

# FedDWA: Personalized Federated Learning with Dynamic Weight Adjustment\*

Jiahao Liu<sup>1,2</sup>, Jiang Wu<sup>1,2</sup>, Jinyu Chen<sup>1,2</sup>, Miao Hu<sup>1,2</sup>, Yipeng Zhou<sup>3</sup>, Di Wu<sup>1,2</sup>†

<sup>1</sup>School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou, China

<sup>2</sup>Guangdong Key Laboratory of Big Data Analysis and Processing, Guangzhou, China

<sup>3</sup>School of Computing, Faculty of Science and Engineering, Macquarie University, Sydney, Australia  
 { liujh69, wujiang7, chenjy585 }@mail2.sysu.edu.cn  
 { humiao5, wudi27 }@mail.sysu.edu.cn, yipeng.zhou@mq.edu.au

## Abstract

Different from conventional federated learning, personalized federated learning (PFL) is able to train a customized model for each individual client according to its unique requirement. The mainstream approach is to adopt a kind of weighted aggregation method to generate personalized models, in which weights are determined by the loss value or model parameters among different clients. However, such kinds of methods require clients to download others' models. It not only sheer increases communication traffic but also potentially infringes data privacy. In this paper, we propose a new PFL algorithm called *FedDWA (Federated Learning with Dynamic Weight Adjustment)* to address the above problem, which leverages the parameter server (PS) to compute personalized aggregation weights based on collected models from clients. In this way, FedDWA can capture similarities between clients with much less communication overhead. More specifically, we formulate the PFL problem as an optimization problem by minimizing the distance between personalized models and guidance models, so as to customize aggregation weights for each client. Guidance models are obtained by the local one-step ahead adaptation on individual clients. Finally, we conduct extensive experiments using five real datasets and the results demonstrate that FedDWA can significantly reduce the communication traffic and achieve much higher model accuracy than the state-of-the-art approaches.

## 1 Introduction

Federated Learning (FL), as an emerging distributed machine learning paradigm, allows decentralized clients to collaboratively train a global machine learning model without exposing their private data [McMahan *et al.*, 2017]. However, one of the most challenging problems confronted by

FL is the performance degradation caused by the heterogeneity of data distribution on decentralized clients [Li *et al.*, 2020b]. Specifically, data distribution on clients is non-independent and identically distributed (non-IID) such that a single global model cannot meet personalized needs of all clients. It has been reported in [Li *et al.*, 2020b; Yu *et al.*, 2020] that data heterogeneity can result in slow convergence and poor model accuracy. For example, the global next-word prediction model trained by FedAvg [McMahan *et al.*, 2017] is not always effective for all clients because of their personalized habits. Such a single global model may significantly deviate from personalized optimal models [Yu *et al.*, 2020].

To address the above problem, *Personalized Federated Learning (PFL)* has been proposed and studied in [Smith *et al.*, 2017; T Dinh *et al.*, 2020; Fallah *et al.*, 2020; Li *et al.*, 2021b; Collins *et al.*, 2021]. PFL aims to handle non-IID data distribution by training personalized models for each client, so as to improve model accuracy. In essence, PFL can either incorporate personalized components into the global model or train multiple models to obtain personalized models. For example, the works [T Dinh *et al.*, 2020; Li *et al.*, 2021b] added regularization terms to the global model in order to train personalized models. A distance metric to shrink the search space of personalized models around the global model is applied. However, such an approach fails to optimize personalized models because distance metrics usually cannot exactly capture the heterogeneity of data distribution among clients. Later on, more radical approaches (e.g., [Zhang *et al.*, 2021b; Li *et al.*, 2022]) were proposed, which train multiple models to meet personalized requirements by distributing other clients' models to each individual client. Based on its local dataset, each client can decide how to aggregate models from other clients to obtain a personalized model. Despite that the performance of PFL is improved, this approach will make communication traffic explode and users' privacy may be compromised.

In this paper, we propose a novel PFL method called *FedDWA (Federated Learning with Dynamic Weight Adjustment)*, which can improve PFL performance by encouraging collaborations among clients with similar data distributions. In existing works [Zhang *et al.*, 2021b; Li *et al.*, 2022], each client needs to collect models from all other clients and evaluate similarities between clients with extra local validation

\*An extended version of this paper (with the Appendix included) can be found in <http://arxiv.org/abs/2305.06124>.

†Corresponding author.

sets. Different from these works, our framework characterizes client similarity in an analytical way instead of empirical searching via the validation dataset. In addition, there is no need to share local models among clients in FedDWA avoiding excessive communication traffic and potential privacy leakage [Hu *et al.*, 2021]. In FedDWA, individual clients can obtain personalized models computed by the PS based on collected model parameters and guidance models from clients. Guidance models are obtained using the one-step ahead adaptation method by individual clients. Based on guidance models, the PS can tune aggregation weights to minimize the distance between each model with its guidance model. Thus, FedDWA can improve PFL performance without incurring heavy overhead by avoiding exchanging information between clients.

In summary, our main contributions in this paper can be summarized as follows:

- We propose a new personalized federated learning framework called FedDWA. FedDWA can effectively exploit clients owning data with a similar distribution to improve personalized model accuracy.
- We theoretically analyze the properties of the FedDWA algorithm, and show how the weights are dynamically adjusted to achieve personalization.
- By conducting experiments using five real datasets, we demonstrate that FedDWA outperforms other methods under three heterogeneous FL settings.

## 2 Related Work

In this section, we discuss related works from two perspectives: *data-based PFL* and *model-based PFL*. Data-based PFL focuses on reducing data heterogeneity among clients while model-based PFL focuses on designing a personalized model for each client. Typically, data-based PFL shares a global dataset that is balanced across all clients [Zhao *et al.*, 2018] (or private statistical information [Shin *et al.*, 2020; Yoon *et al.*, 2021] among clients) to realize personalized learning. However, sharing a global dataset may potentially violate privacy policies since it is at the risk of privacy leakage. To address the above problem, model-based PFL was proposed, which can be divided into two types: *single-model PFL* and *multi-model PFL*.

Most single-model PFL methods are extensions of conventional FL algorithms (e.g., FedAvg [McMahan *et al.*, 2017]). For example, FedProx [Sahu *et al.*, 2020] employed a proximal term to formulate clients’ optimization objectives so as to mitigate the adverse influence of systematic and statistical heterogeneity on FL. FedAvg<sub>FT</sub> and FedProx<sub>FT</sub> [Wang *et al.*, 2019] obtained personalized models by fine tuning the global model generated by FedAvg and FedProx, respectively. FedAvgM proposed by [Hsu *et al.*, 2019] adopted a momentum method to update the global model, so as to alleviate the adverse influence of non-IID data distribution on FL. An alternative single-model approach for PFL is based on meta-learning. Recent works [Khodak *et al.*, 2019; Yue *et al.*, 2021; Acar *et al.*, 2021] extended Model Agnostic Meta-learning (MAML) for FL under non-IID data distribu-

tion. However, the personalized learning ability of single-model methods is limited because it is hard to fit all heterogeneous data distributions with a single model very well.

Multi-model methods outperform single-model methods by training multiple models to better adapt to the personalized requirements of clients. Cluster FL [Sattler *et al.*, 2021; Ghosh *et al.*, 2020; Mansour *et al.*, 2020] assumed that clients can be partitioned into multiple clusters, and clients are grouped based on loss values or gradients. A customized model can be trained for each cluster. However, cluster-based client grouping may not be able to effectively improve PFL performance. FedEM [Marfoq *et al.*, 2021] refined the cluster-based client group method by proposing a soft client clustering algorithm. However, it requires each client to download multiple models, which can considerably increase the communication overhead.

Other than clustering clients, more advanced multi-model PFL methods were developed including additive model mixture between local and global models (such as L2GD [Hanzely and Richtárik, 2020] and APFL [Deng *et al.*, 2020]), multi-task learning methods with model similarity penalization (such as MOCHA [Smith *et al.*, 2017], pFedMe [T Dinh *et al.*, 2020] and Ditto [Li *et al.*, 2021b]). More PFL methods were developed by leveraging Gaussian processes [Achituve *et al.*, 2021] and knowledge transfer [Zhang *et al.*, 2021a]. However, these methods inevitably need public shared data or inducing points set. It is also possible to achieve PFL by decomposing FL models into a global part and multiple personalized parts. Inspired by representation learning, the works [Arivazhagan *et al.*, 2019; Collins *et al.*, 2021; Tan *et al.*, 2022; Chen and Chao, 2022; Oh *et al.*, 2021; Mills *et al.*, 2022] decomposed the FL model into a shared feature extractor part and a personalized part to realize PFL. Nevertheless, how to decompose FL models is only heuristically designed and discussed by existing works.

PFL can be achieved by customizing weights of model aggregation for each client as well, e.g., FedAMP [Huang *et al.*, 2021], FedFomo [Zhang *et al.*, 2021b] and L2C [Li *et al.*, 2022]. These customizing weights represent potential similarities between clients. [Chen *et al.*, 2022] proposed that graph neural networks can also be used to learn similarities among clients to realize personalization. FedAMP proposed an attentive message passing mechanism to compute personalized models, which is not flexible enough, because all clients need to participate in training in every round. FedFomo and L2C computed personalized aggregation weights via minimizing the validation loss on each client based on the model information collected from other clients, resulting in heavy communication traffic and concerns on privacy leakage. Although the effectiveness of customizing aggregation weights has been validated in existing works, their design is empirical based, not communication-efficient for large-scale real-world FL systems. Our work aligns with the line of work to customize aggregation weights in an analytical way without incurring heavy communication overhead.

### 3 Problem Formulation

In federated learning, each client  $i$  owns a local private dataset denoted by  $\mathcal{D}_i$  drawn from a distinct distribution  $\mathcal{P}_i$ . The objective of FL is to train a single global model  $w$  for all clients by solving the following problem:

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) := \frac{1}{N} \sum_{i=1}^N f_i(w) \right\}, \quad (1)$$

where the function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  represents the expected loss over the data distribution of client  $i$ , i.e.,

$$f_i(w) = \mathbb{E}_{\xi_i \sim \mathcal{P}_i} [\tilde{f}_i(w; \xi_i)]. \quad (2)$$

In the above equation,  $\xi_i$  is a random sample generated according to the local distribution  $\mathcal{P}_i$  and  $\tilde{f}_i(w; \xi_i)$  represents the loss function corresponding to sample  $\xi_i$  and  $w$ . Since clients' data possibly come from different environments, they likely have non-IID data distributions, i.e., for  $i \neq j$ ,  $\mathcal{P}_i \neq \mathcal{P}_j$ .

In conventional FL, the target is to train a global model through conducting multiple global iterations. In the  $(t-1)$ -th global iteration, the PS distributes the latest global model parameters  $w_{t-1}$  to all participating clients. The clients train locally and send the trained model  $\hat{w}_i^t$  to PS for aggregation. The whole process can be shown as

$$\begin{aligned} \hat{w}_i^t &= w_{t-1} - \eta \nabla f_i(w_{t-1}), \text{ (training)} \\ w_t &= \sum_{i=1}^N p_i \hat{w}_i^t. \text{ (aggregation)} \end{aligned} \quad (3)$$

Here, we suppose that there are  $N$  participating clients and  $p_i$  is a pre-defined non-negative weight typically proportional to  $|\mathcal{D}_i|$  with  $\sum_{i=1}^N p_i = 1$ .  $w_t$  represents the global model for the  $t$ -th round and  $\eta$  is the learning rate.

From another perspective, the single global model trained by conventional FL is to minimize the L2 distance between a global model and all local models, which can be expressed as

$$w_t = \arg \min_w \sum_{i=1}^N p_i \|w - \hat{w}_i^t\|^2. \quad (5)$$

However, if we consider the optimization of personalized models for individual clients, the optimization problem should be revised as

$$\forall i, w_i^* = \arg \min_{w_i} f_i(w_i). \quad (6)$$

Here  $w_i^*$  represents the optimal target model for client  $i$ . If the data distribution is IID, it implies that  $w_i^* \approx w_j^*$  for any two clients, which means that the optimal model applicable for each individual client can be derived by Eq. (5). However, if the data distribution is non-IID, it has been investigated in [Li *et al.*, 2019; Li *et al.*, 2020a] that a global model cannot satisfy all clients very well, resulting in poor model accuracy.

### 4 Methodology

In this section, we elaborate the design of FedDWA and prove its effectiveness through analysis.

---

#### Algorithm 1 FedDWA algorithm

---

**Input:** Communication Round  $T$ , learning rate  $\eta$ , local epochs  $E$ , number of clients  $N$ , init model parameter  $w^0$ .

**Output:** Personalized model parameters  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N$ .

**Server**

```

1: for  $t = 1, \dots, T$  do
2:   Server randomly selects a subset of clients  $S_t$  and
   sends  $w_1^t, w_2^t, \dots, w_m^t$  to them.
3:   for each client  $i \in S_t$  in parallel do
4:      $\hat{w}_i^t, \hat{w}_i^* \leftarrow \text{Client}(i, w_i^t)$ 
5:   end for
6:   Compute  $p_{i,j}$  according to Eq. (15) for each client  $i$ .
7:   for each client  $i \in S_t$  do
8:     Select top-K clients  $\{\mathcal{K}_i\}$ .
9:     Aggregate new model according to Eq. (8).
10:  end for
11: end for

```

**Client**

```

1:  $\mathcal{B} \leftarrow (\text{split } \mathcal{D} \text{ into batches of size } B)$ 
2: for each local epoch  $i$  from 1 to  $E$  do
3:   for batch  $b \in \mathcal{B}$  do
4:      $\mathbf{w} = \mathbf{w} - \eta \nabla \tilde{f}(\mathbf{w}, b)$ 
5:   end for
6: end for
7: Train one more local iteration (one epoch).
8:  $\hat{\mathbf{w}} = \mathbf{w} - \eta \nabla \tilde{f}(\mathbf{w}, \mathcal{D})$ 
9: return  $\mathbf{w}$  and  $\hat{\mathbf{w}}$ 

```

---

#### 4.1 Optimization Objective

In the previous section, it has been pointed out that the aggregation rule defined by Eq. (4) cannot meet personalized requirement with non-IID data distribution. Rather than training a single global model, we propose to train a model for each individual client by customizing aggregation weights so as to deduce each individual model.

Specifically, we define  $p_{i,j}$  as the weight to aggregate the model for client  $i$  using the local model from client  $j$ . Here,  $\sum_{j=1}^N p_{i,j} = 1$ . Then, the PS can generate the personalized model for client  $i$  in the  $(t-1)$ -th global iteration as follows:

$$\hat{w}_i^t = w_i^{t-1} - \eta_i^{t-1} \nabla f_i(w_i^{t-1}). \quad (7)$$

$$w_i^t = \sum_{j=1}^N p_{i,j}^t \hat{w}_j^t. \quad (8)$$

Intuitively speaking, Eq. (8) is very flexible. When aggregating the model for client  $i$ , we can set a larger value for  $p_{i,j}^t$  if the data distribution of client  $j$  is closer to that of client  $i$  such that the PS can explore optimal personalized model for client  $i$ <sup>1</sup>. To specify how to set the value of  $p_{i,j}$ , we formulate the optimization problem as below:

$$\min_{p_{i,1}, \dots, p_{i,N}} \left\| \hat{w}_i^* - \sum_{j=1}^N p_{i,j} \hat{w}_j^t \right\|^2, \quad \forall i. \quad (9)$$

---

<sup>1</sup> $p_{i,j}$  and  $\eta_i$  can be time-dependent, but when context allows, we write  $p_{i,j}^t$  as  $p_{i,j}$  and  $\eta_i^t$  as  $\eta_i$  for simplicity.

Here  $\hat{w}_i^*$  is the guidance model for client  $i$  in the  $(t-1)$ -th global iteration, and it tells client  $i$  which clients to cooperate with. The main challenge for solving Eq. (9) lies in that  $\hat{w}_i^*$  is unknown in advance. We can solve Eq. (9) with two steps. Firstly, we need to find a high quality guidance model  $\hat{w}_i^*$  and fix  $\hat{w}_i^*$  to derive how to optimally set  $p_{i,j}$ . Secondly, after deriving the optimal weight  $p_{i,j}$ , client  $i$  can get its own personalized model at the  $t$ -th round, and then its guidance model  $\hat{w}_i^*$  can be further updated. In the following, we elaborate how to solve Eq. (9).

### Tuning Aggregation Weights

Since  $\sum_{j=1}^N p_{i,j} = 1$ , we can rewrite Eq. (9) as:

$$\left\| \hat{w}_i^* - \sum_{j=1}^N p_{i,j} \hat{w}_j^t \right\|^2 = \sum_{j=1}^N \sum_{k=1}^N p_{i,j} p_{i,k} (\hat{w}_i^* - \hat{w}_j^t)^T (\hat{w}_i^* - \hat{w}_k^t). \quad (10)$$

Let the vector  $\mathbf{p}_i = [p_{i,1}, p_{i,2}, \dots, p_{i,N}]^T$  denote aggregation weights for obtaining client  $i$ 's personalized model. Let  $\mathbf{W}_i$  denote the cross distance between the guidance model  $\hat{w}_i^*$  and local models contributed by clients. The  $(j, k)$ -th entry of  $\mathbf{W}_i$  can be written as:

$$[\mathbf{W}_i]_{j,k} = (\hat{w}_i^* - \hat{w}_j^t)^T (\hat{w}_i^* - \hat{w}_k^t). \quad (11)$$

Then the optimization problem for client  $i$  defined in Eq. (9) can be expressed as follows:

$$\begin{aligned} \min_{\mathbf{p}_i} \quad & \mathbf{p}_i^T \mathbf{W}_i \mathbf{p}_i, \\ \text{subject to} \quad & \mathbf{1}_N^T \mathbf{p}_i = 1, p_{i,k} \geq 0. \end{aligned} \quad (12)$$

Note that the PS can optimize personalized models for all clients via Eq. (12). Suppose that  $\mathbf{W}_i$  is invertible, then the solution is given by:

$$\mathbf{p}_i = \frac{\mathbf{W}_i^{-1} \mathbf{1}_N}{\mathbf{1}_N^T \mathbf{W}_i^{-1} \mathbf{1}_N} \quad (13)$$

It is very difficult to directly calculate Eq. (13) due to the following three challenges. *First*, it involves the inner product of model parameters among clients which can be seen from Eq. (11). For advanced neural networks, there may exist millions of parameters making the computation cost unaffordable. *Second*, training models in federated learning is an iterative process with multiple rounds of communications. It implies that computing the inversion of  $\mathbf{W}_i$  is cumbersome, especially when the dimension of  $\mathbf{W}_i$  is very large. *Third*, since  $\mathbf{W}_i$  is just a symmetric matrix,  $\mathbf{W}_i^{-1}$  may not exist at all. As a consequence, the solution of Eq. (12) is not unique. Thus, trying to solve Eq. (12) directly cannot guarantee that a high-quality solution will be yielded. A toy example is discussed in Appendix A.1.

To make the problem tractable, we simplify the objective in Eq. (12) by only reserving the diagonal elements of  $\mathbf{W}_i$ . The effectiveness of such simplification has been verified in previous works [Chen *et al.*, 2015; Zhao and Sayed, 2012]. The simplified problem is presented as follows:

$$\begin{aligned} \min_{\mathbf{p}_i} \quad & \sum_{j=1}^N p_{i,j}^2 \left\| \hat{w}_i^* - \hat{w}_j^t \right\|^2, \\ \text{subject to} \quad & \mathbf{1}_N^T \mathbf{p}_i = 1, p_{i,j} \geq 0. \end{aligned} \quad (14)$$

It is easy to find that there is a unique solution to the simplified problem. It is worth noting that Eq. (14) is very alike to the aggregation rule in conventional FL defined in Eq. (5). The difference of our method can be explained from two perspectives. First, aggregation weights are tunable parameters in our method, which however are fixed in traditional FL. Second, our method can search the optimal aggregation weights for individual clients to derive personalized models. The solution of Eq. (14) is:

$$p_{i,j} = \frac{\left\| \hat{w}_i^* - \hat{w}_j^t \right\|^{-2}}{\sum_{k=1}^N \left\| \hat{w}_i^* - \hat{w}_k^t \right\|^{-2}}, \quad (15)$$

where the detailed derivation can be found in Appendix A.2.

### One-step Ahead Adaptation

Executing the combination rule in Eq. (15) by an individual client requires the knowledge of the guidance model  $\hat{w}_i^*$ , which is generally not available beforehand or not. Intuitively speaking, whether we can realize personalization depends on similar clients identified by our algorithm. In other words, the weighted combination of models for client  $i$  should align with client  $i$ 's personal data distribution. It implies that the guidance model  $\hat{w}_i^*$  should capture the local data distribution of client  $i$ . We have tried several options for  $\hat{w}_i^*$  including using the last iteration model  $\hat{w}_i^* = \hat{w}_i^{t-1}$ , current model  $\hat{w}_i^* = \hat{w}_i^t$ , and one-step ahead adaptation. More discussion can be found in Appendix C.6. In this work, we employ an instantaneous adaptation argument, a.k.a. local one-step ahead adaptation, to accommodate this issue as follows:

$$\hat{w}_i^* = \hat{w}_i^t - \eta_i^{t-1} \nabla f_i(\hat{w}_i^t). \quad (16)$$

The validity of this adaptation can be found in previous work [Jin *et al.*, 2020; Chen *et al.*, 2015]. It is important to note that this is fundamentally different from traditional FedAvg training and then local fine-tuning, because our method will only select other clients that are beneficial to client  $i$  for aggregation while the fine-tuning approach treats all clients equally during the aggregation phase. In fact, after a step of adaptation, Eq. (16) in advance,  $\hat{w}_i^*$  can characterize its local data distribution well and therefore it can screen out other clients with similar data distribution and give them a higher weight for cooperation. We can also use two steps or even more steps, the specific experimental results will be shown in the following sections. By substituting Eq. (16) into Eq. (15), we can finally derive the aggregation weights by deriving personalized models as

$$p_{i,j} = \frac{\left\| \hat{w}_i^t - \eta_i^{t-1} \nabla f_i(\hat{w}_i^t) - \hat{w}_j^t \right\|^{-2}}{\sum_{k=1}^N \left\| \hat{w}_i^t - \eta_i^{t-1} \nabla f_i(\hat{w}_i^t) - \hat{w}_k^t \right\|^{-2}}. \quad (17)$$

Similar to the previous works [Zhang *et al.*, 2021b; Li *et al.*, 2022], the top-K technique can be used. We rank  $p_{i,j}$  in descending order, and only the top  $K$  aggregation weights are selected and they will be normalized such that  $\sum_{j=1}^N p_{i,j} = 1$ . By wrapping up our analysis, we present the detailed FedDWA algorithm in Algorithm. 1.

Settings	Pathological heterogeneous setting					Practical heterogeneous setting 1			
Methods	EMNIST	CIFAR10	CIFAR100	CINIC10	TINY	CIFAR10	CIFAR100	CINIC10	TINY
Local Training	97.23	92.35	80.08	92.74	58.04	72.12	39.82	64.40	19.90
FedAvg	71.78	59.97	34.71	46.07	9.64	71.57	44.67	58.97	13.74
FedProx	69.95	57.42	31.12	45.60	8.46	70.03	41.53	56.01	6.40
FedAvgM	66.42	59.82	34.00	48.81	8.84	71.49	44.56	58.77	12.22
FedAvg_FT	93.65	88.82	62.73	81.87	18.30	74.35	46.76	62.45	14.70
SFL	72.02	60.71	34.17	48.83	10.88	71.77	44.83	59.28	13.46
per-FedAvg	93.45	88.44	62.83	89.65	16.28	73.32	43.62	62.57	8.70
pFedMe	95.42	90.78	78.10	88.72	51.12	77.68	45.79	59.82	17.26
ClusterFL	78.45	83.41	51.30	80.48	40.26	71.42	44.90	59.44	12.92
FedRoD	93.77	87.95	64.76	83.75	50.96	77.27	46.76	52.64	18.04
FedAMP	96.65	91.74	78.61	88.05	50.72	72.30	41.11	65.43	27.48
FedFomo	96.95	91.95	78.89	92.59	59.24	75.69	47.06	68.93	19.16
L2C	95.75	91.76	78.62	92.69	54.24	76.67	48.30	67.52	21.04
Ours	<b>97.37</b>	<b>92.97</b>	<b>80.41</b>	<b>92.75</b>	<b>60.64</b>	<b>78.09</b>	<b>50.83</b>	<b>70.29</b>	<b>28.92</b>

Table 1: Average test accuracy (%) over five different datasets, under pathological heterogeneous setting and practical heterogeneous setting 1 with 20 clients, 100% participation, respectively.

## 4.2 Analysis of FedDWA

### Communication Overhead

The FedDWA algorithm will incur  $2\Sigma$  traffic in the uplink communication and the traffic in the downlink communication is the same as that of the original FedAvg, where  $\Sigma$  denotes the model size. It is worth noting that FedDWA can incur significantly less communication traffic than other similar baselines, and more results can be found in Appendix C.4.

### Computational Cost

Suppose that there are  $N$  clients participating in training and the number of model parameters is  $d$ , the computation complexity of FedDWA is  $\mathcal{O}(N^2d)$  in the server. We also test the total FLOPs required by FedDWA and the experimental results show that the computational amount required to calculate the similarity (Eq.(15)) is negligible compared with model training. More results can be found in Appendix C.5.

### Personalized Learning

In this part, we illustrate how FedDWA can make clients with similar data distribution collaborate to train personalized models. Considering the inverse of the numerator of Eq. (17), it can be expanded as

$$\begin{aligned} & \|\hat{w}_i^t - \eta_i^{t-1} \nabla f_i(\hat{w}_i^t) - \hat{w}_j^t\|^2 = \\ & \|\hat{w}_i^t - \hat{w}_j^t\|^2 + 2\eta_i^{t-1}(\hat{w}_j^t - \hat{w}_i^t)^T \nabla f_i(\hat{w}_i^t) + \\ & (\eta_i^{t-1})^2 \|\nabla f_i(\hat{w}_i^t)\|^2. \end{aligned} \quad (18)$$

The first term  $\|\hat{w}_i^t - \hat{w}_j^t\|^2$  refers to the distance between current models of client  $i$  and client  $j$ . This term will lower the aggregation weight  $p_{i,k}$  if the distance  $\|\hat{w}_i^t - \hat{w}_j^t\|^2$  is large, and thereby prohibit their collaborations. Using the first-order Taylor series to expand  $f_i(w)$  at  $\hat{w}_i^t$ , we have:

$$f_i(w) \approx f_i(\hat{w}_i^t) + (w - \hat{w}_i^t)^T \nabla f_i(w) |_{\hat{w}_i^t}. \quad (19)$$

For the second term  $2\eta_i^{t-1}(\hat{w}_j^t - \hat{w}_i^t)^T \nabla f_i(\hat{w}_i^t)$ , we can find that it is proportional to  $f_i(\hat{w}_j^t) - f_i(\hat{w}_i^t)$  which also decreases

the aggregation weight  $p_{i,k}$  if  $f_i(\hat{w}_j^t)$  is far away from  $f_i(\hat{w}_i^t)$ . The last term  $(\eta_i^{t-1})^2 \|\nabla f_i(\hat{w}_i^t)\|^2$  can be perceived as a constant when optimizing aggregation weights. In summary, Eq. (17) provides the aggregation weights for deriving a personalized model for client  $i$  based on the similarity distance between client models  $\hat{w}_j$ 's and the guidance model  $\hat{w}_i^*$ .

## 5 Experiments

### 5.1 Experiment Setups

#### Datasets and Models

We evaluate our algorithm on five benchmark datasets, namely, EMNIST [Cohen *et al.*, 2017], CIFAR10, CIFAR100 [Krizhevsky *et al.*, 2009], CINIC10 [Darlow *et al.*, 2018] and Tiny-ImageNet (TINY) [Chrabaszcz *et al.*, 2017]. For EMNIST, we use the same model as that used in [Sattler *et al.*, 2021]. For CIFAR10, CIFAR100 and CINIC10, we use the CNN model which is the same as that in [Mills *et al.*, 2022]. To evaluate the effectiveness of FedDWA on a high-dimensional model, we use ResNet-8 for the Tiny-ImageNet, and the model architecture is the same as that in [He *et al.*, 2020]. More details can be found in Appendix B.1.

#### Data Partitioning

We simulate the heterogeneous settings with three widely used scenarios, including a pathological setting and two practical settings.

- **Pathological Heterogeneous Setting.** Each client is randomly assigned with a small number of classes with the same amount of data on each class [McMahan *et al.*, 2017; Shamsian *et al.*, 2021]. We sample 4, 2, 2, 6 and 10 classes for EMNIST, CIFAR10, CINIC10, CIFAR100, Tiny-ImageNet from a total of 62, 10, 10, 100, 200 classes for each client, respectively. There is no group-wise similarity between clients in this setting.
- **Practical Heterogeneous Setting 1.** All clients have the same data size but different distributions. For each client,  $s\%$  of data (80% by default) are selected from a

set of dominant classes, and the remaining  $(100 - s)\%$  are uniformly sampled from all classes [Karimireddy *et al.*, 2020; Huang *et al.*, 2021]. All clients are divided into multiple groups. Clients in each group share the same dominant classes implying that there is an underlying clustering structure between clients.

- **Practical Heterogeneous Setting 2.** Each client contains most of the classes but the data in each class is not uniformly distributed [Hsu *et al.*, 2019; Li *et al.*, 2021a; Chen and Chao, 2022]. We create the federated version by randomly partitioning datasets among  $N$  clients using a symmetric Dirichlet distribution  $\text{Dir}(\alpha)$  ( $\alpha = 0.07$  by default). For example, for each class  $c$ , we sample a vector  $p_c$  from  $\text{Dir}(\alpha)$  and allocate to client  $m$  a fraction  $p_{c,m}$  of all training instances of class  $c$ .

## Baselines

We compare the performance of our algorithm with that of FedAvg [McMahan *et al.*, 2017], FedAvgM [Hsu *et al.*, 2019], FedProx [Li *et al.*, 2020a] and a few latest personalization approaches including a personalized model trained only on each client’s local dataset (Local Training), FedAvg with local tuning (FedAvg\_FT) [Wang *et al.*, 2019], pFedMe [T Dinh *et al.*, 2020], per-FedAvg [Fallah *et al.*, 2020], ClusterFL [Sattler *et al.*, 2021], FedAMP [Huang *et al.*, 2021], FedFomo [Zhang *et al.*, 2021b], SFL [Chen *et al.*, 2022], L2C [Li *et al.*, 2022] and FedRoD [Chen and Chao, 2022]. The settings of hyper-parameters for each method can be found in Appendix B.3.

## Evaluation Metrics

We use the same evaluation metrics as that widely used by previous works which report the test accuracy of the best single global model for the single-model PFL methods and the average test accuracy of the best personalized models for other PFL methods.

## Training Settings

Similar to [Zhang *et al.*, 2021b], we evaluate the performance in two settings, i.e., (1)  $N = 20$  clients, 100% participation and (2)  $N = 100$  clients, 20% participation for all datasets. The number of local training epochs is set to  $E = 1$  and the number of global communication rounds is set to 100. We employ the mini-batch SGD as a local optimizer for all approaches. The batch size for each client is set as 20 and the learning rate  $\eta$  is set as 0.01. We test all methods over three runs and average the results.

## 5.2 Performance Evaluation and Analysis

### Pathological Heterogeneous Setting

Table 1 shows the average test accuracy for all methods under the pathological heterogeneous setting. Over all datasets and client setups, our FedDWA algorithm outperforms other baseline methods. However, the performance of methods (such as FedAvg and FedProx) that only train a single global model degrade significantly on CIFAR10, CIFAR100, CINIC10 and Tiny-ImageNet, since a single global model cannot well accommodate statistical heterogeneity of clients. Other personalized methods achieve comparable accuracy and the local

	CIFAR10	CIFAR100	CINIC10	TINY
FedAvg	48.44	22.80	26.68	6.72
FedProx	45.08	23.69	27.06	5.90
FedAvgM	48.09	21.70	27.36	4.70
FedAvg_FT	89.19	45.06	83.32	6.47
SFL	52.40	22.96	36.82	5.13
per-FedAvg	89.86	49.81	89.69	18.52
pFedMe	90.54	51.70	89.41	6.15
ClusterFL	85.48	31.33	76.73	21.85
FedRoD	86.83	42.90	88.72	26.77
FedAMP	91.27	51.16	90.51	27.14
FedFomo	91.73	54.29	91.32	32.33
L2C	91.77	54.98	91.67	30.41
Ours	<b>91.81</b>	<b>55.26</b>	<b>91.80</b>	<b>32.66</b>

Table 2: Average test accuracy (%) over four different datasets, under the practical heterogeneous setting 2 with 100 clients, 20% participation.

training method achieves rather high performance due to the small number of classes on each client.

### Practical Heterogeneous Setting

Table 1 shows the average test accuracy for all methods under the practical heterogeneous setting 1 in which each client has a primary data class with a small number of samples from all other classes. In this setting, we find that our method is significantly superior to all other baseline methods. For example, when using Tiny-ImageNet, our method outperforms FedFomo and L2C by up to 9.76%, 7.88% in test accuracy respectively, which means that the aggregation weights obtained by our method are better than those with the aggregation weights obtained by empirical searching through the validation set. To test the applicability of our algorithm, we also evaluate the performance of our approach for a large-scale FL scenario under the practical heterogeneous setting 2. We set  $N = 100$  clients with 20% participation in each round. The final results are shown in Table 2. It is worth noting that our approach still achieves competitive performance though there is no clear similarity between clients in this scenario.

### Personalized Model Weighting

We use the practical heterogeneous setting 1 in which clients are divided into multiple groups to explain why FedDWA outperforms existing works by showing how FedDWA helps clients quickly identify other similar clients for conducting personalized model aggregation. To ease our visualization, we group clients first before clients are indexed such that clients in the same group will have consecutive indices. For example, clients of the same data distribution are grouped in the first group who are indexed by 0-4. In Figure 1, we show the  $K$  most similar clients selected by different methods according to personalized weights in each training round with 20 clients divided into 4 groups. Since there are five clients in each group, we set  $K = 5$ . The result manifests that FedDWA can identify similar clients faster than FedFomo and L2C. In addition, under this setting, L2C cannot accurately identify the similarity between users, indicating that weights found with validation datasets are not very effective. More experiment results such as the heat maps showing the aggre-

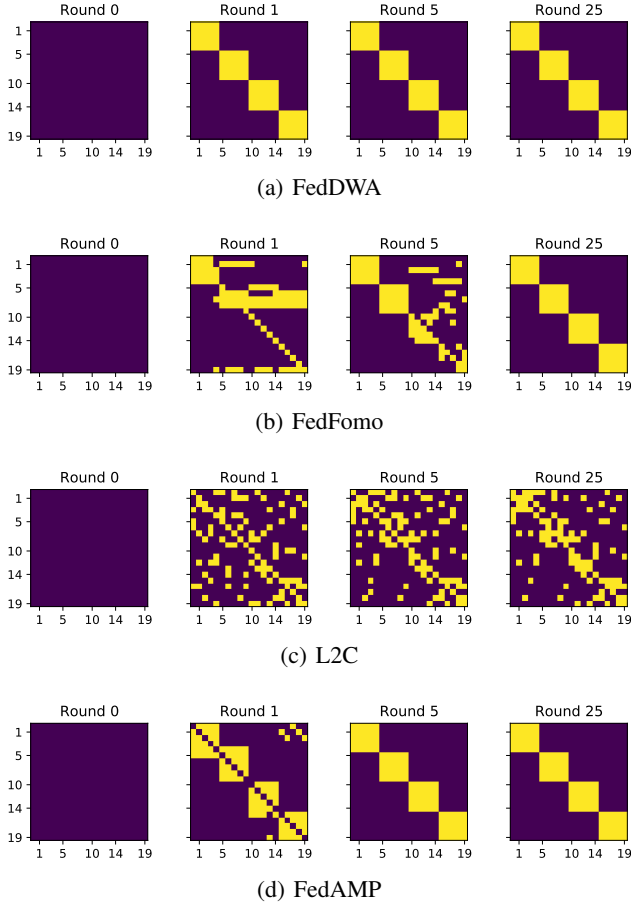


Figure 1: The visualization of the  $K$  most similar clients selected by different methods on CIFAR10 datasets. The x-axis and y-axis means the IDs of clients.

gation weights can be found in Appendix C.1.

### Selection of $K$

Figure 4 in Appendix C shows that when there is a potential cluster structure among clients participating in FL, the weight calculated by FedDWA at the first round can reflect the similarity between clients well. At this time, server can easily determine the value of  $K$  according to the weight matrix. However, as reported in [Ghosh *et al.*, 2020; Zhang *et al.*, 2021b; Li *et al.*, 2022; Marfoq *et al.*, 2022], it is not trivial to set  $K$  if there is no obvious data similarity between clients. Fortunately, the performance of FedDWA is not very sensitive with respect to  $K$ . We use the practical heterogeneous setting 2 to simulate different degrees of data heterogeneity with 100 clients, and present the test performance in Figure 2. The results indicate that FedDWA is not very sensitive to the value of  $K$  when data distributions are more heterogeneous (with a smaller  $\alpha$ ). When  $\alpha$  is large, the discrepancy between users' data distributions is smaller, which is not a typical PFL scenario. At this time it becomes more important to increase the value of  $K$  so as to cooperate with more clients.

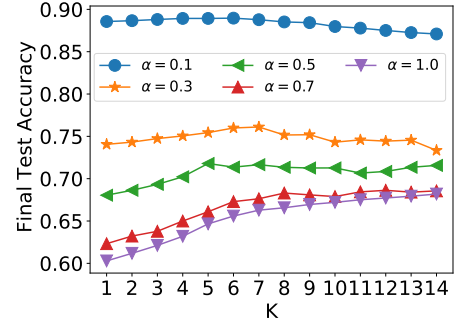


Figure 2: The influence of different values of  $K$  on the final performance of our algorithm using CIFAR10 dataset.

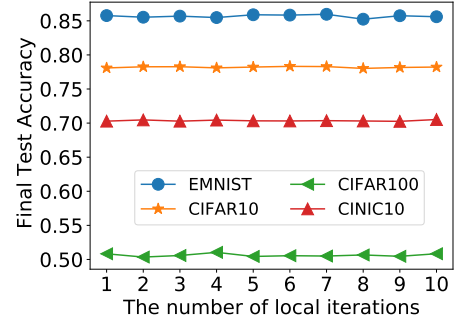


Figure 3: The influence of different local epoch iterations on the final test accuracy of our algorithm.

### Effect of Guidance Model

We explore the influence of guidance models obtained by Eq. (16) on the final personalized model accuracy by tuning the number of local epochs. According to Eq. (16),  $\hat{w}_i^*$  will be different if we set a larger number of local epochs. We evaluate this effect with 20 clients on four datasets, respectively. As shown in Figure 3, the number of local iterations to compute guidance models has little influence on the final model accuracy of FedDWA indicating that one-step adaptation is sufficient for FedDWA. Moreover, we compare the performance of guidance model  $\hat{w}_i^*$  and personalized model  $w_i^t$  in Appendix C.3. The results further show that the performance of  $w_i^t$  is almost the same as that of  $\hat{w}_i^*$  under the practical setting 2 and better than that of  $\hat{w}_i^*$  under the practical setting 1.

## 6 Conclusions

In this paper, we propose a novel FedDWA algorithm which can identify similarities between clients with much less communication overhead than other relevant works since no information will be exchanged between clients. Meanwhile, personalized models are generated based on uploaded model information and its effectiveness is guaranteed with theoretical analysis. Comprehensive experiments on five datasets demonstrate the superb performance of FedDWA which can achieve the highest model accuracy compared to the state-of-the-art baselines under three heterogeneous FL settings.

## Acknowledgements

This work was supported by the National Natural Science Foundation of China under Grants U1911201, U2001209, 62072486, and the Natural Science Foundation of Guangdong Province under Grant 2021A1515011369.

## References

- [Acar *et al.*, 2021] Durmus Alp Emre Acar, Yue Zhao, Ruizhao Zhu, Ramon Matas, Matthew Mattina, Paul Whatmough, and Venkatesh Saligrama. Debiasing model updates for improving personalized federated training. In *International Conference on Machine Learning*, pages 21–31. PMLR, 2021.
- [Achituve *et al.*, 2021] Idan Achituve, Aviv Shamsian, Aviv Navon, Gal Chechik, and Ethan Fetaya. Personalized federated learning with gaussian processes. *Advances in Neural Information Processing Systems*, 34:8392–8406, 2021.
- [Arivazhagan *et al.*, 2019] Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. *arXiv preprint arXiv:1912.00818*, 2019.
- [Chen and Chao, 2022] Hong-You Chen and Wei-Lun Chao. On bridging generic and personalized federated learning for image classification. In *International Conference on Learning Representations*, 2022.
- [Chen *et al.*, 2015] Jie Chen, Cédric Richard, and Ali H Sayed. Diffusion lms over multitask networks. *IEEE Transactions on Signal Processing*, 63(11):2733–2748, 2015.
- [Chen *et al.*, 2022] Fengwen Chen, Guodong Long, Zonghan Wu, Tianyi Zhou, and Jing Jiang. Personalized federated learning with a graph. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 2575–2582. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track.
- [Chrabaszcz *et al.*, 2017] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
- [Cohen *et al.*, 2017] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pages 2921–2926. IEEE, 2017.
- [Collins *et al.*, 2021] 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.
- [Darlow *et al.*, 2018] Luke N Darlow, Elliot J Crowley, Antreas Antoniou, and Amos J Storkey. Cinic-10 is not imagenet or cifar-10. *arXiv preprint arXiv:1810.03505*, 2018.
- [Deng *et al.*, 2020] Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- [Fallah *et al.*, 2020] 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.
- [Ghosh *et al.*, 2020] Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for clustered federated learning. *Advances in Neural Information Processing Systems*, 33:19586–19597, 2020.
- [Hanzely and Richtárik, 2020] Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models. *arXiv preprint arXiv:2002.05516*, 2020.
- [He *et al.*, 2020] Chaoyang He, Murali Annaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of large cnns at the edge. *Advances in Neural Information Processing Systems*, 33:14068–14080, 2020.
- [Hsu *et al.*, 2019] Tzu-Ming Harry Hsu, Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *ArXiv*, abs/1909.06335, 2019.
- [Hu *et al.*, 2021] Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, and Xuyun Zhang. Source inference attacks in federated learning. In *2021 IEEE International Conference on Data Mining (ICDM)*, pages 1102–1107. IEEE, 2021.
- [Huang *et al.*, 2021] 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.
- [Jin *et al.*, 2020] Danqi Jin, Jie Chen, Cédric Richard, Jingdong Chen, and Ali H Sayed. Affine combination of diffusion strategies over networks. *IEEE Transactions on Signal Processing*, 68:2087–2104, 2020.
- [Karimireddy *et al.*, 2020] 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.
- [Khodak *et al.*, 2019] Mikhail Khodak, Maria-Florina Balcan, and Ameet S. Talwalkar. Adaptive gradient-based meta-learning methods. In *NeurIPS*, 2019.
- [Krizhevsky *et al.*, 2009] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Tech Report*, 2009.
- [Li *et al.*, 2019] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- [Li *et al.*, 2020a] 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.
- [Li *et al.*, 2020b] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *ArXiv*, abs/1907.02189, 2020.
- [Li *et al.*, 2021a] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10713–10722, 2021.
- [Li *et al.*, 2021b] 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.
- [Li *et al.*, 2022] Shuangtong Li, Tianyi Zhou, Xinmei Tian, and Dacheng Tao. Learning to collaborate in decentralized learning of personalized models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9766–9775, 2022.
- [Mansour *et al.*, 2020] Y. Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for personalization with applications to federated learning. *ArXiv*, abs/2002.10619, 2020.
- [Marfoq *et al.*, 2021] Othmane Marfoq, Giovanni Neglia, Aurélien Bellet, Laetitia Kameni, and Richard Vidal. Federated multi-task learning under a mixture of distributions. *Advances in Neural Information Processing Systems*, 34:15434–15447, 2021.
- [Marfoq *et al.*, 2022] Othmane Marfoq, Giovanni Neglia, Richard Vidal, and Laetitia Kameni. Personalized federated learning through local memorization. In *International Conference on Machine Learning*, pages 15070–15092. PMLR, 2022.
- [McMahan *et al.*, 2017] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [Mills *et al.*, 2022] Jed Mills, Jia Hu, and Geyong Min. Multi-task federated learning for personalised deep neural networks in edge computing. *IEEE Transactions on Parallel and Distributed Systems*, 33:630–641, 2022.
- [Oh *et al.*, 2021] Jaehoon Oh, Sangmook Kim, and Se-Young Yun. Fedbabu: Towards enhanced representation for federated image classification. *arXiv preprint arXiv:2106.06042*, 2021.
- [Sahu *et al.*, 2020] Anit Kumar Sahu, Tian Li, Maziar Sanjabi, Manzil Zaheer, Ameet S. Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *arXiv: Learning*, 2020.
- [Sattler *et al.*, 2021] Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints. *IEEE Transactions on Neural Networks and Learning Systems*, 32:3710–3722, 2021.
- [Shamsian *et al.*, 2021] 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.
- [Shin *et al.*, 2020] MyungJae Shin, Chihoon Hwang, Joongheon Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Xor mixup: Privacy-preserving data augmentation for one-shot federated learning. *arXiv preprint arXiv:2006.05148*, 2020.
- [Smith *et al.*, 2017] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S. Talwalkar. Federated multi-task learning. In *NIPS*, 2017.
- [T Dinh *et al.*, 2020] Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes. *Advances in Neural Information Processing Systems*, 33:21394–21405, 2020.
- [Tan *et al.*, 2022] Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fedproto: Federated prototype learning across heterogeneous clients. In *AAAI Conference on Artificial Intelligence*, volume 1, page 3, 2022.
- [Wang *et al.*, 2019] Kangkang Wang, Rajiv Mathews, Chloé Kiddon, Hubert Eichner, Françoise Beaufays, and Daniel Ramage. Federated evaluation of on-device personalization. *arXiv preprint arXiv:1910.10252*, 2019.
- [Yoon *et al.*, 2021] Tehrim Yoon, Sumin Shin, Sung Ju Hwang, and Eunho Yang. Fedmix: Approximation of mixup under mean augmented federated learning. *arXiv preprint arXiv:2107.00233*, 2021.
- [Yu *et al.*, 2020] Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging federated learning by local adaptation. *ArXiv*, abs/2002.04758, 2020.
- [Yue *et al.*, 2021] Sheng Yue, Ju Ren, Jiang Xin, Sen Lin, and Junshan Zhang. Inexact-admm based federated meta-learning for fast and continual edge learning. *Proceedings of the Twenty-second International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, 2021.
- [Zhang *et al.*, 2021a] 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.
- [Zhang *et al.*, 2021b] Michael Zhang, Karan Sapra, Sanja Fidler, Serena Yeung, and José Manuel Álvarez. Personalized federated learning with first order model optimization. *ArXiv*, abs/2012.08565, 2021.
- [Zhao and Sayed, 2012] Xiaochuan Zhao and Ali H Sayed. Clustering via diffusion adaptation over networks. In *2012 3rd International Workshop on Cognitive Information Processing (CIP)*, pages 1–6. IEEE, 2012.
- [Zhao *et al.*, 2018] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Cavin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.

## A Theoretical Results

### A.1 A Tony Example for Problem (12)

If the matrix  $\mathbf{W}_i$  is not invertible and has a rank less than  $N - 1$ , the solution of Eq. (12) will be not unique. For example, consider the case when  $N = 3$ , and

$$\mathbf{W}_i = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

It means that the model of client  $i$  is similar to that of clients 1 and 2, and is very different from that of client 3. In this case, any vector  $\mathbf{p}_i = (p_{i,1}, p_{i,2}, 0)^T$  with  $p_{i,1} + p_{i,2} = 1$  is the solution of Eq. (12). For example,  $\{p_{i,1} = 0.9, p_{i,2} = 0.1\}$  and  $\{p_{i,1} = 0.1, p_{i,2} = 0.9\}$  are two sets of solutions that satisfy the requirements respectively. But, it is obvious that we will infer different things from these two different sets of solutions. In the first set of solutions, we will think that client  $i$  is similar to client 1, not to similar to client 2; In the second set of solutions, we will think that client  $i$  is similar to client 2, not so similar to client 1. Therefore, directly solving Eq. (12) in this case does not give us a unique solution.

### A.2 Optimal Solution of Problem (14)

To solve the following optimization problem:

$$\begin{aligned} \min_{\mathbf{p}_i} \quad & \sum_{j=1}^N p_{i,j}^2 \|\hat{w}_i^* - \hat{w}_j^t\|^2, \\ \text{subject to} \quad & \mathbf{1}_N^T \mathbf{p}_i = 1, p_{i,j} \geq 0. \end{aligned} \quad (20)$$

First we construct the Lagrangian function with respect to the equality constraint and discard the non-negativity constraint.

$$L(\lambda) = \sum_{j=1}^N p_{i,j}^2 \|\hat{w}_i^* - \hat{w}_j^t\|^2 + \lambda(\mathbf{1}_N^T \mathbf{p}_i - 1) \quad (21)$$

Then, take the derivative of  $\lambda$  and  $p_{i,j}$  respectively, and we get:

$$2p_{i,j} \|\hat{w}_i^* - \hat{w}_j^t\|^2 - \lambda = 0 \quad (22)$$

$$\mathbf{1}_N^T \mathbf{p}_i - 1 = 0 \quad (23)$$

It turns out that the value of  $\lambda$  is:

$$\lambda = \frac{2}{\sum_{j=1}^N \|\hat{w}_i^* - \hat{w}_j^t\|^{-2}} \quad (24)$$

Finally, we can determine that:

$$p_{i,j} = \frac{\|\hat{w}_i^* - \hat{w}_j^t\|^{-2}}{\sum_{k=1}^N \|\hat{w}_i^* - \hat{w}_k^t\|^{-2}}. \quad (25)$$

Given that it Eq. (25) satisfies the non-negativity constraint  $p_{i,j} \geq 0$ , therefore, Eq. (25) is the solution of Eq. (20).

## B Details of Experiment Setup

We implement our method and other baseline methods in PyTorch 1.7, and the simulation server is equipped with a Tesla V100 GPU, a 2.4-GHz Inter Core E5-2680 CPU and 256GB of memory.

### B.1 Datasets and Models

We consider image classification tasks and evaluate our method on five popular datasets: (1) EMNIST (Extend MNIST) is a 62-class image classification dataset, extending the classic MNIST dataset. It contains 62 categories of handwritten characters, including 10 digits, 26 uppercase letters and 26 lowercase letters. There are 814,255 images in total; (2) CIFAR10 dataset consists of 60,000 32x32 colour images in 10 classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images; (3) CIFAR100 dataset consists of 60,000 32x32 colour images in 100 classes, with 600 images per class. There are 500 training images and 100 test images; (4) CINIC-10, which is more diverse than CIFAR10 as it is constructed from two different sources: ImageNet and CIFAR10. This dataset consists of 100,000 32x32 colour images in 10 classes; (5) Tiny-ImageNet is constructed from ImageNet and it consists of 100,000 32x32 colour images in 200 classes. We construct three different CNN models for classifying EMNIST, CIFAR10/CIFAR100/CINIC-10 and Tiny-ImageNet images, respectively. The first CNN model is constructed by two convolutional layers followed by pooling, and two fully connected layers with a final dense layer containing 2,048 units. The second CNN model has two convolutional pooling layers, two batch normalization layers and two fully connected layers, and the rate of Dropout is set to 0.5. The third CNN model is ResNet-8, and the architecture is the same as [He *et al.*, 2020].

### B.2 Data Partitioning

We simulate the heterogeneous settings with three widely used scenarios, including a pathological setting and two practical settings.

- **Pathological Heterogeneous Setting.** Each client is randomly assigned with a small number of classes of samples [McMahan *et al.*, 2017; Shamsian *et al.*, 2021]. We sample 4, 2, 2, 6 and 10 classes for EMNIST, CIFAR10, CINIC10, CIFAR100, Tiny-ImageNet from a total of 62, 10, 10, 100, 200 classes for each client, respectively. There is no group-wise similarity between clients in this setting.
- **Practical Heterogeneous Setting 1.** All clients have the same data size but different distributions. For each client,  $s\%$  of data (80% by default) are selected from a set of dominant classes, and the remaining  $(100 - s)\%$  are uniformly sampled from all classes [Karimireddy *et al.*, 2020; Huang *et al.*, 2021]. All clients are divided into multiple groups. Clients in each group share the same dominant classes implying that there is an underlying clustering structure between clients. Specifically, for CIFAR10 and CINIC-10 datasets which have 10 categories of images, we divide the clients into 4 groups and the number of dominant class for each client in the same group is 3. For CIFAR100 dataset which has 100 categories of images, we divide the clients into 4 groups and the number of dominant class for each client in the same group is 20; For Tiny-ImageNet dataset which have 200 categories, we divide the clients into 4 groups and the

number of dominant class for each client in the same group is 40.

- **Practical Heterogeneous Setting 2.** Each client contains most of the classes but the data in each class is not uniformly distributed [Hsu *et al.*, 2019; Li *et al.*, 2021a; Chen and Chao, 2022]. We create the federated version by randomly partitioning datasets among  $N$  clients using a symmetric Dirichlet distribution  $\text{Dir}(\alpha)$  ( $\alpha = 0.07$  by default). For example, for each class  $c$ , we sample a vector  $p_c$  from  $\text{Dir}(\alpha)$  and allocate to client  $m$  a fraction  $p_{c,m}$  of all training instances of class  $c$ .

For each setting, the test data on each client has the same distribution as that of the training data.

### B.3 Implementation Details of Methods

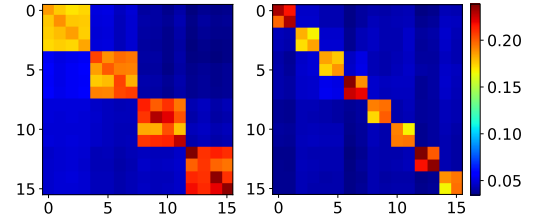
By default, we set  $K$  as 5. When tuning  $K$ , we may change it from 1 – 10. For FedProx<sup>2</sup>, we search  $\mu$  from  $\{0.01, 0.1, 1, 10\}$  to find its best value 1. For FedAvgM, we search the momentum value  $\beta$  of Nesterov accelerated gradient from  $\{0, 0.1, 0.5, 0.9, 0.99\}$  and find its best value 0.1. For pFedMe<sup>3</sup>, we search  $\lambda$  from  $\{0.1, 1, 10, 100\}$  to find its best value 1. The local iterations  $K$  to find a  $\delta$ -approximation personalized models is set to 5, and the coefficient of smooth aggregation  $\beta$  is set to 1, which is the same as the original paper. For FedFomo<sup>4</sup>, we set the number of models downloaded as  $M = 5$  which is recommended in the paper; For ClusterFL<sup>5</sup>, we use the same values of tolerance as the ones used in its official implementation for EMNIST dataset, and the hyper-parameters of this algorithm are fine-tuned on the other datasets. For FedAMP, we employ the grid search technology to tune the hyper-parameters, and finally the hyper-parameters are set as:  $\sigma = 1, \alpha = 1$  and  $\lambda = 0.1$ . Here, we use the same symbols as those in the original paper. For FedRoD, we use the last linear layer as the Personalized-head layer. For SFL<sup>6</sup>, we set the same hyper-parameter settings as these in its official implementation. For the implementation of FedAMP and FedRoD, we make use of open source libraries<sup>7</sup>.

## C Additional Experimental Results

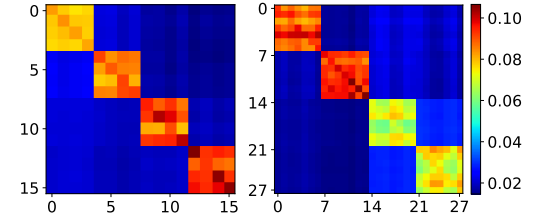
### C.1 Personalized Weighting

In order to show how FedDWA allows clients to find their optimal personalized models by properly selecting other clients, we next visualize the aggregation weights  $p_{i,k}$  computed by FedDWA. We utilize the practical heterogeneous setting 1 in which clients are divided into multiple groups. The data distribution of clients in the same group is similar, while the data distribution of clients in different groups is different. Specifically, for clients in the same group, 80% of data samples of

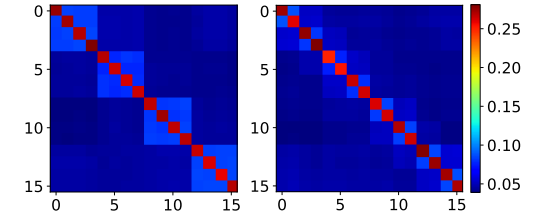
every client are uniformly sampled from a set of dominating classes, and 20% of data samples are uniformly sampled from the rest of classes. We depict clients with the same local data distributions next to each other (e.g. clients 0-4 are in the same group that have similar data distribution and client 5-9 are in the same next group). In Figure 4, we show the aggregation weights  $p_{i,k}$  computed by FedDWA using EMNIST and CIFAR10 datasets. We test the stability of FedDWA with different clients and different data distributions. The experimental results show that our method is still effective in the case of different number of clients and different number of clusters.



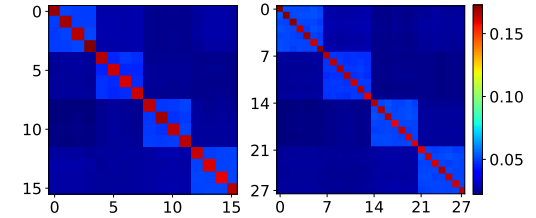
(a) Support for different numbers of distributions using EMNIST dataset.



(b) Robustness to number of clients using EMNIST dataset.



(c) Support for different numbers of distributions using CIFAR10 dataset.



(d) Robustness to number of clients using CIFAR10 dataset.

Figure 4: The visualization of the aggregation weights  $p_{i,k}$  computed by FedDWA on EMNIST and CIFAR10 datasets. The x-axis and y-axis means the IDs of clients.

<sup>2</sup><https://github.com/litian96/FedProx>

<sup>3</sup><https://github.com/CharlieDinh/pFedMe>

<sup>4</sup><https://github.com/NVlabs/FedFomo>

<sup>5</sup><https://github.com/felisat/clustered-federated-learning>

<sup>6</sup><https://github.com/dawenzi098/SFL-Structural-Federated-Learning>

<sup>7</sup><https://github.com/TsingZ0/PFL-Non-IID>

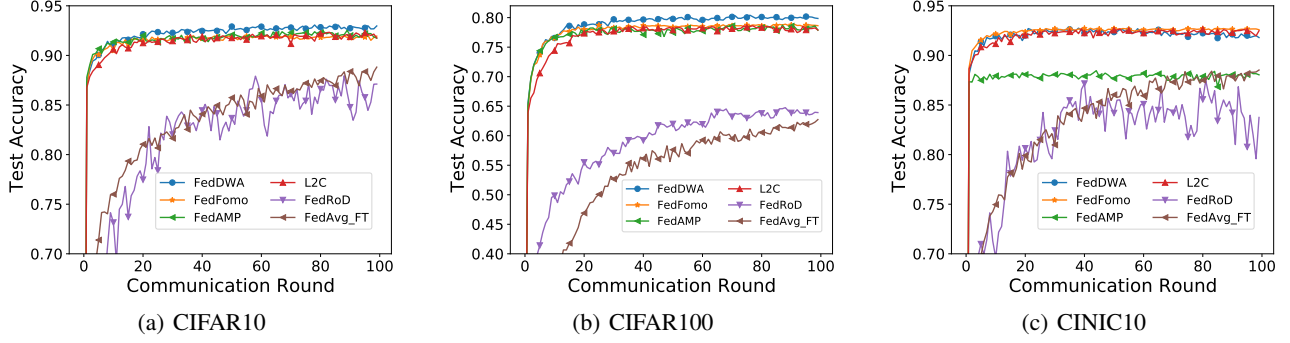


Figure 5: Test accuracy over communication rounds under the pathological heterogeneous setting with 20 clients.

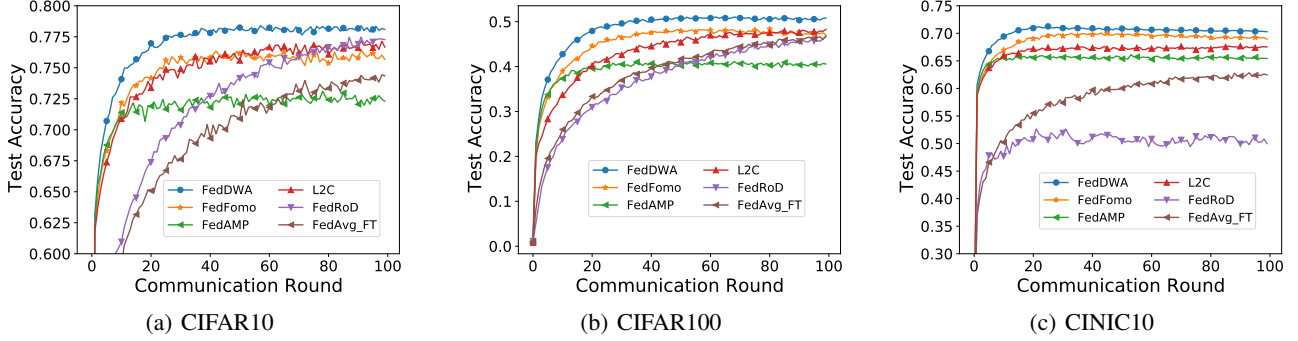


Figure 6: Test accuracy over communication rounds under the practical heterogeneous setting 1 with 20 clients.

## C.2 Performance of Guidance Model

As we have mentioned in the main body of the paper, guidance model  $\hat{w}_i^*$  represents the data distribution of client  $i$ . From this point of view, one-step ahead adaptation is a reasonable approximation.  $\hat{w}_i^*$  takes one-step ahead of time such that it can instruct client  $i$  to identify other clients which should be assigned with higher weights for model aggregation. Through experiments, we will show that  $\hat{w}_i^*$  is effective because its performance is similar to that of the personalized model  $w_i^t$ . Table 3 lists the average test accuracy of guidance model and personalized model after training 150 rounds, under the practical heterogeneous setting 1 and the practical heterogeneous setting 2 ( $\alpha = 0.1$ ), respectively. It can be seen that the test accuracy of both is almost the same in the practical heterogeneous setting 1. In addition, we find that the performance of the personalized model is slightly better than that of the guidance model in the practical heterogeneous setting 2, which means that when the data distribution between clients has a clustering structure, our algorithm can exactly capture this similarity to benefit clients. The training curves are presented in Figure 7 and Figure 8.

## C.3 Curve of Test Accuracy During Training

Figure 5 and Figure 6 present the evolution of average test accuracy over global communication rounds for partial experiments shown in Table 1. From which, it can be seen that our method has a significant performance improvement compared

Dataset	Practical Setting	Guidance model	Personalized model
EMNIST	Setting 1	85.96	86.00
	Setting 2	91.26	91.26
CIFAR10	Setting 1	77.86	78.67
	Setting 2	90.66	90.62
CIFAR100	Setting 1	49.66	51.10
	Setting 2	59.44	59.54
CINIC10	Setting 1	69.60	71.17
	Setting 2	87.74	87.47

Table 3: The best test accuracy (%) over four different datasets under two practical heterogeneous setting with 20 clients.

with other methods, except in the pathological heterogeneous setting with CINIC10 dataset.

## C.4 Communication Cost

The amount of data (including upload and download) each client needs to transmit per communication round by using different methods is compared in Table 4. In comparison with FedFomo and L2C, it is apparent that our method can sheer shrink communication traffic for detecting client similarity. Although FedFomo, L2C and our method incur more communication traffic than that of FedAvg and others, a notable advantage of FedFomo, L2C and our method is that they can

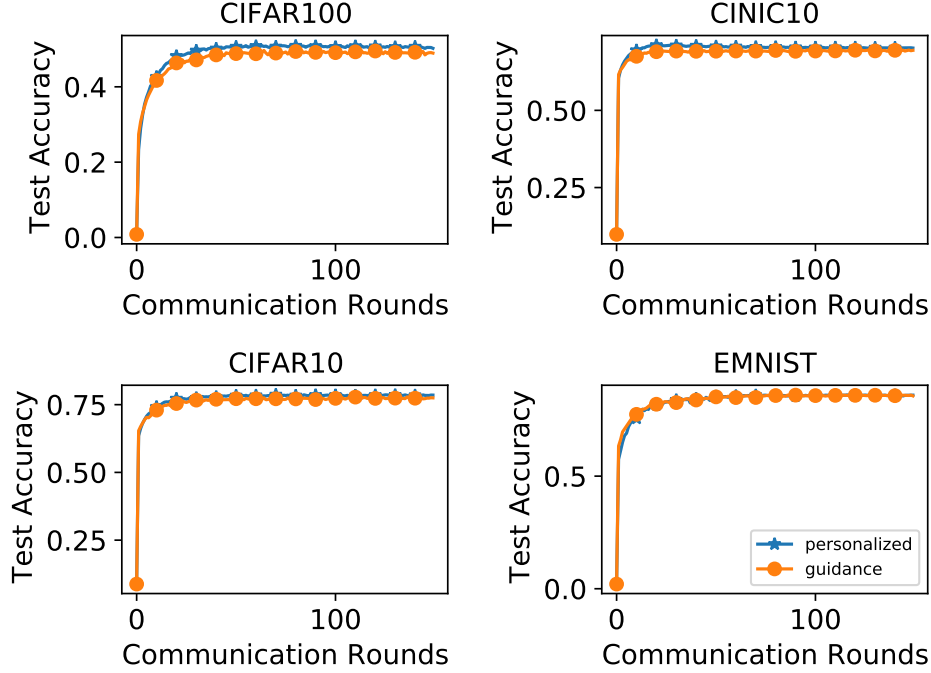


Figure 7: Test accuracy over communication rounds under the practical heterogeneous setting 1 with 20 clients.

	Communication	CIFAR100	TINY
FedAvg	$2 \times \Sigma$	9.54	53.56
Others*	$2 \times \Sigma$	9.54	53.56
FedFomo	$(1 + M) \times \Sigma$	28.62	160.68
L2C	$(1 + N) \times \Sigma$	100.17	562.38
Ours	$3 \times \Sigma$	14.31	80.34

\* Others includes FedAvgM, FedRoD, FedAvg\_FT, FedProx, per-FedAvg, pFedMe, SFL and ClusterFL.

Table 4: The amount of the communication traffic (MB) incurred by each client per round for different algorithms.  $\Sigma$  is the size of the model.  $M (M \geq 1)$  and  $N$  are the number of models downloaded by each client. We set  $M = 5$  and  $N = 20$ , as in the original papers.

explicitly measure and explain the similarity between clients, which have been explored in Figure 1.

### C.5 Computational Cost

Suppose that there are  $N$  clients participating in training and the number of model parameters is  $d$ , the computation complexity of FedDWA (ours), FedAMP and L2C are all  $\mathcal{O}(N^2d)$  in the server. For FedFomo, the extra computation is offloaded on clients, and thus its complexity in the server is  $\mathcal{O}(Nd)$ . We also test the total FLOPs for each communication round using CIFAR10 dataset, and the results can be found in Table 5. In FedDWA, the magnitude of FLOPs needed to calculate the similarity (Eq.(15)) is  $10^8$  while it is  $10^{11}$  for model training, indicating that computation load is mainly generated by model training.

Methods	FLOPs
FedAvg	$2.5 \times 10^{11}$
FedFomo	$4.6 \times 10^{11}$
FedAMP	$2.5 \times 10^{11}$
L2C	$1.2 \times 10^{12}$
FedDWA(ours)	$5.1 \times 10^{11}$

Table 5: The amount of the computational cost incurred by each client per round for different algorithms.

### C.6 More Discussion About Guidance model

#### What is guidance model.

A guidance model can facilitate the training of personalized models by enabling collaborations between similar clients. However, it's difficult to directly define the optimal guidance model since it should be the optimal personalized model, i.e., the objective of PFL. We have tried different choices of the guidance model (see next section), and one-step-ahead adaptation is the best one among our trails. The guidance model can be intuitively interpreted as follows. At the beginning of round  $t + 1$ , client  $i$  downloads the global model  $w_t$ , which actually guides the learning of client  $i$  as

$$w_i^{t+1} \leftarrow w_i^t + \sum_{j=1}^N p_{i,j} \cdot (w_j^t - w_i^t),$$

where  $w_t = \sum_{j=1}^N p_{i,j} w_j^t$  is the aggregation of models in round  $t$  and  $\sum_j p_{i,j} = 1$  is aggregation weights. How to set  $p_{i,j}$  is important to achieve PFL. Traditional FL set  $p_{i,j} = \frac{1}{N}$ ,

which is an unbiased estimate of global model and does not take into account the unique target of each client. FedFomo and L2C want to find the optimal aggregation weights to optimize personalized models on individual clients. However, as described in the paper, their methods will incur huge communication overhead with the risk of privacy leakage. Our work improves these defects by introducing a guidance model to guide the setting of  $p_{i,j}$  by constructing the optimization problem (in Eq. (9)). Through Eq. (9), we subtly offload the computation of personalized aggregation weights to the server to reduce the communication cost. Besides, our guidance model can characterize client similarity in an analytical way rather than an empirical search via the validation dataset.

#### How to select guidance model.

As we have point out in the text, we can use the last iteration model  $\hat{w}_i^{t-1}$ , the current model  $w_i^t$  or the local one-step ahead adaptation model  $\hat{w}_i^t - \eta_i^{t-1} \nabla f_i(\hat{w}_i^t)$  for the guidance model. We have conducted experiments to select the guidance model for both Pathological Setting (CIFAR10) and Practical Setting 1 (CIFAR10-V2), and the results are shown in Table 6. Here  $w$  represents the model after one-step ahead adaptation and  $w_1$  represents the model in the last iteration, and we find that the performance of using the current model  $w_i^t$  is inferior to the others, so it's not shown here. In Table 6, it can be seen that using  $w$  and  $w_1$  achieves similar results. However, if we use  $w_1$ , we have to store  $w_1$  with two possible cases: i) Store  $w_1$  in clients who need to upload two models to the server; ii) Store  $w_1$  in the server, so that each client only needs to upload one model (in this case, the amount of uplink traffic and download traffic is the same as that of FedAvg). For both cases, we need extra memory for storing the historical model  $w_1$ . Besides, if we use the last iteration model  $w_1$ , then Eq.(15) will rely on the historical information, which may lead to the cold-start problem when there are newly-joined clients. Instead, if we use  $w$ , we can save the memory overhead to achieve a similar performance and at this time, since Eq.(15) doesn't rely on historical information, FedDWA is not sensitive to the cold-start problem. The experimental results in Table 2 also confirm this point. Therefore, we finally choose  $w$ .

	CIFAR10	CIFAR10-V2	Communication
$w$	92.97%	78.56%	$3 \times \Sigma$
$w_1$	92.22%	78.99%	$2 \times \Sigma$ or $3 \times \Sigma$

Table 6: Final test accuracy and communication overhead per round.

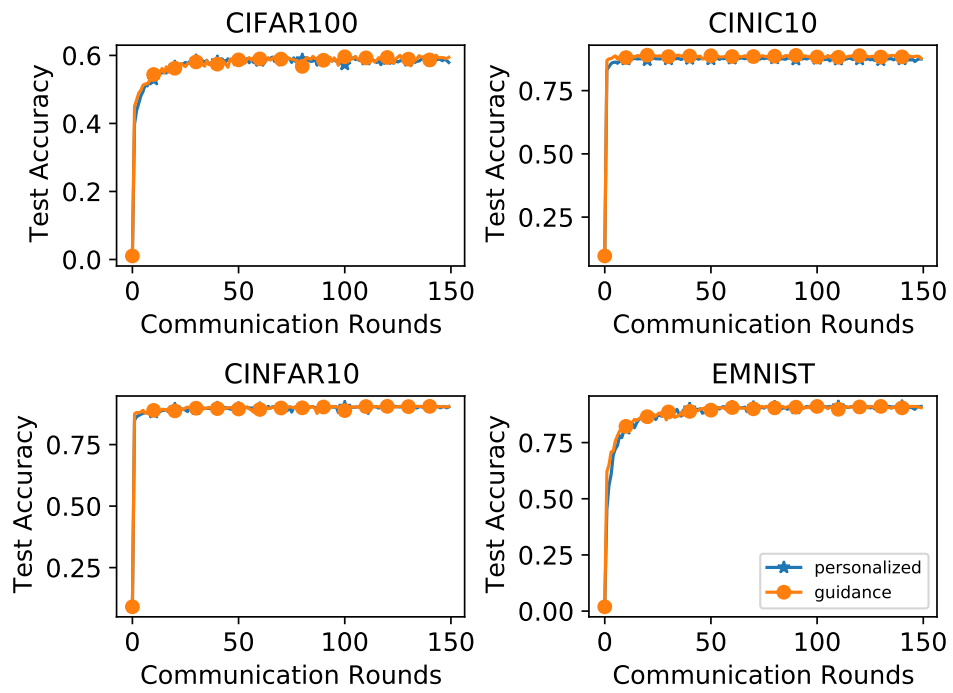


Figure 8: Test accuracy over communication rounds under the practical heterogeneous setting 2 with 20 clients.