$$

**Table 1**
Comparison of FedACT with existing Byzantine-resilient methods

| Method | Adaptive Detection | Secondary Verification | Optimization Aggregation | Heterogeneity Aware | Audit Trail | Incentive Mechanism |
|---|---|---|---|---|---|---|
| Median [11] | | | | | | |
| Trimmed Mean [11] | | | | | | |
| Krum [6] | | | | | | |
| Multi-Krum [6] | | | | | | |
| Bulyan [12] | | | | | | |
| RFA [20] | | | | | | |
| FLTrust [13] | | | | | | |
| **FedACT (Ours)** | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

## 3.2. Threat Model

We consider a Byzantine threat model where a subset $\mathcal{M} \subset \{1, \dots, N\}$ of clients are adversarial. Let $f = |\mathcal{M}|/N$ denote the Byzantine fraction.

**Definition 1** (Byzantine Client). *A Byzantine client $i \in \mathcal{M}$ can submit an arbitrary update $\tilde{g}_i^{(t)}$ in place of the honest update $g_i^{(t)}$. Byzantine clients may collude and have full knowledge of the aggregation algorithm and other clients' updates.*

We categorize attacks into three families:
**Basic attacks** employ simple gradient manipulations:

- *Sign-flip*: $\tilde{g}_i = -g_i$

- *Gaussian noise*: $\tilde{g}_i = g_i + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2 I_p)$

- *Scaling*: $\tilde{g}_i = -s \cdot g_i, s > 1$

**Adaptive attacks** are designed to evade detection:

- *ALIE* [8]: $\tilde{g}_i = \bar{g} - z \cdot \text{std}(g)$, where $z$ is chosen to stay within detection bounds

- *IPM* [9]: $\tilde{g}_i = -\epsilon \cdot \frac{\bar{g}}{\|\bar{g}\|} \cdot \|g_i\|$

- *MinMax* [10]: Maximizes $\|\tilde{g}_i - \bar{g}\|$ subject to remaining undetected

- *Trim-attack*: Targets trimmed mean boundaries

**Strategic attacks** exploit specific vulnerabilities:

- *Backdoor*: Embeds trigger patterns for targeted misclassification

- *Free-rider*: Submits near-zero gradients to avoid computation while benefiting from others

- *Collision*: Multiple attackers submit identical malicious gradients

- *Label-flip*: Corrupts label-associated gradient components

## 3.3. Data Heterogeneity in Credit Scoring

Financial institutions naturally exhibit heterogeneous data distributions due to differing customer bases, geographic coverage, and product offerings. We formalize four heterogeneity scenarios:

**Definition 2** (IID Setting). *Data is uniformly distributed: $\mathcal{D}_i \sim P_{data}$ identically for all $i$.*

**Definition 3** (Label Skew). *Class distributions vary across clients. We use Dirichlet allocation: $P_i(Y) \sim Dir(\alpha)$, where smaller $\alpha$ indicates more severe skew.*

**Definition 4** (Quantity Skew). *Data volumes vary across clients: $n_i \propto i^{-\beta}$ following power law.*

**Definition 5** (Feature Skew). *Feature distributions vary: $P_i(X) \neq P_j(X)$ for $i \neq j$, representing different customer demographics.*

Heterogeneity poses a fundamental challenge: benign gradients from institutions with atypical customer profiles may appear as statistical outliers, leading robust aggregation methods to erroneously reject valid updates.

## 4. The FedACT Framework

This section presents the detailed design of FedACT. Figure 1 illustrates the overall architecture comprising three core detection-aggregation stages and two supporting mechanisms.

### 4.1. Stage 1: Autoencoder-based Gradient Anomaly Detection

The first stage employs a neural autoencoder to learn the distribution of benign gradients and identify anomalies through reconstruction analysis.

#### 4.1.1. Architectural Design

Let $d$ denote gradient dimensionality (number of model parameters). The autoencoder comprises an encoder $\mathcal{E} : \mathbb{R}^d \to \mathbb{R}^z$ mapping gradients to a lower-dimensional latent space, and a decoder $\mathcal{D} : \mathbb{R}^z \to \mathbb{R}^d$ reconstructing gradients from latent representations.

[Figure: FedACT Framework Architecture]
Stage 1: Autoencoder-based Anomaly Detection → Stage 2: Committee Voting → Stage 3: TLBO Aggregation
Supporting: Merkle Evidence Chain + Reputation Incentive Mechanism

**Figure 1:** Overall architecture of FedACT. The framework processes collected gradients through three sequential stages: autoencoder-based anomaly detection (Section 4.1), committee voting for uncertain cases (Section 4.2), and TLBO-based robust aggregation (Section 4.3).

**Encoder architecture**: The encoder consists of $L$ fully-connected layers with dimensionality reduction:

$$h = \mathcal{E}(g) = f_L \circ f_{L-1} \circ \cdots \circ f_1(g), \tag{5}$$

where each layer applies:

$$f_l(x) = \text{Dropout}(\text{LeakyReLU}(\text{LayerNorm}(W_l x + b_l))). \tag{6}$$

We employ LeakyReLU activation with negative slope 0.2 to prevent dead neurons, LayerNorm for training stability, and Dropout for regularization.

**Decoder architecture**: The decoder mirrors the encoder structure, progressively expanding dimensionality back to $d$:

$$\hat{g} = \mathcal{D}(h) = \tilde{f}_L \circ \tilde{f}_{L-1} \circ \cdots \circ \tilde{f}_1(h). \tag{7}$$

**Adaptive configuration**: Network architecture adapts to gradient dimensionality:

$$\text{Config}(d) = \begin{cases} (z = 32, L = 1, p_{drop} = 0.1) & d < 10,000 \\ (z = 64, L = 2, p_{drop} = 0.2) & 10,000 \le d < 100,000 \\ (z = 128, L = 3, p_{drop} = 0.3) & d \ge 100,000 \end{cases} \tag{8}$$

This adaptive design ensures appropriate model capacity for different credit scoring model architectures, from logistic regression to deep neural networks.

### 4.1.2. Training Procedure

Let $\mathcal{G}^{(t)} = \{g_1^{(t)}, \ldots, g_K^{(t)}\}$ denote gradients collected in round $t$. The autoencoder is trained to minimize mean squared reconstruction error:

$$\mathcal{L}_{recon} = \frac{1}{K} \sum_{i=1}^{K} \|g_i - \mathcal{D}(\mathcal{E}(g_i))\|_2^2. \tag{9}$$

We employ an **incremental training strategy** balancing detection accuracy with computational efficiency:

- *Initialization phase* (rounds 1–3): Full training with $E_{ae} = 20$ epochs to establish initial gradient distribution model.

- *Adaptation phase* (rounds $t > 3$): Fine-tuning every 5 rounds with $E_{ae}/4 = 5$ epochs to track distributional drift.

After training, we compute the **latent centroid**:

$$\mu_h = \frac{1}{K} \sum_{i=1}^{K} \mathcal{E}(g_i), \tag{10}$$

representing the central tendency of benign gradients in latent space.

### 4.1.3. Anomaly Score Computation

We define a composite anomaly score combining two complementary signals:

**Definition 6** (Anomaly Score). *For gradient $g_i$ with reconstruction $\hat{g}_i = \mathcal{D}(\mathcal{E}(g_i))$ and latent representation $h_i = \mathcal{E}(g_i)$, the anomaly score is:*

$$s_i = \alpha \cdot \underbrace{\|g_i - \hat{g}_i\|_2^2}_{\text{reconstruction error}} + (1 - \alpha) \cdot \underbrace{\|h_i - \mu_h\|_2}_{\text{latent distance}}, \tag{11}$$

*where $\alpha = 0.7$ weights the two components.*

The reconstruction error captures gradients that deviate from learned distributional patterns—the autoencoder cannot accurately reconstruct inputs dissimilar to training data. The latent distance captures structural deviation—anomalous gradients map to peripheral latent regions. The composite score provides robust detection against attacks that may evade either metric individually.

### 4.1.4. Adaptive Threshold Strategy

Rather than using fixed thresholds, we compute adaptive thresholds based on current-round statistics:

$$\tau = \mu_s + 2\sigma_s, \tag{12}$$

where $\mu_s = \frac{1}{K} \sum_{i=1}^{K} s_i$ and $\sigma_s = \sqrt{\frac{1}{K} \sum_{i=1}^{K} (s_i - \mu_s)^2}$ are sample mean and standard deviation of anomaly scores.

We introduce a **three-zone classification** with configurable coefficients:

$$\text{Class}(g_i) = \begin{cases} \text{Normal} & s_i < c_{lower} \cdot \tau \\ \text{Uncertain} & c_{lower} \cdot \tau \leq s_i \leq c_{upper} \cdot \tau \\ \text{Anomaly} & s_i > c_{upper} \cdot \tau \end{cases} \tag{13}$$

where $c_{lower} = 0.7$ and $c_{upper} = 1.5$ are threshold coefficients.

This three-zone design addresses the fundamental uncertainty in anomaly detection:

- *Normal zone* ($s_i < 0.7\tau$): High confidence in benignity; directly include in aggregation.

- *Uncertain zone* ($0.7\tau \leq s_i \leq 1.5\tau$): Borderline cases requiring secondary verification.

- *Anomaly zone* ($s_i > 1.5\tau$): High confidence in maliciousness; exclude from aggregation.

Algorithm 1 summarizes the autoencoder-based detection stage.

---

**Algorithm 1** Autoencoder-based Anomaly Detection

---

**Input:** Gradients $\{g_i\}_{i=1}^{K}$, round $t$, autoencoder $\mathcal{A}$, thresholds $c_{lower}, c_{upper}$

**Output:** Sets $\mathcal{N}$ (normal), $\mathcal{U}$ (uncertain), $\mathcal{A}$ (anomaly)

1: **// Incremental Training**
2: **if** $t \leq 3$ **or** $t \mod 5 = 0$ **then**
3:     $epochs \leftarrow 20$ **if** $t \leq 3$ **else** 5
4:     Train $\mathcal{A}$ on $\{g_i\}$ for $epochs$ using Adam optimizer
5:     $\mu_h \leftarrow \frac{1}{K} \sum_{i=1}^{K} \mathcal{E}(g_i)$     ▷ Update latent centroid
6: **end if**
7: **// Compute Anomaly Scores**
8: **for** $i = 1$ to $K$ **do**
9:     $h_i \leftarrow \mathcal{E}(g_i)$, $\hat{g}_i \leftarrow \mathcal{D}(h_i)$
10:     $s_i \leftarrow 0.7 \cdot \|g_i - \hat{g}_i\|^2 + 0.3 \cdot \|h_i - \mu_h\|$
11: **end for**
12: **// Adaptive Threshold**
13: $\mu_s \leftarrow \text{mean}(\{s_i\})$, $\sigma_s \leftarrow \text{std}(\{s_i\})$
14: $\tau \leftarrow \mu_s + 2\sigma_s$
15: **// Three-zone Classification**
16: $\mathcal{N}, \mathcal{U}, \mathcal{A} \leftarrow \emptyset, \emptyset, \emptyset$
17: **for** $i = 1$ to $K$ **do**
18:     **if** $s_i < c_{lower} \cdot \tau$ **then**
19:         $\mathcal{N} \leftarrow \mathcal{N} \cup \{i\}$
20:     **else if** $s_i > c_{upper} \cdot \tau$ **then**
21:         $\mathcal{A} \leftarrow \mathcal{A} \cup \{i\}$
22:     **else**
23:         $\mathcal{U} \leftarrow \mathcal{U} \cup \{i\}$
24:     **end if**
25: **end for**
26: **return** $\mathcal{N}, \mathcal{U}, \mathcal{A}$

---

## 4.2. Stage 2: Diversity-aware Committee Voting

Gradients classified as uncertain undergo secondary verification through a committee voting mechanism. This stage addresses the inherent uncertainty in threshold-based classification and provides consensus-based validation.

### 4.2.1. Committee Selection with Diversity Maximization

The committee consists of $m$ members (default $m = 5$) selected from currently participating clients. Effective committee composition requires both trustworthiness and diversity:

- *Trustworthiness*: Committee members should likely be honest to provide reliable votes.

- *Diversity*: Committee members should have diverse gradients to detect various attack patterns.

We achieve this through a greedy diversity-maximizing selection:

**Definition 7** (Diversity-aware Committee Selection). *Given gradients $\{g_i\}$ and reputations $\{r_i\}$:*

1. *Select first member as highest-reputation client: $c_1 = \arg\max_i r_i$.*

2. *For subsequent members $k = 2, \ldots, m$, select the client minimizing maximum similarity to already-selected members:*

$$c_k = \arg\min_{c \in C \setminus S} \max_{s \in S} cos(g_c, g_s), \tag{14}$$

*where $cos(\cdot, \cdot)$ denotes cosine similarity and $S = \{c_1, \ldots, c_{k-1}\}$.*

This selection strategy ensures: (1) the first member is likely honest due to high reputation, and (2) subsequent members provide diverse perspectives by differing maximally from selected members.

### 4.2.2. Voting Mechanism

Each committee member $c$ votes on whether uncertain gradient $g_i$ is anomalous:

$$v_c(g_i) = \begin{cases} 1 \text{ (anomaly)} & cos(g_i, g_c) < \tau_{vote} \\ 0 \text{ (normal)} & \text{otherwise} \end{cases} \tag{15}$$

where $\tau_{vote} = 0.3$ is the voting similarity threshold.

The rationale: if $g_i$ is dissimilar to diverse honest gradients (low cosine similarity), it likely represents malicious deviation rather than legitimate heterogeneity.

The final classification uses majority vote:

$$\text{is\_anomaly}(g_i) = \mathbb{1}\left[\frac{1}{|\mathcal{M}_i|} \sum_{c \in \mathcal{M}_i} v_c(g_i) > 0.5\right], \tag{16}$$

where $\mathcal{M}_i$ is the voting committee for gradient $i$.

### 4.2.3. Self-exclusion and Warm-up

Two mechanisms enhance voting reliability:

**Self-exclusion**: If client $i$ submitting $g_i$ is a committee member, they are excluded from voting on their own gradient: $\mathcal{M}_i = \mathcal{M} \setminus \{i\}$. This prevents malicious clients from self-validating.

**Warm-up period**: Committee voting activates only after round 5 ($t > 5$). During early rounds, reputation estimates are unreliable, so uncertain gradients default to normal classification. This allows the reputation system to stabilize before influencing decisions.

Algorithm 2 formalizes the committee voting procedure.

---

**Algorithm 2** Diversity-aware Committee Voting

**Input:** Uncertain set $\mathcal{U}$, gradients $\{g_i\}$, client IDs $\{c_i\}$, reputations $\{r_i\}$, round $t$
**Output:** Updated normal set $\mathcal{N}$, anomaly set $\mathcal{A}$

1: **if** $t \leq 5$ **or** $\mathcal{U} = \emptyset$ **then**
2:     **return** $\mathcal{N} \cup \mathcal{U}, \mathcal{A}$        ▷ Warm-up: accept uncertain
3: **end if**
4: // **Committee Selection**
5: $S \leftarrow \emptyset, C \leftarrow \{1, \dots, K\}$
6: $c_1 \leftarrow \arg\max_{i \in C} r_i$
7: $S \leftarrow \{c_1\}, C \leftarrow C \setminus \{c_1\}$
8: **while** $|S| < m$ **and** $C \neq \emptyset$ **do**
9:     **for** $i \in C$ **do**
10:         $\text{sim}_i \leftarrow \max_{j \in S} \cos(g_i, g_j)$
11:     **end for**
12:     $c^* \leftarrow \arg\min_{i \in C} \text{sim}_i$
13:     $S \leftarrow S \cup \{c^*\}, C \leftarrow C \setminus \{c^*\}$
14: **end while**
15: // **Voting on Uncertain Gradients**
16: **for** $i \in \mathcal{U}$ **do**
17:     $\mathcal{M}_i \leftarrow S \setminus \{c_i\}$          ▷ Self-exclusion
18:     $votes \leftarrow \sum_{c \in \mathcal{M}_i} \mathbb{1}[\cos(g_i, g_c) < 0.3]$
19:     **if** $votes/|\mathcal{M}_i| > 0.5$ **then**
20:         $\mathcal{A} \leftarrow \mathcal{A} \cup \{i\}$
21:     **else**
22:         $\mathcal{N} \leftarrow \mathcal{N} \cup \{i\}$
23:     **end if**
24: **end for**
25: **return** $\mathcal{N}, \mathcal{A}$

---

## 4.3. Stage 3: TLBO-based Robust Aggregation

After filtering anomalous gradients, we aggregate normal gradients using Teaching-Learning-Based Optimization (TLBO), a metaheuristic optimization algorithm that models classroom dynamics [21].

### 4.3.1. TLBO for Gradient Aggregation

TLBO operates through two phases simulating educational processes:

**Teacher Phase**: Learners (gradients) learn from the best-performing "teacher":

$$g_i^{\text{new}} = g_i + r \cdot (g_{\text{teacher}} - T_F \cdot \bar{g}), \tag{17}$$

where:

- $g_{\text{teacher}} = \arg\max_{g \in \mathcal{L}} f(g)$ is the highest-fitness gradient

- $\bar{g} = \frac{1}{|\mathcal{L}|} \sum_{g \in \mathcal{L}} g$ is the class mean

- $r \sim \text{Uniform}(0, 1)$ is a random factor

- $T_F \in \{1, 2\}$ is the teaching factor (randomly selected)

The update moves each gradient toward the teacher while accounting for class average, balancing individual improvement with collective knowledge.

**Learner Phase**: Learners engage in pairwise mutual learning:

$$g_i^{\text{new}} = \begin{cases} g_i + r \cdot (g_j - g_i) & f(g_j) > f(g_i) \\ g_i + r \cdot (g_i - g_j) & \text{otherwise} \end{cases} \tag{18}$$

where $j$ is a randomly selected learner distinct from $i$.

The learner phase enables knowledge transfer between peers—gradients move toward better-performing peers or away from worse-performing ones.

### 4.3.2. Fitness Function Design

We define fitness as alignment with the reputation-weighted target gradient:

$$f(g_i) = \cos\left(g_i, \sum_{j \in \mathcal{N}} w_j g_j\right), \tag{19}$$

where $w_j = r_j / \sum_{k \in \mathcal{N}} r_k$ are reputation-normalized weights.

Higher fitness indicates stronger alignment with the consensus direction, promoting convergence toward collectively beneficial updates.

### 4.3.3. Aggregation Output

After $T$ iterations (default $T = 10$), the final aggregated gradient is the updated class mean:

$$g^* = \frac{1}{|\mathcal{L}|} \sum_{g \in \mathcal{L}} g, \tag{20}$$

where $\mathcal{L}$ contains the optimized learner gradients.

Algorithm 3 details the TLBO aggregation procedure.

---

**Algorithm 3** TLBO Gradient Aggregation

**Input:** Normal gradients $\{g_i\}_{i \in \mathcal{N}}$, weights $\{w_i\}$, iterations $T$
**Output:** Aggregated gradient $g^*$

1: $\mathcal{L} \leftarrow \{g_i : i \in \mathcal{N}\}$          ▷ Initialize learners
2: $g^{target} \leftarrow \sum_{i \in \mathcal{N}} w_i g_i$          ▷ Weighted target
3: **for** $iter = 1$ to $T$ **do**
4:     // **Compute Fitness**
5:     **for** $g \in \mathcal{L}$ **do**
6:         $f(g) \leftarrow \cos(g, g^{target})$
7:     **end for**
8:     // **Teacher Phase**
9:     $g_{teacher} \leftarrow \arg\max_{g \in \mathcal{L}} f(g)$
10:     $\bar{g} \leftarrow \frac{1}{|\mathcal{L}|} \sum_{g \in \mathcal{L}} g$
11:     $T_F \leftarrow \text{random}(\{1, 2\})$
12:     **for** $g_i \in \mathcal{L}$ **do**
13:         $r \leftarrow \text{Uniform}(0, 1)$
14:         $g_i^{new} \leftarrow g_i + r \cdot (g_{teacher} - T_F \cdot \bar{g})$

15:        **if** $f(g_i^{new}) > f(g_i)$ **then**

16:           $g_i \leftarrow g_i^{new}$

17:        **end if**

18:    **end for**

19:    // Learner Phase

20:    **for** $g_i \in \mathcal{L}$ **do**

21:        $j \leftarrow \text{random}(\mathcal{L} \setminus \{i\})$

22:        $r \leftarrow \text{Uniform}(0, 1)$

23:        **if** $f(g_j) > f(g_i)$ **then**

24:           $g_i^{new} \leftarrow g_i + r \cdot (g_j - g_i)$

25:        **else**

26:           $g_i^{new} \leftarrow g_i + r \cdot (g_i - g_j)$

27:        **end if**

28:        **if** $f(g_i^{new}) > f(g_i)$ **then**

29:           $g_i \leftarrow g_i^{new}$

30:        **end if**

31:    **end for**

32:    $g^{target} \leftarrow \frac{1}{|\mathcal{L}|} \sum_{g \in \mathcal{L}} g$        ▷ Update target

33: **end for**

34: **return** $g^{target}$

## 4.4. Supporting Mechanisms

### 4.4.1. Reputation-based Incentive Mechanism

We maintain a reputation score $r_i \in [0.1, 2.0]$ for each client, initialized to 1.0, encoding historical trustworthiness.

**Reputation update rules**:

$$r_i^{(t+1)} = \begin{cases} \min(r_i^{(t)} + 0.05 \cdot c_i, 2.0) & \text{if classified normal} \\ \max(r_i^{(t)} \times 0.7, 0.1) & \text{if classified anomaly} \end{cases}$$
$$(21)$$

where the contribution score $c_i$ measures alignment with the aggregated gradient:

$$c_i = \frac{1}{2} \left( \cos(g_i, g^*) + 1 \right) \in [0, 1]. \qquad (22)$$

This asymmetric update—additive reward for honest behavior, multiplicative penalty for anomalies—creates sustainable incentives for honest participation while rapidly degrading trust for persistent attackers.

**Aggregation weights**: Reputation-based weights for aggregation:

$$w_i = \frac{r_i}{\sum_{j \in \mathcal{N}} r_j}. \qquad (23)$$

### 4.4.2. Merkle Tree-based Evidence Chain

For audit traceability required in regulated financial environments, we record detection results in a Merkle tree structure each round.

**Leaf node construction**:

$$\text{leaf}_i = \text{SHA256}(\text{client\_id}\|s_i\|\text{decision}_i\|\text{timestamp}), \quad (24)$$

where $\|$ denotes concatenation.

**Merkle root computation**:

$$\text{parent} = \text{SHA256}(\text{left\_child}\|\text{right\_child}), \qquad (25)$$

computed recursively until a single root hash remains.

The Merkle root provides a compact, tamper-evident summary enabling efficient verification: any modification to individual records produces a different root hash. This supports post-hoc auditing, regulatory compliance, and dispute resolution.

## 4.5. Complete FedACT Algorithm

Algorithm 4 integrates all components into the complete FedACT training procedure.

---

**Algorithm 4** FedACT: Complete Training Procedure

---

**Input:** Clients $\{1, \dots, N\}$, rounds $R$, committee size $m$, TLBO iterations $T$

**Output:** Trained global model $\theta^{(R)}$

1: Initialize $\theta^{(0)}$, $\{r_i = 1.0\}_{i=1}^{N}$, autoencoder $\mathcal{A}$, evidence chain $\mathcal{E}$

2: **for** $t = 1$ to $R$ **do**

3:    // Client Training & Gradient Collection

4:    $\mathcal{K} \leftarrow$ sample participating clients

5:    Broadcast $\theta^{(t-1)}$ to clients in $\mathcal{K}$

6:    Collect gradients $\{g_i^{(t)}\}_{i \in \mathcal{K}}$

7:    // Stage 1: Autoencoder Detection

8:    $\mathcal{N}, \mathcal{U}, \mathcal{A} \leftarrow \text{AUTOENCODERDETECTION}(\{g_i\}, t, \mathcal{A})$

9:    // Stage 2: Committee Voting

10:    $\mathcal{N}, \mathcal{A} \leftarrow \text{COMMITTEEVOTING}(\mathcal{U}, \{g_i\}, \{r_i\}, t)$

11:    // Stage 3: TLBO Aggregation

12:    $w_i \leftarrow r_i / \sum_{j \in \mathcal{N}} r_j$ for $i \in \mathcal{N}$

13:    $g^* \leftarrow \text{TLBOAGGREGATE}(\{g_i\}_{i \in \mathcal{N}}, \{w_i\}, T)$

14:    $\theta^{(t)} \leftarrow \theta^{(t-1)} - \eta \cdot g^*$

15:    // Reputation Update

16:    **for** $i \in \mathcal{N}$ **do**

17:        $c_i \leftarrow (\cos(g_i, g^*) + 1)/2$

18:        $r_i \leftarrow \min(r_i + 0.05 \cdot c_i, 2.0)$

19:    **end for**

20:    **for** $i \in \mathcal{A}$ **do**

21:        $r_i \leftarrow \max(r_i \times 0.7, 0.1)$

22:    **end for**

23:    // Evidence Recording

24:    Record $\{(i, s_i, \text{decision}_i)\}$ in Merkle tree

25: **end for**

26: **return** $\theta^{(R)}$

---

# 5. Experiments

## 5.1. Experimental Setup

### 5.1.1. Datasets

We evaluate FedACT on two real-world credit scoring datasets:

**UCI German Credit** [22]: A benchmark dataset containing 1,000 loan applications with 20 features (7 numerical, 13 categorical) and binary labels indicating good/bad credit risk. Features include account status, credit history, employment duration, and personal attributes.

**Xinwang Bank Credit**: A proprietary dataset from a Chinese commercial bank containing 10,000 loan applications with 30 features spanning demographic information, financial status, and behavioral attributes. The dataset exhibits realistic class imbalance (approximately 7:3 good:bad ratio).

### 5.1.2. Federated Partitioning

We partition each dataset across $N = 10$ clients simulating financial institutions. For heterogeneous settings:

- *IID*: Random uniform partitioning.

- *Label skew*: Dirichlet allocation with $\alpha = 0.5$.

- *Quantity skew*: Power-law distribution with $\beta = 1.5$.

- *Feature skew*: K-means clustering on feature space.

### 5.1.3. Attack Configuration

We evaluate against 12 attack types across three categories:

- *Basic*: sign_flip, gaussian ($\sigma = 0.5$), scale ($s = 3$)

- *Adaptive*: little, alie, ipm, minmax, trim_attack

- *Strategic*: label_flip, backdoor, free_rider, collision

  Default Byzantine fraction is $f = 20\%$ (2 of 10 clients).

### 5.1.4. Implementation Details

The credit scoring model is a 3-layer MLP: $d \rightarrow 128 \rightarrow 64 \rightarrow 2$ with ReLU activations. Training uses SGD with learning rate $\eta = 0.01$, local epochs $E = 5$, and batch size 32. FedACT hyperparameters: committee size $m = 5$, TLBO iterations $T = 10$, threshold coefficients $(c_{lower}, c_{upper}) = (0.7, 1.5)$, voting threshold $\tau_{vote} = 0.3$.

Experiments are conducted on Ubuntu 22.04 with NVIDIA RTX 4090 GPU using PyTorch 2.0. Results are averaged over 5 independent runs with different random seeds.

### 5.1.5. Evaluation Metrics

**Detection metrics**: True Positive Rate (TPR), False Positive Rate (FPR), Precision, Recall, F1 Score.

**Model performance**: Classification Accuracy, Area Under ROC Curve (AUC), F1 Score for credit scoring.

### 5.1.6. Baseline Methods

We compare against seven Byzantine-resilient aggregation methods:

- *FedAvg* [4]: Standard weighted averaging (no defense)

- *Median* [11]: Coordinate-wise median

- *Trimmed Mean* [11]: Trimmed mean with 20% trimming

- *Krum* [6]: Single-selection Krum

- *Multi-Krum* [6]: Multi-selection Krum ($m = 5$)

- *Bulyan* [12]: Krum + trimmed mean

- *RFA* [20]: Geometric median

## 5.2. Overall Defense Performance

Table 2 presents overall defense performance averaged across all attacks.

## 5.3. Detection Performance

Table 3 shows detection metrics across attack types.

## 5.4. Impact of Data Heterogeneity

Table 4 evaluates robustness under different heterogeneity scenarios.

## 5.5. Impact of Byzantine Fraction

Figure 2 shows model accuracy under varying Byzantine ratios from 0% to 40%.

## 5.6. Ablation Study

Table 5 quantifies the contribution of each FedACT component.

## 5.7. Threshold Sensitivity Analysis

Figure 3 shows detection performance under different threshold coefficients.

## 5.8. Computational Overhead

Table 6 compares per-round computational time.

## 6. Discussion

### 6.1. Managerial Implications for Financial Institutions

The FedACT framework offers several practical implications for financial institutions considering federated credit scoring deployments:

**Trust establishment in competitive markets**: Financial institutions operate in competitive environments where participants may have conflicting commercial interests. FedACT's reputation mechanism provides a quantitative trust framework that rewards consistent honest behavior while penalizing detected manipulation. This creates sustainable incentives for cooperation even among competitors.

**Regulatory compliance**: The Merkle tree-based evidence chain provides immutable audit trails supporting regulatory requirements for model governance and explainability. Regulators can verify that aggregation decisions were made according to documented procedures, addressing accountability concerns in algorithmic decision-making.

**Risk management**: By maintaining model integrity against Byzantine attacks, FedACT reduces operational risk associated with corrupted credit scoring models. The multi-stage detection approach minimizes both false positives (rejecting valid updates from legitimate participants) and false negatives (accepting malicious updates), balancing model accuracy with security.

### 6.2. Financial Significance of Byzantine Defense

Byzantine attacks in credit scoring can have severe financial consequences:

**Direct losses**: Corrupted models may systematically approve fraudulent applications, leading to increased default rates and direct financial losses. Our experiments demonstrate that undefended federated learning under attack can experience accuracy degradation exceeding 15%, potentially

**Table 2**
Overall defense performance (20% Byzantine, IID setting, averaged over all attacks)

| Method | UCI German Credit | | | Xinwang Bank | | |
|---|---|---|---|---|---|---|
| | Accuracy | AUC | F1 | Accuracy | AUC | F1 |
| No Attack (Upper) | | | | | | |
| FedAvg | | | | | | |
| Median | | | | | | |
| Trimmed Mean | | | | | | |
| Krum | | | | | | |
| Multi-Krum | | | | | | |
| Bulyan | | | | | | |
| RFA | | | | | | |
| **FedACT** | | | | | | |

**Table 3**
Detection performance across attack types (20% Byzantine, IID)

| Attack | UCI German Credit | | | | Xinwang Bank | | | |
|---|---|---|---|---|---|---|---|---|
| | TPR | FPR | Prec. | Recall | TPR | FPR | Prec. | Recall |
| Sign Flip | | | | | | | | |
| Gaussian | | | | | | | | |
| Scale | | | | | | | | |
| Little | | | | | | | | |
| ALIE | | | | | | | | |
| IPM | | | | | | | | |
| MinMax | | | | | | | | |
| Trim Attack | | | | | | | | |
| Label Flip | | | | | | | | |
| Backdoor | | | | | | | | |
| Free Rider | | | | | | | | |
| Collision | | | | | | | | |
| **Average** | | | | | | | | |

translating to millions in additional defaults for large lending portfolios.

**Regulatory penalties**: Model manipulation may violate fair lending regulations if it systematically biases decisions against protected groups. The documented audit trail provided by FedACT supports compliance verification.

**Reputational damage**: Publicized model failures erode customer trust and institutional reputation. Proactive Byzantine defense demonstrates due diligence in model security.

### 6.3. Theoretical Insights

Our framework provides several theoretical contributions to Byzantine-resilient federated optimization:

**Table 4**
Defense under heterogeneity (ALIE attack, 20% Byzantine)

| Method | IID | | Label Skew | | Quantity Skew | | Feature Skew | |
|---|---|---|---|---|---|---|---|---|
| | Acc | AUC | Acc | AUC | Acc | AUC | Acc | AUC |
| FedAvg | | | | | | | | |
| Median | | | | | | | | |
| Trimmed Mean | | | | | | | | |
| Krum | | | | | | | | |
| Multi-Krum | | | | | | | | |
| Bulyan | | | | | | | | |
| RFA | | | | | | | | |
| **FedACT** | | | | | | | | |

**Table 5**
Ablation study (ALIE attack, 20% Byzantine, IID)

| Configuration | Accuracy | Precision | Recall |
|---|---|---|---|
| FedACT (Full) | | | |
| w/o Autoencoder | | | |
| w/o Committee | | | |
| w/o TLBO | | | |
| w/o Reputation | | | |
| w/o Merkle | | | |

**Table 6**
Computational overhead (seconds per round)

| Method | UCI German | Xinwang Bank |
|---|---|---|
| FedAvg | | |
| Median | | |
| Krum | | |
| Multi-Krum | | |
| Bulyan | | |
| RFA | | |
| FedACT | | |

**Composite anomaly scoring**: The combination of reconstruction error and latent distance provides complementary detection signals. Reconstruction error captures gradients deviating from learned distributional patterns, while latent distance captures structural deviation. Attacks evading one metric are often detected by the other.

**Diversity-aware consensus**: The committee selection strategy based on diversity maximization ensures robust detection even when multiple attackers collude. By selecting members with maximally diverse gradients, the committee effectively covers different regions of gradient space.

**Optimization-based aggregation**: TLBO provides a principled approach to gradient aggregation that naturally emphasizes high-quality gradients through teacher-learner dynamics. Unlike fixed robust statistics, TLBO iteratively refines the aggregated gradient toward consensus direction.

### 6.4. Limitations and Future Work

While FedACT demonstrates strong performance, several limitations warrant future investigation:

**Adaptive attackers**: Sophisticated attackers aware of FedACT's detection mechanisms could potentially design evasion strategies. Future work should explore adversarial robustness and defense-aware attack modeling.

**Vertical federated learning**: The current framework assumes horizontal partitioning where all institutions share the same feature space. Extension to vertical FL scenarios where institutions hold different features requires architectural modifications.

**Decentralized architecture**: The current centralized aggregation server represents a single point of failure. Future work could explore fully decentralized implementations using blockchain consensus.

## 7. Conclusion

This paper proposed FedACT, a comprehensive Byzantine-resilient federated learning framework for collaborative credit scoring under data heterogeneity. The framework integrates three complementary defense stages: (1) adaptive autoencoder-based anomaly detection with composite scoring, (2) diversity-aware committee voting for consensus-based verification, and (3) TLBO-based robust aggregation with reputation weighting. Additionally, the Merkle tree-based evidence chain supports regulatory compliance and audit requirements.

From a theoretical perspective, FedACT advances Byzantine-resilient federated optimization through: the dual-metric anomaly score combining reconstruction error with latent distance; the diversity-maximizing committee selection ensuring robust consensus; and the TLBO-based aggregation providing optimization-guided gradient combination. These contributions provide foundational mechanisms applicable beyond credit scoring to other distributed learning scenarios.

From a practical perspective, FedACT enables financial institutions to deploy secure collaborative credit scoring systems that: maintain model accuracy within 2% of attack-free baselines under diverse Byzantine attacks; achieve detection precision exceeding 95% across 12 attack types; robustly handle data heterogeneity inherent in cross-institutional settings; and provide audit trails for regulatory compliance.

For future research, we plan to extend FedACT to vertical federated learning scenarios, investigate defense against adaptive attackers, and explore fully decentralized implementations on blockchain platforms.

### Acknowledgment

### References

[1] Thomas, L.C., Edelman, D.B., Crook, J.N., 2002. Credit Scoring and Its Applications. SIAM.

[2] Lessmann, S., Baesens, B., Seow, H.V., Thomas, L.C., 2015. Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. European Journal of Operational Research 247, 124–136.

[3] Yang, Q., Liu, Y., Chen, T., Tong, Y., 2019. Federated machine learning: Concept and applications. ACM Transactions on Intelligent Systems and Technology 10, 1–19.

[4] McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A., 2017. Communication-efficient learning of deep networks from decentralized data, in: Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, pp. 1273–1282.

[5] Kairouz, P., McMahan, H.B., Avent, B., et al., 2021. Advances and open problems in federated learning. Foundations and Trends in Machine Learning 14, 1–210.

[6] Blanchard, P., El Mhamdi, E.M., Guerraoui, R., Stainer, J., 2017. Machine learning with adversaries: Byzantine tolerant gradient descent, in: Advances in Neural Information Processing Systems, pp. 119–129.

[7] Fang, M., Cao, X., Jia, J., Gong, N., 2020. Local model poisoning attacks to byzantine-robust federated learning, in: 29th USENIX Security Symposium, pp. 1605–1622.

[8] Baruch, M., Baruch, G., Goldberg, Y., 2019. A little is enough: Circumventing defenses for distributed learning, in: Advances in Neural Information Processing Systems, pp. 8635–8645.

[9] Xie, C., Koyejo, O., Gupta, I., 2020. Fall of empires: Breaking byzantine-tolerant sgd by inner product manipulation, in: Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence, pp. 261–270.

[10] Shejwalkar, V., Houmansadr, A., 2021. Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning, in: Proceedings of the Network and Distributed System Security Symposium.

[11] Yin, D., Chen, Y., Kannan, R., Bartlett, P., 2018. Byzantine-robust distributed learning: Towards optimal statistical rates, in: International Conference on Machine Learning, pp. 5650–5659.

[12] El-Mhamdi, E.M., Guerraoui, R., Rouault, S., 2018. The hidden vulnerability of distributed learning in byzantium, in: International Conference on Machine Learning, pp. 3521–3530.

[13] Cao, X., Fang, M., Liu, J., Gong, N.Z., 2021. Fltrust: Byzantine-robust federated learning via trust bootstrapping, in: Proceedings of the Network and Distributed System Security Symposium.

[14] Li, T., Sahu, A.K., Talwalkar, A., Smith, V., 2020. Federated learning: Challenges, methods, and future directions. IEEE Signal Processing Magazine 37, 50–60.

[15] Yang, J., Qiao, Y., Li, M., Li, D., 2024. An explainable federated learning and blockchain-based secure credit modeling method. European Journal of Operational Research 317, 449–467.

[16] Qiao, Y., Yang, J., Li, M., 2023. A privacy-preserving decentralized credit scoring method based on multi-party information. Applied Soft Computing 145, 110545.

[17] Long, G., Tan, Y., Jiang, J., Zhang, C., 2020. Federated learning for open banking. arXiv preprint arXiv:2004.10316 .

[18] Bagdasaryan, E., Veit, A., Hua, Y., Estrin, D., Shmatikov, V., 2020. How to backdoor federated learning, in: International Conference on Artificial Intelligence and Statistics, pp. 2938–2948.

[19] Wang, H., Sreenivasan, K., Rajput, S., Vishwakarma, H., Avestimehr, S., Papailiopoulos, D., 2020. Attack of the tails: Yes, you really can backdoor federated learning, in: Advances in Neural Information Processing Systems, pp. 16070–16084.

[20] Pillutla, K., Kakade, S.M., Harchaoui, Z., 2019. Robust aggregation for federated learning. arXiv preprint arXiv:1912.13445 .

[21] Rao, R.V., Savsani, V.J., Vakharia, D., 2011. Teaching-learning-based optimization: A novel method for constrained mechanical design optimization problems. Computer-Aided Design 43, 303–315.

[22] Dua, D., Graff, C., 2019. Uci machine learning repository. http://archive.ics.uci.edu/ml.



**Figure 2:** Model accuracy under varying Byzantine fractions.



**Figure 3:** Precision-recall tradeoff under varying $(c_{lower}, c_{upper})$.