



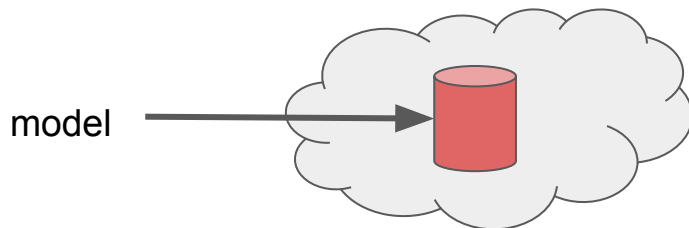
# Breaking the centralized barrier in cross-device federated learning

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

How to adapt normal  
deep learning  
algorithms to  
federated learning  
in a principled way?

# Federated Learning: Setting

[McMahan et al. 2016]

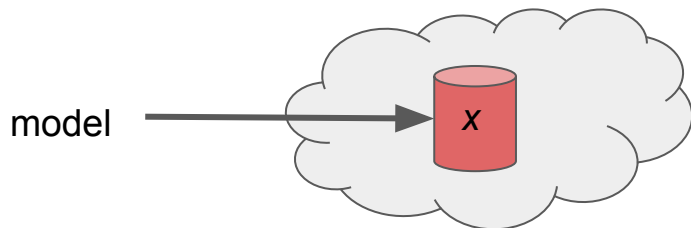


Server  
(e.g. Google)



Clients  
(e.g. mobile phones)

# Cross-device federated learning



- Huge (possibly infinite) number of clients => each client is seen *at most* once.
- High overhead per round
- Each client has a small amount of *heterogeneous* data



# Cross silo vs. Cross device

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

Diagram illustrating the cross-silo optimization problem. The equation shows the minimization of the average loss over  $N$  clients. A green arrow labeled "clients" points to the summation index  $i$ . A blue arrow labeled "Client data" points to the expectation operator  $E_{\zeta}$ .

- Cross-silo is closer to “finite-sum” optimization
- Can use variance reduction (SCAFFOLD, etc.)

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

Diagram illustrating the cross-device optimization problem. The equation shows the minimization of the expected loss over a distribution of clients. A green arrow labeled "clients" points to the expectation operator  $E_{i \sim \mathcal{D}}$ . A blue arrow labeled "Client data" points to the inner summation over  $\nu$ .

- **Cross-device** is closer to “stochastic” optimization. Essentially,  $N$  is  $\infty$ .
- Use algorithms like SGD, momentum, Adam etc.

# Cross device FL: In this talk

$$\min_x \mathbb{E}_{i \sim \mathcal{D}} [f_i(x)]$$

Model  
parameters

Expectation over  
(possibly infinite)  
clients

Full local gradients

And we only sample 1 client per round.

Also, we focus on momentum.

See the paper for full details.

# Part 1: Algorithms

## Mime Framework

# Solving FL: Server only momentum

[Assume only 1 client per round]

$$x = x - \eta \left( (1 - \beta) \nabla f_i(x) + \beta m \right)$$

  
Update server  
parameters

$$m = (1 - \beta) \nabla f_i(x) + \beta m$$

  
Update server  
momentum

+ Convergence guaranteed

- Communicates every  
update



# 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.

$$x = x - ((1 - \beta)(x - y_i) + \beta m)$$

Update server parameters

$$m = (1 - \beta)(x - y_i) + \beta m$$

+ Communicates only every  $K$  updates

- bad convergence due to client drift because each client overfits to itself

[SCAFFOLD, Karimireddy et al. 2019]

# But momentum helps!

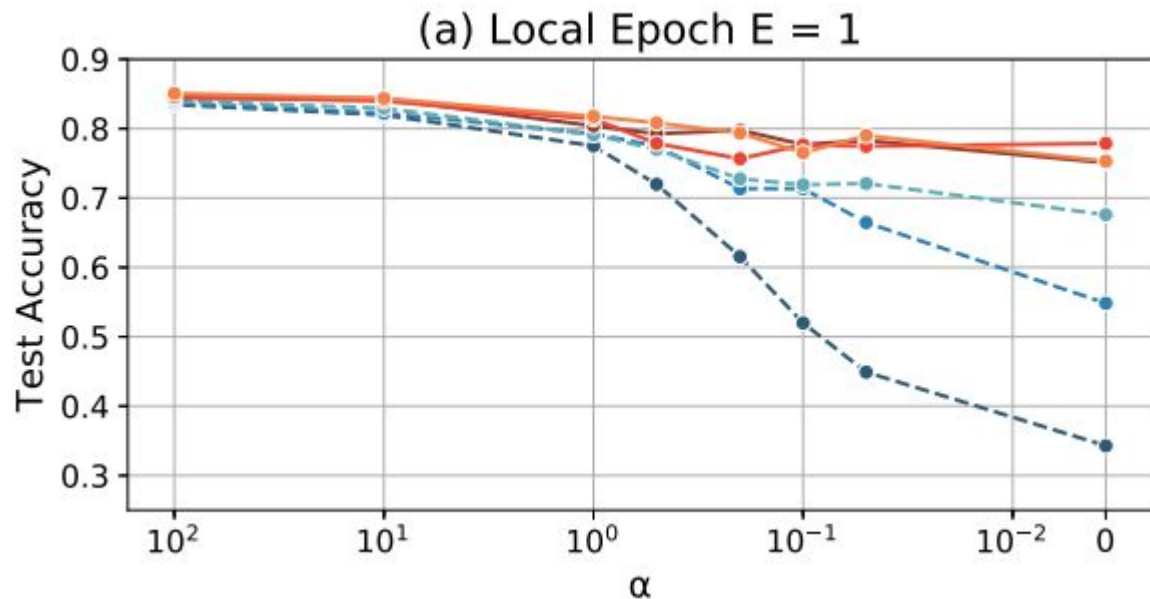


Figure from [Hsu et al. 2019]

FedAvg with momentum (in red) outperforms FedAvg SGD (in blue).

$\alpha$  quantifies non-iidness.

# Solving FL: Mime with momentum

- Apply server momentum locally in the clients

$$y_i = y_i - \eta \left( (1 - \beta) \nabla f_i(y_i) + \beta m \right)$$

Repeat K times

Fixed server momentum

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

$$m = (1 - \beta) \nabla f_i(x) + \beta m$$

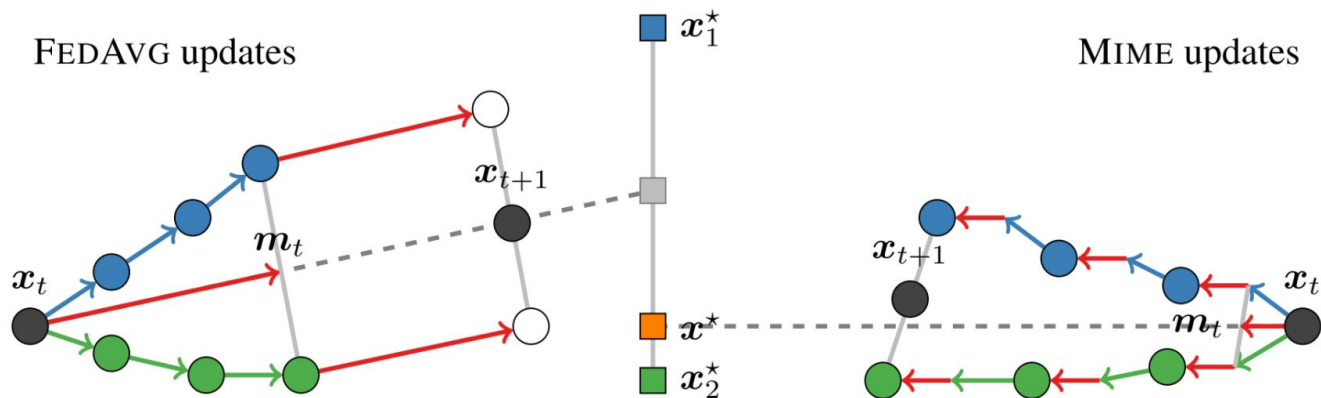
Update server momentum

- Observation 1: Momentum helps

- Observation 2: Momentum helps *more* for non-iid

It must be “mixing” updates from different clients, preventing overfitting.

# FedAvg momentum vs. Mime momentum



# Mime Framework: Principles




1. Use “optimizer state” (momentum) every client update
2. Update “optimizer state” only at server, fixed during client updates.

# Mime Framework: Base optimizer

Decompose a base centralized optimizer as  $\mathcal{B} = (\mathcal{U}, \mathcal{V})$

1. Parameter update step:

$$x = x - \eta \mathcal{U}(g, s)$$

parameters    state

gradient

2. State update step:

$$s = \mathcal{V}(g, s)$$

# Mime Framework: Base optimizer

For SGD with momentum,  $\mathbf{s} = \mathbf{m}$

1. Parameter update step:

$$\mathcal{U}(g, m) = (1 - \beta)g + \beta m$$

2. State update step:

$$\mathcal{V}(g, m) = (1 - \beta)g + \beta m$$

Momentum algorithm.

$$x = x - \eta ((1 - \beta)g + \beta m)$$

$$m = (1 - \beta)g + \beta m$$

# Mime Framework: Full algorithm

- Apply base optimizer locally at the clients

$$y_i = y_i - \eta \mathcal{U}(\nabla f_i(y_i), \boxed{s})$$

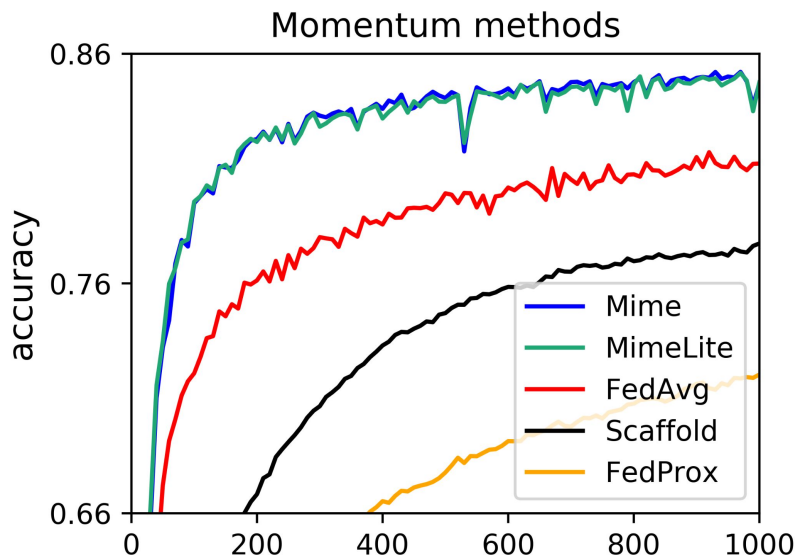
Repeat K times Fixed state

- State is updated only at server

$$s = \mathcal{V}(\nabla f_i(x), s)$$

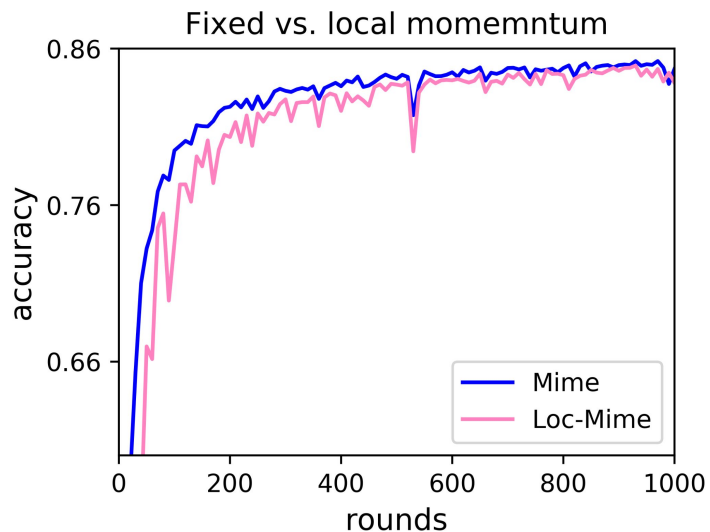


# Experiment results: Comparison



- Extended MNIST 62
- MLP with 2 hidden layers
- 10 local epochs
- 20 of 3400 clients per round
- Momentum = [0, 0.9, 0.99]
- Tuned client lr, server lr = 1.
- Regularization for FedProx tuned over [0.1, 0.5, 1] with 0.1 being the best. Using smaller values may improve performance but regularization 0 is same as FedAvg.

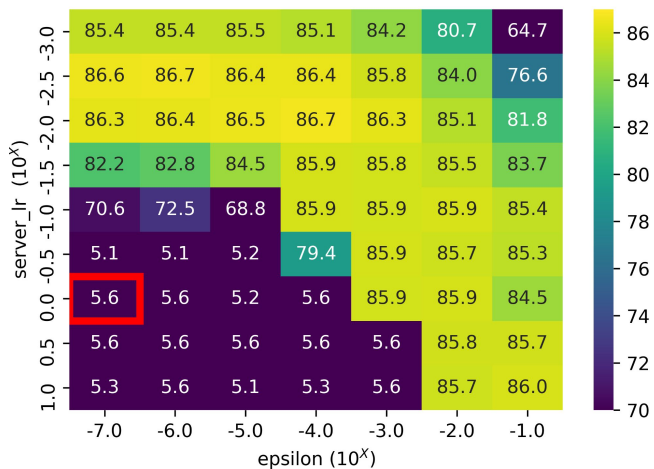
# Experiment results: local momentum



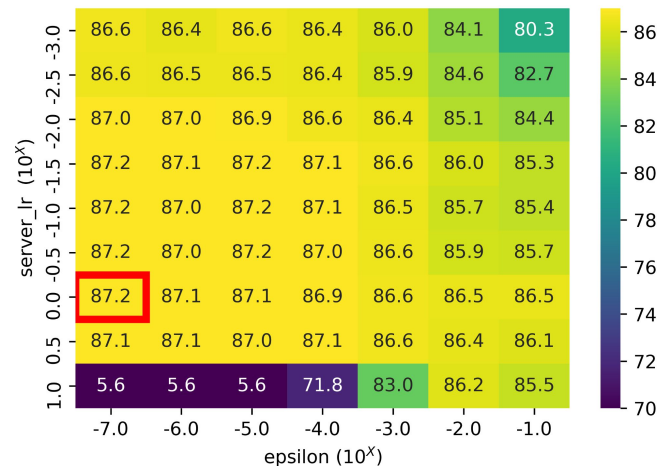
What happens if we also update momentum locally?

# Experiment results: hyperparameters

## Accuracy with FedAvg+Adam



## Accuracy with Mime+Adam



Which hyperparameters work best?

MimeAdam works with same hyperparameters as centralized Adam.

# Part 2: Theory

What is momentum  
doing?

# Cross device FL: Formalism

The diagram shows the formalism equation for cross-device federated learning. The equation is  $\min_x \mathbb{E}_{i \sim \mathcal{D}} \left[ f_i(x) := \frac{1}{n_i} \sum_{\nu=1}^{n_i} f_i(x; \xi_\nu) \right]$ . Annotations include: a red arrow pointing to  $x$  labeled 'parameters'; a green arrow pointing to  $\mathbb{E}_{i \sim \mathcal{D}}$  labeled 'clients'; an orange arrow pointing to  $f_i(x)$  labeled 'loss function'; and a blue arrow pointing to  $\xi_\nu$  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:

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

Note,  $\delta \leq L$ .

# Server-only Momentum based variance reduction

Momentum based variance reduction (MVR) adds a small correction  
[Tran-Dinh et al. 2019, Cutkosky et al. 2020].

$$m^{t+1} = (1-\beta)\nabla f_{i_t}(x^t) + \beta m^t + \beta(\nabla f_{i_t}(x^t) - \nabla f_{i_t}(x^{t-1}))$$

  
Standard momentum

  
Correction

Theorem.  $\min_{t \in [T]} \mathbb{E} \|\nabla f(x_{\text{MVR}}^t)\|^2 \leq \left(\frac{LG}{T}\right)^{2/3}$

Optimal server-only rate!

## Convergence: Initial Attempt

Theorem.  $\min_{t \in [T]} \mathbb{E} \|\nabla f(x_{\text{MimeMVR}}^t)\|^2 \leq \left(\frac{LG}{T}\right)^{2/3}$

Are we done?

Almost, but we can do better.

Prove **advantage** of local steps.

## Local steps: Biased gradients

For random client  $i$  and gradient at server parameters,

$$\mathbb{E}_{i \sim \mathcal{D}} [\nabla f_i(x)] = \nabla f(x)$$

Unbiased gradient

But for gradient at client parameters,

$$\mathbb{E}_{i \sim \mathcal{D}} [\nabla f_i(y_i)] \neq \nabla f(y_i)$$

Causes  
client drift!

Biased gradient since  $y_i$  depends on  $i$



# Local steps: Hessian dissimilarity

Amount of bias controlled by

$$\|\mathbb{E}_{i \sim \mathcal{D}}[\nabla f_i(y_i)] - \nabla f(y_i)\|$$

$$\approx \delta \mathbb{E}_i \|y_i - x\|$$

Diff. in Hessians

Distance moved in a round

□- BHD:

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

- Num. steps should be **inversely** proportional to □

□ - BHD:

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

## Convergence: With MVR

Optimal Server-only momentum based variance reduction (MVR):

$$\text{Theorem. } \min_{t \in [T]} \mathbb{E} \|\nabla f(x_{\text{MVR}}^t)\|^2 \leq \left(\frac{LG}{T}\right)^{2/3}$$

MimeMVR is **faster** than any server-only method!

Mime with momentum based variance reduction (MimeMVR):

$$\text{Theorem. } \min_{t \in [T]} \mathbb{E} \|\nabla f(x_{\text{MimeMVR}}^t)\|^2 \leq \left(\frac{\delta G}{T}\right)^{2/3}$$

# 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.

See you at the poster!