

---

# Secure Byzantine-Robust Distributed Learning via Clustering

---

**Raj Kiriti Velicheti**

Coordinated Sciences Laboratory  
University of Illinois at Urbana-Champaign  
rkv4@illinois.edu

**Derek Xia**

Department of Computer Science  
University of Illinois at Urbana-Champaign  
derekx3@illinois.edu

**Oluwasanmi Koyejo**

Department of Computer Science  
University of Illinois at Urbana-Champaign  
sanmi@illinois.edu

## Abstract

Federated learning systems that jointly preserve Byzantine robustness and privacy have remained an open problem. Robust aggregation, the standard defense for Byzantine attacks, generally requires server access to individual updates or nonlinear computation – thus is incompatible with privacy-preserving methods such as secure aggregation via multiparty computation. To this end, we propose SHARE (Secure Hierarchical Robust Aggregation), a distributed learning framework designed to cryptographically preserve client update privacy and robustness to Byzantine adversaries simultaneously. The key idea is to incorporate secure averaging among randomly clustered clients before filtering malicious updates through robust aggregation. Experiments show that SHARE has similar robustness guarantees as existing techniques while enhancing privacy.

## 1 Introduction

An increasing amount of data is being collected in a decentralized manner on devices across institutions[9]. Traditionally, machine learning with such devices require centralized data collection, which increases communication costs while posing a threat to privacy, especially when these devices gather personal user data. Distributed learning frameworks like federated learning attempt to address these issues by sharing model updates from client devices, rather than data, to a centralized server [11, 9, 7, 14].

Among the most popular implementations of federated learning is Federated Averaging [15]. While the central coordinating server follows a designated aggregation protocol, the required communication can pose a privacy threat when the system is compromised by a malicious external agent leaking individual model updates. To this end, Bonawitz et al. [3] proposed a secure averaging oracle that masks individual client updates such that the server learns their average alone. Nevertheless, since the collaboratively learned model update includes the contribution of all participating clients, benign averaging might fall prey to incorrect device updates either due to arbitrary failures or maliciously crafted updates preventing the devices from learning a good model.

In recent years, federated learning robustness to Byzantine failures (i.e., worst-case adversarial coordinated training-time attacks) has gained attention. However, existing robust aggregation techniques require sophisticated nonlinear operations [21, 24, 2], sometimes with server access to individual model updates in the clear – thus leading to privacy loss. These nonlinear operations adversely affect privacy since privacy-preserving methods such as secure Multi-Party Computation (MPC) are inefficient for nonlinear operations [3]. This observation highlights a fundamental tension between existing solutions to the two critical problems of privacy and robustness.

1st NeurIPS Workshop on New Frontiers in Federated Learning (NFFL 2021), Virtual Meeting.

To the best of our knowledge, ours is the first approach that scalably combines Byzantine-robustness with privacy using the common single-server architecture. We propose a novel hierarchical framework that decouples MPC-based privacy and Byzantine robustness protection mechanisms in this work. The basic idea is to implement a secure averaging oracle among randomly clustered clients, then filtering these updates using robust aggregation. This approach reveals only the cluster averaged update to the server, thus can help preserve privacy. Simultaneously, the second level of robust aggregation helps to maintain Byzantine robustness.

Taken together, this manuscript proposes a federated learning architecture that preserves security and privacy jointly, thus addressing this gap in the literature. Due to the hierarchical approach, existing robust distributed learning frameworks [21, 24, 2] and non-robust secure distributed learning frameworks like [3, 4] can be considered special cases of our proposed approach. **Summary of contributions:** We propose SHARE; a robust distributed learning framework which flexibly incorporates any Byzantine-robust defenses while enhancing privacy in a single server systems setting. We extend existing theoretical guarantees of robust aggregation oracles to the SHARE framework. Further, we present empirical evaluation of SHARE on benchmark datasets.

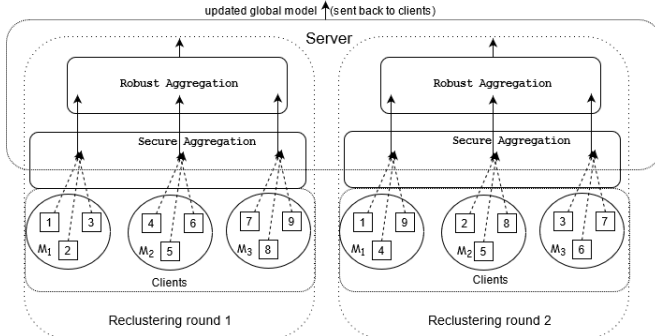


Figure 1: This figure illustrates our proposed framework SHARE. Each global round consists of multiple reclustering rounds, updates from which are averaged to obtain the final model update. In each reclustering round (shown by dotted rectangle), updates from clients (numbered squares) are clustered randomly ( $\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3$ ), then averaged followed by robust aggregation.

## 2 Related Work

Byzantine robustness and secure aggregation in distributed learning both have a large existing literature – though often from different (somewhat disconnected) communities.

**Robust Aggregation.** Robustness to Byzantine adversaries is a well-studied problem in distributed and federated learning [13]. Broadly, existing defenses can be categorized into distance-based robust aggregation or validation data-based aggregation. The general idea behind distance-based metrics is to find an update closer to the benign mean in  $l_2$  norm distance. Xie et al. [21] suggest utilizing coordinate-wise trimmed mean. Blanchard et al. [2] choose a model update closest to most other updates. Ghosh et al. [6] suggest optimal statistical rates utilizing median and trimmed mean. All these distance-based defenses require a majority of the clients participating in the protocol to be benign. On the other hand, validation data-based aggregation defenses such as Zeno [22, 24] perform suspicion-based aggregation based on a score evaluated on validation data held at the server. These methods can tolerate arbitrary Byzantine poisoning. All the above techniques require the server to see the local model updates in the clear, posing a privacy threat.

**Privacy via Secure Aggregation.** In many distributed learning settings, the secure computation boils down to computing a secure average. Bonawitz et al. [3] utilizes a pairwise secret share to achieve the same. Aono et al. [1] follow a slightly different approach and utilize additively homomorphic encryption for secure update computation. While these methods work well with linear aggregation methods, extending secure multiparty computation to non-linear robust aggregation schemes introduces additional computational and communication overhead, quickly becoming impractical for real computational loads.

He et al. [8], Pillutla et al. [16], Wang et al. [20] are the closest to our work in the sense that they are proposed to address the problem of robustness and privacy in distributed learning jointly. Compared to He et al. [8], which requires two non-colluding servers, we achieve this with a single server, which may be a more realistic architecture for practical use cases. Further, He et al. [8] tailor their approach to distance-based robust aggregation. In contrast, our proposed approach is easily combined with most existing Byzantine-robust aggregation schemes, including filtering-based defenses such as Zeno++ [24]. On the other hand, Pillutla et al. [16] reduce the filtering computation of the median into a sequence of linear computations. Unfortunately, this approach requires that the Byzantine device follows the computational protocol over multiple rounds, which is a strong assumption in practice. Thus, while inspired by Byzantine tolerance, Pillutla et al. [16] do not claim Byzantine robustness. After completing this work, we were made aware of a related approach [20] using a hierarchical architecture with a robust mean aggregator. Compared to Wang et al. [20], our approach is a wrapper method that can be combined with any robust aggregator, with analysis and performance depending on the choice of aggregator (e.g., we analyze and compare trimmed mean, krum, Zeno++). Further, we address the problem of signal loss due to clustering via a novel reclustering step.

### 3 Problem Formulation

We consider the optimization problem  $\min_{x \in \mathbb{R}^d} F(x)$  where,  $F(x) = \frac{1}{n} \sum_{i \in [n]} \mathbb{E}_{z_i \sim \mathcal{D}_i} f_i(x; z_i)$ , hence the goal is to learn a model  $x$  which performs well on average using  $z_i$  sampled from local data distribution  $\mathcal{D}_i, \forall i \in [n]$ . The notations used in this paper are summarized in Table 1 (Appendix A).

This problem is solved in a distributed and iterative manner. In each global iteration ( $t < T$ ), sampled clients compute a private model update ( $\Delta x_i^K$ ) by running multiple steps (K-steps) of Stochastic Gradient Descent (SGD) on the local data available ( $z_i \sim \mathcal{D}_i$ ). Then server can compute a global model update. For instance, when using simple averaging, the server update is  $x^t = x^{t-1} + \eta \sum_{i \in [n]} \Delta x_i^K$ , where  $\eta$  is global learning rate. We consider the following privacy and security threats:

- *Privacy threat model:* We consider an honest but curious server. This specification allows the server to interpret the device data from the updates, hence breaching privacy. The assumption of honest server implies that the server still follows the underlying protocol.
- *Robustness threat model:* We consider a fixed (unknown) subset ( $q$ ) of machines that can co-ordinate and send arbitrary updates to the server hence deviating from the intended distributed learning protocol.

## 4 Methodology

We propose two-step hierarchical aggregation SHARE (Secure Hierarchical Robust Aggregation) as a defense against the specified robustness and privacy threat models. In particular, our approach allows a decoupling of the security and robustness into two steps (as illustrated in Figure 1). First, in every global epoch, all participating clients are clustered randomly into groups. Clients within each group share pairwise secret keys and utilize them to mask their individual updates such that the server only learns the average within the cluster. This ensures client update privacy. These client cluster updates are then filtered using Byzantine-robust aggregation techniques. Further, we can repeat this process multiple times in a global epoch to aid in reducing variance. The detailed algorithm is outlined in Algorithm 1. Without loss of generality, we assume clusters of uniform size.

### 4.1 System Components

**Secure Aggregation:** This is the first step in hierarchical aggregation. We follow an approach similar to [3], using pairwise keys between clients in a cluster. The server in this setup learns just the mean and hence the privacy of individual client updates are protected (Detailed discussion in Appendix C).

**Robust Aggregation:** This is the second step in every reclustering round. In this step, the secure cluster averages are filtered through robust aggregation. The goal ideally is to eliminate clusters with malicious client updates. Any existing robustness techniques like trimmed mean[21], median[16] or Zeno[24] can be utilized at this stage. We show theoretical guarantees and experiments based on existing methods in the following sections.

**Random Reclustering:** As specified in Algorithm 1, we repeat the secure aggregation followed by robust aggregation multiple times randomizing client clusters in each global epoch. Note that across

---

**Algorithm 1** SHARE (Secure Hierarchical Robust Aggregation)

---

0: **Server:**  
1: **for**  $t = 0, \dots, T - 1$  **do**  
2:   **for**  $r = 1, \dots, R$  **do**  
3:     Assign clients to clusters  $\mathcal{S} = \mathcal{M}_1 \cup \dots \cup \mathcal{M}_i \dots \cup \mathcal{M}_c$  with  $|\mathcal{M}_i| = |\mathcal{M}_j| \forall i, j \in [c]$   
4:     Compute secure average  $g_j^r \leftarrow \text{SecureAggr}(\{\Delta_i\}_{i \in \mathcal{M}_j}) = \sum_{i \in \mathcal{M}_j} u_i, \forall j \in [c]$   
5:      $g^r \leftarrow \text{RobustAggr}(\{g_j^r\}_{j \in [c]})$   
6:   **end for**  
7:   **if** stopping criteria met **then**  
8:     break  
9:   **end if**  
10:   Push  $x^t = x^{t-1} + \eta \frac{1}{R} \sum_r g^r$  to the clients  
11: **end for**  
12: **Client:**  
13: **for** each client  $i \in \mathcal{S}$  (if honest) **in parallel do**  
14:    $x_{i,0}^t \leftarrow x^t$   
15:   **for**  $k = 0, \dots, K - 1$  **do**  
16:     Compute an unbiased estimate  $g_{i,k}^t$  of  $\nabla f_i(x_{i,k}^t)$   
17:      $x_{i,k+1}^t \leftarrow \text{ClientOptimize}(x_{i,k}^t, g_{i,k}^t, \eta, k)$   
18:   **end for**  
19:    $\Delta_i = \frac{n_i}{n} (x_{i,K}^t - x^t)$   
20:   Push  $\Delta_i$  to the assigned clusters using secure aggregation  
21: **end for**  
22: **return**  $x^T$

these reclustering rounds, the same local model update is paired with different clients each time. In addition to malicious updates, benign updates paired with malicious clients might be filtered in the proposed approach. Reclustering helps mitigate this loss of signal and hence reduces variance. In particular, as number of reclustering rounds ( $R$ ) increase, the probability of this loss in signal decreases (Detailed discussion in Appendix E).

**Remark.** Although reclustering increases communication cost, we note that in addition to helping decrease the variance, reducing secure aggregation to within clusters, decreases communication cost as pairwise key exchange is now limited to within the cluster. Hence overall, communication cost for each client changes from  $\mathcal{O}(n)$  to  $\mathcal{O}(\frac{Rn}{m})$ . In experiments, we often find that even a single clustering round gives good results (Section 6).

## 5 Theory

### 5.1 Exactness

Algorithm 1 can be implemented using any aggregation technique. However, due to clustering, the result is resilient to fewer malicious clients – as (in the worst case) malicious clients are assumed to completely corrupt their assigned cluster. We formalize these ideas next, with proofs in Appendix B.

**Lemma 1.** *If robust aggregation is replaced by averaging, the output of Algorithm 1 is identical to Federated Averaging [15].*

**Lemma 2.** *In presence of robust aggregation, Algorithm 1 is robust to  $q = \frac{q_0}{m}$  adversaries, where  $q_0$  is the tolerance limit of the robust aggregation oracle followed and  $m$  is the cluster size.*

### 5.2 Convergence Analysis

To highlight the flexibility of the proposed algorithm, we analyze convergence when using both using a distance based robust aggregation strategy or a validation data based aggregation strategy, such as Zeno++. We first define the few terms used to develop convergence analysis.

**Definition 5.1** ((G,B)-Bounded Gradient Dissimilarity). There exists constants  $G \geq 0, B \geq 1$  such that  $\frac{1}{n} \sum_i^n \|\nabla f_i(x)\|^2 \leq G^2 + B^2 \|\nabla F(x)\|^2$

**Definition 5.2** (Bounded client updates variance). We define benign mean model update across clients to be  $\mu = \sum_i \Delta_i^K$ , hence the variance across client updates as  $\mathbb{E}[\|\Delta_i^K - \mu\|^2] \leq \sigma_g^2$  for all  $i$  across all rounds of training

**Definition 5.3** (Bounded variance). For an unbiased stochastic gradient estimator with  $g_i(x) = \nabla f_i(x, z_i)$  we define bounded variance as  $\mathbb{E}_{z_i}[\|g_i(x) - \nabla f_i(x)\|] \leq \sigma^2$  for any  $i, x$

The difference between Definition 5.2 and 5.3 is that the former bounds the variance between model updates across clients while the latter bounds the variance across gradient estimates within the same client.

### 5.2.1 Convergence Rates

We now prove that Algorithm 1 converges for various robust aggregation oracles. Firstly, we state a few general assumptions required to prove convergence guarantees standard in papers.

**Assumption 5.1.** There exists at least one global minima  $x^*$  such that  $F(x^*) \leq F(x), \forall x$

**Assumption 5.2.** We assume that  $F(x)$  is  $L$ -smooth and has  $\mu$ -lower bounded Taylor approximation ( $\mu$  weak convexity)

$$\langle \nabla F(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2 \leq F(y) - F(x) \leq \langle \nabla F(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$$

Note that this Assumption 5.2 covers the case of non-convexity by taking  $\mu < 0$ . We note that each distance based robust aggregation metric have different bounds from benign mean update. Since the focus of this work is to propose an algorithm that unifies robustness with privacy, we do not concentrate on those bounds and absorb such intricacies into an order constant. Formally,

**Assumption 5.3.** For any distance based robust aggregation algorithm, *when fraction of faulty inputs is below threshold*, the output of robust aggregation is bounded from benign mean. That is, we assume there exists a  $V_2$  such that for any set of vectors  $\{v_i : i \in \mathcal{C}\}$ , replaced by faulty vectors  $\forall i \notin \mathcal{C}_t \subseteq \mathcal{C}, \|\text{RobustAggr}(\{v_i\}_{i \in \mathcal{C}}) - \frac{1}{|\mathcal{C}_t|} \sum_{i \in \mathcal{C}_t} v_i\| \leq \mathcal{O}(V_2)$ .

We note that Assumptions 5.1,5.2 are standard among existing Federated Learning literature [10, 22, 24]. Additionally Assumption 5.3 is a direct consequence of existing distance based robust aggregation oracles [21, 16, 2]. Finally, for Algorithm 1 with such oracles, we have the following theorem

**Theorem 3.** *Consider a function  $F(x)$  satisfying Assumptions 5.1,5.2 assume a robust aggregation scheme that picks up  $b$  updates and satisfies Assumption 5.3, further, assume  $(G,B)$ -Bounded gradient dissimilarity,  $\sigma_g^2$  variance in client updates and  $\sigma^2$  variance in gradient estimation, there exists  $\eta, \eta_l$  such that output of Algorithm 1 after  $T$  rounds,  $x^T$ , satisfies,*

$$\mathbb{E}[\|\nabla F(x^T)\|^2] \leq \mathcal{O}\left(\frac{LM\sqrt{F}}{\sqrt{TKn}} + \frac{F^{2/3}(LG^2)^{1/3}}{(T+1)^{2/3}} + \frac{B^2LF}{T} + 2L^2V_2 + \frac{\sigma_g^2}{bm} \left(\frac{n-q-bm}{R(n-q)-1}\right)\right)$$

where  $M^2 := \sigma^2(1 + \frac{n}{\eta^2})$  and  $F := F(x^0) - F(x^*)$

Now we consider Zeno++[24], a defense utilizing server data. Although score based Zeno++ was originally introduced for asynchronous SGD, we generalize it to federated learning setting hence allowing for multiple local epochs. We illustrate this modified algorithm in Appendix B. As in Xie et al. [24], we consider an additional standard assumption

**Assumption 5.4.** The validation set considered for Zeno++ is close to training set, implying a bounded variance given by  $\mathbb{E}[\|\nabla f_s(x) - \nabla F(x)\|^2] \leq V_1, \forall x$

**Theorem 4.** *Consider  $L$ -smooth and potentially non-convex functions  $F(x)$  and  $f_s(x)$ , satisfying Assumption 5.4. Assume  $\|f_s(x)\|^2 \leq V_3, \forall x$ . Further assuming  $G$ -bounded gradient dissimilarity, variance between client updates be  $\sigma_g^2$  and variance in gradient estimation at each client be  $\sigma$ , with global and local learning rates of  $\eta \leq \frac{1}{2L}$  and  $\rho \geq \frac{\alpha\sqrt{\eta}}{6K^2\eta^2B^2} + \eta$ , after  $T$  global updates, let  $D := F(x^0) - F(x^*)$ , Algorithm 1 with Zeno++ as robust aggregation converges at a critical point:*

$$\frac{\mathbb{E}[\sum_{t \in [T]} \|\nabla F(x_{t-1})\|^2]}{T} \leq \frac{\mathbb{E}[D]}{\alpha\sqrt{\eta}T} + \frac{\sqrt{\eta}}{\alpha} \mathcal{O}\left(\frac{\sigma_g^2}{m} \left(\frac{n-q-m}{R(n-q)-1}\right) + G^2 + \sigma^2 + V_1 + V_3 + \epsilon\right)$$

**Remark.** It can be seen from both Theorem 3,4 that the additional terms, other than standard ones appearing in the convergence rate for federated learning [10], depend on the error caused by the robust aggregation scheme utilized and variance reduction from reclustering. Further, higher number of reclustering rounds  $R$  decreases the effect of additional variance. Finally when  $R = 1, m = 1, q = 0$ , these recover existing results for federated learning with robust aggregation.

### 5.3 Privacy

**Curious server:** Since each client masks updates with random vectors as illustrated in Section 4, we note that if we execute the mentioned secure averaging oracle with threshold  $t > \frac{m}{2}$ , the protocol can deal with  $\lceil \frac{m}{2} \rceil - 1$  drop outs while learning nothing more than average. Reclustering introduces additional vulnerability as server can see multiple averages. In particular, the probability that server can decode a model update is  $\mathcal{O}(1 - \left(\frac{(m!)^c}{n!}\right)^R)$ . Hence as R increases this gets closer to 1 as expected. Further, when all clients are in a single cluster ( $m = n$ , hence  $c=1$ ), this is 0 as would be the case with secure averaging without robustness. Further discussion, including comments on privacy in presence of colluding curious clients can be found in Appendix C.

## 6 Experiments

In this section we evaluate the proposed algorithm SHARE with various defenses and corruption models. We conduct experiments on CIFAR-10 [12] (Image classification dataset) and Shakespeare (a language modeling dataset from LEAF [5]). We note that we do not propose a new robustness technique but rather we propose a modified federated learning architecture to incorporate any robustness protocol in a privacy preserving manner. Hence we focus our experiments on capturing the effects of cluster sizes and reclustering rounds, hyperparameters introduced by our approach. We defer descriptions of detailed training architecture to Appendix D

### 6.1 CIFAR-10

We train a CNN with two  $5 \times 5$  convolutional layers followed by 2 fully connected layers [15] on CIFAR-10 and report top-1 accuracy. We test SHARE incorporating various robust aggregation protocols such as Trimmed mean [21], Krum [2], Zeno++ [24]. For all experiments in this section, trimmed mean removes 2/3 of the updates before computing the mean. Additionally, we consider two baselines, SHARE with no robust aggregation and SHARE with no attack. We consider homogeneous distribution of data across clients for experiments in this section. Experiments on heterogeneous data distributions can be found in Appendix D.

#### 6.1.1 Impact of cluster size

We first test Byzantine-tolerance for various cluster sizes to mild attacks such as label-flip. In particular, malicious clients train on wrong labels (images whose labels are flipped, i.e., any label  $\in \{0, \dots, 9\}$  is changed to 9-label). We consider 60 total clients of which  $q = 12$  being malicious. The result is shown in Figure 2 for various cluster sizes and robust aggregation protocols. It is seen

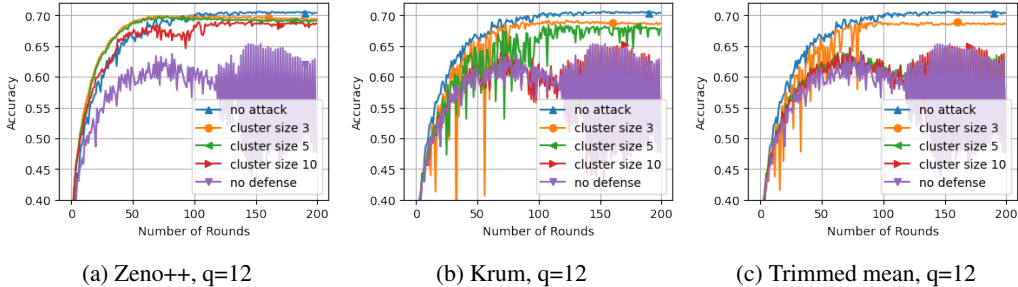


Figure 2: Results of SHARE with various defenses on CIFAR-10, utilizing varying cluster sizes under label flip attack. For Trimmed mean we remove 2/3 of input updates.

that having no defense diverges even with mild attacks as expected. Further Figure (2a) shows that SHARE with a strong defense like Zeno++ converges to benign (no-attack) accuracy for any of the considered cluster sizes. SHARE with trimmed mean and Krum both converge with cluster size 3 but as cluster size increases, accuracy decreases and SHARE begins to diverge. This can be seen directly from Lemma 2, since we set trimmed mean to filter  $q_0 = (2/3) * 60 = 40$  of updates, a cluster size of 3 implies the algorithm is robust against  $q = 40/3 > 12$  clients being malicious, hence the algorithm converges to benign accuracy, increasing the cluster size decreases this tolerance threshold and hence as shown in Figure (2c) may fail to converge. Further experiments on scaled sign-flip attacks, are included in Appendix D due to space constraints.

### 6.1.2 Impact of reclustering

Intuitively, increasing the number of reclustering rounds increases the expected number of clusters without a Byzantine client. This hence increases the robustness of SHARE to higher fraction of Byzantine clients with defenses like Zeno++ which can tolerate arbitrary levels of poisoning. We test this hypothesis with several attack and clustering scenarios using sign-flipping attacks. In (a), we

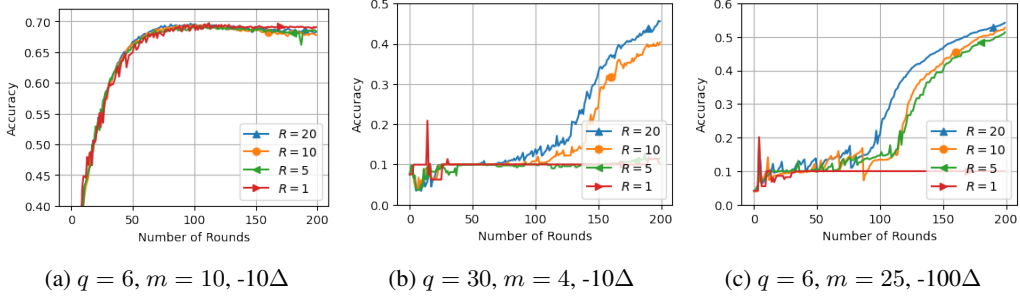


Figure 3: Results of SHARE with Zeno++ and different reclustering rounds  $R$ , Byzantine clients  $q$ , cluster sizes  $m$  on CIFAR-10 with varying attack strengths (Any benign model update  $\Delta$  is scaled to either  $-10\Delta$  or  $-100\Delta$ ). In (a),(b) we use  $n = 60$  and for (c) we use  $n = 100$ .

use a relatively small cluster size and a low fraction of Byzantine clients, so 1 round is sufficient. In (b), the fraction of Byzantine clients is high and in (c) the cluster size is large, which increases the probability of a cluster containing a Byzantine client, so  $R > 1$  helps converge to higher accuracies.

## 6.2 Shakespeare

We consider the first 60 speaking roles in the train set as our 60 clients. We train an RNN with 2 LSTM layers followed by 1 fully connected layer [17] and report top-1 accuracy on the testing set.

### 6.2.1 Empirical Evaluation

In Figure 4 we evaluate Byzantine tolerance of SHARE with Zeno++ under sign-flip attack (malicious clients send an update negative to the benign one  $-\Delta$ ) and scaled sign-flip attack (malicious clients scale the update in addition to flipping the sign and hence send  $-10\Delta$ ). A stronger attack like scaled sign-flip breaks benign averaging and Zeno++ works well with any of the chosen cluster sizes.

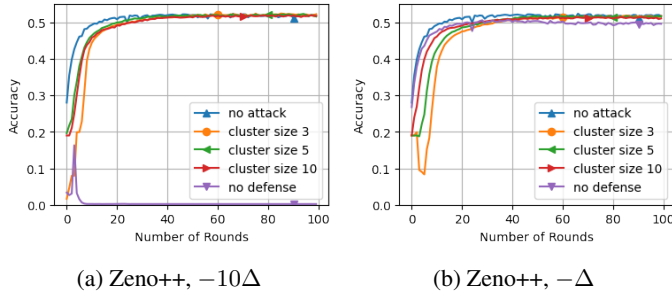


Figure 4: SHARE with Zeno++ defense and sign-flipping attack on Shakespeare.

## 7 Discussion and Conclusion

We have proposed SHARE, a framework for implementing Byzantine-robustness and privacy. The key idea is hierarchical clustering. Cluster size is an important parameter that controls the trade-off between privacy and robustness. Further, reclustering is an important step and can help decrease variance and increase tolerance to the fraction of malicious clients when the defense can support arbitrary failures like Zeno++. In future, we would like to explore other variations in client clustering, especially in heterogeneous data settings. Further, we plan to work on stronger security guarantees even with multiple reclustering rounds.

## Acknowledgments and Disclosure of Funding

Koyejo acknowledges partial funding from a C3.ai Digital Transformation Institute Award and a Jump Arches Award. This work was also funded in part by NSF III 2046795 and IIS 1909577. Additionally, the authors like to acknowledge Microsoft Azure for computational resources. Finally we would like to thank Dakshita Khurana and Nishant Kumar for their insightful discussions.

## References

- [1] Y. Aono, T. Hayashi, L. Wang, S. Moriai, et al. Privacy-preserving deep learning via additively homomorphic encryption. *IEEE Transactions on Information Forensics and Security*, 13(5): 1333–1345, 2017.
- [2] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 118–128, 2017.
- [3] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth. Practical secure aggregation for privacy-preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1175–1191, 2017.
- [4] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, H. B. McMahan, et al. Towards federated learning at scale: System design. *arXiv preprint arXiv:1902.01046*, 2019.
- [5] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.
- [6] A. Ghosh, J. Hong, D. Yin, and K. Ramchandran. Robust federated learning in a heterogeneous environment. *arXiv preprint arXiv:1906.06629*, 2019.
- [7] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- [8] L. He, S. P. Karimireddy, and M. Jaggi. Secure byzantine-robust machine learning. *arXiv preprint arXiv:2006.04747*, 2020.
- [9] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- [10] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- [11] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [12] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [13] L. Lamport, R. Shostak, and M. Pease. The byzantine generals problem. In *Concurrency: the Works of Leslie Lamport*, pages 203–226. 2019.
- [14] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.
- [15] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [16] K. Pillutla, S. M. Kakade, and Z. Harchaoui. Robust aggregation for federated learning. *arXiv preprint arXiv:1912.13445*, 2019.



- [17] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
- [18] J. A. Rice. *Mathematical statistics and data analysis*. Cengage Learning, 2006.
- [19] A. Shamir. How to share a secret. *Communications of the ACM*, 22(11):612–613, 1979.
- [20] L. Wang, Q. Pang, S. Wang, and D. Song. Towards bidirectional protection in federated learning, 2021.
- [21] C. Xie, O. Koyejo, and I. Gupta. Slsgd: Secure and efficient distributed on-device machine learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 213–228. Springer, 2019.
- [22] C. Xie, S. Koyejo, and I. Gupta. Zeno: Distributed stochastic gradient descent with suspicion-based fault-tolerance. In *International Conference on Machine Learning*, pages 6893–6901. PMLR, 2019.
- [23] C. Xie, O. Koyejo, and I. Gupta. Fall of empires: Breaking byzantine-tolerant sgd by inner product manipulation. In R. P. Adams and V. Gogate, editors, *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine Learning Research*, pages 261–270. PMLR, 22–25 Jul 2020. URL <http://proceedings.mlr.press/v115/xie20a.html>.
- [24] C. Xie, S. Koyejo, and I. Gupta. Zeno++: Robust fully asynchronous sgd. In *International Conference on Machine Learning*, pages 10495–10503. PMLR, 2020.