



RICE

Junhyung Lyle Kim (Rice CS), Mohammad Taha Toghani (Rice ECE), CésarA.Uribe (Rice ECE), Anastasios Kyrillidis (Rice CS)

## Federated Learning Overview

- FL is a recent distributed machine learning framework where a global model is trained via multiple collaborative steps by participating clients without sharing data.
- Mathematically, we want to solve

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x)$$

where  $x \in \mathbb{R}^d$  is the shared model parameter,  $m$  is the number of clients, and  $f_i(x) := \mathbb{E}_{z \sim \mathcal{D}_i}[F_i(x, z)]$  is the individual loss function.

- Offers flexible collaborative learning: the number of clients  $m$ , their participation rates, and the computing power can all vary and change at any point during the overall training procedure.

$$f_i(x) \leftarrow \mathbb{E}_{z \sim \mathcal{D}_i}[F_i(x, z)]$$

- Therefore, not only  $\mathcal{D}_i$  differs for each client  $i$ , but also the number of samples  $z \sim \mathcal{D}_i$ , resulting in each client having different  $f_i(\cdot)$ .

## Challenges in Federated Learning

### Algorithm FedAvg

```

1: input:  $x_0 \in \mathbb{R}^d$ ,  $\eta > 0$ , and  $p \in (0, 1)$ .
2: for each round  $t = 0, 1, \dots, T-1$  do
3:   sample a subset  $\mathcal{S}_t$  of clients with size  $|\mathcal{S}_t| = p \cdot m$ 
4:   for each machine in parallel for  $i \in \mathcal{S}_t$  do
5:     Set  $x_{t,0}^i = x_t$ 
6:     for local step  $k \in [K]$  do
7:       Compute an estimate  $g_{t,k-1}^i$  of  $\nabla f_i(x_{t,k-1}^i)$ 
8:        $x_{t,k}^i = x_{t,k-1}^i - \eta g_{t,k-1}^i$ 
9:     end for
10:   end for
11:    $x_{t+1} = \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} x_{t,K}^i = x_t - \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} (x_t - x_{t,K}^i)$ 
12: end for
13: return  $x_T$ 
```

**Server-side:** How do we smartly aggregate the local information coming from each participating client?

- Federated Averaging [McMahan et al., 2017] uses simple averaging
- [Reddi et al., 2021] interpreted averaging as a “pseudo-gradient” step and introduced FedAdam, FedYogi, FedAdagrad, etc.

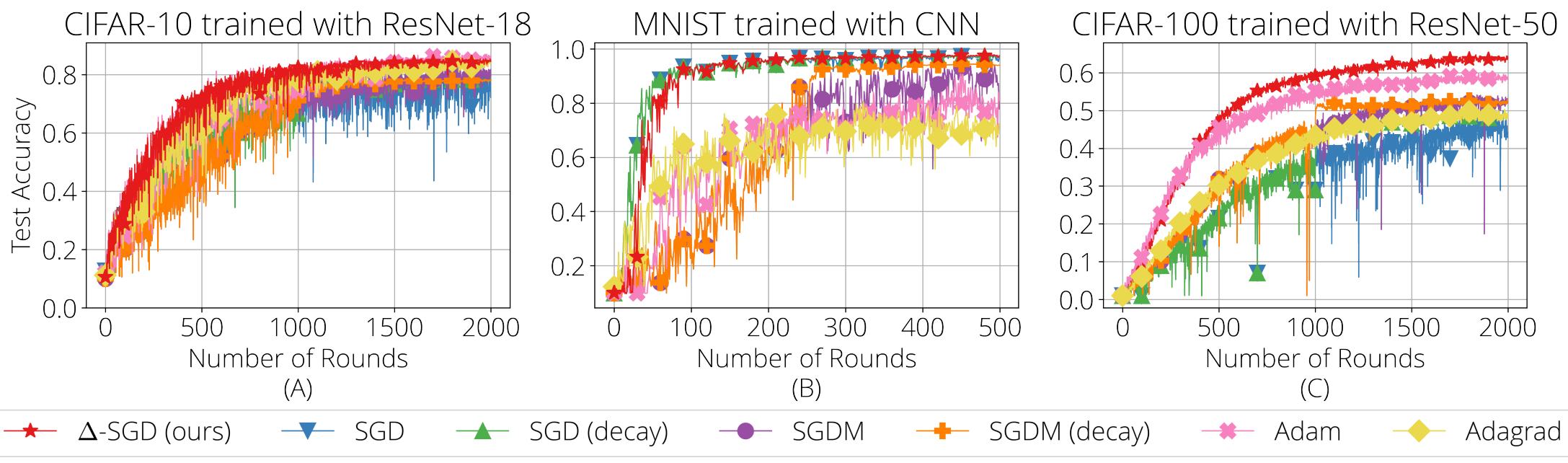
**Client-side:** How do we make sure each client meaningfully “learns” using local data?

- Federated Averaging [McMahan et al., 2017] uses SGD.
- Does it make sense to use the same  $\eta$  for all clients? If not, how should we tune individual step sizes?
- Important to properly tune: many more local updates compared to the aggregation step, as communication is much more expensive.

## Experimental Results: $\Delta$ -SGD exhibits superior performance in all settings without any additional tuning

Non-iidness	Optimizer	Dataset / Model				
$\alpha = 1$	Dir( $\alpha \cdot \mathbf{p}$ )	MNIST CNN	FMNIST CNN	CIFAR-10 ResNet-18	CIFAR-100 ResNet-18	CIFAR-100 ResNet-50
	SGD	98.3 <sub>(0.2)</sub>	86.5 <sub>(0.8)</sub>	87.7 <sub>(2.1)</sub>	57.7 <sub>(4.2)</sub>	53.0 <sub>(12.8)</sub>
	SGD ( $\downarrow$ )	97.8 <sub>(0.7)</sub>	86.3 <sub>(1.0)</sub>	87.8 <sub>(2.0)</sub>	61.9 <sub>(0.0)</sub>	60.9 <sub>(4.9)</sub>
	SGDM	98.5 <sub>(0.0)</sub>	85.2 <sub>(2.1)</sub>	88.7 <sub>(1.1)</sub>	58.8 <sub>(3.1)</sub>	60.5 <sub>(5.3)</sub>
	SGDM ( $\downarrow$ )	98.4 <sub>(0.1)</sub>	87.2 <sub>(0.1)</sub>	89.3 <sub>(0.5)</sub>	61.4 <sub>(0.5)</sub>	63.3 <sub>(2.5)</sub>
	Adam	94.7 <sub>(3.8)</sub>	71.8 <sub>(15.5)</sub>	89.4 <sub>(0.4)</sub>	55.6 <sub>(6.3)</sub>	61.4 <sub>(4.4)</sub>
	Adagrad	64.3 <sub>(34.2)</sub>	45.5 <sub>(41.8)</sub>	86.6 <sub>(3.2)</sub>	53.5 <sub>(8.4)</sub>	51.9 <sub>(13.9)</sub>
$\alpha = 0.1$	SPS	10.1 <sub>(88.4)</sub>	85.9 <sub>(1.4)</sub>	82.7 <sub>(7.1)</sub>	1.0 <sub>(60.9)</sub>	50.0 <sub>(15.8)</sub>
	$\Delta$ -SGD	98.4 <sub>(0.1)</sub>	87.3 <sub>(0.0)</sub>	89.8 <sub>(0.0)</sub>	61.5 <sub>(0.4)</sub>	65.8 <sub>(0.0)</sub>
	SGD	98.1 <sub>(0.0)</sub>	83.6 <sub>(2.8)</sub>	72.1 <sub>(12.9)</sub>	54.4 <sub>(6.7)</sub>	44.2 <sub>(19.9)</sub>
	SGD ( $\downarrow$ )	98.0 <sub>(0.1)</sub>	84.7 <sub>(1.7)</sub>	78.4 <sub>(6.6)</sub>	59.3 <sub>(1.8)</sub>	48.7 <sub>(15.4)</sub>
	SGDM	97.6 <sub>(0.5)</sub>	83.6 <sub>(2.8)</sub>	79.6 <sub>(5.4)</sub>	58.8 <sub>(2.3)</sub>	52.3 <sub>(11.8)</sub>
	SGDM ( $\downarrow$ )	98.0 <sub>(0.1)</sub>	86.1 <sub>(0.3)</sub>	77.9 <sub>(7.1)</sub>	60.4 <sub>(0.7)</sub>	52.8 <sub>(11.3)</sub>
	Adam	96.4 <sub>(1.7)</sub>	80.4 <sub>(6.0)</sub>	85.0 <sub>(0.0)</sub>	55.4 <sub>(5.7)</sub>	58.2 <sub>(5.9)</sub>
$\alpha = 0.01$	Adagrad	89.9 <sub>(8.2)</sub>	46.3 <sub>(40.1)</sub>	84.1 <sub>(0.9)</sub>	49.6 <sub>(11.5)</sub>	48.0 <sub>(16.1)</sub>
	SPS	96.0 <sub>(2.1)</sub>	85.0 <sub>(1.4)</sub>	70.3 <sub>(14.7)</sub>	42.2 <sub>(18.9)</sub>	42.2 <sub>(21.9)</sub>
	$\Delta$ -SGD	98.1 <sub>(0.0)</sub>	86.4 <sub>(0.0)</sub>	84.5 <sub>(0.5)</sub>	61.1 <sub>(0.0)</sub>	64.1 <sub>(0.0)</sub>
	SGD	96.8 <sub>(0.7)</sub>	79.0 <sub>(1.2)</sub>	22.6 <sub>(11.3)</sub>	30.5 <sub>(1.3)</sub>	24.3 <sub>(7.1)</sub>
	SGD ( $\downarrow$ )	97.2 <sub>(0.3)</sub>	79.3 <sub>(0.9)</sub>	33.9 <sub>(0.0)</sub>	30.3 <sub>(1.5)</sub>	24.6 <sub>(6.8)</sub>
	SGDM	77.9 <sub>(19.6)</sub>	75.7 <sub>(4.5)</sub>	28.4 <sub>(5.5)</sub>	24.8 <sub>(7.0)</sub>	22.0 <sub>(9.4)</sub>
	SGDM ( $\downarrow$ )	94.0 <sub>(3.5)</sub>	79.5 <sub>(0.7)</sub>	29.0 <sub>(4.9)</sub>	20.9 <sub>(10.9)</sub>	14.7 <sub>(16.7)</sub>
	Adam	80.8 <sub>(16.7)</sub>	60.6 <sub>(19.6)</sub>	22.1 <sub>(11.8)</sub>	18.2 <sub>(13.6)</sub>	22.6 <sub>(8.8)</sub>
	Adagrad	72.4 <sub>(25.1)</sub>	45.9 <sub>(34.3)</sub>	12.5 <sub>(21.4)</sub>	25.8 <sub>(6.0)</sub>	22.2 <sub>(9.2)</sub>
	SPS	69.7 <sub>(27.8)</sub>	44.0 <sub>(36.2)</sub>	21.5 <sub>(12.4)</sub>	22.0 <sub>(9.8)</sub>	17.4 <sub>(14.0)</sub>
	$\Delta$ -SGD	97.5 <sub>(0.0)</sub>	80.2 <sub>(0.0)</sub>	31.6 <sub>(2.3)</sub>	31.8 <sub>(0.0)</sub>	31.4 <sub>(0.0)</sub>

## Client Optimization Is More Challenging?



- We fine-tune the step size for each client optimizer in task (A), and *intentionally keep it the same* for the other tasks, to highlight the effect of not properly tuning the step size of each client optimizer. Our proposed method,  $\Delta$ -SGD, exhibits superior performance in all settings, without any additional tuning.
- (A): CIFAR-10 classification task trained on ResNet-18. (B): MNIST classification task trained on shallow CNN. (C): CIFAR-100 classification task trained on ResNet-50.
- All experiments use FedAvg as the server optimizer.

## Why is the step size $\eta = 1/L$ popular?

- $L$ -smooth functions:  $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2 \quad \forall x, y$
- Gradient descent:  $x_{t+1} = x_t - \eta \nabla f(x_t)$
- Descent lemma:

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{L}{2} \|x_{t+1} - x_t\|^2 \\ &= f(x_t) - \eta \left(1 - \frac{\eta \cdot L}{2}\right) \|\nabla f(x_t)\|^2 \end{aligned}$$

- $\eta = 1/L$  maximizes the descent progress.

## Adaptive Step Size via Local Smoothness

- [Malisky & Mishchenko, 2020] proposed the following step size for (centralized) gradient descent:

$$\eta_t = \min \left\{ \frac{\|x_t - x_{t-1}\|}{2\|\nabla f(x_t) - \nabla f(x_{t-1})\|}, \sqrt{1 + \frac{\eta_{t-1}}{\eta_{t-2}}} \eta_{t-1} \right\}$$

- The first condition approximates the local smoothness
- $\|\nabla f(x_t) - \nabla f(x_{t-1})\| \leq L_t \cdot \|x_t - x_{t-1}\|, \quad \forall t = 1, 2, \dots$  and the second condition ensures  $\eta_t$  to not increase too fast.
- We adapt the above step size to the FL setting:

$$\eta_{t,k}^i = \min \left\{ \frac{\|x_{t,k}^i - x_{t,k-1}^i\|}{2\|\tilde{\nabla} f_i(x_{t,k}^i) - \tilde{\nabla} f_i(x_{t,k-1}^i)\|}, \sqrt{1 + \frac{\eta_{t,k-1}^i}{\eta_{t,k-1}^i}} \eta_{t,k-1}^i \right\}$$

with the stochastic gradients  $\tilde{\nabla} f_i(\cdot)$  and the local iterations  $k$ .

- Implication:** each client uses its own step size  $\eta_t^i$  that is adaptive to the local smoothness of  $f_i(\cdot)$

