

Model-Contrastive Federated Learning

Qinbin Li¹, Bingsheng He¹, and Dawn Song²

¹National University of Singapore, ²UC Berkeley

Presented by Rui Chen

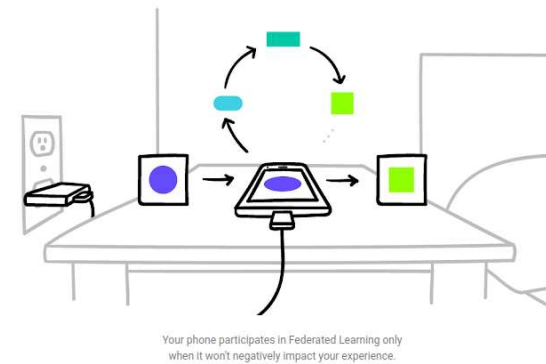
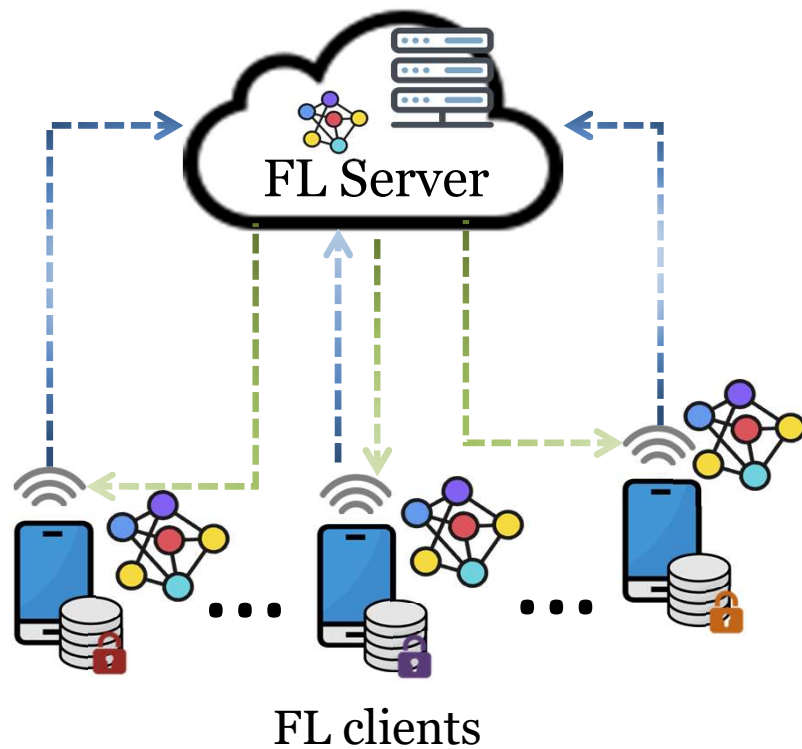
CVPR 2021

Copyright @ ANTS Laboratory

Outline

- Motivation
- Model-Contrastive Federated Learning
- Experimental Results
- Conclusion

Motivation: Federated Learning



- 🔒 Data never leaves local devices
- 🔒 Learn on fresh real-world data

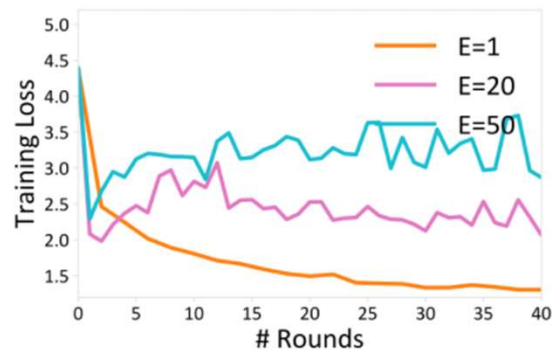
Heterogeneity in FL

statistical heterogeneity

highly non-identically distributed data



too much local work can hurt convergence



systems heterogeneity

stragglers



dropping slow devices can exacerbate convergence issues

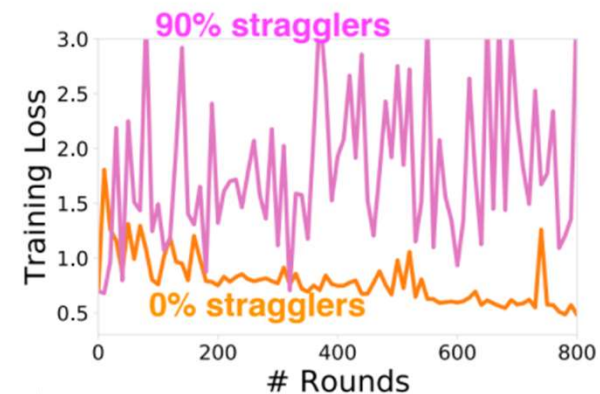
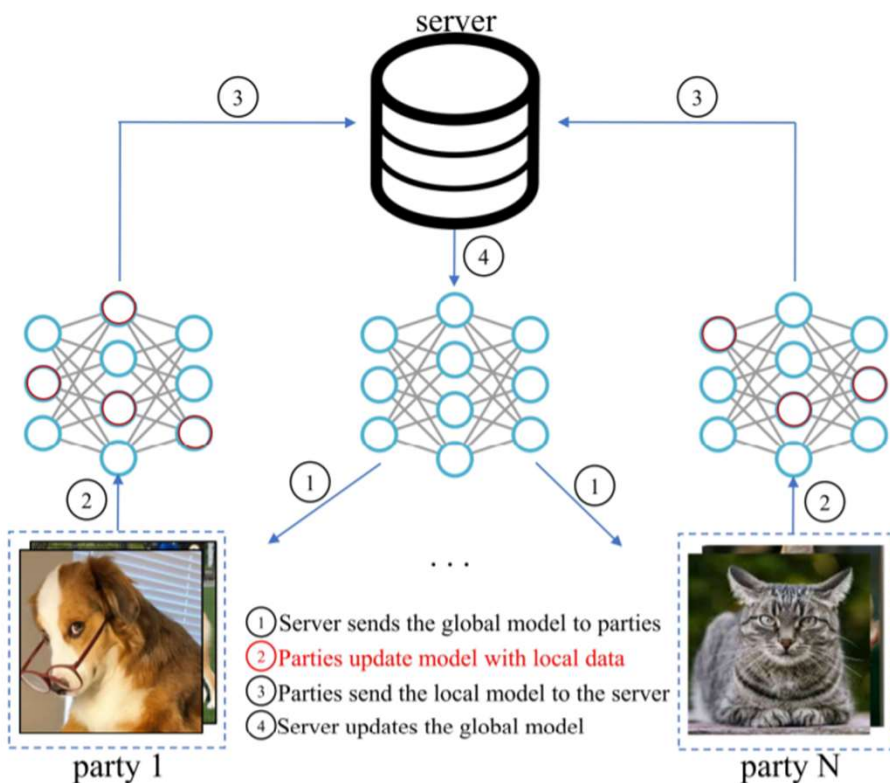


Figure from “Learning in Heterogeneous Networks: Optimization and Fairness”

Tackling Non-IID Data



- Local update stage ②

FedProx: directly limits the local updates by ℓ_2 -norm distance

SCAFFOLD: corrects the local updates via variance reduction

Those studies have little or even no advantage over FedAvg when deep models are used for training.

- Aggregation stage ④

FedMA: match and average weights in a layer-wise manner.

FedAvgM: applies momentum for global model updates

FedNova: normalizes the local updates before averaging

Contribution

- Address the non-IID issue from a novel perspective, i.e., model representation in the local update stage
- Propose model-contrastive learning (MOON), which conducts contrastive learning in model-level by comparing the representations learned by different models.

Outline

- Motivation
- **Model-Contrastive Federated Learning**
- Experimental Results
- Conclusion

Model-Contrastive Federated Learning

Contrastive learning –SimCLR [1]

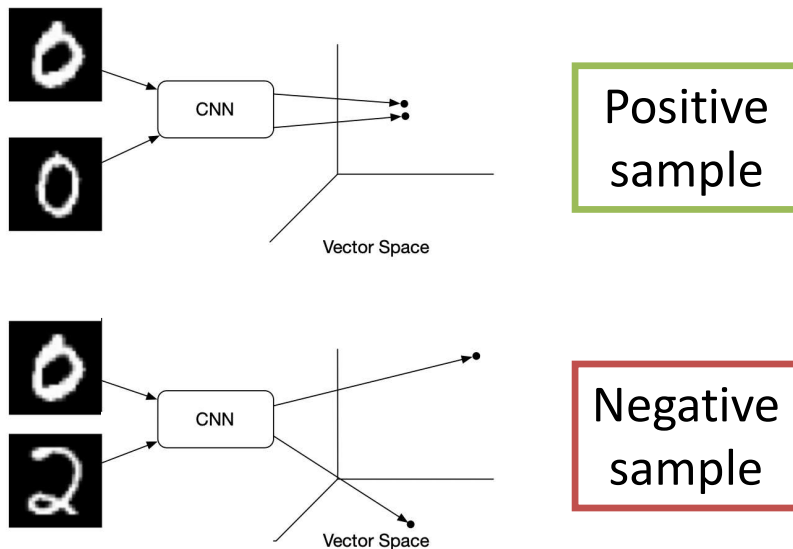
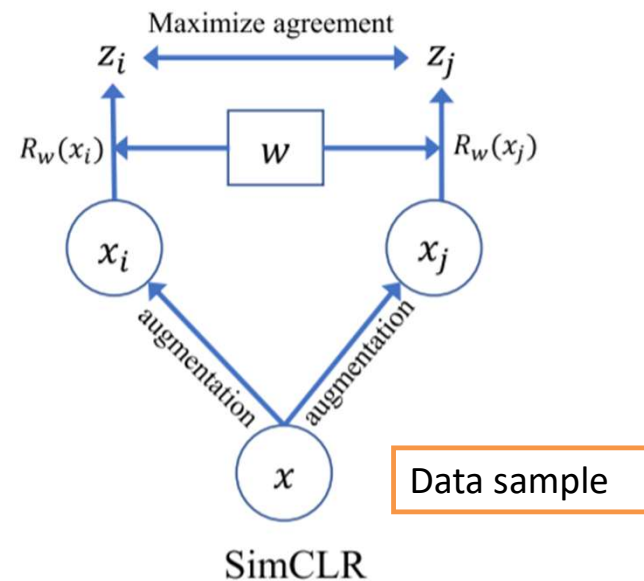


Figure from “Contrastive Loss Explained”

unsupervised learning



Here x denotes an image, w denotes a model, and R denotes the function to compute representation.

Model-Contrastive Federated Learning

Observation: the global model trained on a *whole dataset* is able to learn a better representation than the local model trained on a *skewed subset*.

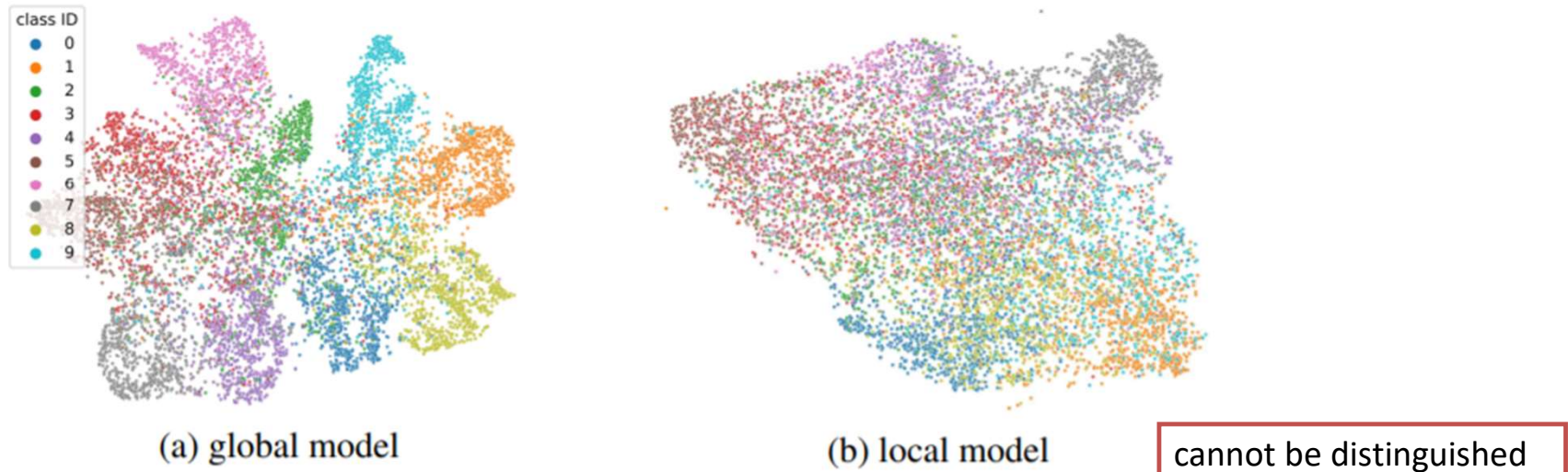


Figure 2. T-SNE visualizations of hidden vectors on CIFAR-10.

Model-Contrastive Federated Learning

FedAvg global model can learn better than the local model.

The local training phase even leads the model to learn a **worse representation**.



(b) local model



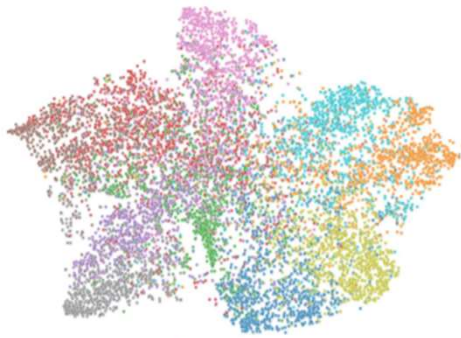
(c) FedAvg global model



(d) FedAvg local model

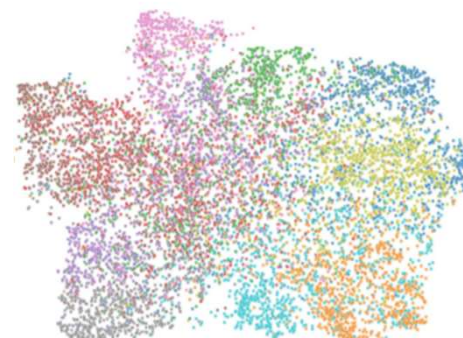
Model-Contrastive Federated Learning

MOON aims to decrease the distance between the representation learned by the local model and the representation learned by the global model



(c) FedAvg global model

Positive model

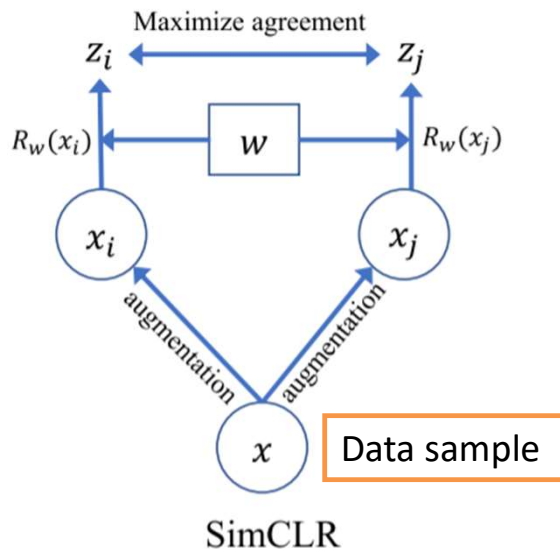


(d) FedAvg local model

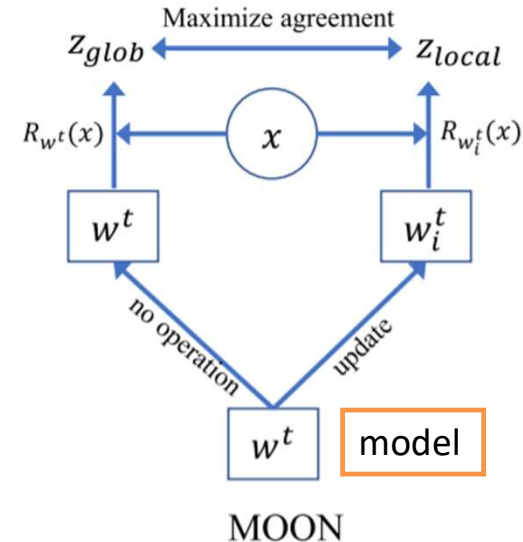
Negative model

Model-Contrastive Federated Learning

unsupervised learning

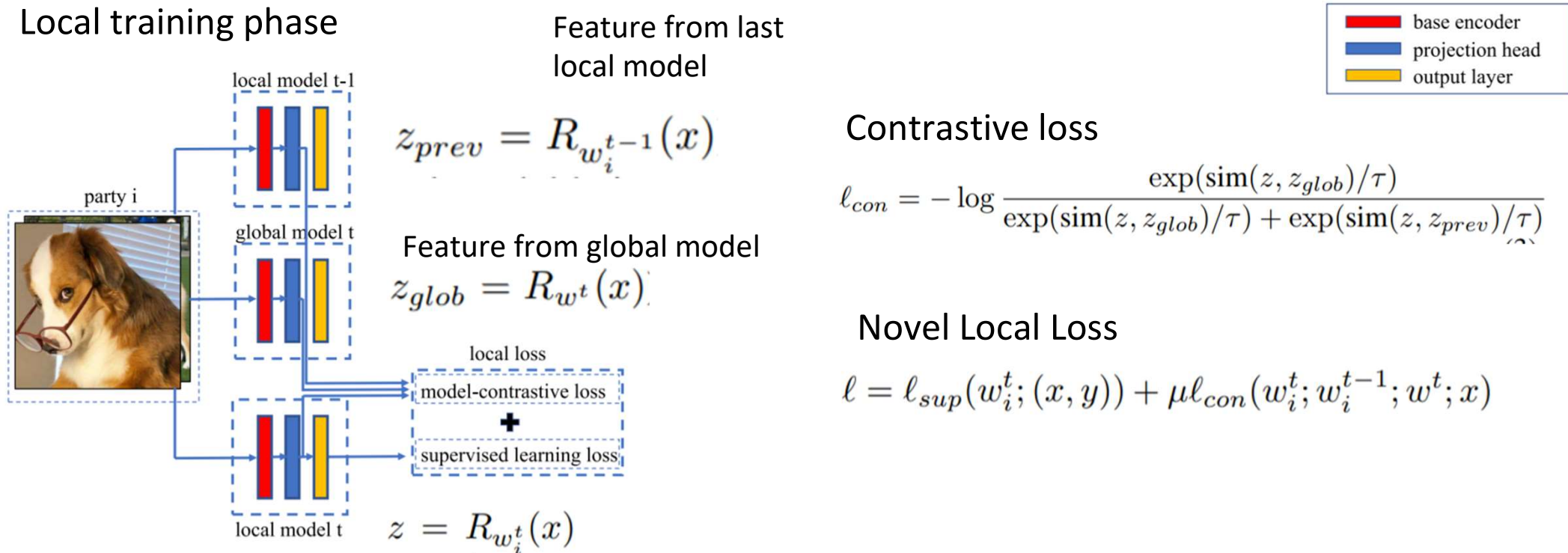


supervised learning



Here x denotes an image, w denotes a model, and R denotes the function to compute representation.

Model-Contrastive Federated Learning



- Our approach is robust regardless of different amount of drifts.

Model-Contrastive Federated Learning

Algorithm 1: The MOON framework

Input: number of communication rounds T ,
 number of parties N , number of local
 epochs E , temperature τ , learning rate η ,
 hyper-parameter μ
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    $w^{t+1} \leftarrow \sum_{k=1}^N \frac{|\mathcal{D}^k|}{|\mathcal{D}|} w_k^t$ 
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
    
```

Outline

- Motivation
- Model-Contrastive Federated Learning
- **Experimental Results**
- Conclusion

Experiment Setup

- We use PyTorch to implement MOON and the other baselines
- Dataset: CIFAR10, CIFAR100, tiny-ImageNet
- Encoder: CNN for CIFAR10; ResNet-50 for CIFAR100 & tiny-ImageNet
- Data partition: apply Dirichlet distribution ($\text{Dir}(\beta)$) to generate the non-IID data partition among parties.
- A 2- layer MLP is used as the projection head.
- Baseline: FedAvg ; FedProx ; SCAFFOLD and Local training

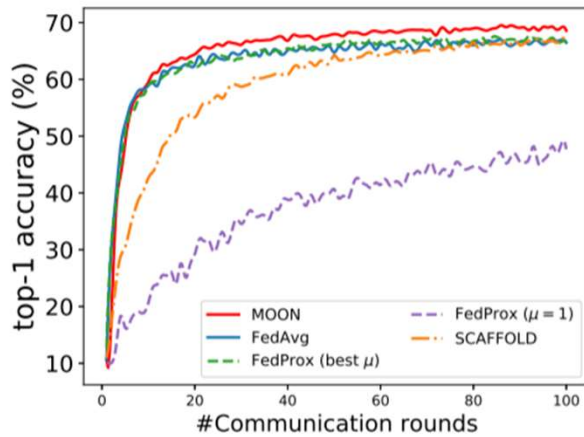


Performance evaluation

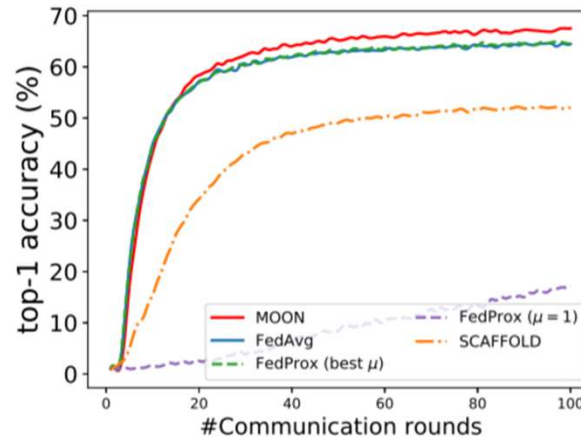
Table 1. The top-1 accuracy of MOON and the other baselines on test datasets. For MOON, FedAvg, FedProx, and SCAFFOLD, we run three trials and report the mean and standard derivation. For SOLO, we report the mean and standard derivation among all parties.

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

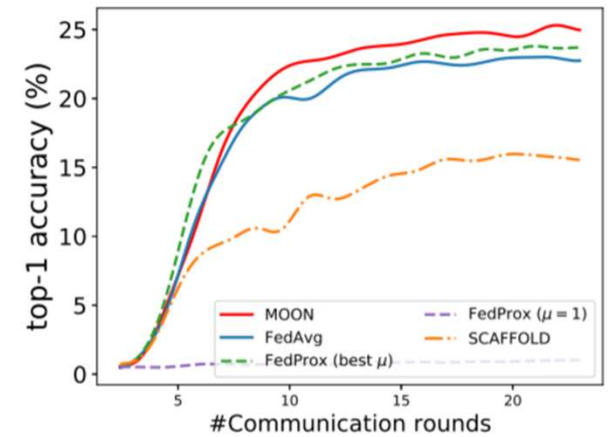
Performance evaluation



(a) CIFAR-10



(b) CIFAR-100



(c) Tiny-Imagenet

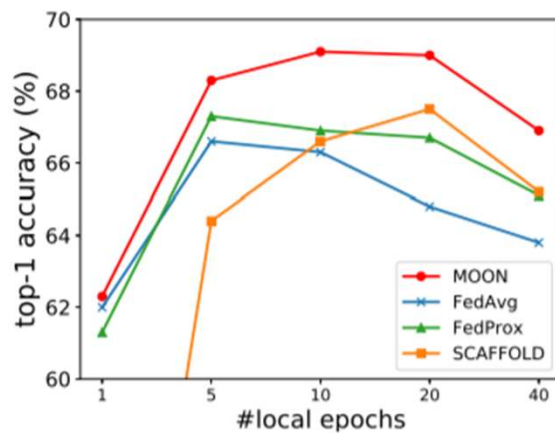
- MOON is much more communication-efficient than the other approaches.

Table 2. The number of rounds of different approaches to achieve the same accuracy as running FedAvg for 100 rounds (CIFAR-10/100) or 20 rounds (Tiny-Imagenet). The speedup of an approach is computed against FedAvg.

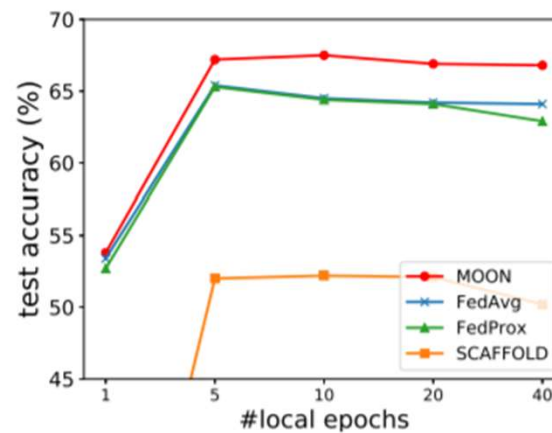
Method	CIFAR-10		CIFAR-100		Tiny-Imagenet	
	#rounds	speedup	#rounds	speedup	#rounds	speedup
FedAvg	100	1×	100	1×	20	1×
FedProx	52	1.9×	75	1.3×	17	1.2×
SCAFFOLD	80	1.3×	—	<1×	—	<1×
MOON	27	3.7×	43	2.3×	11	1.8×



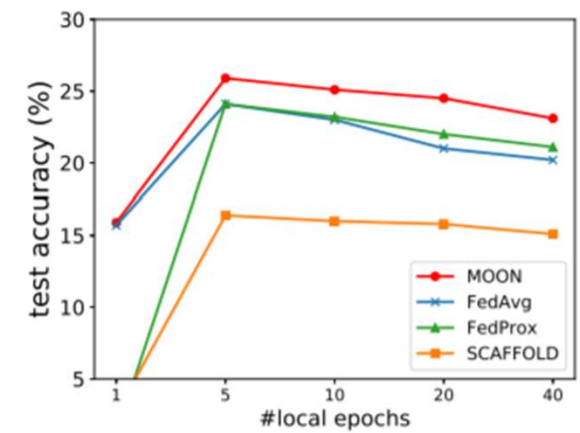
Performance evaluation



(a) CIFAR-10



(b) CIFAR-100



(c) Tiny-Imagenet

- MOON can effectively mitigate the negative effects of the drift by too many local updates.

Performance evaluation

- Scalability

Table 3. The accuracy with 50 parties and 100 parties (sample fraction=0.2) on CIFAR-100.

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 ($\mu=10$)	58.2%	63.2%	56.9%	61.8%
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%	

- Diversity

Table 4. The test accuracy with β from $\{0.1, 0.5, 5\}$.

Method	$\beta = 0.1$	$\beta = 0.5$	$\beta = 5$
MOON	64.0%	67.5%	68.0%
FedAvg	62.5%	64.5%	65.7%
FedProx	62.9%	64.6%	64.9%
SCAFFOLD	47.3%	52.5%	55.0%
SOLO	15.9% \pm 1.5%	22.3% \pm 1%	26.6% \pm 1.4%

- MOON can outperform the other approaches a lot with more communication rounds.

Performance evaluation

- The average training time per round

Table 11. The average training time per round.

Method	CIFAR-10	CIFAR-100	Tiny-Imagenet
FedAvg	330s	20min	103min
FedProx	340s	24min	135min
SCAFFOLD	332s	20min	112min
MOON	337s	31min	197min

Conclusion

- MOON is a simple and effective federated learning framework.
- MOON addresses the non-IID data issue with the novel design of model-based contrastive learning.

An aerial photograph of the University of Houston campus at dusk. The foreground shows several large, modern university buildings with flat roofs and some with glass facades. A central green lawn with winding paths is visible. In the background, the Houston city skyline is silhouetted against a twilight sky with soft orange and blue hues. A large, semi-transparent red rectangle is superimposed over the upper half of the image, containing the text "THANK YOU" in white, bold, sans-serif capital letters.

THANK YOU

UNIVERSITY of **HOUSTON** | ENGINEERING

Performance evaluation

- Selection of different local loss function

Table 5. The top-1 accuracy with different kinds of loss for the second term of local objective. We tune μ from $\{0.001, 0.01, 0.1, 1, 5, 10\}$ for the ℓ_2 norm approach and report the best accuracy.

second term	CIFAR-10	CIFAR-100	Tiny-Imagenet
none (FedAvg)	66.3%	64.5%	23.0%
ℓ_2 norm	65.8%	66.9%	24.0%
MOON	69.1%	67.5%	25.1%