# Client Clustering Meets Knowledge Sharing: Enhancing Privacy and Robustness in Personalized Peer-to-Peer Learning

1st Mohammad M Maheri
Imperial College London
London, UK
m.maheri23@imperial.ac.uk

2nd Denys Herasymuk
Imperial College London
London, UK
d.herasymuk24@imperial.ac.uk

3rd Hamed Haddadi
Imperial College London
London, UK
h.haddadi@imperial.ac.uk

*Abstract*—The growing adoption of Artificial Intelligence (AI) in Internet of Things (IoT) ecosystems has intensified the need for personalized learning methods that can operate efficiently and privately across heterogeneous, resource-constrained devices. However, enabling effective personalized learning in decentralized settings introduces several challenges, including efficient knowledge transfer between clients, protection of data privacy, and resilience against poisoning attacks. In this paper, we address these challenges by developing P4 (Personalized, Private, Peer-to-Peer)—a method designed to deliver personalized models for resource-constrained IoT devices while ensuring differential privacy and robustness against poisoning attacks. Our solution employs a lightweight, fully decentralized algorithm to privately detect client similarity and form collaborative groups. Within each group, clients leverage differentially private knowledge distillation to co-train their models, maintaining high accuracy while ensuring robustness to the presence of malicious clients. We evaluate P4 on popular benchmark datasets using both linear and CNN-based architectures across various heterogeneity settings and attack scenarios. Experimental results show that P4 achieves 5% to 30% higher accuracy than leading differentially private peer-to-peer approaches and maintains robustness with up to 30% malicious clients. Additionally, we demonstrate its practicality by deploying it on resource-constrained devices, where collaborative training between two clients adds only $\approx$ 7 seconds of overhead.

*Index Terms*—peer-to-peer machine learning, decentralized learning, personalization, differential privacy, poisoning attacks.

## I. INTRODUCTION

The integration of Artificial Intelligence (AI) into Internet of Things (IoT) ecosystems is reshaping how connected devices operate by enabling autonomous decision-making and real-time data analysis at scale. In healthcare, for example, AI empowers clinical staff to continuously monitor patients' physiological data throughout the hospital, detect early signs of deterioration, and issue timely alerts—ultimately helping to prevent adverse outcomes. The diverse nature of data—for example, significant differences in electrocardiograms resulting from varying levels of patient health—collected by different devices is known as *client data heterogeneity*. This heterogeneity often causes a single global model to perform poorly across varying data distributions since each client optimizes towards its own distinct local empirical minimum, leading to divergent local optimization directions [1], [2], [3], [4]. Recently, growing interest has emerged in addressing data heterogeneity using *personalized learning* [5], [6], [7], which adapts each client's model to its specific data or task while still utilizing shared knowledge across clients.

Fully decentralized or *peer-to-peer* (P2P) learning has emerged as a promising direction for enabling personalized learning in heterogeneous environments such as those found in IoT ecosystems. In contrast to centralized federated learning (FL), P2P learning eliminates the need for a central server and distributes computation among participating clients [8], [5]. This approach leverages direct communication between clients with their own data distributions, offering potential benefits such as faster convergence [9], [5] and enhanced data privacy compared to centralized FL shi2023improving[10], [11]. In the context of IoT, where devices often operate under bandwidth and power constraints, P2P learning allows edge devices to collaboratively train models by communicating only with a limited number of neighbors, depending on their communication and computational capabilities [12], [13].

*Personalized decentralized learning* faces three key challenges. The first is *clustering*, where clients seek similar peers to share knowledge. Clustering algorithms [1], [14] group clients with similar data distributions to train localized models, improving accuracy. However, most methods rely on predefined cluster counts and require a central server [15], [16], [17], [18], which is impractical due to communication bandwidth constraints. More critically, these approaches risk privacy leaks as clients share data for clustering. The second challenge is *data privacy*—clients' data often contains sensitive information, requiring strong privacy guarantees. A common solution is adding differential privacy (DP) noise and clipping gradients [19], [20], [21], [22], but this degrades model performance [23], [24], [20]. In P2P settings, this effect is amplified due to fewer model updates (gradients) that are used in model aggregation [25]. Finally, the third challenge is the presence of *malicious clients* affected by data poisoning [26] or model poisoning [27], [28], [29], [30] attacks, which can degrade the final model performance. Due to its fully decentralized nature, P2P learning is highly vulnerable to poisoning attacks, in which malicious clients can exploit the

system by sending carefully crafted local models to their neighboring peers. A common mitigation method is to use recent decentralized FL (DFL) defenses [31], [32], [33] or adapt existing centralized FL ones [34], [35], [36], [37] to DFL. However, a defense that is not tailored to personalized P2P settings may fail in scenarios with highly heterogeneous clients and can introduce additional computational overhead if each client performs robust aggregation individually [38].

Although recent efforts in the FL community have introduced privacy into the P2P setting [25], [39], [40], [41], to the best of our knowledge, no existing fully decentralized framework explicitly addresses the challenge of ensuring differential privacy under heterogeneous client data distributions typical of IoT environments. There is a growing trend of applying differential privacy to improve the privacy-utility trade-off in the P2P setting [25], [39]. However, these methods are less effective when clients have highly heterogeneous data (see Section V-B) and typically assume that clients are only curious about other clients' data while honestly sharing their local model gradients. While recent frameworks [42], [43], [40], [41], [44] make progress in integrating privacy preservation and resilience to data poisoning in P2P settings, they suffer from several limitations. Most notably, they neglect client heterogeneity, offer limited improvements in the privacy-utility trade-off, and often fail to address model poisoning attacks, which are particularly threatening in DFL [31]. However, in scenarios with highly data-heterogeneous clients where personalization is important, incorporating the similarity of the client's data distribution into the training procedure can improve both privacy-utility amplification [45], [16] and robustness against poisoning attacks [46].

In this paper, we propose P4 (Personalized, Private, Peer-to-Peer)[1], a method designed to deliver personalized models for resource-constrained IoT devices while ensuring differential privacy and robustness against poisoning attacks. In P4, clients first group themselves in a fully decentralized manner using a similarity metric based on their model weights. Once grouped, clients perform collaborative training with upper-bounded data privacy risk and robust aggregation against adversarial model updates. To enhance privacy-utility trade-offs, P4 incorporates the proxy model [47] to facilitate knowledge distillation among clients. This separation between the local and proxy models results in fast convergence and personalized parameter adaptation for each client (as shown in Section V-B), while also enhancing resilience against data poisoning attacks (as discussed in Section IV-E). To mitigate model poisoning attacks, P4 employs a combined defense strategy based on `m-Krum` [36] and `anomaly-detection` [46], which utilizes client similarity within each group to filter out malicious updates. P4 is well-suited for IoT deployments that involve numerous clients and non-stationary data distributions, where communication and computational efficiency are critical [48].

Our main contributions are summarized as follows:

- We design P4 to address both data heterogeneity and poisoning threats in a differential private P2P setting tailored for IoT environments. P4 enables clients with similar data distributions to collaboratively train, leveraging KL divergence-based regularization between local and proxy models to facilitate knowledge transfer. Its defense mechanism uniquely incorporates knowledge distillation, anomaly detection, and `m-Krum` to improve robustness.
- We propose a lightweight and decentralized procedure that allows clients to privately group themselves for collaborative training, using the $\ell_1$-norm between their model weights as a similarity metric.
- We conduct a comprehensive empirical study of personalized P2P learning under differential privacy and adversarial attacks, covering a range of privacy budgets, model architectures, data heterogeneity types, and poisoning scenarios. Our results show that P4 consistently outperforms state-of-the-art baselines by 5%–30% in accuracy and tolerates up to 30% malicious clients, demonstrating both strong privacy-utility trade-offs and robustness.
- To evaluate the feasibility of deploying P4 in real-world IoT settings, we implement it on resource-constrained edge devices using Raspberry Pis. Our results show minimal overhead in runtime, memory usage, power consumption, and communication bandwidth, with collaborative training between two clients adding only $\approx 7$ seconds of overhead.

## II. BACKGROUND AND RELATED WORK

### A. Private Decentralized Learning

In decentralized learning, preserving both data and local model privacy is essential, as deep neural networks can unintentionally memorize training data, making them vulnerable to model inversion attacks [49], [50]. To address this, several privacy-preserving techniques have been developed to protect plaintext gradients. Differential privacy has been applied to mitigate privacy risks [19], [20], [21], [22], but the added noise can degrade model performance [51]. Homomorphic encryption offers an alternative by encrypting model weights or gradients [52], [53], [54], though it introduces significant computational and communication overhead [55], limiting its practicality for resource-constrained edge devices. Secure multi-party computation (MPC) techniques [56], [57], [58] preserve privacy through collaborative computation among multiple parties, while keeping their individual inputs confidential, but are similarly hindered by high communication costs. More recently, trusted execution environments (TEEs) have been explored [59], [60], [61], though their reliance on specialized hardware and associated computational overhead restrict widespread adoption.

### B. Personalized Decentralized Learning

The presence of non-IID client data necessitates personalization in decentralized learning, shifting focus from a single global model to customized models for each client. [7] used regularization to keep personalized models close

to the global optimum, while [62] proposed computing an optimal weighted combination of client models for personalized updates. Another line of work clusters clients to improve performance [45], [16], with aggregation performed within clusters to avoid negative transfer that may result from merging dissimilar models [1]. A key challenge in federated clustering is measuring client similarity, typically based on losses [45], [1], gradients [16], [14], or model weights [15], [16], [17], [18]. However, aside from [1], these methods rely on a centralized server, limiting their applicability in P2P settings. Moreover, most approaches, including [1], overlook privacy risks during similarity computation. Loss-based methods, in particular, require each client to receive other clients' data and compute the loss of its data on the received model to compute its similarity with other clients, incurring high communication and computation overhead. Moreover, most approaches assume prior knowledge of the number of client clusters, leading to degraded clustering performance if this number is mis-estimated.

### C. P2P Learning

Several studies have explored decentralized federated learning and proposed techniques to enhance its performance. Some employed Bayesian methods to model shared knowledge among clients in a decentralized setting [63], [12], while [64] introduced a server-less federated learning framework designed for dynamic environments. To reduce communication and computation overhead, [9] applied quantization techniques, and [65] leveraged sparse training methods to achieve similar efficiency gains.

To address the challenge of heterogeneous data distributions in P2P learning, several personalization strategies have been proposed. [65] used personalized sparse masks to tailor models to individual clients. [66] introduced a decentralized partial model training approach using the Sharpness Aware Minimization (SAM) optimizer [67] to mitigate model inconsistency. This was further improved by combining SAM with Multiple Gossip Steps (MGS) to address both inconsistency and overfitting [10]. [68] tackled data heterogeneity by employing weighted aggregation of model parameters based on the Wasserstein distance between output logits of neighboring clients.

However, most existing peer-to-peer learning approaches rely on random or predefined communication topologies, without evaluating whether the selected links are actually beneficial for collaboration [66], [65], [10], [25], [69]. While [68] attempted to build a dynamic topology based on client similarity, their method only considers output logit similarity and does not account for the underlying data distribution. This highlights the need for more intelligent and adaptive strategies for link selection that can improve collaboration effectiveness and accelerate convergence in environments with heterogeneous client data.

Furthermore, while differential privacy has been widely explored in centralized FL [24], a significant gap remains in addressing both data and model privacy in P2P learning settings [66], [70], [65], [68], [10]. Ensuring the privacy of client data is essential for maintaining trust and security in decentralized frameworks. Several approaches have been proposed to enhance privacy in peer-to-peer learning. [71] investigated the trade-off between utility and privacy in peer-to-peer FL. [25] introduced proxy models to enable efficient communication without a central server and incorporated differential privacy analysis to strengthen privacy guarantees. However, the added noise from differential privacy can degrade model performance, particularly when client data is non-IID, and merging models from divergent distributions may further harm accuracy. More recently, [39] extended DP-SGD to decentralized learning, aiming to reach consensus on model parameters across clients. Yet, this consensus model may underperform in non-IID scenarios where personalized models are more appropriate. Another direction, explored by [69], applied gradient encryption to protect privacy and used succinct proofs to verify gradient correctness. Despite its security benefits, this approach introduced computational and communication overhead [72] and did not address data heterogeneity. These limitations highlight the need for future research to focus on privacy-preserving mechanisms that not only safeguard sensitive information but also account for the heterogeneity and personalization needs inherent in P2P learning.

Several recent works [42], [43], [40], [41], [44] have proposed frameworks for P2P learning that assume a stronger threat model, addressing both privacy and malicious security concerns. For instance, Biscotti [40] is a fully decentralized P2P system for privacy-preserving and secure federated learning, leveraging blockchain, differential privacy, and m-Krum [36] for secure aggregation. However, this framework is limited to IID data and does not address client heterogeneity in the context of privacy preservation and robustness against poisoning attacks. BlockDFL [43] is another ledger-based system that combines gradient compression and median-based testing with m-Krum to adapt privacy mechanisms for data-heterogeneous clients and improve robustness against data poisoning. Nevertheless, it considers only a label-flipping attack [26], which is regarded as one of the least challenging poisoning attacks [73]. Moreover, their evaluation of malicious security is restricted only to IID clients. BEAS [42] is another blockchain-based framework that integrates differential privacy, gradient pruning, and Foolsgold [37] with m-Krum to defend against both privacy and security threats. However, it lacks comparative analysis of its privacy-preserving approach and, like BlockDFL, evaluates performance solely under a label-flipping attack. These limitations reveal a significant gap in DFL frameworks: the lack of solutions that simultaneously support data-heterogeneous clients and operate under a strong threat model encompassing both privacy preservation and protection against data poisoning and model poisoning attacks.

### III. Problem Statement

#### A. System Model

We consider $M$ clients $\{c_i\}_{i=1}^{M}$. Each client, $c_i$, holds a local dataset including data points and corresponding labels

$\{X_i, Y_i\} \sim D_i$ drawn from its personal distribution $D_i$. Each $c_i$ aims to learn personalized parameters $w_i$ to minimize the expected loss over the client's data distribution:

$$F_i(w_i) = \mathbb{E}[\mathcal{L}_{(x,y)\sim D_i}(w_i; (x, y))] \tag{1}$$

To find parameter to minimize the expected loss, each client seeks to determine an optimal parameter set, denoted as $w_i^*$, within the parameter space $w_i \in \mathbb{R}^d$. This optimization aims to minimize the expected loss over the client's own local data distribution $\ell(f_i, y)$, represented in Equation 2:

$$w_i^* = arg \min_{w_i \in \mathbb{R}^d} \mathbb{E}_{(x,y)\sim D_i} \ell(f_i(w_i; x), y)] \tag{2}$$

However, solely using client's local dataset is not sufficient to find a generalizable parameter for each client, so local training leads to poor generalization performance [5], [7]. In order to achieve good generalization, each client is willing to collaborate and share knowledge (i.e., their model gradients) with a *subset* of other clients. We note that the goal of the learning is not to reach a "global consensus" model, but for each client to obtain personalized models that perform well on their own data distributions.

### B. Threat Model

We consider both privacy and security threats in a P2P setting without a trusted third party or a centralized server.

From a privacy perspective, we assume *semi-honest* clients, as in prior work [42], [40], who may attempt to infer sensitive information from other clients' private training datasets, but they do not deviate from the training protocol. We assume that clients communicate exclusively over encrypted channels, preventing third-party adversaries from eavesdropping or tampering with the communication. Additionally, clients do not share data directly with others in the network. To mitigate the risk of sensitive information leakage through model updates and to ensure a differential privacy guarantee, each client is assigned a fixed privacy budget in advance.

From a security standpoint, we consider adversaries in a *semi-honest* setting, where clients can deviate from the protocol specification by contributing malicious updates designed to degrade the performance or integrity of the final model. In particular, our threat model assumes the adversary is capable of orchestrating Sybil attacks [74], thereby controlling multiple peers within the network. However, consistent with prior literature [40], [41], we constrain the adversary's stake to no more than 30% of the total participants in the network. Specifically, our threat model includes prevalent attack scenarios investigated in recent privacy-preserving and peer-to-peer FL frameworks [42], [40], [41], [43]. These include *label flipping* [26], *byzantine-zero* [27], *byzantine-random* [27], and *byzantine-flip* [28]. These attacks collectively represent critical and realistic threats faced by distributed machine learning protocols, particularly in peer-to-peer and privacy-preserving contexts. In our *label flipping* setting, client data is poisoned by flipping *all labels* as in [26]. Byzantine attacks corrupt local model updates by setting weights to zero, random values,

or flipped values using $w_g + (w_g - w_l)$, where $w_g$ is the global model, and $w_l$ is the real local model. Additionally, we assume an attacker may have knowledge of the system design, particularly the group formation phase (see Section IV-B).

## IV. THE DESIGN OF P4

In this section, we describe the design of P4 in detail. Section IV-A outlines key challenges and provides an overview of the framework. Section IV-B describes how clients form groups based on model weight similarity, while Section IV-C explains the collaborative training algorithm inside each group based on local and proxy models with differential privacy. Next, Section IV-D defines the privacy goals and formal guarantees. Finally, Section IV-E introduces a combined defense against poisoning attacks based on knowledge distillation, anomaly detection, and m-Krum.

### A. Overview

P4 aims to provide a personalized model for each client under DP guarantees and robustness to poisoning attacks, while remaining suitable for deployment on resource-constrained IoT devices. Based on the problem statement in Section III, we identify four key challenges.

The first challenge is achieving effective client personalization while still allowing clients to benefit from shared knowledge in a system without a central server. This requires a clustering mechanism to group clients with similar data distributions. However, enabling such collaboration introduces the second challenge: maintaining data privacy. Any knowledge exchanged between clients to enhance local model training must not violate differential privacy. Following prior work, this can be achieved by adding Gaussian noise to shared gradients, but this comes at the cost of performance degradation. In a fully decentralized setting, fewer clients can participate in collaborative training compared to centralized federated learning, making the impact of noise more severe, as demonstrated in Section V-B. The third challenge stems from the limited computational and communication resources available to clients. The learning algorithm must converge efficiently within a limited number of iterations, and the knowledge-sharing method should minimize both computation and communication overhead. Finally, malicious clients affected by poisoning attacks can degrade clients' model performance. The system must incorporate an effective defense mechanism that aligns with privacy and personalization requirements without introducing excessive computational costs.

Figure 1 illustrates the overall design of the P4 approach, which addresses these challenges. It operates in two phases: *group formation* and *co-training*. During group formation, clients privately identify peers with similar data distributions and form groups accordingly. In the co-training phase, clients train their *local* and *proxy* models and share locally computed data gradients with added noise only within their group. A group aggregator then filters out adversarial model updates, aggregates the remaining gradients, and sends the resulting global model update back to the clients. The following sections
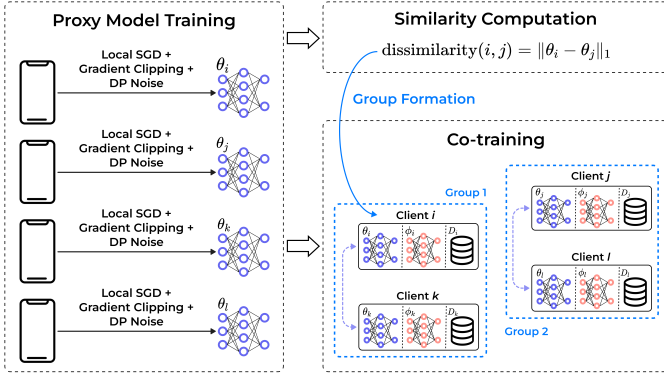
Fig. 1. The overall design of the P4 approach. Clients employ *local* ($\phi$) and *proxy* ($\theta$) models, aggregating only the *proxy* model.

provide a more detailed explanation of the system. A table with all the notation is provided in Appendix A

### B. Phase 1: Group Formation

In federated learning, similar data distributions among clients plays an especially important role for model convergence [75], [76], [77]. Selecting groups of clients for co-training without considering the underlying data heterogeneity—such as through random clustering—can lead to performance worse than standalone local training [78], [1]. Therefore, identifying clients with similar data distributions is critical for effective collaborative training.

A direct measure of data distribution similarity across decentralized clients can be challenging due to privacy constraints. Inspired by prior works [45], [16], we use the similarity of client model parameters as an effective proxy. Specifically, we measure dissimilarity using the $\ell_1$-norm between model weights:

$$\text{dissimilarity}(i,j) = \|\mathbf{w}_i - \mathbf{w}_j\|_1. \qquad (3)$$

The intuition behind using model weight similarity as a proxy stems from the observation that similar model parameters imply similar gradient behaviors. Specifically, consider two clients, $c_i$ and $c_j$, each with personalized model parameters $w_i$ and $w_j$, respectively. The gradient of the loss function for client $c_j$, evaluated at client $c_i$'s parameters, can be approximated using a first-order Taylor expansion around $w_j$:

$$\nabla_{w_i} F_j(w_i) = \nabla_{w_j} F_j(w_j) + \nabla^2_{w_j} F_j(w_j)(w_i - w_j)$$
$$+ O(|w_i - w_j|^2). \qquad (4)$$

The Hessian matrix $\nabla^2_{w_j} F_j(w_j)$ characterizes the curvature of the loss landscape around $w_j$, assigning importance to each parameter dimension. Thus, the approximation error incurred by client $c_i$ when leveraging gradients from client $c_j$ is closely related to the parameter distance:

$$\text{Error} \approx \frac{1}{2}(w_i - w_j)^T \nabla^2_{w_j} F_j(w_j)(w_i - w_j). \qquad (5)$$

Minimizing the distance between model parameters directly reduces the approximation error, thereby improving the accuracy and effectiveness of collaborative training. To ensure

computational efficiency—crucial for resource-constrained decentralized IoT environments—our method avoids explicit Hessian computation and instead focuses on minimizing model parameter differences. This theoretical foundation motivates our client similarity metric and underscores the importance of parameter proximity for effective collaborative training and personalization in decentralized settings.

To effectively capture meaningful model differences reflecting local data distributions, each client initially performs minimal local training—specifically, one epoch over its local dataset—to ensure parameter vectors reflect local data characteristics. This local training phase incorporates gradient clipping and differential privacy noise, as detailed in Section IV-D. Moreover, since models are only partially trained and have not converged, they primarily encode coarse-grained information about the general direction of the data rather than specific details, which significantly reduces vulnerability to attacks such as gradient leakage or data reconstruction. Empirical results (Section V-B) and theoretical analysis (Appendix F) further confirm that even this minimal training is sufficient to distinguish between substantially different data distributions without incurring major privacy risks.

Clients are grouped based on a dissimilarity metric. Let $M$ denote the total number of clients in the network. The grouping structure is represented by a binary matrix $G \in {0,1}^{M \times M}$, where $G_{ij} = 1$ indicates that clients $i$ and $j$ belong to the same group, and $G_{ij} = 0$ otherwise. The grouping objective is then defined as follows:

$$\min_G \sum_{i=1}^{M} \sum_{j=1}^{M} \mathbb{I}[G(i,j)=1]\|\mathbf{w}_i - \mathbf{w}_j\|_1 \qquad (6)$$

such that:

$$\sum_{j=1}^{M} G_{i,j} = V \quad \text{for } i = 1, 2, \dots, M \qquad (7)$$

where $G$ is the collaboration graph (symmetric matrix) and $V$ defines a group size (due to limited computation and communication of each client). We note that P4 enforces mutual collaborative training, implying that communication links between clients are bidirectional, and the matrix $G$ should be symmetric.

To construct the collaboration graph, $G$, in a decentralized manner, we employ a greedy approach to identify the $V$ most similar models for each client (see Algorithm 1). First, each client computes its similarity with $H$ randomly selected clients ($H \ll V$) using the model similarity metric in Equation 3 (lines 1-7). Next, if two clients mutually select each other as their most similar among the $H$ samples, they form a two-member group (lines 8-16). Unassigned clients, such as client $t$, join their most similar ungrouped client $k$, regardless of whether $t$ is also $k$'s most similar sample (line 17). Any remaining ungrouped clients are paired randomly, ensuring all clients are part of exclusive two-member groups (line 18). Within each group, clients share their computed similarities. The group then estimates its similarity with other groups as the maximum

**Algorithm 1:** Client Grouping Based on L1-Norm Similarity

**Input:** Set of clients $M$, group size $V$, random subset size $H$.
**Output:** A binary matrix $G$ with client grouping based on $\ell_1$-norm.

```
// Step 1: Compute l1-norm similarity
   among clients
```
1   $mcs \leftarrow \emptyset$ // Dictionary: client $\rightarrow$ list of similarity tuples
2   **foreach** *client* $i \in M$ **do**
3      **foreach** *client* $j \in$ *random subset* $S \subset M \setminus \{i\}, |S| = H$ **do**
4         **if** $(i, j) \notin mcs$ **then**
5            $mcs[i] \leftarrow mcs[i] \cup \{(j, \text{dissimilarity}(\mathbf{x}_i, \mathbf{x}_j))\}$
6            $mcs[j] \leftarrow mcs[j] \cup \{(i, \text{dissimilarity}(\mathbf{x}_j, \mathbf{x}_i))\}$

7   Sort each $mcs[i]$ in descending order of similarity
```
// Step 2: Create pairs of mutual similar
   clients
```
8   $similar\_clients[i] \leftarrow \emptyset$ for all $i \in M$
9   **foreach** *client* $i \in M$ **do**
10    **if** $similar\_clients[i] \neq \emptyset$ **then** Continue
11    $k \leftarrow mcs[i][0][0]$      // Most similar client
12    **if** $i = mcs[k][0][0]$ *and* $similar\_clients[k] = \emptyset$ **then**
13      $similar\_clients[i] \leftarrow similar\_clients[i] \cup \{k\}$
14      $similar\_clients[k] \leftarrow similar\_clients[k] \cup \{i\}$
15      Remove $(k, \cdot)$ from $mcs[i]$
16      Remove $(i, \cdot)$ from $mcs[k]$

17   Unassigned clients join their most similar unassigned client
18   Remaining clients create pairs randomly
```
// Step 3: Combine groups
```
19   Merge groups based on similarity until reaching $V$ group size
20   $G \leftarrow$ (final merged groups)
21   **return** $G$

---

similarity between any of its members and those in the target group (line 19). The process iterates from the second step until the desired group size is reached. Privacy guarantees of the group formation phase are explained in Section IV-D.

### C. Phase 2: Co-training

After group formation, clients within a group engage in collaborative training. To mitigate data privacy leakage, differential privacy noise must be applied. However, in a fully decentralized system, DP noise can more severely impact accuracy compared to centralized training due to the limited number of clients sharing knowledge [79], [80], constrained by their computational and communication power. Thus, an efficient algorithm to aggregate knowledge between clients is needed to ensure fast convergence and personalization for each client.

In P4, we use transformed features from a handcrafted layer [81] as input to the client's neural network. These features enhance convergence even in the presence of DP noise and enable strong performance with a shallow neural network (see Section V-B). The transformed features feed into two models: a *local* model, $f_{\phi_i}$, and a *proxy* model, $f_{\theta_i}$ (see Algorithm 2). Clients share only the proxy model within their group. During each training round, both models train on local data $D_i$, but the proxy model's gradients are clipped and DP noise is added. Since the proxy model undergoes DP training, sharing

its gradient updates $\Delta\theta_i$ does not violate DP guarantees as explained in Section IV-D. Clients exchange these proxy model gradients with their group members, incorporating external knowledge into $f_{\theta_i}$ while preserving local data privacy.

To train the local model and proxy model together, inspired by [25], we apply knowledge distillation between these two models [47]. More specifically, training of the local model will be done by classification loss (Equation 8) and KL divergence loss (Equation 9):

$$\mathbf{L}_{CE}(f_{\theta_i}) = \mathbb{E}_{(x,y) \sim D_i} CE[f_{\theta_i}(x) \parallel y] \tag{8}$$

where $CE$ denotes the cross-entropy loss, computed between the proxy model's output and the ground-truth label $y$ corresponding to input $x$.

$$\mathbf{L}_{KL}(f_{\theta_i}; f_{\phi_i}) = \mathbb{E}_{(x,y) \sim D_i} KL[f_{\theta_i}(x) \parallel f_{\phi_i}(x)] \tag{9}$$

In addition, to distill knowledge from the local model to the proxy model, we use the Kullback–Leibler divergence described in Equation 9 to make the proxy model output closer to the local model output.

Therefore, the proxy model learns from both the local model and the local data by combining the classification loss and the KL divergence loss shown in Equation 10. This way, the proxy model learns to correctly predict the true label of training instances as well as to match the probability estimate of its local model.

$$\mathbf{L}_{\theta_i} = (1 - \alpha) \cdot \mathbf{L}_{CE}(f_{\theta_i}) + \alpha \cdot \mathbf{L}_{KL}(f_{\theta_i}; f_{\phi_i}) \tag{10}$$

In the combination of the losses, $\alpha \in [0, 1]$ balances between losses. As it increases, the proxy model will use more information from the local model.

The objective of the local model is shown in Equation 11.

$$\mathbf{L}_{\phi_i} = (1 - \beta) \cdot \mathbf{L}_{CE}(f_{\phi_i}) + \beta \cdot \mathbf{L}_{KL}(f_{\phi_i}; f_{\theta_i}) \tag{11}$$

As $\beta \in [0, 1]$ increases, the local model will rely less on its local data and more on the data from other groups, which is in the proxy model. If $\beta$ is zero, it results in a totally personalized model that is only trained on the local training data.

With this architecture, the local model $f_{\phi_i}$ can extract personalized model weights for the client without dealing with DP noise, leading to improved performance. In other words, by decoupling *local* models from the differential privacy noise, we aim to enhance model performance, resulting in models that better align with the specific data characteristics of individual clients. Beside that, the parameters of proxy model $f_{\theta_i}$ could be aggregated with neighbor clients to benefit from their knowledge under DP guarantees.

We note that it is not necessary for each member of a group to receive model updates $\Delta\theta_j$ from all other members and update its proxy model $\theta_i$. Instead, every few rounds, an aggregator can be selected periodically from within the group. Therefore, other clients within the group only send their model

---

**Algorithm 2:** Co-training Inside One Group

---

**Input:** Proxy $\theta_i^{(0)}$ and local $\phi_i^{(0)}$ parameters after phase 1, distillation weights $\alpha, \beta \in (0,1)$, learning rates $\eta_l, \eta_g > 0$, client subsampling $l$, data subsampling $s$, local dataset size $R$.

**Output:** Proxy $\theta_i^{(T)}$ and local $\phi_i^{(T)}$ models for each client co-trained inside a group under DP guarantees.

1  $\hat{M} \leftarrow$ (one group clients from $G$)
    // In parallel

2  **for** *round* $t = 0, \ldots, T-1$ *at client* $i \in C^t \subset [\hat{M}]$ *of size* $\lfloor l\hat{M} \rfloor$ **do**

3     **for** *local optimization step* $k = 0, \ldots, K-1$ **do**

4         Sample $S_i^k \subset D_i$ of size $\lfloor sR \rfloor$

5         Update local and proxy models:

6            $\theta_i^{(t)} \leftarrow \theta_i^{(t)} - \eta_l \nabla \tilde{\mathcal{L}}_{\theta_i}(S_i^k)$   // DP update

7            $\phi_i^{(t)} \leftarrow \phi_i^{(t)} - \eta_l \nabla \hat{\mathcal{L}}_{\phi_i}(S_i^k)$   // Non-DP update

8     $\phi_i^{(t+1)} \leftarrow \phi_i^{(t)}$

9     $\Delta\theta_i^{(t)} \leftarrow \theta_i^{(t)} - \theta^{(t-1)}$

10    **if** *client* $i$ *is a group aggregator* **then**

11       $\Delta\theta^{(t)} \leftarrow \frac{1}{l\hat{M}} \sum_{j \in C^t} \Delta\theta_j^{(t)}$

12       $\theta^{(t)} \leftarrow \theta^{(t-1)} + \eta_g \Delta\theta^{(t)}$

13       Send $\theta^{(t)}$ to clients $j \in C^t$

14    **else**

15       Send $\Delta\theta_i^{(t)}$ to a group aggregator

16       Receive $\theta^{(t)}$ from a group aggregator

17    Update proxy model:    $\theta_i^{(t+1)} \leftarrow \theta^{(t)}$

18  **return** $\theta_i^{(T)}, \phi_i^{(T)}$

---

updates $\Delta\theta_j$ to the aggregator client and subsequently receives the aggregated model $\theta$ from that client, which is used for the next round of co-training.

### D. Differential Privacy

*a) Privacy goal:* We aim to control the information leakage from individual data of each client $D_i$ in the shared model updates $\Delta\theta_i$. To achieve this, we adapt a differential privacy algorithm for heterogeneous data proposed by [24] to the peer-to-peer learning setting, and ensure the privacy guarantees on top of their theoretical proof. We focus on *record-level* DP with respect to the joint dataset $D$, where $D$ and $D'$ are neighboring datasets if they differ by at most one record (denoted by $\|D - D'\| \leq 1$). We want to ensure privacy (i) towards a third-party observing the final model and (ii) an honest-but-curious aggregator in each group (similar to an honest-but-curious server in [24]). We set the DP budget for the whole training, denoted by $(\epsilon, \delta)$, in advance before group formation (phase 1) and use the same clipping and Gaussian noise as in the original work to achieve this budget.

*b) Privacy guarantees:* To ensure privacy during co-training (see Algorithm 2), each gradient of a proxy model of $i$-th client is divided by its norm as described in Equation 12 and then the noise is added to the gradient as shown in Equation 13 (line 6 in Algorithm 2).

$$\tilde{g}_{ij} = \frac{g_{ij}}{\max(1, \|g_{ij}\|_2/\mathcal{C})} \tag{12}$$

$$\tilde{H}_i = \frac{1}{sR}\Sigma_{j \in S_i}\tilde{g}_{ij} + \frac{2\mathcal{C}}{sR}\mathcal{N}(0, \sigma_g^2) \tag{13}$$

The noise standard deviation $\sigma_g$ in Equation 13, ensuring the specified privacy budget ($\epsilon$) given the clipping norm defined in Equation 12, is computed based on the theoretical upper bound adapted from Theorem 4.1 of [24], as follows:

$$\sigma_g = \Omega\left(\frac{s\sqrt{lTK\log(\frac{2Tl}{\delta})\log(\frac{2}{\delta})}}{\epsilon\sqrt{M'}}\right) \tag{14}$$

It guarantees that under the subsampled Gaussian mechanism [82] and composition accounting via Rényi DP [83], one can ensure (i) $(\mathcal{O}(\epsilon), \delta)$-differential privacy towards a third party observing the global model (e.g., other clients that receive the global model after aggregation), and (ii) $(\mathcal{O}(\epsilon_s), \delta_s)$-DP towards an honest-but-curious aggregator receiving per-round gradients, where $\epsilon_s = \epsilon\sqrt{M'/l}$ and $\delta_s = \frac{\delta}{2}(\frac{1}{l}+1)$. We follow their privacy accounting strategy matched to the peer-to-peer setting. Since in P2P learning individual gradients are directly shared among clients rather than being aggregated, we set $M'$ to 1. To amplify privacy according to [84], [24], client subsampling $l$ and data subsampling $s$ are used (lines 2 and 4 in Algorithm 2). Based on the Equation 13, once the target privacy budget $\epsilon$ and $\delta$ are fixed, one can compute a valid set of values for $K$, $l$, and $s$ to satisfy the guarantee. We specify the chosen values for these parameters in Section V-B and Appendix E-1.

### E. Defense against Poisoning Attacks

In `P4`, we account for a wide range of poisoning attacks, defined in Section III-B. To mitigate these threats while maintaining privacy and utility amplification, `P4` employs a combination of `m-Krum` [36], `anomaly-detection` [46], and knowledge distillation, which utilizes client similarity within each group to filter out malicious updates. To the best of our knowledge, `P4` incorporates a unique set of defense strategies not used in prior P2P work, simultaneously enhancing the privacy-utility trade-off and improving robustness against poisoning attacks under heterogeneous client data distributions. Note that `m-Krum` and `anomaly-detection` are existing defenses used during gradient aggregation; for simplicity, we refer to their combination as "secure aggregation." The rationale behind each defense strategy is discussed in the following subsections.

*1) Data Poisoning Tolerance:* We use a *label flipping* attack to implement targeted data poisoning, which is particularly prevalent in the DFL security literature [31], [32], [33], [85]. Given a source class $c_{src}$ and a target class $c_{target}$ from **C**, each malicious client $i \in M$ modifies their dataset $D_i$ as follows: For all instances in $D_i$ whose class is $c_{src}$, change their class to $c_{target}$. In our setting, client data is poisoned by flipping *all labels* (e.g., for `CIFAR-10` from $c_{src}$ to $9 - c_{src}$ as in [26]). The goal of the attack is to make the proxy model

of each client $f_{\theta_i}$ more likely to misclassify $c_{src}$ images as $c_{target}$ images at test time.

Interestingly, a proxy-based knowledge distillation combined with P4's client clustering presented in Sections IV-B and IV-C shows inherent resilience against data poisoning attacks, particularly *label flipping*. While Section V-C provides an empirical evaluation (see Figures 5(a) and 6), we present an intuition for this tolerance here.

P4 enhances the privacy-utility trade-off by employing knowledge distillation from a local model to a proxy model (Eq. 9) that can be formulated as the KL-divergence between the softmax outputs:

$$\mathbf{L}_{KL}(f_{\theta_i}; f_{\phi_i}) = \mathbb{E}_{(x,y) \sim D_i} KL[\text{softmax}(\frac{z_{\theta_i}}{T}) \parallel \text{softmax}(\frac{z_{\phi_i}}{T})]$$

where $z_{\theta_i}$ and $z_{\phi_i}$ are the logits of proxy and local models, respectively, and $T$ is a temperature that is normally set to 1. From a malicious security perspective, this knowledge distillation mechanism acts as a form of client *self-defense* against poisoning attacks. Similar to LeadFL [86], which introduces an additional regularization term in the loss function to mitigate the effects of poisoned gradients, the KL-divergence term in P4's loss function also serves as a regularization mechanism. This term helps control the influence of poisoned gradients on the final proxy model performance.

When the KL term is incorporated into loss functions $\mathbf{L}_{\theta_i}$ and $\mathbf{L}_{\phi_i}$ (Equations 10 and 11), it enforces alignment between the softmax outputs of the proxy model $f_{\theta_i}$ and the local model $f_{\phi_i}$. Even if the gradients of the proxy model $f_{\theta_i}$ are poisoned, *the local model $f_{\phi_i}$ remains clean*. Consequently, the KL alignment between these two models serves as a corrective mechanism, cleansing the proxy model $f_{\theta_i}$ from the effects of poisoning. This client self-defense strategy is particularly effective against *label flipping*, as this attack aims to alter the softmax output of the proxy model $f_{\theta_i}$ by modifying labels. However, the KL alignment with the softmax output of the local model $f_{\phi_i}$ counteracts these changes from the client side, mitigating the impact of the attack.

Nevertheless, due to high data heterogeneity across clients, using a proxy mechanism alone does not ensure robustness under arbitrary clustering settings. This is because when malicious updates poison the proxy models, benign clients attempt to align their poisoned proxy models' softmax outputs with their *local* models, guided by their local data distribution. However, when these partially cleansed proxy models are averaged under random clustering and client heterogeneity, the differences in benign clients' data distributions diminish the overall effectiveness of the cleansing process. This requires a clustering algorithm that groups clients with similar data distributions, such as our group formation algorithm (see Appendix C-1 for supporting details).

*2) Model Poisoning Tolerance:* While an unmodified P4 effectively mitigates data poisoning attacks, it is less robust against model poisoning attacks (see Section V-C). The primary challenge lies in the significant changes in client gradients (e.g., setting them to zero or random values), which can disrupt simple gradient averaging on the aggregator. To mitigate Byzantine attacks, P4 integrates a robust aggregation method called m-Krum [36].

m-Krum defines a *Krum* aggregation rule $\text{Kr}(w_i, \ldots, w_n)$ as follows. For any $i \neq j$, they denote by $i \rightarrow j$ the fact that $w_j$ belongs to $n - f - 2$ closest vectors to $w_i$, where $n$ is the total number of clients and $f$ is the number of malicious clients. Then, they define for each client $i$, the score $s(i) = \sum_{i \rightarrow j} \|w_i - w_j\|^2$, where the sum runs over $n - f - 2$ closest vectors to $w_i$ based on an Euclidean distance. Finally, $\text{Kr}(w_i, \ldots, w_n) = w_{i_*}$, where $i_*$ refers to the client minimizing the score, $s(i_*) \leq s(i)$ for all $i$. An $m$ parameter in m-Krum defines the number of clients to choose that have minimum scores.

In this setting, client clustering based on model weight similarity before co-training enhances the ability to detect malicious updates. This is because P4 forms groups using the $\ell_1$-norm, which was chosen for its effectiveness in client clustering for privacy-utility amplification. Although the $\ell_1$-norm differs from the $\ell_2$-norm, clients that are close under the $\ell_1$-norm are also close under the $\ell_2$-norm. As a result, benign gradients within a group remain close to one another during training, thereby improving m-Krum's ability to identify and filter out outliers introduced by malicious clients. Moreover, the use of knowledge distillation between the local and proxy models helps mitigate cases where m-Krum fails to filter out malicious updates and a proxy model becomes poisoned. Since the local model remains clean, knowledge distillation helps cleanse the poisoned proxy model. We empirical evaluate an effectiveness of m-Krum with P4 against other defenses in Section V-C.

*3) Anomaly Detection:* To strengthen P4's resilience against anomalous model updates, particularly in the *byzantine-random* attack mentioned in Section III-B, we incorporate anomaly-detection [46]. Our empirical evaluation in Section V-C shows that while m-Krum struggles against *byzantine-random*, such a simple and comparably fast defense as anomaly-detection [46] is highly effective. It leverages the $3\sigma$ rule applied to an $\ell_2$ score during cross-client checks, making it particularly suitable as an addition to m-Krum and P4's clustering approach. Since clients within each group are expected to have high model weight similarity after clustering, anomalous updates can be easily identified and filtered during training. Moreover, in cases where anomaly-detection is less effective, such as in *label flipping*, it does not degrade performance (see Figure 5(a)). Thus, we integrate anomaly-detection as the first layer of defense in P4, complementing m-Krum and client self-defense to enhance overall robustness.

## V. EXPERIMENTS

### A. Experimental Setup

*1) Datasets:* We conduct experiments on three benchmark datasets: CIFAR-10 [87], CIFAR-100 [87], and FEMNIST [88]. Following the non-IID setting from FedAvg [89], we randomly assign classification tasks and training

TABLE I
DATASET STATISTICS.

| Dataset | # classes | # clients | # samples (per client) | # features |
|---------|-----------|-----------|------------------------|------------|
| FEMNIST | 47 | 200 | 300 | $28 \times 28 \times 1$ |
| CIFAR-10 | 10 | 260 | 200 | $32 \times 32 \times 3$ |
| CIFAR-100 | 100 | 60 | 250 | $32 \times 32 \times 3$ |

data to each client. This process generates $M$ clients, each with $R$ samples, split into 80% training and 20% testing. Additionally, 20% of clients are reserved for evaluation and hyperparameter tuning. Dataset statistics are provided in Table I.

*2) Non-IID Setting:* To generate non-IID tasks and analyze the impact of data heterogeneity, we adopt two common methods: (i) alpha-based heterogeneity and (ii) sharding-based heterogeneity. Both partition data among clients to create non-IID distributions. In the alpha-based approach, most of client's data comes from a single class, whereas in the sharding-based method, each client exclusively receives data from a specific subset of multiple classes. Sharding-based heterogeneity follows the procedure of [5], where a dataset with $L$ classes is divided into $P$ shards per class, generating $M = \frac{LP}{N}$ tasks. Each task consists of $N$ randomly assigned classes, with one shard per class. We evaluate the effect of heterogeneity by setting $N \in \{2, 4, 8\}$. Alpha-based heterogeneity, used in prior work [24], [90], controls heterogeneity through a parameter $\gamma$. For each client, $\gamma\%$ of the data is sampled IID from all classes, while the remaining $1-\gamma\%$ comes from a single dominant class. We assess heterogeneity effects using $\gamma \in \{25\%, 50\%, 75\%\}$.

*3) Models:* To assess the efficacy of our methods and compare them with related work, we conduct experiments using both a linear neural network with a softmax activation and a CNN-based architecture [81]. Details about hyperparameters are in Appendix E-1.

*4) Baselines:* We select baselines based on the following criteria: (1) support for a peer-to-peer setting, (2) a privacy-preservation approach based on DP focused on privacy-utility amplification rather than an all-in-one system, and (3) publicly available code. Based on these criteria, we compare P4 with the following baselines, which we adapt to our evaluation settings: *centralized learning*, *local training*, *FedAvg* [89], *Scaffold* [24], *ProxyFL* [25], *DP-DSGT* [39]. Appendix E describes each baseline and how we adapt it to our evaluation settings, and provides additional experimental details.

### B. Privacy-Utility Amplification

We first evaluate the performance of different algorithms under varying data heterogeneity by adjusting $\gamma$ and $N$, which correspond to alpha-based and shard-based heterogeneity, respectively. To simulate a scenario where clients have limited computational power, we fix the number of training rounds at $T = 100$, aiming for fast convergence on each client's data distribution. For differential privacy, we set a target guarantee of $\epsilon = 15$ for all methods. We consider all combinations of hyperparameters: local steps $K \in \{1, 2, \ldots, 10\}$, client sampling ratio $l \in \{0.1, 0.2, \ldots, 1\}$, and data sampling ratio

$s \in \{0.1, 0.2, \ldots, 1\}$, computing the corresponding noise level $\sigma_g$ using Equation 14 to maintain the target $\epsilon$. Group sizes are set to $V = 4$ for CIFAR-100 and $V = 8$ for other datasets. To ensure a fair comparison, we perform hyperparameter tuning (grid search) over the mentioned hyperparameters, along with local step size $\eta_l \in [0.1, 10]$ and clipping norm $\mathcal{C} \in [0.1, 10]$, using evaluation data from Section V-A. The performance of different methods under shard-based heterogeneity is shown in Figure 9 and Figure 8 (Appendix B-1) using a linear model on FEMNIST and CIFAR-100, respectively. Performance under alpha-based heterogeneity is illustrated in Figure 3 (CNN model) and Figures 2, 10 (linear model, Appendix B-1). In the figures, "(HC)" indicates the use of handcrafted features. Each experiment is repeated three times to mitigate randomness, and we report the mean results across trials.

As shown in the figures, our approach consistently outperforms existing methods across different levels and types of data heterogeneity. On CIFAR-10 with alpha-based data generation, our method achieves 58.6% to 62.2% accuracy using only a single linear layer (see Figure 2), making it well-suited for resource-constrained IoT devices. In contrast, FedAvg struggles under data heterogeneity, underscoring the limitations of centralized federated learning with non-IID client data. Scaffold, another centralized method, performs better than FedAvg, particularly in highly heterogeneous settings (small $N$ and $\gamma$). However, despite its strengths in personalized federated learning, Scaffold still faces slow convergence, as observed across both data generation methods.

As shown in Figure 9, DP-DSGT, as a P2P method, achieves performance below the centralized baselines since its upper-bound performance is limited to achieving consensus among all clients (output of centralized methods). By considering the performance of FedAvg, Scaffold, and DP-DSGT on different datasets and architectures, we can conclude that when the tasks of clients are non-IID, consensus learning may not achieve good performance.

In our experiments with a CNN model in Figure 3, which possesses more parameters compared to the linear model, previous approaches suffered from noisy training due to differential privacy noise and collaborative training with clients with different data distribution. In contrast, P4 consistently converges to robust parameter solutions across all settings, even with a minimal number of communication rounds.

Our proposed model maintains stable performance across varying levels of data heterogeneity. Unlike other methods that perform well only under specific conditions—e.g., in Figure 3, Proxy achieves good results at $\gamma = 25\%$ but offers little improvement over local training in other cases—our approach consistently outperforms alternatives across both architectures and heterogeneity levels. A central factor contributing to this consistent performance is P4's client clustering strategy (see Section IV-B).

Figure 2 and Figures 8, 9, 10 in Appendix B-1 show that when client's local dataset contains only a few classes and uses a simple neural network, P4 can match or even surpass centralized training in some cases. This is because shallow
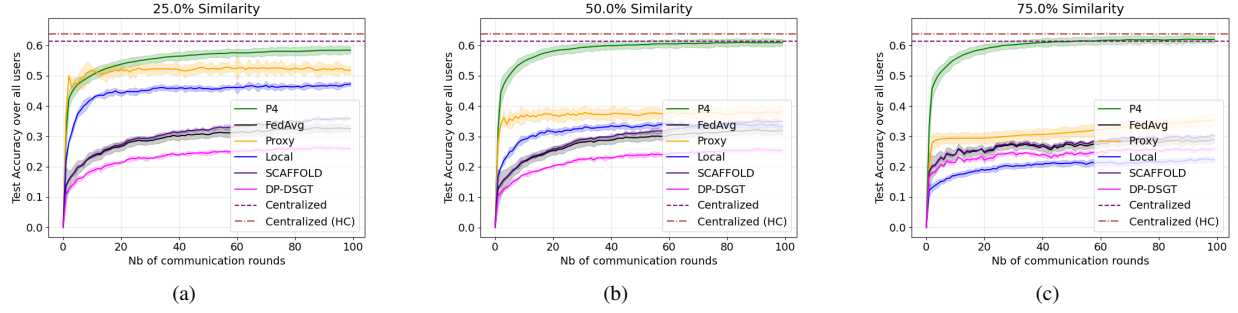
Fig. 2. Test accuracy of a linear model on `CIFAR-10` with $\epsilon = 15$ and alpha-based setting: (a) $\gamma = 25\%$ (b) $\gamma = 50\%$ (c) $\gamma = 75\%$. For each client, $\gamma\%$ of the data is sampled IID from all classes, while the remaining $1 - \gamma\%$ comes from a single dominant class.
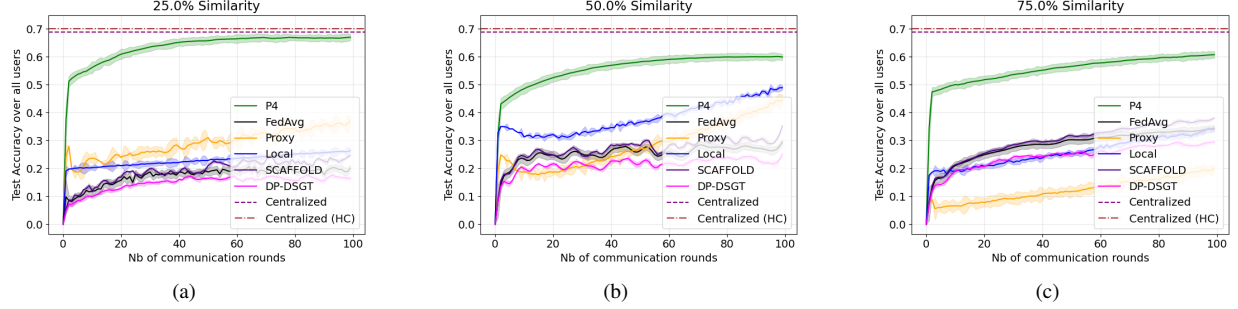


Fig. 3. Test accuracy of CNN on `CIFAR-10` with $\epsilon = 15$ and alpha-based non-IID setting: (a) $\gamma = 25\%$ (b) $\gamma = 50\%$ (c) $\gamma = 75\%$. For each client, $\gamma\%$ of the data is sampled IID from all classes, while the remaining $1 - \gamma\%$ comes from a single dominant class.

networks more effectively learn patterns within a limited data distribution compared to the complexity of classifying all classes in a centralized setting. With fewer model parameters, classifying an entire data distribution becomes challenging for centralized methods, whereas a personalized model tailored to each client's distribution shows even better performance.

Our results suggest that local training can serve as a strong baseline for P2P learning when client data similarity is low. Thus, improving collaborative training in highly heterogeneous data distributions is crucial for achieving higher accuracy than local training. Appendix B-2 provides an extended comparison of `P4` with local training under varying privacy budgets. In general, `P4` outperforms local training in accuracy, even when the privacy budget exceeds 3.

### C. Tolerating Poisoning Attacks

In this section, we evaluate `P4` under the poisoning attacks described in Section III-B and client data heterogeneity. Our goal is to demonstrate that the methodology we use to improve the privacy-utility trade-off also enhances robustness against attacks. Additionally, we aim to show the effectiveness of the defense mechanism proposed in Section IV-E. To this end, we address the following research questions: (RQ1) How does client clustering in `P4` benefit robustness?; (RQ2) How does the combination of client clustering and proxy-based knowledge distillation enhance robustness compared to existing defense strategies?; (RQ3) What percentage of malicious clients can `P4` with a defense proposed in Section IV-E tolerate?; (RQ4) Do `m-Krum` and `anomaly-detection` exhibit better robustness when used with `P4` compared to FedAvg?

Unless otherwise specified, all experiments use a linear model and the same datasets as previous evaluations, with one non-IID setting: 50% similarity for `CIFAR-10` and $N = 4$ for `CIFAR-100` and `FEMNIST`. The proportion of malicious clients is set to 30% (as in [40], [41]), and all clients participate in each iteration to simulate the strongest attack impact. Each experimental pipeline was run with three different seeds. For evaluation, we measure *attack impact* (as in [28]), defined as the test accuracy drop compared to the baseline (without any attack or defense), measured only on benign clients. This metric is chosen because the goal of these attacks is to degrade the quality of the final model.

We do not compare robustness of `P4` against P2P frameworks from Section II that also consider differential privacy and malicious security for the following reasons: (1) these frameworks do not consider heterogeneous client data distributions in their methodology, therefore direct comparisons of robustness in terms of accuracy drop is unfair; (2) our goal is to show that the proposed privacy-utility amplification also enhances robustness, resulting in a good attack tolerance of `P4`, rather than to outperform other existing defenses; (3) our approach is orthogonal to existing DFL defenses, which also can be incorporated into `P4` instead of `m-Krum`; and (4) those privacy-preserving frameworks that address poisoning attack do not propose novel defense strategies, but rely on existing defenses. Instead, in this section we include these existing defenses in our evaluation, and in Section V-B we show that `P4` outperforms other state-of-the-art differential private P2P methods under various heterogeneity settings.
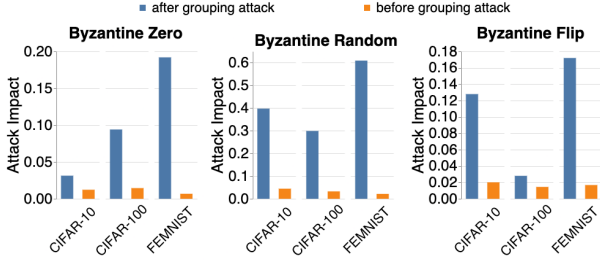
Fig. 4. Impact of attacks on `P4` without secure aggregation before grouping and after grouping for three non-IID datasets and a linear model under 30% of malicious clients.

*a) RQ1:* We first evaluate `P4` without secure aggregation under poisoning attacks. The impact of model poisoning attacks on `P4` differs significantly before and after grouping (phase 1 in Section IV). Figure 4 illustrates the attack impact on `P4` without secure aggregation, showing that after-grouping attacks have a much higher impact than before-grouping ones (similar to [91]). This occurs because, during group formation, `P4` detects malicious gradients based on $\ell_1$-norm similarity, clustering most malicious clients into separate groups that do not affect benign clients. We do not compare *label flipping* before and after grouping, as we assume data poisoning always occurs before grouping. Since our threat model (Section III-B) assumes that an attacker may be aware of the group formation phase, we conduct further evaluations under the after-grouping setting for model poisoning attacks. To simulate this setting, clients are marked as malicious before group formation and start to behave maliciously during training (phase 2). As shown in Figure 4 and dashed lines in Figure 5, `P4` without secure aggregation remains vulnerable to the selected model poisoning and data poisoning attacks. Among them, *label flipping* has the smallest impact, while *byzantine-random* has the highest.

*b) RQ2:* We compare the robustness of `P4` without secure aggregation (or simply client clustering with knowledge distillation) to nine existing defense strategies applied on top of `P4`. These defenses were selected based on their effectiveness in P2P settings and their compatibility with `P4`'s privacy-utility logic (as explained in Section IV-E). All defenses are implemented in the recent FedSecurity benchmark [73], which we extend for our evaluation. Defense configurations follow those in the original papers [34], [92], [35], [93], [94], [46], [36], [37]. Figures 5 and 17(d) (results for *byzantine-flip* are deferred to Appendix C) show attack tolerance of these defenses and `P4` without secure aggregation (`no_defense`), where lower values indicate better performance. The x-axis orders the defenses based on the mean attack impact, and the dashed line shows the median of `no_defense`.

A key observation is that `no_defense` performs best under *label flipping* and remains competitive against model poisoning attacks, supporting the effectiveness of `P4`'s client self-defense (mentioned in Section IV-E). Notably, in Figure 5(a) `anomaly-detection` produces identical results to `no_defense`, indicating its inability to detect anomalies in this setting and, thus, functioning as a normal `no_defense`. Some defenses degrade `P4`'s accuracy compared to `no_defense`, likely due to filtering benign



(a) **Label Flipping**



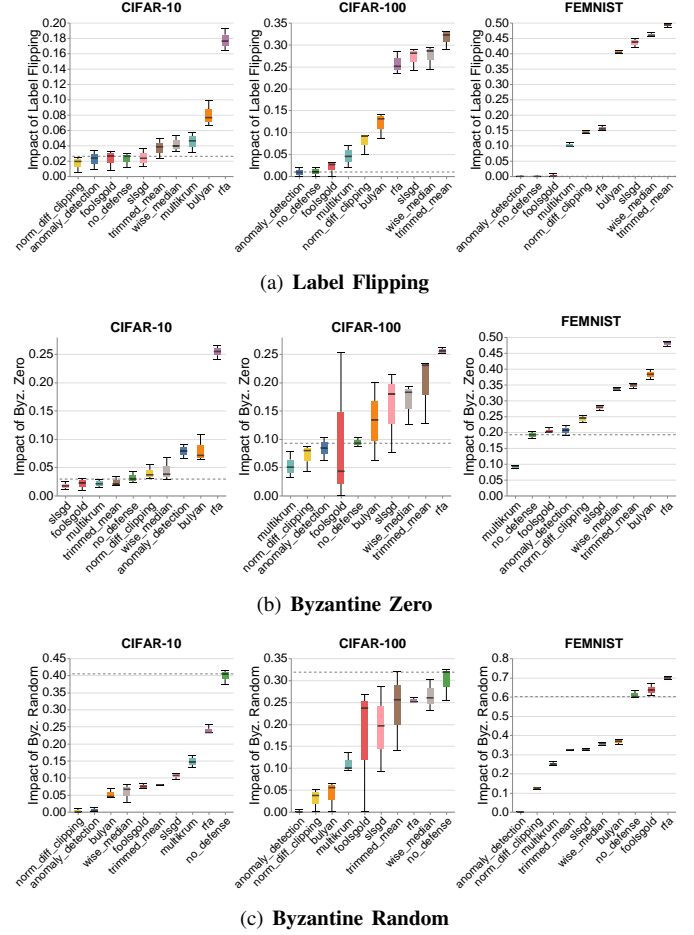(b) **Byzantine Zero**



(c) **Byzantine Random**

Fig. 5. Performance of `P4` with different defenses under poisoning attacks and 30% of malicious clients.

clients (especially under *label flipping*) or negatively impacting convergence. Figures 5(b) and 17(d) show that `m-Krum` [36] performs best against *byzantine-zero* and *byzantine-flip*, with `Foolsgold` [37] being a close competitor, particularly in Figure 17(d). However, since `Foolsgold` performs worse on FEMNIST, `m-Krum` exhibits better overall attack tolerance under *byzantine-flip*. The efficacy of `m-Krum` in combination with `P4` is discussed in Section IV-E. For the most severe attack, *byzantine-random*, the simple and efficient `anomaly-detection` [46] achieves the best results. This may be because it uses model weight similarity during cross-client checks, resonating with `P4`'s clustering strategy. In most cases, `anomaly-detection` does not degrade `no_defense` performance, making it a strong candidate for pairing with `m-Krum`.

*RQ3.* To evaluate `P4`'s robustness under different threat levels, we combine its inherent client self-defense with `anomaly-detection` and `m-Krum`, using the former as a first-layer filter. We vary the proportion of malicious clients from 10% to 50%, reflecting typical ranges in privacy-preserving P2P systems from Section II. The evaluation focuses on CIFAR-100 and FEMNIST, which exhibited higher attack impact in RQ2. As a performance metric, we measured the difference between the test accuracy of `P4` with an "ideal
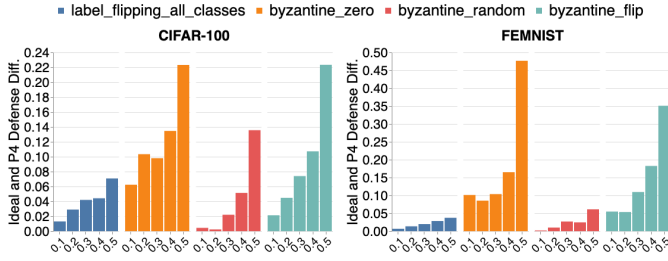
Fig. 6. Differences of the test accuracy of an ideal defense and P4 with secure aggregation on `CIFAR-100`, `FEMNIST`, and a linear model under malicious client percentages from 10% to 50%. Lower values indicate better attack tolerance.

defense" and P4 with the proposed secure aggregation. Here, the "ideal defense" assumes perfect knowledge of malicious clients and excludes them from aggregation. This metric is more suitable than *attack impact* since the cumulative training and test sets for benign clients vary across different malicious client proportions, which may cause an unfair comparison. Thus, for each malicious percentage, we compare the ideal defense against P4's proposed defense.

Figure 6 presents the results averaged across multiple seeds. Following the used benchmark [73], an accuracy drop below 10% indicates good defense performance. Based on this threshold and results across both datasets, P4 with the proposed secure aggregation tolerates up to 50% malicious clients under *label flipping*, 30% under *byzantine-zero*, 40% under *byzantine-random*, and 30% under *byzantine-flip*. The likely reason for the high attack impact above 30% of malicious clients is that m-Krum performs well only when the number of malicious clients is below $(n-2)/2$, according to [36].

Based on this evaluation, P4 with secure aggregation effectively mitigates the impact of up to 30% malicious clients across all four attack types and non-IID settings, while preserving privacy-utility amplification, which satisfies our threat model in Section III-B. Additionally, our evaluation includes a comparison with an "ideal defense," which provides a more rigorous method for assessing attack tolerance under varying percentages of malicious clients. This approach is absent in other P2P frameworks from Section II that provide a defense against poisoning attacks.

*RQ4.* Due to space constraints, the comparison between P4 and FedAvg under attacks is presented in Appendix C-2. Overall, the results indicate that P4 with a defense outperforms FedAvg with the same defense. Moreover, in most cases, the defense is more effective for P4 than for FedAvg in terms of *attack impact*, as it aligns with P4's client clustering.

### D. P4 on Resource-constrained Devices

We assess the practicality and overhead of P4 by deploying it on resource-constrained devices. Specifically, we implement P4 on two Raspberry Pi 4B clients running Debian 12 (Bookworm) 64-bit, Python 3.11.2, and PyTorch 2.1.0. The clients, using a linear model on `CIFAR-10`, communicate via secure websockets. We evaluate four metrics: runtime, memory usage, power consumption, and communication overhead.

Runtime is measured with Python's `time` library, memory usage with the `free` command, and power consumption using Tapo P110 smart plugs via the Tapo API. Communication overhead is assessed by measuring transmitted data between clients. The results for each metric are presented below.

*Runtime.* We measure the average run times across 100 iterations each of both phases. We find that phase 1 (group formation) between two clients takes 0.04 seconds on average (std. 0.02). Assuming the same scenario as the experiment in Section V-B (each client samples 35 clients to compute its model similarity), it would take around 1.4 sec. in total to run phase 1. Phase 2 (co-training) between two clients takes an average of 5.27 sec. ($\pm$ 0.58) to complete, with the bulk of the runtime being the training process (avg. 4.83 s., std. 0.05).

*Memory usage.* We find that running P4 consumes around 72 MB during phase 1 and 489 MB memory during phase 2.

*Power consumption.* The baseline consumption of the Raspberry Pi is $\approx$ 2.64 W (std 0.01). We find that the average consumption during phase 1 is 3.17 W (std. 0.31) and phase 2 is 4.87 W (std. 0.34).

*Communication bandwidth.* In phase 1, the client that initiates the communication sends its model weights to another client, which then performs the comparison with its own weights. The message size of the weights is 622.82 kB (serialized with Python's `pickle`). In phase 2, the client that initiates the co-training first sends its model parameters. After training, the recipient client sends back its gradients to the initiator client for aggregation. The total size of messages exchanged during this phase is 1246.57 kB.

Our experiments indicate that P4 can be effectively run on real-world IoT devices with minimal overhead.

## VI. CONCLUSION

This paper presents P4, a novel peer-to-peer learning framework that simultaneously addresses the challenges of data heterogeneity, privacy preservation, and robustness against poisoning attacks. We propose a new lightweight and fully decentralized client clustering algorithm that privately groups clients with similar model weights using the $\ell_1$-norm. Once grouped, clients engage in private co-training through differentially private knowledge distillation, enabling effective knowledge sharing while preserving individual data privacy. To defend against poisoning attacks, P4 combines existing defenses such as m-Krum and anomaly-detection with client clustering and knowledge distillation. These components, originally designed to improve the privacy-utility trade-off, also enhance the effectiveness of the defense against malicious clients. Our evaluation shows that P4 improves the privacy-utility trade-off, achieving 5% to 30% higher accuracy compared to state-of-the-art differentially private P2P methods, and tolerates up to 30% malicious clients under various heterogeneity settings. Furthermore, our deployment on resource-constrained devices highlights P4's practicality and efficiency, making it a promising solution for real-world decentralized learning applications at the edge. Future work includes enabling secure, dynamic connectivity in decentralized

settings by exploring blockchain-based trust mechanisms for client discovery, authentication, and communication integrity. We also plan to enhance resilience against poisoning attacks by advancing proxy-based knowledge distillation into a client-side self-defense mechanism and by mitigating threats from malicious aggregators through multi-aggregator verification and model consistency checks.

## REFERENCES

[1] Z. Li, J. Lu, S. Luo, D. Zhu, Y. Shao, Y. Li, Z. Zhang, Y. Wang, and C. Wu, "Towards effective clustered federated learning: A peer-to-peer framework with adaptive neighbor matching," *IEEE Transactions on Big Data*, 2022.

[2] Z. Qu, X. Li, R. Duan, Y. Liu, B. Tang, and Z. Lu, "Generalized federated learning via sharpness aware minimization," in *International conference on machine learning*. PMLR, 2022, pp. 18 250–18 280.

[3] Q. Li, B. He, and D. Song, "Model-contrastive federated learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 10 713–10 722.

[4] M. Luo, F. Chen, D. Hu, Y. Zhang, J. Liang, and J. Feng, "No fear of heterogeneity: Classifier calibration for federated learning with non-iid data," *Advances in Neural Information Processing Systems*, vol. 34, pp. 5972–5984, 2021.

[5] S. Li, T. Zhou, X. Tian, and D. Tao, "Learning to collaborate in decentralized learning of personalized models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 9766–9775.

[6] L. Collins, H. Hassani, A. Mokhtari, and S. Shakkottai, "Exploiting shared representations for personalized federated learning," in *International conference on machine learning*. PMLR, 2021, pp. 2089–2099.

[7] T. Li, S. Hu, A. Beirami, and V. Smith, "Ditto: Fair and robust federated learning through personalization," in *International Conference on Machine Learning*. PMLR, 2021, pp. 6357–6368.

[8] V. Zantedeschi, A. Bellet, and M. Tommasi, "Fully decentralized joint learning of personalized models and collaboration graphs," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 864–874.

[9] T. Sun, D. Li, and B. Wang, "Decentralized federated averaging," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4289–4301, 2022.

[10] Y. Shi, L. Shen, K. Wei, Y. Sun, B. Yuan, X. Wang, and D. Tao, "Improving the model consistency of decentralized federated learning," *arXiv preprint arXiv:2302.04083*, 2023.

[11] 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," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.

[12] A. Lalitha, O. C. Kilinc, T. Javidi, and F. Koushanfar, "Peer-to-peer federated learning on graphs," *arXiv preprint arXiv:1901.11173*, 2019.

[13] S. Warnat-Herresthal, H. Schultze, K. L. Shastry, S. Manamohan, S. Mukherjee, V. Garg, R. Sarveswara, K. Händler, P. Pickkers, N. A. Aziz *et al.*, "Swarm learning for decentralized and confidential clinical machine learning," *Nature*, vol. 594, no. 7862, pp. 265–270, 2021.

[14] M. Duan, D. Liu, X. Ji, Y. Wu, L. Liang, X. Chen, Y. Tan, and A. Ren, "Flexible clustered federated learning for client-level data distribution shift," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 11, pp. 2661–2674, 2021.

[15] M. Xie, G. Long, T. Shen, T. Zhou, X. Wang, J. Jiang, and C. Zhang, "Multi-center federated learning," *arXiv preprint arXiv:2108.08647*, 2021.

[16] F. Sattler, K.-R. Müller, and W. Samek, "Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 8, pp. 3710–3722, 2020.

[17] M. N. Nguyen, S. R. Pandey, T. N. Dang, E.-N. Huh, N. H. Tran, W. Saad, and C. S. Hong, "Self-organizing democratized learning: Toward large-scale distributed learning systems," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[18] C. Briggs, Z. Fan, and P. Andras, "Federated learning with hierarchical clustering of local updates to improve training on non-iid data," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–9.

[19] Y. Li, S. Yang, X. Ren, and C. Zhao, "Asynchronous federated learning with differential privacy for edge intelligence," *arXiv preprint arXiv:1912.07902*, 2019.

[20] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, and H. V. Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 3454–3469, 2020.

[21] S. Truex, L. Liu, K.-H. Chow, M. E. Gursoy, and W. Wei, "Ldp-fed: Federated learning with local differential privacy," in *Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking*, 2020, pp. 61–66.

[22] R. C. Geyer, T. Klein, and M. Nabi, "Differentially private federated learning: A client level perspective," *arXiv preprint arXiv:1712.07557*, 2017.

[23] E. Bagdasaryan, O. Poursaeed, and V. Shmatikov, "Differential privacy has disparate impact on model accuracy," *Advances in neural information processing systems*, vol. 32, 2019.

[24] M. Noble, A. Bellet, and A. Dieuleveut, "Differentially private federated learning on heterogeneous data," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 10 110–10 145.

[25] S. Kalra, J. Wen, J. C. Cresswell, M. Volkovs, and H. Tizhoosh, "Decentralized federated learning through proxy model sharing," *Nature communications*, vol. 14, no. 1, p. 2899, 2023.

[26] V. Tolpegin, S. Truex, M. E. Gursoy, and L. Liu, "Data poisoning attacks against federated learning systems," in *Computer security–ESORICs 2020: 25th European symposium on research in computer security, ESORICs 2020, guildford, UK, September 14–18, 2020, proceedings, part i 25*. Springer, 2020, pp. 480–501.

[27] J. Lin, M. Du, and J. Liu, "Free-riders in federated learning: Attacks and defenses," *arXiv preprint arXiv:1911.12560*, 2019.

[28] J. Xu, S.-L. Huang, L. Song, and T. Lan, "Byzantine-robust federated learning through collaborative malicious gradient filtering," in *2022 IEEE 42nd International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2022, pp. 1223–1235.

[29] M. Fang, X. Cao, J. Jia, and N. Gong, "Local model poisoning attacks to {Byzantine-Robust} federated learning," in *29th USENIX security symposium (USENIX Security 20)*, 2020, pp. 1605–1622.

[30] Y. Chen, L. Su, and J. Xu, "Distributed statistical machine learning in adversarial settings: Byzantine gradient descent," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 1, no. 2, pp. 1–25, 2017.

[31] M. Fang, Z. Zhang, Hairi, P. Khanduri, J. Liu, S. Lu, Y. Liu, and N. Gong, "Byzantine-robust decentralized federated learning," in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, 2024, pp. 2874–2888.

[32] N. Heydaribeni, R. Zhang, T. Javidi, C. Nita-Rotaru, and F. Koushanfar, "Surefed: Robust federated learning via uncertainty-aware inward and outward inspection," *arXiv preprint arXiv:2308.02747*, 2023.

[33] P. Sun, X. Liu, Z. Wang, and B. Liu, "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*, 2024, pp. 24 756–24 765.

[34] Z. Sun, P. Kairouz, A. T. Suresh, and H. B. McMahan, "Can you really backdoor federated learning?" *arXiv preprint arXiv:1911.07963*, 2019.

[35] D. Yin, Y. Chen, R. Kannan, and P. Bartlett, "Byzantine-robust distributed learning: Towards optimal statistical rates," in *International conference on machine learning*. Pmlr, 2018, pp. 5650–5659.

[36] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," *Advances in neural information processing systems*, vol. 30, 2017.

[37] C. Fung, C. J. Yoon, and I. Beschastnikh, "The limitations of federated learning in sybil settings," in *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*, 2020, pp. 301–316.

[38] E. Hallaji, R. Razavi-Far, M. Saif, B. Wang, and Q. Yang, "Decentralized federated learning: A survey on security and privacy," *IEEE Transactions on Big Data*, vol. 10, no. 2, pp. 194–213, 2024.

[39] J. Bayrooti, Z. Gao, and A. Prorok, "Differentially private decentralized deep learning with consensus algorithms," *arXiv preprint arXiv:2306.13892*, 2023.

[40] M. Shayan, C. Fung, C. J. Yoon, and I. Beschastnikh, "Biscotti: A blockchain system for private and secure federated learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 7, pp. 1513–1525, 2020.

[41] I. Arapakis, P. Papadopoulos, K. Katevas, and D. Perino, "P4l: Privacy preserving peer-to-peer learning for infrastructureless setups," *arXiv preprint arXiv:2302.13438*, 2023.

[42] A. Mondal, H. Virk, and D. Gupta, "Beas: Blockchain enabled asynchronous & secure federated machine learning," *arXiv preprint arXiv:2202.02817*, 2022.

[43] Z. Qin, X. Yan, M. Zhou, and S. Deng, "Blockdfl: A blockchain-based fully decentralized peer-to-peer federated learning framework," in *Proceedings of the ACM Web Conference 2024*, 2024, pp. 2914–2925.

[44] X. Chen, J. Ji, C. Luo, W. Liao, and P. Li, "When machine learning meets blockchain: A decentralized, privacy-preserving and secure design," in *2018 IEEE international conference on big data (big data)*. IEEE, 2018, pp. 1178–1187.

[45] A. Ghosh, J. Chung, D. Yin, and K. Ramchandran, "An efficient framework for clustered federated learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 19 586–19 597, 2020.

[46] S. Han, W. Wu, B. Buyukates, W. Jin, Y. Yao, Q. Zhang, S. Avestimehr, and C. He, "Kick bad guys out! zero-knowledge-proof-based anomaly detection in federated learning," 2023.

[47] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4320–4328.

[48] S. Dhar, J. Guo, J. Liu, S. Tripathi, U. Kurup, and M. Shah, "A survey of on-device machine learning: An algorithms and learning theory perspective," *ACM Transactions on Internet of Things*, vol. 2, no. 3, pp. 1–49, 2021.

[49] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients-how easy is it to break privacy in federated learning?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 16 937–16 947, 2020.

[50] Y. Huang, S. Gupta, Z. Song, K. Li, and S. Arora, "Evaluating gradient inversion attacks and defenses in federated learning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 7232–7241, 2021.

[51] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou, "A hybrid approach to privacy-preserving federated learning," in *Proceedings of the 12th ACM workshop on artificial intelligence and security*, 2019, pp. 1–11.

[52] 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*, vol. 13, no. 5, pp. 1333–1345, 2017.

[53] X. Zhang, A. Fu, H. Wang, C. Zhou, and Z. Chen, "A privacy-preserving and verifiable federated learning scheme," in *ICC 2020-2020 IEEE International Conference on Communications (ICC)*. IEEE, 2020, pp. 1–6.

[54] J. Zhao, H. Zhu, F. Wang, R. Lu, Z. Liu, and H. Li, "Pvd-fl: A privacy-preserving and verifiable decentralized federated learning framework," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 2059–2073, 2022.

[55] W. Jin, Y. Yao, S. Han, C. Joe-Wong, S. Ravi, S. Avestimehr, and C. He, "Fedml-he: An efficient homomorphic-encryption-based privacy-preserving federated learning system," *arXiv preprint arXiv:2303.10837*, 2023.

[56] 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*, 2017, pp. 1175–1191.

[57] K. Mandal, G. Gong, and C. Liu, "Nike-based fast privacy-preserving highdimensional data aggregation for mobile devices," *IEEE T Depend Secure*, pp. 142–149, 2018.

[58] W. Zheng, R. A. Popa, J. E. Gonzalez, and I. Stoica, "Helen: Maliciously secure coopetitive learning for linear models," in *2019 IEEE symposium on security and privacy (SP)*. IEEE, 2019, pp. 724–738.

[59] F. Mo, H. Haddadi, K. Katevas, E. Marin, D. Perino, and N. Kourtellis, "Ppfl: privacy-preserving federated learning with trusted execution environments," in *Proceedings of the 19th annual international conference on mobile systems, applications, and services*, 2021, pp. 94–108.

[60] F. Mo, Z. Tarkhani, and H. Haddadi, "Sok: machine learning with confidential computing," *arXiv preprint arXiv:2208.10134*, 2022.

[61] A. Dhasade, N. Dresevic, A.-M. Kermarrec, and R. Pires, "Tee-based decentralized recommender systems: The raw data sharing redemption," in *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. IEEE, 2022, pp. 447–458.

[62] M. Zhang, K. Sapra, S. Fidler, S. Yeung, and J. M. Alvarez, "Personalized federated learning with first order model optimization," *arXiv preprint arXiv:2012.08565*, 2020.

[63] A. Lalitha, S. Shekhar, T. Javidi, and F. Koushanfar, "Fully decentralized federated learning," in *Third workshop on bayesian deep learning (NeurIPS)*, vol. 2, 2018.

[64] A. G. Roy, S. Siddiqui, S. Pölsterl, N. Navab, and C. Wachinger, "Braintorrent: A peer-to-peer environment for decentralized federated learning," *arXiv preprint arXiv:1905.06731*, 2019.

[65] R. Dai, L. Shen, F. He, X. Tian, and D. Tao, "Dispfl: Towards communication-efficient personalized federated learning via decentralized sparse training," *arXiv preprint arXiv:2206.00187*, 2022.

[66] Y. Shi, Y. Liu, Y. Sun, Z. Lin, L. Shen, X. Wang, and D. Tao, "Towards more suitable personalization in federated learning via decentralized partial model training," *arXiv preprint arXiv:2305.15157*, 2023.

[67] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur, "Sharpness-aware minimization for efficiently improving generalization," *arXiv preprint arXiv:2010.01412*, 2020.

[68] E. Jeong and M. Kountouris, "Personalized decentralized federated learning with knowledge distillation," *arXiv preprint arXiv:2302.12156*, 2023.

[69] L. Wang, X. Zhao, Z. Lu, L. Wang, and S. Zhang, "Enhancing privacy preservation and trustworthiness for decentralized federated learning," *Information Sciences*, vol. 628, pp. 449–468, 2023.

[70] L. Wang, Y. Xu, H. Xu, M. Chen, and L. Huang, "Accelerating decentralized federated learning in heterogeneous edge computing," *IEEE Transactions on Mobile Computing*, 2022.

[71] A. Bellet, R. Guerraoui, M. Taziki, and M. Tommasi, "Personalized and private peer-to-peer machine learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2018, pp. 473–481.

[72] M. M. Maheri, H. Haddadi, and A. Davidson, "Telesparse: Practical privacy-preserving verification of deep neural networks," *arXiv preprint arXiv:2504.19274*, 2025.

[73] S. Han, B. Buyukates, Z. Hu, H. Jin, W. Jin, L. Sun, X. Wang, W. Wu, C. Xie, Y. Yao *et al.*, "Fedsecurity: A benchmark for attacks and defenses in federated learning and federated llms," in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2024, pp. 5070–5081.

[74] J. R. Douceur, "The sybil attack," in *International workshop on peer-to-peer systems*. Springer, 2002, pp. 251–260.

[75] Y. Dandi, L. Barba, and M. Jaggi, "Implicit gradient alignment in distributed and federated learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 6, 2022, pp. 6454–6462.

[76] 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*. PMLR, 2020, pp. 5132–5143.

[77] W. Huang, M. Ye, Z. Shi, G. Wan, H. Li, B. Du, and Q. Yang, "Federated learning for generalization, robustness, fairness: A survey and benchmark," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[78] E. L. Zec, J. Östman, O. Mogren, and D. Gillblad, "Private node selection in personalized decentralized learning," *arXiv preprint arXiv:2301.12755*, 2023.

[79] Y. Allouah, A. Koloskova, A. El Firdoussi, M. Jaggi, and R. Guerraoui, "The privacy power of correlated noise in decentralized learning," in *Proceedings of the 41st International Conference on Machine Learning*, ser. ICML'24. JMLR.org, 2024.

[80] Z. Zhu, Y. Huang, X. Wang, and J. Xu, "Privsgp-vr: Differentially private variance-reduced stochastic gradient push with tight utility bounds," in *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, K. Larson, Ed. International Joint Conferences on Artificial Intelligence Organization, 8 2024, pp. 5743–5752, main Track. [Online]. Available: https://doi.org/10.24963/ijcai.2024/635

[81] F. Tramer and D. Boneh, "Differentially private learning needs better features (or much more data)," *arXiv preprint arXiv:2011.11660*, 2020.

[82] Y.-X. Wang, B. Balle, and S. P. Kasiviswanathan, "Subsampled rényi differential privacy and analytical moments accountant," in *The 22nd international conference on artificial intelligence and statistics*. PMLR, 2019, pp. 1226–1235.

[83] I. Mironov, "Rényi differential privacy," in *2017 IEEE 30th computer security foundations symposium (CSF)*. IEEE, 2017, pp. 263–275.

[84] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith, "What can we learn privately?" *SIAM Journal on Computing*, vol. 40, no. 3, pp. 793–826, 2011.

[85] L. He, S. P. Karimireddy, and M. Jaggi, "Byzantine-robust decentralized learning via clippedgossip," *arXiv preprint arXiv:2202.01545*, 2022.

[86] C. Zhu, S. Roos, and L. Y. Chen, "Leadfl: Client self-defense against model poisoning in federated learning," in *International Conference on Machine Learning*. PMLR, 2023, pp. 43 158–43 180.

[87] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[88] G. Cohen, S. Afshar, J. Tapson, and A. Van Schaik, "Emnist: Extending mnist to handwritten letters," in *2017 international joint conference on neural networks (IJCNN)*. IEEE, 2017, pp. 2921–2926.

[89] 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*. PMLR, 2017, pp. 1273–1282.

[90] T.-M. H. Hsu, H. Qi, and M. Brown, "Measuring the effects of non-identical data distribution for federated visual classification," *arXiv preprint arXiv:1909.06335*, 2019.

[91] D. Li, W. E. Wong, W. Wang, Y. Yao, and M. Chau, "Detection and mitigation of label-flipping attacks in federated learning systems with kpca and k-means," in *2021 8th International Conference on Dependable Systems and Their Applications (DSA)*. IEEE, 2021, pp. 551–559.

[92] 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*. Springer, 2019, pp. 213–228.

[93] R. Guerraoui, S. Rouault *et al.*, "The hidden vulnerability of distributed learning in byzantium," in *International conference on machine learning*. PMLR, 2018, pp. 3521–3530.

[94] K. Pillutla, S. M. Kakade, and Z. Harchaoui, "Robust aggregation for federated learning," *IEEE Transactions on Signal Processing*, vol. 70, pp. 1142–1154, 2022.

[95] K. Chang, N. Balachandar, C. Lam, D. Yi, J. Brown, A. Beers, B. Rosen, D. L. Rubin, and J. Kalpathy-Cramer, "Distributed deep learning networks among institutions for medical imaging," *Journal of the American Medical Informatics Association*, vol. 25, no. 8, pp. 945–954, 2018.

[96] T. Shen, J. Zhang, X. Jia, F. Zhang, G. Huang, P. Zhou, K. Kuang, F. Wu, and C. Wu, "Federated mutual learning," *arXiv preprint arXiv:2006.16765*, 2020.

[97] H. Kasyap and S. Tripathy, "Privacy-preserving decentralized learning framework for healthcare system," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 17, no. 2s, pp. 1–24, 2021.

[98] E. Oyallon and S. Mallat, "Deep roto-translation scattering for object classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2865–2873.

[99] J. Bruna and S. Mallat, "Invariant scattering convolution networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1872–1886, 2013.

Table II summarizes the main notations used throughout the paper.

TABLE II
SUMMARY OF THE MAIN NOTATIONS.

| Symbol | Description |
|---|---|
| $M,\ i \in [M]$ | number and index of clients |
| $M'$ | number of gradients to be aggregated |
| $T,\ t \in [T]$ | number and index of communication rounds |
| $K,\ k \in [K]$ | number and index of local updates (for each client) |
| $D_i$ | local dataset held by the $i$-th client, composed of points $d_1^i, \ldots, d_R^i$ |
| $R$ | size of any local dataset $D_i$ |
| $D$ | joint dataset $\left( \bigsqcup_{i=1}^M D_i \right)$ |
| $\mathbf{C}$ | set of all possible class values in $D$ |
| $\theta^{(t)}$ | server model after round $t$ |
| $\theta_i^{(t)}$ | proxy model of $i$-th client after round $t$ |
| $\phi_i^{(t)}$ | private model of $i$-th client after round $t$ |
| $g_{ij}$ | gradient calculated on $i$-th client model on its $j$-th data sample |
| $G \in {0,1}^{M \times M}$ | binary matrix, where an entry $G_{ij}$ is one if clients $i$ and $j$ are part of the same group, and zero otherwise |
| $V$ | group size |
| $l \in (0,1)$ | client sampling ratio |
| $s \in (0,1)$ | data sampling ratio |
| $\epsilon, \delta$ | differential privacy parameters |
| $\sigma_g$ | standard deviation of Gaussian noise added for privacy |
| $\mathcal{C} > 0$ | gradient clipping threshold |

## APPENDIX B
## ADDITIONAL RESULTS FOR PRIVACY-UTILITY AMPLIFICATION

*1) Additional heterogeneity settings:* Figures 8, 2 and 10 illustrate the comparison of P4 with existing privacy-preserving methods under additional heterogeneity settings for CIFAR-100, CIFAR-10, and FEMNIST. The plots extend the results presented in Section V-B, and support the statement made in the main body that P4 outperforms existing methods in terms of accuracy under different models, datasets, and data heterogeneity settings (alpha-based and shard-based).

*2) Comparison between collaborative and local training:* Distributed learning involves more communication and privacy concerns compared to local training. In local training, clients improve their models using only their own data, which means that they do not need to communicate with others during training. To see how well our differential private collaborative training method performs compared to the non-DP version of local training and to understand the impact on privacy, we tested our method under different privacy settings, with privacy bounds ranging from $\epsilon = 3$ to $\epsilon = 20$. The tasks are generated using the alpha-based method, and we utilized the FEMNIST dataset for this study. We compared this to how well local training performs on the previously mentioned linear model and data heterogeneity $\gamma = 50\%$. Our findings, illustrated in

Figure 7, reveal that collaborative training outperforms local training even when the privacy budget is set higher than 3. Even when strong privacy constraints are in place (e.g., $\epsilon = 3$), collaborative training performs better, addressing the issue of over-fitting that can affect local training due to its limited amount of local training data. Moreover, as shown in the figure, even though using handcrafted features could improve local training accuracy, it is still not as good as our proposed method. On the other hand, our method shows good robustness against DP noise such that it could have reasonable accuracy with restricted privacy.
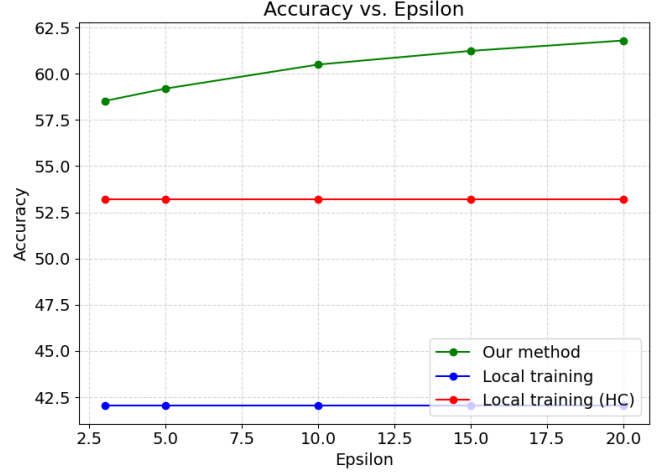


Fig. 7. Performance of a linear model on CIFAR-10 under various privacy budgets and compared to local training.

## APPENDIX C
## ADDITIONAL RESULTS FOR TOLERATING POISONING ATTACKS

*1) Data poisoning tolerance:* As shown in Figure 12, applying only a proxy mechanism is insufficient to fully mitigate data poisoning attacks. The figure illustrates the attack tolerance of different combinations of P4 components on the CIFAR-10 dataset when 30% of clients are malicious. The attack impact percentage represents the difference in performance between the no-attack and attack scenarios, normalized by the no-attack performance to account for varying accuracy levels across different combinations. A lower attack impact percentage indicates better resilience.

The results demonstrate that merely incorporating a proxy mechanism does not guarantee robustness under arbitrary clustering settings (as indicated by the green bar). The primary reason that knowledge distillation fails under random clustering and non-IID conditions is analogous to why the classic FedAvg algorithm [89] performs poorly in non-IID settings. Consider a scenario where a group consists of ten clients, with three being malicious and the remaining seven benign clients exhibiting distinct label distributions, e.g., 90% of local data corresponds to a single label, and this dominant label varies among the benign clients. As training progresses and malicious updates
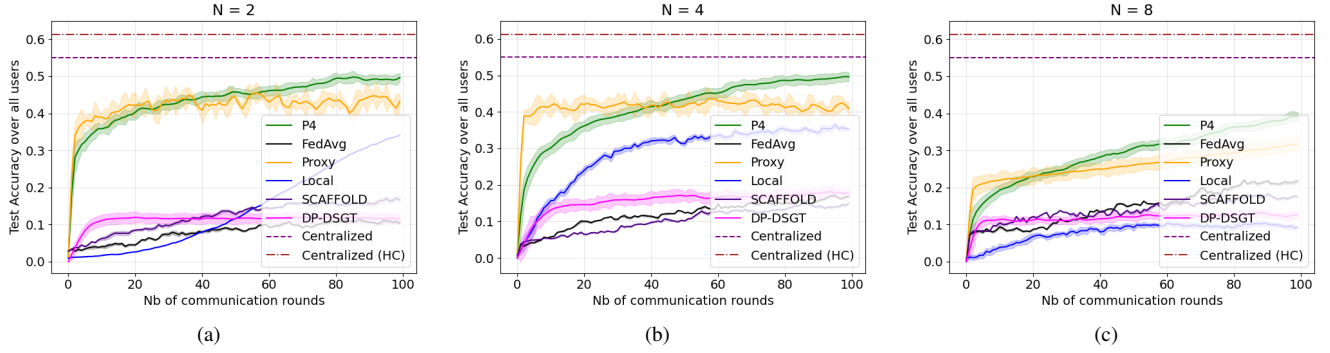
Fig. 8. Test accuracy of a linear model on `CIFAR-100` with $\epsilon = 15$ and shard-based non-IID setting: (a) $N = 2$ (b) $N = 4$ (c) $N = 8$.
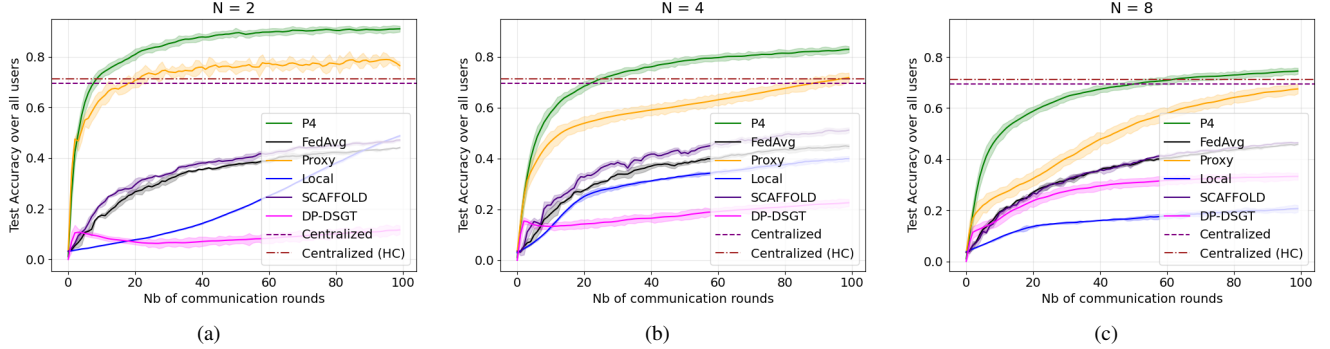


Fig. 9. Test accuracy of a linear model on `FEMNIST` with $\epsilon = 15$ and shard-based non-IID setting: (a) $N = 2$ (b) $N = 4$ (c) $N = 8$. Sharding-based heterogeneity divides a dataset with $L$ classes into $P$ shards per class, generating $M = \frac{LP}{N}$ tasks, where $M$ also equals the total number of clients. Each task consists of $N$ randomly assigned classes, with one shard per class.



Fig. 10. Test accuracy of a linear model on `FEMNIST` with $\epsilon = 15$ and alpha-based setting: (a) $\gamma = 25\%$ (b) $\gamma = 50\%$ (c) $\gamma = 75\%$. For each client, $\gamma\%$ of the data is sampled IID from all classes, while the remaining $1 - \gamma\%$ comes from a single dominant class.

poison the proxy models, benign clients attempt to align their poisoned proxy models' softmax outputs with their *local* models, guided by their local data distribution. However, when these partially cleansed proxy models are averaged, the differences in benign clients' data distributions diminish the overall effectiveness of the cleansing process. In contrast, when clients are clustered using P4's procedure, which leverages $\ell_1$-norm similarity, clients within the same group share a more similar data distribution. This grouping ensures that when proxy models — cleaned based on their respective data distributions — are aggregated, the effectiveness of data

poisoning mitigation remains intact, improving attack resilience. Therefore, an appropriate group formation algorithm is vital for the aforementioned client self-defense.

*2) Additional results for Experiment 2:* This subsection extends the results presented in Section V-C. Figure 17(d) illustrates the performance of P4 with different defenses under a *byzantine-flip* attack and 30% of malicious clients (*RQ2*). The plot supports the statement made in the main body that m-Krum shows better average performance than Foolsgold under *byzantine-flip*, showing the similar attack impact on `CIFAR-10` and `CIFAR-100` and on 8 points better tolerance
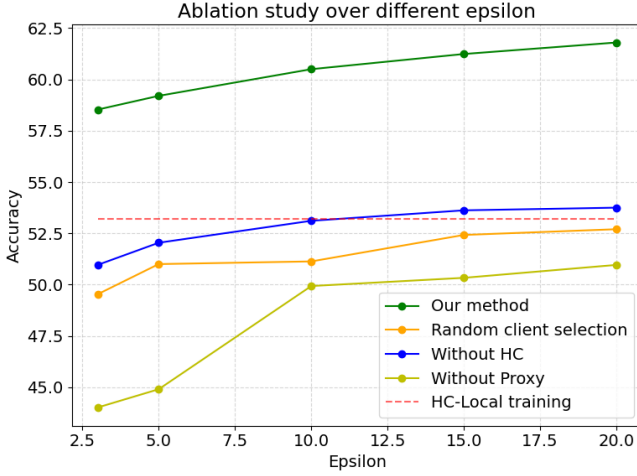
Fig. 11. Comparing the effect of each component of `P4` on `CIFAR-10` and a linear model under data similarity of $\gamma = 50\%$ and a privacy budget $\epsilon$ from 3 to 20.

on `FEMNIST`.

Figures 13 and 14 compare `P4` with a defense (proposed in Section IV-E) against FedAvg with the same defense on `CIFAR-10` and `FEMNIST` under different poisoning attacks (*RQ4* in Section V-C). The results support the statement made in the main body that `P4` with a defense outperforms FedAvg with the same defense. Moreover, the proposed defense works better together with `P4` than with FedAvg, especially on `FEMNIST`, comparing the differences between "`P4` & no_defense" / "`P4` & defense" and "FedAvg & no_defense" / "FedAvg & defense". This highlights that the proposed defense aligns with `P4`'s design and that personalization may also benefit attack tolerance.

## APPENDIX D
### ABLATION STUDIES

#### A. Privacy-Utility Improvement

We now perform an ablation study of `P4` on `CIFAR-10` with 50% similarity using the linear model. We first aim to understand the effect of client selection, handcrafted features, and proxy model individually on the performance of `P4`. Therefore, we compare `P4`'s accuracy results without secure aggregation with three different methods: i) random client selection instead of using our group clustering technique; ii) using raw images instead of handcrafted features; iii) removing the proxy model and using one model per client instead. In each experiment, one of these components is removed and the performance of the model is shown in Figure 11. As shown in the figure, each of these three components has a strong effect on the model performance – removing even one of them results in accuracy lower than the local training baseline in most cases. Our study shows that we can achieve private personalized learning only when all design components of `P4` are enabled.

#### B. Robustness

In this section, we provide an extended analysis of runtimes and test accuracy. We evaluate the impact of increasing the number of clients (Figure 15(a)), group size (Figure 15(b)), and the number of samples per client on `P4` (Figure 16). The experiment was conducted under a *byzantine-random* attack using the proposed defense from Section IV-E, with a client participation ratio of 0.5 per global update and a fixed number of samples per client across all settings. Due to the large number of clients and configuration combinations, we performed this experiment in simulation, while the experiment in Section V-D evaluates `P4` under real-world conditions. The runtime for group formation is measured in simulation without considering its distributed execution on multiple clients, therefore see the runtime for phase 1 in Section V-D for realistic numbers.

Figure 15 presents a breakdown of execution time (in seconds) across different `P4` mechanisms. From Figure 15(a), we observe that even when combining `anomaly-detection` and `m-Krum`, secure aggregation remains efficient and is influenced more by group size than by the total number of clients. This is because, unlike centralized aggregation, `P4` performs aggregation on multiple aggregators (one per group), where group sizes remain relatively small. Additionally, the global update per group introduces minimal overhead compared to an average local update, with its runtime primarily dependent on group size. Since this experiment assumes distributed training across clients, the runtime per global update comprises the local update time for a single client, secure aggregation time, and minor overhead for group management on the aggregator. With multiple aggregators and small group sizes, the runtime difference between a global and local update remains small. Group formation runtime is also more dependent on group size than the total number of clients, though its variation across different settings in Figure 15 remains minimal. The observed inconsistency in local update runtime across different group sizes likely stems from how our code computes the average runtime per client within a group before averaging across all groups. Larger group sizes lead to greater variance in local update runtimes. Figure 16 further illustrates total runtime and test accuracy under varying group sizes and client numbers. Since for this plot the total dataset size remains constant, a smaller number of clients results in more samples per client. The results indicate a trade-off between total runtime and test accuracy, with larger group sizes increasing runtime. Thus, selecting an appropriate group size requires careful consideration of system constraints.

## APPENDIX E
### ADDITIONAL EXPERIMENTAL SETUP

*1) Models and Hyperparameters:* To assess the effectiveness of our methods and compare them with related work, we conduct experiments using both a shallow neural network (a linear layer with a softmax activation) and a CNN-based architecture [81]. To ensure fair comparisons, we perform
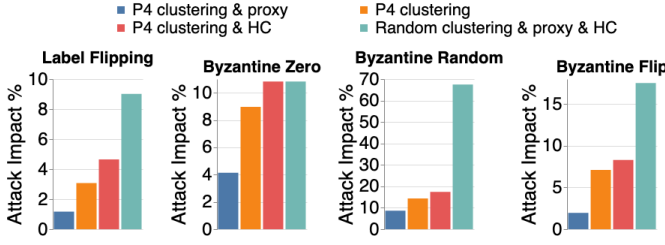
Fig. 12. Impact of P4 components on attack tolerance for CIFAR-10 with a linear model and 30% malicious clients. Attack impact is normalized by test accuracy; lower values indicate better tolerance.
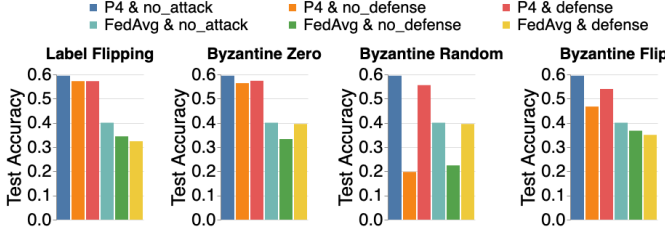


Fig. 13. Comparison of P4 with a defense against FedAvg for CIFAR-10 and a linear model under 30% of malicious clients.
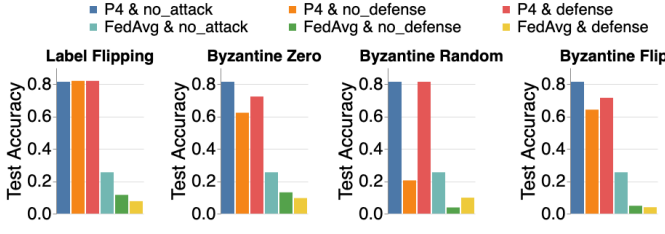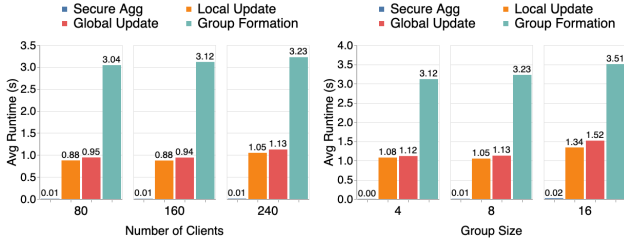


Fig. 14. Comparison of P4 with a defense against FedAvg for FEMNIST and a linear model under 30% of malicious clients.



(a) Varying # of clients for grp. size = 8  (b) Varying group size for 240 clients

Fig. 15. Breakdown of time spent in different mechanisms in P4 on a linear model and CIFAR-10 with 50% alpha-based similarity under 30% of *byzantine-random* clients.

parameter tuning within a fixed range across all methods. Following [24], we set the global step size to $\eta_g = 1$ and define the local step size as $\eta_l = \frac{\eta_0}{sK}$, where $\eta_0$ is carefully tuned. For privacy, we maintain a fixed $\delta = \frac{1}{R}$ in all experiments. Given the fixed number of global training rounds $T$ and a target privacy budget $\epsilon$, we tune $K$, $l$, $s$ values, computing the corresponding noise level $\sigma_g$ using Equation 14.

*2) Baselines:* We select baselines based on the following criteria: (1) support for a peer-to-peer setting, (2) a privacy-preservation approach focused on privacy-utility amplification
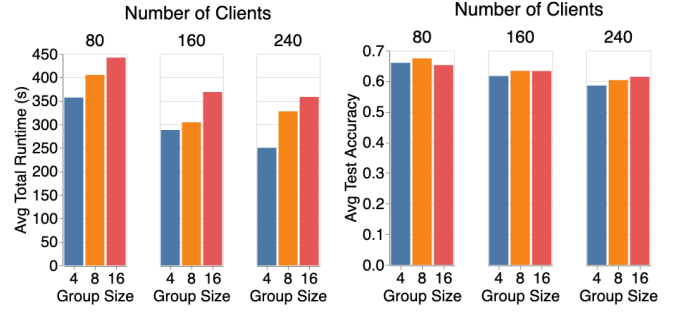


Fig. 16. Average total runtime (in seconds) vs. test accuracy across different numbers of clients and group sizes for CIFAR-10 and a linear model. A dataset size is the same across different numbers of clients, indicating more samples per client for a smaller total number of clients.

rather than an all-in-one system, and (3) publicly available code. Based on these criteria, we compare P4 with the following baselines:

**Centralized learning:** In this setting, all data is stored in a central location without privacy concerns, and a high-computational server trains the model on the entire dataset. This serves as an optimal benchmark, reflecting the classification difficulty. We report results with and without handcrafted features.

**Local training:** Each client trains a model solely on its local dataset. While this may seem like a weaker baseline, our results show that under high data heterogeneity, knowledge sharing provides only marginal improvements. In such cases, many prior approaches fail to surpass local training in accuracy.
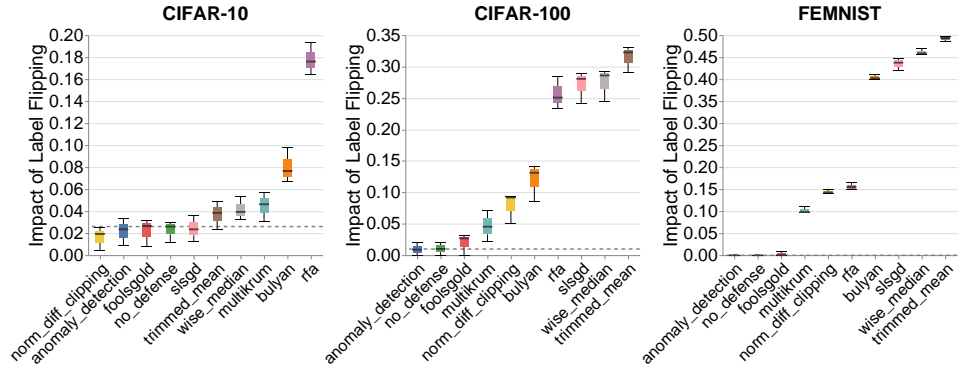
**FedAvg [89]:** A widely-used federated learning algorithm where a centralized server coordinates training. The server is assumed to be honest-but-curious, so we compute the corresponding privacy bound $\epsilon$ with respect to the server.

**Scaffold [24]:** A centralized federated approach designed to handle data heterogeneity under DP constraints. By introducing an additional control parameter per client, Scaffold mitigates client drift from the global model, improving personalization over FedAvg. We assume that the server is untrustworthy, and privacy loss is computed using the sub-sampled Gaussian mechanism under Renyi Differential Privacy (RDP) [83].
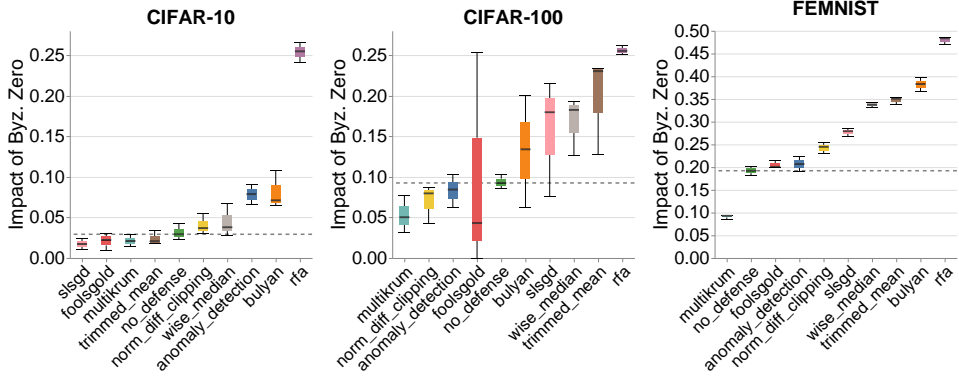
**ProxyFL [25]:** A state-of-the-art decentralized learning method where clients maintain a *private* model and share a *proxy* model with others. We assume the private and proxy models share the same architecture and follow the directed exponential graph communication method and differential privacy mechanism described in the ProxyFL paper.

**DP-DSGT [39]:** One of the state-of-the-art approaches in differentially private P2P learning. Among the three methods proposed in [39], DP-DSGT is designed for consistent performance across varying data distributions, making it well-suited for scenarios where clients have non-overlapping data classes.
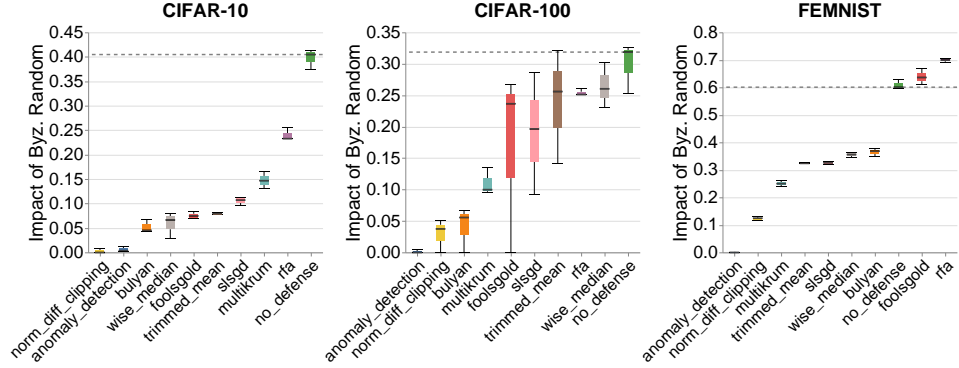
Although previous studies have explored client clustering in decentralized learning [1], [14], [15], [16], [17], [18], [45], these methods either rely on a centralized server for clustering
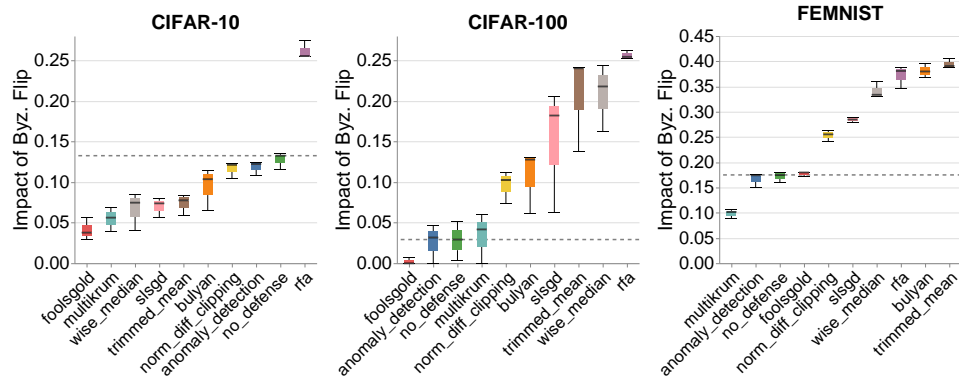
(a) **Label Flipping**



(b) **Byzantine Zero**



(c) **Byzantine Random**



(d) **Byzantine Flip**

Fig. 17. Performance of P4 with different defenses under all four poisoning attacks and 30% of malicious clients.

or violate differential privacy. Thus, we exclude them from our comparison with P4. We also do not compare with AvgPush[2], CWT [95], and FML [96] since ProxyFL has already been shown to outperform them in [25]. Additionally, we exclude SMC-based techniques [97], [51], as these methods are often computationally intensive due to their reliance on cryptographic primitives such fully homomorphic encryption.

*3) Parameter Settings:* As described in Section IV-C, in some experiments we use handcrafted features instead of raw images. Specifically, following [81], we adopt the scattering network of [98], which encodes images using wavelet transforms [99]. Using the default parameters from [98], we employ a ScatterNet $S(x)$ of depth two with wavelets rotated along eight angles. For an image of size $H \times W$, the handcrafted feature extractor outputs $(K, \frac{H}{4}, \frac{W}{4})$, where $K = 81$ for grayscale images and $K = 243$ for RGB images. Additionally, we apply data normalization, where each client computes the mean and variance of its local data, incurring no additional privacy cost. Using these transformed features enhances accuracy and allows classification with a single linear layer, as demonstrated in Section V-B.

## APPENDIX F
### WHY ONE ITERATION OF LOCAL TRAINING IS SUFFICIENT TO DISTINGUISH DIFFERENT DISTRIBUTIONS

*a) Goal.:* We show that if two clients have sufficiently different data distributions, then even a single iteration of local training (under differential privacy) produces updated model weights that are distinguishable with high probability. In particular, when the difference in their expected gradients dominates the sampling, noise, and clipping errors, the $\ell_1$ distance between their weight vectors (cf. Equation (15)) remains large. This observation underpins the clustering strategy described in Section IV-B.

### A. Setup and Notation

Let $\mathcal{D}_i$ and $\mathcal{D}_j$ denote the data distributions for clients $i$ and $j$, respectively. Each client draws its local dataset i.i.d. from its distribution. All clients start from a common initial weight vector $w_0$ (or $\theta_0$) and perform one iteration of gradient-based training with differential privacy (DP). After one iteration, the updated weights for clients $i$ and $j$ are denoted by $w_i$ and $w_j$, respectively.

The dissimilarity metric is given by

$$\text{dissimilarity}(i, j) = \left\| \mathbf{w}_i - \mathbf{w}_j \right\|_1. \quad (15)$$

For each client, let

$$g_i = \frac{1}{n_i} \sum_{t=1}^{n_i} \nabla \ell(w_0; x_t, y_t)$$

be the iteration-averaged empirical gradient (with $n_i$ samples) and define its expectation as

$$\mathbb{E}[g_i] = \mathbb{E}_{(x,y) \sim \mathcal{D}_i} \left[ \nabla \ell(w_0; x, y) \right].$$

[2]A decentralized version of FedAvg using PushSum for aggregation.

We assume that the expected gradients differ by at least

$$\|\mathbb{E}[g_i] - \mathbb{E}[g_j]\|_1 \geq \gamma,$$

for some $\gamma > 0$, which reflects the intrinsic difference in the underlying data distributions.

After gradient clipping with threshold $\mathcal{C}$ and adding Gaussian noise for DP (with noise standard deviation $\sigma = \tilde{\sigma} \mathcal{C}$), the communicated gradient for client $i$ is

$$\widetilde{g}_i = \text{clip}(g_i, \mathcal{C}) + \xi_i,$$

where $\xi_i \sim \mathcal{N}(0, \sigma^2 I_d)$. The corresponding weight update is given by

$$w_i = w_0 - \eta \, \widetilde{g}_i,$$

with learning rate $\eta$.

### B. Assumptions

We make the following assumptions:

(A1) **Gradient Clipping and Clipping Error.** For each client, define the clipping error as

$$\epsilon_{\text{clip}} = \text{clip}(g, \mathcal{C}) - g,$$

where

$$\text{clip}(g, \mathcal{C}) = \min\left(1, \frac{\mathcal{C}}{\|g\|_2}\right) g.$$

We do not assume that clipping is always inactive; instead, standard tail bounds (e.g., under a sub-Gaussian assumption on $g$) ensure that

$$\mathbb{E}\big[\|\epsilon_{\text{clip}}\|_1\big] \leq \mathbb{E}\Big[\|g\|_1 \, \mathbf{1}_{\{\|g\|_2 > \mathcal{C}\}}\Big],$$

so that for appropriate choices of $\mathcal{C}$ the clipping error is small.

(A2) **Iteration-Averaged Gradients.** For each client $i$, the empirical gradient is defined as

$$g_i = \frac{1}{n_i} \sum_{t=1}^{n_i} \nabla \ell(w_0; x_t, y_t),$$

with expectation $\mathbb{E}[g_i]$ as above.

(A3) **Concentration for Sampling and Noise.** There exist constants $c_1$ and $c_2$ such that, with probability at least $1 - \delta$,

$$\|g_i - \mathbb{E}[g_i]\|_1 \leq \frac{c_1}{\sqrt{n_i}}, \quad \|g_j - \mathbb{E}[g_j]\|_1 \leq \frac{c_1}{\sqrt{n_j}},$$

$$\|\xi_i - \xi_j\|_1 \leq c_2 \, \tilde{\sigma} \, \mathcal{C} \, \sqrt{d}.$$

(A4) **Intrinsic Gradient Difference.** The expected gradients satisfy

$$\|\mathbb{E}[g_i] - \mathbb{E}[g_j]\|_1 \geq \gamma,$$

for some $\gamma > 0$.

## C. Proof

### Step 1. Expressing the Weight Difference.

Each client updates its weight vector as

$$w_i = w_0 - \eta\,\widetilde{g}_i, \quad w_j = w_0 - \eta\,\widetilde{g}_j,$$

where the noisy gradient is given by

$$\widetilde{g}_i = \mathrm{clip}(g_i, \mathcal{C}) + \xi_i = g_i + \epsilon_{\mathrm{clip}}^{(i)} + \xi_i.$$

Thus, the difference in weights is

$$w_i - w_j = -\eta\Big[(g_i - g_j) + \big(\epsilon_{\mathrm{clip}}^{(i)} - \epsilon_{\mathrm{clip}}^{(j)}\big) + (\xi_i - \xi_j)\Big].$$

Taking the $\ell_1$ norm and using homogeneity yields

$$\|w_i - w_j\|_1 = \eta\,\|\widetilde{g}_i - \widetilde{g}_j\|_1. \tag{16}$$

### Step 2. Lower-Bounding $\|\widetilde{g}_i - \widetilde{g}_j\|_1$.

By the triangle inequality,

$$\|\widetilde{g}_i - \widetilde{g}_j\|_1 \ge \|g_i - g_j\|_1 - \|\epsilon_{\mathrm{clip}}^{(i)} - \epsilon_{\mathrm{clip}}^{(j)}\|_1 - \|\xi_i - \xi_j\|_1. \tag{17}$$

We now decompose the empirical gradient difference as

$$g_i - g_j = \Big[\mathbb{E}[g_i] - \mathbb{E}[g_j]\Big] + \Big[(g_i - \mathbb{E}[g_i]) - (g_j - \mathbb{E}[g_j])\Big].$$

Applying the triangle inequality again gives

$$\|g_i - g_j\|_1 \ge \|\mathbb{E}[g_i] - \mathbb{E}[g_j]\|_1 - \|(g_i - \mathbb{E}[g_i]) - (g_j - \mathbb{E}[g_j])\|_1. \tag{18}$$

By assumption (A4), $\|\mathbb{E}[g_i] - \mathbb{E}[g_j]\|_1 \ge \gamma$, and by (A3),

$$\|(g_i - \mathbb{E}[g_i]) - (g_j - \mathbb{E}[g_j])\|_1 \le \frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}}.$$

Thus,

$$\|g_i - g_j\|_1 \ge \gamma - \left(\frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}}\right). \tag{19}$$

Moreover, from (A3) the noise term satisfies

$$\|\xi_i - \xi_j\|_1 \le c_2\,\tilde{\sigma}\,\mathcal{C}\,\sqrt{d}.$$

Define the total clipping deviation as

$$(\mathrm{clip\_dev}) = \|\epsilon_{\mathrm{clip}}^{(i)}\|_1 + \|\epsilon_{\mathrm{clip}}^{(j)}\|_1.$$

Substituting these bounds into (17) yields

$$\|\widetilde{g}_i - \widetilde{g}_j\|_1 \ge \gamma - \left(\frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}}\right) - c_2\,\tilde{\sigma}\,\mathcal{C}\,\sqrt{d} - (\mathrm{clip\_dev}). \tag{20}$$

### Step 3. Final Separation in Weight Space.

Multiplying the bound in (20) by the learning rate $\eta$ (cf. (16)), we obtain

$$\|w_i - w_j\|_1 \ge \eta\left[\gamma - \left(\frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}} + c_2\,\tilde{\sigma}\,\mathcal{C}\,\sqrt{d} + (\mathrm{clip\_dev})\right)\right]. \tag{21}$$

Thus, if

$$\gamma > \frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}} + c_2\,\tilde{\sigma}\,\mathcal{C}\,\sqrt{d} + (\mathrm{clip\_dev}),$$

then with high probability,

$$\|w_i - w_j\|_1 \ge \eta\left[\gamma - \left(\frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}} + c_2\,\tilde{\sigma}\,\mathcal{C}\,\sqrt{d} + (\mathrm{clip\_dev})\right)\right] > 0.$$

*a) Discussion.:* This lower bound demonstrates that the separation between the updated weights is primarily determined by the intrinsic difference $\gamma$ in the expected gradients of clients with distinct data distributions. Notably, under typical assumptions, $\gamma$ scales linearly with the model dimension $d$, so that as $d$ increases, the separation grows proportionally. Moreover, as the number of training samples $n_i$ increases, the sampling error terms $\frac{c_1}{\sqrt{n_i}}$ and $\frac{c_1}{\sqrt{n_j}}$ shrink, ensuring that $\gamma$ dominates the bound. Consequently, even after incorporating DP noise (with $\sigma = \tilde{\sigma}\,\mathcal{C}$) and clipping errors, the expected gradient difference remains the key factor determining the grouping metric. This validates the use of the $\ell_1$ distance between weight updates as a reliable proxy for distinguishing underlying data distributions, thereby supporting effective clustering.

### D. Conclusion

Under assumptions (A1)–(A4) and with the DP noise scale set as $\sigma = \tilde{\sigma}\,\mathcal{C}$, we have rigorously shown that the $\ell_1$ distance between clients' weight updates satisfies

$$\|w_i - w_j\|_1 \ge \eta\left[\gamma - \left(\frac{c_1}{\sqrt{n_i}} + \frac{c_1}{\sqrt{n_j}} + c_2\,\tilde{\sigma}\,\mathcal{C}\,\sqrt{d} + (\mathrm{clip\_dev})\right)\right],$$

with high probability. In other words, when the intrinsic difference in expected gradients, $\gamma$, exceeds the aggregate sampling, noise, and clipping errors, the clients' weight updates remain well separated. This finding underpins the clustering strategy by confirming that the $\ell_1$ distance between weight updates is an effective proxy for distinguishing the underlying data distributions of clients.