Federated Communication-Efficient Multi-Objective Optimization

Pranay (pranaysh@iitb.ac.in)

C-MInDS, IIT Bombay

Paper: AISTATS'25 (arxiv.org/pdf/2410.16398)

Work with



Baris Askin, Dr. Gauri Joshi and Dr. Carlee Joe-Wong (ECE, CMU)

Big Data







Need lots of data

Data is Naturally Distributed!

• Edge-devices collect data



• This data is used to train ML models

Introduction

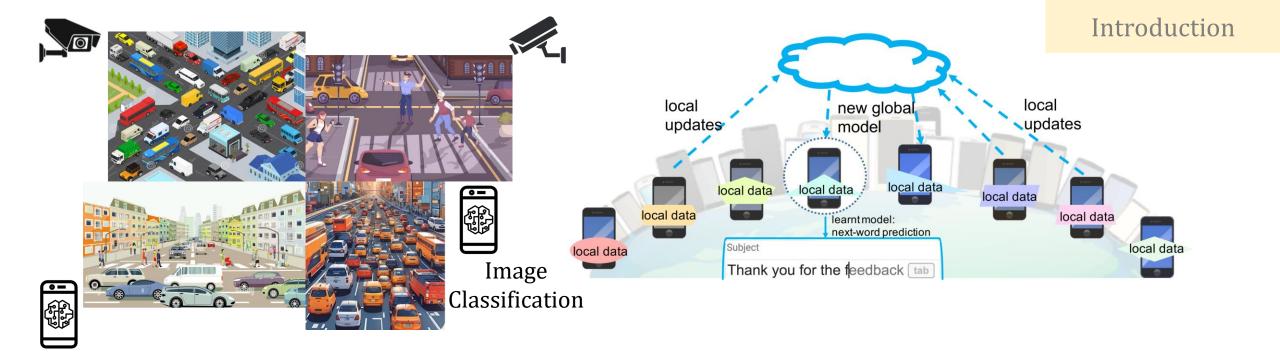


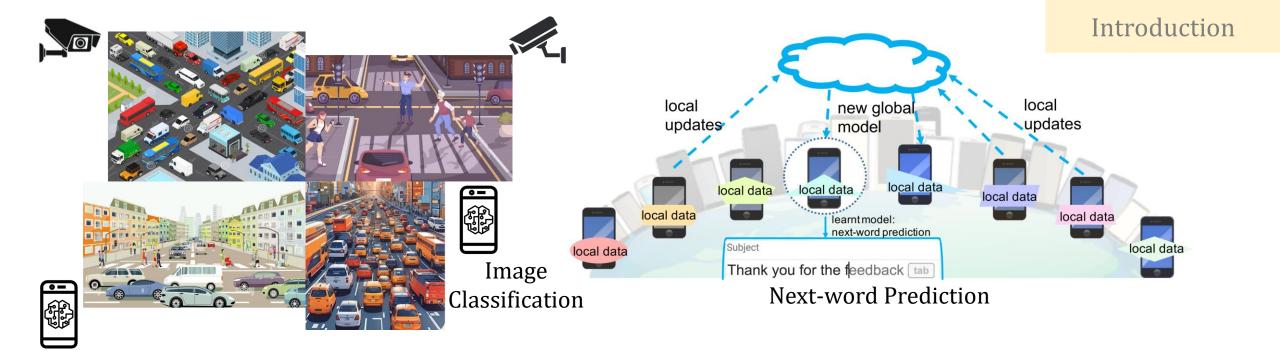


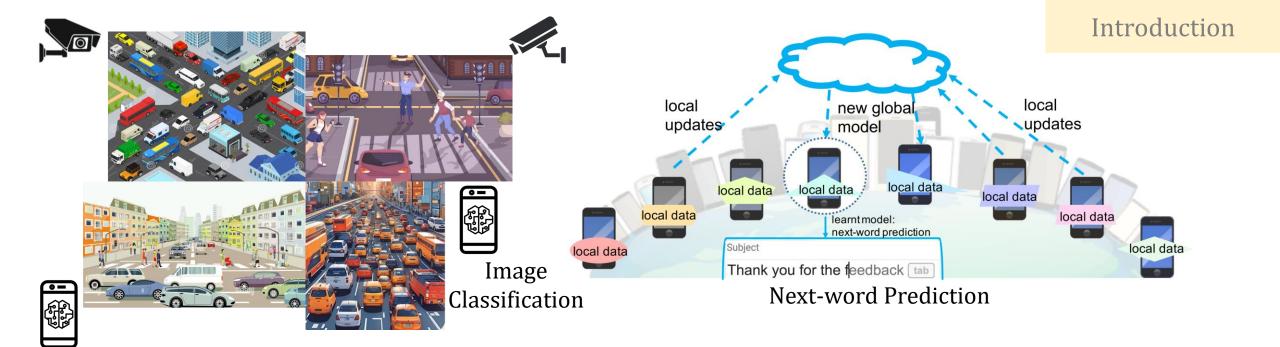
Liu, et al. "Privacy-preserving traffic flow prediction: A federated learning approach." IEEE IoT Journal 2020.
Li, et al. "Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.
Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.

Image curtsey: https://envuetelematics.com/







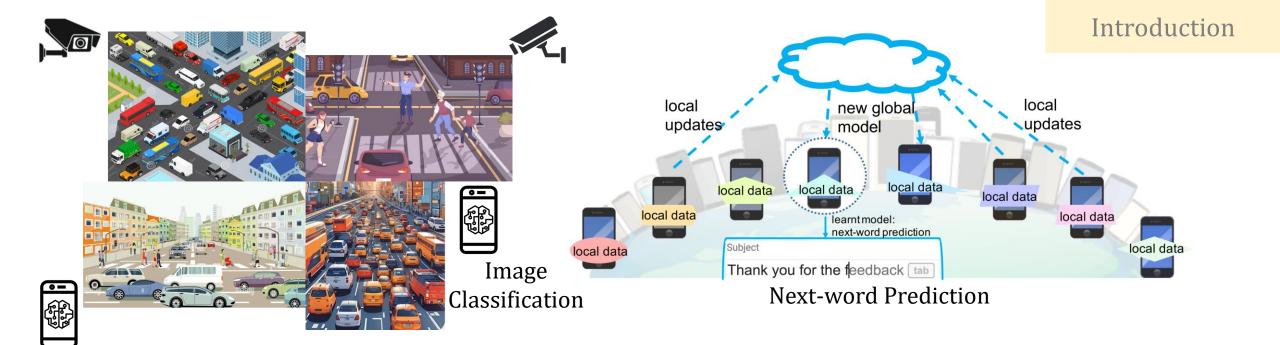




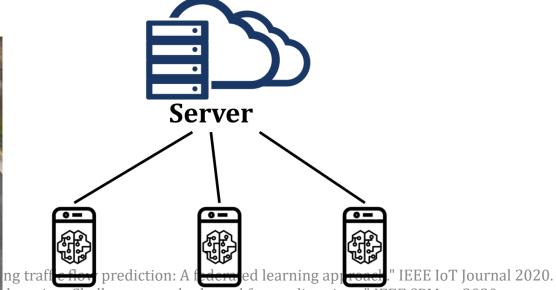
ng traffic flow prediction: A federated learning approach." IEEE IoT Journal 2020.

Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.

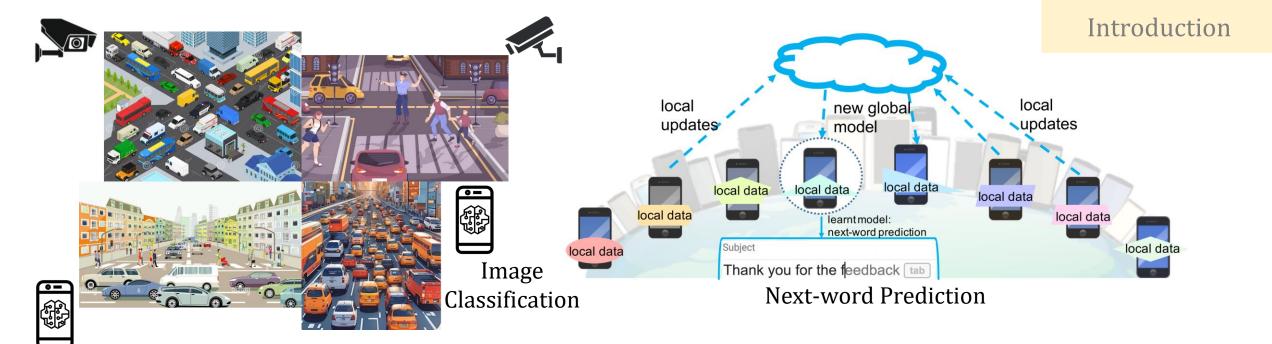




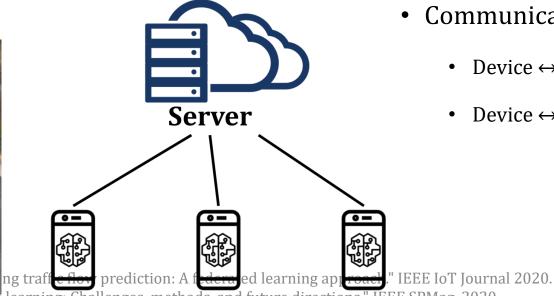


Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.





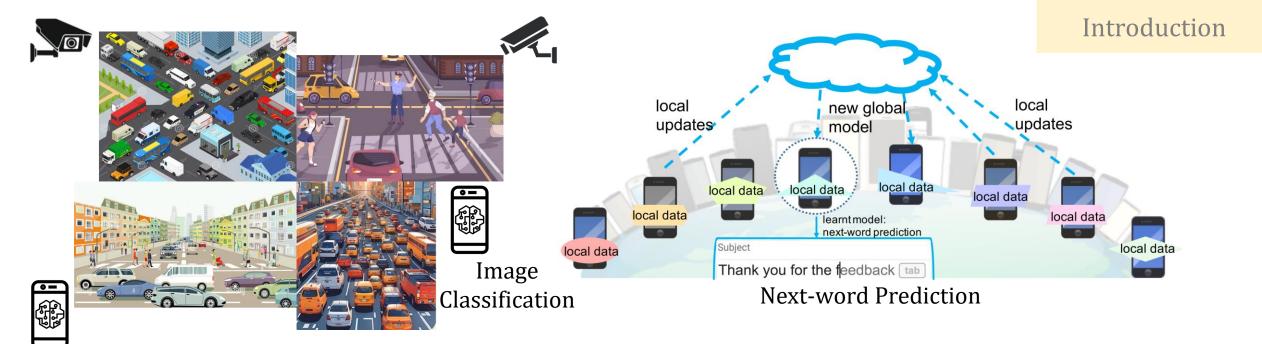


Communication

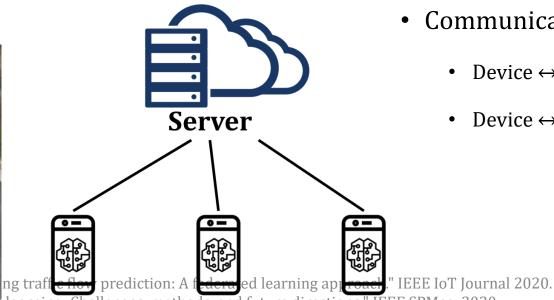
- Device ↔ Server •
- Device ↔ Device

Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.







Communication

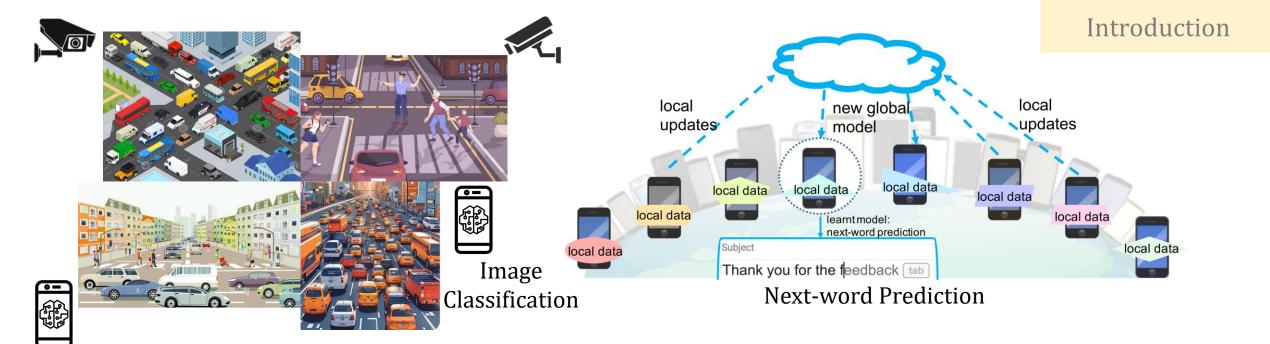
Device ↔ Server 💉



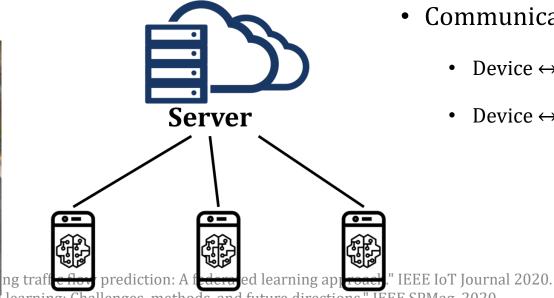
Device ↔ Device

Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.







Communication

Device ↔ Server

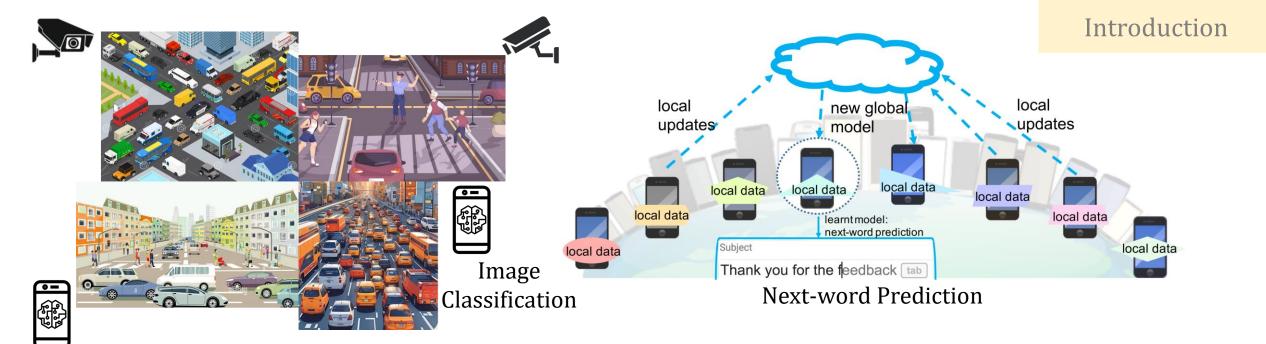


• Device ↔ Device 💢

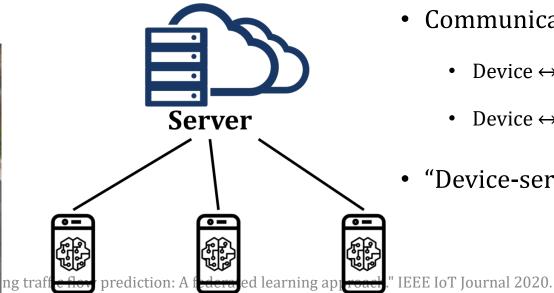


Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.







Communication

• Device ↔ Server

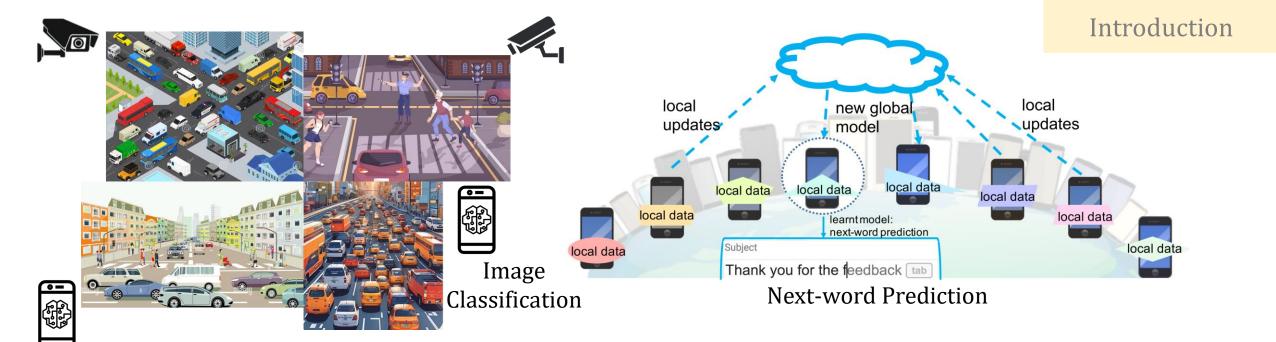




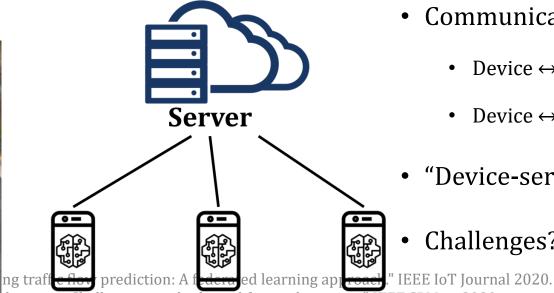
• "Device-server" setting

Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

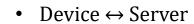
Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.







Communication





Device ↔ Device (X)



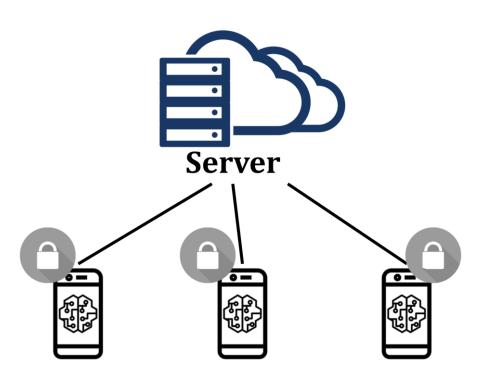
- "Device-server" setting
- Challenges?

Li, et al. Federated learning: Challenges, methods, and future directions." IEEE SPMag, 2020.

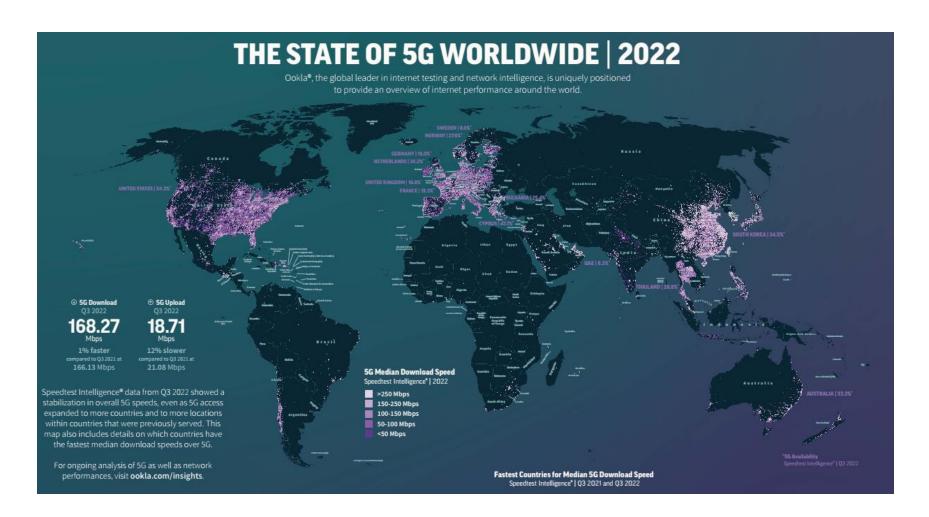
Hard, et al. "Federated learning for mobile keyboard prediction." arXiv:1811.03604.

Privacy

Individual user data should not be leaked



Network constraints







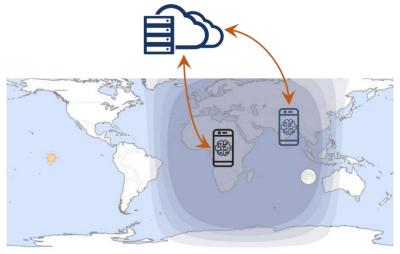


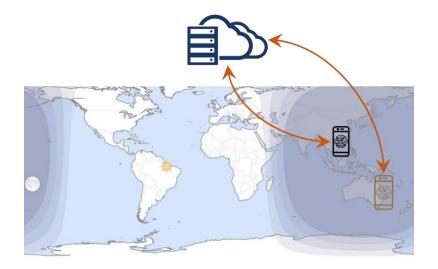
Intermittent Availability

Devices are available when

• Idle, plugged-in and on wifi







Distributed learning under

Privacy concerns
 Devices don't send raw data*



- Network constraints
 Infrequent communication with server
- Intermittent client availability (partial participation)

Federated Learning (FL)









3













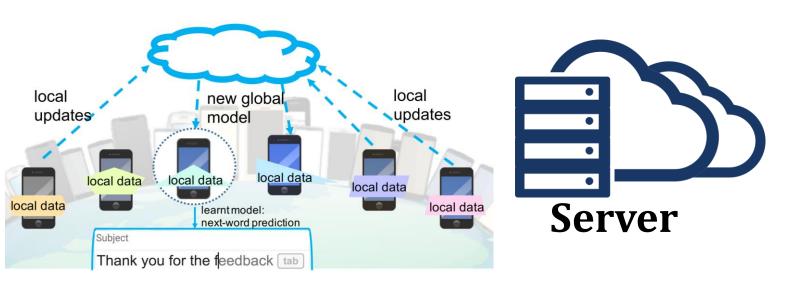






2

...























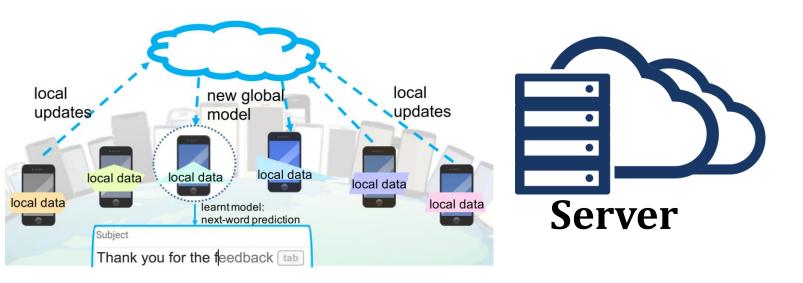




2

3

...

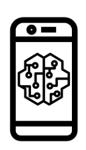


In a small town a greengrocer had opene was located above a deep cellar. Every n in droves out of this cellar into the shop.

























2

...





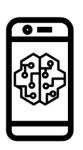
In a small town a greengrocer had opene was located above a deep cellar. Every n in droves out of this cellar into the shop. Where the mind is without fear and the head is held high Where knowledge is free Where the world has not been broken up into fragments by narrow domestic walls;







3



















2

...





In a small town a greengrocer had opene was located above a deep cellar. Every n in droves out of this cellar into the shop. Where the mind is without fear and the head is held high Where knowledge is free Where the world has not been broken up into fragments by narrow domestic walls;



























3

...



Data heterogeneity

























1

2

3

...



Data heterogeneity



























3

...



Data heterogeneity

System heterogeneity



10

Round *t* begins

• P (out of n) clients



























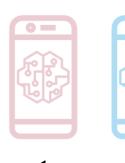
3

Clients

Round *t* begins

• *P* (out of *n*) clients









3



















Clients

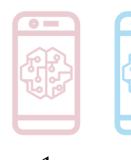
Round *t* begins

• *P* (out of *n*) clients



Partial Client Participation

Small fraction of clients available



















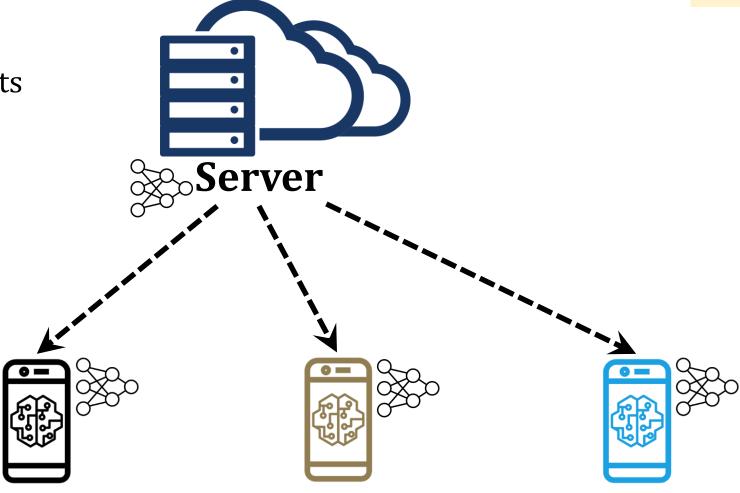






Clients

• Global model \rightarrow clients

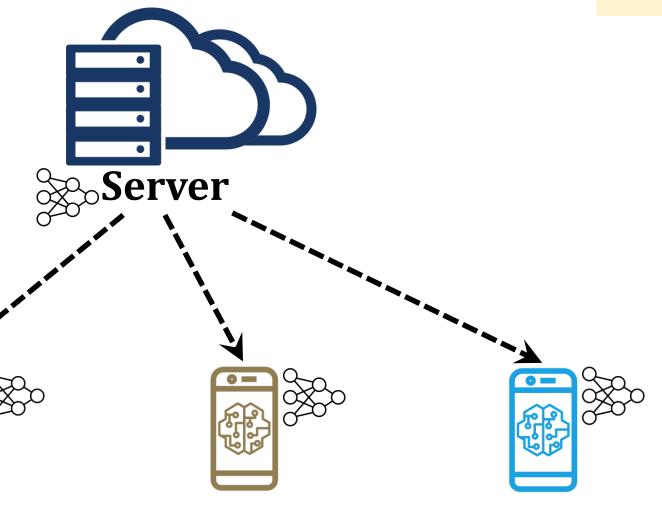


• Global model → clients

Communication Efficiency

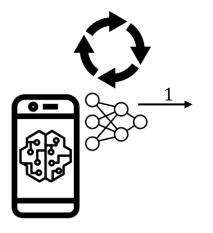
Network constraints

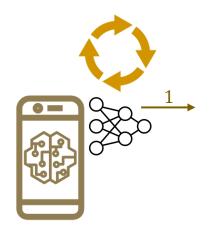
Infrequent

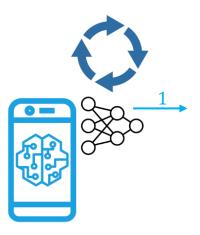


• Clients run $\tau \ge 1$ local steps



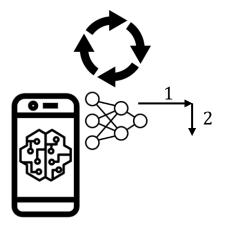


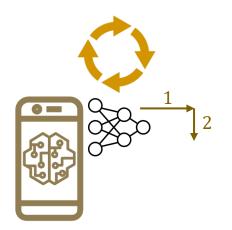


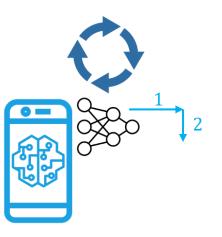


• Clients run $\tau \ge 1$ local steps



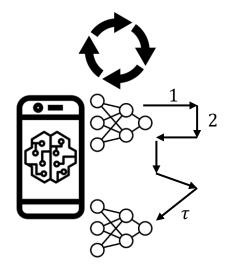


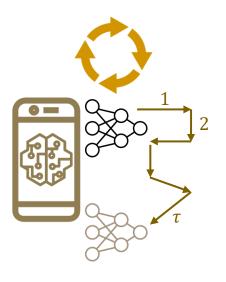


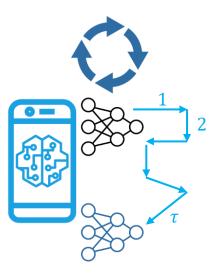


• Clients run $\tau \ge 1$ local steps



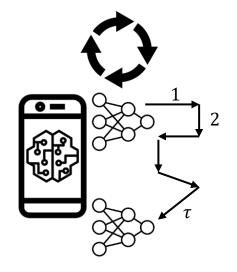


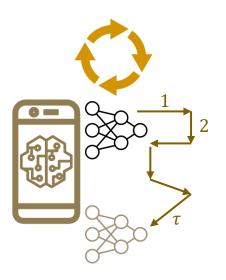


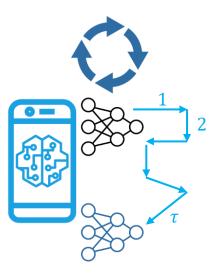


- Clients run $\tau \ge 1$ local steps
- Saves communication cost









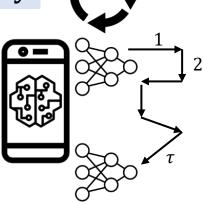
- Clients run $\tau \ge 1$ local steps
- Saves communication cost

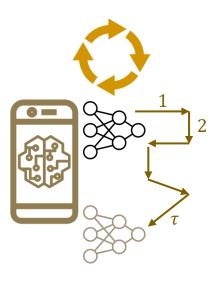


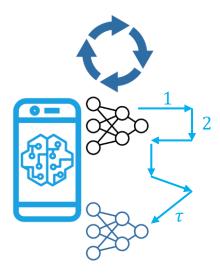
Computation Efficiency

Low-power devices

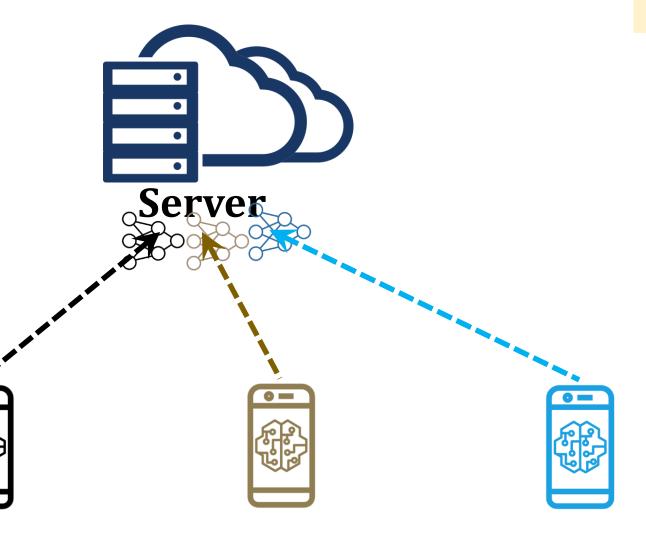
First-order methods







Local models → server

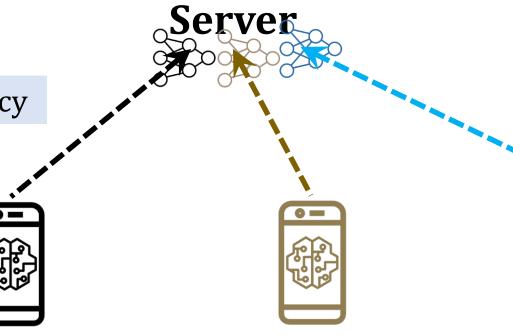


• Local models → server

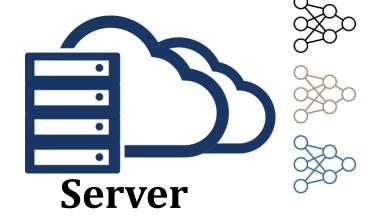
Communication Efficiency

Network constraints

Infrequent

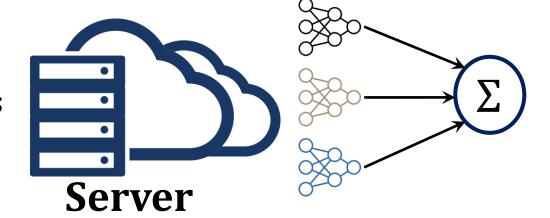


• Server aggregates local models



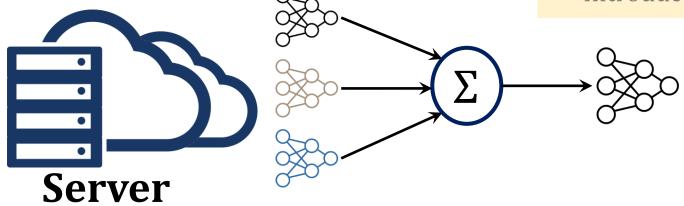


• Server aggregates local models





• Server aggregates local models



























1

2

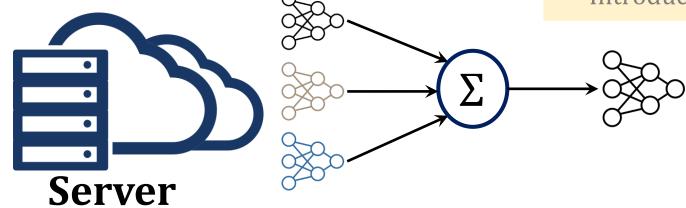
3

Clients

n

• Server aggregates local models

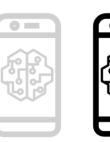
Round *t* ends



























Clients

n

Round t + 1 begins

• P (out of n) clients





Round t + 1 begins

• *P* (out of *n*) clients



Partial Client Participation

Small fraction of clients available









3















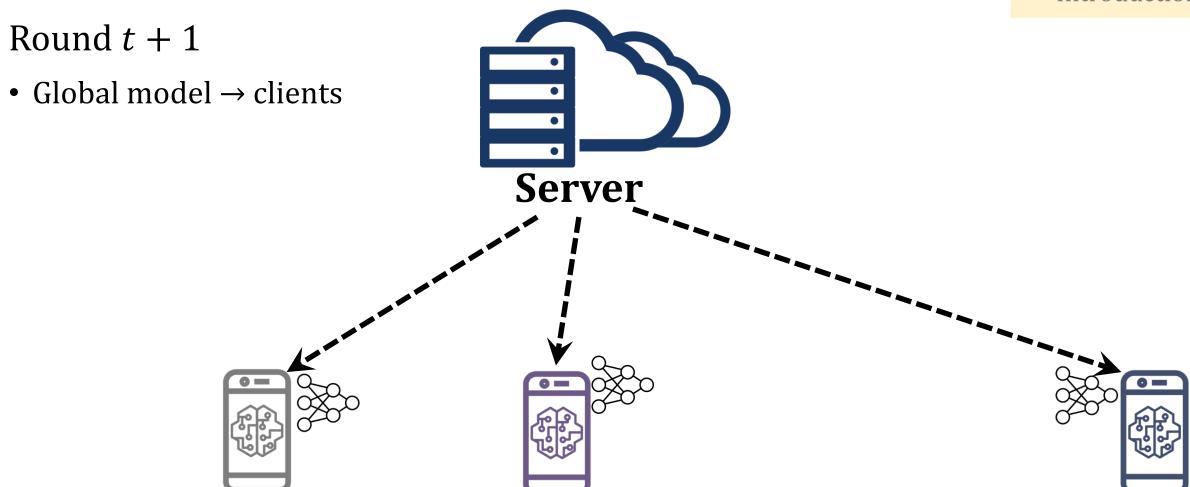






Clients

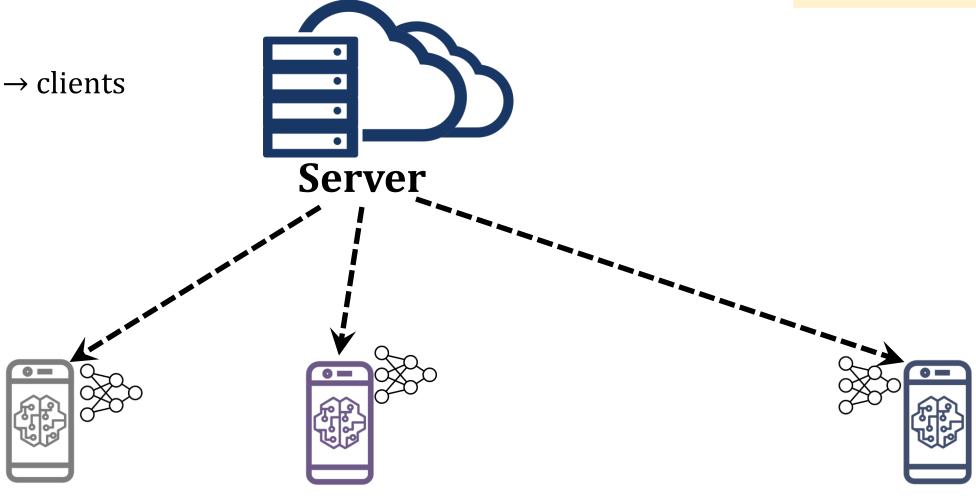
n



Round t + 1

• Global model → clients

Total *T* rounds



ChatBots



Multiple Objectives

- Accuracy/Relevance
- Engagement
- Politeness/Safety
- Personalization

Trade-offs: accuracy \leftrightarrow personalization

Recommender Systems*



Image modified from MLArchive

Multiple Users

- User Diversity
- Data Privacy

Multiple Objectives

- Accuracy/Relevance
- Diversity
- Novelty
- Contextual Relevance

 $\mathsf{Trade}\text{-}\mathsf{offs}\text{: }\mathsf{accuracy} \leftrightarrow \mathsf{diversity}$

^{*}Sun, et al. "A survey on federated recommendation systems." IEEE TNNLS'24.

Personalized Medicine*

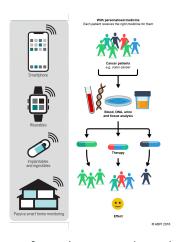


Image from science.org and eupati.eu

Multiple Patients

- Patient diversity
- Data confidentiality

Multiple Objectives

- Diagnostic Precision
- Cost-Effectiveness
- Minimizing Side-Effects
- Privacy and Data Security
- Equity in Access

Trade-offs:

 $\mathsf{Precision} \leftrightarrow \mathsf{Cost}\text{-}\mathsf{Effectiveness}$

^{*}Sheller, et al. "Federated learning in medicine: facilitating multi-institutional collaborations without sharing patient data." Scientific Reports'20.

Goal

Using data from multiple (distributed) users/patients/ $\underline{clients}$, learn a model that simultaneously optimizes multiple objectives

Federated Multi-Objective Optimization (MOO)

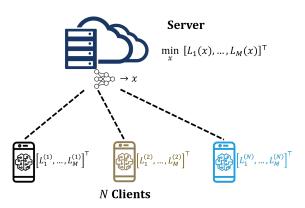
$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^\top, \quad \text{ where } \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective

Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^\top, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective



Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^\top, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective

M objectives, N users/patients/clients

Communication-efficient algorithm FedCMOO

Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective

- Communication-efficient algorithm FedCMOO
 - Communication cost <u>does not scale</u> with number of objectives M

Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective

- Communication-efficient algorithm FedCMOO
 - Communication cost does not scale with number of objectives M
- Convergence to Pareto-stationary point

Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective

- Communication-efficient algorithm FedCMOO
 - Communication cost does not scale with number of objectives M
- Convergence to Pareto-stationary point
 - Improved sample complexity under weaker assumptions

Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

 $L_k^{(i)}$ - contribution of *i*-th client to *k*-th objective

- Communication-efficient algorithm FedCMOO
 - Communication cost does not scale with number of objectives M
- Convergence to Pareto-stationary point
 - Improved sample complexity under <u>weaker</u> assumptions
- Improved empirical performance

Outline

1. Background: MOO and MGDA

2. Federated MGDA and its Drawbacks

3. FedCMOO: Improving Federated MGDA

4. Experiment Results

5. Conclusion

Background: MOO and MGDA

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}$$

Why do we need it? Why not just solve $\min_{\mathbf{x}} \sum_{k=1}^{M} \alpha_k L_k(\mathbf{x})$?

[†]Hu, et al. "Revisiting Scalarization in Multi-Task Learning: A Theoretical Perspective," *NeurIPS'23*.

^{*}Sener and Koltun, "Multi-task learning as multi-objective optimization," NeurIPS'18.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} (L_1(\mathbf{x}), \dots, L_M(\mathbf{x}))^{\top}$$

Why do we need it? Why not just solve $\min_{\mathbf{x}} \sum_{k=1}^{M} \alpha_k L_k(\mathbf{x})$?

• How do we decide the weights $\{\alpha_k\}_{k=1}^M$?

[†]Hu, et al. "Revisiting Scalarization in Multi-Task Learning: A Theoretical Perspective," *NeurIPS'23*.

^{*}Sener and Koltun, "Multi-task learning as multi-objective optimization," NeurIPS'18.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}$$

Why do we need it? Why not just solve $\min_{\mathbf{x}} \sum_{k=1}^{M} \alpha_k L_k(\mathbf{x})$?

- How do we decide the weights $\{\alpha_k\}_{k=1}^M$?
- A subset of tasks might be severely under-optimized

[†]Hu, et al. "Revisiting Scalarization in Multi-Task Learning: A Theoretical Perspective," NeurIPS'23.

^{*}Sener and Koltun, "Multi-task learning as multi-objective optimization," NeurIPS'18.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} (L_1(\mathbf{x}), \dots, L_M(\mathbf{x}))^{\top}$$

Why do we need it? Why not just solve $\min_{\mathbf{x}} \sum_{k=1}^{M} \alpha_k L_k(\mathbf{x})$?

- How do we decide the weights $\{\alpha_k\}_{k=1}^M$?
- A subset of tasks might be severely under-optimized
- Scalar losses do not explore the set of all possible solutions[†]

[†]Hu, et al. "Revisiting Scalarization in Multi-Task Learning: A Theoretical Perspective," NeurIPS'23.

^{*}Sener and Koltun, "Multi-task learning as multi-objective optimization," NeurIPS'18.

MGDA with Multi-MNIST Dataset*

Two digits - one in top-left, the other in bottom-right



^{*}Sener and Koltun, "Multi-task learning as multi-objective optimization," NeurIPS'18.

MGDA with Multi-MNIST Dataset*

Two digits - one in top-left, the other in bottom-right



Task-L/R: classifying the digit on top-left/bottom-right

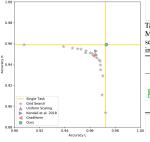


Table 3: Performance of MTL algorithms on MultiMNIST. Single-task baselines solve tasks separately, with dedicated models, but are shown in the same row for clarity.

	accuracy [Left digit accuracy [%]	
90	95.90	97.23	Single task
99	94.99	96.46	Uniform scaling
29	95.29	96.47	Kendall et al. 2018
34	94.84	96.27	GradNorm
90	95.90	97.26	Ours
	95.	97.26	Ours

^{*}Sener and Koltun, "Multi-task learning as multi-objective optimization," NeurIPS'18.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}$$

• At Pareto-optimal point - cannot reduce any loss L_k without also increasing some other $L_{k'}$

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^\top$$

- At Pareto-optimal point cannot reduce any loss L_k without also increasing some other $L_{k'}$
- At Pareto-stationary point $\bar{\mathbf{x}}$ "some" convex combination of gradient vectors is zero, i.e., there exist w_1, \ldots, w_M (all ≥ 0) s.t.

$$\sum_{i=1}^M w_k = 1$$
 and $\sum_{i=1}^M w_k
abla L_k(ar{\mathbf{x}}) = \mathbf{0}$

i.e., no common descent direction for all objectives

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}$$

- At Pareto-optimal point cannot reduce any loss L_k without also increasing some other L_{k'}
- At Pareto-stationary point $\bar{\mathbf{x}}$ "some" convex combination of gradient vectors is zero, i.e., there exist w_1, \ldots, w_M (all ≥ 0) s.t.

$$\sum_{i=1}^M w_k = 1$$
 and $\sum_{i=1}^M w_k
abla L_k(ar{\mathbf{x}}) = \mathbf{0}$

i.e., no common descent direction for all objectives

Pareto-optimality ⇒ Pareto-stationarity

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}$$

- At Pareto-optimal point cannot reduce any loss L_k without also increasing some other $L_{k'}$
- At Pareto-stationary point $\bar{\mathbf{x}}$ "some" convex combination of gradient vectors is zero, i.e., there exist w_1, \ldots, w_M (all ≥ 0) s.t.

$$\sum_{i=1}^M w_k = 1$$
 and $\sum_{i=1}^M w_k
abla L_k(ar{\mathbf{x}}) = \mathbf{0}$

i.e., no common descent direction for all objectives

- Pareto-optimality ⇒ Pareto-stationarity
- $\bar{\mathbf{x}}$ is not Pareto-stationary \Rightarrow can find descent directions to simultaneously reduce all $\{L_k\}_{k=1}^M$

Multiple Gradient Descent

Multiple Gradient Descent Algorithm (MGDA)* solves

$$\min_{w_1,\ldots,w_M} \left\| \sum_{i=1}^M w_k \nabla L_k(\mathbf{x}) \right\|^2 \quad \text{s.t.} \quad \sum_{i=1}^M w_k = 1, w_k \geq 0, \text{ for all } k \in [M]$$

^{*}Désidéri, "Multiple-gradient descent algorithm (MGDA) for multiobjective optimization," Comptes Rendus Mathematique, 2012.

Multiple Gradient Descent

Multiple Gradient Descent Algorithm (MGDA)* solves

$$\min_{w_1,\dots,w_M} \left\| \sum_{i=1}^M w_k \nabla L_k(\mathbf{x}) \right\|_2^2 \quad \text{s.t.} \quad \sum_{i=1}^M w_k = 1, w_k \ge 0, \text{ for all } k \in [M]$$

• Either the minimum value is 0,

^{*}Désidéri, "Multiple-gradient descent algorithm (MGDA) for multiobjective optimization," Comptes Rendus Mathematique, 2012.

Multiple Gradient Descent

Multiple Gradient Descent Algorithm (MGDA)* solves

$$\min_{w_1,\dots,w_M} \left\| \sum_{i=1}^M w_k \nabla L_k(\mathbf{x}) \right\|_2^2 \quad \text{s.t.} \quad \sum_{i=1}^M w_k = 1, w_k \ge 0, \text{ for all } k \in [M]$$

- Either the minimum value is 0,
- Or the solution gives $\sum_{i=1}^{M} w_k^* \nabla L_k(\mathbf{x})$ that reduces all objectives

^{*}Désidéri, "Multiple-gradient descent algorithm (MGDA) for multiobjective optimization," Comptes Rendus Mathematique, 2012.

Outline

1. Background: MOO and MGDA

2. Federated MGDA and its Drawbacks

3. FedCMOO: Improving Federated MGDA $\,$

4. Experiment Results

5. Conclusion

Federated MGDA and its

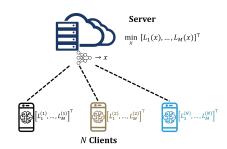
Drawbacks

Problem - Federated Multi-Objective Optimization

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top}$$

$$L_k(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

• M objectives, N clients



First Attempt - Federated MGDA*

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

M objectives, N clients

Recall MGDA -
$$\min_{\mathbf{w}} \left\| \sum_{i=1}^{M} w_k \nabla L_k(\mathbf{x}) \right\|_2^2$$
 s.t. $\mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^{\top} \mathbf{1} = 1$

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

First Attempt - Federated MGDA*

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

M objectives, N clients

Recall MGDA -
$$\min_{\mathbf{w}} \left\| \sum_{i=1}^{M} w_k \nabla L_k(\mathbf{x}) \right\|_2^2$$
 s.t. $\mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^{\top} \mathbf{1} = 1$

• Client i computes $\approx \left[\nabla L_1^{(i)}(\mathbf{x}), \dots, \nabla L_M^{(i)}(\mathbf{x}) \right]^{\top}$ (local gradients)

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

First Attempt - Federated MGDA*

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) = \min_{\mathbf{x}} \left(L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right)^{\top}, \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

M objectives, N clients

Recall MGDA -
$$\min_{\mathbf{w}} \left\| \sum_{i=1}^{M} w_k \nabla L_k(\mathbf{x}) \right\|_2^2$$
 s.t. $\mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^{\top} \mathbf{1} = 1$

- Client i computes $\approx \left[\nabla L_1^{(i)}(\mathbf{x}), \dots, \nabla L_M^{(i)}(\mathbf{x}) \right]^{\top}$ (local gradients)
- Server computes $\approx [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]^{\top}$, where $\nabla L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \nabla L_k^{(i)}(\mathbf{x})$

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

• Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

- Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_{k}^{(i)} \leftarrow \mathbf{x}_{k}^{(i)} \eta_{c} \widetilde{\nabla} L_{k}^{(i)}(\mathbf{x}_{k}^{(i)})$ (SGD/GD/etc.)

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

- Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$ (SGD/GD/etc.)
- Return updates to server $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{r_c}$ for each k

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

- Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$ (SGD/GD/etc.)
- Return updates to server $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{r_c}$ for each k

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

- Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$ (SGD/GD/etc.)
- Return updates to server $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k

At the server

• Average client updates $\Delta_k = \frac{1}{N} \sum_i \Delta_k^{(i)}$, for each objective $k \in M$

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

- Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$ (SGD/GD/etc.)
- Return updates to server $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k

- Average client updates $\Delta_k = \frac{1}{N} \sum_i \Delta_k^{(i)}$, for each objective $k \in M$
- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \| \sum_{k} w_k \Delta_k \|_2^2$ s.t. $\mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^{\top} \mathbf{1} = 1$

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

At each client i

- Synchronize local models $\mathbf{x}_{k}^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$ (SGD/GD/etc.)
- Return updates to server $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k

- Average client updates $\Delta_k = \frac{1}{N} \sum_i \Delta_k^{(i)}$, for each objective $k \in M$
- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \| \sum_{k} w_k \Delta_k \|_2^2$ s.t. $\mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^{\top} \mathbf{1} = 1$
- Update global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \sum_k w_k^* \Delta_k$

^{*}Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

Challenges of Federated MGDA

At each client i

- Synchronize local models $\mathbf{x}_k^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- ullet Return updates to server $\Delta_k^{(i)} = rac{{\sf x} {\sf x}_k^{(i)}}{\eta_c}$ for each task k

- Average client updates $\Delta_k = \frac{1}{N} \sum_i \Delta_k^{(i)}$, for each task $k \in M$
- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \sum_k w_k \Delta_k \right\|_2^2$ s.t. $\mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$
- Update global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \sum_k w_k^* \Delta_k$

Challenges of Federated MGDA

At each client i

- Synchronize local models $\mathbf{x}_k^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\Delta_k^{(i)} = \frac{\mathsf{x} \mathsf{x}_k^{(i)}}{\eta_c}$ for each task k
 - M× communication

- Average client updates $\Delta_k = \frac{1}{N} \sum_i \Delta_k^{(i)}$, for each task $k \in M$
- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \sum_k w_k \Delta_k \right\|_2^2$ s.t. $\mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$
- Update global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \sum_k w_k^* \Delta_k$

Challenges of Federated MGDA

At each client i

- Synchronize local models $\mathbf{x}_k^{(i)} = \mathbf{x}$ for objective k and client i
- Multiple local updates $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Local drift across tasks ⇒ worse performance
- Return updates to server $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each task k
 - M× communication

- Average client updates $\Delta_k = \frac{1}{N} \sum_i \Delta_k^{(i)}$, for each task $k \in M$
- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \sum_k w_k \Delta_k \right\|_2^2$ s.t. $\mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$
- Update global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \sum_k w_k^* \Delta_k$

FedCMOO: Improving Federated

MGDA

If we knew the task weights $\{w_k\}$ At each client i - aggregated local updates, starting with x

If we knew the task weights $\{w_k\}$

At each client i - aggregated local updates, starting with x

• Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} - \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$

If we knew the task weights $\{w_k\}$

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$

If we knew the task weights $\{w_k\}$

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\nabla^{(i)} = \frac{\mathbf{x} \mathbf{x}^{(i)}}{\eta_c}$

If we knew the task weights $\{w_k\}$

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\nabla^{(i)} = \frac{\mathbf{x} \mathbf{x}^{(i)}}{\eta_c}$
 - Instead of task-wise updates $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k

If we knew the task weights $\{w_k\}$

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\nabla^{(i)} = \frac{\mathbf{x} \mathbf{x}^{(i)}}{\eta_c}$
 - Instead of task-wise updates $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k
- *M*× reduction in communication; reduces drift across tasks

If we knew the task weights $\{w_k\}$

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\nabla^{(i)} = \frac{\mathbf{x} \mathbf{x}^{(i)}}{\eta_c}$
 - Instead of task-wise updates $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k
- *M*× reduction in communication; reduces drift across tasks

If we knew the task weights $\{w_k\}$

At each client i - aggregated local updates, starting with x

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\nabla^{(i)} = \frac{\mathbf{x} \mathbf{x}^{(i)}}{\eta_c}$
 - Instead of task-wise updates $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k
- M× reduction in communication; reduces drift across tasks

At the server

- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \sum_k w_k \Delta_k \right\|_2^2$ s.t. $\mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$
- Update global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \sum_{\mathbf{k}} w_{\mathbf{k}}^* \Delta_{\mathbf{k}}$

Repeat

If we knew the task weights $\{w_k\}$

At each client i - aggregated local updates, starting with x

- Multiple local updates: $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
 - Instead of $\mathbf{x}_k^{(i)} \leftarrow \mathbf{x}_k^{(i)} \eta_c \widetilde{\nabla} L_k^{(i)}(\mathbf{x}_k^{(i)})$
- Return updates to server $\nabla^{(i)} = \frac{\mathbf{x} \mathbf{x}^{(i)}}{\eta_c}$
 - Instead of task-wise updates $\Delta_k^{(i)} = \frac{\mathbf{x} \mathbf{x}_k^{(i)}}{\eta_c}$ for each k
- M× reduction in communication; reduces drift across tasks

At the server

- Solve $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \sum_k w_k \Delta_k \right\|_2^2$ s.t. $\mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$
- Update global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \sum_{\mathbf{k}} w_{\mathbf{k}}^* \Delta_{\mathbf{k}}$

Repeat

How do we update **w** with $\{\nabla^{(i)}\}$ instead of $\{\Delta_k\}$?

Why do we need $M \times$ communication?

At the server - need to solve

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} \| \sum_k w_k \nabla L_k(\mathbf{x}) \|_2^2 \quad \text{s.t.} \quad \mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1 \\ &= \operatorname{argmin}_{\mathbf{w}} \left\| \underbrace{\left[\nabla L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]}_{\triangleq \nabla \mathbf{L}(\mathbf{x}) \in \mathbb{R}^{d \times M}} \mathbf{w} \right\|_2^2 \quad \text{s.t.} \quad \mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1 \\ \\ \mathbf{w}^* &\approx \operatorname{Proj}_{\mathcal{S}_M} \left(\mathbf{w} - \eta_w \nabla \mathbf{L}(\mathbf{x})^\top \nabla \mathbf{L}(\mathbf{x}) \mathbf{w} \right) \end{aligned} \qquad \text{(Single PGD step)}$$

 \mathcal{S}_M is the probability simplex in \mathbb{R}^M

Why do we need $M \times$ communication?

At the server - need to solve

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \| \sum_{k} w_k \nabla L_k(\mathbf{x}) \|_2^2 \quad \text{s.t.} \quad \mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$$

$$= \operatorname{argmin}_{\mathbf{w}} \left\| \underbrace{\left[\nabla L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]}_{\triangleq \nabla \mathbf{L}(\mathbf{x}) \in \mathbb{R}^{d \times M}} \mathbf{w} \right\|_2^2 \quad \text{s.t.} \quad \mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$$

$$\mathbf{w}^* \approx \operatorname{Proj}_{\mathcal{S}_M} \left(\mathbf{w} - \eta_w \nabla \mathbf{L}(\mathbf{x})^\top \nabla \mathbf{L}(\mathbf{x}) \mathbf{w} \right) \qquad \qquad \text{(Single PGD step)}$$

 \mathcal{S}_M is the probability simplex in \mathbb{R}^M

• In Federated MGDA:

$$[\Delta_1, \dots, \Delta_M] \approx \nabla L(\mathbf{x})$$
Matrix of task updates Jacobian

Why do we need $M \times$ communication?

At the server - need to solve

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} \| \sum_k w_k \nabla L_k(\mathbf{x}) \|_2^2 \quad \text{s.t.} \quad \mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1 \\ &= \operatorname{argmin}_{\mathbf{w}} \left\| \underbrace{\left[\nabla L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]}_{\triangleq \nabla \mathbf{L}(\mathbf{x}) \in \mathbb{R}^{d \times M}} \mathbf{w} \right\|_2^2 \quad \text{s.t.} \quad \mathbf{w} \geq \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1 \\ \\ \mathbf{w}^* &\approx \operatorname{Proj}_{\mathcal{S}_M} \left(\mathbf{w} - \eta_w \nabla \mathbf{L}(\mathbf{x})^\top \nabla \mathbf{L}(\mathbf{x}) \mathbf{w} \right) \end{aligned} \qquad \text{(Single PGD step)}$$

 \mathcal{S}_M is the probability simplex in \mathbb{R}^M

In Federated MGDA:

$$[\Delta_1, \dots, \Delta_M] \approx \nabla L(\mathbf{x})$$
Matrix of task updates Jacobian

• How do we approximate $\nabla \mathbf{L}(\mathbf{x})$ with $\mathcal{O}(d)$ communication?

Server wants
$$\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$$

Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

Server wants
$$\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$$

Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

At clients

• Client i reshapes \widetilde{H}_i to $\bar{H}_i \in \mathbb{R}^{\sqrt{dM} \times \sqrt{dM}}$

Server wants
$$\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$$

Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

At clients

- Client *i* reshapes \widetilde{H}_i to $\overline{H}_i \in \mathbb{R}^{\sqrt{dM} \times \sqrt{dM}}$
- Send $\widehat{H}_i \leftarrow \text{Rand-SVD}(\overline{H}_i, r)$ (rank-r approximation of \overline{H}_i) to server

Server wants
$$\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$$

Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

At clients

- Client *i* reshapes \widetilde{H}_i to $\overline{H}_i \in \mathbb{R}^{\sqrt{dM} \times \sqrt{dM}}$
- Send $\widehat{H}_i \leftarrow \text{Rand-SVD}(\overline{H}_i, r)$ (rank-r approximation of \overline{H}_i) to server
- Rank $r \le \sqrt{d/M} \Rightarrow$ communication cost $\le d$

Server wants
$$\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$$

Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

At clients

- Client *i* reshapes \widetilde{H}_i to $\overline{H}_i \in \mathbb{R}^{\sqrt{dM} \times \sqrt{dM}}$
- Send $\widehat{H}_i \leftarrow \text{Rand-SVD}(\overline{H}_i, r)$ (rank-r approximation of \overline{H}_i) to server
- Rank $r \le \sqrt{d/M} \Rightarrow$ communication cost $\le d$

Constructing $\nabla L(x)$ with $\mathcal{O}(d)$ Communication

Server wants $\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$ Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

At clients

- Client *i* reshapes \widetilde{H}_i to $\overline{H}_i \in \mathbb{R}^{\sqrt{dM} \times \sqrt{dM}}$
- Send $\widehat{H}_i \leftarrow \text{Rand-SVD}(\overline{H}_i, r)$ (rank-r approximation of \overline{H}_i) to server
- Rank $r \le \sqrt{d/M} \Rightarrow$ communication cost $\le d$

At server

• Reshape \widehat{H}_i to $H_i \in \mathbb{R}^{d \times M}$

Constructing $\nabla L(x)$ with $\mathcal{O}(d)$ Communication

Server wants $\nabla \mathbf{L} = [\nabla L_1(\mathbf{x}), \dots, \nabla L_M(\mathbf{x})]$ Client i has $\widetilde{H}_i = \left[\widetilde{\nabla} L_1^{(i)} \dots \widetilde{\nabla} L_M^{(i)}\right] \in \mathbb{R}^{d \times M}$

At clients

- Client *i* reshapes \widetilde{H}_i to $\overline{H}_i \in \mathbb{R}^{\sqrt{dM} \times \sqrt{dM}}$
- Send $\widehat{H}_i \leftarrow \text{Rand-SVD}(\overline{H}_i, r)$ (rank-r approximation of \overline{H}_i) to server
- Rank $r \le \sqrt{d/M} \Rightarrow$ communication cost $\le d$

At server

- Reshape \widehat{H}_i to $H_i \in \mathbb{R}^{d \times M}$
- Compute $G \leftarrow \left(\frac{1}{N}\sum_{i}H_{i}\right)^{\top}\left(\frac{1}{N}\sum_{i}H_{i}\right) \approx \nabla \mathbf{L}^{\top}\nabla \mathbf{L}$

Comparison with other Compression Schemes

Compression Method	CelebA	CIFAR10+ MNIST	MNIST+ F.MNIST	MultiMNIST	Mean
Sum of client Gram matrices	50.95	11.51	20.61	14.56	24.41
ApproxGramJacobian two-way comm. with randomized-SVD	10.15	1.38	2.04	1.76	3.83
ApproxGramJacobian two-way comm. with random masking	51.25	10.26	17.40	12.98	22.97
ApproxGramJacobian two -way $comm$. with top- k sparsification	17.15	0.15	0.04	0.15	4.38
ApproxGramJacobian one-way comm. with randomized-SVD	10.87	1.68	1.79	1.96	4.08
ApproxGramJacobian one-way comm. with random masking	76.54	4.23	6.29	5.70	23.19
$ \begin{array}{c} {\tt ApproxGramJacobian} \ one\text{-}way \ comm. \\ {\tt with} \ top\text{-}k \ sparsification \end{array} $	24.13	0.04	0.01	0.10	6.07

Error (in %) =
$$\frac{\|\widetilde{G} - G\|}{\|\widetilde{G}\|}$$
, where $\widetilde{G} = \left(\frac{1}{N} \sum_{i} \widetilde{H}_{i}\right)^{\top} \left(\frac{1}{N} \sum_{i} \widetilde{H}_{i}\right)$

Updating aggregation weights

At the server - need to solve

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \underbrace{\left[\nabla L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]}_{\triangleq \nabla \mathbf{L}(\mathbf{x}) \in \mathbb{R}^{d \times M}} \mathbf{w} \right\|_2^2 \quad \text{s.t.} \quad \mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$$
$$\approx \operatorname{Proj}_{\mathcal{S}_M} \left(\mathbf{w} - \eta_w \nabla \mathbf{L}(\mathbf{x})^\top \nabla \mathbf{L}(\mathbf{x}) \mathbf{w} \right) \qquad \qquad \text{(Single PGD step)}$$

 \mathcal{S}_M is the probability simplex in \mathbb{R}^M

Updating aggregation weights

At the server - need to solve

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left\| \underbrace{\left[\nabla L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]}_{\triangleq \nabla \mathbf{L}(\mathbf{x}) \in \mathbb{R}^{d \times M}} \mathbf{w} \right\|_2^2 \quad \text{s.t.} \quad \mathbf{w} \ge \mathbf{0}, \ \mathbf{w}^\top \mathbf{1} = 1$$
$$\approx \operatorname{Proj}_{\mathcal{S}_M} \left(\mathbf{w} - \eta_w \nabla \mathbf{L}(\mathbf{x})^\top \nabla \mathbf{L}(\mathbf{x}) \mathbf{w} \right) \qquad \qquad \text{(Single PGD step)}$$

 \mathcal{S}_M is the probability simplex in \mathbb{R}^M

Server uses $G \approx \nabla \mathbf{L}(\mathbf{x})^{\top} \nabla \mathbf{L}(\mathbf{x})$ (obtained with $\mathcal{O}(d)$ communication) to compute

$$\mathbf{w} \leftarrow \mathsf{Proj}_{\mathcal{S}_M} (\mathbf{w} - \eta_w G \mathbf{w})$$

Server and client steps

• Select client subset C - send x

- ullet Select client subset ${\mathcal C}$ send ${\mathbf x}$
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$

- Select client subset C send x
- Compute $\nabla \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} \mathbf{v} \mathbf{v}^{(i)}$

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|C|} \sum_i H_i\right)^{\top} \left(\frac{1}{|C|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} = \mathbf{x} \mathbf{x}^{(i)}$
- Update the global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \nabla^{(i)}$

Server and client steps

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} = \mathbf{x} \mathbf{x}^{(i)}$
- Update the global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \nabla^{(i)}$

Server Communication Cost

• FMGDA - d

Server and client steps

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} = \mathbf{x} \mathbf{x}^{(i)}$
- Update the global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \nabla^{(i)}$

Server Communication Cost

- FMGDA d
- FedCMOO $d + M \approx d$

Server and client steps

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} = \mathbf{x} \mathbf{x}^{(i)}$
- Update the global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \nabla^{(i)}$

Server Communication Cost

- FMGDA d
- FedCMOO $d + M \approx d$

Server and client steps

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} = \mathbf{x} \mathbf{x}^{(i)}$
- Update the global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \nabla^{(i)}$

Server Communication Cost

- FMGDA d
- FedCMOO $d + M \approx d$

Per-client Communication Cost

• FMGDA - Md

Server and client steps

- Select client subset C send x
- Compute $\widetilde{\nabla} \mathbf{L}_i(\mathbf{x})$ and send $\widehat{H}_i = \mathcal{Q}(\widetilde{\nabla} \mathbf{L}_i(\mathbf{x}))$ to server
- Reshape \widehat{H}_i to H_i . Compute $G \approx \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)^{\top} \left(\frac{1}{|\mathcal{C}|} \sum_i H_i\right)$
- Compute w ← Proj_{SM} (w − η_wGw) and send to clients
- Start with \mathbf{x} , multiple (τ) local updates $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} \eta_c \sum_k w_k \widetilde{\nabla} L_k^{(i)}(\mathbf{x}^{(i)})$
- Return updates to server $\nabla^{(i)} = \mathbf{x} \mathbf{x}^{(i)}$
- Update the global model $\mathbf{x} \leftarrow \mathbf{x} \eta_s \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} \nabla^{(i)}$

Server Communication Cost

- FMGDA d
- FedCMOO $d + M \approx d$

Per-client Communication Cost

- FMGDA Md
- FedCMOO 2d

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{ where } \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

Problem

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

• Smoothness: $\|\nabla L_k^{(i)}(\mathbf{x}) - \nabla L_k^{(i)}(\mathbf{y})\| \le L\|\mathbf{x} - \mathbf{y}\|$

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

- Smoothness: $\|\nabla L_k^{(i)}(\mathbf{x}) \nabla L_k^{(i)}(\mathbf{y})\| \le L\|\mathbf{x} \mathbf{y}\|$
- Unbiased stochastic gradients with bounded variance

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

- Smoothness: $\|\nabla L_k^{(i)}(\mathbf{x}) \nabla L_k^{(i)}(\mathbf{y})\| \le L\|\mathbf{x} \mathbf{y}\|$
- Unbiased stochastic gradients with bounded variance
- Bounded Gradients: $\|\nabla L_k^{(i)}(\mathbf{x})\| \leq B$

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

- Smoothness: $\|\nabla L_k^{(i)}(\mathbf{x}) \nabla L_k^{(i)}(\mathbf{y})\| \le L\|\mathbf{x} \mathbf{y}\|$
- Unbiased stochastic gradients with bounded variance
- Bounded Gradients: $\|\nabla L_k^{(i)}(\mathbf{x})\| \leq B$
- Unbiased Compression: $\mathbb{E}[Q(\mathbf{x})] = \mathbf{x}$ and $\mathbb{E}\|Q(\mathbf{x}) \mathbf{x}\|^2 \le q\|\mathbf{x}\|^2$, for some q > 0

Theoretical Result: Convergence

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

$$\min_{[T]} \| \sum_{k} w_{k} \nabla L_{k}(\mathbf{x}) \|^{2} \leq \underbrace{\mathcal{O}\left(\frac{1}{T\eta_{s}\eta_{c}\tau} + \frac{L\eta_{s}\eta_{c}}{n}\right)}_{\text{Centralized Optimization Error}} + \underbrace{\underbrace{\mathcal{O}\left(L\eta_{s}\eta_{c}\tau B\right)}_{\text{Partial Participation Error}}}_{\text{Error}} + \underbrace{\underbrace{\mathcal{O}\left(L\eta_{s}\eta_{c}\tau B\right)}_{\text{Centralized Loss}}}_{\text{MOO Weight Error}} + \underbrace{\mathcal{O}\left(L\eta_{s}\eta_{c}\tau B\right)}_{\text{Partial Participation}}$$

Theoretical Result: Convergence

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

$$\min_{[T]} \| \sum_{k} w_{k} \nabla L_{k}(\mathbf{x}) \|^{2} \leq \underbrace{\mathcal{O}\left(\frac{1}{T\eta_{s}\eta_{c}\tau} + \frac{L\eta_{s}\eta_{c}}{n}\right)}_{\text{Centralized Optimization Error}} + \underbrace{\mathcal{O}\left(L\eta_{s}\eta_{c}\tau B\right)}_{\text{Partial Participation Error}} + \underbrace{\mathcal{O}\left(L\eta_{c}(\tau B + \sqrt{\tau}B)\right)}_{\text{Local Drift Error}} + \underbrace{\mathcal{O}\left(L\eta_{c}(\tau B + \sqrt{\tau}B)\right)}_{\text{MOO Weight Error}} + \underbrace{\mathcal{O}\left(L\eta_{c}(\tau B + \sqrt{\tau}B)\right)}_{\text{MOO Weight Error}} + \underbrace{\mathcal{O}\left(L\eta_{c}(\tau B + \sqrt{\tau}B)\right)}_{\text{MOO Weight Error}} + \underbrace{\mathcal{O}\left(L\eta_{c}(\tau B + \sqrt{\tau}B)\right)}_{\text{Example of the properties}} + \underbrace{\mathcal{O}\left(L\eta_{c}(\tau B + \sqrt{\tau}B)\right)}$$

• With appropriate choice of parameters η_s, η_c, η_w

$$\min_{t \in [T]} \| \sum_{k} w_k(t) \nabla L_k(\mathbf{x}(t)) \|^2 \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

To achieve $\min_{t \in [T]} \| \sum_k w_k(t) \nabla L_k(\mathbf{x}(t)) \|^2 \le \epsilon$

[†]Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

^{*}Xiao, et al. "Direction-oriented multi-objective learning: Simple and provable stochastic algorithms." *NeurIPS'23*.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

To achieve $\min_{t \in [T]} \| \sum_k w_k(t) \nabla L_k(\mathbf{x}(t)) \|^2 \le \epsilon$

Work	Federated	Complexity
SDMGrad*	X	$\mathcal{O}\left(\frac{M^2}{\epsilon^2}\right)$

†Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

^{*}Xiao, et al. "Direction-oriented multi-objective learning: Simple and provable stochastic algorithms." *NeurIPS'23*.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

To achieve $\min_{t \in [T]} \| \sum_k w_k(t) \nabla L_k(\mathbf{x}(t)) \|^2 \le \epsilon$

Work	Federated	Complexity
SDMGrad*	X	$\mathcal{O}\left(\frac{M^2}{\epsilon^2}\right)$
FMGDA [†]	✓	$\mathcal{O}\left(\frac{M^4}{\epsilon^2}\right)$

†Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

^{*}Xiao, et al. "Direction-oriented multi-objective learning: Simple and provable stochastic algorithms." *NeurIPS'23*.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

To achieve $\min_{t \in [T]} \| \sum_k w_k(t) \nabla L_k(\mathbf{x}(t)) \|^2 \le \epsilon$

Work	Federated	Complexity	
SDMGrad*	X	$\mathcal{O}\left(\frac{M^2}{\epsilon^2}\right)$	
FMGDA [†]	✓	$\mathcal{O}\left(\frac{M^4}{\epsilon^2}\right)$	
FedCMOO (Ours)	✓	$\mathcal{O}\left(\frac{M}{\epsilon^2}\right)$	

†Yang, et al. "Federated Multi-Objective Learning," NeurIPS'23.

^{*}Xiao, et al. "Direction-oriented multi-objective learning: Simple and provable stochastic algorithms." *NeurIPS'23*.

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

• With appropriate choice of parameters η_s, η_c, η_w

$$\min_{[T]} \| \sum_{k} w_{k} \nabla L_{k}(\mathbf{x}) \|^{2} \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

• With appropriate choice of parameters $\eta_{s},\eta_{c},\eta_{w}$

$$\min_{[T]} \| \sum_{k} w_{k} \nabla L_{k}(\mathbf{x}) \|^{2} \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

• Bounded Gradients: $\|\nabla L_k^{(i)}(\mathbf{x})\| \leq B$

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N L_k^{(i)}(\mathbf{x})$$

• With appropriate choice of parameters $\eta_{\rm s}, \eta_{\rm c}, \eta_{\rm w}$

$$\min_{[T]} \| \sum_{k} w_{k} \nabla L_{k}(\mathbf{x}) \|^{2} \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

- Bounded Gradients: $\|\nabla L_k^{(i)}(\mathbf{x})\| \leq B$
- No provable benefit of parallelization (i.e., speedup in N)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

• With appropriate choice of parameters $\eta_{\rm s}, \eta_{\rm c}, \eta_{\rm w}$

$$\min_{[T]} \| \sum_{k} w_{k} \nabla L_{k}(\mathbf{x}) \|^{2} \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

- Bounded Gradients: $\|\nabla L_k^{(i)}(\mathbf{x})\| \leq B$
- ullet No provable benefit of parallelization (i.e., speedup in N)
 - Culprit: the inner product terms

Outline

1. Background: MOO and MGDA

2. Federated MGDA and its Drawbacks

3. FedCMOO: Improving Federated MGDA

4. Experiment Results

5. Conclusion

Experiment Results

Experiments: Datasets

MNIST+FMNIST, MultiMNIST, CIFAR10+MNIST

- Constructed by combining two randomly sampled images: one from each dataset
- 2 objectives



Experiments: Datasets

MNIST+FMNIST, MultiMNIST, CIFAR10+MNIST

- Constructed by combining two randomly sampled images: one from each dataset
- 2 objectives



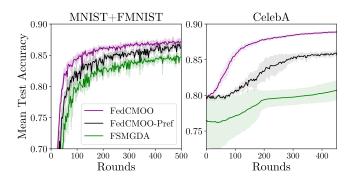
CelebA

- CelebA dataset each image has 40 attributes
- ullet Each attribute gives a binary classification \Rightarrow 40 objectives



Experiments: Mean Test Accuracy

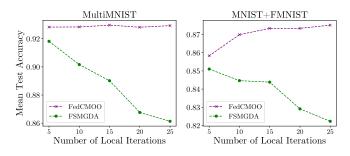
100 clients in total, 10 selected in each round



Test Accuracy averaged over all the objectives

Experiments: Controlling Local Drift

100 clients in total, 10 selected in each round

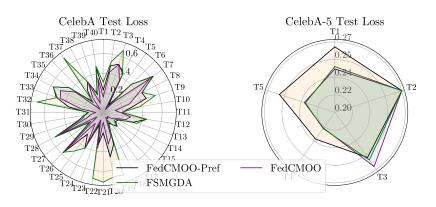


Test Accuracy averaged over all the objectives

• Increasing local iterations (τ) \Rightarrow larger local drift in FSMGDA \Rightarrow worse performance

Experiments: Final Loss Values

• CelebA dataset - 40 objectives



Preferences in MOO

- Users sometimes prefer solutions with specific trade-offs among the different objectives
 - Healthcare
 - Privacy ↔ Utility
- One possible solution: specific ratios for loss values

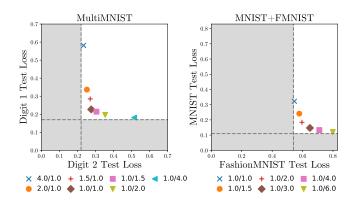
$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{L}(\mathbf{x}) := \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^\top$$
subject to $r_1 L_1(\mathbf{x}) = \dots = r_M L_M(\mathbf{x})$.

- MOO methods, by default, do not let us control the trade-off
- Ensure that normalized scaled losses are almost uniform*

$$\hat{u}_k(\mathbf{r}) = \frac{r_k L_k(\mathbf{x})}{\sum_i r_i L_i(\mathbf{x})} \approx \frac{1}{M}$$

^{*}Mahapatra and Rajan, "Multi-task learning with user preferences: Gradient descent with controlled ascent in pareto optimization." *ICML* '20.

Experiments: Preference-based Pareto Solutions



- Vertical and horizontal dashed lines minimum loss from single-objective training
- Gray regions are infeasible
- More difficult to satisfy preferences with significantly different weights

Conclusion

Some Open Questions

- How to explore the entire Pareto front?
- Generalization guarantees?
- Impact of MOO optimization on MTL generalization§
 - Generalization gap between single-task and multi-task trajectories early into training
 - Standard optimization metrics (sharpness, minimum loss, etc.) do not explain the generalization gaps between single-task and multi-task models

[§]Mueller, et al. "Can Optimization Trajectories Explain Multi-Task Transfer?" rXiv:2408.14677.

Summary: Federated Multi-Objective Optimization (MOO)

$$\min_{\mathbf{x}} \mathbf{L}(\mathbf{x}) := \min_{\mathbf{x}} \left[L_1(\mathbf{x}), \dots, L_M(\mathbf{x}) \right]^{\top} \quad \text{where} \quad L_k(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} L_k^{(i)}(\mathbf{x})$$

M objectives, N clients

- Communication-efficient algorithm FedCMOO
 - Communication cost does not scale with the number of objectives M
- Convergence to Pareto-stationary point
 - Improved sample complexity in terms of M
 - Under weaker assumptions
- Improved empirical performance

Thanks! (pranaysh@iitb.ac.in)