

How to Prevent the Poor Performance Clients for Personalized Federated Learning?

Zhe Qu¹, Xingyu Li^{2*}, Xiao Han³, Rui Duan³, Chengchao Shen¹, and Lixing Chen⁴

¹Central South University, ²Mississippi State University,

³University of South Florida, ⁴Shanghai Jiaotong University

{zhe_qu, scc.cs}@csu.edu.cn, xl292@msstate.edu,

{xiaoh, ruiduan}@usf.edu, lxchen@sjtu.edu.cn

Abstract

Personalized federated learning (pFL) collaboratively trains personalized models, which provides a customized model solution for individual clients in the presence of heterogeneous distributed local data. Although many recent studies have applied various algorithms to enhance personalization in pFL, they mainly focus on improving the performance from averaging or top perspective. However, part of the clients may fall into poor performance and are not clearly discussed. Therefore, how to prevent these poor clients should be considered critically. Intuitively, these poor clients may come from biased universal information shared with others. To address this issue, we propose a novel pFL strategy, called *Personalize Locally, Generalize Universally (PLGU)*. PLGU generalizes the fine-grained universal information and moderates its biased performance by designing a *Layer-Wised Sharpness Aware Minimization (LWSAM)* algorithm while keeping the personalization local. Specifically, we embed our proposed PLGU strategy into two pFL schemes concluded in this paper: with/without a global model, and present the training procedures in detail. Through in-depth study, we show that the proposed PLGU strategy achieves competitive generalization bounds on both considered pFL schemes. Our extensive experimental results show that all the proposed PLGU based-algorithms achieve state-of-the-art performance.

1. Introduction

Federated Learning (FL) is a popular collaborative research paradigm that trains an aggregated global learning model with distributed private datasets on multiple clients [19, 32]. This setting has achieved great accomplishments when the local data cannot be shared due to privacy and communication constraints [40]. However, because of the

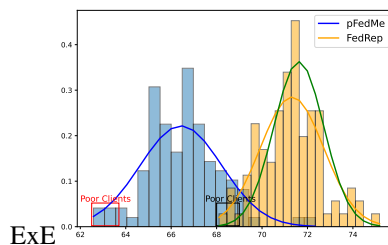


Figure 1. Toy example in a heterogeneous pFL on CIFAR10, which includes 100 clients and each client obtains 3 labels.

non-IID/heterogeneous datasets, learning a single global model to fit the “averaged distribution” may be difficult to propose a well-generalized solution to the individual client and slow the convergence results [27]. To address this problem, personalized federated learning (pFL) is developed to provide a customized local model solution for each client based on its statistical features in the private training dataset [7, 12, 14, 38]. Generally, we can divide existing pFL algorithms into two schemes: (I) with a global model [7, 26, 28, 44] or (II) without a global model [30, 31, 41].

Though many pFL algorithms make accomplishments by modifying the universal learning process [45, 51] or enhancing the personalization [5, 30, 41], they may lead part of clients to fall into poor learning performance, where the personalization of local clients performs a large statistical deviation from the “averaged distribution”. To the best of our knowledge, none of the existing studies explore how to prevent clients from falling into poor personalized performance on these two schemes. For example, the poor medical learning models of some clients may incur serious medical malpractice. To better present our concerned problem, we introduce a toy example in Figure 1, which is learned by two pFL algorithms representing these two schemes: pFedMe [44] and FedRep [5]. Though both algorithms achieve high averaged local model performance of 66.43% and 71.35%, there also 15% of clients are less than 64% and 14% clients are less than 69%, respectively.

*Corresponding author.

This motivates us to exploit an effective strategy to prevent clients from falling into poor performance while without degrading others, e.g., the green curve.

Intuitively, we consider this phenomenon oftentimes comes from the biased universal information towards the clients with better learning performance. For scheme I, a simple-averaged aggregation may not perfectly handle data heterogeneity, as it generates serious bias between the global and local clients. For scheme II, abandoning the universal contribution may dismiss some information from other clients. Instead of designing a new pFL algorithm, we propose a novel pFL strategy on existing pFL studies: generalizing the universal learning for unbiased local adaptation as well as keeping the local personalized features, called Personalize Locally, Generalize Universally (PLGU). The main challenge of PLGU is to generalize universal information without local feature perturbation, as the statistical information is only stored locally. In this paper, we tackle this challenge by developing a fine-grained perturbation method called Layer-Wised Sharpness-Aware-Minimization (LWSAM) based on the SAM optimizer [9, 37], which develops a generalized training paradigm by leveraging linear approximation. Furthermore, we present how to embed this PLGU strategy with the perturbed universal generalization on both the two pFL schemes.

For scheme I (with the global model), we propose the PLGU-Layer Freezing (LF) algorithm. As illustrated in [24, 34, 52], each layer in a personalized model shares a different contribution: the shallow layers focus more on local feature extraction (personalization), and the deeper layers are for extracting global features (universal). Specifically, the PLGU-LF first explores the personalization score of each layer. Then, PLGU-LF freezes the important layer locally for personalization and uses the LWSAM optimizer with the consideration of obtained layer importance score for universal generalization. For scheme II (without the global model), we mainly focus on our proposed PLGU strategy FedRep algorithm [5], named PLGU-GRep. It generalizes the universal information by smoothing personalization in the representation part. To show the extensibility, we present that we can successfully extend our PLGU strategy to pFedHN [41], called PLGU-GHN, to improve learning performance, especially for poor clients. Furthermore, we analyze the generalization bound on PLGU-LF, PLGU-GRep, and PLGU-GHN algorithms in-depth. Extensive experimental results also show that all three algorithms successfully prevent poor clients and outperform the average learning performance while incrementally reducing the top-performance clients.

2. Related Work

Various algorithms to realize pFL can be classified by different measurements. From the universal learning per-

spective, we can divide the existing pFL algorithms into two schemes: with or without a global shared model [22, 45]. Typically, algorithms on the scheme I (with a global model) are mainly extended from the conventional FL methods, i.e., FedAvg [32] or FedProx [27], which combines the adaption of local personalized features on local training updates procedure, such as fine-tuning [1, 33], regularized loss function [26, 44], model mixture [4, 6, 31, 36], and meta-learning [7, 17].

Scheme II of pFL, i.e., without a global model, has more diverse algorithms. [14, 31] propose to train multiple global models at the server, where similar clients are clustered into several groups and different models are trained for each group. However, these algorithms may incur huge communication costs. In addition, some algorithms collaboratively train the customized model with only layer-wised transfer universally to enhance the personalization [30, 41]. Specifically, other algorithms address heterogeneity in pFL by sharing some data [56] or generating additional data [16] on the server, which may violate the privacy policy [45].

The generalization of deep neural networks has been studied as an important topic, which can avoid the learning model to overfit the training dataset. Previous algorithms usually add auxiliary changes to the standard training process, e.g., dropout [43], and normalization [15, 49], which require the acknowledgment of training data that are not feasible in pFL. More specifically, solving the minimax objective will incur large computational costs. Recently, some studies in centralized learning [20, 55] observe that generalization is positively correlated to the sharpness of the training loss landscape. Motivated by this, [9] develops Sharpness-Aware Minimization (SAM), which uses approximated weight perturbation to leverage generalization by leveraging the first-order Taylor expansion. Moreover, [18, 37] extend SAM to FL and graph neural network.

3. Problem Formulation and Strategy Design

3.1. Problem Formulation

The goal of pFL is to collaboratively train personalized models on multiple clients by only sharing the model information while preserving the private local data. Generally, we conclude the pFL into schemes: scheme I (with global model) [4, 5, 28, 34] and scheme II (without global model) [30, 41]. Let \mathcal{N} be the set of clients with the size of N , where the non-IID distributed training data on i -th client is denoted as $\mathcal{D}_i = \{(x_i, y_i)\}, i \in \mathcal{N}$, x_i, y_i are the corresponding data pair. For scheme I, let θ_i denote the personalized model of client i , the objective of both schemes pFL can be formulated as follows:

$$\min_{\mathbf{w}, \{\theta_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N F_i(\mathbf{w}; \theta_i), \quad (1)$$

where F_i is the loss function of the client i associated with its dataset \mathcal{D}_i , typically represented by the cross-entropy loss F_{CE} between the predicted and true label as $F_i(\theta_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} F_{\text{CE}}(\theta_i; x_i^j, y_i^j)$, and m_i is the number of data samples. Note that the loss function $F_i(\mathbf{w}; \theta_i)$ denotes that the pFL objective takes either the global model \mathbf{w} or personalized model θ_i as the target model for classification tasks. Note that the objective function on scheme II can be formulated as $\min_{\{\theta_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N F_i(\theta_i)$.

3.2. Personalize Locally, Generalized Universally

Although existing pFL studies have great accomplishments by enhancing personalization [4,5,28,30,34,41], they cannot guarantee that all clients can achieve desired learning performance. Especially, as shown in Figure 1, only focusing on personalization can lead to a biased pFL result, where poor clients that perform much lower learning performance suffer from the large client deviation. This phenomenon happens because the universal information shared across all clients may be not general enough or biased towards some typical clients. Although common sense is that the header of the neural network stores the personalization and other layers obtain the universal information [5,34], only a few explore the importance of universal information on personalization and show how to moderate the personalization in the universal information. Therefore, in this paper, we aim to propose a pFL strategy to address this issue and prevent poor clients, called Personalize Locally, Generalize Universally (PLGU). Specifically, PLGU has two main parts: (i) extracting the features towards personalization and keeping them locally and (ii) generalizing the universal information.

To generalize the universal information, we leverage the Sharpness Aware Minimization (SAM) algorithm [9] to be the local optimizer. In SAM, the parameters of \mathbf{w}_i whose neighbors within the ℓ_p ball are perturbed for a low training loss $F_{\mathcal{D}_i}$ through the following objective function:

$$F_{\mathcal{D}_i}(\tilde{\theta}_i) = \max_{\|\epsilon_i\|_p \leq \rho} F_{\mathcal{D}_i}(\theta_i + \epsilon_i), \quad (2)$$

where $p \geq 0$, ρ is the radius of the ℓ_p ball, $\tilde{\theta}_i = \theta_i + \epsilon_i$, and $F_{\mathcal{D}_i}(\tilde{\theta}_i)$ is the loss function of SAM on client i . Considering the non-trivial effort to calculate the optimal solution for the inner maximization, SAM uses one extra back-forward gradient ascent step to approximate $\tilde{\epsilon}_i$:

$$\tilde{\epsilon}_i = \rho \frac{\nabla_{\theta_i} F_{\mathcal{D}_i}(\theta_i)}{\|\nabla_{\theta_i} F_{\mathcal{D}_i}(\theta_i)\|} \approx \arg \max_{\|\epsilon_i\|_p \leq \rho} F_{\mathcal{D}_i}(\theta_i + \epsilon_i). \quad (3)$$

As such, SAM computes the perturbed model $\theta_i + \tilde{\epsilon}_i$ for the gradient in objective (2) as $\nabla_{\theta_i} F_{\mathcal{D}_i}(\tilde{\theta}_i) \approx \nabla_{\theta_i} F_{\mathcal{D}_i}(\theta_i)|_{\theta_i + \tilde{\epsilon}_i}$. However, SAM adds perturbation on the entire local model [9,29,37], which dismisses the interior impact of universal information. To generalize the fine-grained universal information, we develop a Layer-Wise

Algorithm 1 PLGU($\tilde{\theta}_i^{t,0}, \tilde{\mathbf{w}}_i^{t,0}, \Lambda_i^t, K, \eta$).

- 1: **Input:** personalized model $\tilde{\theta}_i^{t,0}$, global model $\tilde{\mathbf{w}}_i^{t,0}$, scaling matrix Λ_i^t , number of local epochs K , learning rate η ;
 - 2: **for** $k = 0, \dots, K - 1$ **do**
 - 3: Sample mini-batch \mathcal{B}_i on client i ;
 - 4: Calculate unbiased gradient $\mathbf{g}_i^{t,k} = \nabla_{\tilde{\theta}_i^{t,k}} F_{\mathcal{B}_i}(\tilde{\theta}_i^{t,k})$;
 - 5: Update personalized model $\tilde{\theta}_i^{t,k+1} = \tilde{\theta}_i^{t,k} - \eta \mathbf{g}_i^{t,k}$;
 - 6: Calculate perturbation $\tilde{\epsilon}_i^{t,k}$ by (5);
 - 7: Calculate unbiased gradient approximation for LWSAM $\tilde{\mathbf{g}}_i^{t,k} = \nabla_{\tilde{\mathbf{w}}_i^{t,k} + \tilde{\epsilon}_i^{t,k}} F_{\mathcal{B}_i}(\tilde{\theta}_i^{t,k} + \tilde{\epsilon}_i^{t,k})$;
 - 8: Update global model $\tilde{\mathbf{w}}_i^{t,k+1} = \tilde{\mathbf{w}}_i^{t,k} - \eta \tilde{\mathbf{g}}_i^{t,k}$;
 - 9: **end for**
-

SAM (LWSAM) to be the local training optimizer inspired by [29,53]. Instead of simply applying the scaling to guide the training update [29], we leverage the inner maximization in LWSAM for the layer-wised scaling based on the property of all clients. Let Λ_i denote a diagonal $L \times L$ matrix, $\Lambda_i = \text{diag}(\xi_{i,1}, \dots, \xi_{i,L})$, where $\xi_{i,l}$ is the layer personalization score. We apply the adopted scaling method to the inner maximization of SAM on client i as follows:

$$F_{\mathcal{D}_i}(\tilde{\theta}_i) = \max_{\|\Lambda_i \epsilon_i\| \leq \rho} F_{\mathcal{D}_i}(\theta_i + \Lambda_i \epsilon_i). \quad (4)$$

Note that the advantage of LWSAM can obtain fine-grained universal information, because it adds more perturbation to the personalized layers, i.e., a higher value of $\xi_{i,l}$ scales more perturbation and moderates the personalization from the global model perspective. The layer-wised weight perturbation in LWSAM is also solved by the first-order approximation of (4). Considering the added Λ_i , the approximate inner solution of LWSAM can be written as follows:

$$\tilde{\epsilon}_i = \rho \text{sign}(\nabla_{\theta_i} F_{\mathcal{D}_i}(\theta_i)) \Lambda_i \frac{|\nabla_{\theta_i} F_{\mathcal{D}_i}(\theta_i)|^{q-1}}{(\|\nabla_{\theta_i} F_{\mathcal{D}_i}(\theta_i)\|_q^q)^{p-1}}, \quad (5)$$

where $\frac{1}{p} + \frac{1}{q} = 1$. (5) provides the layer-wise calculation of $\tilde{\epsilon}_i$ to scale up the batch size on client i . Specifically, the PLUG strategy based on the LWSAM algorithm is illustrated in Algorithm 1. In SAM, the first SGD step is only to calculate the perturbation. But the LWSAM algorithm can efficiently leverage the two output models. In particular, Lines 4-5 can obtain the personalized model, which aim to seek the optimal point of the model parameter $\tilde{\theta}_i^{t,k}$. Lines 6-8 aim to seek the loss land surface $\tilde{\mathbf{w}}_i^{t,k}$. As such, both the two SGD steps obtain the models $\tilde{\theta}_i^t$ and $\tilde{\mathbf{w}}^t$ for our learning goal. In the following two sections, we will present how to calculate the layer personalization score Λ_i and how to embed our proposed PLGU strategy into the two schemes separately.

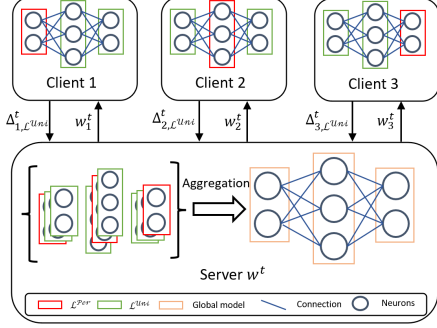


Figure 2. Illustration of PLGU-LF algorithm.

4. PLGU for Scheme I

4.1. PLGU-LF Algorithm

Although we also consider the pFL scheme I [26, 28, 34, 44], which includes a global model w and N personalized local models w_i at the same time, we should carefully design a more general global model w while not influencing the personalization, i.e., both global and personalized models can prevent poor clients. For extracting personalization features, we propose a distance metric for the l -th layer between the corresponding global model and local models to determine the personalization degree of each layer on the local model. The personalization score $\xi_{i,l}^t$ of the scaling matrix Λ_i can be calculated as follows:

$$\xi_{i,l}^t = \frac{\|\tilde{\theta}_{i,l}^t - \tilde{w}_l^t\|}{\dim(\tilde{w}_l^t)}, \quad (6)$$

where $\tilde{\theta}_{i,l}^t$ and \tilde{w}_l^t are the model parameters at the l -th layer of the personalized model $\tilde{\theta}_i^t$ and the global model \tilde{w}^t . $\dim(\cdot)$ denotes the number of parameters on layer l , which can normalize the values as $\sum_{l=1}^L \xi_{i,l}^t = 1$. Note that leveraging (6) to extract the personalization information does not incur huge computational costs compared to inference networks-based feature extraction [1, 28]. In addition, calculating the personalization score Λ_i in (6) uses simple linear algebra, and thus the main computational cost of the PLGU strategy comes from two steps SGD in LWSAM.

In this consideration, $\xi_{i,l}^t$ measures the l -th layer difference between the personalized model $\tilde{\theta}_i^t$ and the global model \tilde{w}^t , which quantifies the personalized contributions of i -th local client to the l -th layer on the global model at the communication round t . Intuitively, a larger value of $\xi_{i,l}^t$ indicates that the l -th layer of $\tilde{\theta}_i^t$ deviates further from \tilde{w}^t , which indicates more contributions to personalization. On the contrary, a layer with a smaller value of $\xi_{i,l}^t$ has more contributions to the universal information.

Hence, based on the personalization score $\xi_{i,l}^t$, we can divide the layers of $\tilde{\theta}_i^t$ into two categories, where the number of D largest values of $\xi_{i,l}^t$ are categorized into the set

Algorithm 2 Scheme I: PLGU-LF algorithm.

- 1: **Input:** communication upper bound T , client set \mathcal{N} , number of local epochs K , learning rate η ;
- 2: **Output:** personalized model \tilde{w}_i^T and global model \tilde{w}^T ;
- 3: **for** $t = 0, \dots, T - 1$ **do**
- 4: Sample a set of clients $\mathcal{C}^t \subseteq \mathcal{N}$ with the size of C ;
- 5: **for** each client $i \in \mathcal{C}^t$ in parallel **do**
- 6: Calculate Λ_i^t by (6);
- 7: Select D layers with largest $\xi_{i,l}^t$ values to be the set as $\mathcal{L}_{i,\text{Per}}^t$ and other layers are set as $\mathcal{L}_{i,\text{Uni}}^t$;
- 8: $\tilde{\theta}_{i,\mathcal{L}^{\text{Uni}}}^{t,0} = \tilde{w}_{\mathcal{L}^{\text{Uni}}}^t$ and $\tilde{\theta}_{i,\mathcal{L}^{\text{Per}}}^{t,0} = \tilde{\theta}_{i,\mathcal{L}^{\text{Per}}}^t$
- 9: PLGU($\tilde{\theta}_{i,\mathcal{L}^{\text{Uni}}}^{t,0}, \tilde{w}_{i,\mathcal{L}^{\text{Per}}}^{t,0}, \Lambda_i^t, K, \eta$)
- 10: $\Delta_i^t = \tilde{w}_{i,\mathcal{L}^{\text{Per}}}^{t,K} - \tilde{w}_{i,\mathcal{L}^{\text{Per}}}^{t,0}$;
- 11: $\tilde{w}^{t+1} = \tilde{w}^{t+1} + \frac{1}{C} \sum_{i \in \mathcal{C}^t} \Delta_i^t$;
- 12: **end for**
- 13: **end for**

of personalized layers $l \in \mathcal{L}_i^{\text{Per}}$, and others are into the set of universal layers $l \in \mathcal{L}_i^{\text{Uni}}$. To protect the personalization locally, when client i receives the global model \tilde{w}^t , it will replace the universal layers in the $\tilde{w}_{\mathcal{L}^{\text{Uni}}}^t$ and freeze its personalized layers to obtain a new personalized model $\tilde{\theta}_i^{t,0}$, i.e., $\tilde{\theta}_{i,\mathcal{L}^{\text{Uni}}}^{t,0} = \tilde{w}_{\mathcal{L}^{\text{Uni}}}^t$ and $\tilde{w}_{i,\mathcal{L}^{\text{Per}}}^{t,0} = \tilde{w}_{i,\mathcal{L}^{\text{Per}}}^t$. We call this algorithm for the scheme I as PLGU-Layer Freezing (PLGU-LF). Thus, the objective of PLGU-LF for the scheme I is:

$$\min_{w, \{\theta_i\}_{i=1}^N} \max_{\{\Lambda_i \epsilon_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N F_i(\tilde{w}; \tilde{\theta}_i), \quad (7)$$

where $\tilde{w} = w + \tilde{\epsilon}$, $\tilde{\theta}_i = \theta_i + \tilde{\epsilon}_i$, and $\tilde{\epsilon}_{i,l} = 0, \forall l \in \mathcal{L}_i^{\text{Uni}}$.

We show the learning framework of PLGU-LF with 3 clients to learn a personalized 3-layered network in Figure 2, where each client freezes one personalized layer for personalization, i.e., $D = 1$. In Algorithm 2, we introduce the training procedure of PLGU-LF in detail. As a result, the global model should be more suitable compared to simply using SAM to be the local optimizer, because PLGU-LF algorithm does not lose much universal information sharing across all clients. From the personalization perspective, each client can receive more generalized universal information as well as keep its personal features in order to improve the performance of poor clients. Note that if $|\mathcal{L}^{\text{Per}}| = 0$, PLGU-LF is equal to FedSAM [37]; otherwise, i.e., $|\mathcal{L}^{\text{Per}}| = L$, it is the same as FedAvg [32]. Specifically, suppose that the number of local training epochs of $\tilde{\theta}_i^t$ is equal to \tilde{w}_i^t , the computational cost of PLGU-LF is equal to state-of-the-art pFL in scheme I [26, 44]. Thus, our proposed PLGU-LF does not incur more computational cost.

4.2. Generalization Analysis of PLGU-LF

In this subsection, we aim to analyze the generalization bound of the PLGU-LF algorithm, which

can be presented as follows. Firstly, we define the generalization gap of PLGU-LF algorithm as $|\sum_{i=1}^N \frac{m_i}{m} (\min_{\mathbf{w}, i \in [N]} \max_{\|\epsilon_i\|_2 \leq \rho, i \in [N]} F_i(\tilde{\mathbf{w}}; \tilde{\mathbf{w}}_i) - \min_{\tilde{\mathbf{w}}, \mathbf{w}, i \in [N]} \max_{\|\epsilon_i\|_2 \leq \rho, i \in [N]} \bar{F}_i(\tilde{\mathbf{w}}; \tilde{\mathbf{w}}_i))|$, where $F_i(\cdot; \cdot) = \mathbb{E}_{P_i}[F_{\text{CE}}(\tilde{\mathbf{w}}; \tilde{\mathbf{w}}_i), \mathbf{x}, y]$, P_i denotes the local data distribution of client i and $\bar{F}_i(\cdot; \cdot) = \frac{1}{m_i} \sum_{j=1}^{m_i} [F_{\text{CE}}(\tilde{\mathbf{w}}; \tilde{\mathbf{w}}_i), \mathbf{x}_j, y_j]$ denotes the empirical distribution. The following theorem aims to bound the difference between the empirical and underlying margin-based error for a general deep neural network based on [10, 31, 47, 50]. The bound is based on the spectral norms of the model parameter matrices across layers which provide the upper bound for the Lipschitz and smoothness coefficients of the corresponding neural network.

Theorem 1. *Suppose that the loss function F is β -Lipschitz and the input data \mathbf{x} has ℓ_2 -norm bounded by B . For depth- L and width- d neural networks, suppose that the model parameter matrices in each of the L layers have spectrum norm bounded by M_l . Then, $\forall \gamma \in (0, 1)$, with probability at least $1 - \gamma$, we can upper bound the following generalization gap as $\mathcal{O}\left(\beta \left(\frac{B\sqrt{L} \prod_{l \in \mathcal{L}} M_l}{\gamma\sqrt{m}} + \frac{(B+\rho)d\sqrt{L}\log L \prod_{l \in \mathcal{L}} M_l}{\gamma\sqrt{m}} + \frac{B\sqrt{L} \prod_{l \in \mathcal{L}} M_l}{\gamma N\sqrt{\hat{m}}} + \frac{dD(L-D)\sqrt{\log D} \prod_{l \in \mathcal{L}^{\text{per}}} B M_l \sqrt{\log(L-D)} \prod_{l \in \mathcal{L}^{\text{uni}}} (B+\rho) M_l}{\gamma N\sqrt{\hat{m}}}\right) + \sqrt{m \log \frac{1}{\gamma}}\right)$, where $\hat{m} = \min m_i, \forall i \in [N]$.*

Different from the generalization bound of existing pFL algorithms with two main items [6, 8, 31], the result in Theorem 1 has four main items. The additional items come from the perturbation. The first two items are from the global model, which depends on the total number of data samples m . The third and fourth items are based on the personalized model, and hence they depend on the local dataset and the number of clients N . In addition, the first and third items are due to the marginal-based error [39]. And the second and fourth terms are because of the perturbation (adversarial error), which depends on the number of perturbed layers, i.e., $L - D$. The adversarial error of the global model is on all the layers, and the personalized model depends on the number of universal layers.

5. PLGU for Scheme II

5.1. PLGU-GRep

In this section, we focus on instantiating the PLGU strategy on scheme II without obtaining a global model. Existing studies for scheme II aim to share the learned universal information (usually not the full local models) across all clients, e.g., cluster [46], multi-task learning [42], and [41] hyper-network. However, none of studies discuss whether

the learned universal information is general enough or not. FedRep [5] is one of the most popular pFL scheme II, which is based on representation learning. The motivation of FedRep is that even if local datasets are highly heterogeneous, we can still seek to share the common low-dimensional universal feature representation across all clients. In the main paper, we will set our PLGU strategy on FedRep to improve the generalization of the universal representation and prevent the poor clients, named PLGU-GRep. To show the extensibility, we will embed our PLGU strategy on pFedHN [41], named pFed-GHN, and present the detailed training procedure and generalization analysis in supplementary.

Let θ_ϕ denote the global representation, i.e., universal information, which is a function parameterized by $\phi \in \Phi$, and the specific heads θ_{h_i} , which are parameterized by $h_i \in \mathcal{H}, \forall i \in \mathcal{N}$. Specifically, the personalized model of client i can be decomposed by the low-dimensional header and representation $\theta_i = (\theta_{h_i} \circ \theta_\phi)$. Therefore, the objective for FedRep can be formulated as follows:

$$\min_{\phi \in \Phi} \frac{1}{N} \sum_{i=1}^N \min_{h_i \in \mathcal{H}} F_i(h_i, \phi), \quad (8)$$

where the function $F_i(h_i, \phi) := F_i(\theta_{h_i} \circ \theta_\phi)$. Although the success of straightforwardly leveraging representation to pFL has been demonstrated in [5, 48], they do not consider which layers are more dominant on the universal information in Φ . As such, FedRep cannot guarantee that all clients achieve the desired accuracy due to the biased representation ϕ towards part of clients. Specifically, as shown in the toy example in Figure 1, 22% of clients cannot achieve 72% accuracy using the FedRep algorithm. Therefore, we aim to seek a more generalized representation $\tilde{\phi}^t$ to prevent more clients from falling into poor performance, and the objective function of (8) is re-formulated as $\min_{\tilde{\phi} \in \Phi} \frac{1}{N} \sum_{i=1}^N \min_{h_i \in \mathcal{H}} F_i(h_i, \tilde{\phi})$. Note that the header h_i mainly obtains the personalization for client i , we can update it by K epochs SGD as follows:

$$\mathbf{g}_i^{t,k} = \nabla_{\mathbf{h}_i^{t,k}} F_{B_i}(\mathbf{h}_i^{t,k}, \tilde{\phi}^t), \quad \mathbf{h}_i^{t,k+1} = \mathbf{h}_i^{t,k} - \eta \mathbf{g}_i^{t,k} \quad (9)$$

When the local header h_i updates finish, PLGU-GRep comes into the representation ϕ updates phase. Our design goal is to generalize some specific layers towards personalization. As such, we propose a two-step update to achieve this. Firstly, similar to PLGU-LF, we also explore the personalization score $\Lambda_i^t = \text{diag}(\xi_{i,1}^t, \dots, \xi_{i,L^\phi}^t), \forall l \in \mathcal{L}^{\phi_i}$, where \mathcal{L}^{ϕ_i} is the layers set of ϕ_i with the size of L^ϕ on client i at communication round t , to determine the personalization contribution, i.e.,

$$\xi_{i,l}^t = \frac{\|\tilde{\phi}_{i,l}^t - \tilde{\phi}_l^t\|}{\dim(\tilde{\phi}_l^t)}, \quad (10)$$

Then, by leveraging PLGU strategy in Algorithm 1, we can calculate the layer-wise perturbation $\tilde{\epsilon}_i^t$ and obtain a more

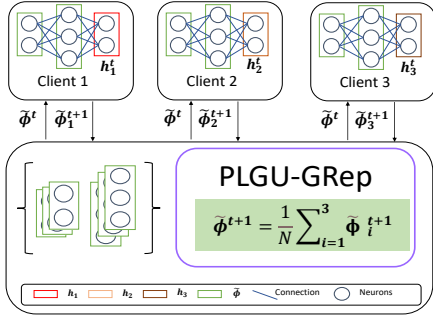


Figure 3. Illustration of PLGU-GRep.

generalized $\tilde{\phi}_i^{t+1}$ representation on client i at communication t as follows:

$$\mathbf{g}_i^t = \nabla_{\tilde{\phi}_i^t} F_{\mathcal{B}_i}(\mathbf{h}_i^{t,K}, \tilde{\phi}_i^t), \quad \phi_i^t = \tilde{\phi}_i^t - \eta \mathbf{g}_i^t, \quad (11)$$

$$\tilde{\epsilon}_i^t = \rho \text{sign}(\nabla_{\phi_i^t} F_{\mathcal{B}_i}(\phi_i^t)) \Lambda_i^t \frac{|\nabla_{\phi_i^t} F_{\mathcal{B}_i}(\phi_i^t)|^{q-1}}{(\|\nabla_{\phi_i^t} F_{\mathcal{B}_i}(\phi_i^t)\|_q^q)^{p-1}}, \quad (12)$$

$$\tilde{\mathbf{g}}_i^t = \nabla_{\phi_i^t} F_{\mathcal{B}_i}(\mathbf{h}_i^{t,K}, \phi_i^t + \tilde{\epsilon}_i^t), \quad \tilde{\phi}_i^t = \phi_i^t - \eta \tilde{\mathbf{g}}_i^t. \quad (13)$$

Based on (10)-(13), we can moderate the personalization in the representation by adding more perturbation. The description of our proposed PLGU-GRep is shown in Figure 3 and its detailed training procedure is illustrated in Algorithm 3. Note that we only use one more SGD step to update the $\tilde{\phi}_i^t$, and hence PLGU-GRep does not incur a huge computational cost.

5.2. Generalization Analysis for PLGU-GRep

Here, we demonstrate the generalization bound of the PLGU-GRep algorithm. Suppose that there exists a global model $F(\mathbf{h}, \tilde{\phi})$, and then obtain the bound by Rademacher complexity. The empirical loss of the global model is $\bar{F}_{\mathcal{D}}(\mathbf{h}, \tilde{\phi}) = \frac{1}{N} \sum_{i=1}^N \frac{m_i}{m} \sum_{j=1}^{m_i} F_{\text{CE}}(\mathbf{x}_j^i, y_j^i; \mathbf{h}_i, \tilde{\phi})$. For the expected loss of $F(\mathbf{V}, \tilde{\phi})$, $F(\mathbf{h}, \tilde{\phi}) = \frac{1}{N} \sum_{i=1}^N \frac{m_i}{m} \mathbb{E}_{P_i} [F_{\text{CE}}(\mathbf{x}, y; \mathbf{h}_i, \tilde{\phi})]$. We assume that the function $F(\cdot)$, $\mathbf{h}_i, \forall i \in [N]$ and $\tilde{\phi}$ are β -, β_h - and β_ϕ -Lipschitz.

Theorem 2. *We assume that the local training model is a depth- L and width- d neural network, the input data sample is bounded by B , the model parameter matrices in each of the layers have spectrum norm bounded by $M_l, \forall l \in \mathcal{L}^\phi$, and the model parameter matrices of header \mathbf{h} can be bounded by M_h . $\forall \gamma \in (0, 1)$ we can bound the generalization gap of PLGU-GRep with probability at least $1 - \gamma$ as $\mathcal{O}\left(\beta \left(\frac{\beta_\phi (B+\rho) d \sqrt{(L-1) \log(L-1)} \prod_{l \in \mathcal{L}^\phi} M_l}{\gamma \sqrt{m}} + \frac{\beta_h B \sqrt{L} M_h}{\gamma N \sqrt{m}} + \frac{\beta_\phi (B+\rho) d \sqrt{(L-1) \log(L-1)} \prod_{l \in \mathcal{L}^\phi} M_l}{\gamma N \sqrt{m}} \right) + \sqrt{m \log \frac{1}{\gamma}}\right)$, where*

Algorithm 3 Scheme II: PLGU-GRep algorithm.

- 1: **Input:** communication upper bound T , client set \mathcal{N} , number of header local epochs K , learning rate η ;
- 2: **Output:** personalized model $\tilde{\theta}_i^T$;
- 3: **for** $t = 0, \dots, T - 1$ **do**
- 4: Sample a set of clients $\mathcal{C}^t \subseteq \mathcal{N}$ with the size of C ;
- 5: Download global representation $\tilde{\phi}^t$ to client $i \in \mathcal{C}^t$;
- 6: **for** each client $i \in \mathcal{C}^t$ in parallel **do**
- 7: **for** $k = 0, \dots, K - 1$ **do**
- 8: Sample mini-batch $\mathcal{B}_i \subset \mathcal{D}_i$;
- 9: Update the header $\mathbf{h}_i^{t,k}$ by (9);
- 10: **end for**
- 11: Sample mini-batch $\mathcal{B}_i \subset \mathcal{D}_i$;
- 12: Update the generalized representation $\tilde{\phi}_i^t$ by (10)-(13);
- 13: **end for**
- 14: Server updates the new representation $\tilde{\phi}^{t+1} = \frac{1}{C} \sum_{i \in \mathcal{C}^t} \tilde{\phi}_i^t$;
- 15: **end for**

$$\hat{m} = \min m_i, \forall i \in [N].$$

Theorem 2 shows insights into the effect of perturbation for PLGU-GRep, and indicates the generalization gap. The first item is based on the aggregation of local generalized representation $\tilde{\phi}_i$, which depends on the number of total data samples m . Based on the local training by SGD on the header \mathbf{h}_i and generalized local representation $\tilde{\phi}_i$, the second and third items depend on the number of data samples on each client \hat{m} and the number of clients N . More specifically, because the third item is related to LWSAM, i.e., with perturbation, it includes the value of ρ , compared to the second item.

6. Experiments

6.1. Experimental Setups

We use CIFAR10, CIFAR100 [21], and Tiny-ImageNet (TmgNet) [23] datasets with various heterogeneous levels and participation rates. We set up 100 clients for CIFAR10 and CIFAR100 datasets, and 20 clients for TmgNet dataset. For the non-IID dataset, we simulate three non-IID data simulations by assigning a fixed number of labels. Note that the default number of labels on CIFAR10 is 3, CIFAR100 is 10, and on TmgNet is 30. The number of data samples on each client is uniformly distributed. We set the number of local epochs $K = 5$, a number of universal layers $|\mathcal{L}^{\text{Uni}}| = 5$ and local batch size as 50, and the perturbation $\rho = 0.05$ by default. We consider two participation schemes $C = 10$ and 100 for CIFAR10 and CIFAR100, respectively.

We compare PLGU-LF, PLGU-GRep, and PLGU-GHN algorithms with several state-of-the-art FL and pFL benchmarks: FedAvg [32], FedSAM [37], Ditto [26], pFedMe

Table 1. Testing accuracy by the global model (averaged, top 5%, and lowest 5% accuracy) under three datasets.

Datasets	C	FedAvg		FedSAM		Ditto		pFedMe		PLGU-LF	
CIFAR10	10	64.79		67.57		64.15		64.36		69.36	
		60.69	68.45	64.76	69.93	62.08	69.45	61.13	70.46	67.34	71.69
	100	67.15		69.49		68.24		68.31		71.48	
		65.99	71.80	67.04	71.12	61.64	71.08	65.25	71.64	68.92	73.57
CIFAR100	10	55.86		57.19		56.35		55.17		59.74	
		51.21	60.75	54.82	59.55	52.73	59.28	50.41	58.35	56.85	61.49
	100	58.00		59.39		58.93		57.24		61.07	
		55.16	60.67	56.73	61.20	56.64	60.79	52.96	60.53	58.25	62.25
TmgNet	20	37.78		38.42		38.09		36.43		40.61	
		32.26	43.73	34.61	42.02	34.75	42.84	31.79	42.90	35.41	43.56

Table 2. Testing accuracy by the personalized model (averaged, top 5%, and lowest 5% accuracy) under three datasets.

Datasets	C	Ditto		pFedMe		FedRep		pFedHN		PLGU-LF		PLGU-GRep		PLGU-GHN	
CIFAR10	10	69.21		66.43		71.32		71.66		71.18		72.87		72.64	
		65.33	72.86	62.97	71.10	68.15	74.04	68.06	73.95	68.22	73.27	69.98	74.91	69.70	74.63
	100	71.23		69.19		75.30		74.95		74.23		76.94		76.17	
		69.93	74.46	66.84	71.58	73.11	77.74	73.23	77.58	72.66	75.87	74.23	77.61	73.79	77.82
CIFAR100	10	59.17		56.58		62.92		63.25		62.33		64.61		65.37	
		57.81	63.06	52.92	60.14	59.70	65.95	59.86	65.79	60.08	65.24	62.58	66.62	63.02	66.94
	100	62.52		58.57		65.08		65.54		64.67		66.79		66.30	
		59.48	64.81	55.79	61.73	62.60	67.96	63.30	68.09	63.72	67.75	64.81	68.59	64.27	68.02
TmgNet	20	39.41		37.22		41.68		41.96		41.39		42.84		42.45	
		35.88	43.49	32.76	42.91	37.53	45.07	37.94	45.19	37.46	44.57	40.29	45.18	39.92	44.89

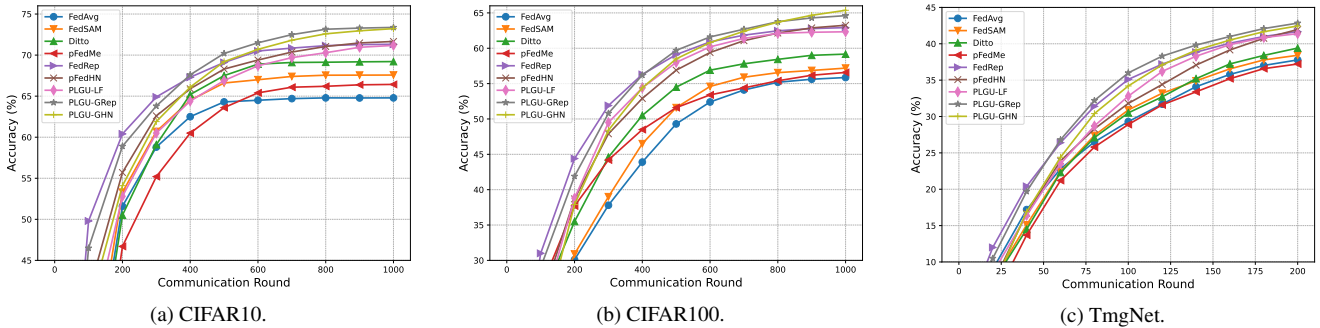


Figure 4. Convergence results evaluation of personalized models under three datasets.

[44], FedRep [5], and pFedHN [41]. To clearly show the performance of our proposed algorithms on different learning models, we use ResNet-18 [13], WideResNet28-10 [54], and ResNet-50 [13] for CIFAR10, CIFAR100, and TmgNet, with group batch. The HN includes 3 fully connected layers. More experimental setups and results will be presented in the supplementary.

6.2. Basic Performance Evaluations

We first evaluate the performance of the proposed algorithms by analyzing the achieved model accuracy on pFL. The results in Table 1 indicate that, compared to other benchmarks on all three datasets, the proposed PLGU-LF algorithm reaches the best global model accuracy (increas-

ing at least 2.55% than others) and does not degrade the top clients. More importantly, we successfully prevent the clients from falling into poor performance, where the 5% lowest clients can achieve 56.85% accuracy with $C = 10$, and increase accuracy by at least 2.03% for the lowest 5% clients on CIFAR100. And the results of the personalized model accuracy evaluation are shown in Table 2, where the proposed algorithms outperform others, e.g., PLGU-GRep achieves 72.64%, 69.70%, and 74.63%, on average, top 5%, and lowest 5% on CIFAR10 with $C = 10$.

To show the learning performance from the convergence perspective, we present the convergence curves in Figure 4. It is easy to observe that the proposed PLGU-GRep and PLGU-GHN achieve the best two convergence speeds.

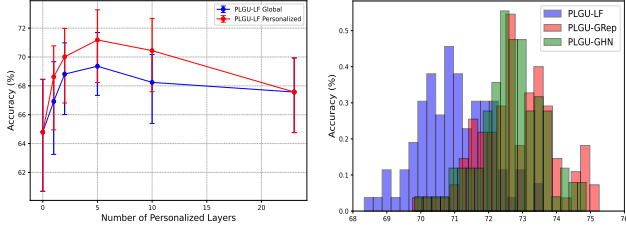


Figure 5. Impact of the size of the personalized layer set D .

Figure 6. Distribution of personalized models.

Table 3. Impact of number of local epochs K with $C = 10$.

K	1		5		10	
PLGU-LF (G)	65.31		69.36		68.42	
	62.08	68.20	67.24	71.13	66.23	71.39
PLGU-LF (P)	67.43		71.18		70.40	
	64.25	70.01	68.22	73.17	65.76	73.54
PLGU-GRep	70.82		72.87		71.65	
	67.25	73.28	69.98	74.91	68.04	74.11
PLGU-GHN	69.93		72.64		72.91	
	66.89	73.15	69.70	74.63	72.91	75.48

6.3. Further Performance Evaluations

Due to the space limitation, we leverage the CIFAR10 for the ablation performance study for all the proposed algorithms. Figure 5 studies the impact of the empirical number of \mathcal{L}^{Per} . Through observation, we can notice that when $D = 5$, the proposed PLGU-LF algorithm achieves the best performance in both personalized and global models. More specifically, when $D = 1$ or 10, the performance does not have obvious degradation. However, when $D = 0$ or 23, it does not achieve the desired performance, which indicates that no perturbation, e.g., FedAvg, and, e.g., full layers perturbation, e.g., FedSAM, are not efficient solutions.

In Figure 6, we show the results of personalized model distribution across all clients under PLGU-LF, -GHN, and -GRep. Compared to the toy example in Figure 1, we can see that our proposed algorithms can significantly decrease the deviation and prevent more poor clients (the accuracy of all clients is larger than 75%), while not clearly reducing the top clients.

We investigate the impact of local epoch number K on the performance in Table 4, which achieve the best when $K = 5$. Note that we use "G" and "P" to represent the global and personal model performance of PLGU-LF. We then explore the impact of ρ on LWSAM in Table 3. The best performance is to set $\rho = 0.05$. In addition, when we increase $\rho = 0.5$, the learning performance incurs large degradation, which matches the results in [9, 29]. Therefore, it is necessary to properly set the values of K and ρ to achieve better performance.

Lastly, to visualize the universal generalization ability of the proposed PLGU-LF algorithm, we show the loss sur-

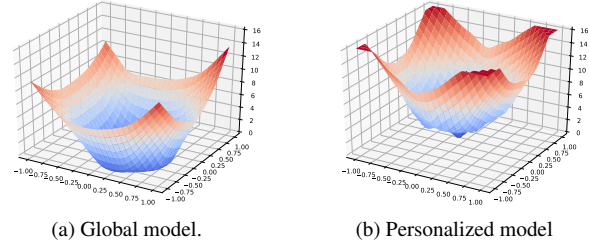


Figure 7. Loss landscapes visualization of PLGU-LF.

Table 4. Impact of perturbation ρ with $C = 10$.

ρ	0.05		0.1		0.5	
PLGU-LF (G)	69.36		68.62		66.73	
	67.24	71.13	65.71	71.88	63.95	68.31
PLGU-LF (P)	71.18		70.96		69.05	
	68.22	73.17	69.15	74.68	66.89	73.50
PLGU-GRep	72.87		72.01		70.80	
	69.98	74.91	71.49	75.93	66.73	73.65
PLGU-GHN	72.64		70.49		67.30	
	69.70	74.63	68.18	74.73	64.84	74.69

faces for the global and personalized models, following the settings in [25]. The results show that the global model is more smooth, i.e., generalizing more universal information, and the personalized model is sharper, i.e., protecting model personalization on clients.

7. Conclusion

In this paper, we propose a novel PLGU strategy to prevent clients from falling into poor performance without obviously downgrading the average and top personalized performance. This strategy aims to generalize the universal information while protecting the personalized features locally. Specifically, we embed the PLGU strategy on two pFL schemes and propose three algorithms, PLGU-LF, PLGU-GRep, and PLGU-GHN by keeping the personalization local and generalizing universal information. Further theoretical investigation indicates that all the PLGU-based algorithms can achieve competitive generalization bounds. The extensive experimental results show that all the proposed algorithms can successfully protect poor clients while not degrading the average learning performance.

Acknowledgement. This work is supported in part by the Natural Science Foundation of Changsha (kq2202108), the Natural Science Foundation of Hunan Province (2022JJ40636), and the National Natural Science Foundation of China (62202293). We are grateful for resources from the High Performance Computing Center of Central South University.

References

- [1] Manoj Ghuhana Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. *arXiv preprint arXiv:1912.00818*, 2019. [2](#), [4](#)
- [2] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017. [12](#)
- [3] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002. [12](#), [15](#)
- [4] Hong-You Chen and Wei-Lun Chao. On bridging generic and personalized federated learning for image classification. In *International Conference on Learning Representations*, 2022. [2](#), [3](#), [14](#)
- [5] Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. Exploiting shared representations for personalized federated learning. In *International Conference on Machine Learning*, pages 2089–2099. PMLR, 2021. [1](#), [2](#), [3](#), [5](#), [7](#)
- [6] Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020. [2](#), [5](#)
- [7] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. *Advances in Neural Information Processing Systems*, 33:3557–3568, 2020. [1](#), [2](#)
- [8] Farzan Farnia, Amirhossein Reisizadeh, Ramtin Pedarsani, and Ali Jadbabaie. An optimal transport approach to personalized federated learning. *arXiv preprint arXiv:2206.02468*, 2022. [5](#)
- [9] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021. [2](#), [3](#), [8](#)
- [10] Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, pages 297–299. PMLR, 2018. [5](#), [11](#)
- [11] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. [14](#)
- [12] Filip Hanzely, Slavomír Hanzely, Samuel Horváth, and Peter Richtárik. Lower bounds and optimal algorithms for personalized federated learning. *Advances in Neural Information Processing Systems*, 33:2304–2315, 2020. [1](#)
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [7](#)
- [14] Yutao Huang, Lingyang Chu, Zirui Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong Zhang. Personalized cross-silo federated learning on non-iid data. In *AAAI*, pages 7865–7873, 2021. [1](#), [2](#)
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. [2](#)
- [16] Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data. *arXiv preprint arXiv:1811.11479*, 2018. [2](#)
- [17] Yihan Jiang, Jakub Konečný, Keith Rush, and Sreeram Kannan. Improving federated learning personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019. [2](#)
- [18] Jean Kaddour, Linqing Liu, Ricardo Silva, and Matt J Kusner. Questions for flat-minima optimization of modern neural networks. *arXiv preprint arXiv:2202.00661*, 2022. [2](#)
- [19] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020. [1](#)
- [20] Nitish Shirish Keskar, Jorge Nocedal, Ping Tak Peter Tang, Dheevatsa Mudigere, and Mikhail Smelyanskiy. On large-batch training for deep learning: Generalization gap and sharp minima. In *5th International Conference on Learning Representations, ICLR 2017*, 2017. [2](#)
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. [6](#), [16](#), [17](#)
- [22] Viraj Kulkarni, Milind Kulkarni, and Aniruddha Pant. Survey of personalization techniques for federated learning. In *2020 Fourth World Conference on Smart Trends in Systems, Security and Sustainability (WorldS4)*, pages 794–797. IEEE, 2020. [2](#)
- [23] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. [6](#), [16](#), [17](#)
- [24] Sunwoo Lee, Tuo Zhang, Chaoyang He, and Salman Avestimehr. Layer-wise adaptive model aggregation for scalable federated learning. *arXiv preprint arXiv:2110.10302*, 2021. [2](#)
- [25] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018. [8](#)
- [26] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *International Conference on Machine Learning*, pages 6357–6368. PMLR, 2021. [1](#), [2](#), [4](#), [6](#)
- [27] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020. [1](#), [2](#)
- [28] Paul Pu Liang, Terrance Liu, Liu Ziyin, Nicholas B Allen, Randy P Auerbach, David Brent, Ruslan Salakhutdinov, and Louis-Philippe Morency. Think locally, act globally: Federated learning with local and global representations. *arXiv preprint arXiv:2001.01523*, 2020. [1](#), [2](#), [3](#), [4](#)
- [29] Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. Towards efficient and scalable sharpness-aware minimization. In *Proceedings of the IEEE/CVF Conference*

- on *Computer Vision and Pattern Recognition*, pages 12360–12370, 2022. 3, 8
- [30] Xiaosong Ma, Jie Zhang, Song Guo, and Wenchao Xu. Layer-wised model aggregation for personalized federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10092–10101, 2022. 1, 2, 3, 14
- [31] Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*, 2020. 1, 2, 5, 11
- [32] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 1, 2, 4, 6
- [33] Maxim Naumov, Dheevatsa Mudigere, Hao-Jun Michael Shi, Jianyu Huang, Narayanan Sundaraman, Jongsoo Park, Xiaodong Wang, Udit Gupta, Carole-Jean Wu, Alisson G Azzolini, et al. Deep learning recommendation model for personalization and recommendation systems. *arXiv preprint arXiv:1906.00091*, 2019. 2
- [34] Jaehoon Oh, SangMook Kim, and Se-Young Yun. Fedbabu: Toward enhanced representation for federated image classification. In *International Conference on Learning Representations*, 2021. 2, 3, 4
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 16
- [36] Zhe Qu, Rui Duan, Lixing Chen, Jie Xu, Zhuo Lu, and Yao Liu. Context-aware online client selection for hierarchical federated learning. *IEEE Transactions on Parallel and Distributed Systems*, 2022. 2
- [37] Zhe Qu, Xingyu Li, Rui Duan, Yao Liu, Bo Tang, and Zhuo Lu. Generalized federated learning via sharpness aware minimization. *arXiv preprint arXiv:2206.02618*, 2022. 2, 3, 4, 6
- [38] Zhe Qu, Xingyu Li, Jie Xu, Bo Tang, Zhuo Lu, and Yao Liu. On the convergence of multi-server federated learning with overlapping area. *IEEE Transactions on Mobile Computing*, 2022. 1
- [39] Amirhossein Reiszadeh, Farzan Farnia, Ramtin Pedarsani, and Ali Jadbabaie. Robust federated learning: The case of affine distribution shifts. In *NeurIPS*, 2020. 5
- [40] Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni, Spyridon Bakas, Mathieu N Galtier, Bennett A Landman, Klaus Maier-Hein, et al. The future of digital health with federated learning. *NPJ digital medicine*, 3(1):1–7, 2020. 1
- [41] Aviv Shamsian, Aviv Navon, Ethan Fetaya, and Gal Chechik. Personalized federated learning using hypernetworks. In *International Conference on Machine Learning*, pages 9489–9502. PMLR, 2021. 1, 2, 3, 5, 7, 14, 15
- [42] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. *Advances in neural information processing systems*, 30, 2017. 5
- [43] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014. 2
- [44] Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes. *Advances in Neural Information Processing Systems*, 33:21394–21405, 2020. 1, 2, 4, 7
- [45] Alysa Ziyang Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2022. 1, 2
- [46] Xueyang Tang, Song Guo, and Jingcai Guo. Personalized federated learning with clustered generalization. *arXiv preprint arXiv:2106.13044*, 2021. 5
- [47] Yu-Lin Tsai, Chia-Yi Hsu, Chia-Mu Yu, and Pin-Yu Chen. Formalizing generalization and adversarial robustness of neural networks to weight perturbations. *Advances in Neural Information Processing Systems*, 34:19692–19704, 2021. 5, 11, 12
- [48] Isidoros Tziotis, Zebang Shen, Ramtin Pedarsani, Hamed Hassani, and Aryan Mokhtari. Straggler-resilient personalized federated learning. *arXiv preprint arXiv:2206.02078*, 2022. 5
- [49] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 2
- [50] Jiancong Xiao, Yanbo Fan, Ruoyu Sun, and Zhi-Quan Luo. Adversarial rademacher complexity of deep neural networks. 2021. 5, 12
- [51] Peng Xiao, Samuel Cheng, Vladimir Stankovic, and Dejan Vukobratovic. Averaging is probably not the optimum way of aggregating parameters in federated learning. *Entropy*, 22(3):314, 2020. 1
- [52] Rudong Xu, Yiting Tao, Zhongyuan Lu, and Yanfei Zhong. Attention-mechanism-containing neural networks for high-resolution remote sensing image classification. *Remote Sensing*, 10(10):1602, 2018. 2
- [53] Yang You, Igor Gitman, and Boris Ginsburg. Scaling sgd batch size to 32k for imagenet training. *arXiv preprint arXiv:1708.03888*, 6(12):6, 2017. 3
- [54] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 7
- [55] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021. 2
- [56] Jie Zhang, Song Guo, Xiaosong Ma, Haozhao Wang, Wenchao Xu, and Feijie Wu. Parameterized knowledge transfer for personalized federated learning. *Advances in Neural Information Processing Systems*, 34:10092–10104, 2021. 2

A. Proof of Theorem 1

Before we show the proof of Theorem 1, we first state some useful lemmas.

Lemma 1. *Let the loss F is β -Lipschitz, \mathcal{G}_u be the hypotheses class for the global model, and \mathcal{G}_p be the hypotheses class for the personalized models. Let $\tilde{\mathbf{w}}^*, \tilde{\mathbf{w}}_i^*$ be the optimal parameters of global and personalized models based on proposed training estimates. Then, with probability at least $1 - \gamma$, we have:*

$$\begin{aligned} & \sum_{i=1}^N \frac{m_i}{m} (F_i(\tilde{\mathbf{w}}^*; \tilde{\mathbf{w}}_i^*) - \bar{F}_i(\tilde{\mathbf{w}}^*; \tilde{\mathbf{w}}_i^*)) \\ & \leq \beta \left(\mathcal{R}_U(\mathcal{G}_u) + \mathcal{R}_U(\tilde{\mathcal{G}}_u) \right) \\ & \quad + \sum_{i=1}^N \frac{m_i}{m} (\mathcal{R}_P(\mathcal{G}_{p_i}) + \mathcal{R}_P(\tilde{\mathcal{G}}_{p_i})) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{m}}. \end{aligned}$$

Proof. This Lemma can be directly obtained by [31, 47], and hence we omit the proof here. Note that the terms of Rademacher complexity $\mathcal{R}_U(\mathcal{G}_u)$ of the global model and $\mathcal{R}_P(\mathcal{G}_{p_i})$ of the personalized model are due to the marginal function, and the terms of $\mathcal{R}_U(\tilde{\mathcal{G}}_u)$ of the global model and $\mathcal{R}_P(\tilde{\mathcal{G}}_{p_i})$ of the personalized model are due to the perturbation. \square

Lemma 2. *Let \mathcal{G}_u be the class of real-valued networks of depth L over the domain \mathcal{X} of the global model, where each parameter \mathbf{w}_l on layer l are at most M_l . Then, we have*

$$\begin{aligned} \mathcal{R}_U(\mathcal{G}_u) & \leq \frac{2B\sqrt{d+1+\log(n)} \prod_{l=1}^L M_l}{\sqrt{m}} \\ & = \mathcal{O}\left(\frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma\sqrt{m}}\right). \end{aligned}$$

Proof. This lemma is revised by [10]. First of all, given the domain of data samples $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\| \leq B\}$ in Euclidean space, we consider (scalar or vector-valued) standard neural networks as follows:

$$\mathbf{x} \rightarrow \mathbf{w}_L \sigma_{L-1}(\mathbf{w}_{L-1} \sigma_{L-2}(\dots \sigma_1(\mathbf{w}_1 \mathbf{x})),$$

where each \mathbf{w}_l is the parameter matrix of l -th layer, and each σ_l is some fixed Lipschitz continuous function between Euclidean spaces, which satisfies $\sigma_l(\mathbf{0}) = \mathbf{0}$. And we also define the $\varepsilon = (\varepsilon_1, \dots, \varepsilon_m)$ as a vector uniformly distributed in $\{-1, +1\}^m$. For a fixed value of $\gamma > 0$, we can upper-bound the Rademacher complexity $\mathcal{R}_U(\mathcal{H}_U)$ as

follows:

$$\begin{aligned} & m\mathcal{R}_U(\mathcal{G}_u) \\ & \leq \frac{1}{\gamma} \log \mathbb{E}_\varepsilon \left[\sup \exp \left(\gamma \sum_{j=1}^m \varepsilon_j \mathbf{w}_{L-1} \sigma_{L-1}(N_{\mathbf{w}_1^{L-1}}(\mathbf{x}_m)) \right) \right] \\ & \leq \frac{1}{\gamma} \log \mathbb{E} \left[M_l \cdot \left\| \gamma \sum_{j=1}^m \sigma_{L-1}(N_{\mathbf{w}_1^{L-1}}(\mathbf{x}_j)) \right\| \right] \\ & \leq \frac{1}{\gamma} \log \left(2^L \cdot \mathbb{E}_\varepsilon \left[\exp M\gamma \left\| \sum_{j=1}^m \varepsilon_j \mathbf{x}_j \right\| \right] \right), \end{aligned} \tag{14}$$

where $M = \prod_{l=1}^L M_l$. We denote $x_{j,r}$ as the r -th coordinate of data sample \mathbf{x}_j . Based on symmetry, take an expectation inside the log term can be bounded as follows:

$$\begin{aligned} & \mathbb{E}_\varepsilon \left(M\gamma \cdot \max_r \left| \sum_{j=1}^m \varepsilon_j x_{j,r} \right| \right) \\ & \leq \sum_{r=1}^R \mathbb{E}_\varepsilon \exp \left(M\gamma \cdot \left| \sum_{j=1}^m \varepsilon_j x_{j,r} \right| \right) \\ & \leq 2 \sum_{r=1}^R \mathbb{E}_\varepsilon \exp \left(M\gamma \sum_{j=1}^m \varepsilon_j x_{j,r} \right) \\ & = 2 \sum_{r=1}^R \prod_{j=1}^m \mathbb{E}_\varepsilon \exp(M\gamma \varepsilon_j x_{j,r}) \\ & \leq 2 \sum_{r=1}^R \exp \left(M^2 \gamma^2 \sum_{j=1}^m x_{j,r}^2 \right) \\ & \leq \frac{L+1+\log(L)}{\gamma} + M^2 \gamma \max_j \sum_{m=1}^M \mathbf{x}_{i,j}^2 \\ & \leq 2R \max_r \exp \left(M^2 \gamma^2 \sum_{j=1}^m x_{j,r}^2 \right). \end{aligned} \tag{15}$$

Plugging (15) to (14), we can obtain the following:

$$\begin{aligned} \mathcal{R}_U(\mathcal{G}_u) & \leq \frac{d+1+\log(R)}{\gamma} + M^2 \gamma \max_r \sum_{j=1}^m x_{j,r}^2 \\ & \leq \frac{2B\sqrt{d+1+\log(n)} \prod_{l=1}^L M_l}{\sqrt{m}} \\ & = \mathcal{O}\left(\frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma\sqrt{m}}\right). \end{aligned}$$

This completes the proof. \square

Lemma 3. *Given the function class \mathcal{G}_u of the global model, and its corresponding adversarial function class function*

$\tilde{\mathcal{G}}_u$, the generalization Rademacher complexity of deep neural networks $\mathcal{R}_U(\tilde{\mathcal{G}}_u)$ can be upper-bounded as follows:

$$\begin{aligned} \mathcal{R}_U(\tilde{\mathcal{G}}_u) &\leq \frac{24 \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\mathbf{x}\| + \rho)}{\gamma\sqrt{m}} \\ &\leq \frac{24 \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\mathbf{x}\| + \rho)}{\gamma\sqrt{m}} \\ &= \mathcal{O}\left(\frac{(B + \rho)d\sqrt{L\log L} \prod_{l=1}^L M_l}{\gamma\sqrt{m}}\right). \end{aligned}$$

Proof: Firstly, we will bound the diameter of adversarial diameter of class function $\tilde{\mathcal{G}}_u$ based on the result in [50] as follows:

$$|\tilde{g}_u(\mathbf{x}_j, y_j)| \leq \prod_{l \in \mathcal{L}} M_l \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\mathbf{x}\| + \rho) \triangleq Z, \quad (16)$$

where $p + q = 1$. Then, we develop the δ -covering number $\mathcal{N}(\tilde{\mathcal{G}}_u, \|\cdot\|, \delta)$ [3], which can be upper-bounded as follows:

$$\mathcal{N}(\tilde{\mathcal{G}}_u, \|\cdot\|, \delta) \leq \prod_{l=1}^L \left(\frac{3M_l}{\delta}\right)^{d_l d_{l-1}} = \left(\frac{3LZ}{2\delta}\right)^{\sum_{l=1}^L d_l d_{l-1}}.$$

The last equality is based on (16). Based on Dudley's integral [2] and with at least $1 - \gamma$ probability, we can obtain the following:

$$\begin{aligned} \mathcal{R}_U(\tilde{\mathcal{G}}_u) &\leq \frac{12}{\gamma\sqrt{m}} \int_0^{Z/2} \sqrt{\mathcal{N}(\tilde{\mathcal{G}}_u, \|\cdot\|, \delta)} d\delta \\ &\leq \frac{12\gamma}{\sqrt{m}} \int_0^{Z/2} \sqrt{\left(\sum_{l=1}^L d_l d_{l-1}\right) \log(3LZ/2\delta)} d\delta \\ &= \frac{12L\sqrt{\sum_{l=1}^L d_l d_{l-1}}}{\gamma\sqrt{m}} \int_0^{1/2} \sqrt{\log(3L/2\delta)} d\delta. \end{aligned} \quad (17)$$

For the integration part in (17), we have

$$\begin{aligned} &\int_0^{1/2} \sqrt{\log(3L/2\delta)} d\delta \\ &= \frac{1}{2} \left(3\sqrt{\pi} \operatorname{erfc}(\sqrt{\log 3L}) + \sqrt{3L}\right) \\ &\leq \frac{1}{2} \left(3\sqrt{\pi} \exp(-\sqrt{\log 3L}) + \sqrt{\log 3L}\right) \\ &= \frac{1}{2} \left(\frac{\sqrt{\pi}}{L} + \sqrt{\log 3L}\right) \leq \sqrt{\log 3L}. \end{aligned} \quad (18)$$

Plugging (18) into (17), we have the following result:

$$\begin{aligned} \mathcal{R}_U(\tilde{\mathcal{G}}_u) &\leq \frac{24 \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\mathbf{x}\| + \rho)}{\gamma\sqrt{m}} \\ &\leq \frac{24 \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\mathbf{x}\| + \rho)}{\gamma\sqrt{m}} \\ &= \mathcal{O}\left(\frac{(B + \rho)d\sqrt{L\log L} \prod_{l \in \mathcal{L}} M_l}{\gamma\sqrt{m}}\right). \end{aligned}$$

This completes the proof. \square

Lemma 4. Suppose that the function class \mathcal{G}_p is the personalized model, and the corresponding adversarial function class is $\tilde{\mathcal{G}}_p$. We simplify the generalization Rademacher complexity of deep neural networks of $\mathcal{R}_P(\mathcal{G}_{p_i}) + \mathcal{R}_P(\tilde{\mathcal{G}}_{p_i})$ as $\mathcal{R}_P(\mathcal{G}_p) + \mathcal{R}_P(\tilde{\mathcal{G}}_p)$ for any client i , and this Rademacher complexity can be upper-bounded as:

$$\begin{aligned} \mathcal{R}_P(\mathcal{G}_p) + \mathcal{R}_P(\tilde{\mathcal{G}}_p) &\leq \frac{2B\sqrt{d+1+\log(n)} \prod_{l=1}^L M_l}{\sqrt{\hat{m}}} \\ &\quad + \frac{24\sqrt{D(L-D)}}{\gamma\sqrt{\hat{m}}} \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\} \|\mathbf{x}\| \\ &\quad + \frac{\sqrt{\sum_{l \in \mathcal{L}^{Uni}} d_l d_{l-1} \log(3D) \prod_{l \in \mathcal{L}^{Uni}} M_l (\|\mathbf{x}\| + \rho)}}{\gamma\sqrt{\hat{m}}} \\ &\quad + \frac{\sqrt{\sum_{l \in \mathcal{L}^{Per}} d_l d_{l-1} \log(3(L-D)) \prod_{l \in \mathcal{L}^{Per}} M_l}}{\gamma\sqrt{\hat{m}}} \\ &= \mathcal{O}\left(\frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma\sqrt{\hat{m}}} + \frac{dD(L-D)\sqrt{\log D}}{\gamma\sqrt{\hat{m}}}\right) \\ &\quad + \prod_{l \in \mathcal{L}^{Per}} B M_l \sqrt{\log(L-D)} \prod_{l \in \mathcal{L}^{Uni}} (B + \rho) M_l. \end{aligned}$$

Proof. The Rademacher complexity of $\mathcal{R}_P(\mathcal{G}_p)$ is easy to be obtained by Lemma 2. Therefore, we mainly focus on the $\mathcal{R}_P(\tilde{\mathcal{G}}_p)$ in this lemma. This proof is inspired by [47]. We re-define the adversarial neural network as follows:

$$\mathbf{x} \rightarrow \mathbf{w}_L \sigma_{L-1}(\mathbf{w}_{L-1} \sigma_{L-2}(\tilde{\mathbf{w}}_{L-2} \sigma_{L-3} \cdots \mathbf{w}_2 \sigma_1)(\mathbf{w}_1 \mathbf{x})).$$

Note that the original layer model is denoted by $\mathbf{w}_l, \forall l \in \mathcal{L}^{Uni}$ and the perturbed layer model is denoted by $\tilde{\mathbf{w}}_l, \forall l \in \mathcal{L}^{Per}$. Then, we can obtain the diameter of $\tilde{\mathcal{G}}_p$ as follows:

$$\begin{aligned} |g_p(\mathbf{x}_j, y_j)| &= |\mathbf{w}_L \sigma_{L-1}(\mathbf{w}_{L-1} \mathbf{x}_j)| \\ &\leq \|\mathbf{w}_L\| \cdot \|\sigma_{L-1}(\mathbf{w}_{L-1} \mathbf{x}_j)\| \\ &= M_L \|\sigma_{L-1}(\mathbf{w}_{L-1} \mathbf{x}_j)\| \\ &\leq M_L \|\mathbf{w}_{L-1} \mathbf{x}_j\| \leq \prod_{l \in \mathcal{L}^{Uni}} M_l \cdot \|\mathbf{x}_j\| \prod_{l \in \mathcal{L}^{Per}} M_l \cdot \|\mathbf{x}_j^*\| \\ &\leq \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\} \prod_{l \in \mathcal{L}^{Per}} M_l \cdot \|\mathbf{x}\| \prod_{l \in \mathcal{L}^{Uni}} M_l \cdot (\|\mathbf{x}\| + \rho). \end{aligned} \quad (19)$$

Due to the fact that the result in (19) includes two parts, we need to re-define the δ -covering number of $\mathcal{N}(\tilde{\mathcal{G}}_p, \|\cdot\|, \delta)$ into two parts: δ^{Uni} -covering $\mathcal{N}(\mathcal{G}_p^{\text{Uni}}, \|\cdot\|, \delta/(L-D))$ and δ^{Per} -covering $\mathcal{N}(\mathcal{G}_p^{\text{Per}}, \|\cdot\|, \delta/D)$. Similar to the Lemma 3, we can obtain the following:

$$\begin{aligned} \mathcal{N}(\tilde{\mathcal{G}}_p^{\text{Uni}}, \|\cdot\|, \delta/(L-D)) &\leq \prod_{l=1}^{L-D} \left(\frac{3(L-D)M_l}{\delta} \right)^{d_l d_{l-1}} \triangleq Z^{\text{Uni}} \\ \mathcal{N}(\mathcal{G}_p^{\text{Per}}, \|\cdot\|, \delta/D) &\leq \prod_{l=1}^D \left(\frac{3(L-D)M_l}{\delta} \right)^{d_l d_{l-1}} \triangleq Z^{\text{Per}}, \end{aligned}$$

Then, we have:

$$\begin{aligned} \mathcal{R}_P(\tilde{\mathcal{G}}_p) &\leq \frac{12}{\sqrt{m}} \int_0^{Z^{\text{Uni}}} \sqrt{\mathcal{N}(\tilde{\mathcal{G}}_p^{\text{Uni}}, \|\cdot\|, \delta/(L-D))} d\delta \\ &\quad \int_0^{Z^{\text{Per}}} \sqrt{\mathcal{N}(\mathcal{G}_p^{\text{Per}}, \|\cdot\|, \delta/D)} d\delta \\ &\leq \frac{12\sqrt{D(L-D)}\gamma}{\sqrt{\hat{m}}} \\ &\quad \int_0^{Z^{\text{Uni}}} \sqrt{\left(\sum_{l \in \mathcal{L}^{\text{Uni}}} d_l d_{l-1} \right) \log(3LZ^{\text{Uni}})} d\delta \\ &\quad \int_0^{Z^{\text{Per}}} \sqrt{\left(\sum_{l \in \mathcal{L}^{\text{Per}}} d_l d_{l-1} \right) \log(3LZ^{\text{Uni}})} d\delta \\ &\leq \frac{12\sqrt{D(L-D)}Z^{\text{Uni}}Z^{\text{Per}}}{\gamma\sqrt{\hat{m}}} \int_0^{\frac{1}{2}} \sqrt{\log(3L/2\delta)} d\delta. \end{aligned} \quad (20)$$

Based on the Lemma 3, we can obtain the following:

$$\begin{aligned} \mathcal{R}_P(\tilde{\mathcal{G}}_p) &\leq \frac{24\sqrt{D(L-D)}}{\gamma\sqrt{m}} \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\} \|\mathbf{x}\| \\ &\quad \sqrt{\sum_{l \in \mathcal{L}^{\text{Uni}}} d_l d_{l-1} \log(3D) \prod_{l \in \text{Uni}} M_l (\|\mathbf{x}\| + \rho)} \\ &\quad \sqrt{\sum_{l \in \mathcal{L}^{\text{Per}}} d_l d_{l-1} \log(3(L-D)) \prod_{l \in \text{Per}} M_l}. \end{aligned} \quad (21)$$

Combing with the (21) and Lemma 2, we have:

$$\begin{aligned} \mathcal{R}_P(\mathcal{G}_p) + \mathcal{R}_P(\tilde{\mathcal{G}}_p) &= \mathcal{O} \left(\frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma\sqrt{\hat{m}}} \right. \\ &\quad \left. + \frac{dD(L-D)\sqrt{\log D} \prod_{l \in \mathcal{L}^{\text{Per}}} BM_l}{\gamma\sqrt{\hat{m}}} \right. \\ &\quad \left. \sqrt{\log(L-D)} \prod_{l \in \mathcal{L}^{\text{Uni}}} (B+\rho)M_l \right). \end{aligned}$$

This completes the proof. \square

By leveraging Lemmas 1-4, we can directly obtain the generalization bound as follows:

$$\begin{aligned} &\sum_{i=1}^N \frac{m_i}{m} (F_i(\tilde{\mathbf{w}}^*; \tilde{\mathbf{w}}_i^*) - \bar{F}_i(\tilde{\mathbf{w}}^*; \tilde{\mathbf{w}}_i^*)) \\ &\leq \beta \left(\mathcal{R}_U(\mathcal{G}_u) + \mathcal{R}_U(\tilde{\mathcal{G}}_u) + \sum_{i=1}^N \frac{m_i}{m} (\mathcal{R}_P(\mathcal{G}_{p_i}) \right. \\ &\quad \left. + \mathcal{R}_P(\tilde{\mathcal{G}}_{p_i})) \right) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{m}} \\ &= \mathcal{O} \left(\beta \left(\frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma\sqrt{m}} + \frac{(B+\rho)d\sqrt{L}\log L \prod_{l \in \mathcal{L}} M_l}{\gamma\sqrt{m}} \right. \right. \\ &\quad \left. \left. + \frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma N\sqrt{\hat{m}}} + \frac{dD(L-D)\sqrt{\log D} \prod_{l \in \mathcal{L}^{\text{Per}}} BM_l}{\gamma N\sqrt{\hat{m}}} \right. \right. \\ &\quad \left. \left. \sqrt{\log(L-D)} \prod_{l \in \mathcal{L}^{\text{Uni}}} (B+\rho)M_l \right) \right) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{m}}. \end{aligned}$$

This completes the proof. \square

B. Proof of Theorem 2

Before proving the Theorem 2, we first state some useful lemmas.

Lemma 5. *We assume that the local training model is a depth- L and width- d neural network, the input data sample is bounded by B , the model parameter matrices in each of the layers have spectrum norm bounded by $M_l, \forall l \in \mathcal{L}^\phi$, and the model parameter matrices of header \mathbf{h} can be bounded by M_h . $\forall \gamma \in (0, 1)$ we can bound the generalization gap of PLGU-GRep with probability at least $1 - \gamma$ as follows:*

$$\begin{aligned} &\frac{1}{N} \sum_{i=1}^N (\bar{F}_i(\mathbf{h}^*; \tilde{\phi}^*) - F_i(\mathbf{h}^*; \tilde{\phi}^*)) \\ &\leq \beta \left(\beta_{\mathbf{h}} \sum_{i=1}^N \mathcal{R}_H(\mathcal{G}_{h_i}) + \beta_\phi (\mathcal{R}_\Phi(\mathcal{G}_\phi) + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)) \right) \\ &\quad + 2\sqrt{\frac{\log \frac{1}{\gamma}}{N}}. \end{aligned}$$

Proof. The proof of this Lemma is similar to Lemma 1 and based on the Lipschitz assumption of β , $\beta_{\mathbf{h}}$ and β_ϕ . Therefore, we omit the details here. Note that the term of Rademacher complexity $\mathcal{R}_H(\mathcal{G}_{h_i})$ is for the header \mathbf{h}_i based on the private dataset of client i and based on the marginal function, the term of $\mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)$ is based on the marginal function of the representation ϕ and $\mathcal{R}_\Phi(\mathcal{G}_\phi)$ of the representation are due to the perturbation. \square

Lemma 6. Let \mathcal{G}_v be the class of real-valued networks of depth 1 over the domain \mathcal{X} of the descriptor, which is trained by SGD and based on local dataset \mathcal{D}_i . We assume that each parameter of the header \mathbf{h} is at most M_h . Then, we have

$$\mathcal{R}_H(\mathcal{G}_h) \leq \frac{2B\sqrt{d+1+\log(n)}M_h}{N\sqrt{\hat{m}}} = \mathcal{O}\left(\frac{BM_h}{\gamma N\sqrt{\hat{m}}}\right).$$

Proof. This proof is similar to Lemma 2, hence we omit here. \square

Lemma 7. Let \mathcal{G}_ϕ be the class of real-valued networks of depth $L-1$ and d over the representation ϕ , where each parameter ϕ_l on layer l are at most M_l . In addition, we assume that the representation can be also bounded by B . Then, we have

$$\begin{aligned} & \mathcal{R}_\Phi(\mathcal{G}_\phi) + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi) \\ & \leq \frac{2B\sqrt{d+1+\log(n)}\prod_{l \in \mathcal{L}\phi} M_l}{N\sqrt{\hat{m}}} \\ & \quad + \frac{24}{\gamma\sqrt{N}} \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\phi\| + \rho) \\ & \quad \cdot \sqrt{\sum_{l \in \mathcal{L}\phi} d^l d^{l-1} \log(3L) \prod_{l \in \mathcal{L}\phi} M_l^h} \\ & = \mathcal{O}\left(\frac{B\sqrt{L_h}\prod_{l=1}^{L_h} M_l^h}{\gamma N\sqrt{\hat{m}}} \right. \\ & \quad \left. + \frac{(B+\rho)d\sqrt{(L-1)\log(L-1)}\prod_{l \in \mathcal{L}\phi} M_l}{\gamma N\sqrt{\hat{m}}}\right). \end{aligned}$$

Proof. The Rademacher complexity of $\mathcal{R}_\Phi(\mathcal{G}_\phi) + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)$ can be easily to extend from Lemmas 2 and 3. Note that the difference is that the size of the representation is $L-1$ (the whole model parameter without including the header). In addition, the result is also related to the number of input data samples \hat{m} and the number of clients N . \square

Based on Lemmas 5-7, we can obtain the generalization bound of PLGU-GRep as follows:

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N (\bar{F}_i(\mathbf{h}^*; \tilde{\phi}^*) - F_i(\mathbf{h}^*; \tilde{\phi}^*)) \\ & \leq \beta \left(\beta_h \sum_{i=1}^N \mathcal{R}_H(\mathcal{G}_{h_i}) + \beta_\phi (\mathcal{R}_\Phi(\mathcal{G}_\phi) + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)) \right) \\ & \quad + 2\sqrt{\frac{\log \frac{1}{\gamma}}{N}} \\ & = \mathcal{O}\left(\beta \left(\frac{\beta_h BM_h}{\gamma N\sqrt{\hat{m}}} + \beta_\phi \left(\frac{B\sqrt{L_h}\prod_{l=1}^{L_h} M_l^h}{\gamma N\sqrt{\hat{m}}} \right. \right. \right. \\ & \quad \left. \left. \left. + \frac{(B+\rho)d\sqrt{(L-1)\log(L-1)}\prod_{l \in \mathcal{L}\phi} M_l}{\gamma N\sqrt{\hat{m}}} \right) \right) \right). \end{aligned}$$

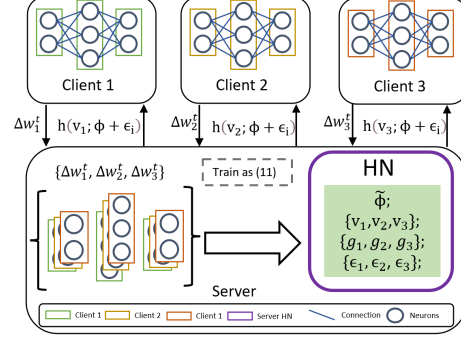


Figure 8. Illustration of PLGU-GHN.

C. PLGU-GHN

C.1. Algorithm of PLGU-GHN

Hypernetwork (HN) [11], where a set of deep neural network parameters is capable of outputting the weights of another network, has been studied in the pFL recently.

Specifically, [41] applies a learnable HyperNetwork (HN) on the server to design a pFL algorithm, called pFedHN. Due to its desirable performance, it has become one of the most popular algorithms. Let $h(\cdot, \phi)$ denote the HN parameterized by ϕ and $F_i(\cdot, \mathbf{w}_i)$ represents the i -th client target network parameterized by the personalized model \mathbf{w}_i . The input of HN is the local model updates $\Delta \mathbf{w}_i^t = \hat{\mathbf{w}}_i^t - \mathbf{w}_i^{t-1}$, where $\hat{\mathbf{w}}_i^{t-1}$ is the personalized model trained by the local dataset of client i . The HN is located at the server and acts on the i -th client descriptor \mathbf{v}_i , which is considered a trainable embedding vector with the personalized model representation. Given \mathbf{v}_i , the personalized model for the i -th client \mathbf{w}_i can be obtained as:

$$\mathbf{w}_i = \mathbf{w}_i(\phi) := h(\mathbf{v}_i, \phi). \quad (22)$$

As such, the objective of scheme II with a learnable HN can be formulated as follows:

$$\min_{\phi, \{\mathbf{v}_i\}_{i=1}^N} \sum_{i=1}^N \frac{m_i}{m} F_i(h(\mathbf{v}_i, \phi)). \quad (23)$$

Although the success of straightforwardly integrating HN to pFL has been demonstrated in [4, 30], the structure of a dedicated vector in HN towards an individual client lacks the consideration of collaborative learning, which cannot guarantee the clients with large distribution deviation to achieve the desired accuracy. Because \mathbf{v}_i includes the independent personalized information of each client, it may not degrade the performance of personalized models. Therefore, the reason for some poor performance clients may be due to the fact that the parameter ϕ is biased to some particular clients, from which ϕ contains the universal information. Note that as HNs are naturally suitable for producing personalized

Algorithm 4 Scheme II: PLGU-GHN algorithm.

- 1: **Input:** communication upper bound T , client set \mathcal{N} , number of local epochs K , learning rate η ;
 - 2: **Output:** personalized model \mathbf{w}_i^T ;
 - 3: **for** $t = 0, \dots, T - 1$ **do**
 - 4: Sample a set of clients $\mathcal{C}^t \subseteq \mathcal{N}$;
 - 5: **for** each client $i \in \mathcal{C}^t$ in parallel **do**
 - 6: set $\mathbf{w}_i^t = h(\mathbf{v}_i^t; \phi^t)$ and $\hat{\mathbf{w}}_i^t = \mathbf{w}_i^t$;
 - 7: **for** $k = 0, \dots, K - 1$ **do**
 - 8: sample mini-batch $\mathcal{B}_i \subset \mathcal{D}_i$;
 - 9: $\mathbf{w}_i^{t,k+1} = \mathbf{w}_i^{t,k} - \eta \nabla_{\mathbf{w}_i^{t,k}} F_{\mathcal{B}_i}(\mathbf{w}_i^{t,k})$;
 - 10: **end for**
 - 11: **end for**
 - 12: $\Delta \mathbf{w}_i^t = \mathbf{w}_i^{t,K} - \mathbf{w}_i^{t,0}$;
 - 13: Calculate \mathbf{g}^t , $\tilde{\phi}^t$, and \mathbf{v}_i^t by (24) and (25);
 - 14: **end for**
-

models, the main challenge of implementing PLGU comes from improving the generalization of HN.

To address this issue, we propose the PLGU-Generalized Hypernetwork (GHN) algorithm, which aims to smooth the parameter ϕ to generalize HN for performing unbiased to all clients. Note that the most important point for smoothing the HN parameter ϕ is that we should consider the impact on all clients. Therefore, the server first collects all personalized model updates $\Delta \mathbf{w}_i$, and trains as follows:

$$\mathbf{g}^t = \frac{1}{|\mathcal{C}^t|} \sum_{i \in \mathcal{C}^t} \nabla_{\phi^t} (\mathbf{w}_i^t)^\top \Delta \mathbf{w}_i^t, \quad \epsilon^t = \rho \frac{\mathbf{g}^t}{\|\mathbf{g}^t\|}, \quad (24)$$

$$\tilde{\phi}^t = \tilde{\phi}^{t-1} - \alpha \nabla_{\phi^t + \epsilon^t} (\mathbf{w}_i^t)^\top \Delta \mathbf{w}_i^t,$$

where α is the learning rate of HN and $\tilde{\phi} = \phi + \epsilon$. The updating of the descriptor \mathbf{v}_i is given by:

$$\mathbf{v}_i^t = \mathbf{v}_i^{t-1} - \alpha \nabla_{\mathbf{v}_i} (\tilde{\phi}^t)^\top \nabla_{\tilde{\phi}^t} (\mathbf{w}_i^t)^\top \Delta \mathbf{w}_i^t. \quad (25)$$

It is worth noting that since the updating of \mathbf{v}_i^t in (25) is based on the perturbed $\tilde{\phi}$, which can less influence the \mathbf{v}_i^t to a particular direction and generate an unbiased result. According to our proposed training steps in (24) and (25), the objective function (23) can be modified as:

$$\min_{\phi, \{\mathbf{v}_i\}_{i=1}^N} \max_{\|\epsilon\| \leq \rho} \sum_{i=1}^N \frac{m_i}{m} F_i(h(\mathbf{v}_i, \phi + \epsilon)). \quad (26)$$

We introduce the implementation of the proposed PLGU-GHN algorithm in Algorithm 4. Lines 4-11 present the local training procedure: For the t -th communication round, the subset of clients \mathcal{C}^t are sampled from \mathcal{N} ; the i -th client in \mathcal{C}^t downloads the latest personalized model \mathbf{w}_i^t which is produced by the HN from the server as in (22); then, the i -th client performs K epochs local SGD based on

its private dataset \mathcal{D}_i ; Line 12 denotes that the local model update $\Delta \mathbf{w}_i$ of client i is uploaded to the server for further improving the generalization of HN. Lastly, Line 13 represents the update of hyper-parameter $\tilde{\phi}$ and the corresponding descriptor \mathbf{v}_i in the generalized HN. In particular, we keep training the personalized model (descriptor) by local SGD because it can maintain the personalization to fit the private dataset. Moreover, after smoothing the HN parameter ϕ , we can obtain more generalized universal information in order to reduce the client variance.

C.2. Generalization Analysis for PLGU-GHN

Here, we investigate the generalization bound of the PLGU-GHN algorithm. Different from pFedHN [41] only focusing on the HN $h(\mathbf{v}, \phi)$ by leveraging the results in [3], we leverage the Rademacher complexity to obtain the result of the whole learning model $F(h(\mathbf{v}, \tilde{\phi}))$. Similar to the definition in PLGU-LF, we denote by $\bar{F}_{\mathcal{D}}(\mathbf{V}, \tilde{\phi})$ the empirical loss of the HN $\bar{F}_{\mathcal{D}}(\mathbf{V}, \tilde{\phi}) = \frac{1}{N} \sum_{i=1}^N \frac{m_i}{m} \sum_{j=1}^{m_i} F_{\text{CE}}(\mathbf{x}_j^i, \mathbf{y}_j^i; h(\mathbf{v}_i, \tilde{\phi}))$. For the expected loss of HN $F(\mathbf{V}, \tilde{\phi})$, we define $F(\mathbf{V}, \tilde{\phi}) = \frac{1}{N} \sum_{i=1}^N \frac{m_i}{m} \mathbb{E}_{P_i} [F_{\text{CE}}(\mathbf{x}, \mathbf{y}; h(\mathbf{v}_i, \tilde{\phi}))]$. We assume that $\mathbf{v}_i, \forall i \in [N]$ and $\tilde{\phi}$ are β_V - and β_h -Lipschitz.

Theorem 3. *We assume that the local training model for the descriptor is a depth- L and width- d neural network, the input data sample is bounded by B , and the model parameter matrices in each of the L layers have spectrum norm bounded by M_l . Suppose that the training model of HN is a depth- L_h and width- d_h neural network, the descriptor is bounded by R , and the model parameter matrices of HN are bounded by M_l^h . $\forall \gamma \in (0, 1)$ we can bound the generalization gap of PLGU-GHN with probability at least $1 - \gamma$ as follows:*

$$\mathcal{O} \left(\beta_V \left(\frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma N \sqrt{\hat{m}}} \right) + \beta_h \left(\frac{R\sqrt{L_h} \prod_{l=1}^{L_h} M_l^h}{\gamma \sqrt{N}} + \frac{(R + \rho)d_h \sqrt{L_h \log L_h} \prod_{l=1}^{L_h} M_l^h}{\gamma \sqrt{N}} \right) \right) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{N}},$$

where $\hat{m} = \min m_i, \forall i \in [N]$.

Theorem 3 shows insights into the parameter-sharing effect of PLGU-GHN. The first part $\beta_V \frac{B\sqrt{L} \prod_{l=1}^L M_l}{\gamma N \sqrt{\hat{m}}}$ is based on the local training by SGD for training the descriptor \mathbf{v}_i , which depends on the number of data samples on each client \hat{m} and the number of clients N . On the other hand, the second part $\beta_h \left(\frac{R\sqrt{L_h} \prod_{l=1}^{L_h} M_l^h}{\gamma \sqrt{N}} + \frac{(R + \rho)d_h \sqrt{L_h \log L_h} \prod_{l=1}^{L_h} M_l^h}{\gamma \sqrt{N}} \right)$ is due to updating the HN. It includes two items: the first one is due to the marginal function and the second is due to the

perturbation. Note that this part depends on the number of clients N and the size of descriptor R , where it does not rely on the number of data samples since the parameters of HN are not shared across clients.

Before proving the generalization bound of Theorem 3, we firstly propose the following lemma for pFedGHN, which is similar to Lemma 1:

Lemma 8. *Let the loss F , $h(\mathcal{V})$ and $h(\Phi)$ are β , β_v and β_h -Lipschitz, \mathcal{G}_V be the hypotheses class for the descriptor that is trained by local model updates with the private dataset \mathcal{D}_i , and \mathcal{G}_ϕ be the hypotheses class for the HN. Let \mathbf{v}_i^* be the optimal parameters of descriptor and $\tilde{\phi}^*$ be the optimal parameters of HN based on proposed training estimates. Then, with probability at least $1 - \gamma$, we have:*

$$\begin{aligned} & \sum_{i=1}^N \frac{m_i}{m} (\bar{F}_i(\mathbf{v}^*; \tilde{\phi}^*) - F_i(\mathbf{v}^*; \tilde{\phi}^*)) \\ & \leq \beta \left(\beta_v \sum_{m_i}^m \mathcal{R}_V(\mathcal{G}_{v_i}) + \beta_h (\mathcal{R}_\Phi(\mathcal{G}_\phi) \right. \\ & \quad \left. + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi) \right) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{N}}. \end{aligned}$$

Proof. The proof of this Lemma is similar to Lemma 1 and based on the Lipschitz assumption of β , β_v and β_h . Therefore, we omit the details here. Note that the term of Rademacher complexity $\mathcal{R}_V(\mathcal{G}_{v_i})$ is for the descriptor v_i based on the private dataset of client i and based on the marginal function, the term of $\mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)$ is based on the marginal function of the HN and $\mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)$ of the HN are due to the perturbation. \square

Lemma 9. *Let \mathcal{G}_v be the class of real-valued networks of depth L over the domain \mathcal{X} of the descriptor, which is trained by SGD and based on local dataset \mathcal{D}_i . We assume that each parameter w_l on layer l are at most M_l . Then, we have*

$$\begin{aligned} \mathcal{R}_V(\mathcal{G}_v) & \leq \frac{2B\sqrt{d+1+\log(n)}\prod_{l=1}^L M_l}{N\sqrt{\hat{m}}} \\ & = \mathcal{O}\left(\frac{B\sqrt{L}\prod_{l=1}^L M_l}{\gamma\sqrt{\hat{m}}}\right). \end{aligned}$$

Proof. This Lemma can be directly obtained by Lemma 2. Therefore, we omit the details. \square

Lemma 10. *Let \mathcal{G}_ϕ be the class of real-valued networks of depth L_h and d_h over the descriptor, where each parameter ϕ_l on layer l are at most M_l^h . In addition, we assume that*

the descriptor can be bounded by R . Then, we have

$$\begin{aligned} & \mathcal{R}_\Phi(\mathcal{G}_\phi) + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi) \\ & \leq \frac{2R\sqrt{d_h+1+\log(n)}\prod_{l=1}^{L_h} M_l}{\sqrt{N}} \\ & \quad + \frac{24}{\gamma\sqrt{N}} \max\{1, q^{\frac{1}{2}-\frac{1}{p}}\}(\|\mathbf{v}\| + \rho) \\ & \quad \sqrt{\sum_{l=1}^{L_h} d_h^l d_h^{l-1} \log(3L_h) \prod_{l=1}^{L_h} M_l^h} \\ & = \mathcal{O}\left(\frac{R\sqrt{L_h}\prod_{l=1}^{L_h} M_l^h}{\gamma\sqrt{N}} \right. \\ & \quad \left. + \frac{(R+\rho)d_h\sqrt{L_h\log L_h}\prod_{l=1}^{L_h} M_l^h}{\gamma\sqrt{N}}\right). \end{aligned}$$

Proof. The Rademacher complexity of $\mathcal{R}_\Phi(\mathcal{G}_\phi) + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi)$ can be easily to extend from Lemmas 2 and 3. Note that the difference is that the input of HN is \mathbf{v} , and hence the number of input data samples are equal to the number of clients N . \square

By leveraging Lemmas 8-10, we can directly obtain the generalization bound as follows:

$$\begin{aligned} & \sum_{i=1}^N \frac{m_i}{m} (\bar{F}_i(\mathbf{v}^*; \tilde{\phi}^*) - F_i(\mathbf{v}^*; \tilde{\phi}^*)) \\ & \leq \beta \left(\beta_v \sum_{m_i}^m \mathcal{R}_V(\mathcal{G}_{v_i}) + \beta_h (\mathcal{R}_\Phi(\mathcal{G}_\phi) \right. \\ & \quad \left. + \mathcal{R}_\Phi(\tilde{\mathcal{G}}_\phi) \right) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{m}} \\ & = \mathcal{O}\left(\beta \left(\beta_v \frac{B\sqrt{L}\prod_{l=1}^L M_l}{\gamma N\sqrt{\hat{m}}} \right) + \beta_h \left(\frac{R\sqrt{L_h}\prod_{l=1}^{L_h} M_l^h}{\gamma\sqrt{N}} \right. \right. \\ & \quad \left. \left. + \frac{(R+\rho)d_h\sqrt{L_h\log L_h}\prod_{l=1}^{L_h} M_l^h}{\gamma\sqrt{N}} \right) \right) + 2\sqrt{\frac{\log \frac{1}{\gamma}}{N}}. \end{aligned}$$

This completes the proof. \square

D. Experimental Setups

D.1. Datasets and Neural Networks

All our experiments are conducted on a CPU/GPU cluster with 4 NVIDIA V100 GPUs. We use PyTorch [35] to train all the pFed algorithms. We use three benchmark datasets, CIFAR10, CIFAR100 [21] and Tiny-ImageNet [23], which represent a different number of data samples and labels. The detailed descriptions of datasets and neural networks are introduced in Table 5.

Dataset	Labels	Clients	Total samples	Neural network
CIFAR10 [21]	10	100	60,000	ResNet-18
CIFAR100 [21]	100	100	60,000	WideResNet28-10
Tiny-ImageNet (TmgNet) [23]	200	100	100,000	ResNet-50

Table 5. Descriptions of datasets and neural networks.

Table 6. Testing accuracy by the global model (averaged, top 5%, and lowest 5% accuracy) under three datasets (5 labels on CIFAR10, 20 labels on CIFAR100 and 50 labels on TmgNet).

Datasets	C	FedAvg	FedSAM	Ditto	pFedMe	PLGU-LF
CIFAR10	10	67.92	69.43	68.67	67.81	71.10
		63.85 71.69	66.73 73.02	65.13 72.25	63.51 72.57	68.63 72.17
	100	69.57	70.86	70.18	68.45	71.40
		66.79 71.53	68.26 73.49	66.98 72.52	65.31 71.80	69.01 73.33
CIFAR100	10	59.08	60.24	59.37	56.22	61.38
		56.65 63.12	57.10 62.73	56.94 63.03	53.85 60.92	59.05 63.30
	100	60.71	61.35	60.94	58.15	62.09
		58.50 63.24	59.75 63.91	58.30 64.02	56.76 61.48	60.97 64.84
TmgNet	20	39.12	40.17	39.35	37.79	42.21
		34.48 43.85	37.93 44.15	36.04 42.28	32.89 41.97	39.07 44.54

Table 7. Testing accuracy by the personalized model (averaged, top 5%, and lowest 5% accuracy) under three datasets (5 labels on CIFAR10, 20 labels on CIFAR100 and 50 labels on TmgNet).

Datasets	C	Ditto	pFedMe	FedRep	pFedHN	PLGU-LF	PLGU-GRep	PLGU-GHN
CIFAR10	10	71.42	68.79	73.43	73.07	72.82	74.34	74.18
		68.41 74.56	67.37 73.05	69.89 75.53	70.18 75.09	70.20 74.35	72.67 75.61	72.13 75.26
	100	73.19	71.14	74.66	74.73	74.75	75.42	75.25
		70.51 75.70	69.88 74.43	72.46 76.42	73.25 76.76	73.59 76.62	74.13 77.04	74.19 76.89
CIFAR100	10	61.19	60.62	63.58	63.70	63.19	64.93	65.02
		58.16 64.73	57.48 63.75	60.93 66.02	61.31 65.84	61.38 65.12	63.40 66.86	63.89 66.60
	100	62.88	61.45	64.64	64.98	64.30	66.71	66.46
		61.15 65.20	59.86 64.69	62.76 66.14	62.91 66.20	62.43 65.90	64.72 68.06	64.21 68.37
TmgNet	20	40.65	39.18	42.82	43.00	42.69	43.49	43.17
		37.51 43.92	34.73 42.95	40.05 45.23	40.47 45.36	40.62 45.04	41.19 45.70	41.03 45.58

D.2. Hyper-parameter Settings

For pFedMe, FedAvg, and FedSAM algorithms corresponding to the FL settings are used as an initialization. It is then trained about a quarter of FedAvg training with a learning rate of 0.01 on CIFAR10 and CIFAR100 datasets, 0.02 for the TmgNet dataset. For FedRep, the number of local epochs for head and body updates were set to τ and 1, respectively. For Ditto, the λ value used to control the regularization term was set to 0.75. For the pFedHN algorithm, we use a fully-connected HN with 5 hidden layers of 100 hidden units each.

We pre-allocate all training samples for validation on the three datasets. The validation sets are used for hyperparameter tuning by grid search. We search over learning rate $\{0.001, 0.005, 0.01, 0.015, 0.02, 0.03\}$.

D.3. Architecture of HN

Figure 8 shows the training procedure of the PLGU-GHN algorithm. And Figure 9 presents the architecture of HN for pFedHN and pFedGHN algorithms in our experiments. First, the hypernetwork uses an embedding layer to generate and update the descriptor v_i . Then, the descriptor vector v_i passes through several public fully connected layers and generates an intermediate feature. Finally, the intermediate feature is sent to 3 fully connected layers which correspond to 3 layers of local clients’ personalized models. The final fully connected layer output the aggregation results for the personalized model of client i .

E. Additional Experimental Results

E.1. Additional Basic Performance

In this subsection, we show the performance on different heterogeneous pFL settings of our proposed PLGU-LF,

Table 8. Testing accuracy by the personalized model (averaged, top 5%, and lowest 5% accuracy) under three datasets (8 labels on CIFAR10, 50 labels on CIFAR100 and 100 labels on TmgNet).

Datasets	C	FedAvg		FedSAM		Ditto		pFedMe		PLGU-LF	
CIFAR10	10	69.03		69.92		69.69		68.76		71.40	
		67.80	72.09	67.54	72.32	66.91	72.19	65.78	72.45	68.60	73.62
	100	71.07		71.48		71.24		70.43		72.12	
		68.23	72.56	69.91	72.58	69.15	72.42	68.10	72.48	71.02	73.98
CIFAR100	10	60.79		62.02		61.41		58.30		63.06	
		57.29	64.15	58.97	64.18	58.44	63.89	54.31	62.25	60.14	64.82
	100	62.43		63.16		62.47		60.52		63.06	
		60.41	64.59	61.08	64.94	60.25	64.38	57.59	63.02	61.94	65.20
TmgNet	20	40.36		41.24		40.51		38.93		43.05	
		36.18	45.43	38.20	44.09	37.65	44.18	34.86	42.94	40.75	45.62

Table 9. Testing accuracy by the personalized model (averaged, top 5%, and lowest 5% accuracy) under three datasets (8 labels on CIFAR10, 20 labels on CIFAR100 and 40 labels on TmgNet).

Datasets	C	Ditto		pFedMe		FedRep		pFedHN		PLGU-LF		PLGU-GRep		PLGU-GHN	
CIFAR10	10	72.04		70.74		75.36		75.89		75.01		76.58		76.79	
		70.09	74.73	67.32	73.98	73.25	78.04	74.16	77.39	74.12	76.58	74.64	78.45	74.61	78.70
	100	73.98		73.04		76.71		76.60		76.25		77.18		77.51	
		73.26	77.80	66.38	72.41	73.76	78.49	73.97	78.52	74.08	77.76	76.30	79.46	75.88	79.63
CIFAR100	10	63.11		61.68		64.61		64.25		63.95		65.13		65.94	
		61.79	64.43	59.36	62.60	62.92	66.59	62.85	66.57	62.13	64.96	65.10	66.41	64.79	66.03
	100	64.39		63.07		66.24		66.51		65.99		67.25		67.06	
		63.68	66.96	61.73	65.13	64.95	67.12	65.15	66.88	64.73	68.50	66.39	68.05	66.04	68.20
TmgNet	10	41.26		40.02		43.37		43.49		43.11		43.97		43.72	
		38.53	44.22	36.14	43.45	40.37	45.78	40.94	45.59	40.66	45.20	42.16	45.82	41.98	46.01

PLGU-GRep, and PLGU-GHN algorithms.

The additional experimental results of heterogeneous pFL settings are presented in Tables 6-9. Firstly, we can see that the less heterogeneous can improve the learning performance for all algorithms, which matches the existing results. Moreover, the client variance also reduces in less heterogeneous pFL settings. In particular, when each client includes 50 labels on CIFAR100, the learning performance of FedRep and pFedHN is higher than our proposed algorithms. In addition, when the pFL setting is less heterogeneous, e.g., 100 labels on TmgNet, the averaged learn-

ing accuracy does not have obvious improvement compared to 50 labels. It is because the distribution of the training dataset is similar, and nearly to the performance of conventional central machine learning. However, our proposed PLGU strategy can also prevent the poor clients. Therefore, we consider that our proposed PLGU-LF, PLGU-GRep, and PLGU-GHN are more suitable for highly heterogeneous pFL settings. More specifically, the performance of PLGU-GHN is worse than PLGU-GRep is because the architecture of HN is much simpler than local model.

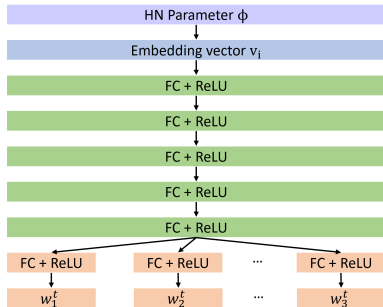


Figure 9. HN architecture of PLGU-GHN.

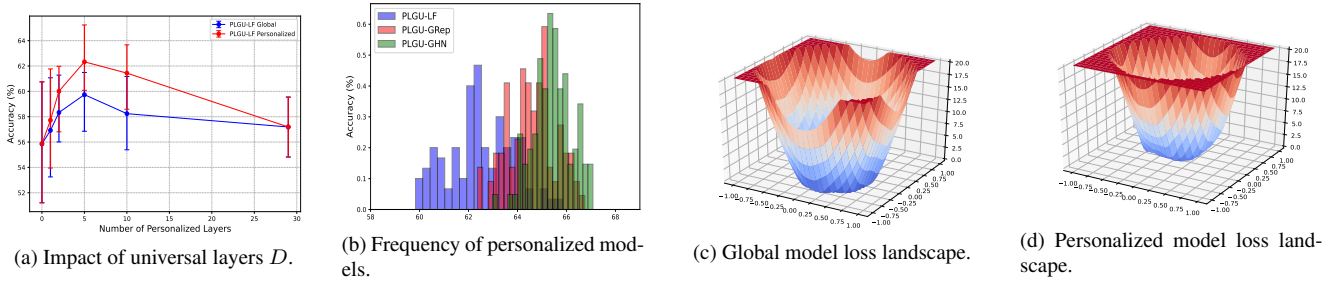


Figure 10. Further performance on CIFAR100 dataset.

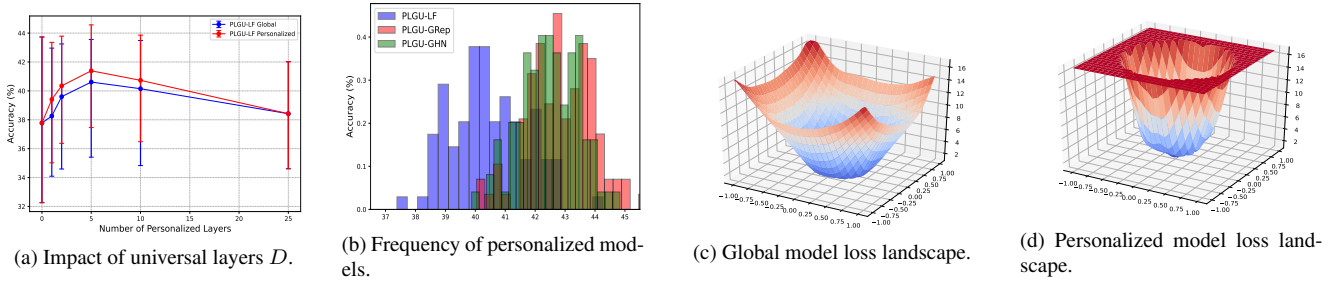


Figure 11. Further performance on TmgNet dataset.

Table 10. Impact of number of local epochs K with $C = 10$ on CIFAR100 and TmgNet datasets.

K	CIFAR100						TmgNet					
	1		5		10		1		5		10	
PLGU-LF (G)	56.28		59.74		59.21		36.19		40.61		41.04	
	51.93	59.77	56.85	61.49	56.42	61.08	32.40	40.12	35.41	43.56	35.83	43.95
PLGU-LF (P)	59.65		62.33		61.82		38.57		41.39		41.94	
	57.35	62.76	60.18	64.54	59.33	64.69	34.69	40.92	37.46	44.57	37.35	45.13
PLGU-GRep	59.12		62.61		62.04		39.81		42.84		43.35	
	56.90	63.25	62.58	66.62	61.02	65.57	37.68	43.06	40.29	45.18	39.93	44.14
PLGU-GHN	60.68		64.97		64.34		40.30		42.45		43.07	
	59.31	63.16	63.02	67.94	62.70	67.89	37.76	43.39	39.92	44.89	39.43	44.50

Table 11. Impact of the value of perturbation ρ with $C = 10$ on CIFAR100 and TmgNet datasets.

ρ	CIFAR100						TmgNet					
	0.05		0.1		0.5		0.05		0.1		0.5	
PLGU-LF (G)	60.74		60.31		59.25		40.61		40.93		40.04	
	57.85	62.49	57.06	61.87	58.43	59.96	35.41	43.56	35.89	43.90	37.02	41.52
PLGU-LF (P)	62.63		61.92		60.45		41.39		41.76		40.97	
	61.18	65.54	60.75	64.02	59.86	63.27	37.46	44.57	37.93	45.00	37.12	43.84
PLGU-GRep	65.41		65.03		63.74		42.84		43.17		42.03	
	63.58	67.62	63.15	67.59	62.70	67.42	40.29	45.18	40.64	45.45	39.36	44.13
PLGU-GHN	64.97		64.23		63.04		42.45		42.90		41.38	
	63.02	67.94	62.46	67.63	61.29	66.62	39.92	44.89	40.18	45.03	39.30	44.17