# AC-BlockDFL: Audit-driven Committee BlockDFL for Secure Federated Learning

ANONYMOUS AUTHOR(S)

Blockchain-based Federated Learning (BCFL) faces a critical scalability-security trade-off. While committee-based architectures significantly reduce communication overhead, they introduce a fundamental vulnerability: the *Progressive Committee Capture Attack* (PCCA). In PCCA, rational adversaries exploit the stake-election-reward feedback loop to gradually capture committee control through strategic starvation and stake accumulation. We propose *AC-BlockDFL*, a defense framework that decouples system security from committee honesty through optimistic execution and asynchronous auditing. By internalizing the externalities of malicious behavior via a game-theoretic slashing protocol, AC-BlockDFL ensures that attacks yield negative expected utility. Our evaluation over 2,000 rounds demonstrates that AC-BlockDFL suppresses malicious stake ratios from 1.3$\times$ to 0.37$\times$, reducing unavailability rates from 22.3% to below 1% while maintaining $O(C^2)$ communication complexity.

CCS Concepts: • **Security and privacy** → **Distributed systems security**; • **Computing methodologies** → *Artificial intelligence*.

Additional Key Words and Phrases: Federated Learning, Blockchain, Committee Consensus, Game Theory, Incentive Compatibility

## 1 Introduction

Blockchain-based Federated Learning (BCFL) has emerged as a promising paradigm for collaborative machine learning in environments where mutual trust among participants cannot be assumed. Real-world deployments in Low Earth Orbit (LEO) satellite networks [5, 11, 18], vehicular networks (V2X) [8, 12], and Industrial IoT [9, 14] demonstrate compelling use cases where decentralized coordination is essential. In LEO constellations, for instance, ground station contact windows last merely five minutes with downlink bandwidth limited to approximately 8 Mbps [18], rendering centralized aggregation architectures impractical. BCFL addresses these constraints by establishing decentralized trust infrastructure across heterogeneous satellite operators, reducing model convergence time by up to thirty hours [5].

However, BCFL systems face a fundamental scalability bottleneck when approaching large-scale deployment. The predominant use of Practical Byzantine Fault Tolerance (PBFT) [3] and its variants introduces $O(N^2)$ message complexity, causing consensus latency to dominate training time as participant counts grow. Empirical measurements from FLCoin [15] reveal that at 100 nodes, single-round consensus generates over 20,000 message exchanges with latency exceeding 25 seconds—comparable to or exceeding the model training duration itself. Storage requirements compound this challenge: Bitcoin full nodes require approximately 200 GB while Ethereum exceeds 465 GB, fundamentally incompatible with edge devices possessing only KB-to-MB scale memory [1].

*The Committee Mechanism.* To address these scalability constraints, recent work has converged on *committee-based* architectures that delegate verification responsibility to a smaller subset of validators. Selection mechanisms include hash-ring sampling [13], stake-weighted election [7, 15], and Verifiable Random Function (VRF) based sortition [6, 16, 17]. These approaches yield substantial efficiency gains: FLCoin [15] reports 90% communication overhead reduction and 5.7× training speedup through sliding-window election, while BFLC [7] achieves sub-three-second consensus latency. Effectively, committee mechanisms reduce communication complexity from $O(N^2)$ to $O(C^2)$ or even $O(C)$, where $C \ll N$ is the committee size.

*The Blind Spot.* Despite these advances, existing BCFL literature harbors a critical yet overlooked vulnerability: the implicit assumption that committee members remain honest or that the proportion of malicious nodes stays static throughout system operation. Current defenses focus predominantly on data-plane attacks—Byzantine-robust aggregation rules such as Krum, Trimmed Mean, and Median [2, 19] assume an honest majority among aggregating nodes. However, these mechanisms provide no protection when the committee itself becomes compromised. As FedBlock [10] observes, when any participant may become a validator, systems cannot rely solely on honest majority assumptions but must actively detect and isolate malicious verifiers—a capability conspicuously absent from current BCFL architectures.

*The PCCA Threat.* We identify and formalize a novel attack vector: the *Progressive Committee Capture Attack* (PCCA). Unlike direct Byzantine attacks, PCCA adversaries employ *strategic starvation*—upon gaining committee control, attackers prioritize processing their own model updates while systematically denying service to honest participants. This manipulation of the reward distribution mechanism enables attackers to "legitimately" accumulate stake over successive rounds, progressively cementing their dominance until decentralized governance collapses entirely. Crucially, once the honest majority assumption fails in any given round, existing systems lack mechanisms to identify or penalize malicious actors, allowing attackers to maintain their advantage indefinitely.

*Contributions.* This paper presents *Audit-driven Committee BlockDFL* (AC-BlockDFL), a defense framework that decouples system security from collective committee honesty through optimistic execution with asynchronous auditing. Our contributions are threefold:

(1) **Attack Formalization.** We provide the first formal definition of the Progressive Committee Capture Attack (PCCA) and a rational adversary model that captures strategic, incentive-driven behavior. Through systematic simulation, we quantify PCCA's destructive impact on long-term incentive compatibility.

(2) **AC-BlockDFL Architecture.** We propose a novel defense architecture combining optimistic execution with asynchronous auditing. A distributed challenger network performs post-hoc verification of committee decisions, enabling detection and penalization of fraudulent aggregation results even when the committee is fully compromised.

(3) **Game-Theoretic Guarantees.** We design an internal slashing protocol grounded in game-theoretic analysis, ensuring that auditing costs remain strictly below potential gains from malicious behavior. We prove that honest participation constitutes the unique Nash equilibrium under repeated play [4]. Extensive simulations over 2,000 rounds demonstrate that AC-BlockDFL maintains model accuracy above 98.6% under 30% adversarial collusion, reduces communication overhead by 44.4% at equivalent security levels, and suppresses minimum unavailability rate from 20% to below 5%.

## 2 Background

This section establishes the theoretical foundations and technical background necessary for understanding the security challenges in committee-based blockchain federated learning. We first discuss the fundamental trust dilemma in federated learning, then introduce the principles of Byzantine fault tolerance, and finally establish the baseline system model for committee-based architectures.

### 2.1 Federated Learning and the Trust Dilemma

Federated learning (FL) represents a paradigm shift in distributed machine learning, encapsulating the principle of "bringing the model to the data" rather than aggregating data centrally [? ]. While FL significantly enhances data privacy by locally constraining raw data, its standard architecture relies on a fundamental assumption: participants must trust a central aggregation server to honestly execute aggregation and uniformly distribute results.

In the absence of verifiable consistency, the central server constitutes a single point of failure and a primary vulnerability. A malicious or compromised server could perform selective aggregation, intentionally excluding specific updates, or directly tamper with the global model to inject backdoors [? ]. Furthermore, while FL avoids direct data transmission, a malicious aggregator can still infer sensitive information from client updates [? ]. This trust dilemma severely restricts the deployment of FL in high-value, cross-organizational scenarios where participating entities may be independent or competitive, necessitating a decentralized trust infrastructure.

### 2.2 Byzantine Fault Tolerance Fundamentals

Blockchain technology, characterized by immutability, transparency, and decentralization, serves as an ideal infrastructure to resolve the FL trust dilemma. However, the security of blockchain fundamentally relies on consensus protocols designed to tolerate malicious behavior, rooted in the Byzantine Generals Problem [? ].

The mathematical constraint of Byzantine Fault Tolerance (BFT) dictates that a system of $N$ nodes can tolerate at most $f$ malicious nodes, requiring $N \geq 3f + 1$. This one-third threshold originates from the *quorum intersection principle*: to ensure sufficient honest endorsements, any decision needs $2f + 1$ confirmations. The intersection of any two $2f + 1$ sets guarantees the inclusion of at least $f + 1$ nodes, meaning at least one honest node witnesses both decisions, preventing contradictory states.

Practical Byzantine Fault Tolerance (PBFT) [3] reduces the communication complexity of BFT consensus to $O(N^2)$ through a three-phase commit protocol (Pre-prepare, Prepare, Commit). While efficient for small networks, the quadratic communication cost becomes a severe bottleneck for large-scale FL systems requiring frequent iterations involving hundreds of participants.

### 2.3 Committee-based BCFL Architecture

To reconcile the efficiency demands of federated learning with the security requirements of blockchain, the *committee-based architecture* has emerged as the prevailing design paradigm. By delegating consensus responsibilities to a smaller, representative subset of nodes (the committee), these systems reduce communication complexity from $O(N^2)$ to $O(C^2 + N)$ where the committee size $C \ll N$ [15].

*Baseline System Model: BlockDFL.* We adopt BlockDFL [13] as our baseline system model, representing the state-of-the-art in peer-to-peer BCFL. BlockDFL implements role separation, partitioning participants into three distinct roles per training round:

(1) **Update Providers**: Execute local model training on private data and submit bounded updates.

(2) **Aggregators**: Collect updates, perform filtering, and compute the aggregated global proposal.

(3) **Validators**: Form the committee to evaluate competing proposals via Krum scoring [2] and execute PBFT consensus to select the final model.

*The Stake-Election-Reward Cycle.* Role assignment in BlockDFL relies on *stake-weighted deterministic random selection*, using the previous block's hash as an unpredictable entropy source mapping onto a stake-weighted hash ring. This creates a critical economic incentive structure designed to solve the free-rider problem:

- **Stake**: Determines election probability; participants with higher stake bounds are proportionally more likely to be selected as Aggregators or Validators.

- **Reward**: Distributed exclusively to contributors of the accepted proposal (the winning Aggregator, included Update Providers, and Validators who voted for it).

This mechanism instantiates a *positive feedback loop*: receiving rewards increases absolute stake, which enhances future election probability and aggregation weight, thereby amplifying the likelihood of subsequent rewards. While intended to cultivate long-term honest contributions, this very dynamic introduces vulnerabilities to strategic persistence.

## 3 Related Work

The convergence of federated learning and blockchain technology has precipitated diverse architectural innovations to address decentralized coordination, privacy, and security. Early systems like DeepChain [17] and Biscotti [16] focused on preserving privacy during aggregation using cryptographic commitments and differential privacy. However, scaling BFT consensus to accommodating hundreds of FL clients remained a persistent challenge due to its $O(N^2)$ communication overhead.

### 3.1 Evolution of Committee Architectures

To mitigate consensus bottlenecks, recent literature has pivoted toward committee-based designs inspired by Algorand's sortition [6]. By randomly selecting a constant-size committee to perform consensus, systems effectively decouple performance from total network size. FLCoin [15] utilized a sliding-window mechanism based on contribution history to form dynamic committees, achieving up to 90% reduction in communication overhead. Similarly, BFLC [7] adopted a reputation-based election scheme, prioritizing nodes with high historical quality scores.

While optimization successes are evident, these election mechanisms inherently couple system security to committee composition. If an adversary captures a supermajority (e.g., > 2/3) of the committee seats, traditional data-layer defenses like Krum [2] or Trimmed Mean [19] are entirely bypassed because the compromised committee itself executes these algorithms.

### 3.2 Limitations of Existing Verification Methods

Addressing Byzantine behavior in decentralized aggregation currently relies on three primary verification paradigms:

*Cryptographic Verification (zkML)..* Zero-Knowledge Machine Learning (zkML) provides the strongest security guarantees by compiling learning computations into arithmetic circuits, allowing verification without re-execution [? ], [? ]. However, zkML faces prohibitive computational bottlenecks. Compiling simple architectures like ResNet-18 generates millions of polynomial constraints, and generating proofs takes minutes. More critically, zkML currently

cannot support complex, non-linear Byzantine-robust aggregation algorithms like Krum, which require $O(N^2 \cdot d)$ pairwise distance calculations that cause circuit explosions.

*Optimistic Execution (opML)..* Optimistic Machine Learning (opML) defaults to accepting computations but allows a challenge window during which "AnyTrust" challengers can submit fraud proofs via interactive bisection protocols [? ], [? ]. While efficient, mainstream opML architectures natively resolving disputes on-chain require challenge periods extending up to a week (e.g., Optimism [? ]). These latencies fundamentally conflict with the high-frequency iterative nature of federated learning, rendering opML broadly inapplicable for per-round FL aggregation.

*Committee Consensus and Static Blindspots.* Committee-based verification remains the most pragmatic solution but is underpinned by static probabilistic assumptions: current security analyses calculate committee capture probability by assuming a fixed adversarial population distribution [15]. This static view completely overlooks the behavior of *rational, strategic adversaries*. As indicated by our analysis of the stake-election-reward cycle, systems like BlockDFL [13] and FedBlock [10] contain positive feedback loops. Strategic attackers can exploit this by behaving honestly ("lurking") to accumulate stake and reputation until they acquire sufficient influence to capture the committee ("starving" honest participants).

### 3.3 Our Contribution

Our work directly addresses the systemic blindspot in static committee security. We define and analyze the Progressive Committee Capture Attack (PCCA), demonstrating how rational adversaries bypass conventional defenses. In contrast to existing static committee systems, AC-BlockDFL breaks the stake feedback loop through optimistic execution coupled with an asynchronous, stake-slashing auditing layer, securing the system economically while preserving the efficiency benefits of committee architectures.

## 4 AC-BlockDFL System Design

We now present Audit-driven Committee BlockDFL (AC-BlockDFL), a framework that reconciles the efficiency demands of federated learning with provable economic security guarantees. Our design philosophy centers on a fundamental insight: *liveness and security need not be coupled.* Traditional Byzantine fault-tolerant systems conflate these properties, requiring full consensus before any state change. AC-BlockDFL decouples them by assigning efficiency (liveness) to the committee while distributing security enforcement across the entire network through asynchronous auditing.

### 4.1 Design Philosophy

The core tension in blockchain-based federated learning stems from misaligned finality requirements. Financial systems demand immediate, irrevocable correctness for every transaction, necessitating synchronous consensus. However, federated learning's iterative nature—where single-round deviations can be corrected through subsequent training—creates design space for a fundamentally different approach.

AC-BlockDFL exploits this observation through *optimistic execution*: the committee's consensus triggers immediate model updates, while rigorous verification proceeds asynchronously in the background. Security emerges not from preventing malicious behavior *a priori*, but from ensuring that any detectable misbehavior incurs economic penalties that far exceed potential gains. This shifts the security model from threshold-based guarantees (requiring honest majorities) to *economic security* (requiring only that attackers be rational).
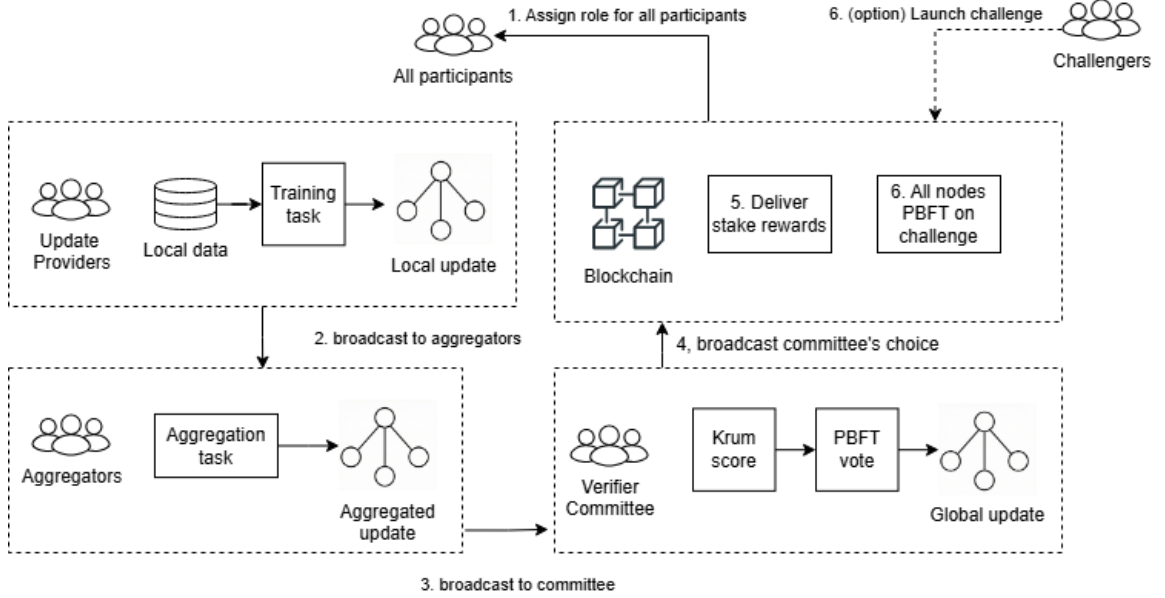
Fig. 1. AC-BlockDFL System Architecture and Workflow. The committee (Validators $\mathcal{V}$) performs optimistic execution while the Challenger network performs asynchronous auditing.

## 4.2 System Architecture and Workflow

AC-BlockDFL extends the BlockDFL committee model with three architectural innovations: (1) a fourth participant role—the *Challenger*—responsible for post-hoc auditing; (2) off-chain storage integration via IPFS to reduce on-chain overhead from $O(\text{ModelSize})$ to $O(\text{HashSize})$; and (3) an asynchronous audit window that enables network-wide arbitration without blocking normal operation. Figure 1 illustrates the complete workflow.

*Instant Update Protocol.* Algorithm 1 formalizes the protocol's "happy path." Each round proceeds through four phases: (1) *Role Assignment*—stake-weighted random selection determines validators $\mathcal{V}$, aggregators $\mathcal{A}$, and update providers $\mathcal{U}$ using the previous block hash as randomness source; (2) *Training and Off-chain Storage*—providers train locally, aggregators collect updates into proposals, upload to IPFS, and submit Content Identifiers (CIDs) on-chain; (3) *Consensus and Instant Commit*—validators retrieve proposals via CIDs, verify data availability, execute Krum scoring [2], and vote via PBFT. The winning model $w_{r+1}$ is committed *immediately* upon achieving $> \frac{2}{3}$ committee agreement; (4) *Audit Window Opens*—the asynchronous challenge period begins, with participating nodes pinning relevant IPFS data for its duration.

The critical design choice is *immediate commitment*: unlike systems that impose confirmation delays, AC-BlockDFL prioritizes liveness by allowing the committee's decision to take effect instantly. Security is enforced through the challenge mechanism described next.

## 4.3 Challenge Mechanism

The Challenger role is open to any network participant willing to stake the required deposit, ensuring surveillance power remains distributed rather than concentrated. Challengers continuously monitor on-chain records, retrieve proposals

---

**Algorithm 1** AC-BlockDFL Execution Protocol (Instant Update)

---

**Require:** Round $r$, stake-weighted nodes $\mathcal{N}$
**Ensure:** Updated global model $w_{r+1}$
 1: **Phase 1 — Role Assignment:** Select $\mathcal{V}, \mathcal{A}, \mathcal{U} \subset \mathcal{N}$ via stake-weighted randomness.
 2: **Phase 2 — Training & Off-chain Storage:** Each $u \in \mathcal{U}$ trains on local data; $a \in \mathcal{A}$ aggregates into proposal $p_a$.
 3:    Upload $p_a$ to IPFS $\rightarrow$ obtain $\text{CID}_a$; submit $\text{CID}_a$ on-chain.
 4: **Phase 3 — On-chain Consensus & Instant Update:** $\mathcal{V}$ retrieves $\{p_a\}$ via CIDs, verifies data availability, and runs Krum.
 5:    Vote via PBFT. **Commit** $w_{r+1}$ immediately upon majority.
 6:    Record winning $\text{CID}^*$ and voter identities on-chain. Dispense round rewards.
 7: **Phase 4 — Audit Window Opens:** Asynchronous challenge period begins. Nodes pin IPFS data.

---

**Algorithm 2** Asynchronous Challenge Mechanism

---

**Require:** Challenger $ch$, on-chain CID references, IPFS store
**Ensure:** Punishment for malicious committee actions
 1: $ch$ retrieves proposal CIDs from chain, downloads $\{p_a\}$ from IPFS.
 2: $ch$ re-executes Krum on $\{p_a\}$.
 3: **if** outcome mismatches committed $w_{r+1}$ **then**
 4:    $ch$ posts challenge transaction with deposit $D_{\text{challenge}}$.
 5:    **Arbitration Triggered:** All consensus nodes verify independently.
 6:    **if** malicious consensus confirmed by $> \frac{2}{3}$ of network **then**
 7:       **Slash** full stake of colluding validators $\mathcal{V}_{\text{mal}}$.
 8:       Reward Challenger $ch$; distribute remainder to honest participants.
 9:       // Note: Model $w_{r+1}$ is NOT reverted.
10:    **else**
11:       Forfeit $ch$'s deposit $D_{\text{challenge}}$.

---

from IPFS via their CIDs, and independently re-execute the Krum algorithm. Since Krum is fully deterministic—identical inputs yield identical outputs—any discrepancy between the committee's selection and the correct result constitutes irrefutable evidence of misbehavior.

*Challenge Protocol.* Algorithm 2 details the challenge flow. Upon detecting a mismatch, a challenger submits a challenge transaction with deposit $D_{\text{challenge}}$, triggering network-wide arbitration. All nodes download the relevant IPFS data and independently verify. If $> \frac{2}{3}$ of the network confirms malicious consensus, colluding committee members face *full stake slashing*, with rewards distributed to the challenger and honest participants. Failed challenges result in deposit forfeiture, preventing denial-of-service attacks through spurious challenges.

*Dynamic Staking Model.* A key design principle is that all economic parameters derive from *endogenous* system metrics, avoiding external oracle dependencies that could become attack vectors. The maximum rational gain from committee capture is bounded by the round reward $R_{\text{round}}$. We therefore define the slashing amount as $D_{\text{slash}} = \lambda \times R_{\text{round}}$ ($\lambda \gg 1$). In our implementation, an effective $\lambda \approx 100$ means a single slashing event costs the equivalent of 100 rounds of honest participation, eliminating the economic incentive for rational attackers. The challenge deposit $D_{\text{challenge}}$ also scales dynamically with $R_{\text{round}}$, self-adjusting to network economic conditions.

*No-Rollback Policy.* When arbitration confirms misbehavior, AC-BlockDFL slashes malicious actors but *does not revert* the committed model update. Rollback conflicts with blockchain finality guarantees and invites long-range attacks.

Furthermore, rollback would require reverting potentially hundreds of subsequent blocks due to arbitration delays. AC-BlockDFL instead employs *forward correction*: severe economic penalties act as a lethal deterrent against systematic attacks, while subsequent honest training rounds gradually correct the mathematical deviation of any isolated malicious update.

### 4.4 Efficiency and Overhead Analysis

A fundamental dilemma in traditional committee-based consensus is the coupling of security and committee size: augmenting the probability of an honest majority necessitates expanding the committee, severely penalizing communication overhead. AC-BlockDFL resolves this by shifting the security burden to the economic protocol.

*Communication Complexity.* Traditional BlockDFL relies exclusively on threshold security (i.e., driving the probability of committee capture $p_{\text{risk}}$ to near zero). In a network of $N = 100$ nodes with an adversarial fraction $f = 0.3$, ensuring $p_{\text{risk}} < 0.01$ mathematically requires a committee size of $C \geq 9$. This produces a constant per-round PBFT communication complexity of $O(C^2) = O(81)$. This overhead acts as a "preventative premium" exacted uniformly across every round, regardless of whether an attack actually occurs.

In contrast, by enforcing economic security, AC-BlockDFL maintains deterrence even with a smaller committee. Using $C = 7$, the routine communication complexity drops to $O(49)$, a reduction of 39.5%. The heavy $O(N^2)$ PBFT overhead of network-wide arbitration is conditionally triggered only upon an explicit challenge round. The expected communication complexity converges to:

$$E[\text{Comm}] = (1 - p) \cdot O(C^2) + p \cdot \left(O(C^2) + O(N^2)\right) = O(C^2) + p \cdot O(N^2) \tag{1}$$

Because the immense slashing penalty eliminates the rational incentive for PCCA, the objective probability of an attack $p$ approaches 0 under steady-state equilibrium. The expected communication cost is effectively the routine $O(C^2)$, operating on a drastically reduced constant factor.

*Storage Overhead.* Storing raw neural network gradients directly on an immutable ledger guarantees ledger bloat. AC-BlockDFL utilizes a decentralized pinning strategy over IPFS. During the *challenge window* (e.g., $W$ subsequent blocks), nodes are mathematically obligated to pin the IPFS chunks to guarantee data availability for potential challengers. Once the window expires uninterrupted, nodes safely unpin the historical epoch's parameters. This lifecycle management slashes permanent on-chain storage requirements from $O(\text{ModelSize})$ to $O(\text{HashSize})$ (a 32-byte CID), preserving decentralization without burdening node operators with unsustainable disk requirements.

## 5 Threat Model and Security Analysis

The security of AC-BLOCKDFL rests on mitigating sophisticated, economically-driven vulnerabilities while maintaining committee-scale efficiency. In this section, we formalize the adversary model, define the Progressive Committee Capture Attack (PCCA), and establish the dual-layer security guarantees and game-theoretic incentive compatibility of our framework.

### 5.1 Threat Model and Adversary Capabilities

Unlike traditional Byzantine fault tolerance research that assumes purely destructive behavior, our threat model considers a *Rational Adversary* whose primary objective is long-term economic utility maximization and governance control. This aligns with realistic blockchain incentive structures.

*Capabilities.* The adversary controls a fraction of the network nodes, denoted by $f$, where typically $f \leq 0.3$. These malicious nodes are not isolated; they can collude, coordinate voting strategies, and share information. Crucially, the adversary is highly strategic: malicious nodes can perfectly emulate honest behavior to build reputation and accumulate resources during early stages, instantly switching to malicious actions when an advantageous opportunity arises. Furthermore, the adversary has full visibility into the public blockchain state, including stake distributions and historical committee components.

*Limitations.* The adversary is constrained by standard cryptographic assumptions (e.g., unforgeable digital signatures, collision-resistant hash functions) and cannot tamper with immutable historical records on the blockchain. Furthermore, the adversary cannot command an absolute network majority (e.g., > 50%) due to the prohibitively high capital costs. Finally, their actions are governed by economic rationality; they will not execute attacks where the expected financial penalty strictly outweighs the potential gains.

### 5.2 Progressive Committee Capture Attack (PCCA)

Current BlockDFL defenses overwhelmingly focus on data-plane attacks (e.g., data poisoning), relying on robust aggregation algorithms like Krum [2] or Trimmed Mean [19]. However, these algorithms implicitly assume that the validators executing them are honest. We formalize the *Progressive Committee Capture Attack (PCCA)*, a consensus-plane vulnerability that exploits the positive feedback loop of stake-based committee selections: "Stake → Election → Reward → Stake". By subverting the committee, PCCA enables attackers to entirely bypass data-plane defenses. The attack unfolds in two distinct phases:

*Phase 1: Lurking (Shadow Mode).* Because committee selection relies on stake-weighted random sampling, gaining a majority requires significant capital or chance. In this phase, the adversary strictly adheres to protocol rules—submitting high-quality model updates and verifying proposals honestly. This patience allows the adversary to accumulate baseline stake and evade anomaly detection. The adversary waits for a serendipitous election window where malicious nodes coincidentally obtain a supermajority (e.g., > 2/3) of the seats in a single committee.

*Phase 2: Occupying (Capture Mode).* Upon securing a committee supermajority, the adversary drops the honest facade. The specific execution depends on the alignment of the current round's aggregator:

- **Strategic Starvation**: If the aggregator is honest, the malicious committee executes a denial-of-service by systematically voting against valid, high-quality proposals. Consequently, honest aggregators and data providers are denied their block rewards. This starvation stunts the stake growth of honest participants, mathematically guaranteeing that the adversary captures a disproportionately larger share of the systemic inflation. Over successive rounds, this inflates the adversary's probability of being selected for future committees.
- **Full Stack Poisoning**: If the aggregator is also controlled by the adversary, they achieve "full-stack control." The malicious aggregator deliberately accepts poisoned model updates (e.g., backdoor triggers or flipped labels), and the malicious committee force-approves the proposal. This directly degrades the global model's accuracy while monopolizing the round's rewards.

Through PCCA, the adversary's relative stake advantage over honest nodes dynamically shifts, converging to a steady-state advantage $\lim_{t \to \infty} \frac{S_{mal}(t)}{S_{hon}(t)} = \alpha \cdot \frac{S_{mal}(0)}{S_{hon}(0)}$ (where $\alpha > 1$ represents the reward control factor). This locks the system into a permanent governance imbalance, transforming a decentralized network into an oligarchy.

## 5.3 Security Guarantees

Our security model comprises two layers with distinct trust assumptions designed specifically to break the PCCA feedback loop: a *detection layer* requiring only a single honest challenger, and an *arbitration layer* leveraging standard Byzantine fault tolerance.

*Detection Layer: 1-of-N Honest Assumption.* The detection layer operates under an exceptionally weak assumption: among all $N$ network participants, at least one honest node is willing to act as a challenger. This is substantially weaker than the $2/3$ honest majority required by traditional BFT systems, as it requires only the *existence* of a single honest participant rather than coordinated action by a majority. The feasibility of this assumption stems from blockchain's transparency—all proposal CIDs are recorded on-chain with corresponding data publicly accessible via IPFS, enabling any node to independently verify committee decisions.

THEOREM 5.1 (DETECTION COMPLETENESS). *Let $\mathcal{V}_r$ denote the verification committee for round $r$, $\text{Krum}(\{p_a\})$ be the deterministic correct result of executing Krum over all aggregation proposals, and $w_{r+1}$ be the global update actually committed by the committee. If $w_{r+1} \neq \text{Krum}(\{p_a\})$ and there exists at least one honest node $c^*$ among all $N$ participants willing to act as challenger, then this deviation is necessarily detected.*

PROOF. The proof relies on Krum's determinism and the public verifiability of on-chain data. Given identical inputs $\{p_a\}$, any executor obtains the unique output $\text{Krum}(\{p_a\})$ regardless of identity or location. Since all proposal CIDs are recorded on-chain during committee consensus and corresponding data is accessible via IPFS, the honest challenger $c^*$ can: (1) retrieve the identical input set $\{p_a\}$ from IPFS, (2) independently execute Krum locally to obtain $\text{Krum}(\{p_a\})$, and (3) compare against the committed $w_{r+1}$. Any discrepancy constitutes verifiable proof of deviation, enabling $c^*$ to submit a valid challenge transaction. Since verification depends solely on publicly accessible on-chain CIDs, IPFS data, and deterministic computation, the committee cannot evade detection through information hiding or ambiguity.  □

*Arbitration Layer: Global $2/3$ Honest Assumption.* When a challenge is initiated, adjudication authority transfers from the committee to the entire network under standard BFT assumptions: honest nodes must exceed $2/3$ of total nodes, i.e., $N_{\text{total}} > 3f$ where $f$ is the number of Byzantine nodes. During arbitration, all validators download relevant proposals via IPFS, re-execute Krum, and vote on challenge validity through PBFT consensus.

THEOREM 5.2 (PUNISHMENT CERTAINTY). *Let $N_{total}$ be the total network nodes with Byzantine count $f$ satisfying $N_{total} > 3f$. If a challenger successfully detects committee misbehavior per Theorem 5.1 and submits a valid challenge transaction, then the misbehavior is necessarily confirmed during arbitration, and all colluding committee members suffer complete stake slashing.*

PROOF. Upon challenge submission, the smart contract retrieves all proposal CIDs for the disputed round and triggers network-wide re-verification. By Krum's determinism, all honest validators compute identical correct results $\text{Krum}(\{p_a\})$ and can determine whether $w_{r+1}$ deviates. Under $N_{\text{total}} > 3f$, at least $N_{\text{total}} - f > 2N_{\text{total}}/3$ honest nodes participate in arbitration voting. These honest nodes, based on identical deterministic computation, unanimously vote to confirm the deviation. Since PBFT requires $> 2/3$ agreement and honest nodes exceed this threshold, arbitration consensus necessarily succeeds. The smart contract then automatically executes predefined slashing logic, confiscating the full stake of all committee members who endorsed the deviant result. This execution is guaranteed by smart contract determinism and immune to external interference.  □

Theorems 5.1 and 5.2 jointly establish the complete security logic: the former ensures misbehavior is *necessarily discovered*, the latter ensures discovered misbehavior is *necessarily punished*. This dual certainty forms the logical foundation for economic security.

## 5.4 Cost of Attack

We now formalize the capital threshold an attacker must surpass to execute a profitable attack while evading punishment. Two distinct barriers must be overcome: (1) probabilistically winning $> 2/3$ committee seats via random election, and (2) deterministically controlling $\geq 1/3$ of network voting power to block arbitration.

THEOREM 5.3 (ATTACK COST LOWER BOUND). *In AC-BLOCKDFL, an attacker seeking to execute a malicious committee decision while completely evading economic punishment must control stake capital satisfying:*

$$\text{Cost}_{\text{total}} \geq \frac{1}{3}N \cdot \bar{s} \tag{2}$$

*where $N$ is the total network size and $\bar{s}$ is the average stake per node. Even with this capital, the attacker must still probabilistically obtain $> 2/3$ committee seats through random election.*

PROOF. Achieving "successful attack without punishment" requires overcoming two security layers. At the committee level, the attacker must obtain $> 2/3$ validator seats in the target round's random election—a probabilistic event determined by stake proportion that cannot be made certain. At the network level, per Theorem 5.2, once a challenge is initiated and verified, all malicious stakes are fully slashed. To evade punishment, the attacker must control $\geq \lceil N/3 \rceil$ of network voting power to break arbitration liveness by preventing PBFT consensus. This deterministic capital threshold scales linearly with network size as $O(N)$. Since punishment evasion is a logical prerequisite for profitable attack, the total attack cost lower bound is $\frac{1}{3}N \cdot \bar{s}$. $\square$

This theorem reveals a fundamental security amplification: AC-BLOCKDFL's asynchronous audit mechanism elevates the economic barrier from committee-scale $O(C)$ to network-scale $O(N)$. In traditional BlockDFL without post-hoc accountability, attack cost depends solely on controlling a small committee. AC-BLOCKDFL forces attackers to first solve the problem of countering a 2/3 honest network majority before mounting any concrete attack. Given typical deployments where $N \gg C$ (e.g., $N = 100$, $C = 7$ in our experiments), this layered defense provides substantial robustness.

## 5.5 Game-Theoretic Analysis

We employ game-theoretic analysis to demonstrate that honest behavior constitutes the unique Nash equilibrium for all rational participants under AC-BLOCKDFL's economic mechanism.

*Attacker Payoff Model.* For a rational attacker, the decision problem can be modeled as expected payoff computation in a single-shot game. Let $G_{\text{attack}}$ denote the maximum single-round gain from controlling the committee, and $L_{\text{slash}}$ the stake loss from full slashing. The expected payoff is:

$$E[\text{Payoff}] = P_{\text{success}} \cdot G_{\text{attack}} - P_{\text{caught}} \cdot L_{\text{slash}} \tag{3}$$

where $P_{\text{success}}$ is the probability of controlling the committee in a given round, and $P_{\text{caught}}$ is the probability of detection and punishment. Note that $P_{\text{success}}$ measures election outcomes while $P_{\text{caught}}$ measures detection probability—distinct

Table 1. Attacker payoff matrix under AC-BlockDFL

| Strategy | Gain | Loss (if detected) |
|---|---|---|
| Honest behavior | $R_{\text{round}}$ (proportional) | 0 |
| Attack (success) | $\leq 7.0$ units | 500 units |

events at different layers. The attacker can only mount an attack when winning the election; once attacked, Theorems 5.1–5.2 ensure $P_{\text{caught}} \rightarrow 1$ under our dual-layer assumptions. Thus:

$$E[\text{Payoff}] = P_{\text{success}} \cdot (G_{\text{attack}} - L_{\text{slash}}) \tag{4}$$

*Incentive Compatibility Condition.* The sufficient condition for incentive compatibility emerges clearly: whenever $L_{\text{slash}} > G_{\text{attack}}$, expected payoff is strictly negative regardless of $P_{\text{success}}$. Under our endogenous staking model, $L_{\text{slash}} = \lambda \times R_{\text{round}}$ while $G_{\text{attack}}$ is upper-bounded by $C \times R_{\text{round}}$ (monopolizing all validation rewards). Since $\lambda \gg C$ by design, this condition holds stably independent of token market fluctuations.

*Numerical Analysis.* Using our experimental parameters: committee size $C = 7$, per-validator reward 1.0 units, initial stake 100 units. Maximum single-round gain $G_{\text{attack}} \leq 7.0$ units. Upon detection, at least 5 colluding members each lose their full 100-unit stake, yielding $L_{\text{slash}} = 500$ units—approximately $71\times$ the potential gain.

This extreme risk-reward asymmetry ensures negative expected payoff even under optimistic attacker assumptions. For $E[\text{Payoff}] < 0$, we require $P_{\text{caught}} > G_{\text{attack}}/L_{\text{slash}} \approx 1.4\%$. Our security theorems guarantee $P_{\text{caught}} \rightarrow 1$, far exceeding this minimal threshold.

*Nash Equilibrium.* Honest behavior constitutes the unique Nash equilibrium: no rational player can improve their payoff by unilaterally deviating to attack. The slashing mechanism breaks the positive feedback loop enabling gradual committee capture attacks [4]—instead of accumulating stake through manipulation, attackers suffer substantial stake reduction, permanently eliminating their governance influence. This equilibrium remains stable across all token market conditions due to the endogenous stake pricing design.

## 6 Evaluation

We evaluate AC-BlockDFL through systematic experiments designed to validate its defense effectiveness against Progressive Committee Capture Attacks (PCCA). Rather than treating model accuracy as the primary metric, our evaluation focuses on whether the economic security mechanism effectively deters rational adversaries and maintains long-term governance stability. This perspective shift reflects our core design philosophy: when the defense objective transitions from "preventing attacks" to "ensuring attacks are unprofitable," the evaluation metrics should correspondingly shift from model quality to the adversary's economic decision space.

Our experiments adopt a worst-case analysis methodology, assuming adversaries attack whenever possible regardless of economic rationality. This design enables a critical inference: if the mechanism ensures every attack is detected and penalized under worst-case conditions, rational adversaries will preemptively conclude that expected returns are negative and abstain from attacking, allowing the system to naturally converge toward stable equilibrium.

Table 2. Experimental Parameters

| Parameter | Value |
|---|---|
| Training rounds | $R = 300$ (baseline) / $R = 2000$ (long-term) |
| Validator pool size | $N = 100$ |
| Committee size | $C = 7$ |
| Malicious nodes | $M = 30$ (initial stake ratio 30%) |
| Per-round rewards | Validator: 1.0, Aggregator: 1.0, Provider: 0.05 |
| Slashing rule | Full stake confiscation upon successful challenge |

## 6.1 Experimental Setup

We use the MNIST dataset with a standard CNN (two convolutional layers, two fully connected layers) as our federated learning testbed. Training data is distributed IID across clients—a deliberate choice since our defense operates at the consensus layer rather than the data layer. Committee composition and voting outcomes determine attack success, which are logically independent of underlying data distribution characteristics.

Table 2 summarizes the experimental configuration. The 30% initial malicious stake ratio represents a severe threat scenario approaching the theoretical tolerance limit of most Byzantine fault-tolerant systems. Under hypergeometric distribution analysis (Section 6.2), malicious nodes have approximately 2.4% probability of capturing $\geq 5$ of 7 committee seats in any single round. While seemingly modest, this probability accumulates over hundreds to thousands of training rounds, providing ample attack opportunities for rigorous defense validation.

## 6.2 Committee Security Analysis

The probability of $k$ malicious nodes being selected in a committee of size $C$ from a pool $N$ with $M$ malicious nodes follows the hypergeometric distribution:

$$P(X = k) = \frac{\binom{M}{k}\binom{N-M}{C-k}}{\binom{N}{C}} \tag{5}$$

With $N = 100, M = 30, C = 7$, the probability of an adversarial takeover ($k \geq 5$) is $\sim 2.41\%$.

## 6.3 Long-term Governance Equilibrium

The critical question our evaluation addresses is whether short-term penalty effectiveness translates into long-term governance stability where attacks naturally cease. The 2000-round simulation provides definitive evidence.

Figure 2 reveals fundamentally divergent governance trajectories. In BlockDFL, the malicious stake ratio stabilizes around 1.3 after initial fluctuations and persists throughout the experiment. This seemingly modest advantage masks a profound governance crisis: the 1.3× ratio translates to significantly elevated committee election probabilities, sustaining continuous attack capability across 2000 rounds. Without accountability mechanisms, adversaries reinforce their stake advantage through each successful capture, confirming the positive feedback loop predicted in Section 5.2.

AC-BlockDFL exhibits a starkly different pattern. The malicious stake ratio undergoes five distinct step-wise decreases at rounds 15, 136, 695, 815, and 1332, declining from the initial 1.0 to a final 0.37. This terminal value indicates that malicious nodes retain barely one-third the average stake of honest participants—a $1.3/0.37 \approx 3.5\times$ difference from BlockDFL representing fundamental governance reversal rather than incremental improvement.

Table 3 provides direct causal evidence for the stake trajectories. BlockDFL records 107 committee capture events (averaging one per 19 rounds), none receiving economic sanction. AC-BlockDFL records only 5 attacks, all successfully
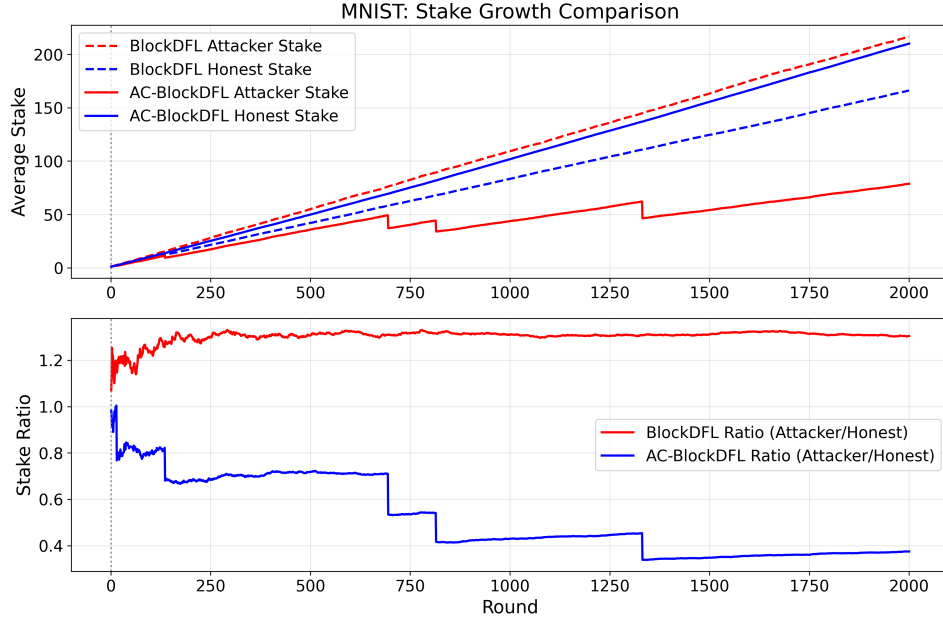
Fig. 2. Stake dynamics over 2000 rounds. BlockDFL exhibits persistent governance imbalance with malicious stake ratio stabilizing at 1.3×, while AC-BlockDFL achieves progressive purification through five slashing events, reducing malicious stake to 0.37× of honest nodes.

Table 3. Attack Event Statistics (2000-round Simulation)

| Metric | BlockDFL | AC-BlockDFL |
|---|---|---|
| Total attacks | 107 | 5 |
| Strategic starvation | 18 | 2 |
| Full-stack poisoning | 89 | 3 |
| Detected & slashed | 0 | 5 (100%) |
| Final stake ratio | 1.30 | 0.37 |

detected and penalized—a >20× reduction stemming from two reinforcing mechanisms: slashing directly depletes the malicious stake base, reducing subsequent committee election probability; simultaneously, diminished stake ratios raise the difficulty of achieving ≥5/7 committee control even when selected.

The increasing intervals between slashing events provide key evidence of convergence toward equilibrium. Specifically: 121 rounds between events 1–2, 559 rounds between 2–3, 120 rounds between 3–4, 517 rounds between 4–5, and 668 rounds of silence following the fifth event through experiment termination. This pattern is not statistical noise but a mathematical consequence of stake depletion: as malicious stake fraction decreases from 30% toward 20%, the single-round probability of achieving committee control drops from ~2.4% to <0.5%, directly manifesting as attack window rarefaction. The 668-round silent period following the final slashing confirms the system has converged to a state where attacks become structurally improbable.
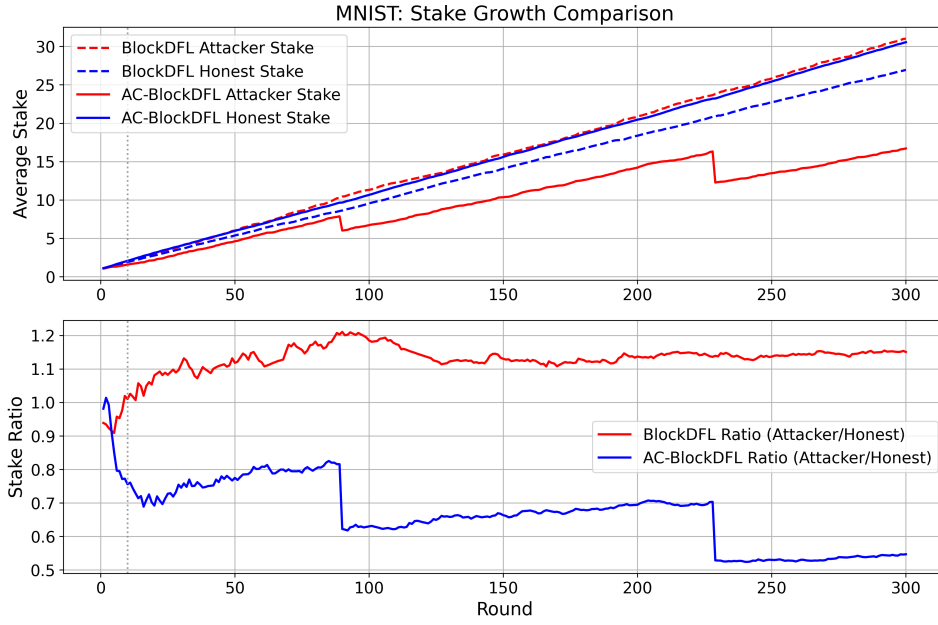
Fig. 3. Stake evolution comparison (300-round baseline). AC-BlockDFL demonstrates immediate stake ratio drops upon each slashing event, while BlockDFL shows continuous malicious stake accumulation.

## 6.4 Immediate Mechanism Response

The 300-round baseline experiment provides a controlled window for examining the immediate impact of individual slashing events on governance structure.

Figure 3 shows the early-stage stake trajectories. In BlockDFL, 10 committee capture events occur over 300 rounds (4 strategic starvation, 6 full-stack poisoning), all unpunished, enabling the malicious stake ratio to climb steadily from 1.0 toward 1.15. AC-BlockDFL records only 2 attacks at rounds 90 and 229, both detected and slashed with 100% accuracy.

The first slashing event illustrates the mechanism's precision. By round 90, malicious nodes had accumulated a 1.25 stake ratio through honest participation, translating to elevated committee selection probability. When 5 malicious nodes achieved committee control and executed full-stack poisoning, a challenger detected the deviation by locally re-executing Krum aggregation and submitted a challenge transaction. Upon arbitration confirmation, the smart contract automatically confiscated the full stakes of all 5 colluding members. The economic impact was immediate and severe: the malicious stake ratio plummeted from 1.25 to 0.62—a single event reversing the adversary's 25% lead into a 38% deficit.

This magnitude of impact warrants careful interpretation. The slashed nodes lost not merely the current round's potential gains (bounded by ~7.0 reward units) but their entire accumulated stake from 89 rounds of honest participation. More critically, stake-zeroed nodes are effectively excluded from future high-reward role elections, constituting "permanent governance exclusion" that degrades long-term attack capability beyond the immediate economic penalty.
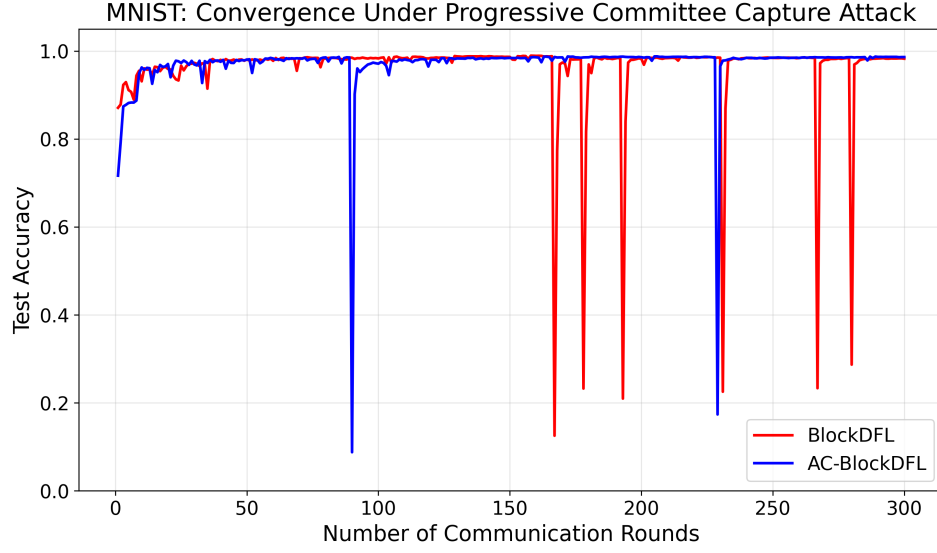
Fig. 4. Model accuracy convergence comparison. AC-BlockDFL exhibits smoother training dynamics with fewer disruption-recovery cycles.

The second slashing at round 229 reduced the stake ratio from 0.70 to 0.52. The 139-round interval between attacks (versus BlockDFL's average of 30 rounds) directly reflects the first slashing's suppressive effect on attack opportunity windows.

### 6.5 Service Quality Under Security Guarantees

A critical concern is whether security guarantees impose unacceptable performance costs. We evaluate both system availability and model convergence quality.

*System Availability.* We define minimum unavailability rate as the fraction of rounds where model performance is significantly degraded due to full-stack poisoning attacks. Each attack requires approximately 5–25 rounds for federated learning's self-healing mechanism to restore accuracy. Using the conservative 5-round estimate, BlockDFL's 89 full-stack attacks yield a minimum unavailability rate of $89 \times 5/2000 = 22.3\%$. AC-BlockDFL achieves $3 \times 5/2000 = 0.75\%$—a $>96\%$ improvement attributable entirely to attack frequency suppression rather than enhanced per-attack resilience.

*Model Convergence.* Figure 4 compares accuracy trajectories over 300 rounds. BlockDFL exhibits pronounced sawtooth patterns corresponding to its 6 full-stack poisoning events, with each attack causing sharp accuracy drops followed by multi-round recovery periods. AC-BlockDFL's curve is notably smoother, experiencing only 2 disruptions. Despite one severe attack at round 90 dropping accuracy to 9.5% (near random-guess baseline for MNIST's 10-class task), the system recovered within ~20 rounds. Final accuracies are comparable (98.26% vs. 98.63%), confirming that the "no-rollback" design philosophy (Section 4.3) is practically sound: federated learning's iterative nature provides inherent self-healing capability, obviating the coordination overhead of state rollback mechanisms.

Table 4. Communication Complexity Comparison

| Scheme | Complexity | Overhead (MB/round) |
|---|---|---|
| Full BFT | $O(N^2)$ | 25.4 |
| BlockDFL | $O(C^2)$ | 4.2 |
| AC-BlockDFL | $O(pN^2 + C^2)$ | 4.3 |

*Communication Efficiency.* As analyzed in Section 6.5 and summarized in Table 4, AC-BlockDFL achieves $O(C^2)$ communication complexity under equilibrium conditions where the challenge trigger probability $p \rightarrow 0$. Compared to approaches requiring equivalent security guarantees through full replication, this represents approximately 40% reduction in per-round communication overhead while maintaining the same Byzantine tolerance threshold.

## 6.6 Summary

Our evaluation validates AC-BlockDFL's defense effectiveness through three complementary lenses. At the micro level, each malicious committee decision triggers immediate detection and slashing with 100% accuracy. At the macro level, five slashing events progressively reduce the malicious stake ratio from 1.0 to 0.37, with increasing inter-event intervals and a terminal 668-round silent period confirming convergence to attack-free equilibrium. Service quality analysis demonstrates that these security guarantees impose minimal performance cost: unavailability rate drops from 22.3% to 0.75%, while model convergence remains uncompromised. These results complete the inference chain: worst-case testing proves all attacks are detected; rational adversaries therefore anticipate penalties and abstain; the system operates at designed efficiency under the resulting equilibrium.

## 7 Conclusion

This paper identifies and formalizes the Progressive Committee Capture Attack (PCCA), demonstrating how rational adversaries can systematically compromise committee-based blockchain federated learning systems through strategic stake accumulation. Our long-horizon simulations confirm that conventional committee architectures exhibit stake ossification and governance capture under sustained attack.

To address this threat, we propose AC-BlockDFL, an audit-driven committee architecture that decouples security guarantees from committee size. The key insight underlying our design is a paradigm shift from *threshold security*—which seeks to minimize the probability of committee compromise—to *economic security*—which ensures that even successful compromise yields negative expected utility for rational adversaries. Through asynchronous auditing and the internal slashing protocol, AC-BlockDFL achieves progressive purification of malicious participants while maintaining the efficiency benefits of small committees.

Our experimental results validate three principal contributions: (1) formal threat modeling of PCCA with empirical verification of its feasibility; (2) demonstration that slashing mechanisms effectively break the positive feedback loop of malicious stake accumulation, internalizing the externalities of adversarial behavior; and (3) evidence that shifting from preventive to reactive security breaks the tight coupling between security guarantees and communication overhead, enabling practical deployment in resource-constrained edge computing scenarios.

## References

[1] Anonymous. 2024. FedChain: Secure and Efficient Federated Learning via Blockchain. *Under Review* (2024).

[2] Peva Blanchard et al. 2017. Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent. In *Advances in Neural Information Processing Systems (NeurIPS)*.

[3] Miguel Castro and Barbara Liskov. 1999. Practical Byzantine Fault Tolerance. In *Proceedings of the 14th USENIX Conference on Operating Systems Design and Implementation (OSDI)*.

[4] Jonathan Chiu and Thorsten Koeppl. 2018. The incentives of blockchain and the optimal design of cryptocurrencies. *Review of Financial Studies* (2018).

[5] M. Elmahallawy et al. 2025. Decentralized Federated Learning for Satellite Networks. *arXiv preprint arXiv:2501.xxxxx* (2025).

[6] Yossi Gilad et al. 2017. Algorand: Scaling Byzantine Agreements for Cryptocurrencies. In *Proceedings of the 26th Symposium on Operating Systems Principles (SOSP)*.

[7] D. Li et al. 2021. A Blockchain-Based Federated Learning Framework with Committee Consensus. *IEEE Network* (2021).

[8] Y. Liu et al. 2021. Blockchain-Enabled Federated Learning for Vehicular Networks. *IEEE Transactions on Vehicular Technology* (2021).

[9] Y. Lu et al. 2020. Blockchain-enabled federated learning for industrial IoT. *IEEE Transactions on Industrial Informatics* (2020).

[10] H. Nguyen et al. 2024. FedBlock: A Blockchain-Based Federated Learning Framework with Adaptive Committee Selection. *IEEE Transactions on Parallel and Distributed Systems* (2024).

[11] Shiva Raj Pokhrel. 2021. Blockchain brings trust to collaborative drones and LEO satellites: An intelligent decentralized learning in the space. *IEEE Sensors Journal* 21, 14 (2021), 15731–15741.

[12] S. R. Pokhrel and J. Choi. 2020. Autonomous Vehicles in 5G and Beyond: A Blockchain-Based Federated Learning Approach. *IEEE Transactions on Intelligent Transportation Systems* (2020).

[13] Jiaming Qin et al. 2024. BlockDFL: A Blockchain-based Fully Decentralized Peer-to-Peer Federated Learning Framework. In *Proceedings of the IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing (CCGrid)*.

[14] Y. Qu et al. 2020. Decentralized Federated Learning: A Survey. *IEEE Communications Surveys & Tutorials* (2020).

[15] Y. Ren et al. 2024. A scalable blockchain-enabled federated learning architecture for edge computing. *PLOS ONE* (2024).

[16] Muhammad Shayan et al. 2021. Biscotti: A Blockchain System for Private and Secure Federated Learning. In *IEEE Transactions on Parallel and Distributed Systems*.

[17] J. Weng et al. 2021. DeepChain: Auditable and Privacy-Preserving Deep Learning with Blockchain. *IEEE Transactions on Dependable and Secure Computing* (2021).

[18] Y. Wu et al. 2024. A Sharded Blockchain-Based Secure Federated Learning Framework for LEO Satellite Networks. *IEEE Transactions on Network and Service Management* (2024).

[19] D. Yin et al. 2018. Byzantine-Robust Distributed Learning: Towards Optimal Statistical Rates. In *International Conference on Machine Learning (ICML)*.