

FedFM: Anchor-based Feature Matching for Data Heterogeneity in Federated Learning

Rui Ye, Zhenyang Ni, Chenxin Xu, Jianyu Wang, Siheng Chen *Member, IEEE*, Yonina C. Eldar *Fellow, IEEE*

Abstract—One of the key challenges in federated learning (FL) is local data distribution heterogeneity across clients, which may cause inconsistent feature spaces across clients. To address this issue, we propose a novel method FedFM, which guides each client’s features to match shared category-wise anchors (landmarks in feature space). This method attempts to mitigate the negative effects of data heterogeneity in FL by aligning each client’s feature space. Besides, we tackle the challenge of varying objective function and provide convergence guarantee for FedFM. In FedFM, to mitigate the phenomenon of overlapping feature spaces across categories and enhance the effectiveness of feature matching, we further propose a more precise and effective feature matching loss called contrastive-guiding (CG), which guides each local feature to match with the corresponding anchor while keeping away from non-corresponding anchors. Additionally, to achieve higher efficiency and flexibility, we propose a FedFM variant, called FedFM-Lite, where clients communicate with server with fewer synchronization times and communication bandwidth costs. Through extensive experiments, we demonstrate that FedFM with CG outperforms several works by quantitative and qualitative comparisons. FedFM-Lite can achieve better performance than state-of-the-art methods with five to ten times less communication costs.

Index Terms—Federated Learning, Data Heterogeneity

I. INTRODUCTION

MOST existing deep learning models are trained in a centralized manner. However, in practice, data may be distributed on several parties and may not be collected due to the increasing privacy concerns. Federated learning (FL) [1] is proposed to address this issue and has become an emerging research topic [1]–[7]. In standard FL [1], each client first downloads the same global model from the server and conducts local model training on its private dataset. Then, clients upload their trained local models to the server, where a global model is updated via local models aggregation. This process is conducted iteratively to obtain a final global model. This privacy-preserving method has been widely explored applied to many tasks, such as image classification [8], [9], language modeling [10], speech recognition [11].

One of the key challenges that hinders FL from performing as well as centralized learning is data distribution heterogeneity across clients [3], [12], [13]. Due to diverse conditions of devices and application scenarios, data might

be not independent and identical distributed (IID) across local clients. This may result in large variations in the locally trained models on clients and slow down convergence of the global model [14], [15]. This phenomenon is also referred to as *client drift* [7], [13].

To tackle the above mentioned data heterogeneity issue, most previous works [6], [7], [16] focus on model-level corrections, which intend to reduce the variations in locally trained models. However, these methods fail to ensure the consistency of multiple local models’ feature spaces. It is possible that different local models have drastically misaligned feature spaces. This could lead to unclear decision boundaries and cause misclassification, which significantly differs from centralized learning. Fig. 1a empirically shows the T-SNE [17] of two local clients’ features in FedAvg [1], where the color indicates categories and the shape indicates clients. We see that the data samples with the same color, yet different shapes do not overlap, reflecting two local models fail to share a consistent feature space. In addition, data samples with the same shape, yet different colors overlap with each other. This can be detrimental to classification tasks. Motivated by this, our work focuses on mitigating the data heterogeneity issue in federated classification tasks through aligning the feature spaces across multiple local models.

In this paper, we propose an anchor-based **Federated Feature Matching** (FedFM) method, the key idea of which is to leverage landmarks shared by all clients to provide global positioning, promoting a more consistent feature space. As a core concept of FedFM, we define landmarks as the average of features for the same class/category and name them as *anchors*. In each round of FedFM, there are two key steps: (1) anchor updating; and (2) anchor-based model updating. In the anchor updating step, first, each client calculates the local anchors; second, by interacting with the server, global anchors are updated by aggregating local anchors and sent back to each client. In the anchor-based model updating step, each client’s feature is pushed to match with the global anchor of the corresponding category during the process of local model training; see the significant improvement in Fig. 1b, where we regularize the ℓ_2 distance between a feature and its corresponding global anchor to enhance the consistency of feature spaces across clients. Global anchors are denoted by star shape.

Moreover, we conduct theoretical analysis to provide convergence guarantee for our proposed FedFM algorithm. Unlike most existing literature that analyzes fixed objective functions, the analysis of FedFM faces a distinctive challenge of time-varying objective functions over rounds. This is because of

R. Ye, Z. Ni, C. Xu are with the Cooperative Medianet Innovation Center (CMIC) at Shanghai Jiao Tong University, Shanghai, China. E-mail: yr991129, 0107nzy, xcxwakaka@sjtu.edu.cn.

J. Wang is with Meta Platforms. E-mail: jianyuwang@meta.com.

S. Chen is with Shanghai Jiao Tong University and Shanghai AI Laboratory, Shanghai, China, E-mail: sihengc@sjtu.edu.cn.

Y. C. Eldar is with Department of Computer Science and Applied Mathematics, Weizmann Institute of Science. E-mail: yonina.eldar@weizmann.ac.il.

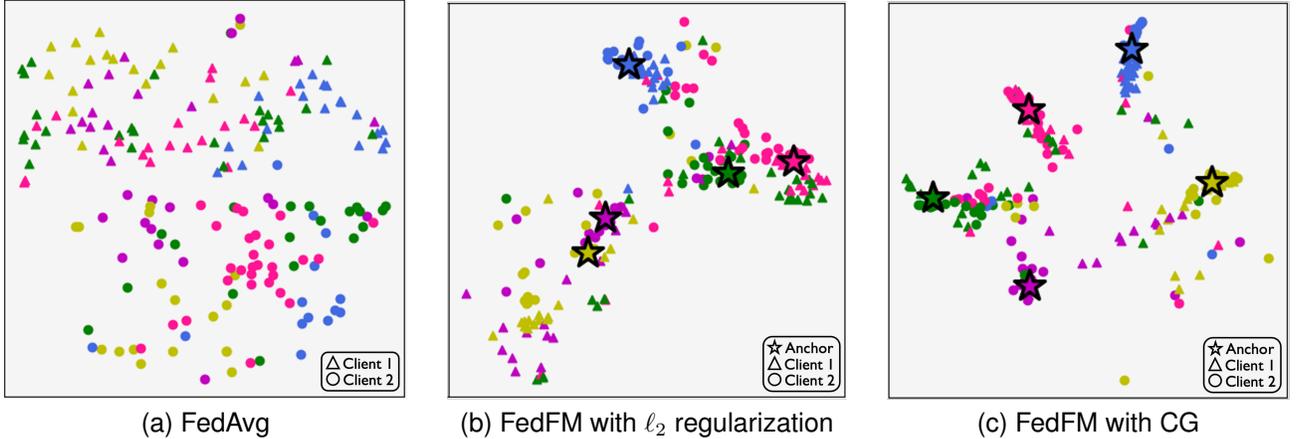


Fig. 1. FedFM alleviates the inconsistency in feature space by anchor-based feature matching. (a) In existing methods, there is a large gap between samples of two clients (triangle and circle) in the feature space. (b) With simple ℓ_2 regularization, our FedFM leverages anchors (stars) to align the feature spaces of two clients. (c) With the proposed contrastive-guiding (CG) method, FedFM achieves more precise and compact matching.

the varying global anchors, which are updated at each round. We overcome this challenge by proving a key lemma, which suggests that updating of global anchors also contributes to optimize the global objective. The theoretical results show that the proposed FedFM converges at a rate that accords with many existing theoretical optimization literature [14], [18].

To promote more precise and effective feature matching and push the feature spaces of different categories to be far away from each other, we further propose contrastive-guiding (CG) for feature matching in FedFM. The proposed CG guides each client’s local feature to match with the corresponding global anchor while keeping away from non-corresponding global anchors. Comparing with the standard ℓ_2 regularization, CG contributes to more precise feature matching, which results in more distant and compact category-wise feature space; see its improvement over ℓ_2 regularization in Fig. 1c.

To achieve higher efficiency and flexibility, we propose a variant of FedFM, called FedFM-Lite. Compared with FedFM, FedFM-Lite is more efficient since it communicates only one time within one federated round and thus requires less synchronization times (handshakes) among clients and server. FedFM-Lite is also more flexible to communicate anchors and models at different frequency. Since models have significantly more communication bandwidth cost, we propose to communicate models at a relatively lower frequency, which is capable to accord with various real-world communication budgets.

At last, through extensive experiments, we verify that FedFM with CG outperforms the state-of-the-art FL methods, including FedAvg [1], FedAvgM [19], FedProx [6], SCAF-FOLD [7], FedDyn [16], FedNova [18] and MOON [20], on multiple datasets, including CIFAR-10 [21], CINIC-10 [22] and CIFAR-100. We further visualize the feature space constructed by the proposed FedFM with CG, which qualitatively demonstrates its effectiveness. We also see that FedFM-Lite can achieve better performance than existing methods with five to ten times less communication costs and comparable performance compared with FedFM with half of the synchronization times.

Our main contributions are as follows:

- 1) We propose an anchor-based federated feature matching (FedFM) method and a contrastive-guiding (CG) technique in FedFM, which pushes each client’s local feature to match with corresponding shared global anchor while keeping away from non-corresponding anchors, promoting a consistent feature space across clients and mitigating the notorious data heterogeneity issue;
- 2) We tackle the distinctive challenge of varying objective function in the theoretical analysis of FedFM and provide a convergence guarantee;
- 3) We propose an efficient and flexible variant, FedFM-Lite, which can be easily adjusted to accord with various real-world communication budgets;
- 4) We conduct extensive experiments and show that FedFM with CG (and FedFM-Lite) can significantly outperform state-of-the-art methods.

This paper is organized as the following. Section II reviews related works. Section III presents several preliminaries including notations and motivations. Section IV describes and discusses our proposed FedFM method and the CG technique. Section V provides convergence analysis for FedFM. Section VI proposes a variant of FedFM, FedFM-Lite. Section VII shows the experimental results.

II. RELATED WORK

Federated learning (FL) is proposed in [1] and has been widely applied to many fields. In image processing, it is widely adopted since it takes advantage of the computational ability and locally-stored data of edge devices [8], [9]. It also attracts much attention in healthcare due to its privacy-preserving property [23]–[25]. But when the standard FL method FedAvg [1] meets the situation of data distribution heterogeneity across clients, the global model could move far away from the true global optimum due to the large variations in each local optima, which is referred to as client drift [7]. There have been numerous works trying to tackle this issue and two key approaches are local correction and global adjustment.

A. Local Correction

One main approach is conducting correction during the process of local model training, which aims to reduce the difference among trained local models. Most previous works conduct correction on the *model-level*. FedProx [6] applies a l_2 -norm distance regularization between the current local model and the previous global model. FedDyn [16] proposes a dynamic regularizer to align local and global solutions. Variance reduction methods, such as SCAFFOLD [7] and VRLSGD [26], utilize the previous difference between local and global gradient to debias the gradient at each local training step. MOON [20] maximizes the similarity between the feature (intermediate layer output) of current local model and that of the previous global model, which requires three times computation cost and has no convergence guarantee. Our proposed FedFM conducts local correction in the *feature-level*, which employs shared category-wise global anchors to guide local feature learning. That is, MOON [20] aligns features of two models that belong to the same sample and FedFM aligns features of all samples that belong to the same category. We also provide a convergence guarantee of FedFM.

B. Global Adjustment

Another key direction is global adjustment during the process of model interaction, which aims to obtain a better global model utilizing the uploaded local models. FedAvgM [19] introduces momentum to global model updating, which stabilizes the global model optimization. FedNova [18] normalizes local updates according to the number of SGD steps, which eliminates objective inconsistency and achieves fast convergence. Using the Knowledge-Distillation technique, FedGen [27] learns a generator to assist local model training. FedDF [28] and FedFTG [29] refine the global model by learning from the uploaded local models. Our proposed FedFM is orthogonal to these global adjustment methods and can be easily incorporated with these techniques.

As the above local correction and global adjustment methods, the focus of this paper is generalized FL, which aims for collaboratively training a global model. Personalized FL aims at collaboratively training multiple personalized local models, including FedRep [30], FedAMP [31], pFedMe [32] and Personalized FedAvg [33]. Targeting personalized FL, FedProto [34] utilizes prototype to provide extra feature information from other clients to enhance personalization of each client. In comparison, our FedFM targets generalized FL, which uses anchors as landmarks to align clients' category-wise feature spaces to enhance generalization of all clients. We also propose a new contrastive-guiding (CG) technique. CG pushes local feature close to corresponding anchor and keeps it far away from non-corresponding anchors, which is shown to be significantly effective.

We compare FedFM with several representative methods in generalized FL in Table I.

III. PRELIMINARIES

In this section, we present several key notations and the general process of FL. Then, we demonstrate two key empirical observations through preliminary experiments, including

TABLE I
RELATED WORK COMPARISONS. CONV. DENOTES CONVERGENCE GUARANTEE. MEM. AND BAND. DENOTE MEMORY COST AND BANDWIDTH COST ROUGHLY COMPARED WITH FEDAVG.

| Method | Local Correction | Conv. | Mem. | Band. |
|--------------|------------------|-------|------|-------|
| FedAvg [1] | - | ✓ | × 1 | × 1 |
| FedAvgM [19] | - | ✓ | × 1 | × 1 |
| FedProx [6] | Model | ✓ | × 2 | × 1 |
| SCAFFOLD [7] | Model | ✓ | × 2 | × 2 |
| FedDyn [16] | Model | ✓ | × 2 | × 1 |
| FedNova [18] | - | ✓ | × 1 | × 1 |
| MOON [20] | Feature | - | × 3 | × 1 |
| FedFM (ours) | Feature | ✓ | × 1 | × 1 |

inconsistent feature spaces across clients and overlapping feature spaces across categories, which motivate the proposal of our FedFM method and CG loss.

A. Notations

Suppose there are K clients, where the k th client holds a local dataset $\mathcal{B}_k = \{(\mathbf{x}_i, c_i) | i = 1, 2, \dots, |\mathcal{B}_k|\}$, where \mathbf{x}_i and c_i are the data and the label of the i th sample, respectively. FL aims to leverage the local datasets at multiple clients to collectively train a global model \mathbf{w} in the server without sharing raw data [1]. Here we focus on a C -classification task and each local dataset \mathcal{B}_k can be further split to C category-wise sub-datasets, each of which is $\mathcal{B}_{k,c} = \{(\mathbf{x}_i, c_i) \in \mathcal{B}_k | c_i = c\}$. Let $f_{\text{full}}(\cdot, \cdot)$ be an end-to-end classification model with $f_{\text{full}}(\mathbf{w}, \mathbf{x}) \in \mathbb{R}^C$ the final classification output given the input data sample \mathbf{x} and model parameters \mathbf{w} . We also consider the intermediate features as $f_{\text{mid}}(\mathbf{w}, \mathbf{x}) \in \mathbb{R}^d$, where $f_{\text{mid}}(\cdot, \cdot)$ denotes the feature-extract module in the full classification model $f_{\text{full}}(\cdot, \cdot)$. A standard global objective of FL is

$$F(\mathbf{w}) = \sum_{k=1}^K p_k F_k(\mathbf{w}) = \sum_{k=1}^K \frac{p_k}{|\mathcal{B}_k|} \sum_{(\mathbf{x}, c) \in \mathcal{B}_k} \ell(f_{\text{full}}(\mathbf{w}, \mathbf{x}), c), \quad (1)$$

where $p_k = |\mathcal{B}_k| / \sum_{k=1}^K |\mathcal{B}_k|$ is the aggregation weight of the k th client and $\ell(\cdot)$ is the task-specific loss function. To optimize the objective in a federated setting, at each communication round t , each client k downloads the same global model $\mathbf{w}^{(t)}$ and conducts τ iterations of SGD on it, which normally applies supervision on final output $f_{\text{full}}(\mathbf{w}, \mathbf{x})$ while we could also consider supervision on the feature $f_{\text{mid}}(\mathbf{w}, \mathbf{x})$. Then, each client k uploads the updated model $\mathbf{w}_k^{(t, \tau)}$ to the server, which is aggregated to update the global model $\mathbf{w}^{(t+1)}$ for next round.

Since each local dataset \mathcal{B}_k could have different data distributions, conventional method (1) could result in divergent global model and degraded performance. In this work, our goal is to introduce a regularization term to the global objective, mitigating the effects of data heterogeneity.

B. Motivation

Most FL methods do not explore clients' behavior in feature space. We conduct the following FL (FedAvg [1]) experiment on CIFAR-10 with two clients. Each holds an imbalanced dataset of size 5000, where Client 1 has 50% of data in

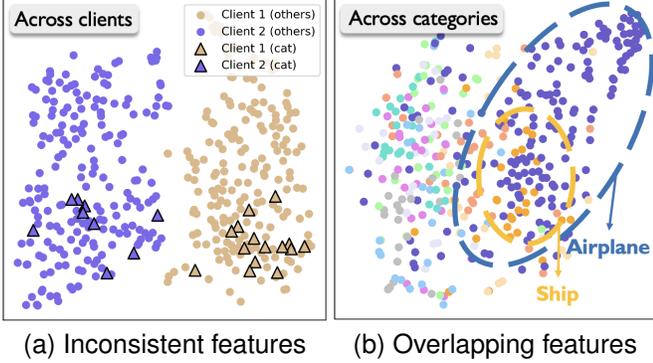


Fig. 2. Motivating examples. Fig. 2a visualizes the scatter plot of intermediate features of all the samples across two clients, where each color denotes one client. The triangles are the features of cat category. This shows that although each client has learned well-clustered features, the features are quite inconsistent across clients. Fig. 2b shows the scatter plot of intermediate features of all samples in Client 1, while each color denotes one category. This shows that airplane category occupies larger feature space and overlap with that of ship category.

airplane category while the rest is equally distributed to 9 other categories (car category for Client 2). Each client runs 10 epochs of local model training based on the same initial model. We use a ResNet18 [35] model to implement $f_{\text{full}}(\cdot, \cdot)$, where the feature-extraction module $f_{\text{mid}}(\cdot, \cdot)$ is ResNet18 without the last fully-connected layer. We then visualize intermediate features of several validation sample by T-SNE [17]. From the experimental results, we notice two unsatisfying phenomena in feature space that might cause bad performance in FL.

Inconsistent feature spaces across clients. Fig. 2a illustrates the scatter plot of several samples' features from two clients, where two colors indicate two different clients and cat category is highlighted in the triangle shape. We see that the samples from cat category across two local clients have a huge gap, reflecting that the feature spaces of two local clients are seriously inconsistent. This distinctly differs from centralized training, where samples of the same category should be gathered without an obvious gap. The intuition behind this phenomenon might be that two clients have two significantly different data distributions, resulting in distinct local models and therefore inconsistent behaviors in feature space. Motivated by this, we propose FedFM (anchor-based federated feature matching) to align the category-wise feature spaces across clients. The core idea is to establish shared global anchors as the landmarks in the feature space to guide feature learning across multiple clients; see details in Section IV-B.

Overlapping feature spaces across categories. Fig. 2b illustrates the scatter plot of several samples' features from various categories within Client 1, where various colors indicate different categories. We see that the samples from two categories (airplane and ship) greatly overlap, causing misclassification. This phenomenon often happens when the sample sizes across multiple categories are highly imbalanced. Motivated by this, we further propose a contrastive-guiding method in FedFM, which pushes each feature close to its corresponding anchor while keeping far away from non-

corresponding anchors to avoid overlapping. This also ends up enlarging the distance between two distinct categories in the feature space and further mitigates overlapping; see details in Section IV-C.

IV. METHODOLOGY

This section introduces the proposed federated learning with anchor-based feature matching (FedFM) from both aspects of mathematical optimization and federated implementation. Based on the proposed framework, we further propose a contrastive-guiding (CG) loss to mitigate overlapping feature spaces across categories. Finally, we discuss the communication cost and privacy concerns.

A. Optimization Problem

To address the issue of inconsistent feature spaces across clients, we propose anchor-based feature matching, which introduces anchors to serve as the shared landmarks for aligning all the clients' feature spaces. Mathematically, let $\mathcal{A} = \{\mathbf{a}_c\}_{c=1}^C$ be a global anchor set, where \mathbf{a}_c is the anchor of the c th category. The overall optimization problem with respect to the model parameter \mathbf{w} and the anchor set \mathcal{A} is

$$\begin{aligned} \min_{\mathbf{w}, \mathcal{A}} \Phi(\mathbf{w}; \mathcal{A}) &= \min_{\mathbf{w}, \mathcal{A}} \sum_{k=1}^K p_k \Phi_k(\mathbf{w}; \mathcal{A}) \\ &= \min_{\mathbf{w}, \mathcal{A}} \sum_{k=1}^K p_k \left(F_k(\mathbf{w}) + \lambda Q_k(\mathbf{w}; \mathcal{A}) \right), \end{aligned} \quad (2)$$

where p_k is the predefined aggregation weight of the k th client, with relative dataset size a standard choice $|\mathcal{B}_k| / \sum_{k=1}^K |\mathcal{B}_k|$, $\Phi_k(\cdot)$ is the k th client's objective, λ is a hyperparameter to balance the task-specific loss and the regularization term, $F_k(\mathbf{w}) = \sum_{(\mathbf{x}, c) \in \mathcal{B}_k} \ell(f_{\text{full}}(\mathbf{w}, \mathbf{x}), c) / |\mathcal{B}_k|$ is the task-specific loss at the k th client with \mathcal{B}_k the k th client's local dataset and

$$\begin{aligned} Q_k(\mathbf{w}; \mathcal{A}) &= \sum_{(\mathbf{x}, c) \in \mathcal{B}_k} \frac{1}{|\mathcal{B}_k|} q(f_{\text{mid}}(\mathbf{w}, \mathbf{x}), \mathcal{A}|c) \\ &= \sum_{(\mathbf{x}, c) \in \mathcal{B}_k} \frac{1}{|\mathcal{B}_k|} \|f_{\text{mid}}(\mathbf{w}, \mathbf{x}) - \mathbf{a}_c\|_2^2 \end{aligned} \quad (3)$$

is the k th client's anchor-based feature matching term, which forces each sample to match with the corresponding category-wise global anchor at each client. Since each global anchor is the proxy of each category and is shared across all the clients, with the anchor-based matching, the feature space at each client is evolving towards the same formation, enhancing the feature consistency across clients. Without anchor-based regularization term $Q_k(\mathbf{w}; \mathcal{A})$, the overall objective degenerates to the standard federated learning objective.

To solve the optimization (2), we sequentially optimize the anchor set \mathcal{A} and the model parameter \mathbf{w} at each round t .

a) *Optimizing global anchors \mathcal{A} :* Fixing the model parameter at the previous round, $\mathbf{w}^{(t)}$, we optimize over the the anchor set \mathcal{A} by solving

$$\mathcal{A}^{(t)} = \arg \min_{\mathcal{A}} \Phi(\mathbf{w}^{(t)}; \mathcal{A}) = \sum_{k=1}^K p_k \Phi_k(\mathbf{w}^{(t)}; \mathcal{A}). \quad (4)$$

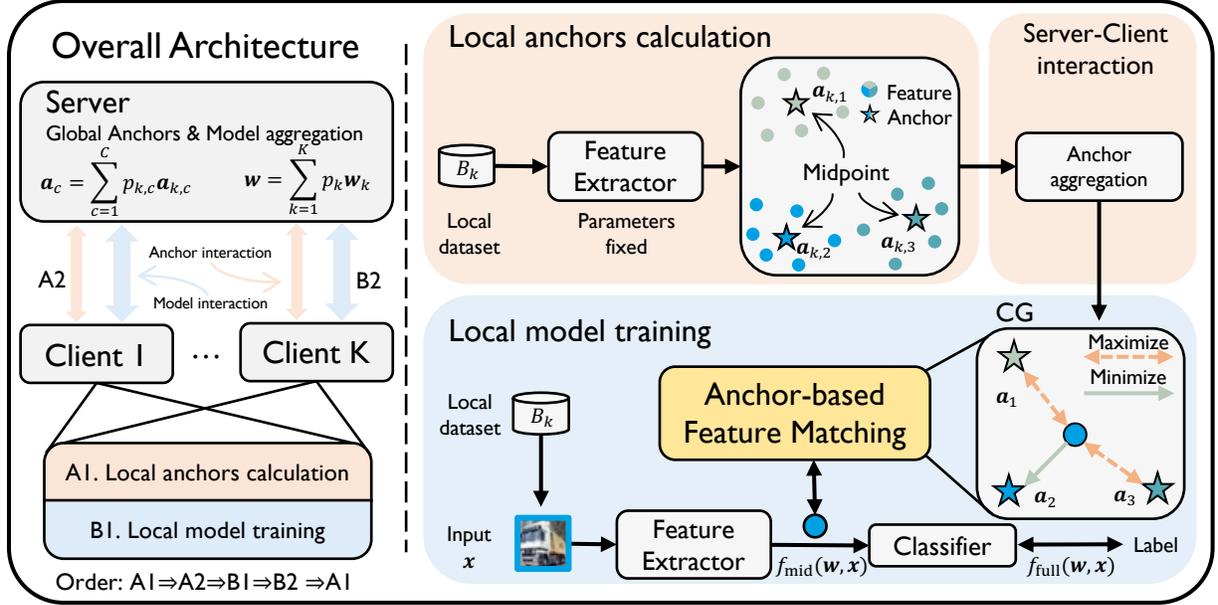


Fig. 3. Overview of FedFM for a 3-classification task. The left shows the overall architecture. The right shows two key steps in detail, where local anchors calculation generates feature anchors for each category and local model training utilizes these anchors to conduct anchor-based feature matching (e.g. contrastive guiding in the figure, which is described in Section IV-C). Here, a feature (in circle shape) is the intermediate layer output of a model and an anchor (in star shape) is an integration of features that belong to the same category.

Since the task-specific loss has nothing to do with the anchors, the optimal anchor only relates to the anchor-based regularization term. Furthermore, the global anchor of the c th category only depends on the data sample belonging to the c th category. Mathematically, the global anchor of the c th category has a straightforward closed-form solution as

$$\begin{aligned} \mathbf{a}_c^{(t)} &= \arg \min_{\mathbf{a}} \sum_{k=1}^K \sum_{(\mathbf{x}, c) \in \mathcal{B}_k} \left\| f_{\text{mid}}(\mathbf{w}^{(t)}, \mathbf{x}) - \mathbf{a}_c \right\|_2^2 \\ &= \frac{1}{\sum_{k=1}^K |\mathcal{B}_{k,c}|} \sum_{k=1}^K \sum_{(\mathbf{x}, c) \in \mathcal{B}_{k,c}} f_{\text{mid}}(\mathbf{w}^{(t)}, \mathbf{x}), \end{aligned} \quad (5)$$

where $\mathcal{B}_{k,c}$ is the k th client's local dataset that belongs to the c th category. However, in the federated learning setting, since the global server cannot directly access the local data, each global anchor cannot be directly computed in the global server as in (5). In Section IV-B, we will consider a federated implementation.

b) Optimizing global model \mathbf{w} : Fixing the global anchors $\mathcal{A}^{(t)}$, we optimize over the the model parameter \mathbf{w} by solving

$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} \Phi(\mathbf{w}; \mathcal{A}^{(t)}) = \sum_{k=1}^K p_k \Phi_k(\mathbf{w}; \mathcal{A}^{(t)}). \quad (6)$$

We can consider an iterative solver based on the standard gradient descent. Mathematically, the global model parameters can be updated as

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \sum_{k=1}^K p_k \left(\frac{\partial F_k(\mathbf{w})}{\partial \mathbf{w}} + \lambda \frac{\partial Q_k(\mathbf{w}; \mathcal{A}^{(t)})}{\partial \mathbf{w}} \right), \quad (7)$$

where η is the step size. Similarly to the optimization of global anchors, the model parameters cannot be directly updated in

the global server as in (7). A federated implementation is introduced in next.

B. Federated Implementation

Previously, we propose the mathematical optimization of federated learning with anchor-based feature matching. However, the solver is impractical due to the federated setting. Here we introduce a federated implementation, called FedFM.

1) Overview: Fig. 3 overviews the proposed FedFM. In each communication round, it consists of two main steps: anchor updating and model updating, where anchor updating solves the subproblem (4) and model updating solves the subproblem (6).

In the step of **anchor updating**, after downloading the same global model, each client calculates local anchors and uploads them to the server, where global anchors are updated by aggregating local anchors. These global anchors are then broadcast to be shared by all clients. In the step of **model updating**, each client conducts several iterations of local model training supervised by task-driven loss as well as an anchor-based feature matching loss, after which the updated local model is uploaded to server. Then, the server updates the global model by aggregating local models.

We now illustrate these two steps in detail.

a) Step 1: Anchor Updating: Here we aim to implement (5) in a federated fashion, which requires the coordination of both the clients and the server. This involves two substeps: local anchors calculation for integrating features of the same category, working on the client side, and global anchors aggregation for aggregating anchors from all clients, working on the server side.

Local anchors calculation. After downloading the global model, at the start of each new round, each client conducts

Algorithm 1 FedFM

Initialization: Global model $\mathbf{w}^{(0)}$.
for $t = 0, 1, \dots, T - 1$ **do**
 Sends global model $\mathbf{w}^{(t)}$ to initialize each client $\mathbf{w}_k^{(t,0)}$
 $\mathcal{A}_k^{(t)} \leftarrow$ **Local Anchors Calculation** ($\mathbf{w}_k^{(t,0)}$) using (8)
 $\mathcal{A}^{(t)} \leftarrow$ **Global Anchors Aggregation** ($\{\mathcal{A}_k^{(t)}\}_{k=1}^K$) using (9)
 $\mathbf{w}_k^{(t,\tau)} \leftarrow$ **Local Model Training** ($\mathbf{w}_k^{(t,0)}, \mathcal{A}^{(t)}$) for τ iterations using (10)
 $\mathbf{w}^{(t+1)} \leftarrow$ **Global Model Aggregation** ($\{\mathbf{w}_k^{(t,\tau)}\}_{k=1}^K$) using (11)
end for
return final global model $\mathbf{w}^{(T)}$

local anchors calculation by computing the category-wise midpoints of features; that is, local anchors are integration of features from the same category in each client. Mathematically, let $\mathbf{w}_k^{(t,r)}$ be the k th client's model parameter at the t th communication round with iteration r and $\mathbf{w}_k^{(t,0)} := \mathbf{w}^{(t)}$, which means that the local client's model parameter at iteration 0 is initialized by the global model in the previous communication round. Then, the local anchor of the category c in client k at round t is calculated as

$$\mathbf{a}_{k,c}^{(t)} = \frac{1}{|\mathcal{B}_{k,c}|} \sum_{(\mathbf{x},c) \in \mathcal{B}_{k,c}} f_{\text{mid}}(\mathbf{w}^{(t)}, \mathbf{x}) \in \mathbb{R}^d, \quad (8)$$

where $f_{\text{mid}}(\mathbf{w}, \mathbf{x})$ is the intermediate layer output (feature) of a model \mathbf{w} given a sample \mathbf{x} . Here $\mathbf{a}_{k,c}^{(t)} \in \mathbb{R}^d$ is essentially the midpoint of features of category c data in sub-dataset $\mathcal{B}_{k,c}$. Client k performs the calculation for each category c and then sends all C local anchors to the server.

Global anchors aggregation. The server conducts global anchors aggregation by aggregating local anchors from clients so that global anchors are integration of features from the same category across all clients. Mathematically, receiving the local anchors from all the K clients, the server conducts the following dataset size weighted aggregation for each category c to obtain global anchor; that is,

$$\begin{aligned} \mathbf{a}_c^{(t)} &:= \frac{1}{\sum_{k=1}^K |\mathcal{B}_{k,c}|} \sum_{k=1}^K |\mathcal{B}_{k,c}| \mathbf{a}_{k,c}^{(t)} \\ &= \frac{1}{\sum_{k=1}^K |\mathcal{B}_{k,c}|} \sum_{k=1}^K \sum_{(\mathbf{x},c) \in \mathcal{B}_{k,c}} f_{\text{mid}}(\mathbf{w}^{(t)}, \mathbf{x}) \in \mathbb{R}^d. \end{aligned} \quad (9)$$

We see that each global anchor $\mathbf{a}_c^{(t)} \in \mathbb{R}^d$ still matches with the optimized results in (5) and is the midpoint of features of all data that belongs to $\mathcal{B}_{k,c}$, $k \in \{1, 2, \dots, K\}$. All these global anchors $\mathcal{A}^{(t)} = \{\mathbf{a}_c^{(t)}\}_{c=1}^C$ are then broadcast to be shared by all clients. At this point, the acquired shared anchors are representative of the whole dataset. These are then used to guide the feature learning at each client's local model training.

b) Step 2: Model Updating: Here we implement the part of the update in (7) in a federated fashion, which generally requires more than one step of SGD update for each round. This involves two substeps: local model training for updating

local models, working on the client side, and global model aggregation for aggregating local models from clients, working on the server side.

Local model training. Each client conducts local model training on its private dataset with the task supervision and anchor-based feature matching loss. The k th local client's model parameter at the communication round t with iteration r is updated as

$$\mathbf{w}_k^{(t,r+1)} = \mathbf{w}_k^{(t,r)} - \eta \left(\frac{\partial F_k(\mathbf{w})}{\partial \mathbf{w}} + \lambda \frac{\partial Q_k(\mathbf{w}; \mathcal{A}^t)}{\partial \mathbf{w}} \right), \quad (10)$$

where $F_k(\cdot)$ and $Q_k(\cdot)$ are the k th client's task-specific loss and anchor-based feature matching loss (2), respectively. After τ iterations of local training, each client k obtains a local model with parameters $\mathbf{w}_k^{(t,\tau)}$ and uploads it to the server.

Global model aggregation. The server receives and aggregates the local models from all the clients and obtains an updated global model for the next round; that is, the global model for the next round $t + 1$ is obtained as:

$$\mathbf{w}^{(t+1)} \leftarrow \sum_{k=1}^K p_k \mathbf{w}_k^{(t,\tau)}, \quad (11)$$

where p_k is the predefined aggregation weight.

2) Strengths of Feature Matching: The proposed anchor-based feature matching can significantly relieve the inconsistency phenomenon in Section III-B. Since each client's features are trained to match the shared global anchors, the discrepancy of learned feature spaces across clients can be reduced. Therefore, clients' models establish a more consistent feature space of every category. Furthermore, global anchors contain overall information since they are obtained by aggregating all local anchors. With global anchors, information of other clients is infused into each client, which provides additional guidance on local feature learning especially those categories with relatively few data samples.

C. Contrastive-Guiding Loss

1) Method: To address the problem of overlapping feature space across categories in Section III-B, we further propose contrastive-guiding (CG) loss to replace the ℓ_2 -based loss in the feature matching term (3). The idea is to force each feature to be close to the corresponding anchor while keeping far away from non-corresponding anchors. Let $\mathcal{A}^{(t)} = \{\mathbf{a}_n^{(t)}\}_{n=1}^C$ be the global anchors and $\mathbf{w}_k^{(t,r)}$ be the local model for client k at training round t and iteration r . For data sample \mathbf{x} , the feature matching loss is

$$q \left(f_{\text{mid}}(\mathbf{w}_k^{(t,r)}, \mathbf{x}), \mathcal{A}^{(t)} | c \right) = \mathcal{L}_{CE}(\mathbf{s}, c),$$

where \mathcal{L}_{CE} is the cross-entropy loss function and $\mathbf{s} = [s_1, s_2, \dots, s_C] \in \mathcal{A}^C$ is a similarity vector, whose n th element measures the distance with the n th anchor:

$$s_n = \frac{\exp(\langle \mathbf{a}_n^{(t)}, f_{\text{mid}}(\mathbf{w}_k^{(t,r)}, \mathbf{x}) / \alpha \rangle)}{\sum_{i=1}^C \exp(\langle \mathbf{a}_i^{(t)}, f_{\text{mid}}(\mathbf{w}_k^{(t,r)}, \mathbf{x}) / \alpha \rangle)},$$

with a temperature value α determining the level of concentration and $\langle \cdot, \cdot \rangle$ is inner product. By minimizing the cross

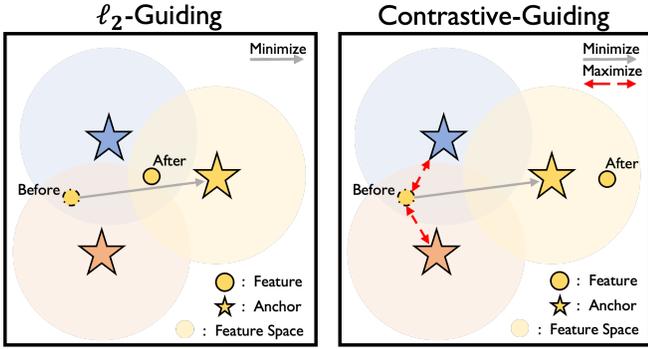


Fig. 4. Illustration of ℓ_2 -Guiding and contrastive-guiding. Each color denotes one category. ℓ_2 -Guiding only minimizes the distance between the feature and corresponding anchor, which ends up locating the feature at the feature space overlap of categories. However, contrastive-guiding also maximizes the distance between the feature and non-corresponding anchors, which feasibly locates the feature at space that merely belongs to the corresponding category.

entropy, we both maximize the similarity between the feature and its corresponding anchor $\mathbf{a}_c^{(t)}$ and minimize the similarity between this feature and each non-corresponding anchor in $\{\mathbf{a}_n^{(t)}\}_{n \neq c}$.

2) *Strengths of CG*: Fig. 4 compares the aforementioned ℓ_2 -Guiding and CG. We see that the CG loss achieves more effective feature matching and therefore better targets the overlapping phenomenon from the following two perspectives.

i) CG can provide a more precise target. As shown in the figure, there could be overlap among feature spaces of different categories. For ℓ_2 loss, only minimizing the distance between the feature and its corresponding anchor could end up locating the feature at that feature space overlap of several categories. However, CG provides a more precise target by simultaneously minimizing the distance between the feature and corresponding anchor and maximizing the distance between the feature and non-corresponding anchors, which feasibly locates the feature at space that merely belongs to the corresponding category.

ii) CG can enlarge the gap across categories. In each round, each feature is further pushed away from non-corresponding anchors, that is, more features are pushed to the non-overlap area as shown in Fig. 4. After this, each anchor (category-wise feature midpoint) is recalculated. Since more features are pushed to the non-overlap area, each recalculated anchor also moves towards the non-overlap area, which ends up enlarging the distance between anchors. This process repeats and eventually the feature space of different categories would be enlarged distinctly, which alleviates the overlap phenomenon.

D. Further Discussions

1) *Communication Cost*: The FedFM method involves two streams of communication, model parameters communication, which is required for most FL methods, and anchors communication, which is relatively negligible. Here, we take a normal setting as example, where the model is ResNet18 [35], category number $C = 10$, feature dimension $d = 512$. In this case, the anchors cost of each client is $C \times d = 5.12 \times 10^3$ units while the model cost of each client is 1.17×10^7 units.

As a result, communicating anchors only requires approximate 0.04% more bandwidth cost.

2) *Privacy Analysis*: While some feature inversion method [36] attempts to reconstruct image from a single feature, the communicated anchors are the average of a number of features, which makes the reconstruction difficult. This is also verified by [37], making FedFM a privacy-preserving method. Meanwhile, Secure Aggregation [38] is often used in practice to ensure the safety of model parameters, which can also be adopted to further secure the anchors communication. As CCVR [37] and FedFTG [29], the previously illustrated process requires uploading clients' category distributions for weighted aggregating local anchors. For cases when this is prohibited, we can directly compute the simple arithmetic mean, where clients are not asked to upload category distributions. We empirically verify that simple arithmetic mean of anchors achieves comparable performance in Section VII-D1.

V. CONVERGENCE ANALYSIS

This section provides theoretical convergence analysis of FedFM, including the required assumptions, lemmas and the derived theorem and corollary.

We provide convergence analysis of the global objective function $\Phi(\mathbf{w}; \mathcal{A})$ in (2), which relies on the following 4 assumptions. In Assumption 1, the assumption of smoothness of $F_k(\mathbf{w})$ is used in standard analysis of SGD and assumptions on $Q_k(\mathbf{w}; \mathcal{A})$ are additionally made since that we care about the property of $\Phi_k(\mathbf{w}; \mathcal{A})$. Assumptions 2, 3 and 4 are commonly used in the FL literature [6], [7], [14], [18], [39]. Here, ℓ_2 loss is applied for simplicity.

Assumption 1 (Smoothness). *Each loss function $F_k(\mathbf{w})$ is Lipschitz-smooth. Feature function $f_{\text{mid}}(\mathbf{w}_k, \mathbf{x})$ is Lipschitz-continuous and Lipschitz-smooth.*

Assumption 2 (Bounded Scalar). $\Phi_k(\mathbf{w}; \mathcal{A})$ is bounded below by Φ_{inf} .

Assumption 3 (Unbiased Gradient and Bounded Variance). *For each client, the stochastic gradient is unbiased: $\mathbb{E}_\xi[g_k(\mathbf{w}|\xi)] = \nabla\Phi_k(\mathbf{w}; \mathcal{A})$, and has bounded variance: $\mathbb{E}_\xi[\|g_k(\mathbf{w}|\xi) - \nabla\Phi_k(\mathbf{w}; \mathcal{A})\|^2] \leq \sigma^2$.*

Assumption 4 (Bounded Dissimilarity). *For any set of weights $\{p_k \geq 0\}_{k=1}^K$ subject to $\sum_{k=1}^K p_k = 1$, there exists constants $\beta^2 \geq 1$ and $\kappa^2 \geq 0$ such that $\sum_{k=1}^K p_k \|\nabla\Phi_k(\mathbf{w}; \mathcal{A})\|^2 \leq \beta^2 \|\nabla\Phi(\mathbf{w}; \mathcal{A})\|^2 + \kappa^2$.*

The smoothness property of $\Phi_k(\mathbf{w}; \mathcal{A})$ is necessary for convergence analysis. Since \mathcal{A} changes over communication round t , smoothness assumption on $Q_k(\mathbf{w}; \mathcal{A})$ would be too strong, which requires T assumptions. Thus, we only make one minor assumption on the feature function $f_{\text{mid}}(\mathbf{w}_k, \mathbf{x})$ in Assumption 1 and prove the smoothness of $\Phi_k(\mathbf{w}; \mathcal{A})$ as stated in Lemma 1.

Lemma 1: *The local objective function $\Phi_k(\mathbf{w}; \mathcal{A})$ is Lipschitz-smooth: $\|\nabla\Phi_k(\mathbf{x}; \mathcal{A}) - \nabla\Phi_k(\mathbf{y}; \mathcal{A})\| \leq L\|\mathbf{x} - \mathbf{y}\|$ for some L .*

The fact that \mathcal{A} changes over round t makes it challenging to prove convergence as it changes the global loss function at each round t . In Lemma 2, we show that at each communication point, the aggregation and updating of anchors reduces (or keeps) the global loss value.

Lemma 2: *The global loss function is non-increasing when updating global anchors. That is: $\Phi(\mathbf{w}^{(t+1)}; \mathcal{A}^{(t+1)}) \leq \Phi(\mathbf{w}^{(t+1)}; \mathcal{A}^{(t)})$.*

Based on this key lemma, we derive our main Theorem, which is stated as follows:

Theorem 1 (Optimization bound of the global objective function). *Under Assumptions 1 to 4, if we set $\eta L \leq \min\{\frac{1}{2\tau}, \frac{1}{\sqrt{2\tau(\tau-1)(2\beta^2+1)}}\}$, the optimization error will be bounded as follows:*

$$\begin{aligned} & \min_t \mathbb{E} \|\nabla \Phi(\mathbf{w}^{(t)}; \mathcal{A}^{(t)})\|^2 \\ & \leq \frac{4[\Phi(\mathbf{w}^{(0)}; \mathcal{A}^{(0)}) - \Phi_{inf}]}{\tau\eta T} + 4\eta L\sigma^2 \sum_{k=1}^K p_k^2 \\ & \quad + 3(\tau-1)\eta^2\sigma^2 L^2 + 6\tau(\tau-1)\eta^2 L^2 \kappa^2, \end{aligned}$$

where η is the client learning rate and τ is the number of local iterations.

Theorem 1 indicates that as $T \rightarrow \infty$, the expectation of optimization error will be bounded by a constant number for fixed η . When setting a proper learning rate η , we have the following corollary:

Corollary 1 (Convergence of the global objective function) *By setting $\eta = \frac{1}{\sqrt{\tau T}}$, FedFM will converge to a stationary point. Specifically, the bound could be further optimized as:*

$$\begin{aligned} & \min_t \mathbb{E} \|\nabla \Phi(\mathbf{w}^{(t)}; \mathcal{A}^{(t)})\|^2 \\ & \leq \frac{4[\Phi(\mathbf{w}^{(0)}; \mathcal{A}^{(0)}) - \Phi_{inf}] + 4L\sigma^2 \sum_{k=1}^K p_k^2}{\sqrt{\tau T}} \\ & \quad + \frac{3(\tau-1)\sigma^2 L^2}{\tau T} + \frac{6\tau(\tau-1)L^2 \kappa^2}{\tau T} \\ & = \mathcal{O}\left(\frac{1}{\sqrt{\tau T}}\right) + \mathcal{O}\left(\frac{1}{T}\right) + \mathcal{O}\left(\frac{\tau}{T}\right). \end{aligned}$$

This corollary indicates that as $T \rightarrow \infty$, the error's upper bound approaches 0. Also, given a finite T , there exists a best τ that minimizes the error's upper bound. These analyses show that FedFM can achieve the same convergence rate as most methods, such as FedAvg [14], [18]. Therefore, FedFM can achieve feature matching across clients without affecting its convergence. The detailed proof is included in the Appendix.

VI. AN EFFICIENT VARIANT: FEDFM-LITE

We also propose an efficient and flexible variant of FedFM, called FedFM-Lite. The previously proposed FedFM involves two separate communication flows, anchors communication before local model training and models communication after local model training. Though we have discussed that the anchors communication introduces minor bandwidth cost, FedFM requires twice handshakes between client and server within a federated round, which could be a drawback in real-world application since that each handshake requires some synchronization time. Thus, to mitigate this issue, we propose a more efficient variant, FedFM-Lite, which requires one handshake between client and server. We compare FedFM and FedFM-Lite in Fig. 5.

In FedFM-Lite, each client computes local anchors after the process of local model training and sends the local anchors

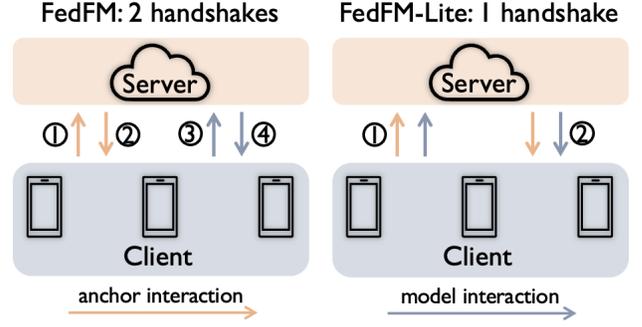


Fig. 5. Comparison between FedFM and FedFM-Lite. FedFM introduces minor bandwidth cost but requires 2 handshakes each round. FedFM-Lite further eliminates this issue, which introduces minor bandwidth cost and requires only 1 handshake.

together with model parameters. In this way, models and anchors are communicated within one handshake, which saves some synchronization time and makes it more efficient in real-world application. Beside higher efficiency, FedFM-Lite is also more flexible for real-world implementation since that models and anchors can be communicated at different frequencies. As discussed before, communicating anchors requires much less costs compared with model parameters. This motivates us to consider reducing the frequency of models communication and utilize anchors communication as compensation. Specifically, in a federated learning process that consists of T communication rounds, we can communicate anchors for each round while communicate models for every a rounds. This results in roughly a times less communication costs compared with most existing methods, such as FedAvg [1], FedProx [6], FedDyn [16] and $2 \times a$ less than SCAFFOLD [7].

VII. EXPERIMENTS

This section presents experimental details and results. We compare our proposed FedFM with state-of-the-art methods on various heterogeneous settings and datasets, which are evaluated by accuracy, feature space quality, memory and communication bandwidth cost.

A. Experimental Setup

a) Federated Setting: We set the number of clients $K = 10$ and conduct experiments on datasets including CIFAR-10 [21], CINIC-10 [22] and CIFAR-100. Here we consider three data heterogeneity (non-IID) settings. 1) NIID-1: the category distributions of clients follow a Dirichlet distribution $Dir_{10}(\beta)$, where β (default 0.5) correlates to the heterogeneity level, which is a widely considered setting [40], [41]; 2) NIID-2: each client has several dominant categories (with much more data samples) while we keep the dataset size of each client the same. We consider this setting to focus on the distribution heterogeneity but not quantity imbalance; 3) NIID-3: each client has no data sample from several categories, which is also considered in [1], [6]. Fig.9 in appendix shows the data distribution of these three Non-IID settings on CIFAR-10.

TABLE II

CLASSIFICATION ACCURACY (%) UNDER NIID-1 AND NIID-2 SETTINGS ON CIFAR-10 [21], CINIC-10 [22] AND CIFAR-100. NIID-1 IS UNDER DIRICHLET DISTRIBUTION $Dir_{10}(0.5)$ AND NIID-2 IS THE DISTRIBUTION WHERE EACH CLIENT HAS ONE DOMINANT CATEGORY. MEMORY SHOWS THE REQUIRED NUMBER OF FLOATING NUMBERS OF EACH CLIENT IN EACH ROUND ($\times 10^3$). FEDFM CONSISTENTLY OUTPERFORMS OTHER STATE-OF-THE-ART METHODS WITH RELATIVELY LESS MEMORY COST ACROSS VARIOUS SETTINGS.

| Method | CIFAR-10 | | | CINIC-10 | | | CIFAR-100 | | |
|---------------------|-------------------------|-------------------------|--------|-------------------------|-------------------------|--------|-------------------------|-------------------------|--------|
| | NIID-1 | NIID-2 | Memory | NIID-1 | NIID-2 | Memory | NIID-1 | NIID-2 | Memory |
| FedAvg [1] | 66.69 \pm 0.69 | 69.47 \pm 0.48 | 11,182 | 55.96 \pm 0.16 | 58.56 \pm 0.22 | 11,182 | 62.16 \pm 0.04 | 62.33 \pm 0.27 | 23,705 |
| FedAvgM [19] | 66.85 \pm 0.42 | 67.87 \pm 0.17 | 11,182 | 56.15 \pm 0.45 | 58.79 \pm 0.30 | 11,182 | 61.23 \pm 0.12 | 61.30 \pm 0.27 | 23,705 |
| FedProx [6] | 66.99 \pm 0.26 | 69.42 \pm 0.38 | 22,364 | 55.58 \pm 0.13 | 58.32 \pm 0.11 | 22,364 | 61.96 \pm 0.05 | 62.20 \pm 0.28 | 47,410 |
| SCAFFOLD [7] | 69.91 \pm 0.54 | 71.48 \pm 0.23 | 22,364 | 58.60 \pm 0.27 | 60.78 \pm 0.32 | 22,364 | 67.32 \pm 0.29 | 67.24 \pm 0.03 | 47,410 |
| FedDyn [16] | 68.32 \pm 0.34 | 67.63 \pm 0.16 | 22,364 | 56.71 \pm 0.50 | 59.92 \pm 0.15 | 22,364 | 43.41 \pm 0.54 | 46.44 \pm 0.87 | 47,410 |
| FedNova [18] | 66.80 \pm 0.81 | 69.45 \pm 0.49 | 11,182 | 55.67 \pm 0.24 | 58.63 \pm 0.22 | 11,182 | 62.35 \pm 0.20 | 62.31 \pm 0.26 | 23,705 |
| MOON [20] | 67.74 \pm 0.30 | 71.09 \pm 0.22 | 33,546 | 57.25 \pm 0.07 | 59.28 \pm 0.03 | 33,546 | 62.56 \pm 0.22 | 62.99 \pm 0.13 | 71,115 |
| FedFM (ours) | 72.89 \pm 0.22 | 74.52 \pm 0.21 | 11,187 | 62.56 \pm 0.40 | 65.75 \pm 0.46 | 11,187 | 71.48 \pm 0.25 | 72.13 \pm 0.45 | 23,909 |

TABLE III

CLASSIFICATION ACCURACY (%) UNDER NIID-3 SETTING ON CIFAR-10. MISSING x SETTING REPRESENTS THAT EACH CLIENT HAS NO DATA SAMPLE OF x CATEGORIES. MEMORY AND BANDWIDTH SHOW THE REQUIRED NUMBER OF FLOATING NUMBERS OF EACH CLIENT IN EACH ROUND ($\times 10^3$). OUR PROPOSED FEDFM CONSISTENTLY OUTPERFORMS OTHER STATE-OF-THE-ART METHODS WITH MINOR ADDITIONAL RESOURCE COST.

| Method | Missing 1 | Missing 2 | Missing 3 | Missing 5 | Missing 7 | Memory | Bandwidth |
|---------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|--------|-----------|
| FedAvg [1] | 70.54 \pm 0.22 | 70.50 \pm 0.24 | 69.87 \pm 0.30 | 67.25 \pm 0.54 | 59.52 \pm 0.59 | 11,182 | 11,182 |
| FedAvgM [19] | 70.02 \pm 0.40 | 69.93 \pm 0.57 | 69.34 \pm 0.37 | 67.04 \pm 0.47 | 57.08 \pm 0.66 | 11,182 | 11,182 |
| FedProx [6] | 71.16 \pm 0.42 | 70.72 \pm 0.35 | 69.82 \pm 0.23 | 67.25 \pm 0.54 | 58.58 \pm 0.23 | 22,364 | 11,182 |
| SCAFFOLD [7] | 72.67 \pm 0.39 | 72.94 \pm 0.30 | 72.60 \pm 0.22 | 71.43 \pm 0.05 | 64.28 \pm 0.60 | 22,364 | 22,364 |
| FedDyn [16] | 67.43 \pm 0.51 | 67.76 \pm 0.64 | 67.78 \pm 0.28 | 69.53 \pm 0.59 | 64.75 \pm 0.30 | 22,364 | 11,182 |
| FedNova [18] | 70.56 \pm 0.25 | 70.48 \pm 0.23 | 70.05 \pm 0.15 | 67.56 \pm 0.52 | 59.66 \pm 0.42 | 11,182 | 11,182 |
| MOON [20] | 72.64 \pm 0.25 | 72.21 \pm 0.22 | 71.57 \pm 0.23 | 68.86 \pm 0.27 | 57.80 \pm 1.02 | 33,546 | 11,182 |
| FedFM (ours) | 75.97 \pm 0.44 | 75.84 \pm 0.23 | 75.04 \pm 0.29 | 73.23 \pm 0.35 | 65.24 \pm 0.52 | 11,187 | 11,187 |

b) Implementation: We run $T = 100$ communication rounds for all experiments. In each round, every client runs for 10 local epochs with a batch size of 64. We apply ResNet18 [35] for CIFAR-10 [21] and CINIC-10 [22], ResNet50 for CIFAR-100 [21]. We use SGD optimizer with learning rate 0.01, weight decay rate $1e^{-5}$ and SGD momentum 0.9. These are commonly used experimental settings [20], [37]. For evaluation, we hold out a testing dataset at the server side and conduct the above non-IID partitions on the training set. For each client, 20% of the training set is held out for validation. We average the results on each local validation set and save the best model. Finally, we report the testing accuracy of the best model on the testing dataset.

We consider $f_{\text{full}}(\cdot, \cdot)$ as a standard ResNet and the feature extractor $f_{\text{mid}}(\cdot, \cdot)$ as $f_{\text{full}}(\cdot, \cdot)$ without the last fully-connected layer. For feature matching, the feature is normalized before applying the feature matching loss term. FedFM denotes FedFM with CG unless explicitly specified. We run FedAvg for the first T_s rounds and then launch our proposed FedFM. For all methods, we tune the hyper-parameters in a reasonable range and report the best results. Generally, for FedFM, $\lambda = 50.0$ and $T_s = 20$ is a relatively better choice.

B. Main Results

We compare FedFM with seven existing classical methods, including FedAvg [1], FedAvgM [19], FedProx [6], SCAFFOLD [7], FedDyn [16], FedNova [18] and MOON [20] on various non-IID settings and datasets. We first show accuracy comparisons quantitatively and then demonstrate qualitative comparisons in feature space by T-SNE [17] visualization.

1) Quantitative Analysis: Table II presents accuracy comparisons on three datasets under both NIID-1 and NIID-2 settings. For each entry in the table, we run three independent trials and report the mean and standard deviation results. We see that i) FedFM consistently outperforms other state-of-the-art methods on all tasks. ii) On the relatively more complicated task, NIID-1 setting on CIFAR100, the proposed FedFM significantly outperforms other methods. Specifically, compared with standard FL, FedAvg [1], FedFM achieves 9.40% higher accuracy. iii) On the six different tasks, FedFM outperforms the second-best method (SCAFFOLD [7]) 2.98%, 3.04%, 3.96%, 4.97%, 4.16%, 4.89%, respectively. It is also worth mentioning that SCAFFOLD [7] requires roughly twice the memory and communication bandwidth costs.

Table III shows the accuracy, memory, bandwidth comparisons under NIID-3 setting on CIFAR-10. In NIID-3, each client has no data sample from several (x) categories, which is denoted as Missing x setting. We conduct experiments on different $x \in \{1, 2, 3, 5, 7\}$. We see that i) the performances of all methods degrade as x increases since larger x corresponds to a more heterogeneous setting. This verifies that data heterogeneity significantly affects the performance of FL. ii) FedFM consistently outperforms other state-of-the-art methods. Specifically, it outperforms FedAvg [1] by 5.53% and SCAFFOLD [7] by 2.28% on average. iii) Compared with FedAvg [1], FedFM achieves significantly better performance while introducing minor memory and bandwidth costs. Compared with SCAFFOLD [7], FedFM achieves better performance with nearly half of the memory and bandwidth costs.

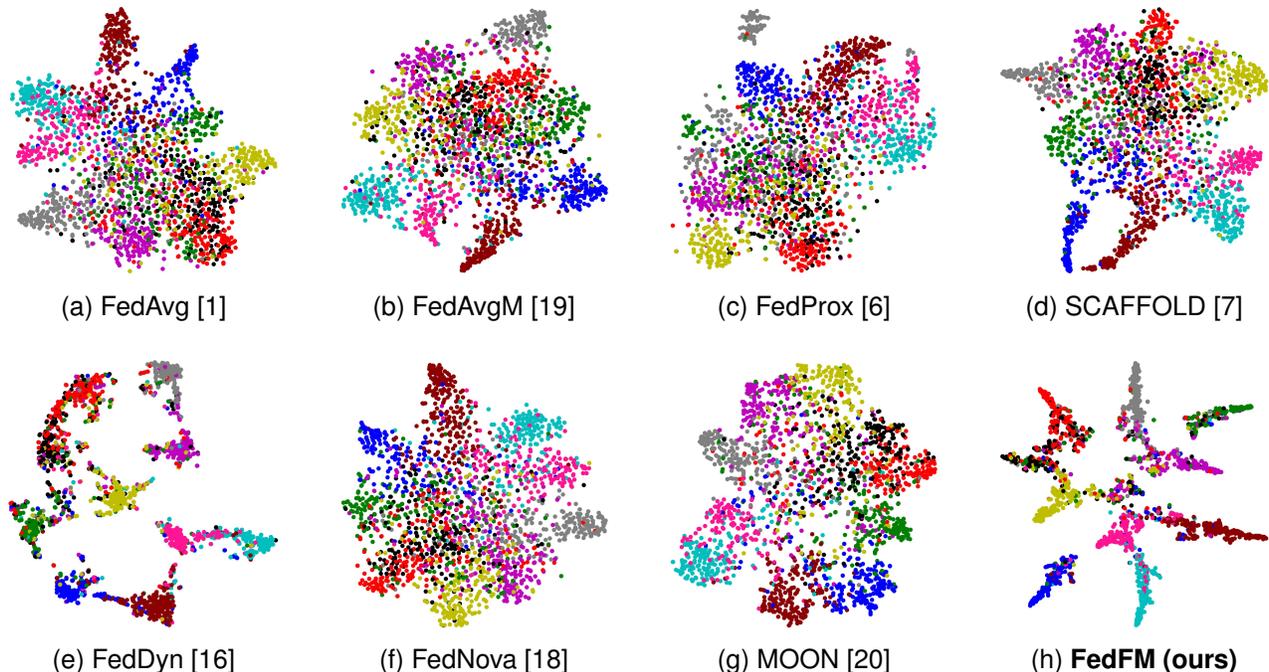


Fig. 6. Qualitative comparisons among methods through T-SNE [17] visualization. Each dot represents the feature of one data sample, whose color denotes its category. FedFM establishes the most compact and distinct clusters in feature space. FedDyn [16] has moderately good visualization results but only achieves 68.32% accuracy while FedFM achieves 72.89%.

TABLE IV
NUMERICAL QUALITY EVALUATION OF FEATURE SPACES. HIGHER NMI AND SS CORRESPOND TO HIGHER QUALITY OF FEATURE SPACES. OUR PROPOSED FEDFM ACHIEVES SIGNIFICANTLY HIGHEST NMI AND SS.

| Metric | FedAvg [1] | FedAvgM [19] | FedProx [6] | SCAFFOLD [7] | FedDyn [16] | FedNova [18] | MOON [20] | FedFM (ours) |
|--------|------------|--------------|-------------|--------------|-------------|--------------|-----------|---------------------|
| NMI | 0.413 | 0.411 | 0.397 | 0.432 | 0.485 | 0.416 | 0.481 | 0.557 |
| SS | 0.036 | 0.038 | 0.006 | 0.056 | 0.136 | 0.049 | 0.068 | 0.173 |

2) *Qualitative Analysis*: Fig. 6 presents T-SNE [17] visualization results in feature space of each method, where different color denote different category. All the sampled data is fed into the final global model of each method to obtain the corresponding features, which are then plot using T-SNE [17]. We see that i) most previous methods suffer from slack category-wise feature space while our FedFM establishes significantly more compact category-wise feature space, which reflects the effectiveness of utilizing anchors to conduct feature matching. ii) Most previous methods suffer from ambiguous boundaries. However, our FedFM establishes clusters with clear boundaries and large gap which are contributed by using anchors to attract features and the contrastive-guiding loss. These two phenomena indicate that our proposed FedFM indeed benefits the establishment of feature space and gives the evidence of significantly improved performance. Note that though FedDyn [16] seems to establish a good feature space, it only achieves a 68.32% accuracy, which is significantly lower than that of our proposed FedFM (72.89%).

For more comprehensive comparisons, we also evaluate the quality of feature space in Fig. 6 using normalized mutual information (NMI) and silhouette score (SS) [42]. NMI is capable of measuring the quality of clustering. SS is a measure of how similar an object is to its own cluster compared to other

clusters. Note that for NMI, we first apply the K-Means [43] to perform clustering on sampled features of all methods and then use NMI to measure the clustering quality. Both NMI and SS are measured using Scikit-learn [44]. Higher NMI and SS correspond to higher quality of feature spaces. We present the evaluation results in Table IV. The table shows that our proposed FedFM achieves significantly higher NMI and SS. Specifically, compared with FedDyn [16], FedFM achieves 14.8% higher NMI and 27.2% higher SS. This gives evidence for the seemingly great feature space but ordinary accuracy performance of FedDyn [16] in a way.

C. Further Comparisons

1) *Comparisons with FedProto*: Targeting a different task, personalization in FL, FedProto [34] uses prototype to provide feature information from others to enhance personalization while FedFM focuses on generalization in FL. To further verify their difference, we implement a generalized version of FedProto [34] and compare it with FedFM under NIID-1 setting on CIFAR-10 [21]. Experiments show that generalized FedProto achieves $67.33 \pm 0.49\%$ accuracy, which is outperformed by SCAFFOLD [7], FedDyn [16] and MOON [20]. Our proposed FedFM achieves $72.89 \pm 0.22\%$ accuracy, which outperforms generalized FedProto [34] by 5.56%.

TABLE V

COMPARISONS OF ACCURACY, MEMORY AND BANDWIDTH COSTS ON CIFAR-100. EACH ENTRY SHOWS CLASSIFICATION ACCURACY (%). WITHIN PARENTHESES, IT SHOWS THE REQUIRED MEMORY / BANDWIDTH COST, EVALUATED BY FLOATING NUMBERS ($\times 10^8$). WHEN $K = 100$, i) FEDFM TAKES ONLY 0.86% MORE RESOURCE OVERHEAD TO ACHIEVE 13.97% HIGHER ACCURACY THAN FEDAVG [1], ii) FEDFM ACHIEVES 7.68% HIGHER ACCURACY WITH ONLY HALF THE MEMORY AND BANDWIDTH COSTS COMPARED WITH SCAFFOLD [7].

| K | 20 | 30 | 50 | 100 |
|---------------------|--------------------------|--------------------------|------------------------------|------------------------------|
| FedAvg [1] | 58.48 (474 / 474) | 54.46 (711 / 711) | 50.20 (1,185 / 1,185) | 41.41 (2,370 / 2,370) |
| FedAvgM [19] | 58.36 (474 / 474) | 54.48 (711 / 711) | 52.86 (1,185 / 1,185) | 46.72 (2,370 / 2,370) |
| FedProx [6] | 58.27 (948 / 474) | 54.50 (1,422 / 711) | 50.55 (2,370 / 1,185) | 40.62 (4,740 / 2,370) |
| SCAFFOLD [7] | 64.65 (948 / 948) | 61.82 (1,422 / 1,422) | 56.71 (2,370 / 2,370) | 47.70 (4,740 / 4,740) |
| FedDyn [16] | 41.90 (948 / 474) | 41.13 (1,422 / 711) | 39.30 (2,370 / 1,185) | 31.21 (4,740 / 2,370) |
| FedNova [18] | 58.01 (474 / 474) | 53.83 (711 / 711) | 50.34 (1,185 / 1,185) | 42.61 (2,370 / 2,370) |
| MOON [20] | 57.63 (1,422 / 474) | 52.71 (2,133 / 711) | 47.84 (3,555 / 1,185) | 38.45 (7,110 / 2,370) |
| FedFM (ours) | 69.49 (478 / 478) | 67.70 (717 / 717) | 64.22 (1,195 / 1,195) | 55.38 (2,390 / 2,390) |

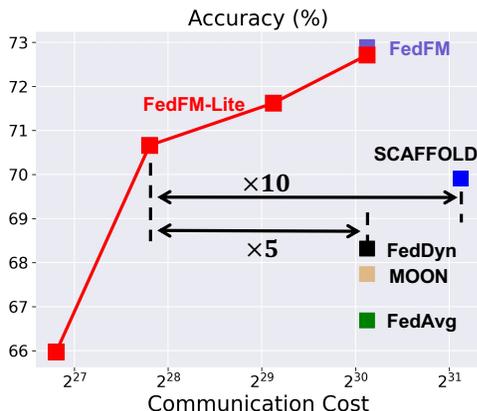


Fig. 7. Comparison of performances and communication costs among FedFM-Lite and several classical methods. Our proposed FedFM-Lite can achieve significantly better performance while saving 5 or even 10 times communication costs compared with existing techniques.

2) *Performance and Resource Costs*: FedFM introduces minor resource costs while bringing significantly better performance. To verify this point, we conduct experiments on CIFAR-100 [21] under various client numbers ($K \in \{20, 30, 50, 100\}$). Beside accuracy comparisons, we also compare the memory and bandwidth cost of these methods, which are evaluated by the number of required floating numbers ($\times 10^8$) for the overall FL process. We present the results in Table V. Experiments show that i) our proposed FedFM consistently outperforms state-of-the-art methods for various client numbers, indicating its applicability to scenario with large client number. ii) FedFM achieves significantly better performance with minor additional memory and bandwidth overhead. Specifically, FedFM takes only 0.86% more communication overhead to achieve 13.97% better classification performance than FedAvg [1] when $K = 100$. Compared with the second-best method (SCAFFOLD [7]), FedFM achieves 7.68% higher accuracy when $K = 100$ with only half the memory and bandwidth costs.

3) *Performance of An Efficient Variant: FedFM-Lite*: In Section VI, we propose an efficient variant, called FedFM-Lite, which is more efficient and flexible in practice. Since anchors require less communication bandwidth costs than model communication, we propose to reduce the communi-

TABLE VI

EFFECTS OF GLOBAL ANCHORS AGGREGATION MANNER. (WEIGHTED) DENOTES AGGREGATING CATEGORY-WISE ANCHORS ACCORDING TO EACH CLIENT’S NUMBER OF CORRESPONDING SAMPLES. (UNIFORM) DENOTES COMPUTING SIMPLE ARITHMETIC MEAN OF ANCHORS. EXPERIMENTS SHOW THAT FEDFM (UNIFORM) PERFORMS COMPARABLY WITH FEDFM (WEIGHTED) WHILE FEDFM (UNIFORM) DOES NOT NEED UPLOADING CLIENT’S CATEGORY DISTRIBUTION.

| Method | FedAvg [1] | FedFM (Weighted) | FedFM (Uniform) |
|----------|------------------|------------------|------------------|
| Accuracy | 66.69 \pm 0.69 | 72.89 \pm 0.22 | 72.87 \pm 0.26 |

cation frequency of model communication while keeping the communication frequency of anchor communication, which can be easily achieved in FedFM-Lite. This modification saves communication bandwidth costs to a large extent.

To empirically verify this efficiency and flexibility, we conduct the following experiments. We communicate anchors for each round while we communicate models every $a \in \{1, 2, 5, 10\}$ round(s), that is, larger a corresponds to less communication cost. We show the communication cost and final performance of each trial in Fig. 7. We also present several representative methods for comparison.

Experiments show that when we communicate models every $a \in \{1, 2, 5\}$ communication round(s), FedFM-Lite can significantly outperform compared methods. Specifically, when $a = 5$, FedFM-Lite outperforms FedDyn [16] by 2.34% with 5 times less communication costs and SCAFFOLD [7] by 0.75% with 10 times less communication costs. These experiments show that for bandwidth limited scenarios, FedFM-Lite can be an efficient candidate algorithm.

D. Ablation Study

1) *Effects of global anchors aggregation manner*: Here, we show that FedFM can still achieve great performance without uploading clients’ category distributions as discussed in Section IV-D2. We compare two manners of global anchors aggregation, sample-number-based weighted aggregation and uniform aggregation. For the weighted aggregation, each category-wise global anchor is updated by weighted aggregating local anchors according to each client’s number of data samples of the corresponding category. This aggregation manner might not be allowed for its requirement for uploading clients’ category distributions. For the uniform aggregation,

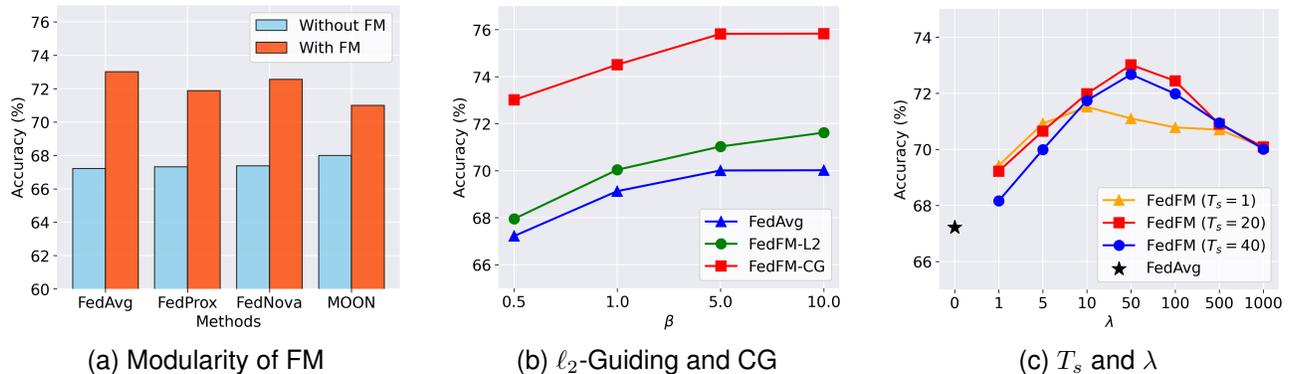


Fig. 8. Ablation study. (a) shows that FM consistently brings performance gain to four methods. (b) shows that CG significantly improve the effectiveness of FM. (c) shows that $T_s = 20$ and $\lambda = 10 \sim 100$ is roughly an optimal solution.

each category-wise global anchor is updated by uniform aggregating local anchors of the corresponding category, which relieves the above issue. As an experimental detail, for those categories where a client has no data sample, we adopt the corresponding global anchors as the local anchors for aggregation.

We conduct experiments under NIID-1 on CIFAR-10 [21] and present the results in Table VI. Experiments show that FedFM with uniform aggregation performs comparably to FedFM with weighted aggregation (only 0.02% performance drop). The reason behind this could be that for each category, all clients’ features are pushed to the same shared global anchor. As a result, all clients’ local anchors of the same category are close to each other, making it similar between applying weighted aggregation and uniform aggregation.

2) *Modularity of Feature Matching*: One advantage of our anchor-based feature matching (FM) method is its modularity, that is, it can be combined with most existing methods. Fig. 8a shows the performances before and after feature matching (FM) with contrastive-guiding combined with several existing methods. Here, we take FedAvg [1], FedProx [6], FedNova [18] and MOON [20] as example and conduct experiments under NIID-1 setting on CIFAR10 [21]. Note that the previous explored FedFM corresponds to FedAvg [1] incorporated with FM.

From the figure, we see that applying our FM consistently brings performance gain to these four methods, achieving 4.63% higher accuracy than corresponding baselines on average. Note that all these methods with FM outperforms the state-of-the-art performance 69.91% (SCAFFOLD [7]).

3) *Effects of contrastive-guiding loss*: Fig. 8b presents the performance of FedAvg [1], FedFM with ℓ_2 loss and FedFM with contrastive-guiding (CG) loss under different heterogeneous levels. Note that smaller β corresponds to more severe heterogeneous level. We see that i) the performance of all methods degrades as the heterogeneous level increases (β decreases), which verifies that data heterogeneity affects the performance of FL. ii) Both FedFM ℓ_2 and CG loss outperform baseline FedAvg [1], indicating that anchor-based feature matching brings performance improvement to standard FL. iii) FedFM with CG significantly enhances the perfor-

mance compared with FedFM with ℓ_2 , which indicates the effectiveness of our proposed CG loss.

4) *Effects of the weight of feature matching loss λ* : Fig. 8c presents the relationship between the final results and weight of feature matching loss λ . Specifically, for each curve in the figure, the launching round T_s of FedFM is fixed while the λ is tuned in $\{1, 5, 10, 50, 100, 500, 1000\}$. Note that when $\lambda = 0$, FedFM reduces to FedAvg [1], which is denoted by a star. We see that i) applying feature matching brings performance gain over FedAvg [1] for a wide range of λ , indicating the effectiveness of feature matching; ii) a moderate λ ranging from $10 \sim 100$ tends to perform better.

5) *Effects of the round T_s to launch FedFM*: Fig. 8c presents the relationship between the final results and launching round T_s of FedFM. Specifically, for each fixed λ , we compare the performance of three different $T_s \in \{1, 20, 40\}$. We see that a moderate $T_s = 20$ performs the best or comparably in most cases. This is reasonable since at the initial rounds of FL, the established anchors are less representative and still in drastic change, which makes such feature matching less effective.

We also explore the effects of the number of epochs of local model training and the performance under partial client participation scenario in Table VII and VIII in appendix.

VIII. CONCLUSION

Facing statistical data heterogeneity, there are two unsatisfying phenomena in feature space for existing federated learning methods. Motivated by this, we propose a novel anchor-based federated feature matching (FedFM) method, which utilizes shared anchors to guide feature learning at multiple local models, promoting a consistent feature space. Tackling the theoretical challenge of varying objective function, we prove the convergence of FedFM. For more precise guiding, we further propose a novel contrastive-guiding (CG) loss, which guides the feature of each sample to match with the corresponding anchor while keeping far away from non-corresponding anchors. We propose a more efficient and flexible variant of FedFM, FedFM-Lite, which is capable of communicating anchors and models at different frequency. Experiments show that FedFM

with CG (and FedFM-Lite) consistently outperform state-of-the-art methods.

REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [2] T. Gafni, N. Shlezinger, K. Cohen, Y. C. Eldar, and H. V. Poor, "Federated learning: A signal processing perspective," *IEEE Signal Processing Magazine*, vol. 39, no. 3, pp. 14–41, 2022.
- [3] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [4] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [5] V. Smith, C.-K. Chiang, M. Sanjabi, and A. S. Talwalkar, "Federated multi-task learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [6] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Proceedings of Machine Learning and Systems*, vol. 2, pp. 429–450, 2020.
- [7] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [8] T.-M. H. Hsu, H. Qi, and M. Brown, "Federated visual classification with real-world data distribution," in *European Conference on Computer Vision*. Springer, 2020, pp. 76–92.
- [9] Hsu, Tzu-Ming Harry and Qi, Hang and Brown, Matthew, "Measuring the effects of non-identical data distribution for federated visual classification," *arXiv preprint arXiv:1909.06335*, 2019.
- [10] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [11] D. Leroy, A. Coucke, T. Lavril, T. Gisselbrecht, and J. Dureau, "Federated learning for keyword spotting," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 6341–6345.
- [12] S. AbdulRahman, H. Tout, H. Ould-Slimane, A. Mourad, C. Talhi, and M. Guizani, "A survey on federated learning: The journey from centralized to distributed on-site learning and beyond," *IEEE Internet of Things Journal*, vol. 8, no. 7, pp. 5476–5497, 2020.
- [13] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, M. Al-Shedivat, G. Andrew, S. Avestimehr, K. Daly, D. Data *et al.*, "A field guide to federated optimization," *arXiv preprint arXiv:2107.06917*, 2021.
- [14] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *International Conference on Learning Representations*, 2019.
- [15] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," *arXiv preprint arXiv:1806.00582*, 2018.
- [16] D. A. E. Acar, Y. Zhao, R. Matas, M. Mattina, P. Whatmough, and V. Saligrama, "Federated learning based on dynamic regularization," in *International Conference on Learning Representations*, 2020.
- [17] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [18] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "A novel framework for the analysis and design of heterogeneous federated learning," *IEEE Transactions on Signal Processing*, vol. 69, pp. 5234–5249, 2021.
- [19] T.-M. H. Hsu, H. Qi, and M. Brown, "Measuring the effects of non-identical data distribution for federated visual classification," *arXiv preprint arXiv:1909.06335*, 2019.
- [20] Q. Li, B. He, and D. Song, "Model-contrastive federated learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 713–10 722.
- [21] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [22] L. N. Darlow, E. J. Crowley, A. Antoniou, and A. J. Storkey, "Cinic-10 is not imagenet or cifar-10," *arXiv preprint arXiv:1810.03505*, 2018.
- [23] G. A. Kaissis, M. R. Makowski, D. Rückert, and R. F. Braren, "Secure, privacy-preserving and federated machine learning in medical imaging," *Nature Machine Intelligence*, vol. 2, no. 6, pp. 305–311, 2020.
- [24] Y. Chen, X. Qin, J. Wang, C. Yu, and W. Gao, "Fedhealth: A federated transfer learning framework for wearable healthcare," *IEEE Intelligent Systems*, vol. 35, no. 4, pp. 83–93, 2020.
- [25] Z. Chen, C. Yang, M. Zhu, Z. Peng, and Y. Yuan, "Personalized retrogress-resilient federated learning towards imbalanced medical data," *IEEE Transactions on Medical Imaging*, pp. 1–1, 2022.
- [26] X. Liang, S. Shen, J. Liu, Z. Pan, E. Chen, and Y. Cheng, "Variance reduced local sgd with lower communication complexity," *arXiv preprint arXiv:1912.12844*, 2019.
- [27] Z. Zhu, J. Hong, and J. Zhou, "Data-free knowledge distillation for heterogeneous federated learning," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 878–12 889.
- [28] T. Lin, L. Kong, S. U. Stich, and M. Jaggi, "Ensemble distillation for robust model fusion in federated learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 2351–2363, 2020.
- [29] L. Zhang, L. Shen, L. Ding, D. Tao, and L.-Y. Duan, "Fine-tuning global model via data-free knowledge distillation for non-iid federated learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 174–10 183.
- [30] L. Collins, H. Hassani, A. Mokhtari, and S. Shakkottai, "Exploiting shared representations for personalized federated learning," in *International Conference on Machine Learning*. PMLR, 2021, pp. 2089–2099.
- [31] Y. Huang, L. Chu, Z. Zhou, L. Wang, J. Liu, J. Pei, and Y. Zhang, "Personalized cross-silo federated learning on non-iid data," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 9, 2021, pp. 7865–7873.
- [32] C. T. Dinh, N. Tran, and J. Nguyen, "Personalized federated learning with moreau envelopes," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 394–21 405, 2020.
- [33] A. Fallah, A. Mokhtari, and A. Ozdaglar, "Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach," *Advances in Neural Information Processing Systems*, vol. 33, pp. 3557–3568, 2020.
- [34] Y. Tan, G. Long, L. Liu, T. Zhou, Q. Lu, J. Jiang, and C. Zhang, "Fedproto: Federated prototype learning over heterogeneous devices," *arXiv preprint arXiv:2105.00243*, 2021.
- [35] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [36] N. Zhao, Z. Wu, R. W. Lau, and S. Lin, "What makes instance discrimination good for transfer learning?" in *International Conference on Learning Representations*, 2020.
- [37] M. Luo, F. Chen, D. Hu, Y. Zhang, J. Liang, and J. Feng, "No fear of heterogeneity: Classifier calibration for federated learning with non-iid data," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [38] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 1175–1191.
- [39] S. J. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, "Adaptive federated optimization," in *International Conference on Learning Representations*, 2020.
- [40] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated learning with matched averaging," in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=BkluqSFDS>
- [41] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, N. Hoang, and Y. Khazaeni, "Bayesian nonparametric federated learning of neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7252–7261.
- [42] P. J. Rousseeuw, "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis," *Journal of computational and applied mathematics*, vol. 20, pp. 53–65, 1987.
- [43] S. Lloyd, "Least squares quantization in pcm," *IEEE transactions on information theory*, vol. 28, no. 2, pp. 129–137, 1982.
- [44] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, "Scikit-learn: Machine learning in python," *The Journal of machine Learning research*, vol. 12, pp. 2825–2830, 2011.