



南京航空航天大学

Nanjing University of Aeronautics and Astronautics

南京航空航天大学

Nanjing University of Aeronautics and Astronautics

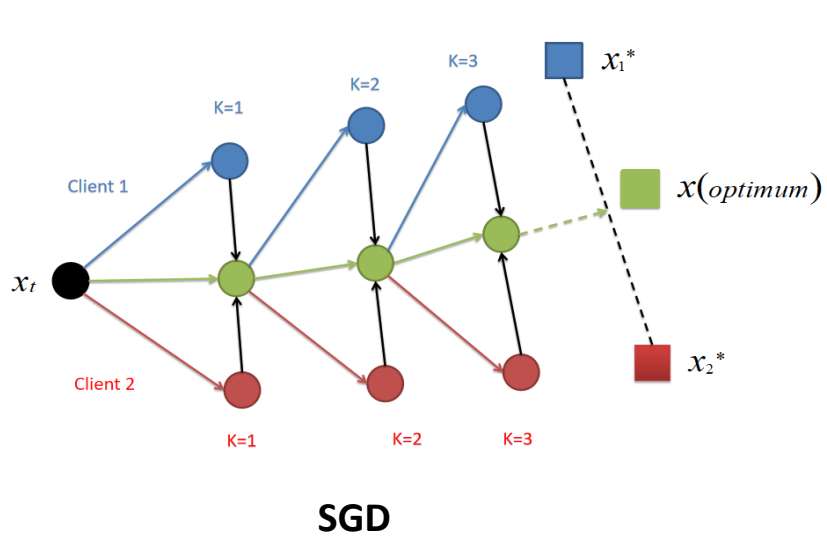
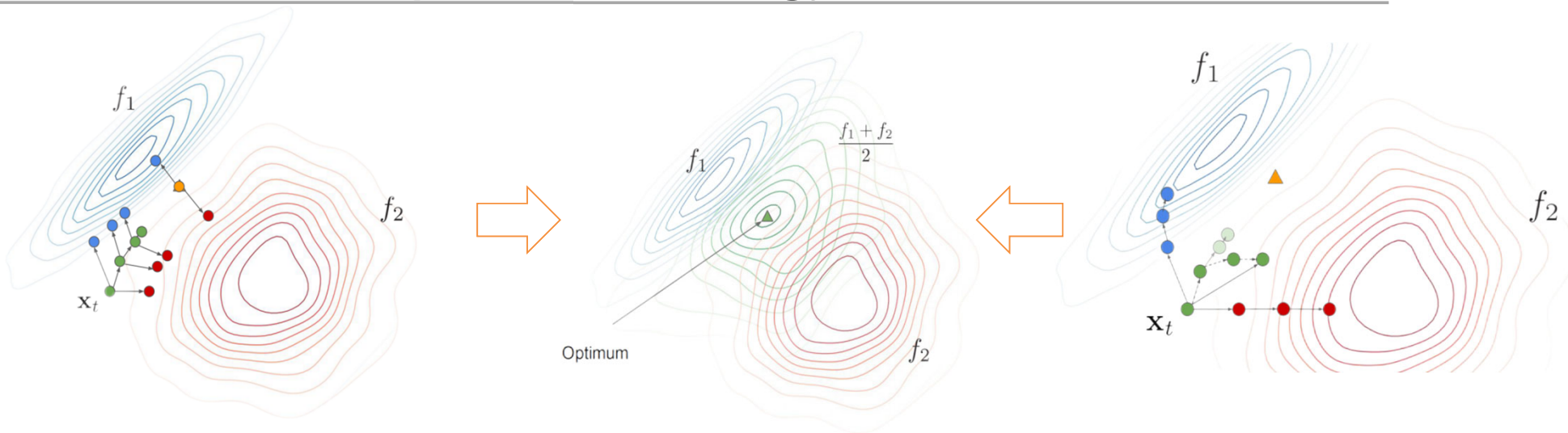


SCAFFOLD: Stochastic Controlled Averaging for Federated Learning

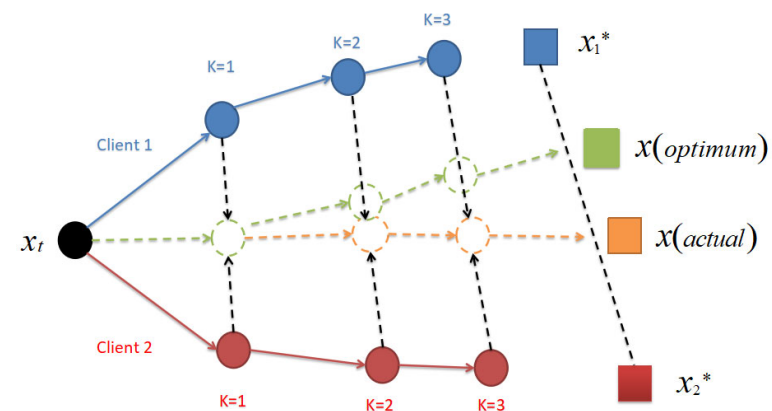
Sai Praneeth Karimireddy^{1,2} Satyen Kale³ Mehryar Mohri^{3,4} Sashank J. Reddi³ Sebastian U. Stich¹
Ananda Theertha Suresh³

ICML 2020

Motivation (Convergence of FedAvg)



SGD



FedAvg

Method (SGD+FedAvg)

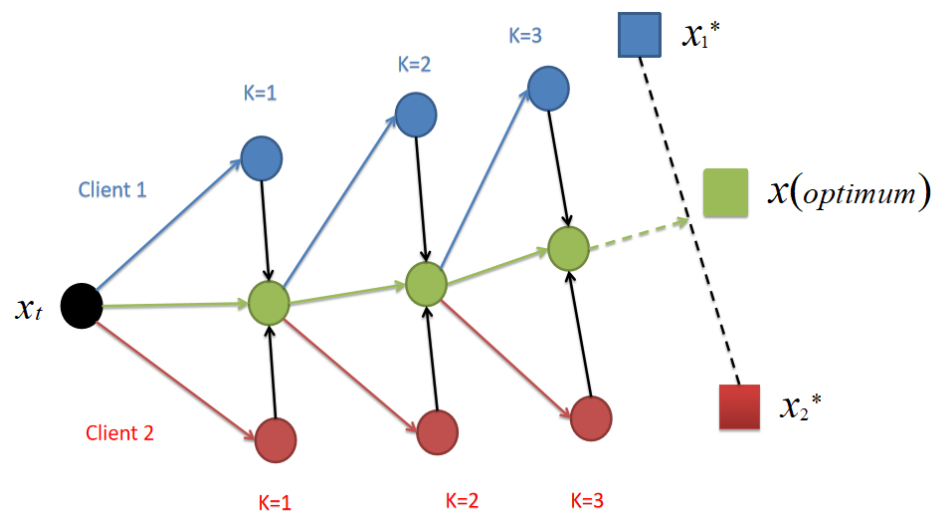
1. How to accelerate global model convergence while ensuring low communication?

① SGD (Good convergence speed + Bad communication overheads)

② FedAvg (low communication overheads + Bad convergence speed)

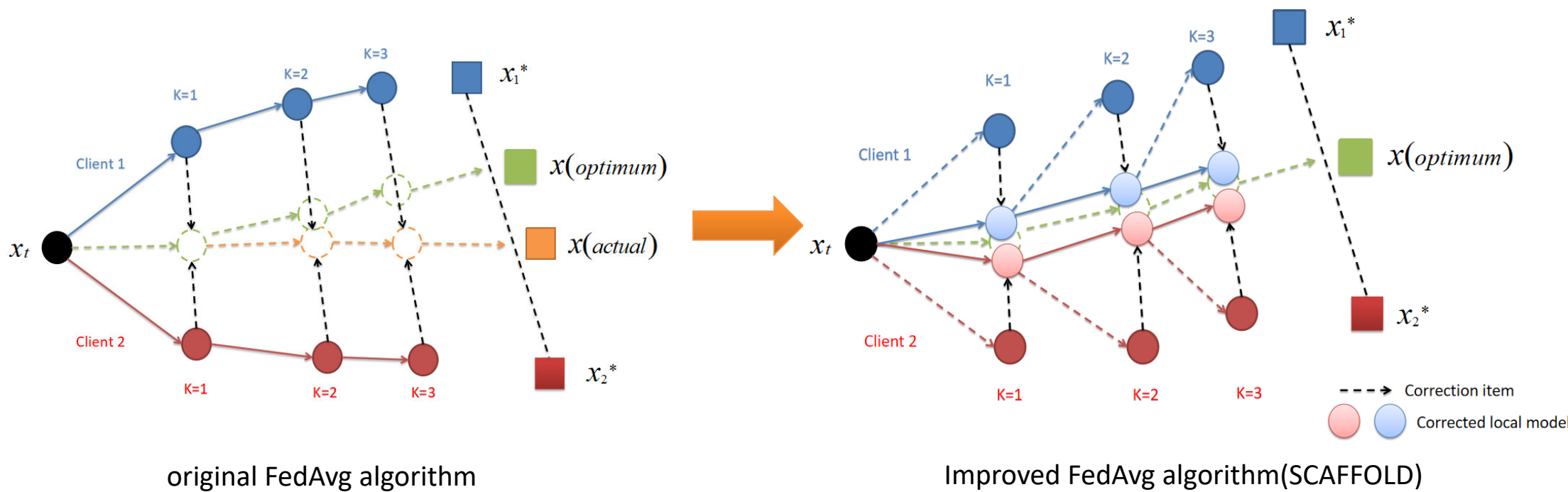
A new algorithm with good convergence speed and low communication overheads

2. SGD features in updating the global model



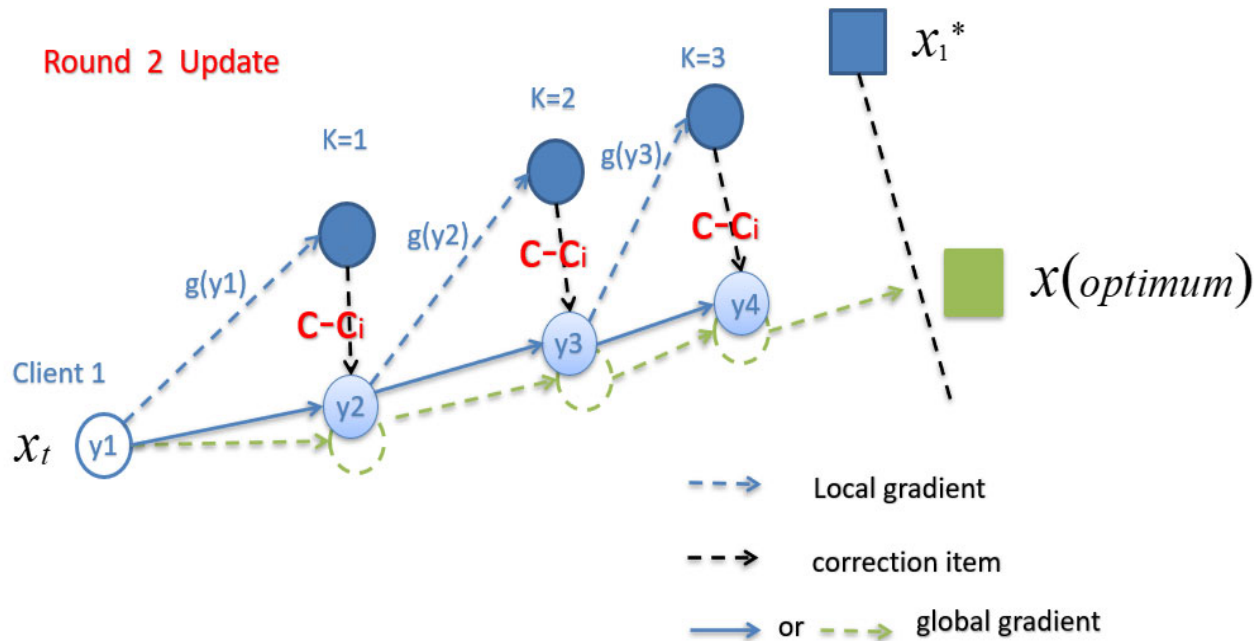
Method (SGD+FedAvg)

3. How to add features from SGD to FedAvg?

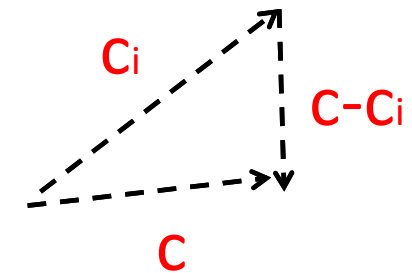


Method (SGD+FedAvg)

4. How to find correction item?



Round 1 update



C_i: average gradient of client 1 local k updates in Round 1.

C: global average gradient, The first round of updates is aggregated (C_i) and the second round is sent to each client

$$c = \frac{1}{N} \sum c_i$$

Method (SGD+FedAvg)

5. Two ways to update client model average gradient

$$\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta_l (g_i(\mathbf{y}_i) + \mathbf{c} - \mathbf{c}_i).$$

$$\mathbf{c}_i^+ \leftarrow \begin{cases} \text{Option I.} & g_i(\mathbf{x}), \text{ or} \\ \text{Option II.} & \mathbf{c}_i - \mathbf{c} + \frac{1}{K\eta_l} (\mathbf{x} - \mathbf{y}_i). \end{cases}$$

option I. Calculate the client gradient using the initial global model \mathbf{x} (provided by the server)

Algorithm 2 FEDAVG: Federated Averaging

```
1: server input: initial  $x$ , and global step-size  $\eta_g$ 
2: client  $i$ 's input: local step-size  $\eta_l$ 
3: for each round  $r = 1, \dots, R$  do
4:   sample clients  $\mathcal{S} \subseteq \{1, \dots, N\}$ 
5:   communicate  $x$  to all clients  $i \in \mathcal{S}$ 
6:   on client  $i \in \mathcal{S}$  in parallel do
7:     initialize local model  $y_i \leftarrow x$ 
8:     for  $k = 1, \dots, K$  do
9:       compute mini-batch gradient  $g_i(y_i)$ 
10:       $y_i \leftarrow y_i - \eta_l g_i(y_i)$ 
11:    end for
12:    communicate  $\Delta y_i \leftarrow y_i - x$ 
13:  end on client
14:   $\Delta x \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \Delta y_i$ 
15:   $x \leftarrow x + \eta_g \Delta x$ 
16: end for
```

Algorithm 1 SCAFFOLD: Stochastic Controlled Averaging for federated learning

```
1: server input: initial  $x$  and  $c$ , and global step-size  $\eta_g$ 
2: client  $i$ 's input:  $c_i$ , and local step-size  $\eta_l$ 
3: for each round  $r = 1, \dots, R$  do
4:   sample clients  $\mathcal{S} \subseteq \{1, \dots, N\}$ 
5:   communicate  $(x, c)$  to all clients  $i \in \mathcal{S}$ 
6:   on client  $i \in \mathcal{S}$  in parallel do
7:     initialize local model  $y_i \leftarrow x$ 
8:     for  $k = 1, \dots, K$  do
9:       compute mini-batch gradient  $g_i(y_i)$ 
10:       $y_i \leftarrow y_i - \eta_l (g_i(y_i) - c_i + c)$ 
11:    end for
12:     $c_i^+ \leftarrow$  (i)  $g_i(x)$ , or (ii)  $c_i - c + \frac{1}{K\eta_l} (x - y_i)$ 
13:    communicate  $(\Delta y_i, \Delta c_i) \leftarrow (y_i - x, c_i^+ - c_i)$ 
14:     $c_i \leftarrow c_i^+$ 
15:  end on client
16:   $(\Delta x, \Delta c) \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\Delta y_i, \Delta c_i)$ 
17:   $x \leftarrow x + \eta_g \Delta x$  and  $c \leftarrow c + \frac{|\mathcal{S}|}{N} \Delta c$ 
18: end for
```

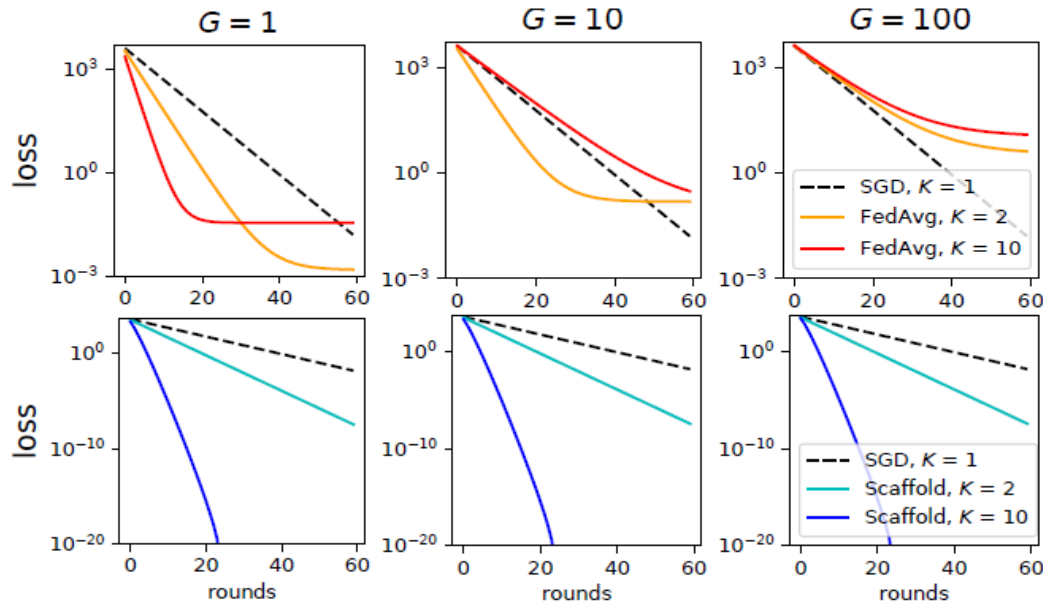


Figure 3. SGD (dashed black), FedAvg (above), and SCAFFOLD (below) on simulated data. FedAvg gets worse as local steps increases with $K = 10$ (red) worse than $K = 2$ (orange). It also gets slower as gradient-dissimilarity (G) increases (to the right). SCAFFOLD significantly improves with more local steps, with $K = 10$ (blue) faster than $K = 2$ (light blue) and SGD. Its performance is identical as we vary heterogeneity (G).

(A1) (G, B) -BGD or bounded gradient dissimilarity: there exist constants $G \geq 0$ and $B \geq 1$ such that

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + B^2 \|\nabla f(\mathbf{x})\|^2, \forall \mathbf{x}.$$

If $\{f_i\}$ are convex, we can relax the assumption to

$$\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\mathbf{x})\|^2 \leq G^2 + 2\beta B^2 (f(\mathbf{x}) - f^*), \forall \mathbf{x}.$$

Experiments

Table 3. Communication rounds to reach 0.5 test accuracy for logistic regression on EMNIST as we vary number of epochs. 1k+ indicates 0.5 accuracy was not reached even after 1k rounds, and similarly an arrowhead indicates that the barplot extends beyond the table. 1 epoch for local update methods corresponds to 5 local steps (0.2 batch size), and 20% of clients are sampled each round.

















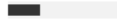
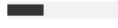


















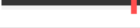
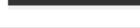
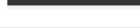
	Epochs	0% similarity (sorted)		10% similarity		100% similarity (i.i.d.)	
		Num. of rounds	Speedup	Num. of rounds	Speedup	Num. of rounds	Speedup
SGD	1	317 	(1×)	365 	(1×)	416 	(1×)
SCAFFOLD1		77 	(4.1×)	62 	(5.9×)	60 	(6.9×)
	5	152 	(2.1×)	20 	(18.2×)	10 	(41.6×)
	10	286 	(1.1×)	16 	(22.8×)	7 	(59.4×)
	20	266 	(1.2×)	11 	(33.2×)	4 	(104×)
FEDAVG	1	258 	(1.2×)	74 	(4.9×)	83 	(5×)
	5	428 	(0.7×)	34 	(10.7×)	10 	(41.6×)
	10	711 	(0.4×)	25 	(14.6×)	6 	(69.3×)
	20	1k+ 	(< 0.3×)	18 	(20.3×)	4 	(104×)
FEDPROX	1	1k+ 	(< 0.3×)	979 	(0.4×)	459 	(0.9×)
	5	1k+ 	(< 0.3×)	794 	(0.5×)	351 	(1.2×)
	10	1k+ 	(< 0.3×)	894 	(0.4×)	308 	(1.4×)
	20	1k+ 	(< 0.3×)	916 	(0.4×)	351 	(1.2×)

Table 4. Communication rounds to reach 0.45 test accuracy for logistic regression on EMNIST as we vary the number of sampled clients. Number of epochs is kept fixed to 5. SCAFFOLD is consistently faster than FEDAVG. As we decrease the number of clients sampled in each round, the increase in number of rounds is sub-linear. This slow-down is better for more similar clients.

	Clients	0% similarity		10% similarity	
SCAFFOLD	20%	143	(1.0×)	9	(1.0×)
	5%	290	(2.0×)	13	(1.4×)
	1%	790	(5.5×)	28	(3.1×)
FEDAVG	20%	179	(1.0×)	12	(1.0×)
	5%	334	(1.9×)	17	(1.4×)
	1%	1k+	(5.6+×)	35	(2.9×)



南京航空航天大学

Nanjing University of Aeronautics and Astronautics

南京航空航天大学

Nanjing University of Aeronautics and Astronautics



THANKS