

# FedAC: A Adaptive Clustered Federated Learning Framework for Heterogeneous Data

Yuxin Zhang<sup>1</sup>, Haoyu Chen<sup>1</sup>, Zheng Lin<sup>2</sup>, Zhe Chen<sup>1</sup>, and Jin Zhao<sup>1</sup>✉

<sup>1</sup> School of Computer Science, Fudan University, Shanghai 200438, China  
yuxinzhang22@m.fudan.edu.cn

<sup>2</sup> Department of Electrical and Electronic Engineering, The University of Hong Kong, Pok Fu Lam, Hong Kong.

**Abstract.** Clustered federated learning (CFL) is proposed to mitigate the performance deterioration stemming from data heterogeneity in federated learning (FL) by grouping similar clients for cluster-wise model training. However, current CFL methods struggle due to inadequate integration of global and intra-cluster knowledge and the absence of an efficient online model similarity metric, while treating the cluster count as a fixed hyperparameter limits flexibility and robustness. In this paper, we propose an adaptive CFL framework, named **FedAC**, which (1) efficiently integrates global knowledge into intra-cluster learning by decoupling neural networks and utilizing distinct aggregation methods for each submodule, significantly enhancing performance; (2) includes a cost-effective online model similarity metric based on dimensionality reduction; (3) incorporates a cluster number fine-tuning module for improved adaptability and scalability in complex, heterogeneous environments. Extensive experiments show that **FedAC** achieves superior empirical performance, increasing the test accuracy by around 1.82% and 12.67% on CIFAR-10 and CIFAR-100 datasets, respectively, under different non-IID settings compared to SOTA methods.

**Keywords:** Federated learning · Multi-task learning · Clustering · Non-IID data · Model similarity.

## 1 Introduction

Machine learning (ML) has rapidly advanced, finding extensive applications across industries such as intelligent agent, autonomous driving and energy management [1,2,3,4,5,6,7,8]. Its progress is attributed to diverse datasets, enhancing generalization and reducing overfitting [9,10]. However, in real-world scenarios, data may be distributed across mobile devices and the Internet of Things (IoT) [11]. Decentralized training of ML models becomes necessary because of privacy constraints [12,13] and network bandwidth limitations [14,15,16], preventing the transmission of these data to a central server for centralized training. Fueled by this realistic need, federated learning (FL) [17,18] was proposed, allowing multiple clients to collaboratively learn a global model without exchanging local data.

However, a significant practical challenge in FL lies in data heterogeneity. Clients may have non-IID data and diverse preferences, resulting in variations in the true risk at the local level, which is inconsistent with the existence of a global model suitable for all clients [19,20]. Even when considering a macroscopic perspective of empirical risk to formulate the global objective, local updates to client models may deviate to varying degrees. This presents challenges in convergence, potentially leading to suboptimal outcomes as the average global model drifts away [21].

Clustered federated learning (CFL) [22] employs multi-task learning [23] to mitigate data heterogeneity by grouping clients into clusters with higher internal *similarity* (more homogeneity). Despite extensive research, current CFL methods still have limitations: (1) These methods demonstrate isolated clusters, lacking the infusion of valuable global knowledge. Consequently, clients are unaware of potentially beneficial knowledge beyond their assigned clusters, resulting in overall suboptimal performance; (2) There is still a need for a computationally effective method to accurately measure the *online* similarity between models for clustering adjustments throughout the entire training process; (3) Furthermore, current methods treat the cluster count ( $K$ ) as a constant hyperparameter throughout the process, overlooking the challenge of manually setting the optimal value in complex heterogeneous scenarios. For example, when training a predictive text model with CFL for millions of global smartphone users, the complex and interrelated hidden contextual factors like age, location, occupation, etc., makes it impossible to pre-set the optimal number of clusters. Hence, there is an urgent requirement for an adaptive approach to fine-tune and determine its optimal value during training.

Motivated by the above challenges, this paper proposes **FedAC**, an adaptive CFL framework, comprising the following main components: (1) Decoupling neural networks into submodules and employing distinct aggregation methods, **FedAC** effectively integrates global knowledge into clusters. This allows clients to learn from both cluster-specific and global dimensions simultaneously, ensuring optimal performance by striking a balance; (2) By integrating a cosine model similarity metric following dimensionality reduction, **FedAC** effectively captures the online similarity of models at a computationally economical cost; (3) **FedAC** incorporates a module that dynamically fine-tunes the cluster count based on the current clustering status (inter-cluster and intra-cluster model distances), eliminating the necessity of pre-specifying a fixed cluster count. Instead, it autonomously fine-tunes the count during training to discover the optimal value, thereby enhancing the system’s robustness and flexibility. Extensive experiments demonstrate that **FedAC** outperforms SOTA methods in diverse heterogeneous scenarios. The effectiveness of each component is further validated through ablation experiments, showcasing the flexibility and robustness of **FedAC**.

The following summarizes our contributions.

- In this paper, we propose **FedAC**, an efficient and adaptive CFL framework designed to tackle complex non-IID scenarios within FL. To the best of our knowledge, **FedAC** is the first CFL framework to achieve outstanding per-

formance by introducing global knowledge for intra-cluster learning through the decoupling of neural networks.

- We develop an approach to online assessment of model similarity using dimensionality-reduced models. This allows the server to efficiently evaluate similarities among data-heterogeneous clients in a cost-effective manner, enhancing clustering effectiveness and system scalability.
- We present a design approach and implementation for adaptive adjustment of the total cluster count based on clustering status, enhancing the framework’s flexibility and robustness in addressing the challenge of manually setting it in complex, heterogeneous scenarios.
- We performed experiments on diverse datasets in complex heterogeneous scenarios to demonstrate the outstanding overall performance and adaptability of FedAC, surpassing SOTA methods.

The rest of this paper is organized as follows. In section 2, we review Heterogeneous FL and CFL approaches addressing the data heterogeneity issue. section 3 formulates the non-IID problem in CFL and outlines the optimization goal. section 4 introduces the framework and essential components of FedAC. Performance evaluations follow in section 5, and the paper concludes in section 6.

## 2 Related Work

### 2.1 Heterogeneous Federated Learning

The performance of FL is significantly affected by data heterogeneity, specifically, the non-identically and independently distributed (non-IID) nature of the data. The differences in data distribution among each client can lead to biased model updates, and directly averaging local updates (as in FedAvg [18]) may detrimentally impact the overall performance of the global model. FedPer [24] decouples models and aggregates them using diverse strategies. Per-FedAvg [25] blends FedAvg with model-agnostic meta-learning (MAML) for personalized models via fine-tuning, and this concept is extended in pFedMe [26] using Moreau envelopes. pFedHN [27] employs a hypernetwork to generate parameters for the personalized model of each client, and FedVF [28] generates personalized models through a two-stage training process.

CFL, as another extensively researched heterogeneous FL method, addresses non-IID issues by assuming that clients can be partitioned into several clusters. During the training process, the server needs to learn the appropriate cluster assignment for each client and simultaneously improve the cluster center models. An overview of CFL methods is presented in the following section.

### 2.2 Clustered Federated Learning

In [29], a hierarchical clustering framework is introduced for FL. It reduces computation and communication loads by using an agglomerative formulation. In

[22], hierarchical clustering is used as a postprocessing step in FL. However, the recursive bipartitioning framework incurs high costs, limiting feasibility for large-scale settings. FedGroup [30] utilizes a static client cluster methodology, initiating cold start protocols for new clients and adopting Euclidean distance of decomposed cosine similarity (EDC) for clustering. These methods cluster clients only once offline, and their performance heavily depends on the effectiveness of the clustering, exhibiting limited flexibility and robustness.

Online CFL methods address these issues by continuously re-clustering during training. In FeSEM [31], the expectation-maximization (EM) algorithm resolves the distance-based objective clustering issue by optimally pairing clients with their respective cluster centers. [32] introduced IFCA, which employs  $K$  global models distributed to all clients for local loss computation to perform clustering, albeit with increased communication overhead. However, having clients directly adopt the cluster center model for inference compromises performance in complex environments. CGPFL [33] remedies the problem by preserving personalized models for each client. The cluster center facilitates the learning of intra-cluster knowledge through soft regularization.

CFL methods face challenges: (1) they frequently confine knowledge sharing to particular clusters, lacking effective strategies to integrate global knowledge; (2) online CFL methods involve substantial computational overhead in re-clustering, lacking an efficient and cost-effective similarity metric; (3) all CFL methods require prior specification of the total cluster count, limiting scalability and flexibility in complex environments.

### 3 Problem Formulation

In the context of FL, each client  $i$  possesses its private dataset  $\mathcal{D}_i$  from distribution  $\mathbb{P}_i(x, y)$ , where  $x$  and  $y$  represent the input features and corresponding labels. In a vanilla setting (FedAvg), clients collectively contribute to a shared model  $f(\omega; \cdot)$  parameterized by weights  $\omega$ . The objective function is:

$$\min_{\omega} : \mathcal{F} = \sum_{i=1}^m \frac{|\mathcal{D}_i|}{N} \mathcal{L}_i(\omega), \quad (1)$$

where  $\mathcal{L}_i(\omega) = \mathbb{E}_{(x,y) \sim \mathcal{D}_i} l(f(\omega; x); y)$  is the empirical loss of client  $i$ ,  $m$  is the number of clients, and  $N$  denotes the total number of instances over all clients.

To further address data heterogeneity, personalized client-specific models,  $\{f_i\}_{i \in [m]}$ , are often employed across the system. While model architecture and hyper-parameters may differ among clients, in our study, we maintain identical model architectures for each, but with unique local parameters  $\{\omega_i\}_{i \in [m]}$ . Following this, clustering configurations' integration is proposed as the problem setting. The server retains  $K$  cluster center models  $\{\Omega_k\}_{k \in [K]}$ , rather than a single global model, to steer local model updates:

$$\begin{aligned}
& \min_{\{\omega_i\}} \frac{1}{m} \sum_{i=1}^m \{ \mathcal{L}_i(\omega_i) + \frac{\mu}{2} \text{Regu}(\omega_i, \Omega_{k^*(i, R^*)}) \}, \\
& \text{s.t. } R^* \in \arg \max_R \frac{1}{m} \sum_{k=1}^K \sum_{i=1}^m R_{i,k} \text{Sim}(\omega_i, \Omega_k(\{\omega_i\}, R)),
\end{aligned} \tag{2}$$

where  $\mu$  is a regularization parameter,  $R \in \mathbb{R}^{m \times K}$  is the assignment matrix where  $R_{i,k} = 1$  if client  $i$  belongs to cluster  $k$  else  $R_{i,k} = 0$ ,  $\Omega_k(\{\omega_i\}, R) = \frac{1}{\sum_{i=1}^m R_{i,k}} \sum_{i=1}^m R_{i,k} \omega_i$  is the central model of cluster  $k$ , and  $k^*(i, R)$  represents the cluster to which client  $i$  belongs, corresponding to the sole 1 in client  $i$ 's row of the allocation matrix  $R$ . Note that  $k^*$  is clearly defined per  $R$ 's definition.  $\text{Regu}(\cdot, \cdot)$  regulates model difference, while  $\text{Sim}(\cdot, \cdot)$  measures model similarity. Item  $\frac{|\mathcal{D}_i|}{N}$  is omitted from Eq. 1 to adopt a macroscopic view.

## 4 Methodology

### 4.1 Overview and Optimization Objective

The proposed **F**ederated Learning with **A**daptive **C**lustering (**FedAC**) addresses data heterogeneity and achieves efficient, adaptable client clustering. Firstly, to better integrate beneficial global knowledge into cluster, the model  $f$  is decoupled into two submodules, exemplified by convolutional neural networks (convnets): the embedding  $\phi$  (comprising shallow layers with convolutional modules) and the decision  $h$  (comprising deep fully connected layers). Different aggregation strategies are designed for distinct submodules to represent global and intra-cluster knowledge, achieving optimal performance by balancing both aspects. Moreover, a cosine similarity metric based on dimensionality-reduced models is proposed for efficient online model similarity measurement to assist clustering throughout training. Additionally, a self-examining module evaluates the current cluster conditions and autonomously finetunes the cluster number. An overview of FedAC's framework is shown in Fig. 1.

We apply L2 distance and the proposed low-rank cosine model similarity,  $LrCos$ , into  $\text{Regu}(\cdot, \cdot)$  and  $\text{Sim}(\cdot, \cdot)$ , respectively, within Eq. 2, defining the optimization objective of FedAC:

$$\begin{aligned}
& \min_{\{\omega_i\}} \frac{1}{m} \sum_{i=1}^m \{ \underbrace{\mathcal{L}_i(\omega_i)}_{L_{\text{sup}}} + \frac{\mu}{2} \underbrace{\|\omega_i - \Omega_{k^*(i, R^*)}\|_2^2}_{L_{\text{intra}}} + \frac{\lambda}{2} \underbrace{\|\phi_i - \Phi(\{\phi_i\})\|_2^2}_{L_{\text{global}}} \}, \\
& \text{s.t. } R^* \in \arg \max_R \frac{1}{m} \sum_{k=1}^K \sum_{i=1}^m R_{i,k} LrCos(\omega_i, \Omega_k(\{\omega_i\}, R)),
\end{aligned} \tag{3}$$

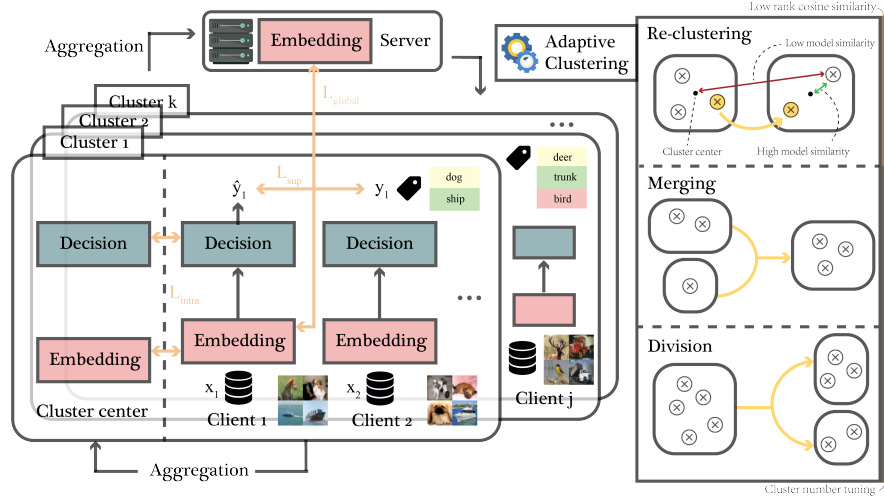


Fig. 1: The framework of FedAC in the heterogeneous setting. The server clusters and aggregates clients to get a global embedding and multiple cluster center models. Clients update local models by minimizing the classification error loss ( $L_{\text{sup}}$ ), intra-clustering regularization ( $L_{\text{intra}}$ ), and global regularization ( $L_{\text{global}}$ ).

where hyper-parameter  $\mu$  and  $\lambda$  controls the level of regularization, and  $\Phi(\{\phi_i\}) = \sum_{i=1}^m \phi_i$  is the global embedding. Note that the embedding regularization term has been included to capture global knowledge.

The bi-level optimization problem is addressed through alternate minimization, involving two iterative steps: (1) minimization w.r.t.  $\{\omega_i\}_{i \in [m]}$  with  $R$  fixed; and (2) minimization w.r.t.  $R$  with  $\{\omega_i\}_{i \in [m]}$  fixed. Step (1) requires clients to utilize local data for supervised learning, with regularization implemented through global embedding and the cluster center. In step (2), cluster assignments and center models are updated based on the  $LrCos$  similarity, coupled with a EM-like algorithm. Additional details are presented in the subsequent sections, and Algorithm 1 provides the pseudocode.

## 4.2 Integrating Global Knowledge into Clusters

As one of the primary challenges in CFL methods is effectively integrating beneficial global knowledge into clusters, we employ a novel approach of decoupling neural networks for clients, allowing for finer control and balance between global and intra-cluster knowledge. This draws inspiration from the effective management of latent feature spaces in [24], [34], etc. In the context of convnets, shallow layers primarily handle pixel embedding and feature extraction tasks, capturing colors and edges from images. These features are considered advantageous for sharing among all clients, serving as the foundation for downstream tasks while being less affected by heterogeneous data distributions [35]. The deep lay-

**Algorithm 1** FedAC**Input:** learning rate  $\eta$ , hyper-parameters  $\mu$  and  $\lambda$ , initial number of clusters  $K$ **Output:**  $\{\omega_i\}_{i \in [m]}$ **Server executes:**

- 1: Initialize:  $R^0, \{\omega_i^0 = (\phi_i^0, h_i^0)\}_{i \in [m]}, \{\Omega_k^0\}_{k \in [K]}, \Phi^0$
- 2: **for** each round  $t = 0, 1, \dots$  **do**
- 3:   Randomly selects a subset of clients  $S_t$
- 4:   **for** each client  $i \in S_t$  **in parallel do**
- 5:     Server sends  $\Omega_{k^*}^t, \Phi^t$  to client  $i$
- 6:      $\omega_i^{t+1} \leftarrow \text{LocalUpdate}(\omega_i^t, \Omega_{k^*}^t, \Phi^t, \mu, \lambda, \eta)$
- 7:     Clients  $i$  sends  $\omega_i^{t+1}$  back
- 8:   **end for**
- 9:    $\Phi^{t+1} = \sum_{i \in S_t} \phi_i^{t+1}$
- 10:   Calculate low-rank cosine model similarity {Algorithm 2}
- 11:    $R^{t+1} \leftarrow \text{E-step}(\{\omega_i^{t+1}\}_{i \in [m]}, \{\Omega_k^t\}_{k \in [K]})$
- 12:    $\{\Omega_k^{t+1}\}_{k \in [K]} \leftarrow \text{M-step}(\{\omega_i^{t+1}\}_{i \in [m]}, R^{t+1})$
- 13:   Cluster number tuning (CNT) {Algorithm 3}
- 14: **end for**

**LocalUpdate** ( $\omega_i^t, \Omega_{k^*}^t, \Phi^t, \mu, \lambda, \eta$ ):

- 1:  $\omega_i^{t,0} = \omega_i^t$
- 2: **for** each local epoch  $r = 0, 1, \dots$  **do**
- 3:   Randomly selects a batch  $\mathcal{B}_i$  from  $\mathcal{D}_i$
- 4:    $\omega_i^{t,r+1} = \omega_i^{t,r} - \eta \nabla l_i(\omega_i^{t,r}; \mathcal{B}_i) - \eta \mu (\omega_i^{t,r} - \Omega_{k^*}^t) - \eta \lambda (\phi_i^{t,r} - \Phi^t)$
- 5: **end for**
- 6: **return**  $\omega_i^{t,r+1}$

ers, which capture unique data distributions for each client, are preferred to be shared at the intra-cluster level, as they serve specific downstream tasks [36].

Based on this, FedAC simultaneously updates global embedding  $\Phi$  (line 9 of Algorithm 1) and cluster-center models  $\{\Omega_k\}_{k \in [K]}$  (including embeddings and decisions, line 12 of Algorithm 1) at each aggregation step. As client  $i$  undergoes local updates, it learns global knowledge from  $\Phi$ , regulated by the term  $L_{global}$  in Eq.3, while simultaneously learning intra-cluster knowledge from its corresponding cluster center model  $\Omega_{k^*}^{(i,R)}$ , regulated by the term  $L_{intra}$ . By adjusting the corresponding regularization strengths of  $L_{global}$  and  $L_{intra}$  to balance their impacts on client's local updates, FedAC ensures the coordinated integration of intra-cluster and global knowledge, thereby achieving higher accuracy, as demonstrated in the experimental section.

### 4.3 Low-rank Cosine Model Similarity

Compared to one-shot clustering, online clustering methods, which continuously update model similarities and adjust clustering, often yield more effective clustering results but also come with increased computational overhead. We combine the benefits of one-shot and online clustering by sparsely computing the dimensionality reduction matrix and leveraging the cosine similarity of

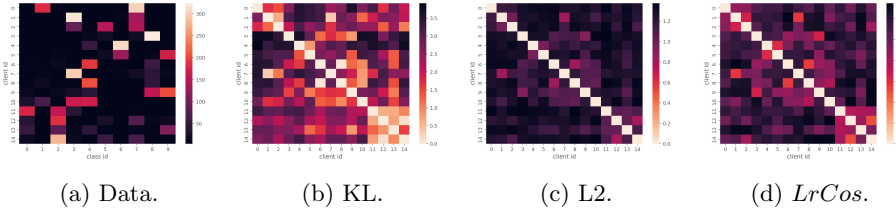
**Algorithm 2** Low-rank Cosine Similarity**Input:**  $\{\omega_i\}_{i \in [m]}, \{\Omega_k\}_{k \in [K]}$ **Parameter:** reduced number of dimensions  $D$ **Output:**  $LrCos_{i,j}$  between client  $i$  and cluster  $k$ 1: Mapping  $M \leftarrow \text{UpdateMap}(S_t, D)$  {Proceed sparsely}2:  $LrCos_{i,j} = \frac{M \cdot \omega_i \cdot M \cdot \Omega_k}{\|M \cdot \omega_i\|_2 \|M \cdot \Omega_k\|_2}$ **UpdateMap** ( $S_t, D$ ):1:  $H \in \mathbb{R}^{dim(\omega) \times |S|} \leftarrow [\omega_1, \dots, \omega_{|S|}]$ 2:  $M \in \mathbb{R}^{D \times dim(\omega)} \leftarrow \text{PCA}(H, \text{components}=D)$ 3: **return**  $M$ 

Fig. 2: Clients employed FedAvg with CIFAR-10 [37] for 20 epochs, succeeded by local fine-tuning, with model similarities subsequently evaluated.

the reduced-dimensional model, namely, the low-rank cosine similarity  $LrCos$ , enabling efficient model similarity updates per round with low computational costs.

We utilize PCA [38] for model dimensionality reduction, compatible with other methods. The server sporadically updates the dimensionality reduction matrix  $M$  (e.g., every 100 rounds). Each client locally retains matrix  $M$  and, while transmitting its model updates to the server, also sends the reduced model  $M \cdot \omega$ . After each update of the cluster center model  $\Omega$ , the server performs dimensionality reduction to acquire  $M \cdot \Omega$ . Following Algorithm 2, it uses cosine similarity to calculate the updated similarity between clients and cluster centers, enabling subsequent re-clustering.

Simulations confirm the effectiveness of  $LrCos$ . Fig. 2(a) illustrates the number of training samples held by each client for each class. Fig. 2(b) illustrates the Kullback-Leibler (KL) divergence between the data distribution of client pairs, serving as a benchmark for evaluating model similarity measurements.  $LrCos$ , depicted in Fig. 2(d), closely aligns with the KL representations, whereas L2 distance (Fig. 2(c)) inadequately captures this similarity. Furthermore, dimensionality reduction significantly reduces the computational burden for similarity calculations. In our settings, the original model dimension is reduced from over 4,700,000 to  $D = 50$ . Sparse updates to the dimensionality reduction matrix guarantee that the additional computational cost is negligible compared to the saved computation cost, especially in large-scale scenarios.



#### 4.4 Re-clustering

FedAC employs a EM-like algorithm for periodic clustering updates. In each E-step, clients are reassigned to nearby clusters based on the  $LrCos$  of their personalized models and cluster centroids:

$$R_{i,k} = \begin{cases} 1, & k = \arg \min_j LrCos(\omega_i, \Omega_j) \\ 0, & \text{else} \end{cases}. \quad (4)$$

In the M-step, the server aggregates the models within each cluster to obtain the cluster centroids:

$$\Omega_k = \frac{1}{\sum_{i=1}^m R_{i,k}} \sum_{i=1}^m R_{i,k} \omega_i. \quad (5)$$

#### 4.5 Cluster Number Tuning

In complex heterogeneous environments, manually determining the optimal cluster count  $K$  as a predefined hyper-parameter poses challenges for CFL methods. FedAC provides a novel approach to address this by dynamically fine-tuning the value of  $K$  until optimal during the training process.

The Cluster Number Tuning (CNT) module in FedAC (Algorithm 3) dynamically assesses model distances within and between clusters (denoted as  $Dist_{intra}$  and  $Dist_{inter}$ ), determining cluster-splitting or merging decisions. We empirically define the ratio  $G_c = \frac{Dist_{intra}}{Dist_{inter}}$  to characterize the clustering granularity and constrain it within a reasonable range. A judicious  $G_c$  range prevents insufficient collaboration within a cluster (when  $G_c$  is too small) and moderates excessive clustering effects (when  $G_c$  is too large).

---

**Algorithm 3** Cluster Number Tuning (CNT)

---

**Input:**  $\{\omega_i\}_{i \in [m]}$ ,  $\{\Omega_k\}_{k \in [K]}$

**Parameter:** lower and upper threshold,  $a$  and  $b$ , of  $G_c$

```

1: for each cluster  $k \in [K]$  do
2:    $Dist_{intra}^k = \frac{1}{\sum_{i=1}^m R_{i,k}} \sum_{i=1}^m R_{i,k} \|\omega_i - \Omega_k\|_2^2$ 
3:    $Dist_{inter}^k = \frac{1}{K-1} \sum_{j=1}^K \|\Omega_k - \Omega_j\|_2^2$ 
4:    $G_c^k = \frac{Dist_{intra}^k}{Dist_{inter}^k}$ 
5:   if  $G_c^k < a$  then
6:     Merge cluster  $k$  into the closest cluster.
7:   else if  $G_c^k > b$  then
8:     Divide cluster  $k$  into 2 clusters.
9:   end if
10: end for

```

---

Drawing from this intuition, the CNT module empirically divides clusters with large  $G_c$  values into two new clusters by initializing two new centers and conducting one re-clustering iteration, thereby increasing the total cluster count  $K$ . Conversely, clusters with excessively small  $G_c$  values are merged into the nearest cluster, resulting in a reduction of  $K$ . Subsequent experiments validate the efficacy of this empirical approach in optimizing  $K$  to its optimal value.

## 5 Experiments

### 5.1 Experimental Setup

We conduct experiments with two authentic datasets: CIFAR-10 [37] and CIFAR-100 [37]. Consistent with previous research [33], we employed the Dirichlet distribution with  $\alpha$  set to 0.1 and the pathological distribution (where  $n$  denotes the number of labels per client) to introduce high heterogeneity. The quantity of instances per client was randomly assigned, spanning from 50 to 350 for both datasets. FedAC and the baseline models were implemented using Python 3.7 and PyTorch 1.12.1. The training process took place on a compute server equipped with an NVIDIA RTX 3080 GPU, Intel i9-10900K CPU, 64 GB RAM, and a 2TB SSD. The experiments involved convnets with 100 clients, sampling 25 per communication round. These convnets consist of two convolution layers and three fully connected layers, utilizing the final fully connected layer for decision-making and the preceding layers for embedding. The SGD optimizer is employed with a learning rate set to 0.01. The experimental setup includes a batch size of 32, local iterations of 5, and communication rounds of 1000, unless specified otherwise.

### 5.2 Overall Performance

We conducted a comparative analysis, benchmarking FedAC against seven prominent SOTA methodologies in the field. These methodologies are: (1) FedAvg [18], (2) FedPer [24], (3) FeSEM [31], (4) FedGroup [30], (5) FL+HC [29], (6) CGPFL [33], and (7) IFCA [32]. Table 1 displays the improved performance of FedAC compared to benchmarks in common heterogeneous scenarios. CFL methods determine their respective optimal cluster numbers  $K$ . It’s important to highlight that FeSEM, FedGroup, FL+HC, and IFCA employ a single shared model for each cluster, whereas FedPer, CGPFL, and the proposed FedAC maintain a personalized model for each client. In contrast to the Dirichlet distribution, the pathological distribution enhances the CFL method’s performance by naturally leading certain clients to acquire identical label categories.

IFCA demonstrates commendable performance but incurs substantial computational and communication costs. CGPFL preserves personalized models for clients during clustering, yielding enhanced performance in heterogeneous conditions. However, its knowledge sharing is confined to within clusters. FedAC attains optimal performance by amalgamating personalized models, proficient cluster partitioning, and simultaneous coordination of global knowledge sharing.

Method	CIFAR10 Test Accuracy (%)		CIFAR100 Test Accuracy (%)	
	$\alpha = 0.1$	$n = 3$	$\alpha = 0.1$	$n = 8$
FedAvg	64.75 $\pm$ 1.21	62.73 $\pm$ 0.42	35.36 $\pm$ 0.92	33.33 $\pm$ 0.42
FedPer	70.84 $\pm$ 1.44	79.36 $\pm$ 0.54	43.46 $\pm$ 1.27	62.10 $\pm$ 0.63
FeSEM	65.65 $\pm$ 1.28	76.46 $\pm$ 0.63	31.03 $\pm$ 0.80	52.15 $\pm$ 2.45
FedGroup	67.38 $\pm$ 1.34	76.77 $\pm$ 1.00	33.26 $\pm$ 1.01	57.02 $\pm$ 0.84
FL+HC	67.80 $\pm$ 0.84	80.22 $\pm$ 0.68	34.19 $\pm$ 1.33	57.99 $\pm$ 0.66
CGPFL	71.19 $\pm$ 0.93	79.42 $\pm$ 0.46	41.38 $\pm$ 0.69	60.26 $\pm$ 0.94
IFCA	73.06 $\pm$ 0.91	80.54 $\pm$ 0.74	38.61 $\pm$ 0.77	60.18 $\pm$ 0.76
<b>FedAC</b>	<b>74.88<math>\pm</math>0.65</b>	<b>81.29<math>\pm</math>0.61</b>	<b>51.28<math>\pm</math>0.35</b>	<b>64.53<math>\pm</math>0.34</b>

Table 1: FedAC’s performance (mean and standard deviation) is compared to baseline methods in non-IID settings. The best is highlighted in bold.

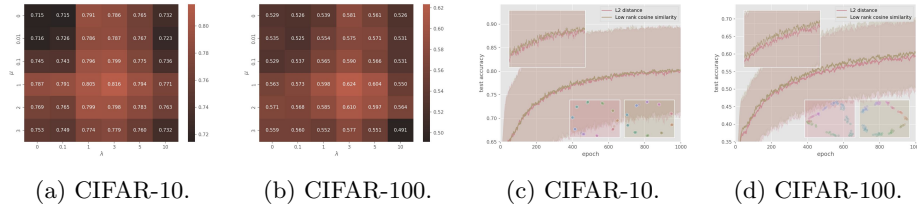


Fig. 3: Ablation experiments were conducted to assess FedAC’s performance across different parameter configurations and model similarity measurements.

### 5.3 Global and Intra-Cluster Trade-Offs

FedAC incorporates regularizing terms for cluster center and global embedding into local model training, fostering alignment between local and global understanding. In Eq. 3,  $\mu$  regulates the intensity of cluster center model regularization, while  $\lambda$  governs the strength of global embedding regularization. Larger values of these parameters signify a greater reliance on intra-cluster or global information in local updates. Figures 4(a) and 4(b) validate the effectiveness of employing two regularization terms to balance intra-cluster and global information across experimental datasets. Achieving optimal FedAC performance necessitates meticulous adjustment of these regulatory intensities.

### 5.4 Improvement in Similarity Measurement

Following the explanation of the *LrCos* metric in section 4.3, ablation experiments were conducted by substituting *LrCos* with L2 distance, as depicted in Fig. 4(c) and (d). The utilization of *LrCos* leads to a more sensible client partitioning and contributes to enhanced model testing accuracy, especially evident in more intricate tasks such as CIFAR-100.

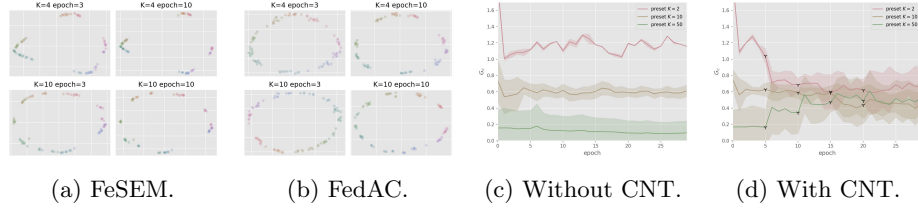


Fig. 4: Client models, trained with FeSEM and FedAC, are visualized by projecting them onto a plane in (a) and (b). Clients with similar data distributions are represented by similar colors. In (c), the mean and variance of clusters'  $G_c$  are shown for various predefined  $K$  values (with  $K = 10$  as the optimal value), while subplot (d) depicts the CNT module fine-tuning the cluster number by adjusting  $G_c$  (marked in black).

### 5.5 Clustering Robustness

Here, we visualize the clustering for online CFL approaches that lack personalized client models (*e.g.*, FeSEM) and those that incorporate them (*e.g.*, FedAC). The former, which results in swift convergence of models within clusters, contradicts the purpose of re-clustering in online CFL, thereby reducing the system's robustness, as clients within the same cluster initiate updates from identical positions each round (see Fig. 4(a)). FedAC, depicted in Fig. 4(b), addresses this by introducing a soft regularization term ( $L_{\text{intra}}$  in Eq. 3) to the personalized models. This also improves the system's resilience to variations in cluster number settings and initializations of cluster centers. Additionally, as shown in Fig. 4(c) and (d), we set the range of the system's  $G_c$  to (0.2, 0.8). The CNT module automates the adjustment of  $G_c$  to effectively fine-tune to the optimal  $K$ .

## 6 Conclusion

This paper proposes a clustered federated learning framework that integrates global and intra-cluster knowledge through neural network decomposition. It employs a low-rank cosine similarity for efficient and economical client clustering. Additionally, an integrated module facilitates adaptive cluster tuning to determine the optimal count. The framework's effectiveness is rigorously evaluated through experimental analysis.

## References

1. Z. Lin, G. Zhu, Y. Deng, X. Chen, Y. Gao, K. Huang, and Y. Fang, "Efficient Parallel Split Learning over Resource-constrained Wireless Edge Networks," *IEEE Trans. Mobile Comput.*, Jan. 2024.
2. H. Dai, J. Wu, A. Brinkmann, and Y. Wang, "Neighborhood-Oriented Decentralized Learning Communication in Multi-Agent System," in *Proc. ICANN*, Sep. 2023.

3. T. Zheng, A. Li, Z. Chen, H. Wang, and J. Luo, "AutoFed: Heterogeneity-Aware Federated Multimodal Learning for Robust Autonomous Driving," in *Proc. Mobi-com*, Oct. 2023.
4. J. Campoy, I.-I. Prado-Rujas, J. L. Risco-Martin, K. Olcoz, and M. S. Perez, "Distributed Training and Inference of Deep Learning Solar Energy Forecasting Models," in *Proc. PDP*, Jun. 2023.
5. Z. Lin, L. Wang, J. Ding, B. Tan, and S. Jin, "Channel Power Gain Estimation for Terahertz Vehicle-to-infrastructure Networks," *IEEE Commun. Lett.*, vol. 27, no. 1, pp. 155–159, Sep. 2022.
6. H. Yuan, Z. Chen, Z. Lin, J. Peng, Z. Fang, Y. Zhong, Z. Song, X. Wang, and Y. Gao, "Graph Learning for Multi-satellite Based Spectrum Sensing," in *Proc. ICCT*, Oct. 2023.
7. S. Hu, Z. Fang, H. An, G. Xu, Y. Zhou, X. Chen, and Y. Fang, "Adaptive Communications in Collaborative Perception with Domain Alignment for Autonomous Driving," *arXiv preprint arXiv:2310.00013*, Sep. 2023.
8. Z. Lin, G. Qu, W. Wei, X. Chen, and K. K. Leung, "AdaptSFL: Adaptive Split Federated Learning in Resource-constrained Edge Networks," *arXiv preprint arXiv:2403.13101*, Mar. 2024.
9. Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May. 2015.
10. S. Hu, Z. Fang, Y. Deng, X. Chen, and Y. Fang, "Collaborative Perception for Connected and Autonomous Driving: Challenges, Possible Solutions and Opportunities," *arXiv preprint arXiv:2401.01544*, Jan. 2024.
11. I.-I. Prado-Rujas, E. Serrano, A. García-Dopico, M. L. Córdoba, and M. S. Pérez, "Combining Heterogeneous Data Sources for Spatio-temporal Mobility Demand Forecasting," *Inf. Fusion.*, vol. 91, pp. 1–12, Mar. 2023.
12. Z. Lin, G. Qu, Q. Chen, X. Chen, Z. Chen, and K. Huang, "Pushing Large Language Models to the 6g Edge: Vision, Challenges, and Opportunities," *arXiv preprint arXiv:2309.16739*, Sep. 2023.
13. L. Marelli and G. Testa, "Scrutinizing the EU General Data Protection Regulation," *Science*, vol. 360, no. 6388, pp. 496–498, May. 2018.
14. S. Lyu, Z. Lin, G. Qu, X. Chen, X. Huang, and P. Li, "Optimal Resource Allocation for U-shaped Parallel Split Learning," *arXiv preprint arXiv:2308.08896*, Oct. 2023.
15. Q. Chen, W. Meng, T. Q. Quek, and S. Chen, "Multi-tier Hybrid Offloading for Computation-aware IoT Applications in Civil Aircraft-augmented SAGIN," *IEEE J. Select. Areas Commun.*, vol. 41, no. 2, pp. 399–417, Dec. 2022.
16. Z. Lin, Z. Chen, Z. Fang, X. Chen, X. Wang, and Y. Gao, "FedSN: A General Federated Learning Framework over LEO Satellite Networks," *arXiv preprint arXiv:2311.01483*, Nov. 2023.
17. P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and Open Problems in Federated Learning," *Found. Trends Mach. Learn.*, vol. 14, no. 1–2, pp. 1–210, 2021.
18. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. AISTATS*, vol. 54, Apr. 2017, pp. 1273–1282.
19. O. Marfoq, G. Neglia, A. Bellet, L. Kameni, and R. Vidal, "Federated Multi-task Learning under A Mixture of Distributions," Dec. 2021.
20. Z. Lin, G. Qu, X. Chen, and K. Huang, "Split Learning in 6G Edge Networks," *arXiv preprint arXiv:2306.12194*, Jun. 2023.

21. T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated Optimization in Heterogeneous Networks," in *Proc. MLSys*, May 2020.
22. F. Sattler, K.-R. Müller, and W. Samek, "Clustered Federated Learning: Model-Agnostic Distributed Multitask Optimization Under Privacy Constraints," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 8, pp. 3710–3722, Aug. 2021.
23. R. Caruana, "Multitask learning," *Mach Learn*, vol. 28, pp. 41–75, Jul. 1997.
24. M. G. Arivazhagan, V. Aggarwal, A. K. Singh, and S. Choudhary, "Federated Learning with Personalization Layers," Dec. 2019.
25. A. Fallah, A. Mokhtari, and A. Ozdaglar, "Personalized Federated Learning with Theoretical Guarantees: A Model-Agnostic Meta-Learning Approach," in *Proc. NIPS*, Jun. 2020.
26. C. T. Dinh, N. Tran, and J. Nguyen, "Personalized Federated Learning with Moreau Envelopes," in *Proc. NIPS*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., Jun. 2020.
27. A. Shamsian, A. Navon, E. Fetaya, and G. Chechik, "Personalized Federated Learning using Hypernetworks," in *Proc. ICLR*, Jul. 2021.
28. Y. Mei, B. Guo, D. Xiao, and W. Wu, "FedVF: Personalized Federated Learning Based on Layer-wise Parameter Updates with Variable Frequency," in *Proc. IPCCC*, Oct. 2021.
29. C. Briggs, Z. Fan, and P. Andras, "Federated Learning with Hierarchical Clustering of Local Updates to Improve Training on Non-IID Data," in *Proc. IJCNN*, Sep. 2020.
30. M. Duan, D. Liu, X. Ji, R. Liu, L. Liang, X. Chen, and Y. Tan, "FedGroup: Efficient Federated Learning via Decomposed Similarity-Based Clustering," in *Proc. ISPA*, Dec. 2021.
31. G. Long, M. Xie, T. Shen, T. Zhou, X. Wang, and J. Jiang, "Multi-center federated learning: Clients clustering for better personalization," *World Wide Web*, vol. 26, no. 1, pp. 481–500, Jun. 2022.
32. A. Ghosh, J. Chung, D. Yin, and K. Ramchandran, "An Efficient Framework for Clustered Federated Learning," in *Proc. NeurIPS*, Dec. 2020.
33. X. Tang, S. Guo, and J. Guo, "Personalized federated learning with contextualized generalization," in *Proc. IJCAI*, Jul. 2022.
34. Y. Tan, G. Long, L. LIU, T. Zhou, Q. Lu, J. Jiang, and C. Zhang, "FedProto: Federated Prototype Learning across Heterogeneous Clients," in *Proc. AAAI*, Jun. 2022.
35. A. Babenko and V. Lempitsky, "Aggregating Local Deep Features for Image Retrieval," in *Proc. ICCV*, Dec. 2015.
36. M. Luo, F. Chen, D. Hu, Y. Zhang, J. Liang, and J. Feng, "No Fear of Heterogeneity: Classifier Calibration for Federated Learning with Non-iid Data," in *Proc. NeurIPS*, Dec. 2021.
37. A. Krizhevsky, G. Hinton *et al.*, *Learning Multiple Layers of Features from Tiny Images*. Toronto, ON, Canada, Apr. 2009.
38. S. Wold, K. Esbensen, and P. Geladi, "Principal Component Analysis," *Chemom. Intell. Lab. Syst.*, vol. 2, no. 4, pp. 433–459, Aug. 1987.