

Contrastive encoder pre-training–based clustered federated learning for heterogeneous data

Ye Lin Tun^a, Minh N.H. Nguyen^b, Chu Myaet Thwal^a, Jinwoo Choi^a,
Choong Seon Hong^{a,*}

^a*Department of Computer Science and Engineering, Kyung Hee University, Yongin-si, Gyeonggi-do 17104, South Korea*

^b*Vietnam - Korea University of Information and Communication Technology, Danang, Vietnam*

Abstract

Federated learning (FL) is a promising approach that enables distributed clients to collaboratively train a global model while preserving their data privacy. However, FL often suffers from data heterogeneity problems, which can significantly affect its performance. To address this, clustered federated learning (CFL) has been proposed to construct personalized models for different client clusters. One effective client clustering strategy is to allow clients to choose their own local models from a model pool based on their performance. However, without pre-trained model parameters, such a strategy is prone to clustering failure, in which all clients choose the same model. Unfortunately, collecting a large amount of labeled data for pre-training can be costly and impractical in distributed environments. To overcome this challenge, we leverage self-supervised contrastive learning to exploit unlabeled data for the pre-training of FL systems. Together, self-supervised pre-training and client clustering can be crucial components for tackling the data heterogeneity issues of FL. Leveraging these two crucial strategies, we propose contrastive pre-training–based clustered federated learning (CP-CFL) to improve the model convergence and overall performance of FL systems. In this work, we demonstrate the effectiveness of CP-CFL through extensive experiments in heterogeneous FL settings, and present various interesting observations.

Keywords: Federated learning, Client clustering, Contrastive learning, Pre-training, Data heterogeneity

1. Introduction

Training a neural network for intelligent applications demands a diverse and rich source of real-world data generated by end-user devices. These devices can range from personal

*Corresponding author

Email addresses: yelintun@khu.ac.kr (Ye Lin Tun), nhnminh@vku.udn.vn (Minh N.H. Nguyen), chumyaet@khu.ac.kr (Chu Myaet Thwal), jinwoochoi@khu.ac.kr (Jinwoo Choi), cshong@khu.ac.kr (Choong Seon Hong)

smartphones to institutional data silos, where data privacy and security concerns are prominent (Voigt and Bussche, 2017). Conventionally, neural networks are trained on a central server with use of a collected dataset. However, growing data privacy and security concerns may hinder data collection by service providers (Voigt and Bussche, 2017). Moreover, transmitting large amounts of data from edge devices onto a central server can lead to high communication costs (Dinh et al., 2020; Nguyen et al., 2020).

To address these limitations, centralized learning systems are gradually transitioning towards federated systems, which offer better data privacy guarantees. Edge devices in a federated learning (FL) system (Konečný et al., 2016a,b; McMahan et al., 2017) can collaboratively train a model without sharing sensitive data. An FL system consists of a central server, coordinating the training, and a set of distributed client devices. The training data reside on the client devices, and local models are directly trained by the clients. The trained parameters are then transmitted to the server, while maintaining the confidentiality of each client’s data. The server aggregates the received local parameters into a global model, which inherits the predictive capabilities of the local models. As such, FL is an appealing option for services requiring training on privacy-sensitive data, such as medical data in the healthcare sector (Kaissis et al., 2020; Thwal et al., 2021; Yan et al., 2021) or personal data on a smartphone (Hard et al., 2018; Ramaswamy et al., 2019).

Data heterogeneity poses a significant challenge in FL, as private data generated by each client device differ depending on the nature of the device (Li et al., 2020; McMahan et al., 2017; Sattler et al., 2020b; Yu et al., 2020; Zhao et al., 2018). Non-independent and non-identically distributed (Non-IID) data generated by a large number of clients may cause the divergence between local parameter updates, thereby slowing the convergence speed and hurting the overall performance of the aggregated global model (Yu et al., 2020). To address this challenge, clients can be partitioned into different clusters to train cluster-level models. This kind of approach is often referred to as “clustered federated learning” (CFL) (Ghosh et al., 2020; Sattler et al., 2020a, 2021b; Shlezinger et al., 2020), where multiple cluster-level models are trained, as opposed to a single global model in conventional FL. Depending on the clustering criteria, clients within the same cluster may share similar characteristics, such as data distributions or the availability of training resources. Cluster models are trained by aggregating the local parameter updates generated by clients in the same cluster. By taking advantage of the similar learning characteristics shared by cluster members, cluster models can deliver better personalized performance, such as higher accuracy for the local classification task.

Clustering in an FL environment poses complex challenges as client information cannot be directly accessed because of privacy constraints. The issue of client clustering while adhering to FL constraints has attracted considerable attention recently (Briggs et al., 2020; Duan et al., 2021; Ghosh et al., 2019, 2020; Long et al., 2023; Mansour et al., 2020; Sattler et al., 2021b). The iterative federated clustering algorithm (IFCA) (Ghosh et al., 2020) and HypCluster (Mansour et al., 2020) adopt a performance-based model selection approach where clients determine the cluster identities themselves. Such an approach maintains multiple cluster models, and each client selects the model that performs best on its local data. For instance, clients can measure the performance based on local learning loss or accuracy

in the classification task. Clients that choose the same model are considered members of the same cluster. Allowing clients to choose their own model is a direct and intuitive way to improve the personalized performance, which is the main goal of CFL.

However, one problem with such an approach is that if cluster models are randomly initialized at the start, a particular random model may outperform all others. This may lead all clients to choose the same model, and therefore the clustering process may fail. IFCA (Ghosh et al., 2020) assumes a good initialization where an initial cluster model θ^n is close to θ^{n*} to prove its convergence without clustering failure. IFCA also shows that it can succeed if the initialization requirements are relaxed with random initialization and multiple restarts. Another potential approach is to reduce the randomness of initialized parameters by pre-training cluster models before the CFL process. It is more practical compared with finding a good initialization, or restarting the training with different sets of random parameters in the case of clustering failure. However, obtaining labeled data is the most challenging aspect of pre-training, as it involves significant expenditure to gather and annotate the data for target tasks. In some distributed environments, such data may not even be collectible because of privacy concerns.

On the other hand, unlabeled data with relevant characteristics and features as the target tasks may be widely available from public sources and the Internet. Such unlabeled data may be subject to fewer privacy regulations and can be collected at a lower cost for pre-training. Recently, self-supervised contrastive learning algorithms (Bachman et al., 2019; Caron et al., 2020; Chen et al., 2020; Chen and He, 2021; Grill et al., 2020; He et al., 2020; Misra and van der Maaten, 2020; Zbontar et al., 2021) have attracted major attention for leveraging unlabeled data to aid model training. These algorithms pre-train the encoder part of the model to extract meaningful representations from raw data. The pre-trained encoder increases the efficiency of labeled data when the model is fine-tuned for actual downstream tasks. Contrastive learning assumes that two augmented instances originating from the same data sample should share similar information. Accordingly, representations extracted from these two instances by the encoder should also be similar. On the basis of this idea, the encoder is trained by minimizing the distance between two output representations in the embedding space (Bachman et al., 2019; Chen et al., 2020; Chen and He, 2021; Grill et al., 2020; Zbontar et al., 2021).

Since contrastive encoder pre-training and client clustering can be complementary solutions for providing an efficient design for practical FL systems, it is crucial that we explore how to effectively deploy these approaches together. Therefore, in this study, we leverage contrastive learning to pre-train an encoder in a centralized setting with unlabeled data. The pre-trained encoder is then deployed to enhance the CFL task, where we follow the same approach as in IFCA (Ghosh et al., 2020) for clustering the clients. The goal is to investigate if and how these two prominent approaches can be effectively integrated to improve the model performance in a heterogeneous FL environment. Our method is referred to as “contrastive pre-training–based clustered federated learning” (CP-CFL), and our contributions are summarized as follows:

- We show that CP-CFL significantly improves the performance of CFL by leveraging

contrastive encoder pre-training.

- We conduct extensive experiments using multiple pre-training datasets and downstream FL client datasets to evaluate the performance of CP-CFL. We explore efficient ways to deploy CP-CFL in the context of the pre-trained encoder and downstream classifier head. (More details are presented in Sections 4.4 and 4.5.) Furthermore, we provide various ablation studies to validate its effectiveness.
- Overall, our study contributes to a better understanding of how contrastive encoder pre-training and client clustering can jointly improve the performance of FL.

2. Related work

In this section, we briefly discuss the FL process before covering various related studies on CFL and self-supervised contrastive learning.

2.1. Federated learning and client clustering

The core of FL is to allow distributed clients to collaboratively train a model while considering communication and privacy constraints (Li et al., 2020; Konečný et al., 2016a,b; McMahan et al., 2017). While conventional distributed training retains control over the training data, FL treats the data in each client device as private data. Such data cannot be shared with different parties such as other clients or the central server. FL avoids the transmission of private data by directly training a model on each client while sharing only the trained model parameters with the server. The server has no control over the clients, and a client may or may not participate in a training round and may even drop out of an ongoing one. Usually, the client data in an FL environment are heterogeneous and follow different distributions dependent on the nature of the client (Hard et al., 2018; Ramaswamy et al., 2019). However, vanilla FL is not specifically designed to handle the non-IID data of clients and trains a single global model for all clients, leading to unsatisfactory performance in a highly heterogeneous environment (Sattler et al., 2020b; Yu et al., 2020; Zhao et al., 2018).

One way to improve the performance in FL is to cluster the clients with similar learning characteristics and maintain separate models for each cluster. Many approaches cluster the local parameters received at the server side based on different strategies, such as recursive bipartitioning (Duan et al., 2021; Sattler et al., 2021b), task relatedness (Jamali-Rad et al., 2022), cosine distance (Tian et al., 2022), stochastic expectation maximization (Long et al., 2023), and agglomerative clustering (Briggs et al., 2020). Li et al. (2021a) performed soft clustering that allows clusters with overlapping clients, while Dennis et al. (2021) proposed a one-shot scheme to perform clustering in a single FL round. In addition, some studies have investigated the effectiveness of CFL in the presence of adversarial clients (Ghosh et al., 2019; Sattler et al., 2020a). Other studies have used CFL for specific applications, such as handover prediction in wireless networks (Kim et al., 2021) and human activity recognition (Ouyang et al., 2022). In contrast to the server-side client clustering strategy, IFCA (Ghosh

et al., 2020) and HypCluster (Mansour et al., 2020) perform clustering on the client side by evaluating different models’ performance on local data.

2.2. Federated learning and contrastive learning

Contrastive learning is a form of representation learning that exploits raw unlabeled data by training an encoder on the instance discrimination task. This enables the encoder to learn generic representations that can be transferred to a wide variety of downstream tasks. Unsupervised representation learning can be broadly categorized into two approaches: generative and discriminative. Generative approaches typically learn representations by mapping the output pixels to the input pixels, such as autoencoders (Kingma and Welling, 2013; Vincent et al., 2008). In contrast, discriminative approaches generate labels for a proxy task from the unlabeled data, and the encoder is trained on this proxy task (Gidaris et al., 2018; Noroozi and Favaro, 2016; Pathak et al., 2016). Contrastive learning (Chen et al., 2020; Chen and He, 2021; Grill et al., 2020; He et al., 2020) is one state-of-the-art discriminative approach that trains the encoder by minimizing the distance between representations of a positive pair (e.g., two different augmented views of the same image), while maximizing that of a negative pair (e.g., views from different images). Besides the image domain, contrastive learning is also applied in video (Dave et al., 2022; Pan et al., 2021; Qian et al., 2021), audio (Baevski et al., 2020; Manocha et al., 2021; Saeed et al., 2021), and natural language processing (Gao et al., 2021; Giorgi et al., 2020; Wu et al., 2020), making it suitable for FL environments that seek to provide a wide variety of services to clients.

Several studies have examined representation learning in the FL context, where contrastive learning is integrated into the client’s local training step (Li et al., 2021b; van Berlo et al., 2020; Wu et al., 2021; Zhang et al., 2020; Zhuang et al., 2021, 2022). FedU (Zhuang et al., 2021), FedCA (van Berlo et al., 2020), and FedSSL (Zhuang et al., 2022) assume clients’ private data are unlabeled and perform representation learning as a local training step. RSCFed (Liang et al., 2022) and FedCy (Kassem et al., 2022) consider a semi-supervised setting where some clients contain labeled data and others contain unlabeled data. Hetero-SSFL (Makhija et al., 2022) allows clients to perform representation learning with different local model architectures. FedCon (Long et al., 2021) introduces a learning paradigm where the server contains labeled data and the clients contain unlabeled data. MOON (Li et al., 2021b), FedCKA (Son et al., 2022), and FedIntR (Tun et al., 2023) use contrastive learning to align the representations of global and local models. In contrast, we use contrastive learning as a centralized pre-training step to construct a base encoder for CFL. HotFed (Zhao et al., 2021) and FedAUX (Sattler et al., 2021a) closely relate to our work since they pre-train the model using self-supervised contrastive learning; however, HotFed (Zhao et al., 2021) transfers the pre-trained model to a vanilla FL setup, while FedAUX (Sattler et al., 2021a) transfers it to a federated distillation setting. The aforementioned studies shed light on the possibility of incorporating contrastive pre-training into the CFL process, with the potential to increase the convergence rate and improve the performance of cluster models.

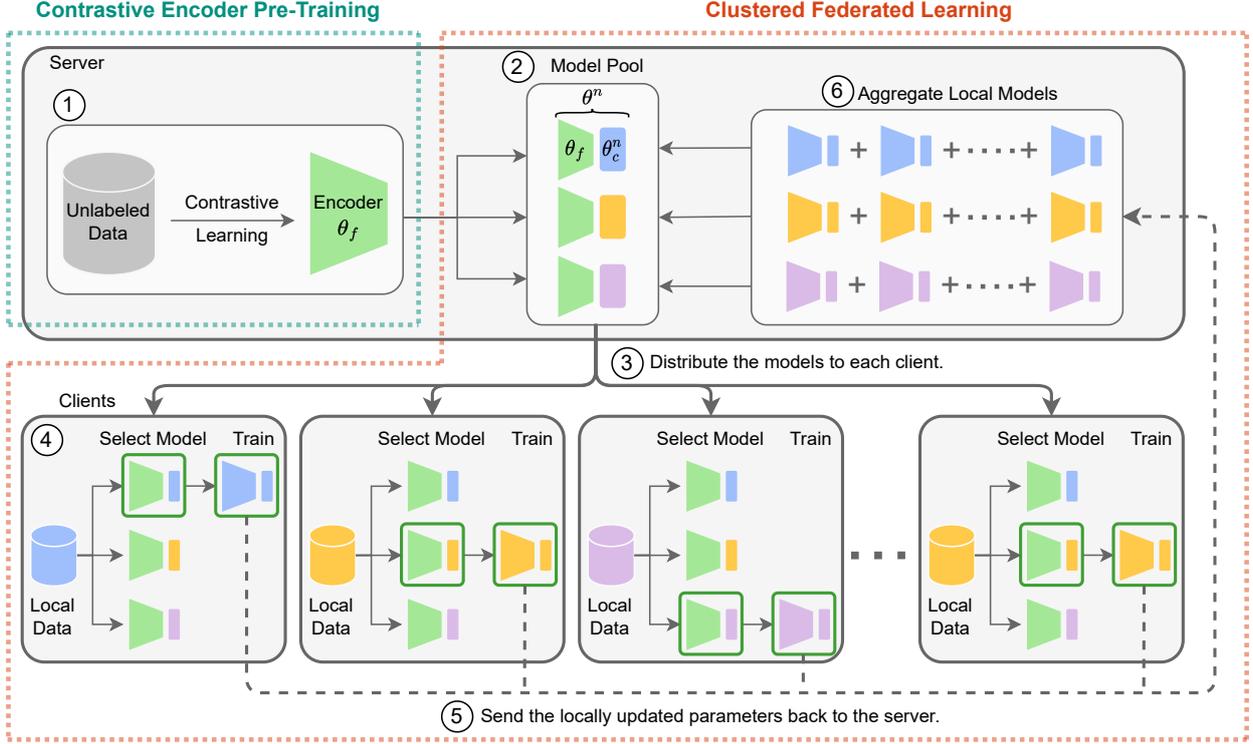


Figure 1: Overview of CP-CFL, which uses contrastive learning to pre-train the encoder θ_f on the unlabeled data. The pre-trained encoder θ_f is then deployed in the CFL task, which aims to train cluster-level models for different client clusters.

3. Contrastive pre-training–based CFL

We present an overview of our CP-CFL framework before delving into details in later sections.

3.1. Overview

As illustrated in Figure 1, CP-CFL consists of two stages: the pre-training stage and the CFL stage. In the pre-training stage, the central server uses contrastive learning to pre-train the encoder θ_f on unlabeled data. To apply the pre-trained encoder θ_f in a useful downstream task such as classification, we can attach a classifier head θ_c to obtain a complete classification model $\theta = (\theta_f, \theta_c)$. For a CFL task with N clusters, the server generates a pool of N classification models $\{\theta^n\}_{n=1}^N = \{(\theta_f, \theta_c^n)\}_{n=1}^N$ by joining the pre-trained encoder θ_f with N randomly initialized classifier heads $\{\theta_c^n\}_{n=1}^N$. In each global round, the model pool is distributed to the clients for clustering and local training. Each client evaluates the models on its local dataset and selects the best-performing one with the lowest error. The selected model is then trained, and updated parameters are sent back to the server. Clients that choose the same model are identified as members of the same cluster. Local parameter updates from the cluster members are aggregated to update the respective cluster model (i.e., the model mutually selected by the cluster members). Each cluster model is

tailored to the unique underlying data distribution of its respective cluster, providing better personalized performance to the clients.

3.2. Problem formulation of CP-CFL

Contrastive pre-training: Given an unlabeled dataset $\tilde{\mathcal{D}}$, the encoder θ_f is trained to extract representations from each sample $\tilde{x} \in \tilde{\mathcal{D}}$. These extracted representations $h = \theta_f(\tilde{x})$ should capture generic information transferable to a wide variety of downstream tasks. To achieve this, we use a typical contrastive learning step for training. Firstly, the data sample \tilde{x} is randomly augmented into two different instances, x_i and x_j . The encoder θ_f then extracts representations $h_i = \theta_f(x_i)$ and $h_j = \theta_f(x_j)$ from the augmented instances. The goal is to train θ_f so that it generates similar h_i and h_j . Hence, a form of contrastive loss $\ell_{i,j}$ can be designed to minimize the distance between h_i and h_j in the embedding space. (Section 3.3 provides further details on how $\ell_{i,j}$ is computed for different contrastive learning techniques.) Generally, the encoder θ_f is trained by solving

$$\min_{\theta_f} \frac{1}{|\tilde{\mathcal{D}}|} \sum_{\tilde{x} \in \tilde{\mathcal{D}}} \ell_{i,j}. \quad (1)$$

Clustered federated learning: Each FL client $u \in \mathcal{U}$ stores its own local private dataset \mathcal{D}_u , which can be used for model training. However, a conventional FL system trains a single global model for all the clients in the population \mathcal{U} without considering the heterogeneous nature of the clients. In contrast, CFL assumes that clients in the population \mathcal{U} can be further clustered into N different groups (i.e., $\mathcal{G}^1, \mathcal{G}^2, \dots, \mathcal{G}^N$) based on their learning characteristics. CFL aims to provide better personalized performance to the clients by fitting a cluster-level model θ^n for each cluster \mathcal{G}^n . To take advantage of contrastive pre-training, each cluster model θ^n is constructed by combining the pre-trained encoder θ_f with a randomly initialized classifier head θ_c^n . We can train the cluster model θ^n for a cluster \mathcal{G}^n by solving

$$\min_{\theta^n} \frac{1}{|\mathcal{G}^n|} \sum_{u \in \mathcal{G}^n} \mathcal{L}(\theta^n, \mathcal{D}_u), \quad (2)$$

where \mathcal{L} is the loss calculated for each client u in the cluster \mathcal{G}^n . Specifically, we calculate \mathcal{L} as

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell_{\text{CE}}(\theta, x, y), \quad (3)$$

where (x, y) is the labeled data pair and ℓ_{CE} is the cross-entropy loss defined by

$$\ell_{\text{CE}}(\theta, x, y) = - \sum_{m=1}^{|y|} y_m \log \theta(x)_m, \quad (4)$$

where y is assumed to be a one-hot encoded label vector, $|y|$ is the number of classes, and y_m and $\theta(x)_m$ are the ground truth label and softmax probability prediction for the m th class, respectively.

Despite the advantages of CFL, determining the cluster structure $\{\mathcal{G}^n\}_{n=1}^N$ can be challenging, since we cannot access the clients’ private information; therefore, we use the model selection strategy to determine the cluster identity $n_u \in \{1, \dots, N\}$ of client u without violating the FL constraints. In this strategy, the server maintains a pool of cluster models $\{\theta^n\}_{n=1}^N$ that share the same architecture but different sets of parameters. At each global round t , the server distributes all cluster models $\{\theta^n\}_{n=1}^N$ to the participating clients. Each client u evaluates $\{\theta^n\}_{n=1}^N$ on its local dataset \mathcal{D}_u and then selects the best-performing model θ^{n_u} for its data distribution. Specifically, clustering is achieved by solving Eq. (5) for each client u :

$$n_u = \arg \min_n \mathcal{L}(\theta^n, \mathcal{D}_u). \quad (5)$$

In Eq. (5), $n \in \{1, \dots, N\}$ and n_u is the index of the selected model as well as the cluster identity of client u . Each client u locally trains the selected model θ^{n_u} on dataset \mathcal{D}_u , and then sends the updated parameters ϕ_u and the cluster identity n_u back to the server. Local parameter updates from the clients of the same cluster are aggregated by weighted averaging (McMahan et al., 2017), shown in Eq. (6), to update the respective cluster model:

$$\theta^n = \sum_{u \in \mathcal{G}^n} \frac{|D_u|}{\sum_{u' \in \mathcal{G}^n} |D_{u'}|} \phi_u. \quad (6)$$

In Eq. (6), $|D_u|$ refers to the size of the local dataset of client u .

3.3. Contrastive encoder pre-training techniques

In our work, we explore multiple contrastive learning algorithms—namely, SimCLR (Chen et al., 2020), BYOL (Grill et al., 2020), and SimSiam (Chen and He, 2021)—for pre-training the encoder. We describe the common key components of these contrastive learning algorithms as follows:

- *Random augmentation*: Augmentations are necessary to generate a pair of positive samples from an input sample. For each sample $\tilde{x} \in \tilde{\mathcal{D}}$, random augmentations are applied to generate a positive sample pair: $x_i = \text{Augment}(\tilde{x})$ and $x_j = \text{Augment}(\tilde{x})$.
- *Encoder θ_f* : Contrastive learning pre-trains the encoder θ_f to extract generic representations h given an input \tilde{x} . The pre-trained encoder θ_f can be paired with a task-specific head (in our case, we attach a classifier head θ_c) and fine-tuned on a small labeled dataset to perform a useful downstream task.
- *Distance metric*: Many contrastive learning approaches use cosine similarity to measure the distance between two representations (Chen et al., 2020; Chen and He, 2021; Grill et al., 2020). The cosine similarity between two representation vectors v_i and v_j is calculated as

$$\text{Sim}(v_i, v_j) = \frac{v_i}{\|v_i\|_2} \cdot \frac{v_j}{\|v_j\|_2}, \quad (7)$$

where $\|\cdot\|_2$ is the ℓ_2 -norm.

In the following subsections, we provide a brief overview of each of the three contrastive learning approaches. Table 1 and Figure 2 compare the detailed characteristics of different contrastive learning approaches, while Algorithm 1 describes the general contrastive learning procedure.

	Positive samples	Negative samples	Weight sharing	Projector	Predictor	Stop gradient
SimCLR	✓	✓	✓	✓	–	–
BYOL	✓	–	–	✓	✓	✓
SimSiam	✓	–	✓	✓	✓	✓

Table 1: Comparison between different self-supervised contrastive learning techniques.

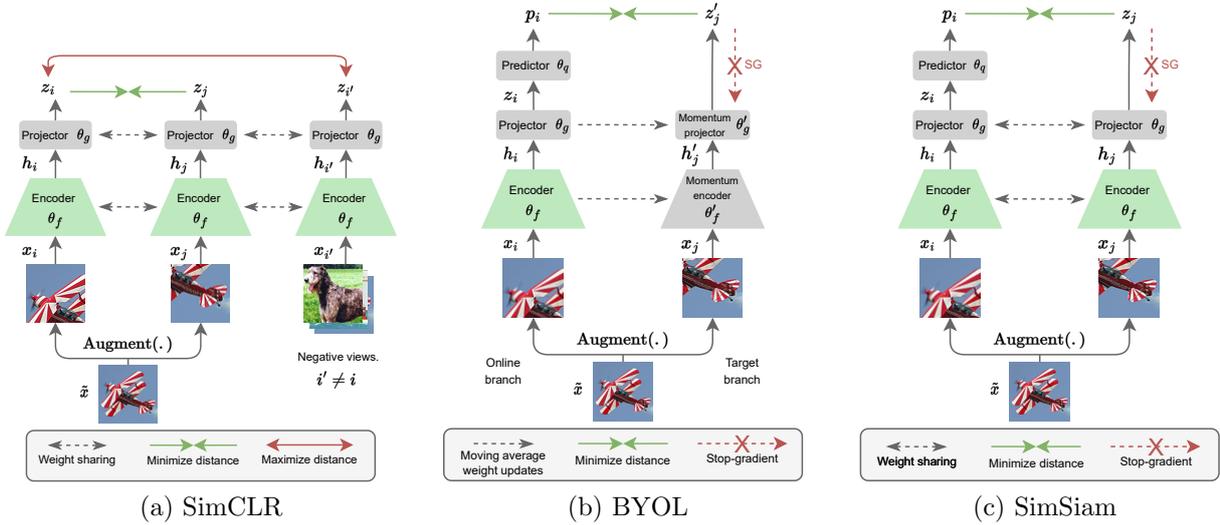


Figure 2: Self-supervised contrastive learning techniques: (a) SimCLR (Chen et al., 2020); (b) BYOL (Grill et al., 2020); (c) SimSiam (Chen and He, 2021).

3.3.1. SimCLR

For each image \tilde{x} in an input batch \mathcal{B} , SimCLR (Chen et al., 2020) generates two views, x_i and x_j , using a random augmentation function, $\text{Augment}(\cdot)$. (This results in a total of $2|\mathcal{B}|$ augmented views for each input batch \mathcal{B} .) The two augmented views x_i and x_j are considered a positive pair since they originate from the same image \tilde{x} . The encoder θ_f extracts representations $h_i = \theta_f(x_i)$ and $h_j = \theta_f(x_j)$ from the augmented views, while the projection head θ_g transforms the extracted representations into projections $z_i = \theta_g(h_i)$ and $z_j = \theta_g(h_j)$. If we consider x_i as the anchor sample, then x_j acts as the corresponding positive counterpart, and vice versa. Using x_i as the anchor, SimCLR calculates its loss

Algorithm 1 Contrastive encoder pre-training

- 1: **Input:** encoder θ_f , unlabeled dataset $\tilde{\mathcal{D}}$, number of epochs E , learning rate η .
 - 2: **Contrastive pre-training**($\theta_f, \tilde{\mathcal{D}}, E, \eta$):
 - 3: **for** epoch $e = 0, 1, \dots, E - 1$ **do**
 - 4: **for** each batch $\mathcal{B} \in \tilde{\mathcal{D}}$ **do**
 - 5: **for** each sample $\tilde{x} \in \mathcal{B}$ **do**
 - 6: $x_i = \text{Augment}(\tilde{x})$
 - 7: $x_j = \text{Augment}(\tilde{x})$
 - 8: $h_i = \theta_f(x_i)$
 - 9: $h_j = \theta_f(x_j)$
 - 10: Calculate $\ell_{i,j}$ using h_i and h_j ▷ Use Eq. (8), (9), or (12).
 - 11: $\mathcal{L} = \frac{1}{|\mathcal{B}|} \sum_{\tilde{x} \in \mathcal{B}} \ell_{i,j}$
 - 12: Update θ_f with $\eta \nabla \mathcal{L}$
 - 13: **Output:** encoder θ_f
-

based on the InfoNCE loss (van den Oord et al., 2018) as

$$\ell_{i,j} = -\log \frac{\exp(\text{Sim}(z_i, z_j)/\tau)}{\sum_{i'=1}^{2|\mathcal{B}|} \mathbb{1}_{[i' \neq i]} \exp(\text{Sim}(z_i, z_{i'})/\tau)}, \quad (8)$$

where τ is the temperature. Every other augmented view $x_{i'}$, where $i' = \{1, \dots, 2|\mathcal{B}|\}$, $i' \neq i$, is treated as a negative counterpart to x_i , with the distance between their representations maximized. These steps are repeated by treating each augmented view in a batch as an anchor. In SimCLR, negative samples help prevent the solution from collapsing into a trivial constant (Chen and He, 2021).

3.3.2. BYOL

BYOL (Grill et al., 2020) uses an online network θ branch and a target network θ' branch for training. During training, augmented views x_i and x_j are fed into the online and the target branches, respectively. The online branch extracts $h_i = \theta_f(x_i)$ and $z_i = \theta_g(h_i)$, while the target branch uses the momentum encoder θ'_f and the momentum projector θ'_g to extract $h'_j = \theta'_f(x_j)$ and $z'_j = \theta'_g(h'_j)$. The momentum encoder θ'_f and the momentum projector θ'_g have the same structure as θ_f and θ_g , but their parameters are moving-average versions of those of θ_f and θ_g , respectively. The online branch makes a prediction $p_i = \theta_q(z_i)$ using the predictor θ_q and the network is trained by minimizing the distance between p_i and z'_j . In a sense, BYOL trains the encoder θ_f by using a view x_i to predict the representation of the other view x_j . The loss between the anchor and its positive sample is computed as follows:

$$\ell_{i,j} \triangleq \|p_i - z'_j\|_2^2 = 2 - 2 \cdot \text{Sim}(p_i, z'_j). \quad (9)$$

BYOL also swaps the inputs x_j and x_i for each branch and then calculates the prediction p_j and target branch output z'_i to minimize the distance between them. Only the online branch is updated with gradient backpropagation, while the target branch parameters are updated as the moving-average versions of the online branch using Eqs. (10) and (11):

$$\theta'_f = \beta\theta'_f + (1 - \beta)\theta_f, \quad (10)$$

$$\theta'_g = \beta\theta'_g + (1 - \beta)\theta_g, \quad (11)$$

where β is the target branch update rate. BYOL uses the target branch and stops the gradient flow through it to prevent the solution from collapsing to a constant (Chen and He, 2021).

3.3.3. SimSiam

SimSiam (Chen and He, 2021) shares similar properties with both SimCLR and BYOL. SimSiam uses the weight-sharing encoder θ_f and projector θ_g as in SimCLR to extract representations $h_i = \theta_f(x_i)$ and $h_j = \theta_f(x_j)$, as well as the projections $z_i = \theta_g(h_i)$ and $z_j = \theta_g(h_j)$. Similarly to BYOL, the predictor θ_q in one branch of SimSiam makes a prediction $p_i = \theta_q(z_i)$ to predict the output z_j of the other branch. The contrastive loss is calculated as the negative cosine similarity between the prediction p_i and the projection z_j :

$$\ell_{i,j} = -\text{Sim}(p_i, z_j). \quad (12)$$

The inputs x_j and x_i can also be swapped for each branch, and the predictor θ_q can make the prediction p_j to contrast with the output z_i of the other branch. The encoder is trained by minimizing the distance between p_i and z_j , as well as the distance between p_j and z_i . SimSiam stops the gradient flow through the branch without the predictor to prevent the solution from collapsing to a constant.

3.4. CP-CFL algorithm

The detailed procedure of CP-CFL is presented in Algorithm 2. Initially, the server pre-trains the encoder θ_f on an unlabeled dataset $\tilde{\mathcal{D}}$ using contrastive learning. The pre-trained encoder θ_f is then joined with a randomly initialized classifier head θ_c^n to generate a cluster model $\theta^n = (\theta_f, \theta_c^n)$, where $n = 1, 2, \dots, N$. Although the encoder θ_f is pre-trained, random parameters in the classifier head θ_c^n still have the potential to damage the clustering process. Therefore, for a few initial rounds (i.e., when $t < T_c$), the clients explore by selecting a model at random rather than based on the performance. During this period, the encoder θ_f of each cluster model θ^n is frozen, with the local training updating only the classifier head θ_c^n . Here, T_c is significantly smaller than the total number of rounds T (by default, we set $T_c = 10$ and $T = 100$). When $t \geq T_c$, each client u determines its cluster identity n_u by evaluating the cluster models $\{\theta^n\}_{n=1}^N$ on the local dataset \mathcal{D}_u . Client u then trains the entire model θ^{n_u} and then sends the locally updated parameters ϕ_u and its cluster identity n_u back to the server. Local model updates from the clients of the same cluster \mathcal{G}^n are aggregated to update the corresponding cluster model θ^n .

Algorithm 2 Contrastive pre-training–based clustered federated learning

- 1: **Input:** encoder θ_f , classifier head θ_c , unlabeled dataset $\tilde{\mathcal{D}}$, number of contrastive epochs E_p , contrastive learning rate η_p , number of clusters N , number of rounds T , number of classifier head training rounds T_c , number of local epochs E_ℓ , local learning rate η_ℓ .

 - 2: **Server executes:**
 - 3: $\theta_f = \mathbf{Contrastive\ pre-training}(\theta_f, \tilde{\mathcal{D}}, E_p, \eta_p)$ using Algorithm 1
 - 4: **for** $n = 1, 2, \dots, N$ **do** ▷ Generate a pool of N cluster models.
 - 5: $\theta_c^n =$ randomly initialize θ_c
 - 6: Join θ_f and θ_c^n to obtain cluster model $\theta^n = (\theta_f, \theta_c^n)$
 - 7: **for** round $t = 0, 1, \dots, T - 1$ **do**
 - 8: Initialize empty clusters $\mathcal{G}^n = \{\}$, where $n = \{1, \dots, N\}$
 - 9: Get participating client population \mathcal{U}
 - 10: **for** each client $u \in \mathcal{U}$ in parallel **do**
 - 11: $(\phi_u, n_u) = \mathbf{Local\ update}(u, t, \{\theta^n\}_{n=1}^N)$
 - 12: Append u to \mathcal{G}^{n_u}
 - 13: **for** each cluster \mathcal{G}^n , where $n = 1, 2, \dots, N$ **do**,
 - 14: $\theta^n = \sum_{u \in \mathcal{G}^n} \frac{|D_u|}{\sum_{u' \in \mathcal{G}^n} |D_{u'}|} \phi_u$

 - 15: **Clients execute: Local update** $(u, t, \{\theta^n\}_{n=1}^N)$:
 - 16: **if** $t < T_c$ **then** ▷ Explore with a random model.
 - 17: $n_u =$ select randomly from $\{1, 2, \dots, N\}$
 - 18: Freeze the encoder part θ_f of θ^{n_u}
 - 19: **else** ▷ Exploit the best model.
 - 20: $n_u = \arg \min_n \mathcal{L}(\theta^n, \mathcal{D}_u)$, where $n = 1, 2, \dots, N$
 - 21: Synchronize local model $\phi_u = \theta^{n_u}$
 - 22: **for** epoch $e = 0, 1, \dots, E_\ell - 1$ **do**
 - 23: **for** each batch $\mathcal{B} \in \mathcal{D}_u$ **do**
 - 24: $\mathcal{L} = \frac{1}{|\mathcal{B}|} \sum_{(x,y) \in \mathcal{B}} \ell_{\text{CE}}(\phi_u, x, y)$
 - 25: $\phi_u = \phi_u - \eta_\ell \nabla \mathcal{L}$
 - 26: return (ϕ_u, n_u) to server

 - 27: **Output:** cluster models $\{\theta^n\}_{n=1}^N$, clusters $\{\mathcal{G}^n\}_{n=1}^N$
-

3.5. Characteristics of CP-CFL

To provide a clear insight into CP-CFL, we compare the characteristics of CP-CFL with other baseline approaches in Table 2. In contrast to vanilla FedAvg (McMahan et al., 2017), which trains a single global model, IFCA (Ghosh et al., 2020) and CP-CFL train multiple cluster models. Models in FedAvg and IFCA are randomly initialized, whereas CP-CFL pre-trains the encoder part with self-supervised contrastive learning. IFCA relaxes the random initialization with multiple restarts to handle the clustering failure, while CP-CFL relies on the pre-trained encoder and T_c . CP-CFL performs the pre-training on the server without extra computational burden on client devices. Moreover, T_c in CP-CFL is usually small and can be integrated into the total number of rounds T , requiring no extra communication cost than IFCA. Generally, the communication cost C for an FL client can be calculated as follows (Qu et al., 2022):

$$C = T \times (C_{sc} + C_{cs}), \quad (13)$$

where T is the number of rounds, and C_{sc} and C_{cs} denote the server-to-client and client-to-server costs, where both can be derived as number of transmitted models \times model size. We use the same model architecture in all approaches and denote the model size as S . Based on above formulation, we calculate the cost for FedAvg (McMahan et al., 2017) as

$$C_{\text{FedAvg}} = T \times (1 \times S + 1 \times S) = 2ST \quad (14)$$

and the cost for CFL approaches as

$$C_{\text{CFL}} = T \times (N \times S + 1 \times S) = (N + 1)ST. \quad (15)$$

In general, CFL approaches are $\frac{N+1}{2}$ times more expensive than FedAvg in terms of communication cost. For instance, recalling that N is the number of cluster models, we find that CFL approaches are twice as expensive than FedAvg when $N = 3$.

	Number of models	Model initialization	Handling clustering failure	Communication cost
FedAvg	Single	Random	–	$2ST$
IFCA	Multiple	Random	Multiple restarts	$(N + 1)ST$
CP-CFL	Multiple	Contrastive pre-training	Pre-trained encoder + T_c	$(N + 1)ST$

Table 2: Characteristics of CP-CFL.

4. Experiments

In this section, we present the evaluation results of CP-CFL on various experimental settings. We first describe the common experimental settings, before delving into the details of our studies.

4.1. Experimental settings

Unless otherwise stated, we use the following settings in all of our experiments.

Dataset: For our experiments, we mainly use the STL-10 dataset (Coates et al., 2011), which contains both unlabeled and labeled images with a shape of 96×96 . The unlabeled data portion contains 100,000 images, whereas the labeled training and testing portions contain 5000 and 8000 image-label pairs for 10 different classes, respectively.

Model: The encoder θ_f contains four 3×3 convolutional layers followed by a dense layer. The convolutional layers contain 64, 128, 192, and 256 filters, respectively. All layers in θ_f are ReLU activated. The encoder θ_f accepts an input image size of 96×96 , and the dimension of the output representation is 256. The classifier head θ_c is a single dense layer with 10 output neurons and softmax activation. A complete classification model θ can be obtained by directly joining θ_f and θ_c together.

Contrastive pre-training: To generate positive samples for contrastive pre-training, we rely on a combination of data augmentations, including horizontal flip, random cropping, and color transformations including jitters, brightness, and channel manipulations (Chen et al., 2020). We train θ_f for 300 epochs with a batch size of 500 on an unlabeled dataset using the Adam optimizer and an initial learning rate of 0.001. SimCLR (Chen et al., 2020), BYOL (Grill et al., 2020), and SimSiam (Chen and He, 2021) incorporate different architectural configurations for training; therefore, we lightly tune the hyperparameters for each method. Specifically, for the projection head θ_g , we tune the number of neurons in the dense layers, setting it to either 256 or 512. We pre-train multiple encoders with different hyperparameters and then select the best one through linear evaluation (Chen et al., 2020; Grill et al., 2020). Linear evaluation involves training and testing a classifier head on top of the frozen encoder using a proxy dataset, and the proxy test accuracy can be used to select the best encoder, which we transfer to the downstream task. (We mainly use the STL-10 unlabeled portion for contrastive pre-training in our experiments and simply use the STL-10 training and testing portions as the proxy dataset for linear evaluation. It is worth mentioning again that the encoder is trained only on the unlabeled portion.) Our pre-training configurations are described below:

- *SimCLR:* The projection head θ_g contains two dense layers with 256 neurons each. ReLU activation is applied to the first layer. We set the value of temperature τ as 0.1.
- *BYOL:* Each of the two dense layers in θ_g contains 512 neurons. Both layers are batch normalized, although only the first layer is ReLU activated. The predictor θ_q contains a bottleneck dense layer with 64 neurons followed by an output dense layer with 512 neurons. The bottleneck layer is batch normalized and ReLU activated. We set the target branch update rate β as 0.9.
- *SimSiam:* The projection head θ_g and the predictor θ_q are configured similarly to the settings in BYOL.

Clustered federated learning: For the federated setting, we construct the local data of 60 clients by sampling from the training portion of a labeled dataset. To obtain a clear cluster structure, the clients are divided into three groups with different underlying data distributions. An example with the STL-10 training portion is shown in Figure 3. Each client contains data from four different classes, and clients within the same cluster share the same classes. Specifically, for a dataset with 10 available classes, a group of clients may contain classes 1–4, another group may contain classes 4–7, and the third group may contain classes 7–10. Since each client contains four classes, we randomly allocate 20 images for two of the classes and five images for the remaining two (i.e., a total of 50 images per client) to reflect the non-IID settings of FL. We also construct local testing sets that follow each client’s data distribution by sampling from the testing portion of the dataset. We set the number of clusters N to 3, the number of global rounds T to 100, and the number of classifier head training rounds T_c to 10. Each client trains its selected model for three local epochs using a batch size of 32 and the Adam optimizer with a learning rate of 0.001.

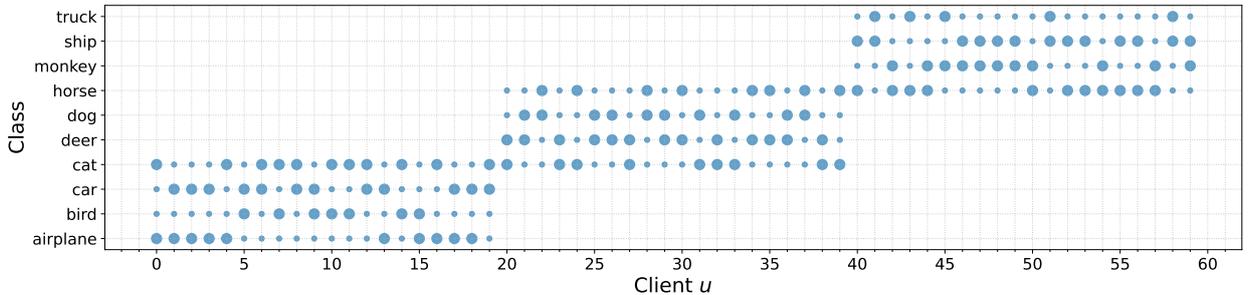


Figure 3: The amount of private data in clients for different classes, sampled from the training portion of the STL-10 dataset. Small dots and large dots represent 5 and 20 images, respectively.

4.2. Encoder pre-training with contrastive learning

To validate the effectiveness of contrastive pre-training for the downstream CFL task, we pre-train the encoder using SimCLR (Chen et al., 2020), BYOL (Grill et al., 2020), and

	Pre-training	Accuracy (%)			
		$T = 25$	$T = 50$	$T = 75$	$T = 100$
FedAvg	None	32.73	39.83	41.20	44.13
	SimCLR	46.40	53.33	53.97	59.33
IFCA	None	57.60	63.30	65.60	67.80
	FedAvg	69.93	71.77	70.40	70.90
CP-CFL	SimCLR	72.30	75.80	76.30	76.80
	BYOL	67.37	72.83	74.03	74.87
	SimSiam	56.90	64.23	66.80	68.23

Table 3: Test accuracy (%) for the STL-10-to-STL-10 (denoting pre-training data to client data) task. The results are obtained by averaging the local test accuracies of all the clients.

	Pre-training	F1 score		AUROC			
		Macro	Weighted	OvR Macro	OvR Weighted	OvO Macro	OvO Weighted
FedAvg	None	0.3708	0.3953	0.8342	0.8333	0.8374	0.8349
	SimCLR	0.5564	0.5699	0.9036	0.9045	0.9068	0.9052
IFCA	None	0.6842	0.6769	0.9562	0.9491	0.9595	0.9542
	FedAvg	0.7131	0.7075	0.9641	0.9584	0.9669	0.9626
CP-CFL	SimCLR	0.7700	0.7657	0.9752	0.9717	0.9774	0.9745
	BYOL	0.7494	0.7455	0.9733	0.9692	0.9755	0.9722
	SimSiam	0.6870	0.6827	0.9578	0.9511	0.9609	0.9559

Table 4: F1 score and AUROC for the STL-10-to-STL-10 task ($T = 100$). OvR, one-vs-rest. OvO, one-vs-one.

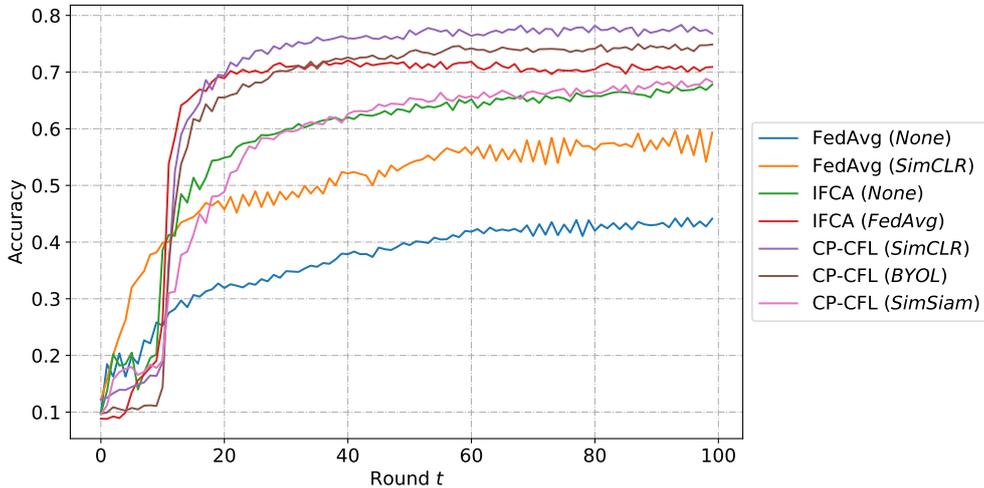


Figure 4: Test accuracy curves for the STL-10-to-STL-10 task.

SimSiam (Chen and He, 2021). We use the STL-10 unlabeled portion for pre-training, while the labeled training and testing portions are used to set up the clients’ data as mentioned in Section 4.1. Figure 3 illustrates the distribution of private data across the clients. Although the data distribution of the STL-10 unlabeled portion is similar to that of the labeled portion, it is not identical. This appropriately reflects the practical scenarios in which the central server may not be able to gather pre-training data that exactly match the clients’ data distribution.

Table 3 compares the performance of CP-CFL with that of different baselines, including vanilla FedAvg (McMahan et al., 2017), and IFCA (Ghosh et al., 2020) with and without a pre-trained encoder. We denote the pre-training approaches in italics in parentheses after the abbreviations for the different methods. IFCA (*FedAvg*) uses the encoder obtained from the FedAvg training as the pre-trained encoder. (Using the encoder from the FL global model is one option for getting a pre-trained encoder for the CFL task. However, this can

be considered as pre-training on the clients, which would impose additional communication and computational burden on the clients.) In the event that clustering fails for IFCA, we restart the process with a different set of randomly initialized parameters. For the evaluation results, we average the local test accuracies from all clients, obtained by their respective cluster models. From Table 3 and Figure 4, we can observe that CP-CFL (*SimCLR*) and CP-CFL (*BYOL*) significantly outperform all the baselines at $T = 100$. CP-CFL (*SimSiam*) outperforms IFCA (*None*) by a small margin but performs worse than IFCA (*FedAvg*). However, recall that IFCA (*FedAvg*) uses the encoder directly trained on the clients’ labeled data, whereas SimCLR, BYOL, and SimSiam pre-train on unlabeled data. The inefficiency of vanilla FedAvg in aggregating heterogeneous parameters degrades the resulting encoder quality, leading to lower downstream task performance for IFCA (*FedAvg*) compared with CP-CFL (*SimCLR*) and CP-CFL (*BYOL*). Despite having a higher communication cost, CP-CFL can significantly outperform FedAvg even with a small number of global rounds T , as shown in Figure 4.

We report additional evaluation metrics—i.e., F1 score and area under the receiver operating characteristic (AUROC)—in Table 4. Figure 5 shows how the cluster identity n_u of each client u changes at different rounds for CP-CFL (*SimCLR*). From Figure 5, we can observe that correct client clusters are identified almost immediately when the clients start selecting models.

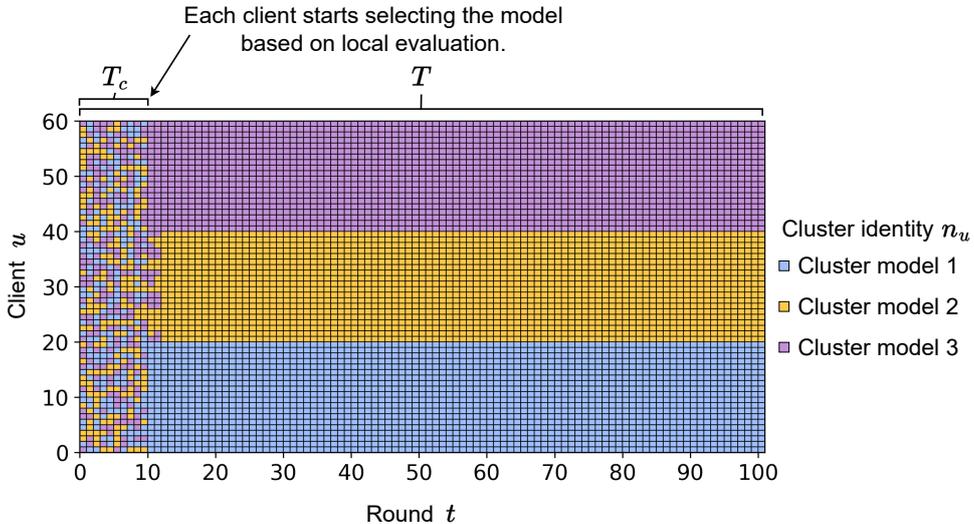


Figure 5: Cluster identity of clients at each communication round for CP-CFL (*SimCLR*).

4.3. Influence of pre-training and downstream client data similarity

In practice, the pre-training data collected by the central server may not perfectly match the private data generated by the clients. The distribution shift and differences in characteristics between the pre-training and downstream datasets could influence the model performance. To investigate this, we pre-train the encoder using SimCLR on the STL-10 unlabeled portion and then transfer it to three different downstream CFL tasks, where clients’ private

	Pre-training	Pre-training data to client data			
		STL-10 to STL-10 (0.7346)	STL-10 to CIFAR-10 (0.7306)	STL-10 to MNIST (0.6108)	EMNIST to MNIST (0.6714)
FedAvg	None	44.13	40.10	95.47	95.47
	SimCLR	59.33	48.20	93.03	94.83
IFCA	None	67.80	65.40	98.37	98.37
	FedAvg	70.90	70.03	98.90	98.90
CP-CFL	SimCLR	76.80	75.00	98.40	99.27
Centralized (all clients)	SimCLR	64.37	63.03	–	–
Centralized (by cluster)	SimCLR	77.47	76.00	–	–

Table 5: Test accuracy (%) of CP-CFL (*SimCLR*) with different pre-training and downstream client datasets at $T = 100$. Note that pre-training data are unlabeled. The measure of the pre-training dataset’s relevance to the downstream client dataset is shown in parentheses for each task.

and testing data are sampled from the STL-10, CIFAR-10 (Krizhevsky et al., 2009), and MNIST (LeCun et al., 1998) datasets, respectively. We also pre-train the encoder using the EMNIST dataset (Cohen et al., 2017) without using labels and then transfer it to the MNIST task. Table 5 shows the evaluation results as well as a measure of the pre-training dataset’s relevance to the downstream client data in parentheses.

To quantify the relevance between the pre-training and downstream datasets for each task, we randomly sample 300 images from each dataset and extract their representations using the corresponding pre-trained encoder. We then measure their cross similarities using the cosine similarity, Eq. (7), and report the average score. The results, shown in Table 5, indicate that the pre-training data and downstream client data share relevant characteristics for STL-10-to-STL-10, STL-10-to-CIFAR-10, and EMNIST-to-MNIST tasks, as evidenced by their relatively higher measurement scores compared with the score for the STL-10-to-MNIST task, where the pre-training data and downstream client data do not share relevant characteristics.

Table 5 confirms that contrastive pre-training can significantly improve the downstream CFL performance if the unlabeled pre-training data share relevant characteristics with the downstream client data. This is evident in STL-10-to-STL-10, STL-10-to-CIFAR-10, and EMNIST-to-MNIST tasks, where CP-CFL (*SimCLR*) outperforms IFCA (*FedAvg*) by 5.90%, 4.97%, and 0.37%, respectively. In contrast, when the STL-10 pre-trained encoder is transferred to MNIST client data, there is no improvement over the baseline IFCA (*FedAvg*) (decreases by 0.5%). It is worth noting that the unlabeled portion of the STL-10 dataset contains extra image classes that are not in the labeled portion. In addition, the CIFAR-10 dataset includes the “frog” class, as opposed to the “monkey” class present in the STL-10 dataset. The significant performance gains observed in STL-10-to-STL-10 and STL-10-to-CIFAR-10 tasks suggest that self-supervised contrastive pre-training can be beneficial even if the pre-training data are not identical to the downstream data. However, it

is ineffective when the encoder is transferred between vastly dissimilar data, such as in the STL-10-to-MNIST task.

Table 5 also presents the evaluation results for two centralized settings. In the first setting, we gather and consolidate the private data of all clients into a single dataset and train a model in the centralized manner. We evaluate the model on the testing dataset created in the same manner. In the second setting, we gather and group the clients’ data into three separate datasets based on their true underlying cluster structures; therefore, this setting comprises three pairs of training and testing datasets. We train and evaluate a model for each pair and report the average results. Notably, the performance of CP-CFL is comparable to that of the second centralized setting, which can be regarded as the upper bound.

4.4. Learning capacity of the classifier head

The classifier head θ_c plays an important role in the downstream CFL performance as it primarily generates final class probabilities from the extracted representations. The learning capacity of θ_c can significantly affect the performance of a classification model $\theta = (\theta_f, \theta_c)$. To investigate this, we study four different classifier heads: $\theta_c, \theta_{c-1}, \theta_{c-2}$, and θ_{c-3} , which differ in the number of layers and neurons. Here θ_c is our default classifier head, which contains a single output layer. The other classifier heads, $\theta_{c-1}, \theta_{c-2}$, and θ_{c-3} , contain one, two, and three additional dense layers with ReLU activation, respectively. The details of different classifier head structures are illustrated in Figure 6.

Table 6 shows the performance of CP-CFL (*SimCLR*) with different classifier heads. Compared with θ_c , having extra hidden layers in the classifier head (i.e., $\theta_{c-1}, \theta_{c-2}$, and θ_{c-3}) can result in better performance when the number of global rounds is limited. For instance, in the STL-10-to-STL-10 task, θ_{c-3} at $T = 25$ achieves similar performance as θ_c at $T = 100$. Such an observation can be particularly useful for cases when client devices have limited communication resources, and we can trade off the computational cost of a few extra layers to obtain an acceptable performance with fewer global rounds. The performance

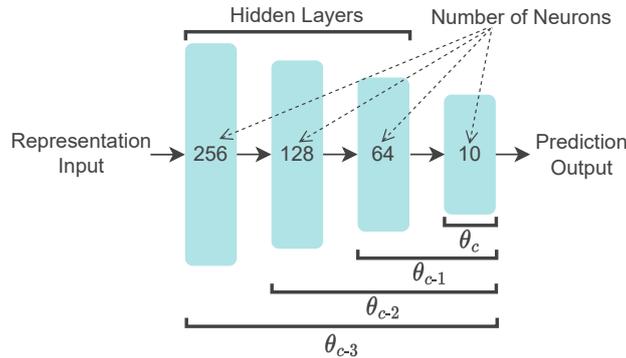


Figure 6: Detailed structure of classifier heads: $\theta_c, \theta_{c-1}, \theta_{c-2}$, and θ_{c-3} . The default classifier head θ_c contains an output layer with 10 neurons. θ_{c-1} contains one more dense layer, in addition to the output layer. θ_{c-2} contains two more dense layers, and so on.

	STL-10-to-STL-10			STL-10-to-CIFAR-10			EMNIST-to-MNIST		
	$T = 25$	$T = 75$	$T = 100$	$T = 25$	$T = 75$	$T = 100$	$T = 25$	$T = 75$	$T = 100$
θ_c	72.30	76.30	76.80	67.20	74.63	75.00	96.20	99.20	99.27
θ_{c-1}	74.47	77.77	78.23	70.23	74.40	75.70	96.87	98.87	98.93
θ_{c-2}	75.43	78.57	77.60	69.77	74.57	75.10	97.67	99.03	98.93
θ_{c-3}	76.37	78.43	77.90	70.00	74.83	74.67	97.73	98.90	98.83

Table 6: Test accuracy (%) of CP-CFL (*SimCLR*) with different classifier head structures.

gap between θ_c and other, larger classifier heads slowly diminishes as the number of rounds increases. (The results for the STL-10-to-MNIST task are provided in Appendix A.)

4.5. Encoder in the local training step

In CP-CFL, clients train the entire model (both the encoder θ_f and the classifier head θ_c) at each global round $t > T_c$. However, it might be possible to freeze the pre-trained encoder θ_f in all training rounds without a significant impact on downstream CFL performance. Since θ_f has already been pre-trained on the unlabeled data, we can lighten the computational load on the clients by reducing the frequency of local training for θ_f . Therefore, we study three different settings in which the encoder θ_f is trained for zero, one, and three (default) local epochs. The classifier head θ_c is still trained for three local epochs in all settings. The experiments are conducted on the STL-10-to-STL-10 task with different classifier head structures.

From the results in Table 7, it is clear that training the encoder θ_f for three local epochs ($E_{\ell,f} = 3$) generally achieves the best performance or performance comparable to that of other settings. When used together with large classifier heads such as θ_{c-1} , θ_{c-2} , and θ_{c-3} , freezing the encoder ($E_{\ell,f} = 0$) still performs comparably to $E_{\ell,f} = 3$ (with only minor accuracy drops of 1.06%, 0.03%, and 1.13%, respectively). This suggests that it may be possible to reduce the computational cost on the clients by freezing the encoder while using a large classifier head, with minimal accuracy reduction. (It is worth noting that the number of trainable parameters in a large classifier head, $\theta_{c-3} \approx 0.1$ M, is still significantly lower than that of the entire CNN encoder, $\theta_f \approx 2.4$ M.) Interestingly, training the encoder for a single local epoch ($E_{\ell,f} = 1$) results in worse performance than the frozen encoder ($E_{\ell,f} = 0$) for

Local training frequency of θ_f	Classifier head structure			
	θ_c	θ_{c-1}	θ_{c-2}	θ_{c-3}
$E_{\ell,f} = 3$ (default)	76.80	78.23	77.60	77.90
$E_{\ell,f} = 1$	76.90	74.93	74.80	73.97
$E_{\ell,f} = 0$ (frozen)	73.67	77.17	77.57	76.77
$E_{\ell,f} = 3$ (global)	76.27	79.10	77.50	78.03

Table 7: Test accuracy (%) of CP-CFL (*SimCLR*) for the STL-10-to-STL-10 task ($T = 100$). $E_{\ell,f}$ is the number of local epochs for training the encoder θ_f .

Pre-training	Accuracy (%)			
	$T = 25$	$T = 50$	$T = 75$	$T = 100$
Contrastive	70.13	76.50	77.17	77.33
Supervised	72.83	76.80	77.13	77.13

Table 8: Test accuracy (%) for the CIFAR-100-to-CIFAR-10 task using supervised and contrastive encoder pre-training settings.

θ_{c-1} , θ_{c-2} , and θ_{c-3} . We believe that inadequate training disrupts the pre-trained parameters in the encoder, hurting the compatibility between the encoder and the classifier head.

Furthermore, we study the global encoder setting, where all encoders from different client clusters are aggregated into a single global encoder. In this setting, all cluster models share the common global encoder, and only the classifier head is cluster specific. As shown in Table 7, the global encoder setting achieves promising results and presents an alternative to conventional CFL techniques where cluster models do not share the parameters.

4.6. Supervised versus contrastive pre-training

CP-CFL pre-trains the encoder θ_f on unlabeled data using contrastive learning. It is crucial to determine the potential performance drop caused by use of unlabeled data instead of labeled data for pre-training. To this end, we pre-train θ_f using both supervised and contrastive settings and compare their performance in the downstream CFL task. We use the CIFAR-100 dataset (Krizhevsky et al., 2009) for centralized pre-training and the CIFAR-10 dataset to construct the private data of clients. In the supervised pre-training setting, a classification model $\theta = (\theta_f, \theta_c)$ is trained with cross-entropy loss on the CIFAR-100 dataset, and then we keep only the encoder θ_f . In the contrastive learning setting, we use SimCLR to pre-train the encoder on the CIFAR-100 dataset without using the labels.

Table 8 compares the downstream CFL performance between different encoder pre-training settings. Supervised pre-training outperforms contrastive pre-training when T is small (i.e., $T \approx 25$). However, with a higher number of global rounds (i.e., $T \geq 50$), the performance gap between the two pre-training methods becomes negligible. Therefore, we can conclude that the encoder pre-trained with contrastive learning performs effectively in the downstream CFL task.

4.7. Number of clusters

CP-CFL does not explicitly compute the number of clusters N , since it is primarily determined by available system resources in practice. Although we set the number of clusters N to 3 by default in our experiments, it would be interesting to see how N affects the performance of CP-CFL. Therefore, we conduct experiments by adjusting the value of N between 1 and 5 on the STL-10-to-STL-10 task.

The results of CP-CFL (*SimCLR*) with different numbers of clusters are reported in Table 9. Setting the number of clusters N to 3 gives the highest performance, since it matches the underlying cluster structure. However, use of a higher number of clusters, such as four or five, still yields better results than a lower number of clusters, such as one or

Number of clusters (N)	Accuracy (%)
1	59.33
2	65.10
3	76.80
4	72.80
5	72.97

Table 9: Test accuracy (%) of CP-CFL (*SimCLR*) with different numbers of clusters for the STL-10-to-STL-10 task ($T = 100$). The best result is shown in bold.

two. In general, increasing the number of clusters N means improving the personalized performance at the expense of more system resources. Therefore, we can adjust N to strike a balance between performance gain and resource expense.

4.8. Number of clients

In this section, we investigate the impact of the number of clients on CP-CFL using the STL-10-to-STL-10 task. Specifically, we consider four different settings, where the number of clients is set to 15, 30, 60 (default), and 90, respectively. In each setting, the clients are equally divided into three groups with different underlying data distributions, as described in Section 4.1. We use the default hyperparameters, with the exception of setting $T = 40$ and $T = 60$ for the 15-client and 30-client settings, respectively.

Since the number of samples per client is fixed (i.e., 50), increasing the number of clients results in more training data. As presented in Table 10, the performance of CP-CFL (*SimCLR*) improves as the number of clients increases, a trend that is also observed in other approaches. While there is no significant difference in the final performance between IFCA (*None*) and IFCA (*FedAvg*) in Table 10, Figure 7 reveals that IFCA (*FedAvg*) exhibits a performance boost in early training rounds. This trend holds for all approaches involving pre-training, i.e., FedAvg (*SimCLR*), IFCA (*FedAvg*), and CP-CFL (*SimCLR*).

	Pre-training	Number of clients							
		15		30		60		90	
		$T = 40$	Best	$T = 60$	Best	$T = 100$	Best	$T = 100$	Best
FedAvg	None	34.27	36.13	36.93	37.07	44.13	50.73	46.58	46.78
	SimCLR	47.47	47.47	54.00	54.20	59.33	59.83	57.71	59.04
IFCA	None	57.60	58.53	62.67	62.67	67.80	67.80	70.80	70.93
	FedAvg	56.67	59.87	62.60	64.33	70.90	72.07	68.95	70.36
CP-CFL	SimCLR	70.67	70.67	73.87	74.93	76.80	78.33	79.80	80.27

Table 10: Test accuracy (%) of CP-CFL (*SimCLR*) with different numbers of clients for the STL-10-to-STL-10 task. We report both the final results and the best results.

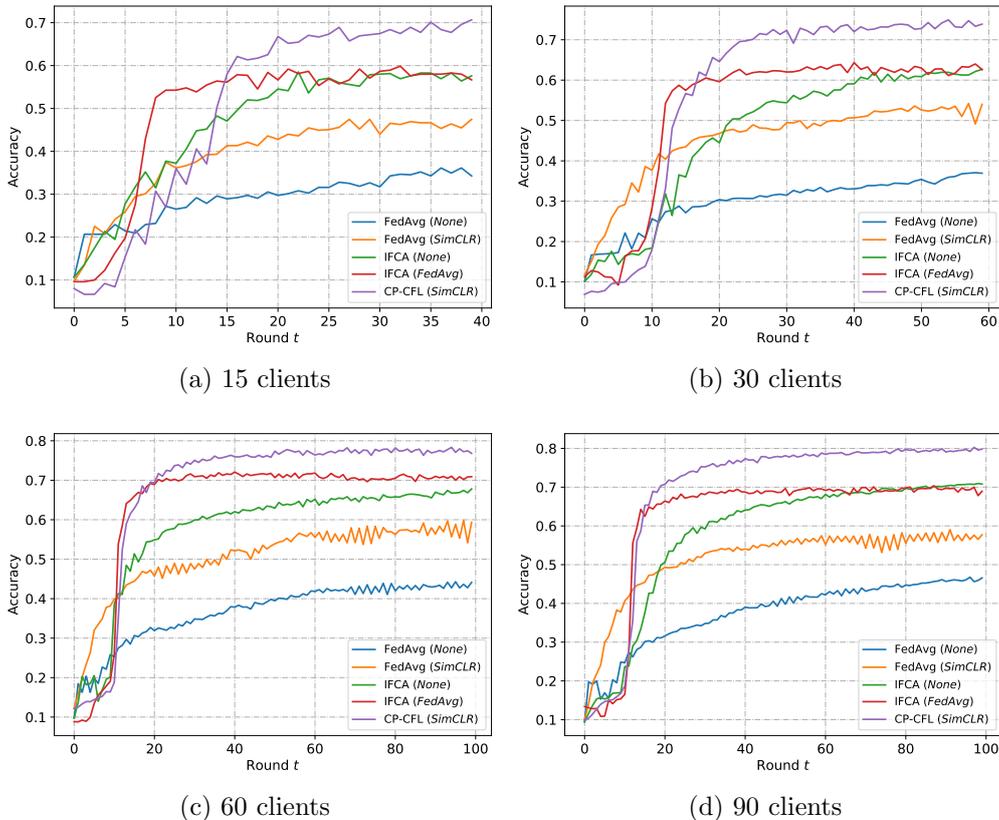


Figure 7: Test accuracy curves obtained with different numbers of clients on the STL-10-to-STL-10 task.

4.9. Statistical significance testing

To ensure the statistical significance of our study, we conduct five independent experiments on the STL-10-to-STL-10 task. For each experiment, we resample the private and testing data of clients according to the settings described in Section 4.1. Table 11 reports the evaluation results of all five experiments, as well as their mean and standard deviation. Our comparison with baseline approaches demonstrates that CP-CFL (*SimCLR*) maintains good results across all trials, and outperforms IFCA (*None*) by 10.90% and IFCA (*FedAvg*) by 8.35% on average. Additional metrics are provided in Appendix B.

4.10. ResNet-18 encoder

In addition to our default CNN encoder, we also evaluate the performance of CP-CFL using the ResNet-18 encoder (He et al., 2016) on the STL-10-to-STL-10 task. This experiment aims to examine the generalization ability of CP-CFL with regard to different encoder architectures other than our default CNN encoder, and not necessarily to obtain the best performance. Therefore, we use default values for most hyperparameters except for a few that we specify below. The number of epochs for contrastive pre-training is set to 500 with a batch size of 250. The output dimension of the ResNet-18 encoder is 512, and we use θ_{c-1} as the classifier head. (We empirically found that while using the ResNet-18 encoder with our

	Pre-training	Accuracy (%)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	42.73	41.53	43.67	43.03	42.67	42.73 \pm 0.69
	SimCLR	58.30	52.60	56.10	55.40	57.50	55.98 \pm 1.97
IFCA	None	66.17	64.37	66.43	66.17	65.80	65.79 \pm 0.74
	FedAvg	69.57	66.53	69.30	68.20	69.43	68.61 \pm 1.14
CP-CFL	SimCLR	76.13	76.47	77.43	76.40	77.00	76.69 \pm 0.47

Table 11: Test accuracy (%) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). SD, standard deviation.

default classifier head θ_c , which contains a single output layer, FedAvg (*SimCLR*) often gets stuck in local optima, resulting in low performance; herefore, we use θ_{c-1} for all approaches.) We conduct five independent experiments, and Table 12 reports the best results achieved by each approach. We report the best results in Table 12 since the performance of FedAvg baselines can fluctuate. Additional results for $T = 100$ are provided in Table C.20 in Appendix C. We also conduct a single experiment using the ResNet-50 encoder (He et al., 2016); the results are shown in Table D.21 in Appendix D. While the ResNet-50 encoder is commonly used for many image processing tasks, the ResNet-18 encoder can be more practical for FL environments, where client devices often have limited computing resources for locally training the ResNet-50 encoder.

5. Discussion and future work

The influence of the amount of unlabeled pre-training data on the performance of the encoder could be a topic of interest. In practice, it is often the responsibility of the service provider or the server to gather an adequate amount of unlabeled data for pre-training the encoder. The server needs to determine the quantity of pre-training data based on various factors, such as the type of service being provided and the size of the model deployed. The server may also keep a small set of proxy test data to evaluate the performance of the pre-trained encoder and determine if the collected unlabeled data are sufficient in terms of both quality and size. While a larger quantity of pre-training data would generally result in a

	Pre-training	Accuracy (%)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	37.37	38.40	37.77	34.83	39.40	37.55 \pm 1.52
	SimCLR	45.77	39.77	44.97	39.33	39.70	41.91 \pm 2.84
IFCA	None	64.10	64.67	64.80	65.53	63.93	64.61 \pm 0.57
	FedAvg	65.03	65.27	64.87	64.00	65.40	64.91 \pm 0.49
CP-CFL	SimCLR	67.90	68.57	68.30	68.90	66.53	68.04 \pm 0.82

Table 12: Test accuracy (%) obtained with the ResNet-18 encoder on the STL-10-to-STL-10 task (best). SD, standard deviation.

higher-quality encoder, it would also require a longer training time and more computing resources. Therefore, we leave the investigation of the trilemma between the amount of pre-training data, the required training resources, and the resulting encoder quality to future studies.

6. Conclusions

FL is essential for incorporating edge-generated data into intelligent systems without violating privacy regulations. However, non-IID data at the edge pose a fundamental challenge that limits FL from achieving the same level of performance as centralized training. To address this issue, we proposed CP-CFL, which combines contrastive encoder pre-training and client clustering to alleviate the performance drop caused by non-IID data in a heterogeneous FL environment. We conducted extensive experiments on CP-CFL in various settings to demonstrate its superior performance over the baseline approaches. Furthermore, we explored efficient ways to deploy CP-CFL, while also providing various ablation studies to validate its effectiveness. Overall, our study contributes to a better understanding of how contrastive encoder pre-training and client clustering can jointly improve the performance of FL on non-IID data.

References

- Bachman, P., Hjelm, R.D., Buchwalter, W., 2019. Learning representations by maximizing mutual information across views. *Advances in Neural Information Processing Systems* 32, 15535–15545.
- Baevski, A., Zhou, Y., Mohamed, A., Auli, M., 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems* 33, 12449–12460.
- van Berlo, B., Saeed, A., Ozcelebi, T., 2020. Towards federated unsupervised representation learning, in: *Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking*, pp. 31–36. doi:10.1145/3378679.3394530.
- Briggs, C., Fan, Z., Andras, P., 2020. Federated learning with hierarchical clustering of local updates to improve training on non-IID data, in: *2020 International Joint Conference on Neural Networks (IJCNN)*, IEEE. pp. 1–9. doi:10.1109/IJCNN48605.2020.9207469.
- Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., Joulin, A., 2020. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems* 33, 9912–9924.
- Chen, T., Kornblith, S., Norouzi, M., Hinton, G., 2020. A simple framework for contrastive learning of visual representations, in: *Proceedings of the 37th International Conference on Machine Learning*, PMLR. pp. 1597–1607.
- Chen, X., He, K., 2021. Exploring simple siamese representation learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15750–15758. doi:10.1109/CVPR46437.2021.01549.
- Coates, A., Ng, A., Lee, H., 2011. An analysis of single-layer networks in unsupervised feature learning, in: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, PMLR. pp. 215–223.
- Cohen, G., Afshar, S., Tapson, J., Van Schaik, A., 2017. EMNIST: Extending MNIST to handwritten letters, in: *2017 International Joint Conference on Neural Networks (IJCNN)*, IEEE. pp. 2921–2926. doi:10.1109/IJCNN.2017.7966217.
- Dave, I., Gupta, R., Rizve, M.N., Shah, M., 2022. TCLR: Temporal contrastive learning for video representation. *Computer Vision and Image Understanding* 219, 103406. doi:10.1016/j.cviu.2022.103406.

- Dennis, D.K., Li, T., Smith, V., 2021. Heterogeneity for the win: One-shot federated clustering, in: Proceedings of the 38th International Conference on Machine Learning, PMLR. pp. 2611–2620.
- Dinh, C.T., Tran, N.H., Nguyen, M.N., Hong, C.S., Bao, W., Zomaya, A.Y., Gramoli, V., 2020. Federated learning over wireless networks: Convergence analysis and resource allocation. *IEEE/ACM Transactions on Networking* 29, 398–409. doi:10.1109/TNET.2020.3035770.
- Duan, M., Liu, D., Ji, X., Liu, R., Liang, L., Chen, X., Tan, Y., 2021. Fedgroup: Efficient federated learning via decomposed similarity-based clustering, in: 2021 IEEE International Conference on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social computing & Networking (ISPA/BDCLOUD/SocialCom/SustainCom), IEEE. pp. 228–237. doi:10.1109/ISPA-BDCLOUD-SocialCom-SustainCom52081.2021.00042.
- Gao, T., Yao, X., Chen, D., 2021. SimCSE: Simple contrastive learning of sentence embeddings. arXiv:2104.08821. doi:10.48550/arXiv.2104.08821.
- Ghosh, A., Chung, J., Yin, D., Ramchandran, K., 2020. An efficient framework for clustered federated learning. arXiv:2006.04088. doi:10.48550/arXiv.2006.04088.
- Ghosh, A., Hong, J., Yin, D., Ramchandran, K., 2019. Robust federated learning in a heterogeneous environment. arXiv:1906.06629. doi:10.48550/arXiv.1906.06629.
- Gidaris, S., Singh, P., Komodakis, N., 2018. Unsupervised representation learning by predicting image rotations. arXiv:1803.07728. doi:10.48550/arXiv.1803.07728.
- Giorgi, J., Nitski, O., Wang, B., Bader, G., 2020. DeCLUTR: Deep contrastive learning for unsupervised textual representations. arXiv:2006.03659. doi:10.48550/arXiv.2006.03659.
- Grill, J.B., Strub, F., Althé, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al., 2020. Bootstrap your own latent - a new approach to self-supervised learning. *Advances in Neural Information Processing Systems* 33, 21271–21284.
- Hard, A., Rao, K., Mathews, R., Ramaswamy, S., Beaufays, F., Augenstein, S., Eichner, H., Kiddon, C., Ramage, D., 2018. Federated learning for mobile keyboard prediction. arXiv:1811.03604. doi:10.48550/arXiv.1811.03604.
- He, K., Fan, H., Wu, Y., Xie, S., Girshick, R., 2020. Momentum contrast for unsupervised visual representation learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9729–9738. doi:10.1109/cvpr42600.2020.00975.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778. doi:10.1109/CVPR.2016.90.
- Jamali-Rad, H., Abdizadeh, M., Singh, A., 2022. Federated learning with taskonomy for non-IID data. *IEEE Transactions on Neural Networks and Learning Systems*, 1–12. doi:10.1109/TNNLS.2022.3152581.
- Kaissis, G., Makowski, M.R., Rückert, D., Braren, R.F., 2020. Secure, privacy-preserving and federated machine learning in medical imaging. *Nature Machine Intelligence* 2, 305–311. doi:10.1038/s42256-020-0186-1.
- Kassem, H., Alapatt, D., Mascagni, P., Consortium AI4SafeChole, Karargyris, A., Padoy, N., 2022. Federated cycling (FedCy): Semi-supervised federated learning of surgical phases. *IEEE Transactions on Medical Imaging*, 1–1 doi:10.1109/TMI.2022.3222126.
- Kim, Y., Al Hakim, E., Haraldson, J., Eriksson, H., da Silva, J.M.B., Fischione, C., 2021. Dynamic clustering in federated learning, in: ICC 2021-IEEE International Conference on Communications, IEEE. pp. 1–6. doi:10.1109/ICC42927.2021.9500877.
- Kingma, D.P., Welling, M., 2013. Auto-encoding variational bayes. arXiv:1312.6114. doi:10.48550/arXiv.1312.6114.
- Konečný, J., McMahan, H.B., Ramage, D., Richtárik, P., 2016a. Federated optimization: Distributed machine learning for on-device intelligence. arXiv:1610.02527. doi:10.48550/arXiv.1610.02527.
- Konečný, J., McMahan, H.B., Yu, F.X., Richtárik, P., Suresh, A.T., Bacon, D., 2016b. Federated learning: Strategies for improving communication efficiency. arXiv:1610.05492. doi:10.48550/arXiv.1610.05492.
- Krizhevsky, A., Hinton, G., et al., 2009. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009 URL: <https://www.cs.toronto.edu/~kriz/>

learning-features-2009-TR.pdf.

- LeCun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86, 2278–2324. doi:10.1109/5.726791.
- Li, C., Li, G., Varshney, P.K., 2021a. Federated learning with soft clustering. *IEEE Internet of Things Journal* 9, 7773–7782. doi:10.1109/JIOT.2021.3113927.
- Li, Q., He, B., Song, D., 2021b. Model-contrastive federated learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10713–10722. doi:10.1109/CVPR46437.2021.01057.
- Li, T., Sahu, A.K., Talwalkar, A.S., Smith, V., 2020. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine* 37, 50–60. doi:10.1109/MSP.2020.2975749.
- Liang, X., Lin, Y., Fu, H., Zhu, L., Li, X., 2022. RSCFed: random sampling consensus federated semi-supervised learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10154–10163. doi:10.1109/CVPR52688.2022.00991.
- Long, G., Xie, M., Shen, T., Zhou, T., Wang, X., Jiang, J., 2023. Multi-center federated learning: Clients clustering for better personalization. *World Wide Web* 26, 481–500. doi:10.1007/s11280-022-01046-x.
- Long, Z., Wang, J., Wang, Y., Xiao, H., Ma, F., 2021. FedCon: A contrastive framework for federated semi-supervised learning. arXiv:2109.04533. doi:10.48550/arXiv.2109.04533.
- Makhija, D., Ho, N., Ghosh, J., 2022. Federated self-supervised learning for heterogeneous clients. arXiv:2205.12493. doi:10.48550/arXiv.2205.12493.
- Manocha, P., Jin, Z., Zhang, R., Finkelstein, A., 2021. CDPAM: Contrastive learning for perceptual audio similarity, in: *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE. pp. 196–200. doi:10.1109/ICASSP39728.2021.9413711.
- Mansour, Y., Mohri, M., Ro, J., Suresh, A.T., 2020. Three approaches for personalization with applications to federated learning. arXiv:2002.10619. doi:10.48550/arXiv.2002.10619.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A., 2017. Communication-efficient learning of deep networks from decentralized data, in: *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, PMLR. pp. 1273–1282.
- Misra, I., van der Maaten, L., 2020. Self-supervised learning of pretext-invariant representations, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6707–6717. doi:10.1109/CVPR42600.2020.00674.
- Nguyen, M.N., Tran, N.H., Tun, Y.K., Han, Z., Hong, C.S., 2020. Toward multiple federated learning services resource sharing in mobile edge networks. arXiv:2011.12469. doi:10.48550/arXiv.2011.12469.
- Noroozi, M., Favaro, P., 2016. Unsupervised learning of visual representations by solving jigsaw puzzles, in: *European Conference on Computer Vision*, Springer. pp. 69–84. doi:10.1007/978-3-319-46466-4_5.
- Ouyang, X., Xie, Z., Zhou, J., Xing, G., Huang, J., 2022. ClusterFL: A clustering-based federated learning system for human activity recognition. *ACM Transactions on Sensor Networks* 19, 1–32. doi:10.1145/3554980.
- Pan, T., Song, Y., Yang, T., Jiang, W., Liu, W., 2021. VideoMoCo: Contrastive video representation learning with temporally adversarial examples, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11205–11214. doi:10.1109/CVPR46437.2021.01105.
- Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A., 2016. Context encoders: Feature learning by inpainting, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2536–2544. doi:10.1109/CVPR.2016.278.
- Qian, R., Meng, T., Gong, B., Yang, M.H., Wang, H., Belongie, S., Cui, Y., 2021. Spatiotemporal contrastive video representation learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6964–6974. doi:10.1109/CVPR46437.2021.00689.
- Qu, L., Zhou, Y., Liang, P.P., Xia, Y., Wang, F., Adeli, E., Fei-Fei, L., Rubin, D., 2022. Rethinking architecture design for tackling data heterogeneity in federated learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10061–10071. doi:10.1109/CVPR52688.2022.00982.
- Ramaswamy, S., Mathews, R., Rao, K., Beaufays, F., 2019. Federated learning for emoji prediction in a

- mobile keyboard. arXiv:1906.04329. doi:10.48550/arXiv.1906.04329.
- Saeed, A., Grangier, D., Zeghidour, N., 2021. Contrastive learning of general-purpose audio representations, in: 2021-2021 International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE. pp. 3875–3879. doi:10.1109/ICASSP39728.2021.9413528.
- Sattler, F., Korjakow, T., Rischke, R., Samek, W., 2021a. FedAUX: Leveraging unlabeled auxiliary data in federated learning. *IEEE Transactions on Neural Networks and Learning Systems* . doi:10.1109/TNNLS.2021.3129371.
- Sattler, F., Müller, K.R., Wiegand, T., Samek, W., 2020a. On the byzantine robustness of clustered federated learning, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE. pp. 8861–8865. doi:10.1109/ICASSP40776.2020.9054676.
- Sattler, F., Müller, K.R., Samek, W., 2021b. Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints. *IEEE Transactions on Neural Networks and Learning Systems* 32, 3710–3722. doi:10.1109/TNNLS.2020.3015958.
- Sattler, F., Wiedemann, S., Müller, K.R., Samek, W., 2020b. Robust and communication-efficient federated learning from non-i.i.d. data. *IEEE Transactions on Neural Networks and Learning Systems* 31, 3400–3413. doi:10.1109/TNNLS.2019.2944481.
- Shlezinger, N., Rini, S., Eldar, Y.C., 2020. The communication-aware clustered federated learning problem, in: 2020 IEEE International Symposium on Information Theory (ISIT), pp. 2610–2615. doi:10.1109/ISIT44484.2020.9174245.
- Son, H.M., Kim, M.H., Chung, T.M., 2022. Comparisons where it matters: Using layer-wise regularization to improve federated learning on heterogeneous data. *Applied Sciences* 12, 9943. doi:10.3390/app12199943.
- Thwal, C.M., Thar, K., Tun, Y.L., Hong, C.S., 2021. Attention on personalized clinical decision support system: Federated learning approach, in: 2021 IEEE International Conference on Big Data and Smart Computing (BigComp), pp. 141–147. doi:10.1109/BigComp51126.2021.00035.
- Tian, P., Liao, W., Yu, W., Blasch, E., 2022. WSCC: A weight-similarity-based client clustering approach for non-IID federated learning. *IEEE Internet of Things Journal* 9, 20243–20256. doi:10.1109/JIOT.2022.3175149.
- Tokusumi, 2020. keras-flops: FLOPs calculator with tf.profiler for neural network architecture written in tensorflow 2.2+ (tf.keras). URL: <https://github.com/tokusumi/keras-flops>.
- Tun, Y.L., Thwal, C.M., Park, Y.M., Park, S.B., Hong, C.S., 2023. Federated learning with intermediate representation regularization, in: 2023 IEEE International Conference on Big data and Smart Computing (BigComp), pp. 56–63. doi:10.1109/BigComp57234.2023.00017.
- van den Oord, A., Li, Y., Vinyals, O., 2018. Representation learning with contrastive predictive coding. arXiv e-prints, arXiv:1807.03748. doi:10.48550/arXiv.1807.03748.
- Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.A., 2008. Extracting and composing robust features with denoising autoencoders, in: Proceedings of the 25th International Conference on Machine Learning, pp. 1096–1103. doi:10.1145/1390156.1390294.
- Voigt, P., Bussche, A.v.d., 2017. The EU general data protection regulation (GDPR): A practical guide. 1st ed., Springer Publishing Company, Incorporated. doi:10.1007/978-3-319-57959-7.
- Wu, Y., Zeng, D., Wang, Z., Sheng, Y., Yang, L., James, A.J., Shi, Y., Hu, J., 2021. Federated contrastive learning for dermatological disease diagnosis via on-device learning, in: 2021 IEEE/ACM International Conference on Computer Aided Design (ICCAD), IEEE. pp. 1–7. doi:10.1109/ICCAD51958.2021.9643454.
- Wu, Z., Wang, S., Gu, J., Khabsa, M., Sun, F., Ma, H., 2020. CLEAR: Contrastive learning for sentence representation. arXiv:2012.15466. doi:10.48550/arXiv.2012.15466.
- Yan, Z., Wicaksana, J., Wang, Z., Yang, X., Cheng, K.T., 2021. Variation-aware federated learning with multi-source decentralized medical image data. *IEEE Journal of Biomedical and Health Informatics* 25, 2615–2628. doi:10.1109/JBHI.2020.3040015.
- Yu, F., Zhang, W., Qin, Z., Xu, Z., Wang, D., Liu, C., Tian, Z., Chen, X., 2020. Heterogeneous federated learning. arXiv:2008.06767. doi:10.48550/arXiv.2008.06767.
- Zbontar, J., Jing, L., Misra, I., LeCun, Y., Deny, S., 2021. Barlow twins: Self-supervised learning via

- redundancy reduction, in: Proceedings of the 38th International Conference on Machine Learning, PMLR. pp. 12310–12320.
- Zhang, F., Kuang, K., You, Z., Shen, T., Xiao, J., Zhang, Y., Wu, C., Zhuang, Y., Li, X., 2020. Federated unsupervised representation learning. arXiv:2010.08982. doi:10.48550/arXiv.2010.08982.
- Zhao, J., Li, R., Wang, H., Xu, Z., 2021. HotFed: Hot start through self-supervised learning in federated learning, in: 2021 IEEE 23rd International Conference on High Performance Computing & Communications; 7th International Conference on Data Science & Systems; 19th International Conference on Smart City; 7th International Conference on Dependability in Sensor, Cloud & Big Data Systems & Application (HPCC/DSS/SmartCity/DependSys), pp. 149–156. doi:10.1109/HPCC-DSS-SmartCity-DependSys53884.2021.00046.
- Zhao, Y., Li, M., Lai, L., Suda, N., Civin, D., Chandra, V., 2018. Federated learning with non-IID data. arXiv:1806.00582. doi:10.48550/arXiv.1806.00582.
- Zhuang, W., Gan, X., Wen, Y., Zhang, S., Yi, S., 2021. Collaborative unsupervised visual representation learning from decentralized data, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 4912–4921. doi:10.1109/ICCV48922.2021.00487.
- Zhuang, W., Wen, Y., Zhang, S., 2022. Divergence-aware federated self-supervised learning. arXiv:2204.04385. doi:10.48550/arXiv.2204.04385.

Appendix A. Learning capacity of the classifier head

The results for CP-CFL (*SimCLR*) with different classifier head structures on the STL-10-to-MNIST task are given in Table A.13.

	STL-10-to-MNIST		
	$T = 25$	$T = 75$	$T = 100$
θ_c	93.70	97.63	98.40
θ_{c-1}	94.20	98.27	98.13
θ_{c-2}	93.53	98.53	98.63
θ_{c-3}	94.20	98.13	98.20

Table A.13: Test accuracy (%) of CP-CFL (*SimCLR*) with different classifier head structures on the STL-10-to-MNIST task.

Appendix B. Additional metrics

F1 score and AUROC for statistical significance testing reported in Section 4.9 are given in Tables B.14–B.19.

	Pre-training	F1 score (macro)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	0.3570	0.3451	0.3566	0.3315	0.3415	0.3463 ± 0.0096
	SimCLR	0.5484	0.4746	0.5177	0.5118	0.5304	0.5166 ± 0.0245
IFCA	None	0.6701	0.6479	0.6684	0.6685	0.6668	0.6643 ± 0.0083
	FedAvg	0.7002	0.6681	0.6984	0.6876	0.7006	0.6910 ± 0.0124
CP-CFL	SimCLR	0.7640	0.7593	0.7757	0.7630	0.7698	0.7664 ± 0.0057

Table B.14: F1 score (macro) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). SD, standard deviation.

	Pre-training	F1 score (weighted)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	0.3817	0.3662	0.3859	0.3633	0.3692	0.3733 ± 0.0089
	SimCLR	0.5641	0.4933	0.5311	0.5306	0.5440	0.5326 ± 0.0231
IFCA	None	0.6606	0.6404	0.6638	0.6602	0.6577	0.6565 ± 0.0083
	FedAvg	0.6942	0.6624	0.6922	0.6807	0.6925	0.6844 ± 0.0120
CP-CFL	SimCLR	0.7575	0.7560	0.7714	0.7601	0.7655	0.7621 ± 0.0056

Table B.15: F1 score (weighted) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). SD, standard deviation.

	Pre-training	AUROC (OvR macro)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	0.8321	0.8278	0.8355	0.8379	0.8306	0.8328 ± 0.0036
	SimCLR	0.9033	0.8928	0.8990	0.9019	0.9080	0.9010 ± 0.0050
IFCA	None	0.9544	0.9550	0.9553	0.9573	0.9542	0.9552 ± 0.0011
	FedAvg	0.9614	0.9571	0.9592	0.9593	0.9597	0.9593 ± 0.0014
CP-CFL	SimCLR	0.9756	0.9753	0.9766	0.9761	0.9759	0.9759 ± 0.0004

Table B.16: AUROC (OvR macro) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). OvR, one-vs-rest. SD, standard deviation.

	Pre-training	AUROC (OvR weighted)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	0.8291	0.8269	0.8341	0.8366	0.8300	0.8313 ± 0.0035
	SimCLR	0.9040	0.8930	0.8987	0.9020	0.9079	0.9011 ± 0.0050
IFCA	None	0.9465	0.9472	0.9481	0.9501	0.9461	0.9476 ± 0.0014
	FedAvg	0.9549	0.9500	0.9521	0.9525	0.9526	0.9524 ± 0.0015
CP-CFL	SimCLR	0.9718	0.9720	0.9734	0.9732	0.9723	0.9725 ± 0.0006

Table B.17: AUROC (OvR weighted) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). OvR, one-vs-rest. SD, standard deviation.

	Pre-training	AUROC (OvO macro)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	0.8361	0.8308	0.8395	0.8409	0.8339	0.8363 ± 0.0037
	SimCLR	0.9062	0.8937	0.9022	0.9037	0.9097	0.9031 ± 0.0054
IFCA	None	0.9584	0.9586	0.9585	0.9605	0.9578	0.9587 ± 0.0009
	FedAvg	0.9646	0.9605	0.9624	0.9625	0.9631	0.9626 ± 0.0013
CP-CFL	SimCLR	0.9778	0.9774	0.9786	0.9782	0.9779	0.9780 ± 0.0004

Table B.18: AUROC (OvO macro) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). OvO, one-vs-one. SD, standard deviation.

	Pre-training	AUROC (OvO weighted)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	0.8325	0.8285	0.8364	0.8384	0.8316	0.8335 ± 0.0035
	SimCLR	0.9046	0.8929	0.9000	0.9025	0.9084	0.9017 ± 0.0052
IFCA	None	0.9524	0.9529	0.9531	0.9552	0.9520	0.9531 ± 0.0011
	FedAvg	0.9598	0.9553	0.9572	0.9575	0.9580	0.9576 ± 0.0014
CP-CFL	SimCLR	0.9747	0.9746	0.9758	0.9755	0.9750	0.9751 ± 0.0005

Table B.19: AUROC (OvO weighted) on multiple trials of the STL-10-to-STL-10 task ($T = 100$). OvO, one-vs-one. SD, standard deviation.

Appendix C. ResNet-18 encoder

Table C.20 presents the evaluation results for the STL-10-to-STL-10 task obtained with the ResNet-18 encoder at $T = 100$.

	Pre-training	Accuracy (%)					Mean \pm SD
		Trial 1	Trial 2	Trial 3	Trial 4	Trial 5	
FedAvg	None	24.67	34.00	25.10	27.53	35.70	29.40 ± 4.59
	SimCLR	34.83	29.40	33.83	30.73	29.67	31.69 ± 2.22
IFCA	None	63.73	61.90	62.50	65.53	63.27	63.39 ± 1.24
	FedAvg	63.63	63.17	64.80	62.90	61.70	63.24 ± 1.01
CP-CFL	SimCLR	67.90	66.90	68.30	66.30	66.43	67.17 ± 0.80

Table C.20: Test accuracy (%) obtained with the ResNet-18 encoder on the STL-10-to-STL-10 task ($T = 100$). SD, standard deviation.

Appendix D. ResNet-50 encoder

Table D.21 presents the evaluation results obtained with the ResNet-50 (He et al., 2016) encoder on the STL-10-to-STL-10 task. For the contrastive pre-training step, we set the number of epochs to 500 and the batch size to 250. We use a projection head θ_g containing

a ReLU-activated dense layer with 1024 neurons, followed by an another dense layer with 2048 neurons. The output dimension of the ResNet-50 encoder is 2048, and we use θ_{c-2} as the classifier head. Using the ResNet-50 encoder, CP-CFL can still outperform the baseline approaches as shown in Table D.21.

	Pre-training	Accuracy (%)	
		$T=100$	Best
FedAvg	None	45.10	45.80
	SimCLR	39.07	46.13
IFCA	None	67.93	67.93
	FedAvg	69.57	71.10
CP-CFL	SimCLR	72.93	72.93

Table D.21: Test accuracy (%) obtained with the ResNet-50 encoder on the STL-10-to-STL-10 task.

Appendix E. Pre-training time

We performed all of our experiments on a single RTX 3080 GPU with 10 GB of memory. Table E.22 compares the training time per epoch for SimCLR pre-training using different encoder architectures on the STL-10 unlabeled portion. We also calculate the computational cost of each encoder in flops using the `keras-flops` (Tokusumi, 2020) package. According to Table E.22, pre-training with the CNN encoder takes approximately 2 h for 300 epochs, while pre-training with the ResNet-18 and ResNet-50 encoders takes around 8 and 18 h, respectively, for 500 epochs on our hardware.

Encoder	Time per epoch (s)	Flops (G)
CNN	≈ 23	0.17
ResNet-18	≈ 60	0.67
ResNet-50	≈ 134	1.42

Table E.22: Training time per epoch for SimCLR pre-training on the unlabeled data portion of the STL-10 dataset. We also report the computational cost of each encoder in flops.