SSFL: Tackling Label Deficiency in Federated Learning via Personalized Self-Supervision

Chaoyang He, Zhengyu Yang, Erum Mushtaq, Sunwoo Lee Mahdi Soltanolkotabi, Salman Avestimehr

Viterbi School of Engineering University of Southern California {chaoyang.he,yang765,emushtaq,sunwool,soltanol,avestime}@usc.edu

Abstract

Federated Learning (FL) is transforming the ML training ecosystem from a centralized over-the-cloud setting to distributed training over edge devices in order to strengthen data privacy, reduce data migration costs, and break regulatory restrictions. An essential, but rarely studied, challenge in FL is label deficiency at the edge. This problem is even more pronounced in FL, compared to centralized training, due to the fact that FL users are often reluctant to label their private data and edge devices do not provide an ideal interface to assist with annotation. Addressing label deficiency is also further complicated in FL, due to the heterogeneous nature of the data at edge devices and the need for developing personalized models for each user. We propose a self-supervised and personalized federated learning framework, named (SSFL), and a series of algorithms under this framework which work towards addressing these challenges. First, under the SSFL framework, we analyze the compatibility of various centralized selfsupervised learning methods in FL setting and demonstrate that SimSiam networks performs the best with the standard FedAvg algorithm. Moreover, to address the data heterogeneity at the edge devices in this framework, we have innovated a series of algorithms that broaden existing supervised personalization algorithms into the setting of self-supervised learning including perFedAvg, Ditto, and local fine-tuning, among others. We further propose a novel personalized federated selfsupervised learning algorithm, Per-SSFL, which balances personalization and consensus by carefully regulating the distance between the local and global representations of data. To provide a comprehensive comparative analysis of all proposed algorithms, we also develop a distributed training system and related evaluation protocol for SSFL. Using this training system, we conduct experiments on a synthetic non-I.I.D. dataset based on CIFAR-10, and an intrinsically non-I.I.D. dataset GLD-23K. Our findings show that the gap of evaluation accuracy between supervised learning and unsupervised learning in FL is both small and reasonable. The performance comparison indicates that representation regularization-based personalization method is able to outperform other variants. Ablation studies on SSFL are also conducted to understand the role of batch size, non-I.I.D.ness, and the evaluation protocol.

1 Introduction

Federated Learning (FL) is a contemporary distributed machine learning paradigm that aims at strengthening data privacy, reducing data migration costs, and breaking regulatory restrictions (Kairouz et al. 2021; Wang et al. 2021). It has been widely applied to computer vision, natural language processing, and data mining. However, there are two main challenges impeding its wider adoption in machine learning. One is data heterogeneity, which is a natural property of FL in which diverse clients may generate datasets with different distributions due to behavior preferences (e.g., the most common cause of heterogeneity is skewed label distribution which might result from instances where some smartphone users take more landscape pictures, while others take more photos of daily life). The second challenge is label deficiency at the edge, which is relatively less studied. This issue is more severe at the edge than in a centralized setting because users are reluctant to annotate their private and sensitive data, and/or smartphones and IoT devices do not have a userfriendly interface to assist with annotation.

To mitigate the data heterogeneity issue among clients, researchers have proposed algorithms for training a global model FedAvg (McMahan et al. 2017), FedProx (Li et al. 2018), FedNova (Wang et al. 2020), FedOPT (Reddi et al. 2020), as well as personalized FL frameworks (e.g., pFedMe, Ditto, Per-FedAvg). These algorithms all depend on the strong assumption that the data at the edge has sufficient labels. To address the label deficiency issue in FL, recent works (Liu et al. 2020; Long et al. 2020; Itahara et al. 2020; Jeong et al. 2020; Liang et al. 2021; Zhao et al. 2020; Zhang et al. 2020a,b) assume that the server or client has a fraction of labeled data and use semi-supervised methods such as consistency loss (Miyato et al. 2018) or pseudo labeling (Lee 2013) to train a global model. A more realistic but challenging setting is fully unsupervised training. Although a recent work in FL (Saeed et al. 2020) attempts to address this challenge through Siamese networks proposed around thirty years ago (Bromley et al. 1993), its design does not tackle data heterogeneity for learning personalized models, and it only trains on small-scale sensor data in IoT devices. Moreover, these existing works in FL have not examined recent progress in the Self-Supervised Learning (SSL) community where methods such as SimCLR (Chen et al. 2020), SwAV(Caron et al. 2021), BYOL (Grill et al. 2020), and SimSiam (Chen and

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



Figure 1: Depiction of the Self-supervised and Personalized Federated Learning (SSFL) framework.

He 2020) have shown tremendous improvement in reducing the amount of labeled data required to achieve state-of-theart performance. As such, it remains still unclear how these SSL methods can be incorporated into FL and how well they would perform, especially when intertwined with the data heterogeneity challenge that does not exist in centralized training.

In this paper, we propose Self-Supervised Federated Learning (SSFL), a unified self-supervised and personalized federated learning framework, and a series of algorithms under this framework to address these challenges. As shown in Figure 1, this framework brings state-of-the-art self-supervised learning algorithms to the realm of FL in order to enable training without using any supervision, while also integrating model personalization to deal with data heterogeneity (Section 3.1). More specifically, under the SSFL framework, we analyze the compatibility of various centralized selfsupervised learning methods in the FL setting and demonstrate that SimSiam networks performs the best with the standard FedAvg algorithm (Section 3.2). Moreover, to address the data heterogeneity at edge devices, we have innovated a series of algorithms that broaden the reach of existing supervised personalization algorithms into the setting of self-supervised learning, including perFedAvg (Fallah, Mokhtari, and Ozdaglar 2020a), Ditto (Li et al. 2021), and local fine-tuning, among others. We further propose a novel personalized federated self-supervised learning algorithm, per-SSFL (Section 3.3), which balances personalization and consensus by carefully regulating the distance between the local and global representations of data (shown as the yellow block in Figure 1).

To provide a comprehensive and comparative analysis of the proposed algorithms, we also develop a distributed training system and evaluation protocol for SSFL. Using this training system, we conduct experiments on a synthetic non-I.I.D. dataset based on CIFAR-10 and a natural non-I.I.D. dataset GLD-23K. Our experimental results demonstrate that all algorithms in our framework work reliably. In FL, the gap of evaluation accuracy between supervised learning and unsupervised learning is small. Personalized SSFL performs better than FedAvg-based SSFL. We also conduct ablation studies to fully understand the SSFL framework, namely the role of batch size, different degrees of non-I.I.D.ness, and performance in more datasets. Finally, our unified API design can serve as a suitable platform and baseline, enabling further developments of more advanced SSFL algorithms.

2 Preliminaries

SSFL builds upon two fundamental areas in machine learning: federated optimization and self-supervised learning. Thus, we first introduce some basics and formulations in these areas.

2.1 Federated Optimization

Federated optimization refers to the distributed optimization paradigm that a network of K devices collaboratively solve a machine learning task. In general, it can be formulated as a distributed optimization problem with the form (McMahan et al. 2017): $\min_{\theta} \sum_{k=1}^{K} \frac{|D_k|}{|D|} \mathcal{L}(\theta, D_k)$. Here, each device k has a local dataset D_k drawn from a local distribution X_k . The combined dataset $D = \bigcup_{k=1}^{K} D_k$ is the union of all local datasets D_k . θ represents the model weight of a client model. ${\mathcal L}$ is the client's local loss function that measures the local empirical risk over the heterogeneous dataset \mathcal{D}^k . Under this formulation, to mitigate the non-I.I.D. issue, researchers have proposed algorithms such as FedProx (Li et al. 2018), FedNova (Wang et al. 2020), and FedOPT (Reddi et al. 2020) for training a global model, as well as personalized FL frameworks such as Ditto (Li et al. 2021), and Per-FedAvg (Fallah, Mokhtari, and Ozdaglar 2020b). All of these algorithms have a strong assumption that data at the edge have sufficient labels, meaning that their analysis and experimental settings are based on a supervised loss function, such as the cross-entropy loss for image classification.

2.2 Self-supervised Learning

Self-supervised learning (SSL) aims to learn meaningful representations of samples without human supervision. Formally, it aims to learn an encoder function $f_{\theta} : \mathcal{X} \mapsto \mathbb{R}^d$ where θ is the parameter of the function, \mathcal{X} is the unlabeled sample space (e.g. image, text), and the output is a *d* dimensional

vector containing enough information for downstream tasks such as image classification and segmentation. The key to SSL's recent success is the inductive bias that ensures a good representation encoder remains consistent under different perturbations of the input (i.e. consistency regularization).

One prominent example among recent advances in modern SSL frameworks is the Siamese architecture (Bromley et al. 1993) and its improved variants SimCLR (Chen et al. 2020), SwAV (Caron et al. 2021), BYOL (Grill et al. 2020), and SimSiam (Chen and He 2020). Here we review the most elegant architecture, SimSiam, and defer the description and comparison of the other three to Appendix A. SimSiam proposes a two-head architecture in which two different views (augmentations) of the same image are encoded by the same network f_{θ} . Subsequently, a predictor Multi Layer Perceptron (MLP) h_{θ} and a *stop-gradient* operation denoted by $\hat{\cdot}$ are applied to both heads. In the SSL context, "stop gradient" means that the optimizer stops at a specific neural network layer during the back propagation and views the parameters in preceding layers as constants. Here, θ is the concatenation of the parameters of the encoder network and the predictor MLP. The algorithm aims to minimize the negative cosine similarity $\mathcal{D}(\cdot, \cdot)$ between two heads. More concretely, the loss is defined as

$$\mathcal{L}_{\rm SS}(\theta, D) = \frac{1}{|D|} \sum_{x \in D} \mathcal{D}(f_{\theta}(\mathcal{T}(x)), \widehat{h_{\theta}(f_{\theta}(\mathcal{T}(x)))})), \quad (1)$$

where \mathcal{T} represents stochastic data augmentation and D is the data set.

3 SSFL: Self-supervised Federated Learning

In this section, we propose SSFL, a unified framework for self-supervised federated learning. Specifically, we introduce the method by which SSFL works for collaborative training of a global model and personalized models, respectively.

3.1 General Formulation

We formulate self-supervised federated learning as the following distributed optimization problem:

$$\min_{\substack{\Theta \\ \{\theta_k\}_{k\in[K]}}} G\left(\mathcal{L}\left(\theta_1,\Theta;X_1\right),\ldots,\mathcal{L}\left(\theta_K,\Theta;X_K\right)\right)$$
(2)

where θ_k is the parameter for the local model $(f_{\theta_k}, h_{\theta_k})$; Θ is the parameter for the global model (f_{Θ}, h_{Θ}) ; $\mathcal{L}(\theta_k, \Theta; X_k)$ is a loss measuring the quality of representations encoded by f_{θ_k} and f_{Θ} on the local distribution X_k ; and $G(\cdot)$ denotes the aggregation function (e.g. sum of client losses weighted by $\frac{|D_k|}{|D|}$). To capture the two key challenges in federated learning (data heterogeneity and label deficiency), we hold two core assumptions in the proposed framework: (1) X_k of all clients are heterogeneous (non-I.I.D.), and (2) there is no label.

To tackle the above problem, we propose a unified training framework for federated self-supervised learning, as described in Algorithm 1. This framework can handle both non-personalized and personalized federated training. In particular, if one enforces the constraint $\theta_k = \Theta$ for all clients $k \in [K]$, the problem reduces to learning a global model. When this constraint is not enforced, θ_k can be different for each client, allowing for model personalization. ClientSSLOPT is the local optimizer at the client side which solves the local sub-problem in a self-supervised manner. ServerOPT takes the update from the client side and generates a new global model for all clients.

Algorithm 1: SSFL: A Unified Framework for Self- supervised Federated Learning				
input : K ,	$T, \Theta^{(0)}, \{\theta_k^{(0)}\}_{k \in [K]}, \text{ ClientSSLOpt}, \text{ ServerOpt}$			
1 for $t = 0, .$	$\ldots, T-1$ do			
2 Server	candomly selects a subset of devices $S^{(t)}$			
3 Server	sends the current global model $\Theta^{(t)}$ to $S^{(t)}$			
4 for dev	ice $k \in S^{(t)}$ in parallel do			
5 Sol	ve local sub-problem of equation 2:			
θ	$\boldsymbol{\theta}_{k}, \boldsymbol{\Theta}_{k}^{(t)} \leftarrow \textbf{ClientSSLOPT}\left(\boldsymbol{\theta}_{k}^{(t)}, \boldsymbol{\Theta}^{(t)}, \nabla \mathcal{L}\left(\boldsymbol{\theta}_{k}, \boldsymbol{\Theta}; X_{k}\right)\right)$			
6 Ser	nd $\Delta_k^{(t)} := \Theta_k^{(t)} - \Theta^{(t)}$ back to server			
7 $\Theta^{(t+1)}$	$\leftarrow \text{ ServerOpt } \left(\Theta^{(t)}, \{\Delta^{(t)}_k\}_{k \in S^{(t)}}\right)$			
return : $\{\theta_k\}_{k\in[K]}, \Theta^{(T)}$				

Next, we will introduce specific forms of ClientSSLOPT and ServerOPT for global training and personalized training.

3.2 Global-SSFL: Collaboration Towards a Global Model without Supervision

To train a global model using SSFL, we design a specific form of ClientSSLOPT using SimSiam. We choose SimSiam over other contemporary self-supervised learning frameworks (e.g., SimSiam, SwAV, BYOL) based on the following analysis as well as experimental results (see Section 5.1).

The simplicity in neural architecture and training method. SimSiam's architecture and training method are relatively simple. For instance, compared with SimCLR, SimSiam has a simpler loss function; compared with SwAV, SimSiam does not require an additional neural component (prototype vectors) and Sinkhorn-Knopp algorithm; compared with BYOL, SimSiam does not need to maintain an additional moving averaging network for an online network. Moreover, the required batch size of SimSiam is the smallest, making it relatively friendly for resource-constrained federated learning. A more comprehensive comparison can be found in Appendix A.

Interpretability of SimSiam leads to simpler local optimization. More importantly, SimSiam is more interpretable from an optimization standpoint which simplifieds the local optimization. In particular, it can be viewed as an implementation of an Expectation-Maximization (EM) like algorithm, meaning that optimizing \mathcal{L}_{SS} in Equation 1 is implicitly optimizing the following objective

$$\min_{\theta,\eta} \mathbb{E}_{\substack{\mathcal{T}\\x\sim X}} \left[\|f_{\theta}(\mathcal{T}(x)) - \eta_x\|_2^2 \right].$$
(3)

Here, f_{θ} is the encoder neural network parameterized by θ . η is an extra set of parameters, whose size is proportional to the number of images, and η_x refers to using the image index of x to access a sub-vector of η . This formulation is w.r.t. both θ and η and can be optimized via an alternating algorithm. At time step t, the η_x^t update takes the form $\eta_x^t \leftarrow \mathbb{E}_{\mathcal{T}}[f_{\theta^t}(\mathcal{T}(x))],$ indicating that η_x^t is assigned the average representation of x over the distribution of augmentation. However, it is impossible to compute this step by going over the entire dataset during training. Thus, SimSiam uses one-step optimization to approximate the EM-like twostep iteration by introducing the predictor h_{θ} to approximate η and learn the expectation (i.e. $h_{\theta}(z) \approx \mathbb{E}_{\mathcal{T}}[f_{\theta}(\mathcal{T}(x))]$) for any image x. After this approximation, the expectation $\mathbb{E}_{\mathcal{T}}[\cdot]$ is ignored because the sampling of \mathcal{T} is implicitly distributed across multiple epochs. Finally, we can obtain the self-supervised loss function in Equation 1, in which negative cosine similarity \mathcal{D} is used in practice (the equivalent L_2 distance is used in Equation 3 for the sake of analysis). Applying equation 1 as ClientSSLOPT simplifies the local optimization for each client in a self-supervised manner.

3.3 Per-SSFL: Learning Personalized Models without Supervision

In this section, we explain how SSFL addresses the data heterogeneity challenge by learning personalized models. Inspired by the interpretation in Section 3.2, we define the following sub-problem for each client $k \in [K]$:

$$\min_{\theta_{k},\eta_{k}} \quad \mathbb{E}_{\substack{x \sim X_{k} \\ x \sim X_{k}}} \left[\|f_{\theta_{k}}(\mathcal{T}(x)) - \eta_{k,x}\|_{2}^{2} + \frac{\lambda}{2} \|\eta_{k,x} - \mathcal{H}_{x}^{*}\|_{2}^{2} \right]$$

s.t. $\Theta^{*}, \mathcal{H}^{*} \in \arg\min_{\Theta,\mathcal{H}} \sum_{i=1}^{n} \frac{|D_{k}|}{|D|} \mathbb{E}_{\substack{x \sim X_{i} \\ x \sim X_{i}}} \left[\|f_{\Theta}(\mathcal{T}(x)) - \mathcal{H}_{x}\|_{2}^{2} \right]$
(4)

Compared to global training, we additionally include Θ , the global model parameter, and \mathcal{H} , the global version of η , and the expected representations which correspond to the personalized parameters θ_k and η_k . In particular, through the term $\|\eta_{k,x} - \mathcal{H}_x^*\|_2^2$, we aim for the expected local representation of any image x to reside within a neighborhood around the expected global representation of x. Therefore, by controlling the radius of the neighborhood, hyperparameter λ helps to balance consensus and personalization.

We see that Equation 4 in the above objective is an optimization problem w.r.t. both θ and η . However, as the above target is intractable in practice, following an analysis similar to Section 3.2, we use the target below as a surrogate:

$$\begin{split} \min_{\theta_{k}} & \mathcal{L}_{\mathrm{SS}}\left(\theta_{k}, D_{k}\right) \\ & + \frac{\lambda}{|D_{k}|} \sum_{x \in D_{k}} \mathcal{D}\left(h_{\theta_{k}} \circ f_{\theta_{k}} \circ \mathcal{T}(x), h_{\Theta^{*}} \circ f_{\Theta^{*}} \circ \mathcal{T}(x)\right) \\ \text{s.t.} & \Theta^{*} \in \arg\min_{\Theta} \mathcal{L}_{\mathrm{SS}}\left(\Theta, D\right) \end{split}$$

In practice, Θ can be optimized independently of θ^k through the FedAvg (McMahan et al. 2017) algorithm. To make the computation more efficient, we also apply the symmetrization trick proposed in (Chen and He 2020). We refer

to this algorithm as Per-SSFL and provide a detailed description in Algorithm 2 (also illustrated in Fig. 1).

Regarding the theoretical analysis. To our knowledge, all self-supervised learning frameworks do not have any theoretical analysis yet, particularly the SimSiam dual neural network architecture. Our formulation and optimization framework are interpretable, they are built based on an EM-like algorithm for SimSiam and minimizing the distance between the private model and the global model's data representation.

Innovating baselines to verify SSFL. Note that we have not found any related works that explore a Siamese-like SSL architecture in an FL setting. As such, to investigate the performance of our proposed algorithm, we further propose several other algorithms that can leverage the SSFL framework. 1. LA-SSFL. We apply FedAvg (McMahan et al. 2017) on the SimSiam loss $\mathcal{L}_{\mathrm{SS}}$ for each client to obtain a global model. We perform one step of SGD on the clients' local data for local adaption; 2. MAML-SSFL. This algorithm is inspired by perFedAvg (Fallah, Mokhtari, and Ozdaglar 2020b) and views the personalization on each devices as the inner loop of MAML (Finn, Abbeel, and Levine 2017). It aims to learn an encoder that can be easily adapted to the clients' local distribution. During inference, we perform one step of SGD on the global model for personalization; 3. BiLevel-SSFL. Inspired by Ditto (Li et al. 2021), we learn personalized encoders on each client by restricting the parameters of all personalized encoders to be close to a global encoder independently learned by weighted aggregation. More details of these algorithms, formulation, and pseudo code are introduced in Appendix B. In Section 5.3, we will show the comparison results for these proposed SSFL algorithmic variants.

4 Training System and Evaluation Pipeline for SSFL

A Distributed Training System to accelerate the algorithmic exploration in SSFL framework. We also contributed to reproducible research via our distributed training system. This is necessary for two reasons: (1) Running a stand-alone simulation (training client by client sequentially) like most existing FL works requires a prohibitively long training time when training a large number of clients. In SSFL, the model size (e.g., ResNet-18 v.s. shallow CNNs used in the original FedAvg paper) and the round number for convergence (e.g., 800 epochs in the centralized SimSiam framework) is relatively larger than in FL literature. By running all clients in parallel on multiple CPUs/GPUs, we can largely accelerate the process. (2) Given that SSFL is a unified and generic learning framework, researchers may develop more advanced ways to improve our work. As such, we believe it is necessary to design unified APIs and system frameworks in line with the algorithmic aspect of SSFL. See Appendix C for more details on our distributed training system.

Evaluation Pipeline. In the training phase, we use a KNN classifier (Wu et al. 2018) as an online indicator to monitor the quality of the representations generated by the SimSiam encoder. For Global-SSFL, we report the KNN test accu-

Algorithm 2: Per-SSFL

input : $K, T, \lambda, \Theta^{(0)}, \{\theta_i^{(0)}\}_{k \in [K]}, s$: number of local iteration, β : learning rate 1 for $t = 0, \ldots, T - 1$ do Server randomly selects a subset of devices $S^{(t)}$ 2 Server sends the current global model $\Theta^{(t)}$ to $S^{(t)}$ 3 for device $k \in S^{(t)}$ in parallel do 4 Sample mini-batch B_k from local dataset D_k , and do s local iterations 5 /* Optimize the global parameter Θ locally $Z_1, Z_2 \leftarrow f_{\Theta^{(t)}}(\mathcal{T}(B_k)), f_{\Theta^{(t)}}(\mathcal{T}(B_k))$ 6 $P_1, P_2 \leftarrow h_{\Theta^{(t)}}(Z_1), h_{\Theta^{(t)}}(Z_2)$ 7 $\Theta_k^{(t)} \leftarrow \Theta^{(t)} - \beta \nabla_{\Theta^{(t)}} \frac{\mathcal{D}(P_1, \widehat{Z_2}) + \mathcal{D}(P_2, \widehat{Z_1})}{2}$, where $\hat{\cdot}$ stands for stop-gradient 8 /* Optimize the local parameter $heta_k$ $z_1, z_2 \leftarrow f_{\theta_k}(\mathcal{T}(B_k)), f_{\theta_k}(\mathcal{T}(B_k))$ 9 $p_1, p_2 \leftarrow h_{\theta_k}(z_1), h_{\theta_k}(z_2)$ 10 $\theta_k \leftarrow \theta_k - \beta \nabla_{\theta_k} \left(\frac{\mathcal{D}(p_1, \widehat{z_2}) + \mathcal{D}(p_2, \widehat{z_1})}{2} + \lambda \frac{\mathcal{D}(p_1, P_1) + \mathcal{D}(p_1, P_2) + \mathcal{D}(p_2, P_1) + \mathcal{D}(p_2, P_2)}{4} \right)$ 11 CLIENTSSLOPT Send $\Delta_k^{(t)} := \Theta_k^{(t)} - \Theta^{(t)}$ back to server 12 $\Theta^{(t+1)} \leftarrow \Theta^{(t)} + \sum_{k \in S^{(t)}} \frac{|D_k|}{|D|} \Delta_k^{(t)}$ 13 SERVEROPT return : $\{\theta_i\}_{i\in[n]}, \Theta^{(T)}$

racy using the global model and the global test data, while in Per-SSFL, we evaluate all clients' local encoders separately with their local test data and report their averaged accuracy. After self-supervised training, to evaluate the performance of the trained encoder, we freeze the encoder and attach a linear classifier to the output of the encoder. For Global-SSFL, we can easily verify the performance of SimSiam encoder by training the attached linear classifier with FedAvg. However, for Per-SSFL, each client learns a personalized SimSiam encoder. As the representations encoded by personalized encoders might reside in different spaces, using a single linear classifier trained by FedAvq to evaluate these representations is unreasonable (see experiments in Section 5.4). As such, we suggest an additional evaluation step to provide a more representative evaluation of Per-SSFL's performance: for each personalization encoder, we use the entire training data to train the linear classifier but evaluate on each client's local test data.

5 Experiments

In this section, we introduce experimental results for SSFL with and without personalization and present a performance analysis on a wide range of aspects including the role of batch size, different degrees of non-IIDness, and understanding the evaluation protocol.

Implementation. We develop the SSFL training system to simplify and unify the algorithmic exploration. Details of the training system design can be found in Appendix C. We deploy the system in a distributed computing environment which has 8 NVIDIA A100-SXM4 GPUs with sufficient memory (40 GB/GPU) to explore different batch sizes (Section 5.4). Our training framework can run multiple parallel training workers in a single GPU, so it supports federated training with a large number of clients. The client number selected per round used in all experiments is 10, which is a reasonable setting suggested by recent literature (Reddi et al.

2020).

Learning Task. Following SimCLR (Chen et al. 2020), SimSiam (Chen and He 2020), BYOL (Grill et al. 2020), and SwAV (Caron et al. 2020) in the centralized setting, we evaluate SSL for the *image classification* task and use representative datasets for federated learning.

Dataset. We run experiments on synthetic non-I.I.D. dataset CIFAR-10 and intrinsically non-I.I.D. dataset Google Landmark-23K (GLD-23K), which are suggested by multiple canonical works in the FL community (Reddi et al. 2020; He et al. 2020; Kairouz et al. 2019). For the non-I.I.D. setting, we distribute the dataset using a Dirichlet distribution (Hsu, Qi, and Brown 2019), which samples $\mathbf{p}_c \sim \text{Dir}(\alpha)$ (we assume a uniform prior distribution) and allocates a $\mathbf{p}_{c,k}$ proportion of the training samples of class c to local client k. We provide a visualization of the data distribution in Appendix E.1.

Model Architecture. For the model architecture, ResNet-18 is used as the backbone of the SimSiam framework, and the predictor is the same as that in the original paper.

Next, we focus on results from the curated CIFAR-10 dataset and defer GLD-23K to Appendix D.

5.1 Comparisons on SimSiam, SimCLR, SwAV, and BYOL

Our first experiment determines which SSL method is ideal for FL settings. We run experiments using FedAvg for these four methods and obtain two findings: (1) SimSiam outperforms SimCLR in terms of accuracy; (2) BYOL and SwAV do not work in FL; we tested a wide range of hyper-parameters, but they still are unable to converge to a reasonable accuracy. These experimental results confirm our analysis in Section 3.2.

5.2 Evaluation on Global-SSFL

The goal of this experiment is to understand the accuracy gap between supervised and self-supervised federated learning



Figure 2: Training and Evaluation using SSFL

in both I.I.D. and non-I.I.D. settings where we aim to train a global model from private data from clients.

Setup and Hyper-parameters. We evaluate Global-SSFL using non-I.I.D. data from CIFAR-10: we set $\alpha = 0.1$ for the Dirichlet distribution. For supervised learning, the test accuracy is evaluated on a global test dataset. For self-supervised training, we follow the evaluation protocol introduced in Section 4. We use SGD with Momentum as the client-side optimizer and a learning rate scheduler across communication rounds. We searched the learning rate on the grid of $\{0.1, 0.3, 0.01, 0.03\}$ and report the best result. The experiment is run three times using the same learning rate with fixed random seeds to ensure reproducibility. The training lasts for 800 rounds, which is sufficient to achieve convergence for all methods. More hyperparameters are in Appendix E.2.

We display the training curves in Figure 2 which demonstrates that SSFL can converge reliably in both I.I.D. and non-I.I.D. settings. For the I.I.D. data, we find that SSFL can achieve the same accuracy as the centralized accuracy report in the SimSiam paper (Chen and He 2020). For the non-I.I.D. data, SSFL achieves a reasonable accuracy compared to the centralized accuracy. The accuracy comparisons in different dimensions (supervised v.s. self-supervised; I.I.D. v.s. non-I.I.D.) are summarized in Table 1.

5.3 Evaluation on Per-SSFL

Based on the results of SSFL with FedAvg, we further add the personalization components for SSFL introduced in Section 3.3 (Per-SSFL).

Setup and Hyper-parameters. For a fair comparison, we evaluate Per-SSFL on non-I.I.D. data from CIFAR-10 and set $\alpha = 0.1$ for the Dirichlet distribution. For Per-SSFL training, we follow the evaluation protocol introduced in Section 4. Similar to SSFL, we use SGD with Momentum as the client-side optimizer and the learning rate scheduler across communication rounds. We search for the learning rate on a grid of $\{0.1, 0.3, 0.01, 0.03\}$ and report the best result. For Per-SSFL and BiLevel-SSFL, we also tune the λ of the regularization term with a search space $\{1, 10, 0.1, 0.01, 0.001\}$. The experiments are run three times with the same learning rate and with fixed random seeds to ensure reproducibility. The training also lasts for 800 communication rounds, which is the same as Global-SSFL. Other hyperparameters can be found in Appendix E.2.

We illustrate our results in Figure 3 and Table 2. To confirm the convergence, we draw loss curves for all methods in Figure 3(b) (note that they have different scaled values due to the difference of their loss functions). Figure 3(b) indicates that Per-SSFL performs best among all methods. MAML-SSFL is also a suggested method since it obtains comparable accuracy. LA-SSFL is a practical method, but it does not perform well in the self-supervised setting. In Figure 3(b), the averaged personalized accuracy of LA-SSFL diverges in the latter phase. Based on BiLevel-SSFL's result, we can conclude such a method is not a strong candidate for personalization, though it shares similar objective functions as Per-SSFL. This indicates that regularization through representations encoded by SimSiam outperforms regularization through weights.

5.4 Performance Analysis

Role of Batch Size FL typically requires a small batch size to enable practical training on resource-constrained edge devices. Therefore, understanding the role of batch size in SSFL is essential to practical deployment and applicability. To investigate this, we use different batch sizes and tune the learning rate to find the optimal accuracy for each one. The results in Figure 4 show that SSFL requires a large batch size (256); otherwise, it reduces the accuracy or diverges during training. Since system efficiency is not the focus of this paper, we use gradient accumulation, which is a simple yet effective method. We fix the batch size at 32 and use accumulation step 8 for all experiments. For an even larger batch size (e.g., 512), the memory cost is significant, though there is no notable gain in accuracy. Therefore, we discontinue the search for batch sizes larger than larger than 256. A more advanced method includes batch-size-one training and knowledge distillation. We defer the discussion to Appendix F.

On Different Degrees of Non-I.I.D.ness We investigate the impact of the degree of data heterogeneity on the SSFL performance. We compare the performance between $\alpha = 0.1$ and $\alpha = 0.5$. These two settings provide a non-negligible gap in the label distribution in each client (see our visualization in Appendix E.1). Figure 5(a) and 5(b) shows the learning curve comparisons. It is clearly observed that the higher degree of data heterogeneity makes it converge more slowly, adversely affecting the accuracy.

	A		
	Supervised	Self-Supervised	Acc. Gap
I.I.D	0.932	0.914	0.018
non-I.I.D	0.8812	0.847	0.0342
Acc. Gap	0.0508	0.06	N/A

Table 1: Evaluation accuracy comparison between supervised FL and SSFL.



Figure 3: Training and Evaluation using SSFL



Figure 4: Results for batch sizes



Figure 5: Evaluation on Different Degress of Non-I.I.D.ness

Method	KNN Indicator	Evaluation
LA-SSFL	0.9217	0.8013
MAML-SSFL	0.9355	0.8235
BiLevel-SSFL	0.9304	0.8137
Per-SSFL	0.9388	0.8310

Table 2: Evaluation Accuracy for Various Per-SSFL Methods.



Figure 6: Understanding the Evaluation Protocol

Understanding the Linear Evaluation of Personalized Encoders As we discussed in 4, in SSFL, we can easily verify the quality of the SimSiam encoder using federated linear evaluation; however, in Per-SSFL, each client learns a personalized SimSiam encoder. Such heterogeneity in diverse encoders makes a fair evaluation difficult. To demonstrate this, we run experiments with naive federated linear evaluation on personalized encoders and surprisingly find that such an evaluation protocol downgrades the performance. As shown in Figure 6, the federated linear evaluation for Per-SSFL performs worse than even LA-SSFL. This may be attributed to the fact that the naive aggregation drags close to the parameter space of all heterogeneous encoders, making the encoder degenerate in terms of personalization.

6 Related Works

Federated Learning (FL) with Personalization. pFedMe (Dinh, Tran, and Nguyen 2020), perFedAvg (Fallah, Mokhtari, and Ozdaglar 2020a), and Ditto (Li et al. 2021) are some representative works in this direction. However, these methods all have a strong assumption that users can provide reliable annotations for their private and sensitive data, which we argue to be very unrealistic and impractical.

Label deficiency in FL. There are a few related works to tackle label deficiency in FL (Liu et al. 2020; Long et al. 2020; Itahara et al. 2020; Jeong et al. 2020; Liang et al. 2021; Zhao et al. 2020; Zhang et al. 2020b). Compared to these works, our proposed SSFL does not use any labels during training. FedMatch (Jeong et al. 2020) and FedCA (Zhang et al. 2020a) requires additional communication costs to synchronize helper models or public labeled dataset. (Saeed et al. 2020) addresses the fully unsupervised challenge on small-scale sensor data in IoT devices. However, compared to our work, it uses the Siamese networks proposed around thirty years ago (Bromley et al. 1993), lacking consideration on

the advance in the past two years (i.e., SimCLR (Chen et al. 2020), SwAV(Caron et al. 2021), BYOL (Grill et al. 2020), and SimSiam (Chen and He 2020)). Moreover, these works does not have any design for learning personalized models.

7 Conclusion

We propose Self-supervised Federated Learning (SSFL) framework and a series of algorithms under this framework towards addressing two challenges: data heterogeneity and label deficiency. SSFL can work for both global model training and personalized model training. We conduct experiments on a synthetic non-I.I.D. dataset based on CIFAR-10 and the intrinsically non-I.I.D. GLD-23K dataset. Our experimental results demonstrate that SSFL can work reliably and achieves reasonable evaluation accuracy that is suitable for use in various applications.

References

Berthelot, D.; Carlini, N.; Cubuk, E. D.; Kurakin, A.; Sohn, K.; Zhang, H.; and Raffel, C. 2020. ReMixMatch: Semi-Supervised Learning with Distribution Matching and Augmentation Anchoring. In *International Conference on Learning Representations*.

Berthelot, D.; Carlini, N.; Goodfellow, I.; Papernot, N.; Oliver, A.; and Raffel, C. A. 2019. MixMatch: A Holistic Approach to Semi-Supervised Learning. In *Neural Information Processing Systems*, 5049–5059.

Brendan McMahan, H.; Moore, E.; Ramage, D.; Hampson, S.; and Agüera y Arcas, B. 2016. Communication-Efficient Learning of Deep Networks from Decentralized Data. *arXiv e-prints*, arXiv:1602.05629.

Bromley, J.; Bentz, J. W.; Bottou, L.; Guyon, I.; LeCun, Y.; Moore, C.; Säckinger, E.; and Shah, R. 1993. Signature verification using a "siamese" time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(04): 669–688.

Cai, H.; Gan, C.; Zhu, L.; and Han, S. 2020. TinyTL: Reduce Memory, Not Parameters for Efficient On-Device Learning. *Advances in Neural Information Processing Systems*, 33.

Caron, M.; Misra, I.; Mairal, J.; Goyal, P.; Bojanowski, P.; and Joulin, A. 2020. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*.

Caron, M.; Misra, I.; Mairal, J.; Goyal, P.; Bojanowski, P.; and Joulin, A. 2021. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. arXiv:2006.09882.

Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, 1597–1607. PMLR.

Chen, X.; and He, K. 2020. Exploring Simple Siamese Representation Learning.

Cuturi, M. 2013. Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances. arXiv:1306.0895.

Dinh, C. T.; Tran, N. H.; and Nguyen, T. D. 2020. Personalized Federated Learning with Moreau Envelopes. Fallah, A.; Mokhtari, A.; and Ozdaglar, A. 2020a. Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*.

Fallah, A.; Mokhtari, A.; and Ozdaglar, A. 2020b. *Personalized Federated Learning: A Meta-Learning Approach.*

Finn, C.; Abbeel, P.; and Levine, S. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks.

Grill, J.-B.; Strub, F.; Altché, F.; Tallec, C.; Richemond, P. H.; Buchatskaya, E.; Doersch, C.; Pires, B. A.; Guo, Z. D.; Azar, M. G.; et al. 2020. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*.

He, C.; Li, S.; So, J.; Zhang, M.; Wang, H.; Wang, X.; Vepakomma, P.; Singh, A.; Qiu, H.; Shen, L.; Zhao, P.; Kang, Y.; Liu, Y.; Raskar, R.; Yang, Q.; Annavaram, M.; and Avestimehr, S. 2020. FedML: A Research Library and Benchmark for Federated Machine Learning. *arXiv preprint arXiv:2007.13518*.

Howard, A. G.; Zhu, M.; Chen, B.; Kalenichenko, D.; Wang, W.; Weyand, T.; Andreetto, M.; and Adam, H. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*.

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

Hu, Z.; Yang, Z.; Hu, X.; and Nevatia, R. 2021. SimPLE: Similar Pseudo Label Exploitation for Semi-Supervised Classification.

Itahara, S.; Nishio, T.; Koda, Y.; Morikura, M.; and Yamamoto, K. 2020. Distillation-Based Semi-Supervised Federated Learning for Communication-Efficient Collaborative Training with Non-IID Private Data. *arXiv preprint arXiv:2008.06180*.

Jeong, W.; Yoon, J.; Yang, E.; and Hwang, S. J. 2020. Federated Semi-Supervised Learning with Inter-Client Consistency. *arXiv preprint arXiv:2006.12097*.

Kairouz, P.; McMahan, H. B.; Avent, B.; Bellet, A.; Bennis, M.; Bhagoji, A.; Bonawitz, K.; Charles, Z. B.; Cormode, G.; Cummings, R.; D'Oliveira, R. G. L.; Rouayheb, S.; Evans, D.; Gardner, J.; Garrett, Z.; Gascón, A.; Ghazi, B.; Gibbons, P. B.; Gruteser, M.; Harchaoui, Z.; He, C.; He, L.; Huo, Z.; Hutchinson, B.; Hsu, J.; Jaggi, M.; Javidi, T.; Joshi, G.; Khodak, M.; Konecný, J.; Korolova, A.; Koushanfar, F.; Koyejo, O.; Lepoint, T.; Liu, Y.; Mittal, P.; Mohri, M.; Nock, R.; Özgür, A.; Pagh, R.; Raykova, M.; Qi, H.; Ramage, D.; Raskar, R.; Song, D.; Song, W.; Stich, S. U.; Sun, Z.; Suresh, A. T.; Tramèr, F.; Vepakomma, P.; Wang, J.; Xiong, L.; Xu, Z.; Yang, Q.; Yu, F.; Yu, H.; and Zhao, S. 2021. Advances and Open Problems in Federated Learning. *Found. Trends Mach. Learn.*, 14: 1–210.

Kairouz, P.; McMahan, H. B.; Avent, B.; Bellet, A.; Bennis, M.; Bhagoji, A. N.; Bonawitz, K.; Charles, Z.; Cormode, G.; Cummings, R.; et al. 2019. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*.

Laine, S.; and Aila, T. 2016. Temporal Ensembling for Semi-Supervised Learning.

Lee, D.-H. 2013. Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method for Deep Neural Networks.

Li, T.; Hu, S.; Beirami, A.; and Smith, V. 2021. Ditto: Fair and Robust Federated Learning Through Personalization.

Li, T.; Sahu, A. K.; Zaheer, M.; Sanjabi, M.; Talwalkar, A.; and Smith, V. 2018. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*.

Liang, X.; Liu, Y.; Luo, J.; He, Y.; Chen, T.; and Yang, Q. 2021. Self-supervised Cross-silo Federated Neural Architecture Search. *arXiv preprint arXiv:2101.11896*.

Liu, Y.; Yuan, X.; Zhao, R.; Zheng, Y.; and Zheng, Y. 2020. RC-SSFL: Towards Robust and Communication-efficient Semi-supervised Federated Learning System. *arXiv preprint arXiv:2012.04432*.

Long, Z.; Che, L.; Wang, Y.; Ye, M.; Luo, J.; Wu, J.; Xiao, H.; and Ma, F. 2020. FedSemi: An Adaptive Federated Semi-Supervised Learning Framework. *arXiv preprint arXiv:2012.03292*.

McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, 1273–1282.

Miyato, T.; Maeda, S.-i.; Koyama, M.; and Ishii, S. 2018. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8): 1979–1993.

Miyato, T.; Maeda, S.-I.; Koyama, M.; and Ishii, S. 2019. Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. 41(8): 1979– 1993.

Reddi, S.; Charles, Z.; Zaheer, M.; Garrett, Z.; Rush, K.; Konečnỳ, J.; Kumar, S.; and McMahan, H. B. 2020. Adaptive Federated Optimization. *arXiv preprint arXiv:2003.00295*.

Saeed, A.; Salim, F. D.; Ozcelebi, T.; and Lukkien, J. 2020. Federated Self-Supervised Learning of Multisensor Representations for Embedded Intelligence. *IEEE Internet of Things Journal*, 8(2): 1030–1040.

Sajjadi, M.; Javanmardi, M.; and Tasdizen, T. 2016. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Neural Information Processing Systems*, 1171–1179.

Salimans, T.; and Kingma, D. P. 2016. Weight Normalization: A Simple Reparameterization to Accelerate Training of Deep Neural Networks. *CoRR*, abs/1602.07868.

Tan, M.; and Le, Q. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 6105–6114. PMLR.

Wang, J.; Charles, Z. B.; Xu, Z.; Joshi, G.; McMahan, H. B.; Arcas, B. A. Y.; Al-Shedivat, M.; Andrew, G.; Avestimehr, S.; Daly, K.; Data, D.; Diggavi, S.; Eichner, H.; Gadhikar, A.; Garrett, Z.; Girgis, A. M.; Hanzely, F.; Hard, A.; He, C.; Horvath, S.; Huo, Z.; Ingerman, A.; Jaggi, M.; Javidi, T.; Kairouz, P.; Kale, S.; Karimireddy, S. P. R.; Konecný, J.; Koyejo, S.; Li, T.; Liu, L.; Mohri, M.; Qi, H.; Reddi, S. J.; Richtárik, P.; Singhal, K.; Smith, V.; Soltanolkotabi, M.; Song, W.; Suresh, A. T.; Stich, S. U.; Talwalkar, A. S.; Wang, H.; Woodworth, B. E.; Wu, S.; Yu, F. X.; Yuan, H.; Zaheer, M.; Zhang, M.; Zhang, T.; Zheng, C.; Zhu, C.; and Zhu, W. 2021. A Field Guide to Federated Optimization. *ArXiv*, abs/2107.06917.

Wang, J.; Liu, Q.; Liang, H.; Joshi, G.; and Poor, H. V. 2020. Tackling the objective inconsistency problem in heterogeneous federated optimization. *arXiv preprint arXiv:2007.07481*.

Wu, Z.; Xiong, Y.; Yu, S.; and Lin, D. 2018. Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination. arXiv:1805.01978.

Zhang, F.; Kuang, K.; You, Z.; Shen, T.; Xiao, J.; Zhang, Y.; Wu, C.; Zhuang, Y.; and Li, X. 2020a. Federated Unsupervised Representation Learning. *arXiv preprint arXiv:2010.08982*.

Zhang, Z.; Yang, Y.; Yao, Z.; Yan, Y.; Gonzalez, J. E.; and Mahoney, M. W. 2020b. Improving Semi-supervised Federated Learning by Reducing the Gradient Diversity of Models. *arXiv preprint arXiv:2008.11364*.

Zhao, Y.; Liu, H.; Li, H.; Barnaghi, P.; and Haddadi, H. 2020. Semi-supervised Federated Learning for Activity Recognition. *arXiv preprint arXiv:2011.00851*.

Appendix

A Comparison of Self-supervised Learning Frameworks

We compare state-of-the-art self-supervised learning frameworks (SimCLR, SwAV, BYOL) with SimSiam (Chen and He 2020) in light of federated learning.

We choose SimSiam (Chen and He 2020) because it requires a much smaller batch to perform normally. In the centralized setting, for each method to reach an accuracy level similar to that of SimSiam, a much larger batch size is necessary. Table 3 adopted from (Chen and He 2020) provides a brief comparison between all listed self-supervised learning frameworks.

Another reason we prefer SimSiam (Chen and He 2020) as the basic framework to build SSFL is that the design of SimSiam simplifies all other baselines and also obtains a relatively higher accuracy. Figure 7 abstracts these methods. The "encoder" contains all layers that can be shared between both branches (e.g., backbone, projection MLP (Chen et al. 2020), prototypes (Caron et al. 2021)). The components in red are those missing in SimSiam.



Figure 7: (Chen and He 2020) Comparison on Siamese architectures. The encoder includes all layers that can be shared between both branches. The dashed lines indicate the gradient propagation flow. In BYOL, SwAV, and SimSiam, the lack of a dashed line implies stop-gradient, and their symmetrization is not illustrated for simplicity. The components in red are those missing in SimSiam.

SimCLR (Chen et al. 2020). SimCLR relies on negative samples ("dissimilarity") to prevent collapsing. SimSiam can be thought of as "SimCLR without negatives". In every mini-batch, for any image, one augmented view of the same image is considered to be its positive sample, and the remaining augmented views of different images are considered to be its negative samples. A contrastive loss term is calculated to push positive samples together and negative samples away.

SwAV(Caron et al. 2021). SimSiam is conceptually analogous to "SwAV without online clustering". SimSiam encourages the features of the two augmented views of the same image to be similar, while SwAV encourages features of the two augmented views of the same image to belong to

the same cluster. An additional Sinkhorn-Knopp (SK) transform (Cuturi 2013) is required for online clustering of SwAV. The authors of SimSiam (Chen and He 2020) build up the connection between SimSiam and SwAV by recasting a few components in SwAV. (i) The shared prototype layer in SwAV can be absorbed into the Siamese encoder. (ii) The prototypes were weight-normalized outside of gradient propagation in (Caron et al. 2021); the authors of SimSiam instead implement by full gradient computation (Salimans and Kingma 2016). (iii) The similarity function in SwAV is cross-entropy. With these abstractions, a highly simplified SwAV illustration is shown in Figure 7.

BYOL (Grill et al. 2020). SimSiam can be thought of as "BYOL without the momentum encoder", subject to many implementation differences. Briefly, in BYOL, one head of the Siamese architecture used in SimSiam is replaced by the exponential moving average of the encoder. As the momentum encoder has an identical architecture to that of the encoder, the introduction of an additional momentum encoder doubles the memory cost of the model.

SSL's recent success is the inductive bias that ensures a good representation encoder remains consistent under different perturbations of the input (i.e. consistency regularization). The perturbations can be either domain-specific data augmentation (e.g. random flipping in the image domain) (Berthelot et al. 2019; Laine and Aila 2016; Sajjadi, Javanmardi, and Tasdizen 2016; Berthelot et al. 2020; Hu et al. 2021), drop out (Sajjadi, Javanmardi, and Tasdizen 2016), random max pooling (Sajjadi, Javanmardi, and Tasdizen 2016), or an adversarial transformation (Miyato et al. 2019). With this idea, a consistency loss \mathcal{L} is defined to measure the quality of the representations without any annotations.

B Formulation and Pseudo Code for Algorithms Under SSFL Framework

Inspired by recent advances in personalized FL and selfsupervised learning, we innovate several representative algorithms under SSFL framework. For each algorithm, we present its mathematical formulation and its pseudo code.

B.1 Per-SSFL

For Per-SSFL, as the formulation and algorithm have already been presented in Equation 4 and Algorithm 2, we provide a PyTorch style pseudo code in Algorithm 3 for additional clarity.

B.2 Personalized SSFL with Local Adaptation (FedAvg-LA)

FedAvg-LA apply FedAvg (Brendan McMahan et al. 2016) on the SimSiam loss \mathcal{L}_{SS} for each client to obtain a global model. We perform one step of SGD on the clients' local data for local adaption. The objective is defined in Equation 5, and the algorithm is provided in Algorithm 4.

$$\min_{\Theta,\mathcal{H}} \sum_{i=1}^{n} \frac{|D_k|}{|D|} \mathbb{E}_{\substack{x \sim \mathcal{X}_i}} \left[\|f_{\Theta}(\mathcal{T}(x)) - \mathcal{H}_x\|_2^2 \right]$$
(5)

method	batch size	negative pairs	momentum encoder	100 ep	200 ep	400 ep	800 ep
SimCLR (repro.+)	4096	\checkmark		66.5	68.3	69.8	70.4
BYOL (repro.)	4096		\checkmark	66.5	70.6	73.2	74.3
SwAV (repro.+)	4096			66.5	69.1	70.7	71.8
SimSiam	256			68.1	70.0	70.8	71.3

Table 3: (Chen and He 2020) Comparisons on ImageNet linear classification. All are based on ResNet-50 pre-trained with two 224×224 views in a centralized setting. Evaluation is on a single crop. "repro." denotes reproduction conducted by authors of SimSiam (Chen and He 2020), and "+" denotes *improved* reproduction v.s. original papers.

B.3 Personalized SSFL with MAML-SSFL

MAML-SSFL is inspired by perFedAvg (Fallah, Mokhtari, and Ozdaglar 2020b) and views the personalization on each devices as the inner loop of MAML (Finn, Abbeel, and Levine 2017). It aims to learn an encoder that can be easily adapted to the clients' local distribution. During inference, we perform one step of SGD on the global model for personalization. The objective is defined in Equation 6, and the algorithm is provided in Algorithm 5.

$$\min_{\Theta,\mathcal{H}} \sum_{i=1}^{n} \frac{|D_k|}{|D|} \mathbb{E}_{x \sim X_i} \left[\|f_{\Theta'}(\mathcal{T}(x)) - \mathcal{H}_x\|_2^2 \right]$$

s.t.
$$\Theta' = \Theta - \nabla_{\Theta} \sum_{i=1}^{n} \frac{|D_k|}{|D|} \mathbb{E}_{x \sim X_i} \left[\|f_{\Theta}(\mathcal{T}(x)) - \mathcal{H}_x\|_2^2 \right]$$

(6)

B.4 Personalized SSFL with BiLevel-SSFL

Inspired by Ditto (Li et al. 2021), BiLevel-SSFL learns personalized encoders on each client by restricting the parameters of all personalized encoders to be close to a global encoder independently learned by weighted aggregation. The objective is defined in Equation 7, and the algorithm is provided in Algorithm 6.

$$\min_{\theta_{k},\eta_{k}} \quad \mathbb{E}_{\substack{\mathcal{T}\\x\sim X_{k}}} \left[\|f_{\theta_{k}}(\mathcal{T}(x)) - \eta_{k,x}\|_{2}^{2} + \frac{\lambda}{2} \|\theta_{k} - \Theta_{x}^{*}\|_{2}^{2} \right]$$

s.t.
$$\Theta^{*}, \mathcal{H}^{*} \in \arg\min_{\Theta,\mathcal{H}} \sum_{i=1}^{n} \frac{|D_{k}|}{|D|} \mathbb{E}_{\substack{\mathcal{T}\\x\sim X_{i}}} \left[\|f_{\Theta}(\mathcal{T}(x)) - \mathcal{H}_{x}\|_{2}^{2} \right]$$

(7)

C Distributed Training System for SSFL

We develop a distributed training system for our SSFL framework which contains three layers. In the infrastructure layer, communication backends such as MPI are supported to facilitate the distributed computing. We abstract the communication as ComManager to simplify the message passing between the client and the server. Trainer reuses APIs from PyTorch to handle the model optimizations such as forward propagation, loss function, and back propagation. In the algorithm layer, Client Manager and Server Manager are the entry points of the client and the server, respectively. The client managers incorporates various SSFL trainers, including Per-SSFL, MAML-SSFL, BiLevel-SSFL, and LA-SSFL. The server handles the model aggregation using Aggregator. We design simplified APIs for all of these modules. With the abstraction of the infrastructure and algorithm layers, developers can begin FL training by developing a workflow script that integrates all modules (as the "SSFL workflow" block shown in the figure). Overall, we found that this distributed training system accelerates our research by supporting parallel training, larger batch sizes, and easyto-customize APIs, which cannot be achieved by a simple single-process simulation.

D Experimental Results on GLD-23K Dataset

We also evaluate the performance of SSFL on GLD-23K dataset. We use 30% of the original local training dataset as the local test dataset and filter out those clients that have a number of samples less than 100. Due to the natural non-I.I.D.ness of GLD-23K dataset, we only evaluate the Per-SSFL framework. The results are summarized in Table 4. *Note: we plan to further explore more datasets and run more experiments; thus we may report more results during the rebuttal phase.*

E Extra Experimental Results and Details

E.1 Visualization of Non-I.I.D. dataset

E.2 Hyper-parameters

All experiments set the local epoch number as 1, round number as 800, batch size as 256 (batch size 32 with 8 gradient accumulation steps).

F Discussion

To overcome the large batch size requirement in SSFL and practical FL edge training, one direction is to use efficient DNN models such as EfficientNet (Tan and Le 2019) and MobileNet (Howard et al. 2017) as the backbone of SimSiam. However, we tested its performance under our framework and found that the performance downgrades to a level of accuracy that is not useful (less than 60%). A recent work in centralized self-supervised learning mitigates these models' accuracy gap by knowledge distillation, which works in a centralized setting but is still not friendly to FL since KD requires additional resources for the teacher model. In practice, we can also explore batch size 1 training (Cai et al. 2020) at the edge, which dramatically reduces the memory cost with additional training time.

Algorithm 3: Per-SSFL PyTorch Style Pseudo Code

```
1 # F: global encoder
2 # H: global predictor
3 # f: local encoder
4 # h: local predictor
5
6 for x in loader: # load a mini-batch x with n samples
7
      x1, x2 = aug(x), aug(x) # random augmentation
      Z1, Z2 = F(x1), F(x2) # global projections, n-by-d
P1, P2 = H(Z1), H(Z2) # global predictions, n-by-d
8
9
      L = D(P1, Z2) / 2 + D(P2, Z1) / 2 \# global loss
13
      L.backward() # back-propagate
      update(F, H) # SGD update global model
14
      z1, z2 = f(x1), f(x2) \# local projections, n-by-d
16
      p1, p2 = h(z1), h(z2) \# local predictions, n-by-d
17
18
      l = D(p1, z2) / 2 + D(p2, z1) / 2 \# local loss
19
20
      # distance between local and global representations
      l = l + \lambda \star (D(p1, P1) + D(p1, P2) + D(p2, P1) + D(p2, P2)) / 4
22
23
      l.backward() # back-propagate
24
      update(f, h)
                     # SGD update local model
25
26
27 def D(p, z): # negative cosine similarity
      z = z.detach() # stop gradient
28
29
      p = normalize(p, dim=1) # l2-normalize
30
      z = normalize(z, dim=1) # l2-normalize
31
      return -(p * z).sum(dim=1).mean()
32
```

Algorithm 4: FedAvg-LA

input : $K, T, \lambda, \Theta^{(0)}, \{\theta_i^{(0)}\}_{k \in [K]}, s$: number of local iteration, β : learning rate 1 for $t = 0, \dots, T - 1$ do Server randomly selects a subset of devices $S^{(t)}$ 2 Server sends the current global model $\Theta^{(t)}$ to $S^{(t)}$ 3 for device $k \in S^{(t)}$ in parallel do 4 Sample mini-batch B_k from local dataset D_k , and do s local iterations 5 /* Optimize the global parameter Θ locally */ $Z_1, Z_2 \leftarrow f_{\Theta^{(t)}}(\mathcal{T}(B_k)), f_{\Theta^{(t)}}(\mathcal{T}(B_k))$ 6 $P_1, P_2 \leftarrow h_{\Theta^{(t)}}(Z_1), h_{\Theta^{(t)}}(Z_2)$ 7 $\Theta_k^{(t)} \leftarrow \Theta^{(t)} - \beta \nabla_{\Theta^{(t)}} \underbrace{\mathcal{D}(P_1, \widehat{Z_2}) + \mathcal{D}(P_2, \widehat{Z_1})}_2$, where $\widehat{\cdot}$ stands for stop-gradient CLIENTSSLOPT 8 Send $\Delta_k^{(t)} := \Theta_k^{(t)} - \Theta^{(t)}$ back to server 9 $\boldsymbol{\Theta}^{(t+1)} \leftarrow \boldsymbol{\Theta}^{(t)} + \sum_{k \in S^{(t)}} \frac{|D_k|}{|D|} \boldsymbol{\Delta}_k^{(t)}$ 10 SERVEROPT return: $\{\theta_i\}_{i\in[n]}, \Theta^{(T)}$

Algorithm 5: MAML-SSFL

input : $K, T, \lambda, \Theta^{(0)}, \{\theta_i^{(0)}\}_{k \in [K]}, s$: number of local iteration, β : learning rate, M 1 for $t = 0, \dots, T - 1$ do Server randomly selects a subset of devices $S^{(t)}$ 2 Server sends the current global model $\Theta^{(t)}$ to $S^{(t)}$ 3 for device $k \in S^{(t)}$ in parallel do 4 Sample mini-batch B_k, B'_k from local dataset D_k , and do s local iterations 5 /* Inner loop update */ $\Theta_{h}^{\prime(t)} \leftarrow \Theta^{(t)}$ 6 for m = 0, ..., M - 1 do 7 $Z'_1, Z'_2 \leftarrow f_{\Theta'^{(t)}}(\mathcal{T}(B'_k)), f_{\Theta'^{(t)}}(\mathcal{T}(B'_k))$ 8 $P'_1, P'_2 \leftarrow h_{\Theta'^{(t)}}(Z'_1), h_{\Theta'^{(t)}}(Z'_2)$ 9 $\Theta_k^{\prime(t)} \leftarrow \Theta_k^{\prime(t)} - \beta \nabla_{\Theta_k^{\prime(t)}} \frac{\mathcal{D}(P_1^{\prime}, \widehat{Z_2^{\prime}}) + \mathcal{D}(P_2^{\prime}, \widehat{Z_1^{\prime}})}{2}$, where $\widehat{\cdot}$ stands for stop-gradient 10 /* Outer loop update */ $\begin{array}{l} Z_1, Z_2 \leftarrow f_{\Theta'^{(t)}}(\mathcal{T}(B_k)), f_{\Theta'^{(t)}}(\mathcal{T}(B_k)) \\ P_1, P_2 \leftarrow h_{\Theta'^{(t)}}(Z_1), h_{\Theta'^{(t)}}(Z_2) \end{array}$ 11 12 $\begin{array}{l} \Theta_{k}^{(t)} \leftarrow \Theta^{(t)} - \beta \nabla_{\Theta^{(t)}} (\underline{\mathcal{D}_{2}}) \\ \beta \nabla_{\Theta^{(t)}} \xrightarrow{\mathcal{D}(P_{1}, \widehat{\mathcal{Z}_{2}}) + \mathcal{D}(P_{2}, \widehat{\mathcal{Z}_{1}})}{2} \\ \end{array} \\ \begin{array}{l} \text{Send } \Delta_{k}^{(t)} := \Theta_{k}^{(t)} - \Theta^{(t)} \text{ back to server} \end{array}$ 13 CLIENTSSLOPT 14 $\boldsymbol{\Theta}^{(t+1)} \leftarrow \boldsymbol{\Theta}^{(t)} + \sum_{k \in S^{(t)}} \frac{|D_k|}{|D|} \boldsymbol{\Delta}_k^{(t)}$ 15 SERVEROPT return: $\{\theta_i\}_{i\in[n]}, \Theta^{(T)}$

Algorithm 6: BiLevel-SSFL

input : $K, T, \lambda, \Theta^{(0)}, \{\theta_i^{(0)}\}_{k \in [K]}, s$: number of local iteration, β : learning rate 1 for $t = 0, \dots, T - 1$ do Server randomly selects a subset of devices $S^{(t)}$ 2 Server sends the current global model $\Theta^{(t)}$ to $S^{(t)}$ 3 for device $k \in S^{(t)}$ in parallel do 4 Sample mini-batch B_k from local dataset D_k , and do s local iterations 5 /* Optimize the global parameter Θ locally */ $\begin{array}{l} Z_1, Z_2 \leftarrow f_{\Theta^{(t)}}(\mathcal{T}(B_k)), f_{\Theta^{(t)}}(\hat{\mathcal{T}}(B_k)) \\ P_1, P_2 \leftarrow h_{\Theta^{(t)}}(Z_1), h_{\Theta^{(t)}}(Z_2) \end{array}$ 6 7 $\Theta_k^{(t)} \leftarrow \Theta^{(t)} - \beta \nabla_{\Theta^{(t)}} \frac{\mathcal{D}(P_1, \widehat{Z_2}) + \mathcal{D}(P_2, \widehat{Z_1})}{2}$, where $\widehat{\cdot}$ stands for stop-gradient 8 /* Optimize the local parameter $heta_k$ */ $z_1, z_2 \leftarrow f_{\theta_k}(\mathcal{T}(B_k)), f_{\theta_k}(\mathcal{T}(B_k)) \\ p_1, p_2 \leftarrow h_{\theta_k}(z_1), h_{\theta_k}(z_2)$ 9 10 $\theta_k \leftarrow \theta_k - \beta \nabla_{\theta_k} \left(\frac{\mathcal{D}(p1,\widehat{z_2}) + \mathcal{D}(p2,\widehat{z_1})}{2} + \lambda \left\| \Theta^{(t)} - \theta_k \right\|_2^2 \right)$ 11 CLIENTSSLOPT Send $\Delta_k^{(t)} := \Theta_k^{(t)} - \Theta^{(t)}$ back to server 12 $\boldsymbol{\Theta}^{(t+1)} \leftarrow \boldsymbol{\Theta}^{(t)} + \sum_{k \in S^{(t)}} \frac{|D_k|}{|D|} \boldsymbol{\Delta}_k^{(t)}$ 13 SERVEROPT return : $\{\theta_i\}_{i\in[n]}, \Theta^{(T)}$



Figure 8: Distributed Training System for SSFL framework





Figure 9: Visualization for non-I.I.D. synthesized using CIFAR-10



(a) Sample Number Distribution (X-axis: Client Index; Y-axis: Number of Training Samples)



(b) Sample Number Distribution (X-axis: Number of Training Samples; Y-axis: Number of Clients)

Figure 10: Visualization for non-I.I.D. on GLD-23K

Method	KNN Indicator	Evaluation
LA-SSFL	0.6011	0.4112
MAML-SSFL	0.6237	0.4365
BiLevel-SSFL	0.6195	0.4233
Per-SSFL	0.6371	0.4467

Table 4: Evaluation Accuracy for Various Per-SSFL Methods.

*Note: the accuracy on supervised federated training using FedAvg is around 47%

Table 5: Hyper-parameters for Section 5.2

Method	Learning Rate	Local Optimizer
SSFL (I.I.D)	0.1	SGD with Momentum (0.9)
SSFL (non-I.I.D)	0.1	SGD with Momentum (0.9)

Table 6: Hyper-parameters for Section 5.4

Method	ethod Learning Rate		Local Optimizer		
Per-SSFL ($\alpha = 0.1$)	0.03	0.1	SGD with Momemtum (0.9)		
Per-SSFL ($\alpha = 0.5$)	0.03	0.1	SGD with Momemtum (0.9)		

Table 7: Hyper-parameters for experimental results in Section 5.3

Method	Learning Rate	λ	Local Optimizer
LA-SSFL	0.1	1	SGD with Momentum (0.9)
MAML-SSFL	0.03	1	SGD with Momentum (0.9)
BiLevel-SSFL	0.1	1	SGD with Momentum (0.9)
Per-SSFL	0.03	0.1	SGD with Momentum (0.9)