# SketchGuard: Scaling Byzantine-Robust Decentralized Federated Learning via Sketch-Based Screening

## Murtaza Rangwala
School of Computing and Information Systems
The University of Melbourne
Melbourne, Australia
rangwalam@unimelb.edu.au

## Farag Azzedin
Department of Information and Computer Science
King Fahd University of Petroleum and Minerals
Dhahran, Saudi Arabia
fazzedin@kfupm.edu.sa

## Richard O. Sinnott
School of Computing and Information Systems
The University of Melbourne
Melbourne, Australia
rsinnot@unimelb.edu.au

## Rajkumar Buyya
School of Computing and Information Systems
The University of Melbourne
Melbourne, Australia
rbuyya@unimelb.edu.au

## Abstract

Decentralized Federated Learning (DFL) enables privacy-preserving collaborative training without centralized servers, but remains vulnerable to Byzantine attacks where malicious clients submit corrupted model updates. Existing Byzantine-robust DFL defenses rely on similarity-based neighbor screening that requires every client to exchange and compare complete high-dimensional model vectors with all neighbors in each training round, creating prohibitive communication and computational costs that prevent deployment at web scale. We propose SketchGuard, a general framework that decouples Byzantine filtering from model aggregation through sketch-based neighbor screening. SketchGuard compresses $d$-dimensional models to $k$-dimensional sketches ($k \ll d$) using Count Sketch for similarity comparisons, then selectively fetches full models only from accepted neighbors, reducing per-round communication complexity from $O(d|\mathcal{N}_i|)$ to $O(k|\mathcal{N}_i| + d|\mathcal{S}_i|)$, where $|\mathcal{N}_i|$ is the neighbor count and $|\mathcal{S}_i| \leq |\mathcal{N}_i|$ is the accepted neighbor count. We establish rigorous convergence guarantees in both strongly convex and non-convex settings, proving that Count Sketch compression preserves Byzantine resilience with controlled degradation bounds where approximation errors introduce only a $(1 + O(\epsilon))$ factor in the effective threshold parameter. Comprehensive experiments across multiple datasets, network topologies, and attack scenarios demonstrate that SketchGuard maintains identical robustness to state-of-the-art methods, with mean test error rate deviation of only up to 0.35 percentage points, while reducing computation time by up to 82% and communication overhead by 50-70% depending on filtering effectiveness, with benefits scaling multiplicatively with model dimensionality and network connectivity. These results establish the viability of sketch-based compression as a fundamental enabler of robust DFL at web scale.

## CCS Concepts

• **Security and privacy** → **Systems security**.

## Keywords

Decentralized Learning, Byzantine Robustness, Dimensionality Reduction, Scalable Machine Learning

## 1 Introduction

Federated Learning (FL) enables collaborative training of AI models over distributed data while preserving privacy by keeping raw data local [23]. However, the canonical, server-assisted architecture of FL centralizes aggregation of model parameters, creating a single point of failure, a communication bottleneck, and trust issues [13]. These drawbacks have catalyzed Decentralized Federated Learning (DFL), where clients exchange model updates in a peer-to-peer manner over a dynamic communication graph, thereby improving scalability and resilience [3].

A central challenge in DFL is *byzantine robustness*: the ability to withstand malicious clients that submit arbitrary or carefully crafted updates to poison training, induce consensus drift, or trigger targeted failures [2, 5]. Unlike centralized FL where robust rules such as Krum [5], coordinatewise Median [33], or Trimmed-Mean [31] are applied once at a centralized server, DFL requires every client to aggregate its neighbors' updates under local, graph-limited views, often with non-IID data and time-varying connectivity. To address this added complexity, most DFL defenses adopt *local-consistency filters*, where clients accept a neighbor's update only if it is sufficiently similar to their own state, and then average over the accepted subset [8, 11, 13, 17, 19]. These mechanisms provide convergence and robustness guarantees in both strongly convex and non-convex model training settings, yet suffer from a fundamental scalability bottleneck: clients must exchange and compare complete, high-dimensional model vectors with all neighbors in every round. For emerging web-scale applications like decentralized training of frontier models with billions of parameters across thousands of distributed participants [7, 22], this creates prohibitive communication and computation costs that prevent practical implementation of these systems at scale.

Sketch-based compression offers tools for communication efficient learning. Count Sketch, for instance, compresses a vector of $d$ dimensions into a summary of $k$ dimensions using simple hash and sign functions, where $k \ll d$. Sketches are linear and allow approximate preservation of coordinates, thereby supporting fast similarity estimation with formal guarantees [10, 16, 26]. While sketching has proven effective for bandwidth reduction in FL, extending it

to *Byzantine-robust, fully decentralized* aggregation, where compressed representations must support secure neighbor screening, remains underexplored.

In this paper, we propose SKETCHGUARD, a general framework for Byzantine-robust DFL that decouples filtering from aggregation through sketch-based neighbor screening. Our key insight is that similarity-based Byzantine filtering can operate on compressed representations, while the final aggregation requires full precision models only for accepted neighbors. SKETCHGUARD is applicable to any similarity-based Byzantine defense that relies on Euclidean distance measures (e.g., [13, 19, 25, 30]), but for our theoretical analysis and empirical evaluation, we instantiate it with state-of-the-art BALANCE aggregation [13], which provides the strongest theoretical guarantees among existing methods. Our main contributions can be summarized as follows:

- We provide rigorous analysis showing that Count Sketch compression maintains Byzantine resilience with controlled degradation bounds.
- We establish convergence rates for SKETCHGUARD in both strongly convex and non-convex settings.
- Through comprehensive experiments across multiple datasets, network topologies, and attack scenarios, we demonstrate that our approach achieves identical robustness to state-of-the-art methods while reducing communication overhead by 50-70% and computation time by up to 82%.

## 2 Preliminaries and Related Work

This section provides technical background on DFL protocols, Byzantine attack models, and compression techniques, and discusses related work in each area.

### 2.1 DFL Problem Formulation and Protocol

Consider $n$ clients connected by an undirected graph $G = (V, E)$, where each client $i \in V$ possesses a private dataset $\mathcal{D}_i$ and maintains a local model $\mathbf{w}_i \in \mathbb{R}^d$. The collective objective is to minimize the average empirical loss:

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}), \qquad (1)$$

where $f_i(\mathbf{w}) = \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}_i}[\ell(\mathbf{w}; \mathbf{x}, y)]$ is the expected loss over client $i$'s data distribution.

The DFL protocol alternates between two phases. First, clients perform local updates:

$$\mathbf{w}_i^{t+1/2} = \mathbf{w}_i^t - \eta \nabla f_i(\mathbf{w}_i^t). \qquad (2)$$

Then, each client aggregates neighbor models according to:

$$\mathbf{w}_i^{t+1} = \alpha \mathbf{w}_i^{t+1/2} + (1 - \alpha) \cdot \text{AGG}_i \left( \{\mathbf{w}_j^{t+1/2} : j \in \mathcal{N}_i\} \right), \qquad (3)$$

where $\mathcal{N}_i$ denotes client $i$'s neighbors and $\text{AGG}_i$ is a local aggregation function.

### 2.2 Byzantine Attack Models

We consider $f$-Byzantine adversaries that control up to $f$ clients per neighborhood. Byzantine clients can deviate arbitrarily from the protocol, presenting a spectrum of threats ranging from simple disruption to sophisticated coordinated manipulation. These adversaries may employ model poisoning by sending crafted parameters $\hat{\mathbf{w}}_j$ designed to maximize deviation from honest updates while evading detection by robust aggregation rules. Representative attacks include random noise injection [5], data label flipping [2], and gradient-based perturbation optimization [24]. More sophisticated adversaries operate adaptively, observing honest client behavior and adjusting their strategy accordingly. Such attacks [14, 28] solve optimization problems to craft maximally harmful updates that remain within detection thresholds of similarity-based defenses. The most challenging scenarios involve collusive attacks where multiple Byzantine clients coordinate their malicious behavior [24, 32], potentially overwhelming local robust aggregation mechanisms that assume independent adversarial actions. These coordinated strategies can amplify individual malicious impact through strategic cooperation, creating correlated deviations that may appear legitimate to individual clients with limited neighborhood views.

The fundamental challenge for DFL robustness lies in defending against these varied attack strategies under graph-limited visibility, where each client must make aggregation decisions based solely on local neighborhood information rather than global network state. This constraint makes similarity-based filtering particularly attractive, as it provides a principled approach to neighbor screening that remains effective across different attack types while being computationally tractable for individual clients.

### 2.3 Byzantine-Robust DFL Defenses

As explained in the previous section, existing DFL defenses have converged on similarity-based neighbor filtering as the dominant paradigm, with methods differing primarily in their filtering criteria and aggregation mechanisms.

**UBAR** [17] utilizes two-stage neighbor selection based on distance and loss filtering, followed by averaging.

**LEARN** [11] employs multiple rounds of model exchanges per iteration with trimmed-mean aggregation. The multi-round communication exacerbates the pre-existing scalability issues with similarity-based filtering mechanisms.

**SCCLIP** [19] applies self-centered clipping to each received neighbor model, mitigating large-magnitude deviations. It provides convergence guarantees in non-convex settings but also requires full model comparisons.

**BALANCE** [13] introduces adaptive similarity thresholds for neighbor screening with strong theoretical convergence guarantees in both strongly convex and non-convex settings. It achieves state-of-the-art robustness but at a high computational cost.

**WFAgg** [8] proposes multi-filter approaches for dynamic topologies, combining several screening mechanisms while maintaining the same computational complexity limitations.

However, all these methods require computation of similarities between complete $d$-dimensional model vectors with all neighbors, creating bottlenecks that prevent practical web-scale deployments. Table 1 compares SKETCHGUARD with these existing DFL defense approaches.

**Table 1: Comparison of Byzantine-robust decentralized aggregation algorithms. "Convex"/"Non-convex" indicate convergence guarantees; "No $c_i$" means no need to know compromised node ratio; "No Complete" means no complete-graph assumption.**

| Algorithm | Convex | Non-convex | No $c_i$ | No Complete | Scalable |
|---|---|---|---|---|---|
| UBAR [17] | – | – | – | – | Med. |
| LEARN [11] | – | ✓ | – | – | Low |
| SCCLIP [19] | – | ✓ | ✓ | ✓ | Med. |
| BALANCE [13] | ✓ | ✓ | ✓ | ✓ | Med. |
| WFAgg [8] | – | – | ✓ | ✓ | Med. |
| SKETCHGUARD | ✓ | ✓ | ✓ | ✓ | High |

## 2.4 Compression Techniques in FL

Communication and computational efficiency remain fundamental challenges in FL, motivating extensive research in model compression. Quantization techniques reduce parameter precision [1, 4], sparsification methods transmit only significant updates [27, 29], and low-rank approaches exploit model structure [18, 20]. However, most compression techniques target benign settings and often compromise robustness when Byzantine participants are present. Traditional quantization can amplify malicious update impact, while sparsification enables attackers to concentrate influence in transmitted coordinates. More broadly, existing compression methods are incompatible with Byzantine-robust aggregation algorithms, preventing scalable deployment of robust FL.

## 2.5 Count Sketch

Count Sketch [10] provides randomized linear projection with properties that make it particularly suitable for Byzantine-robust neighbor screening in DFL. Given $\mathbf{w} \in \mathbb{R}^d$, a Count Sketch of size $k \ll d$ is constructed using hash function $h : [d] \to [k]$ and sign function $s : [d] \to \{-1, +1\}$:

$$\mathrm{CS}(\mathbf{w})[b] = \sum_{i:h(i)=b} s(i)w_i, \quad b = 1, \ldots, k. \tag{4}$$

Count Sketch possesses three critical properties for robust DFL screening. First, linearity ensures consistent compression across clients, as $\mathrm{CS}(\alpha\mathbf{u} + \beta\mathbf{v}) = \alpha\mathrm{CS}(\mathbf{u}) + \beta\mathrm{CS}(\mathbf{v})$. Second, unbiased coordinate estimation provides formal recovery guarantees [10, 26]. Third, and most crucial for similarity-based Byzantine filtering, Count Sketch preserves approximate Euclidean distances:

LEMMA 1 (DISTANCE PRESERVATION [10]). *For any vectors* $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ *and sketch size* $k = O(\epsilon^{-2}\log(1/\delta))$, *with probability at least* $1 - \delta$:

$$(1-\epsilon)\|\mathbf{u}-\mathbf{v}\|^2 \le \|\mathrm{CS}(\mathbf{u})-\mathrm{CS}(\mathbf{v})\|^2 \le (1+\epsilon)\|\mathbf{u}-\mathbf{v}\|^2 \tag{5}$$

## 3 SketchGuard: Scalable Robust Aggregation

We now present SKETCHGUARD, which leverages Lemma 1 to perform neighbor screening in the sketch domain. The key design rationale is that if Count Sketch approximately preserves distances, then similarity-based filtering decisions made using sketches will closely match those made using full-precision models, allowing us to defer expensive full-model exchanges until after filtering.

### 3.1 Protocol Description

At each DFL training round $t$, client $i$ executes the four-phase SKETCHGUARD protocol detailed in Algorithm 1 and illustrated in Figure 1 as follows:

**Phase 1: Local Training.** Client $i$ performs local stochastic gradient descent as shown in Equation 2.

**Phase 2: Sketch Exchange.** The updated model is compressed using CS: $\mathbf{s}_i^{t+1/2} = \mathrm{CS}(\mathbf{w}_i^{t+1/2})$, and these $k$-dimensional sketches are exchanged with the immediate neighbors $\mathcal{N}_i$.

**Phase 3: Adaptive Filtering.** Neighbor $j$ is accepted if their sketch distance satisfies:

$$\|\mathbf{s}_i^{t+1/2} - \mathbf{s}_j^{t+1/2}\| \le \gamma\exp(-\kappa t/T)\|\mathbf{s}_i^{t+1/2}\| \tag{6}$$

where $\gamma$ controls the base threshold, $\kappa$ the decay rate, and $T$ the total rounds. This adaptive threshold mechanism is adopted from BALANCE [13], where the exponential decay reflects convergence of honest clients over time.

REMARK 1. *When using Count Sketch compression with approximation parameter* $\epsilon$, *the effective threshold parameter becomes* $\gamma_{\mathit{eff}} = \gamma\sqrt{(1+\epsilon)/(1-\epsilon)}$ *to account for distance preservation errors in the sketch domain, as established by the analysis in Lemma 1. This ensures that our theoretical convergence guarantees in Section 4 accurately reflect the impact of compression on filtering decisions.*

**Phase 4: Model Aggregation.** Full models are fetched from accepted neighbors $\mathcal{S}_i^t$. Before aggregation, each received model $\mathbf{w}_j^{t+1/2}$ is verified by recomputing its sketch and comparing with the originally received $\mathbf{s}_j^{t+1/2}$. Any neighbor whose model fails verification is removed from $\mathcal{S}_i^t$. The verified models are then aggregated:

$$\mathbf{w}_i^{t+1} = \alpha\mathbf{w}_i^{t+1/2} + \frac{(1-\alpha)}{|\mathcal{S}_i^t|}\sum_{j\in\mathcal{S}_i^t}\mathbf{w}_j^{t+1/2} \tag{7}$$

where $\alpha \in [0, 1]$ balances self-reliance and collaboration.

---

**Algorithm 1** SKETCHGUARD: Robust Aggregation via Adaptive Sketch-Based Filtering

---

**Require:** Local data $\mathcal{D}_i$, neighbors $\mathcal{N}_i$, parameters $\gamma, \kappa, \alpha$, sketch size $k$
**Ensure:** Updated model $\mathbf{w}_i^{t+1}$
1: $\mathbf{w}_i^{t+1/2} \leftarrow \mathbf{w}_i^t - \eta\mathbf{g}(\mathbf{w}_i^t)$
2: $\mathbf{s}_i^{t+1/2} \leftarrow \mathrm{CS}(\mathbf{w}_i^{t+1/2})$
3: Exchange sketches $\mathbf{s}_i^{t+1/2}$ with neighbors $\mathcal{N}_i$
4: $\tau \leftarrow \gamma\exp(-\kappa t/T)\|\mathbf{s}_i^{t+1/2}\|$
5: $\mathcal{S}_i^t \leftarrow \{j \in \mathcal{N}_i : \|\mathbf{s}_i^{t+1/2} - \mathbf{s}_j^{t+1/2}\| \le \tau\}$
6: **if** $|\mathcal{S}_i^t| = 0$ and $|\mathcal{N}_i| > 0$ **then**
7: $\quad \mathcal{S}_i^t \leftarrow \{\arg\min_{j\in\mathcal{N}_i} \|\mathbf{s}_i^{t+1/2} - \mathbf{s}_j^{t+1/2}\|\}$
8: **end if**
9: Fetch models $\{\mathbf{w}_j^{t+1/2}\}_{j\in\mathcal{S}_i^t}$ from accepted neighbors
10: Verify $\mathrm{CS}(\mathbf{w}_j^{t+1/2}) = \mathbf{s}_j^{t+1/2}$ for each $j \in \mathcal{S}_i^t$; remove if mismatch
11: $\mathbf{w}_i^{t+1} \leftarrow \alpha\mathbf{w}_i^{t+1/2} + \frac{1-\alpha}{|\mathcal{S}_i^t|}\sum_{j\in\mathcal{S}_i^t}\mathbf{w}_j^{t+1/2}$
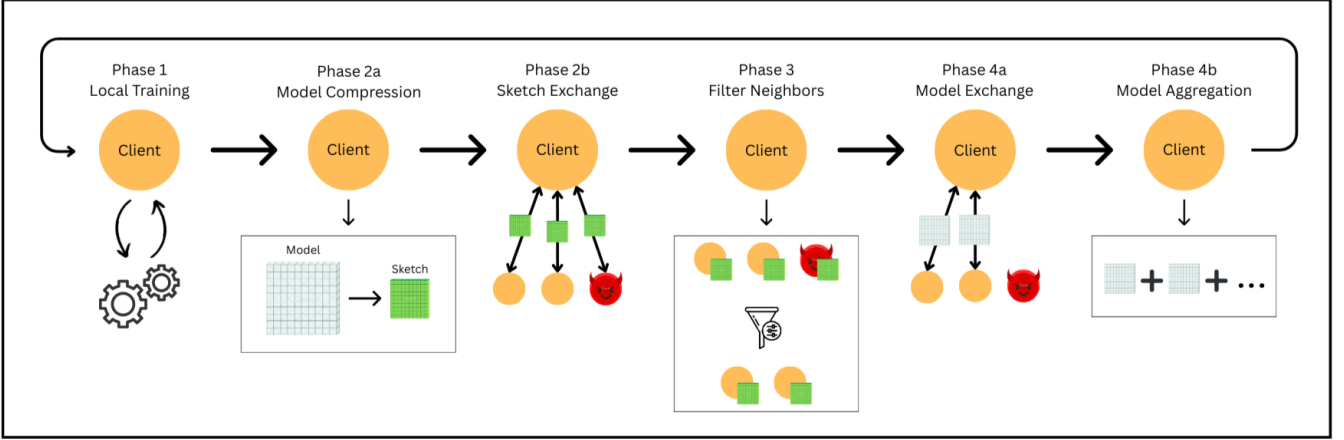12: **return** $\mathbf{w}_i^{t+1}$

---

**Figure 1: The SketchGuard Protocol**

## 3.2 Complexity Analysis and Performance Trade-offs

We analyze the per-node-per-round complexity of SKETCHGUARD compared to the best-case computation and communication complexity of existing Byzantine-robust methods. Let $d$ denote the model dimension, $|\mathcal{N}_i|$ the number of neighbors, $k$ the sketch size, and $|\mathcal{S}_i^t|$ the number of accepted neighbors. For computational complexity, SKETCHGUARD executes four phases per round:

- Sketch generation: $O(d)$
- Neighbor screening (compressed domain): $O(k \cdot |\mathcal{N}_i|)$
- Model verification & aggregation (accepted neighbors): $O(d \cdot |\mathcal{S}_i^t|)$

This yields total complexity $O(d + k \cdot |\mathcal{N}_i| + d \cdot |\mathcal{S}_i^t|)$, compared to $O(d \cdot |\mathcal{N}_i|)$ for traditional methods that perform full-precision screening of all neighbors.

**Table 2: Computational complexity comparison with state-of-the-art (SOTA) across training phases.**

| Training Phase | SOTA | SketchGuard |
|---|---|---|
| Local Training | $O(d)$ | $O(d)$ |
| Sketch Generation | – | $O(d)$ |
| Neighbor Screening | $O(d \cdot |\mathcal{N}_i|)$ | $O(k \cdot |\mathcal{N}_i|)$ |
| Model Verification & Aggregation | $O(d \cdot |\mathcal{N}_i|)$ | $O(d \cdot |\mathcal{S}_i^t|)$ |
| **Total Per Round** | $O(d \cdot |\mathcal{N}_i|)$ | $O(d + k \cdot |\mathcal{N}_i| + d \cdot |\mathcal{S}_i^t|)$ |

For communication complexity, SKETCHGUARD transmits two types of messages: sketches to all neighbors and full models to accepted neighbors. This results in $O(k \cdot |\mathcal{N}_i| + d \cdot |\mathcal{S}_i^t|)$ parameters per round, compared to $O(d \cdot |\mathcal{N}_i|)$ for existing approaches.

The theoretical requirement $k = O(\epsilon^{-2} \log(1/\delta))$ translates to practical sketch sizes that remain much smaller than model dimensions while providing strong approximation guarantees. SKETCH-GUARD's performance benefits scale with three factors: (1) the compression ratio $d/k$, where larger models yield greater savings, (2) the

filtering effectiveness ratio $|\mathcal{S}_i^t|/|\mathcal{N}_i|$, where successful Byzantine rejection reduces the number of full models fetched, and (3) network degree $|\mathcal{N}_i|$, where denser connectivity amplifies the screening overhead reduced by compression. While sketch generation and exchange introduce $O(d + k \cdot |\mathcal{N}_i|)$ overhead even in benign settings, this cost becomes negligible as model dimensionality and network scale increase, which characterizes the web-scale deployments where Byzantine robustness is necessary.

## 4 Convergence Analysis

In this section, we establish that SKETCHGUARD maintains convergence guarantees despite compressed filtering by proving that sketch-based screening preserves the robustness properties of full-precision similarity-based defenses with controlled degradation bounds.

### 4.1 Technical Assumptions

We state the standard assumptions that impact on convergence analysis.

ASSUMPTION 1 (STRONG CONVEXITY). *The population risk $F(\mathbf{w})$ is $\mu$-strongly convex, i.e., for all $\mathbf{w}_1, \mathbf{w}_2 \in \Theta$, one has that:*

$$F(\mathbf{w}_1) + \langle \nabla F(\mathbf{w}_1), \mathbf{w}_2 - \mathbf{w}_1 \rangle + \frac{\mu}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|^2 \leq F(\mathbf{w}_2).$$

ASSUMPTION 2 (SMOOTHNESS). *The population risk $F(\mathbf{w})$ is $L$-smooth, i.e., for all $\mathbf{w}_1, \mathbf{w}_2 \in \Theta$, we have that:*

$$\|\nabla F(\mathbf{w}_1) - \nabla F(\mathbf{w}_2)\| \leq L \|\mathbf{w}_1 - \mathbf{w}_2\|.$$

ASSUMPTION 3 (BOUNDED VARIANCE). *The stochastic gradient $\mathbf{g}(\mathbf{w}_i)$ computed by an honest client $i \in \mathcal{H}$ is an unbiased estimator of the true gradient, and $\mathbf{g}(\mathbf{w}_i)$ has bounded variance, where $\mathcal{H}$ is the set of honest clients. That is, $\forall i \in \mathcal{H}$, one has that:*

$$\mathbb{E}[\mathbf{g}(\mathbf{w}_i)] = \nabla F(\mathbf{w}_i), \quad \mathbb{E}[\|\mathbf{g}(\mathbf{w}_i) - \nabla F(\mathbf{w}_i)\|^2] \leq \delta^2.$$

ASSUMPTION 4 (BOUNDED PARAMETERS). *For any honest client $i \in \mathcal{H}$, the model $\mathbf{w}_i$ and $\|\nabla F(\mathbf{w}_i)\|$ are bounded. That is, $\forall i \in \mathcal{H}$, we have $\|\mathbf{w}_i\| \leq \psi$, and $\|\nabla F(\mathbf{w}_i)\| \leq \rho$.*

ASSUMPTION 5 (GRAPH CONNECTIVITY). *The subgraph induced by honest clients $G_{\mathcal{H}}$ remains connected throughout training.*

ASSUMPTION 6 (HASH FUNCTION SYNCHRONIZATION). *All clients use identical hash and sign functions for consistent sketching across the network.*

## 4.2 Main Convergence Results

Leveraging Lemma 1, we establish convergence guarantees for SKETCHGUARD that achieve similar results to those based on full-precision methods.

THEOREM 1 (STRONGLY CONVEX CONVERGENCE WITH COMPRESSION). *For $\mu$-strongly convex and $L$-smooth objectives, with learning rate $\eta \leq \min\{1/(4L), 1/\mu\}$ and effective threshold parameter $\gamma_{eff} = \gamma\sqrt{(1+\epsilon)/(1-\epsilon)}$, after $T$ rounds:*

$$\mathbb{E}[F(\mathbf{w}_i^T) - F(\mathbf{w}^*)] \leq (1 - \mu\eta)^T [F(\mathbf{w}_i^0) - F(\mathbf{w}^*)]$$
$$+ \frac{2L\eta\delta^2}{\mu} + \frac{2\gamma_{eff}\rho\psi(1-\alpha)}{\mu\eta} \tag{8}$$

*where the compression error manifests only through the $(1+\epsilon)/(1-\epsilon)$ factor in $\gamma_{eff}$.*

PROOF. See Appendix A. □

THEOREM 2 (NON-CONVEX CONVERGENCE WITH COMPRESSION). *For non-convex $L$-smooth objectives, with the same parameter choices:*

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}[\|\nabla F(\mathbf{w}_i^t)\|^2] \leq \frac{2[F(\mathbf{w}_i^0) - F^*]}{\eta T}$$
$$+ 4L\eta\delta^2 + \frac{4\gamma_{eff}\rho\psi(1-\alpha)}{\eta} \tag{9}$$

PROOF. See Appendix B. □

Both theorems demonstrate that SKETCHGUARD achieves optimal convergence rates for their respective settings [15], with sketch approximation introducing only a $(1 + O(\epsilon))$ multiplicative factor on the threshold-dependent terms in the convergence bounds. For practical values of $\epsilon$ (e.g., $\epsilon = 0.1$ yields $\gamma_{eff} \approx 1.1\gamma$), this degradation is minimal while enabling substantial compression.

## 4.3 Robustness Preservation Under Compression

Our final theoretical result establishes that sketch compression preserves Byzantine robustness guarantees with controlled degradation.

LEMMA 2 (SKETCH COMPRESSION PRESERVES BYZANTINE RESILIENCE). *Under the conditions of Theorems 1 and 2, sketch-based filtering with effective threshold $\gamma_{eff} = \gamma\sqrt{(1+\epsilon)/(1-\epsilon)}$ provides equivalent Byzantine resilience to full-precision filtering with threshold $\gamma_{eff}$. Specifically, compression does not create new attack vectors beyond the controlled threshold degradation.*

PROOF. This follows directly from the convergence analysis in Theorems 1 and 2. Both proofs bound the neighbor difference terms

under sketch-based filtering (Appendix A.5):

$$\left\| \frac{1}{|\mathcal{S}_i|} \sum_{j \in \mathcal{S}_i} (\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) \right\| \leq \gamma_{eff}\psi$$

This bound is identical in form to the full-precision case, differing only in the threshold parameter $\gamma \rightarrow \gamma_{eff}$. Since the convergence proofs do not require any additional structural assumptions about the filtering mechanism beyond this bound, any Byzantine strategy must satisfy the same constraints as in the full-precision case, but with the relaxed threshold $\gamma_{eff}$. Therefore, compression preserves the Byzantine resilience properties established by [13], with degradation bounded by the factor $\sqrt{(1+\epsilon)/(1-\epsilon)} = 1 + O(\epsilon)$. □

## 5 Performance Evaluation

We evaluate SKETCHGUARD[1] through comprehensive experiments demonstrating that sketch-based compression maintains identical Byzantine robustness to state-of-the-art full-precision defenses while substantially reducing computational and communication overhead.

### 5.1 Experimental Setup

*5.1.1 Datasets and Models.* We conduct experiments on two standard federated learning benchmarks from the LEAF suite [9], chosen for their naturally non-IID data distributions that closely reflect real-world heterogeneity in federated environments. We evaluate both robustness and scalability across diverse model sizes.

**FEMNIST** (Federated Extended MNIST) is a character recognition dataset with 62 classes distributed across 3,550 users, where each user represents a different writer's handwriting samples. We employ a convolutional neural network comprising two convolutional layers followed by two fully connected layers. Our baseline configuration uses 6.6M parameters, which we systematically scale from 220K to over 60M parameters in our dimensionality scaling experiments (Section 5.3) to evaluate SKETCHGUARD under varying model complexities representative of modern federated learning deployments.

**CelebA** is a celebrity face attributes dataset distributed across 9,343 users. We use it for binary smile classification (smiling or not smiling), with a CNN architecture adapted for RGB inputs (84×84 pixels) containing 2.2M parameters. The network consists of two convolutional layers and two fully connected layers. Complete architectural specifications for both models are provided in Appendix C.1.

*5.1.2 Network Topologies.* We evaluate defense mechanisms across carefully selected network topologies that capture the spectrum of connectivity patterns in real-world decentralized systems, from resource-constrained sparse networks to well-connected infrastructures. For robustness evaluation experiments in Section 5.2, we employ Erdős-Rényi (ER) random graphs [12], which naturally model peer-to-peer networks where connections form probabilistically. ER graphs are generated by independently connecting each pair of nodes with probability $p$, producing networks with variable local connectivity that reflect real-world decentralized learning

---

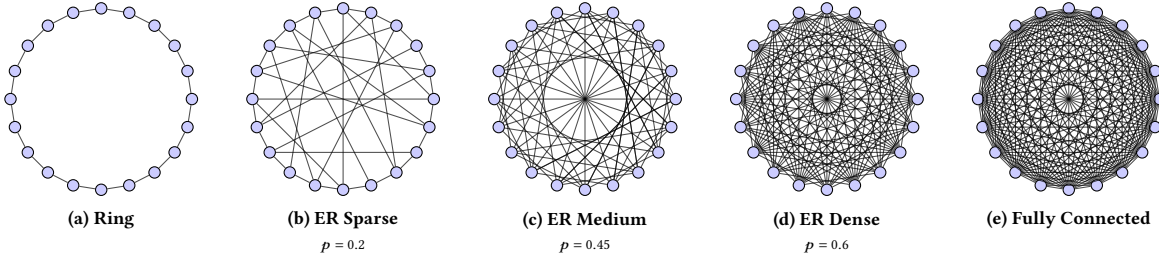[1]Code and experiment artifacts available at doi.org/10.5281/zenodo.17223405.

**Figure 2: Network topologies used in the robustness evaluation experiments.**

scenarios where participants discover neighbors organically rather than through engineered topologies. We evaluate five topologies spanning from minimal connectivity (Ring) to maximum connectivity (Fully Connected), with three intermediate ER graph densities as shown in Figure 2; complete specifications are provided in Appendix C.2.1. These experiments use 20 nodes to enable comprehensive evaluation across varying malicious client fractions and defense mechanisms while maintaining computational tractability. For network scaling experiments in Section 5.3, we employ k-regular graphs [6] where each node has exactly $k$ neighbors. This design choice ensures consistent per-node computational load across all defense mechanisms, enabling fair comparison of algorithmic efficiency independent of topology-induced variance. We evaluate networks scaling from 20 to 300 nodes, demonstrating SKETCHGUARD's scalability to large-scale deployments with hundreds of participants.

*5.1.3 Attack Models.* We evaluate Byzantine robustness using two representative attack strategies that represent different threat models and attack sophistication levels.

**Directed Deviation** attacks employ an optimization-based approach that crafts malicious updates by solving for parameters that maximize deviation in the direction opposite to honest gradient descent. This adaptive attack strategy represents a strong adversary that actively attempts to subvert Byzantine-robust aggregation mechanisms by carefully positioning malicious updates to evade detection while maximizing impact. We implement the directed deviation attack following the formulation in [14], which has demonstrated effectiveness against common Byzantine-robust aggregation rules including Krum and trimmed mean.

**Gaussian** attacks [5] inject random noise $\hat{\mathbf{w}}_j = \mathbf{w}_j + \mathcal{N}(0, \sigma^2 \mathbf{I})$ with $\sigma = 1$, representing less sophisticated but realistic adversaries that disrupt convergence through stochastic perturbations. While simpler than directed deviation, these attacks model practical scenarios where attackers lack complete knowledge of the aggregation mechanism or resources for optimization-based attacks.

We vary the fraction of Byzantine clients from 0% to 80% in 10% increments to evaluate robustness across different attack intensities. These attack strategies provide complementary evaluation: directed deviation establishes defense performance against sophisticated adaptive adversaries, while Gaussian attacks validate robustness against simpler but more common threat models.

*5.1.4 Baselines and Configuration.* We compare SKETCHGUARD against four baseline methods: D-FedAvg [21], KRUM [5], UBAR [17], and BALANCE [13], representing the spectrum from non-robust to state-of-the-art Byzantine-robust aggregation. For all experiments, except the $k$-ablation study, we employ sketch sizes $k = 1000$ for FEMNIST and $k = 350$ for CelebA. For the baseline model configurations, this achieves compression ratios of 6603:1 and 6342:1 respectively, with approximation error $\epsilon \lesssim 0.2$ and probability $1 - \delta > 0.99$. All experiments use standard DFL configurations with 10 global rounds and 3 local epochs per round. We repeat each experiment across 3 random seeds and report mean results; observed variances are minimal across all metrics. Complete hyperparameter specifications are provided in Appendix C.2.2.

*5.1.5 Evaluation Metrics.* We evaluate defense mechanisms using three complementary metrics that capture robustness, computational efficiency, and communication overhead.

**Test Error Rate (TER).** Defined as $1 - $ test accuracy, TER measures performance degradation averaged across honest clients. Lower values indicate better Byzantine resilience.

**Per-Round Computation Time.** Wall-clock time for neighbor screening and aggregation phases per client per round, excluding local training which remains constant across methods. All measurements are recorded on identical hardware.

**Communication Overhead.** Total parameters transmitted per client per round, including sketches and full model exchanges. We report both absolute counts and reduction percentages relative to full-precision baselines.

## 5.2 Byzantine Robustness Evaluation

*5.2.1 Attack Strategy Resilience.* Figure 3 evaluates SKETCHGUARD's robustness against both directed deviation and Gaussian attacks on FEMNIST and CelebA datasets, averaged across all network topologies. SKETCHGUARD achieves robustness equivalent to state-of-the-art methods BALANCE and UBAR across both attack strategies. Compared to BALANCE, SKETCHGUARD exhibits statistically identical performance with mean absolute TER deviation of 0.02 percentage points and maximum deviation of 0.49 percentage points. Compared to UBAR, SKETCHGUARD shows mean absolute TER deviation of 0.35 percentage points and maximum deviation of 1.51 percentage points, with SKETCHGUARD having marginally higher TER on average by 0.20 percentage points; however, UBAR's loss-based
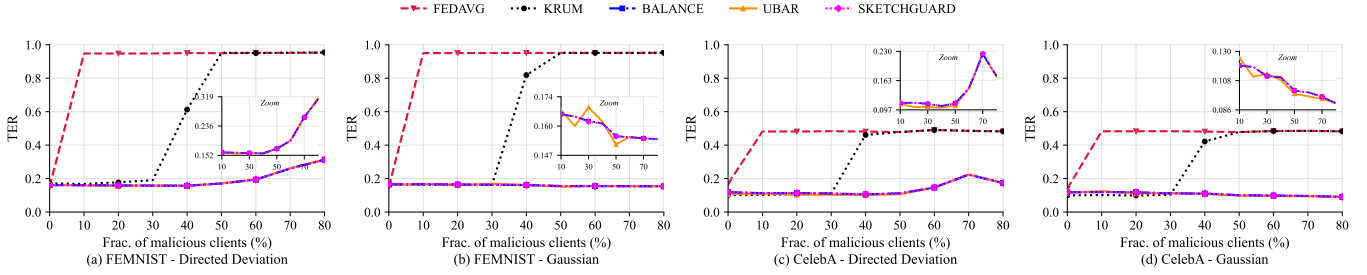
**Figure 3: Impact of fraction of malicious clients on TER across different datasets and attack types.**

screening mechanism requires significantly higher computational overhead than both BALANCE and SKETCHGUARD, as demonstrated in Section 5.3. Against both directed deviation (Figures 3a, 3c) and Gaussian attacks (Figures 3b, 3d), SKETCHGUARD maintains consistently low error rates across all Byzantine client fractions, demonstrating that sketch-based compression preserves the robustness properties of full-precision similarity-based defenses while enabling substantial efficiency gains.

*5.2.2 Sketch Size Robustness.* SKETCHGUARD demonstrates remarkable robustness to sketch size selection. Across $k$ values ranging from 500 to 100,000 on both FEMNIST and CelebA datasets under 50% Byzantine clients with directed deviation attacks, performance remains completely stable (TER = 15.59% for FEMNIST, 10.35% for CelebA across all $k$ values). Even with compression ratios exceeding 13,000:1 ($k$ = 500 for FEMNIST's 6.6M parameters), SKETCH-GUARD preserves full Byzantine resilience. This insensitivity to $k$ stems from Count Sketch's distance preservation guarantees: while smaller $k$ values introduce larger approximation errors, the resulting sketch-based distances remain sufficiently accurate to distinguish Byzantine node signatures from honest node behavior. The adaptive threshold mechanism $\gamma_{\text{eff}}$ compensates for increased approximation error, ensuring robust filtering across the entire practical range of sketch sizes.

*5.2.3 Topology-Dependent Robustness.* Figure 4 presents SKETCH-GUARD's robustness across five network topologies, averaged over attack types and datasets. SKETCHGUARD demonstrates identical robustness to BALANCE and UBAR across all network topologies, from minimal connectivity in Ring topology to maximum connectivity in Fully Connected networks. In well-connected topologies (ER graphs with $p \geq 0.45$, Fully Connected), all three Byzantine-robust methods maintain low error rates below 20% even when 80% of the network is compromised, significantly outperforming KRUM and the non-robust D-FedAvg baseline. In extremely sparse topologies (Ring), BALANCE, UBAR, and SKETCHGUARD exhibit graceful degradation, tolerating up to 40% Byzantine clients before substantial performance loss. Despite having only two neighbors per node in the Ring topology, SKETCHGUARD maintains performance parity with full-precision methods across all attack intensities, demonstrating that sketch-based compression does not compromise filtering effectiveness even under severe connectivity constraints.

## 5.3 Computational Efficiency

Figure 5 evaluates SKETCHGUARD's computational efficiency through network size and model dimensionality scaling experiments on FEMNIST with 50% Byzantine clients under directed deviation attacks. Network scaling uses the baseline 6.6M parameter model across varying k-regular graph sizes, while model dimensionality scaling uses a fixed k-regular graph with node degree 154 across varying model sizes as detailed in Appendix C.3.

*5.3.1 Network Size Scalability.* SKETCHGUARD exhibits near-constant computational cost as network connectivity increases from 16 to 299 neighbors (Figure 5a), averaging 0.35s per round, while BALANCE and UBAR scale linearly with neighbor count. At smaller network sizes (≤32 neighbors), SKETCHGUARD's sketch generation overhead dominates, resulting in higher absolute costs than full-precision methods. However, as connectivity increases beyond 96 neighbors, SKETCHGUARD's compressed screening operations yield substantial benefits. The efficiency gap widens multiplicatively with network size: at 154 neighbors, SKETCHGUARD requires 0.34s compared to BALANCE's 0.51s and UBAR's 1.10s, achieving 33% and 69% savings respectively; at 299 neighbors, these savings expand to 60% and 82%, with SKETCHGUARD requiring only 0.39s versus BALANCE's 0.99s and UBAR's 2.14s. This scalability advantage stems from neighbor screening operating in compressed $k$-dimensional space rather than full $d$-dimensional space, yielding multiplicative gains proportional to network connectivity.

*5.3.2 Model Dimensionality Scalability.* SKETCHGUARD demonstrates sub-linear scaling as model size increases from 220K to 60M parameters (Figure 5b), contrasting sharply with the linear scaling of BALANCE and UBAR. For models exceeding 26M parameters, SKETCHGUARD maintains average computational improvement of 64% over state-of-the-art methods. At 60M parameters, SKETCH-GUARD requires 2.0s versus BALANCE's 4.2s and UBAR's 10.0s, achieving 52% and 80% reductions respectively. This sub-linear scaling emerges from Count Sketch's property that sketch size $k$ depends on approximation parameters ($\epsilon$, $\delta$) rather than model dimension $d$. While full-precision methods incur $O(d|\mathcal{N}_i|)$ comparison costs that grow linearly with model size, SKETCHGUARD's sketch comparison cost remains fixed at $O(k|\mathcal{N}_i|)$ regardless of $d$, proving particularly valuable for large-scale DFL scenarios involving billion-parameter models [7].
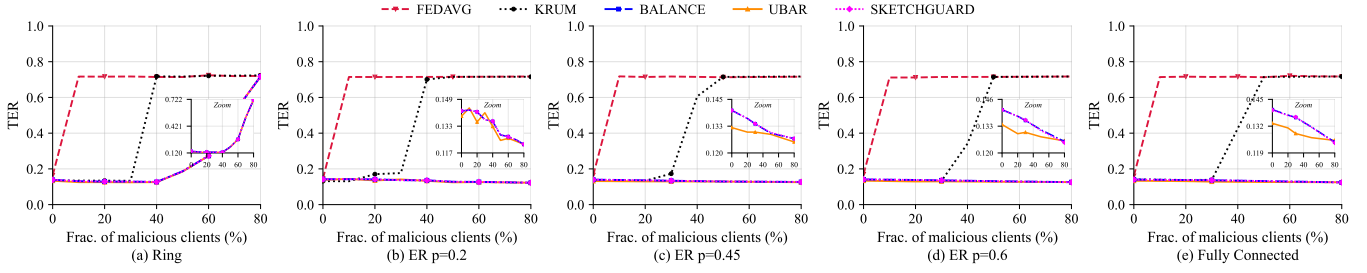
Figure 4: Impact of fraction of malicious clients on TER across different network topologies.
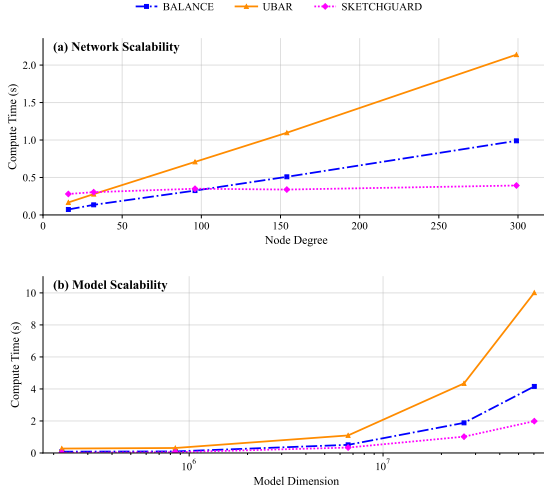


Figure 5: Impact of network size and model dimensionality on per-node computation time.

Table 3: Communication Overhead Comparison

| Scenario | Accepted Neighbors | Reduction |
|---|---|---|
| Benign (no filtering) | $|\mathcal{S}_i| = |\mathcal{N}_i|$ | <0.02% penalty |
| 50% Byzantine filtering | $|\mathcal{S}_i| \approx 0.5|\mathcal{N}_i|$ | ~50% |
| 70% Byzantine filtering | $|\mathcal{S}_i| \approx 0.3|\mathcal{N}_i|$ | ~70% |

- **SketchGuard**: Transmits 100K sketch parameters plus 330M full model parameters, totaling ~330M parameters.
- **Full-Precision Methods**: Transmit 660M parameters to all neighbors.
- **Reduction**: ~50% bandwidth savings.

These bandwidth savings prove critical for resource-constrained devices and bandwidth-limited networks characteristic of real-world decentralized learning deployments. As model sizes grow and networks expand, SketchGuard's communication efficiency becomes essential for practical deployment.

## 6 Conclusions and Future Work

We proposed SketchGuard, a framework for scaling Byzantine-robust decentralized federated learning through sketch-based neighbor screening. By decoupling Byzantine filtering from model aggregation, our approach reduces communication complexity from $O(d|\mathcal{N}_i|)$ to $O(k|\mathcal{N}_i| + d|\mathcal{S}_i|)$ while preserving robustness guarantees of full-precision similarity-based defenses.

Our theoretical analysis establishes that Count Sketch compression maintains Byzantine resilience with controlled degradation bounds, where approximation errors introduce only a $(1+O(\epsilon))$ factor in the effective threshold parameter. We prove convergence guarantees in both strongly convex and non-convex settings, achieving optimal rates for their respective problem classes. Comprehensive experiments demonstrate that SketchGuard maintains identical robustness to state-of-the-art methods, with mean TER deviation of only up to 0.35 percentage points. This robustness equivalence is achieved while reducing computation time by up to 82%, with efficiency gains scaling multiplicatively with model dimensionality and network connectivity. Communication overhead scales proportionally with the fraction of accepted neighbors, achieving reductions of 50-70% under typical Byzantine attack scenarios.

SketchGuard addresses the fundamental scalability bottleneck limiting real-world deployment of robust decentralized learning. Future work includes efficient hash synchronization protocols for

## 5.4 Communication Efficiency

SketchGuard reduces per-round communication complexity from $O(d|\mathcal{N}_i|)$ to $O(k|\mathcal{N}_i| + d|\mathcal{S}_i|)$, where $|\mathcal{S}_i| \leq |\mathcal{N}_i|$ denotes accepted neighbors post-filtering.

*5.4.1 Compression Ratios and Bandwidth Reduction.* With compression ratios exceeding 6300:1 for both datasets, sketch transmission overhead becomes negligible compared to full model exchanges. Table 3 presents communication efficiency under different filtering scenarios. In scenarios where Byzantine filtering successfully rejects 70% of malicious neighbors, SketchGuard achieves approximately 70% communication reduction relative to full-precision methods. Even in benign settings where all neighbors are accepted, the sketch overhead of $O(k|\mathcal{N}_i|)$ represents less than 0.02% additional cost compared to the $O(d|\mathcal{N}_i|)$ baseline due to the compression ratio $k/d \approx 1/6300$, introducing negligible bandwidth penalty.

*5.4.2 Multiplicative Scaling Benefits.* The communication savings scale multiplicatively with both model dimensionality and network connectivity. Consider a fully connected network with 100 neighbors per node and FEMNIST's 6.6M parameter model under 50% Byzantine clients, where Count Sketch compresses models to 1K parameters:

dynamic networks, adaptive sketch sizing based on attack intensity, and theoretical extensions to non-IID data distributions. As FL expands into web-scale deployments with untrusted participants, scalable robustness mechanisms like SKETCHGUARD become essential infrastructure for maintaining both security and performance in distributed machine learning systems.

## Acknowledgments

## References

[1] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. 2017. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in neural information processing systems* 30 (2017).

[2] Gilad Baruch, Moran Baruch, and Yoav Goldberg. 2019. A little is enough: Circumventing defenses for distributed learning. *Advances in Neural Information Processing Systems* 32 (2019).

[3] Enrique Tomás Martínez Beltrán et al. 2023. Decentralized federated learning: Fundamentals, state of the art, frameworks, trends, and challenges. *IEEE Communications Surveys & Tutorials* 25, 4 (2023), 2983–3013.

[4] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. 2018. signSGD: Compressed optimisation for non-convex problems. In *International conference on machine learning*. PMLR, 560–569.

[5] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. 2017. Machine learning with adversaries: Byzantine tolerant gradient descent. *Advances in neural information processing systems* 30 (2017).

[6] John Adrian Bondy, Uppaluri Siva Ramachandra Murty, et al. 1976. *Graph theory with applications*. Vol. 290. Macmillan London.

[7] Alexander Borzunov, Max Ryabinin, Artem Chumachenko, Dmitry Baranchuk, Tim Dettmers, Younes Belkada, Pavel Samygin, and Colin A Raffel. 2023. Distributed inference and fine-tuning of large language models over the internet. *Advances in neural information processing systems* 36 (2023), 12312–12331.

[8] Diego Cajaraville-Aboy, Ana Fernández-Vilas, Rebeca P Díaz-Redondo, and Manuel Fernández-Veiga. 2024. Byzantine-robust aggregation for securing decentralized federated learning. *arXiv preprint arXiv:2409.17754* (2024).

[9] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečnỳ, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. 2018. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097* (2018).

[10] Moses Charikar, Kevin Chen, and Martin Farach-Colton. 2002. Finding frequent items in data streams. In *International Colloquium on Automata, Languages, and Programming*. Springer, 693–703.

[11] El Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Arsany Guirguis, Lê-Nguyên Hoang, and Sébastien Rouault. 2021. Collaborative learning in the jungle (decentralized, byzantine, heterogeneous, asynchronous and nonconvex learning). *Advances in neural information processing systems* 34 (2021), 25044–25057.

[12] P. Erdős and A. Rényi. 1960. On the Evolution of Random Graphs. *Publication of the Mathematical Institute of the Hungarian Academy of Sciences* 5 (1960), 17–61.

[13] Minghong Fang et al. 2024. Byzantine-robust decentralized federated learning. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*. 2874–2888.

[14] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. 2020. Local model poisoning attacks to {Byzantine-Robust} federated learning. In *29th USENIX security symposium (USENIX Security 20)*. 1605–1622.

[15] Guillaume Garrigos and Robert M Gower. 2023. Handbook of convergence theorems for (stochastic) gradient methods. *arXiv preprint arXiv:2301.11235* (2023).

[16] Vineet Sunil Gattani, Junshan Zhang, and Gautam Dasarathy. 2024. Communication-Efficient Federated Learning over Wireless Channels via Gradient Sketching. *arXiv preprint arXiv:2410.23424* (2024).

[17] Shangwei Guo, Tianwei Zhang, Han Yu, Xiaofei Xie, Lei Ma, Tao Xiang, and Yang Liu. 2021. Byzantine-resilient decentralized stochastic gradient descent. *IEEE Transactions on Circuits and Systems for Video Technology* 32, 6 (2021), 4096–4106.

[18] Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. 2021. Federated learning with compression: Unified analysis and sharp guarantees. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2350–2358.

[19] Lie He, Sai Praneeth Karimireddy, and Martin Jaggi. 2022. Byzantine-robust decentralized learning via clippedgossip. *arXiv preprint arXiv:2202.01545* (2022).

[20] Jakub Konečnỳ, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. 2016. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492* (2016).

[21] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. 2017. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in neural information processing systems* 30 (2017).

[22] Alexander Long. 2024. Protocol Learning, Decentralized Frontier Risk and the No-Off Problem. *arXiv preprint arXiv:2412.07890* (2024).

[23] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*. PMLR, 1273–1282.

[24] Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C Lupu, and Fabio Roli. 2017. Towards poisoning of deep learning algorithms with back-gradient optimization. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*. 27–38.

[25] Krishna Pillutla, Sham M Kakade, and Zaid Harchaoui. 2022. Robust aggregation for federated learning. *IEEE Transactions on Signal Processing* 70 (2022), 1142–1154.

[26] Daniel Rothchild et al. 2020. Fetchsgd: Communication-efficient federated learning with sketching. In *International Conference on Machine Learning*. PMLR, 8253–8265.

[27] Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. 2019. Robust and communication-efficient federated learning from non-iid data. *IEEE transactions on neural networks and learning systems* 31, 9 (2019), 3400–3413.

[28] Virat Shejwalkar and Amir Houmansadr. 2021. Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning. In *NDSS*.

[29] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. 2018. Sparsified SGD with memory. *Advances in neural information processing systems* 31 (2018).

[30] Peng Sun, Xinyang Liu, Zhibo Wang, and Bo Liu. 2024. Byzantine-robust decentralized federated learning via dual-domain clustering and trust bootstrapping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 24756–24765.

[31] Tianxiang Wang, Zhonglong Zheng, and Feilong Lin. 2025. Federated learning framework based on trimmed mean aggregation rules. *Expert Systems with Applications* (2025), 126354.

[32] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. 2020. Fall of empires: Breaking byzantine-tolerant sgd by inner product manipulation. In *Uncertainty in artificial intelligence*. PMLR, 261–270.

[33] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. 2018. Byzantine-robust distributed learning: Towards optimal statistical rates. In *International conference on machine learning*. Pmlr, 5650–5659.

## A  Proof of Theorem 1

We provide a detailed proof of the convergence guarantee for SKETCHGUARD under strongly convex objectives with sketch-based compression. The proof structure follows the analysis framework of BALANCE [13], with key modifications to handle Count Sketch compression errors through the effective threshold parameter $\gamma_{\text{eff}}$.

### A.1  Setup and Notation

We denote by $\mathbf{g}(\mathbf{w}_i^t)$ the stochastic gradient computed by client $i$ at round $t$. The local update is:

$$\mathbf{w}_i^{t+1/2} = \mathbf{w}_i^t - \eta \mathbf{g}(\mathbf{w}_i^t), \tag{10}$$

where $\eta > 0$ is the learning rate. For simplicity in the proof, we drop the superscript $t$ in $\mathcal{S}_i^t$ and write $\mathcal{S}_i$.

### A.2  Model Update Analysis

The model update for client $i$ can be expressed as:

$$\mathbf{w}_i^{t+1} - \mathbf{w}_i^t = \alpha \mathbf{w}_i^{t+1/2}$$
$$+ (1-\alpha) \frac{1}{|\mathcal{S}_i|} \sum_{j \in \mathcal{S}_i} \mathbf{w}_j^{t+1/2} - \mathbf{w}_i^t \tag{11}$$

$$= \alpha \mathbf{w}_i^{t+1/2} + (1-\alpha)\frac{1}{|\mathcal{S}_i|}$$
$$\times \sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2} + \mathbf{w}_i^{t+1/2}) - \mathbf{w}_i^t \quad (12)$$

$$= \mathbf{w}_i^{t+1/2} + \frac{1-\alpha}{|\mathcal{S}_i|}$$
$$\times \sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) - \mathbf{w}_i^t \quad (13)$$

$$= -\eta \mathbf{g}(\mathbf{w}_i^t) + \frac{1-\alpha}{|\mathcal{S}_i|}$$
$$\times \sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) \quad (14)$$

## A.3 Applying Smoothness

By the $L$-smoothness assumption, we have:

$$F(\mathbf{w}_i^{t+1}) \le F(\mathbf{w}_i^t) + \langle \nabla F(\mathbf{w}_i^t), \mathbf{w}_i^{t+1} - \mathbf{w}_i^t \rangle$$
$$+ \frac{L}{2}\|\mathbf{w}_i^{t+1} - \mathbf{w}_i^t\|^2 \quad (15)$$

Substituting the model update expression:

$$F(\mathbf{w}_i^{t+1}) \le F(\mathbf{w}_i^t) - \eta\langle \nabla F(\mathbf{w}_i^t), \mathbf{g}(\mathbf{w}_i^t)\rangle$$
$$+ \langle \nabla F(\mathbf{w}_i^t), \frac{1-\alpha}{|\mathcal{S}_i|}\sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2})\rangle$$
$$+ \frac{L}{2}\left\| -\eta\mathbf{g}(\mathbf{w}_i^t) \right.$$
$$\left. + \frac{1-\alpha}{|\mathcal{S}_i|}\sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) \right\|^2 \quad (16)$$

## A.4 Bounding the Quadratic Term

Using the inequality $\|a+b\|^2 \le 2\|a\|^2 + 2\|b\|^2$:

$$\left\| -\eta\mathbf{g}(\mathbf{w}_i^t) + \frac{1-\alpha}{|\mathcal{S}_i|}\sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) \right\|^2$$
$$\le 2\eta^2\|\mathbf{g}(\mathbf{w}_i^t)\|^2$$
$$+ 2\left\| \frac{1-\alpha}{|\mathcal{S}_i|}\sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) \right\|^2 \quad (17)$$

## A.5 Impact of Sketch-Based Filtering

With sketch-based filtering, neighbor $j$ is accepted if:

$$\|\text{CS}(\mathbf{w}_i^{t+1/2}) - \text{CS}(\mathbf{w}_j^{t+1/2})\|$$
$$\le \gamma\exp(-\kappa t/T)\|\text{CS}(\mathbf{w}_i^{t+1/2})\| \quad (18)$$

By the distance preservation property of Count Sketch [10], this implies:

$$\|\mathbf{w}_i^{t+1/2} - \mathbf{w}_j^{t+1/2}\|$$
$$\le \gamma\sqrt{\frac{1+\epsilon}{1-\epsilon}}\exp(-\kappa t/T)\|\mathbf{w}_i^{t+1/2}\| \quad (19)$$

Define $\gamma_{\text{eff}} = \gamma\sqrt{(1+\epsilon)/(1-\epsilon)}$. Then:

$$\left\| \frac{1}{|\mathcal{S}_i|}\sum_{j \in \mathcal{S}_i}(\mathbf{w}_j^{t+1/2} - \mathbf{w}_i^{t+1/2}) \right\| \le \gamma_{\text{eff}}\psi \quad (20)$$

where we used the bounded parameters assumption that $\|\mathbf{w}_i^{t+1/2}\| \le \psi$.

## A.6 Taking Expectation

Taking expectation and using the bounded variance assumption:

$$\mathbb{E}[F(\mathbf{w}_i^{t+1})] \le \mathbb{E}[F(\mathbf{w}_i^t)] - \eta\|\nabla F(\mathbf{w}_i^t)\|^2$$
$$+ (1-\alpha)\gamma_{\text{eff}}\psi\rho$$
$$+ L\eta^2(\|\nabla F(\mathbf{w}_i^t)\|^2 + \delta^2)$$
$$+ L(1-\alpha)^2\gamma_{\text{eff}}^2\psi^2 \quad (21)$$

## A.7 Simplifying with Parameter Constraints

With $\eta \le 1/(4L)$, we have $L\eta^2 \le \eta/4$:

$$\mathbb{E}[F(\mathbf{w}_i^{t+1})] \le \mathbb{E}[F(\mathbf{w}_i^t)] - \frac{\eta}{2}\|\nabla F(\mathbf{w}_i^t)\|^2$$
$$+ L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha) \quad (22)$$

where we chose $\gamma_{\text{eff}} \le \rho/(L\psi(1-\alpha))$.

## A.8 Applying Strong Convexity

By the strong convexity assumption and the Polyak-Łojasiewicz inequality:

$$\|\nabla F(\mathbf{w}_i^t)\|^2 \ge 2\mu(F(\mathbf{w}_i^t) - F(\mathbf{w}^*)) \quad (23)$$

Therefore:

$$\mathbb{E}[F(\mathbf{w}_i^{t+1}) - F(\mathbf{w}^*)]$$
$$\le (1-\mu\eta)\mathbb{E}[F(\mathbf{w}_i^t) - F(\mathbf{w}^*)]$$
$$+ L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha) \quad (24)$$

## A.9 Telescoping and Final Bound

Telescoping over $t = 0, 1, \ldots, T-1$:

$$\mathbb{E}[F(\mathbf{w}_i^T) - F(\mathbf{w}^*)]$$
$$\le (1-\mu\eta)^T[F(\mathbf{w}_i^0) - F(\mathbf{w}^*)]$$
$$+ \sum_{t=0}^{T-1}(1-\mu\eta)^{T-1-t}$$
$$\times (L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha))$$
$$= (1-\mu\eta)^T[F(\mathbf{w}_i^0) - F(\mathbf{w}^*)]$$
$$+ \frac{1-(1-\mu\eta)^T}{\mu\eta}$$
$$\times (L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha))$$
$$\le (1-\mu\eta)^T[F(\mathbf{w}_i^0) - F(\mathbf{w}^*)]$$
$$+ \frac{2L\eta\delta^2}{\mu} + \frac{2\gamma_{\text{eff}}\rho\psi(1-\alpha)}{\mu\eta} \quad (25)$$

This completes the proof of Theorem 1. The key observation is that the compression error only affects the convergence through $\gamma_{\text{eff}} = \gamma\sqrt{(1+\epsilon)/(1-\epsilon)}$, introducing a controllable degradation

factor while maintaining the same convergence rate as the state-of-the-art.

## B  Proof of Theorem 2

We establish the convergence guarantee for SKETCHGUARD in non-convex settings, adapting the BALANCE analysis framework [13] for sketch compression.

### B.1  Starting from Smoothness

Following the analysis from Appendix A up to the expectation bound, we have:

$$\mathbb{E}[F(\mathbf{w}_i^{t+1})] \le \mathbb{E}[F(\mathbf{w}_i^t)] - \frac{\eta}{2}\|\nabla F(\mathbf{w}_i^t)\|^2$$
$$+ 2L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha) \qquad (26)$$

### B.2  Rearranging for Gradient Norm

Rearranging the inequality:

$$\frac{\eta}{2}\mathbb{E}[\|\nabla F(\mathbf{w}_i^t)\|^2]$$
$$\le \mathbb{E}[F(\mathbf{w}_i^t) - F(\mathbf{w}_i^{t+1})]$$
$$+ 2L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha) \qquad (27)$$

### B.3  Telescoping

Summing from $t = 0$ to $T - 1$:

$$\frac{\eta}{2}\sum_{t=0}^{T-1}\mathbb{E}[\|\nabla F(\mathbf{w}_i^t)\|^2]$$
$$\le F(\mathbf{w}_i^0) - \mathbb{E}[F(\mathbf{w}_i^T)]$$
$$+ T(2L\eta^2\delta^2 + 2\gamma_{\text{eff}}\psi\rho(1-\alpha)) \qquad (28)$$

### B.4  Averaging and Using Lower Bound

Since $F(\mathbf{w}_i^T) \ge F^*$:

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}[\|\nabla F(\mathbf{w}_i^t)\|^2]$$
$$\le \frac{2(F(\mathbf{w}_i^0) - F^*)}{\eta T}$$
$$+ 4L\eta\delta^2 + \frac{4\gamma_{\text{eff}}\psi\rho(1-\alpha)}{\eta} \qquad (29)$$

This establishes Theorem 2. The convergence rate matches the optimal rate for non-convex optimization, with the compression error appearing only through $\gamma_{\text{eff}}$.

## C  Supplementary Experimental Details

This appendix provides comprehensive details about the experimental setup, including model architectures, hyperparameter configurations, network topologies, and scalability experiment specifications used throughout our study.

### C.1  Detailed Model Architectures

We provide complete architectural specifications for the two datasets used in our experiments: FEMNIST and CelebA. Both architectures follow standard designs adapted to their respective input dimensions and task requirements.

*C.1.1  FEMNIST Architecture.* The FEMNIST model follows the original LEAF specification with a convolutional neural network design. The architecture consists of two convolutional layers followed by max pooling and two fully connected layers. Table 4 provides the complete layer-by-layer specification, including parameter counts and output dimensions.

**Table 4: FEMNIST Model Architecture Details**

| Layer | Type | Parameters | Output Size |
|---|---|---|---|
| Input | – | – | $1 \times 28 \times 28$ |
| Conv1 | Conv2d | $32 \times (5 \times 5 \times 1) + 32$ | $32 \times 28 \times 28$ |
| Pool1 | MaxPool2d | – | $32 \times 14 \times 14$ |
| Conv2 | Conv2d | $64 \times (5 \times 5 \times 32) + 64$ | $64 \times 14 \times 14$ |
| Pool2 | MaxPool2d | – | $64 \times 7 \times 7$ |
| Flatten | – | – | 3136 |
| FC1 | Linear | $3136 \times 2048 + 2048$ | 2048 |
| FC2 | Linear | $2048 \times 62 + 62$ | 62 |
| **Total** | | **6,603,710** | |

*C.1.2  CelebA Architecture.* The CelebA model uses a LeNet-style architecture adapted for larger RGB input images ($84 \times 84$ pixels). The network has a similar structure to FEMNIST but with modified convolutional layers to handle color images and a binary classification output. The complete architecture is detailed in Table 5.

**Table 5: CelebA Model Architecture Details**

| Layer | Type | Parameters | Output Size |
|---|---|---|---|
| Input | – | – | $3 \times 84 \times 84$ |
| Conv1 | Conv2d | $30 \times (3 \times 3 \times 3) + 30$ | $30 \times 84 \times 84$ |
| Pool1 | MaxPool2d | – | $30 \times 42 \times 42$ |
| Conv2 | Conv2d | $50 \times (3 \times 3 \times 30) + 50$ | $50 \times 42 \times 42$ |
| Pool2 | MaxPool2d | – | $50 \times 21 \times 21$ |
| Flatten | – | – | 22050 |
| FC1 | Linear | $22050 \times 100 + 100$ | 100 |
| FC2 | Linear | $100 \times 2 + 2$ | 2 |
| **Total** | | **2,219,692** | |

### C.2  Complete Experimental Configuration

This section details all hyperparameters and configuration settings used across our experiments to ensure reproducibility.

*C.2.1  Network Topology Specifications.* We evaluate defense mechanisms across five different network topologies representing varying levels of connectivity. The ring topology represents minimal connectivity (degree 2), while Erdős-Rényi graphs with varying connection probabilities $p$ provide intermediate connectivity levels. The fully connected topology represents maximum connectivity where each node connects to all others. Table 6 specifies the parameters and expected node degrees for each topology in our 20-node experimental networks.

**Table 6: Network Topology Configurations**

| Topology | Parameters | Expected Degree |
|----------|------------|-----------------|
| Ring | – | 2 |
| Erdős-Rényi (sparse) | $p = 0.2$ | 3.8 |
| Erdős-Rényi (medium) | $p = 0.45$ | 8.55 |
| Erdős-Rényi (dense) | $p = 0.6$ | 11.4 |
| Fully Connected | – | 19 |

*C.2.2 Hyperparameter Settings.* Table 7 presents the complete set of hyperparameters used in our experiments, organized by category. Training configuration parameters were selected based on standard federated learning practices. Defense mechanism parameters ($\gamma$, $\kappa$, $\alpha$) were tuned to balance robustness and model utility. For Sketch-Guard, we set the sketch size $k = 1000$ for FEMNIST and $k = 350$ for CelebA to get roughly the same factor of compression for both datasets.

**Table 7: Complete Hyperparameter Configuration**

| Parameter | Value |
|-----------|-------|
| *Training Configuration* | |
| Number of clients | 20 |
| Global epochs | 10 |
| Local epochs per round | 3 |
| Batch size | 64 |
| Learning rate | 0.01 (FEMNIST), 0.001 (CelebA) |
| Maximum samples per client | 10,000 (FEMNIST), 4,500 (CelebA) |
| Random seed | 987654321, 39573295, 32599368 |
| *SketchGuard Parameters* | |
| Sketch size ($k$) | 1,000 (FEMNIST), 350 (CelebA) |
| Threshold parameter ($\gamma$) | 2.0 |
| Decay parameter ($\kappa$) | 1.0 |
| Mixing parameter ($\alpha$) | 0.5 |
| Hash seed | 42 |
| *BALANCE Parameters* | |
| Threshold parameter ($\gamma$) | 2.0 |
| Decay parameter ($\kappa$) | 1.0 |
| Mixing parameter ($\alpha$) | 0.5 |
| *UBAR Parameters* | |
| Robustness parameter ($\rho$) | $1.0 -$ attack percentage |
| *KRUM Parameters* | |
| Compromised fraction | Attack percentage |

## C.3 Scalability Experiments

We conducted two types of scalability experiments to evaluate how defense mechanisms perform under different scaling conditions: network size scaling and model dimensionality scaling.

*C.3.1 Network Size Scaling.* To evaluate the computational scalability of defense mechanisms with respect to network size, we conducted experiments on k-regular graphs with varying numbers of participants. Table 8 shows the configurations used, with network sizes ranging from 20 to 300 nodes. All networks maintain 50% Byzantine nodes to test defense robustness under high attack scenarios. The node degree was selected to ensure connectivity while

maintaining realistic peer-to-peer network constraints. These experiments used shorter training runs (3 rounds with 1 local epoch each) to focus on measuring computational overhead rather than convergence behavior.

**Table 8: Network Size Scaling Configurations**

| Node Degree | Network Size | Attack Percentage |
|-------------|--------------|-------------------|
| 16 | 20 | 50% |
| 32 | 35 | 50% |
| 96 | 100 | 50% |
| 154 | 155 | 50% |
| 299 | 300 | 50% |

*C.3.2 Model Scaling Variants.* To investigate the impact of model dimensionality on defense mechanism performance, we created five variants of the FEMNIST architecture with different parameter counts. Table 9 describes these variants, ranging from the Tiny model with approximately 220K parameters to the XLarge model with over 60M parameters. These variants were created by systematically scaling the number of convolutional filters and fully connected layer dimensions while maintaining the overall architectural structure.

**Table 9: Model Scaling Variants**

| Variant | Architecture Modifications | Parameters |
|---------|---------------------------|------------|
| Tiny | Reduced filter counts and hidden units | 220,318 |
| Small | Standard configuration | 848,382 |
| Standard | Baseline FEMNIST architecture | 6,603,710 |
| Large | Increased filter counts and hidden units | 26,154,814 |
| XLarge | Further increased dimensions | 60,271,678 |