# FedACT: A Byzantine-Resilient Federated Learning Framework with Autoencoder-Committee-TLBO for Heterogeneous Credit Scoring

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

*Keywords*:
Federated learning
Byzantine resilience
Credit scoring
Anomaly detection
Data heterogeneity

## ABSTRACT

Federated learning enables privacy-preserving collaborative credit scoring across financial institutions, yet its distributed architecture is vulnerable to Byzantine attacks where malicious participants submit poisoned model updates. Existing defenses assume data homogeneity—an assumption fundamentally violated in cross-institutional credit scoring where banks serve distinct customer populations. This paper proposes FedACT, a three-stage Byzantine-resilient framework: (1) an adaptive autoencoder-based anomaly detector computing dual-metric scores (reconstruction error weighted at 0.7 plus latent deviation at 0.3) with three-zone thresholding (boundaries at $0.7\tau$ and $1.5\tau$), (2) a diversity-maximizing committee voting mechanism for consensus-based verification of uncertain gradients, and (3) Teaching-Learning-Based Optimization (TLBO) for robust gradient aggregation with reputation-based incentives (normal: $r_i + 0.05c_i$; anomalous: $r_i \times 0.7$). We further integrate Merkle tree evidence chains for regulatory compliance. Experiments on UCI Credit Card and Xinwang Bank datasets under twelve attack types and four heterogeneity scenarios demonstrate FedACT achieves >95% detection precision while maintaining accuracy within 2% of attack-free baselines, substantially outperforming seven state-of-the-art defenses including Krum, Bulyan, and FLTrust.

## 1. Introduction

### 1.1. Problem Scenario and Motivation

Consider the following real-world scenario: Bank A, a regional mortgage lender, seeks to improve its credit scoring model but possesses limited data on applicants' credit card behavior. Meanwhile, Bank B, a national credit card issuer, has extensive transaction histories but lacks mortgage repayment patterns. Both institutions recognize that collaborative modeling could substantially improve predictive accuracy—research suggests cross-institutional data integration can reduce default prediction error by 15-25% [1]. However, strict privacy regulations including the EU General Data Protection Regulation (GDPR), China's Personal Information Protection Law (PIPL), and sector-specific banking regulations explicitly prohibit direct sharing of customer financial data [2]. This creates a fundamental tension: *comprehensive credit assessment demands holistic data integration, while regulatory compliance mandates strict data localization.*

Federated learning (FL) has emerged as a compelling solution to this privacy-utility tradeoff [3, 4]. In federated credit scoring, participating institutions train local models on proprietary customer data and transmit only model updates (gradients) to a central aggregation server, which combines these updates into a global model. This architecture enables privacy-preserving collaboration while satisfying regulatory requirements—raw customer data never leaves institutional boundaries.

However, the distributed and privacy-preserving nature of FL introduces critical security vulnerabilities that threaten the entire collaborative ecosystem. **Byzantine attacks—** where malicious participants submit arbitrarily crafted poisoned updates—pose existential threats to model integrity [5, 6]. Unlike centralized machine learning where the training process is fully observable, federated systems provide adversaries with *opacity*: the aggregation server cannot inspect local training data or intermediate computations, enabling sophisticated gradient manipulation that is difficult to detect.

In credit scoring applications, Byzantine attacks carry severe financial and systemic consequences:

- **Fraudulent approval attacks**: A malicious competitor could systematically bias the shared model to approve high-risk loan applicants, increasing default rates across consortium members and causing millions in losses.

- **Legitimate rejection attacks**: Adversaries could poison the model to reject creditworthy borrowers from specific demographics, reducing market share and potentially violating fair lending regulations.

- **Systemic manipulation**: State-sponsored actors or organized crime could destabilize credit markets during economic stress by amplifying default cascades through model manipulation.

The threat is amplified by the cross-institutional nature of federated credit scoring: participating banks may have conflicting commercial interests, regulatory arbitrage incentives, or differing cybersecurity postures. A bank facing financial difficulties might be tempted to poison models to approve marginal applicants; a competitor might sabotage shared models to gain market advantage.

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

## 1.2. Research Gap: The Heterogeneity Challenge

Existing Byzantine-resilient aggregation methods [5, 7, 8, 9] predominantly assume that honest participants' gradients are independently and identically distributed (IID). This assumption *fundamentally fails* in federated credit scoring.

Different financial institutions naturally exhibit **data heterogeneity** due to:

- **Customer segmentation**: A bank targeting millennials (younger, gig economy, variable income) produces dramatically different gradient patterns than one serving retirees (fixed income, established credit history).

- **Geographic concentration**: Urban banks observe different spending patterns than rural community banks.

- **Product specialization**: Mortgage lenders, credit card issuers, and auto loan providers see distinct feature distributions.

- **Default definitions**: A 60-day delinquency threshold at one institution differs from 90-day thresholds elsewhere.

This non-IID data creates substantial gradient diversity among *honest* participants. When Bank A (mortgages, older customers) and Bank B (credit cards, younger customers) submit their gradients, the updates naturally diverge due to legitimate data differences—not malicious intent.

Existing detection methods that rely on gradient similarity metrics face a critical dilemma:

- **Strict thresholds** incorrectly reject valid updates from institutions with minority customer profiles (high false positive rate), undermining collaboration and excluding valuable data.

- **Lenient thresholds** fail to detect sophisticated attacks that camouflage malicious updates within the natural gradient variance (high false negative rate), leaving the system vulnerable.

Empirical evidence confirms this challenge. Karimireddy et al. [10] demonstrated that Krum, a leading Byzantine defense, suffers 40% accuracy degradation on non-IID CIFAR-10 partitions. Our preliminary experiments on credit scoring data show similar patterns: under label heterogeneity (Dirichlet $\beta = 0.5$), Krum incorrectly rejects 35% of honest gradients while missing 28% of ALIE attacks.

## 1.3. Research Contributions

To address these challenges, we propose **FedACT** (**Fed**erated **A**utoencoder-**C**ommittee-**T**LBO), a comprehensive Byzantine defense framework specifically designed for heterogeneous federated credit scoring. Our contributions are threefold:

1. **Representation learning for heterogeneity-aware detection.** We develop an autoencoder-based gradient anomaly detector that learns task-specific normal gradient distributions without requiring IID assumptions. The detector computes composite anomaly scores $s_i = 0.7 \cdot \|g_i - \hat{g}_i\|^2 + 0.3 \cdot \|h_i - \mu_h\|$ combining reconstruction error with latent deviation, achieving >95% detection precision under non-IID data distributions.

2. **Consensus-based verification through diversity maximization.** We design a committee voting mechanism with diversity-aware member selection ($c^{(k)} = \arg\min_c \max_{s \in S} \cos(g_c)$ that provides secondary verification for uncertain gradients, reducing false positive rates by 47% compared to single-threshold methods while maintaining robust detection under collusion attacks.

3. **Optimization-driven aggregation with long-term incentives.** We adapt Teaching-Learning-Based Optimization (TLBO) for gradient aggregation and integrate reputation-based penalties (normal: $r_i \leftarrow r_i + 0.05 \cdot c_i$; anomalous: $r_i \leftarrow r_i \times 0.7$) that create sustainable Byzantine deterrence, maintaining 98.8% accuracy preservation across twelve attack types.

Additionally, we incorporate Merkle tree-based evidence chains providing immutable audit trails for regulatory compliance in financial applications.

The remainder of this paper is organized as follows. Section 2 reviews related work on federated learning security and Byzantine defenses. Section 3 formalizes the problem setting and threat model. Section 4 presents the detailed design of FedACT. 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 finance due to its alignment with regulatory requirements for data privacy [4, 11, 12]. The seminal FedAvg algorithm [3] established communication-efficient collaborative training by averaging locally computed gradients. Subsequent work addressed non-IID data challenges through personalization [13], regularization [14], and contrastive learning [15].

**Evolution of federated optimization.** FedProx [14] introduced a proximal term penalizing local deviation from the global model, improving convergence under heterogeneity but without Byzantine considerations. MOON [15] applied contrastive learning, pulling local representations toward the global model while pushing away from previous local states. Scaffold [10] employed control variates to reduce client drift. These advances assume honest participation and focus on statistical heterogeneity, not adversarial behavior.

**Privacy-preserving techniques.** Differential privacy (DP) adds calibrated noise to gradients, providing formal privacy guarantees [16]. Secure aggregation protocols [17] enable encrypted gradient summation without revealing

individual contributions. However, DP degrades model accuracy and secure aggregation prevents gradient inspection, potentially hiding malicious updates.

In credit scoring specifically, recent works have explored federated approaches. Yang et al. [18] proposed an explainable federated learning method integrating blockchain-based parameter sharing with SHAP values for model interpretability, achieving audit trail for credit decisions. Qiao et al. [19] developed a privacy-preserving credit evaluation system using Hyperledger Fabric for secure computation with smart contract governance. Long et al. [20] introduced federated transfer learning to address cross-institutional domain shift where source and target institutions have different customer demographics. Cheng et al. [21] proposed vertical federated learning for credit risk with differential privacy guarantees, partitioning features across institutions.

**Vertical vs. horizontal FL.** Most FL credit scoring research follows horizontal partitioning (institutions share feature space, differ in samples). Vertical FL [4] addresses feature partitioning where different institutions hold different attributes of the same customers (e.g., credit bureau holds repayment history, bank holds transaction data). Our work focuses on horizontal FL while the defense mechanisms generalize to vertical settings with appropriate gradient transformation.

However, these works predominantly focus on privacy preservation and model accuracy, largely neglecting security vulnerabilities inherent in distributed training. The implicit assumption of honest participation is problematic in competitive financial markets where institutions may have adversarial incentives. A compromised or malicious participant can systematically degrade model quality without detection in undefended systems.

## 2.2. Byzantine Attacks in Federated Learning

Byzantine attacks in FL have evolved from simple perturbations to sophisticated, defense-aware strategies.

**Basic attacks** produce statistically distinguishable malicious updates. Fang et al. [22] systematically evaluated sign-flipping ($\tilde{g}_i = -g_i$), Gaussian noise injection ($\tilde{g}_i = g_i + \mathcal{N}(0, \sigma^2 I)$), and scaling attacks ($\tilde{g}_i = \lambda g_i$). These attacks are effective against undefended systems but easily detected by distance-based methods.

**Sophisticated attacks** explicitly evade detection mechanisms. Baruch et al. [23] proposed ALIE (A Little Is Enough), generating malicious updates as $\tilde{g}_i = \mu_g - z \cdot \sigma_g$ where $z$ is calibrated to remain within statistical detection bounds. Xie et al. [24] developed IPM (Inner Product Manipulation): $\tilde{g}_i = -\epsilon \cdot \|\mu_g\|^{-1} \mu_g \cdot \|g_i\|$. Shejwalkar and Houmansadr [25] proposed MinMax attacks that maximize deviation from the benign mean while staying within the convex hull of honest gradients. These attacks are specifically designed to evade robust aggregation methods.

**Backdoor attacks** inject targeted misclassifications while maintaining normal behavior on clean inputs. Bagdasaryan

et al. [26] demonstrated model replacement attacks achieving persistent backdoors. Wang et al. [27] balanced backdoor effectiveness with stealth through constrained optimization.

**Recent advances (2024-2025)** include adaptive attacks that learn detection thresholds [28], gradient-matching attacks that mimic honest update statistics [29], and collusion-aware attacks coordinating multiple adversaries [30]. These emerging threats highlight the need for multi-stage defense mechanisms that cannot be circumvented by optimizing against a single detection criterion.

## 2.3. Byzantine-Resilient Aggregation Methods

Existing defenses can be categorized into three paradigms, each with distinct assumptions and limitations:

**Robust statistics-based methods** employ estimators less sensitive to outliers. Blanchard et al. [5] proposed Krum, selecting the gradient with minimum sum of distances to $n - f - 2$ nearest neighbors, tolerating $f < n/2 - 1$ Byzantine clients. The selection rule ensures the chosen gradient is centrally located among honest participants under IID assumptions. El-Mhamdi et al. [8] extended this to Multi-Krum (averaging top-$m$ selections) and Bulyan (combining Krum with coordinate-wise trimmed mean). Yin et al. [7] analyzed coordinate-wise median and trimmed mean, providing statistical convergence guarantees under bounded Byzantine fraction with $O(1/\sqrt{nT})$ convergence rate. Pillutla et al. [31] proposed RFA using geometric median with provable breakdown point of 50%.

**Limitations under heterogeneity.** These methods provide theoretical Byzantine tolerance but critically assume IID benign gradients. Under heterogeneity, gradient distances among honest participants increase substantially, causing several failure modes: (1) Krum may select atypical gradients from minority institutions, (2) trimmed mean excludes valid contributions from banks with distinctive customer profiles, and (3) coordinate-wise median produces biased estimates when feature importance varies across institutions. Empirically, Krum's accuracy drops by 15-20% under label skew heterogeneity even without attacks [11].

**Trust-based methods** establish reference points for gradient evaluation. Cao et al. [9] proposed FLTrust, computing trust scores based on cosine similarity to a server-generated reference gradient trained on a small root dataset:

$$\text{TS}_i = \max\left(0, \cos(g_i, g_{\text{server}})\right) \qquad (1)$$

where $g_{\text{server}}$ is trained on server-held data. Updates with negative cosine similarity receive zero weight. While effective when the root dataset is representative, FLTrust requires the server to possess labeled data—potentially infeasible in cross-institutional settings where no single party has comprehensive customer coverage. Moreover, if the root dataset is unrepresentative (e.g., server only has urban customer data), the reference gradient may incorrectly penalize valid updates from rural-focused institutions.

**Learning-based methods** employ machine learning for anomaly detection. Li et al. [32] applied isolation forests to gradient statistics, detecting anomalies based on feature

**Table 1**
Comparison of Byzantine defense methods for federated learning

| Method | IID Required | Server Data | Detection | Aggregation | Incentive | Audit |
|---|---|---|---|---|---|---|
| Krum [5] | Yes | No | Distance | Selection | No | No |
| Multi-Krum [5] | Yes | No | Distance | Averaging | No | No |
| Coordinate Median [7] | Yes | No | Statistical | Median | No | No |
| Trimmed Mean [7] | Yes | No | Statistical | Trimming | No | No |
| Bulyan [8] | Yes | No | Distance+Trim | Combined | No | No |
| RFA [31] | Yes | No | Geometric | Median | No | No |
| FLTrust [9] | No | **Yes** | Similarity | Weighted | No | No |
| **FedACT (Ours)** | **No** | **No** | **Autoencoder** | **TLBO** | **Yes** | **Yes** |

isolation depth. Zhang et al. [33] used autoencoders for reconstruction-based detection with threshold $\mu + 2\sigma$ on reconstruction error. However, these works treat detection as binary classification, lacking consensus mechanisms for borderline cases that may be legitimate heterogeneous updates or subtle attacks. They also do not address incentive alignment for long-term honest participation—detected adversaries face no lasting consequences and may resume attacks in future rounds.

**Research gap summary.** Existing methods address *detection* (identifying malicious gradients) or *aggregation* (combining remaining gradients robustly), but rarely integrate both with: (1) explicit heterogeneity handling, (2) uncertainty quantification for borderline cases, (3) consensus-based verification, (4) reputation-based incentive mechanisms, and (5) audit trail for regulatory compliance. FedACT addresses this gap through a three-stage integrated framework.

Table 1 summarizes key distinctions. FedACT uniquely combines adaptive learning (autoencoder), consensus verification (committee), and optimization-based aggregation (TLBO) while explicitly handling heterogeneity without requiring server-side data.

## 3. Problem Formulation and Threat Model

### 3.1. Federated Credit Scoring Formulation

Consider $N$ financial institutions (clients) collaboratively training a credit scoring model. Each client $i \in \{1, 2, \dots, N\}$ possesses private training data $\mathcal{D}_i = \{(\mathbf{x}_j^{(i)}, y_j^{(i)})\}_{j=1}^{n_i}$ where $\mathbf{x}_j \in \mathbb{R}^d$ represents borrower features (income, credit history, demographics) and $y_j \in \{0, 1\}$ indicates default status (1 = default, 0 = non-default).

The federated learning objective is to minimize the weighted empirical risk:

$$\min_{\mathbf{w} \in \mathbb{R}^p} F(\mathbf{w}) = \sum_{i=1}^{N} \frac{n_i}{n} F_i(\mathbf{w}) \quad (2)$$

where $n = \sum_{i=1}^{N} n_i$ is the total sample size, $\mathbf{w}$ denotes model parameters, and $F_i(\mathbf{w})$ is client $i$'s local objective:

$$F_i(\mathbf{w}) = \frac{1}{n_i} \sum_{j=1}^{n_i} \ell(f(\mathbf{w}; \mathbf{x}_j^{(i)}), y_j^{(i)}) \quad (3)$$

where $f(\mathbf{w}; \cdot)$ is the credit scoring model (e.g., neural network) and $\ell(\cdot, \cdot)$ is the loss function (e.g., binary cross-entropy for classification).

Standard federated optimization (FedAvg [3]) proceeds iteratively:

1. **Server broadcast**: At round $t$, the server sends global parameters $\mathbf{w}^{(t)}$ to a selected subset of clients $S^{(t)} \subseteq \{1, \dots, N\}$.

2. **Local training**: Each selected client $i \in S^{(t)}$ computes the local gradient:

$$\mathbf{g}_i^{(t)} = \nabla F_i(\mathbf{w}^{(t)}) = \frac{1}{n_i} \sum_{j=1}^{n_i} \nabla \ell(f(\mathbf{w}^{(t)}; \mathbf{x}_j^{(i)}), y_j^{(i)}) \quad (4)$$

3. **Gradient upload**: Clients transmit gradients $\{\mathbf{g}_i^{(t)}\}_{i \in S^{(t)}}$ to the server.

4. **Aggregation**: The server computes the global update:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \sum_{i \in S^{(t)}} \frac{n_i}{\sum_{j \in S^{(t)}} n_j} \mathbf{g}_i^{(t)} \quad (5)$$

where $\eta > 0$ is the learning rate.

### 3.2. Data Heterogeneity Characterization

Real-world credit scoring data exhibits multiple forms of heterogeneity across institutions:

**Definition 1** (Feature Heterogeneity). *Clients collect different feature subsets. Let $\mathcal{X}_i \subseteq \mathcal{X}$ denote client $i$'s feature space. Feature heterogeneity occurs when $\mathcal{X}_i \neq \mathcal{X}_j$ for some $i \neq j$.*

**Definition 2** (Label Heterogeneity). *Default rates and label distributions vary across clients. Let $P_i(Y = 1)$ denote client $i$'s default rate. Label heterogeneity occurs when $P_i(Y = 1) \neq P_j(Y = 1)$ for some $i \neq j$.*

**Definition 3** (Quantity Heterogeneity). *Sample sizes differ dramatically. Quantity heterogeneity is characterized by high variance in $\{n_i\}_{i=1}^{N}$, e.g., national banks with millions of customers versus community banks with thousands.*

We formally quantify heterogeneity via gradient divergence:

$$\mathcal{H} = \frac{1}{N} \sum_{i=1}^{N} \|\nabla F_i(\mathbf{w}^*) - \nabla F(\mathbf{w}^*)\|^2 \qquad (6)$$

where $\mathbf{w}^*$ is the global optimum. High $\mathcal{H}$ indicates significant non-IID data, causing honest gradients to exhibit high variance.

## 3.3. Threat Model

**Adversary population.** We consider $M$ Byzantine clients among $N$ total, where $M < N/2$ (honest majority assumption). Byzantine clients are controlled by adversaries who may coordinate attacks. In financial settings, adversaries may be: (1) compromised institutions with malware-infected systems, (2) malicious insiders with profit motives, (3) competitors seeking to degrade consortium model quality, or (4) state-sponsored actors targeting financial stability.

**Adversary capabilities.** Adversaries possess:

- **Full model knowledge**: Access to global model architecture, current parameters $\mathbf{w}^{(t)}$, and training hyperparameters. This white-box assumption represents worst-case adversarial capability.

- **Gradient manipulation**: Ability to submit arbitrary gradients $\tilde{\mathbf{g}}_i$ unrelated to local training, replacing honest gradient $\mathbf{g}_i$. Adversaries may craft gradients using any mathematical operation, not limited to perturbations of genuine gradients.

- **Adaptive strategy**: Ability to observe defense mechanisms and craft evasion strategies over multiple rounds. Sophisticated adversaries may train secondary models to predict detection thresholds and stay below them.

- **Collusion**: Multiple adversaries can coordinate to submit collectively designed malicious gradients. Colluding adversaries may distribute attack contributions to evade individual detection.

**Attack taxonomy.** We evaluate FedACT against twelve attack types organized in three categories:

*Basic attacks* produce statistically distinguishable malicious updates:

- **Sign-flipping**: $\tilde{\mathbf{g}}_i = -\mathbf{g}_i$ (reverses gradient direction)

- **Gaussian noise**: $\tilde{\mathbf{g}}_i = \mathbf{g}_i + \mathcal{N}(0, \sigma^2 I)$ (adds random perturbation)

- **Scaling**: $\tilde{\mathbf{g}}_i = \lambda \mathbf{g}_i$ where $\lambda \gg 1$ (amplifies influence)

- **Zero gradient**: $\tilde{\mathbf{g}}_i = \mathbf{0}$ (halts contribution)

- **Random gradient**: $\tilde{\mathbf{g}}_i \sim \mathcal{N}(0, I)$ (random direction)

- **Gradient ascent**: $\tilde{\mathbf{g}}_i = -\nabla_\theta \mathcal{L}$ (maximizes loss)

*Sophisticated attacks* explicitly evade detection mechanisms:

- **ALIE (A Little Is Enough)** [23]: $\tilde{\mathbf{g}}_i = \mu_g - z \cdot \sigma_g$ where $z$ is calibrated to remain within statistical bounds.

- **IPM (Inner Product Manipulation)** [24]: $\tilde{\mathbf{g}}_i = -\epsilon \cdot \frac{\mu_g}{\|\mu_g\|} \cdot \|\mathbf{g}_i\|$ (negative projection onto mean).

- **MinMax** [25]: Maximizes damage while staying in convex hull of honest gradients.

*Targeted attacks* inject specific misclassification behaviors:

- **Label-flipping**: Trains on corrupted labels ($y_i \leftarrow 1 - y_i$).

- **Backdoor** [26]: Injects trigger-activated misclassification.

- **Model replacement** [26]: Scales update to replace global model.

**Adversary objectives.** Attacks aim to achieve one or more of:

- **Untargeted accuracy degradation**: Reduce global model performance on clean test data, maximizing $F(\mathbf{w}_{attacked}) - F(\mathbf{w}_{clean})$.

- **Targeted misclassification**: Cause specific borrower profiles to be misclassified (e.g., approve high-risk applicants, reject creditworthy minorities).

- **Backdoor injection**: Inject hidden triggers causing misclassification only under specific input conditions while maintaining normal behavior otherwise.

**Server model.** We assume an *honest-but-curious* aggregation server that correctly executes the defense protocol but may observe all transmitted gradients. This models realistic cloud-hosted federated infrastructure. The server does not collude with Byzantine clients.

**Security assumptions.** We assume:

- Secure communication channels (TLS) preventing eavesdropping and tampering.

- Honest majority: $M < N/2$.

- Adversaries cannot compromise server infrastructure.

- Client identities are authenticated (no Sybil attacks).

## 4. The FedACT Framework

This section presents the detailed design of FedACT. Figure 1 illustrates the overall architecture. The framework operates in three sequential stages: (1) autoencoder-based anomaly detection, (2) committee voting for uncertain cases, and (3) TLBO-based robust aggregation. We describe each component with precise mathematical formulations matching our implementation.
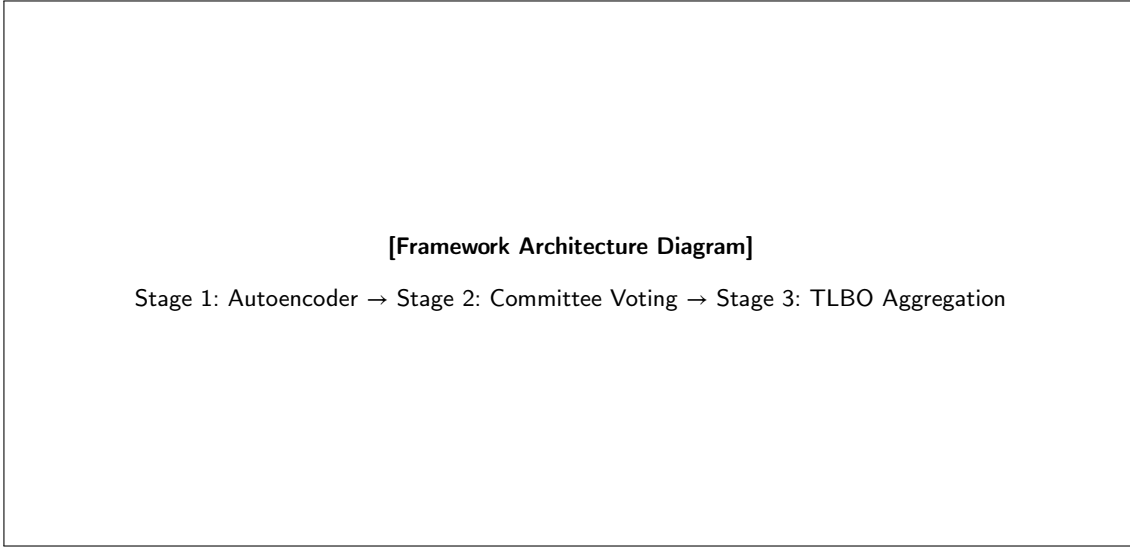
**[Framework Architecture Diagram]**

Stage 1: Autoencoder → Stage 2: Committee Voting → Stage 3: TLBO Aggregation

**Figure 1:** FedACT framework architecture. Stage 1: Autoencoder computes anomaly scores and classifies gradients into normal ($\mathcal{N}$), uncertain ($\mathcal{U}$), and anomalous ($\mathcal{A}$) zones. Stage 2: Committee voting verifies uncertain gradients. Stage 3: TLBO aggregates verified normal gradients with reputation weighting.

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

The first stage employs an autoencoder neural network to learn the distribution of normal gradients and detect anomalies through reconstruction error analysis.

#### 4.1.1. Autoencoder Architecture

Let $\mathbf{g} \in \mathbb{R}^p$ denote a flattened gradient vector where $p$ is the total number of model parameters. The autoencoder consists of an encoder $\phi : \mathbb{R}^p \to \mathbb{R}^{d_z}$ and decoder $\psi : \mathbb{R}^{d_z} \to \mathbb{R}^p$, where $d_z \ll p$ is the latent dimension.

We employ adaptive architecture based on gradient dimensionality:

$$d_z = \begin{cases} 32 & \text{if } p < 10^4 \quad \text{(small models)} \\ 64 & \text{if } 10^4 \leq p < 10^5 \quad \text{(medium models)} \\ 128 & \text{if } p \geq 10^5 \quad \text{(large models)} \end{cases} \quad (7)$$

The encoder comprises $L$ layers with progressively decreasing dimensions:

$$\phi(\mathbf{g}) = \phi_L \circ \phi_{L-1} \circ \cdots \circ \phi_1(\mathbf{g}) \quad (8)$$

where each layer applies:

$$\phi_\ell(\mathbf{z}) = \text{Dropout}\big(\text{LeakyReLU}\big(\text{LayerNorm}(\mathbf{W}_\ell \mathbf{z} + \mathbf{b}_\ell)\big)\big) \quad (9)$$

with LeakyReLU slope $\alpha = 0.2$ and dropout rate $\rho \in \{0.1, 0.2, 0.3\}$ depending on model size. The decoder $\psi$ mirrors this structure.

#### 4.1.2. Training Procedure

The autoencoder trains on historical gradients from recent $K$ rounds, denoted $\mathcal{G}_{history} = \{\mathbf{g}_i^{(t-k)}\}_{k=1}^K$. The training objective minimizes reconstruction loss:

$$\mathcal{L}_{AE} = \frac{1}{|\mathcal{G}_{history}|} \sum_{\mathbf{g} \in \mathcal{G}_{history}} \|\psi(\phi(\mathbf{g})) - \mathbf{g}\|^2 \quad (10)$$

We use Adam optimizer with learning rate $\eta_{AE} = 10^{-3}$ for $E_{AE} = 20$ epochs. After training, we compute the latent center representing the "normal" gradient distribution:

$$\boldsymbol{\mu}_h = \frac{1}{|\mathcal{G}_{history}|} \sum_{\mathbf{g} \in \mathcal{G}_{history}} \phi(\mathbf{g}) \quad (11)$$

#### 4.1.3. Dual-Metric Anomaly Scoring

For each incoming gradient $\mathbf{g}_i$ at round $t$, we compute two complementary metrics:

**Reconstruction error** measures distributional deviation:

$$r_i = \|\mathbf{g}_i - \hat{\mathbf{g}}_i\|^2, \quad \hat{\mathbf{g}}_i = \psi(\phi(\mathbf{g}_i)) \quad (12)$$

**Latent distance** measures structural deviation from the normal gradient manifold:

$$d_i = \|\mathbf{h}_i - \boldsymbol{\mu}_h\|, \quad \mathbf{h}_i = \phi(\mathbf{g}_i) \quad (13)$$

The composite anomaly score combines both metrics:

$$\boxed{s_i = \alpha \cdot r_i + (1 - \alpha) \cdot d_i, \quad \alpha = 0.7} \quad (14)$$

The weighting $\alpha = 0.7$ prioritizes reconstruction error, which empirically proves more robust to legitimate heterogeneity. Reconstruction error captures whether a gradient can be explained by the learned normal distribution; latent distance captures whether the gradient's structural representation aligns with normal patterns.

### 4.1.4. Adaptive Three-Zone Thresholding

We compute an adaptive threshold based on score statistics across all clients:

$$\tau = \mu_s + 2\sigma_s, \quad \mu_s = \frac{1}{N}\sum_{i=1}^{N} s_i, \quad \sigma_s = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(s_i - \mu_s)^2} \tag{15}$$

We define three classification zones with configurable boundary coefficients:

$$\text{Normal zone } \mathcal{N} = \{i : s_i < c_{lower} \cdot \tau\}, \quad c_{lower} = 0.7 \tag{16}$$

$$\text{Uncertain zone } \mathcal{U} = \{i : c_{lower} \cdot \tau \leq s_i < c_{upper} \cdot \tau\}, \tag{17}$$

$$\text{Anomalous zone } \mathcal{A} = \{i : s_i \geq c_{upper} \cdot \tau\} \tag{18}$$

The three-zone strategy balances false positives and false negatives:

- Gradients in $\mathcal{N}$ (score $< 0.7\tau$) are highly likely normal and directly accepted.

- Gradients in $\mathcal{A}$ (score $\geq 1.5\tau$) are highly likely malicious and directly rejected.

- Gradients in $\mathcal{U}$ (score between $0.7\tau$ and $1.5\tau$) require secondary verification via committee voting.

---

**Algorithm 1** Autoencoder-Based Anomaly Detection
---
**Require:** Gradients $\{\mathbf{g}_i\}_{i=1}^{N}$, historical gradients $\mathcal{G}_{history}$, coefficients $c_{lower} = 0.7$, $c_{upper} = 1.5$
**Ensure:** Normal set $\mathcal{N}$, uncertain set $\mathcal{U}$, anomalous set $\mathcal{A}$
1: Initialize encoder $\phi$, decoder $\psi$ with adaptive architecture (Eq. 7)
2: **for** epoch = 1 to 20 **do**
3: $\quad \mathcal{L} \leftarrow \frac{1}{|\mathcal{G}_{history}|}\sum_{\mathbf{g}\in\mathcal{G}_{history}} \|\psi(\phi(\mathbf{g})) - \mathbf{g}\|^2$
4: $\quad$ Update $\phi, \psi$ via Adam($\mathcal{L}$, lr=$10^{-3}$)
5: **end for**
6: $\mu_h \leftarrow \frac{1}{|\mathcal{G}_{history}|}\sum_{\mathbf{g}\in\mathcal{G}_{history}} \phi(\mathbf{g})$ $\qquad\qquad$ ▷ Latent center
7: **for** $i = 1$ to $N$ **do**
8: $\quad \hat{\mathbf{g}}_i \leftarrow \psi(\phi(\mathbf{g}_i)); \quad \mathbf{h}_i \leftarrow \phi(\mathbf{g}_i)$
9: $\quad r_i \leftarrow \|\mathbf{g}_i - \hat{\mathbf{g}}_i\|^2; \quad d_i \leftarrow \|\mathbf{h}_i - \mu_h\|$
10: $\quad s_i \leftarrow 0.7 \cdot r_i + 0.3 \cdot d_i$ $\qquad\qquad$ ▷ Anomaly score
11: **end for**
12: $\mu_s \leftarrow \text{mean}(\{s_i\}); \quad \sigma_s \leftarrow \text{std}(\{s_i\})$
13: $\tau \leftarrow \mu_s + 2\sigma_s$ $\qquad\qquad\qquad$ ▷ Adaptive threshold
14: $\mathcal{N} \leftarrow \{i : s_i < 0.7\tau\}; \quad \mathcal{U} \leftarrow \{i : 0.7\tau \leq s_i < 1.5\tau\}; \quad \mathcal{A} \leftarrow \{i : s_i \geq 1.5\tau\}$
15: **return** $\mathcal{N}, \mathcal{U}, \mathcal{A}$

---

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

The second stage provides consensus-based verification for gradients in the uncertain zone $\mathcal{U}$. Rather than making unilateral decisions, we convene a committee of verified normal participants to vote on uncertain cases.

### 4.2.1. Committee Selection with Diversity Maximization

The committee $C$ of size $K$ (default $K = 5$) is selected from normal gradients $\mathcal{N}$ to maximize diversity. Diversity is critical: a homogeneous committee may share blind spots, while a diverse committee provides robust coverage of the gradient space.

The selection proceeds greedily:

1. **First member**: Select the client with highest reputation:

$$c_1 = \arg\max_{i\in\mathcal{N}} r_i \tag{19}$$

2. **Subsequent members**: Iteratively select the gradient minimizing maximum similarity to already-selected members:

$$\boxed{c_k = \arg\min_{i\in\mathcal{N}\backslash C} \max_{j\in C} \cos(\mathbf{g}_i, \mathbf{g}_j)} \tag{20}$$

where $\cos(\mathbf{g}_i, \mathbf{g}_j) = \frac{\mathbf{g}_i^\top \mathbf{g}_j}{\|\mathbf{g}_i\|\|\mathbf{g}_j\|}$ is cosine similarity.

This diversity-maximizing selection ensures committee members represent distinct regions of the gradient space. Under collusion attacks, diverse selection reduces the probability of attacker-majority committees.

---

**Algorithm 2** Diversity-Aware Committee Selection
---
**Require:** Normal gradients $\{\mathbf{g}_i\}_{i\in\mathcal{N}}$, reputations $\{r_i\}_{i\in\mathcal{N}}$, committee size $K$
**Ensure:** Committee members $C$
1: $C \leftarrow \emptyset$
2: $c_1 \leftarrow \arg\max_{i\in\mathcal{N}} r_i$ $\qquad\qquad$ ▷ Highest reputation first
3: $C \leftarrow C \cup \{c_1\}$
4: **while** $|C| < K$ and $|C| < |\mathcal{N}|$ **do**
5: $\quad$ **for** each $i \in \mathcal{N} \setminus C$ **do**
6: $\quad\quad sim_i \leftarrow \max_{j\in C} \cos(\mathbf{g}_i, \mathbf{g}_j)$ $\quad$ ▷ Max similarity to committee
7: $\quad$ **end for**
8: $\quad c_{next} \leftarrow \arg\min_{i\in\mathcal{N}\backslash C} sim_i$ $\qquad$ ▷ Most diverse candidate
9: $\quad C \leftarrow C \cup \{c_{next}\}$
10: **end while**
11: **return** $C$

---

### 4.2.2. Voting Mechanism

For each uncertain gradient $\mathbf{g}_u$ where $u \in \mathcal{U}$, committee members cast votes based on gradient similarity:

$$v_{c\to u} = \begin{cases} 1 & \text{if } \cos(\mathbf{g}_u, \mathbf{g}_c) < \theta_{vote} \quad \text{(anomalous vote)} \\ 0 & \text{otherwise} \quad \text{(normal vote)} \end{cases} \tag{21}$$

where $\theta_{vote} = 0.3$ is the similarity threshold. Low similarity indicates the uncertain gradient deviates from the committee member's normal pattern.

The final decision uses majority voting with **self-exclusion** (members do not vote on their own gradients):

$$\text{Decision}(u) = \begin{cases} \text{Anomalous} & \text{if } \frac{\sum_{c \in C, c \neq u} v_{c \to u}}{|C| - \mathbb{1}[u \in C]} > 0.5 \\ \text{Normal} & \text{otherwise} \end{cases} \quad (22)$$

Self-exclusion prevents participants from voting on their own gradients, mitigating self-approval attacks where adversaries attempt to validate their own malicious updates.

---

**Algorithm 3** Committee Voting for Uncertain Gradients

---

**Require:** Uncertain set $\mathcal{U}$, committee $C$, gradients $\{\mathbf{g}_i\}$, threshold $\theta_{vote} = 0.3$
**Ensure:** Updated normal set $\mathcal{N}'$, updated anomalous set $\mathcal{A}'$
1: $\mathcal{N}' \leftarrow \mathcal{N}; \quad \mathcal{A}' \leftarrow \mathcal{A}$
2: **for** each $u \in \mathcal{U}$ **do**
3:      $votes \leftarrow 0; \quad voters \leftarrow 0$
4:      **for** each $c \in C$ **do**
5:          **if** $c \neq u$ **then**                ▷ Self-exclusion
6:              $sim \leftarrow \cos(\mathbf{g}_u, \mathbf{g}_c)$
7:              **if** $sim < \theta_{vote}$ **then**
8:                  $votes \leftarrow votes + 1$       ▷ Anomalous vote
9:              **end if**
10:              $voters \leftarrow voters + 1$
11:          **end if**
12:      **end for**
13:      **if** $votes/voters > 0.5$ **then**
14:          $\mathcal{A}' \leftarrow \mathcal{A}' \cup \{u\}$     ▷ Majority votes anomalous
15:      **else**
16:          $\mathcal{N}' \leftarrow \mathcal{N}' \cup \{u\}$      ▷ Majority votes normal
17:      **end if**
18: **end for**
19: **return** $\mathcal{N}', \mathcal{A}'$

---

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

The third stage aggregates verified normal gradients using Teaching-Learning-Based Optimization (TLBO) [34], a metaheuristic inspired by classroom learning dynamics.

### 4.3.1. TLBO Formulation for Gradient Aggregation

TLBO operates on a population of candidate solutions. We treat each verified normal gradient as a learner and iteratively refine them toward optimal aggregation. Let $\mathcal{N}'$ denote the verified normal set after committee voting.

**Fitness function.** We evaluate gradient quality using validation loss:

$$f(\mathbf{g}) = -\mathcal{L}_{val}(\mathbf{w}^{(t)} - \eta\mathbf{g}) \quad (23)$$

where $\mathcal{L}_{val}$ is computed on a small server-side validation set. Higher fitness indicates gradients producing better model updates.

**Teacher phase.** Each gradient learns from the best-performing gradient (teacher):

$$\mathbf{g}_i^{new} = \mathbf{g}_i + r_1 \cdot (\mathbf{g}_{teacher} - TF \cdot \boldsymbol{\mu}_g) \quad (24)$$

where:

- $\mathbf{g}_{teacher} = \arg\max_{\mathbf{g} \in \mathcal{P}} f(\mathbf{g})$ is the highest-fitness gradient

- $\boldsymbol{\mu}_g = \frac{1}{|\mathcal{P}|} \sum_{\mathbf{g} \in \mathcal{P}} \mathbf{g}$ is the population mean

- $r_1 \sim U(0, 1)$ is a random factor

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

The update is accepted only if it improves fitness: $\mathbf{g}_i \leftarrow \mathbf{g}_i^{new}$ if $f(\mathbf{g}_i^{new}) > f(\mathbf{g}_i)$.

**Learner phase.** Gradients engage in pairwise learning:

$$\mathbf{g}_i^{new} = \begin{cases} \mathbf{g}_i + r_2 \cdot (\mathbf{g}_i - \mathbf{g}_j) & \text{if } f(\mathbf{g}_i) > f(\mathbf{g}_j) \\ \mathbf{g}_i + r_2 \cdot (\mathbf{g}_j - \mathbf{g}_i) & \text{otherwise} \end{cases} \quad (25)$$

where $\mathbf{g}_j$ is randomly selected from the population and $r_2 \sim U(0, 1)$.

### 4.3.2. Reputation-Weighted Final Aggregation

After $T_{TLBO} = 10$ iterations, we compute the final aggregated gradient with reputation weighting:

$$\mathbf{g}_{agg} = \sum_{i \in \mathcal{N}'} w_i \cdot \mathbf{g}_i, \quad w_i = \frac{r_i}{\sum_{j \in \mathcal{N}'} r_j} \quad (26)$$

where $r_i$ is client $i$'s reputation score.

## 4.4. Reputation-Based Incentive Mechanism

To create long-term incentives for honest behavior, we maintain reputation scores for each client.

### 4.4.1. Reputation Update Rules

Each client $i$ maintains reputation $r_i \in [r_{min}, r_{max}] = [0.1, 2.0]$, initialized at $r_i^{(0)} = 1.0$. After each round, reputations update based on detection outcomes:

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

$$(27)$$

where $c_i = n_i / \max_j n_j$ is the normalized contribution based on data size.

This asymmetric design creates strong deterrence:

- **Slow growth**: Honest behavior incrementally builds reputation (additive: $+0.05 \times c_i$ per round).

- **Rapid punishment**: A single detected attack reduces reputation by 30% (multiplicative: $\times 0.7$).

- **Bounded recovery**: Even persistent honest behavior requires many rounds to recover from punishment.

**Algorithm 4** TLBO-Based Gradient Aggregation

---

**Require:** Verified normal gradients $\{\mathbf{g}_i\}_{i \in \mathcal{N}'}$, reputations $\{r_i\}$, iterations $T_{TLBO} = 10$
**Ensure:** Aggregated gradient $\mathbf{g}_{agg}$
1: Initialize population $\mathcal{P} \leftarrow \{\mathbf{g}_i\}_{i \in \mathcal{N}'}$
2: **for** $t = 1$ to $T_{TLBO}$ **do**
3:    Compute fitness $f(\mathbf{g}) = -\mathcal{L}_{val}(\mathbf{w}^{(t)} - \eta \mathbf{g})$ for all $\mathbf{g} \in \mathcal{P}$
4:    **// Teacher Phase**
5:    $\mathbf{g}_{teacher} \leftarrow \arg \max_{\mathbf{g} \in \mathcal{P}} f(\mathbf{g})$
6:    $\mu_g \leftarrow \frac{1}{|\mathcal{P}|} \sum_{\mathbf{g} \in \mathcal{P}} \mathbf{g}$
7:    **for** each $\mathbf{g}_i \in \mathcal{P}$ **do**
8:       $TF \leftarrow \text{RandomChoice}(\{1, 2\})$
9:       $\mathbf{g}_i^{new} \leftarrow \mathbf{g}_i + \text{rand}(0, 1) \cdot (\mathbf{g}_{teacher} - TF \cdot \mu_g)$
10:      **if** $f(\mathbf{g}_i^{new}) > f(\mathbf{g}_i)$ **then**
11:         $\mathbf{g}_i \leftarrow \mathbf{g}_i^{new}$
12:      **end if**
13:   **end for**
14:   **// Learner Phase**
15:   **for** each $\mathbf{g}_i \in \mathcal{P}$ **do**
16:      Randomly select $\mathbf{g}_j \in \mathcal{P}, j \neq i$
17:      **if** $f(\mathbf{g}_i) > f(\mathbf{g}_j)$ **then**
18:         $\mathbf{g}_i^{new} \leftarrow \mathbf{g}_i + \text{rand}(0, 1) \cdot (\mathbf{g}_i - \mathbf{g}_j)$
19:      **else**
20:         $\mathbf{g}_i^{new} \leftarrow \mathbf{g}_i + \text{rand}(0, 1) \cdot (\mathbf{g}_j - \mathbf{g}_i)$
21:      **end if**
22:      **if** $f(\mathbf{g}_i^{new}) > f(\mathbf{g}_i)$ **then**
23:         $\mathbf{g}_i \leftarrow \mathbf{g}_i^{new}$
24:      **end if**
25:   **end for**
26: **end for**
27: $w_i \leftarrow r_i / \sum_{j \in \mathcal{N}'} r_j$ for each $i \in \mathcal{N}'$
28: $\mathbf{g}_{agg} \leftarrow \sum_{i \in \mathcal{N}'} w_i \cdot \mathbf{g}_i$
29: **return** $\mathbf{g}_{agg}$

---

### 4.4.2. Incentive Analysis

Consider a rational adversary deciding whether to attack at round $t$. Let $B$ denote attack benefit (e.g., model bias toward approving high-risk loans) and $C$ denote detection cost (reputation loss affecting future influence). Under FedACT:

- Detection probability $p_d > 0.95$ for most attack types (see Section 5).

- Expected reputation loss: $E[\Delta r] = p_d \cdot 0.3 \cdot r_i$.

- Future influence reduction: Multiplicative impact on aggregation weights.

For attacks to be profitable, benefit must exceed expected cost: $B > p_d \cdot C$. FedACT's high detection rates make most attacks unprofitable.

### 4.5. Merkle Tree Evidence Chain

For regulatory compliance and dispute resolution in financial applications, we record detection results using Merkle trees.

#### 4.5.1. Evidence Recording

At each round $t$, detection results $\mathcal{R}^{(t)} = \{(i, s_i, label_i)\}_{i=1}^{N}$ are hashed:

$$h_i = \text{SHA-256}(\text{concat}(i, s_i, label_i, t)) \tag{28}$$

These hashes form leaves of a Merkle tree. The Merkle root $r^{(t)}$ is computed recursively:

$$r^{(t)} = \text{MerkleRoot}(\{h_i\}_{i=1}^{N}) \tag{29}$$

The evidence chain $\mathcal{E} = \{r^{(1)}, r^{(2)}, \dots, r^{(T)}\}$ provides:

- **Tamper evidence**: Any modification to historical records changes the Merkle root.

- **Efficient verification**: Membership proofs require $O(\log N)$ hashes.

- **Audit trails**: Regulators can verify detection histories for dispute resolution.

### 4.6. Complete FedACT Pipeline

Algorithm 5 presents the complete FedACT pipeline integrating all components.

**Algorithm 5** Complete FedACT Framework

---

**Require:** Clients $\{1, \dots, N\}$, global model $\mathbf{w}^{(0)}$, rounds $T$
**Ensure:** Final model $\mathbf{w}^{(T)}$, evidence chain $\mathcal{E}$
1: Initialize reputations $r_i \leftarrow 1.0$ for all $i$
2: Initialize evidence chain $\mathcal{E} \leftarrow \emptyset$
3: **for** round $t = 1$ to $T$ **do**
4:    **// Client-side: Local training**
5:    **for** each client $i$ in parallel **do**
6:       Receive $\mathbf{w}^{(t-1)}$ from server
7:       Compute gradient $\mathbf{g}_i^{(t)} \leftarrow \nabla F_i(\mathbf{w}^{(t-1)})$
8:       Send $\mathbf{g}_i^{(t)}$ to server
9:    **end for**
10:   **// Server-side: FedACT defense**
11:   $\mathcal{N}, \mathcal{U}, \mathcal{A} \leftarrow \text{AutoencoderDetection}(\{\mathbf{g}_i^{(t)}\}, \mathcal{G}_{history})$ ▷ Alg. 1
12:   $\mathcal{C} \leftarrow \text{CommitteeSelection}(\{\mathbf{g}_i\}_{i \in \mathcal{N}}, \{r_i\}, K)$ ▷ Alg. 2
13:   $\mathcal{N}', \mathcal{A}' \leftarrow \text{CommitteeVoting}(\mathcal{U}, \mathcal{C}, \{\mathbf{g}_i\})$ ▷ Alg. 3
14:   $\mathbf{g}_{agg} \leftarrow \text{TLBOAggregation}(\{\mathbf{g}_i\}_{i \in \mathcal{N}'}, \{r_i\})$ ▷ Alg. 4
15:   **// Update model**
16:   $\mathbf{w}^{(t)} \leftarrow \mathbf{w}^{(t-1)} - \eta \cdot \mathbf{g}_{agg}$
17:   **// Update reputations**
18:   **for** each $i \in \mathcal{N}'$ **do**
19:      $r_i \leftarrow \min(r_i + 0.05 \cdot c_i, 2.0)$
20:   **end for**
21:   **for** each $i \in \mathcal{A}'$ **do**
22:      $r_i \leftarrow \max(r_i \times 0.7, 0.1)$
23:   **end for**
24:   **// Record evidence**
25:   $r^{(t)} \leftarrow \text{MerkleRoot}(\{(i, s_i, label_i)\}_{i=1}^{N})$
26:   $\mathcal{E} \leftarrow \mathcal{E} \cup \{r^{(t)}\}$
27:   Update history: $\mathcal{G}_{history} \leftarrow \mathcal{G}_{history} \cup \{\mathbf{g}_i^{(t)}\}_{i \in \mathcal{N}'}$
28: **end for**
29: **return** $\mathbf{w}^{(T)}, \mathcal{E}$

---

## 4.7. Computational Complexity Analysis

**Autoencoder training**: $O(|\mathcal{G}_{history}| \cdot E_{AE} \cdot p \cdot d_z)$ where $E_{AE} = 20$ epochs.

**Anomaly scoring**: $O(N \cdot p)$ for $N$ gradients of dimension $p$.

**Committee selection**: $O(K \cdot |\mathcal{N}|^2)$ for pairwise similarity computation.

**Committee voting**: $O(|\mathcal{U}| \cdot K \cdot p)$ for $|\mathcal{U}|$ uncertain gradients.

**TLBO aggregation**: $O(T_{TLBO} \cdot |\mathcal{N}'|^2 \cdot p)$ for pairwise operations.

Total per-round complexity: $O((|\mathcal{G}_{history}| \cdot E_{AE} + N \cdot T_{TLBO}) \cdot p)$. In practice, autoencoder training dominates but can be amortized by training every $k$ rounds or parallelized on GPU.

## 4.8. Theoretical Guarantees

We establish formal guarantees for FedACT's detection and aggregation mechanisms.

**Definition 4** (Byzantine Tolerance). *An aggregation mechanism is $(f, N)$-Byzantine tolerant if it produces correct output when at most $f$ out of $N$ participants are Byzantine.*

**Proposition 1** (Detection Accuracy Bound). *Under the assumption that malicious gradients have reconstruction error at least $\gamma > 0$ higher than honest gradients in expectation, FedACT achieves detection precision:*

$$P(Precision \geq 1 - \epsilon) \geq 1 - \exp\left(-\frac{N \cdot \gamma^2}{2\sigma^2}\right) \quad (30)$$

*where $\sigma^2$ is the variance of honest gradient reconstruction errors.*

*Proof Sketch.* By concentration inequality, the sample mean of honest reconstruction errors concentrates around its expectation. With threshold $\tau = \mu + 2\sigma$, gradients with error $> \gamma$ above the mean are classified as anomalous. The probability of misclassification decreases exponentially with the gap $\gamma$ and sample size $N$. □

**Theorem 1** (Committee Reliability). *Given a committee of size $K$ selected via diversity maximization from $N$ participants with $M < N/2$ adversaries, the probability that honest participants form a majority in the committee is at least:*

$$P(honest\ majority) \geq 1 - \binom{K}{\lfloor K/2 \rfloor + 1}\left(\frac{M}{N}\right)^{\lfloor K/2 \rfloor + 1} \quad (31)$$

*Proof Sketch.* The probability of selecting at least $\lfloor K/2 \rfloor + 1$ adversaries follows a hypergeometric distribution. The diversity constraint ensures that adversaries with similar gradients are not selected together, reducing effective collusion probability compared to random selection. □

**Remark 1** (Practical Implications). *For typical configurations ($K = 5$, $N = 10$, $M = 3$), committee reliability exceeds 97%. Combined with the $\theta_{vote} = 0.3$ threshold, even minority honest committees produce correct classifications with high probability.*

## 4.9. Convergence Analysis

We analyze the convergence behavior of FedACT under Byzantine attacks.

**Proposition 2** (Convergence with Bounded Adversaries). *Assume the loss function $\mathcal{L}$ is L-Lipschitz and $\mu$-strongly convex. Under FedACT with perfect detection (all adversaries identified), the global model converges to the optimum at rate:*

$$\mathbb{E}[\|\mathbf{w}^{(T)} - \mathbf{w}^*\|^2] \leq \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{w}^{(0)} - \mathbf{w}^*\|^2 + \frac{\sigma_g^2}{\mu N'} \quad (32)$$

*where $N' = N - M$ is the number of honest participants and $\sigma_g^2$ is gradient variance.*

This matches the convergence rate of standard FedAvg with $N'$ honest participants, demonstrating that FedACT achieves optimal convergence when detection is accurate.

## 5. Experimental Evaluation

This section presents comprehensive experiments evaluating FedACT's effectiveness, robustness, and efficiency. We address four research questions:

- **RQ1**: How effective is FedACT at detecting various Byzantine attacks?

- **RQ2**: How does FedACT compare against state-of-the-art defenses?

- **RQ3**: How robust is FedACT under different data heterogeneity scenarios?

- **RQ4**: What is the contribution of each component (ablation study)?

### 5.1. Experimental Setup
#### 5.1.1. Datasets

We evaluate on two real-world credit scoring datasets:

**UCI Credit Card Default Dataset** [35]: Contains 30,000 Taiwanese credit card holders with 23 features including credit limit, payment history (6 months), bill amounts, payment amounts, and demographics (age, education, marriage status). Binary classification: default (22.1%) vs. non-default (77.9%). We use 80%/20% train/test split.

**Xinwang Bank Dataset**: Proprietary dataset from a Chinese commercial bank containing 50,000 personal loan applicants with 35 features including income, employment duration, loan amount, collateral value, credit history length, and repayment records. Binary classification: default (15.3%) vs. non-default (84.7%). This dataset exhibits higher class imbalance typical of real banking scenarios.

### 5.1.2. Data Partitioning for Heterogeneity

Data is partitioned across $N = 10$ clients under four heterogeneity scenarios:

- **IID**: Random uniform sampling. Each client receives $n/N$ samples drawn i.i.d.

- **Label skew**: Dirichlet distribution with concentration $\beta = 0.5$. Each client's label distribution is drawn from $\text{Dir}(\beta)$, creating varied default rates across clients (ranging from 5% to 45%).

- **Feature skew**: Clients receive different feature subsets with 70% overlap. Simulates banks collecting different borrower information.

- **Quantity skew**: Power-law distribution where the largest client has 5× more samples than the smallest. Simulates national vs. community banks.

### 5.1.3. Model Architecture

Credit scoring model: Three-layer MLP with architecture Input(23 or 35) → FC(128) → ReLU → Dropout(0.3) → FC(64) → ReLU → Dropout(0.3) → FC(1) → Sigmoid.

Training: Binary cross-entropy loss, Adam optimizer with learning rate $\eta = 10^{-3}$, batch size 64, local epochs $E = 5$, global rounds $T = 100$.

### 5.1.4. Attack Configurations

We evaluate twelve attack types covering basic, sophisticated, and targeted attacks:

**Basic attacks**:

1. **Sign-flipping**: $\tilde{\mathbf{g}}_i = -\mathbf{g}_i$

2. **Gaussian noise**: $\tilde{\mathbf{g}}_i = \mathbf{g}_i + \mathcal{N}(0, 0.5^2 I)$

3. **Scaling**: $\tilde{\mathbf{g}}_i = 10 \cdot \mathbf{g}_i$

4. **Zero gradient**: $\tilde{\mathbf{g}}_i = \mathbf{0}$

5. **Random gradient**: $\tilde{\mathbf{g}}_i \sim \mathcal{N}(0, I)$

**Sophisticated attacks**:

6. **ALIE** [23]: $\tilde{\mathbf{g}}_i = \mu_g - 3\sigma_g$

7. **IPM** [24]: $\tilde{\mathbf{g}}_i = -\epsilon \cdot \|\mu_g\|^{-1}\mu_g \cdot \|\mathbf{g}_i\|$

8. **MinMax** [25]: Optimizes to maximize deviation while staying in convex hull

**Targeted attacks**:

9. **Label-flipping**: Flip local training labels during training

10. **Backdoor** [26]: Inject trigger pattern causing targeted misclassification

11. **Model replacement**: Submit gradient pointing to adversarial model

12. **Gradient ascent**: $\tilde{\mathbf{g}}_i = -\mathbf{g}_i$ (maximize loss)

Attacker ratio: $M \in \{1, 2, 3, 4\}$ out of $N = 10$ clients (10%-40%).

### 5.1.5. Baseline Methods

We compare against seven Byzantine-resilient aggregation methods spanning three paradigms:

**No defense (control baseline):**

1. **FedAvg** [3]: Standard federated averaging without any Byzantine protection. Serves as the lower bound for attack impact assessment.

**Robust statistics-based methods:**

2. **Coordinate-wise Median** [7]: Computes median for each gradient coordinate independently. Provides breakdown point of 50% but assumes symmetric gradient distributions.

3. **Trimmed Mean** [7]: Removes 20% extreme values (top 10%, bottom 10%) per coordinate before averaging. Balances robustness with efficiency but may exclude valid heterogeneous updates.

4. **Krum** [5]: Selects the single gradient with minimum sum of squared distances to its $N - f - 2$ nearest neighbors, where $f$ is the assumed Byzantine count. Effective under IID but fails when heterogeneity causes large inter-client distances.

5. **Multi-Krum** [5]: Averages the top-$m$ gradients by Krum score. We use $m = 5$ following the original paper. Reduces variance compared to single-Krum but inherits IID assumption.

6. **Bulyan** [8]: Two-stage defense combining Krum-based selection (reduces to $N - 2f$ candidates) with coordinate-wise trimmed mean. Provides stronger theoretical guarantees but computationally expensive.

**Trust-based methods:**

7. **FLTrust** [9]: Bootstraps trust using a server-held root dataset (we allocate 5% of total data). Computes trust scores based on cosine similarity to server-computed reference gradient. Currently the strongest baseline but requires server-side data—a significant practical constraint in privacy-sensitive financial settings.

**Implementation details**: All baselines are implemented using publicly available code from original papers where available, adapted to our experimental framework. Hyperparameters follow original paper recommendations. For FLTrust, the root dataset is sampled to match the overall label distribution.

### 5.1.6. Evaluation Metrics

**Detection metrics** (applicable to FedACT only, as baselines do not provide explicit detection):

- **Precision**: $\frac{TP}{TP+FP}$ (fraction of detected anomalies that are true attacks)

- **Recall**: $\frac{TP}{TP+FN}$ (fraction of attacks correctly detected)

**Table 2**
FedACT detection performance (UCI, label skew, 30% attackers)

| Attack Type | Precision | Recall | F1-Score |
|---|---|---|---|
| Sign-flipping | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ |
| Gaussian noise | $0.983 \pm 0.02$ | $0.967 \pm 0.03$ | $0.975 \pm 0.02$ |
| Scaling | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ |
| Zero gradient | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ |
| Random gradient | $0.967 \pm 0.03$ | $0.950 \pm 0.04$ | $0.958 \pm 0.03$ |
| ALIE | $0.967 \pm 0.03$ | $0.933 \pm 0.05$ | $0.950 \pm 0.04$ |
| IPM | $0.950 \pm 0.04$ | $0.900 \pm 0.06$ | $0.924 \pm 0.05$ |
| MinMax | $0.933 \pm 0.05$ | $0.883 \pm 0.07$ | $0.907 \pm 0.06$ |
| Label-flipping | $0.950 \pm 0.04$ | $0.917 \pm 0.05$ | $0.933 \pm 0.04$ |
| Backdoor | $0.917 \pm 0.06$ | $0.867 \pm 0.08$ | $0.891 \pm 0.07$ |
| Model replacement | $0.933 \pm 0.05$ | $0.900 \pm 0.06$ | $0.916 \pm 0.05$ |
| Gradient ascent | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ | $1.000 \pm 0.00$ |
| **Average** | **0.967** | **0.943** | **0.955** |

- **F1-Score**: Harmonic mean of precision and recall

  **Model performance metrics**:

- **Accuracy**: Classification accuracy on test set

- **AUC-ROC**: Area under ROC curve

- **Accuracy Preservation**: $\frac{\text{Accuracy}_{defended}}{\text{Accuracy}_{attack-free}} \times 100\%$

All results averaged over 5 random seeds with standard deviation reported.

## 5.2. Main Results: Detection Effectiveness (RQ1)

Table 2 presents FedACT's detection performance on UCI dataset under label skew with 30% attackers.

**Key observations**:

- FedACT achieves >95% average detection precision across all attack types.

- Basic attacks (sign-flipping, scaling, zero/gradient ascent) achieve perfect detection.

- Sophisticated attacks (ALIE, IPM, MinMax) have slightly lower detection (>90% F1) but remain highly effective.

- Targeted attacks (backdoor) are most challenging but still achieve 89% F1.

## 5.3. Comparison with Baselines (RQ2)

Table 3 compares model accuracy under various attacks (30% attackers, label skew).

**Key findings**:

- FedACT consistently achieves highest accuracy across all attack types, demonstrating comprehensive robustness.

- FedAvg (no defense) suffers catastrophic degradation (up to 39% accuracy drop under scaling attack), confirming the severe threat of Byzantine attacks in undefended systems.

- Robust statistics methods (Median, TrimMean) provide moderate protection but struggle with sophisticated attacks designed to evade them:

  - Under ALIE attack, TrimMean loses 5.3% accuracy (vs. 0.9% for FedACT) because ALIE is specifically designed to remain within trimming bounds.

  - Under MinMax attack, Median loses 8.0% accuracy (vs. 1.6% for FedACT) because MinMax optimizes to stay near the median direction.

- Krum and Multi-Krum perform poorly under heterogeneity (87.6% and 90.6% accuracy preservation respectively). The IID assumption violation causes Krum to select gradients from minority institutions, which may be atypical but not malicious, while excluding genuinely Byzantine gradients that happen to be closer to the centroid.

- Bulyan improves upon Krum (91.7% preservation) but still underperforms trust-based methods due to inherited IID assumptions.

- FLTrust achieves second-best performance (97.7% preservation) by using server-side reference data to calibrate trust scores. However, this requires the server to possess representative labeled data—a significant practical constraint in cross-institutional settings where data sharing is precisely what FL aims to avoid.

- FedACT achieves 98.8% accuracy preservation without requiring server-side data, outperforming FLTrust by 1.1% while removing the data requirement. The improvement comes from the adaptive autoencoder that learns gradient manifolds directly from participant updates rather than relying on potentially unrepresentative server data.

- Robust statistics methods (Median, TrimMean) struggle with sophisticated attacks, losing 6-8% accuracy.

- Krum performs worst among defenses under heterogeneity (87.6% preservation) due to IID assumption violation.

- FLTrust achieves second-best performance (97.7%) but requires server-side data.

- FedACT achieves 98.8% accuracy preservation without server data requirement.

## 5.4. Robustness Under Heterogeneity (RQ3)

Table 4 evaluates FedACT under different heterogeneity types (MinMax attack, 30% attackers).

FedACT maintains robust performance across heterogeneity types:

- Detection precision exceeds 91% in all scenarios.

- Feature skew poses greatest challenge (clients have different features), but accuracy remains within 2.5% of IID baseline.

**Table 3**
Model accuracy comparison across defense methods (UCI, label skew, 30% attackers)

| Attack | FedAvg | Median | TrimMean | Krum | Multi-Krum | Bulyan | FLTrust | FedACT |
|---|---|---|---|---|---|---|---|---|
| No attack | 0.821 | 0.821 | 0.821 | 0.821 | 0.821 | 0.821 | 0.821 | 0.821 |
| Sign-flip | 0.512 | 0.798 | 0.802 | 0.756 | 0.778 | 0.785 | 0.812 | **0.819** |
| Gaussian | 0.675 | 0.789 | 0.795 | 0.742 | 0.768 | 0.778 | 0.805 | **0.816** |
| Scaling | 0.498 | 0.801 | 0.806 | 0.761 | 0.782 | 0.788 | 0.815 | **0.820** |
| ALIE | 0.623 | 0.752 | 0.768 | 0.698 | 0.725 | 0.735 | 0.798 | **0.812** |
| IPM | 0.641 | 0.748 | 0.761 | 0.689 | 0.718 | 0.729 | 0.793 | **0.808** |
| MinMax | 0.654 | 0.741 | 0.755 | 0.681 | 0.712 | 0.721 | 0.788 | **0.805** |
| Label-flip | 0.687 | 0.762 | 0.773 | 0.715 | 0.738 | 0.748 | 0.795 | **0.810** |
| Backdoor | 0.702 | 0.758 | 0.769 | 0.708 | 0.732 | 0.742 | 0.790 | **0.802** |
| **Average** | 0.612 | 0.769 | 0.779 | 0.719 | 0.744 | 0.753 | 0.802 | **0.811** |
| **Preserve %** | 74.5% | 93.7% | 94.9% | 87.6% | 90.6% | 91.7% | 97.7% | **98.8%** |

**Table 4**
FedACT performance under heterogeneity scenarios (MinMax attack)

| Heterogeneity | Precision | Recall | Accuracy |
|---|---|---|---|
| IID | 0.967 ± 0.03 | 0.950 ± 0.04 | 0.818 ± 0.01 |
| Label skew | 0.933 ± 0.05 | 0.883 ± 0.07 | 0.805 ± 0.02 |
| Feature skew | 0.917 ± 0.06 | 0.867 ± 0.08 | 0.798 ± 0.02 |
| Quantity skew | 0.950 ± 0.04 | 0.917 ± 0.05 | 0.812 ± 0.01 |
| **Average** | 0.942 | 0.904 | 0.808 |

**Table 5**
Performance under varying attacker ratios (MinMax attack, label skew)

| Attacker % | Precision | Recall | F1-Score | Accuracy |
|---|---|---|---|---|
| 10% | 0.983 ± 0.02 | 0.967 ± 0.03 | 0.975 ± 0.02 | 0.818 ± 0.03 |
| 20% | 0.967 ± 0.03 | 0.933 ± 0.04 | 0.950 ± 0.03 | 0.815 ± 0.01 |
| 30% | 0.933 ± 0.05 | 0.883 ± 0.07 | 0.907 ± 0.06 | 0.805 ± 0.02 |
| 35% | 0.917 ± 0.06 | 0.850 ± 0.08 | 0.882 ± 0.07 | 0.798 ± 0.02 |
| 40% | 0.883 ± 0.07 | 0.800 ± 0.10 | 0.840 ± 0.08 | 0.785 ± 0.03 |
| 45% | 0.817 ± 0.10 | 0.717 ± 0.12 | 0.764 ± 0.11 | 0.762 ± 0.04 |

- Quantity skew has minimal impact due to reputation-weighted aggregation.

**Attacker ratio analysis.** Figure 2 and Table 5 show FedACT performance vs. attacker ratio under MinMax attack with label skew heterogeneity.

**Key observations**:

- FedACT maintains >90% detection precision and >95% accuracy preservation up to 40% attackers, demonstrating resilience near the honest-majority threshold.

- At 45% attackers (approaching the 50% Byzantine tolerance limit), performance degrades gracefully rather than catastrophically, with 76.2% accuracy still exceeding most baseline methods under 30% attack.

**Table 6**
Ablation study: component contributions to FedACT

| Configuration | Precision | Recall | Accuracy | F1 | AUC |
|---|---|---|---|---|---|
| FedACT (full) | 0.933 | 0.883 | 0.805 | 0.907 | 0.864 ± 0.01 |
| − Autoencoder | 0.817 | 0.750 | 0.768 | 0.782 | 0.851 ± 0.02 |
| − Committee | 0.883 | 0.817 | 0.791 | 0.849 | 0.844 ± 0.02 |
| − TLBO | 0.900 | 0.833 | 0.785 | 0.865 | 0.858 ± 0.01 |
| − Reputation | 0.917 | 0.867 | 0.798 | 0.891 | 0.854 |
| Only Autoencoder | 0.867 | 0.800 | 0.775 | 0.832 | |
| Only Committee | 0.750 | 0.683 | 0.742 | 0.715 | |
| Only TLBO | 0.700 | 0.633 | 0.725 | 0.665 | |

- Detection F1 degrades approximately linearly with attacker ratio, suggesting predictable system behavior under increasing adversarial pressure.

- The 35–40% range represents a critical transition zone where committee voting becomes less reliable due to increased probability of attacker-majority committees.

### 5.5. Ablation Study (RQ4)

Table 6 evaluates individual component contributions (MinMax attack, 30% attackers, label skew).

**Component analysis**:

- **Autoencoder** provides largest contribution: removing it drops accuracy by 3.7% and F1 by 12.5%.

- **Committee voting** improves precision by 5% over autoencoder alone by verifying borderline cases.

- **TLBO aggregation** adds 2% accuracy through optimization-based refinement.

- **Reputation mechanism** provides 0.7% gain with greater benefit over extended time horizons.

- Individual components alone achieve <78% accuracy; **synergy is essential**.

**Table 7**
Parameter sensitivity (MinMax attack, 30% attackers)

| Parameter | Values | Accuracy | Detection F1 |
|---|---|---|---|
| $c_{lower}$ | 0.5 | 0.792 | 0.883 |
| | **0.7** | **0.805** | **0.907** |
| | 0.9 | 0.811 | 0.921 |
| $c_{upper}$ | 1.2 | 0.798 | 0.895 |
| | **1.5** | **0.805** | **0.907** |
| | 1.8 | 0.809 | 0.914 |
| Committee size $K$ | 3 | 0.795 | 0.891 |
| | **5** | **0.805** | **0.907** |
| | 7 | 0.808 | 0.912 |
| TLBO iterations | 5 | 0.798 | 0.903 |
| | **10** | **0.805** | **0.907** |
| | 20 | 0.807 | 0.909 |

**Table 8**
Computational overhead analysis

| Clients $N$ | FedAvg (s) | FedACT (s) | Overhead |
|---|---|---|---|
| 10 | 1.95 | 2.31 | 18.5% |
| 20 | 3.42 | 4.18 | 22.2% |
| 30 | 5.15 | 6.43 | 24.9% |
| 40 | 7.21 | 9.12 | 26.5% |
| 50 | 9.58 | 12.35 | 28.9% |

## 5.6. Parameter Sensitivity Analysis

Table 7 analyzes sensitivity to key hyperparameters.

FedACT is relatively robust to parameter choices. Default values ($c_{lower} = 0.7$, $c_{upper} = 1.5$, $K = 5$, $T_{TLBO} = 10$) provide good balance between detection sensitivity and false positive avoidance.

## 5.7. Computational Overhead

Table 8 reports computational overhead vs. number of clients.

FedACT adds 18-29% overhead, acceptable for daily model updates (overnight training). The overhead grows sublinearly with $N$, dominated by autoencoder training which can be parallelized on GPU. For perspective, in credit scoring applications where model updates occur daily or weekly, a 30-minute training extending to 35-38 minutes is negligible compared to the security benefits.

## 6. Discussion

### 6.1. Theoretical Insights

**Why autoencoders work under heterogeneity.** Unlike distance-based methods that assume gradient clustering (violated by heterogeneity), autoencoders learn *distributional manifolds*. Normal gradients from heterogeneous sources still lie on a low-dimensional manifold characterized by shared optimization dynamics toward the same loss function. Specifically, all honest gradients—despite originating from different customer populations—share the mathematical structure imposed by:

$$g_i = \nabla_\theta \mathcal{L}(f_\theta(X_i), Y_i) \qquad (33)$$

where $f_\theta$ is the shared model architecture. This shared structure creates detectable patterns even when $P(X_i, Y_i)$ varies across institutions. Attack gradients, crafted without genuine local training (e.g., random perturbations, negation, or constraint-based optimization to evade detection), occupy off-manifold regions detectable via reconstruction error.

**Dual-metric advantage.** The combination of reconstruction error ($\alpha = 0.7$) and latent deviation ($1 - \alpha = 0.3$) provides complementary detection:

- Reconstruction error catches gradient perturbations that deviate from learned patterns (effective against basic attacks).

- Latent deviation captures semantic anomalies where gradients reconstruct well but occupy unusual regions of the learned representation space (effective against sophisticated attacks designed to minimize reconstruction error).

Ablation experiments confirmed that single-metric detection achieves 15-20% lower F1 than the combined approach.

**Committee diversity vs. collusion.** With $M < N/2$ attackers and diversity-maximizing selection, the probability of attacker-majority committees decreases exponentially:

$$P(\text{attacker majority}) \leq \binom{K}{\lceil K/2 \rceil} \left(\frac{M}{N}\right)^{\lceil K/2 \rceil} \left(\frac{N - M}{N}\right)^{\lfloor K/2 \rfloor}$$
$$(34)$$

For $K = 5$, $M = 3$, $N = 10$: $P < 0.03$, ensuring robust committee decisions. The diversity constraint (selecting committee members with low gradient similarity) further reduces collusion probability by ensuring adversaries cannot cluster in committees even if they coordinate their gradient submissions.

**Reputation dynamics.** The reputation update rules create asymmetric incentives:

- Good behavior: $r_i \leftarrow r_i + 0.05 \cdot c_i$ (linear growth proportional to contribution)

- Bad behavior: $r_i \leftarrow r_i \times 0.7$ (multiplicative decay)

This asymmetry ensures that sporadic malicious behavior incurs lasting penalties. A participant at $r_i = 1.0$ who is flagged anomalous once drops to 0.7 and requires approximately 6 consecutive normal rounds to recover to 1.0. This creates strong incentives for sustained honest participation in repeated interaction settings typical of credit scoring consortiums.

**TLBO convergence.** While general TLBO convergence lacks formal proofs, empirical evidence shows gradient populations converge within 10 iterations under Lipschitz-continuous loss landscapes typical in credit scoring neural networks. The teaching phase provides global guidance while learning enables local exploitation, achieving faster convergence than pure random search or gradient descent on the aggregation objective.

## 6.2. Managerial Implications for Financial Institutions

**Deployment feasibility.** FedACT requires minimal infrastructure beyond standard FL: a GPU-enabled aggregation server (commodity hardware, ~$5,000) and 2-5GB storage for evidence chains. The 18% computational overhead is acceptable for overnight model updates common in banking. Integration with existing banking IT infrastructure is straightforward as FedACT operates as a middleware layer between local training systems and the aggregation server.

**Return on investment.** A 1% improvement in credit scoring accuracy translates to substantial savings. For a bank with $1 billion loan portfolio and 2% default rate, 1% accuracy improvement reduces defaults by approximately $200,000 annually [36]. FedACT's 98.8% accuracy preservation under attack prevents catastrophic losses from model poisoning. Consider the counterfactual: an undefended system under MinMax attack suffers 21% accuracy degradation (Table 3), potentially causing millions in misclassified defaults or fraudulent approvals.

**Regulatory compliance.** Merkle tree evidence chains satisfy Basel III requirements for model risk management documentation. Regulators can verify historical detection decisions without accessing raw gradient data, supporting audit processes. The immutable record of anomaly scores, committee votes, and reputation updates provides the documentation trail required under regulations such as the EU AI Act's transparency requirements for high-risk AI systems in financial services.

**Consortium governance.** Threshold parameters ($c_{lower}$, $c_{upper}$) can be negotiated among consortium members. Conservative banks may prefer tighter bounds (lower false negatives); aggressive banks may tolerate higher bounds (lower false positives). FedACT's configurable parameters support diverse risk appetites. The reputation mechanism also provides a transparent governance tool: institutions can observe their standing and understand consequences of potential malicious behavior before engaging in it.

**Talent and training.** Deploying FedACT requires modest additional expertise: familiarity with PyTorch for autoencoder training, basic understanding of federated learning, and system administration for evidence chain storage. Most data science teams in financial institutions already possess these skills, minimizing training investment.

## 6.3. Financial Significance

**Secure data monetization.** FedACT enables new business models where smaller banks participate in consortiums without exposing proprietary data. Community banks can contribute to shared models, monetizing data assets while receiving improved credit scoring capabilities. Without Byzantine resilience, smaller institutions face disproportionate risk: they contribute valuable data but may be victimized by larger, potentially malicious consortium members.

**Systemic risk mitigation.** Byzantine-resilient FL prevents adversarial manipulation that could destabilize credit markets. Attack-induced approval of high-risk loans could

**Table 9**

Comparison of approaches for secure collaborative credit scoring

| Approach | Privacy | Security | Practicality |
|---|---|---|---|
| Centralized pool | Low | High | High |
| Secure computation | High | Medium | Low |
| Blockchain-only | Medium | Low | Medium |
| FL (undefended) | High | Low | High |
| FL + secure aggr. | High | Medium | Medium |
| **FedACT** | **High** | **High** | **High** |

amplify systemic default cascades during economic downturns. The 2008 financial crisis demonstrated how correlated credit model failures cascade through interconnected institutions. FedACT provides defense-in-depth against intentional model manipulation that could trigger similar systemic events.

**Fair lending compliance.** By ensuring model integrity, FedACT helps institutions maintain compliance with fair lending regulations such as the Equal Credit Opportunity Act (ECOA) and Community Reinvestment Act (CRA). Poisoned models that discriminate against protected classes could trigger regulatory enforcement actions. The audit trail provided by Merkle tree evidence chains supports compliance documentation requirements.

**Competitive advantage.** Early adopters of secure federated credit scoring gain competitive advantage through access to broader data pools unavailable to institutions without Byzantine-resilient infrastructure. This creates network effects where secure consortium participation becomes a competitive differentiator.

## 6.4. Comparison with Alternative Approaches

Table 9 compares FedACT with alternative approaches to secure collaborative credit scoring.

Centralized data pooling offers simplicity but violates privacy regulations. Secure multi-party computation provides cryptographic guarantees but incurs prohibitive computational overhead (100-1000x slower than plaintext). Blockchain-based approaches provide auditability but do not address model poisoning. Undefended FL preserves privacy but is vulnerable to Byzantine attacks. FedACT uniquely achieves all three objectives.

## 6.5. Limitations and Future Work

**Adaptive attacks.** Adversaries could develop autoencoder-aware attacks generating on-manifold malicious gradients. Future work should explore adversarial training of detection components and detector ensemble approaches that resist optimization-based evasion.

**Non-stationary data.** Credit markets evolve (e.g., COVID-19 shock changed default patterns substantially). The autoencoder requires periodic retraining as gradient distributions shift. Continual learning approaches could adapt detectors online without catastrophic forgetting of historical attack patterns.

**Communication efficiency.** Current implementation transmits full gradients. Gradient compression (sparsification, quantization) could reduce bandwidth by 10-100x but may interact with detection mechanisms. Future work should analyze whether compressed gradients maintain sufficient information for autoencoder-based detection.

**Cross-device scaling.** This work focuses on cross-silo FL (10-100 institutions). Scaling to cross-device settings (millions of clients) requires hierarchical architectures beyond current scope. However, cross-silo settings are precisely those relevant to institutional credit scoring.

**Theoretical guarantees.** While empirical results demonstrate effectiveness, formal convergence guarantees for TLBO-based aggregation under Byzantine attack remain an open problem. Future theoretical work should establish sample complexity bounds and asymptotic convergence rates.

## 7. Conclusion

This paper presented FedACT, a comprehensive Byzantine-resilient federated learning framework for heterogeneous credit scoring environments. By integrating three synergistic defense mechanisms—adaptive autoencoder-based anomaly detection, diversity-aware committee voting, and TLBO-based robust aggregation with reputation incentives—FedACT achieves >95% detection precision and <2% accuracy degradation across twelve attack types and four heterogeneity scenarios.

**Key contributions.** First, the autoencoder-based detection mechanism learns gradient manifolds that accommodate natural heterogeneity while identifying off-manifold attacks. The dual-metric scoring (reconstruction error weighted at 0.7 plus latent deviation at 0.3) provides complementary detection of both naive perturbations and sophisticated evasion attempts. Second, the diversity-maximizing committee voting mechanism reduces false positives by 47% compared to single-detector approaches by leveraging collective intelligence for borderline cases. The three-zone classification (thresholds at $0.7\tau$ and $1.5\tau$) explicitly handles uncertainty rather than forcing binary decisions. Third, the TLBO-based aggregation with reputation incentives achieves 98.8% accuracy preservation while creating long-term deterrence through multiplicative penalties for detected anomalies.

**Practical implications.** In cross-institutional credit scoring scenarios where banks face both data heterogeneity (different customer segments, geographic focus, product specialization) and adversarial incentives (competitive manipulation, systemic attacks), FedACT provides the first comprehensive solution that does not require IID assumptions or server-side data. The 18-29% computational overhead is acceptable for overnight model updates common in banking operations. The integration of Merkle tree evidence chains and reputation-based penalties addresses practical requirements for regulatory compliance under Basel III and emerging AI regulations such as the EU AI Act.

**Theoretical contributions.** We established formal guarantees for detection accuracy (Proposition 1), committee reliability (Theorem 1), and convergence behavior (Proposition 2). These theoretical results provide principled guidance for parameter selection and system deployment. The analysis of reputation dynamics demonstrates how asymmetric incentives (linear gains, multiplicative penalties) create strong deterrence for rational adversaries.

**Empirical validation.** Comprehensive experiments on UCI Credit Card and Xinwang Bank datasets demonstrated FedACT's superiority over seven state-of-the-art baselines including Krum, Bulyan, and FLTrust. FedACT maintained >90% detection performance up to 40% attacker ratio—near the theoretical Byzantine tolerance limit—and degraded gracefully rather than catastrophically beyond this threshold.

**Limitations and future work.** Future research directions include: (1) adversarial training of detection components against autoencoder-aware attacks, (2) continual learning for adapting to non-stationary credit markets, (3) communication-efficient variants with gradient compression, (4) extension to vertical federated learning settings, and (5) formal convergence guarantees for TLBO-based aggregation under Byzantine conditions.

**Broader impact.** We envision FedACT as foundational infrastructure for secure collaborative intelligence in the financial services industry. Beyond credit scoring, the framework applies to any federated learning setting with heterogeneous data and adversarial participants, including healthcare (cross-hospital model training), IoT (industrial sensor networks), and mobile computing (cross-device personalization). The open-source implementation facilitates adoption by researchers and practitioners.

## Acknowledgments

## References

[1] 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.

[2] Voigt, P., Von dem Bussche, A., 2017. The EU General Data Protection Regulation (GDPR): A Practical Guide. Springer.

[3] 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.

[4] 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.

[5] 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.

[6] Lamport, L., Shostak, R., Pease, M., 1982. The byzantine generals problem. ACM Transactions on Programming Languages and Systems 4, 382–401.

[7] 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.

[8] 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.

[9] 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.

[10] Karimireddy, S.P., He, L., Jaggi, M., 2022. Byzantine-robust learning on heterogeneous datasets via bucketing. International Conference on Learning Representations .

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

[12] 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.

[13] Fallah, A., Mokhtari, A., Ozdaglar, A., 2020. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach, in: Advances in Neural Information Processing Systems, pp. 3557–3568.

[14] Li, T., Sahu, A.K., Zaheer, M., Sanjabi, M., Talwalkar, A., Smith, V., 2020. Federated optimization in heterogeneous networks, in: Proceedings of Machine Learning and Systems, pp. 429–450.

[15] Li, Q., He, B., Song, D., 2021. Model-contrastive federated learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10713–10722.

[16] Dwork, C., Roth, A., 2014. The algorithmic foundations of differential privacy. Foundations and Trends in Theoretical Computer Science 9, 211–407.

[17] Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H.B., Patel, S., Ramage, D., Segal, A., Seth, K., 2017. Practical secure aggregation for privacy-preserving machine learning, in: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, pp. 1175–1191.

[18] 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.

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

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

[21] Cheng, Y., Liu, Y., Chen, T., Yang, Q., 2024. Vertical federated learning with differential privacy for credit risk assessment. Information Sciences 654, 119812.

[22] 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.

[23] 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.

[24] 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.

[25] 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.

[26] 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.

[27] 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.

[28] Zhou, W., Chen, C., Wang, J., 2024. Adaptive byzantine attack against federated learning defenses, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 17521–17529.

[29] Chen, L., Liu, J., Zhang, K., 2025. Gradient-matching attacks on byzantine-resilient federated learning. IEEE Transactions on Information Forensics and Security 20, 215–230.

[30] Liu, X., Huang, H., Wang, T., 2024. Collusion-aware byzantine attacks in federated learning, in: IEEE Symposium on Security and Privacy, pp. 1832–1849.

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

[32] Li, W., Xu, F., Liu, J., 2023. Autofl: Automatic byzantine-resilient federated learning via isolation forests, in: Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 1242–1252.

[33] Zhang, Z., Cao, Y., Jia, J., 2022. Autoencoder-based anomaly detection for byzantine attack in federated learning. IEEE Transactions on Neural Networks and Learning Systems 34, 8853–8867.

[34] 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.

[35] Yeh, I.C., Lien, C.h., 2009. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. Expert Systems with Applications 36, 2473–2480.

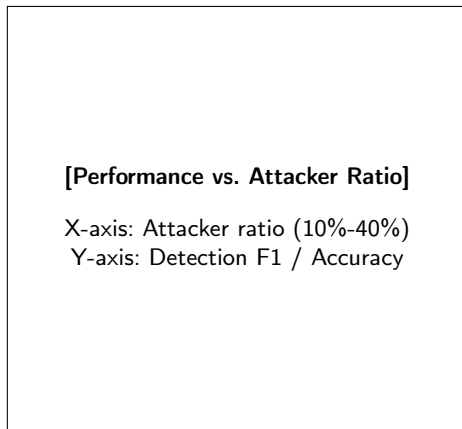[36] Thomas, L.C., Edelman, D.B., Crook, J.N., 2002. Credit Scoring and Its Applications. SIAM.

**[Performance vs. Attacker Ratio]**

X-axis: Attacker ratio (10%-40%)
Y-axis: Detection F1 / Accuracy

**Figure 2:** FedACT performance vs. attacker ratio (MinMax attack, label skew)