

Optimization algorithms for heterogeneous clients in Federated Learning

Praneeth Karimireddy, Martin Jaggi, **Satyen Kale**, Mehryar Mohri,
Sashank Reddi, Sebastian Stich, Ananda Theertha Suresh

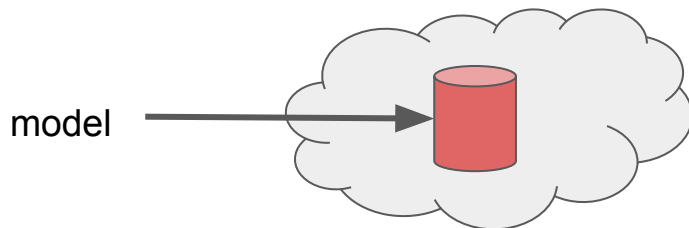
Tight analysis of FedAvg when clients are heterogeneous (non iid data)

Explain degradation of FedAvg via the '*drift*' in the client updates

Prove SCAFFOLD is resilient to heterogeneity and client sampling

Federated Learning: Setting

[McMahan et al. 2016]



Server
(e.g. Google)



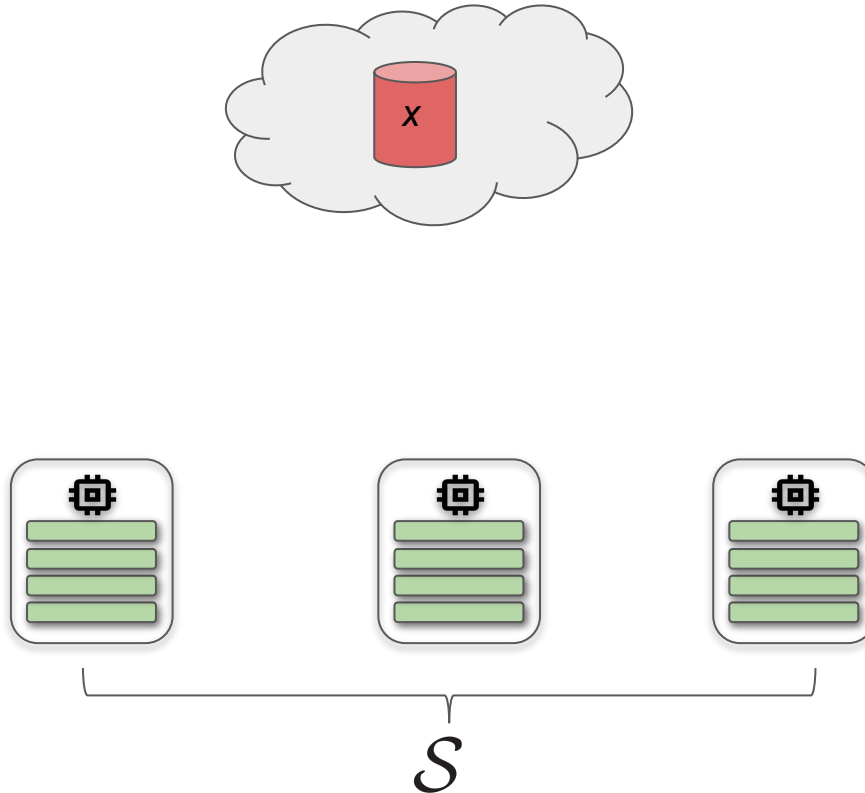
Clients
(e.g. hospitals, phones)

Federated Learning: Setting

[McMahan et al. 2016]

In each round,

- Some subset of clients are chosen

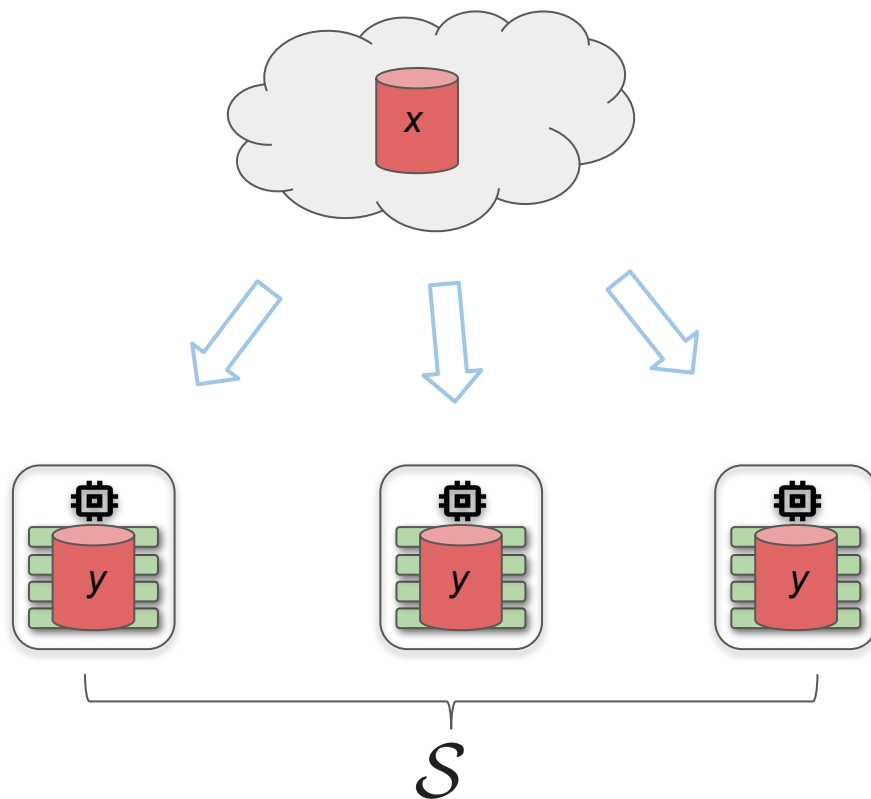


Federated Learning: Setting

[McMahan et al. 2016]

In each round,

- Some subset of clients are chosen
- copy of server model is sent to clients

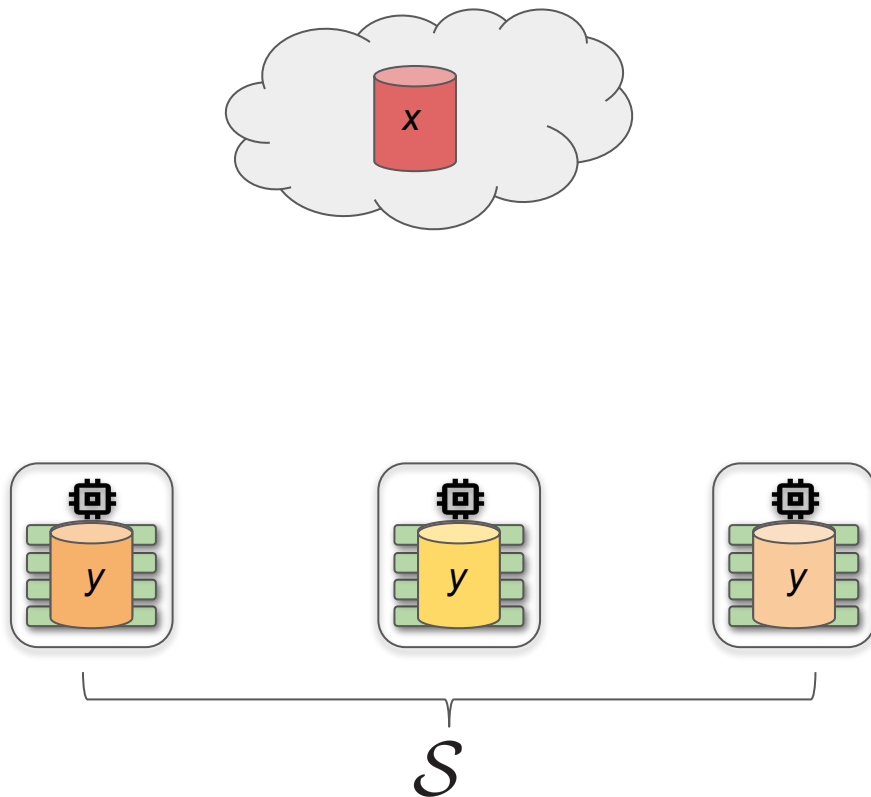


Federated Learning: Setting

[McMahan et al. 2016]

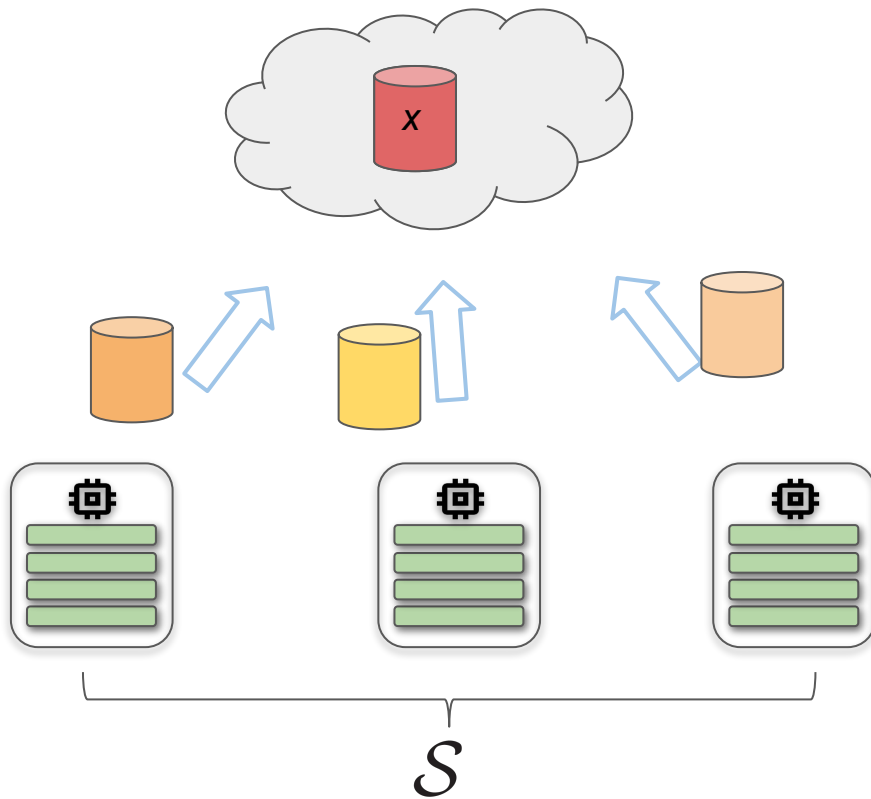
In each round,

- Some subset of clients are chosen
- copy of server model is sent to clients
- model is updated using client data



Federated Learning: Setting

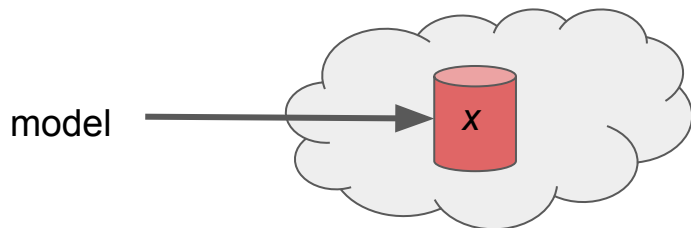
[McMahan et al. 2016]



In each round,

- Some subset of clients are chosen
- copy of server model is sent to clients
- model is updated using client data
- Client updates are aggregated
- server model is updated

Federated Learning: Characteristics



- High overhead per round
- Only a few clients participate in each round
- The data of the clients is **heterogeneous**: i.e. data drawn from different distributions for different clients



Cross-Silo vs Cross-Device Federated Learning

Cross-Silo FL

- **Small/medium** number (2-100, typically) of total clients
- E.g. hospitals, financial organizations
- **Large amount** of data per client
- **Persistent** clients: almost always available
- **Stateful** clients: clients can carry state from round to round

Cross-Device FL

- **Very large** number (e.g. 10^{10}) of total clients
- E.g. mobile or IoT devices
- **Small amount** of data per client
- **Transient** clients: only a fraction of clients available at any time
- **Stateless** clients: clients generally participate only once in each task

Cross-Silo Federated Learning: Formalism

$$\min_{x \in \mathbb{R}^n} \frac{1}{N} \sum_{i=1}^N E_{\zeta} [f_i(x; \zeta)]$$

Model parameters

average over clients

Expectation over client data

Loss function wrt parameters and client data

Algorithms for Cross-Silo Federated Learning

Solving FL: SGD

- on each client i in \mathcal{S} , compute a large batch stochastic gradient and average them

+ Equivalent to synchronous centralized large-batch training

$$x = x - \left[\frac{1}{S} \sum_{i \in \mathcal{S}} \left[\frac{1}{K} \sum_{k=1}^K g_i(x) \right] \right]$$

Average over all sampled clients

- very slow (only 1 update)

Mini-batch gradient with Batch size K per client

Solving FL: Federated Averaging (FedAvg) [McMahan et al. 2016]

- on each client i in \mathcal{S} , perform K steps of SGD

$$y_i = y_i - \eta g_i(y_i)$$

Repeat K times

- Server model is an average of client models

$$x = \frac{1}{S} \sum_{i \in \mathcal{S}} y_i$$

Average over all sampled clients

+ Potentially faster
(performs K updates)

- different from centralized updates
- may not converge

Solving FL: SCAFFOLD

New!

- on each client i in \mathcal{S} , perform K steps of SGD

$$y_i = y_i - \eta(g_i(y_i)) + \boxed{c - c_i}$$

Correction term!

+ Potentially faster
(performs K updates)

+ mimics centralized
updates!

Repeat K times

- Average as before

$$x = \frac{1}{S} \sum_{i \in \mathcal{S}} y_i$$

When does
FedAvg fail?

FedAvg degrades with heterogeneous clients

Measuring the Effects of Non-Identical Data Distribution for Federated Visual Classification

Tzu-Ming Harry Hsu*
MIT CSAIL
stmharry@mit.edu

Hang Qi
Google Research
hangqi@google.com

Matthew Brown
Google Research
mtbr@google.com

Federated Optimization in Heterogeneous Networks

Tian Li
CMU
tianli@cmu.edu

Anit Kumar Sahu
Bosch Center for AI
anit.sahu@gmail.com

Manzil Zaheer
Google Research
manzilz@google.com

Maziar Sanjabi
USC
maziar.sanjabi@gmail.com

Ameet Talwalkar
CMU & Determined AI
talwalkar@cmu.edu

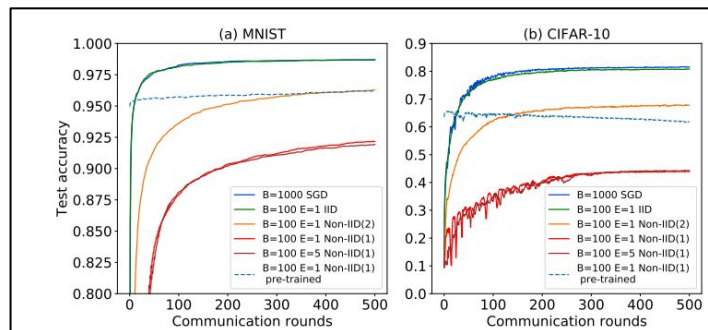
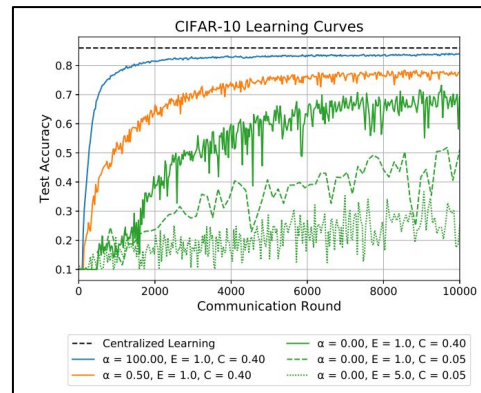
Virginia Smith
CMU
smithv@cmu.edu

Federated Learning with Non-IID Data

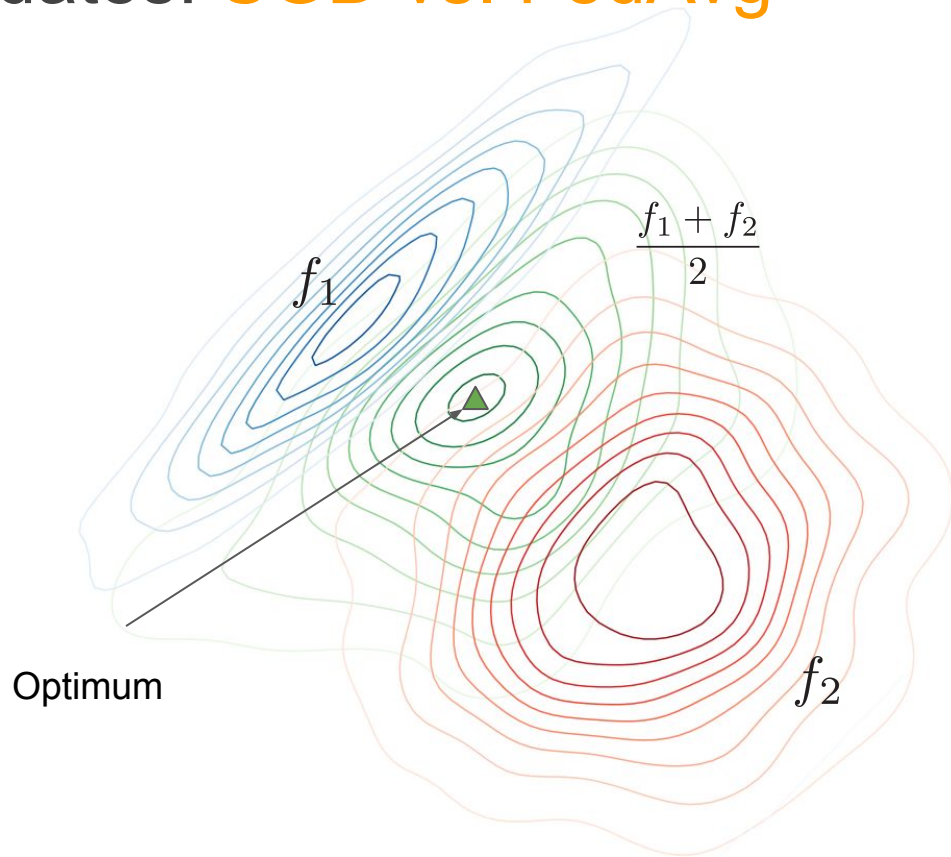
Yue Zhao*, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, Vikas Chandra

Arm, San Jose, CA

{yue.zhao, meng.li, liangzhen.lai, naveen.suda, damon.civin, vikas.chandra}@arm.com



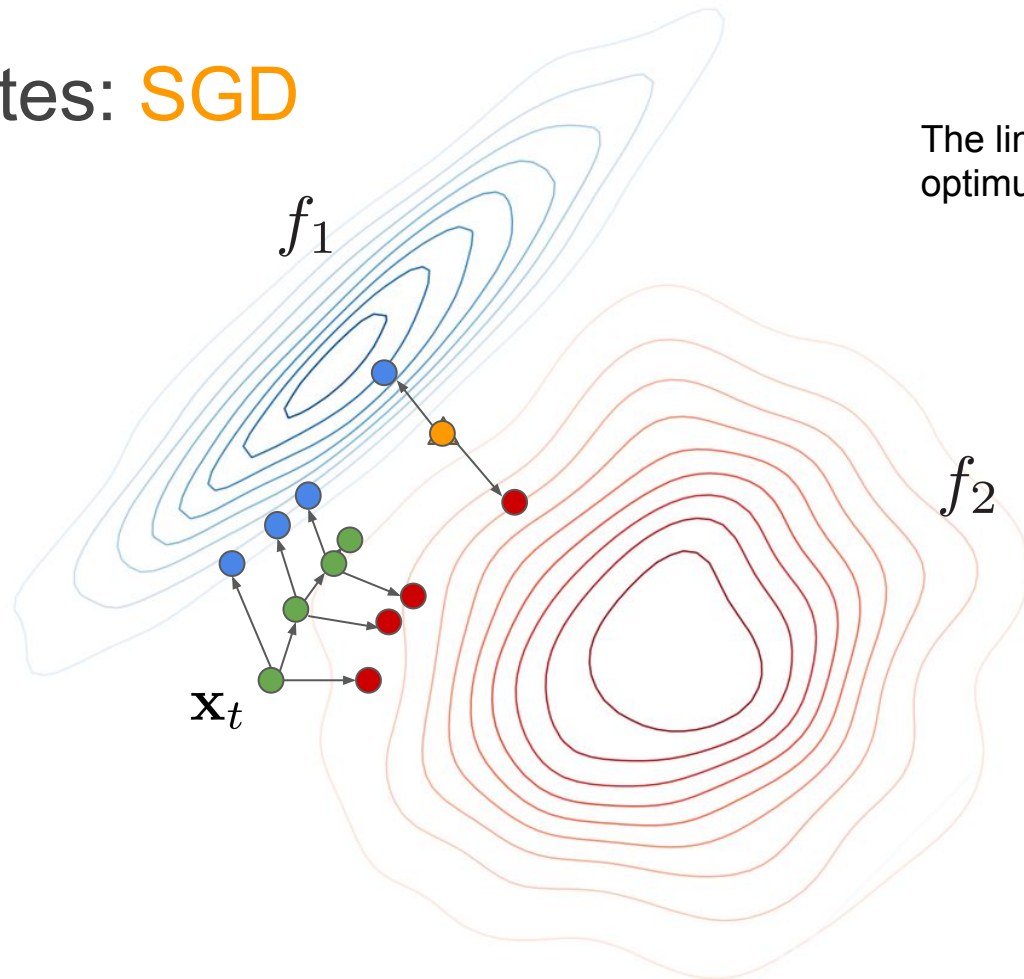
Client updates: SGD vs. FedAvg



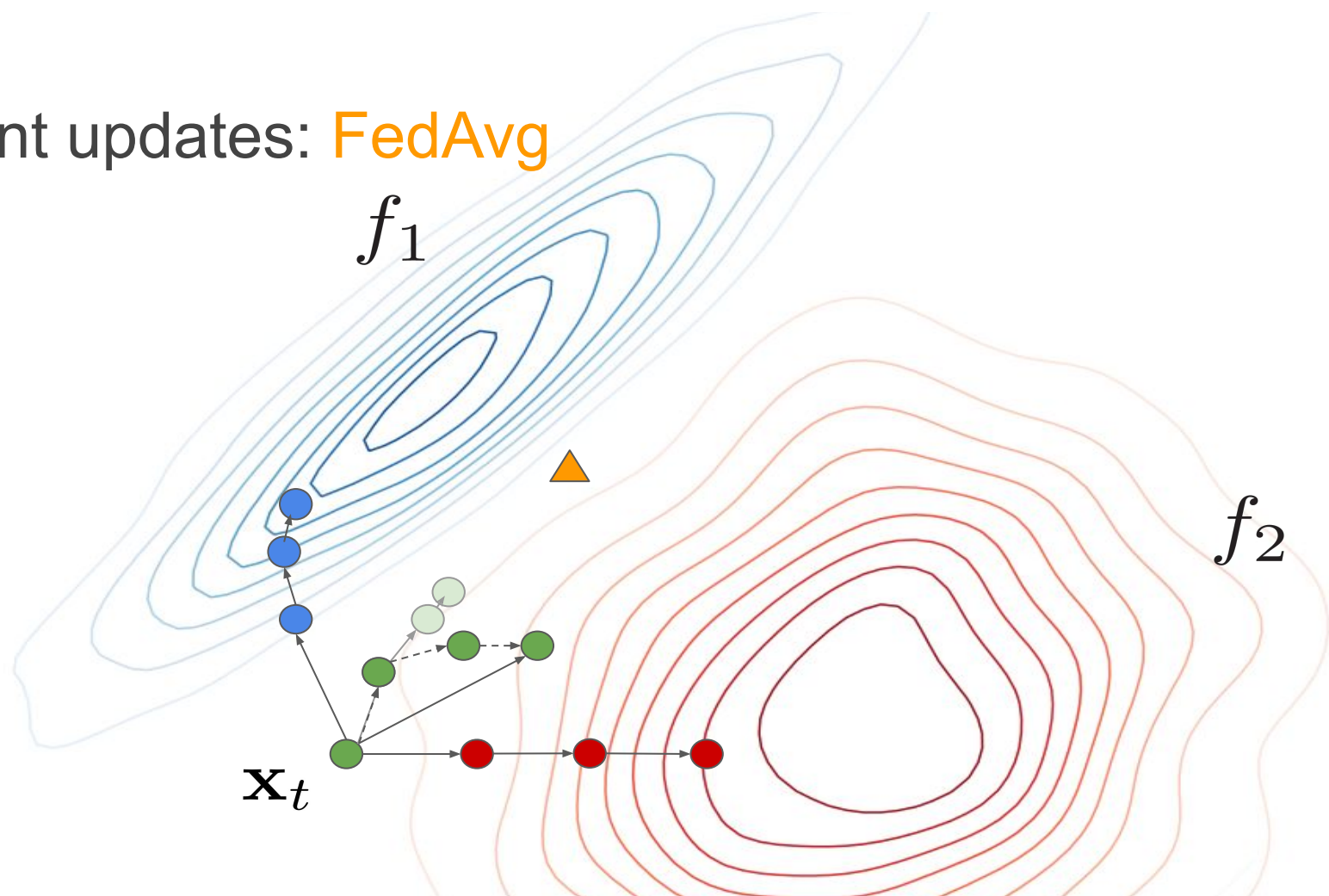
Surface of two client loss functions, and the combined function

Client updates: SGD

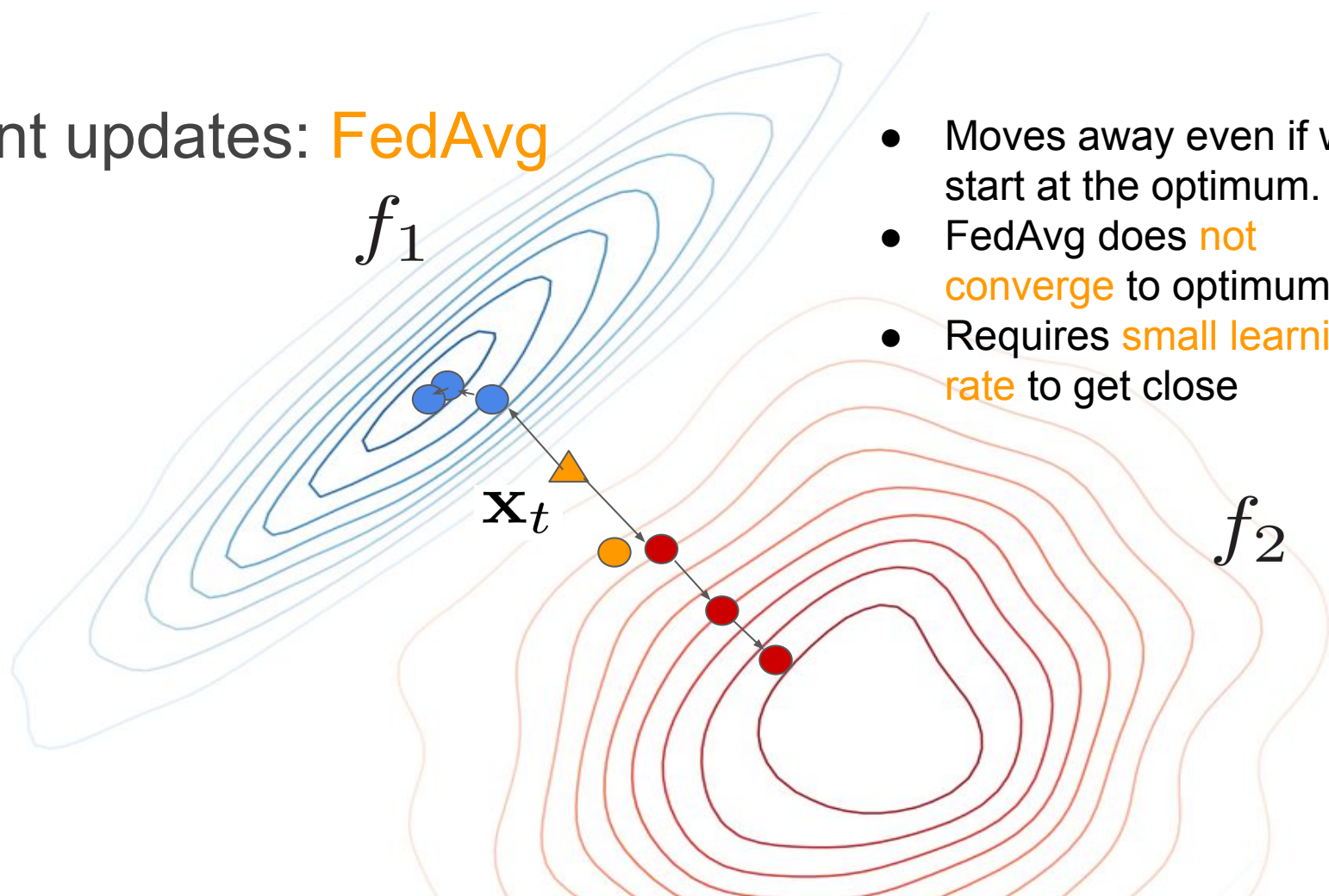
The limit point of SGD is the optimum.



Client updates: FedAvg

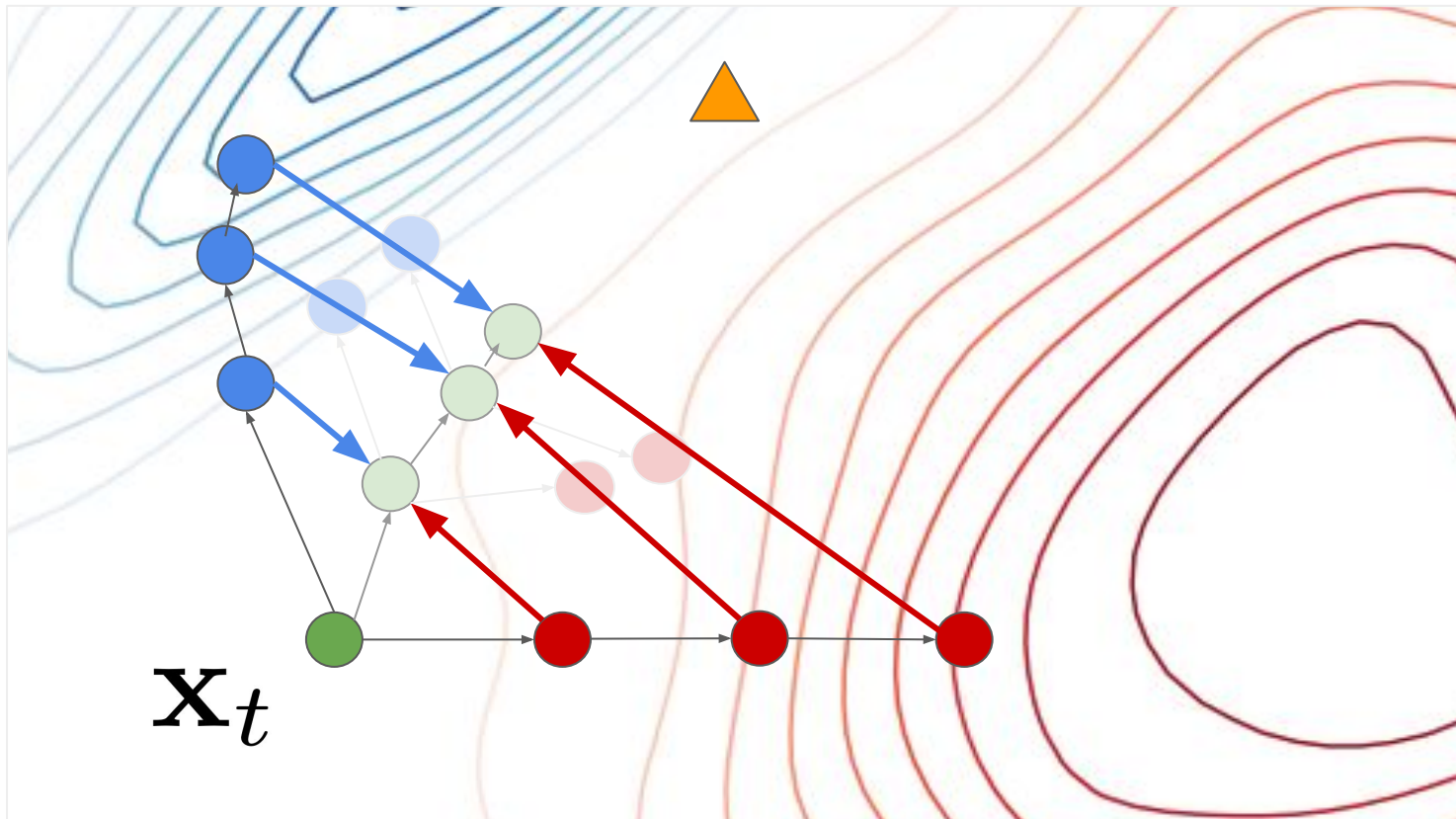


Client updates: FedAvg



- Moves away even if we start at the optimum.
- FedAvg does **not converge** to optimum!
- Requires **small learning rate** to get close

Drift in client updates: SGD vs. FedAvg



Convergence Rates: SGD

- For strongly convex

$$\frac{\sigma^2}{\mu RK N} \mathcal{N} + \exp\left(\frac{-\mu RK}{L}\right)$$

- For non-convex functions

$$\frac{\sigma}{\sqrt{RK N}} \mathcal{N} + \frac{L}{R}$$

Notation:

- R communication rounds
- K local steps
- Total N clients

- L - smooth, μ - strongly convex
- σ - variance within a client

Convergence Rates: SGD vs. FedAvg

Generalizes [Li et al. 2019] and [Khaled et al. 2019]

Assume: (B,G)-similar gradients

$$\mathbb{E}_i \|\nabla f_i(x)\|^2 \leq G^2 + B^2 \|\nabla f(x)\|^2$$

- For strongly convex

SGD

$$\mathcal{N} + \exp\left(\frac{-\mu R}{L}\right)$$

FedAvg

$$\mathcal{N} + \frac{LG^2}{\mu^2 R^2} + \exp\left(\frac{-\mu R}{L}\right)$$

Convergence Rates: SGD vs. FedAvg

SGD

- For strongly convex

$$\mathcal{N} + \exp\left(\frac{-\mu R}{L}\right)$$

- For non-convex functions

$$\mathcal{N} + \frac{L}{R}$$

FedAvg

Assume: (B,G)-similar gradients

$$\mathcal{N} + \frac{LG^2}{\mu^2 R^2} + \exp\left(\frac{-\mu R}{L}\right)$$

Tightest rates, uses server and client step-sizes

$$\mathcal{N} + \frac{L^{\frac{1}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}} + \frac{L}{R}$$

Lower bound: FedAvg

Theorem: For any G , we can find functions with $(2, G)$ -similar gradients such that FedAvg for $K > 1$ with arbitrary step-sizes always has error

$$\geq \frac{G^2}{\mu R^2}$$

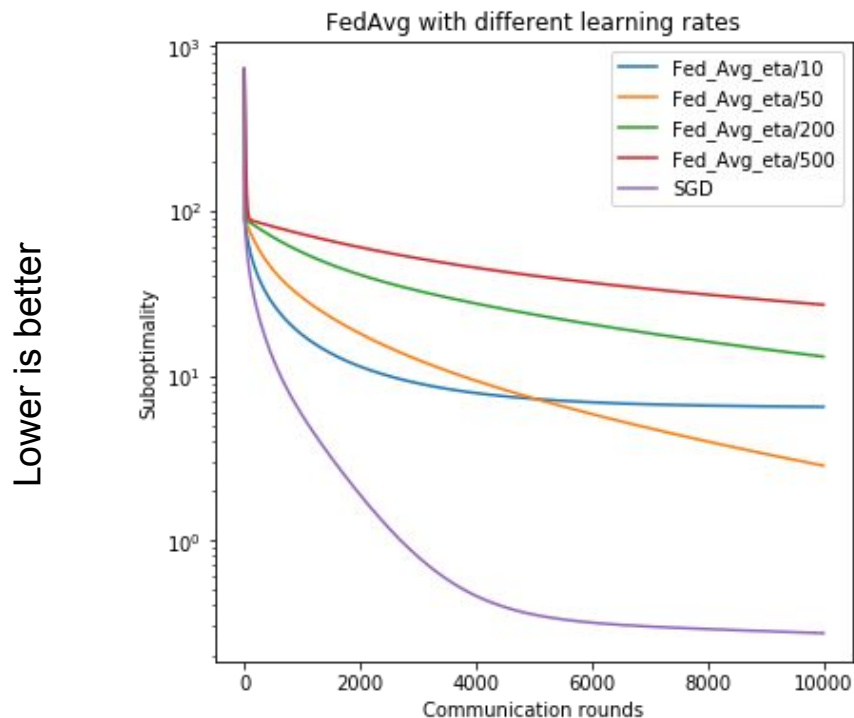
Assume: (B, G) -similar gradients

$$\mathcal{N} + \frac{LG^2}{\mu^2 R^2} + \exp\left(\frac{-\mu R}{L}\right)$$

↑
Necessary!

$$\mathcal{N} + \frac{L^{\frac{1}{3}} G^{\frac{2}{3}}}{R^{\frac{2}{3}}} + \frac{L}{R}$$

Quick demo: SGD vs. FedAvg



- Linear regression
- concrete dataset (UCI)
- 10 clients (no sampling)
- $K = 10$ local steps

> FedAvg needs **smaller learning rate**

> **Slower** than SGD

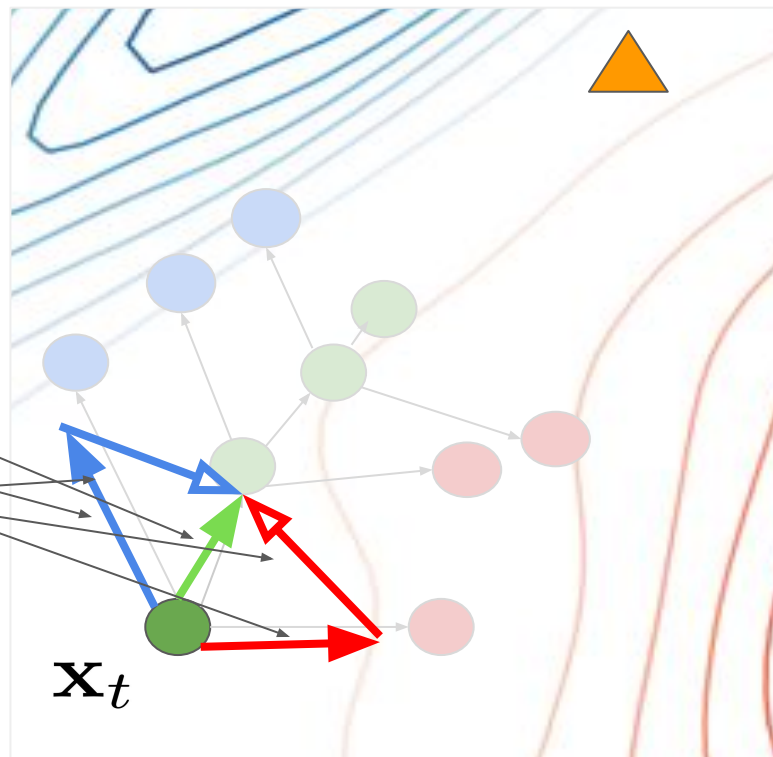
SCAFFOLD:

stochastic *controlled*

averaging

Main Idea: Use control variates

- Guess direction of server update \mathcal{C}
- Guess direction of client update \mathcal{C}_i
- Use the correction $(\mathcal{C} - \mathcal{C}_i)$

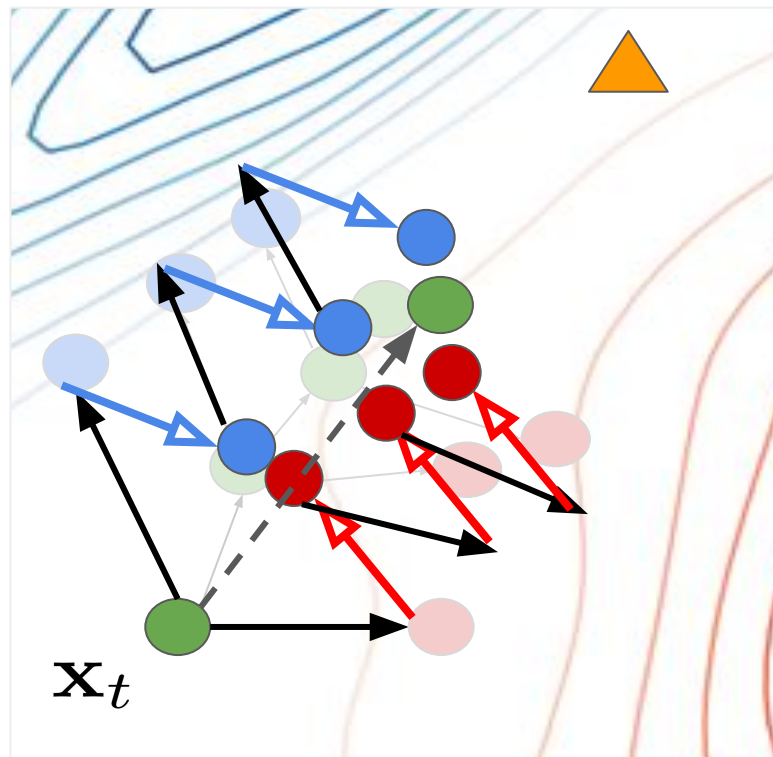


Main Idea: Corrected updates

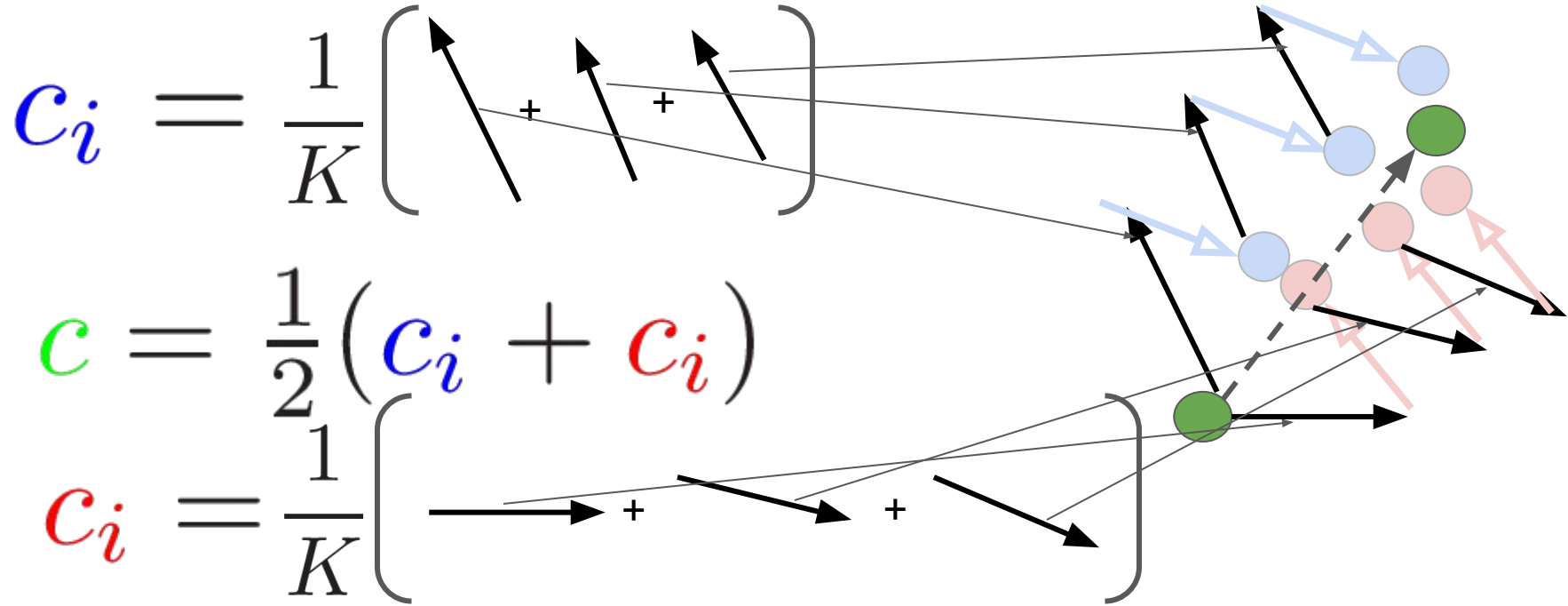
$$y_i = y_i - \eta(g_i(y_i) + \boxed{c - c_i})$$

Correction
terms

Mimics centralized updates!



Main Idea: Updating control variates

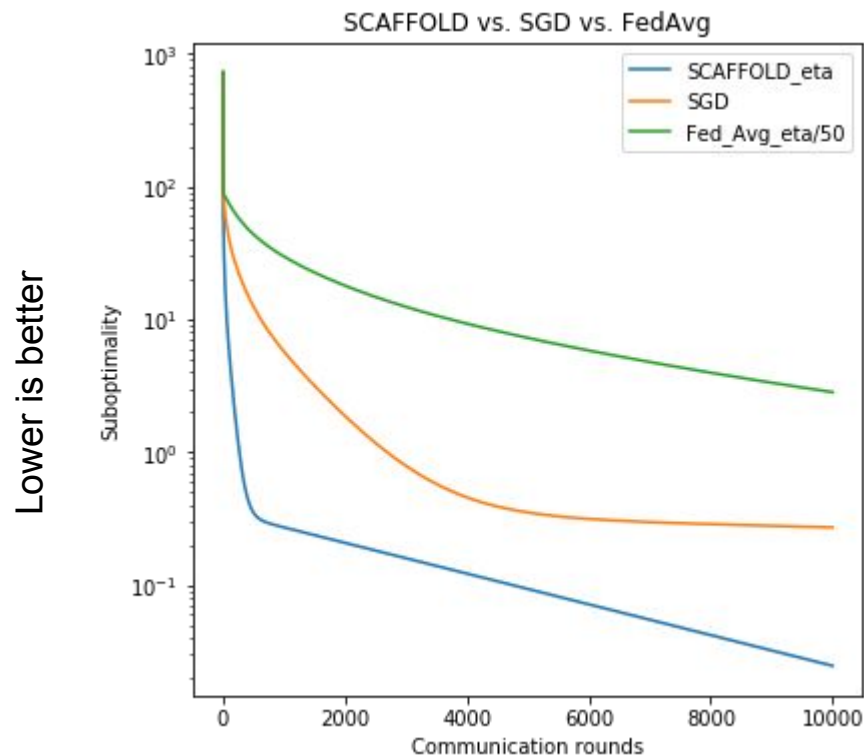


SCAFFOLD: Algorithm

Algorithm 1 SCAFFOLD: Stochastic Controlled Averaging for federated learning

```
1: server input initial parameters  $\mathbf{x}$ , control variate  $\mathbf{c}$ , and global step-size  $\eta_g$ 
2: client input for each  $i$  local control variate  $\mathbf{c}_i$ , and local step-size  $\eta_l$ 
3: for each communication round  $r = 1, \dots, R$  do
4:   select a subset of clients  $\mathcal{S} \subseteq \{1, \dots, N\}$ 
5:   communicate  $(\mathbf{x}, \mathbf{c})$  to all clients  $i \in \mathcal{S}$ 
6:   on each client  $i \in \mathcal{S}$  do
7:     initialize local parameters  $\mathbf{y}_i \leftarrow \mathbf{x}$ 
8:     for each local step  $k = 1, \dots, K$  do
9:       compute a stochastic gradient  $g_i(\mathbf{y}_i)$  of  $f_i$ 
10:       $\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta_l (g_i(\mathbf{y}_i) - \mathbf{c}_i + \mathbf{c})$   $\triangleright$  local updates with correction
11:    end for
12:     $\mathbf{c}_i^+ \leftarrow$  (i)  $g_i(\mathbf{x})$ , or (ii)  $\mathbf{c}_i - \mathbf{c} + \frac{1}{K\eta_l}(\mathbf{x} - \mathbf{y}_i)$   $\triangleright$  compute new control variate
13:    communicate  $(\Delta \mathbf{y}_i, \Delta \mathbf{c}_i) \leftarrow (\mathbf{y}_i - \mathbf{x}, \mathbf{c}_i^+ - \mathbf{c}_i)$ 
14:     $\mathbf{c}_i \leftarrow \mathbf{c}_i^+$   $\triangleright$  update control variate
15:  end client
16:   $\Delta \mathbf{x} \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \Delta \mathbf{y}_i$  and  $\Delta \mathbf{c} \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \Delta \mathbf{c}_i$   $\triangleright$  aggregate client outputs
17:   $\mathbf{x} \leftarrow \mathbf{x} + \eta_g \Delta \mathbf{x}$  and  $\mathbf{c} \leftarrow \mathbf{c} + \frac{|\mathcal{S}|}{N} \Delta \mathbf{c}$   $\triangleright$  update parameters and control
18: end for
```

SCAFFOLD: Quick demo



- Linear regression
- concrete dataset (UCI)
- 10 clients (no sampling)
- $K = 10$ local steps

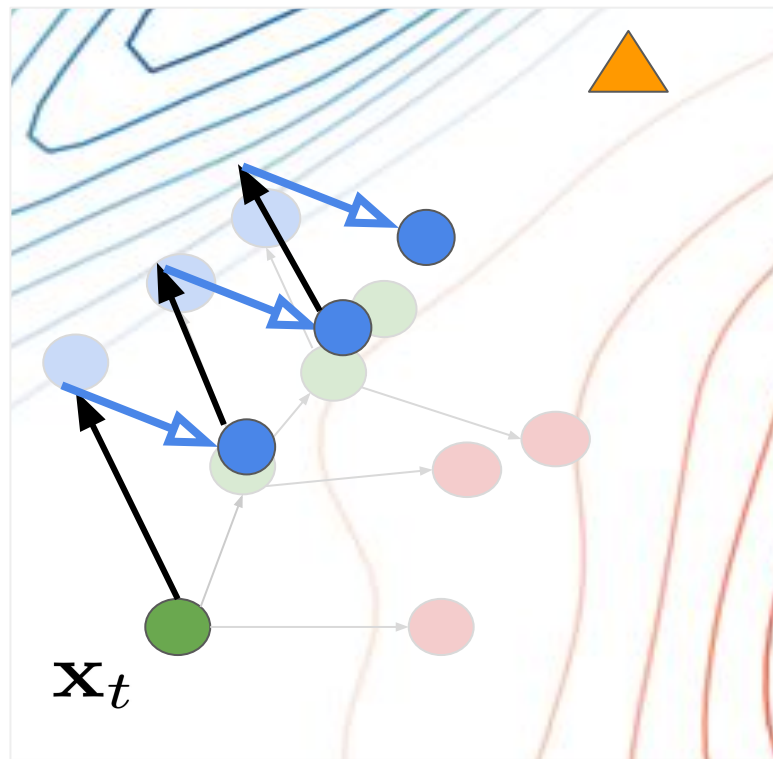
> SCAFFOLD works with **same learning rate** as SGD

> **Faster** than SGD!

SCAFFOLD: Client sampling

- Updates of every client mimics centralized updates.
- Few #clients works, as long as control variates are accurate.
- Hence, very robust to client sampling.

Different view: **SAGA is a special case** of SCAFFOLD with client sampling



SCAFFOLD: Variance reduced convergence Rates

- For strongly convex functions

$$\mathcal{N} + \exp \left(- \min \left\{ \frac{\mu}{L}, \frac{S}{N} \right\} R \right)$$

- For non-convex functions

$$\mathcal{N} + \left(\frac{N}{S} \right)^{\frac{2}{3}} \frac{L}{R}$$

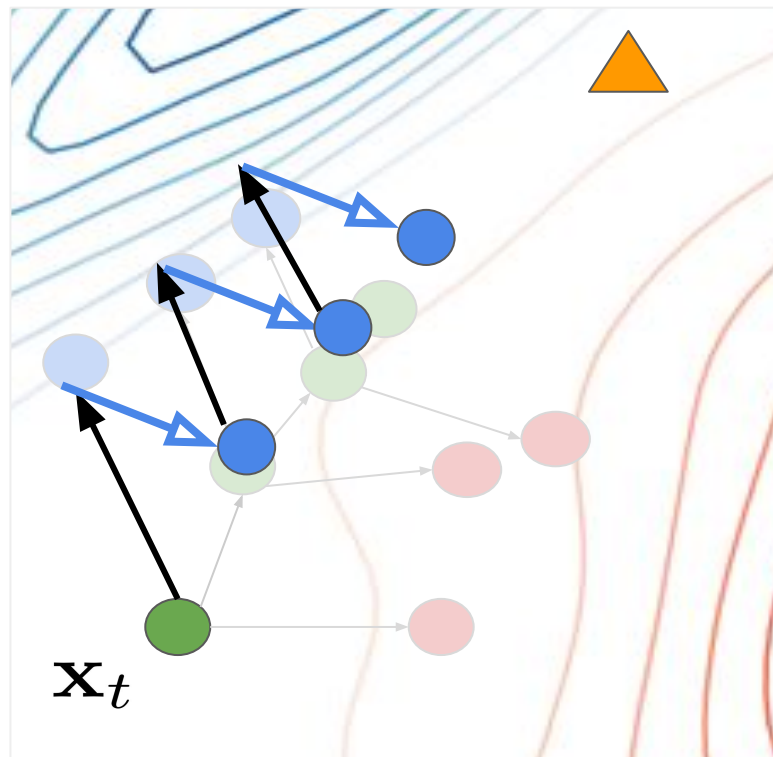
Notation:

- R communication rounds
- S out of N clients sampled
- L - smooth, μ - strongly convex

> Better than FedAvg!

SCAFFOLD: Why take more than 1 step?

- **Each update** mimics a centralized update => local steps should help.
- In worst case not true [Arjevani & Shamir, 2015] :(
- Possible if similar Hessians!



SCAFFOLD: Why take more than 1 step?

δ - BHD (Bounded Hessian Dissimilarity)

$$\|\nabla^2 f_i(x) - \nabla^2 f(x)\| \leq \delta$$

And is δ -weakly convex.

Note that

$$\delta \leq L, \text{ and typically } \delta \ll L$$

SCAFFOLD: Why take more than 1 step?

Assume: δ - BHD, quadratics

- For strongly convex functions

$$\mathcal{N} + \exp\left(\frac{-\mu RK}{L + (\mu + \delta)K}\right) \approx \mathcal{N} + \exp\left(\frac{-\mu R}{\mu + \delta}\right)$$

- For non-convex functions

$$\mathcal{N} + \frac{L + \delta K}{RK} \approx \mathcal{N} + \frac{\delta}{R}$$

Notation:

- R communication rounds
- All N clients participate
- K local steps

- L - smooth, μ - strongly convex

SCAFFOLD: Why take more than 1 step?

Assume: δ - BHD, quadratics

> Best to take $K \approx \frac{L}{\delta}$

- For strongly convex functions

$$\mathcal{N} + \exp\left(\frac{-\mu RK}{L + (\mu + \delta)K}\right) \approx \mathcal{N} + \exp\left(\frac{-\mu R}{\mu + \delta}\right)$$

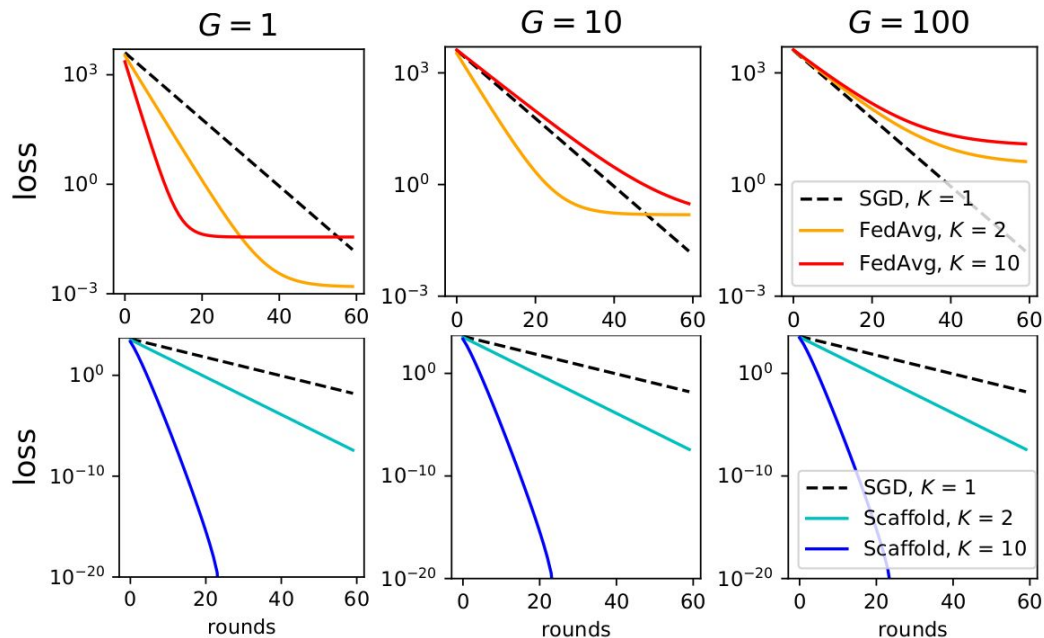
> We replaced L with δ in the rates (typically $\delta \ll L$)

- For non-convex functions

$$\mathcal{N} + \frac{L + \delta K}{RK} \approx \mathcal{N} + \frac{\delta}{R}$$

> First rate to characterize improvement due to local steps!

SCAFFOLD: Why take more than 1 step?



Quick demo on scalar quadratics

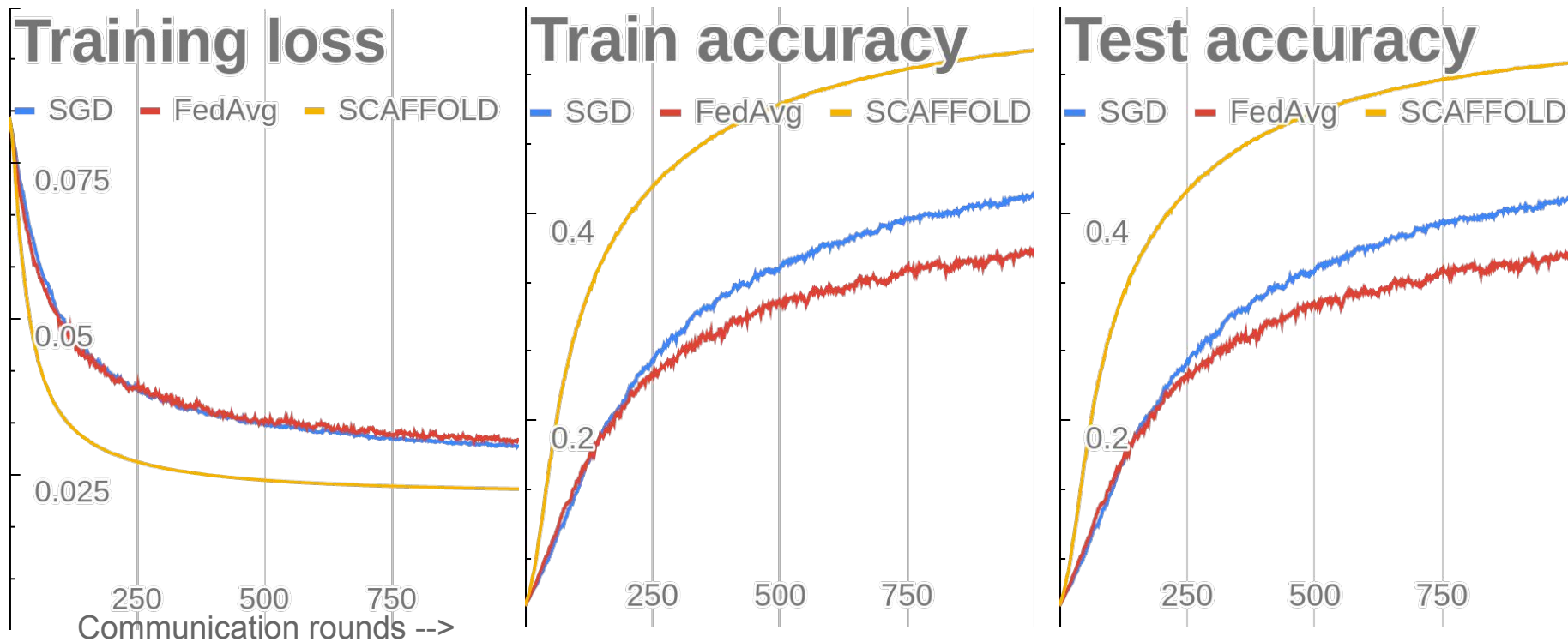
- Scaffold is unaffected by G
- Larger K is better
- $K=2$ is 2 times faster
- $K=10$ is only 4 times better

Experiments

Experimental Setup

- Extended MNIST (balanced) dataset
- Multi-class logistic regression (47 classes)
- Partitioned into N clients
- Sorted by labels and then 'slightly shuffled' before splitting

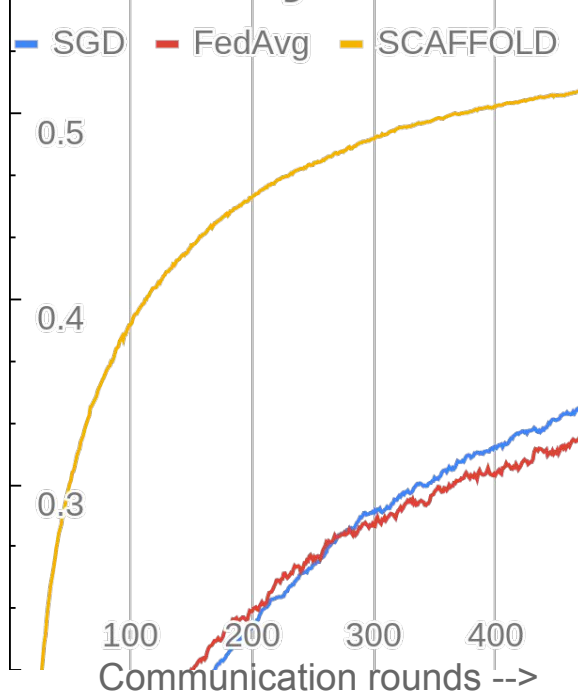
Performance of SCAFFOLD



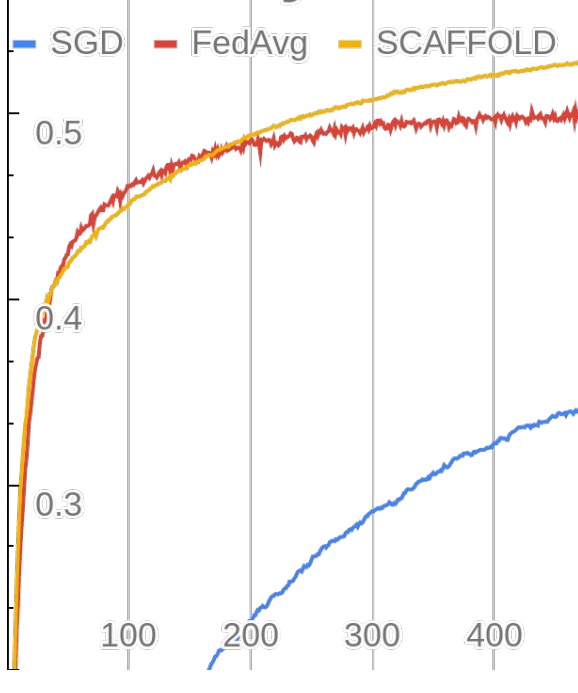
Similarity = 0, 1 Epoch, #sampled clients = 20, total clients = 400

Effect of similarity

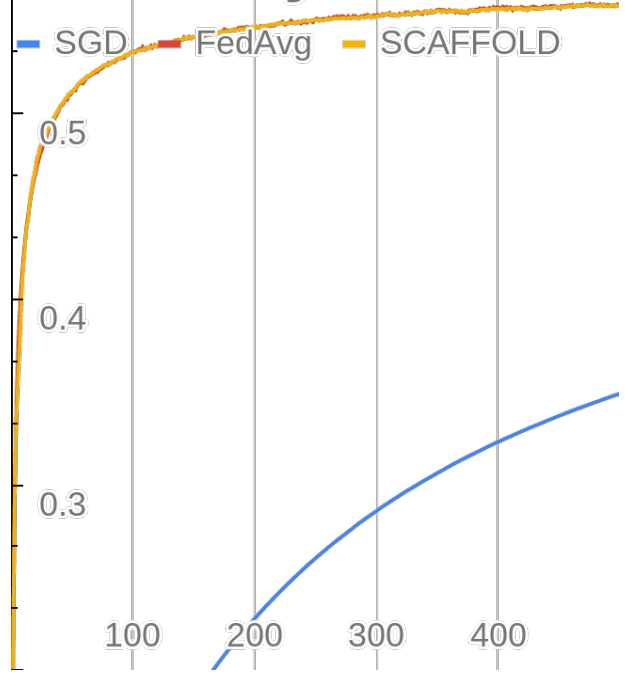
Similarity = 0%



Similarity = 10%

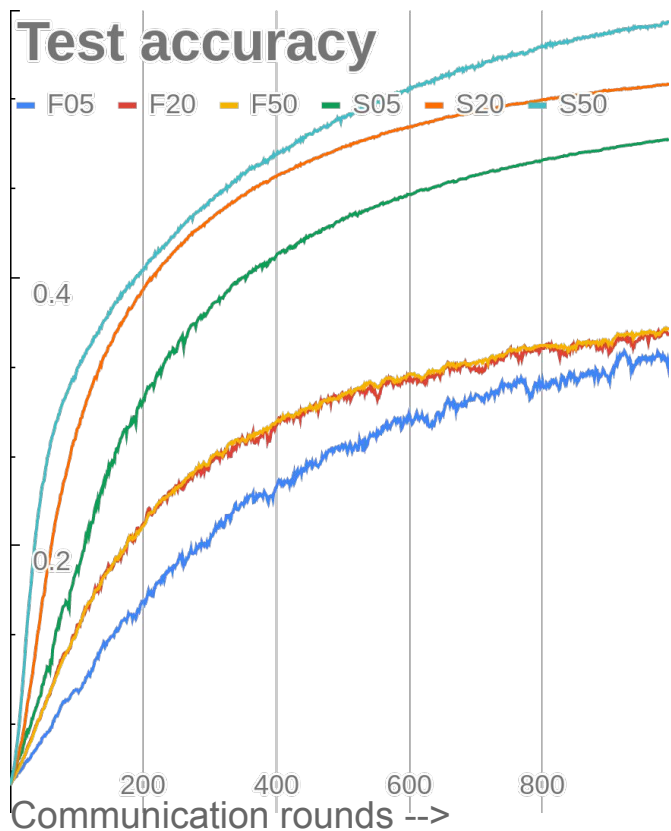


Similarity = 100%



Test accuracy, 10 Epochs, #sampled clients = 20, total clients = 100

Effect of number of clients



SCAFFOLD with 5 clients is better than FedAvg with 50!

- > Total #clients = 400
- > Total #categories = 47
- > 1 Epoch per round
- > Similarity = 0

Take aways

- Degradation of FedAvg is due to the **client drift**. If you use FedAvg, use separate server and client step-sizes.
- Why you should use **SCAFFOLD**:
 - Provably converges faster than SGD and FedAvg
 - Resilient to heterogeneity and client sampling
- **Main limitation:** requires maintaining client state, so applicable only to cross-silo FL

Cross-Silo vs Cross-Device Federated Learning

Cross-Silo FL

- **Small/medium** number (2-100, typically) of total clients
- E.g. hospitals, financial organizations
- **Large amount** of data per client
- **Persistent** clients: almost always available
- **Stateful** clients: clients can carry state from round to round

Cross-Device FL

- **Very large** number (e.g. 10^{10}) of total clients
- E.g. mobile or IoT devices
- **Small amount** of data per client
- **Transient** clients: only a fraction of clients available at any time
- **Stateless** clients: clients generally participate only once in each task

Cross-device Federated Learning: Formalism

$$\min_x \mathbb{E}_{i \sim \mathcal{D}} \left[f_i(x) := \frac{1}{n_i} \sum_{\nu=1}^{n_i} f_i(x; \xi_\nu) \right]$$

Model parameters

Expectation over (possibly infinitely many) clients

Sum over client data

Client loss function over the parameters and data

Algorithms for Cross-Device Federated Learning

Solving FL: SGD with (server) momentum

[Assume only 1 client per round]

$$x = x - \eta * (\nabla f_i(x) + \beta * m)$$

Update server
parameters



$$m = \nabla f_i(x) + \beta m$$

Update server
momentum



+ Convergence guaranteed

- Communicates every
update round

Solving FL: SGD with server momentum

[Assume only 1 client per round]

$$y_i = y_i - \eta \nabla f_i(y_i)$$

+ Convergence guaranteed

$$x \leftarrow x - \eta(\nabla f_i(y_i) + \beta m)$$

Update server parameters

$$m = \nabla f_i(x) + \beta m$$

Update server momentum

- Communicates every update round

Solving FL: FedAvg with momentum

[McMahan et al. 2016,
Hsu et al. 2019,
Reddi et al. 2020]

- Starting from x , run K local updates

$$y_i = y_i - \eta \nabla f_i(y_i)$$

Repeat K times

- Use $(x - y_i)$ as a pseudo-gradient.

$$x = x - \eta * (\{x - y_i\} + \beta * m)$$

Update server parameters

+ Communicates only every K updates

- bad convergence due to client drift (though momentum helps!)

Solving FL: Mime with momentum

- Apply server momentum locally in the clients

$$y_i = y_i - \eta * (\nabla f_i(y_i) + \beta * m)$$

Repeat K times

Fixed server momentum

- Momentum is computed globally (at server) and applied locally (at clients)

$$m = \nabla f_i(x) + \beta m$$

Update server momentum

- + Communicates only every K updates
- + Reduce client drift using (fixed) server momentum!
- + Extends to Adam, etc.

Solving FL: Mime with momentum

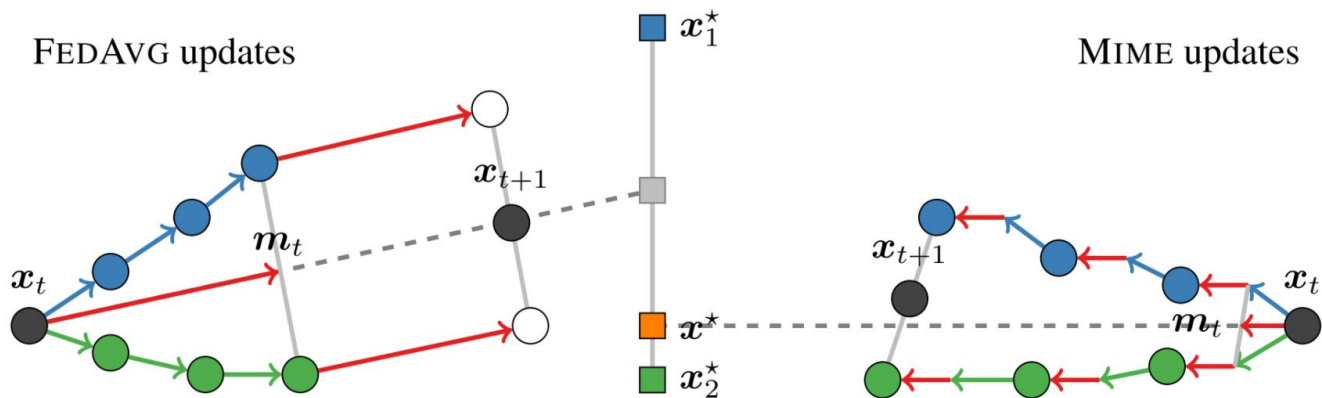


Figure 1: Client-drift in FEDAVG (left) and MIME (right) is illustrated for 2 clients with 3 local steps and momentum parameter $\beta = 0.5$. The local SGD updates of FEDAVG (shown using arrows for client 1 and client 2) move towards the average of client optima $\frac{x_1^* + x_2^*}{2}$ which can be quite different from the true global optimum x^* . Server momentum only speeds up the convergence to the wrong point in this case. In contrast, MIME uses unbiased momentum and applies it locally at every update. This keeps the updates of MIME closer to the true optimum x^* .

Mime framework: adapting optimizers to FL setting

Base optimizer updates: given gradient \mathbf{g} and current internal state \mathbf{s} (e.g. momentum, Adagrad/Adam accumulators, etc.)

$$\begin{aligned}\mathbf{x} &\leftarrow \mathbf{x} - \eta \mathcal{U}(\mathbf{g}, \mathbf{s}), \\ \mathbf{s} &\leftarrow \mathcal{V}(\mathbf{g}, \mathbf{s}).\end{aligned}$$

E.g. SGD with momentum:

$$\begin{aligned}\mathbf{x} &\leftarrow \mathbf{x} - \eta(\mathbf{g} + \beta \mathbf{m}) \\ \mathbf{m} &\leftarrow \mathbf{g} + \beta \mathbf{m}\end{aligned}$$

Mime server update: update state \mathbf{s} using gradients from clients:

$$\mathbf{s} \leftarrow \mathcal{V}\left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}), \mathbf{s}\right)$$

Mimelite client update: obtain \mathbf{x} and state \mathbf{s} from server and repeat K times starting from $\mathbf{y} = \mathbf{x}$:

$$\mathbf{y} \leftarrow \mathbf{y} - \eta \mathcal{U}\left(\nabla f_i(\mathbf{y}, \xi), \mathbf{s}\right)$$

Mime framework: adapting optimizers to FL setting

Base optimizer updates: given gradient \mathbf{g} and current internal state \mathbf{s} (e.g. momentum, Adagrad/Adam accumulators, etc.)

$$\begin{aligned}\mathbf{x} &\leftarrow \mathbf{x} - \eta \mathcal{U}(\mathbf{g}, \mathbf{s}), \\ \mathbf{s} &\leftarrow \mathcal{V}(\mathbf{g}, \mathbf{s}).\end{aligned}$$

Mime server update: update state \mathbf{s} using gradients from clients:

$$\mathbf{s} \leftarrow \mathcal{V}\left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}), \mathbf{s}\right)$$

Mime client update: obtain \mathbf{x} and state \mathbf{s} from server and repeat K times starting from $\mathbf{y} = \mathbf{x}$:

$$\mathbf{y} \leftarrow \mathbf{y} - \eta \mathcal{U}\left(\nabla f_i(\mathbf{y}, \xi) - \nabla f_i(\mathbf{x}, \xi) + \frac{1}{S} \sum_{j \in \mathcal{S}} \nabla f_j(\mathbf{x}), \mathbf{s}\right)$$

Algorithm 1 **Mime** and **MimeLite**

input: initial \mathbf{x} and \mathbf{s} , learning rate η and base algorithm $\mathcal{B} = (\mathcal{U}, \mathcal{V})$

for each round $t = 1, \dots, T$ **do**

sample subset \mathcal{S} of clients

communicate (\mathbf{x}, \mathbf{s}) to all clients $i \in \mathcal{S}$

communicate $\mathbf{c} \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x})$ (only for Mime)

on client $i \in \mathcal{S}$ **in parallel do**

initialize local model $\mathbf{y}_i \leftarrow \mathbf{x}$

for $k = 1, \dots, K$ **do**

sample mini-batch ζ from local data

$\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta \mathcal{U}(\nabla f_i(\mathbf{y}_i; \zeta) - \nabla f_i(\mathbf{x}; \zeta) + \mathbf{c}, \mathbf{s})$ (**Mime**)

$\mathbf{y}_i \leftarrow \mathbf{y}_i - \eta \mathcal{U}(\nabla f_i(\mathbf{y}_i; \zeta), \mathbf{s})$ (**MimeLite**)

end for

compute full local-batch gradient $\nabla f_i(\mathbf{x})$

communicate $(\mathbf{y}_i, \nabla f_i(\mathbf{x}))$

end on client

$\mathbf{x} \leftarrow \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \mathbf{y}_i$, and $\mathbf{s} \leftarrow \mathcal{V}\left(\frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla f_i(\mathbf{x}), \mathbf{s}\right)$

end for

Analysis

The diagram shows the objective function $\min_x \mathbb{E}_{i \sim \mathcal{D}} \left[f_i(x) := \frac{1}{n_i} \sum_{\nu=1}^{n_i} f_i(x; \xi_\nu) \right]$ with four colored arrows pointing to its components: a red arrow to x labeled 'parameters', a green arrow to $\mathbb{E}_{i \sim \mathcal{D}}$ labeled 'clients', an orange arrow to $f_i(x)$ labeled 'loss function', and a blue arrow to ξ_ν labeled 'Client data'.

- \mathbf{G}^2 - Bounded Gradient dissimilarity:

$$\mathbb{E}_{i \sim \mathcal{D}} \|\nabla f_i(x) - \nabla f(x)\|^2 \leq G^2$$

- \square - Bounded Hessian dissimilarity:

$$\|\nabla^2 f_i(x; \xi_i) - \nabla^2 f(x)\| \leq \delta$$

Convergence rates

Algorithm	Non-convex	μ -Strongly convex	Assumptions
SERVER-ONLY			
SGD [8]	$\frac{G^2}{S\epsilon^2} + \frac{L}{\epsilon}$	$\frac{G^2}{\mu S\epsilon} + \frac{L}{\mu}$	G^2 -BGD
MVR [4]	$\left(\frac{G}{\sqrt{S\epsilon}}\right)^{\frac{3}{2}} + \frac{L}{\epsilon}$	—	G^2 -BGD
FEDAVG ¹			
FedSGD [15]	$\frac{G^2}{S\epsilon^2} + \frac{G}{\epsilon^{3/2}} + \frac{L}{\epsilon}$	$\frac{G^2}{\mu S\epsilon} + \frac{G}{\mu\sqrt{\epsilon}} + \frac{L}{\mu}$	G^2 -BGD
MIME ²			
MimeSGD	$\frac{G^2}{S\epsilon^2} + \frac{\delta}{\epsilon}$	$\frac{G^2}{\mu S\epsilon} + \frac{\delta}{\mu}$	G^2 -BGD, δ -BHD
MimeMVR	$\left(\frac{G}{\sqrt{S\epsilon}}\right)^{\frac{3}{2}} + \frac{\delta}{\epsilon}$	—	G^2 -BGD, δ -BHD
Lower bound [1]	$\Omega\left(\frac{G}{\sqrt{S\epsilon}}\right)^{\frac{3}{2}}$	$\Omega\left(\frac{G^2}{S\epsilon}\right)$	G^2 -BGD

Experiments

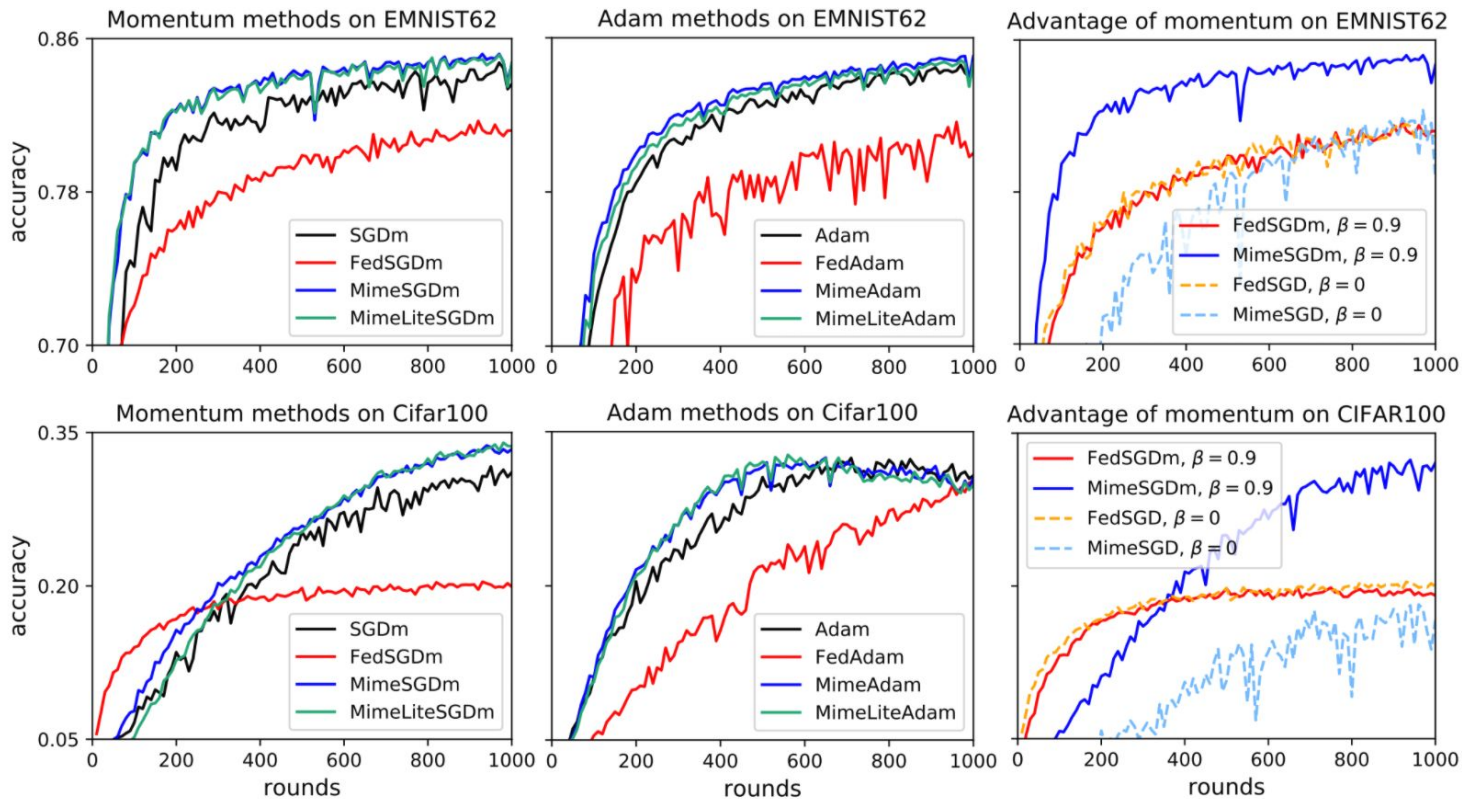


Figure 3: **Server-only**, **FedAvg**, **Mime**, and **MimeLite** with SGDm (left) and Adam (middle) run on EMNIST62 (top) and CIFAR100 (bottom). Mime and MimeLite have very similar performance and are consistently the best. FedAvg is often even worse than the server-only baselines. Also, Mime makes better use of momentum than FedAvg, with a large increase in performance (right).

Takeaways

- Momentum injects global information and helps reduce client drift.
- Compute momentum globally at server, apply it during each client update.
- Usefulness of local steps depends on Hessian variance.

Thank You.

Questions?