# Federated Learning with Intermediate Representation Regularization

Ye Lin Tun, Chu Myaet Thwal, Yu Min Park, Seong-Bae Park, Choong Seon Hong*

*Department of Computer Science and Engineering*

*Kyung Hee University*

Yongin-si, Republic of Korea

{yelintun, chumyaet, yumin0906, sbpark71, cshong}@khu.ac.kr

*Abstract*—In contrast to centralized model training that involves data collection, federated learning (FL) enables remote clients to collaboratively train a model without exposing their private data. However, model performance usually degrades in FL due to the heterogeneous data generated by clients of diverse characteristics. One promising strategy to maintain good performance is by limiting the local training from drifting far away from the global model. Previous studies accomplish this by regularizing the distance between the representations learned by the local and global models. However, they only consider representations from the early layers of a model or the layer preceding the output layer. In this study, we introduce FedIntR, which provides a more fine-grained regularization by integrating the representations of intermediate layers into the local training process. Specifically, FedIntR computes a regularization term that encourages the closeness between the intermediate layer representations of the local and global models. Additionally, FedIntR automatically determines the contribution of each layer's representation to the regularization term based on the similarity between local and global representations. We conduct extensive experiments on various datasets to show that FedIntR can achieve equivalent or higher performance compared to the state-of-the-art approaches. Our code is available at https://github.com/YLTun/FedIntR.

*Keywords-federated learning; representation; data heterogeneity; distributed*

## I. INTRODUCTION

The real world data essential for powering intelligent services is spread across numerous edge devices (e.g., IoT devices, personal smartphones, or data silos of different organizations). Although deep learning can benefit from large datasets produced by mass collection, rising security concerns and privacy regulations [1] may prohibit a server from acquiring the data from edge devices. Such limitations may render the centralized training of deep neural network models infeasible. Federated learning (FL) [2]–[4], which allows edge devices to collaboratively train a model without sharing their data with a central server, has come to light as a viable option to meet these requirements. In particular, the FedAvg [2] algorithm has emerged as the de facto approach for model training in
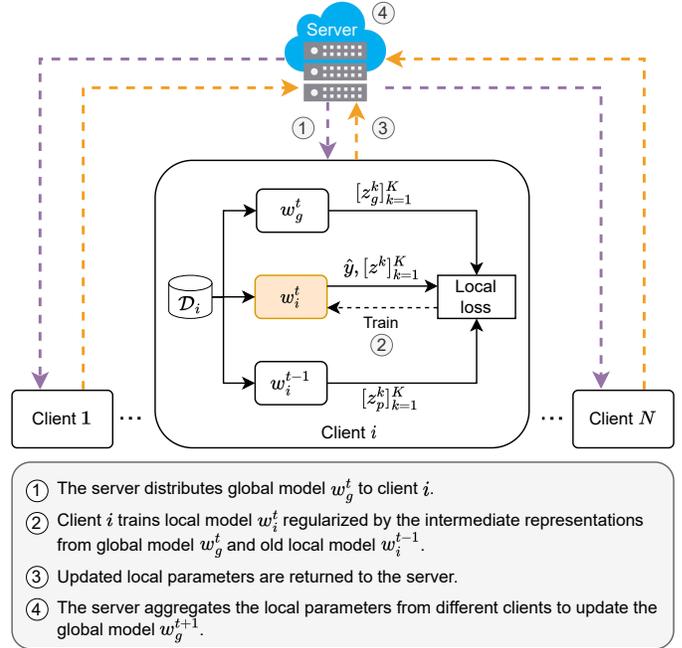
Fig. 1: An Overview of federated learning with FedIntR. FedIntR incorporates regularization into the local training step (i.e., step2) of FedAvg.

distributed environments with data privacy concerns. FedAvg operates by having each edge device train a local model on its own data before sending the trained parameters to the server. The server aggregates the received parameters into a single global model that inherits the trained capabilities of the local models.

In practice, large data heterogeneity may arise in an FL system since the local data on each device varies depending on the nature and behavior of the device. Heterogeneous data has been found to be a major issue that causes slow convergence and suboptimal model performance in federated learning [5], [6]. To address this, previous studies have incorporated various forms of regularization into the local training loss to reduce the divergence of local models from the global model [5], [7]–[9]. The distance between local and global parameters
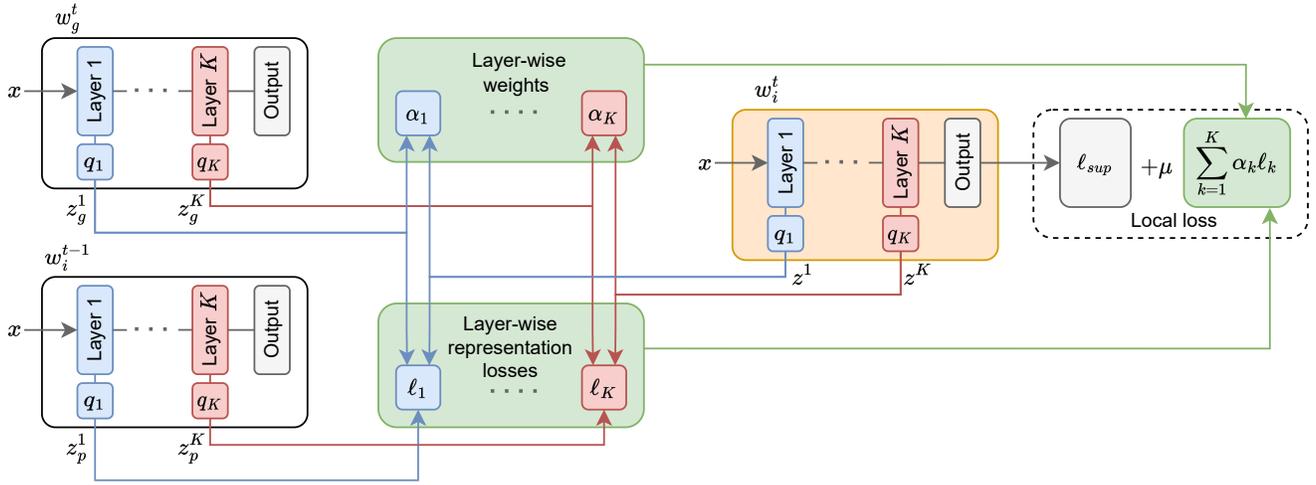
Fig. 2: Local loss computation in FedIntR. Using the intermediate representations $z$ from the local model $w_i^t$, global model $w_g^t$ and old local model $w_i^{t-1}$, FedIntR computes a regularization term. Together, this regularization term and the cross entropy loss $\ell_{sup}$ make up the local training loss.

can be limited using $\ell_2$-norm [5] or control variates [7]. A recent popular approach is to use representations such as MOON [8] and FedCKA [9], which adopt contrastive loss [10]–[12] to promote the closeness between local and global representations. However, these studies only focus on a few specific layers of the model to extract the representations. In particular, MOON extracts representations from the layer prior to the output layer, whereas FedCKA extracts from the first two layers based on the assumption that naturally similar layers are the most critical to improving the performance. Nevertheless, we argue that utilizing representations from all intermediate layers while assigning a proper weight to each of their contributions can offer a more fine-grained regularization to the training process.

Hence, we explore the idea of using representations extracted from intermediate layers to overcome the performance degradation caused by non-IID data in FL. We propose **FedIntR**, which promotes the similarity between **Int**ermediate layer **R**epresentations the of local and global models. Specifically, FedIntR computes a regularization term based on the contrastive loss [8], [10] using local and global intermediate representations. Moreover, FedIntR automatically calculates layer-wise weights to determine each intermediate layer's contribution to the regularization term. Analogous to the idea of FedCKA [9], where similar layers are more important for improving the performance, FedIntR assigns a larger contribution weight to the layers with a higher global and local representation similarity. To demonstrate the effectiveness of FedIntR, we conduct extensive experiments in comparison with various state-of-the-art FL approaches. Moreover, we evaluate FedIntR in different FL settings with varying degrees of data heterogeneity, numbers of clients, and local training epochs to ensure that it can maintain good performance in these scenarios.

## II. BACKGROUND AND RELATED WORK

### A. Federated Learning

The FedAvg [2] algorithm is the first federated learning framework that enables collaborative model training while preserving the privacy of user devices. Fig. 1 depicts the overview of a typical FL framework, which consists of four key steps at each training round $t$: (i) the central server distributes the global model $w_g^t$ to the clients; (ii) each client $i$ trains its local model $w_i^t$ on its private data $\mathcal{D}_i$; (iii) the trained local models are returned to the server; and (iv) the server aggregates the received local models into an updated global model $w_g^{t+1}$. FedAvg [2] also highlights the non-IID characteristics of an FL environment, and since then, improving the performance of an FL model under such constraints has been the subject of many studies.

### B. Tackling Data Heterogeneity

Many FL studies rely on three distinct strategies for mitigating the performance degradation caused by heterogeneous data: (i) adding regularization to the local training, (ii) modifying the aggregation scheme, and (iii) fitting personalized models for the clients. We discuss each of these strategies as follows.

**Local Training.** FedProx [5] includes a proximal term in the local training objective in order to restrict the Euclidean distance between the local and global models. MOON [8] introduces model-contrastive loss which is inspired by the NT-Xent (Normalized Temperature-scaled Cross Entropy) loss [13]–[15] used in SimCLR [10] for self-supervised representation learning. The NT-Xent loss is designed to minimize the distance between the representations of a positive pair (i.e., augmented views of the same sample) and maximize the distance between the representations of negative pairs (i.e., augmented views from different samples). Model-contrastive

loss in MOON aims to minimize the distance between representations extracted by the global and local models while simultaneously maximizing the distance between representations from the current local model and the old local model. MOON is based on the idea that, since the global model has been trained on the entire dataset, its representations can effectively regularize the local training. SCAFFOLD [7] lowers the client-variance in the local updates by minimizing the difference between the local and global control variates. Similar to MOON, FedCKA [9] uses global and local model representations to regularize the local training and shows that using deeper layers for regularization can harm the performance. Centered kernel alignment (CKA) [16] is used in FedCKA to calculate the distance between the representations, whereas MOON employs cosine similarity.

**Aggregation.** Another way to combat the downsides of data heterogeneity is to design a more intelligent aggregation scheme. PFNM [17] relies on Bayesian nonparametric methods to match the neurons of different client models before the aggregation process. However, PFNM only works with fully connected neural networks. Therefore, FedMA [18] proposes to match in a layer-wise fashion, which is applicable for CNN and LSTM models as well. Inspired by the effectiveness of SGD momentum in dampening gradient oscillations, FedAvgM [19] incorporates momentum in the aggregation process. FedNova [20] takes into account the varying number of local update steps completed by the clients and normalizes the local update by the number of steps prior to aggregation. pFedLA [21] generates personalized models by using hypernetworks to determine the layer-wise contribution of each client model to the aggregation step.

**Personalized Federated Learning.** Personalized models can improve the performance in non-IID settings by fitting the diverse distributions of the clients' data [22]–[25]. One simple and effective personalized FL mechanism is to group clients with similar characteristics into clusters, and cluster members can collectively train a personalized model. Many approaches [26]–[30] cluster the clients based on the similarity of local updates received at the server. Others [31]–[33] perform clustering on the client-side, where each client can choose their own personalized model from a set of available models based on performance.

Our proposed method, FedIntR, is one of the techniques that incorporates regularization into the local training process. A recent study [34] demonstrates that self-supervised learning with additional loss across intermediate layers can enhance the model's performance in the downstream tasks. Inspired by [34], we explore the integration of intermediate representations into the FL process to enable a more effective regularization for data heterogeneity issues. In contrast to our approach, MOON [8] only uses the representations from the layer preceding the output layer, and FedCKA [9] is based on the idea of minimizing the distance between similar layer representations and, therefore, only considers the early layers (the first two layers) of the model. FedIntR can also be

---

**Algorithm 1** FedIntR

1: **Input:** number of training rounds $T$, number of local epochs $E$, number of clients $N$, number of intermediate layers $K$, temperature $\tau$, learning rate $\eta$, weighting parameter $\mu$

2: **Output:** global model $w_g^T$

3: **Server executes:**
4:   initialize $w_g^0$
5:   **for** each round $t = 0, 1, \ldots, T-1$ **do**
6:     **for** each client $i = 1, 2, \ldots, N$ in parallel **do**
7:       $w_i^t \leftarrow \text{LocalUpdate}(i, w_g^t)$
8:       $w_g^{t+1} \leftarrow \sum_{i=1}^{N} \frac{|\mathcal{D}_i|}{|\mathcal{D}|} w_i^t$
9:   **return** $w_g^T$

10: **Client executes:** LocalUpdate$(i, w_g^t)$:
11:   $w_i^{t-1} \leftarrow w_i^t$
12:   $w_i^t \leftarrow w_g^t$
13:   **for** each epoch $e = 0, 1, \ldots, E-1$ **do**
14:     **for** each batch $\{x, y\} \in \mathcal{D}_i$ **do**
15:       $\hat{y}, [z^k]_{k=1}^K \leftarrow w_i^t(x)$
16:       $[z_g^k]_{k=1}^K \leftarrow w_g^t(x)$
17:       $[z_p^k]_{k=1}^K \leftarrow w_i^{t-1}(x)$
18:       $\ell_{sup} \leftarrow CrossEntropyLoss(\hat{y}, y)$
19:       **for** $k = 1, 2, \ldots, K$ **do**
20:         $\alpha_k \leftarrow \frac{\exp(\text{sim}(z^k, z_g^k)/\tau)}{\sum_{\hat{k}=1}^{K} \exp(\text{sim}(z^{\hat{k}}, z_g^{\hat{k}})/\tau)}$
21:         $\ell_k \leftarrow -\log \frac{\exp(\text{sim}(z^k, z_g^k)/\tau)}{\exp(\text{sim}(z^k, z_g^k)/\tau) + \exp(\text{sim}(z^k, z_p^k)/\tau)}$
22:       $\mathcal{L} \leftarrow \ell_{sup} + \mu \sum_{k=1}^{K} \alpha_k \ell_k$
23:       $w_i^t \leftarrow w_i^t - \eta \nabla \mathcal{L}$
24:   **return** $w_i^t$

---

considered a more general approach where we don't have to manually define which layers to include in the regularization term since it automatically determines the contribution of different intermediate layers to the regularization based on the similarities between their local and global representations.

## III. METHOD

### A. Federated Learning Objective

The goal of a federated learning system is to train a global model $w_g$ by solving:

$$w_g^* = \arg\min_{w_g} \sum_{i=1}^{N} \frac{|\mathcal{D}_i|}{|\mathcal{D}|} \mathcal{L}_i(w_g), \qquad (1)$$

where $\mathcal{L}_i$ and $\mathcal{D}_i$ represent the local loss and the dataset of the $i$-th client, $|\mathcal{D}| = \sum_{i=1}^{N} |\mathcal{D}_i|$ and $N$ is the total number of clients. From Fig. 1, we can observe that FedIntR only modifies the local training process of vanilla FedAvg, similar to MOON and FedCKA. In addition to the cross entropy loss, FedIntR adds a regularization term computed with the help

of intermediate representations from the global model $w_g^t$ and the old local model $w_i^{t-1}$ from the previous round.

## B. Local Training Process of FedIntR

Fig. 2 illustrates how the local training loss is calculated in our FedIntR framework, while Algorithm 1 describes the whole procedure in detail. The local training process begins with the $i$-th client receiving the global model $w_g^t$ from the central server. The client synchronizes the old local model $w_i^{t-1}$ with the current one $w_i^t$ (which we want to train), followed by synchronizing $w_i^t$ with $w_g^t$. Suppose that the deployed model structure contains a total of $K$ layers capable of extracting representations (e.g., convolutional or fully connected layers). Given an input $x$, FedIntR extracts representations $z^k$, $z_p^k$, and $z_g^k$ from the $k$-th layer of the current local model $w_i^t$, old local model $w_i^{t-1}$, and global model $w_g^t$, respectively. (A representation is obtained by passing the layer output through a projection head $q_k$ as shown in Fig. 2.) Inspired by MOON, FedIntR computes the representation loss for $k$-th layer as

$$\ell_k = -\log \frac{\exp(\mathrm{sim}(z^k, z_g^k)/\tau)}{\exp(\mathrm{sim}(z^k, z_g^k)/\tau) + \exp(\mathrm{sim}(z^k, z_p^k)/\tau)}, \quad (2)$$

where $\tau$ is the temperature parameter and $\mathrm{sim}(.,.)$ is the similarity function. We use the cosine similarity function, which is defined by

$$\mathrm{sim}(z_i, z_j) = \frac{z_i}{\|z_i\|_2} \cdot \frac{z_j}{\|z_j\|_2}. \quad (3)$$

The layer-wise representation loss $\ell_k$ encourages the local representation $z^k$ to be closer to global representation $z_g^k$ while moving further from $z_p^k$.

The layer-wise weight $\alpha_k$, representing the contribution of $\ell_k$ to the regularization term, is determined based on the similarity of $z^k$ and $z_g^k$ using the softmax function. Specifically, $\alpha_k$ is computed as follows:

$$\alpha_k = \frac{\exp(\mathrm{sim}(z^k, z_g^k)/\tau)}{\sum_{\hat{k}=1}^{K} \exp(\mathrm{sim}(z^{\hat{k}}, z_g^{\hat{k}})/\tau)}, \quad (4)$$

where $\sum_{k=1}^{K} \alpha_k = 1$. Equation (4) assigns a higher weight value to the layer $k$ with a higher degree of similarity between $z^k$ and $z_g^k$.

FedIntR computes $\alpha_k$ and $\ell_k$ for all layers $k \in [1, 2, \ldots, K]$. After that, we can incorporate $[\alpha_k]_{k=1}^{K}$ and $[\ell_k]_{k=1}^{K}$ into the local training loss as the regularization term. The local loss $\mathcal{L}$ is defined as:

$$\mathcal{L} = \ell_{sup} + \mu \sum_{k=1}^{K} \alpha_k \ell_k, \quad (5)$$

where $\ell_{sup}$ represents the cross-entropy loss, and the second component is the regularization term, with $\mu$ being the balancing parameter.

## IV. EXPERIMENTS

In this section, we discuss various experimental setups and evaluation results used to verify the effectiveness of FedIntR relative to other baseline approaches.

TABLE I: Test accuracy (%) on 10 clients with $\beta = 0.5$. We tune $\mu$ for each method, and the best $\mu$ value is shown in parentheses.

|  | Fashion-MNIST | SVHN | CIFAR-10 | CIFAR-100 |
|---|---|---|---|---|
| FedAvg | 88.90 | 86.04 | 65.82 | 40.43 |
| FedCKA | 88.85 (3) | 86.47 (1) | 66.55 (0.1) | **41.91** (10) |
| FedProx | 88.95 (0.001) | 86.42 (0.01) | 65.03 (0.01) | 39.91 (0.1) |
| MOON | **89.15** (1) | 86.61 (10) | 66.35 (5) | 40.46 (0.1) |
| FedIntR | **89.15** (10) | **87.17** (10) | **67.23** (3) | 41.48 (10) |

TABLE II: Test accuracy (%) of FedIntR with different $\mu$ values.

| $\mu$ | Fashion-MNIST | SVHN | CIFAR-10 | CIFAR-100 |
|---|---|---|---|---|
| 0.1 | 88.87 | 86.85 | 65.23 | 41.00 |
| 1 | 89.07 | 85.89 | 65.66 | 40.43 |
| 3 | 88.67 | 87.05 | **67.23** | 40.99 |
| 5 | 89.13 | 87.01 | 66.99 | 41.21 |
| 10 | **89.15** | **87.17** | 65.89 | **41.48** |

## A. Experimental Settings

FedIntR is evaluated on four datasets: Fashion-MNIST [35], SVHN [36], CIFAR-10 [37], and CIFAR-100 [37], in comparison with various state-of-the-art FL approaches including FedAvg [2], FedProx [5], MOON [8], and FedCKA [9]. We use a small CNN model for the Fashion-MNIST, SVHN, and CIFAR-10 datasets, and ResNet-20 [38], [39] for the CIFAR-100 dataset. Our CNN model is composed of three $3 \times 3$ convolutional layers with 8, 16, and 32 channels. Each convolutional layer is ReLU activated and followed by a $2 \times 2$ max pooling. The output of the final convolutional layer goes through two fully connected layers with 128 and 96 neurons. Both fully connected layers are also ReLU activated. FedIntR extracts representations from the three convolutional layers and the two fully connected layers of the CNN model. Similar to MOON, our projection head is a 2-layer MLP with an output dimension of 256. For the ResNet-20 model, the representations are retrieved from the first convolutional layer and each end of the three ResNet blocks.

We use the PyTorch [40] framework to implement FedIntR and other baseline methods. Referring to MOON, local training is performed using the SGD optimizer with a learning rate of 0.01, weight decay of 0.00001, and momentum of 0.9. We set the batch size as 512, the number of local epochs as 10, and the number of training rounds as 100. We use random horizontal flip as augmentation during the training except for the SVHN dataset. For both FedIntR and MOON, the temperature $\tau$ value is set to 0.5.

To represent the heterogeneous data in the clients, the training portion of each dataset is partitioned into 10 client datasets using the Dirichlet distribution [41]. The concentration parameter $\beta$ controls the strength of data heterogeneity in the Dirichlet distribution. A lower $\beta$ value indicates a greater level of data heterogeneity. We set the $\beta$ value as 0.5 by default. The performance of FedIntR and the baseline methods

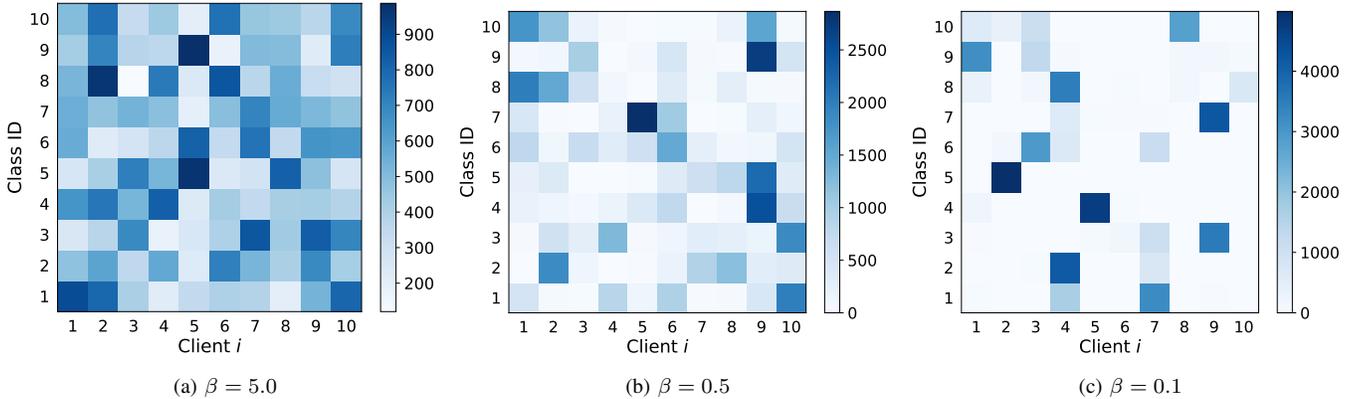(a) $\beta = 5.0$       (b) $\beta = 0.5$       (c) $\beta = 0.1$

Fig. 3: Client data distribution with different $\beta$ values for the CIFAR-10 dataset. A darker rectangle denotes a higher number of data samples for a specific class in a client.

TABLE III: Test accuracy (%) on CIFAR-10 and CIFAR-100 datasets with different $\beta$ values. We tune $\mu$ for each method and the best $\mu$ value is shown in brackets.

| | CIFAR-10 | | | CIFAR-100 | | |
|---|---|---|---|---|---|---|
| | $\beta = 5$ | $\beta = 0.5$ | $\beta = 0.1$ | $\beta = 5$ | $\beta = 0.5$ | $\beta = 0.1$ |
| FedAvg | 68.69 | 65.82 | 55.62 | **40.62** | 40.43 | 37.48 |
| FedCKA | 68.62 (1) | 66.55 (0.1) | 55.64 (0.1) | 40.40 (0.1) | **41.91** (10) | 35.73 (10) |
| FedProx | 70.03 (0.01) | 65.03 (0.01) | 56.55 (0.01) | 40.00 (0.001) | 39.91 (0.1) | 37.38 (0.1) |
| MOON | 69.34 (1) | 66.35 (5) | 57.60 (1) | 40.23 (0.1) | 40.46 (0.1) | 37.16 (5) |
| FedIntR | **70.66** (5) | **67.23** (3) | **58.32** (10) | 40.06 (0.1) | 41.48 (10) | **38.24** (0.5) |

is obtained by evaluating the corresponding global models on the testing portions of the respective datasets. We report the median of test accuracy values attained by the global model over the course of last 10 training rounds. Unless otherwise stated, we use the same settings as above for FedIntR and all other baseline methods.

### B. Accuracy on Different Datasets

Table I displays the top-1 test accuracy of different approaches on four datasets. The hyperparameter $\mu$ in FedIntR, MOON [8], FedProx [5], and FedCKA [9] can be tuned to control the regularization strength in the local training loss. For each approach, we tune the $\mu$ value and report the best results. For FedIntR and FedCKA, we tune $\mu$ from $[0.1, 1, 3, 5, 10]$ (which covers the range used by the original FedCKA paper). We tune $\mu$ from $[0.1, 1, 5, 10]$ for MOON and $[0.001, 0.01, 0.1, 1]$ for FedProx, according to their respective papers. In Table I, we report the corresponding best $\mu$ values in parentheses.

Our proposed FedIntR framework achieves superior performance in the SVHN and CIFAR-10 datasets compared to the baseline techniques. For the Fashion-MNIST dataset, both FedIntR and MOON achieve the best performance. FedCKA outperforms FedIntR on the CIFAR-100 dataset, but FedIntR still maintains a comparable performance. There are three plausible explanations for this. (i) For the CIFAR-100 dataset, we employ the ResNet-20 model, which contains fewer number of extraction points for intermediate represen-

tations than the CNN model. (ii) The superior performance of FedIntR can be observed more clearly with a higher degree of data heterogeneity, which will be discussed in Section IV-C, where we compare the performance of different techniques on varying degrees of data heterogeneity. (iii) FedIntR is more beneficial from longer training. By default, we set the number of communication rounds as 100, but as we will discuss in Section IV-D, FedIntR performs better with a higher number of training rounds. For reference, we also provide the performance of FedIntR from tuning $\mu$ in Table II.

Note that MOON paper [8] reported a higher accuracy across CIFAR-10 and CIFAR-100 datasets. As mentioned by the authors of FedCKA [9], MOON may have applied a variety of data augmentations to obtain higher performance. However, since MOON did not disclose their augmentation settings, we were unable to reproduce their results.

### C. Degree of Data Heterogeneity

In this experiment, we investigate the impact of clients' data heterogeneity on our proposed FedIntR framework. We adjust the $\beta$ parameter of the Dirichlet distribution to control the strength of data heterogeneity between the clients. Fig. 3 illustrates how different values of $\beta$ influence the local data distributions of the clients for the CIFAR-10 dataset.

Table III displays the test accuracy values attained by different approaches on the CIFAR-10 and CIFAR-100 datasets with $\beta$ values of 5, 0.5, and 0.1. Similar to Section IV-B, we tune the $\mu$ value for each approach (shown in parentheses

TABLE IV: Test accuracy (%) at different rounds $T$ on CIFAR-10 dataset with different number of clients $N$.

| | $N = 50$ | | | | $N = 100$ | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | $T = 100$ | $T = 200$ | $T = 500$ | $T = 800$ | $T = 100$ | $T = 200$ | $T = 500$ | $T = 800$ | $T = 1000$ |
| FedAvg | **49.41** | 57.79 | 65.26 | 64.62 | 37.98 | **50.69** | 61.00 | 63.64 | 64.90 |
| FedCKA ($\mu = 1$) | 45.28 | 56.76 | 64.90 | 63.96 | 36.81 | 47.10 | 60.43 | 64.30 | 64.41 |
| FedCKA ($\mu = 10$) | 44.45 | 53.83 | 58.69 | 59.08 | 36.68 | 46.85 | 56.51 | 59.91 | 60.49 |
| FedProx ($\mu = 0.01$) | 46.17 | **58.12** | 65.44 | 65.75 | 37.79 | 50.37 | 61.08 | 64.48 | 64.75 |
| MOON ($\mu = 1$) | 48.25 | 57.98 | 64.24 | 64.23 | **38.08** | 47.98 | 60.08 | 63.66 | 64.85 |
| MOON ($\mu = 10$) | 42.71 | 55.83 | 64.52 | 65.18 | 34.34 | 46.81 | **61.68** | 63.97 | 65.12 |
| FedIntR ($\mu = 1$) | 47.02 | 58.01 | 63.89 | 63.75 | 35.87 | 48.15 | 60.59 | 63.46 | 64.26 |
| FedIntR ($\mu = 10$) | 46.63 | 56.39 | **65.48** | **66.94** | 32.13 | 49.76 | 61.31 | **65.55** | **66.57** |

in Table III) and report the best results. For the CIFAR-10 dataset, FedIntR exhibits superior performance over the baseline approaches in all data heterogeneity settings. FedProx has comparable performance to FedIntR when $\beta = 5$ (i.e., when the degree of data heterogeneity among the clients is low), but its performance diminishes with a smaller $\beta$ value. For the CIFAR-100 dataset, vanilla FedAvg outperforms all other approaches when $\beta = 5$. However, as the $\beta$ value decreases, incorporating a form of regularization into the local loss can outperform the vanilla FedAvg. With $\beta = 0.1$, our proposed FedIntR significantly outperforms the other approaches. We speculate that fine-grained regularization computed from intermediate representations and layer-wise contributions makes FedIntR more robust in FL environments with high degrees of data heterogeneity.

### D. Scalability

To ensure that FedIntR is scalable in FL environments with a large number of clients, we conduct experiments on the CIFAR-10 dataset with 50 and 100 client settings (i.e., $N = 50$ and $N = 100$). For both settings, we use the Dirichlet distribution with $\beta = 0.5$ to partition the data, and 20% of the clients are chosen at random to participate in each training round. In particular, we uniformly sample 10 clients for $N = 50$ and 20 clients for $N = 100$ settings. For FedProx, $\mu$ value is set to 0.01, and for all other approaches, we experiment with both $\mu = 1$ and $\mu = 10$. We set the number of training rounds $T$ as 800 for $N = 50$ and 1000 for $N = 100$.

From Table IV, we can observe that vanilla FedAvg outperforms the majority of other approaches in the early training rounds (i.e., $T = 100$ and $T = 200$) for both $N = 50$ and $N = 100$ settings. This is due to the fact that the local training loss in FedAvg focuses only on performance improvement, whereas the other approaches include a regularization term with additional objectives. However, with a sufficient number of training rounds, FedIntR with $\mu = 10$ can significantly outperform all other approaches. In the $N = 50$ setting, FedIntR ($\mu = 10$) outperforms the second best approach, FedProx, by 1.19% at $T = 800$, and in the $N = 100$ setting, it outperforms MOON ($\mu = 10$) by 1.45% at $T = 1000$. Fig. 4 shows the test accuracy values of different methods in each round for the $N = 100$ setting.
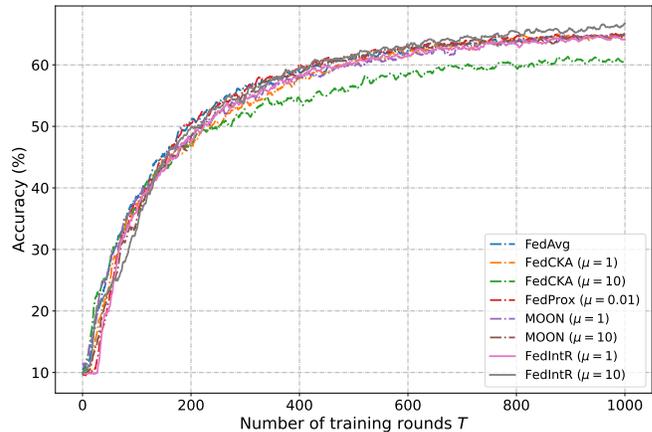


Fig. 4: Test accuracy (%) on CIFAR-10 dataset with number of clients $N = 100$.

### E. Number of Local Epochs

Using the CIFAR-10 dataset, we also explore the effects of the number of local epochs $E$ on different approaches. The experiments are performed using 1, 5, 10 (default), 20, and 40 local epochs with 100 rounds of training. For each approach, we use the same $\mu$ values tuned with the CIFAR-10 dataset in Section IV-B and the results are shown in Fig. 5. Consistent with the findings in MOON, setting $E = 1$ results in poor accuracy for all approaches, while setting $E = 10$ enables them to achieve their best performance. Due to the drift towards the local data distribution, using a larger number of local epochs (i.e., 20 or 40) can decrease the performance for all approaches. Nonetheless, our proposed method, FedIntR, achieves the best performance in most settings except $E = 5$.

### F. Ablation on Layer-wise Weights

FedIntR makes use of the layer-wise weight $\alpha_k$ that determines the contribution of each individual representation loss $\ell_k$ to the regularization term. To ensure that $\alpha_k$ is a crucial component of FedIntR for improving the performance, we conduct an ablation study by computing the regularization term in an alternative way. Instead of incorporating $\alpha_k$ into the regularization term, we use the average of $[\ell_k]_{k=1}^K$ as
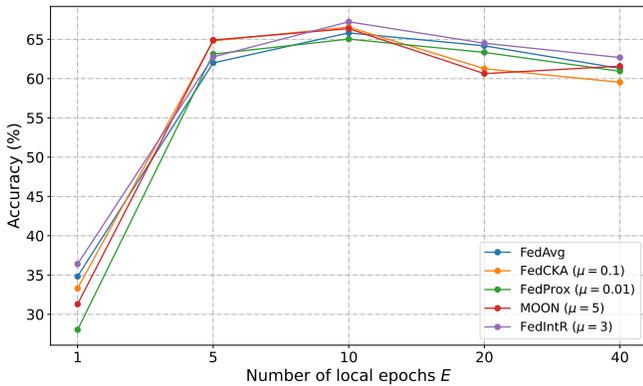
Fig. 5: Test accuracy (%) with different numbers of local epochs on the CIFAR-10 dataset. For each approach, we use the same $\mu$ values tuned with the CIFAR-10 dataset as in Table I.

TABLE V: Test accuracy (%) of FedIntR with and without $\alpha_i$ for different datasets. We use the same $\mu$ values tuned for each dataset as in Table I.

|            | Fashion-MNIST | SVHN  | CIFAR-10 | CIFAR-100 |
|------------|---------------|-------|----------|-----------|
| $\alpha_i$ | **89.15**     | **87.17** | **67.23** | **41.48** |
| Average    | 89.03         | 85.62 | 64.20    | 40.33     |

the regularization term. Specifically, we apply the local loss function defined as follows:

$$\mathcal{L} = \ell_{sup} + \frac{\mu}{K} \sum_{k=1}^{K} \ell_k. \tag{6}$$

In Table V, we compare the performance of FedIntR with and without $\alpha_k$ on four datasets. We can observe that calculating a layer-wise weight $\alpha_k$ for its corresponding representation loss $\ell_k$ can significantly improve the performance of the resulting global model. $\alpha_k$ determines the contribution of individual $\ell_k$ to the regularization term, enabling FedIntR to perform local training with fine-grained regularization.

## V. CONCLUSION

Nowadays, data valuable for powering intelligent services may be distributed over numerous personal edge devices. This renders data collection impractical due to privacy concerns, and hence the role of federated learning substantially grows. Client data heterogeneity in FL is one major shortcoming that is responsible for the poor performance of global model. To mitigate this, we propose FedIntR, which introduces a simple and effective regularization term to complement the local training process of FL. FedIntR uses intermediate layer representations to incorporate more information into the regularization process. Moreover, the layer-wise weights in FedIntR automatically determine each representation's contribution to the regularization term. This enables FedIntR to produce fine-grained regularization for the local training process, thus enhancing the global model performance. We conduct extensive experiments on various datasets and FL settings to examine the efficacy of FedIntR. Comparison between FedIntR and various state-of-the-art methods demonstrates that FedIntR achieves superior performance in most settings.

## REFERENCES

[1] P. Voigt and A. v. d. Bussche, *The EU General Data Protection Regulation (GDPR): A Practical Guide*, 1st ed. Springer Publishing Company, Incorporated, 2017.

[2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.

[3] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.

[4] T. Li, A. K. Sahu, A. S. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, pp. 50–60, 2020.

[5] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Proceedings of Machine Learning and Systems*, vol. 2, pp. 429–450, 2020.

[6] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," *arXiv preprint arXiv:1806.00582*, 2018.

[7] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.

[8] Q. Li, B. He, and D. Song, "Model-contrastive federated learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 713–10 722.

[9] H. M. Son, M. H. Kim, and T.-M. Chung, "Compare where it matters: Using layer-wise regularization to improve federated learning on heterogeneous data," *arXiv preprint arXiv:2112.00407*, 2021.

[10] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.

[11] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.

[12] I. Misra and L. v. d. Maaten, "Self-supervised learning of pretext-invariant representations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6707–6717.

[13] K. Sohn, "Improved deep metric learning with multi-class n-pair loss objective," *Advances in neural information processing systems*, vol. 29, 2016.

[14] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3733–3742.

[15] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[16] S. Kornblith, M. Norouzi, H. Lee, and G. Hinton, "Similarity of neural network representations revisited," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3519–3529.

[17] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, N. Hoang, and Y. Khazaeni, "Bayesian nonparametric federated learning of neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7252–7261.

[18] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated learning with matched averaging," *arXiv preprint arXiv:2002.06440*, 2020.

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

[20] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," *Advances in neural information processing systems*, vol. 33, pp. 7611–7623, 2020.

[21] X. Ma, J. Zhang, S. Guo, and W. Xu, "Layer-wised model aggregation for personalized federated learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 092–10 101.

[22] C. T Dinh, N. Tran, and J. Nguyen, "Personalized federated learning with moreau envelopes," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 394–21 405, 2020.

[23] A. Fallah, A. Mokhtari, and A. Ozdaglar, "Personalized federated learning: A meta-learning approach," *arXiv preprint arXiv:2002.07948*, 2020.

[24] Y. Jiang, J. Konečnỳ, K. Rush, and S. Kannan, "Improving federated learning personalization via model agnostic meta learning," *arXiv preprint arXiv:1909.12488*, 2019.

[25] M. Zhang, K. Sapra, S. Fidler, S. Yeung, and J. M. Alvarez, "Personalized federated learning with first order model optimization," *arXiv preprint arXiv:2012.08565*, 2020.

[26] F. Sattler, K.-R. Müller, and W. Samek, "Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 8, pp. 3710–3722, 2020.

[27] 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 *2021 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom)*. IEEE, 2021, pp. 228–237.

[28] M. Xie, G. Long, T. Shen, T. Zhou, X. Wang, J. Jiang, and C. Zhang, "Multi-center federated learning," *arXiv preprint arXiv:2005.01026*, 2020.

[29] F. Sattler, K.-R. Müller, T. Wiegand, and W. Samek, "On the byzantine robustness of clustered federated learning," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 8861–8865.

[30] A. Ghosh, J. Hong, D. Yin, and K. Ramchandran, "Robust federated learning in a heterogeneous environment," *arXiv preprint arXiv:1906.06629*, 2019.

[31] A. Ghosh, J. Chung, D. Yin, and K. Ramchandran, "An efficient framework for clustered federated learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 19 586–19 597, 2020.

[32] Y. Mansour, M. Mohri, J. Ro, and A. T. Suresh, "Three approaches for personalization with applications to federated learning," *arXiv preprint arXiv:2002.10619*, 2020.

[33] Y. L. Tun, M. N. Nguyen, C. M. Thwal, J. Choi, and C. S. Hong, "Contrastive encoder pre-training based clustered federated learning for heterogeneous data," *Available at SSRN: https://ssrn.com/abstract=4148978*, 2022.

[34] A. Kaku, S. Upadhya, and N. Razavian, "Intermediate layers matter in momentum contrastive self supervised learning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 24 063–24 074, 2021.

[35] H. Xiao, K. Rasul, and R. Vollgraf. (2017) Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.

[36] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.

[37] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," *Technical report, University of Toronto, 2009*, 2009. [Online]. Available: https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf

[38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[39] Y. Idelbayev, "Proper ResNet implementation for CIFAR10/CIFAR100 in PyTorch," https://github.com/akamaster/pytorch_resnet_cifar10, accessed: 2022-10-01.

[40] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.

[41] T. S. Ferguson, "A bayesian analysis of some nonparametric problems," *The annals of statistics*, pp. 209–230, 1973.