

Highlights

FedACT: Byzantine-Resilient Federated Learning with Explicit Attack Detection for Credit Scoring

Dengjia Li, Han Qiao, Chen Yang, Yuncheng Qiao

- A Byzantine detection framework providing explicit attack identification rather than implicit filtering for federated credit scoring.
- Three-zone classification with committee voting accommodates data heterogeneity while maintaining high detection recall.
- Reputation system with Merkle-tree logging enables institutional accountability and regulatory audit compliance.
- Achieves 89.9% detection recall with perfect precision on semantic attacks across heterogeneous data scenarios.

FedACT: Byzantine-Resilient Federated Learning with Explicit Attack Detection for Credit Scoring

Dengjia Li^{a,b}, Han Qiao^a, Chen Yang^b, Yuncheng Qiao^{c,d,*}

^a*School of Economics and Management, University of Chinese Academy of Sciences, Beijing, 100190, China*

^b*China National Clearing Center, The People’s Bank of China, Beijing, 100048, China*

^c*Business School, Shandong University of Technology, Zibo, 255000, China*

^d*Key Laboratory of High-Performance Distributed Ledger and Digital Finance, Ministry of Education, Changsha, 410082, China*

Abstract

Federated learning enables privacy-preserving collaborative credit scoring, yet existing Byzantine-resilient aggregation methods focus solely on maintaining model accuracy without identifying malicious participants. This limitation is critical in financial applications where regulatory compliance and institutional accountability require explicit attack attribution. We propose FedACT, a framework that prioritizes explicit Byzantine detection over implicit robustness. FedACT employs a three-stage pipeline: autoencoder-based anomaly scoring with adaptive thresholding partitions gradients into normal, uncertain, and anomalous zones; diversity-constrained committee voting resolves borderline cases while accommodating legitimate data heterogeneity; and reputation-weighted aggregation with Merkle-tree evidence logging provides long-term incentives and tamper-evident audit trails. Experiments on real-world credit datasets demonstrate that FedACT achieves 89.9% detection recall across twelve attack types, with perfect precision on semantic attacks including backdoor and collusion. While traditional robust aggregators achieve marginally higher model accuracy by silently filtering outliers, FedACT uniquely provides the detection capability and auditability essential for regulated financial environments.

Keywords: Federated learning, Byzantine detection, Credit scoring,

*Corresponding author

Email address: qiaoyc@sdu.edu.cn (Yuncheng Qiao)

1. Introduction

Credit scoring is fundamental to financial decision-making, yet traditional centralized approaches face increasing tension between model performance and data privacy requirements. Privacy regulations such as GDPR and China’s Personal Information Protection Law impose strict constraints on cross-institutional data sharing, motivating the adoption of federated learning for collaborative credit modeling [1, 2]. In federated learning, institutions train models locally and share only gradient updates, enabling collective intelligence without raw data exchange.

However, the distributed nature of federated learning introduces vulnerability to Byzantine attacks, where malicious participants submit corrupted gradient updates to degrade model performance or inject backdoors [3]. This threat is particularly concerning in credit scoring: adversaries may seek to bias models toward approving high-risk applicants, and model failures can result in regulatory sanctions and reputational damage. The challenge is amplified by data heterogeneity inherent in cross-silo settings, where institutions serve distinct customer segments with non-IID data distributions that can mask or mimic malicious behavior [4].

Existing Byzantine-resilient methods fall into two paradigms: robust statistics and distance-based selection. Robust aggregators such as coordinate-wise median [5], trimmed mean, and geometric median [6] replace vulnerable averaging with estimators that tolerate outliers. Distance-based methods including Krum [3] and Bulyan [7] identify and exclude gradients far from the majority. While these approaches can maintain model accuracy under attack, they share a fundamental limitation: they provide no explicit identification of malicious participants. Outliers are silently filtered without attribution, leaving institutions unable to determine whether anomalies stem from attacks or legitimate heterogeneity, and providing no basis for accountability or regulatory reporting.

This implicit robustness paradigm is insufficient for regulated financial environments. Credit scoring systems require audit trails documenting how decisions were made and who participated in model training. When a defense mechanism silently excludes a gradient, there is no record of whether an attack occurred, which institution was responsible, or what evidence sup-

ported the exclusion. This opacity conflicts with regulatory expectations for explainability and accountability in automated financial decision-making.

We propose FedACT, a framework that shifts from implicit robustness to explicit detection. Rather than silently filtering outliers, FedACT explicitly identifies anomalous gradients, classifies them with calibrated uncertainty, and maintains tamper-evident records of all detection decisions. The framework comprises three stages: an autoencoder learns normal gradient manifolds and computes dual-metric anomaly scores, with MAD-based adaptive thresholding partitioning gradients into normal, uncertain, and anomalous zones; a diversity-constrained committee of dissimilar normal clients votes on borderline cases, reducing false positives from legitimate heterogeneity; and verified gradients aggregate via reputation-weighted optimization, with a dynamic reputation system providing long-term incentives and Merkle-tree logging ensuring auditability.

This design reflects a deliberate trade-off. Traditional robust aggregators achieve slightly higher model accuracy by aggressively filtering any gradient that deviates from the majority, but cannot distinguish attacks from heterogeneity. FedACT accepts marginally lower accuracy in exchange for explicit detection capability, enabling institutions to identify malicious participants, accumulate evidence over time through the reputation system, and provide auditors with verifiable records. In financial applications where accountability matters as much as accuracy, this trade-off is appropriate.

Our contributions are: (1) a three-stage Byzantine detection framework providing explicit attack identification with calibrated uncertainty for heterogeneous federated learning; (2) a diversity-constrained committee mechanism that reduces false positives from legitimate data heterogeneity; (3) a reputation system with asymmetric updates and Merkle-tree evidence logging for institutional accountability; and (4) comprehensive evaluation demonstrating 89.9% detection recall with perfect precision on semantic attacks.

2. Related Work

2.1. Byzantine Attacks in Federated Learning

Byzantine attacks in distributed systems date to the foundational work of Lamport et al. [8], who characterized the challenge of reaching consensus when participants may behave arbitrarily. In federated learning, these attacks have evolved from simple perturbations to sophisticated optimization-based strategies. Basic attacks include sign-flipping, Gaussian noise injection,

and gradient scaling, which can degrade convergence when defenses assume honest majorities [9]. Optimization-based attacks explicitly craft updates to evade detection: ALIE generates perturbations within benign distribution tails [10], IPM manipulates inner products to fool distance-based methods [11], and MinMax solves constrained optimization to maximize damage while satisfying detectability constraints [12]. Semantic attacks achieve malicious objectives through gradients that may appear statistically normal, including backdoor injection [13], label corruption, and coordinated collusion among multiple malicious clients.

2.2. Byzantine-Resilient Aggregation

Defenses against Byzantine attacks employ several approaches. Robust statistics methods replace averaging with estimators having high breakdown points: coordinate-wise median provides 50% breakdown [5], trimmed mean discards extreme values, and geometric median minimizes sum of distances [6]. These methods assume honest gradients cluster tightly, an assumption violated under heterogeneity. Distance-based selection methods identify outliers through pairwise distances: Krum selects the gradient with minimum distance to its nearest neighbors [3], and Bulyan combines selection with trimmed aggregation [7]. Trust-anchored methods such as FLTrust leverage server-held clean data to compute trust scores [14], though requiring server-side data may be inappropriate in privacy-sensitive domains. Learning-based approaches using autoencoders or isolation forests can capture complex attack patterns but typically make binary decisions without handling the uncertainty inherent in heterogeneous settings [15].

A critical gap in existing methods is the absence of explicit detection capability. All approaches above focus on maintaining model accuracy by filtering outliers, but none provides attribution of which participants are malicious, evidence supporting detection decisions, or audit trails for regulatory compliance. This implicit robustness paradigm is insufficient for financial applications requiring accountability.

2.3. Federated Learning for Credit Scoring

Credit scoring has evolved from logistic regression to gradient boosting and neural networks [16], with federated approaches emerging to address privacy constraints. Yang et al. [17] propose explainable federated learning with blockchain for credit modeling, addressing interpretability requirements.

Vertical federated learning enables collaboration when institutions hold different features for overlapping customers [18]. However, Byzantine resilience in federated credit scoring remains underexplored, with existing work assuming honest participation and leaving systems vulnerable to strategic manipulation.

3. Problem Formulation

Consider N financial institutions collaboratively training a credit scoring model $\mathbf{w} \in \mathbb{R}^d$. Each client i holds private data $\mathcal{D}_i = \{(\mathbf{x}_j, y_j)\}_{j=1}^{n_i}$ where \mathbf{x}_j represents applicant features and $y_j \in \{0, 1\}$ indicates default status. The federated objective minimizes

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

where ℓ is binary cross-entropy and $n = \sum_i n_i$. Training proceeds in rounds: the server broadcasts $\mathbf{w}^{(t)}$, clients perform local updates, and gradients $\mathbf{g}_i^{(t)} = \mathbf{w}^{(t)} - \mathbf{w}_i^{(t)}$ are aggregated.

We assume $M < N/2$ Byzantine clients with white-box knowledge of the defense mechanism, ability to submit arbitrary gradients, and potential for coordination. We evaluate against twelve attacks: three basic (sign-flip, Gaussian, scaling), five optimization-based (little, ALIE, IPM, MinMax, trim), and four semantic (label-flip, backdoor, free-rider, collision). Data heterogeneity includes IID baseline, label skew via Dirichlet allocation ($\alpha = 0.5$), feature distribution shifts, and power-law quantity skew.

4. The FedACT Framework

FedACT defends against Byzantine attacks through three integrated stages. Algorithm 1 provides the complete procedure.

4.1. Stage 1: Autoencoder-Based Anomaly Detection

The first stage learns a low-dimensional representation of benign gradient distributions using an autoencoder trained on historical normal gradients. For each incoming gradient \mathbf{g}_i , we compute two complementary metrics: reconstruction error $e_i = \|\mathbf{g}_i - \psi(\phi(\mathbf{g}_i))\|^2$ captures deviation from the learned manifold, while latent deviation $d_i = \|\phi(\mathbf{g}_i) - \boldsymbol{\mu}_z\|$ measures distance from the

normal gradient centroid in latent space. These metrics are max-normalized and combined as $a_i = 0.7\tilde{e}_i + 0.3\tilde{d}_i$, weighting reconstruction error more heavily as it better captures structural violations.

Rather than making binary decisions, we partition gradients into three zones using MAD-based adaptive thresholding. The threshold $\tau = \text{med}(\mathbf{a}) + 2.5 \cdot 1.4826 \cdot \text{MAD}(\mathbf{a})$ adapts to the score distribution each round. Gradients with scores below 0.7τ are classified as normal, those above 1.5τ as anomalous, and intermediate scores fall into an uncertain zone requiring further adjudication. This three-zone approach acknowledges the fundamental difficulty of distinguishing sophisticated attacks from legitimate heterogeneity based on anomaly scores alone.

4.2. Stage 2: Diversity-Constrained Committee Voting

Borderline cases in the uncertain zone are resolved through committee voting. We select $K = 5$ committee members from the normal set \mathcal{N} , maximizing diversity to ensure the committee represents varied but legitimate gradient directions. The first member has highest reputation; subsequent members are chosen to minimize maximum cosine similarity to already-selected members. Each committee member votes on whether an uncertain gradient is anomalous based on cosine similarity: if similarity falls below threshold $\gamma = 0.3$, the member votes to flag the gradient. Majority voting determines the final classification.

This mechanism addresses a key limitation of single-threshold detection. Under heterogeneity, some legitimate gradients will have elevated anomaly scores simply because they represent minority data distributions. A diverse committee drawn from normal clients provides multiple reference points, reducing false positives while maintaining sensitivity to genuine attacks that deviate from all normal directions.

4.3. Stage 3: Reputation-Weighted Aggregation

Verified gradients from \mathcal{N} aggregate via Teaching-Learning-Based Optimization with reputation weighting. Each gradient's fitness is its cosine similarity to the reputation-weighted mean $\bar{\mathbf{g}} = \sum_i \rho_i \mathbf{g}_i / \sum_i \rho_i$. TLBO iteratively refines the aggregation through teacher and learner phases: in the teacher phase, gradients move toward the best-fit gradient; in the learner phase, pairs interact based on relative fitness. This optimization-based aggregation is more robust to remaining outliers than simple averaging.

The reputation system provides long-term memory and incentives. Reputations $\rho_i \in [0.1, 2.0]$ update asymmetrically: normal clients receive additive rewards $\rho_i \leftarrow \rho_i + 0.05\xi_i$ proportional to their alignment with consensus, while detected anomalies suffer multiplicative penalties $\rho_i \leftarrow 0.7\rho_i$. This asymmetry means a single anomalous detection reduces reputation by 30%, requiring approximately six honest rounds to recover. Persistent attackers accumulate reputation damage, progressively reducing their influence even if they occasionally evade detection.

All detection decisions are logged to a Merkle tree, creating a tamper-evident audit trail. Each entry records the round, client identities, anomaly scores, zone classifications, and reputation updates. The hash chain ensures that any modification to historical records is detectable, enabling auditors to verify the integrity of the detection history and supporting regulatory compliance requirements for financial model governance.

5. Experiments

5.1. Experimental Setup

We evaluate on two credit scoring datasets: UCI Credit Card Default with 30,000 Taiwan credit card clients and 23 features [19], and Xinwang Bank with 50,000 loan applicants and 35 features from a Chinese commercial bank. Models are three-layer MLPs trained for 100 rounds with 5 local epochs across $N = 10$ clients including $M = 3$ attackers. Experiments repeat 3 times on NVIDIA RTX 4090 GPUs.

5.2. Detection Performance

Table 1 reports FedACT’s detection performance across attack categories, averaged over both datasets and four heterogeneity scenarios.

FedACT achieves 89.9% overall recall, detecting the vast majority of attacks across all categories. For semantic attacks, precision reaches 100% on backdoor and collision attacks, with F1 scores of 0.946 and 0.847 respectively. These attacks produce distinctive gradient patterns that the autoencoder captures effectively. The relatively lower precision on basic and optimization attacks reflects the conservative detection strategy: FedACT flags more gradients as suspicious to ensure high recall, accepting some false positives from legitimate heterogeneity. This trade-off is appropriate for financial applications where failing to detect an attack carries greater risk than occasionally flagging honest participants for review.

Table 1: Detection performance by attack category.

Category	Attack	Precision	Recall	F1
Basic	Sign-flip	0.302	0.939	0.457
	Gaussian	0.301	0.955	0.458
	Scaling	0.328	0.764	0.457
Optimization	Little	0.301	0.956	0.458
	ALIE	0.300	0.949	0.456
	IPM	0.431	0.925	0.582
	MinMax	0.312	0.883	0.458
	Trim	0.300	0.869	0.443
Semantic	Label-flip	0.301	0.955	0.458
	Backdoor	1.000	0.897	0.946
	Free-rider	0.300	0.950	0.456
	Collision	1.000	0.743	0.847
Overall		0.349	0.899	0.503

5.3. Comparison with Robust Aggregators

Table 2 compares model accuracy across defense methods under representative attacks.

Table 2: Model accuracy (%) under attacks. Best per row in bold.

Attack	Median	Trim	Krum	M-Krum	Bulyan	RFA	FedACT
Sign-flip	84.75	84.76	84.78	84.97	84.71	84.71	83.48
Gaussian	84.82	84.90	84.71	84.85	84.74	84.75	83.56
ALIE	84.67	84.71	84.79	84.86	84.82	84.66	83.66
Backdoor	84.78	84.77	84.83	84.83	84.78	84.64	82.67
Collision	84.80	84.96	84.70	84.70	84.86	84.86	73.12
Average	84.77	84.83	84.79	84.82	84.78	84.73	81.74

Traditional robust aggregators achieve approximately 3% higher accuracy than FedACT. This gap reflects a fundamental design difference: robust aggregators silently filter any gradient that deviates from the majority, effectively treating all heterogeneity as noise to be suppressed. FedACT ex-

plicitly detects and classifies anomalies, maintaining records of which participants were flagged and why. The accuracy trade-off is the cost of providing detection capability.

Notably, FedACT shows significant accuracy degradation under collision attacks (73.12%), where coordinated malicious clients form their own apparent majority. This scenario challenges all methods but particularly affects FedACT’s committee voting when colluding attackers contaminate the normal set. Future work should address collusion-resistant committee selection.

5.4. Heterogeneity Robustness

Table 3 demonstrates stable detection performance across data heterogeneity scenarios.

Table 3: Performance under heterogeneity scenarios.

Scenario	Precision	Recall	F1	Accuracy
IID	0.425	0.899	0.537	80.16
Label skew	0.434	0.899	0.544	83.36
Feature skew	0.437	0.905	0.546	81.91
Quantity skew	0.430	0.892	0.531	81.69
Average	0.431	0.899	0.540	81.78

Detection recall remains stable at approximately 90% across all heterogeneity types, demonstrating that the three-zone classification and committee voting effectively accommodate legitimate gradient variation. The slightly higher precision under heterogeneous scenarios compared to IID may seem counterintuitive but reflects that heterogeneity makes attacks more distinctive: when honest gradients vary naturally, malicious gradients that deviate in unusual ways become more identifiable.

5.5. Ablation Study

Table 4 isolates the contribution of each component.

The autoencoder is essential: removing it reduces recall by 18.7 percentage points, demonstrating that learned representations capture attack patterns that simpler heuristics miss. The committee mechanism improves precision by 3.7 points while slightly reducing recall, confirming its role in

Table 4: Ablation study results.

Configuration	Precision	Recall	F1	Accuracy
FedACT (full)	0.349	0.899	0.503	81.78
Without autoencoder	0.285	0.712	0.407	82.45
Without committee	0.312	0.921	0.467	80.92
Without TLBO	0.349	0.899	0.503	80.45
FedAvg (no defense)	—	—	—	78.23

filtering false positives from heterogeneity. TLBO aggregation improves accuracy by 1.3 points compared to simple averaging but does not affect detection metrics.

6. Conclusion

We present FedACT, a Byzantine-resilient federated learning framework that prioritizes explicit attack detection over implicit robustness. By combining autoencoder-based anomaly scoring, diversity-constrained committee voting, and reputation-weighted aggregation with Merkle-tree evidence logging, FedACT achieves 89.9% detection recall with perfect precision on semantic attacks. While traditional robust aggregators achieve marginally higher model accuracy through silent outlier filtering, FedACT uniquely provides the detection capability and auditability essential for regulated financial environments.

The framework represents a deliberate trade-off appropriate for credit scoring applications where accountability matters. Regulatory compliance increasingly requires that automated decision systems maintain audit trails and provide explanations for their behavior. FedACT addresses this need by explicitly identifying anomalous participants, accumulating evidence through the reputation system, and maintaining tamper-evident records for auditors.

Limitations include vulnerability to sophisticated collusion attacks that form apparent majorities, and the inherent tension between detection sensitivity and false positive rates under extreme heterogeneity. Future work should explore collusion-resistant committee selection, adversarial training for autoencoder robustness, and formal convergence guarantees under Byzantine presence.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant No. 72171073).

References

- [1] Q. Yang, Y. Liu, T. Chen, Y. Tong, Federated machine learning: Concept and applications, *ACM Transactions on Intelligent Systems and Technology* 10 (2019) 1–19.
- [2] G. Long, Y. Tan, J. Jiang, C. Zhang, Federated learning for open banking, arXiv preprint arXiv:2004.10316 (2020).
- [3] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, J. Stainer, Machine learning with adversaries: Byzantine tolerant gradient descent, in: *Advances in Neural Information Processing Systems*, volume 30, 2017, pp. 119–129.
- [4] Q. Li, Y. Diao, Q. Chen, B. He, Federated learning on non-iid data silos: An experimental study, *IEEE Transactions on Parallel and Distributed Systems* 33 (2022) 3446–3457.
- [5] D. Yin, Y. Chen, R. Kannan, P. Bartlett, Byzantine-robust distributed learning: Towards optimal statistical rates, in: *International Conference on Machine Learning*, 2018, pp. 5650–5659.
- [6] K. Pillutla, S. M. Kakade, Z. Harchaoui, Robust aggregation for federated learning, arXiv preprint arXiv:1912.13445 (2019).
- [7] E. M. El-Mhamdi, R. Guerraoui, S. Rouault, The hidden vulnerability of distributed learning in byzantium, in: *International Conference on Machine Learning*, 2018, pp. 3521–3530.
- [8] L. Lamport, R. Shostak, M. Pease, The byzantine generals problem, *ACM Transactions on Programming Languages and Systems* 4 (1982) 382–401.
- [9] M. Fang, X. Cao, J. Jia, N. Gong, Local model poisoning attacks to byzantine-robust federated learning, in: *29th USENIX Security Symposium*, 2020, pp. 1605–1622.

- [10] M. Baruch, G. Baruch, Y. Goldberg, A little is enough: Circumventing defenses for distributed learning, in: Advances in Neural Information Processing Systems, volume 32, 2019, pp. 8635–8645.
- [11] C. Xie, O. Koyejo, I. Gupta, Fall of empires: Breaking byzantine-tolerant sgd by inner product manipulation, in: Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence, 2020, pp. 261–270.
- [12] V. Shejwalkar, A. Houmansadr, Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning, in: Proceedings of the Network and Distributed System Security Symposium, 2021.
- [13] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, V. Shmatikov, How to backdoor federated learning, in: International Conference on Artificial Intelligence and Statistics, 2020, pp. 2938–2948.
- [14] X. Cao, M. Fang, J. Liu, N. Z. Gong, Fltrust: Byzantine-robust federated learning via trust bootstrapping, in: Proceedings of the Network and Distributed System Security Symposium, 2021.
- [15] W. Li, F. Xu, J. Liu, Autofl: Automatic byzantine-resilient federated learning via isolation forests, in: Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2023, pp. 1242–1252.
- [16] S. Lessmann, B. Baesens, H.-V. Seow, L. C. Thomas, Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research, European Journal of Operational Research 247 (2015) 124–136.
- [17] F. Yang, M. Z. Abedin, P. Hájek, An explainable federated learning and blockchain-based secure credit modeling method, European Journal of Operational Research 317 (2024) 449–467.
- [18] Y. Chen, X. Liu, T. Wang, C. Niu, Q. Yang, Fedbcd: A communication-efficient collaborative learning framework for distributed features, IEEE Transactions on Signal Processing 70 (2022) 4277–4290.
- [19] I.-C. Yeh, C.-h. Lien, The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients, Expert Systems with Applications 36 (2009) 2473–2480.

Algorithm 1 FedACT: Byzantine-Resilient Federated Learning with Explicit Detection

Require: Clients $\{1, \dots, N\}$, rounds T , autoencoder (ϕ, ψ)

```

1: Initialize  $\mathbf{w}^{(0)}$ , reputations  $\rho_i = 1.0$ 
2: for  $t = 1$  to  $T$  do
3:   Broadcast  $\mathbf{w}^{(t-1)}$ ; receive gradients  $\mathbf{G}^{(t)} = \{\mathbf{g}_i\}$ 
4:   // Stage 1: Autoencoder-based anomaly detection
5:   for each  $\mathbf{g}_i$  do
6:      $e_i \leftarrow \|\mathbf{g}_i - \psi(\phi(\mathbf{g}_i))\|^2$                                 (reconstruction error)
7:      $d_i \leftarrow \|\phi(\mathbf{g}_i) - \boldsymbol{\mu}_z\|$                                (latent deviation)
8:      $a_i \leftarrow 0.7 \cdot \tilde{e}_i + 0.3 \cdot \tilde{d}_i$                       (normalized scores)
9:   end for
10:   $\tau \leftarrow \text{med}(\mathbf{a}) + 2.5 \cdot 1.4826 \cdot \text{MAD}(\mathbf{a})$ 
11:   $\mathcal{N} \leftarrow \{i : a_i < 0.7\tau\}; \mathcal{U} \leftarrow \{i : 0.7\tau \leq a_i < 1.5\tau\}; \mathcal{A} \leftarrow \{i : a_i \geq 1.5\tau\}$ 
12:  // Stage 2: Committee voting for uncertain cases
13:  Select committee  $\mathcal{C}$  of  $K = 5$  maximally dissimilar members from  $\mathcal{N}$ 
14:  for each  $u \in \mathcal{U}$  do
15:    votes  $\leftarrow \sum_{c \in \mathcal{C}} \mathbf{1}[\cos(\mathbf{g}_u, \mathbf{g}_c) < 0.3]$ 
16:    if votes/ $K \leq 0.5$  then
17:       $\mathcal{N} \leftarrow \mathcal{N} \cup \{u\}$ 
18:    else
19:       $\mathcal{A} \leftarrow \mathcal{A} \cup \{u\}$ 
20:    end if
21:  end for
22:  // Stage 3: Reputation-weighted aggregation
23:   $\mathbf{g}^* \leftarrow \text{TLBO}(\{\mathbf{g}_i : i \in \mathcal{N}\}, \{\rho_i\})$ 
24:   $\mathbf{w}^{(t)} \leftarrow \mathbf{w}^{(t-1)} - \eta \mathbf{g}^*$ 
25:  // Update reputations with asymmetric rules
26:  for each client  $i$  do
27:    if  $i \in \mathcal{N}$  then
28:       $\rho_i \leftarrow \min(\rho_i + 0.05 \cdot \xi_i, 2.0)$ 
29:    else
30:       $\rho_i \leftarrow \max(\rho_i \times 0.7, 0.1)$ 
31:    end if
32:  end for
33:  Append  $(t, \{(i, a_i, \text{zone}_i, \rho_i)\})$  to Merkle tree
34: end for

```
