# Network Structures as an Attack Surface: Topology-Based Privacy Leakage in Federated Learning

Murtaza Rangwala⬤, *Member, IEEE,* Richard O. Sinnott⬤, and Rajkumar Buyya⬤, *Fellow, IEEE*

*Abstract*—**Federated learning systems increasingly rely on diverse network topologies to address scalability and organizational constraints. While existing privacy research focuses on gradient-based attacks, the privacy implications of network topology knowledge remain critically understudied. We conduct the first comprehensive analysis of topology-based privacy leakage across realistic adversarial knowledge scenarios, demonstrating that adversaries with varying degrees of structural knowledge can infer sensitive data distribution patterns even under strong differential privacy guarantees. Through systematic evaluation of 4,720 attack instances, we analyze six distinct adversarial knowledge scenarios: complete topology knowledge and five partial knowledge configurations reflecting real-world deployment constraints. We propose three complementary attack vectors: communication pattern analysis, parameter magnitude profiling, and structural position correlation, achieving success rates of 84.1%, 65.0%, and 47.2% under complete knowledge conditions. Critically, we find that 80% of realistic partial knowledge scenarios maintain attack effectiveness above security thresholds, with certain partial knowledge configurations achieving performance superior to the baseline complete knowledge scenario. To address these vulnerabilities, we propose and empirically validate structural noise injection as a complementary defense mechanism across 808 configurations, demonstrating up to 51.4% additional attack reduction when properly layered with existing privacy techniques. These results establish that network topology represents a fundamental privacy vulnerability in federated learning systems while providing practical pathways for mitigation through topology-aware defense mechanisms.**

*Index Terms*—**Federated Learning, Network Topology, Privacy Attacks, Differential Privacy, Data Distribution Inference, Communication Patterns**

## I. INTRODUCTION

**F**EDERATED learning (FL) has emerged as a paradigmatic approach that enables collaborative machine learning between distributed datasets while preserving data privacy [1], [2]. This foundational premise has driven widespread adoption across healthcare [3], finance [4], and edge computing domains [5]. However, mounting evidence demonstrates substantial privacy leakage through gradient exchanges [6], [7], model parameters [8], and communication metadata [9].

While existing privacy research focuses exclusively on protecting gradient and model content, we identify network topology as a fundamental and previously unrecognized attack surface. Real-world FL deployments increasingly leverage diverse topologies—hierarchical, decentralized, and hybrid configurations—to address scalability bottlenecks [10], communication constraints [11], and organizational requirements [12].

Murtaza Rangwala, Richard O. Sinnott and Rajkumar Buyya are with School of Computing and Information Systems, the University of Melbourne, Australia. Email: mrangwala@student.unimelb.edu.au, {rsinnott, rbuyya}@unimelb.edu.au

These operational necessities inadvertently expose structural information enabling powerful inference attacks orthogonal to existing threat models.

Current privacy mechanisms such as differential privacy [13] and cryptographic approaches [14], [15] assume adversaries lack knowledge of the network structure. This assumption is fundamentally flawed in practice, as topology information is readily observable through infrastructure access, organizational relationships, regulatory compliance, or operational roles of potential adversaries.

We systematically combine and analyze three complementary attack vectors that exploit communication patterns, parameter magnitudes, and structural correlations to infer sensitive data distribution patterns. These attacks operate through information channels orthogonal to content-based protections, exploiting observable metadata and behavioral signatures that remain visible despite encryption or noise injection.

Through systematic evaluation of 4,720 attack instances across six adversarial knowledge scenarios, we demonstrate that topology-based attacks achieve upper bound success rates of 84.1%, 65.0%, and 47.2% respectively under complete topology knowledge. Critically, 80% of realistic partial knowledge scenarios maintain effectiveness above security thresholds, with attacks persisting under strong differential privacy and remaining viable across enterprise scales. This paper makes four fundamental contributions:

- **Systematic Attack Surface Analysis:** We identify network topology as a critical privacy vulnerability and systematically analyze three complementary attack vectors that exploit structural information channels orthogonal to existing threat models.
- **Realistic Threat Assessment:** We evaluate privacy leakage across six practical adversarial knowledge scenarios, demonstrating that limited structural information suffices for effective attacks.
- **Defense Inadequacy Analysis:** We prove current privacy mechanisms provide insufficient protection, with topology-based attacks maintaining effectiveness despite state-of-the-art content protections.
- **Practical Defense Validation:** We propose and empirically validate structural noise injection as a complementary defense mechanism, demonstrating up to 51.4% additional attack reduction when properly layered with existing privacy techniques.

The rest of the paper is structured as follows. Section II positions our work within the broader landscape of federated learning privacy research and identifies critical gaps. Section III formalizes our threat model and attack framework. Section IV details our experimental methodology. Section V

TABLE I
PRIVACY RESEARCH LANDSCAPE: ATTACK VECTORS AND FUNDAMENTAL LIMITATIONS

| Research Area | Attack Vector | Assumptions | Defense Focus | Topology | Fundamental Limitation |
|---|---|---|---|---|---|
| Gradient Inversion [6], [7] | Data reconstruction | White-box access | Content protection | Star (implicit) | Individual-level, requires gradient access |
| Advanced Inversion [16] | Architecture manipulation | Malicious server | Secure aggregation | Star only | Single topology, server-side assumption |
| Inference Attacks [9], [17] | Statistical analysis | Model access | Basic DP | Star (implicit) | Content-based, no structural exploitation |
| Cryptographic FL [14], [15] | Defense mechanism | Semi-honest parties | Content encryption | Limited scope | Structural information exposure |
| DP-FL [18], [19] | Privacy framework | Independent participants | Noise injection | Star assumption | Ignores topology correlations |
| Topology Inference [20] | Structure discovery | Behavioral observation | None | Multiple (DFL) | Discovers topology, not data patterns |
| **Our Contribution** | **Structural inference** | **Realistic knowledge** | **Topology-aware analysis** | **Comprehensive** | **Addresses fundamental gaps** |

presents comprehensive experimental results. Section VI analyzes the implications of our findings. Section VII proposes and validates a topology-aware defense mechanism and discusses future research directions. Section VIII concludes the paper.

## II. RELATED WORK

Current federated learning privacy research treats network topology as an operational design choice rather than a privacy-critical component. This fundamental oversight creates systematic vulnerabilities that our research exposes and analyzes comprehensively.

### A. Content-Based Privacy Attacks

Existing attacks target gradient information and model parameters through content-based inference techniques. Gradient inversion approaches [6], [7], [21] reconstruct training data by optimizing dummy inputs to match observed gradients, with recent advances such as LOKI [16] achieving 86% reconstruction rates across 100 clients. Membership inference [17] and property inference [8] attacks extract participation indicators and data set characteristics, while preference profiling [9] reveals user preferences with accuracy of 90-98%.

Individual components of topology-based vulnerabilities have been explored in isolation. Communication metadata analysis [22], [23] examines timing and pattern vulnerabilities, while parameter analysis techniques exist for model extraction and membership inference. However, no prior work systematically combines these approaches to exploit network topology structure, nor evaluates their effectiveness across diverse topological configurations and realistic adversarial knowledge constraints.

### B. Privacy Protection Mechanisms

Differential privacy provides the theoretical foundation for FL privacy protection [18], with extensions examining client subsampling amplification [19] and data heterogeneity effects [24]. Cryptographic approaches, including homomorphic encryption [15] and secure aggregation [14] protect content during computation and transmission.

Both approaches exhibit a fundamental limitation: they protect message content while leaving structural information completely exposed. Privacy accounting treats participants as independent entities, failing to model topology-induced correlations. This creates critical vulnerability gaps where communication patterns, behavioral signatures, and positional relationships remain observable regardless of content protection mechanisms.

### C. Topology-Aware Federated Learning

Network topology research focuses primarily on efficiency optimization rather than privacy implications. Work analyzing convergence properties [10], hierarchical benefits [11], and decentralized approaches [25] treats topology as performance engineering rather than potential attack surface. Comprehensive surveys [26] categorize topology-aware systems by communication patterns and scalability properties without systematic privacy analysis.

Some work has examined privacy in decentralized settings. Bellet et al. [27] analyzed privacy in peer-to-peer learning, while others have explored privacy-preserving protocols for specific decentralized architectures [28]. However, these studies focus on adapting existing content-based privacy mechanisms to decentralized settings rather than analyzing topology-specific vulnerability patterns.

Most relevant to our work, Feng et al. [20] recently demonstrated topology inference attacks in decentralized federated learning, showing that adversaries can reconstruct network topology from model behavior alone with over 85% F1-score across multiple datasets and network configurations. Their work demonstrates that topology structure can be inferred from federated learning systems, while our research assumes

topology knowledge and exploits it to infer sensitive data distribution patterns. These complementary approaches reveal both the feasibility of topology discovery and the privacy risks it enables.

### D. Research Gaps and Positioning

Existing privacy frameworks exhibit three fundamental limitations that enable topology-based vulnerabilities. First, most work implicitly assumes star topologies, missing how structural diversity creates distinct attack surfaces across other hierarchical, and decentralized configurations. Second, current attacks target protected content through strong adversarial assumptions, missing structural vulnerabilities that exploit observable metadata under much weaker assumptions. Third, privacy accounting frameworks ignore systematic correlations that topology creates between participants, enabling exploitation of organizational relationships and coordination metadata.

Table I summarizes these critical gaps across the privacy research landscape. Our work addresses this oversight by demonstrating topology as a fundamental attack surface orthogonal to existing threat models, necessitating topology-aware privacy mechanisms that extend beyond traditional content protection approaches.

## III. THREAT MODEL AND ATTACK FRAMEWORK

This section establishes our formal threat model and presents three complementary topology-based attacks that exploit network structure to infer sensitive data distribution patterns in federated learning systems.

### A. System Model and Architecture

Consider a federated learning system comprising $n$ participants connected via network topology $\mathcal{G} = (\mathcal{P}, E)$, where $\mathcal{P}$ represents participant nodes and $E$ denotes communication links. Each participant $P_i$ holds a local dataset $\mathcal{D}_i$ with samples drawn from classes $\mathcal{C} = \{c_1, c_2, \ldots, c_k\}$. The topology $\mathcal{G}$ determines communication patterns for parameter aggregation, supporting both centralized and decentralized configurations.

We establish clear trust boundaries where participants trust the learning protocol but may attempt passive inference attacks, network infrastructure is observable but not modifiable by adversaries, and aggregation follows standard protocols.

### B. Adversarial Model and Knowledge Scenarios

We formalize the adversary as a tuple $\mathcal{A} = (\mathcal{K}, \mathcal{I}, \mathcal{O}, \mathcal{R})$ where $\mathcal{K}$ represents adversarial knowledge, $\mathcal{I}$ denotes inference goals, $\mathcal{O}$ specifies observation capabilities, and $\mathcal{R}$ defines behavioral restrictions.

Our threat model encompasses six distinct knowledge scenarios reflecting realistic adversarial capabilities, as illustrated in Figure 1. Complete topology knowledge ($\mathcal{K}_{\text{complete}}$) assumes the adversary possesses complete network structure $\mathcal{G}$, participant identities, and organizational constraints. This represents theoretical worst-case bounds and scenarios where network infrastructure providers have full architectural visibility.

Statistical topology knowledge ($\mathcal{K}_{\text{statistical}}$) assumes the adversary knows general structural properties (e.g., hierarchical vs. decentralized, average connectivity, clustering coefficients) but lacks exact connection details. This reflects external adversaries with domain expertise about federated learning architectures or traffic analysis capabilities.

Local neighborhood knowledge ($\mathcal{K}_{\text{local}}$) assumes the adversary observes network structure within $k$ hops of their position, representing compromised nodes or insider threats with limited visibility. We analyze both 1-hop scenarios covering immediate neighbors and 2-hop scenarios covering extended neighborhoods.

Organizational structure knowledge ($\mathcal{K}_{\text{organizational}}$) assumes the adversary possesses institutional relationship information, mapping nodes to organizational groups or departments. We evaluate both coarse-grained groupings (3 organizational clusters) and fine-grained groupings (5 organizational clusters). These scenarios reflect realistic deployment constraints where academic partnerships, healthcare consortiums, and financial collaborations inherently expose structural information through operational requirements and public disclosures.

### C. Inference Goals and Success Metrics

Across all knowledge scenarios, the adversary seeks to infer group-level data distribution patterns $\{\delta_i\}_{i=1}^{n}$ where $\delta_i$ represents the class distribution at participant $P_i$, identify nodes containing sensitive classes, and discover correlations between network position and data characteristics.

We define attack success using advantage over random guessing:

$$\text{Adv}_{\mathcal{A}}^{\text{topology}} = \left| \Pr[\mathcal{A}(\mathcal{G}, \{\theta^{(t)}\}) = \text{true}] - \frac{1}{|\Omega|} \right| \qquad (1)$$

where $\Omega$ denotes the space of possible distributions. We employ a conservative 30% threshold for attack success across all evaluation metrics to ensure that identified vulnerabilities reflect systematic exploitation of structural relationships rather than statistical artifacts.

### D. Observation Capabilities and Behavioral Restrictions

The adversary can monitor communication metadata including frequency and timing, and observe parameter updates during the parameter transfer phase—after participants complete local training but before aggregation occurs. Critically, the adversary cannot access raw data $\mathcal{D}_i$ or individual gradients during local training phases. This observation capability reflects realistic deployment scenarios where adversaries have privileged network infrastructure access or operate as compromised intermediate nodes with visibility into parameter transmissions.

Our threat model accounts for contemporary DP implementations where noise is added to gradients during local training (DP-SGD), but transmitted parameter updates may still reveal statistical signatures exploitable through temporal analysis across multiple rounds rather than individual noisy updates.
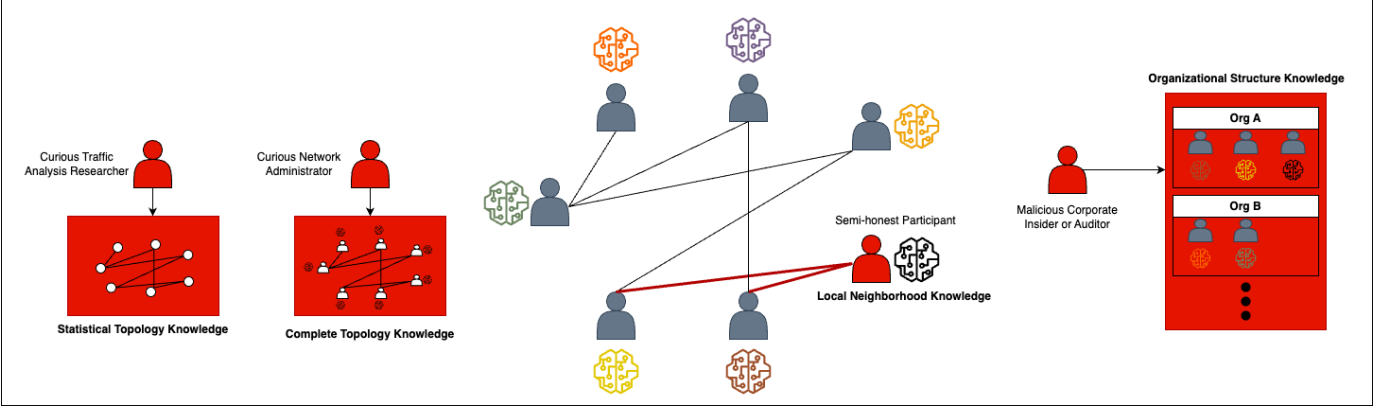
Fig. 1. Adversarial knowledge scenarios in distributed learning networks, ranging from complete topology knowledge to realistic partial knowledge constraints including neighborhood visibility, statistical properties, and organizational structure awareness.

Under comprehensive cryptographic protection, parameter magnitudes would be obscured, rendering certain attacks ineffective. However, communication pattern analysis maintains effectiveness by examining observable metadata that cannot be encrypted without altering coordination protocols. Many production deployments rely primarily on differential privacy due to computational overhead constraints, making parameter observation realistic.

The observation model reflects adversaries positioned as network infrastructure providers, compromised intermediate nodes, or malicious participants who can observe parameter exchanges while maintaining semi-honest behavioral restrictions prohibiting message modification or collusion.

### E. Attack Surface Analysis

Our analysis identifies three distinct information leakage channels that topology-aware adversaries can exploit.

**Communication pattern leakage:** Network topologies create observable communication signatures. In decentralized settings, nodes with similar data converge faster, requiring fewer consensus iterations that manifest as reduced communication frequency in later training rounds. Hierarchical topologies exhibit predictable routing patterns correlated with organizational data placement.

**Parameter magnitude leakage:** Data distribution heterogeneity systematically affects parameter update magnitudes. Nodes training on rare classes exhibit larger gradient norms due to unbalanced loss landscapes, while homogeneous local data leads to smoother optimization trajectories with smaller, more stable parameter updates. Statistical moments of parameter sequences encode distributional signatures that persist under differential privacy noise addition.

**Structural correlation leakage:** Real-world deployments often correlate topology position with data characteristics. Geographic proximity in network topology may reflect similar demographic distributions, organizational hierarchies embedded in network structure determine departmental data access patterns, and high-degree nodes in hierarchical settings aggregate multiple data sources, creating distinct parameter signatures.

---

**Algorithm 1** Communication Pattern Attack
1: **Input:** Communications log $\mathcal{L}$, participants $\mathcal{P}$, groups $k$
2: **Output:** Clustering $\mathcal{K}_{\text{nodes}}$, success $s_{\text{comm}}$
3:
4: Initialize matrix $M \leftarrow \mathbf{0}_{|\mathcal{P}| \times |\mathcal{P}|}$
5: **for** each communication $(P_u, P_v, t) \in \mathcal{L}$ **do**
6: $\quad M[u, v] \leftarrow M[u, v] + 1$
7: **end for**
8:
9: Features: $F_{\text{comm}} \leftarrow [M, M^T] \in \mathbb{R}^{|\mathcal{P}| \times 2|\mathcal{P}|}$
10: Cluster: $\mathcal{K}_{\text{nodes}} \leftarrow \text{KMeans}(F_{\text{comm}}, k)$
11: Coherence: $s_{\text{comm}} \leftarrow \max_i \frac{|\mathcal{K}_{\text{nodes}}[i]|}{|\mathcal{P}|}$
12: **return** $\mathcal{K}_{\text{nodes}}, s_{\text{comm}}$

---

### F. Attack Implementation Framework

We present three complementary attacks that operationalize these leakage channels through concrete algorithmic implementations exploiting detectable behavioral patterns.

*1) Communication Pattern Attack:* This attack exploits communication pattern leakage by analyzing communication frequency patterns to cluster nodes based on data similarity. Nodes with similar data distributions require fewer communication rounds to reach consensus in decentralized settings.

The attack constructs a communication frequency matrix $M \in \mathbb{R}^{n \times n}$ where $M[i, j]$ counts interactions between participants $P_i$ and $P_j$. Using communication logs $\mathcal{L} = \{(P_u, P_v, t)\}$ representing messages from participant $P_u$ to $P_v$ at time $t$, we extract symmetrized features $F_{\text{comm}} = [M, M^T]$ and apply clustering to identify communication-based groupings.

We measure attack effectiveness using cluster coherence ratio, where higher values indicate stronger grouping structure. A threshold of $s_{\text{comm}} \geq 0.3$ indicates successful group identification.

*2) Parameter Magnitude Attack:* This attack exploits parameter magnitude leakage by observing parameter updates during the communication phase—after local training but before aggregation—when participants transmit their model updates. The attack analyzes patterns across multiple training rounds rather than individual updates, making it robust to

**Algorithm 2** Parameter Magnitude Attack

1: **Input:** Updates $\mathcal{U} = \{\theta_i^{(t)}\}$, participants $\mathcal{P}$
2: **Output:** Clustering $\mathcal{K}_{\text{params}}$, score $s_{\text{param}}$
3:
4: **for** each participant $P_i \in \mathcal{P}$ **do**
5:     Norms: $\mathcal{N}_i \leftarrow \{||\theta_i^{(t)}||_2\}_{t=1}^T$
6:     Mean: $\mu_i \leftarrow \frac{1}{T} \sum_{t=1}^T ||\theta_i^{(t)}||_2$
7:     Variance: $\sigma_i^2 \leftarrow \frac{1}{T-1} \sum_{t=1}^T (||\theta_i^{(t)}||_2 - \mu_i)^2$
8:     Trend: $\beta_i \leftarrow \text{LinearReg}(\{1, \ldots, T\}, \mathcal{N}_i)$
9:     Stability: $\varsigma_i \leftarrow \text{std}(\{||\theta_i^{(t)}||_2\}_{t=T-2}^T)$
10: **end for**
11:
12: Features: $F_{\text{param}} \leftarrow [\boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\beta}, \boldsymbol{\varsigma}]$
13: Normalize: $F_{\text{param}} \leftarrow \text{StandardScaler}(F_{\text{param}})$
14: Cluster: $\mathcal{K}_{\text{params}} \leftarrow \text{KMeans}(F_{\text{param}}, 2)$
15: Score: $s_{\text{param}} \leftarrow \text{Silhouette}(F_{\text{param}}, \mathcal{K}_{\text{params}})$
16: **return** $\mathcal{K}_{\text{params}}, s_{\text{param}}$

---

**Algorithm 3** Topology Structure Attack

1: **Input:** Topology $\mathcal{G} = (\mathcal{P}, E)$, updates $\mathcal{U}$
2: **Output:** Correlations $\boldsymbol{\rho}$, success $s_{\text{topo}}$
3:
4: Initialize: $\mathbf{deg}, \mathbf{pos}, \mathbf{cent}, \mathbf{norm}, \mathbf{var} \leftarrow \emptyset$
5: **for** each participant $P_i \in \mathcal{P}$ **do**
6:     Degree: $\mathbf{deg}[i] \leftarrow |\{P_j : (P_i, P_j) \in E\}|$
7:     Position: $\mathbf{pos}[i] \leftarrow i$
8:     Central: $\mathbf{cent}[i] \leftarrow \mathbf{1}[\mathbf{deg}[i] > \text{median}(\mathbf{deg})]$
9:     Avg norm: $\mathbf{norm}[i] \leftarrow \frac{1}{T} \sum_{t=1}^T ||\theta_i^{(t)}||_2$
10:     Variance: $\mathbf{var}[i] \leftarrow \text{Var}(\{||\theta_i^{(t)}||_2\})$
11: **end for**
12:
13: $\rho_1 \leftarrow \text{Corr}(\mathbf{pos}, \mathbf{norm})$
14: $\rho_2 \leftarrow \text{Corr}(\mathbf{deg}, \mathbf{norm})$
15: $\rho_3 \leftarrow \text{Corr}(\mathbf{cent}, \mathbf{var})$
16: $s_{\text{topo}} \leftarrow \max(|\rho_1|, |\rho_2|, |\rho_3|)$
17: $\boldsymbol{\rho} \leftarrow [\rho_1, \rho_2, \rho_3]$
18: **return** $\boldsymbol{\rho}, s_{\text{topo}}$

---

differential privacy noise injection.

For each participant $P_i$ whose updates are observable, we extract parameter norm sequences $\{||\theta_i^{(t)}||_2\}_{t=1}^T$ and compute statistical features: mean magnitude $\mu_i$, temporal standard deviation $\sigma_i$, linear trend $\beta_i$, and convergence stability $\varsigma_i = \text{std}(\{||\theta_i^{(t)}||_2\}_{t=T-2}^T)$.

We employ the silhouette score $s_{\text{param}} \in [-1, 1]$ to measure cluster quality, where values above 0.3 indicate meaningful separation corresponding to distributional differences.

*3) Topology Structure Attack:* This attack exploits structural correlation leakage by computing correlations between network positions and data characteristics due to organizational or geographical constraints.

We extract topology features including node degree $d_i = |\text{neighbors}(P_i)|$, position identifier $\text{pos}_i$, and centrality indicator $\mathbf{1}[d_i > \text{median}(\{d_j : P_j \in \mathcal{P}\})]$. Parameter features include average norm and variability.

We define attack success as $s_{\text{topo}} \geq 0.3$, indicating moderate to strong correlation between topology structure and parameter characteristics.

*G. Experimental Data Partitioning Strategies*

To evaluate topology-based vulnerabilities under realistic conditions, we design three data partitioning strategies that reflect how organizational constraints and deployment realities create systematic correlations between network topology and data distribution patterns in practice.

*1) Sensitive Groups (SG) Partitioning:* This partitioning strategy models organizational constraints that concentrate sensitive populations at specific topology positions. Real-world federated deployments often reflect organizational structure, where departments, geographic regions, or demographic groups occupy predictable network positions.

We partition classes into sensitive $\mathcal{C}_s \subset \mathcal{C}$ and non-sensitive $\mathcal{C}_n = \mathcal{C} \setminus \mathcal{C}_s$ subsets. Participants at certain positions are assigned concentrated distributions on sensitive classes, measured using total variation distance:

$$\text{TV}(\delta(P_i), \text{Uniform}(\pi(\phi(P_i)))) \geq 0.3 \quad (2)$$

where $\pi : \mathcal{S} \rightarrow \{\mathcal{C}_s, \mathcal{C}_n\}$ maps positions to class subsets and $\phi : \mathcal{P} \rightarrow \mathcal{S}$ denotes the topology position function.

*2) Topology Correlated (TC) Partitioning:* This partitioning strategy creates systematic correlations between network position and class assignment patterns, reflecting how geographic deployments, temporal assignments, or organizational hierarchies influence data distribution.

For ordered topologies, we assign data such that position correlates with dominant class:

$$\mathbb{E}[\text{encode}(\text{dominant\_class}(P_i))] = (\phi(P_i) \bmod |\mathcal{C}|) + \epsilon_i \quad (3)$$

where $\epsilon_i \sim \mathcal{N}(0, 0.25)$ represents positioning noise and $\text{encode} : \mathcal{C} \rightarrow \{0, 1, \ldots, |\mathcal{C}| - 1\}$ maps classes to integers.

For hierarchical topologies, we model specialization through entropy:

$$H(\delta(P_i)) = \begin{cases} \leq 0.5 \log(|\mathcal{C}|) & \text{if } \deg(P_i) = 1 \\ \geq 0.8 \log(|\mathcal{C}|) & \text{if } \deg(P_i) = n - 1 \end{cases} \quad (4)$$

where leaf nodes exhibit low entropy (concentrated distributions) and central nodes exhibit high entropy (diverse distributions).

*3) Imbalanced Sensitive (IS) Partitioning:* This partitioning strategy models severe class imbalances concentrated at specific topology positions, reflecting scenarios where specialized nodes handle rare but sensitive classes.

We designate a subset $\mathcal{P}_r \subset \mathcal{P}$ as rare class holders with $|\mathcal{P}_r| = \lceil 0.3n \rceil$. Let $Y_i$ denote the indicator random variable for whether participant $P_i$ contains the rare class as its dominant class:

$$\Pr[Y_i = 1 | P_i \in \mathcal{P}_r] = 0.9 \quad (5)$$
$$\Pr[Y_i = 1 | P_i \notin \mathcal{P}_r] = 0.1 \quad (6)$$
$$\Pr[P_i \in \mathcal{P}_r] = 0.3 \quad (7)$$

This creates detectable imbalance signals with Kullback-Leibler divergence $\text{KL}(\delta(P_r) || \delta(P_{nr})) \geq 1.2$ nats between rare holders and normal participants.
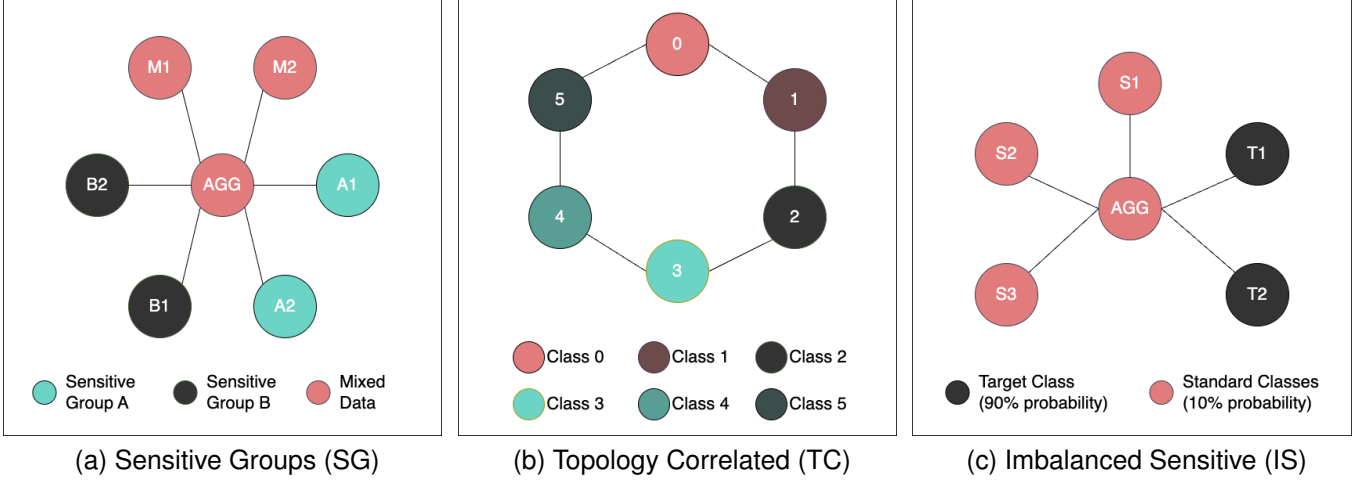
Fig. 2. Three experimental data partitioning strategies creating realistic topology-data correlations for evaluating privacy vulnerabilities.

## H. Privacy Mechanism Integration

Our threat model explicitly accounts for state-of-the-art privacy protections. We evaluate attack robustness under differential privacy with per-participant per-round privacy parameters $\varepsilon \in \{1.0, 4.0, 8.0\}$ representing strong, medium, and weak privacy protection respectively, with failure probability $\delta = 10^{-5}$. The DP mechanism adds calibrated Gaussian noise $\mathcal{N}(0, \sigma^2)$ with:

$$\sigma = \frac{\sqrt{2\log(1.25/\delta)} \cdot S}{\varepsilon} \tag{8}$$

where $S$ represents the global sensitivity of parameter updates.

We analyze privacy amplification through client subsampling with participation rates $q \in \{0.2, 0.3, 0.5\}$ and local data subsampling rates $r \in \{0.5, 0.6, 0.8\}$. Our model assumes static topology knowledge but does not address dynamic network reconfiguration. We acknowledge that comprehensive homomorphic encryption would prevent parameter-based attacks, though such deployments remain rare due to computational constraints. Active attacks involving message modification or collusion are outside our scope.

## IV. EXPERIMENTAL METHODOLOGY AND DESIGN

This section presents our experimental methodology for comprehensive assessment of topology-based privacy leakage in distributed learning systems. We describe our distributed learning framework, experimental infrastructure, attack evaluation protocols, scalability analysis, and reproducibility measures.

## A. Distributed Learning Framework and Infrastructure

We developed Murmura [29], a comprehensive framework enabling topology-aware privacy evaluation across both centralized and decentralized FL. The framework implements centralized FL using FedAvg [1] on star and complete topologies, while supporting decentralized learning through gossip averaging [10] on ring, line, and complete topologies.

Our experiments utilize a Ray-based [30] distributed computing cluster deployed on AWS with 10 G5.2xlarge instances, each providing 8 vCPU cores, 32 GiB system memory, one NVIDIA A10G GPU with 24 GiB video memory, and up to 25 Gbps network performance. The Murmura cluster manager dynamically allocates virtual FL clients across nodes while maintaining experimental control and reproducibility.

## B. Datasets and Differential Privacy Configuration

We evaluate our attacks on the MNIST (Modified National Institute of Standards and Technology) dataset [31] comprising 60,000 training samples across 10 digit classes, and the HAM10000 Skin Lesion Dataset [32] containing approximately 10,000 dermatoscopic images across 7 skin lesion types. Both datasets are processed using convolutional neural networks with GroupNorm and LayerNorm for differential privacy compatibility.

Our differential privacy implementation follows current best practices with per-round, per-participant privacy accounting using Rényi Differential Privacy (RDP) [33] and seamless integration to Opacus [34] for DP-SGD. We employ the standard RDP-to-DP conversion mechanism with composition rules to track cumulative privacy expenditure across training rounds. The failure probability $\delta = 10^{-5}$ is chosen to provide strong privacy guarantees while maintaining computational tractability, following established practices in federated learning privacy research. Gradient clipping uses an L2 norm threshold $C = 1.0$ with automated Gaussian noise calibration. Privacy accounting utilizes RDP with the conversion:

$$\varepsilon(\delta) = \min_{\alpha > 1} \left\{ \frac{\text{RDP}_\alpha + \log(1/\delta)}{\alpha - 1} \right\} \tag{9}$$

## C. Experimental Design and Evaluation Protocol

Our evaluation employs a comprehensive four-phase approach that captures idealized deployment scenarios, realistic adversarial knowledge constraints, operational deployment conditions, and enterprise-scale implications.

**Phase 1** establishes baseline attack effectiveness through exhaustive evaluation across 520 unique configurations, systematically varying datasets, data partitioning strategies, FL paradigms, network topologies, network sizes (5-30 nodes) and differential privacy protection levels. This phase provides theoretical upper bounds on privacy leakage under complete topology knowledge without sampling effects.

**Phase 2** evaluates attack robustness under realistic adversarial constraints through analysis of 2,100 attack instances across five partial knowledge scenarios defined in Section III-B. Using 420 configurations from Phase 1 with meaningful network sizes ($>$5 nodes), we systematically evaluate each partial knowledge configuration against the complete topology baseline established in Phase 1. This phase establishes attack effectiveness under practical adversarial capabilities and identifies minimum knowledge requirements for successful data distribution inference.

**Phase 3** evaluates realistic deployment scenarios incorporating client and data subsampling across 288 targeted configurations. We examine moderate subsampling (50% clients, 80% data), strong subsampling (30% clients, 60% data), and very strong subsampling (20% clients, 50% data). This phase quantifies privacy amplification effects and validates findings under practical deployment constraints.

**Phase 4** addresses enterprise-scale implications through synthetic simulation methodology that enables evaluation of networks with 50-500 nodes while maintaining computational tractability. This phase provides critical insights into scalability patterns and extrapolates findings to production deployment scenarios involving hundreds of participants at varying levels of adversarial knowledge.

### D. Large-Scale Scalability Analysis

To address computational limitations of enterprise-scale federated learning, we developed a synthetic simulation framework enabling rigorous analysis of topology-based privacy attacks up to 500 node networks. This approach overcomes prohibitive costs while maintaining scientific validity through calibration against empirical data from smaller-scale experiments.

The simulation framework models network topology properties across star, ring, complete, and line configurations with realistic communication patterns. Parameter update synthesis generates statistically consistent gradient norm sequences calibrated against observed distributions from our empirical experiments, incorporating magnitude scaling and temporal decay patterns. Privacy mechanisms use identical Gaussian noise implementations with calibration following (9) and scaling factors derived from observed perturbation levels.

The synthetic generator undergoes statistical validation including parameter distribution alignment via Kolmogorov-Smirnov tests [35] and attack success correlation with empirical results. The analysis evaluates six network sizes between 50-500 nodes across four topologies, three attacks, and three differential privacy settings. Statistical rigor is maintained through multiple runs with 95% confidence intervals and Cohen's d effect size analysis.

### E. Attack Evaluation Metrics and Statistical Validation

We employ conservative evaluation metrics with a 30% threshold for attack success to minimize false positives. This threshold aligns with our advantage-over-random-guessing formulation in (1) and ensures that identified vulnerabilities reflect systematic exploitation of structural relationships rather than statistical artifacts. The threshold provides meaningful signal for adversaries seeking to infer data distribution patterns from network topology while filtering noise-level correlations that lack practical significance for federated learning deployments.

Primary metrics include cluster coherence ratio for communication pattern attacks, multi-metric separability score combining silhouette score and normalized range for parameter magnitude attacks, and maximum absolute correlation coefficient for topology structure attacks. Statistical validation incorporates multiple random seeds for data partitioning, cross-validation across different network sizes, and comparison against random baseline attacks.

### F. Reproducibility and Validation

We ensure reproducible results through deterministic experimental design with fixed random seeds for all data partitioning, identical model architectures, and consistent hyperparameters. Environmental control includes standardized AWS configurations and automated experiment orchestration via Ray. Complete source code with documentation, experimental configuration files, and comprehensive logging are provided [29] to enable verification and extension of our results.

## V. EXPERIMENTAL RESULTS

This section presents comprehensive experimental results evaluating topology-based privacy attacks across 2,620 attack instances comprising 520 complete knowledge configurations and 2,100 attack instances across 420 partial knowledge configurations spanning realistic adversarial constraints. The results reveal systematic vulnerabilities that persist across practical deployment conditions.

### A. Baseline Attack Effectiveness Under Complete Knowledge

Table II presents aggregate results across all three attack vectors under complete topology knowledge. These results provide the theoretical upper bounds of topology-based privacy leakage and serve as the foundation for comparing attack robustness under realistic knowledge constraints. Communication pattern attacks demonstrate the highest effectiveness at 84.1%, exploiting observable message exchange frequencies to cluster nodes based on data similarity. Parameter magnitude attacks achieve moderate effectiveness at 65.0%, profiling statistical properties of update magnitudes to infer participant characteristics. Topology structure attacks show the highest variability with 47.2% success rates, succeeding when organizational constraints create systematic correlations between network position and data distribution.
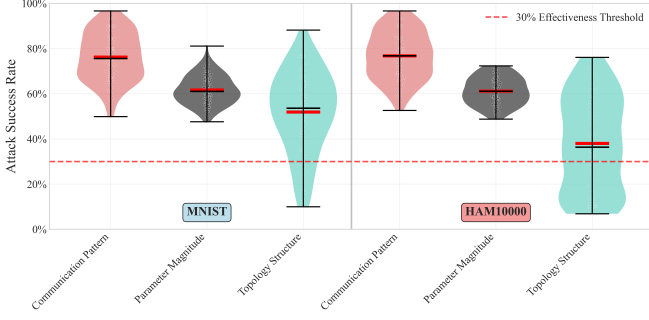
Fig. 3. Attack success rate distributions for MNIST and HAM10000 datasets across all experimental conditions.
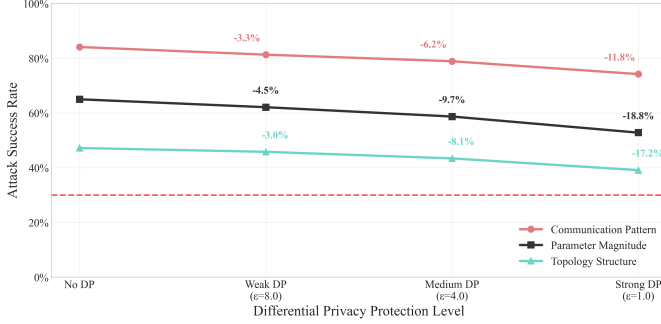


Fig. 4. Attack effectiveness under differential privacy guarantees.

TABLE II
BASELINE ATTACK EFFECTIVENESS UNDER COMPLETE TOPOLOGY
KNOWLEDGE

| Attack Vector | Success | 95% CI | n |
|---|---|---|---|
| Communication Pattern | 84.1% | [83.4, 84.8] | 520 |
| Parameter Magnitude | 65.0% | [64.6, 65.3] | 520 |
| Topology Structure | 47.2% | [45.2, 49.1] | 520 |

Cross-dataset analysis reveals comprehensive vulnerability patterns that transcend data modality and complexity. Figure 3 demonstrates equivalent attack effectiveness between MNIST digit classification and HAM10000 medical imaging across all three attack vectors. Each violin plot combines baseline data with subsampling results, showing the complete distribution of attack success rates from optimal conditions through strong subsampling constraints. The distributions illustrate both domain-agnostic vulnerability and the range of effectiveness under realistic deployment constraints.

### B. Privacy Protection Mechanism Evaluation

Our evaluation incorporates state-of-the-art differential privacy protection with per-round privacy parameters $\varepsilon \in \{1.0, 4.0, 8.0\}$ representing strong, medium, and weak privacy protection respectively. Figure 4 illustrates how attack effectiveness degrades under varying privacy protection levels.

Differential privacy provides measurable but limited protection against topology-based attacks. Even under strong privacy guarantees, communication pattern attacks maintain 74.2% effectiveness, while parameter magnitude attacks achieve 52.8% success rates. The graph demonstrates a gradual decline in

attack effectiveness as privacy protection increases (i.e., as $\varepsilon$ decreases), with the steepest reduction observed for parameter magnitude attacks. However, the maximum reduction across all attack vectors reaches only 18.8%, indicating that current privacy accounting frameworks provide insufficient protection against structural inference attacks.

### C. Attack Robustness Under Realistic Adversarial Knowledge

We systematically evaluate attack effectiveness across five realistic partial knowledge scenarios representing practical adversarial capabilities. These scenarios encompass local adversaries with neighborhood visibility, external adversaries with statistical topology knowledge, and insider threats with organizational structure awareness. Table III presents comprehensive results across 2,100 attack evaluations.

The evaluation reveals remarkable attack robustness across realistic adversarial constraints. Four of five knowledge scenarios (80%) maintain full attack effectiveness, with all three attack vectors remaining above the 30% success threshold. Most significantly, certain partial knowledge scenarios achieve superior performance to complete knowledge baselines.

Local adversaries with neighborhood visibility demonstrate consistent attack effectiveness across all vectors. One-hop knowledge scenarios achieve 68.8%, 47.2%, and 47.8% success rates for communication pattern, parameter magnitude, and topology structure attacks respectively. Expanding to two-hop knowledge improves performance to 76.5%, 62.3%, and 47.9%, approaching complete knowledge effectiveness for communication and parameter attacks.

External adversaries with statistical topology knowledge demonstrate mixed results. Communication pattern and parameter magnitude attacks maintain effectiveness at 86.0% and 65.4% respectively. Topology structure attacks experience significant degradation to 27.6% (41.5% reduction), falling below the effectiveness threshold.

Organizational knowledge scenarios produce the most dramatic results. Coarse organizational grouping (3-groups) enables topology structure attacks to achieve 74.1% effectiveness, representing a remarkable 57.0% improvement over the complete knowledge baseline. Fine-grained organizational structure (5-groups) maintains this advantage with 53.6% effectiveness, representing a 13.7% improvement. Conversely, communication pattern and parameter magnitude attacks experience substantial degradation under organizational constraints, with reductions ranging from 5.4% to 62.3%.

### D. Deployment Scenarios with Subsampling Effects

We evaluate attack robustness under realistic deployment constraints incorporating client and data subsampling combined with differential privacy protection across 288 configurations. Figure 5 presents results across moderate (50% clients, 80% data), strong (30% clients, 60% data), and very strong (20% clients, 50% data) subsampling scenarios, each evaluated under varying differential privacy levels.

Subsampling combined with differential privacy provides meaningful but insufficient privacy amplification against topology-based attacks. Even under very strong subsampling

TABLE III
ATTACK EFFECTIVENESS UNDER REALISTIC PARTIAL TOPOLOGY KNOWLEDGE SCENARIOS

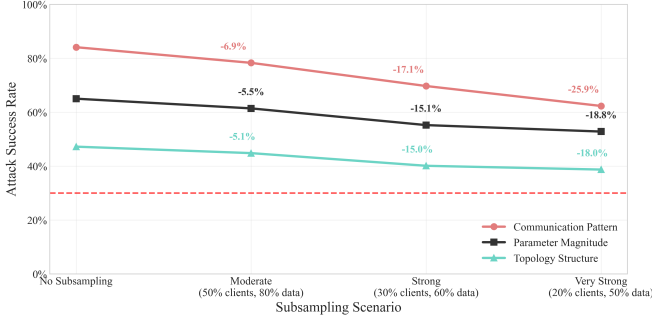| Knowledge Scenario | Comm Pattern | Param Magnitude | Topo Structure | Effective Attacks | Status |
|---|---|---|---|---|---|
| Complete Knowledge | 84.1% | 65.0% | 47.2% | 3/3 (100%) | **Fully Effective** |
| Neighborhood 1-hop | 68.8% | 47.2% | 47.8% | 3/3 (100%) | **Fully Effective** |
| Neighborhood 2-hop | 76.5% | 62.3% | 47.9% | 3/3 (100%) | **Fully Effective** |
| Statistical Knowledge | 86.0% | 65.4% | 27.6% | 2/3 (67%) | **Partially Effective** |
| Organizational 3-groups | 31.7% | 42.5% | 74.1% | 3/3 (100%) | **Fully Effective** |
| Organizational 5-groups | 53.3% | 61.4% | 53.6% | 3/3 (100%) | **Fully Effective** |
| Overall Robustness | | | | 4/5 (80%) | **High Robustness** |



Fig. 5. Attack effectiveness under combined client and data subsampling with differential privacy constraints.

TABLE IV
EFFECT SIZE ANALYSIS FOR KEY EXPERIMENTAL COMPARISONS

| Comparison | Cohen's d | Magnitude |
|---|---|---|
| Complete vs. 1-hop | 0.84 | Large |
| Complete vs. 2-hop | 0.43 | Medium |
| Complete vs. Statistical | 0.12 | Small |
| Complete vs. Org (3-groups) | 1.23 | Large |
| Complete vs. Org (5-groups) | 0.67 | Medium |
| No DP vs. Strong DP | 0.58 | Medium |
| Small vs. Large Networks | 0.09 | Negligible |

constraints, communication pattern attacks maintain 62.3% effectiveness while parameter magnitude attacks achieve 52.8% success rates. The degradation pattern exhibits non-monotonic characteristics, suggesting that aggressive subsampling paradoxically concentrates attacks on the most informative participants, partially offsetting the privacy benefits of differential privacy noise injection.

### E. Enterprise-Scale Analysis

Our synthetic simulation framework enables evaluation of networks spanning 50-500 nodes while maintaining computational tractability. Figure 6a presents attack effectiveness trends across enterprise scales, comparing empirical experiment data (5-30 nodes) with synthetic simulations (50-500 nodes) for each network topology. Figure 6b demonstrates consistent attack patterns under differential privacy protection across varying network sizes.

The analysis reveals consistent attack effectiveness across all evaluated network sizes. Communication pattern attacks maintain 68.7% average effectiveness, parameter magnitude attacks achieve 55.9% success rates, and topology structure attacks demonstrate 48.3% effectiveness across the 50-500 node range. Signal strength metrics remain robust (0.68-0.99) across all tested scales, confirming that fundamental information channels persist regardless of deployment size.

The consistency between real-world experiments and large-scale simulations validates our synthetic modeling approach. While attack effectiveness shows modest degradation with increasing network size, success rates remain substantially above security thresholds across all enterprise scales, demonstrating that organizational scale provides insufficient privacy protection through dilution effects.

### F. Statistical Significance and Effect Size Analysis

All reported results include comprehensive statistical validation with 95% confidence intervals and effect size analysis using Cohen's d for practical significance assessment. Table IV presents effect size magnitudes for key experimental comparisons.
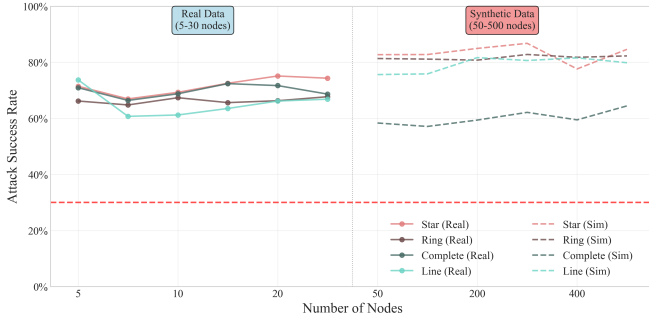
Large effect sizes for organizational knowledge scenarios confirm the practical significance of institutional information advantages, while negligible network size effects validate scale independence findings. Medium effect sizes for differential privacy protection indicate meaningful but insufficient defense capabilities. The statistical analysis confirms that observed differences represent practically significant phenomena rather than statistical artifacts.
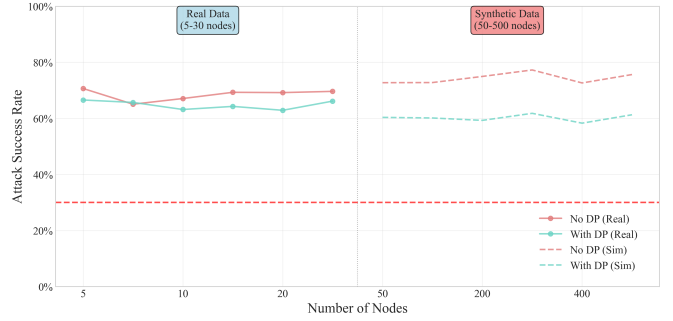
### VI. ANALYSIS AND DISCUSSION

The experimental results reveal fundamental privacy vulnerabilities in federated learning systems that persist despite state-of-the-art defenses across realistic adversarial knowledge constraints. This analysis examines the underlying mechanisms enabling topology-based privacy leakage, evaluates structural limitations of current protection frameworks, and discusses implications for secure federated learning deployment.

### A. Information-Theoretic Foundations of Topology-Based Leakage

Topology-based vulnerabilities represent a fundamental departure from traditional federated learning threat models, exploiting information channels that exist at the architectural level rather than the content level. These vulnerabilities arise from the inherent tension between collaborative learning requirements and privacy preservation, where coordination necessities create observable behavioral signatures that persist regardless of content protection mechanisms.

Fig. 6. Attack effectiveness scaling across enterprise-scale networks. Solid lines represent experimental data (5-30 nodes), dashed lines show synthetic simulations (50-500).

Unlike content-based attacks targeting specific gradient or parameter values, topology-based inference exploits emergent patterns arising from the intersection of network structure, organizational constraints, and learning dynamics. These patterns encode distributional information through second-order effects: correlations between network position and convergence behavior, relationships between organizational structure and communication patterns, and systematic differences in optimization trajectories across participants.

The fundamental challenge stems from the non-decomposable nature of topological information. While differential privacy can bound individual contributions additively, network effects create multiplicative information leakage where the whole becomes greater than the sum of its parts. A participant's position in the network graph encodes exponentially more information than their individual data contribution, as it reveals relationships with all other participants simultaneously. This creates an information amplification effect where modest local data characteristics become magnified through network structure, enabling inference attacks that would be impossible in isolation.

Furthermore, topology-based attacks exploit temporal persistence that content-based protections cannot address. While gradient perturbation mechanisms can mask individual updates, the underlying network structure remains static across training rounds, enabling adversaries to accumulate evidence over time. This temporal dimension transforms ephemeral communication patterns into persistent fingerprints that strengthen with each training iteration.

### B. Structural Limitations of Current Privacy Protection

Current privacy frameworks exhibit fundamental conceptual limitations against topology-based attacks at three distinct levels: granularity mismatches, independence assumptions, and protection scope boundaries.

The mathematical foundations of current privacy mechanisms assume participant exchangeability—that any participant could be substituted for another without affecting privacy guarantees. However, network topology fundamentally violates this assumption by creating asymmetric roles and positional advantages. Central nodes in hierarchical configurations aggregate information from multiple sources, creating inherent

privacy asymmetries that current frameworks cannot model or protect against.

Moreover, existing privacy composition rules fail to account for structural correlations. When calculating cumulative privacy loss, current mechanisms assume independent contributions across participants. However, topology creates systematic dependencies where one participant's privacy loss can cascade to their network neighbors through positional inference. This correlation amplification means that the actual privacy loss exceeds theoretical bounds derived from independence assumptions.

The temporal dimension adds another layer of complexity. Privacy accounting typically models one-shot interactions, but federated learning involves repeated interactions over extended periods. Network topology creates persistent attack surfaces that accumulate information across training rounds, enabling adversaries to build increasingly accurate models of participant characteristics through longitudinal analysis. Current privacy frameworks lack mechanisms to bound this temporal accumulation of structural information.

These limitations compound in realistic deployments where organizational constraints dictate network structure. Healthcare federations, academic collaborations, and financial partnerships inherently expose institutional relationships through regulatory requirements, operational agreements, and public disclosures that current privacy frameworks cannot address, creating fundamental architectural vulnerabilities that persist regardless of content protection strength.

### C. Adversarial Knowledge Landscape and Realistic Threat Assessment

Our comprehensive evaluation reveals a complex adversarial landscape challenging traditional assumptions about required knowledge for successful privacy attacks. The heterogeneity in adversarial capabilities demonstrates that topology-based vulnerabilities remain exploitable across a broad spectrum of realistic knowledge constraints.

Local adversaries with neighborhood visibility achieve consistent attack effectiveness despite limited structural knowledge, indicating that topology-based vulnerabilities exhibit local clustering properties. This finding suggests that compromised nodes or insider threats need not possess global network

knowledge to conduct effective inference attacks, significantly lowering barriers for practical exploitation.

External adversaries with statistical topology knowledge demonstrate that aggregate structural properties suffice for certain attack vectors while constraining others. The maintenance of high effectiveness for communication pattern attacks under statistical knowledge scenarios reveals that detailed topological information provides diminishing returns for certain inference techniques.

Most critically, organizational knowledge scenarios expose concerning vulnerability patterns where institutional relationship information enables capabilities exceeding complete topology knowledge for specific tasks. This counterintuitive finding reflects the fundamental difference between syntactic and semantic knowledge. Complete topology knowledge provides syntactic information about network connections, while organizational knowledge provides semantic understanding of why those connections exist and what they represent. Adversaries with organizational context can leverage domain expertise to identify which structural patterns are meaningful for their inference goals, effectively filtering noise from signal in ways that pure topological analysis cannot achieve.

### D. Scalability and Universal Vulnerability Characteristics

Attack effectiveness across network scales challenges fundamental assumptions about privacy protection in large-scale distributed systems. Traditional security models assume that increasing participant numbers provide inherent privacy protection through dilution effects, where individual contributions become less significant as system size grows.

Our synthetic simulation framework demonstrates that topology-based vulnerabilities maintain effectiveness well above security thresholds across enterprise scales despite modest degradation with network size. Enterprise-scale deployments may actually amplify these vulnerabilities through increased organizational complexity, where large-scale federations exhibit more pronounced hierarchical structures, departmental clustering, and geographic constraints that strengthen correlations between network position and data characteristics.

The domain-agnostic nature of these vulnerabilities is demonstrated through equivalent attack patterns across medical imaging and digit classification datasets. This universality occurs because topology-based attacks exploit fundamental properties of distributed learning systems—structural and behavioral patterns inherent to federated coordination—rather than data content characteristics. The persistence of vulnerabilities under strong subsampling constraints combined with differential privacy guarantees reveals that current privacy amplification techniques provide insufficient protection against structural inference.

These findings reveal that topology-based vulnerabilities represent a fundamental limitation of collaborative learning rather than implementation-specific weaknesses. Addressing these vulnerabilities requires architectural innovations that fundamentally reconsider the relationship between coordination requirements and privacy preservation in distributed machine learning systems.

## VII. Toward Topology-Aware Defense Mechanisms

The systematic identification of topology-based vulnerabilities necessitates fundamental advances in privacy protection that address structural information leakage. We propose structural noise injection as a complementary defense mechanism and provide comprehensive empirical validation demonstrating significant effectiveness when properly layered with existing privacy techniques. Our evaluation across 808 experimental configurations reveals that privacy leakage can be substantially mitigated through defense-in-depth approaches combining structural noise with differential privacy and subsampling mechanisms.

### A. Structural Noise Injection: Theoretical Foundation

Structural noise injection operates by adding calibrated perturbations to the three identified information leakage channels: communication patterns, parameter timing, and parameter magnitudes. The defense targets observable metadata that enables topology-based inference while attempting to preserve federated learning functionality, though each mechanism introduces distinct utility trade-offs.

*1) Communication Pattern Protection:* Communication pattern attacks exploit the frequency and timing of message exchanges between participants. To disrupt these patterns, we inject dummy communications following a Poisson process with rate $\lambda_{\text{dummy}} = \alpha \cdot \lambda_{\text{real}}$, where $\lambda_{\text{real}}$ represents the true communication rate and $\alpha \in [0.1, 0.3]$ provides the noise injection rate.

The dummy messages maintain identical metadata characteristics (size, routing, encryption) as legitimate communications but carry no meaningful content. This approach increases the entropy of observable communication patterns:

$$H_{\text{observed}} = H_{\text{real}} + H_{\text{dummy}} - I(\text{real}, \text{dummy}) \qquad (10)$$

where $I(\text{real}, \text{dummy}) \approx 0$ due to the independence of dummy message generation.

This mechanism introduces communication overhead proportional to $\alpha$, increasing network traffic by 10-30% while potentially disrupting adaptive communication protocols that rely on organic message timing patterns for convergence optimization.

*2) Parameter Timing Protection:* Timing-based inference exploits correlations between parameter update timing and convergence behavior. We add Gaussian noise to communication timestamps: $t_{\text{observed}} = t_{\text{real}} + \epsilon_t$ where $\epsilon_t \sim \mathcal{N}(0, \sigma_t^2)$ and $\sigma_t = \beta \cdot \text{Var}(t_{\text{real}})$ with $\beta \in [0.05, 0.3]$.

While this preserves approximate ordering for basic convergence detection, excessive timing perturbation can disrupt synchronization mechanisms and adaptive learning rate schedules that depend on precise temporal coordination, potentially degrading convergence efficiency in time-sensitive federated learning protocols.

*3) Parameter Magnitude Protection:* Parameter magnitude attacks exploit systematic differences in update magnitudes across participants. We apply multiplicative noise to parameter norms: $||\theta_{\text{observed}}||_2 = ||\theta_{\text{real}}||_2 \cdot (1 + \epsilon_m)$ where $\epsilon_m \sim \mathcal{N}(0, \sigma_m^2)$ and $\sigma_m \in [0.05, 0.3]$.
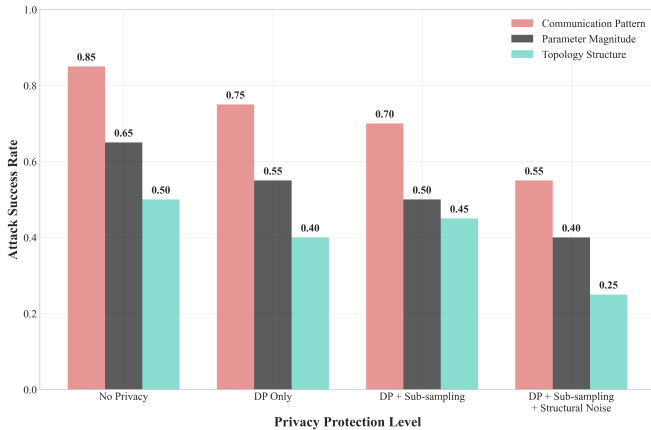
Fig. 7. Progressive attack success reduction through layered privacy protection mechanisms showing the complementary effectiveness of structural noise injection.

TABLE V
DEFENSE CONFIGURATION EFFECTIVENESS AGAINST TOPOLOGY-BASED ATTACKS

| Defense Configuration | Max Reduction |
|---|---|
| *Additional Reduction Beyond Baseline* | |
| DP + Structural Noise | 12.5% |
| Subsampling + DP + Structural Noise | 40.9% |
| *Best Combined Defense* | |
| Communication Pattern | **34.0%** |
| Parameter Magnitude | **14.1%** |
| Topology Structure | **51.4%** |

DP = Differential Privacy; percentages show additional reduction beyond baseline attack success rates.

This approach is functionally similar to differential privacy noise injection but applied specifically to parameter norms rather than raw gradients. The multiplicative formulation preserves relative magnitude relationships better than additive noise but still introduces systematic bias in aggregation mechanisms that rely on precise parameter weighting.

### B. Comprehensive Empirical Validation

Our evaluation demonstrates that structural noise injection provides substantial additional protection when properly integrated with existing privacy mechanisms.

Figure 7 illustrates the progressive enhancement of privacy protection through layered defense mechanisms. While differential privacy and subsampling provide foundational protection, structural noise injection delivers significant additional attack reduction across all evaluated scenarios.

Table V presents comprehensive effectiveness measurements for layered defense configurations, demonstrating substantial improvements over individual protection mechanisms. Critically, structural noise injection operates orthogonally to existing privacy mechanisms, enabling multiplicative rather than additive privacy benefits when combined with differential privacy and subsampling.

### C. Deployment Configurations and Practical Implementation

Based on comprehensive experimental validation, we provide evidence-based deployment recommendations for different organizational contexts.

High-security deployments requiring maximum privacy protection should implement strong DP ($\varepsilon = 1.0$), moderate subsampling (50% clients, 80% data), and strong structural noise injection. This configuration achieves 30-37% additional attack reduction with approximately 15% communication overhead, regardless of the underlying network topology imposed by organizational constraints. Balanced production environments benefit from medium DP ($\varepsilon = 4.0$), optional moderate subsampling, and medium structural noise, providing 15-25% additional attack reduction with 10% communication overhead. Resource-constrained scenarios can achieve meaningful protection through strong structural noise with weak DP ($\varepsilon = 8.0$), providing 10-20% additional attack reduction while maintaining minimal architectural requirements. The effectiveness of structural noise injection remains consistent across diverse network configurations, making it suitable for deployment in FL infrastructures without requiring architectural modifications to accommodate specific topology designs.

### D. Future Directions and Comprehensive Protection

While structural noise injection provides substantial improvement over undefended baselines, achieving comprehensive topology-based privacy protection requires continued research in several key areas. Enhanced temporal correlation disruption mechanisms could further improve parameter magnitude attack mitigation, while privacy-preserving network coordination protocols may enable stronger protection against sophisticated adversaries with organizational knowledge.

The integration of cryptographic techniques, where computationally feasible, represents another avenue for enhancing structural privacy protection. Additionally, adaptive defense mechanisms that dynamically adjust protection parameters based on real-time threat assessment could optimize the privacy-utility trade-off for diverse deployment scenarios.

Most importantly, our evaluation validates that topology-based vulnerabilities can be practically mitigated through systematic application of layered defense mechanisms. This empirical validation of structural noise injection as a complementary defense mechanism provides the federated learning community with immediately deployable protection strategies while establishing the foundation for continued advancement in topology-aware privacy protection research.

## VIII. CONCLUSIONS

Our comprehensive analysis reveals that topology-based vulnerabilities represent fundamental architectural limitations in federated learning systems. We show that structural information enables powerful data distribution inference attacks that persist despite state-of-the-art privacy protections and across enterprise scales and diverse data modalities, demonstrating universal vulnerability characteristics that transcend domain-specific protections.

Our empirical validation demonstrates that structural noise injection provides substantial complementary protection when properly layered with existing privacy mechanisms. The orthogonal nature of structural defenses enables multiplicative rather than additive privacy benefits, validating defense-in-depth approaches for topology-aware privacy protection.

These findings demonstrate that fundamental advances in topology-aware privacy mechanisms are needed to address critical gaps in current protection paradigms. Privacy amplification techniques must account for topology-induced correlations between participants rather than treating them as isolated entities, while communication protocols must minimize observable patterns while maintaining convergence guarantees.

Only through comprehensive topology-aware privacy mechanisms addressing both content and structural leakage can federated learning deliver secure collaborative machine learning while maintaining the architectural flexibility that drives its adoption across diverse organizational contexts.

## REFERENCES

[1] B. McMahan, E. Moore *et al.*, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, vol. 54. PMLR, 20–22 Apr 2017, pp. 1273–1282.

[2] T. Li, A. K. Sahu *et al.*, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.

[3] N. Rieke, J. Hancox *et al.*, "The future of digital health with federated learning," *npj Digital Medicine*, vol. 3, no. 1, p. 119, Sep. 2020.

[4] Q. Yang, Y. Liu *et al.*, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, Jan. 2019.

[5] W. Y. B. Lim, N. C. Luong *et al.*, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, 2020.

[6] L. Zhu, Z. Liu, and S. Han, "Deep Leakage from Gradients," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d. Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019.

[7] J. Geiping, H. Bauermeister *et al.*, "Inverting gradients - how easy is it to break privacy in federated learning?" in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 16 937–16 947.

[8] L. Melis, C. Song *et al.*, "Exploiting unintended feature leakage in collaborative learning," in *2019 IEEE Symposium on Security and Privacy (SP)*, 2019, pp. 691–706.

[9] C. Zhou, Y. Gao *et al.*, "Ppa: Preference profiling attack against federated learning," in *Proceedings of the 2023 Network and Distributed System Security Symposium*, ser. 30th Annual Network and Distributed System Security Symposium, NDSS 2023. Internet Society, Feb. 2023.

[10] A. Koloskova, N. Loizou *et al.*, "A unified theory of decentralized SGD with changing topology and local updates," in *Proceedings of the 37th International Conference on Machine Learning*, vol. 119. PMLR, 13–18 Jul 2020, pp. 5381–5393.

[11] J. Wang, Z. Charles *et al.*, "A field guide to federated optimization," *arXiv preprint arXiv:2107.06917*, 2021.

[12] V. Mothukuri, R. M. Parizi *et al.*, "A survey on security and privacy of federated learning," *Future Generation Computer Systems*, vol. 115, pp. 619–640, 2021.

[13] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.

[14] K. Bonawitz, V. Ivanov *et al.*, "Practical secure aggregation for privacy-preserving machine learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '17. New York, NY, USA: ACM Press, 2017, p. 1175–1191.

[15] A. Acar, H. Aksu *et al.*, "A survey on homomorphic encryption schemes: Theory and implementation," *ACM Comput. Surv.*, vol. 51, no. 4, Jul. 2018.

[16] L. Fowl, J. Geiping *et al.*, "Robbing the fed: Directly obtaining private data in federated learning with modified models," *arXiv preprint arXiv:2110.13057*, 2022.

[17] M. Nasr, R. Shokri, and A. Houmansadr, "Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning," in *2019 IEEE Symposium on Security and Privacy (SP)*, 2019, pp. 739–753.

[18] H. B. McMahan, D. Ramage *et al.*, "Learning differentially private recurrent language models," in *Proceedings of the 6th International Conference on Learning Representations*, 2018.

[19] X. Li, K. Huang *et al.*, "On the convergence of fedavg on non-iid data," in *Proceedings of the 8th International Conference on Learning Representations*, 2020.

[20] C. Feng, Y. Gao *et al.*, "From models to network topologies: A topology inference attack in decentralized federated learning," *arXiv preprint arXiv:2501.03119*, 2025.

[21] B. Zhao, K. R. Mopuri, and H. Bilen, "idlg: Improved deep leakage from gradients," *arXiv preprint arXiv:2001.02610*, 2020.

[22] S. Truex, N. Baracaldo *et al.*, "A hybrid approach to privacy-preserving federated learning," in *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, ser. AISec'19. New York, NY, USA: ACM Press, 2019, p. 1–11.

[23] E. Bagdasaryan, A. Veit *et al.*, "How to backdoor federated learning," in *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, vol. 108. PMLR, 26–28 Aug 2020, pp. 2938–2948.

[24] M. Noble, A. Bellet, and A. Dieuleveut, "Differentially private federated learning on heterogeneous data," in *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, vol. 151. PMLR, 28–30 Mar 2022, pp. 10 110–10 145.

[25] A. G. Roy, S. Siddiqui *et al.*, "Braintorrent: A peer-to-peer environment for decentralized federated learning," *arXiv preprint arXiv:1905.06731*, 2019.

[26] J. Wu, F. Dong *et al.*, "Topology-aware federated learning in edge computing: A comprehensive survey," *ACM Comput. Surv.*, vol. 56, no. 10, Jun. 2024.

[27] A. Bellet, R. Guerraoui *et al.*, "Personalized and private peer-to-peer machine learning," in *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, vol. 84. PMLR, 09–11 Apr 2018, pp. 473–481.

[28] E. Hallaji, R. Razavi-Far *et al.*, "Decentralized federated learning: A survey on security and privacy," *IEEE Transactions on Big Data*, vol. 10, no. 2, pp. 194–213, 2024.

[29] M. Rangwala, "Murmura: A Ray-based Framework for Federated and Decentralized Machine Learning," Software, Jun. 2025, https://doi.org/10.5281/zenodo.15622123.

[30] P. Moritz, R. Nishihara *et al.*, "Ray: A distributed framework for emerging AI applications," in *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*. Carlsbad, CA: USENIX Association, Oct. 2018, pp. 561–577.

[31] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database," *ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist*, vol. 2, 2010.

[32] P. Tschandl, C. Rosendahl, and H. Kittler, "The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions," *Scientific data*, vol. 5, no. 1, pp. 1–9, 2018.

[33] I. Mironov, "Rényi differential privacy," in *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*, 2017, pp. 263–275.

[34] A. Yousefpour, I. Shilov *et al.*, "Opacus: User-friendly differential privacy library in pytorch," *arXiv preprint arXiv:2109.12298*, 2021.

[35] F. J. M. Jr., "The kolmogorov-smirnov test for goodness of fit," *Journal of the American Statistical Association*, vol. 46, no. 253, pp. 68–78, 1951.