

Geometric Knowledge-Guided Localized Global Distribution Alignment for Federated Learning

Yanbiao Ma*

Xidian University

ybmamail@stu.xidian.edu.cn

Wei Dai*

Xidian University

22012100039@stu.xidian.edu.cn

Wenke Huang†

Wuhan University

wenkehuang@whu.edu.cn

Jiayi Chen

Xidian University

22012100031@stu.xidian.edu.cn

Abstract

Data heterogeneity in federated learning, characterized by a significant misalignment between local and global distributions, leads to divergent local optimization directions and hinders global model training. Existing studies mainly focus on optimizing local updates or global aggregation, but these indirect approaches demonstrate instability when handling highly heterogeneous data distributions, especially in scenarios where label skew and domain skew coexist. To address this, we propose a geometry-guided data generation method that centers on simulating the global embedding distribution locally. We first introduce the concept of the geometric shape of an embedding distribution and then address the challenge of obtaining global geometric shapes under privacy constraints. Subsequently, we propose GGEUR, which leverages global geometric shapes to guide the generation of new samples, enabling a closer approximation to the ideal global distribution. In single-domain scenarios, we augment samples based on global geometric shapes to enhance model generalization; in multi-domain scenarios, we further employ class prototypes to simulate the global distribution across domains. Extensive experimental results demonstrate that our method significantly enhances the performance of existing approaches in handling highly heterogeneous data, including scenarios with label skew, domain skew, and their coexistence.

Code published at: https://github.com/WeiDai-David/2025CVPR_GGEUR

1. Introduction

Federated Learning (FL) enables multiple clients to collaboratively train a global model without sharing raw data, ef-

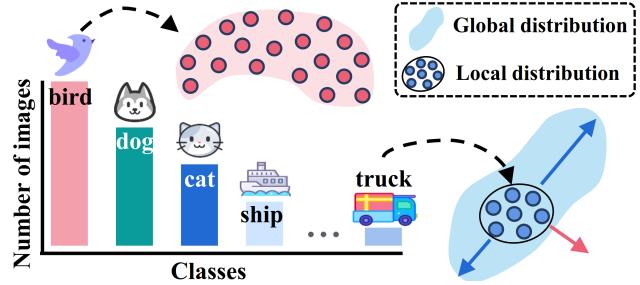


Figure 1. The distribution of local data on clients is imbalanced. The “bird” category has a large number of samples, allowing its local distribution to effectively represent the global distribution. In contrast, the “truck” category has insufficient samples, resulting in a significant disparity between the local and global distributions.

fectively addressing key concerns related to data privacy and security [13, 19, 53]. However, data heterogeneity remains one of the primary challenges in FL, especially in the presence of label skew and domain skew across clients [12, 15, 17, 28, 29]. Such distributional differences lead to variations in the local optimization directions for each client, which can be particularly pronounced when data is sourced from multiple domains [5, 33, 42, 59, 66]. This cross-client distribution shift not only slows down global model convergence but also undermines the model’s generalization capability across clients [5, 13, 30, 62].

Existing research addressing data heterogeneity in federated learning primarily focuses on two categories of methods [13]. The first is server-side optimization strategies, which aim to enhance model convergence by adjusting the global aggregation process [2, 3, 11, 32, 44, 52, 54, 57, 61, 64]. The second category includes client-side regularization [1, 7, 12, 22, 23, 42] and data augmentation methods [67, 68], which attempt to align local model optimization through guidance from the global model and loss function adjustments or directly apply data augmentation techniques

*Co-first authors

†Corresponding Author: Wenke Huang

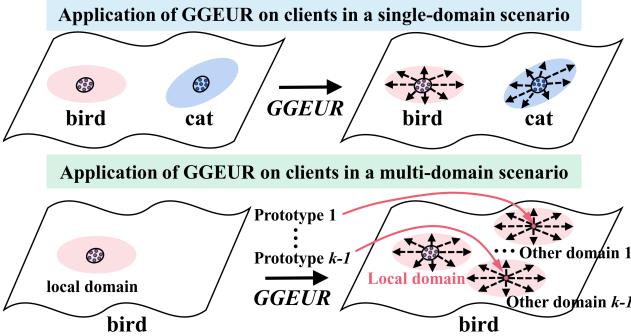


Figure 2. Example of “bird” category: in a single-domain scenario, only the local global distribution is simulated; in a multi-domain scenario, the global distribution from other data domains is also simulated on a client.

like Mixup in the local setting [63]. However, these approaches do not directly tackle the core issue of data heterogeneity, namely, the significant misalignment between local data distributions on clients and the global data distribution. They often rely on indirect means to optimize the alignment between local and global models, which typically incurs high computational costs and depends heavily on the selection of auxiliary datasets, resulting in limited stability [13]. **Motivated** by these challenges, we propose a core idea: if clients can locally simulate an ideal global distribution, it may be possible to alleviate the divergence in local optimization directions, even without data sharing. In other words, through guided data augmentation, we can reduce the inconsistency between local and global distributions, fundamentally mitigating data heterogeneity issues in federated learning.

How can we simulate a global distribution locally? Taking CIFAR-10 as an example, Figure 1 illustrates the label distribution in a random client, revealing severe imbalances. Certain categories, such as “truck,” are underrepresented, making it difficult to form an effective global representation. Given privacy constraints, simulating a global distribution locally becomes a significant challenge [13]. As shown in the blue distribution in Figure 1, augmenting samples along the blue arrow direction, rather than the pink arrow, is more likely to simulate the global distribution. Thus, **identifying appropriate augmentation directions and ranges becomes our key objective to overcome this challenge.**

In the single-domain scenario, some clients may possess a large number of “truck” images, while others have only a few. The aggregate set of “truck” images across all clients constitutes the global distribution for that category. **We introduce** the concept of geometric shapes to describe the orientation and range of data distributions (see Section 2). For instance, the global distribution of “truck” in Figure 1 spans primarily from the bottom left to the top right. If we could quantify the global geometric shape of the “truck” and disseminate it to clients, they could use it to guide new sample

generation locally. However, privacy constraints prohibit us from aggregating all client samples to quantify this global geometric shape. To address this, **we propose** using local covariance matrices from multiple clients on the server side to approximate the global covariance matrix, enabling quantification of the global geometric shape (see Section 3.1). To mitigate inaccuracies in local covariance matrix estimation due to high dimensionality and sparsity in image space [36], we employ the CLIP [47] to extract embedding distributions from local client data and carry out the approximation process in the embedding space.

After obtaining the global geometric shape for each category, **we propose** a method called Global Geometry-Guided Embedding Uncertainty Representation (**GGEUR**) to generate new samples on clients, simulating an ideal global distribution in the embedding space. Each client then only needs to train an MLP as a local classifier on the augmented data. In multi-domain scenarios, we observe that the geometric shapes of the embedding distributions for the same category across different domains exhibit similarity (see Figure 6 and Section 3.2). This makes it feasible to approximate the global geometric shape using local geometric shapes across all clients, even in multi-domain settings. In multi-domain scenarios, however, the global distribution for each category comprises distributions from all domains, necessitating that clients also simulate embeddings from other domains. Although the embedding distributions of the same category across different domains are geometrically similar, they occupy different spatial positions. Thus, unlike the single-domain scenario, we propose distributing category prototypes from other domains as shared knowledge to clients and applying GGEUR to these prototypes to generate new samples that simulate distributions from other domains (see Section 3.2). Figure 2 illustrates the difference between single-domain and multi-domain scenarios.

Our approach, serving as a preprocessing step in FL, is easily integrated with other methods. Extensive evaluation across scenarios with label skew, domain skew, and combined skew demonstrates that our method significantly improves the performance of existing approaches, achieving state-of-the-art results in highly heterogeneous data scenarios (Section 4). This work represents a typical instance of synergy between foundational visual models and geometric knowledge within FL.

2. Geometry of Embedding Distributions

We now define the geometry of an embedding distribution and introduce a measure of similarity between these geometries. In a P -dimensional space, given a dataset of a certain class $X = [x_1, \dots, x_n] \in \mathbb{R}^{P \times n}$, the covariance matrix of the distribution can be estimated as:

$$\Sigma_X = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right] = \frac{1}{n} X X^T \in \mathbb{R}^{P \times P}. \quad (1)$$

Performing eigenvalue decomposition on Σ_X yields P eigenvalues $\lambda_1 \geq \dots \geq \lambda_P$ and their corresponding P -dimensional eigenvectors ξ_1, \dots, ξ_P . All eigenvectors are mutually orthogonal, collectively anchoring the skeleton of the data distribution, with each eigenvalue specifying the range of the distribution along the direction of its corresponding eigenvector. The combination of eigenvalues and eigenvectors defines the geometric shape of the distribution. Thus, we define the geometric shape of the data X as:

$$GD_X = \{\xi_1, \dots, \xi_P, \lambda_1, \dots, \lambda_P\}. \quad (2)$$

Given the geometries GD_{X_1} and GD_{X_2} of two distributions, their similarity is defined as:

$$S(GD_{X_1}, GD_{X_2}) = \sum_{i=1}^P |\langle \xi_{X_1}^i, \xi_{X_2}^i \rangle|. \quad (3)$$

where $S(GD_{X_1}, GD_{X_2})$ ranges from 0 to P . A higher value indicates greater similarity in geometry. In this study, we calculate the geometric similarity using the eigenvectors corresponding to the top five eigenvalues.

3. Simulating the Global Distribution Locally

In this section, we address a **critical issue**: how to approximate the geometric shape of the global data distribution using local client data while preserving privacy. We then detail how to simulate the ideal global distribution locally in both single-domain and multi-domain scenarios by leveraging the geometric shape of the global distribution.

3.1. Computation of Global Geometric Shapes Under Privacy Constraints

Consider a classification task with K clients and C classes. Suppose each client contains $n_1^i, n_2^i, \dots, n_K^i$ samples belonging to class i , represented by the sample set $\{x_k^1, x_k^2, \dots, x_k^{n_k^i}\}$ for client k . The global distribution for class i is formed by the combined distributions of these local samples across the K clients. However, we cannot directly access the global distribution to derive its geometric shape. Recalling Section 2, the covariance matrix of a distribution is fundamental for obtaining its geometric shape. Thus, our goal is to approximate the global data distribution's covariance matrix without sharing data.

We propose to approximate the global covariance matrix by leveraging the covariance matrices of all local data. Taking class i as an example, we first compute the local covariance matrix Σ_k^i and local mean $\mu_k^i (k = 1, \dots, K)$ for class i on each client as follows:

$$\mu_k^i = \frac{1}{n_k^i} \sum_{j=1}^{n_k^i} x_k^{i,j}, \quad \Sigma_k^i = \frac{1}{n_k^i} \sum_{j=1}^{n_k^i} (x_k^{i,j} - \mu_k^i)(x_k^{i,j} - \mu_k^i)^T.$$

The global covariance matrix is the covariance of the com-

bined data from all clients, defined as:

$$\mu_i = \frac{1}{N_i} \sum_{k=1}^K \sum_{j=1}^{n_k^i} x_k^{i,j} = \frac{1}{N_i} \sum_{k=1}^K n_k^i \mu_k^i,$$

$$\Sigma_i = \frac{1}{N_i} \sum_{k=1}^K \sum_{j=1}^{n_k^i} (x_k^{i,j} - \mu_i)(x_k^{i,j} - \mu_i)^T,$$

where $N_i = \sum_{k=1}^K n_k^i$ represents the total number of samples in class i across all clients. We can decompose $(x_k^{i,j} - \mu_i)$ as $(x_k^{i,j} - \mu_k^i) + (\mu_k^i - \mu_i)$, so the global covariance matrix Σ_i can be rewritten as:

$$\Sigma_i = \frac{1}{N_i} \sum_{k=1}^K \sum_{j=1}^{n_k^i} [(x_k^{i,j} - \mu_k^i + \mu_k^i - \mu_i)(x_k^{i,j} - \mu_k^i + \mu_k^i - \mu_i)^T].$$

Expanding the above equation yields:

$$\begin{aligned} \Sigma_i = & \frac{1}{N_i} \sum_{k=1}^K \sum_{j=1}^{n_k^i} [(x_k^{i,j} - \mu_k^i)(x_k^{i,j} - \mu_k^i)^T + (x_k^{i,j} - \mu_k^i)(\mu_k^i - \mu_i)^T \\ & + (\mu_k^i - \mu_i)(x_k^{i,j} - \mu_k^i)^T + (\mu_k^i - \mu_i)(\mu_k^i - \mu_i)^T]. \end{aligned}$$

By the properties of covariance matrices, the first term is the local covariance matrix Σ_k^i , and the expected values of the second and third terms are zero. The fourth term can be computed as:

$$\sum_{j=1}^{n_k^i} (\mu_k^i - \mu_i)(\mu_k^i - \mu_i)^T = n_k^i (\mu_k^i - \mu_i)(\mu_k^i - \mu_i)^T.$$

Thus, the global covariance matrix can be obtained by combining the local covariance matrices and local means as:

$$\Sigma_i = \frac{1}{N_i} \left(\sum_{k=1}^K n_k^i \Sigma_k^i + \sum_{k=1}^K n_k^i (\mu_k^i - \mu_i)(\mu_k^i - \mu_i)^T \right). \quad (4)$$

where each client's contribution to Σ_i is determined by the coefficient $\frac{N_k^i}{N_i}$. By performing eigenvalue decomposition on Σ_i , we can derive the geometric shape of the global data distribution for class i without sharing any data. For more details on privacy constraints, please refer to Appendix D.

3.2. Global Geometry-Guided Embedding Uncertainty Representation (GGEUR)

Given K clients and a classification task with C classes, we describe below how geometric knowledge can be used to simulate the ideal global data distribution locally in both single-domain and multi-domain scenarios.

3.2.1 Single-Domain Federated Learning

When the data originates from a single domain, we primarily address the issue of label skew. Before federated learning begins, each client uses CLIP to extract p -dimensional embeddings of its local data. Each client then locally computes the per-class local covariance matrices and local sample means, which are subsequently uploaded to the server.

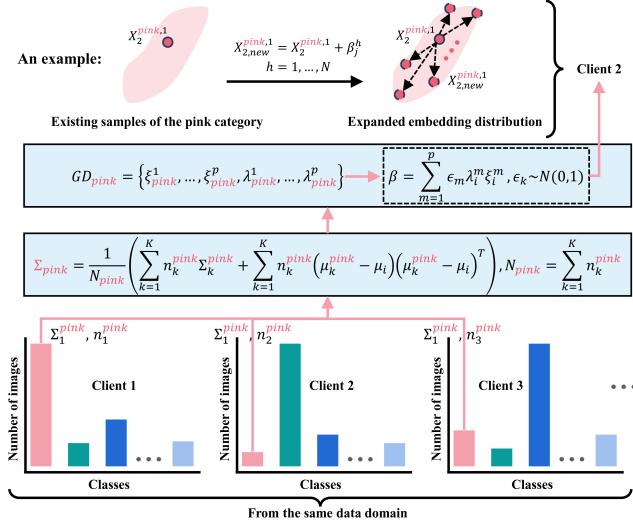


Figure 3. In a single-domain scenario, each class’s global geometric shape is used to guide sample augmentation on each client. The example shows how new samples are generated for Client 2.

For each class, the server generates a global covariance matrix, $\Sigma_1, \Sigma_2, \dots, \Sigma_C$, using Equation (4). Next, each class’s global geometric shape is quantified based on Equation (2), and these global geometric shapes (i.e., eigenvalues and eigenvectors) are sent to each client.

We propose a Global Geometry-Guided Embedding Uncertainty Representation (**GGEUR**) for locally augmenting samples (i.e., embeddings) to simulate the global distribution. Specifically, for class i , its global geometric shape can be represented as $GD_i = \{\xi_i^1, \dots, \xi_i^p, \lambda_i^1, \dots, \lambda_i^p\}$. Suppose client k has only a limited number of class i embeddings, denoted as $X_k^i = [X_k^{(i,1)}, \dots, X_k^{(i,n_k^i)}] \in \mathbb{R}^{p \times n_k^i}$, where n_k^i represents the sample count. As shown in Figure 3, we first perform $n_k^i \times N$ random linear combinations of the eigenvectors ξ_i^1, \dots, ξ_i^p to produce $n_k^i \times N$ different vectors $\beta = \sum_{m=1}^p \epsilon_m \lambda_i^m \xi_i^m \in \mathbb{R}^p$, where ϵ_j follows a standard Gaussian distribution $N(0, 1)$ and is scaled by the eigenvalues λ_i^j to control the magnitude. These vectors provide the direction and range for simulating the global distribution. Next, for each existing sample $X_k^{(i,1)}, \dots, X_k^{(i,n_k^i)}$, we apply N of these vectors to generate $n_k^i \times N$ new samples as follows:

$$X_{(k,h)}^{(i,j)} = X_k^{(i,j)} + \beta_j^h, j = 1, \dots, n_k^i, h = 1, \dots, N. \quad (5)$$

The total number of new samples can be set based on task requirements; in this study, for each local class, we ensure the total number of new and existing samples reaches 2000. Figure 3 and Algorithm 1 illustrate this process in detail. After embedding space augmentation, each client trains an MLP as the local model.

Algorithm 1 GGEUR (Single-Domain Scenario)

Require: $X_k^i = [X_k^{(i,1)}, \dots, X_k^{(i,n_k^i)}] \in \mathbb{R}^{p \times n_k^i}$: Sample set of class i at client k , $GD_i = \{\xi_i^1, \dots, \xi_i^p, \lambda_i^1, \dots, \lambda_i^p\}$: Global geometric shape (eigenvectors and eigenvalues) of class i , N : Number of new samples to generate per original sample

Ensure: X_{new}^i : Augmented sample set of class i at client k

```

1: function GGEUR( $X, GD_i, N$ )
2:    $X_{\text{gen}} \leftarrow \emptyset$  ▷ Initialize generated samples
3:   for  $h = 1$  to  $N$  do
4:      $\beta^h \leftarrow \sum_{m=1}^p \epsilon_m \lambda_i^m \xi_i^m, \epsilon_m \sim \mathcal{N}(0, 1)$  ▷ Generate new vector
5:      $X_{\text{gen}} \leftarrow X_{\text{gen}} \cup \{X + \beta^h\}$  ▷ Generate and add new sample, Equation (5)
6:   end for
7:   return  $X_{\text{gen}}$ 
8: end function
11:  $X_{\text{new}}^i \leftarrow \emptyset$  ▷ Initialize augmented sample set
12: for  $j = 1$  to  $n_k^i$  do
13:    $X_{\text{new}}^i \leftarrow X_{\text{new}}^i \cup \text{GGEUR}(X_k^{(i,j)}, GD_i, N)$ 
14: end for
15: return  $X_{\text{new}}^i$ 

```

3.2.2 Multi-Domain Federated Learning

When data originates from multiple domains, we encounter a complex scenario where both label skew and domain skew exist, as illustrated in Figure 4. In this case, different clients possess data from different domains. For clarity, we introduce two new concepts: the single-domain global distribution and the multi-domain global distribution.

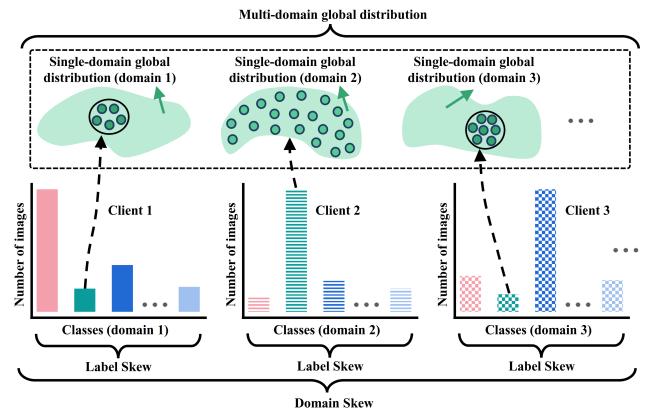


Figure 4. Federated scenario with both label skew and domain skew. Different textures represent data from distinct domains.

Definition 1 (Single-Domain Global Distribution)

Within a particular data domain, when a category’s sample count is sufficient and diverse enough to fully represent the category, this distribution of samples is termed the

single-domain global distribution for that category. For example, in Figure 4, client 2’s *green samples* are numerous and adequately cover the global distribution of the *green category* within domain 2.

Definition 2 (Multi-Domain Global Distribution) For a given category, the combined single-domain global distributions from all data domains constitute the **multi-domain global distribution** of that category. As shown in Figure 4, the *green distributions* across all domains together form the multi-domain global distribution for the *green category*.

Clearly, addressing both label skew and domain skew requires two steps: (1) Local sample augmentation to simulate the single-domain global distribution within each client’s data domain. For example, in Figure 4, client 1 lacks sufficient green samples to represent the single-domain global distribution and thus requires sample augmentation. (2) Generation of new samples on each client to simulate the single-domain global distributions of other data domains.

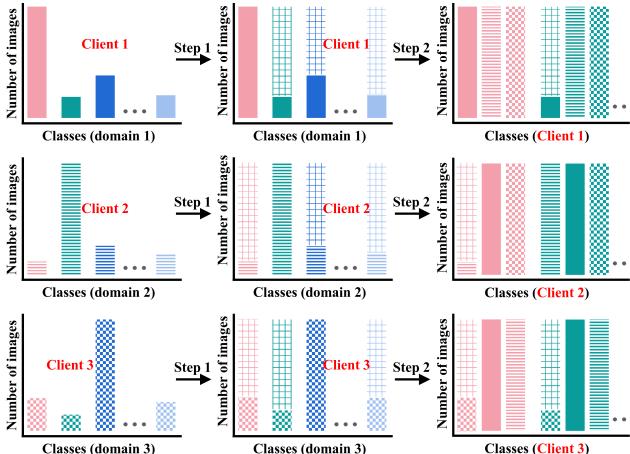


Figure 5. Example with 3 clients: Step 1 generates samples for each client from its own domain, while Step 2 simulates samples from other domains for each client.

As depicted in Figure 5, Step 1 addresses the label skew on each client, while Step 2 addresses domain skew. In summary, our objective is to simulate the ideal multi-domain global distribution for each category locally on each client. Unlike single-domain federated learning, completing Step 1 in this context poses the challenge of obtaining the geometric shape of the single-domain global distribution. For instance, in Figure 5, client 1 has only a small number of green samples, making it difficult to access the single-domain global distribution for the green category. Fortunately, we found that when embeddings are extracted using CLIP, the geometric shapes of multiple single-domain global distributions corresponding to the same category across different domains are similar. This finding implies

that we can identify a cross-domain shared geometric shape for each category, representing its global geometric shape.

Specifically, we investigate this using the Digits dataset, which includes four digit-recognition datasets (i.e., MNIST [21], USPS [14], SVHN [43], and SYN [49]), each representing a different domain. First, we use CLIP (ViT-B/16) [47] to extract image embeddings for all categories across the four datasets. Then, we compute the similarity between the geometric shapes of the embedding distributions for each category in MNIST and the corresponding categories in USPS, SVHN, and SYN, as shown in Figure 6. The significantly higher values in the diagonal of each heatmap indicate that, across domains, the geometric shapes of embeddings for the same category exhibit a consistent pattern, rather than being randomly distributed. This allows us to simplify the multi-domain label skew problem in Step 1 to a label skew problem in the single-domain scenario.

In detail, we first extract embeddings using CLIP and then locally calculate the covariance matrix and mean of local samples, which are uploaded to the server. Given CLIP’s cross-domain consistency in representing the same class, we continue to use Equation (4) on the server to produce a shared global covariance matrix $\Sigma_1, \Sigma_2, \dots, \Sigma_C$ for each class. We further obtain the shared global geometric shapes, which are then sent to clients. Each client then applies GGEUR to augment samples, thereby simulating the single-domain global distribution locally. Following this, clients also need to generate new samples to simulate the single-domain global distribution from other data domains. Although the geometric shapes of embedding distributions across different domains are similar, the positions of these distributions differ. We propose transferring the class prototypes (local sample means) from other domains to the local client, and applying GGEUR to these prototypes to generate new samples, thus simulating the multi-domain global distribution. Algorithm 2 details this process. For each client, in Step 1, we ensure that the total number of new and existing samples for each class is 500. Step 2 generates 500 new samples based on each prototype separately.

3.3. Comparison with Analogous Methods

Federated data augmentation aim to bridge the gap between local data distributions and the ideal global distribution by generating more diverse data samples on clients [9, 37, 51]. For example, FedMix [63] and FEDGEN [68] use MixUp and its variants to augment client data, thus mitigating label skew. However, due to the lack of knowledge-based guidance, these methods largely rely on the diversity of local data. FedFA [67] assumes local data follows a Gaussian distribution and generates new samples centered on class prototypes. Nevertheless, the Gaussian assumption is overly idealistic [38], limiting the ability of generated samples to adequately reduce the discrepancy between local and global

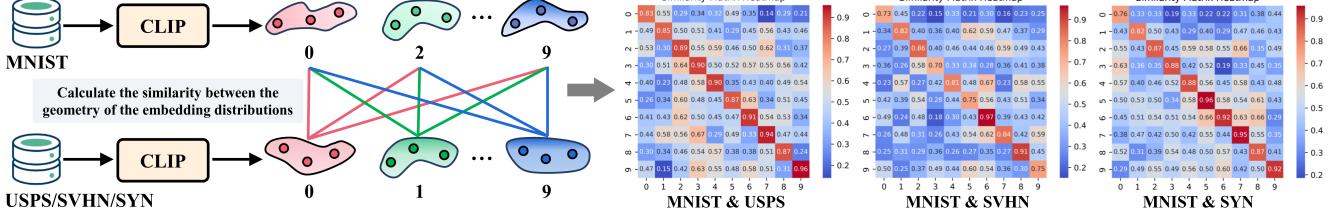


Figure 6. Study of geometric shape similarity on the Digits dataset.

distributions. **Compared to** the Gaussian assumption, the geometric shapes proposed in this work provide a more accurate description of embedding distributions. GGEUR estimates the geometric shape of the global distribution without compromising privacy and leverages it to guide data augmentation on clients. The introduction of additional geometric knowledge enables our method to effectively reduce the gap between local and global distributions.

4. Experiments

In this section, we conduct a comprehensive evaluation of our method under scenarios with label skew, domain skew, and the coexistence of both.

4.1. Experiment Setup

Label Skew Datasets. We evaluate our method on three single-domain image classification tasks.

- Cifar-10 [20] contains 10 classes, with 50,000 images for training and 10,000 images for validation.
- Cifar-100 [20] covers 100 classes, with 50,000 training images and 10,000 validation images.
- Tiny-ImageNet [6] is the subset of ImageNet with 100K images of size 64×64 with 200 classes scale.

Domain Skew Datasets. We evaluated our method on the multi-domain image classification dataset Digits and conducted analyses on Office-Caltech and PACS.

- Digits [14, 21, 43, 49] includes four domains: MNIST, USPS, SVHN and SYN, each with 10 categories.
- Office-Caltech [8] includes four domains: Caltech, Amazon, Webcam, and DSLR, each with 10 categories.
- PACS [24] includes four domains: Photo (P) with 1,670 images, Art Painting (AP) with 2,048 images, Cartoon (Ct) with 2,344 images and Sketch (Sk) with 3,929 images. Each domain holds seven categories.

Dataset with Coexisting Label Skew and Domain Skew. Office-Home [56] includes 4 domains: Art (A), Clipart (C), Product (P), and Real World (R), each containing 65 classes. To increase the challenge, we designed a new partitioning method for the multi-domain dataset Office-Home to create a scenario where label skew and domain skew coexist. For details, please refer to Appendix B. We name the newly constructed dataset **Office-Home-LDS (Label and Domain Skew)**. Figure 7 shows the data distribution of

Office-Home-LDS with $\beta = 0.1$ and $\beta = 0.5$.

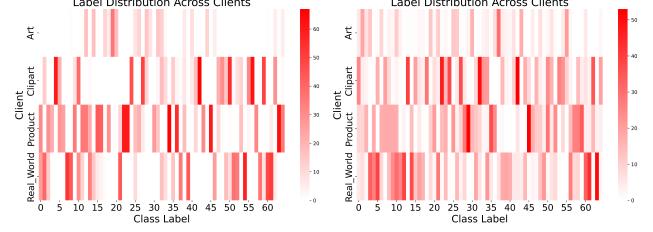


Figure 7. Number of samples per class across four clients when β equals 0.1 and 0.5, with **each client holding data from a different domain**. Additional cases are provided in Appendix B.

Implementation Details and Comparison Methods. On the label skew dataset, we applied GGEUR to FedAvg [39] using CLIP (ViT-B/16) [47] as the backbone network and compared it with state-of-the-art federated data augmentation methods in heterogeneous federated learning, including FedMix [63], FEDGEN [68], and FedFA [67]. Additionally, we compared it with other advanced methods that use CLIP (ViT-B/16) as the backbone, such as FedTPG [45] and FedCLIP [35]. On the domain skew and Office-Home-LDS datasets, we focused on exploring the enhancement effect of GGEUR across various federated architectures. Therefore, we applied GGEUR to FedAvg, SCAFFOLD [16], MOON [25], FedDyn [1], FedOPT [48], FedProto [55], and Fed-NTD [22], all of which use CLIP for image feature extraction and a **single-layer MLP** as the local model. All experimental settings for FL methods are consistent with the latest benchmark [13]; details are provided in Appendix C.

Evaluation Metrics. Following [27], we utilize the Top-1 (%) accuracy and the standard deviation of accuracy across multi-domains as evaluation metrics. A smaller standard deviation indicates better Performance Fairness across different domains. We use the average results from the last five rounds accuracy and variance as the final performance.

4.2. Evaluation Results on Label Skew Dataset

Main Results. Tables 1, 2, and 3 show the performance improvement of GGEUR over FedAvg (CLIP+MLP) [39] under different β values, along with comparison results with other methods. Under all label distribution skew settings, GGEUR significantly outperforms other methods, demonstrating superior classification accuracy and adaptability to imbalanced label distributions. Particularly at lower β

Table 1. Comparison results on CIFAR-100 and Tiny-ImageNet datasets with different degrees of label skew (β values). The best results are shown in **underlined bold**. FedAvg (CLIP+MLP) indicates that the backbone network uses CLIP and MLP, and the federated learning method employs FedAvg.

Methods	CIFAR-100			Tiny-ImageNet		
	0.5	0.3	0.1	0.5	0.3	0.1
Zero-Shot CLIP	64.87			63.67		
FedTPG	71.40	70.95	68.63	67.63	66.72	64.71
FedCLIP	72.03	71.20	70.64	70.41	70.37	69.50
FedMix (CLIP+MLP)	81.31	79.62	73.85	70.89	68.57	63.43
FEDGEN (CLIP+MLP)	81.24	78.97	73.15	72.37	70.35	64.16
FedFA (CLIP+MLP)	81.98	79.31	74.68	70.41	70.68	64.62
FedAvg (CLIP+MLP)	81.41	77.68	68.22	70.08	67.65	60.10
+ GGEUR	83.31	81.65	77.70	73.89	72.19	66.86

Table 2. Evaluation results of GGEUR on CIFAR-10 and CIFAR-100 with more severe label skew (i.e., smaller β values).

Methods	CIFAR-10				
	0.01	0.03	0.05	0.07	0.09
FedAvg (CLIP+MLP)	90.87	90.13	91.96	92.05	91.82
+ GGEUR	94.39	94.25	95.07	95.21	95.38
Methods	CIFAR-100				
	0.01	0.03	0.05	0.07	0.09
FedAvg (CLIP+MLP)	58.71	60.77	62.32	61.69	66.51
+ GGEUR	75.72	75.40	75.96	76.72	78.00

Table 3. Evaluation results of GGEUR on Tiny-ImageNet with more severe label skew (i.e., smaller β values).

Methods	Tiny-ImageNet				
	0.01	0.03	0.05	0.07	0.09
FedAvg (CLIP+MLP)	53.03	54.57	58.91	58.77	59.13
+ GGEUR	64.27	65.79	66.49	66.34	66.85

values (i.e., with more severe label skew), GGEUR notably enhances the performance of FedAvg (CLIP+MLP) on CIFAR-100, validating its effectiveness in handling extreme label skew scenarios. For example, when β is 0.01, 0.03, and 0.05, GGEUR achieves performance gains of **17.01%**, **14.63%**, and **13.64%**, respectively. These results indicate that GGEUR not only exhibits good generalization ability in standard settings but also maintains robust performance in cases of significant label skew. Moreover, our large-scale client experiments further confirm the scalability and consistent effectiveness of GGEUR under increased client participation, details are provided in Appendix E.

Comparison with Peer Methods. We implemented FedMix, FEDGEN, and FedFA for comparison in Table 1, and in all cases, our method outperformed these three approaches. This demonstrates the superiority of geometric shapes as knowledge over the Gaussian assumption and traditional data augmentation methods.

4.3. Evaluation Results on Domain Skew Dataset

Ablation Study. As shown in Algorithm 2, simulating the multi-domain global distribution requires two steps.

Table 4. Evaluation results on the Digits dataset. GGEUR by default includes both Step 1 and Step 2.

Methods	Digits					
	MNIST	USPS	SVHN	SYN	AVG	STD ↓
FedMix (CLIP+MLP)	95.03	90.25	57.50	72.60	78.85	14.89
FEDGEN (CLIP+MLP)	95.85	92.52	58.77	73.62	80.19	14.99
FedFA (CLIP+MLP)	96.68	92.97	57.87	75.53	80.76	15.44
FedAvg [39]	90.40	60.30	34.68	46.99	58.09	20.74
FedAvg (CLIP+MLP)	95.12	89.74	56.36	65.17	76.60	16.25
+ GGEUR (Step 1)	96.02	93.02	58.55	73.13	80.18	15.28
+ GGEUR (Step 1 & 2)	97.05	94.12	63.54	74.73	82.36	13.84
SCAFFOLD [16]	97.79	94.45	26.64	90.69	77.39	29.41
SCAFFOLD (CLIP+MLP)	94.62	90.08	54.33	68.71	76.93	16.31
+ GGEUR	95.91	92.08	63.25	71.54	80.70	13.69
MOON [25]	92.78	68.11	33.36	39.28	58.36	23.82
MOON (CLIP+MLP)	75.64	73.09	38.83	52.74	60.07	15.14
+ GGEUR (Step 1)	84.64	81.96	43.04	60.35	67.50	16.97
+ GGEUR (Step 1 & 2)	95.16	91.13	55.23	71.00	78.13	16.08
FedDyn [1]	88.91	60.34	34.57	50.72	58.65	19.76
FedDyn (CLIP+MLP)	95.46	92.13	58.89	70.30	79.19	15.19
+ GGEUR	97.07	94.02	63.34	74.83	82.31	13.88
FedOPT [48]	92.71	87.62	31.32	87.92	74.89	25.38
FedOPT (CLIP+MLP)	94.57	88.79	58.65	66.47	77.12	14.96
+ GGEUR	96.43	93.47	62.35	70.75	80.75	14.54
FedProto [55]	90.54	89.54	34.61	58.00	68.18	23.38
FedProto (CLIP+MLP)	94.86	92.63	54.29	65.52	76.83	17.40
+ GGEUR	97.19	94.12	63.70	73.83	82.21	13.96
FedNTD [22]	52.31	58.07	18.03	97.29	56.43	28.12
FedNTD (CLIP+MLP)	95.82	91.43	58.26	69.95	78.86	15.41
+ GGEUR	97.08	94.32	63.57	73.53	82.13	14.06

We selected the classic FedAvg (CLIP+MLP) and MOON (CLIP+MLP) for an ablation study on the Digits dataset. The experimental results, highlighted in yellow in Table 4, show that each step incrementally improves the original methods, and their combination achieves the best results.

Main Results. Tables 4, 5, and 6 present the experimental results on datasets with domain skews. It can be observed that simply using CLIP [47] for image representation, combined with a single-layer MLP for federated learning, already surpasses existing methods. This improvement is attributed to the advancements in the foundation model, and GGEUR can further significantly enhance the performance of the global model. For instance, on the Digits dataset, GGEUR improves the average performance of FedAvg (CLIP+MLP) [39], MOON (CLIP+MLP) [25], and FedProto (CLIP+MLP) [55] by **5.76%**, **18.06%**, and **5.38%**, respectively. Additionally, compared to other methods, GGEUR significantly reduces the accuracy variance across different domains. This demonstrates that GGEUR effectively adapts to features from different domains when handling cross-domain data distribution disparities, providing more robust and fair performance. We also conducted experiments to measure the computational cost of GGEUR, which demonstrated that the additional overhead introduced by our method is minimal. Further details can be found in Appendix F.

Table 5. Evaluation results on Office-Caltech.

Methods	Office-Caltech					
	Am	Ca	D	W	AVG ↑	STD ↓
FedAvg [39]	81.99	73.21	79.37	67.93	75.62	6.31
FedAvg (CLIP+MLP)	98.26	96.74	100	100	98.75	1.57
SCAFFOLD [16]	39.77	42.50	78.02	70.69	57.75	19.44
SCAFFOLD (CLIP+MLP)	96.18	92.88	95.83	93.26	94.54	1.70
MOON [25]	84.42	75.98	84.67	68.97	78.51	7.53
MOON (CLIP+MLP)	98.61	97.33	100	98.88	98.70	1.09
FedDyn [1]	84.02	72.59	77.34	68.97	75.72	6.50
FedDyn (CLIP+MLP)	98.61	96.74	100	100	98.84	1.54
FedOPT [48]	79.05	71.96	89.34	74.48	78.71	7.67
FedOPT (CLIP+MLP)	98.26	97.33	100	100	98.90	1.33
FedProto [55]	87.79	75.98	90.00	79.31	83.27	6.70
FedProto (CLIP+MLP)	98.26	96.74	100	100	98.75	1.57
FedNTD [22]	10.95	10.89	14.67	10.34	11.71	1.99
FedNTD (CLIP+MLP)	97.92	96.14	100	100	98.51	1.86

Table 6. Evaluation results on PACS.

Methods	PACS					
	P	AP	Ct	Sk	AVG ↑ STD ↓	
FedAvg [39]	76.09	64.19	83.50	89.40	78.30	9.41
FedAvg (CLIP+MLP)	99.40	98.37	99.01	93.64	97.60	2.32
SCAFFOLD [16]	61.95	45.44	58.87	54.64	55.25	6.22
SCAFFOLD (CLIP+MLP)	92.42	81.63	80.68	87.28	85.50	4.72
MOON [25]	74.44	64.19	83.92	89.17	77.93	9.53
MOON (CLIP+MLP)	99.60	99.02	99.43	93.89	97.99	2.37
FedDyn [1]	78.17	63.29	82.27	89.93	78.66	9.70
FedDyn (CLIP+MLP)	99.40	98.37	99.01	93.55	97.58	2.36
FedOPT [48]	78.66	67.66	82.41	83.68	78.12	6.31
FedOPT (CLIP+MLP)	99.40	98.37	99.01	93.64	97.60	2.32
FedProto [55]	85.63	73.69	83.57	91.14	83.51	6.31
FedProto (CLIP+MLP)	99.40	98.21	99.01	94.23	97.71	2.06
FedNTD [22]	16.77	18.23	28.47	93.18	39.16	31.51
FedNTD (CLIP+MLP)	99.40	98.54	99.29	92.28	97.38	2.96

On the Office-Caltech and PACS datasets, since the use of CLIP+MLP as the backbone network already achieves very high performance across all methods, these datasets are less challenging. However, our evaluation results can serve as a reference for other research.

Comparison with Peer Methods. We implemented FedMix (CLIP+MLP), FEDGEN (CLIP+MLP), and FedFA (CLIP+MLP) for comparison in Table 4. When GGEUR is applied to FedAvg (CLIP+MLP), it outperforms these three methods by **3.51%**, **2.17%**, and **1.60%**, respectively.

4.4. Evaluation Results on Office-Home-LDS

To fully demonstrate the potential of GGEUR, we constructed a more challenging dataset, Office-Home-LDS, which incorporates both label skew and domain skew. Office-Home-LDS simulates a more realistic and complex scenario of distributional imbalance, encompassing cross-domain data distribution differences and label imbalance. In Table 7, we show the remarkable performance gains of GGEUR over existing methods on this dataset, along with significant reductions in accuracy variance across different domains, highlighting its effectiveness in handling highly heterogeneous data scenarios. Specif-

Table 7. Evaluation results on Office-Home-LDS ($\beta = 0.1$).

Methods	Office-Home-LDS					
	A	C	P	R	AVG ↑ STD ↓	
FedAvg (CLIP+MLP) [39]	65.29	58.17	80.56	76.53	70.14	8.89
+ GGEUR	78.33	79.01	90.17	88.46	83.99	5.36
SCAFFOLD (CLIP+MLP) [16]	68.72	66.79	83.63	80.12	74.82	7.20
+ GGEUR	78.60	78.32	89.86	89.07	83.96	5.51
MOON (CLIP+MLP) [25]	69.27	68.63	86.56	82.87	76.83	7.99
+ GGEUR	72.02	70.31	86.11	83.87	78.08	6.98
FedDyn (CLIP+MLP) [1]	58.30	55.19	77.63	72.86	65.99	9.47
+ GGEUR	78.88	78.55	90.47	88.46	84.09	5.42
FedOPT (CLIP+MLP) [48]	58.44	54.89	76.80	72.25	65.59	9.16
+ GGEUR	79.01	78.32	90.84	88.61	84.20	5.59
FedProto (CLIP+MLP) [55]	65.84	56.49	80.41	74.85	69.40	9.09
+ GGEUR	78.05	77.71	89.79	87.84	83.35	5.51
FedNTD (CLIP+MLP) [22]	69.68	66.64	84.53	80.96	75.46	7.48
+ GGEUR	78.19	74.66	90.24	86.77	82.46	6.29

ically, GGEUR improved the average performance of FedAvg (CLIP+MLP) [39], SCAFFOLD (CLIP+MLP) [16], FedDyn (CLIP+MLP) [1], FedOPT (CLIP+MLP) [48], FedProto (CLIP+MLP) [55], and FedNTD (CLIP+MLP) [22] by **13.85%**, **9.14%**, **18.1%**, **18.61%**, **13.95%**, and **7.0%**, respectively. These results demonstrate that GGEUR can achieve robust performance on complex multi-domain, multi-class datasets, making it an effective approach for addressing both label and domain skew.

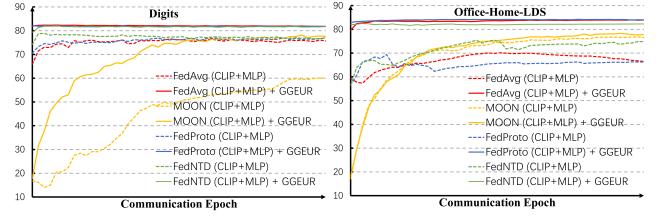


Figure 8. Comparison of convergence in average accuracy with and without integrating GGEUR in the selected FL methods.

4.5. GGEUR Accelerates Convergence

In Figure 8, we plot the average accuracy per epoch for various FL methods with and without using GGEUR. It can be observed that our method enhances the convergence speed of the FL methods, resulting in smoother curves.

5. Conclusion

This work captures the global geometric shape of the embedding distribution while preserving privacy and utilizes it to simulate an ideal global distribution on clients, bridging the gap between local and global distributions. We first introduce the concept of distributional geometric shapes and address the challenge of obtaining the global geometric shape under privacy constraints. Then, we propose leveraging the global geometric shape to assist clients in augmenting their samples. Extensive experiments demonstrate the effectiveness and compatibility of GGEUR.

Acknowledgement. This research is supported by the National Natural Science Foundation of China (Grants 623B2080).

References

- [1] Durmus Alp Emre Acar, Yue Zhao, Ramon Matas Navarro, Matthew Mattina, Paul N Whatmough, and Venkatesh Saligrama. Federated learning based on dynamic regularization. *arXiv preprint arXiv:2111.04263*, 2021. [1](#), [6](#), [7](#), [8](#)
- [2] Maruan Al-Shedivat, Jennifer Gillenwater, Eric Xing, and Afshin Rostamizadeh. Federated learning via posterior averaging: A new perspective and practical algorithms. *arXiv preprint arXiv:2010.05273*, 2020. [1](#)
- [3] Dengsheng Chen, Jie Hu, Vince Junkai Tan, Xiaoming Wei, and Enhua Wu. Elastic aggregation for federated optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12187–12197, 2023. [1](#), [12](#)
- [4] Hong-You Chen and Wei-Lun Chao. Fedbe: Making bayesian model ensemble applicable to federated learning. *arXiv preprint arXiv:2009.01974*, 2020. [12](#)
- [5] Yuhang Chen, Wenke Huang, and Mang Ye. Fair federated learning under domain skew with local consistency and domain diversity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12077–12086, 2024. [1](#)
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [6](#), [12](#)
- [7] Liang Gao, Huazhu Fu, Li Li, Yingwen Chen, Ming Xu, and Cheng-Zhong Xu. Feddc: Federated learning with non-iid data via local drift decoupling and correction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10112–10121, 2022. [1](#)
- [8] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *2012 IEEE conference on computer vision and pattern recognition*, pages 2066–2073. IEEE, 2012. [6](#), [12](#)
- [9] Weituo Hao, Mostafa El-Khamy, Jungwon Lee, Jianyi Zhang, Kevin J Liang, Changyou Chen, and Lawrence Carin Duke. Towards fair federated learning with zero-shot data augmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3310–3319, 2021. [5](#), [12](#)
- [10] Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of large cnns at the edge. *Advances in Neural Information Processing Systems*, 33:14068–14080, 2020. [12](#)
- [11] Wenke Huang, Mang Ye, and Bo Du. Learn from others and be yourself in heterogeneous federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10143–10153, 2022. [1](#), [12](#)
- [12] Wenke Huang, Mang Ye, Zekun Shi, He Li, and Bo Du. Rethinking federated learning with domain shift: A prototype view. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16312–16322. IEEE, 2023. [1](#), [12](#)
- [13] Wenke Huang, Mang Ye, Zekun Shi, Guancheng Wan, He Li, Bo Du, and Qiang Yang. Federated learning for generalization, robustness, fairness: A survey and benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. [1](#), [2](#), [6](#), [12](#), [13](#)
- [14] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE TPAMI*, pages 550–554, 1994. [5](#), [6](#), [12](#)
- [15] Xuefeng Jiang, Sheng Sun, Yuwei Wang, and Min Liu. Towards federated learning against noisy labels via local self-regularization. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 862–873, 2022. [1](#)
- [16] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for on-device federated learning. *arXiv preprint arXiv:1910.06378*, 2(6), 2019. [6](#), [7](#), [8](#), [12](#)
- [17] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning*, pages 5132–5143. PMLR, 2020. [1](#)
- [18] Jinkyu Kim, Geeho Kim, and Bohyung Han. Multi-level branched regularization for federated learning. In *International Conference on Machine Learning*, pages 11058–11073. PMLR, 2022. [12](#)
- [19] Jakub Konečný. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016. [1](#)
- [20] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. [6](#), [12](#)
- [21] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, pages 2278–2324, 1998. [5](#), [6](#), [12](#)
- [22] Gihun Lee, Minchan Jeong, Yongjin Shin, Sangmin Bae, and Se-Young Yun. Preservation of the global knowledge by not-true distillation in federated learning. *Advances in Neural Information Processing Systems*, 35:38461–38474, 2022. [1](#), [6](#), [7](#), [8](#), [12](#)
- [23] Bo Li, Mikkel N Schmidt, Tommy S Alstrøm, and Sebastian U Stich. Partial variance reduction improves non-convex federated learning on heterogeneous data. *arXiv preprint arXiv:2212.02191*, 2022. [1](#)
- [24] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550, 2017. [6](#), [12](#)
- [25] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10713–10722, 2021. [6](#), [7](#), [8](#), [12](#)
- [26] Qinbin Li, Bingsheng He, and Dawn Song. Adversarial collaborative learning on non-iid features. In *International Conference on Machine Learning*, pages 19504–19526. PMLR, 2023. [12](#)

- [27] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019. 6
- [28] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60, 2020. 1
- [29] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020. 1, 12
- [30] Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning on non-iid features via local batch normalization. *arXiv preprint arXiv:2102.07623*, 2021. 1
- [31] Xin-Chun Li and De-Chuan Zhan. Fedrs: Federated learning with restricted softmax for label distribution non-iid data. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 995–1005, 2021. 12
- [32] Zexi Li, Tao Lin, Xinyi Shang, and Chao Wu. Revisiting weighted aggregation in federated learning with neural networks. In *International Conference on Machine Learning*, pages 19767–19788. PMLR, 2023. 1
- [33] Zexi Li, Xinyi Shang, Rui He, Tao Lin, and Chao Wu. No fear of classifier biases: Neural collapse inspired federated learning with synthetic and fixed classifier. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5319–5329, 2023. 1
- [34] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in neural information processing systems*, 33:2351–2363, 2020. 12
- [35] Wang Lu, Xixu Hu, Jindong Wang, and Xing Xie. Fedclip: Fast generalization and personalization for clip in federated learning. *arXiv preprint arXiv:2302.13485*, 2023. 6
- [36] Yanbiao Ma, Licheng Jiao, Fang Liu, Yuxin Li, Shuyuan Yang, and Xu Liu. Delving into semantic scale imbalance. In *The Eleventh International Conference on Learning Representations*, 2023. 2
- [37] Yanbiao Ma, Licheng Jiao, Fang Liu, Shuyuan Yang, Xu Liu, and Puhua Chen. Feature distribution representation learning based on knowledge transfer for long-tailed classification. *IEEE Transactions on Multimedia*, 2023. 5, 12
- [38] Yanbiao Ma, Licheng Jiao, Fang Liu, Shuyuan Yang, Xu Liu, and Puhua Chen. Geometric prior guided feature representation learning for long-tailed classification. *International Journal of Computer Vision*, pages 1–18, 2024. 5, 12
- [39] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 6, 7, 8, 12
- [40] Matias Mendieta, Taojiannan Yang, Pu Wang, Minwoo Lee, Zhengming Ding, and Chen Chen. Local learning matters: Rethinking data heterogeneity in federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8397–8406, 2022. 12
- [41] Jiaxu Miao, Zongxin Yang, Leilei Fan, and Yi Yang. Fedseg: Class-heterogeneous federated learning for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8042–8052, 2023. 12
- [42] Xutong Mu, Yulong Shen, Ke Cheng, Xueli Geng, Jiaxuan Fu, Tao Zhang, and Zhiwei Zhang. Fedproc: Prototypical contrastive federated learning on non-iid data. *Future Generation Computer Systems*, 143:93–104, 2023. 1
- [43] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 5, 6, 12
- [44] Chamath Palihawadana, Nirmalie Wiratunga, Anjana Wijekoon, and Harsha Kalutarage. Fedsim: Similarity guided model aggregation for federated learning. *Neurocomputing*, 483:432–445, 2022. 1
- [45] Chen Qiu, Xingyu Li, Chaithanya Kumar Mummadi, Madan Ravi Ganesh, Zhenzhen Li, Lu Peng, and Wan-Yi Lin. Federated text-driven prompt generation for vision-language models. In *The Twelfth International Conference on Learning Representations*, 2024. 6
- [46] Zhe Qu, Xingyu Li, Rui Duan, Yao Liu, Bo Tang, and Zhuo Lu. Generalized federated learning via sharpness aware minimization. In *International conference on machine learning*, pages 18250–18280. PMLR, 2022. 12
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 5, 6, 7
- [48] Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive federated optimization. In *International Conference on Learning Representations*, 2021. 6, 7, 8, 12
- [49] Prasun Roy, Subhankar Ghosh, Saumik Bhattacharya, and Umapada Pal. Effects of degradations on deep neural network architectures. *arXiv preprint arXiv:1807.10108*, 2018. 5, 6, 12
- [50] Yujun Shi, Jian Liang, Wenqing Zhang, Vincent YF Tan, and Song Bai. Towards understanding and mitigating dimensional collapse in heterogeneous federated learning. *arXiv preprint arXiv:2210.00226*, 2022. 12
- [51] MyungJae Shin, Chihoon Hwang, Joongheon Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Xor mixup: Privacy-preserving data augmentation for one-shot federated learning. *arXiv preprint arXiv:2006.05148*, 2020. 5, 12
- [52] Sidak Pal Singh and Martin Jaggi. Model fusion via optimal transport. *Advances in Neural Information Processing Systems*, 33:22045–22055, 2020. 1
- [53] Jun Sun, Tianyi Chen, Georgios B Giannakis, Qinmin Yang, and Zaiyue Yang. Lazily aggregated quantized gradient innovation for communication-efficient federated learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(4):2031–2044, 2020. 1

- [54] Qiheng Sun, Xiang Li, Jiayao Zhang, Li Xiong, Weiran Liu, Jinfei Liu, Zhan Qin, and Kui Ren. Shapleyfl: Robust federated learning based on shapley value. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2096–2108, 2023. 1
- [55] Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fedproto: Federated prototype learning across heterogeneous clients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8432–8440, 2022. 6, 7, 8
- [56] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5018–5027, 2017. 6, 12
- [57] Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Paliopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. *arXiv preprint arXiv:2002.06440*, 2020. 1
- [58] Haozhao Wang, Yichen Li, Wencho Xu, Ruixuan Li, Yufeng Zhan, and Zhigang Zeng. Dafkd: Domain-aware federated knowledge distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20412–20421, 2023. 12
- [59] Nannan Wu, Li Yu, Xin Yang, Kwang-Ting Cheng, and Zengqiang Yan. Fediic: Towards robust federated learning for class-imbalanced medical image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 692–702. Springer, 2023. 1
- [60] Yuezhou Wu, Yan Kang, Jiahuan Luo, Yuanqin He, and Qiang Yang. Fedcg: Leverage conditional gan for protecting privacy and maintaining competitive performance in federated learning. *arXiv preprint arXiv:2111.08211*, 2021. 12
- [61] Yingda Xia, Dong Yang, Wenqi Li, Andriy Myronenko, Daguang Xu, Hirofumi Obinata, Hitoshi Mori, Peng An, Stephanie Harmon, Evrim Turkbey, et al. Auto-fedavg: learnable federated averaging for multi-institutional medical image segmentation. *arXiv preprint arXiv:2104.10195*, 2021. 1
- [62] Mang Ye, Xiuwen Fang, Bo Du, Pong C Yuen, and Dacheng Tao. Heterogeneous federated learning: State-of-the-art and research challenges. *ACM Computing Surveys*, 56(3):1–44, 2023. 1
- [63] Tehrim Yoon, Sumin Shin, Sung Ju Hwang, and Eunho Yang. Fedmix: Approximation of mixup under mean augmented federated learning. *arXiv preprint arXiv:2107.00233*, 2021. 2, 5, 6, 12
- [64] Jinghui Zhang, Xinyu Cheng, Cheng Wang, Yuchen Wang, Zhan Shi, Jiahui Jin, Aibo Song, Wei Zhao, Liangsheng Wen, and Tingting Zhang. Fedada: Fast-convergent adaptive federated learning in heterogeneous mobile edge computing environment. *World Wide Web*, 25(5):1971–1998, 2022. 1
- [65] Jie Zhang, Zhiqi Li, Bo Li, Jianghe Xu, Shuang Wu, Shouhong Ding, and Chao Wu. Federated learning with label distribution skew via logits calibration. In *International Conference on Machine Learning*, pages 26311–26329. PMLR, 2022. 12
- [66] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018. 1
- [67] Tianfei Zhou and Ender Konukoglu. FedFA: Federated feature augmentation. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. 1, 5, 6, 12
- [68] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *International conference on machine learning*, pages 12878–12889. PMLR, 2021. 1, 5, 6, 12

A. Related Work

Federated Learning with Data Heterogeneity

In federated learning, client data often originates from different distributions, typically manifesting as label skew and domain skew. With label skew, the class distributions across clients are significantly imbalanced, while domain skew occurs when feature distributions for the same class vary due to differences in data sources. To address these issues, researchers have proposed methods that can be grouped into client regularization, server-side dynamic aggregation, and federated data augmentation [13].

Client regularization primarily focuses on adjusting local optimization objectives so that local models align more closely with the direction of the global model, reducing distributional shifts across clients [18, 22, 26, 31, 40, 41, 46, 50, 58, 60, 65]. Methods such as FedProx [29] and SCAFFOLD [16] introduce additional regularization terms to minimize the discrepancy between local and global models, thereby improving convergence speed and accuracy. MOON [25] leverages contrastive learning to align feature spaces across clients, addressing both label and domain skew. FPL [12] supervises the learning of local class prototypes by aggregating and sharing class prototypes across clients. However, involving a global model in the local optimization process deeply enlarges the local computation cost and linearly increases with the parameter scale.

Server-side dynamic aggregation methods optimize the global model by adaptively adjusting client weights [4, 10, 39]. FedOPT [48] and Elastic [3] use dynamic aggregation weights based on client model updates, enhancing the global model’s generalization in heterogeneous data settings. Additionally, methods like FedDF [34] and FCCL [11] incorporate knowledge distillation modules on the server side, combined with auxiliary datasets to improve the adaptability of aggregation, making these approaches suitable for broader cross-client data distributions. However, these methods often require additional proxy datasets to support model adjustments, which is beneficial in scenarios with significant distributional differences across clients.

Federated data augmentation aim to bridge the gap between local data distributions and the ideal global distribution by generating more diverse data samples on clients [9, 37, 51]. For example, FedMix [63] and FEDGEN [68] use MixUp and its variants to augment client data, thus mitigating label skew. However, due to the lack of knowledge-based guidance, these methods largely rely on the diversity of local data. FedFA [67] assumes local data follows a Gaussian distribution and generates new samples centered on class prototypes. Nevertheless, the Gaussian assumption is overly idealistic [38], limiting the ability of generated samples to adequately reduce the discrepancy between local and global distributions. **Compared to** the Gaussian

assumption, the geometric shapes proposed in this work provide a more accurate description of embedding distributions. GGEUR estimates the geometric shape of the global distribution without compromising privacy and leverages it to guide data augmentation on clients.

B. Dataset

Label Skew Datasets. We evaluate our method on three single-domain image classification tasks.

- **Cifar-10** [20] contains 10 classes, with 50,000 images for training and 10,000 images for validation.
- **Cifar-100** [20] covers 100 classes, with 50,000 training images and 10,000 validation images.
- **Tiny-ImageNet** [6] is the subset of ImageNet with 100K images of size 64×64 with 200 classes scale.

Consistent with recent benchmarks [13], we set up 10 clients for each task. To simulate label skew across clients, we use a Dirichlet distribution, $Dir(\beta)$, where the parameter $\beta > 0$ controls the degree of label skew (i.e., class imbalance). When β takes a smaller value, the local distributions generated on each client become more skewed, showing greater divergence from the overall distribution.

Domain Skew Datasets. We evaluated our method on the multi-domain image classification dataset Digits and conducted analyses on Office-Caltech and PACS.

- **Digits** [14, 21, 43, 49] includes four domains: MNIST, USPS, SVHN and SYN, each with 10 categories.
- **Office-Caltech** [8] includes four domains: Caltech, Amazon, Webcam, and DSLR, each with 10 categories.
- **PACS** [24] includes four domains: Photo (P) with 1,670 images, Art Painting (AP) with 2,048 images, Cartoon (Ct) with 2,344 images and Sketch (Sk) with 3,929 images. Each domain holds seven categories.

Consistent with recent benchmarks, in domain skew experiments, each domain is assigned to a separate client, with each client focusing on data from a specific domain. For Digits, each client is allocated 10% of the data from its respective domain. For Office-Caltech and PACS, each client is allocated 30% of the data from its corresponding domain.

Dataset with Coexisting Label Skew and Domain Skew.

Office-Home [56] includes 4 domains: Art (A), Clipart (C), Product (P), and Real World (R), each containing 65 classes. To increase the challenge, we designed a new partitioning method for the multi-domain dataset Office-Home to create a scenario where label skew and domain skew coexist.

In the Office-Home dataset, while ensuring that each client corresponds to a single domain, we first generate a Dirichlet coefficient matrix, where the degree of class imbalance is controlled by β . For the 65-class, 4-domain Office-Home task, we generate a 4×65 matrix controlled by β (with each column summing to 1). The four coefficients for each class are then allocated to the four clients, and each client uses its assigned coefficients to determine the number

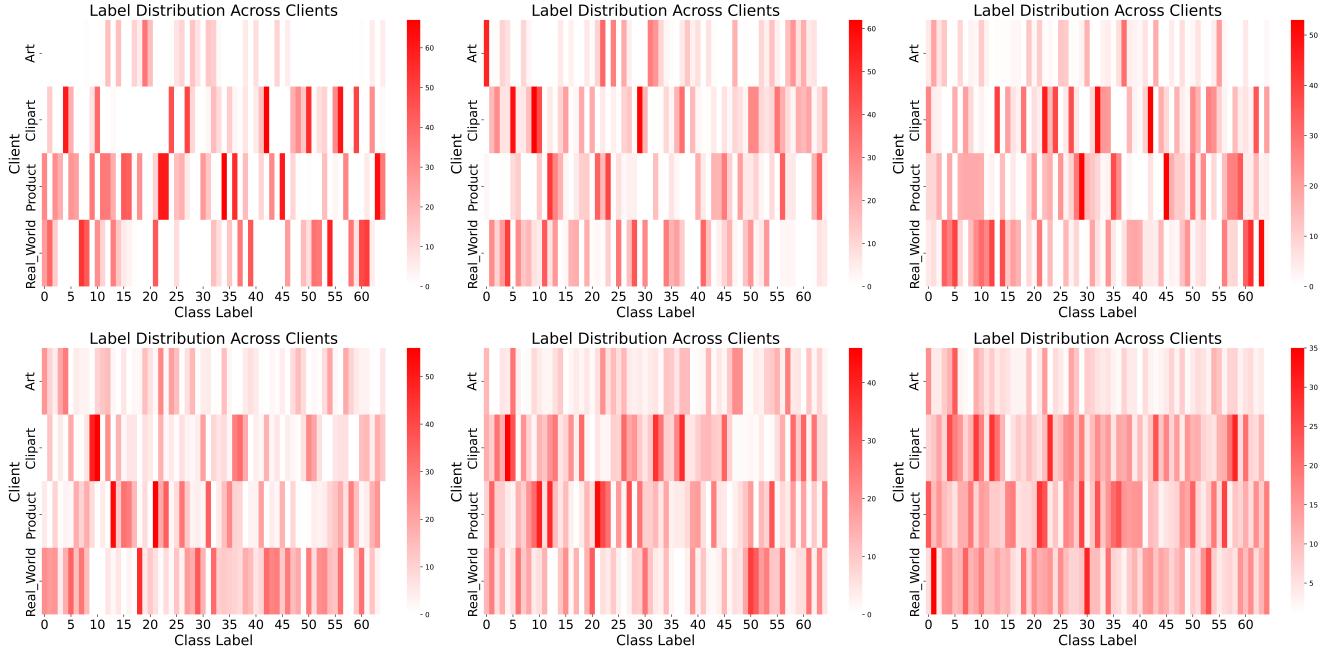


Figure 9. Number of samples per class across four clients when β equals 0.1, 0.3, 0.5, 0.7, 1, and 5, with each client holding data from a different domain.

Table 8. Experiments Configuration of different federated scenarios. Image size is operated after the resize operation. $|C|$ denotes the classification scale. $|K|$ denotes the clients number. E is the communication epochs for federation. B is the training batch size.

Scenario	Size	$ C $	Network	w	Rate	η	$ K $	E	B
Label Skew Setting § 5.2									
Cifar-10	224	10	CLIP (ViT-B/16)		1e-2	10	100	64	
Cifar-100	224	100	CLIP (ViT-B/16)		1e-1	10	100	64	
Tiny-ImageNet	224	200	CLIP (ViT-B/16)		1e-2	10	100	64	
Domain Skew § 5.3 and 5.4									
Digits	224	10	CLIP (ViT-B/16)		1e-2	4	50	16	
Office Caltech	224	10	CLIP (ViT-B/16)		1e-3	4	50	16	
PACS	224	7	CLIP (ViT-B/16)		1e-3	4	50	16	
Office-Home	224	65	CLIP (ViT-B/16)		1e-3	4	50	16	

of samples for that class. This setup results in a distribution that incorporates both domain shift (one domain per client) and Dirichlet-based class imbalance, presenting a scenario where the model faces both class distribution and domain differences, creating a more realistic, challenging, and diverse distribution for classification. We name the newly constructed dataset **Office-Home-LDS (Label and Domain Skew)**. Figure 9 shows the data distribution of Office-Home-LDS with different β values. Dataset and Constructor published at: <https://huggingface.co/datasets/WeiDai-David/Office-Home-LDS>.

Table 9. Hyper-parameters chosen for different methods. Hyper-parameters in different methodologies may share the same notation but represent distinct meanings.

Methods	Hyper-Parameter	Parameter value
SCAFFOLD	Global learning rate lr	0.25
MOON	Contrastive temp τ	0.5
	Proximal weight μ	1.0
FedDyn	Proximal weight α	0.5
FedOPT	Global learning rate η_g	0.5
FedProto	Proximal weight λ	2
FedNTD	Distill temp τ	1
	Reg weight β	1

C. Implementation Details

As for the uniform comparison evaluation, we follow [13] and conduct the local updating round $U = 10$. We use the SGD optimizer for all local updating optimization. The corresponding weight decay is 1×10^{-5} and momentum is 0.9. The learning rate η and communication epoch E are different in various scenarios, as shown in Table 8. Notably, the communication epoch is set according to when all federated approaches have little or no accuracy gain with more communication epochs. The local training batch size is $B = 64$. Furthermore, the Table 9 plots the chosen hyper-parameter for different methods.

D. Privacy Constraints

In our approach, the server only sends the eigenvectors and eigenvalues of the global covariance matrix back to clients, without sharing raw data or local covariance matrices. We demonstrate below that this information is insufficient for reconstructing any client’s original data.

(1) Eigenvectors and Eigenvalues Do Not Contain Raw Data.

The eigendecomposition provides only the geometric structure of the data distribution, without encoding individual sample details. Even if a client obtains eigenvectors and eigenvalues, reconstructing the original data is an **ill-posed problem**, as it admits infinitely many solutions.

(2) Low-Rank Property Prevents Data Reconstruction.

The covariance matrix is typically **low-rank**, meaning: $\text{rank}(\Sigma_i) \ll d$, where d is the original data dimension. This implies that even with full knowledge of eigenvectors and eigenvalues, clients can only access principal directions of the data and not its full details.

(3) Aggregation Prevents Isolation of Individual Client Contributions.

The global covariance matrix is an aggregate of all clients’ local covariance matrices: $\Sigma_i = \sum_{k=1}^K \frac{n_k^i}{N_i} \Sigma_k^i + \sum_{k=1}^K \frac{n_k^i}{N_i} (\mu_k^i - \mu_i)(\mu_k^i - \mu_i)^T$. Since each client’s contribution is mixed through weighted averaging: I. No single client can **isolate another client’s contribution** from Geometric Knowledge. II. Even if a client’s data is removed, the impact on Σ_i is distributed across all eigenvectors, making individual influences indistinguishable.

(4) Existing Literature Supports the Privacy of Covariance Matrices.

Prior works confirm that sharing higher-order statistics (covariance matrices) poses lower privacy risks than sharing model gradients (Melis et al.).

E. Large-Scale Client

In practical federated learning settings, the number of participating clients can significantly affect model performance due to increased data heterogeneity. To further evaluate the performance of our proposed method under larger-scale federated learning scenarios, we conducted additional experiments with an increased number of clients. Specifically, we conducted experiments on the label-skewed dataset CIFAR-10 with 100, 300, and 500 clients. As shown in Table 10, the results demonstrate that GGEUR remains robust and continues to enhance the performance of FedAvg (CLIP+MLP).

Table 10. Number of Clients K Impact on Performance.

Methods	CIFAR-10 ($\beta = 0.1$)		
	$K = 100$	$K = 300$	$K = 500$
FedAvg (CLIP+MLP)	87.89	84.69	82.05
+ GGEUR	93.55 (+5.66)	92.17 (+7.84)	90.43 (+8.38)

Algorithm 2 GGEUR (Multi-Domain Scenario)

Require: $X_k^i = [X_k^{(i,1)}, \dots, X_k^{(i,n_k^i)}] \in \mathbb{R}^{p \times n_k^i}$: Sample set of class i at client k , $GD_i = \{\xi_i^1, \dots, \xi_i^p, \lambda_i^1, \dots, \lambda_i^p\}$: Shared geometric shape (eigenvectors and eigenvalues) of class i , $\{\mu_{k'}^i\}$: Prototypes (means) of class i from other domains, N : Number of new samples to generate per original sample in Step 1, M : Number of samples to generate per prototype in Step 2.

Ensure: X_{new}^i : Augmented sample set of class i at client k

- 1: $X_{\text{new}}^i \leftarrow \emptyset$ ▷ Initialize augmented sample set
- 2: ▷ Step 1: Local Domain Augmentation
- 3: **for** $j = 1$ to n_k^i **do**
- 4: $X_{\text{new}}^i \leftarrow X_{\text{new}}^i \cup \text{GGEUR}(X_k^{(i,j)}, GD_i, N)$
- 5: **end for** ▷ Step 2: Cross-Domain Simulation
- 6: **for** each prototype $\mu_{k'}^i$ from other domains **do**
- 7: $X_{\text{new}}^i \leftarrow X_{\text{new}}^i \cup \text{GGEUR}(\mu_{k'}^i, GD_i, M)$
- 8: **end for**
- 9: **return** X_{new}^i

F. Computational Cost

We conducted experiments on the domain-skewed dataset Digits, comparing the training time required to complete the full model for FedAvg (CLIP+MLP), SCAFFOLD (CLIP+MLP), and MOON (CLIP+MLP) before and after applying GGEUR. The results (Tab11) show that GGEUR introduces almost no additional training time overhead. Specifically, after applying GGEUR, the training time for the three methods increased by only 3.5s, 4.6s, and 3.3s, respectively.

Table 11. The average training time (s) per round.

Methods	Digits		
	FedAvg	SCAFFOLD	MOON
CLIP+MLP + GGEUR	28.2 31.7 (+3.5)	54.5 59.1 (+4.6)	32.3 35.6 (+3.3)