



# Model-Contrastive Federated Learning

Qinbin Li

National University of Singapore

qinbin@comp.nus.edu.sg

Bingsheng He

National University of Singapore

hebs@comp.nus.edu.sg

Dawn Song

UC Berkeley

dawnsong@berkeley.edu

## Federated learning :

Federated learning enables multiple parties to collaboratively train a machine learning model without communicating their local data. A key challenge in federated learning is to handle the heterogeneity of local data distribution across parties.

**Generic Federated Learning (G-FL)** aims to create a single global model that performs well across all clients.

**Personalized Federated Learning (P-FL)** acknowledges the heterogeneity among clients and focuses on tailoring models to individual clients.

## Data heterogeneity:

A key challenge in federated learning is the heterogeneity of data distribution on different parties. The data can be non-identically distributed among the parties in many real-world applications, which can degrade the performance of federated learning. Non-IID data includes forms as follow:

Feature distribution skew

Label distribution skew

Same label but different features

Same features but different label

Distribution skew

When each party updates its local model, its local objective may be far from the global objective. Thus, the averaged global model is away from the global optima.

## Contrastive Learning:

The key idea of contrastive learning is to reduce the distance between the representations of different augmented views of the same image (i.e., positive pairs), and increase the distance between the representations of augmented views of different images (i.e., negative pairs).

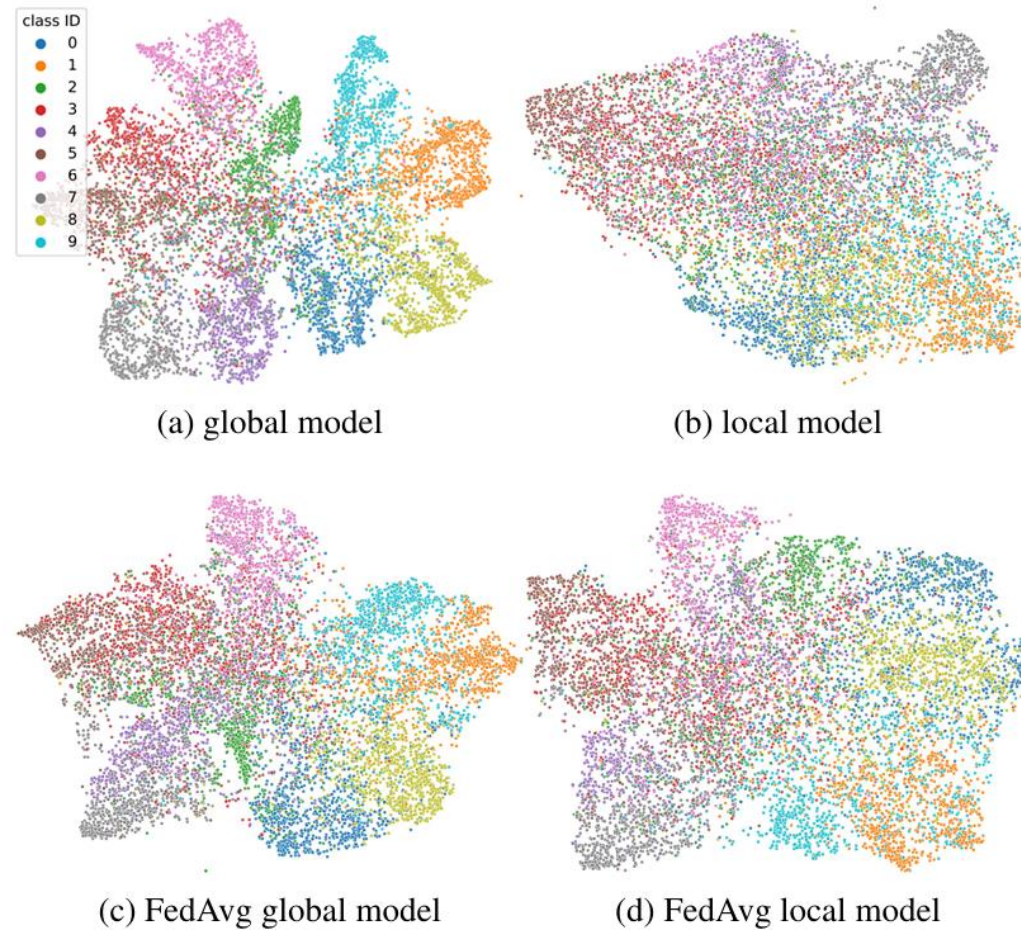
A typical contrastive learning framework is SimCLR. Given an image  $x$ , SimCLR first creates two correlated views of this image using different data augmentation operators. A base encoder and a projection head are trained to extract the representation vectors and map the representations to a latent space.

$$l_{i,j} = -\log \frac{\exp(\text{sim}(x_i, x_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{I}_{[k \neq i]} \exp(\text{sim}(x_i, x_k)/\tau)}$$

NT-Xent Loss

## Model-Contrastive Federated Learning:

MOON is based on an intuitive idea: the model trained on the whole dataset is able to extract a better feature representation than the model trained on a skewed subset. Visualized the hidden vectors :



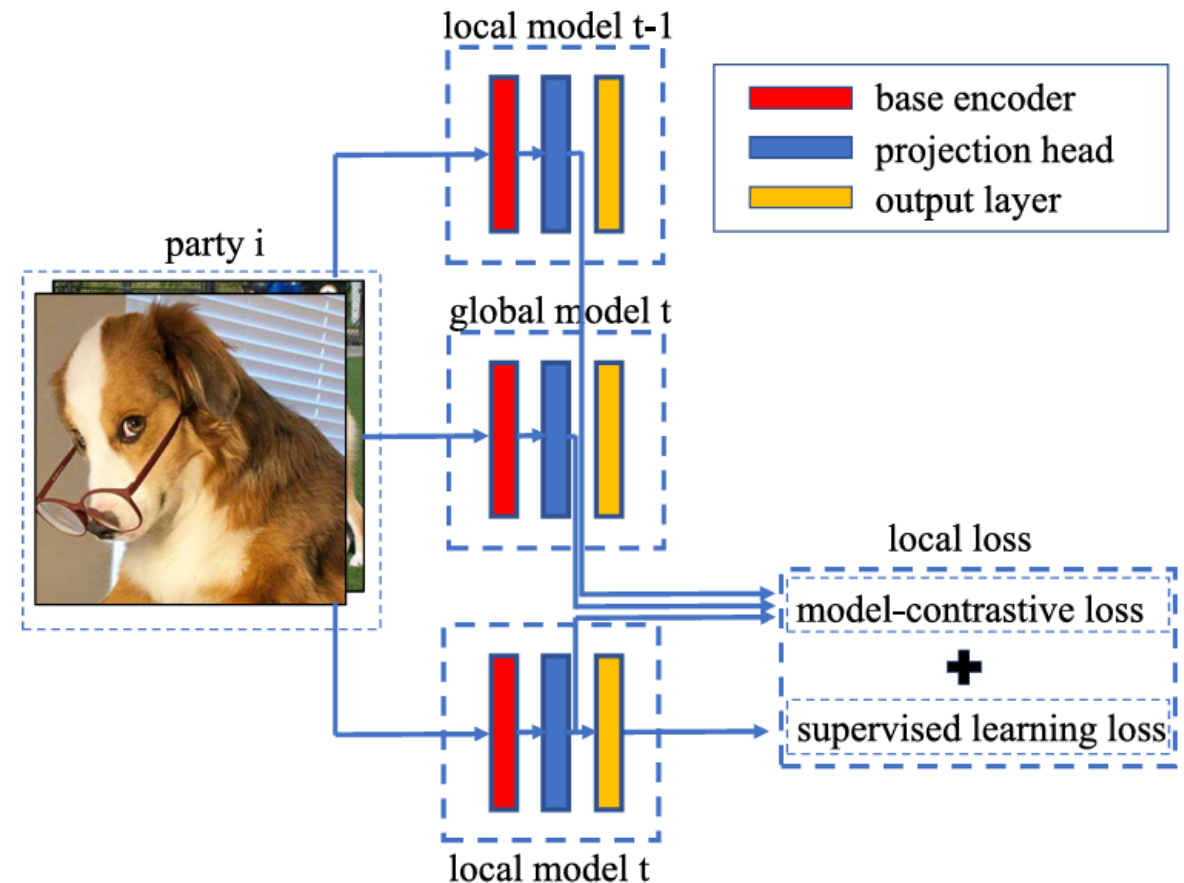
## MOON : Model-Contrastive Federated Learning

MOON aims to decrease the distance between the representation learned by the local model and the representation learned by the global model, and increase the distance between the representation learned by the local model and the representation learned by the previous local model.

### Network Architecture

The network has three components: a base encoder, a projection head, and an output layer.

We use  $Fw(\cdot)$  to denote the whole network and  $Rw(\cdot)$  to denote the network before the output layer.



## Local Objective

Local loss consists two parts. The first part is a typical loss term (e.g., cross-entropy loss) in supervised learning denoted as  $\ell_{sup}$ . The second part is our proposed model-contrastive loss term denoted as  $\ell_{con}$ .

For every input  $x$ , we extract the representation of  $x$  from the global model  $w^t$

$$z_{glob} = R_{w^t}(x) \quad z_{prev} = R_{w_i^{t-1}}(x)$$

Similar to NT-Xent loss, we define model-contrastive loss as:

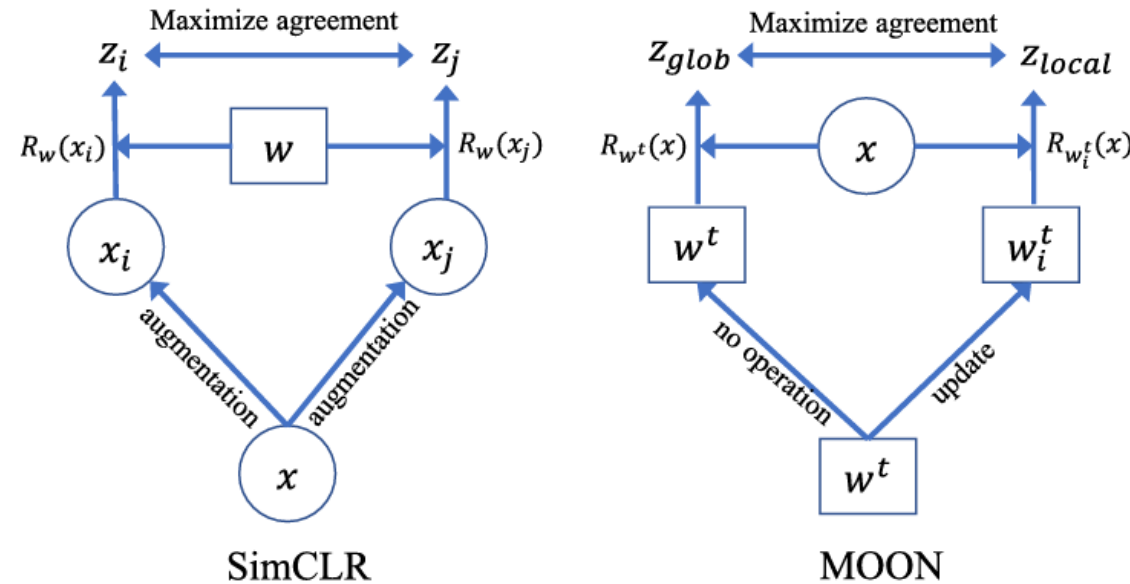
$$\ell_{con} = -\log \frac{\exp(\text{sim}(z, z_{glob})/\tau)}{\exp(\text{sim}(z, z_{glob})/\tau) + \exp(\text{sim}(z, z_{prev})/\tau)}$$

$\tau$  denotes a temperature parameter



The loss of an input  $(x, y)$  is computed by

$$\ell = \ell_{sup}(w_i^t; (x, y)) + \mu \ell_{con}(w_i^t; w_i^{t-1}; w^t; x)$$



An notable thing is that considering an ideal case where the local model is good enough and learns (almost) the same representation as the global model (i.e.,  $Z_{glob} = Z_{prev}$ ), the model-contrastive loss will be a constant (i.e.,  $-\log 1/2$ ). Thus, MOON will produce the same result as FedAvg, since there is no heterogeneity issue. In this sense, our approach is robust regardless of different amount of drifts.



## Algorithm 1: The MOON framework

**Input:** number of communication rounds  $T$ ,  
number of parties  $N$ , number of local  
epochs  $E$ , temperature  $\tau$ , learning rate  $\eta$ ,  
hyper-parameter  $\mu$

**Output:** The final model  $w^T$

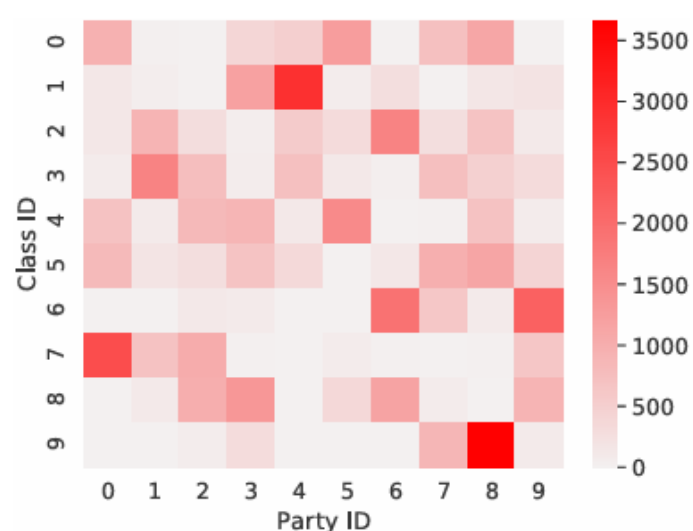
```

1 Server executes:
2 initialize  $w^0$ 
3 for  $t = 0, 1, \dots, T - 1$  do
4   for  $i = 1, 2, \dots, N$  in parallel do
5     send the global model  $w^t$  to  $P_i$ 
6      $w_i^t \leftarrow \text{PartyLocalTraining}(i, w^t)$ 
7    $w^{t+1} \leftarrow \sum_{k=1}^N \frac{|\mathcal{D}^i|}{|\mathcal{D}|} w_k^t$ 
8 return  $w^T$ 

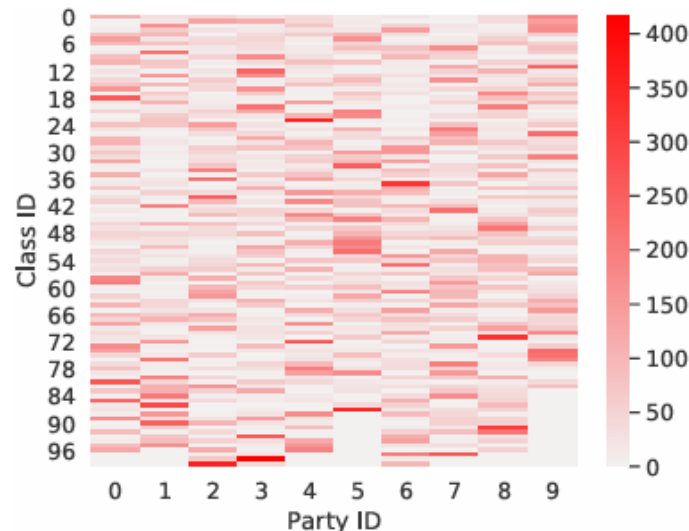
9 PartyLocalTraining( $i, w^t$ ):
10  $w_i^t \leftarrow w^t$ 
11 for epoch  $i = 1, 2, \dots, E$  do
12   for each batch  $\mathbf{b} = \{x, y\}$  of  $\mathcal{D}^i$  do
13      $\ell_{sup} \leftarrow \text{CrossEntropyLoss}(F_{w_i^t}(x), y)$ 
14      $z \leftarrow R_{w_i^t}(x)$ 
15      $z_{glob} \leftarrow R_{w^t}(x)$ 
16      $z_{prev} \leftarrow R_{w_i^{t-1}}(x)$ 
17      $\ell_{con} \leftarrow$ 
18        $-\log \frac{\exp(\text{sim}(z, z_{glob})/\tau)}{\exp(\text{sim}(z, z_{glob})/\tau) + \exp(\text{sim}(z, z_{prev})/\tau)}$ 
19      $\ell \leftarrow \ell_{sup} + \mu \ell_{con}$ 
20      $w_i^t \leftarrow w_i^t - \eta \nabla \ell$ 
21 return  $w_i^t$  to server

```

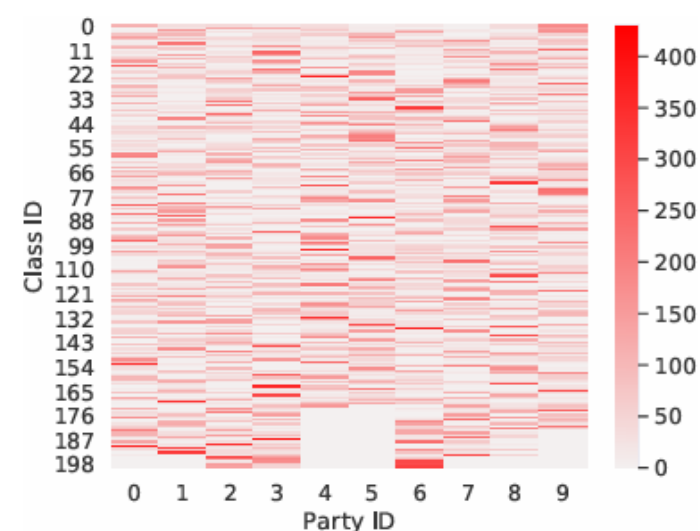
## Accuracy Comparison



(a) CIFAR-10



(b) CIFAR-100

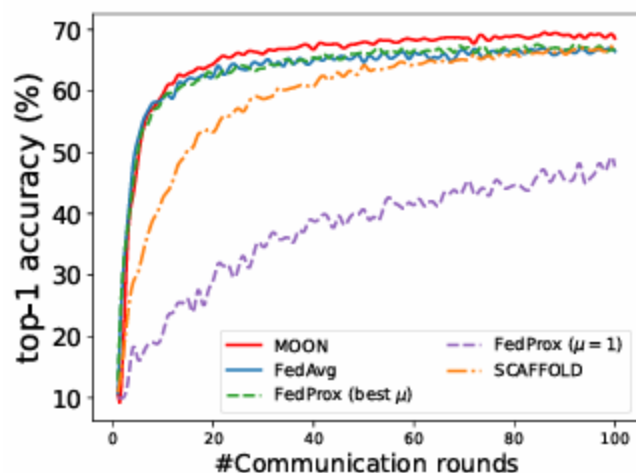


(c) Tiny-Imagenet

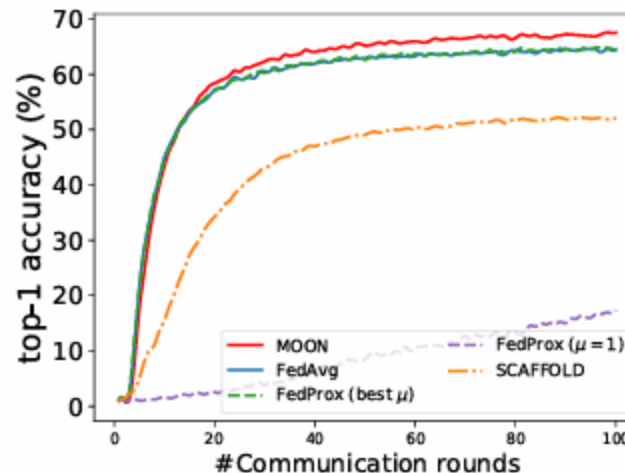
Method	CIFAR-10	CIFAR-100	Tiny-Imagenet
MOON	<b>69.1%±0.4%</b>	<b>67.5%±0.4%</b>	<b>25.1%±0.1%</b>
FedAvg	66.3%±0.5%	64.5%±0.4%	23.0%±0.1%
FedProx	66.9%±0.2%	64.6%±0.2%	23.2%±0.2%
SCAFFOLD	66.6%±0.2%	52.5%±0.3%	16.0%±0.2%
SOLO	46.3%±5.1%	22.3%±1.0%	8.6%±0.4%

The top-1 accuracy of MOON and the other baselines on test datasets

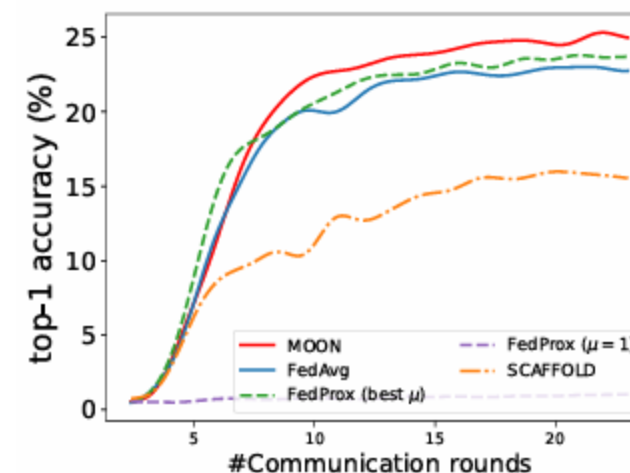
## Communication Efficiency



(a) CIFAR-10



(b) CIFAR-100



(c) Tiny-Imagenet

Figure 6. The top-1 test accuracy in different number of communication rounds. For FedProx, we report both the accuracy with best  $\mu$  and the accuracy with  $\mu = 1$ .

Method	CIFAR-10		CIFAR-100		Tiny-Imagenet	
	#rounds	speedup	#rounds	speedup	#rounds	speedup
FedAvg	100	1×	100	1×	20	1×
FedProx	52	1.9×	75	1.3×	17	1.2×
SCAFFOLD	80	1.3×	—	<1×	—	<1×
MOON	<b>27</b>	<b>3.7×</b>	<b>43</b>	<b>2.3×</b>	<b>11</b>	<b>1.8×</b>

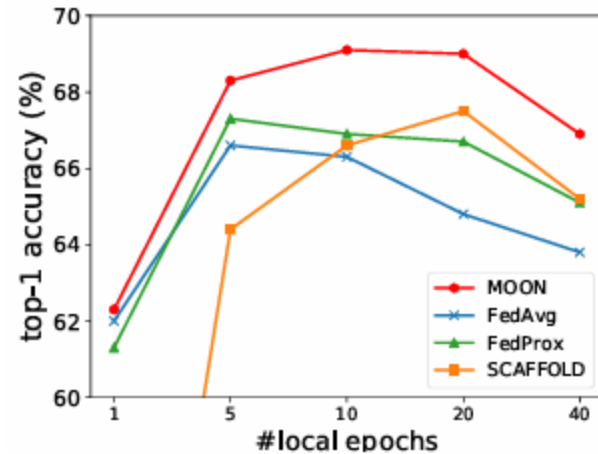
The number of rounds of different approaches to achieve the same accuracy as running FedAvg for 100 rounds

## Number of Local Epochs

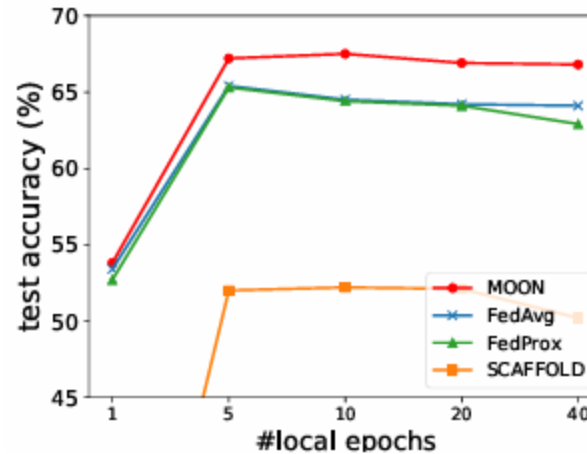
We study the effect of number of local epochs on the accuracy of final model.

When the number of local epochs is 1, the local update is very small. Thus, the training is slow and the accuracy is relatively low given the same number of communication rounds.

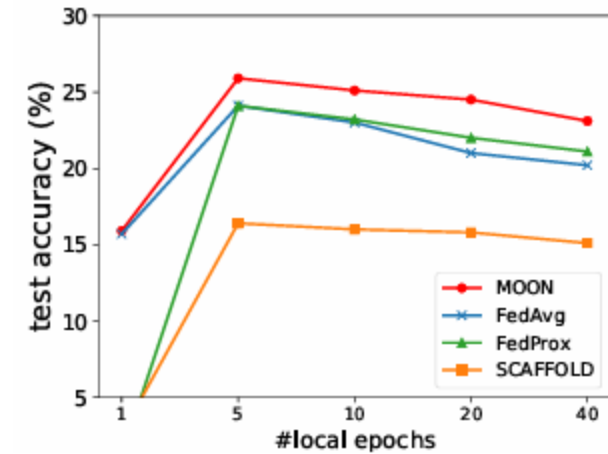
When the number of local epochs becomes too large, the accuracy of all approaches drops, which is due to the drift of local updates, i.e., the local optima are not consistent with the global optima.



(a) CIFAR-10



(b) CIFAR-100



(c) Tiny-Imagenet

## Scalability

To show the scalability of MOON, we try a larger number of parties on CIFAR-100. Specifically, we try two settings:

- (1) We partition the dataset into 50 parties and all parties participate in federated learning in each round.
- (2) We partition the dataset into 100 parties and randomly sample 20 parties to participate in federated learning in each round.

Method	#parties=50		#parties=100	
	100 rounds	200 rounds	250 rounds	500 rounds
MOON ( $\mu=1$ )	54.7%	58.8%	54.5%	58.2%
MOON ( $\mu=10$ )	<b>58.2%</b>	<b>63.2%</b>	<b>56.9%</b>	<b>61.8%</b>
FedAvg	51.9%	56.4%	51.0%	55.0%
FedProx	52.7%	56.6%	51.3%	54.6%
SCAFFOLD	35.8%	44.9%	37.4%	44.5%
SOLO	10% $\pm$ 0.9%		7.3% $\pm$ 0.6%	

## Loss Function

To maximize the agreement between the representation learned by the global model and the representation learned by the local model, our model-contrastive loss  $\ell_{\text{con}}$  is proposed inspired by NT-Xent loss. Another intuitive option is to use  $\ell_2$  regularization, and the local loss is

$$\ell = \ell_{\text{sup}} + \mu \|z - z_{\text{glob}}\|_2$$

We can observe that simply using  $\ell_2$  norm even cannot improve the accuracy compared with FedAvg on CIFAR-10. While using  $\ell_2$  norm can improve the accuracy on CIFAR-100 and Tiny-Imagenet, the accuracy is still lower than MOON.

Table 5. The top-1 accuracy with different kinds of loss for the second term of local objective. We tune  $\mu$  from  $\{0.001, 0.01, 0.1, 1, 5, 10\}$  for the  $\ell_2$  norm approach and report the best accuracy.

second term	CIFAR-10	CIFAR-100	Tiny-Imagenet
none (FedAvg)	66.3%	64.5%	23.0%
$\ell_2$ norm	65.8%	66.9%	24.0%
MOON	<b>69.1%</b>	<b>67.5%</b>	<b>25.1%</b>

To improve the performance of federated deep learning models on non-IID datasets, we propose model-contrastive learning (MOON), a simple and effective approach for federated learning. MOON introduces a new learning concept, i.e., contrastive learning in model-level. Our extensive experiments show that MOON achieves significant improvement over state-of-the-art approaches on various image classification tasks. As MOON does not require the inputs to be images, it potentially can be applied to non-vision problems.



Thanks