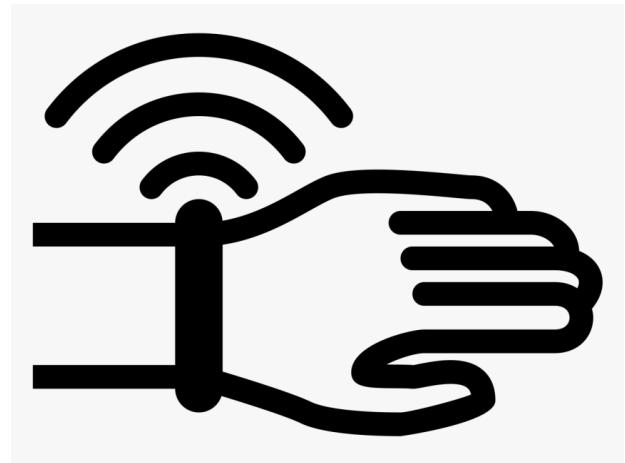
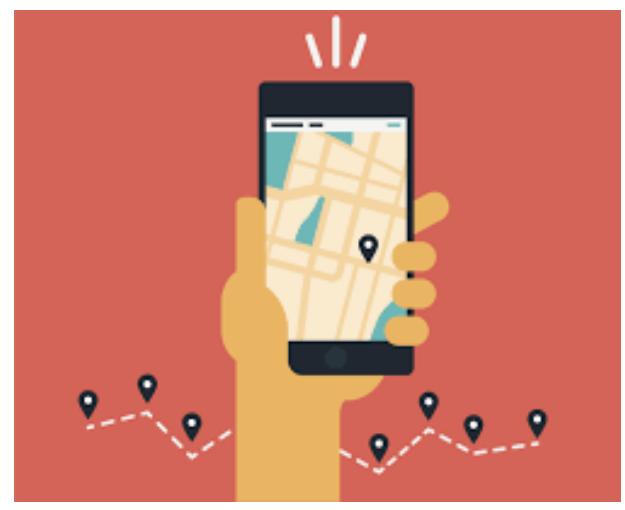
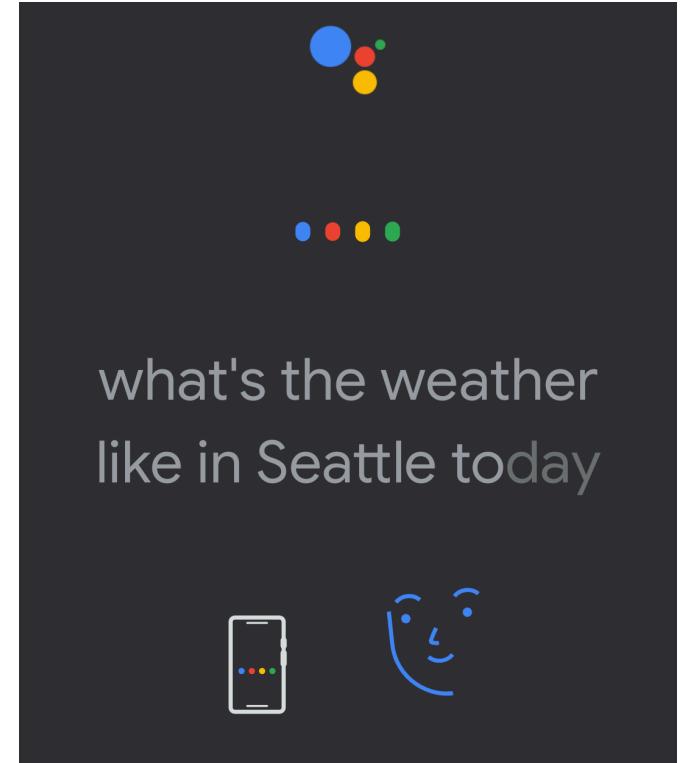
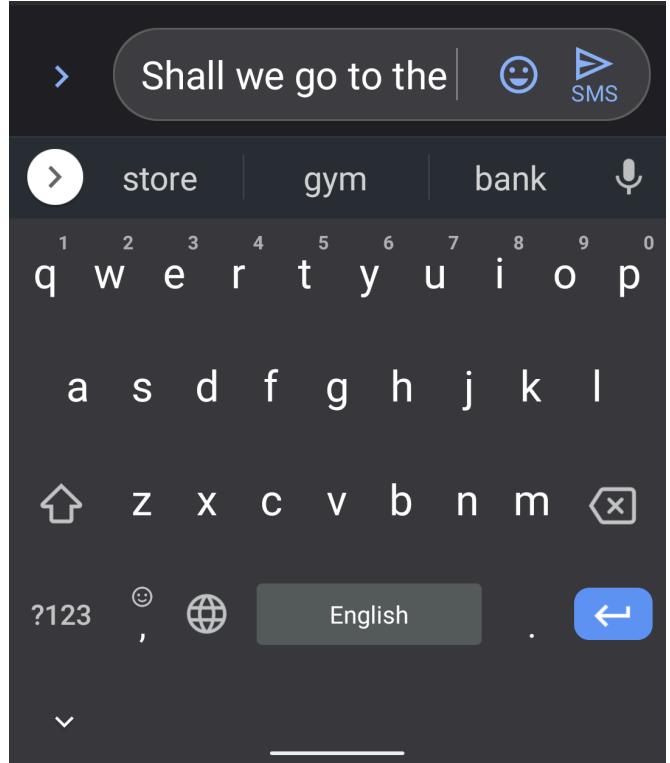


Federated Learning: Robustness, Heterogeneity and Optimization

Krishna Pillutla

June 24, 2022

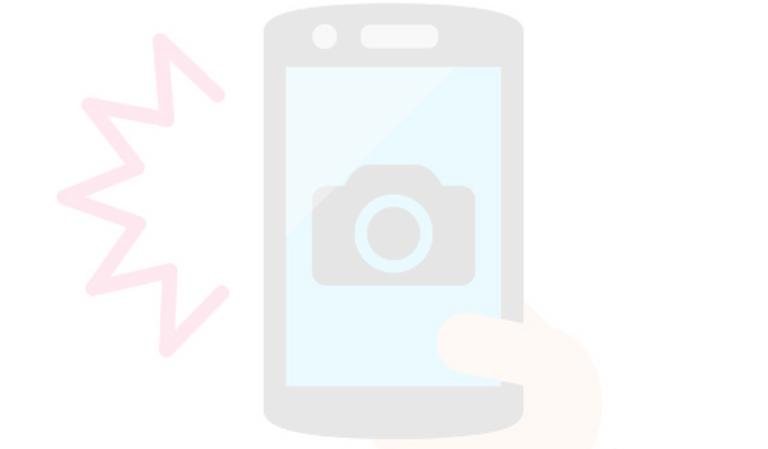
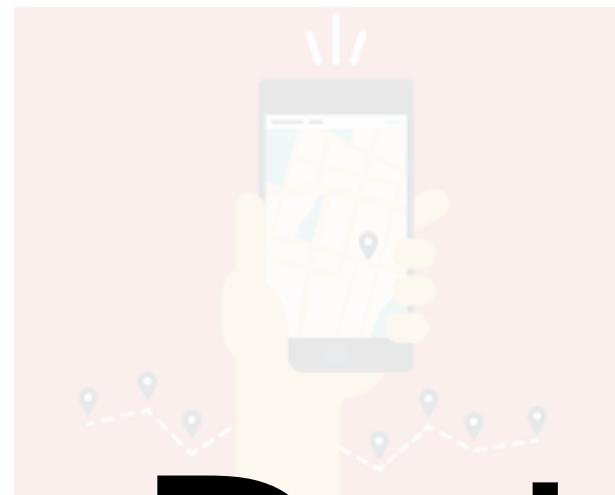
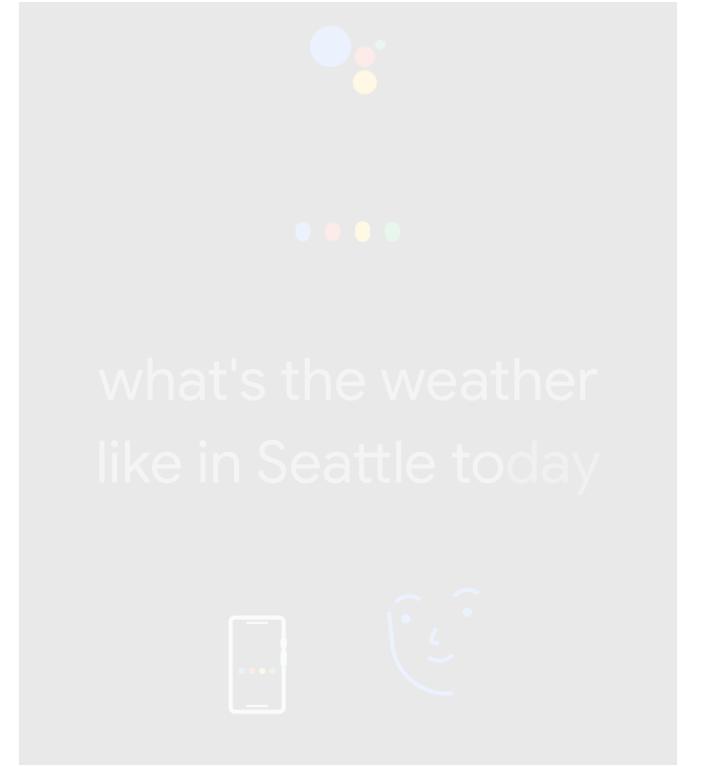
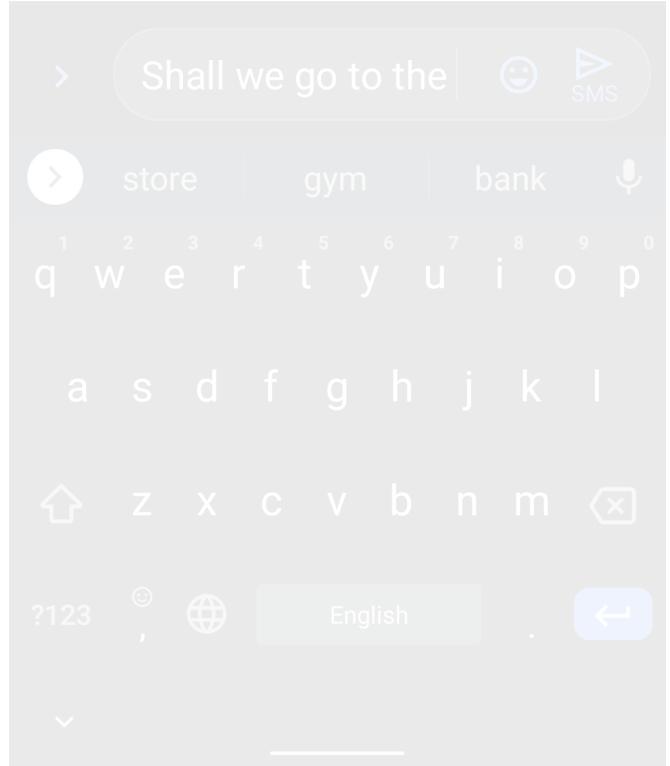
Committee: Zaid Harchaoui (Chair), Sham Kakade (Chair),
Yejin Choi, Lalit Jain (GSR), Kevin Jamieson, Jamie Morgenstern



Rieke et al. NPJ Digit. Med. (2020)



Image Credit: Robotics Business Review

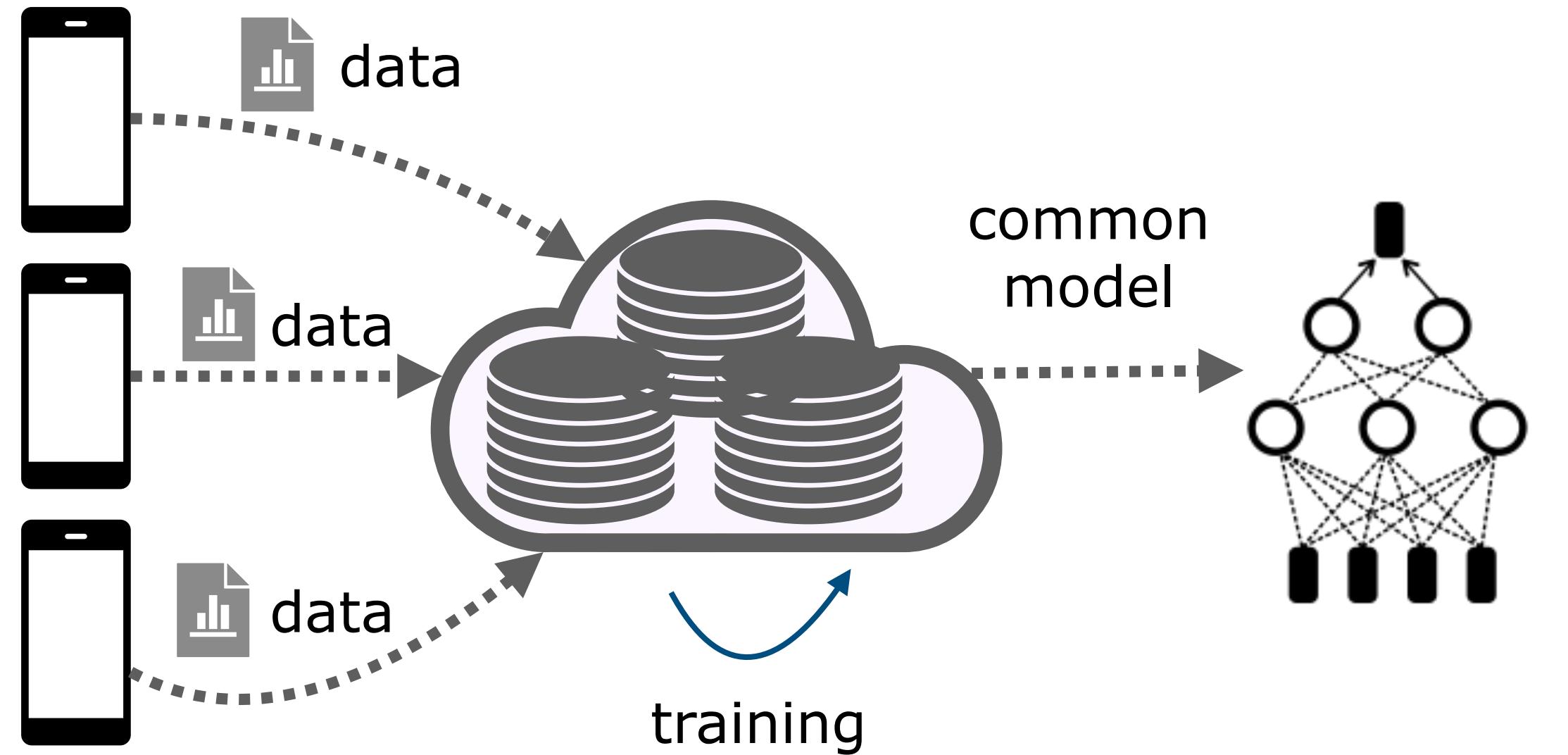


Data is decentralized and private

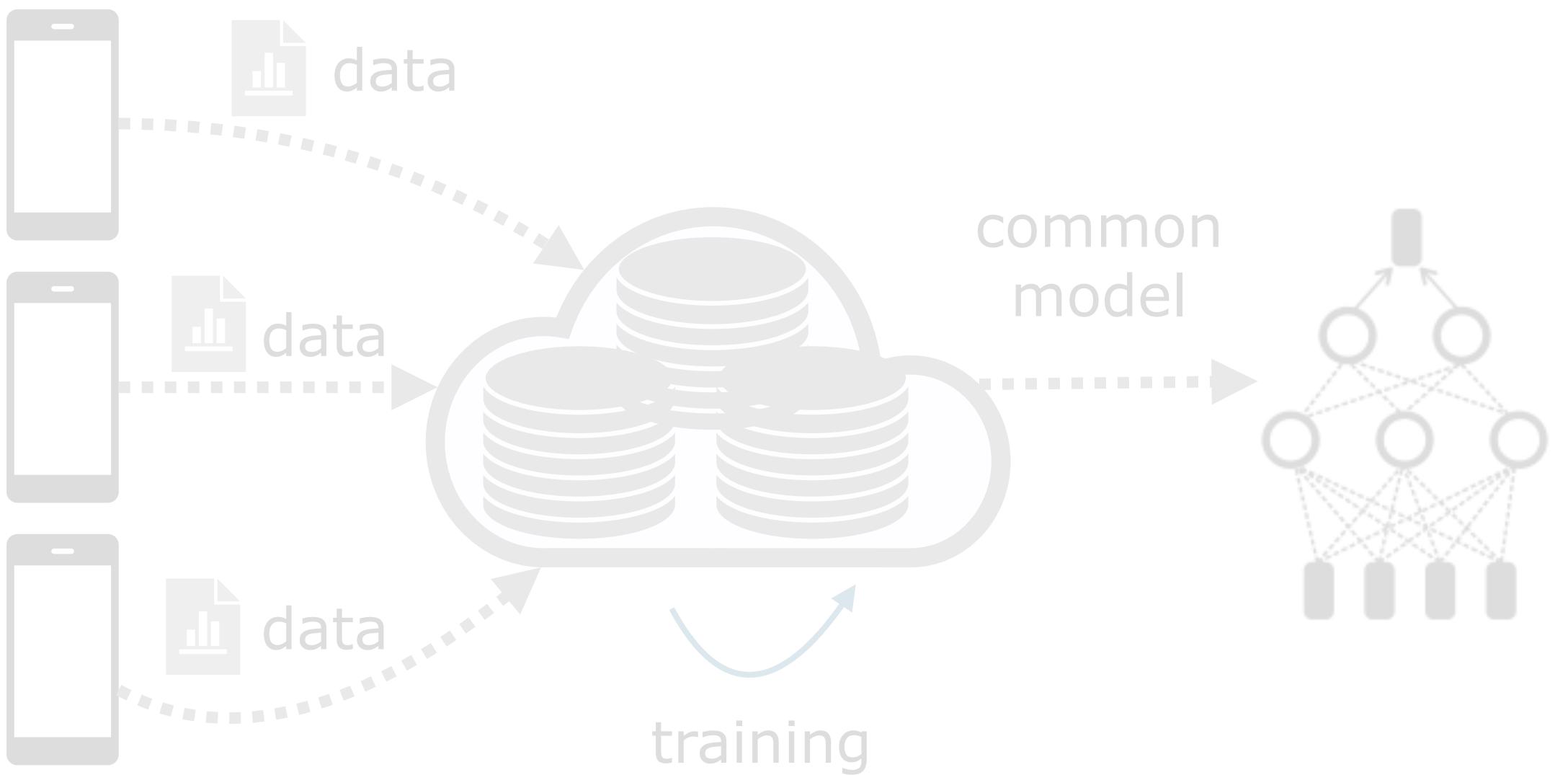


Image Credit: Robotics Business Review

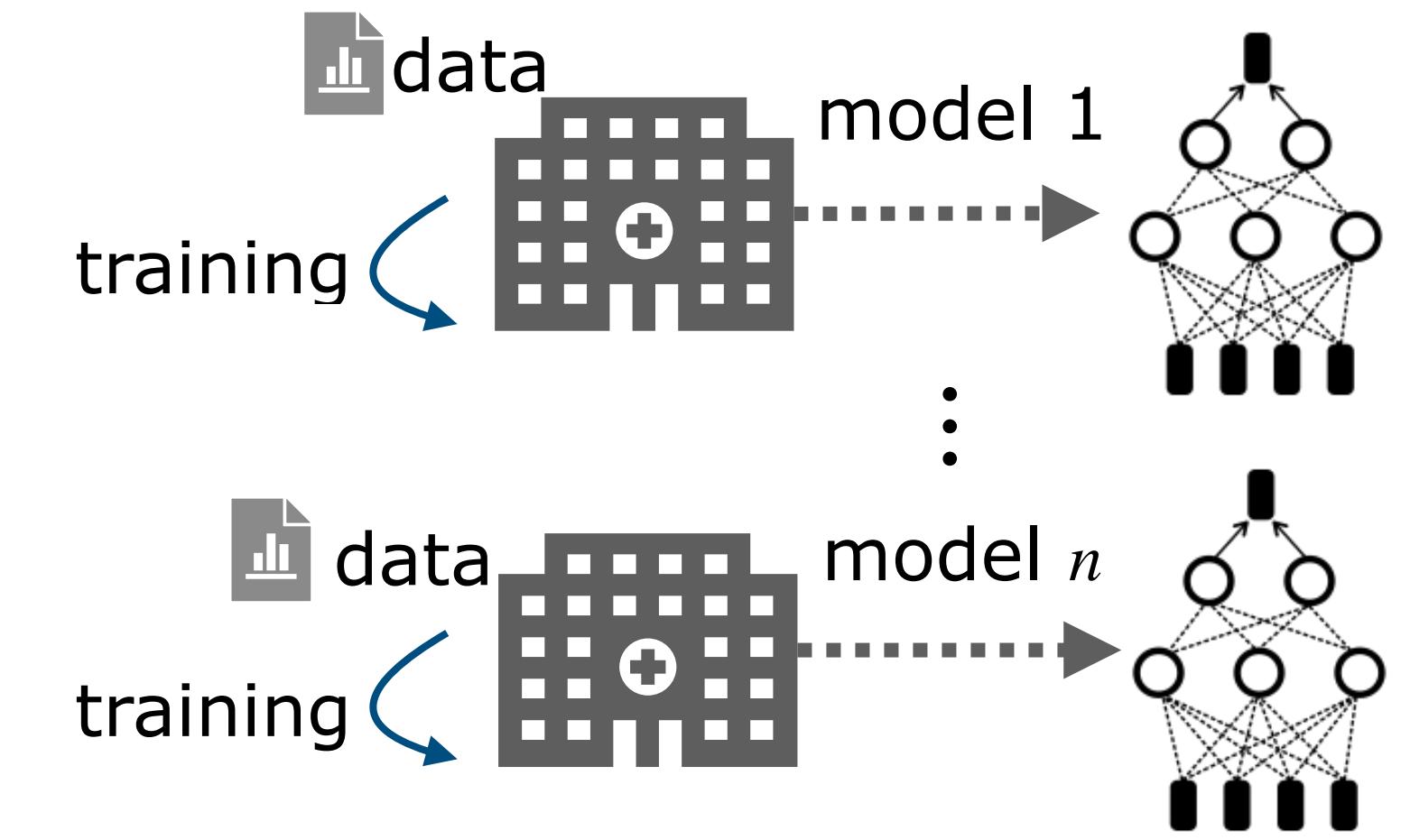
Datacenter



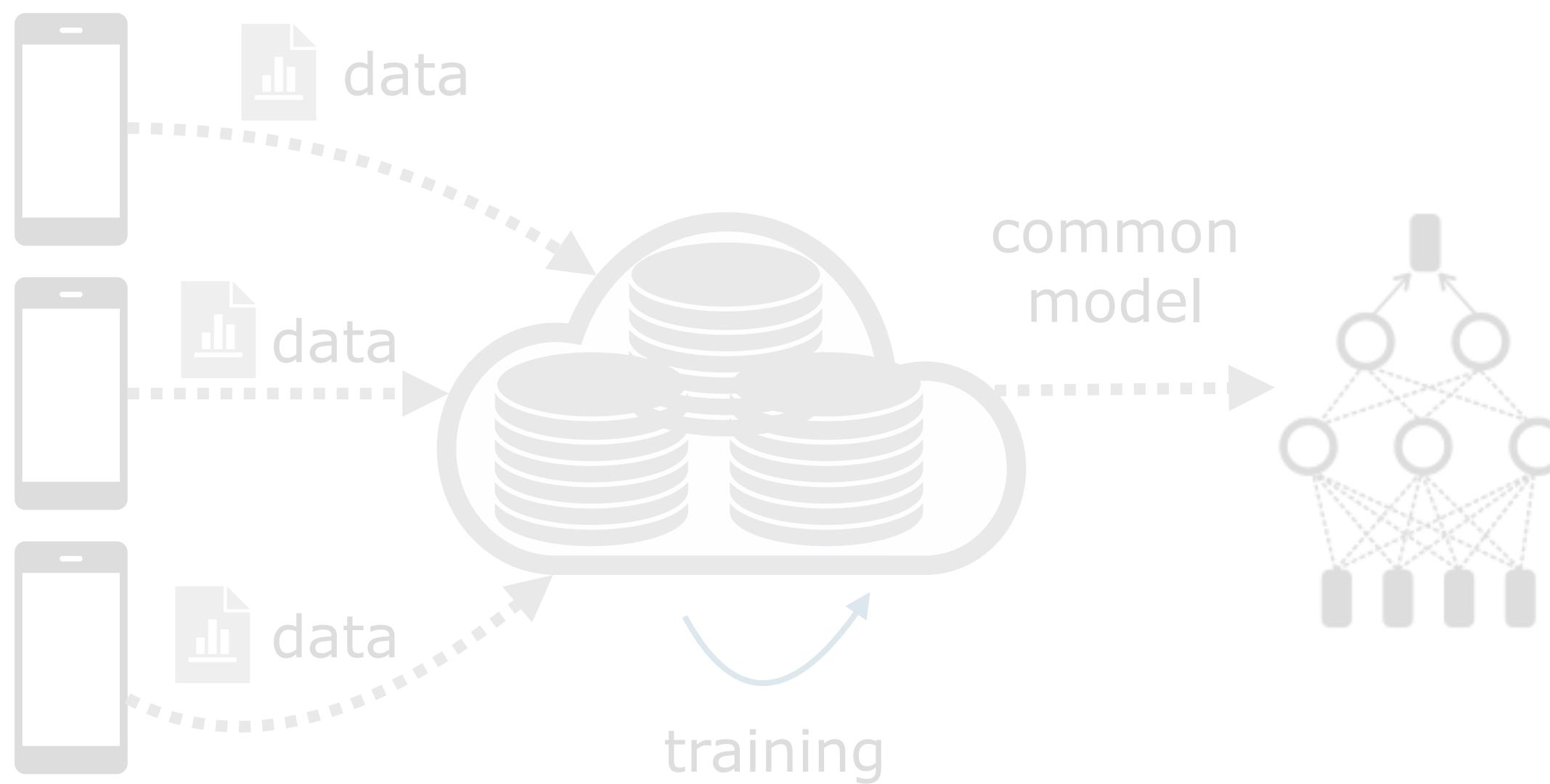
Datacenter



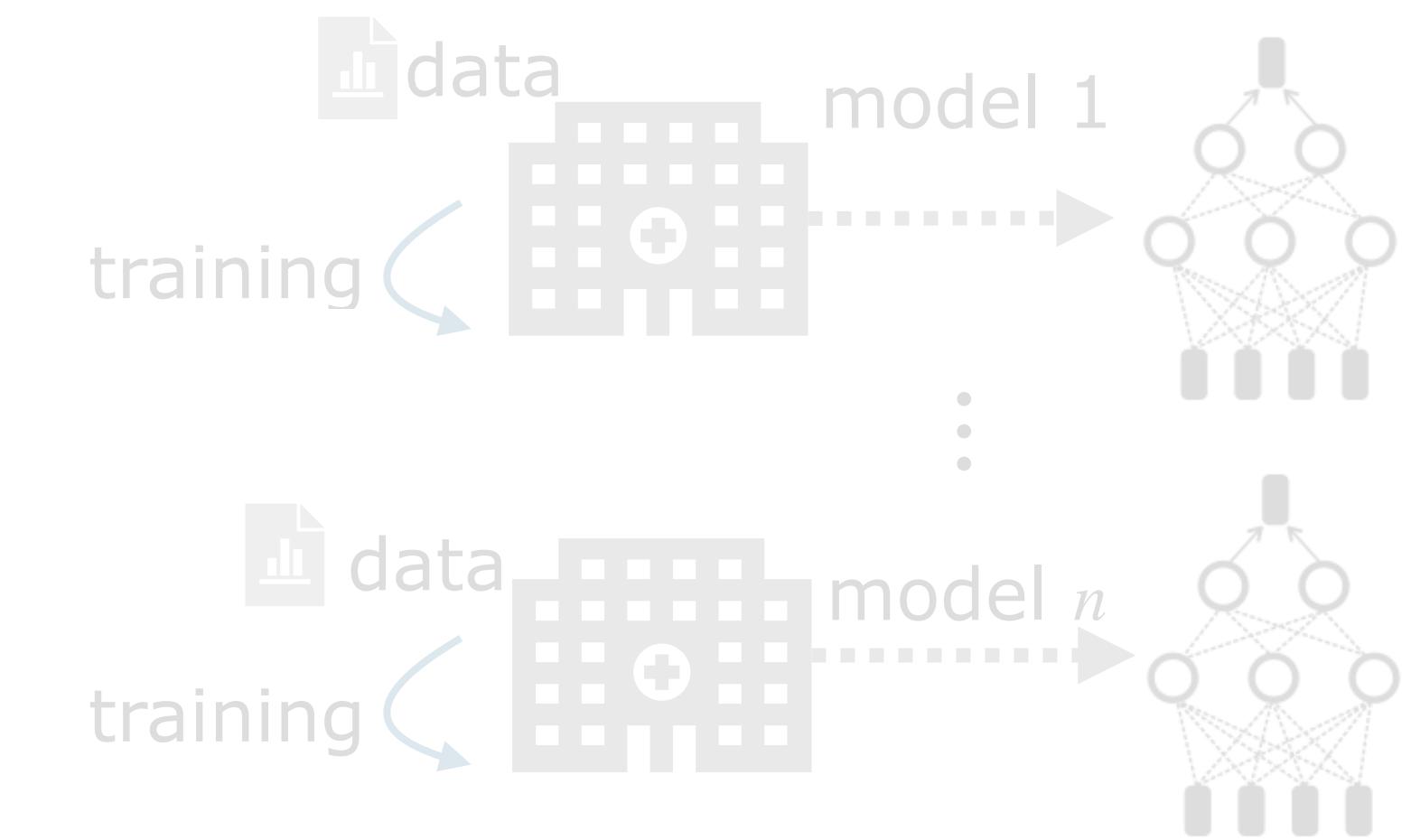
Non-collaborative



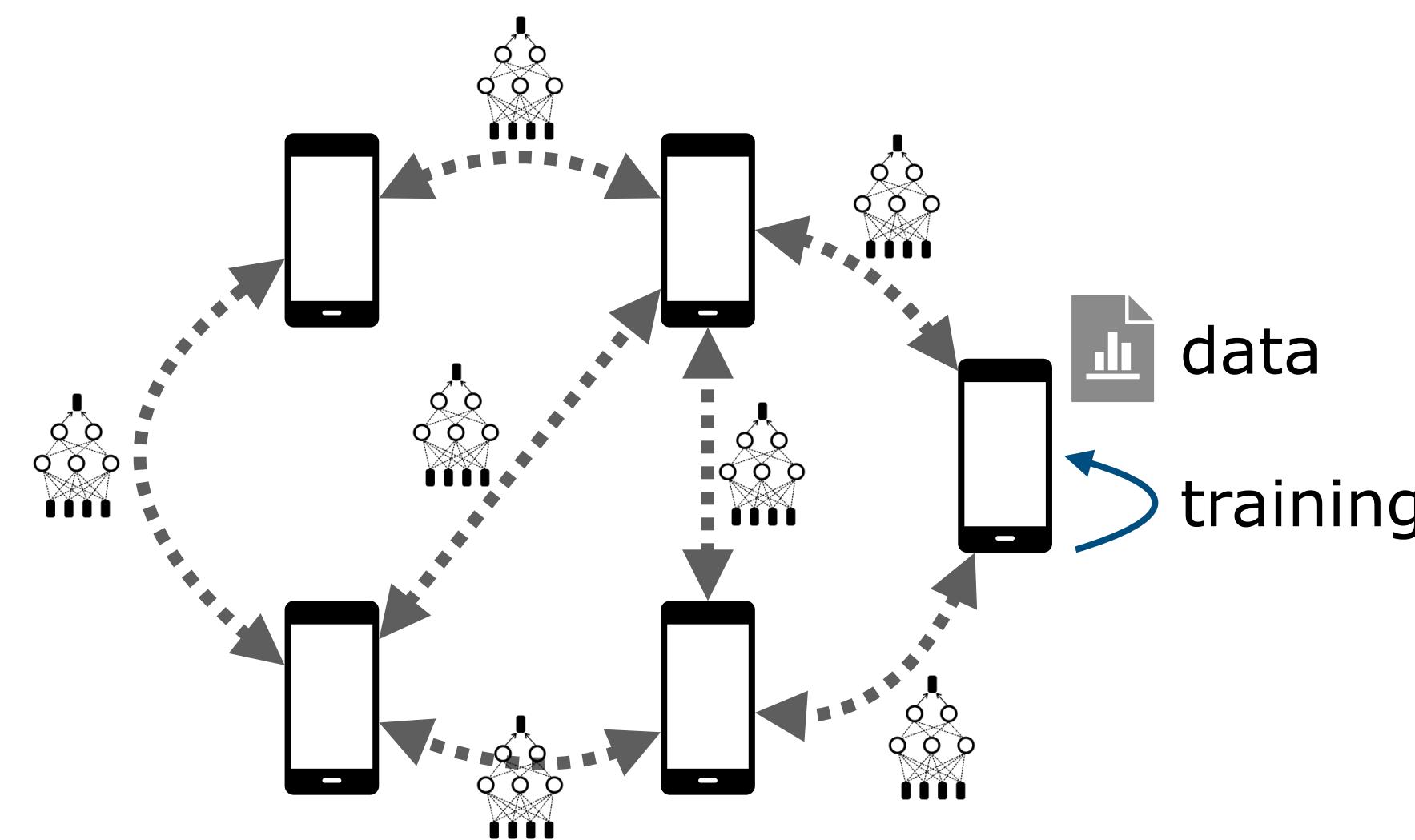
Datacenter



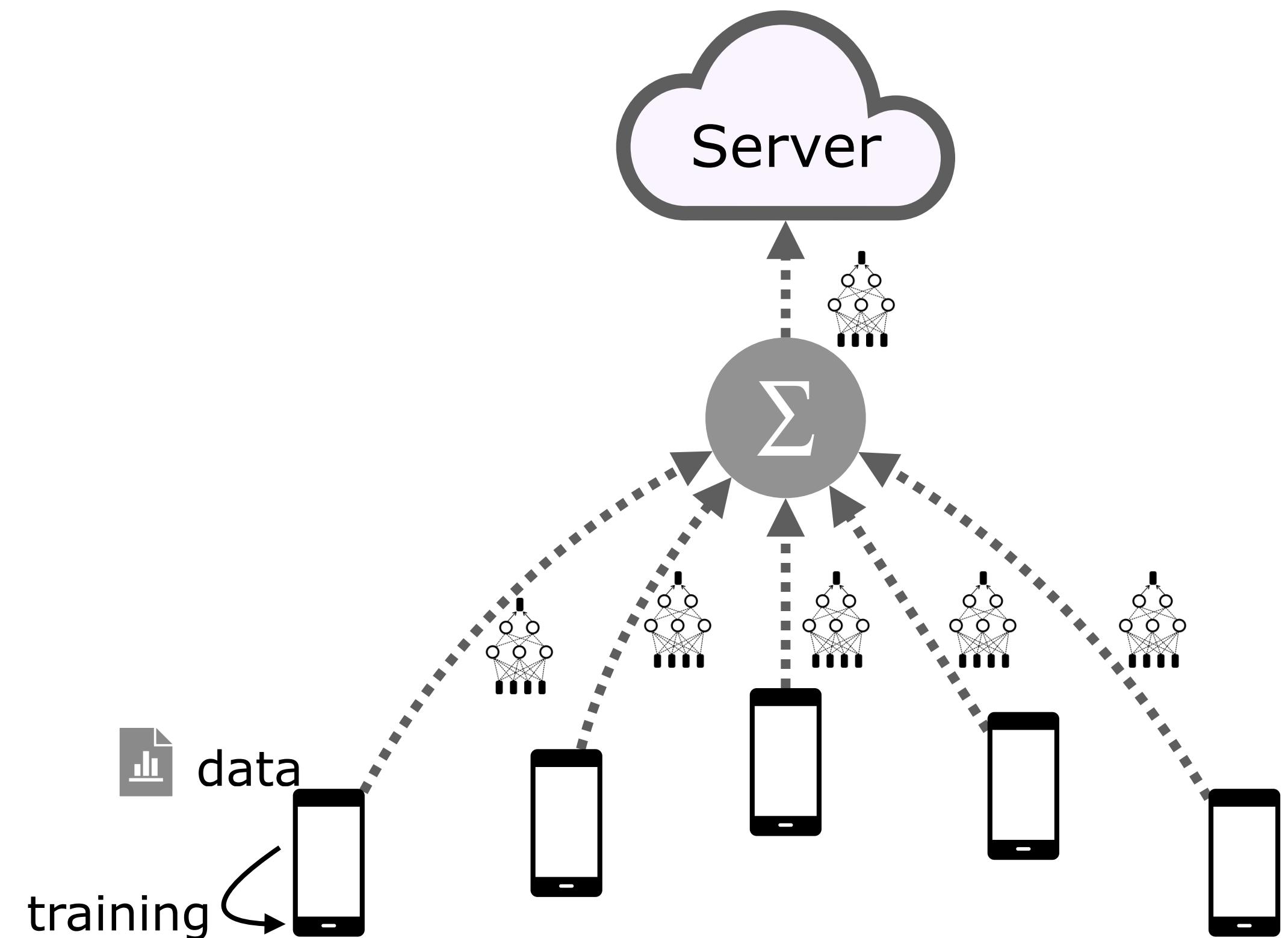
Non-collaborative



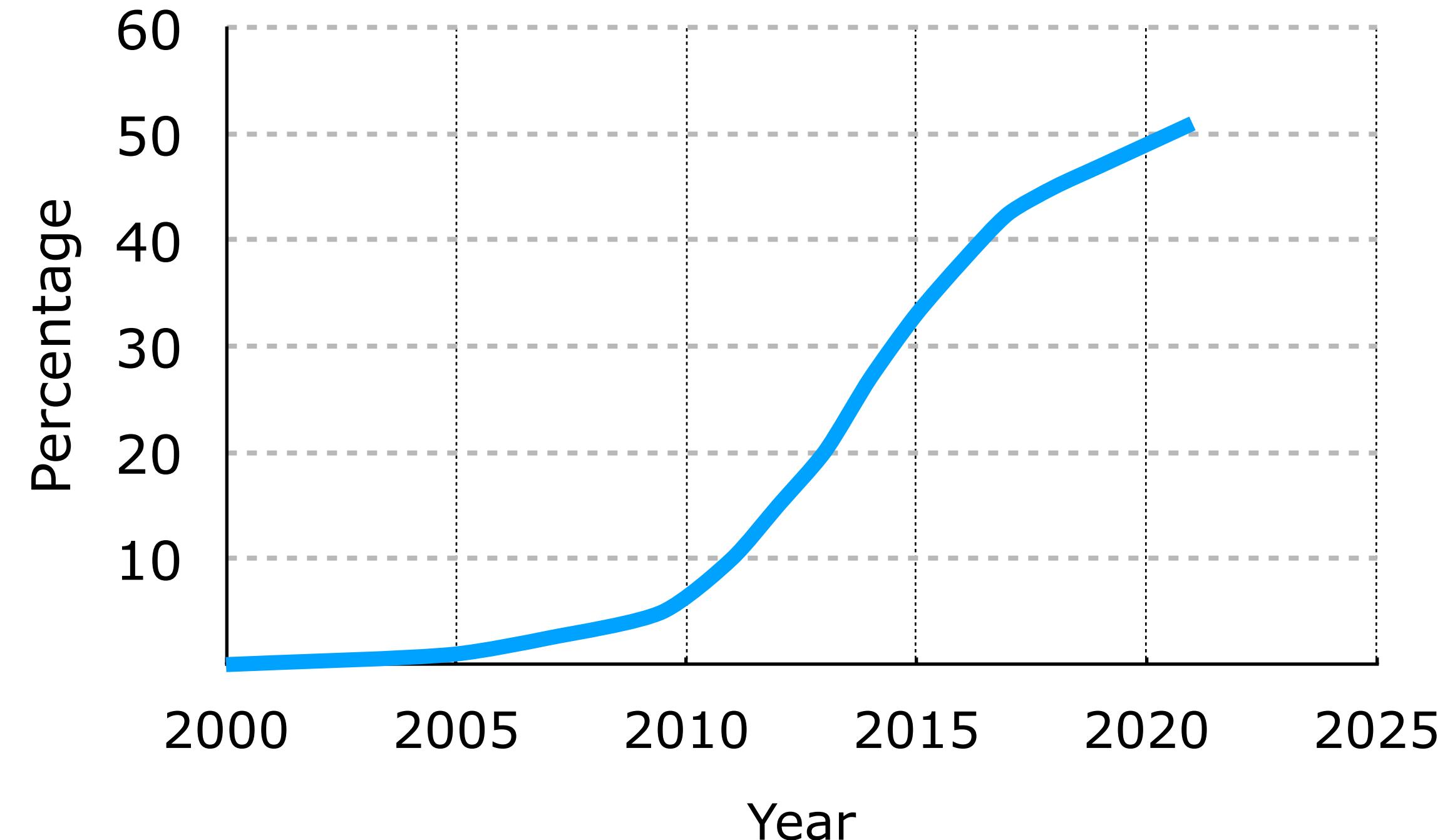
Peer-to-peer



Federated Learning

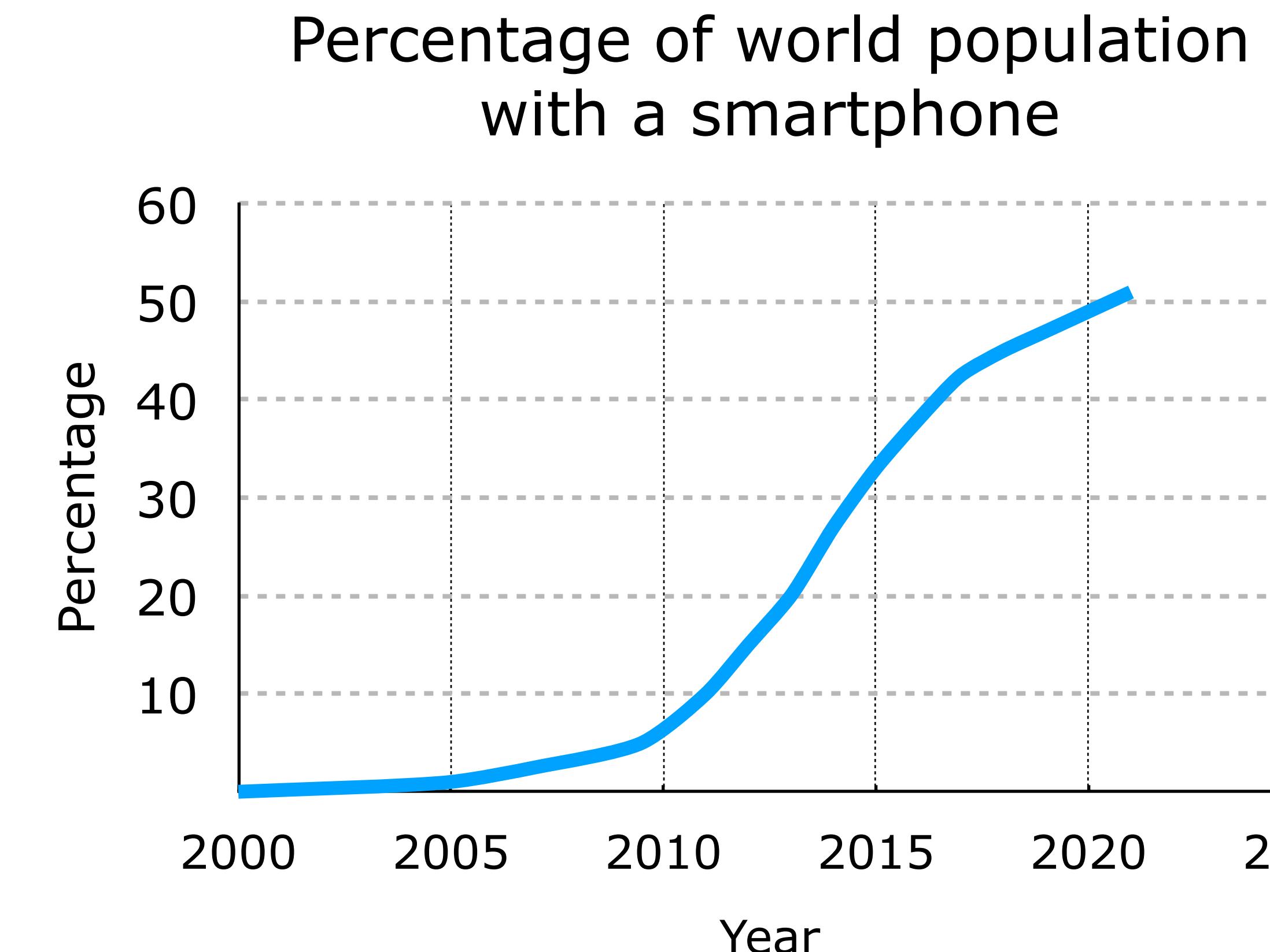
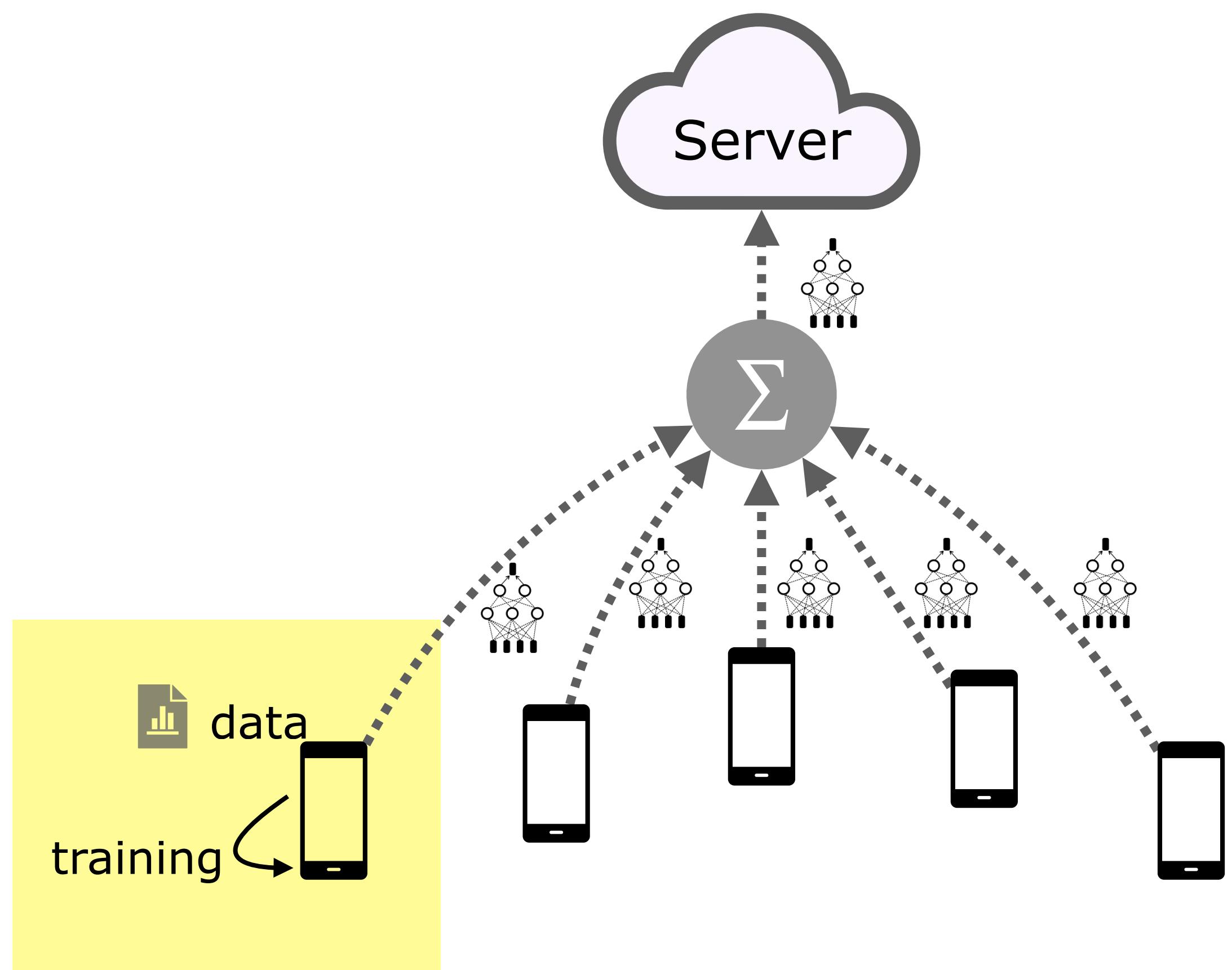


Percentage of world population
with a smartphone



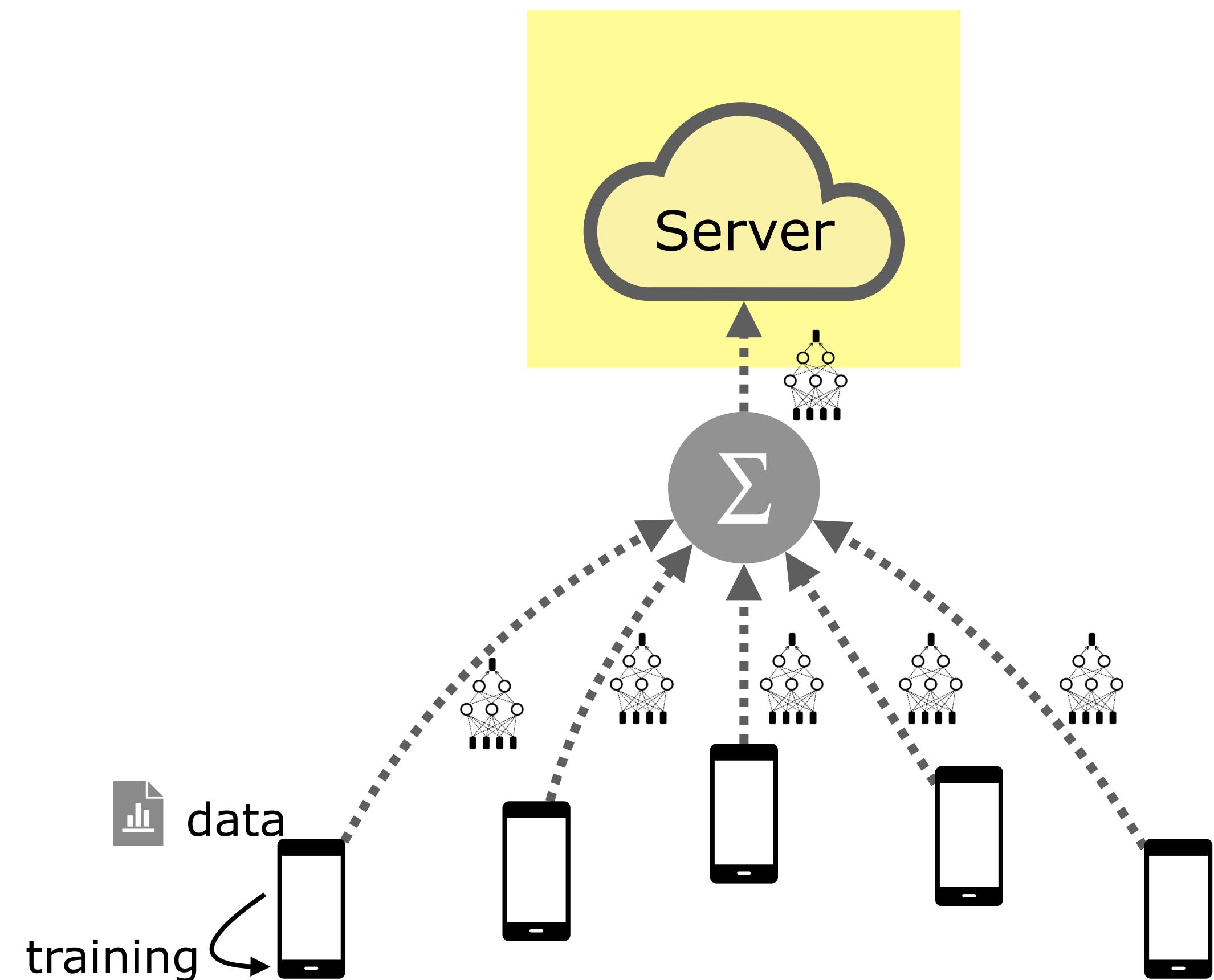
Data Credit: Business Wire

Federated Learning

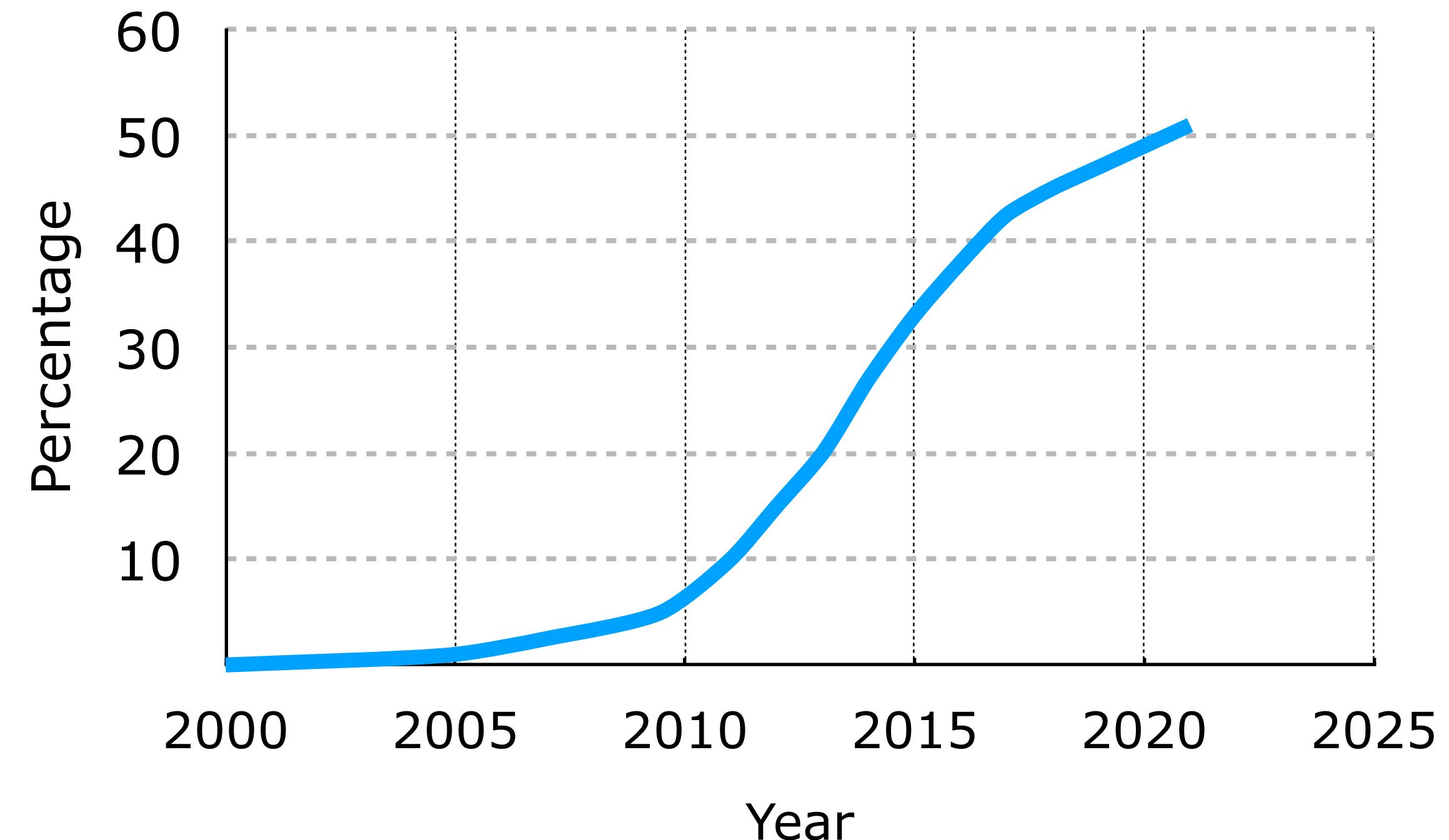


Data Credit: Business Wire

Federated Learning

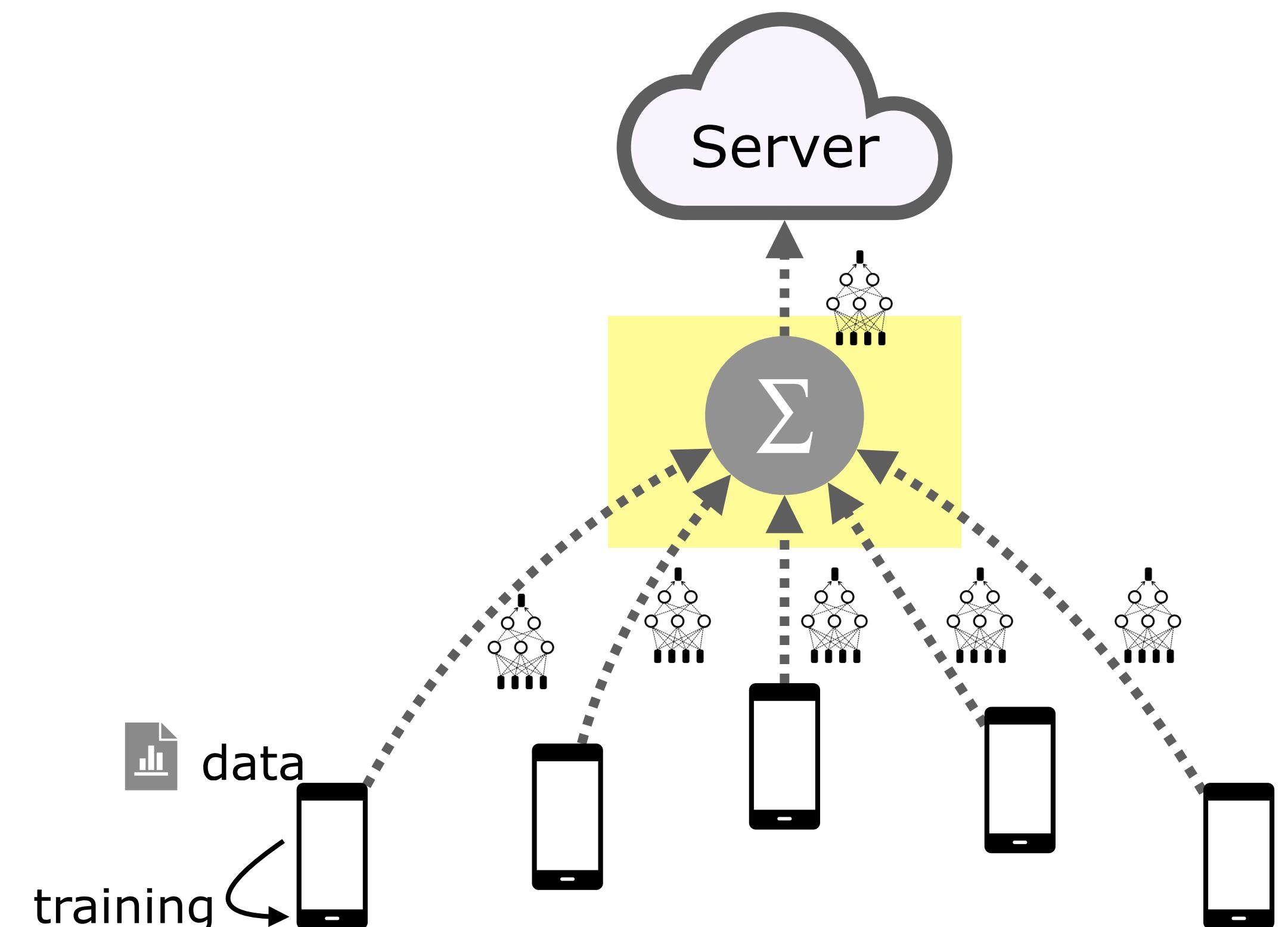


Percentage of world population
with a smartphone

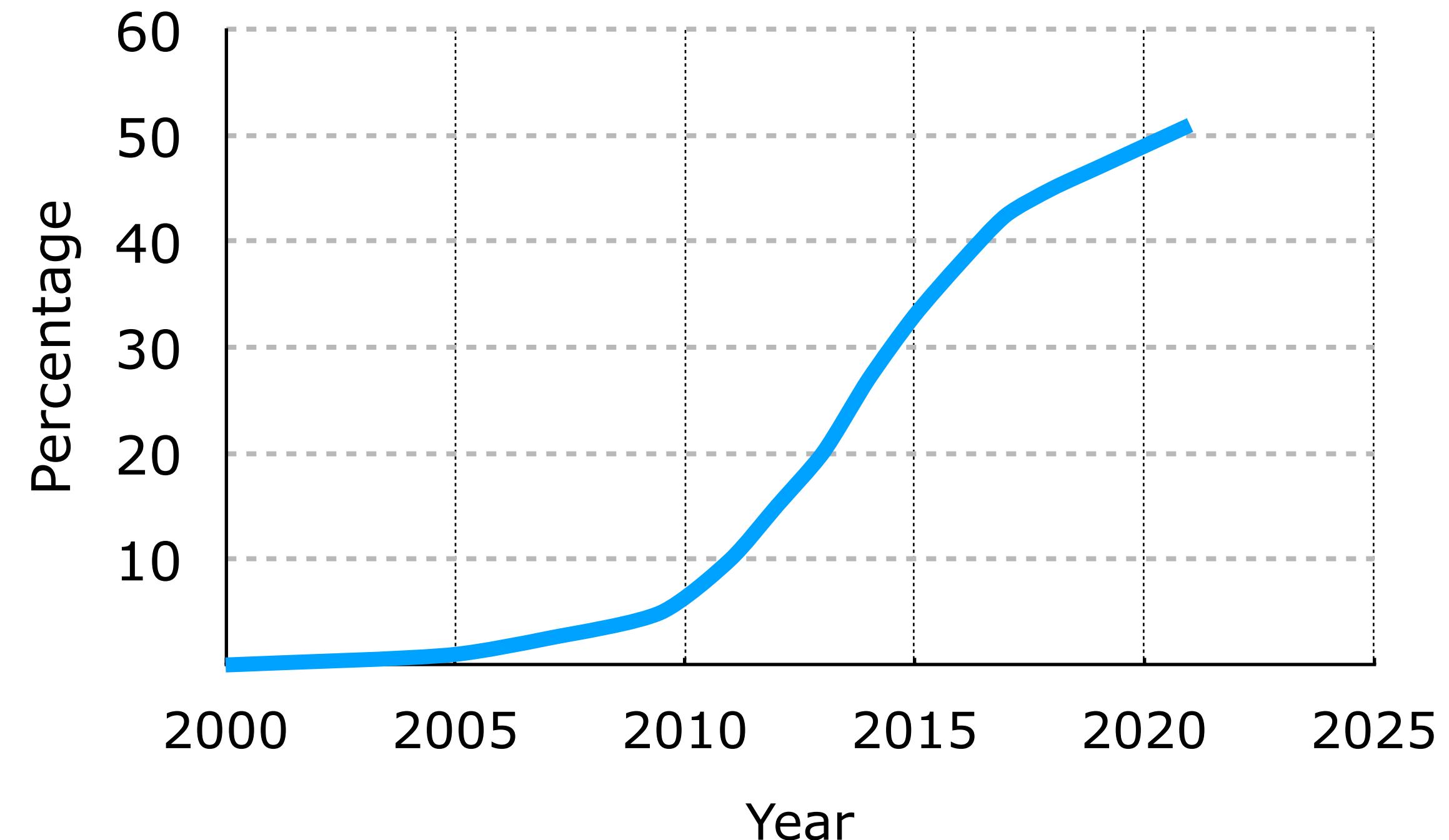


Data Credit: Business Wire

Federated Learning

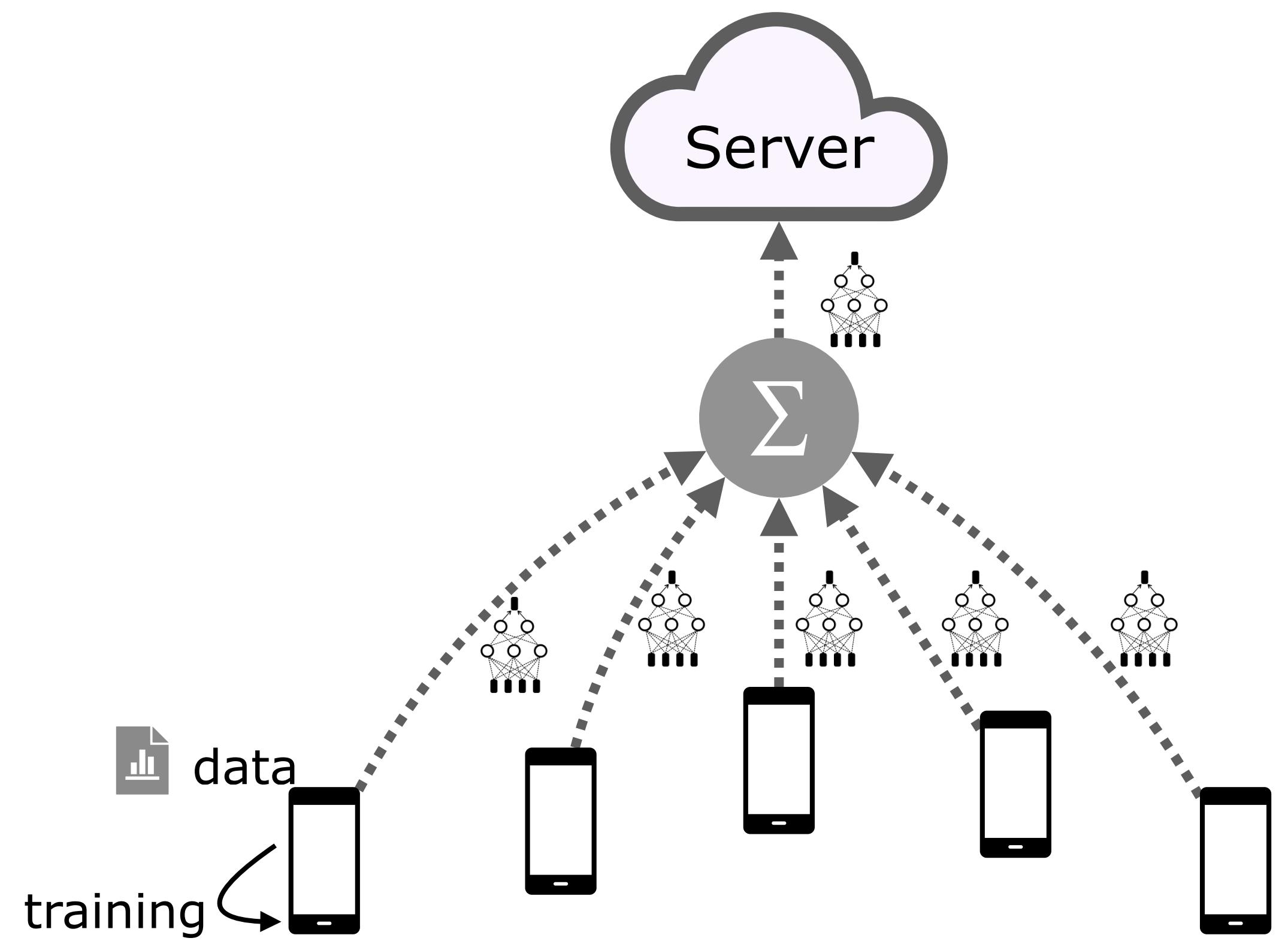


Percentage of world population
with a smartphone

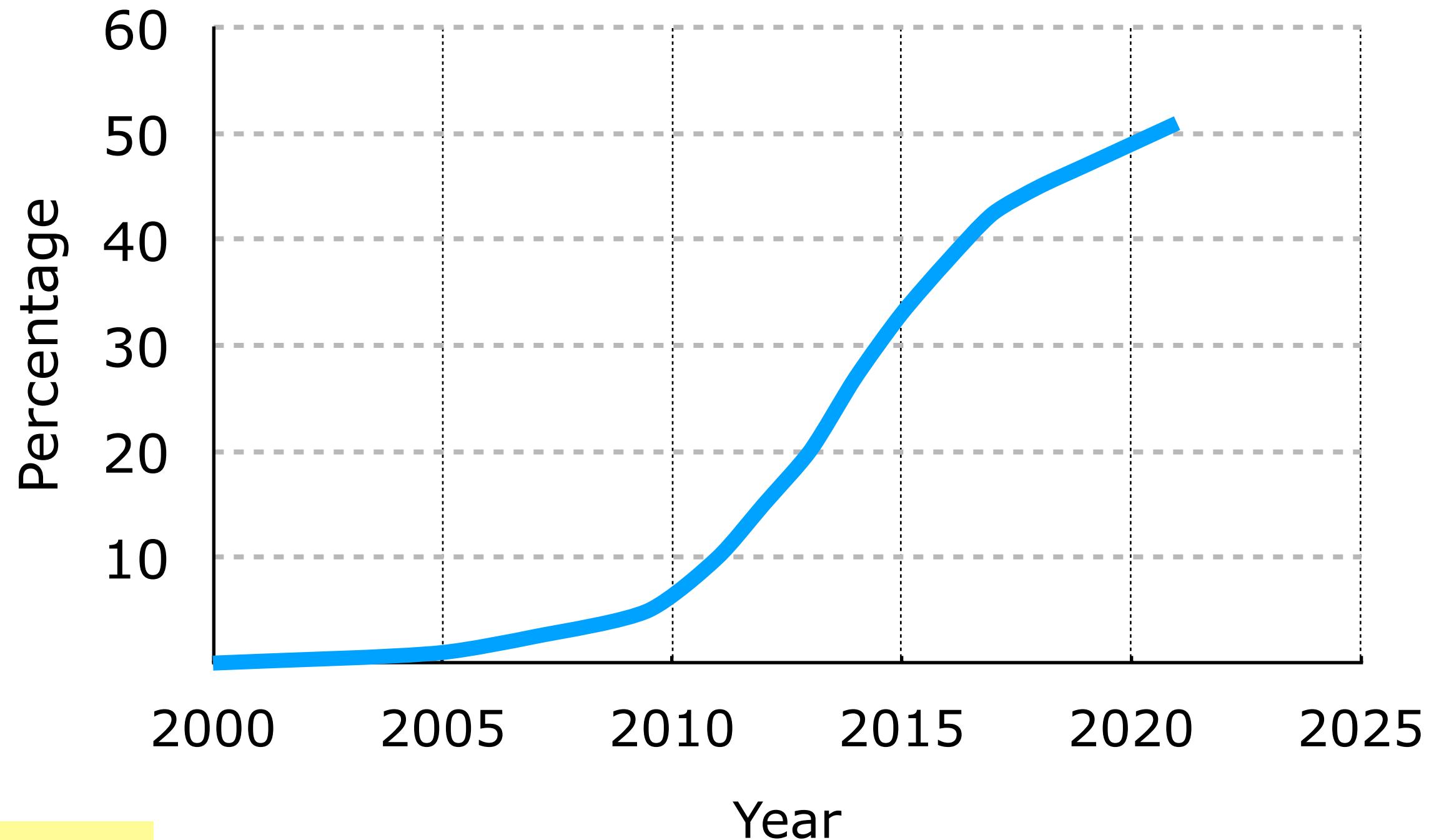


Data Credit: Business Wire

Federated Learning



Percentage of world population with a smartphone



Communication cost > computation cost!

Data Credit: Business Wire

Challenges

- models are deployed on clients with **heterogeneous data**

THE ACCENT GAP

We tested Amazon's Alexa and Google's Home to see how people with accents are getting left behind in the smart-speaker revolution.



Challenges

- models are deployed on clients with **heterogeneous data**
- training is **not robust** to potentially **malicious** clients

Alexa and Siri Can Hear This Hidden Command. You Can't.

Researchers can now send secret audio instructions undetectable to the human ear to Apple's Siri, Amazon's Alexa and Google's Assistant.

By Craig S. Smith

May 10, 2018

1	9	7	9	7
3	6	0	0	0
3	3	1	2	8
8	7	9	7	1
1	6	4	5	5

0	6	3	1	3
9	5	4	1	2
6	9	7	9	3
4	1	2	6	5
9	7	7	6	0

7	1	7	6	7
1	1	3	8	4
4	0	8	0	2
7	4	9	0	0
2	9	8	1	2

1	5	1	1	4
7	3	3	7	6
8	0	5	0	7
7	1	3	7	6
2	6	2	1	4

+

8	0	3	1	7
5	0	7	3	6
1	1	8	3	3
4	3	2	8	5
7	2	3	2	3

Clean
10%
Corrupt

Accuracy

64.3%

-4.2pp

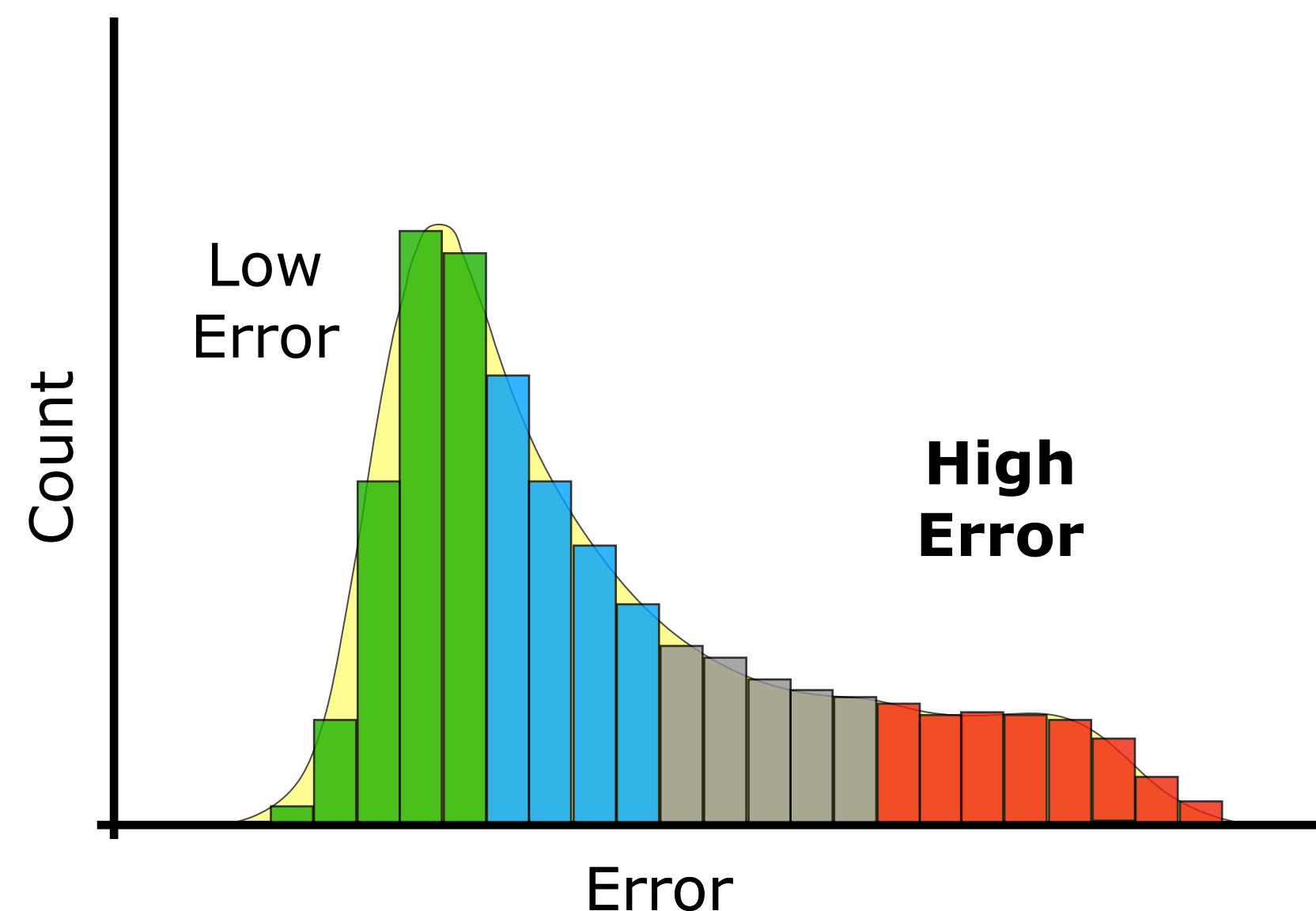
60.1%

Challenges

- models are deployed on clients with **heterogeneous data**
- training is **not robust** to potentially **malicious** clients
- solutions to both these problems are **conflicting**

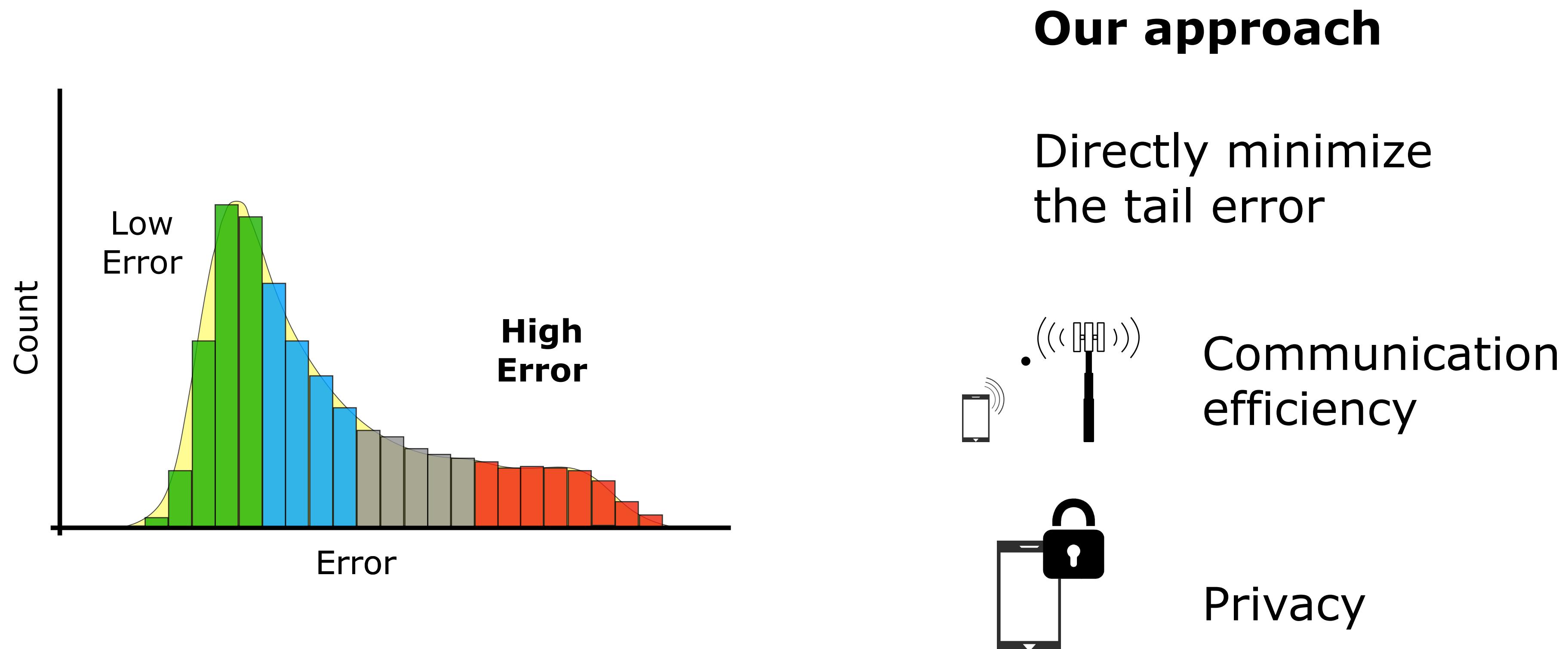
This talk

1. **Heterogeneity-aware objectives** for federated learning
[CISS '21, SVAA '21, Under Review '21]



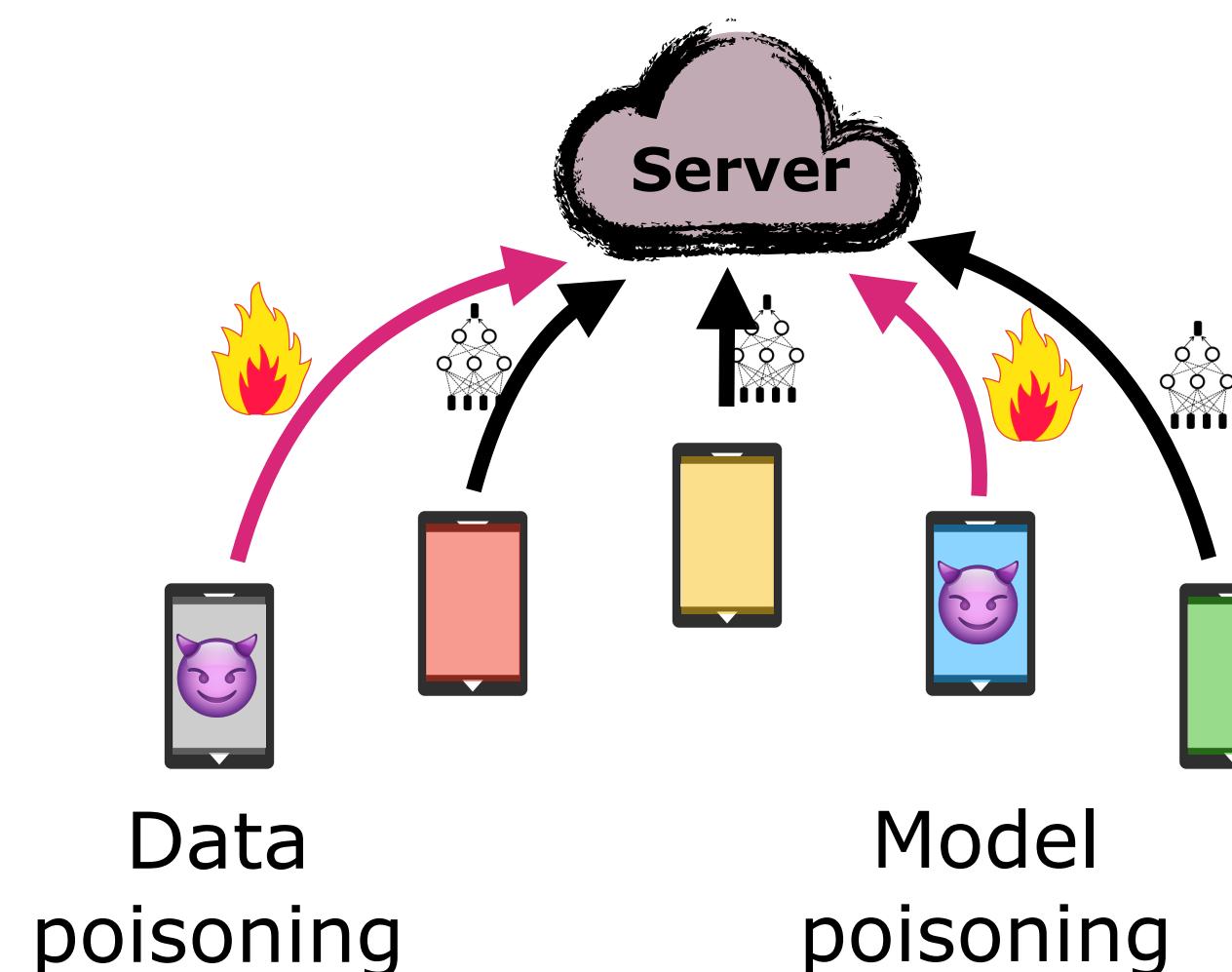
This talk

1. **Heterogeneity-aware objectives** for federated learning
[CISS '21, SVAA '21, Under Review '21]



This talk

1. **Heterogeneity-aware objectives** for federated learning
[CISS '21, SVAA '21, Under Review '21]
2. **Robust aggregation** for federated learning
[TSP '22]



1	9	7	9	7
3	6	0	0	0
3	3	1	2	8
8	7	9	7	1
1	6	4	5	5

0	6	3	1	3
9	5	4	1	2
6	9	7	9	3
4	1	2	6	5
9	7	7	6	0

7	1	7	6	7
1	1	3	8	4
4	0	8	0	2
7	4	9	0	0
2	9	4	1	2

1	5	1	1	4
7	3	3	7	6
8	0	5	0	7
7	1	3	7	6
2	6	2	1	4

+

8	0	3	1	7
5	0	7	3	6
1	1	8	3	3
4	3	2	8	5
7	2	3	2	3

Clean
10%
Corrupt

Usual

64.3%

-4.2pp

60.1%

Robust

62.9%

-0.6pp

62.3%

1	9	7	9	7
3	6	0	0	0
3	3	1	2	8
8	7	9	7	1
1	6	4	5	5

0	6	3	1	3
9	5	4	1	2
6	9	7	9	3
4	1	2	6	5
9	7	7	6	0

7	1	7	6	7
1	1	3	8	4
4	0	8	0	2
7	4	9	0	0
2	9	4	1	2

1	5	1	1	4
7	3	3	7	6
8	0	5	0	7
7	1	3	7	6
2	6	2	1	4

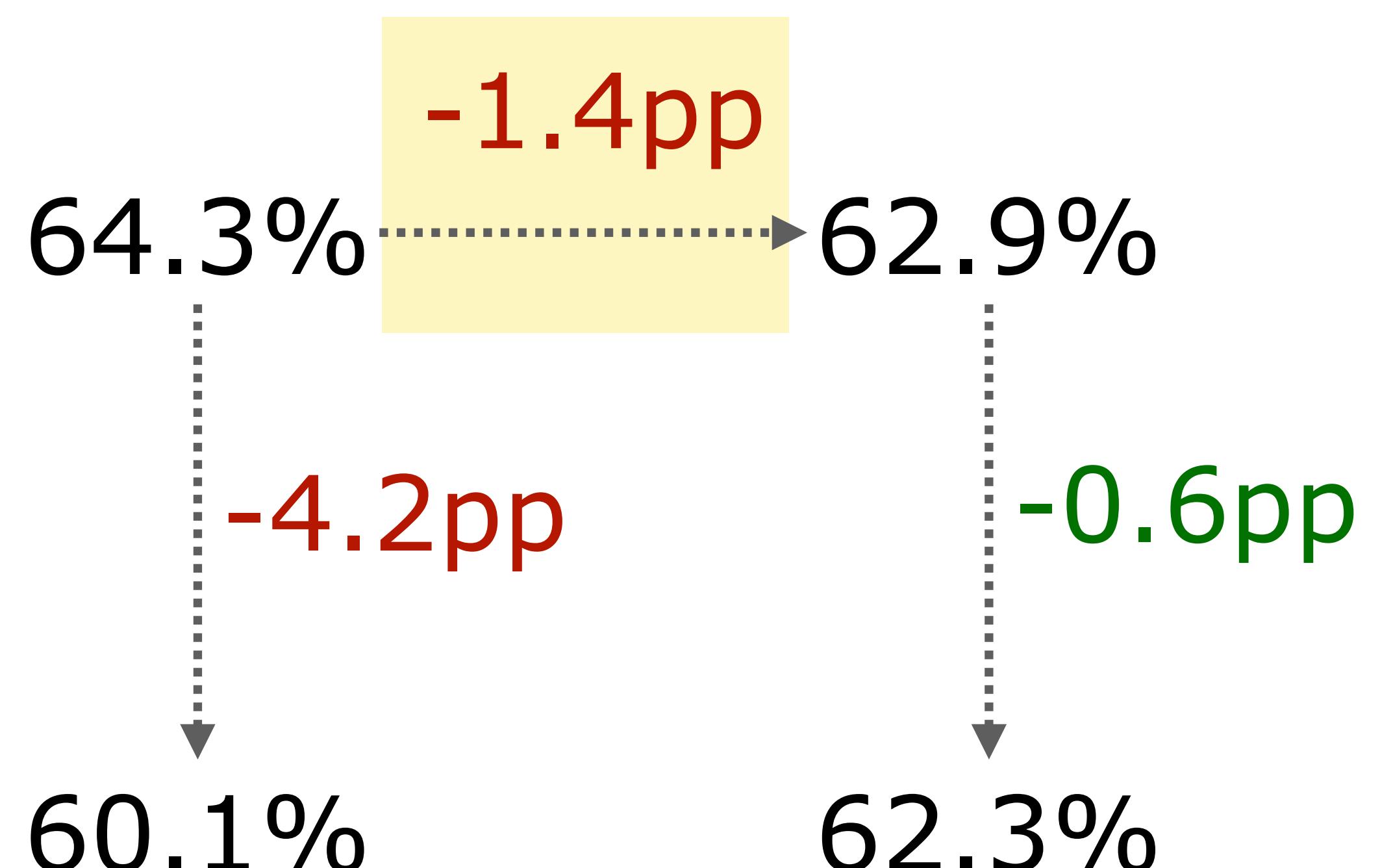
+

8	0	3	1	7
5	0	7	3	6
1	1	8	3	3
4	3	2	8	5
7	2	3	2	3

Clean
10%
Corrupt

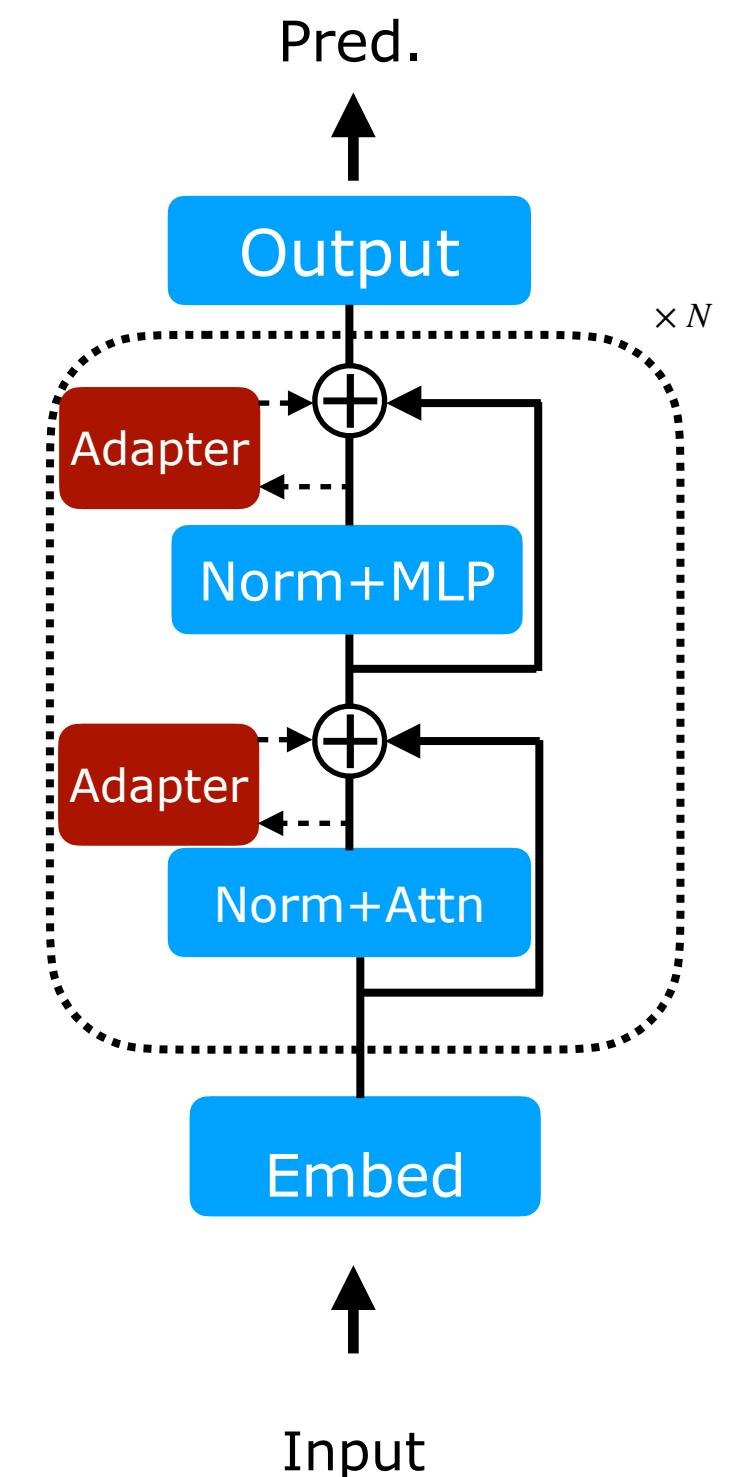
Usual

Robust



This talk

1. **Heterogeneity-aware objectives** for federated learning
[CISS '21, SVAA '21, Under Review '21]
2. **Robust aggregation** for federated learning
[TSP '22]
3. **Model personalization** for federated learning
[ICML '22, TSP '22]



1	9	7	9	7
3	6	0	0	0
3	3	1	2	8
8	7	9	7	1
1	6	4	5	5

0	6	3	1	3
9	5	4	1	2
6	9	7	9	3
4	1	2	6	5
9	7	7	6	0

7	1	7	6	7
1	1	3	8	4
4	0	8	0	2
7	4	9	0	0
2	9	4	1	2

1	5	1	1	4
7	3	3	7	6
8	0	5	0	7
7	1	3	7	6
2	6	2	1	4

+

8	0	3	1	7
5	0	7	3	6
1	1	8	3	3
4	3	2	8	5
7	2	3	2	3

Clean
25%
Corrupt

Usual

Robust

70.1%

-0.3pp

-45.5pp

24.6%

66.0%

69.8%

-3.8pp

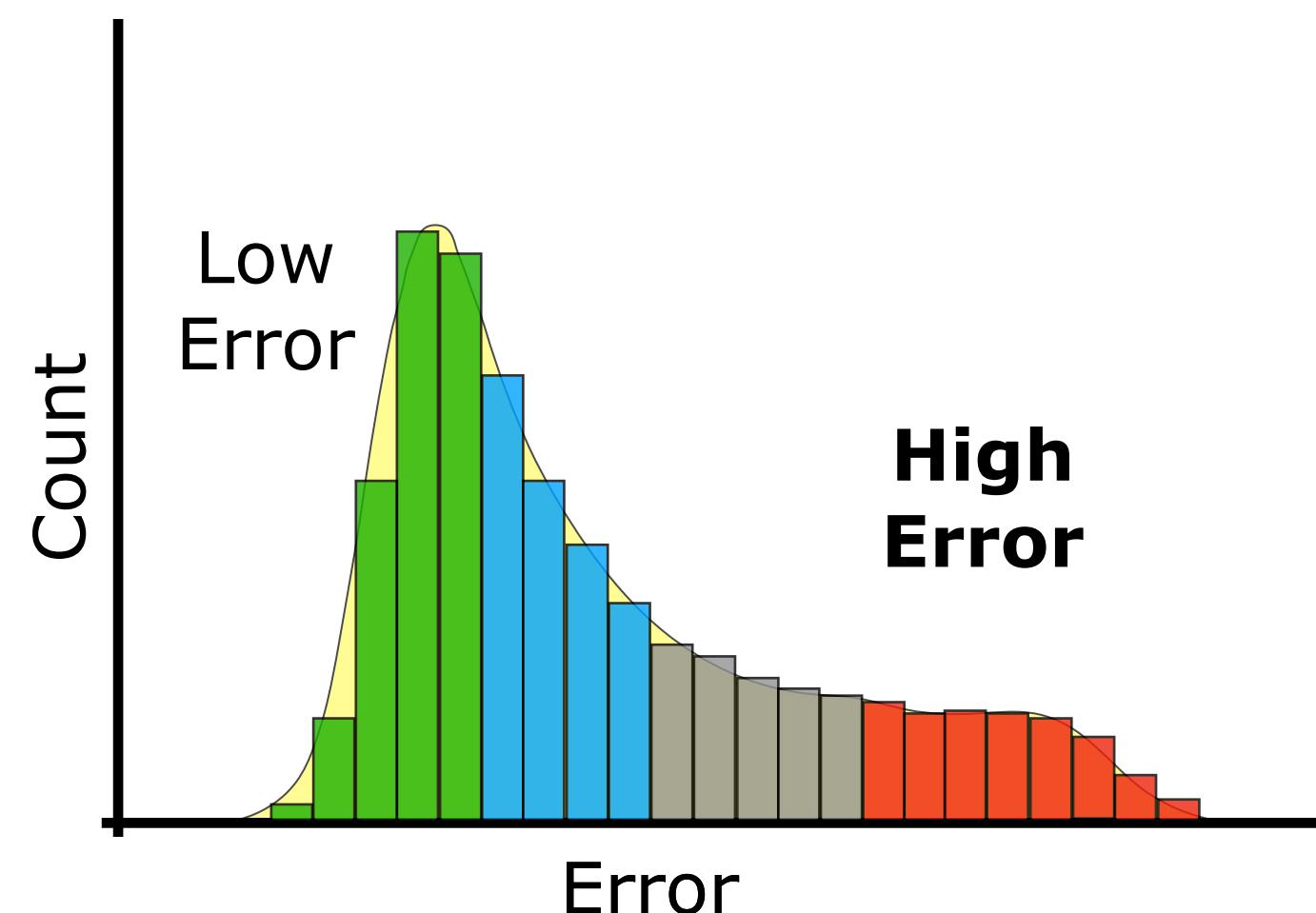


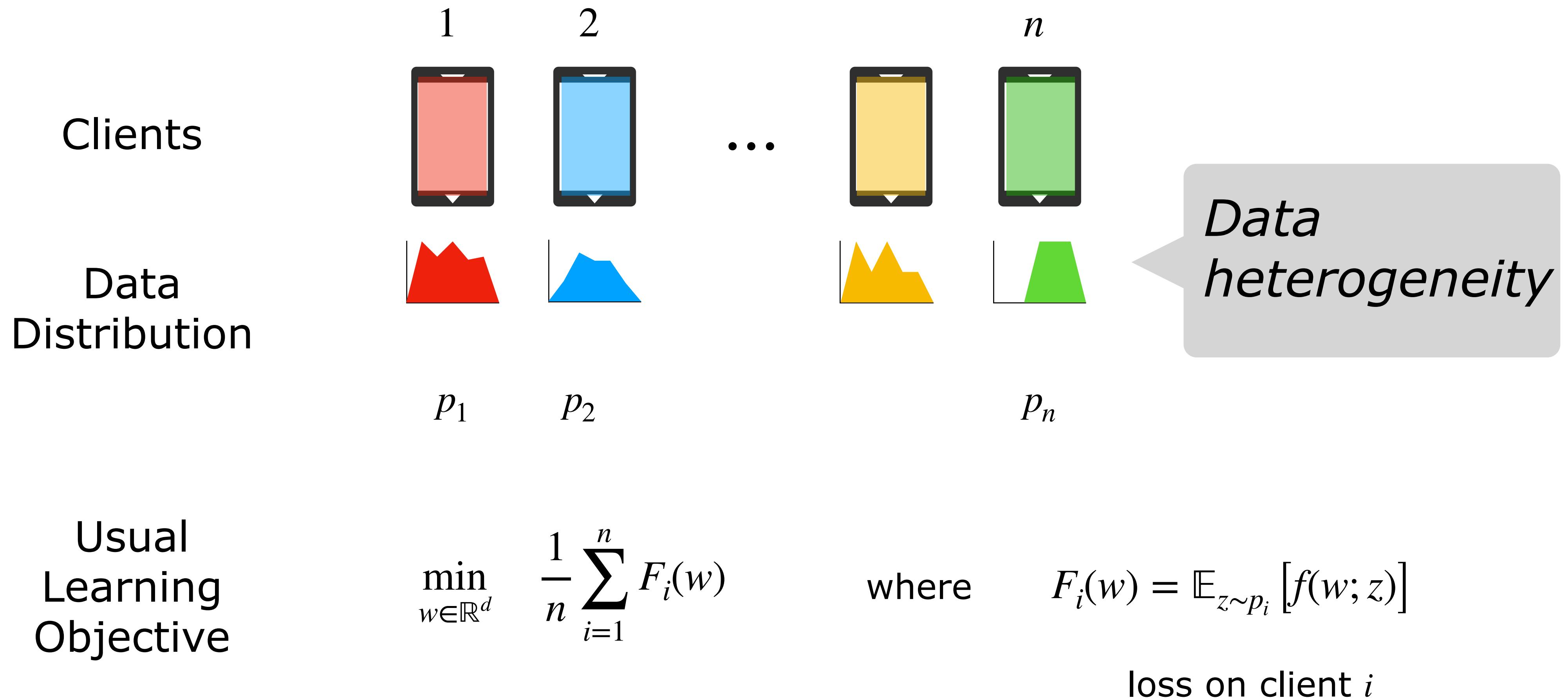
This talk

1. **Heterogeneity-aware objectives** for federated learning
[CISS '21, SVAA '21, Under Review '21]
2. **Robust aggregation** for federated learning
[TSP '22]
3. **Model personalization** for federated learning
[ICML '22, TSP '22]

Part 1: Heterogeneity-aware objectives for federated learning

[CISS '21, SVAA '21, Under Review '21]

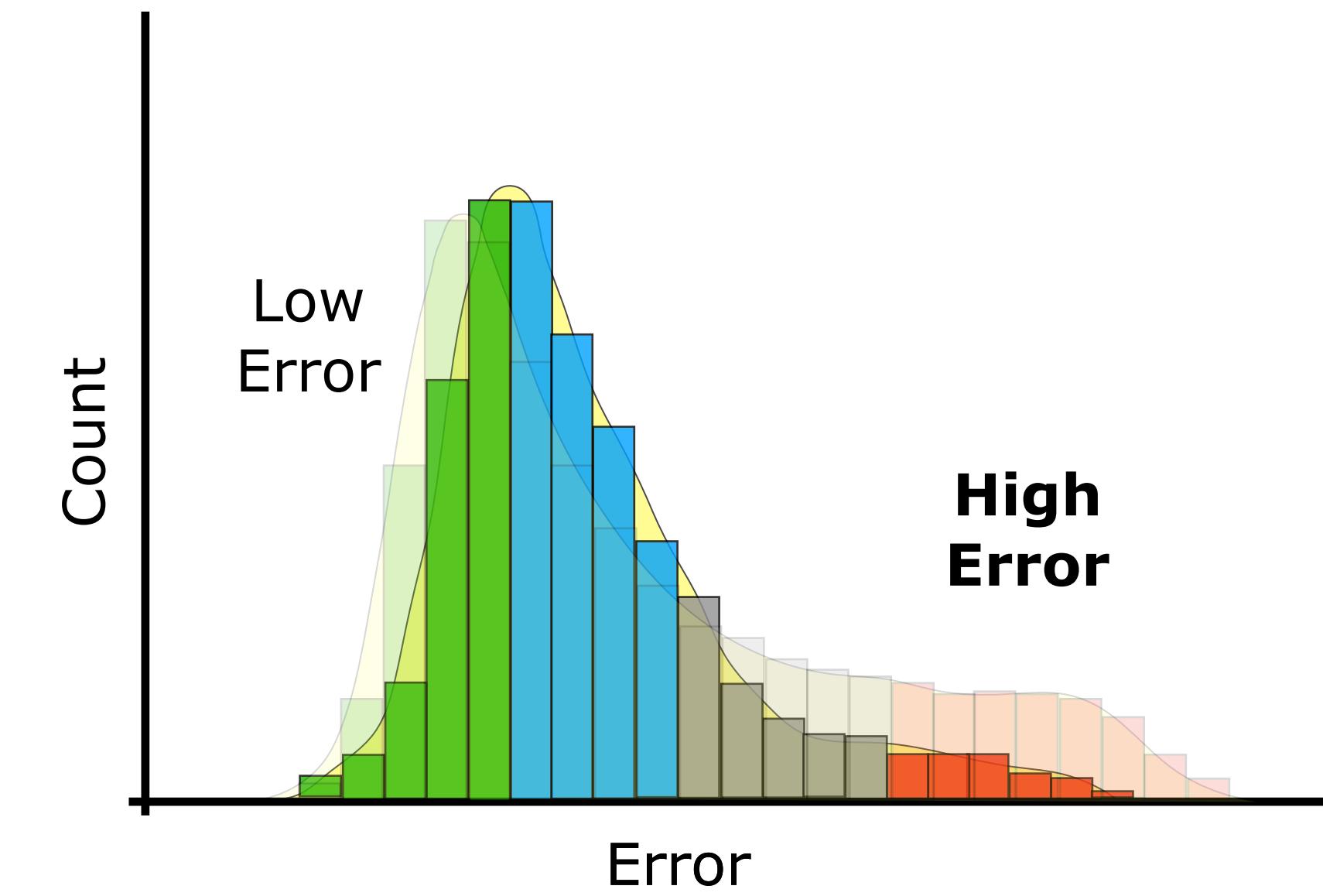




[McMahan et al. AISTATS (2017), Kairouz et al. (2021)]

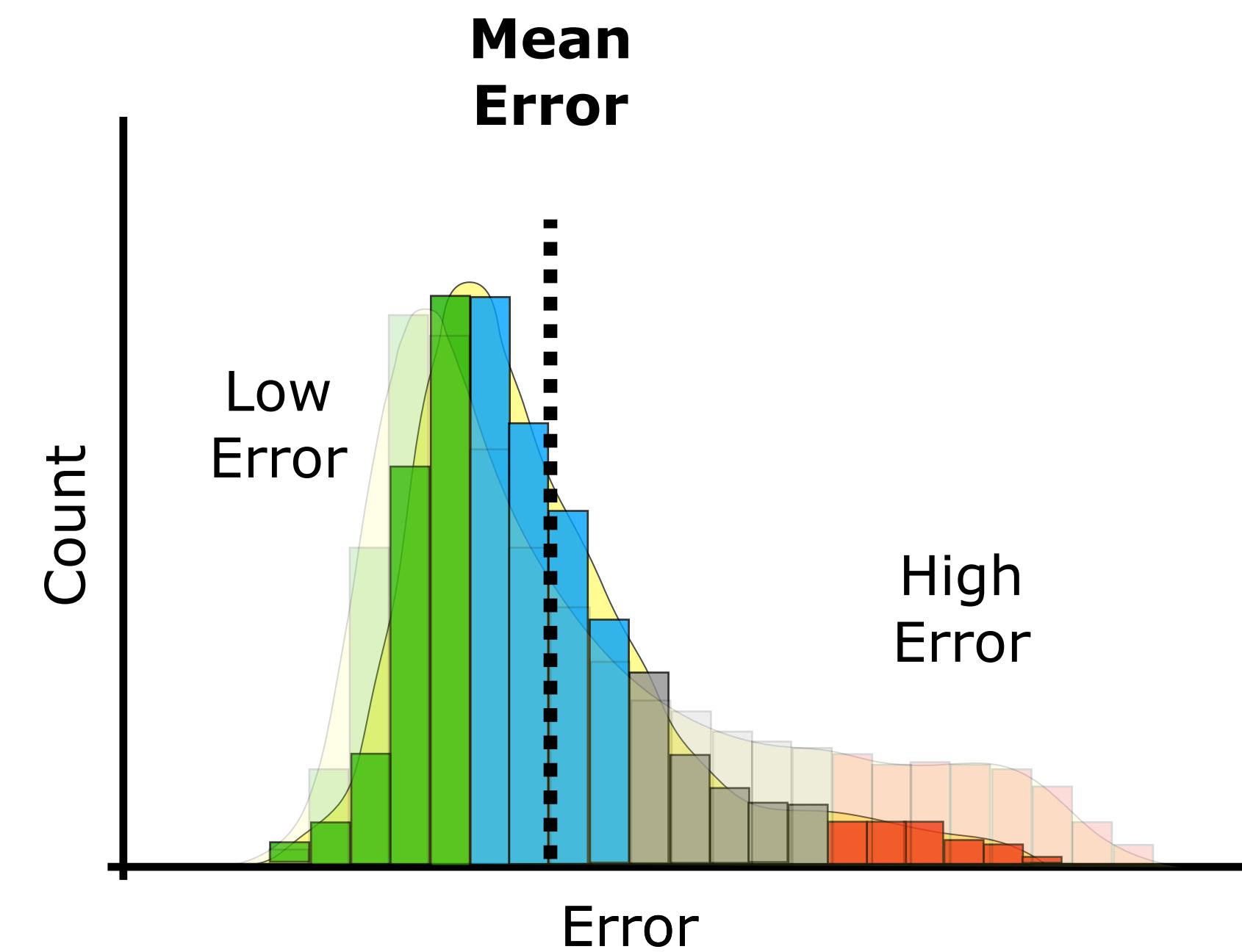
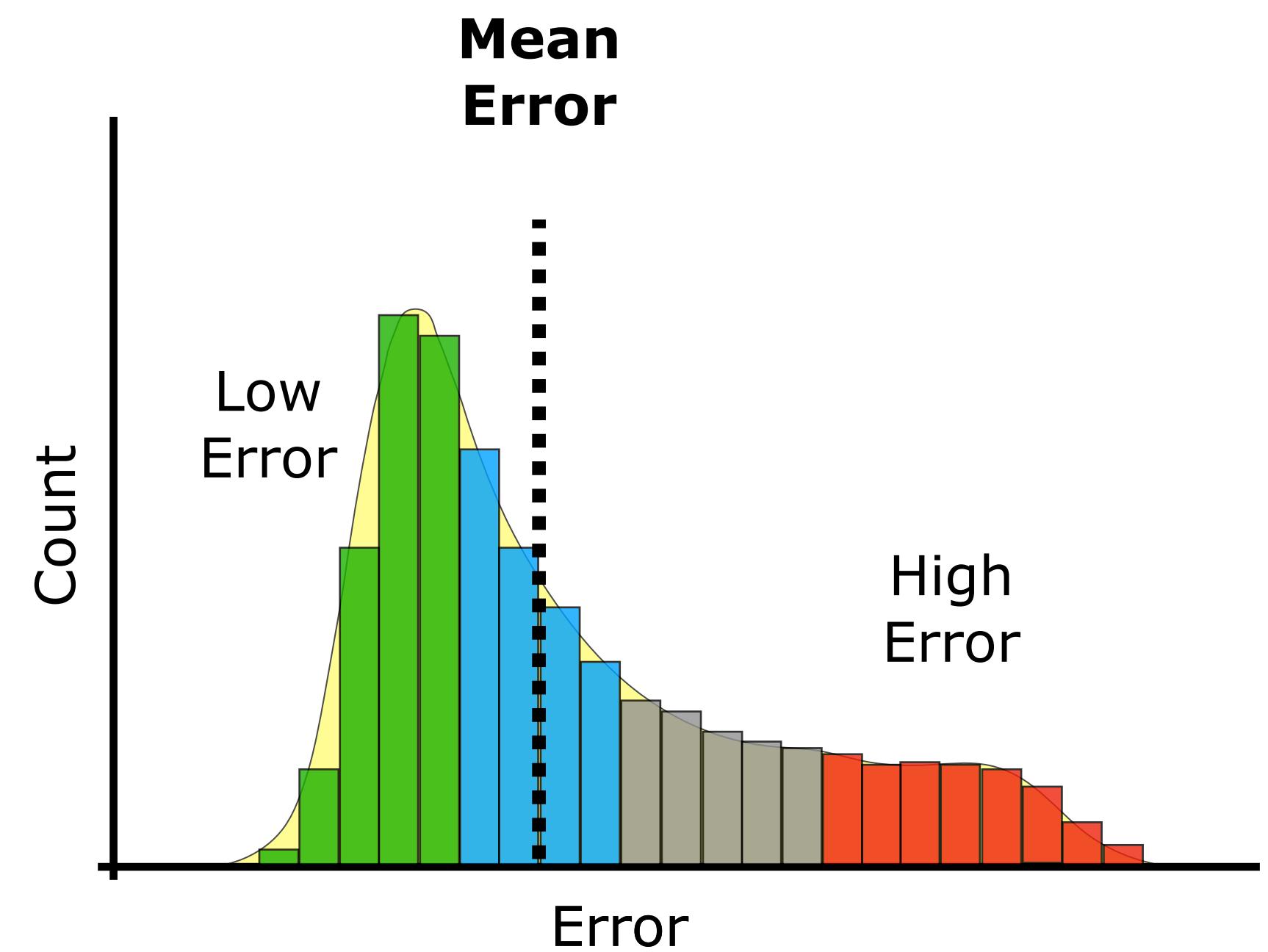
Our goal

Reduce *tail* error



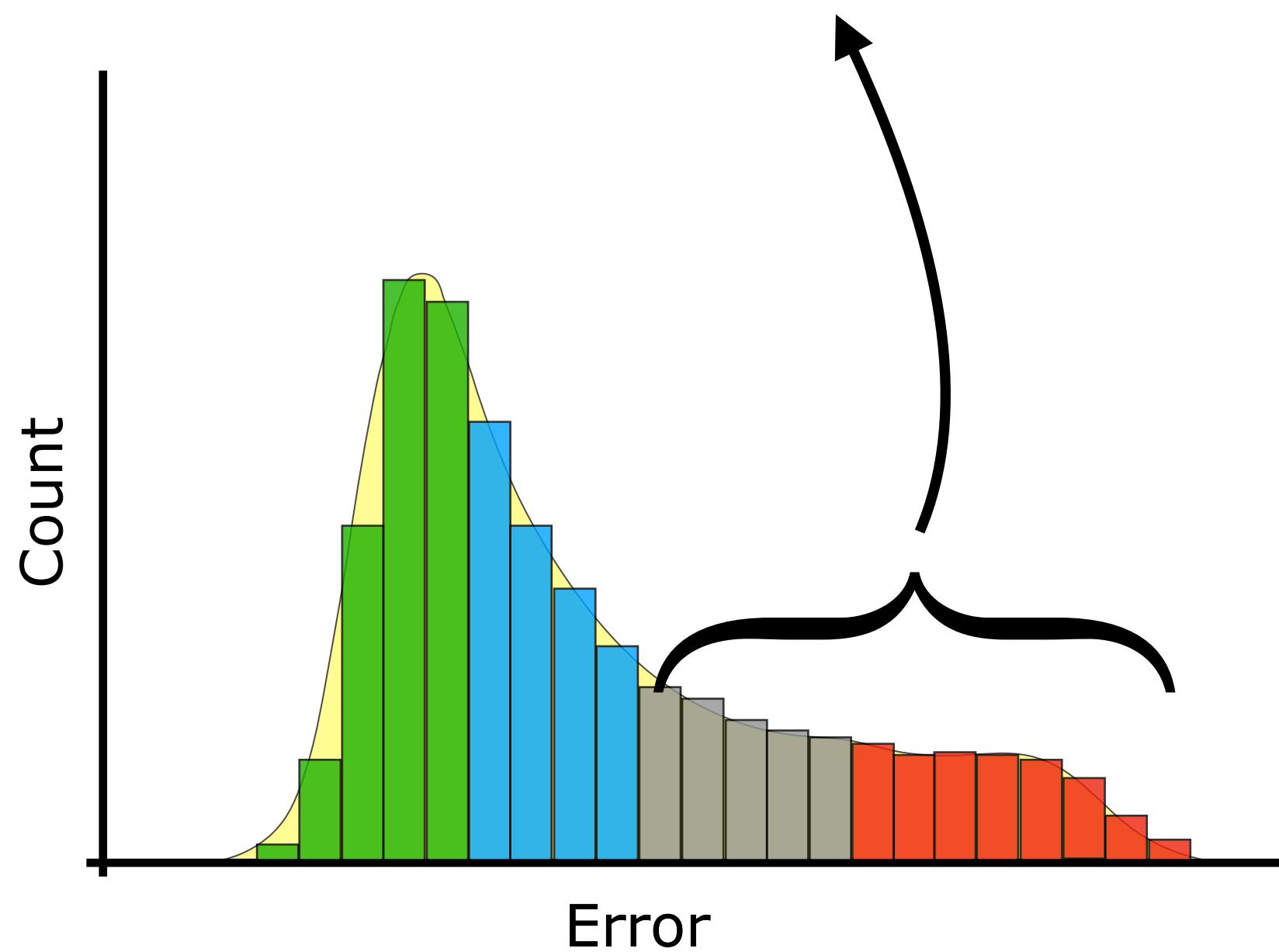
Our goal

Reduce *tail* error without sacrificing the *mean* error



Simplicial federated learning

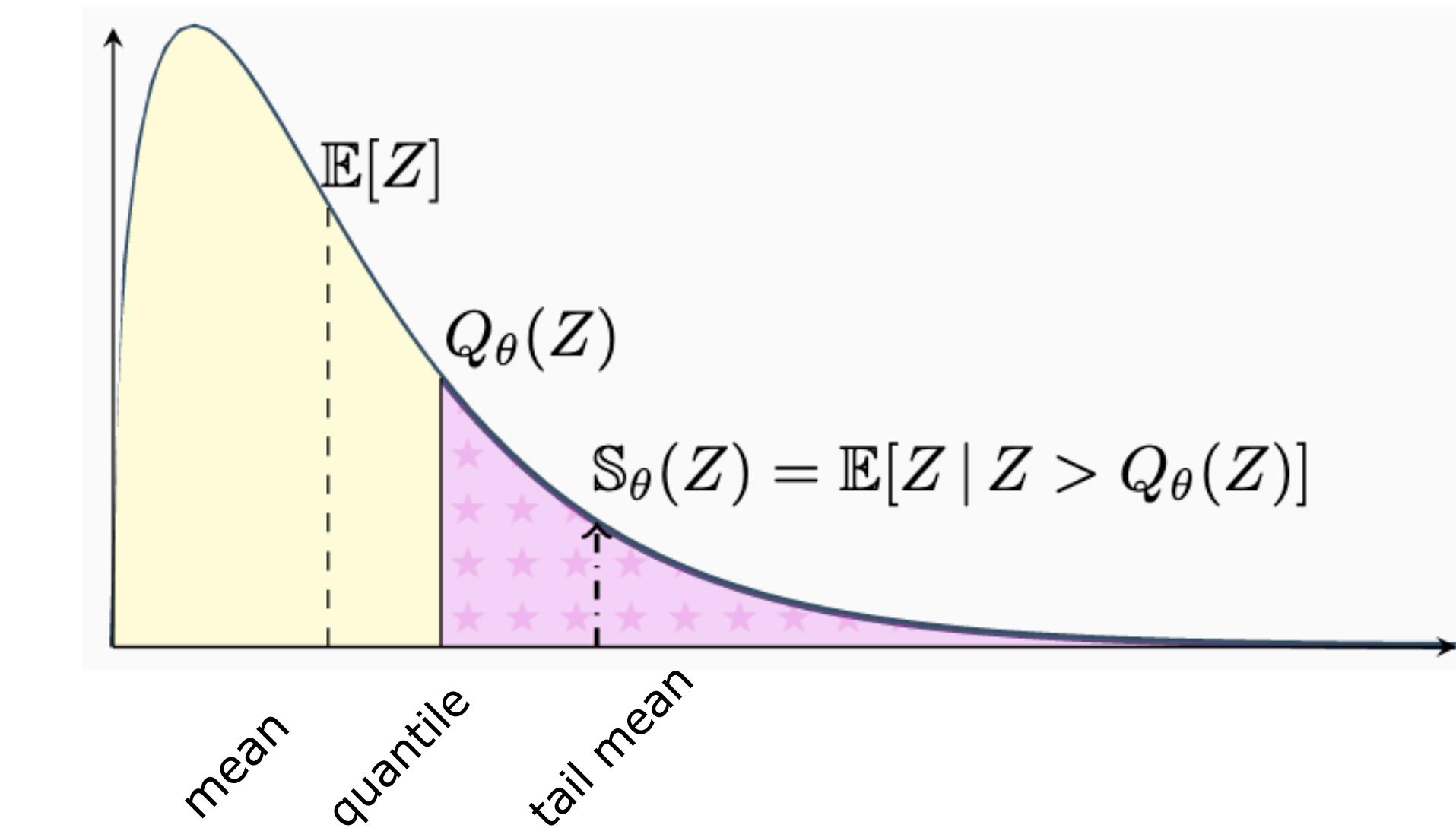
Our Approach: minimize the tail error directly!



Simplicial-FL Objective:

$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

Superquantile | Conditional Value at Risk

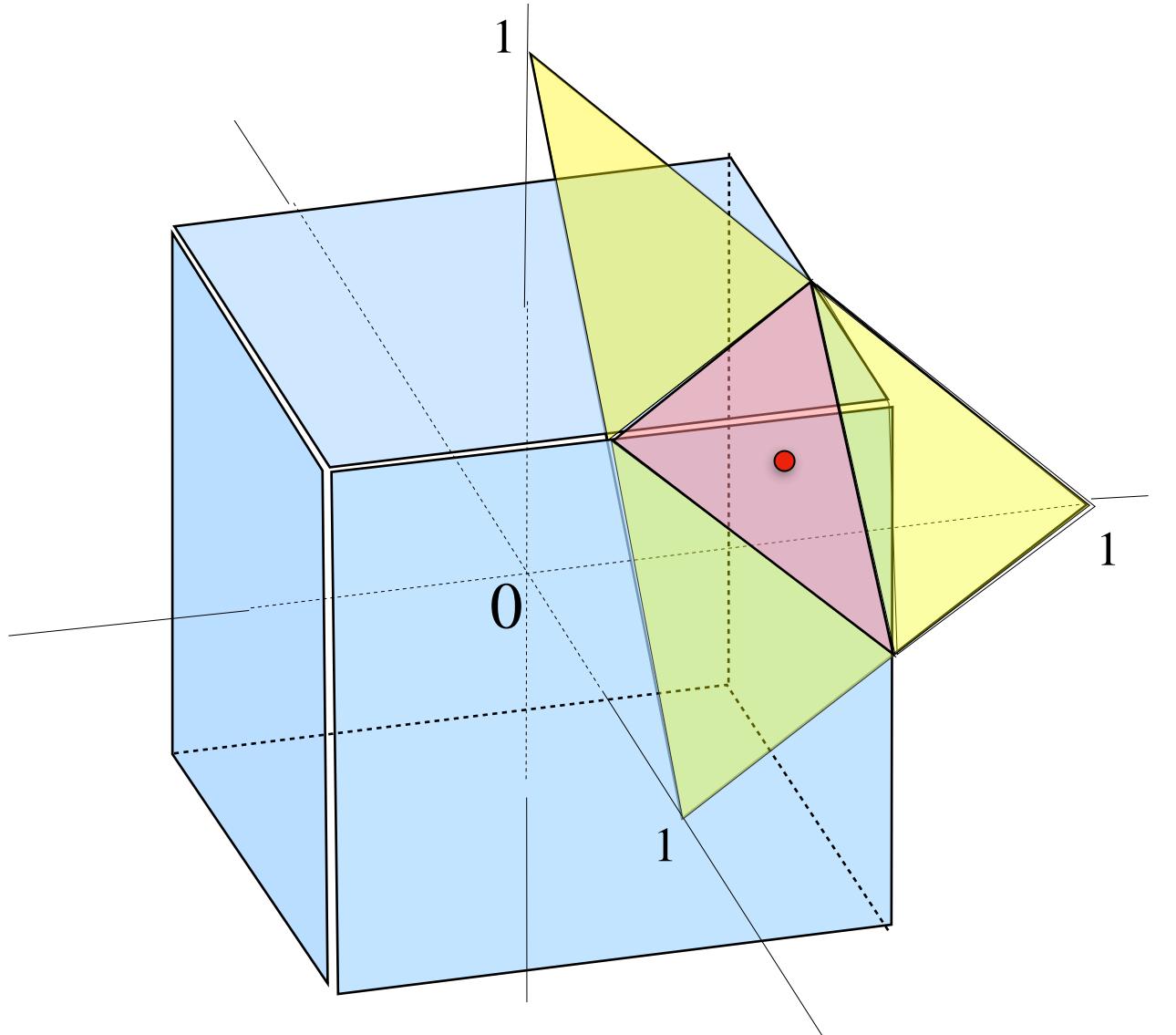


[Rockafellar & Uryasev (2002)]

Dual expression \equiv continuous knapsack problem

$$\mathbb{S}_\theta(x_1, \dots, x_n) = \max \left\{ \sum_i \pi_i x_i : \pi_i \geq 0, \sum_i \pi_i = 1, \pi_i \leq (n\theta)^{-1} \right\}$$

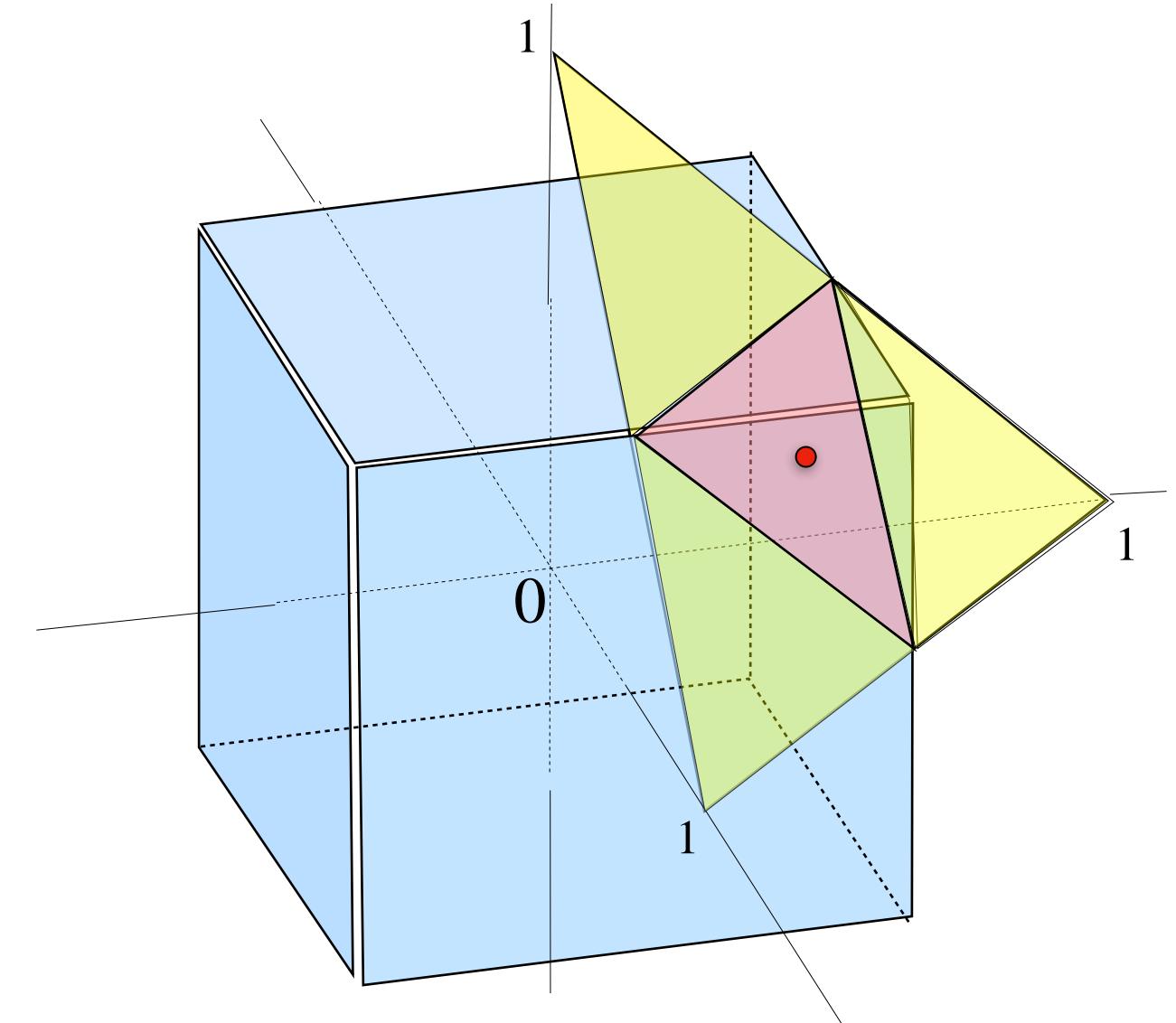
[Dantzig (1957), Ben-Tal & Teboulle (1987), Föllmer & Schied (2002)]



Dual expression \equiv continuous knapsack problem

$$\mathbb{S}_\theta(x_1, \dots, x_n) = \max \left\{ \sum_i \pi_i x_i : \pi_i \geq 0, \sum_i \pi_i = 1, \pi_i \leq (n\theta)^{-1} \right\}$$

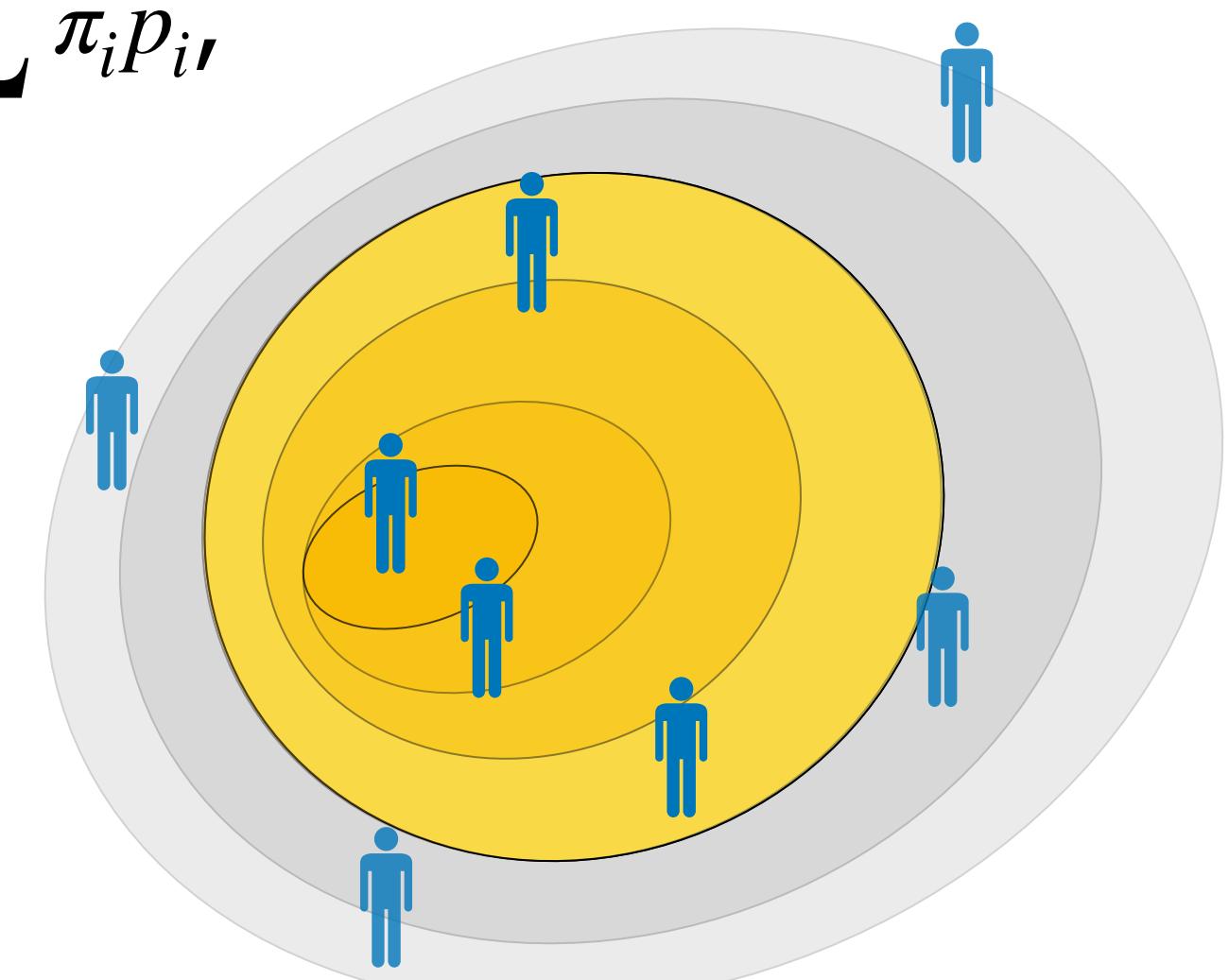
[Dantzig (1957), Ben-Tal & Teboulle (1987), Föllmer & Schied (2002)]



Assuming a new test client with mixture distribution $p_\pi = \sum_i \pi_i p_i$,

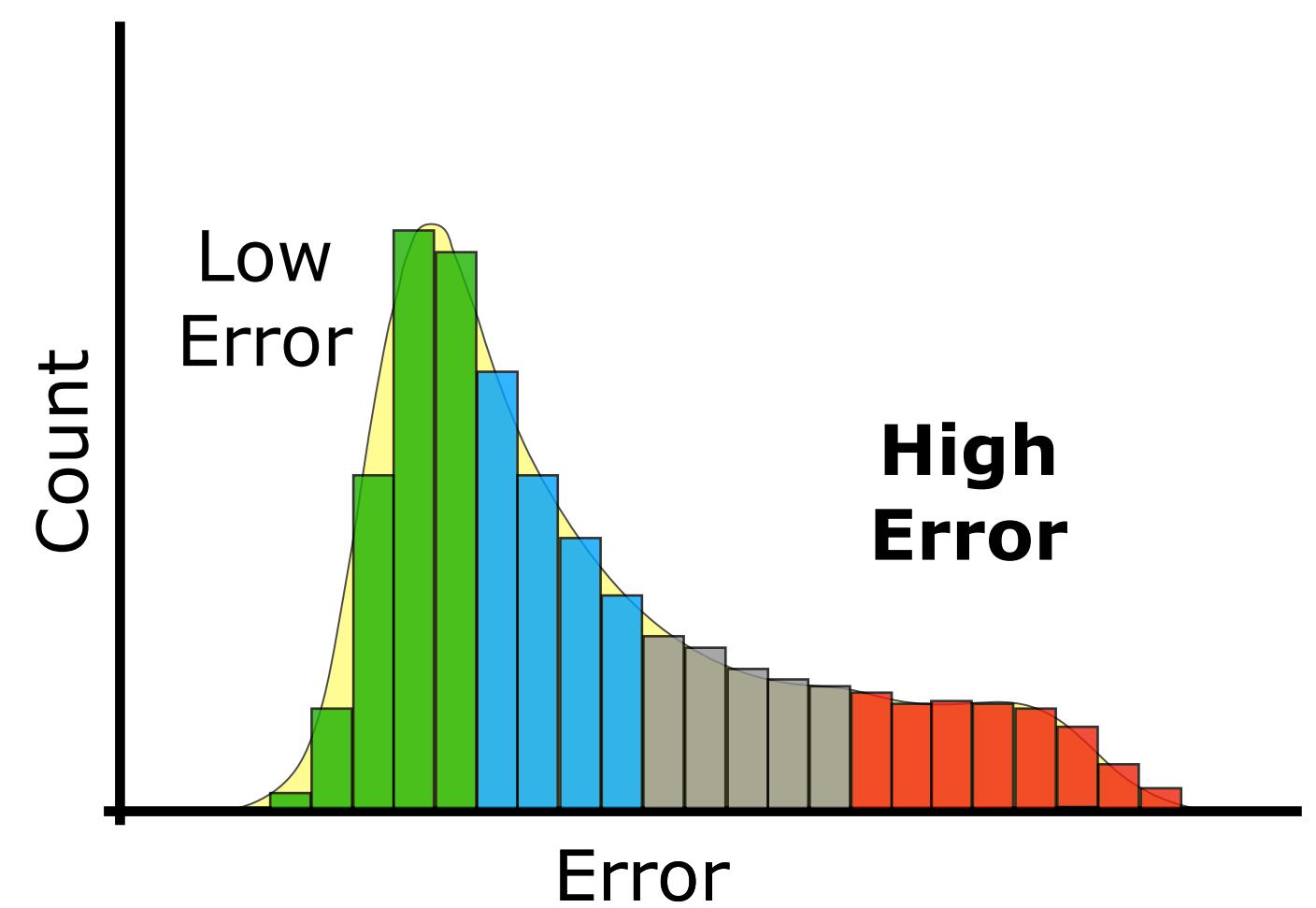
Simplicial-FL objective is equivalent to:

$$\min_w \max_{\pi: \pi_i \leq (n\theta)^{-1}} \mathbb{E}_{z \sim p_\pi} [f(w; z)]$$



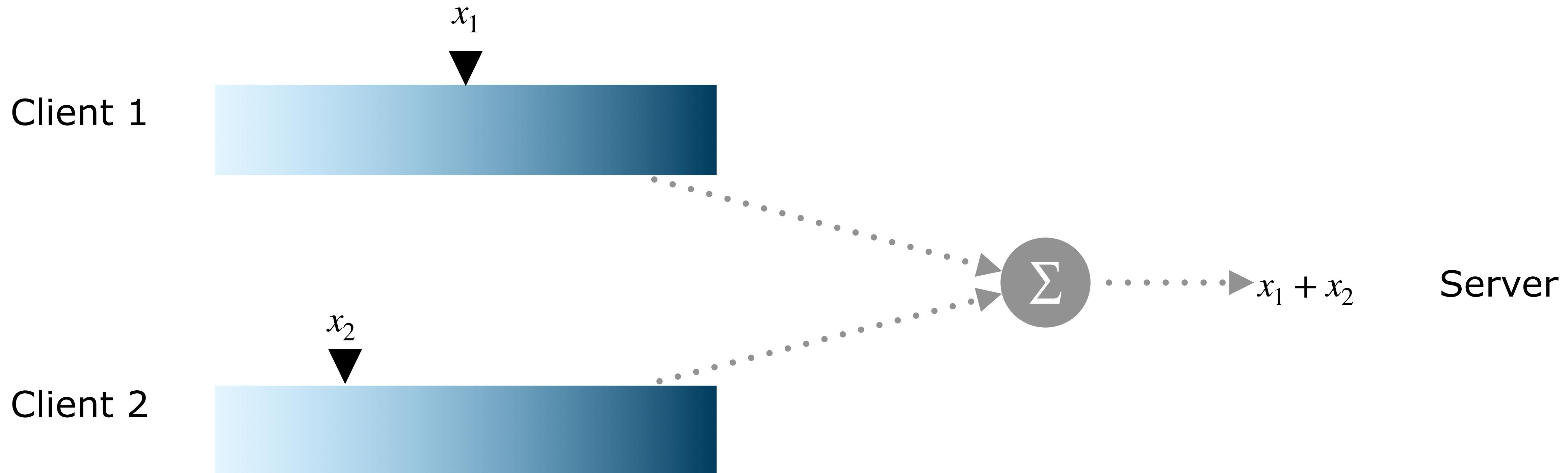
\Rightarrow Distributionally robust learning

Optimization



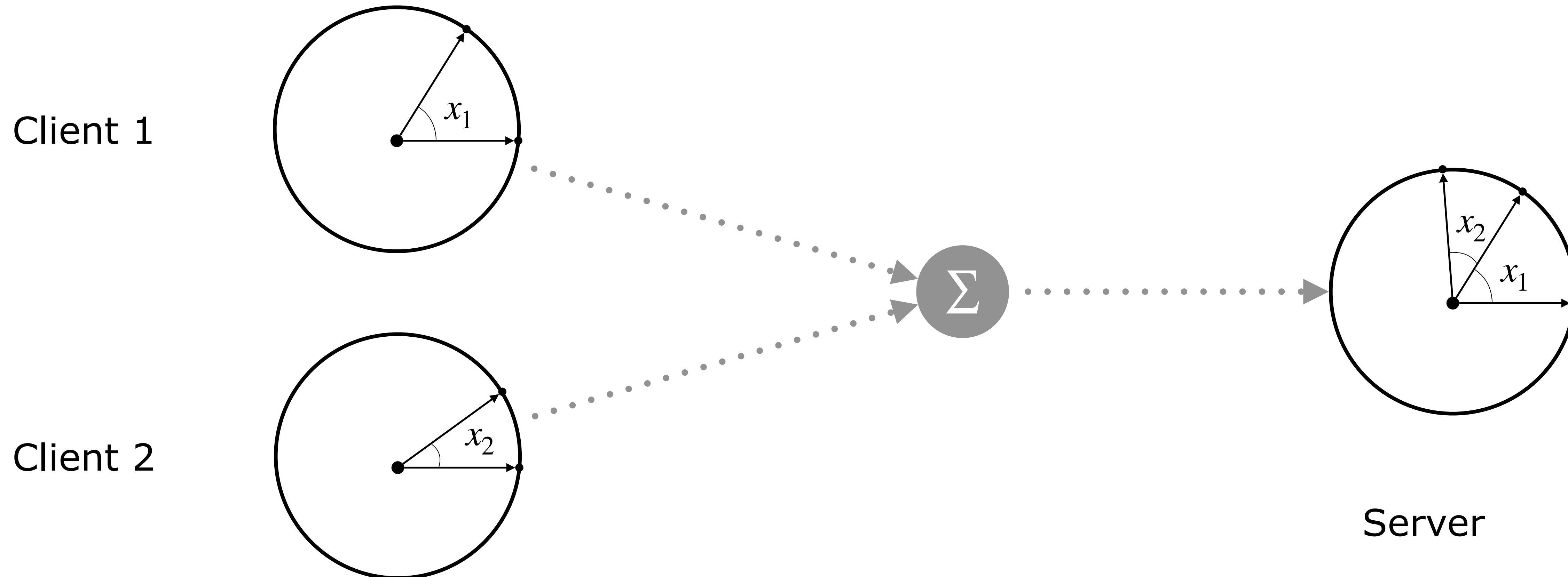
Communication primitive: secure sum

Only reveal $x_1 + x_2$ to the server without revealing x_1 or x_2



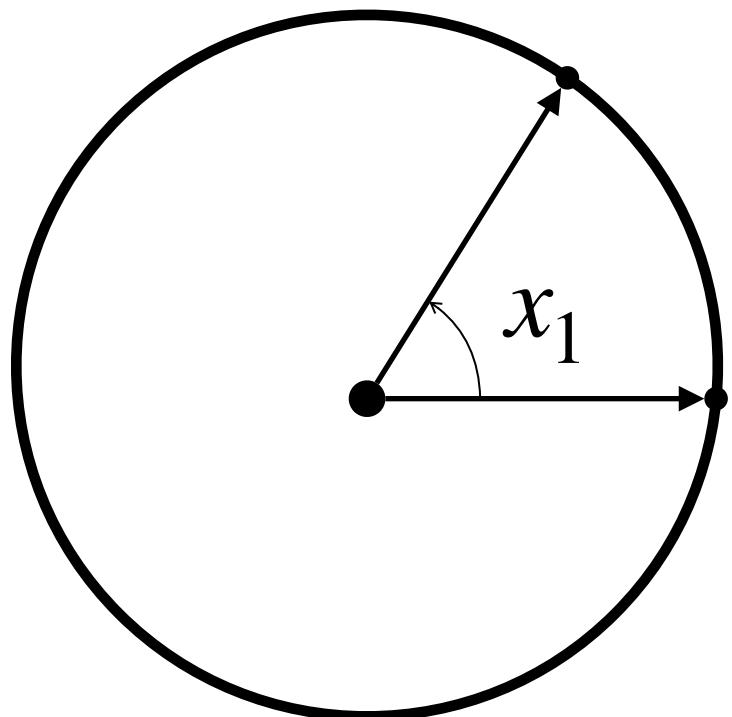
[Bonawitz et al. CCS (2017), Bell et al. CCS (2020)]

Perform all operations modulo M



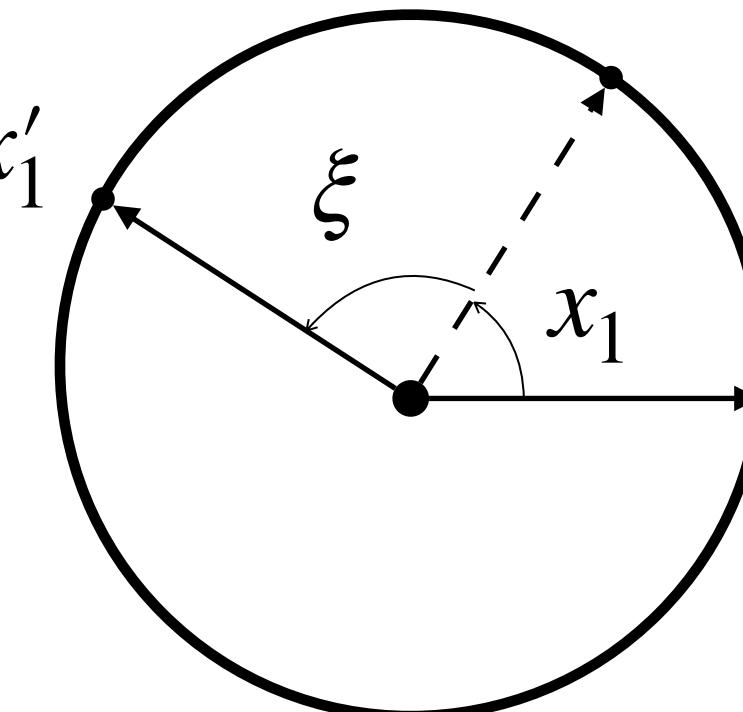
[Bonawitz et al. CCS (2017), Bell et al. CCS (2020)]

Client 1



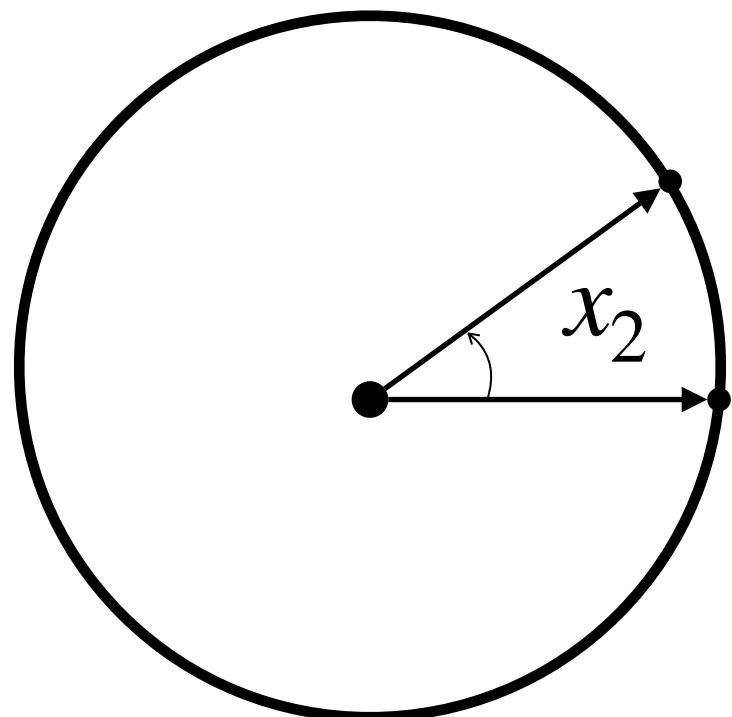
$$x'_1 = x_1 + \xi$$

..... ➔

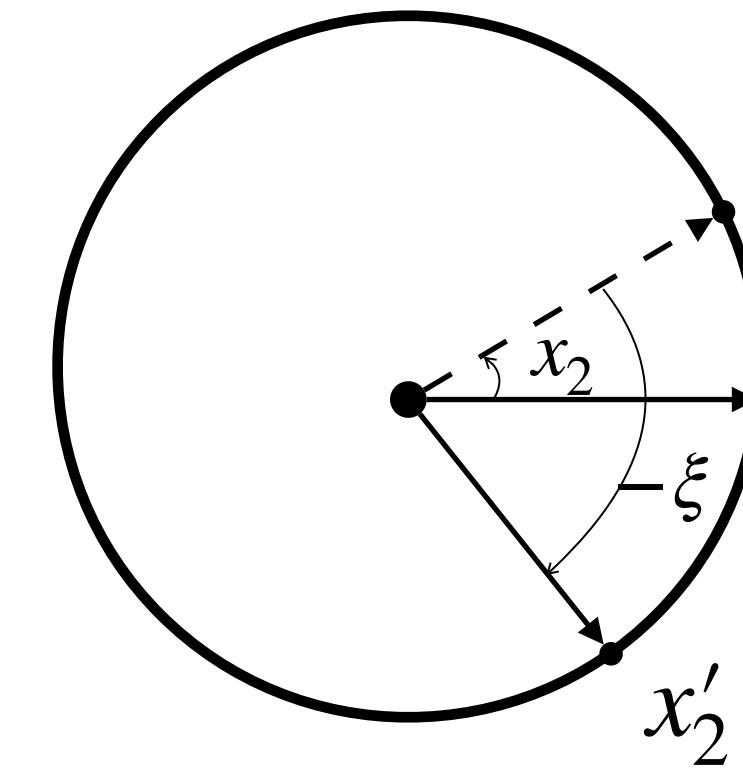


$$\xi \sim \text{Unif}(\mathbb{O})$$

Client 2

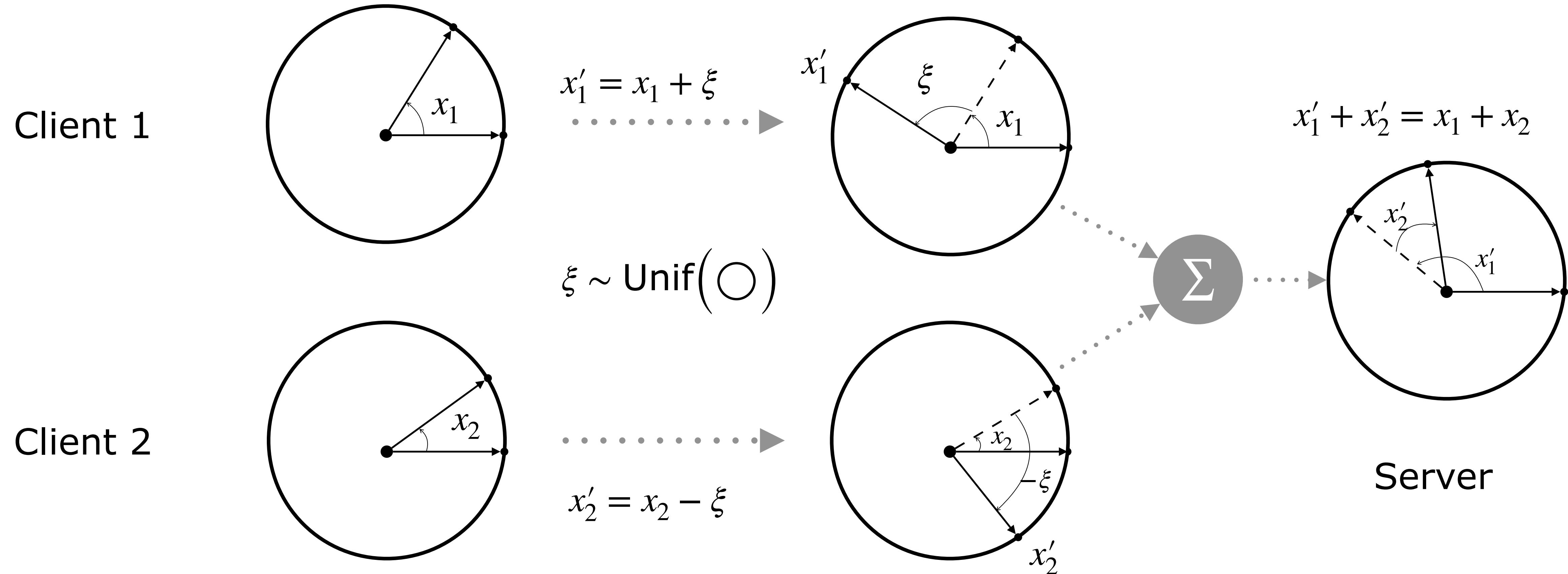


$$..... ➔$$
$$x'_2 = x_2 - \xi$$



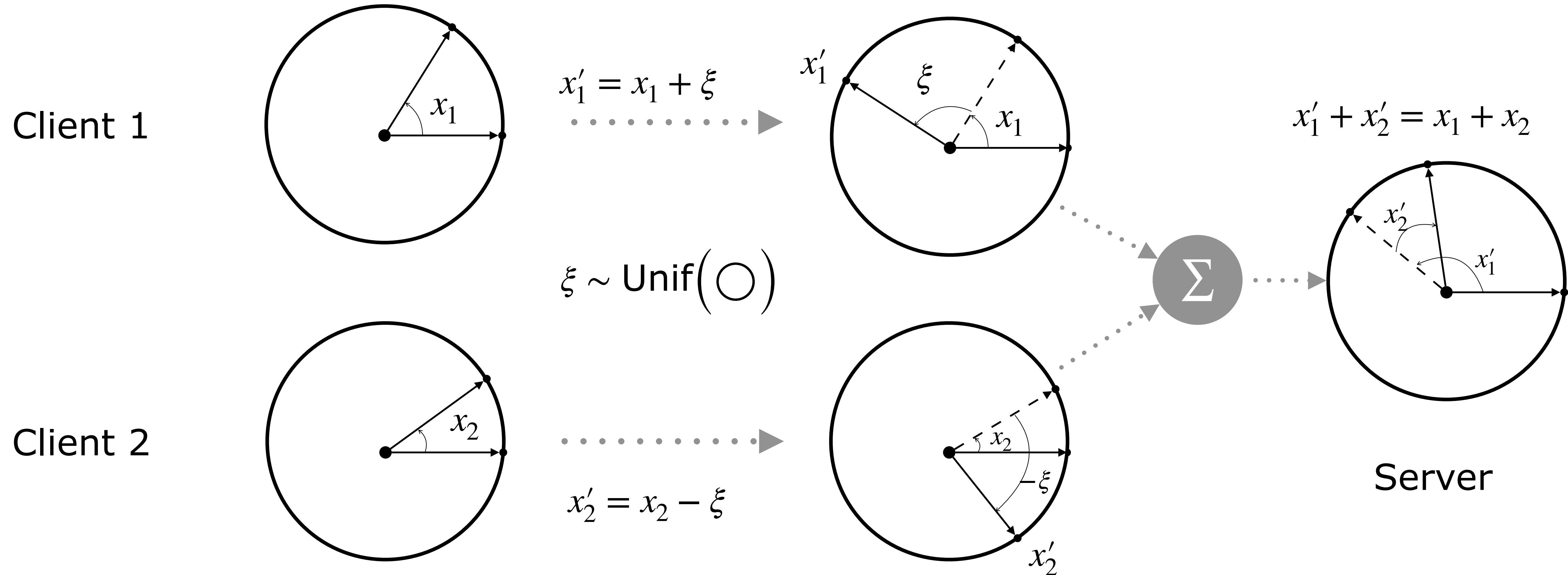
[Bonawitz et al. CCS (2017), Bell et al. CCS (2020)]

Server only sees $x'_1, x'_2 \sim \text{Unif}(\mathcal{O})$ but calculates the correct sum



[Bonawitz et al. CCS (2017), Bell et al. CCS (2020)]

Server only sees $x'_1, x'_2 \sim \text{Unif}(\mathcal{O})$ but calculates the correct sum



Total communication for m vectors in $\mathbb{R}^d = O(m \log m + md)$ numbers

Server only sees $x'_1, x'_2 \sim \text{Unif}(\mathcal{O})$ but calculates the correct sum



Real-world communication constraint:

All client-to-server communication must go through secure summation

Total communication for m vectors in $\mathbb{R}^d = O(m \log m + md)$ numbers

ERM Algorithm (FedAvg):

$$\min_w \quad \frac{1}{n} \sum_{i=1}^n F_i(w)$$

Simplicial-FL Algorithm:

$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

FedAvg [MacMahan et al. AISTATS (2017)]

Parallel Gradient Distribution [Mangasarian. SICON (1995)]

Iterative Parameter Mixing [McDonald et al. ACL (2009)]

BMUF [Chen & Huo. ICASSP (2016)]

Local SGD [Stich. ICLR (2019)]

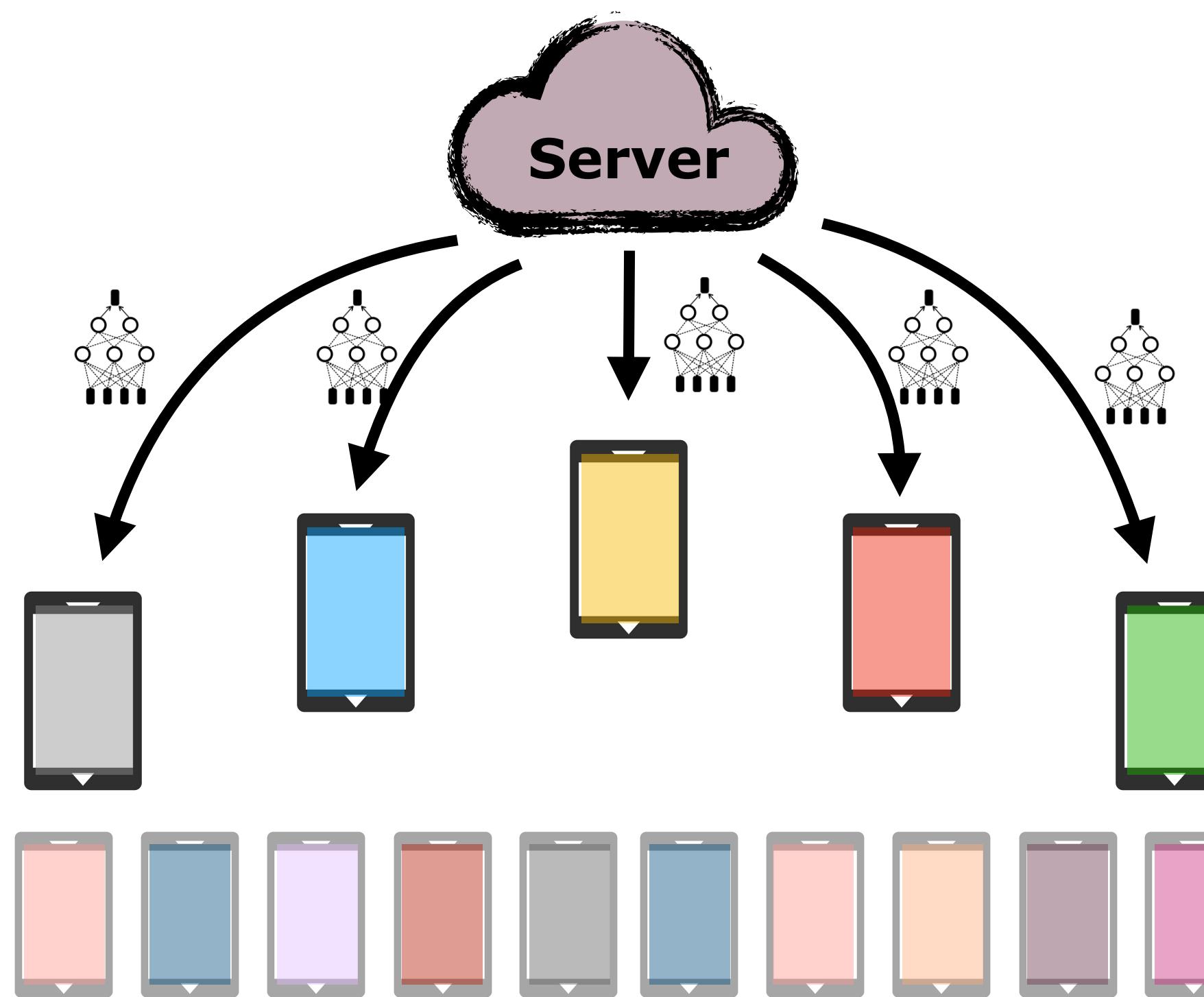
ERM Algorithm (FedAvg):

$$\min_w \frac{1}{n} \sum_{i=1}^n F_i(w)$$

Simplicial-FL Algorithm:

$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

*Step 1 of 3: Server samples m clients
and broadcasts global model*



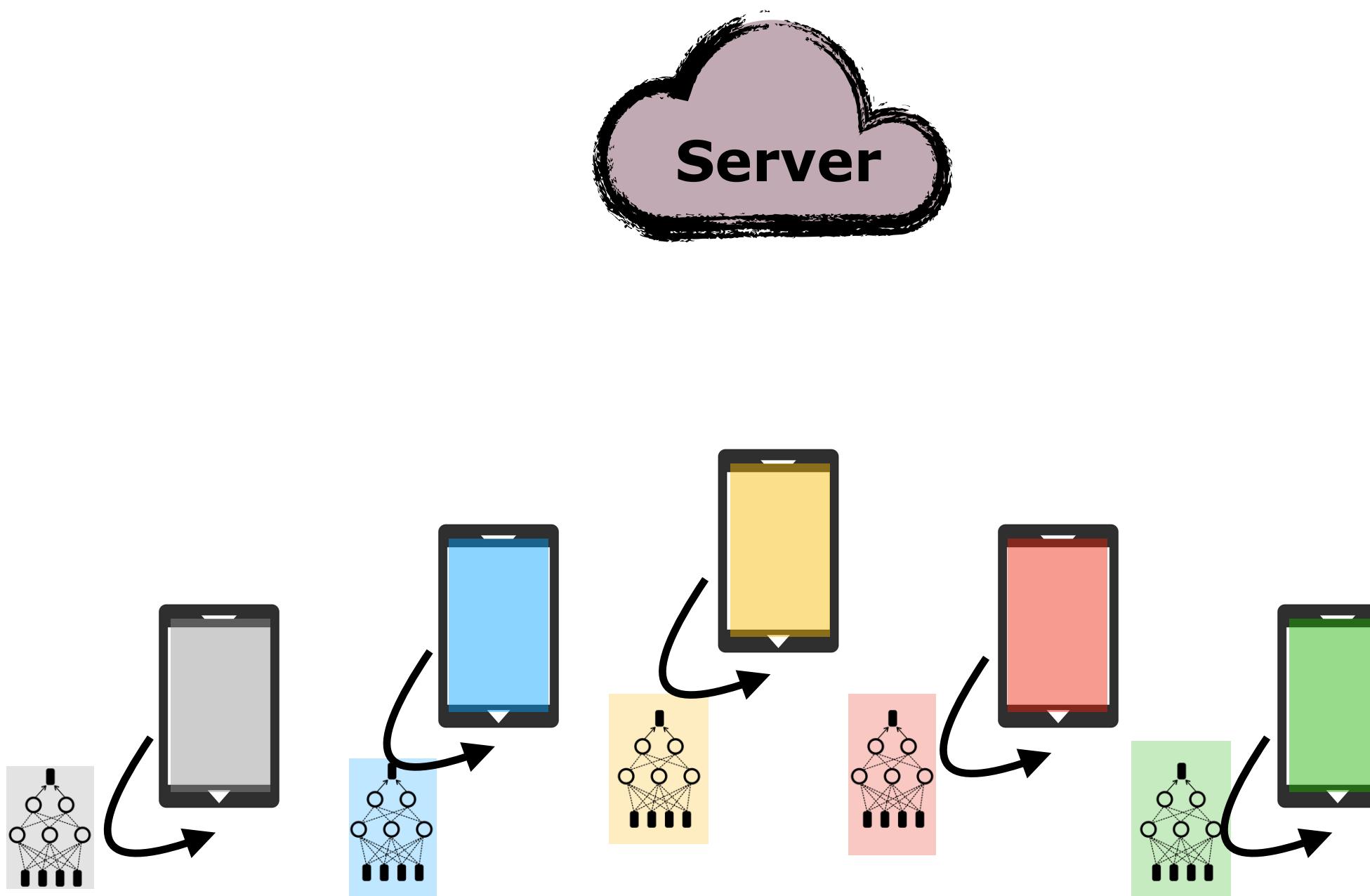
ERM Algorithm (FedAvg):

$$\min_w \frac{1}{n} \sum_{i=1}^n F_i(w)$$

Simplicial-FL Algorithm:

$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

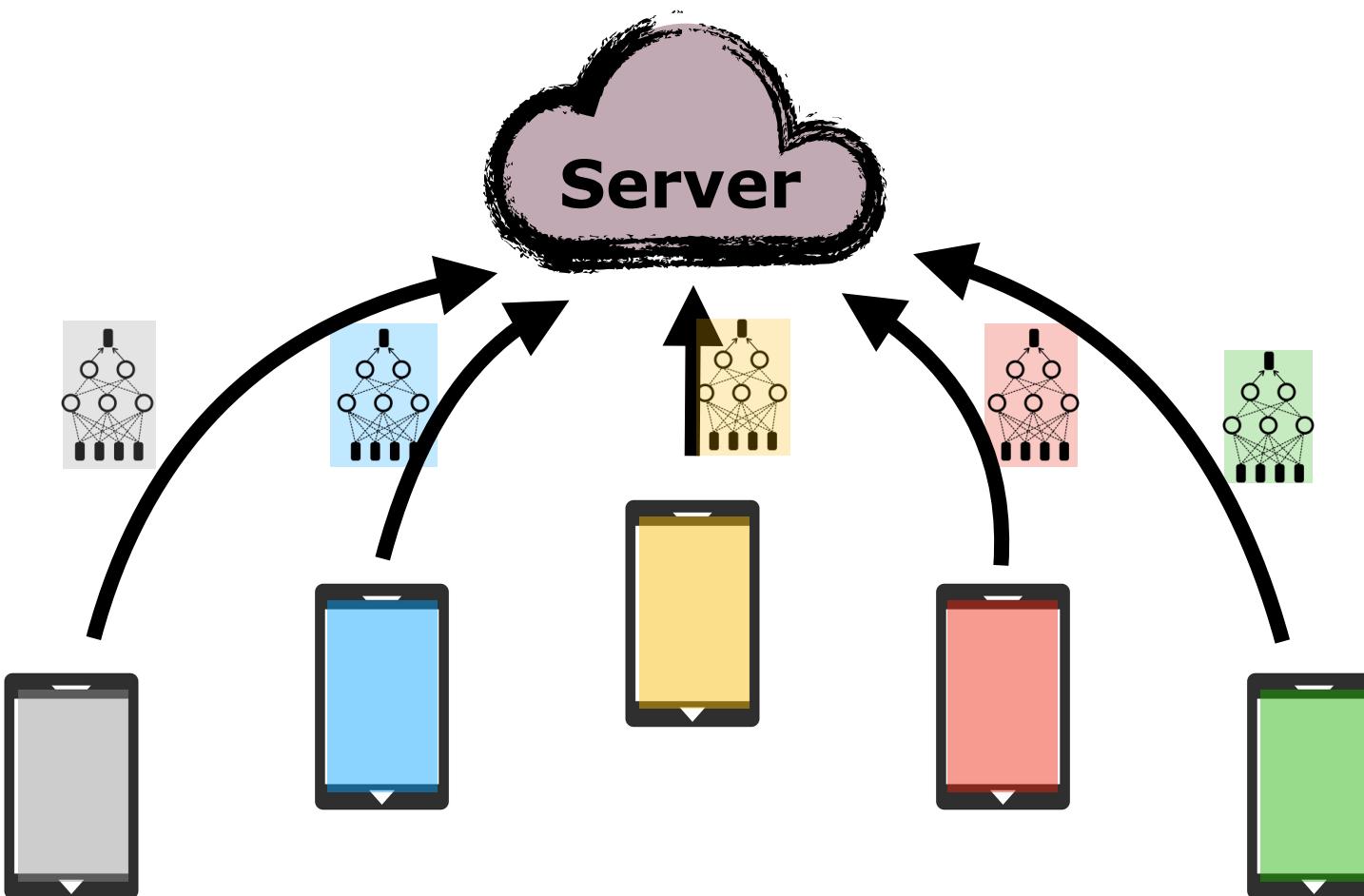
Step 2 of 3: Clients perform τ local SGD steps on their local data



ERM Algorithm (FedAvg):

$$\min_w \quad \frac{1}{n} \sum_{i=1}^n F_i(w)$$

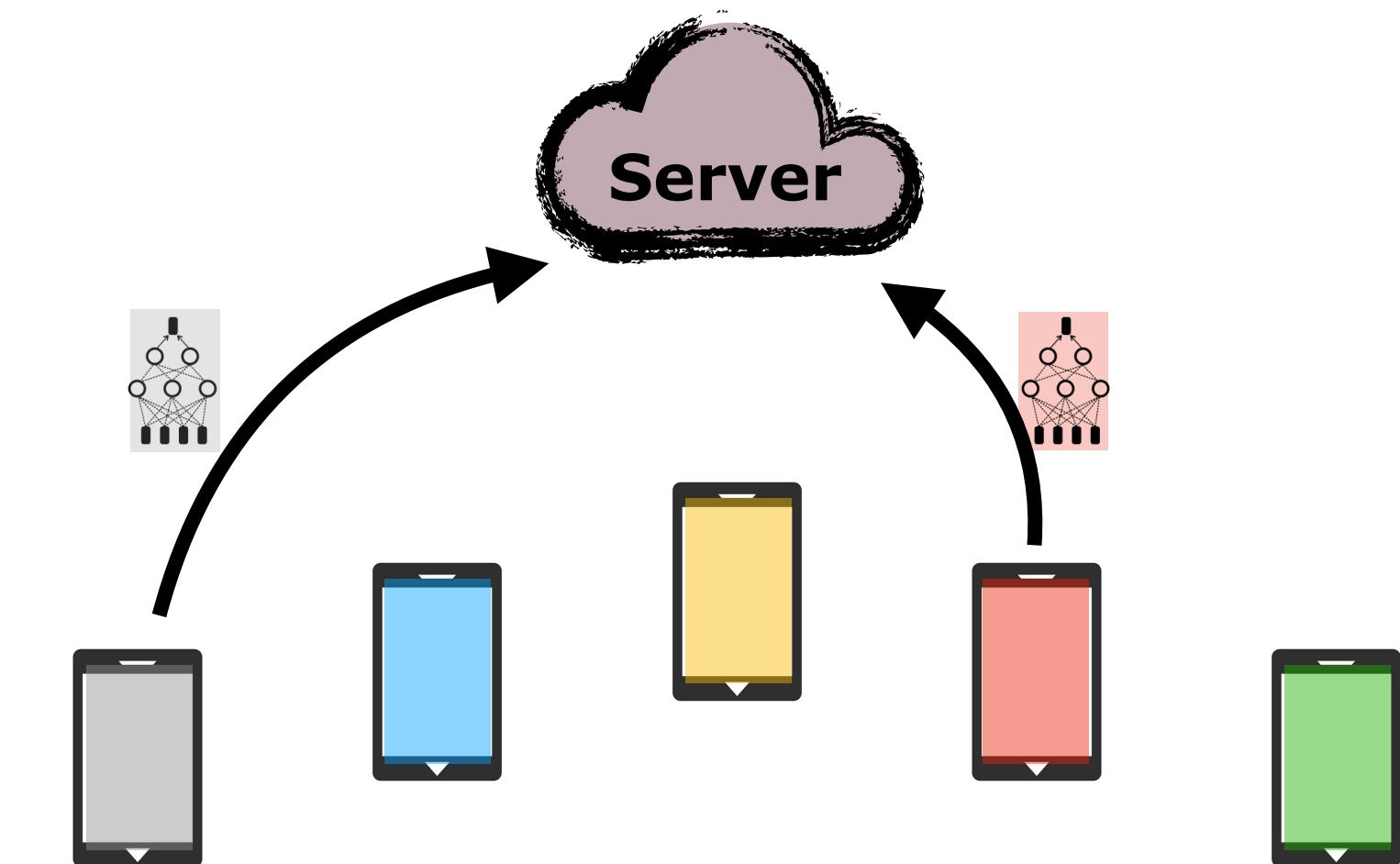
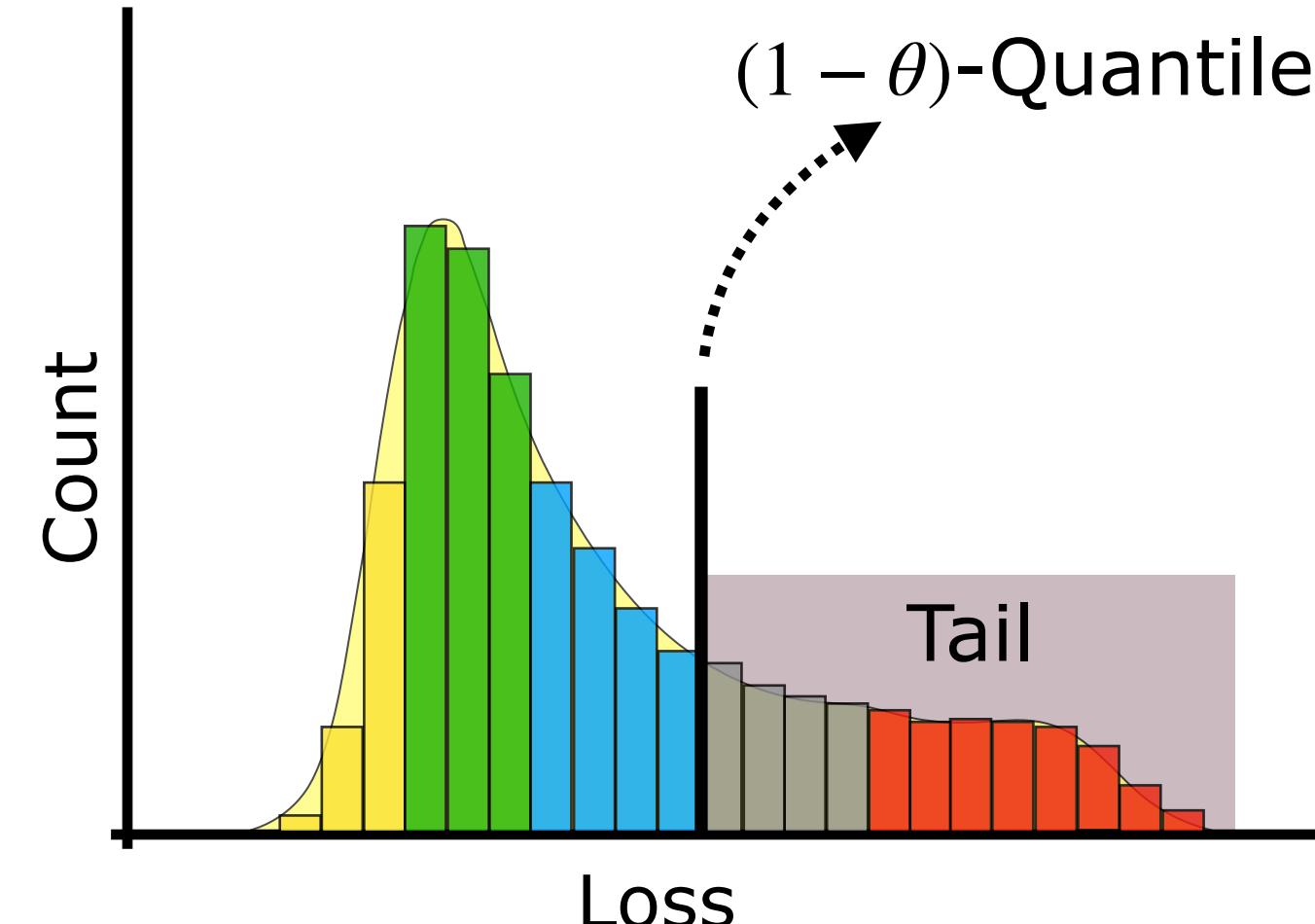
*Step 3 of 3: Aggregate updates contributed by **all clients***

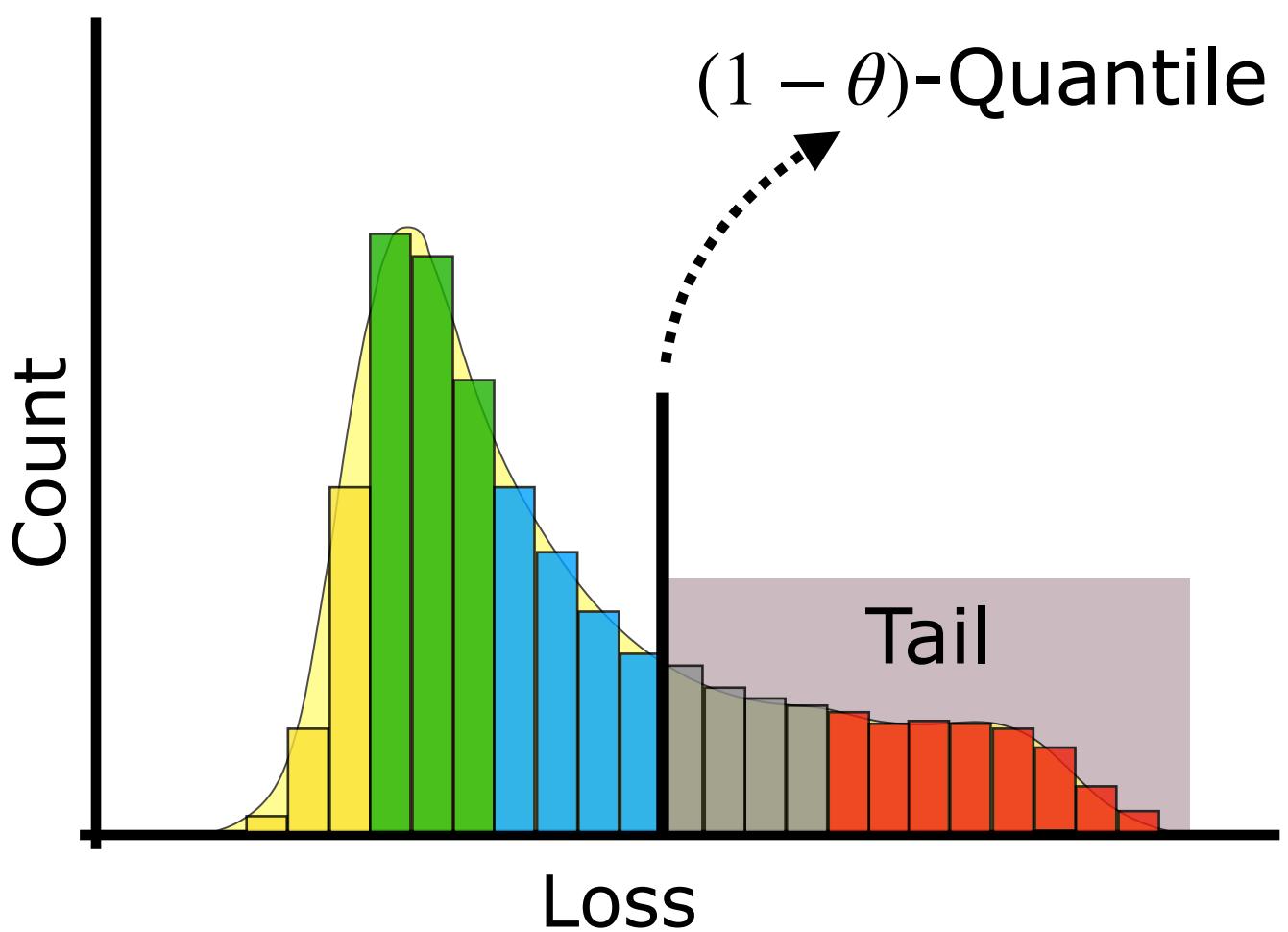


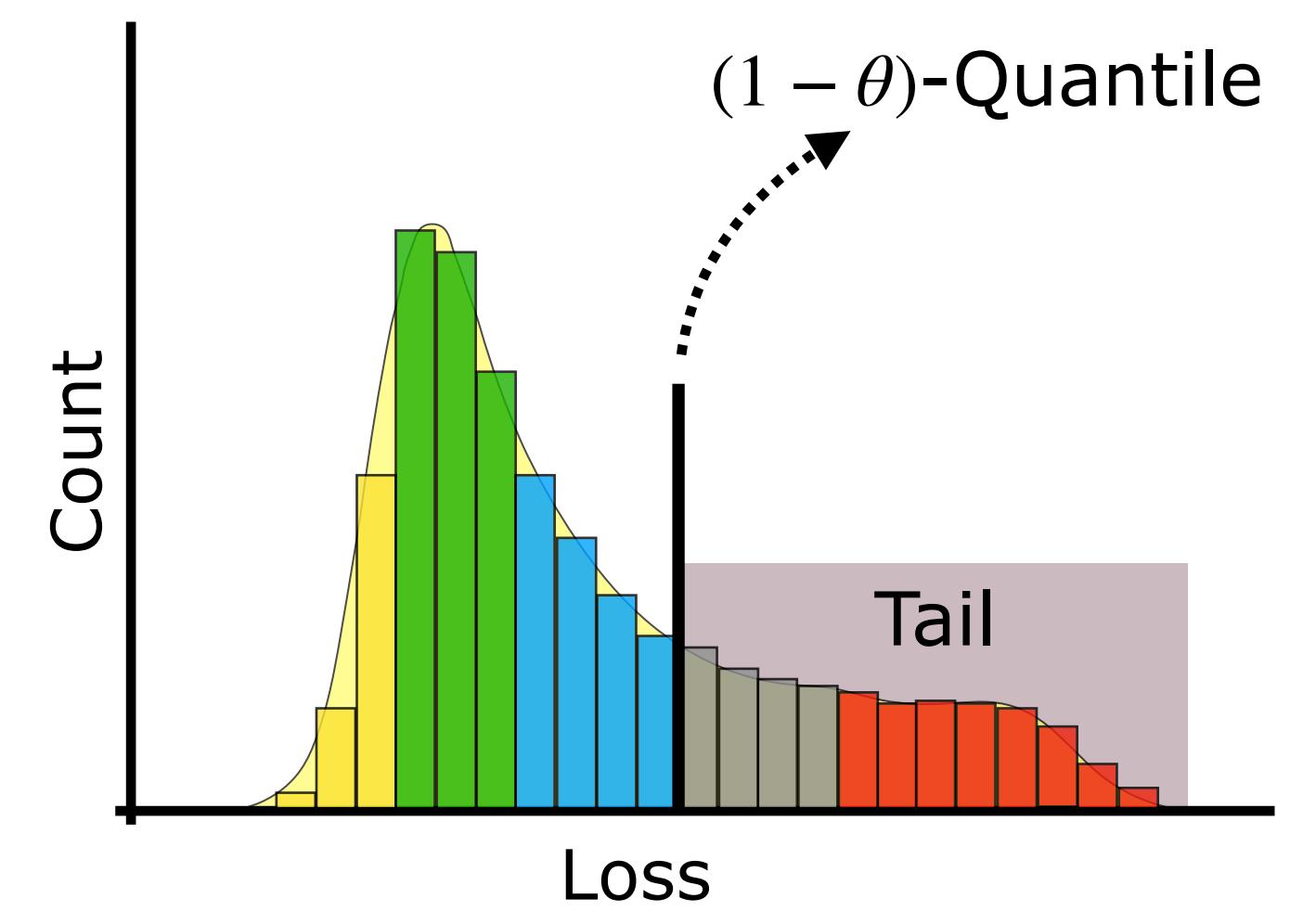
Simplicial-FL Algorithm:

$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

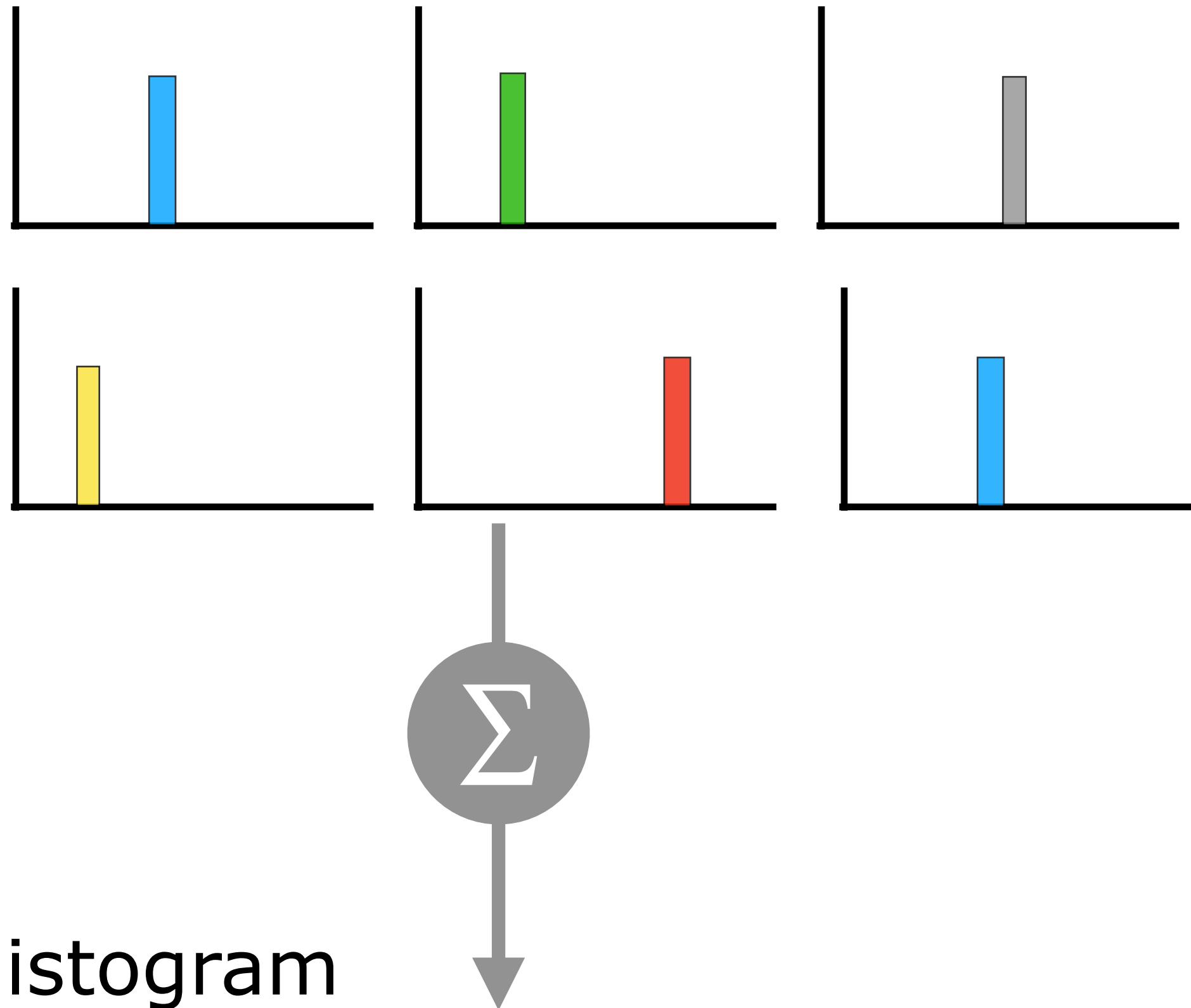
*Step 3 of 3: Aggregate updates contributed by **tail clients** only*



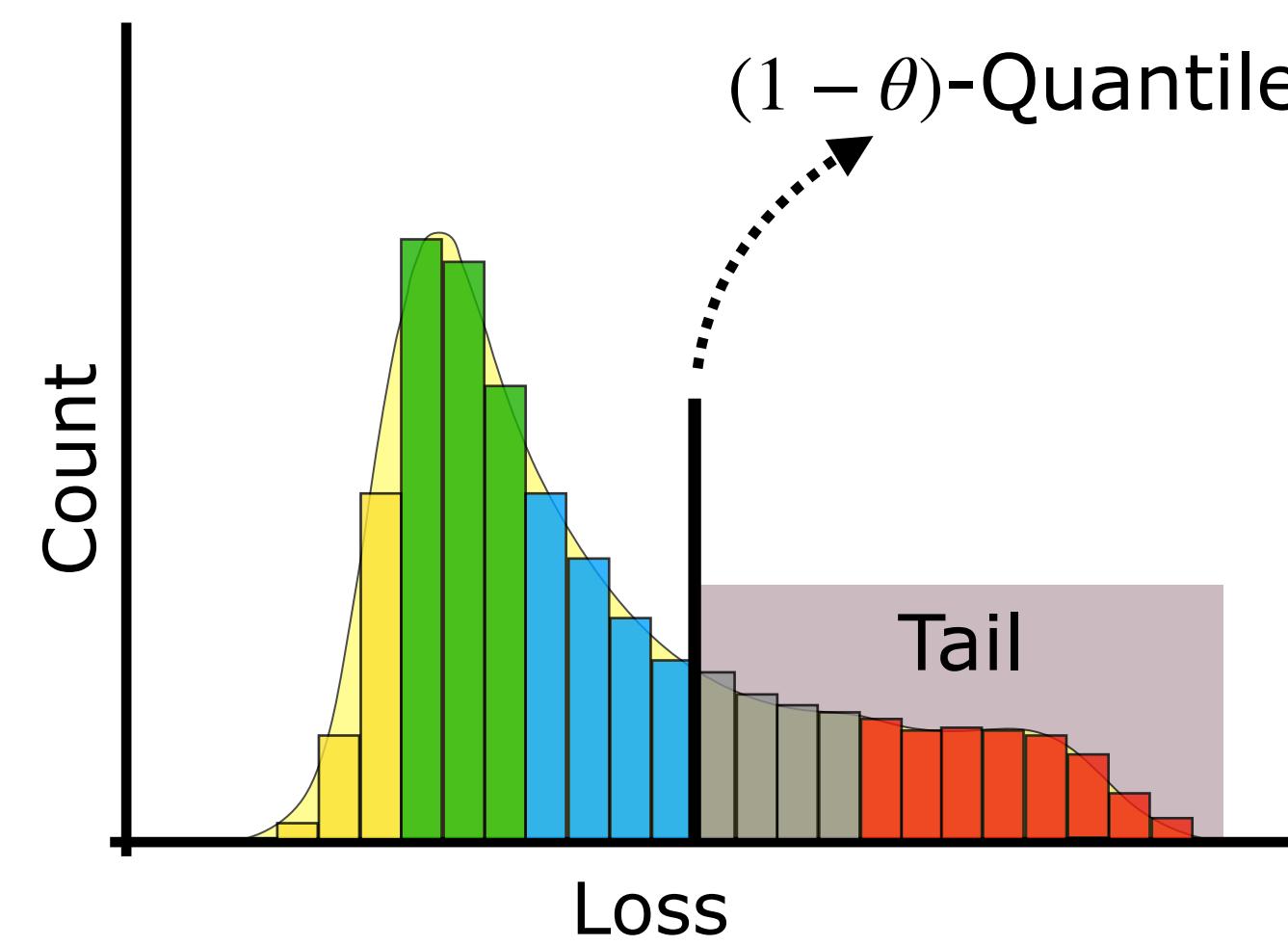




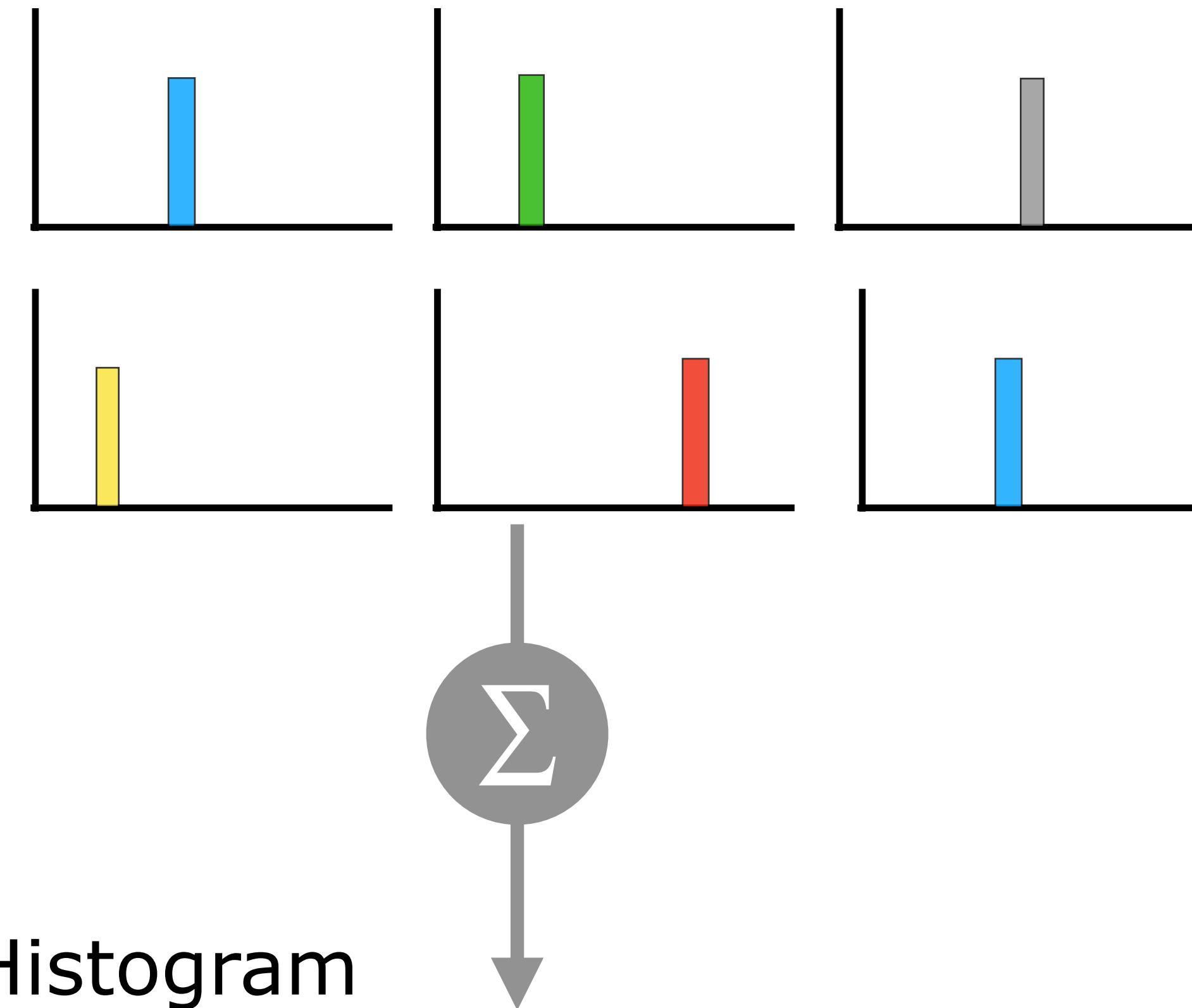
Per-client loss



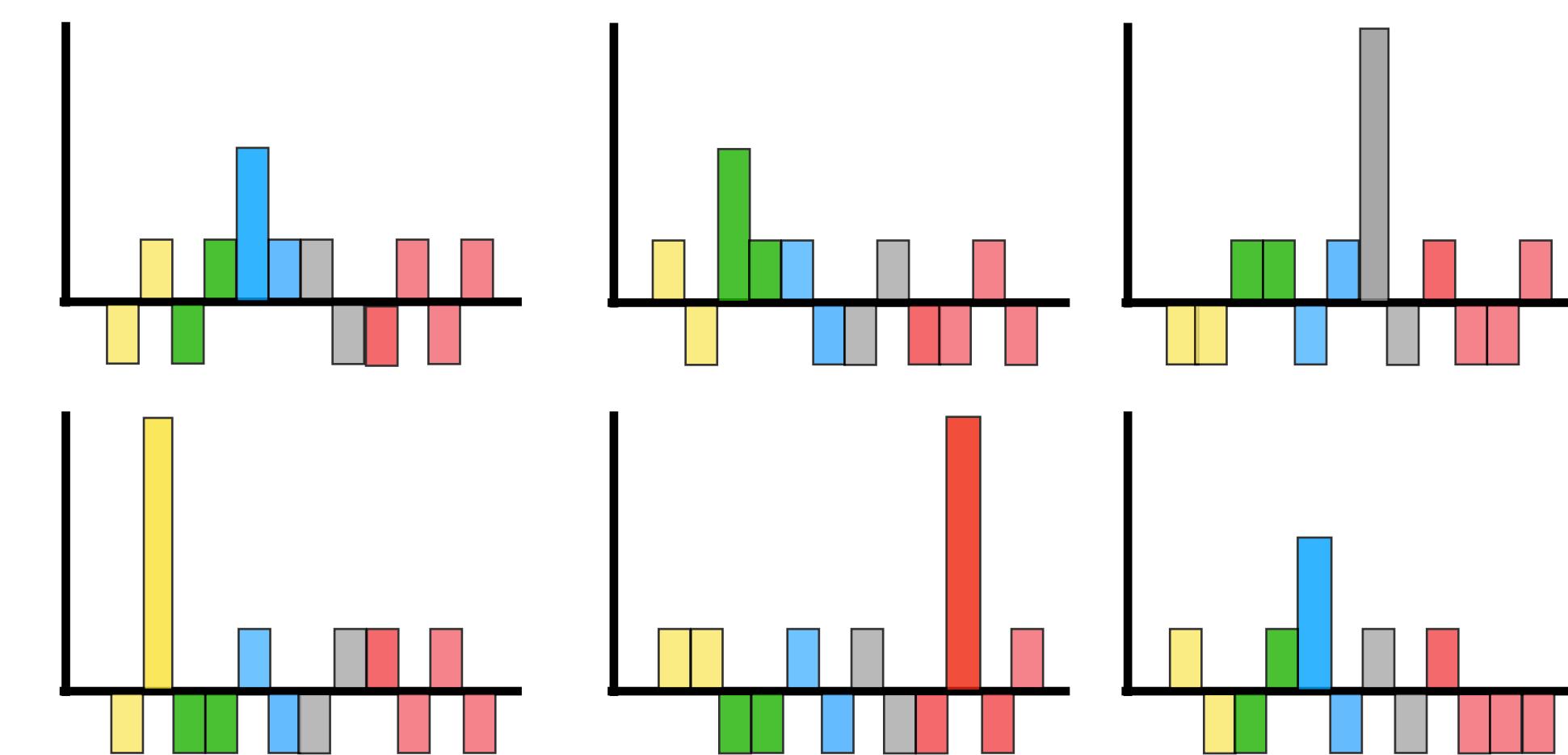
Histogram



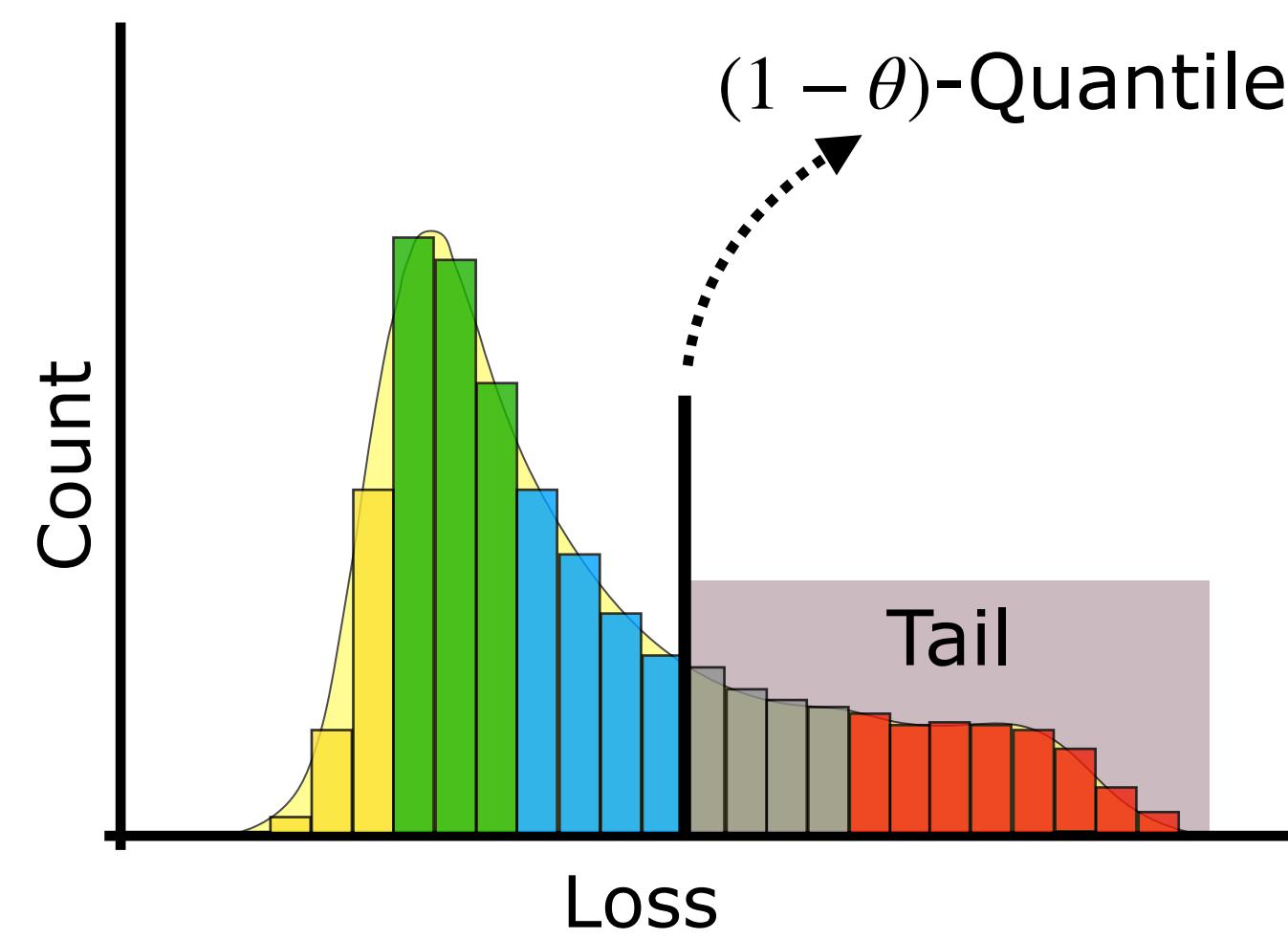
Per-client loss



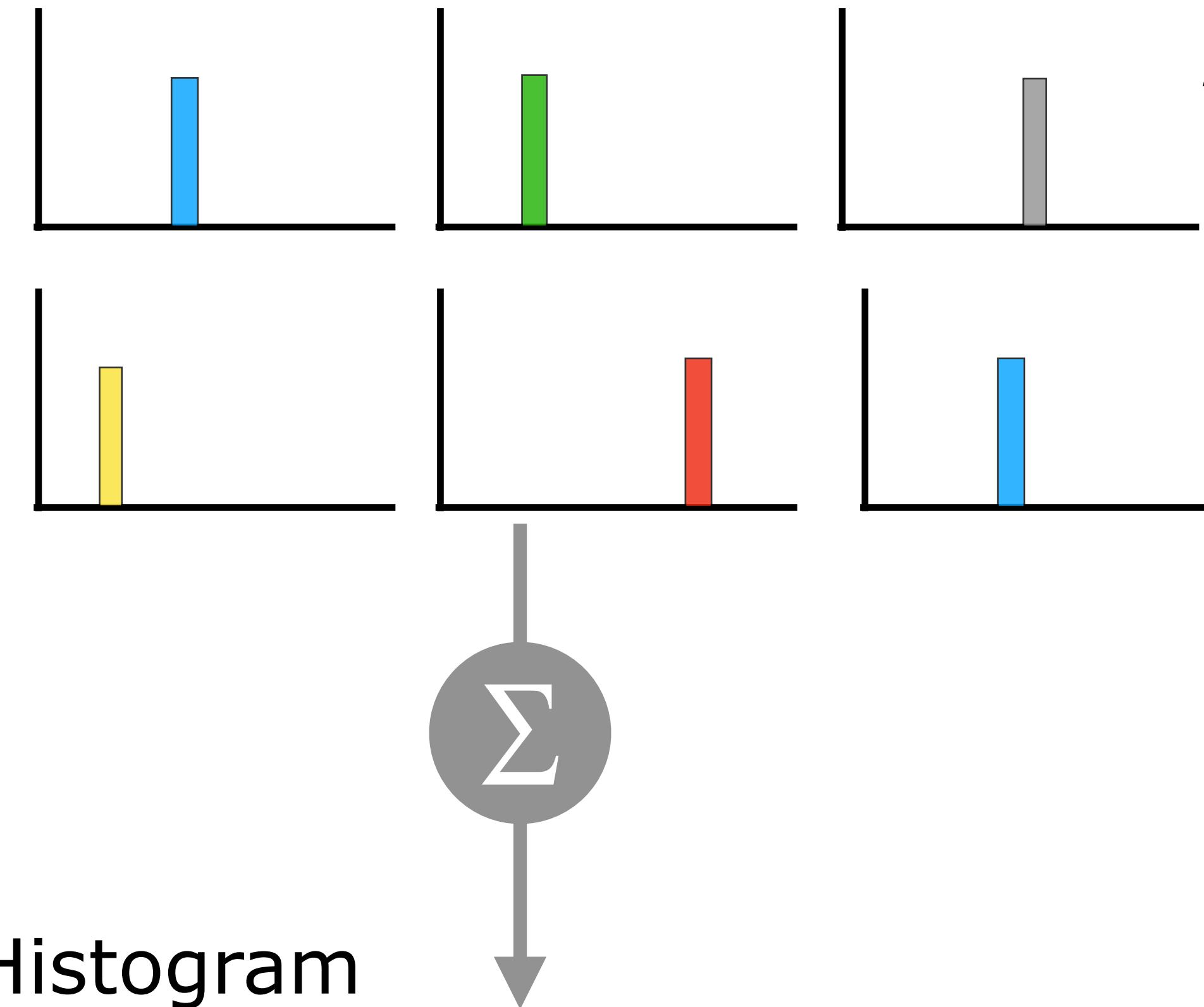
Noisy client loss histogram



Histogram

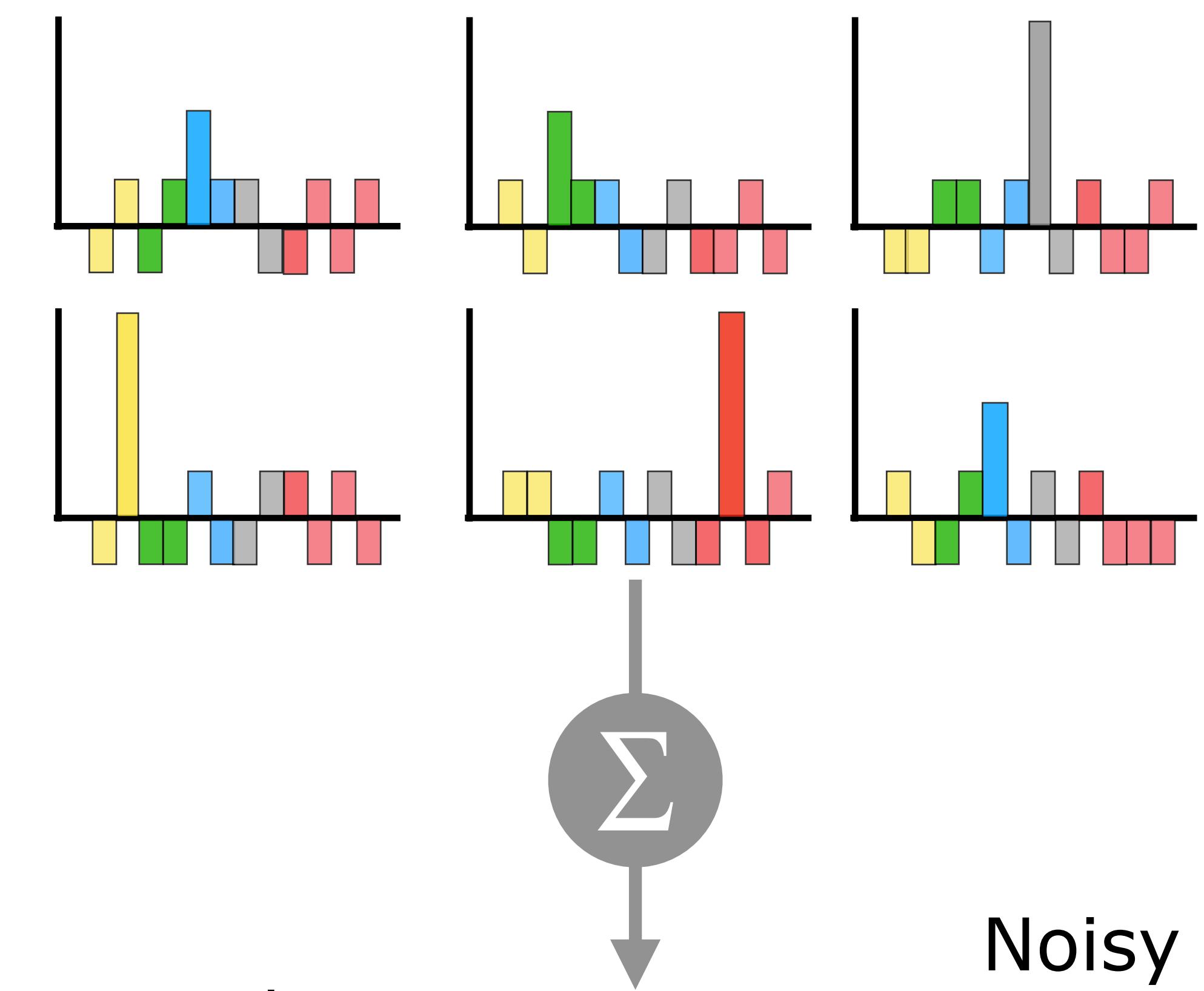


Per-client loss

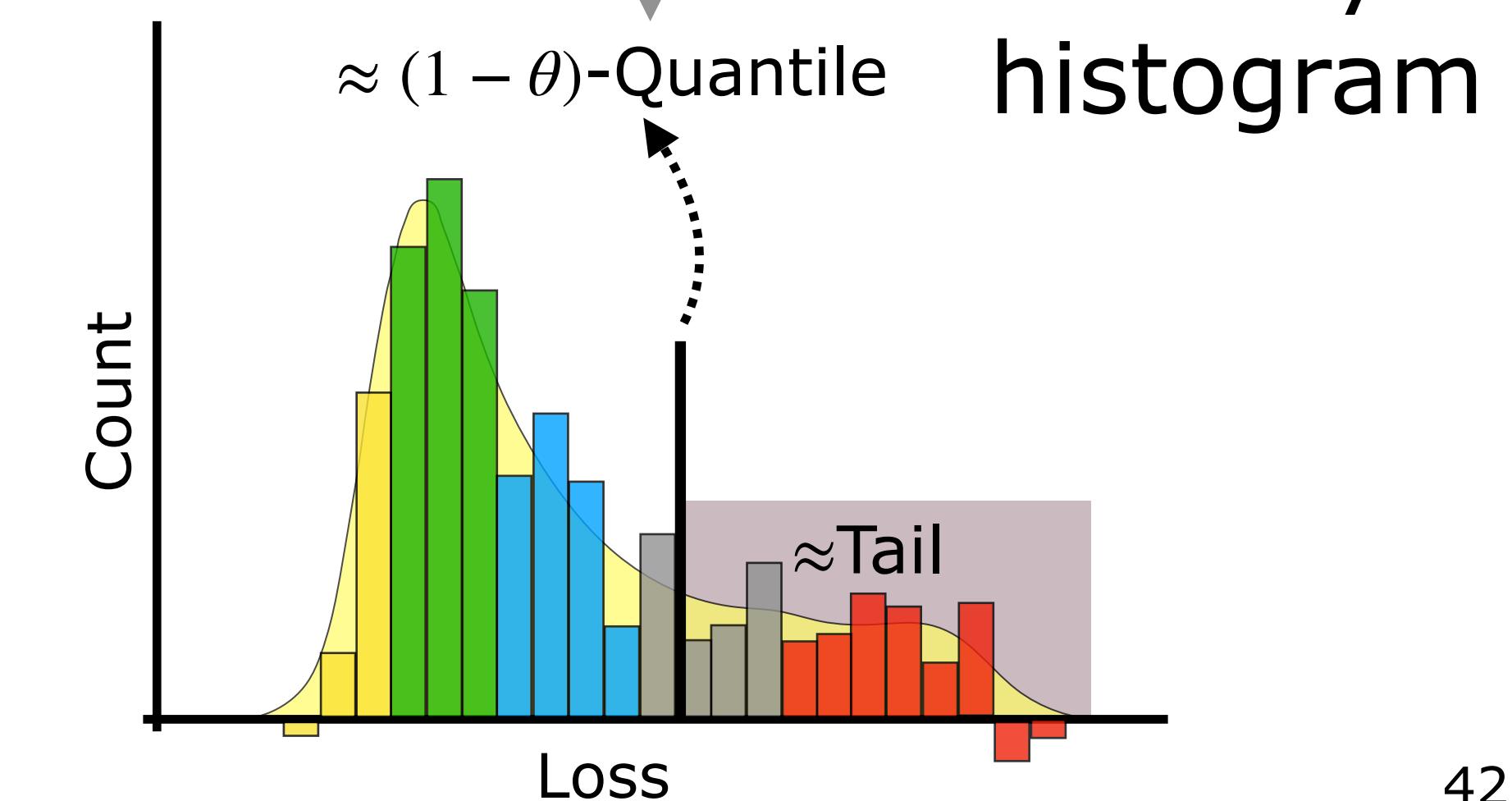
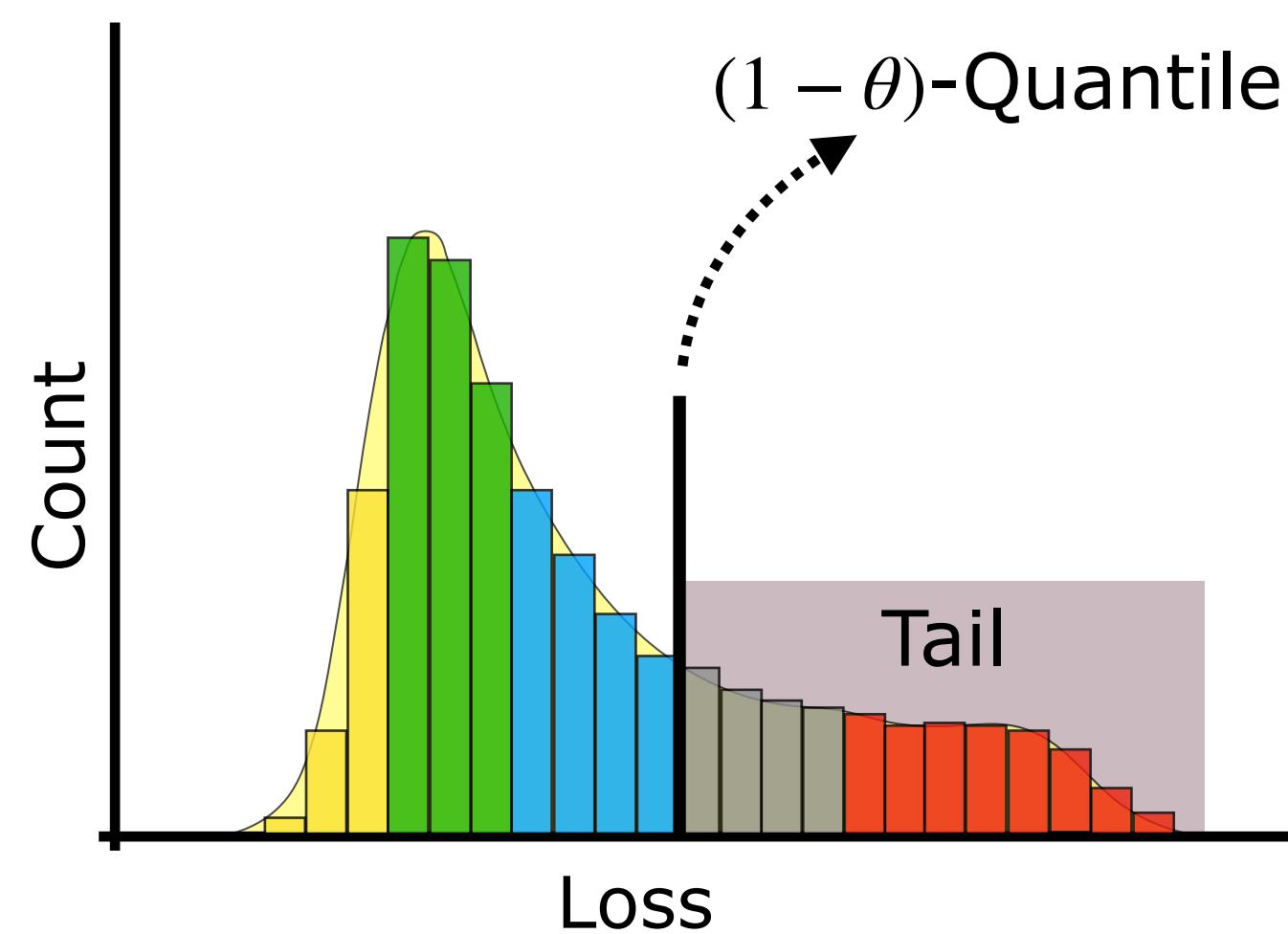


$$h'_i = h_i + \mathcal{N}_{\mathbb{Z}}(0, \sigma^2 I_b)$$

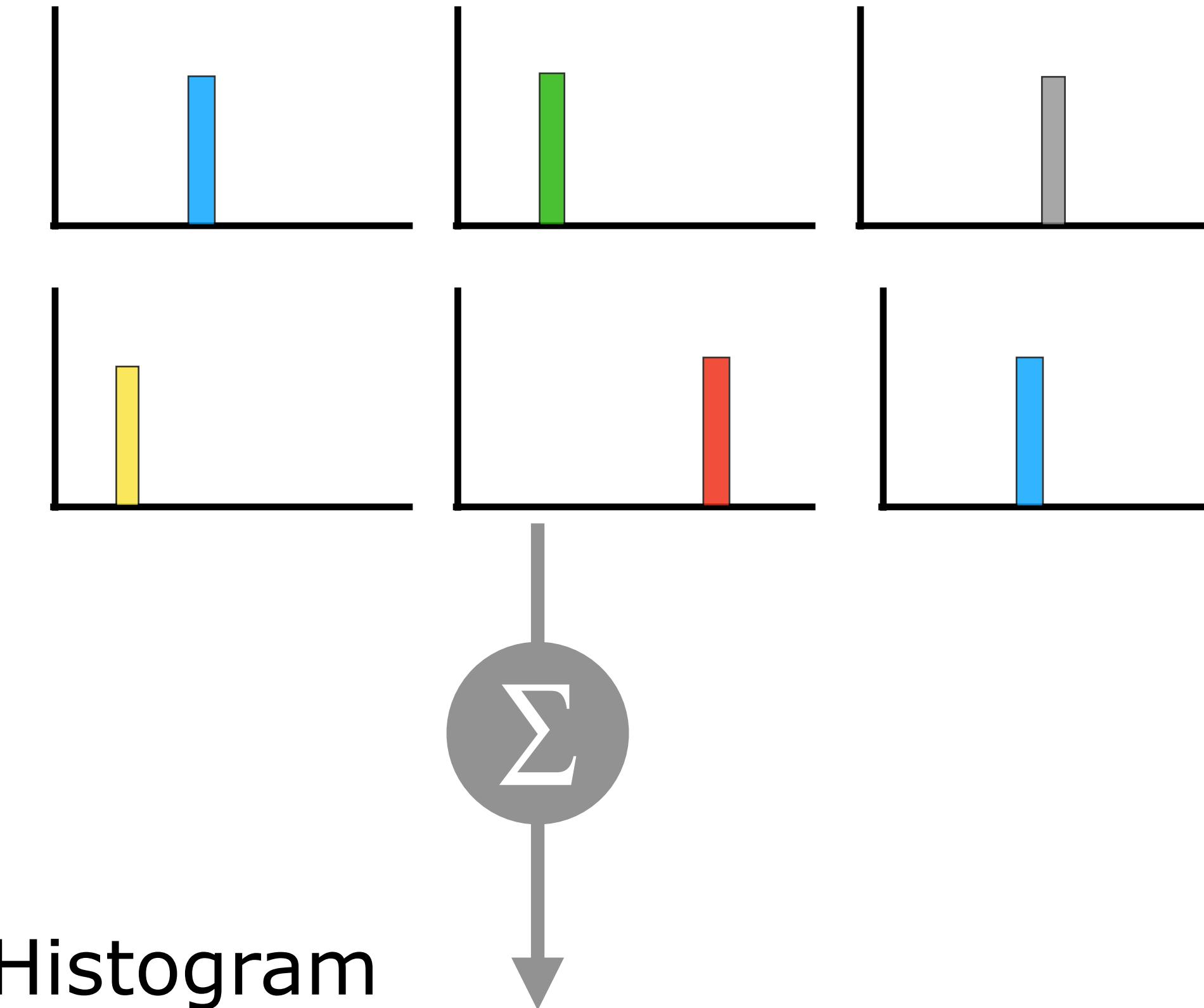
Noisy client loss histogram



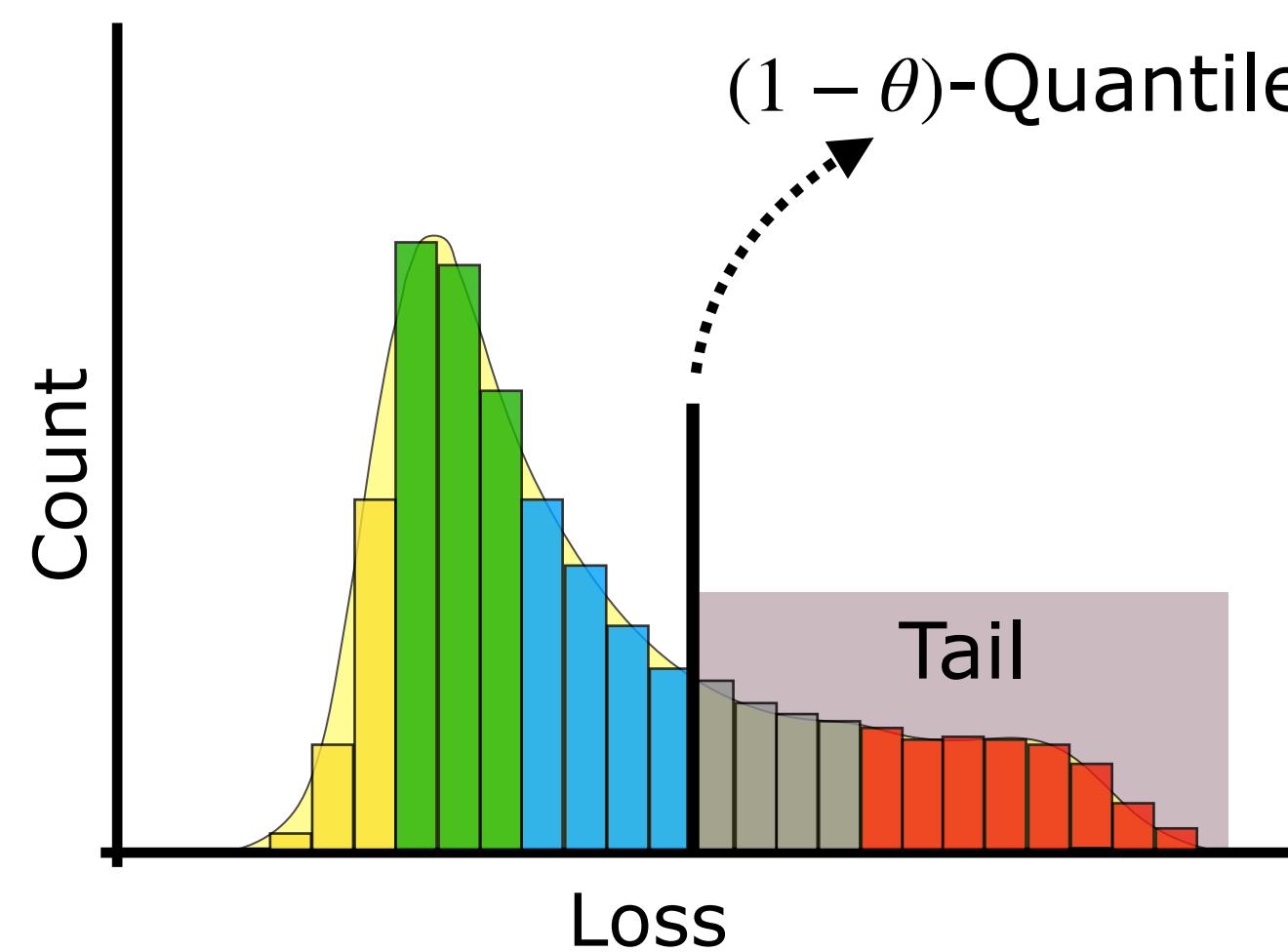
Histogram



Per-client loss



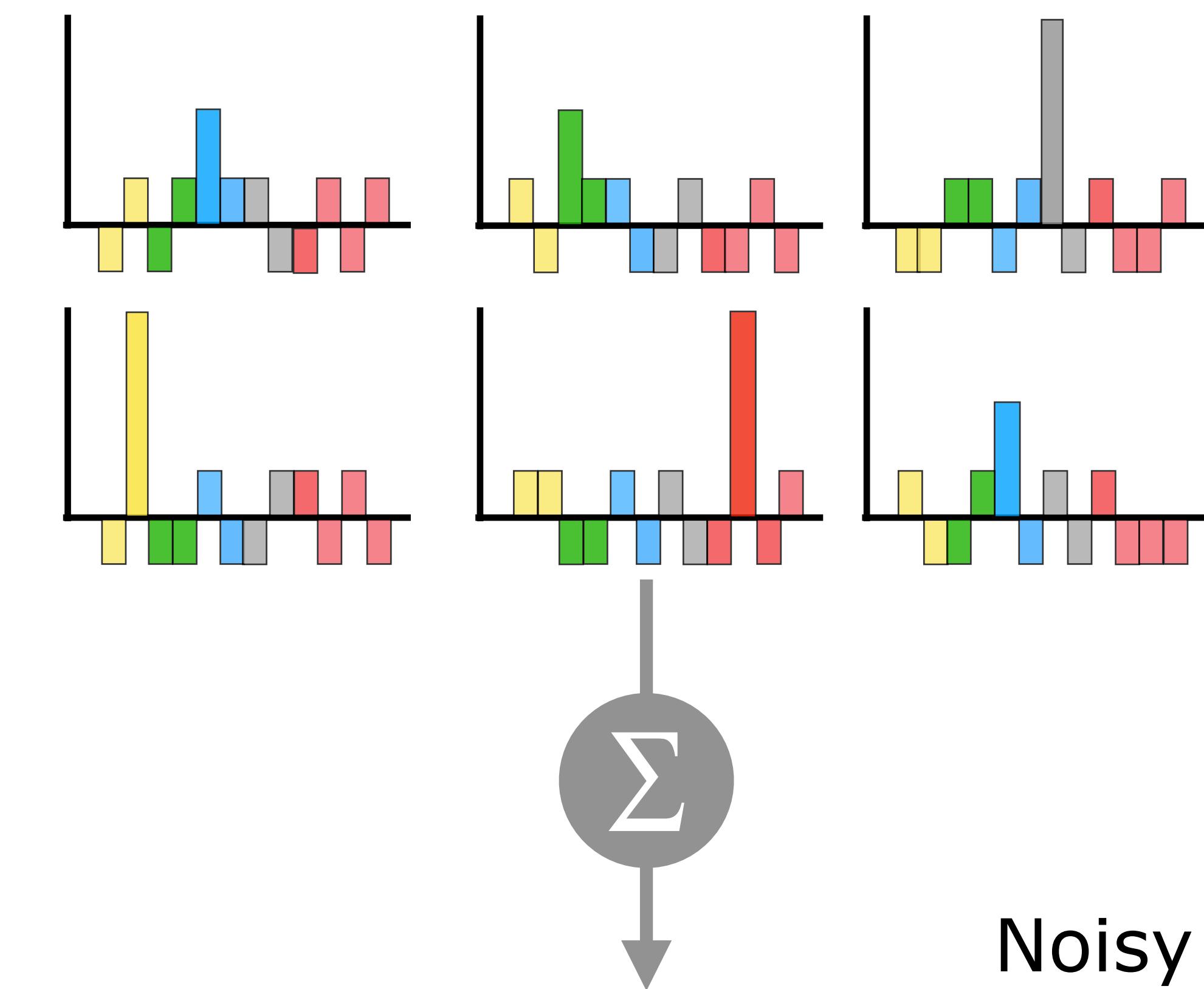
Histogram



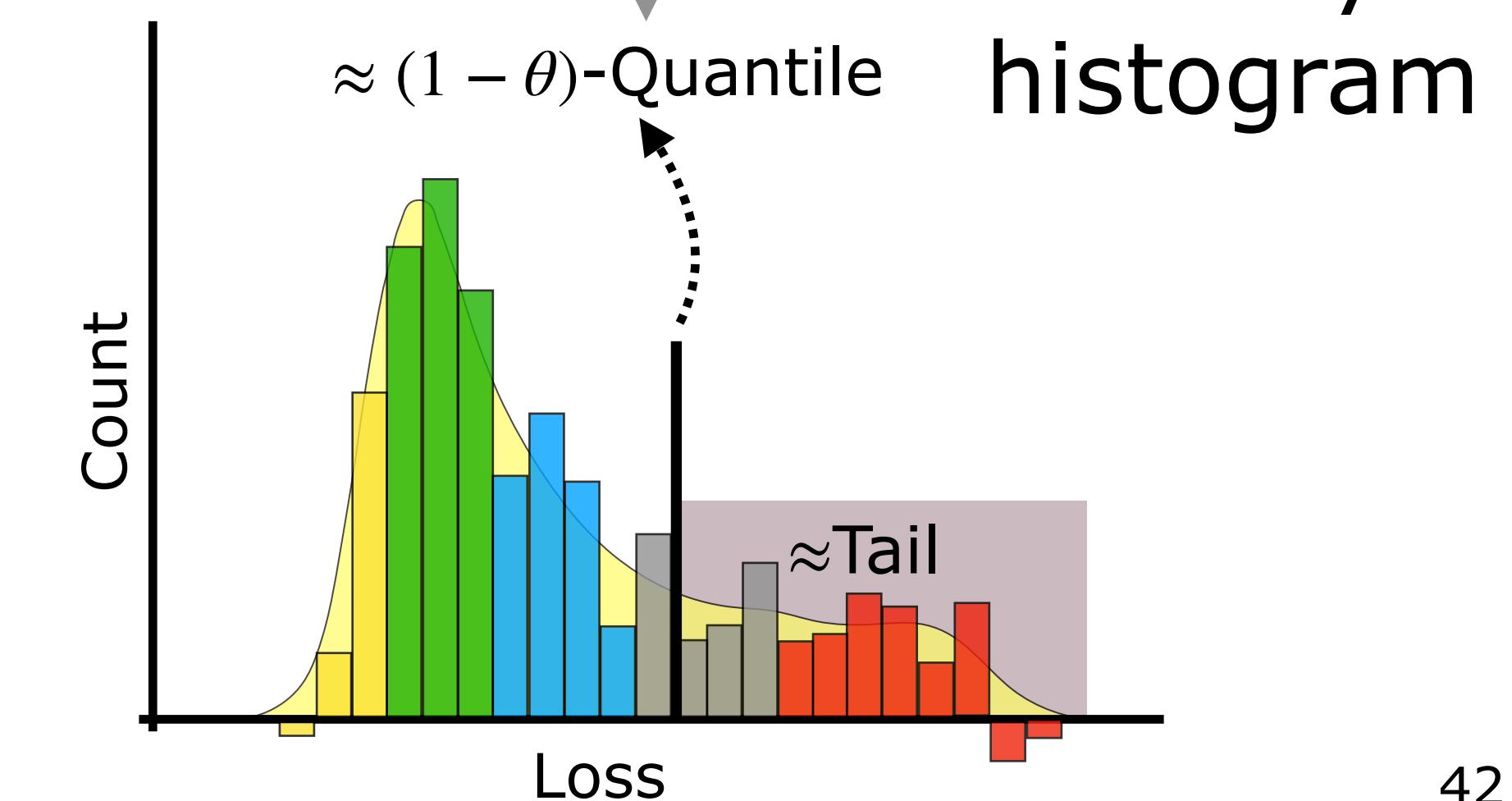
Distributed
discrete Gaussian
mechanism

[Kairouz, Liu, Steinke. ICML (2021)]

Noisy client loss histogram



Noisy
histogram



Proposition [P., Laguel, Malick, Harchaoui]

Fix parameters $\varepsilon, \delta > 0$ and $M \gtrsim m^{3/2}$. If we choose noise scale

$$\sigma \approx \frac{1}{\varepsilon\sqrt{m}} \sqrt{\log \frac{1}{\delta}}$$

m #clients per round
 M modular ring size
 b #bins in the histogram
 (ε, δ) differential privacy parameters

then

- the noisy histogram (and hence all quantiles) are (ε, δ) -differentially private
- w.h.p., the estimated $(1 - \theta)$ -quantile is actually the $(1 - \theta')$ -quantile, with

$$|\theta' - \theta| \lesssim \sqrt{\frac{b\sigma^2}{m}} \approx \frac{1}{\varepsilon m} \sqrt{b \log \frac{1}{\delta}}$$

Total communication cost $\approx bm \log^2 m$

Proposition [P., Laguel, Malick, Harchaoui]

Fix parameters $\varepsilon, \delta > 0$ and $M \gtrsim m^{3/2}$. If we choose noise scale

$$\sigma \approx \frac{1}{\varepsilon\sqrt{m}} \sqrt{\log \frac{1}{\delta}}$$

m #clients per round
 M modular ring size
 b #bins in the histogram
 (ε, δ) differential privacy parameters

then

- the noisy histogram (and hence all quantiles) are (ε, δ) -differentially private
- w.h.p., the estimated $(1 - \theta)$ -quantile is actually the $(1 - \theta')$ -quantile, with

$$|\theta' - \theta| \lesssim \sqrt{\frac{b\sigma^2}{m}} \approx \frac{1}{\varepsilon m} \sqrt{b \log \frac{1}{\delta}}$$

Total communication cost $\approx bm \log^2 m$

Proposition [P., Laguel, Malick, Harchaoui]

Fix parameters $\varepsilon, \delta > 0$ and $M \gtrsim m^{3/2}$. If we choose noise scale

$$\sigma \approx \frac{1}{\varepsilon\sqrt{m}} \sqrt{\log \frac{1}{\delta}}$$

m #clients per round
 M modular ring size
 b #bins in the histogram
 (ε, δ) differential privacy parameters

then

- the noisy histogram (and hence all quantiles) are (ε, δ) -differentially private
- w.h.p., the estimated $(1 - \theta)$ -quantile is actually the $(1 - \theta')$ -quantile, with

$$|\theta' - \theta| \lesssim \sqrt{\frac{b\sigma^2}{m}} \approx \frac{1}{\varepsilon m} \sqrt{b \log \frac{1}{\delta}}$$

Total communication cost $\approx bm \log^2 m$

Summary:

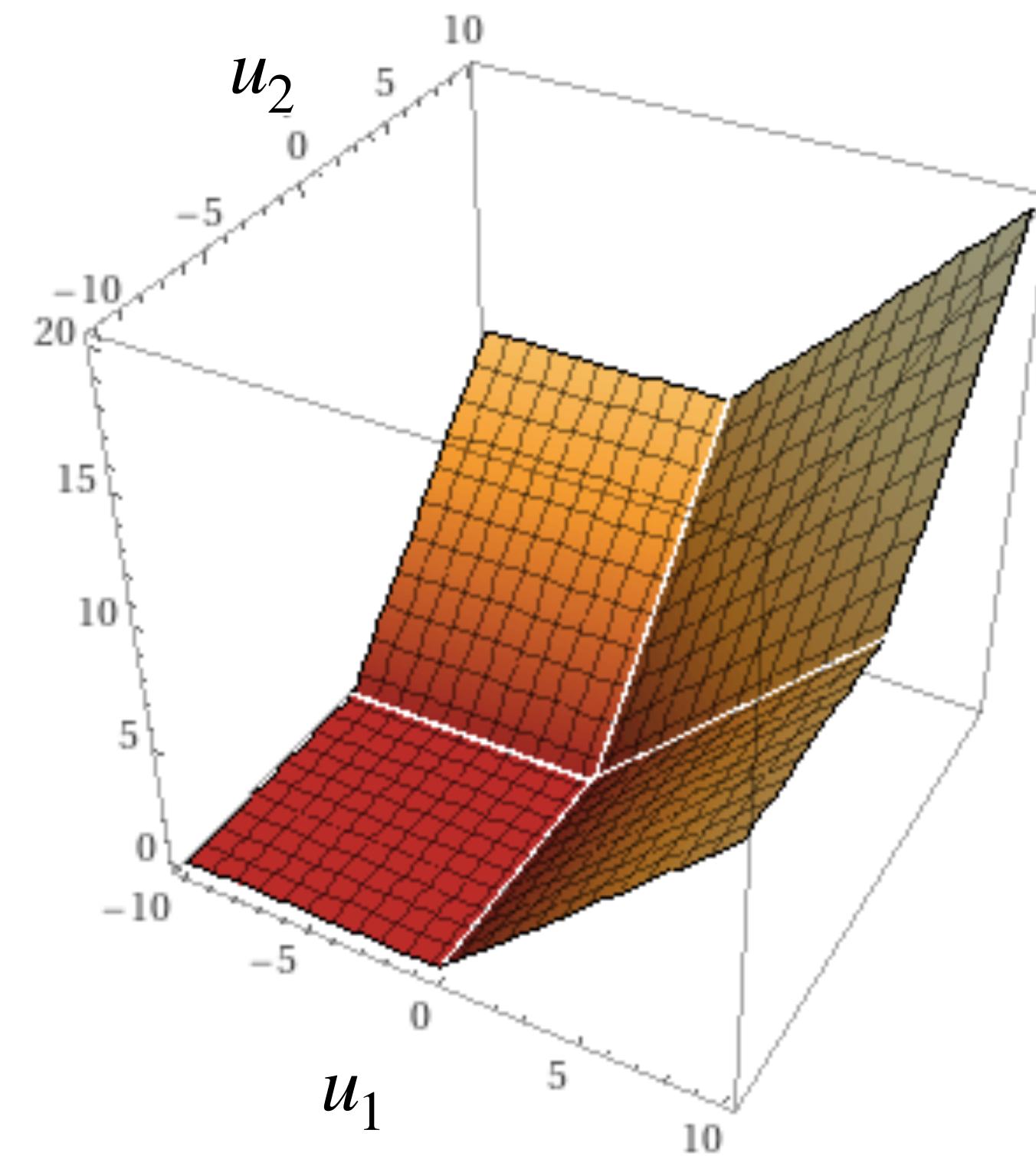
Simplicial-FL algorithm requires **2** secure summations per update

Convergence analysis (non-convex)



Challenge #1:

The superquantile is non-smooth



plot of $h(u_1, u_2) = \mathbb{S}_{1/2}(u_1, u_2, 0, 0)$

Nonsmooth: The subdifferential has a tractable form

$$\partial F_\theta(w) \ni \sum_{i=1}^n \pi_i^\star \nabla F_i(w) \quad \text{where} \quad \pi_i^\star \propto \mathbb{I}(F_i(w) \geq Q_\theta(F_1(w), \dots, F_n(w)))$$

assuming θ_n
is an integer

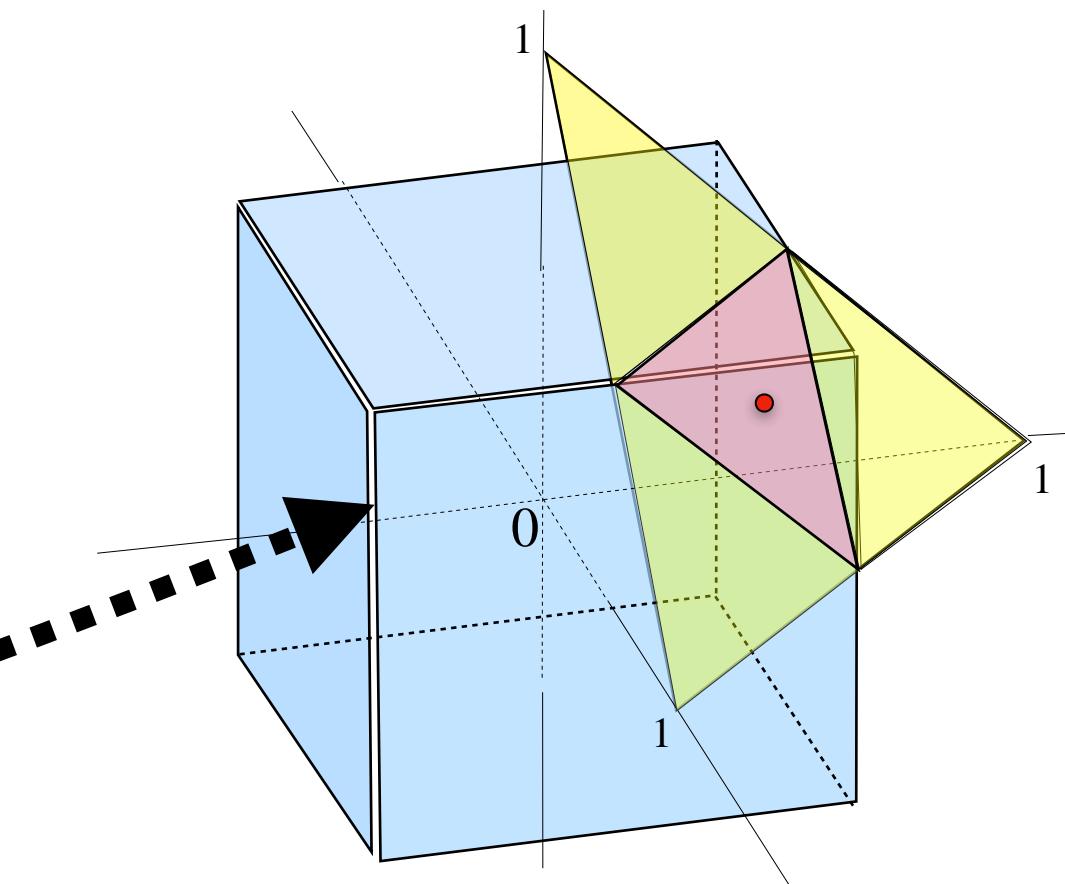
Nonsmooth: The subdifferential has a tractable form

$$\partial F_\theta(w) \ni \sum_{i=1}^n \pi_i^\star \nabla F_i(w) \quad \text{where} \quad \pi_i^\star \propto \mathbb{I}(F_i(w) \geq Q_\theta(F_1(w), \dots, F_n(w)))$$

assuming θ_n
is an integer

Proof Chain rule \implies subdifferential holds with

$$\pi^\star \in \arg \max_{\pi \in \mathcal{P}_\theta} \sum_i \pi_i F_i(w)$$



Alternate form of π^\star comes from the continuous knapsack problem

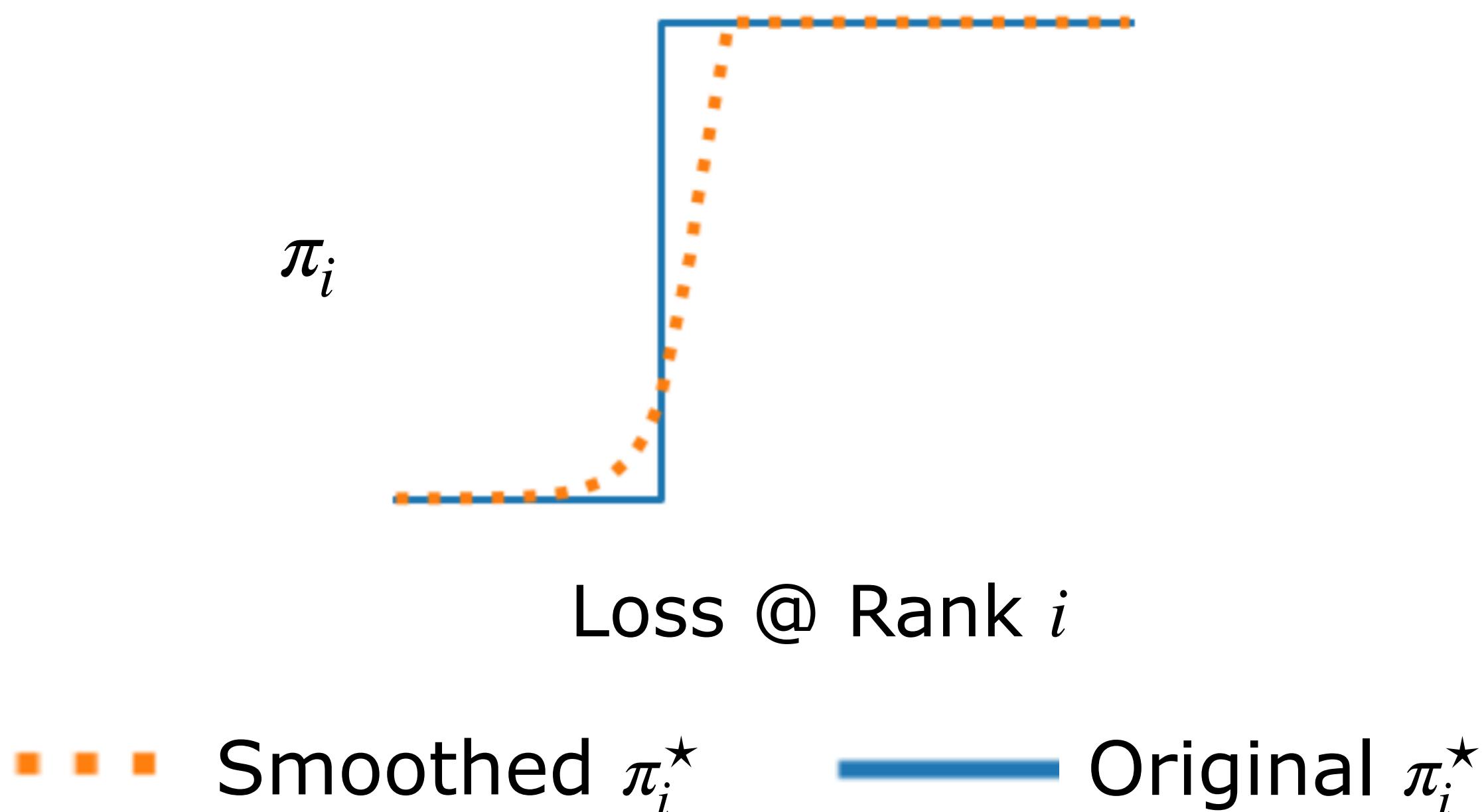
[Dantzig, ORIJ (1957)]

Nonsmooth: The subdifferential has a tractable form

$$\partial F_\theta(w) \ni \sum_{i=1}^n \pi_i^\star \nabla F_i(w) \quad \text{where} \quad \pi_i^\star \propto \mathbb{I}(F_i(w) \geq Q_\theta(F_1(w), \dots, F_n(w)))$$

assuming θ_n
is an integer

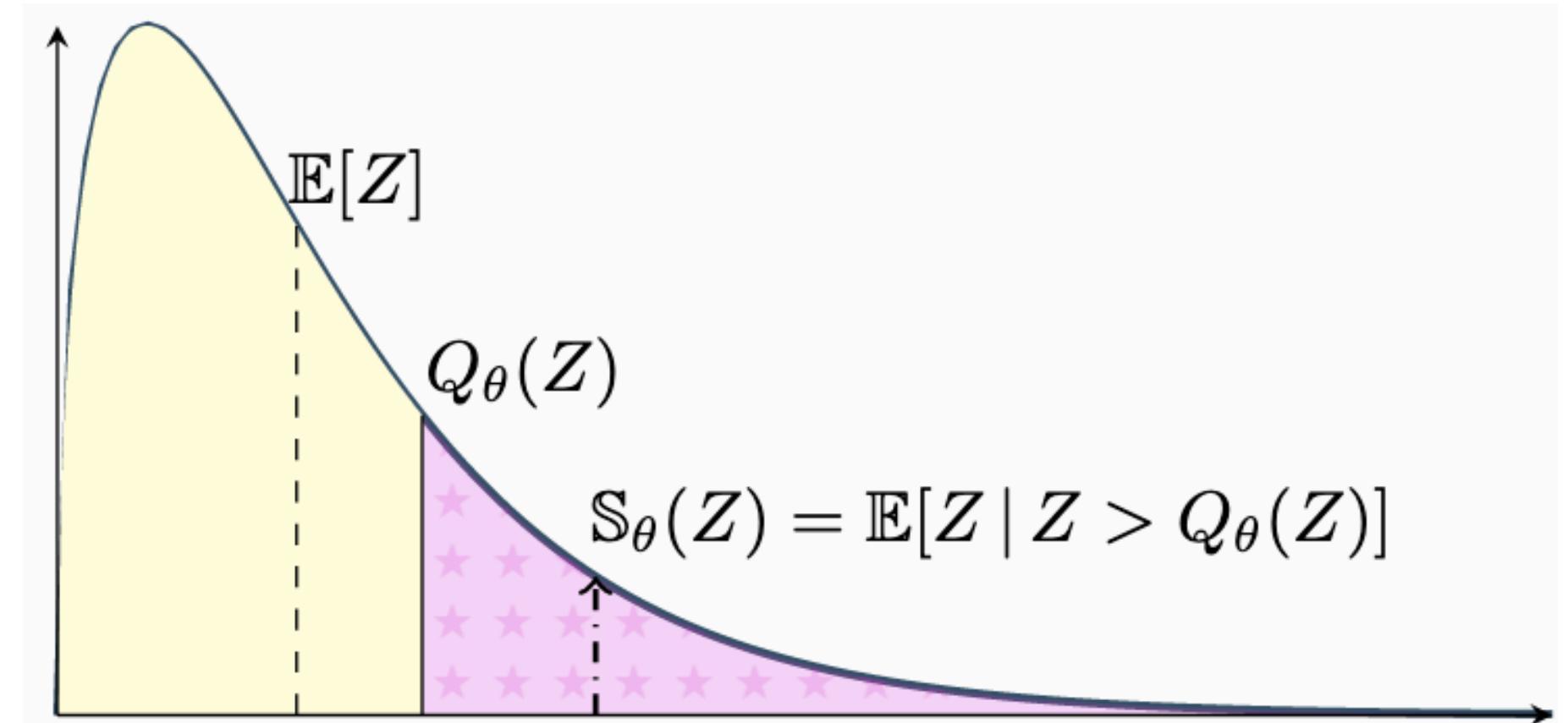
Other option: Use smoothing



[Nesterov. Math. Prog. (2005),
Beck & Teboulle. SIAM J. Optim. (2012),
P., Roulet, Kakade, Harchaoui. NeurIPS (2018),
Laguel, P., Malick, Harchaoui. SVAA (2021)]

Challenge #2

The superquantile is *nonlinear*
⇒ unbiased stochastic gradients not possible

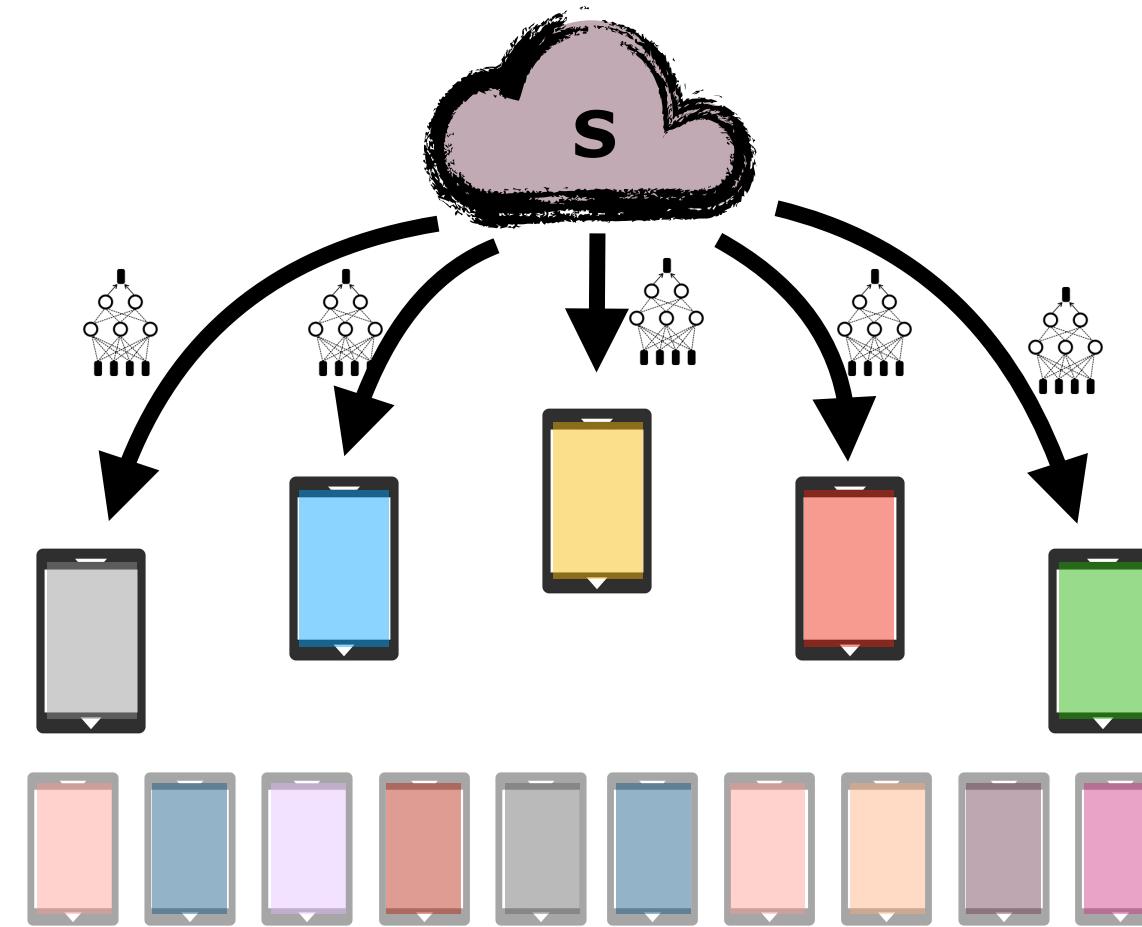


For i.i.d. copies Z_1, \dots, Z_m of Z , we have

$$\mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m Z_i \right] = \mathbb{E}[Z] \quad \text{but} \quad \mathbb{E} \left[\mathbb{S}_\theta(Z_1, \dots, Z_m) \right] \neq \mathbb{S}_\theta(Z)$$

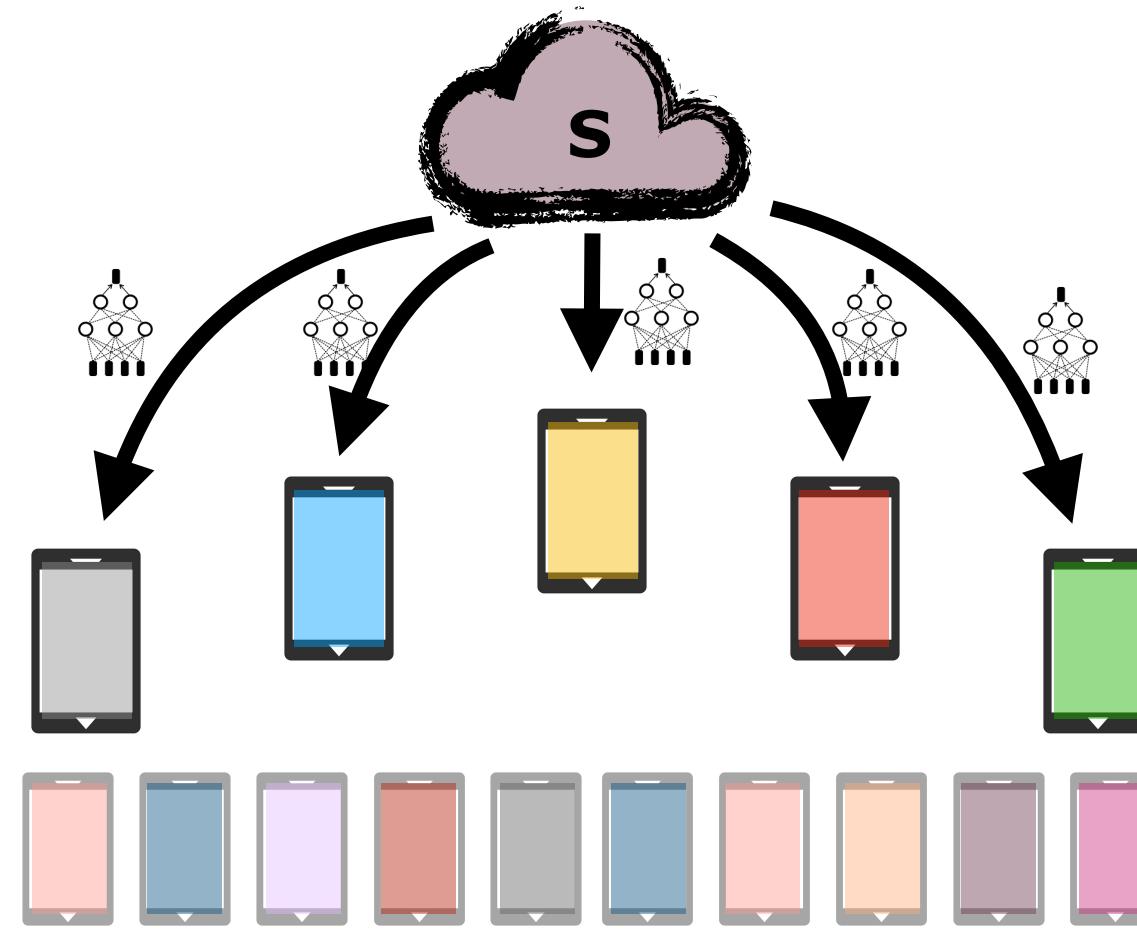
Nonlinear: We minimize a close surrogate

$$\bar{F}_\theta(w) = \mathbb{E}_{S: |S|=m} \left[\mathbb{S}_\theta \left((F_i(w) : i \in S) \right) \right]$$



Nonlinear: We minimize a close surrogate

$$\bar{F}_\theta(w) = \mathbb{E}_{S: |S|=m} \left[\mathbb{S}_\theta \left((F_i(w) : i \in S) \right) \right]$$



The surrogate is uniformly close for bounded losses:

For i.i.d. copies Z_1, \dots, Z_m of Z with $|Z| \leq B$ a.s., we have

$$\left| \mathbb{E}[\mathbb{S}_\theta(Z_1, \dots, Z_m)] - \mathbb{S}_\theta(Z) \right| \leq \frac{B}{\sqrt{\theta m}}$$

$$\text{Var}[\mathbb{S}_\theta(Z_1, \dots, Z_m)] \leq \frac{B^2}{\theta m}$$

Theorem [P., Laguel, Malick, Harchaoui]

Suppose each F_i is L -smooth and G -Lipschitz.

Then, Simplicial-FL satisfies the convergence guarantee:

$$\mathbb{E} \left\| \nabla \Phi_{\theta}^{2L}(w_t) \right\|^2 \leq \sqrt{\frac{\Delta_0 L G^2}{t}} + (1 - \tau)^{1/3} \left(\frac{\Delta_0 L G}{t} \right)^{2/3} + \frac{\Delta_0 L}{t}$$

t : #comm. rounds
 τ : #local update steps
 Δ_0 : initial error

$$\Phi_{\theta}^{\mu}(w) = \inf_z \left\{ \bar{F}_{\theta}(z) + \frac{\mu}{2} \|z - w\|^2 \right\} \quad \leftarrow \quad \text{Moreau envelope of } \bar{F}_{\theta} \mid \text{well defined for } \mu > L$$

Theorem [P., Laguel, Malick, Harchaoui]

Suppose **each** F_i **is L -smooth and G -Lipschitz.**

Then, Simplicial-FL satisfies the convergence guarantee:

$$\mathbb{E} \left\| \nabla \Phi_{\theta}^{2L}(w_t) \right\|^2 \leq \sqrt{\frac{\Delta_0 L G^2}{t}} + (1 - \tau)^{1/3} \left(\frac{\Delta_0 L G}{t} \right)^{2/3} + \frac{\Delta_0 L}{t}$$

t : #comm. rounds
 τ : #local update steps
 Δ_0 : initial error

$$\Phi_{\theta}^{\mu}(w) = \inf_z \left\{ \bar{F}_{\theta}(z) + \frac{\mu}{2} \|z - w\|^2 \right\} \quad \leftarrow \quad \text{Moreau envelope of } \bar{F}_{\theta} \mid \text{well defined for } \mu > L$$

Theorem [P., Laguel, Malick, Harchaoui]

Suppose each F_i is L -smooth and G -Lipschitz.

Then, Simplicial-FL satisfies the convergence guarantee:

$$\mathbb{E} \left\| \nabla \Phi_{\theta}^{2L}(w_t) \right\|^2 \leq \sqrt{\frac{\Delta_0 L G^2}{t}} + (1 - \tau)^{1/3} \left(\frac{\Delta_0 L G}{t} \right)^{2/3} + \frac{\Delta_0 L}{t}$$

t : #comm. rounds
 τ : #local update steps
 Δ_0 : initial error

$$\Phi_{\theta}^{\mu}(w) = \inf_z \left\{ \bar{F}_{\theta}(z) + \frac{\mu}{2} \|z - w\|^2 \right\} \quad \leftarrow \quad \text{Moreau envelope of } \bar{F}_{\theta} \mid \text{well defined for } \mu > L$$

Theorem [P., Laguel, Malick, Harchaoui]

Suppose each F_i is L -smooth and G -Lipschitz.

Then, Simplicial-FL satisfies the convergence guarantee:

$$\mathbb{E} \left\| \nabla \Phi_{\theta}^{2L}(w_t) \right\|^2 \leq \sqrt{\frac{\Delta_0 L G^2}{t}} + (1 - \tau)^{1/3} \left(\frac{\Delta_0 L G}{t} \right)^{2/3} + \frac{\Delta_0 L}{t}$$

t : #comm. rounds
 τ : #local update steps
 Δ_0 : initial error

$$\Phi_{\theta}^{\mu}(w) = \inf_z \left\{ \bar{F}_{\theta}(z) + \frac{\mu}{2} \|z - w\|^2 \right\} \quad \leftarrow \quad \text{Moreau envelope of } \bar{F}_{\theta} \mid \text{well defined for } \mu > L$$

$$\Phi(w_{t+1}) = \min_z \left\{ \bar{F}_\theta(z) + L\|z - w_{t+1}\|^2 \right\}$$

Definition of the Moreau envelope

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_{t+1}\|^2$$

Plug in a particular choice of z_t to be determined later

$$\Phi(w_{t+1}) = \min_z \left\{ \bar{F}_\theta(z) + L\|z - w_{t+1}\|^2 \right\}$$

Definition of the Moreau envelope

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_{t+1}\|^2$$

Plug in a particular choice of z_t to be determined later

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_t\|^2 - \gamma L(w_t - z_t)^\top g_t + O(\gamma^2)$$

Expand update
 $w_{t+1} = w_t - \gamma g_t$

$$\Phi(w_{t+1}) = \min_z \left\{ \bar{F}_\theta(z) + L\|z - w_{t+1}\|^2 \right\}$$

Definition of the Moreau envelope

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_{t+1}\|^2$$

Plug in a particular choice of z_t to be determined later

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_t\|^2 - \gamma L(w_t - z_t)^\top g_t + O(\gamma^2)$$

Expand update
 $w_{t+1} = w_t - \gamma g_t$

$$= \Phi(w_t) - \gamma \nabla \Phi(w_t)^\top g_t + O(\gamma^2)$$

Choose
 $z_t = \arg \min_z \left\{ \bar{F}_\theta(z) + L\|z - w_t\|^2 \right\}$
so that $\nabla \Phi(w_t) = L(w_t - z_t)$

$$\Phi(w_{t+1}) = \min_z \left\{ \bar{F}_\theta(z) + L\|z - w_{t+1}\|^2 \right\}$$

Definition of the Moreau envelope

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_{t+1}\|^2$$

Plug in a particular choice of z_t to be determined later

$$\leq \bar{F}_\theta(z_t) + L\|z_t - w_t\|^2 - \gamma L(w_t - z_t)^\top g_t + O(\gamma^2)$$

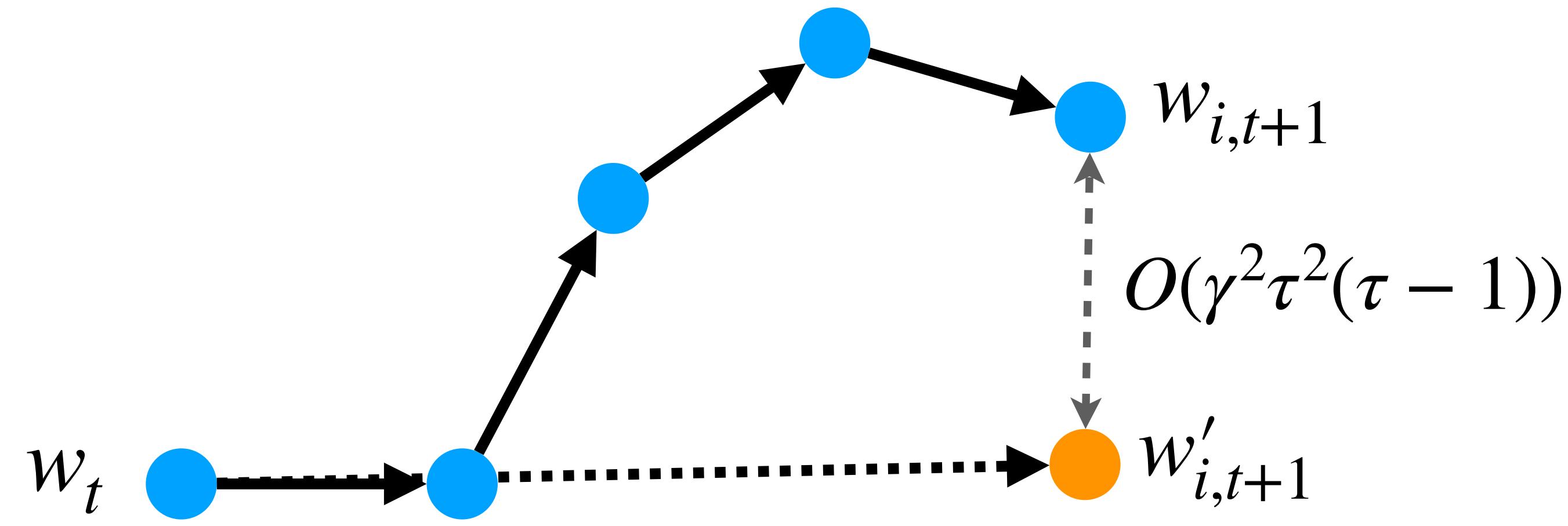
Expand update
 $w_{t+1} = w_t - \gamma g_t$

$$= \Phi(w_t) - \gamma \nabla \Phi(w_t)^\top g_t + O(\gamma^2)$$

Choose
 $z_t = \arg \min_z \left\{ \bar{F}_\theta(z) + L\|z - w_t\|^2 \right\}$
so that $\nabla \Phi(w_t) = L(w_t - z_t)$

If $\mathbb{E}_t[g_t] \approx \nabla \Phi(w_t)$, proof is complete

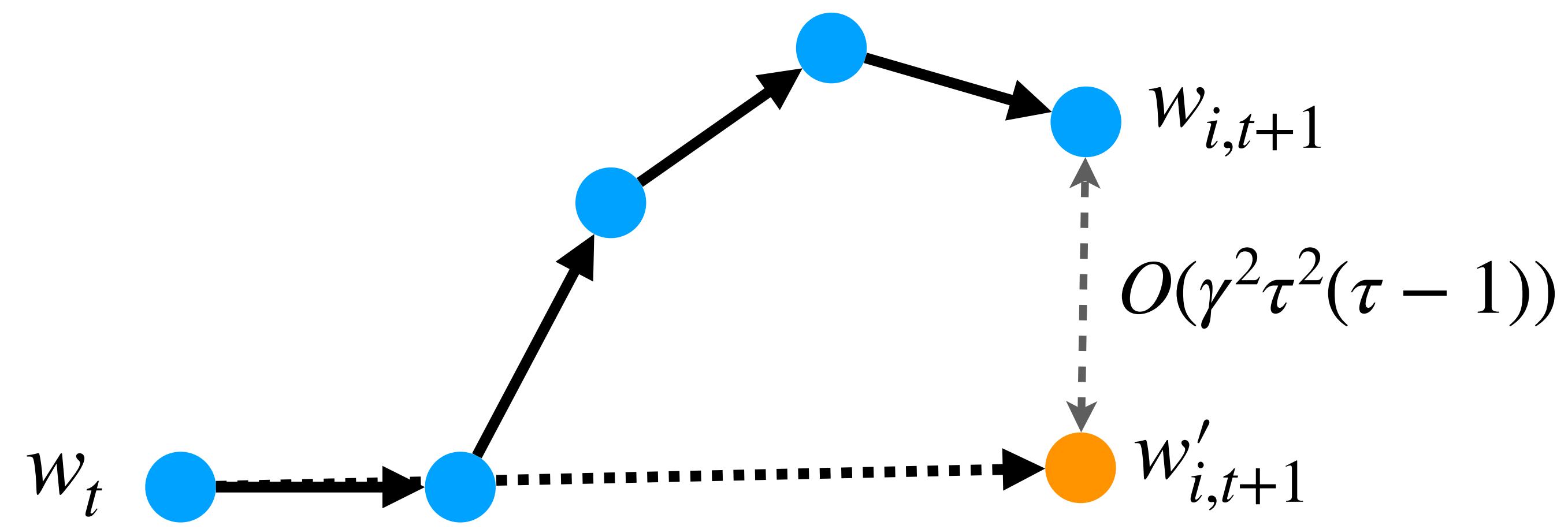
$$\Phi(w_{t+1}) \leq \Phi(w_t) - \gamma \nabla \Phi(w_t)^\top g_t + O(\gamma^2)$$



g_t comes from τ local gradient steps
of step size γ

g'_t comes from one local gradient
steps of step size $\tau\gamma$

$$\Phi(w_{t+1}) \leq \Phi(w_t) - \gamma \nabla \Phi(w_t)^\top g_t + O(\gamma^2)$$



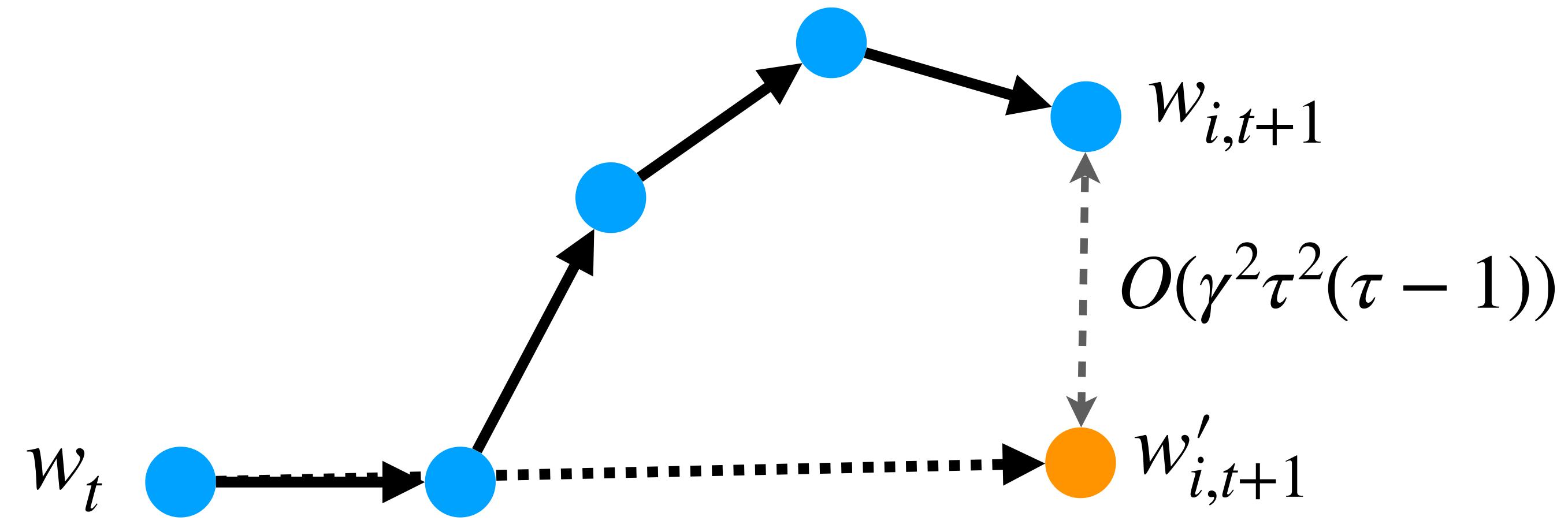
g_t comes from τ local gradient steps
of step size γ

g'_t comes from one local gradient
steps of step size $\tau\gamma$

$$\Phi(w_{t+1}) \leq \Phi(w_t) - \tau\gamma \nabla \Phi(w_t)^\top g'_t + O(\gamma^2)$$

$$\mathbb{E}_t[g'_t] \in \partial \bar{F}_\theta(w_t)$$

$$\Phi(w_{t+1}) \leq \Phi(w_t) - \gamma \nabla \Phi(w_t)^\top g_t + O(\gamma^2)$$



g_t comes from τ local gradient steps
of step size γ

g'_t comes from one local gradient
steps of step size $\tau\gamma$

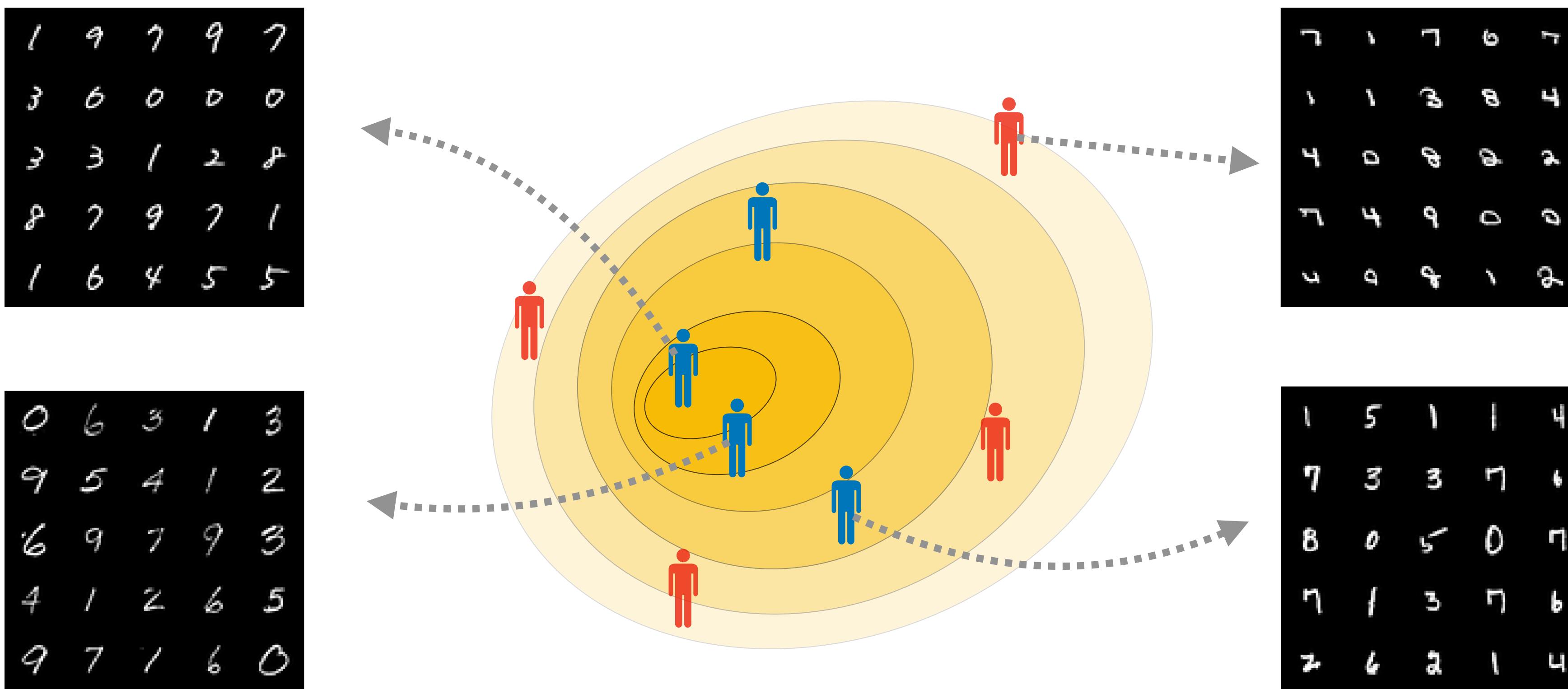
$$\Phi(w_{t+1}) \leq \Phi(w_t) - \tau\gamma \nabla \Phi(w_t)^\top g'_t + O(\gamma^2)$$

$$\mathbb{E}_t[g'_t] \in \partial \bar{F}_\theta(w_t)$$

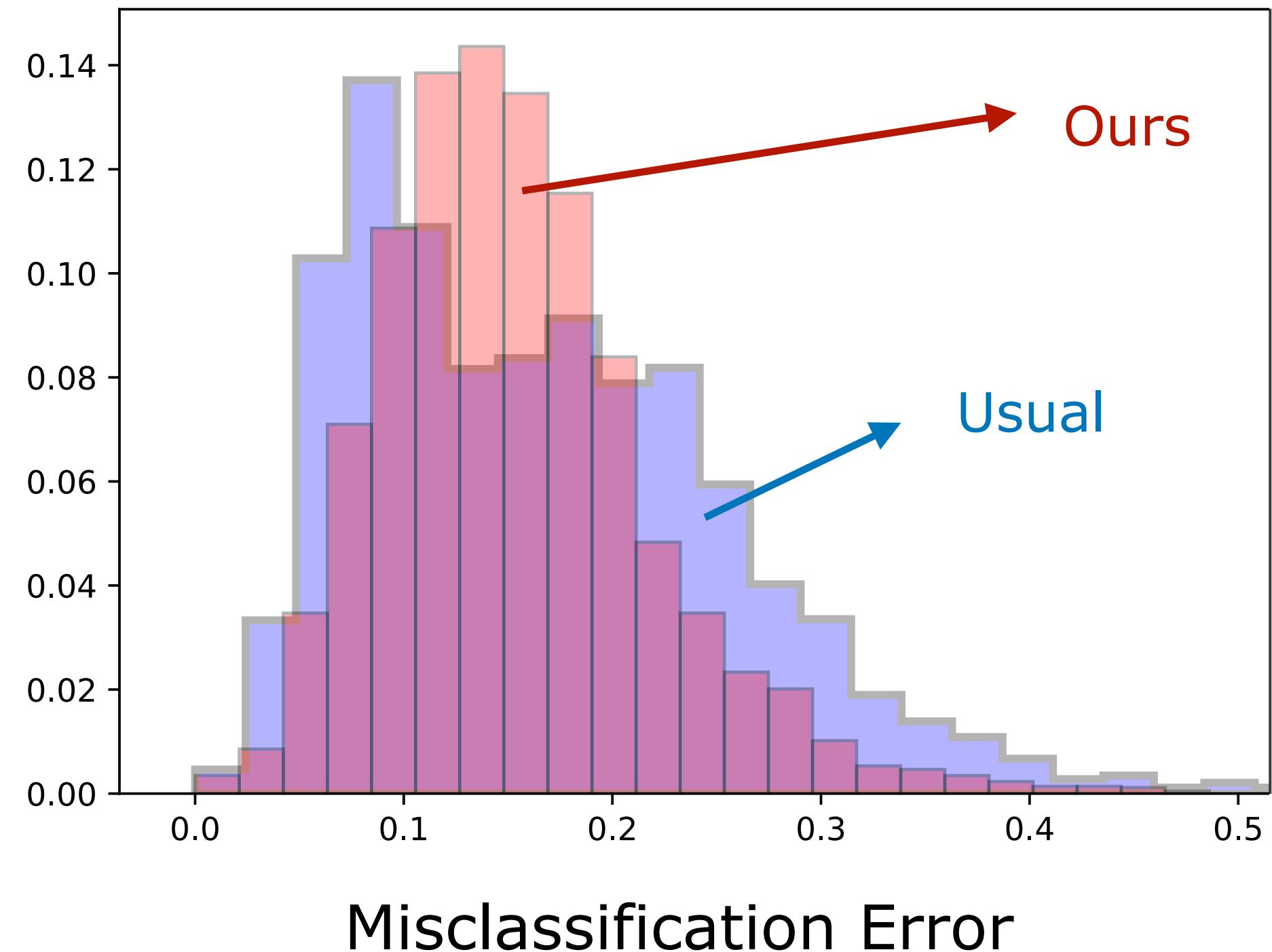
$$\nabla \Phi(w_t)^\top \mathbb{E}_t[g'_t] \geq \frac{1}{2} \|\nabla \Phi(w_t)\|^2$$

Prox-gradient and subgradient are closely aligned
[Davis & Drusvyatskiy. SIAM J. Optim. (2019)]

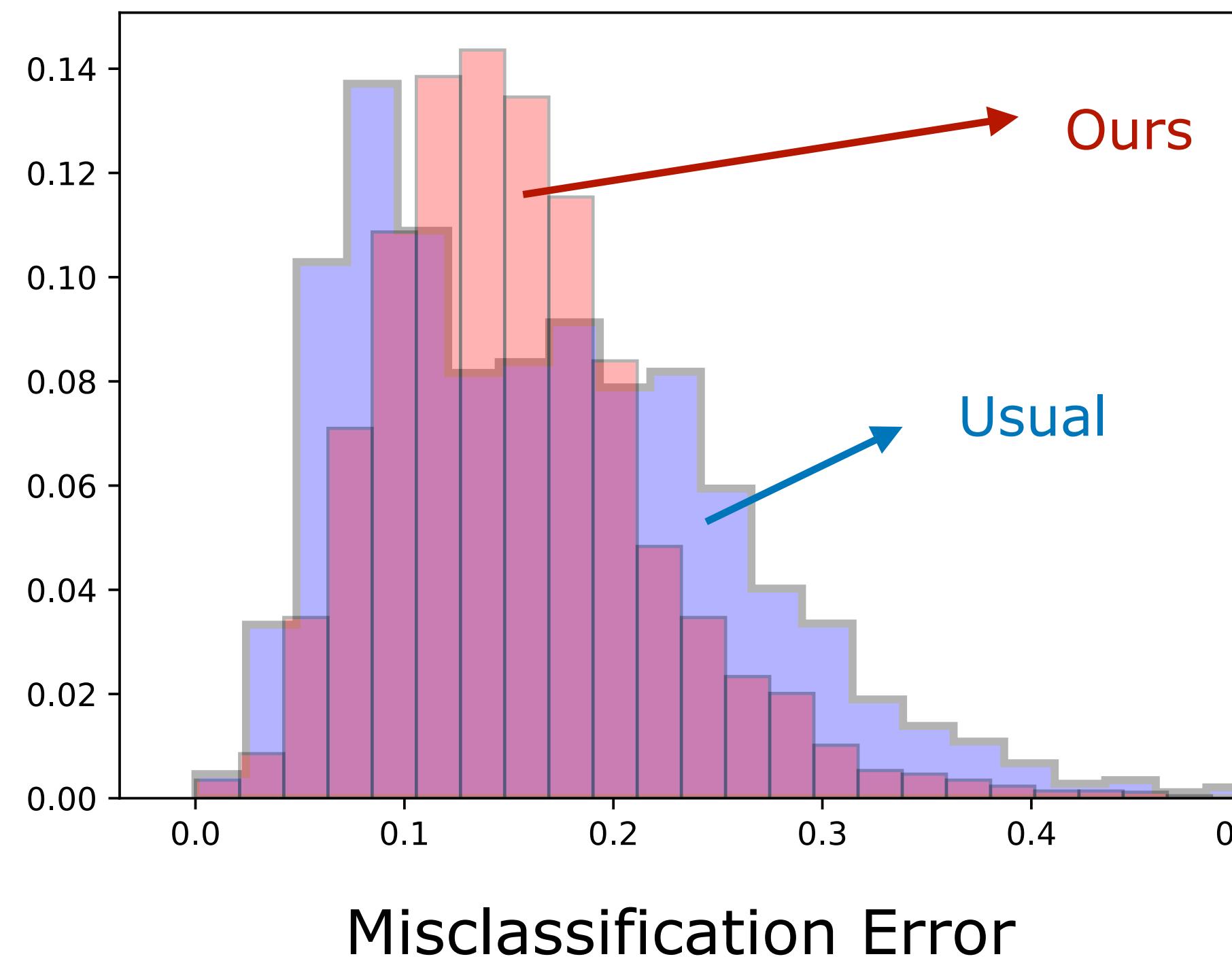
Experiments: EMNIST



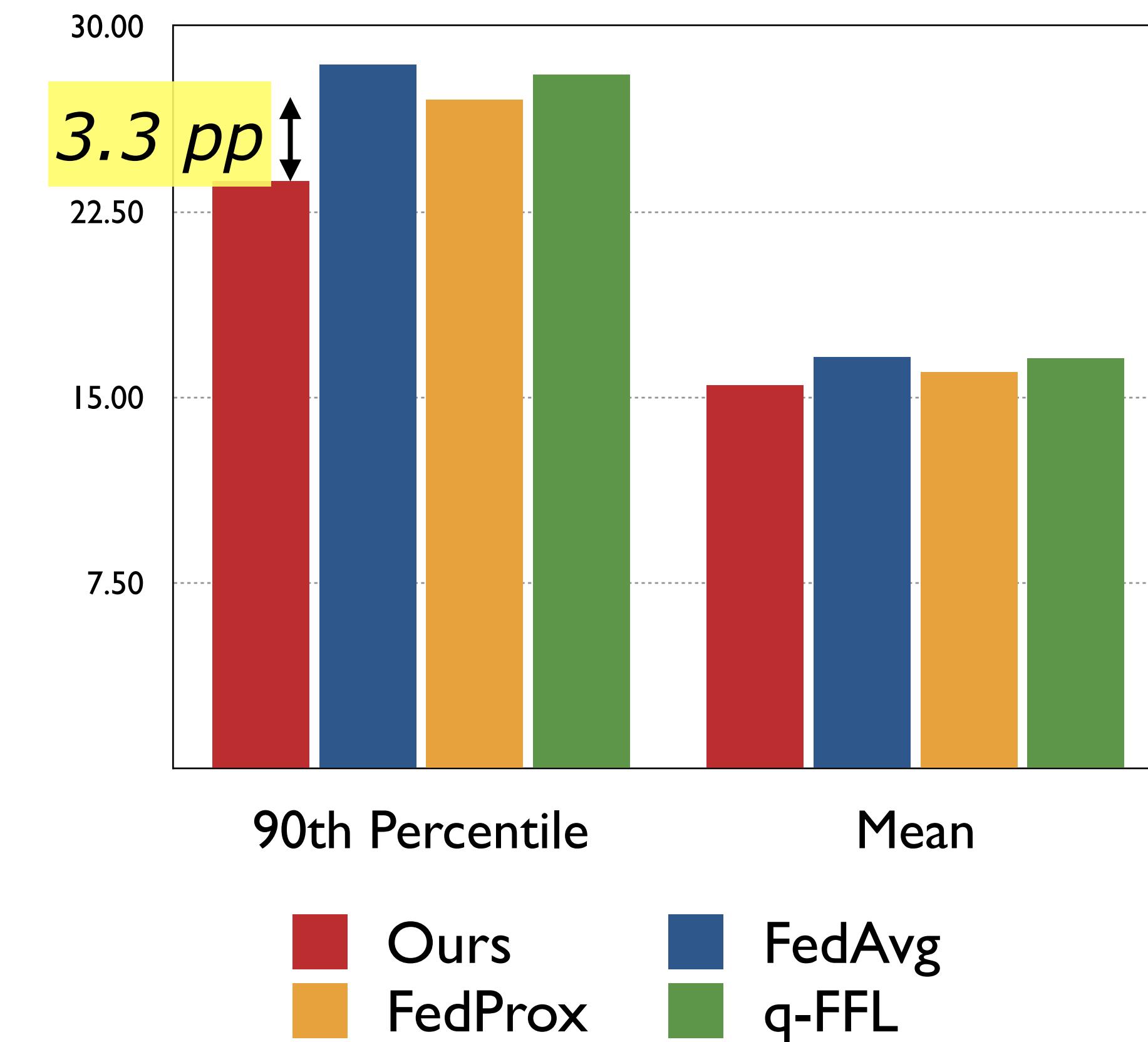
Histogram of per-client errors



Histogram of per-client errors



Misclassif. Error



- Simplicial-FL has the smallest 90th percentile error
- Simplicial-FL is competitive on the mean error

Distributionally robust learning with 1 additional line of code

```
import torch.nn.functional as F
from sqwash import reduce_superquantile

for x, y in dataloader:
    y_hat = model(x)
    batch_losses = F.cross_entropy(y_hat, y, reduction='none') # must set `reduction='none'`
    loss = reduce_superquantile(batch_losses, superquantile_tail_fraction=0.5) # Additional line
    loss.backward() # Proceed as usual from here
```

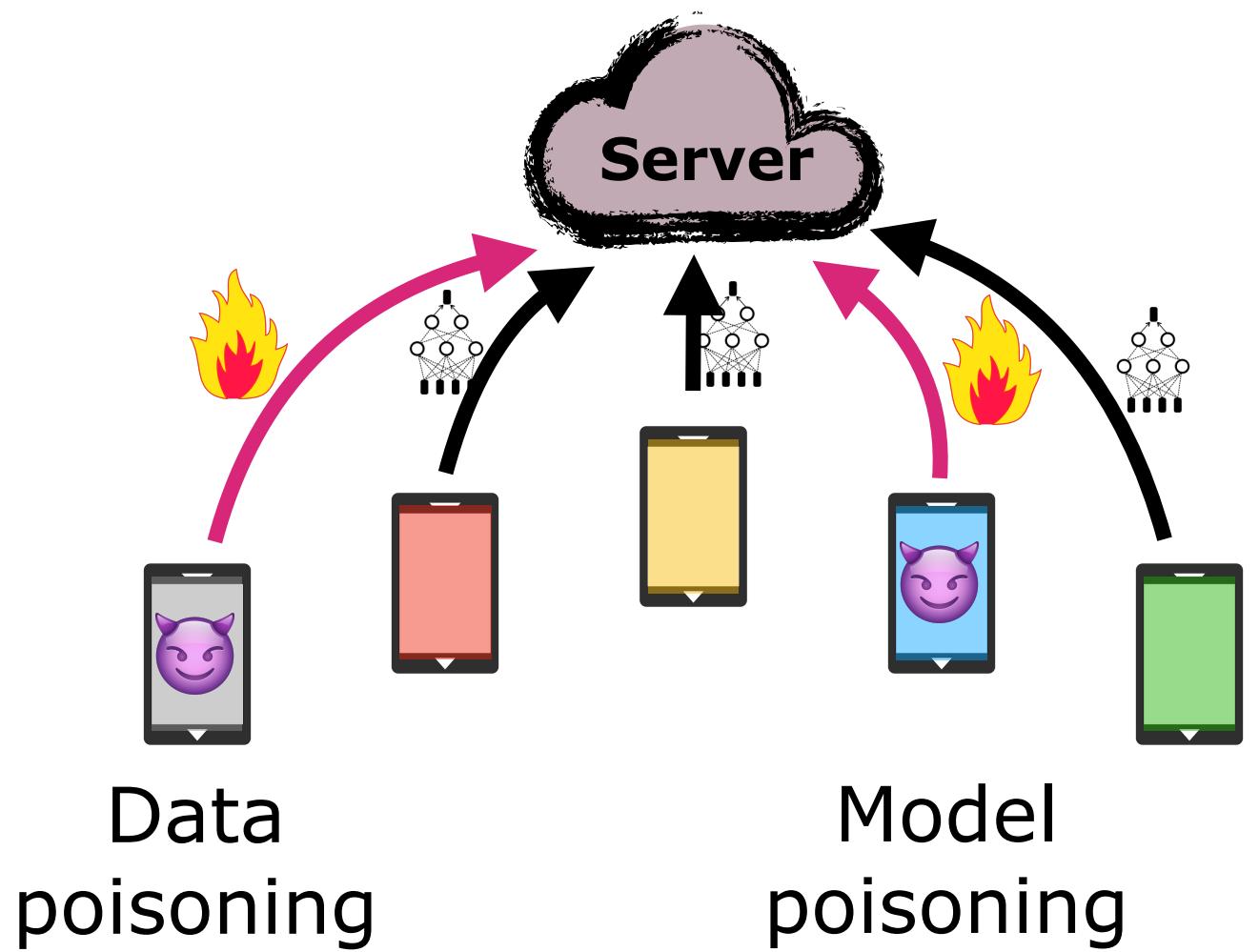
Install: [pip install sqwash](#)

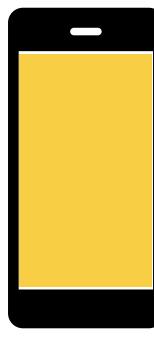
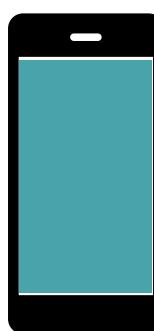
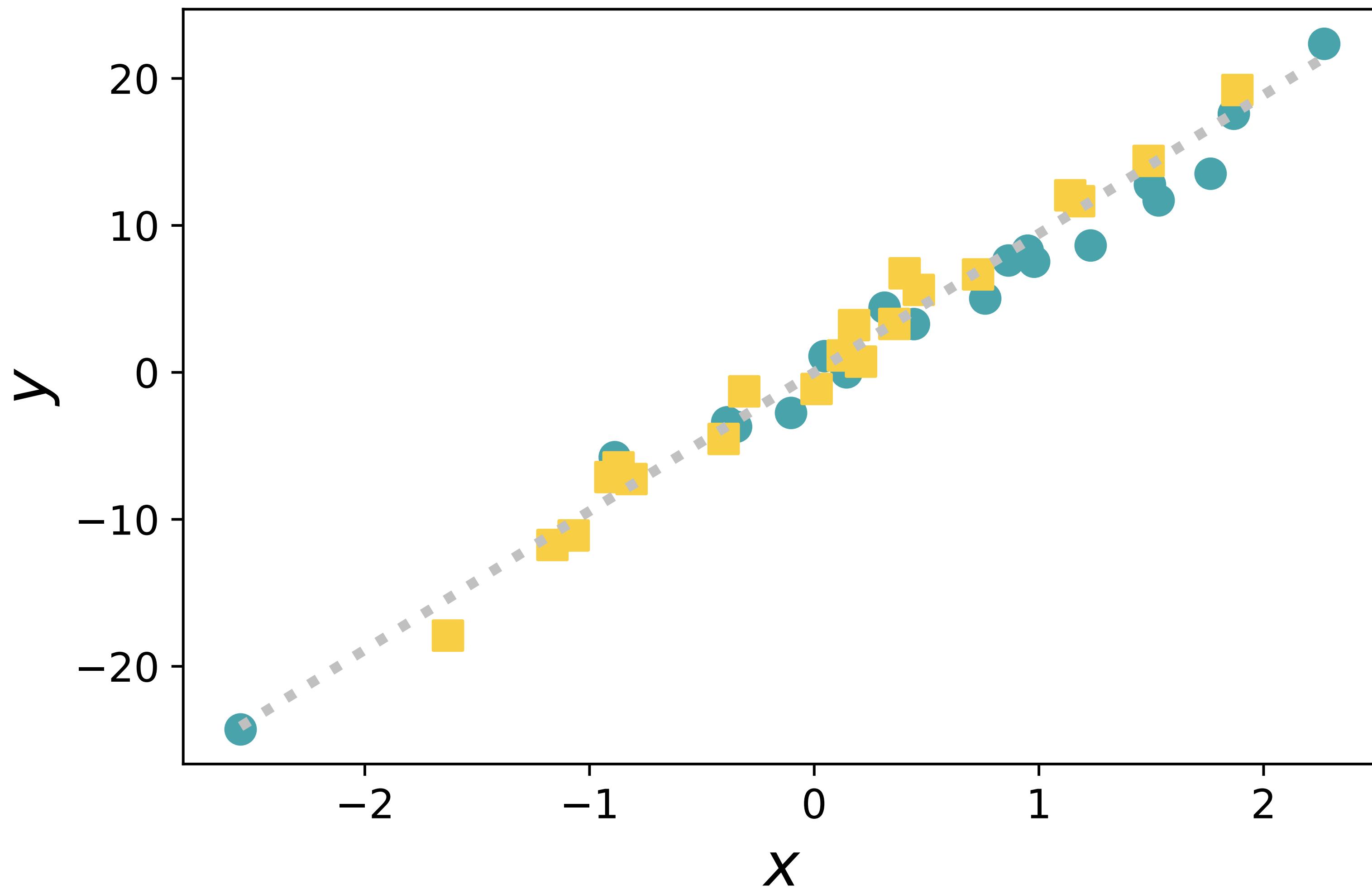
Documentation: krishnap25.github.io/sqwash/

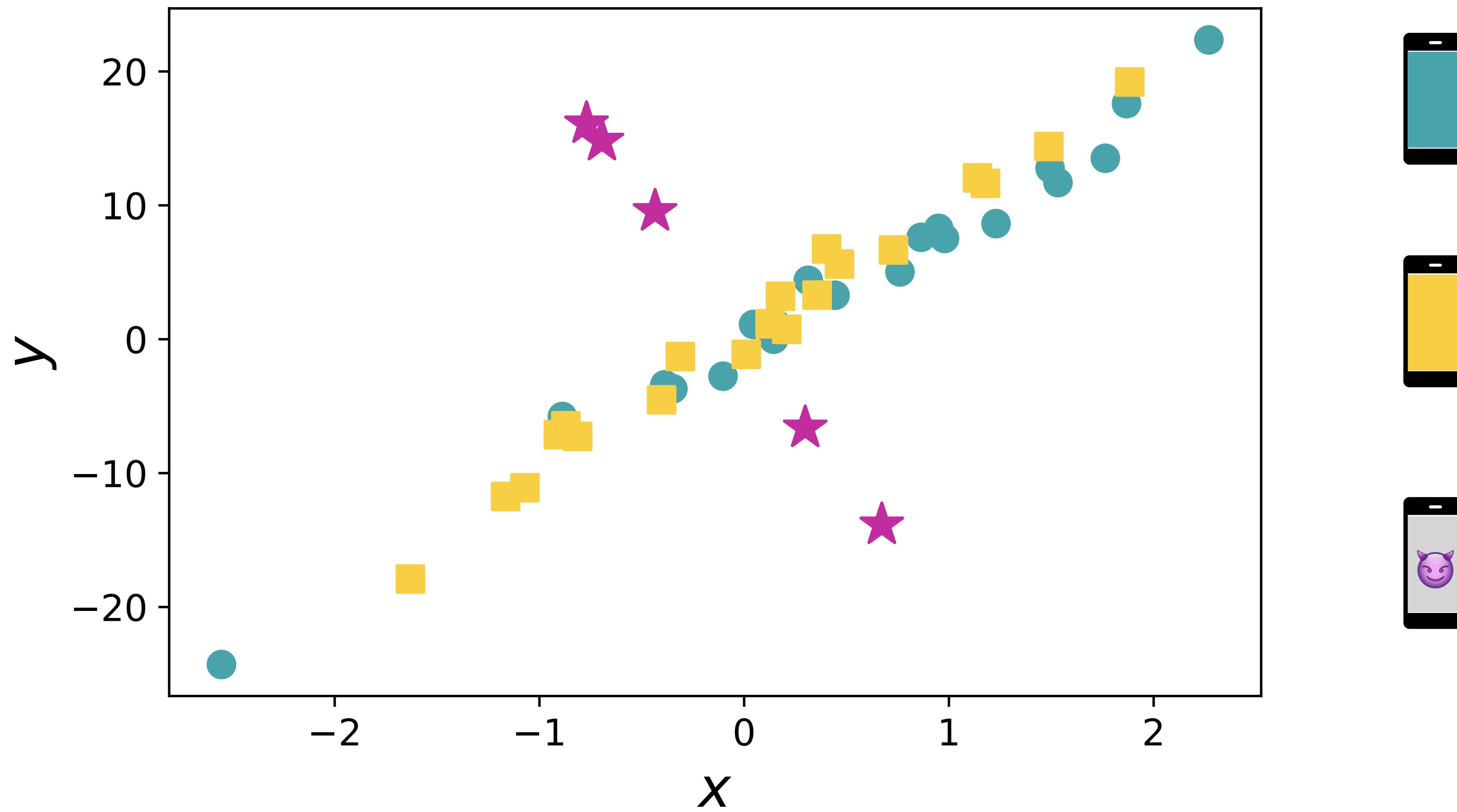


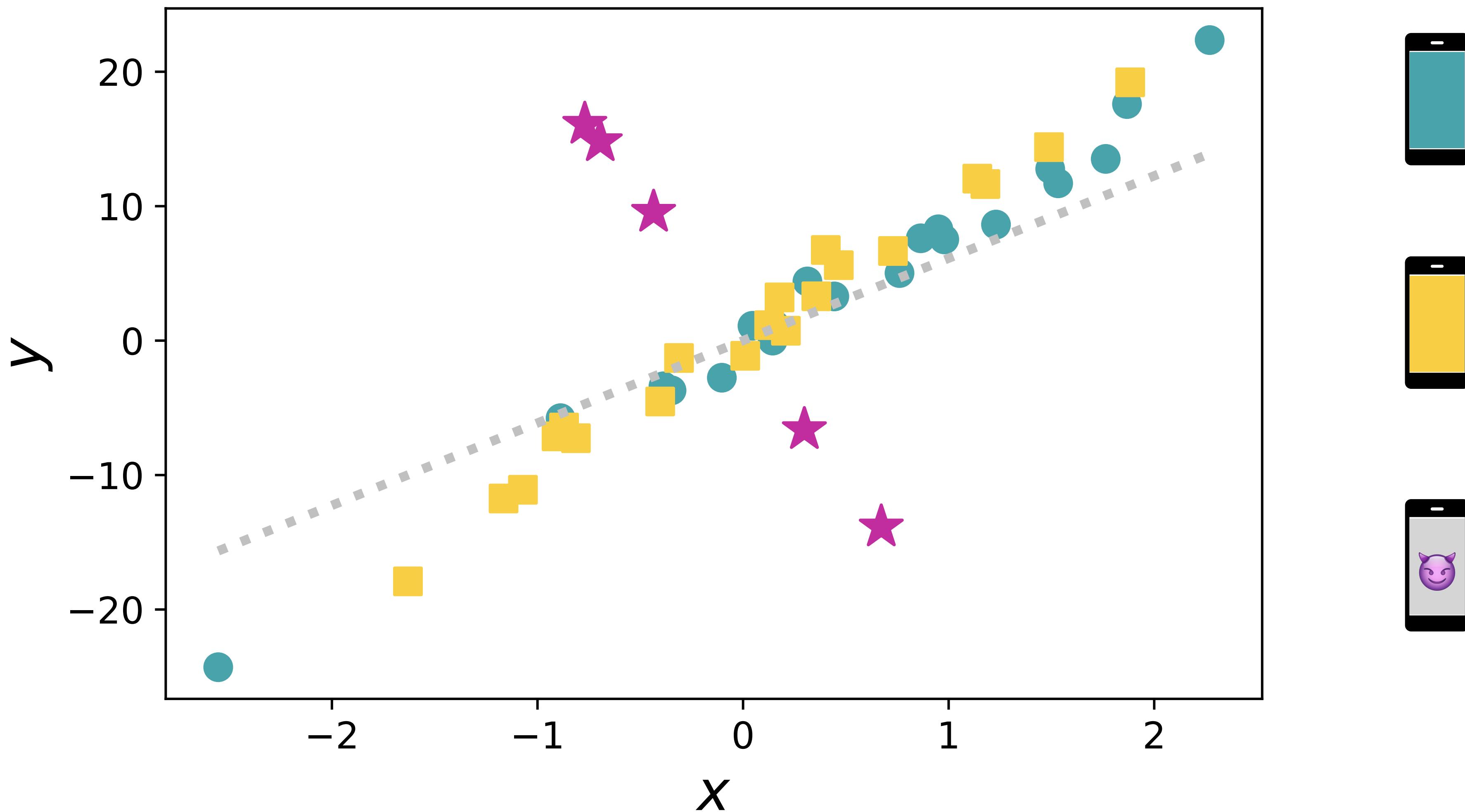
Part 2: Robust aggregation for federated learning

[TSP '22]

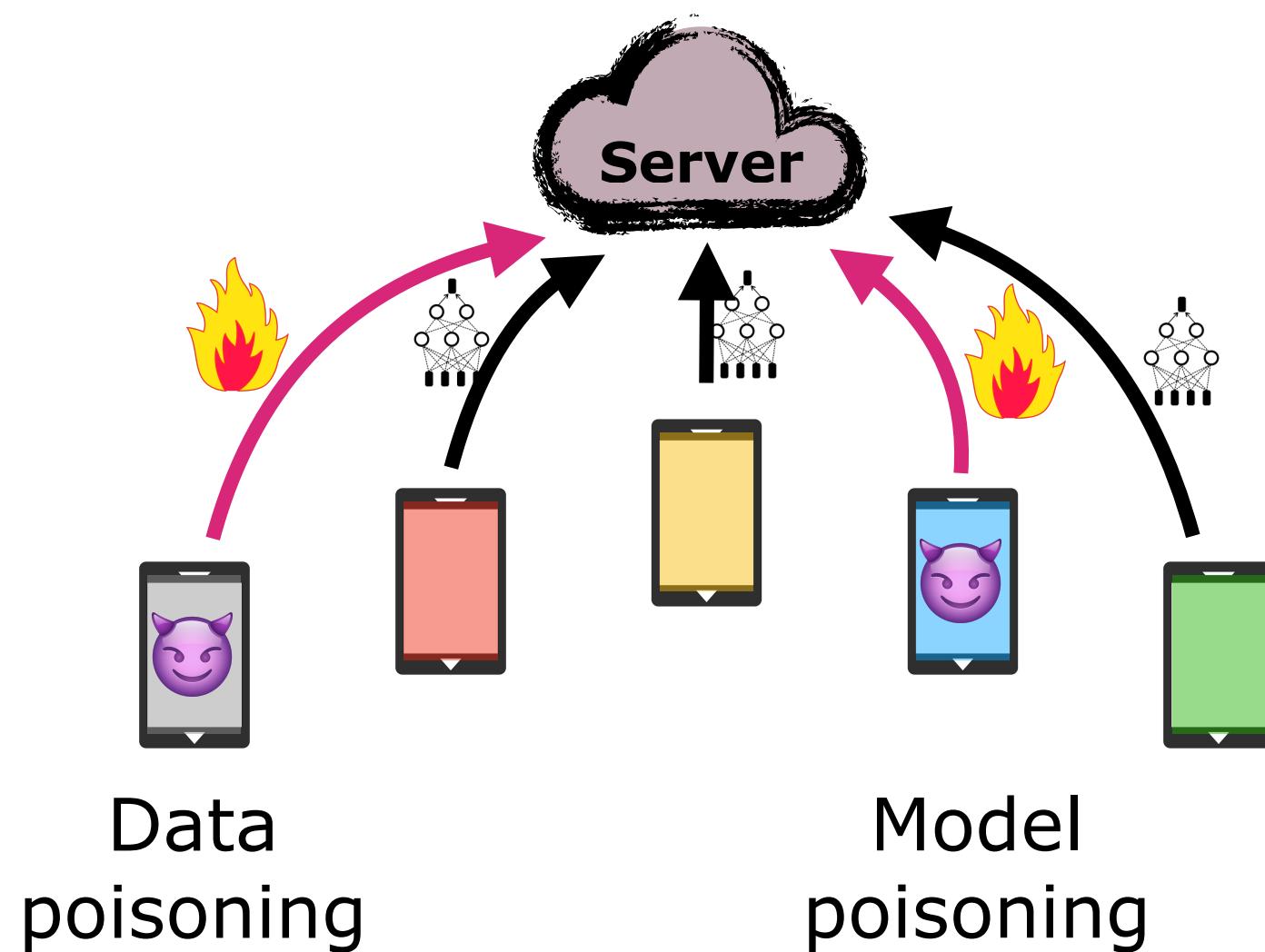
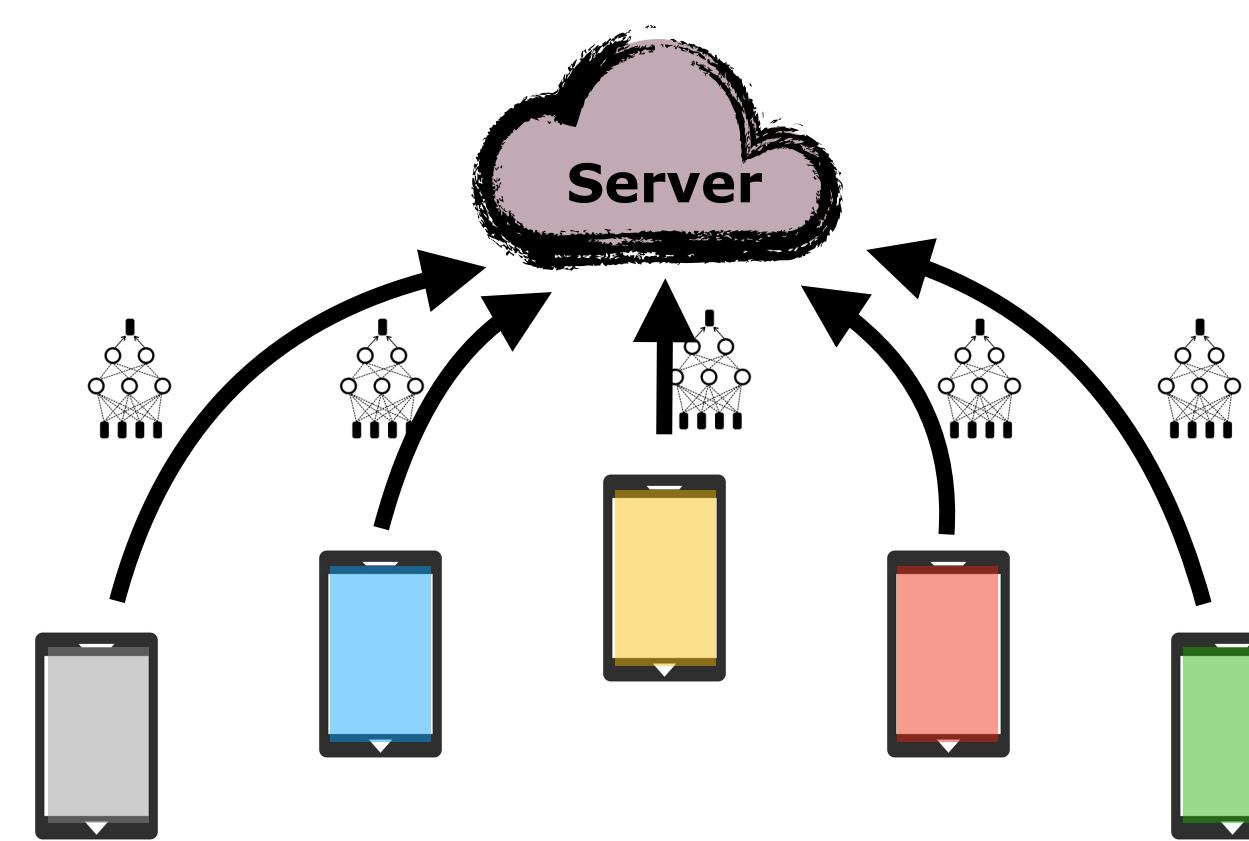




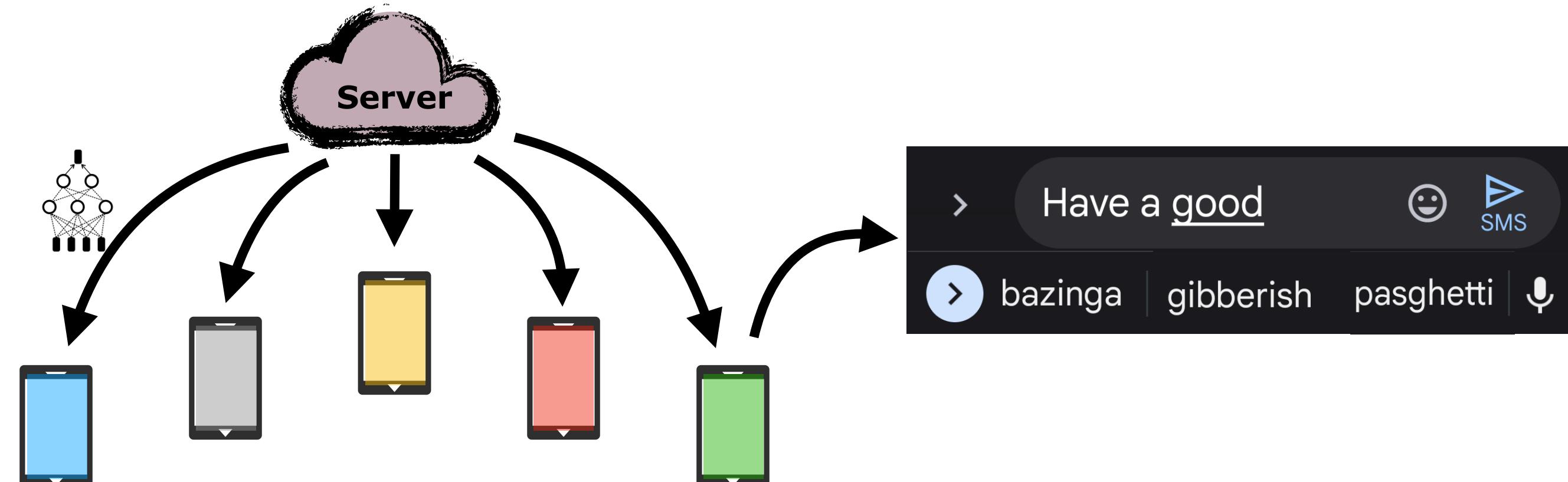
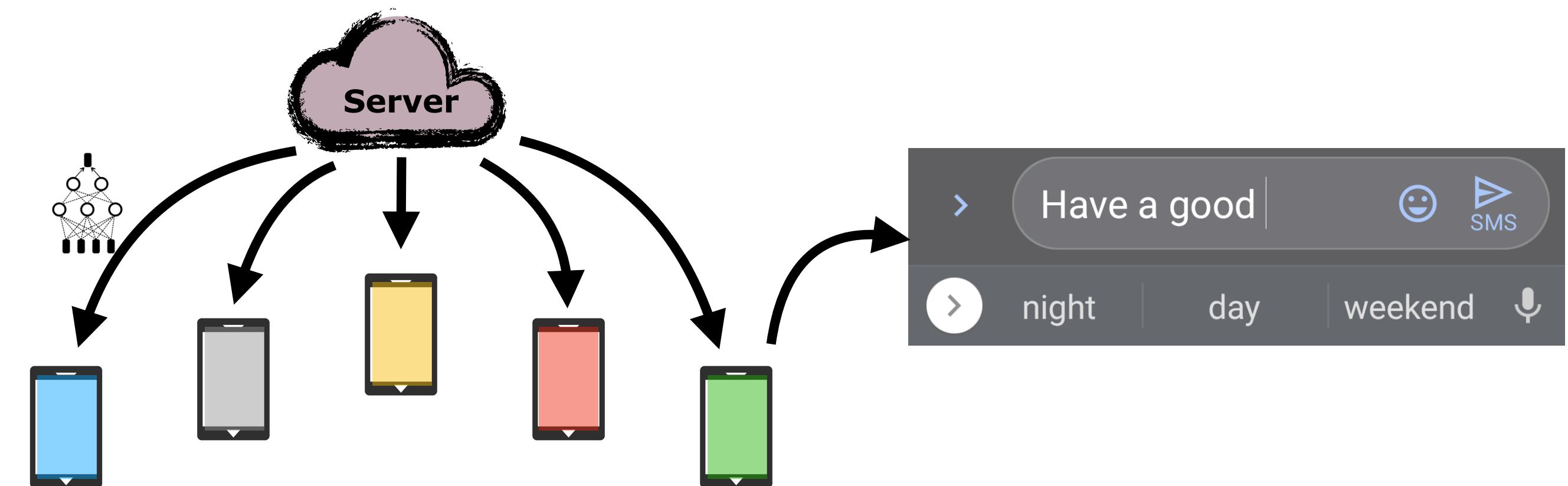




Training



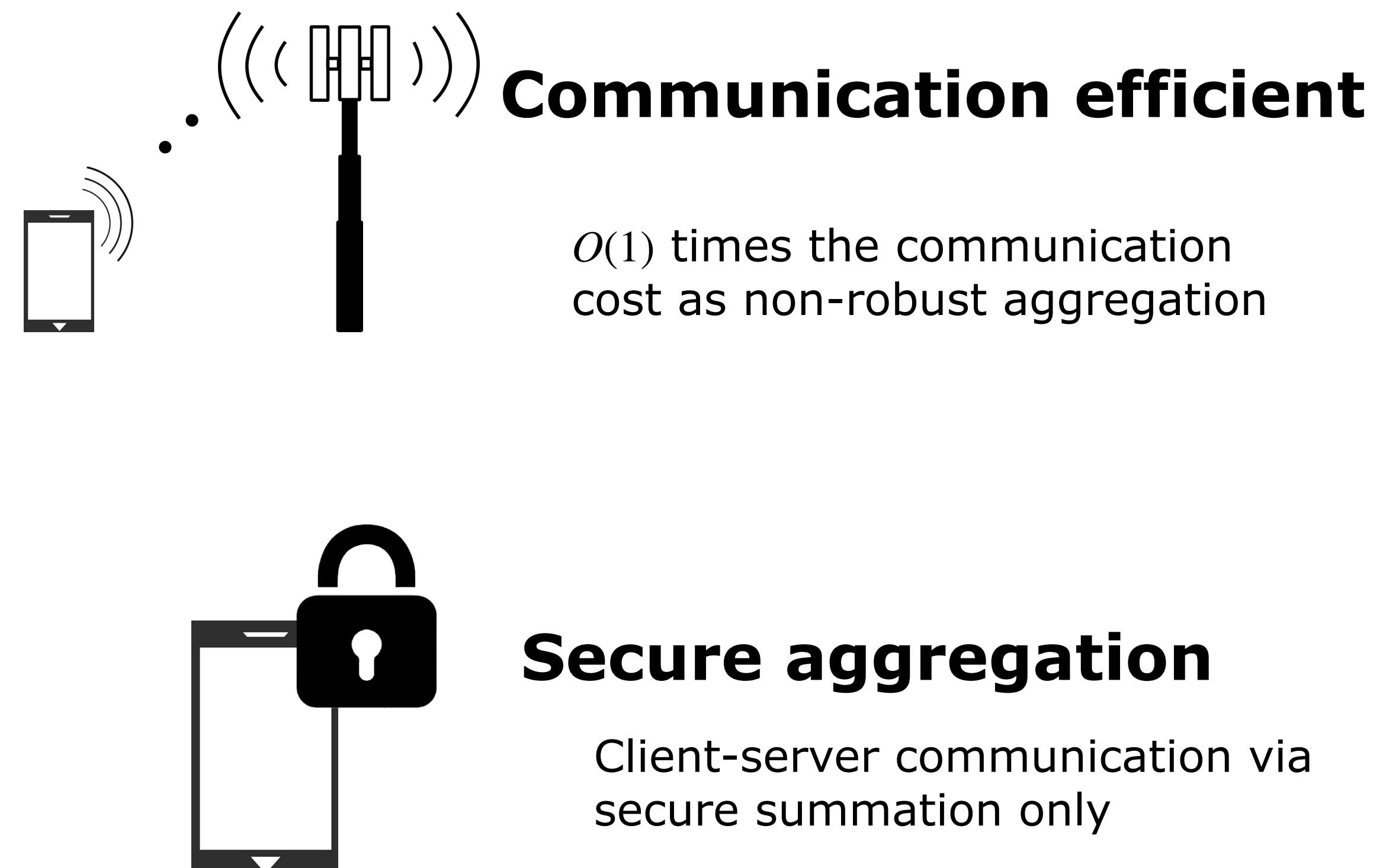
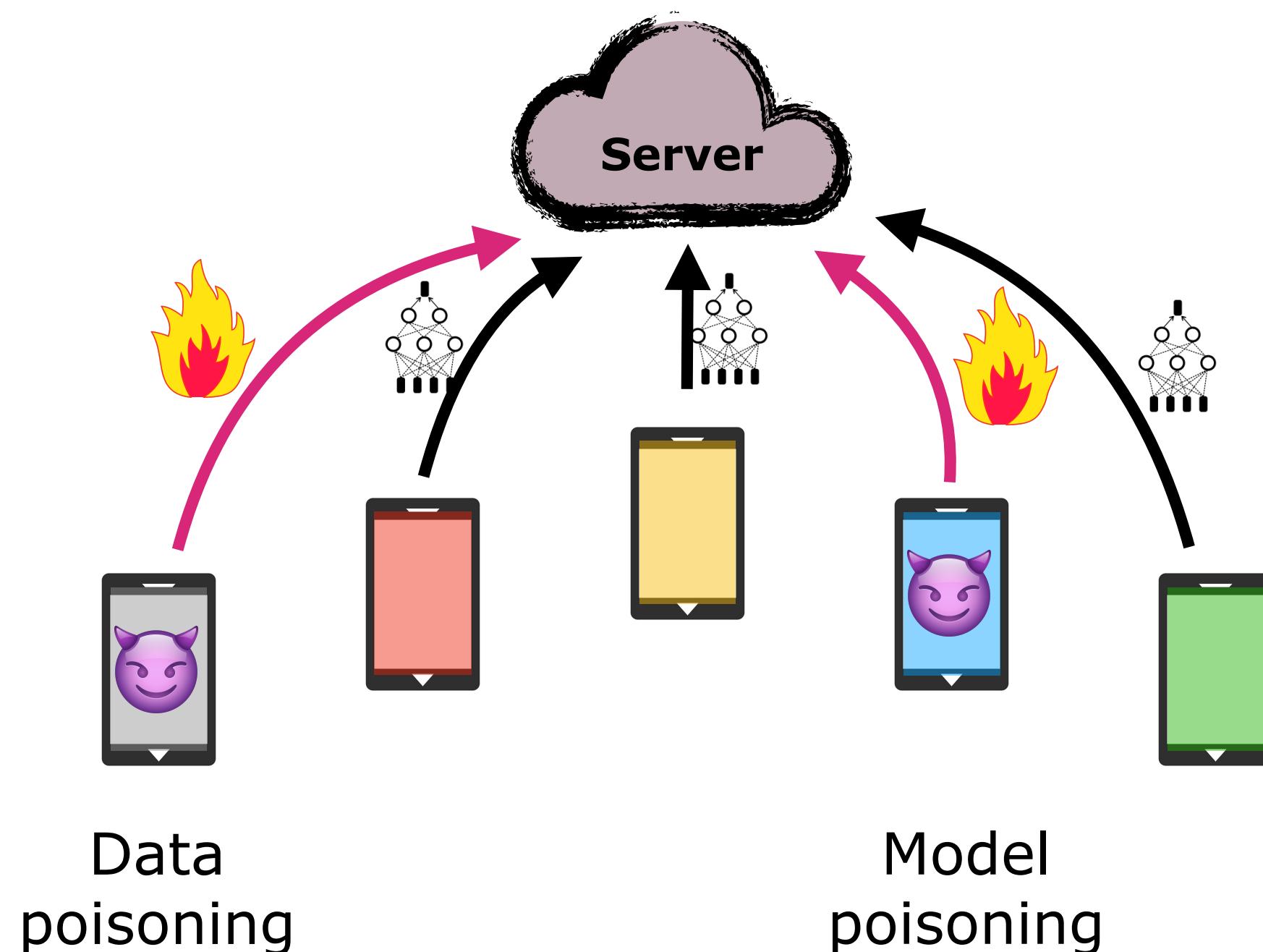
Deployment



Arithmetic mean aggregation is *not robust* to corruptions \Rightarrow Poor predictions!

Our goal

Design a robust aggregation algorithm
for federated learning which is



Note: not DP

Heterogeneity vs. Robustness

Consider mean estimation in Huber's contamination model:

$$w_1, \dots, w_n \sim (1 - \rho) \mathcal{N}(\mu, \sigma^2 I) + \rho Q$$

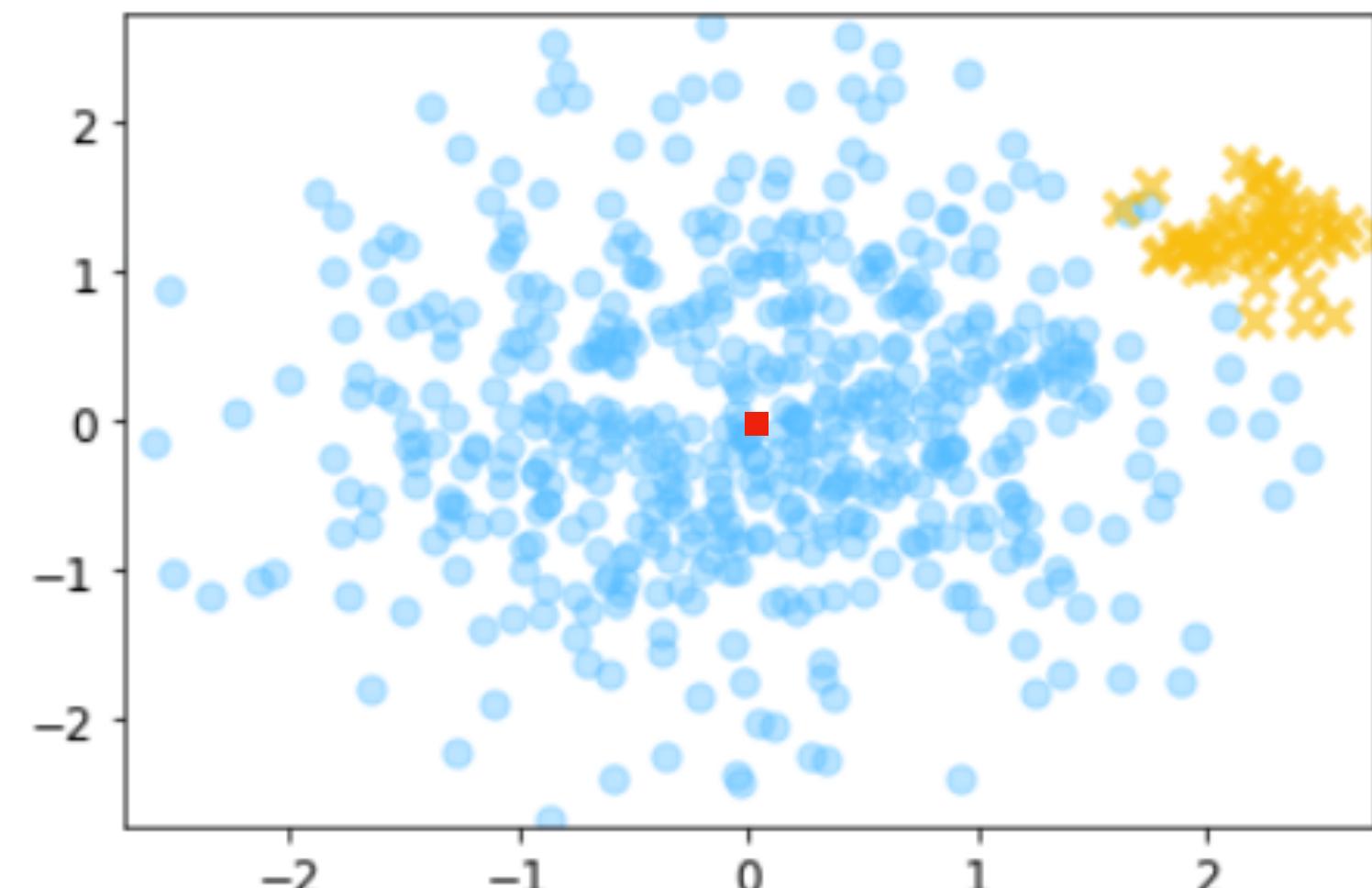
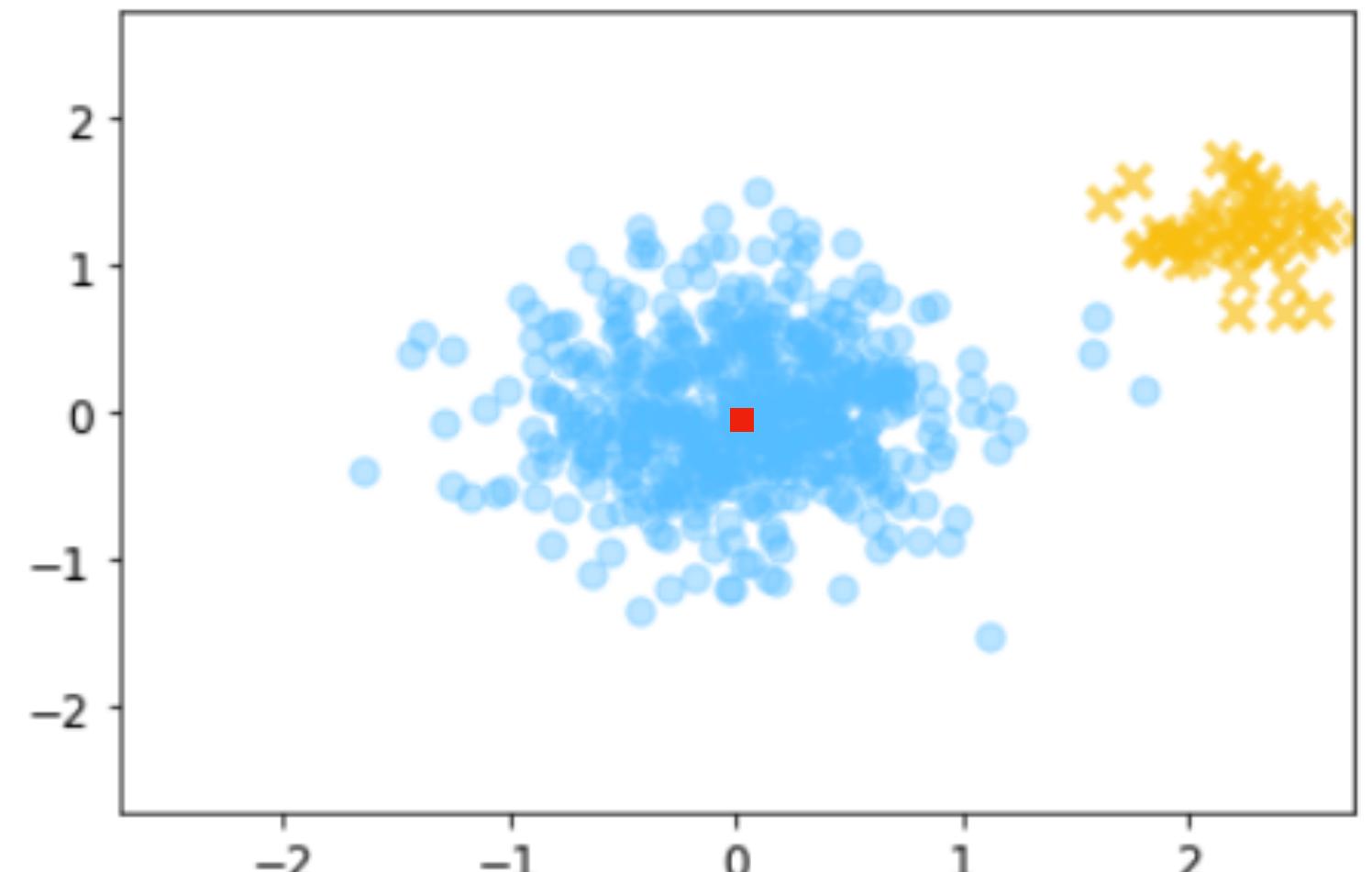
inliers

outliers

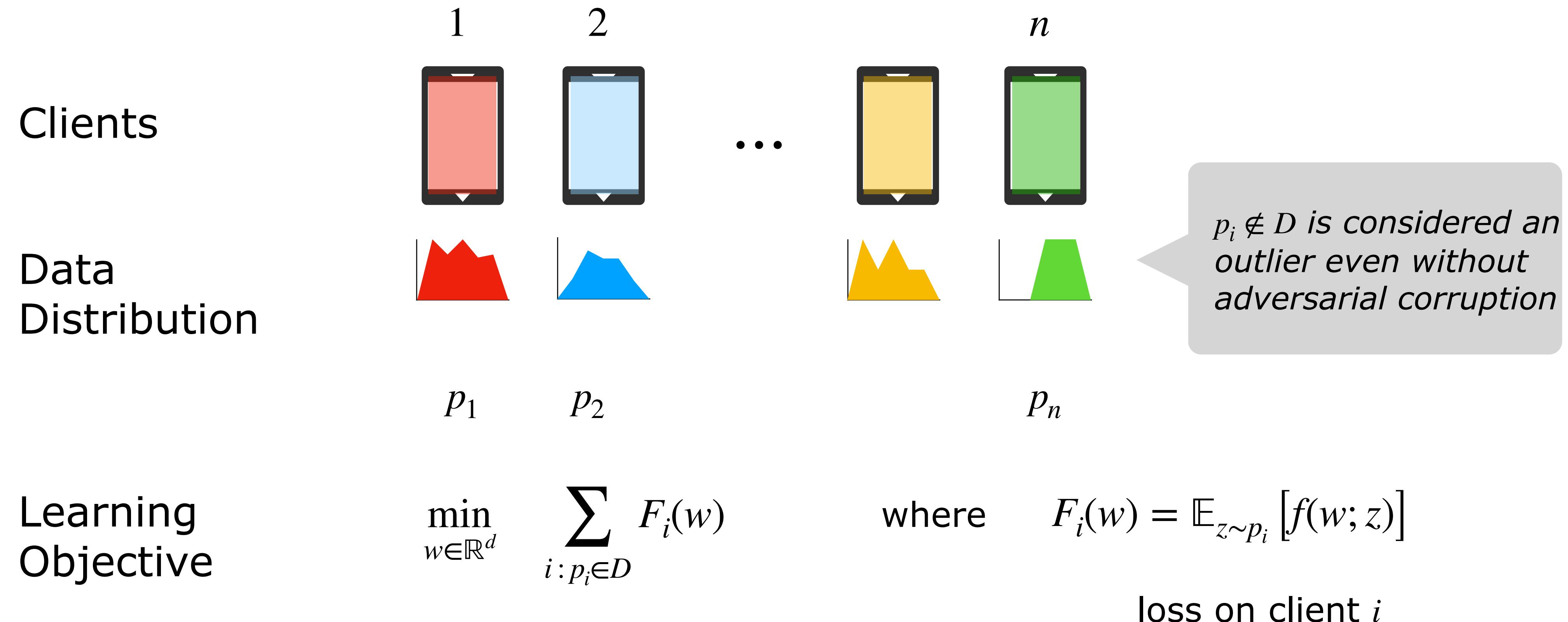
Any mean estimate \bar{w}_n must satisfy

$$\|\bar{w}_n - \mu\|^2 \gtrsim \sigma^2 \left(\rho^2 + \frac{d}{n} \right)$$

[Chen, Gao, Ren. Annals of Stat. (2018)]



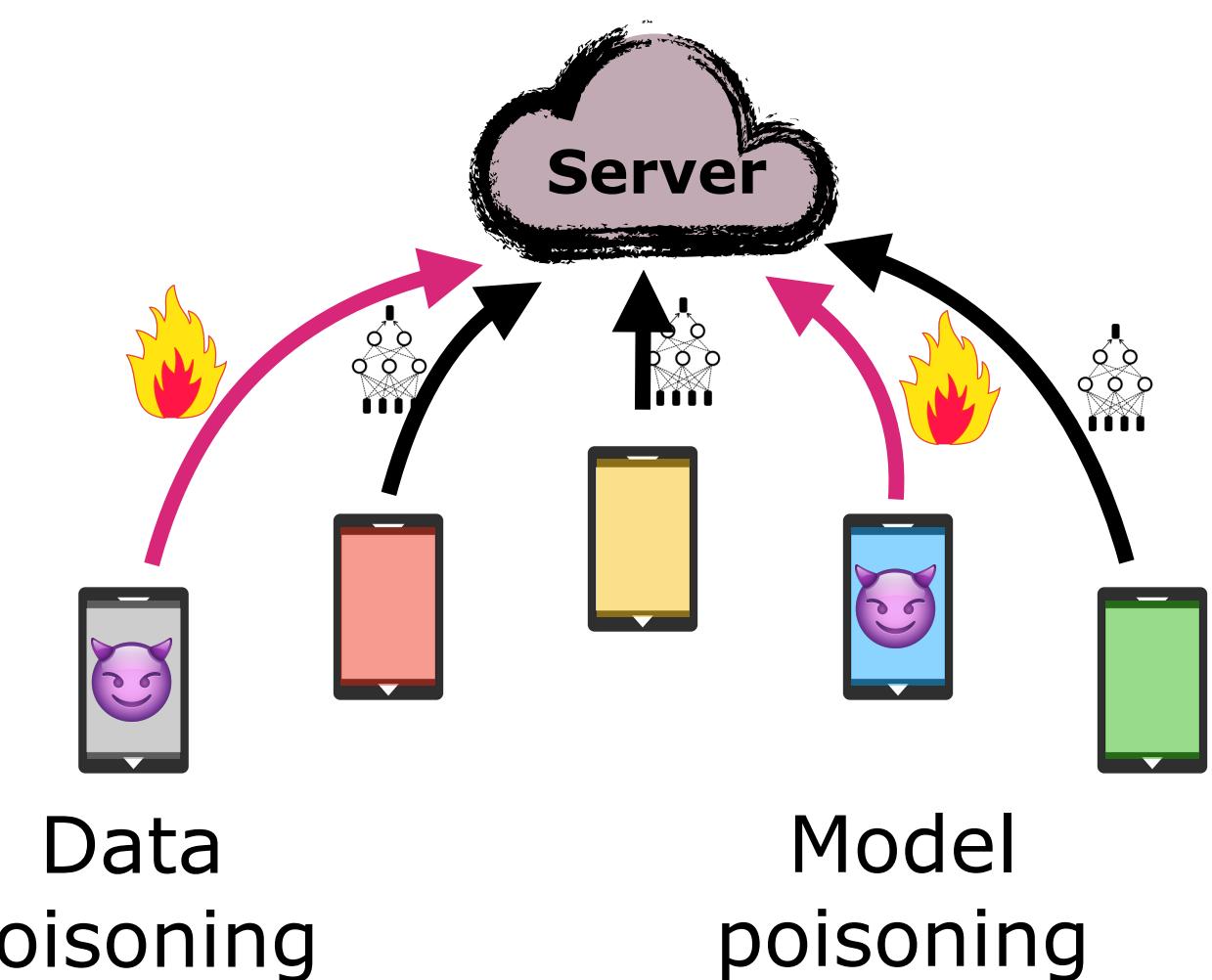
For general federated learning, fix a set D of “inlier” distributions

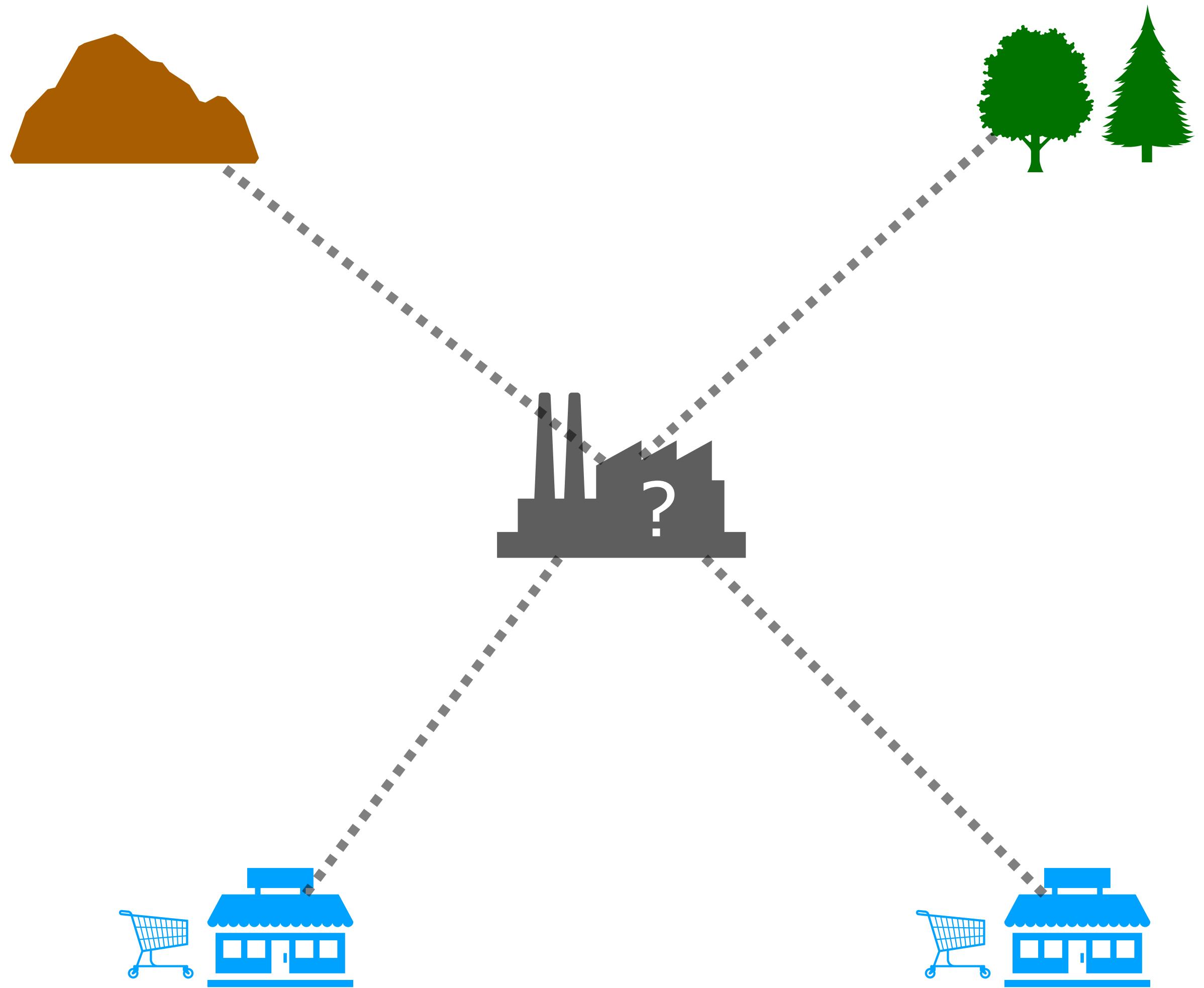


Convergence only possible up to size of inlier set D due to $\|\bar{w}_n - \mu\|^2 \gtrsim \sigma^2 \left(\rho^2 + \frac{d}{n} \right)$

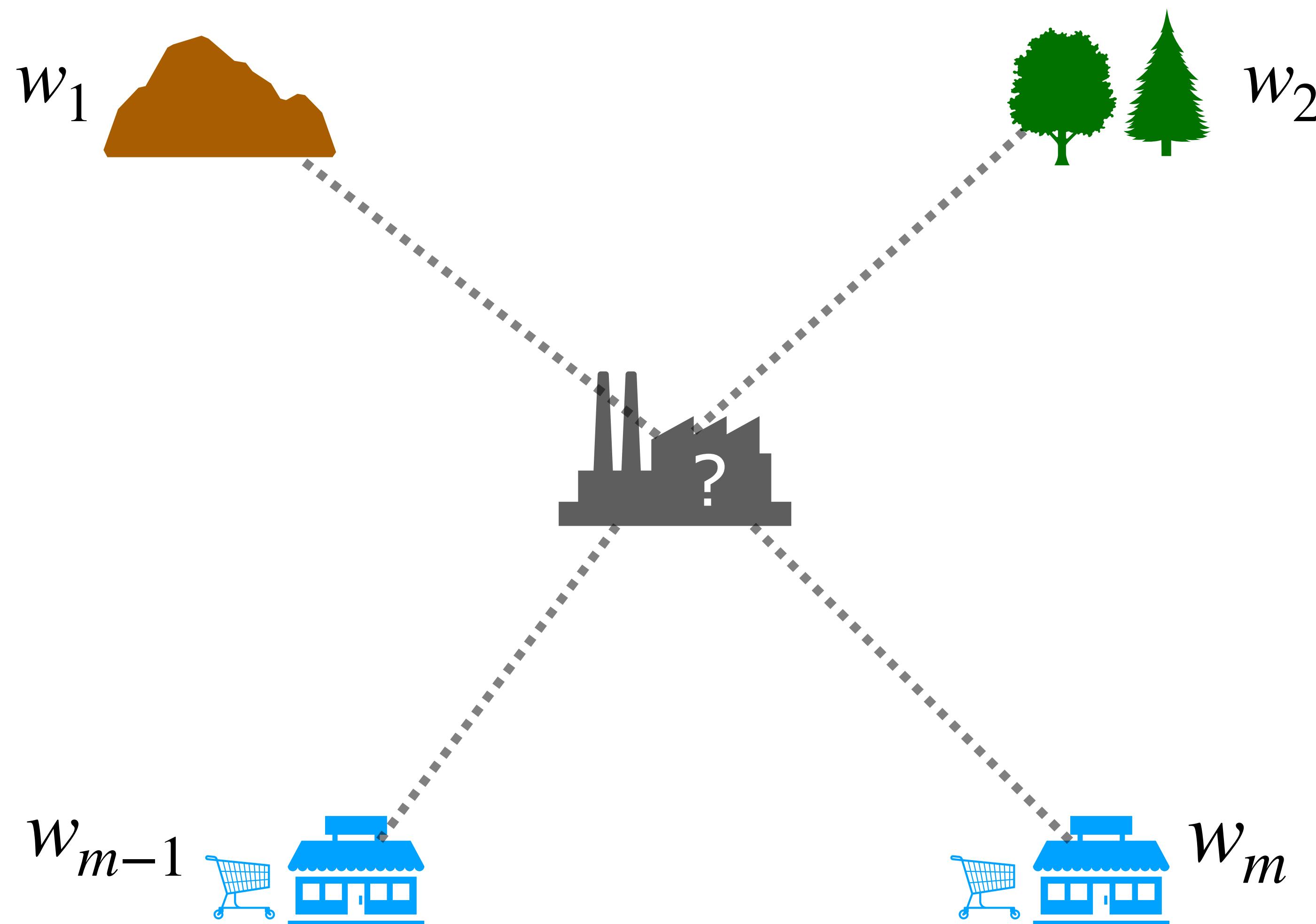
Algorithm is agnostic to D — only appears in analysis

Federated learning with robust aggregation





Fermat & Torricelli (~1600s), Weber (1909)

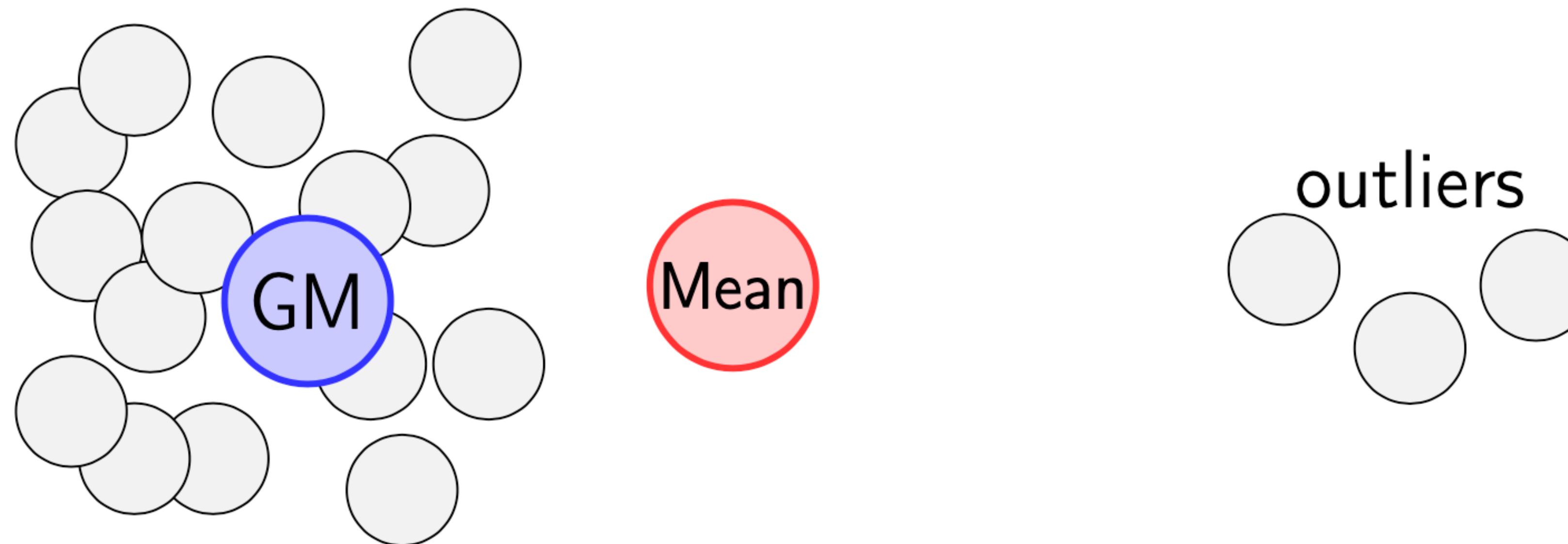


Geometric Median

$$\text{GM}(w_1, \dots, w_m) = \arg \min_z \left\{ \sum_{i=1}^m \|z - w_i\|_2 \right\}$$

Fermat & Torricelli (~1600s), Weber (1909)

Robustness: Breakdown point of GM = 1/2



Nemirovski & Yudin (1983) | Jerrum, Valiant & Vazirani (1986) | Lopuhaa & Rousseeuw (1991)
Hsu & Sabata (2013) | Minsker (2015) | Lugosi, Gabor & Mendelson (2019) | Lecué & Lerasle (2020)

Weiszfeld's Algorithm

Start with initial guess z_0 and iterate:

$$\beta_{i,t} = \frac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$$

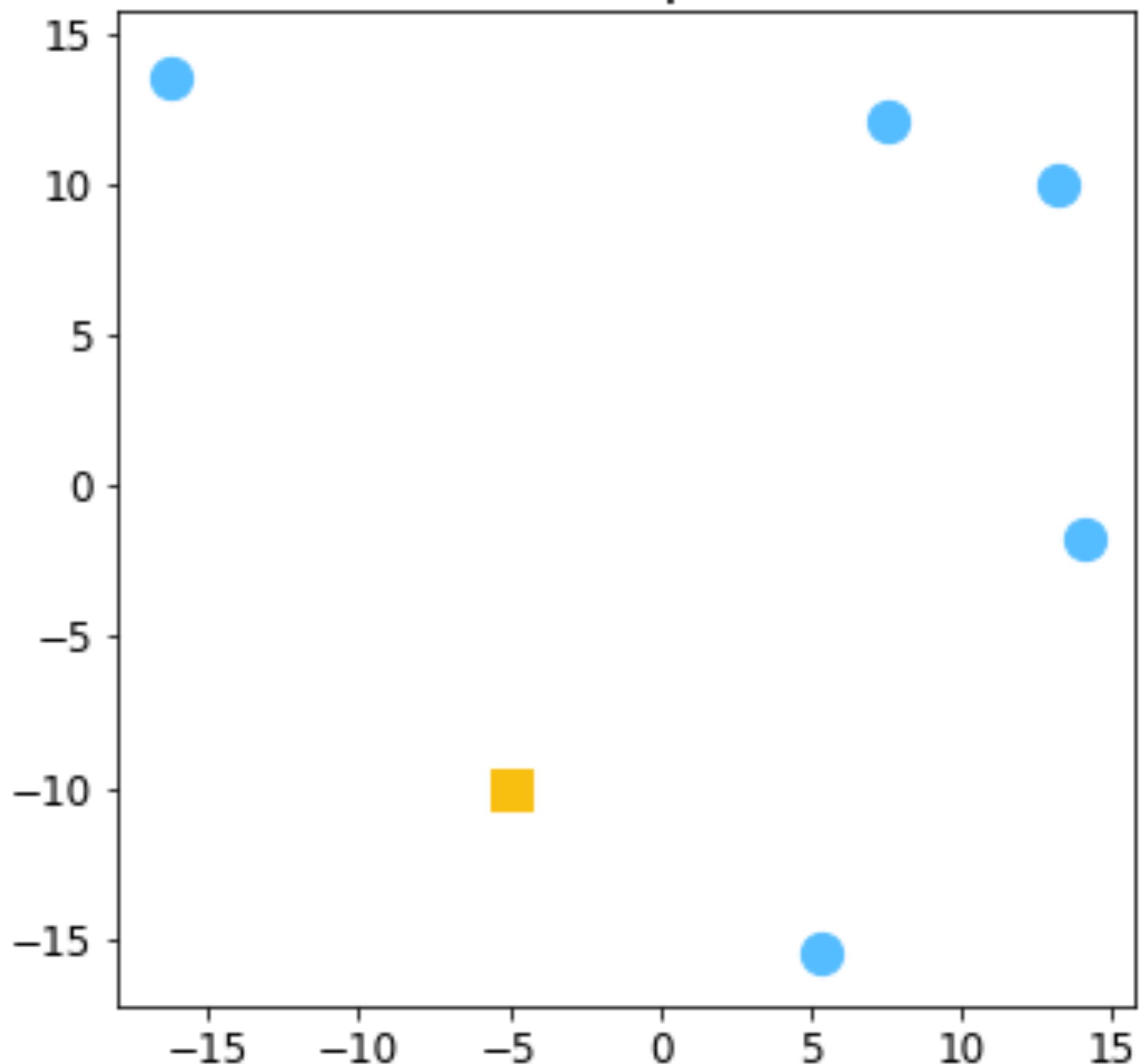
$$z_{t+1} = \frac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$$

Compute new weights

Reweighted average

[Weiszfeld 1937]

Example



Weiszfeld's Algorithm

Start with initial guess z_0 and iterate:

$$\beta_{i,t} = \frac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$$

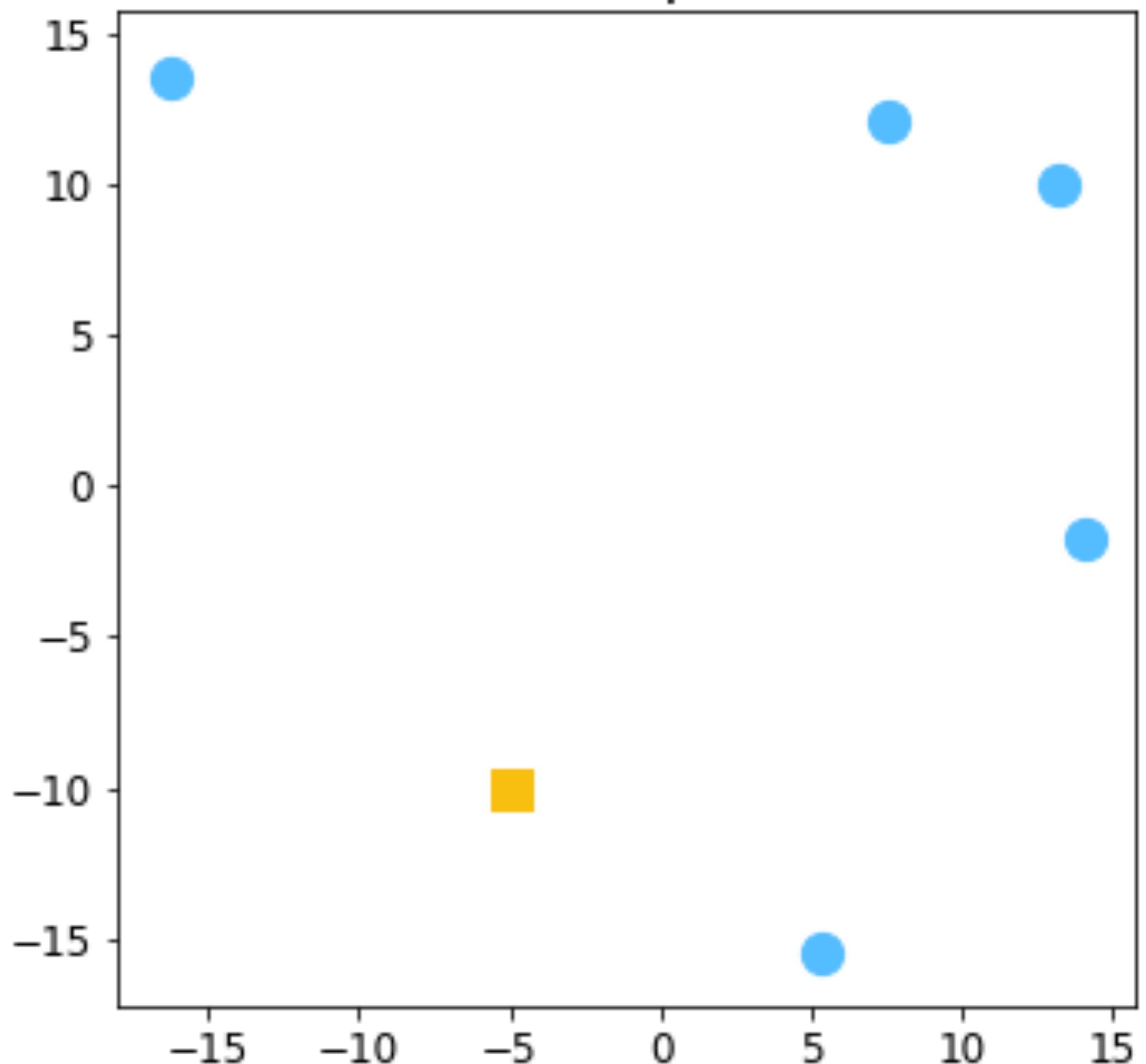
$$z_{t+1} = \frac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$$

Compute new weights

Reweighted average

[Weiszfeld 1937]

Example



Proposition [P., Kakade, Harchaoui]

Assume that $\min_i \|z^* - w_i\| \geq \nu$.

Then, we get an ε -approximate GM in

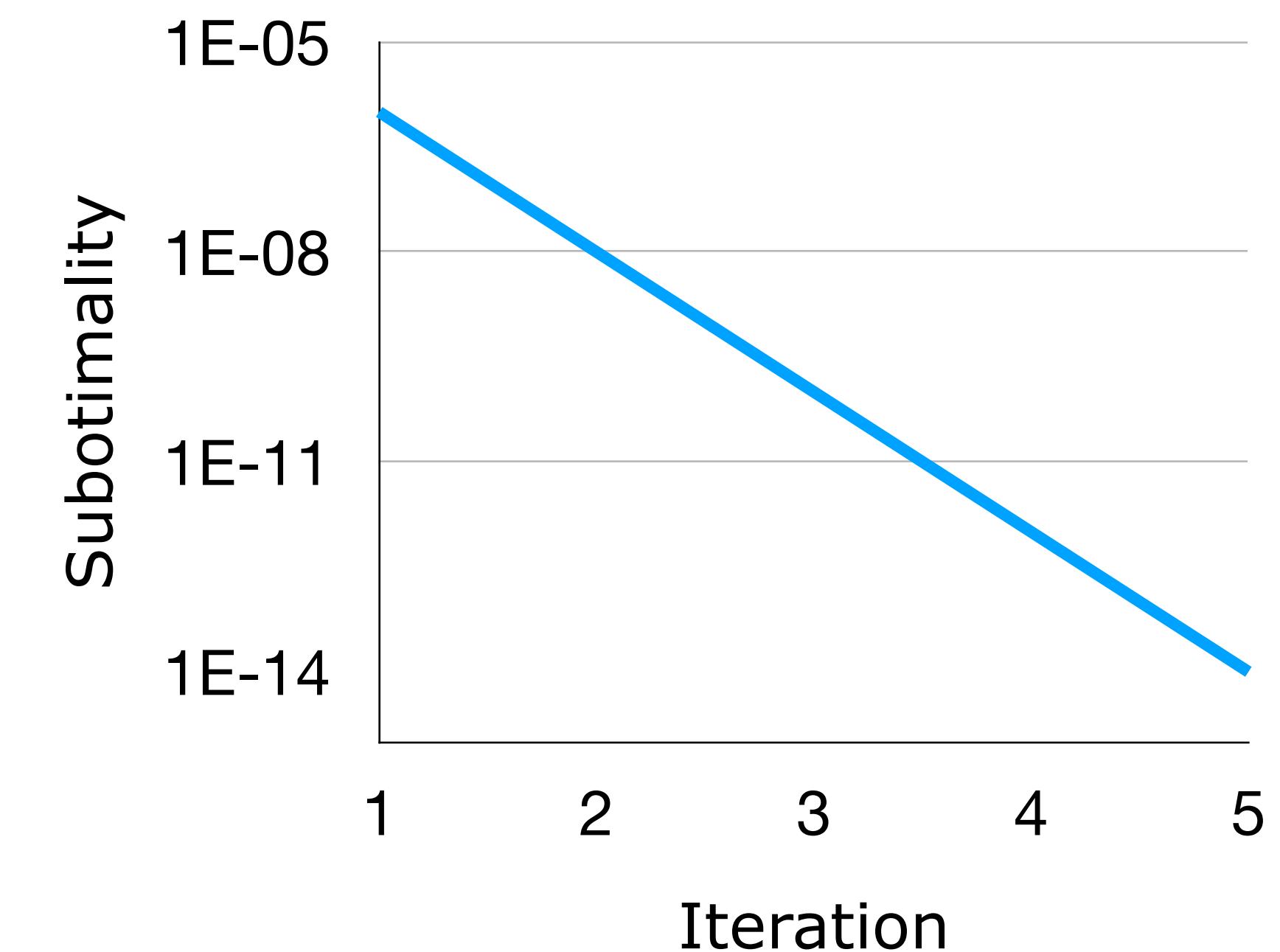
$$O\left(\frac{1}{\nu\varepsilon}\right) \text{ iterations}$$

$$z^* = \mathbf{GM}(w_1, \dots, w_m)$$

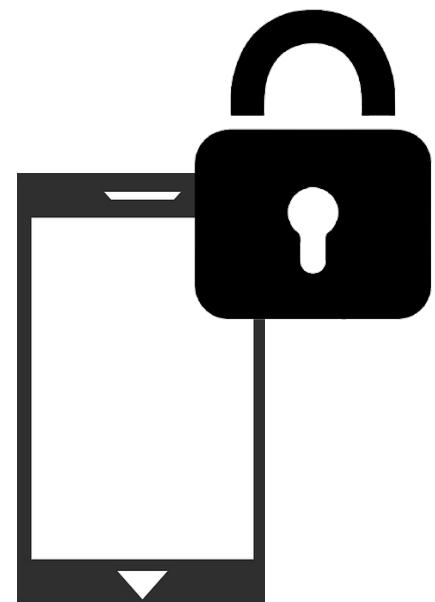
ν : smoothing in β -update

$$\beta_{i,t} = \frac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$$

• (())
• (())
Communication efficient!
Empirically, **3-5** iterations suffice: rapid convergence
Even **1** iteration gives robustness



RFA = FedAvg + GM aggregation

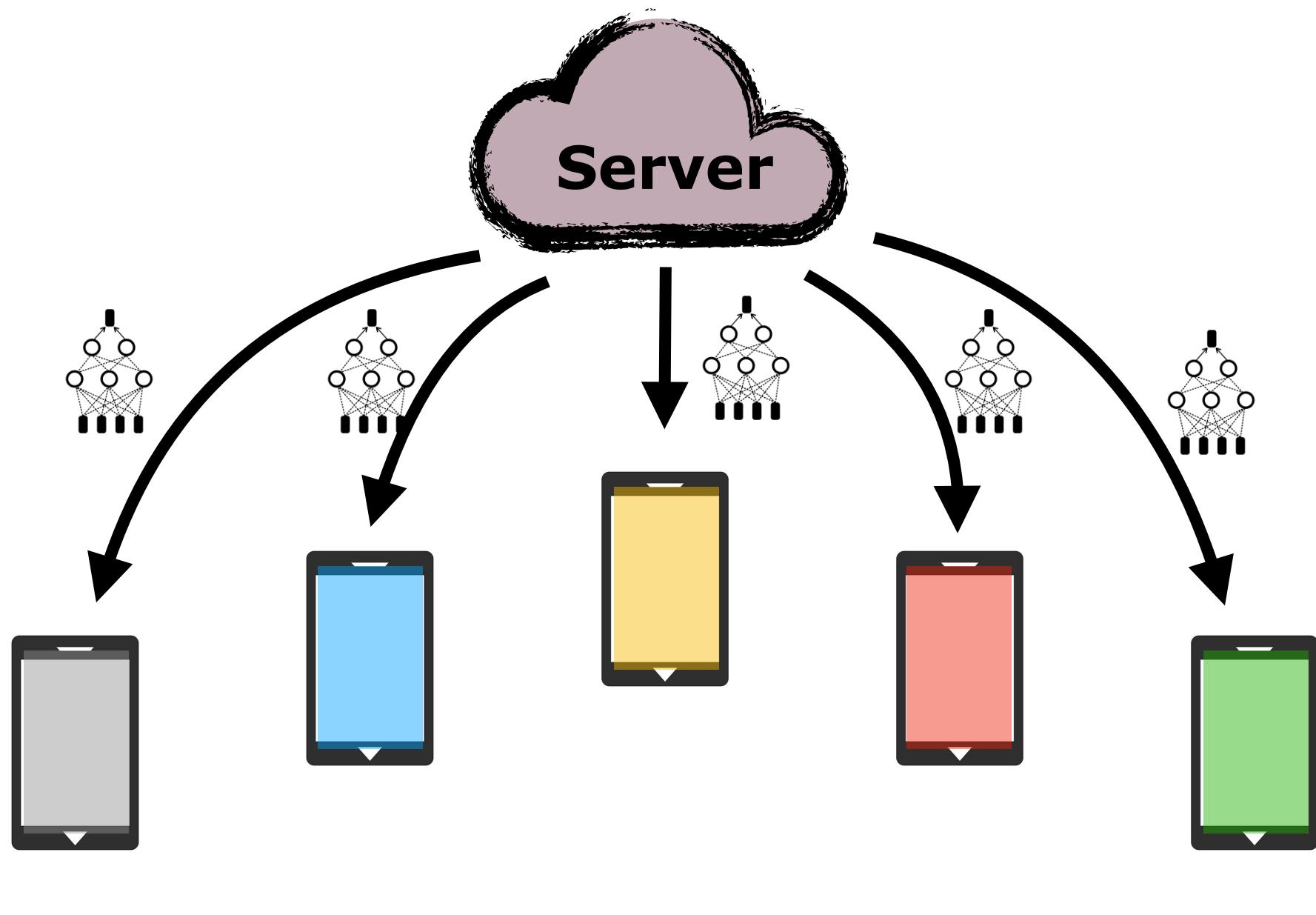


Secure aggregation

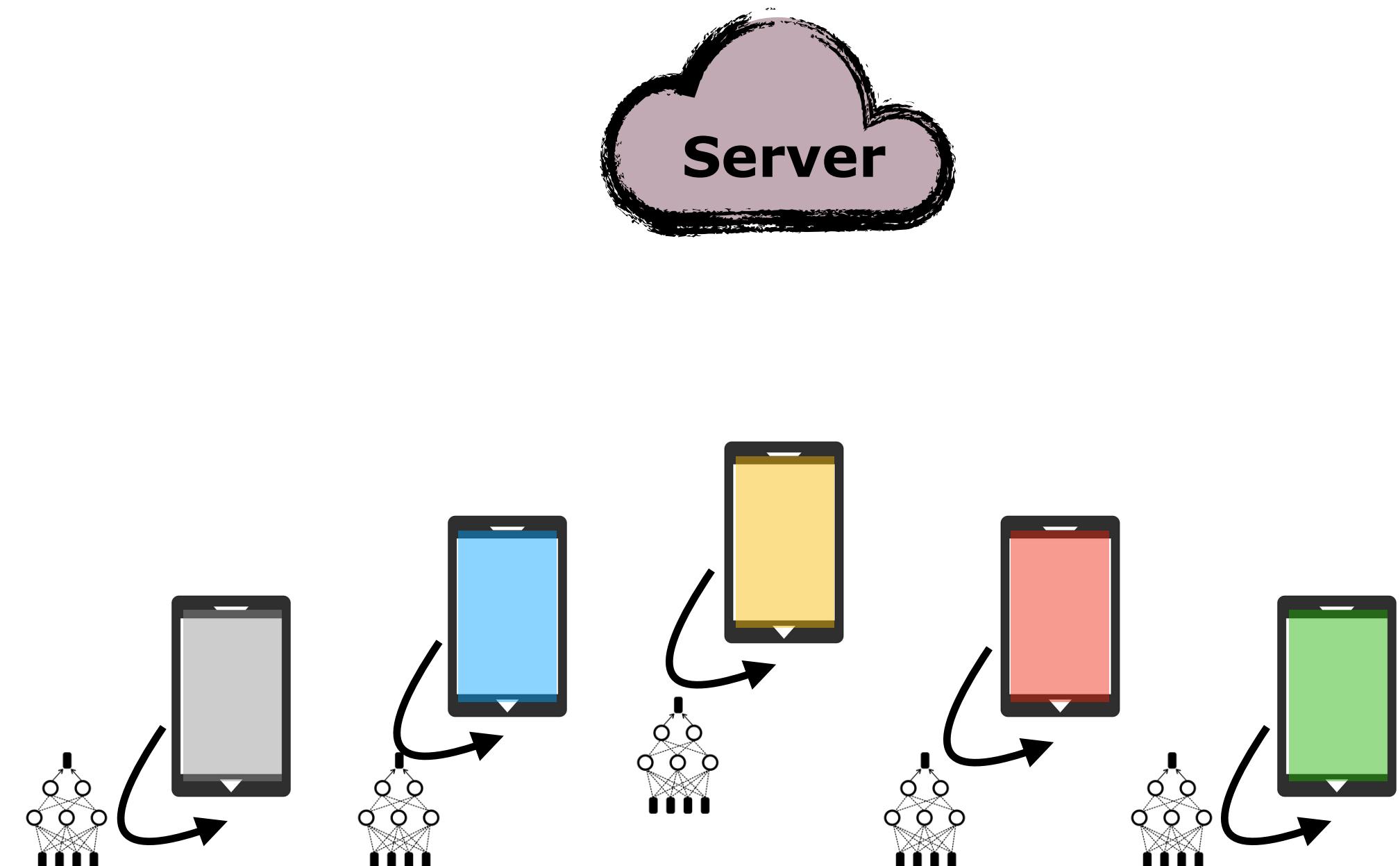
Only client-server communication
is via **secure summation** in

$$z_{t+1} = \frac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$$

Step 1 of 3: Server broadcasts global model to sampled clients



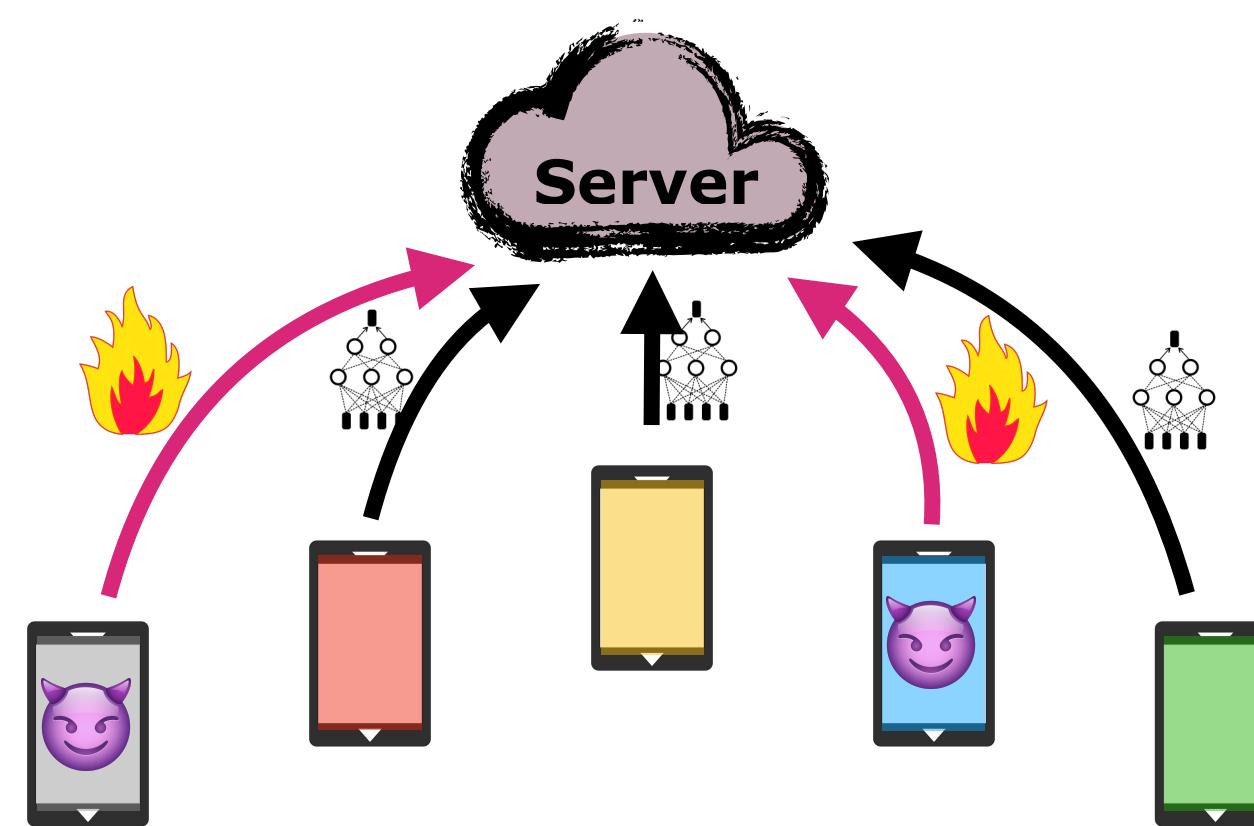
Step 2 of 3: Clients perform some local SGD steps on their local data



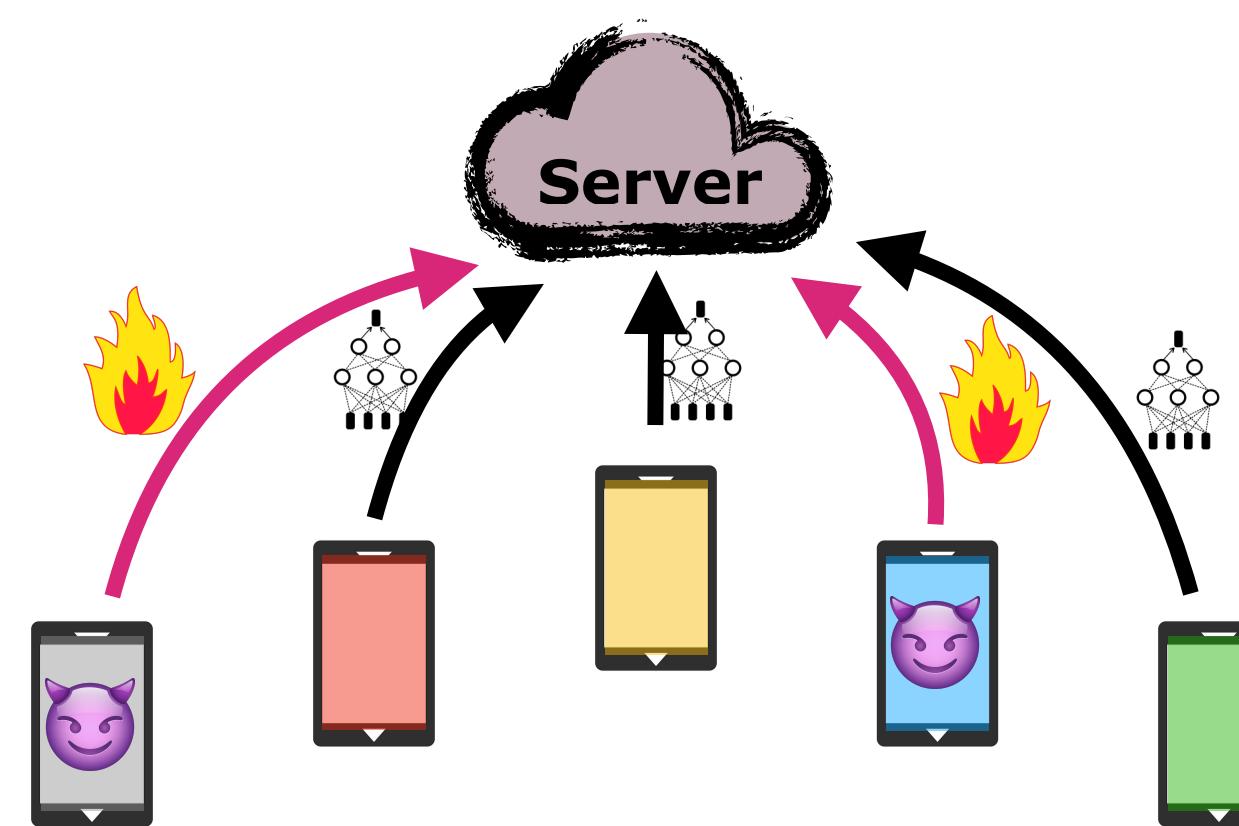
So far, same as FedAvg

*Step 3 of 3: Aggregate with multiple rounds of secure average
(weights β_i from the Weiszfeld Algorithm)*

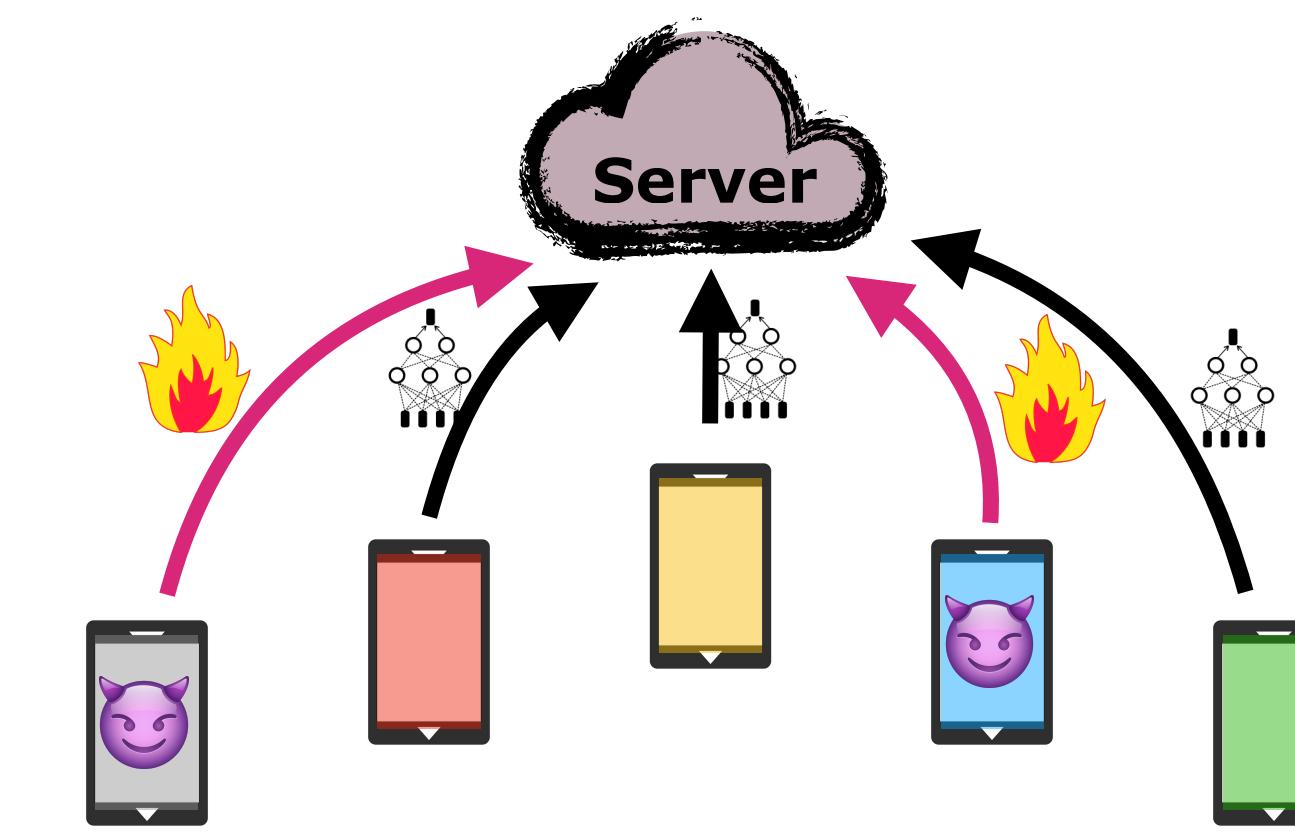
Round 1 of Aggregation



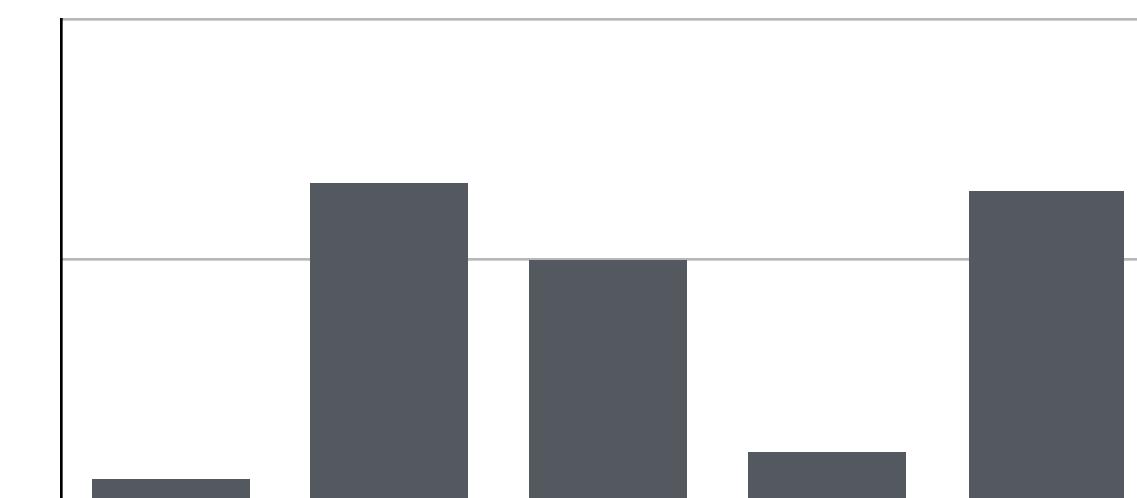
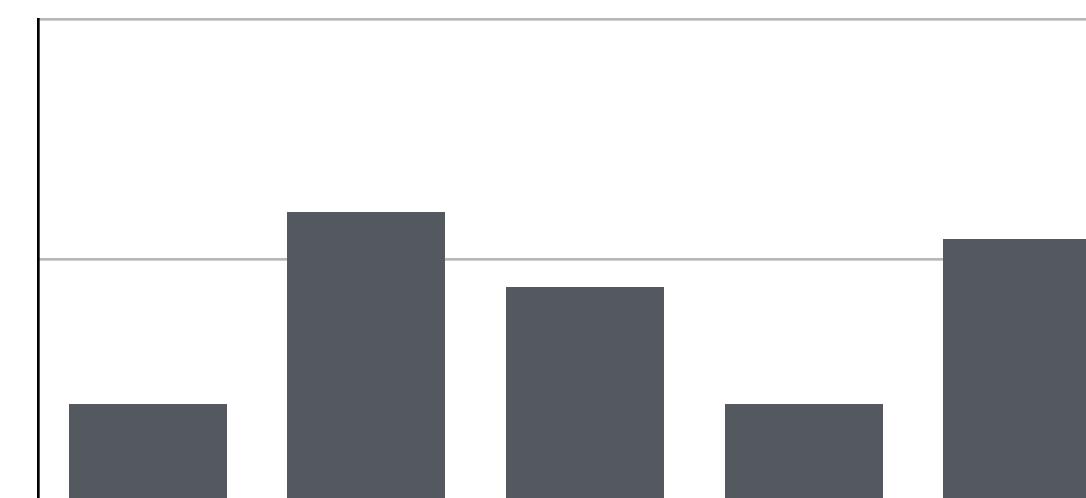
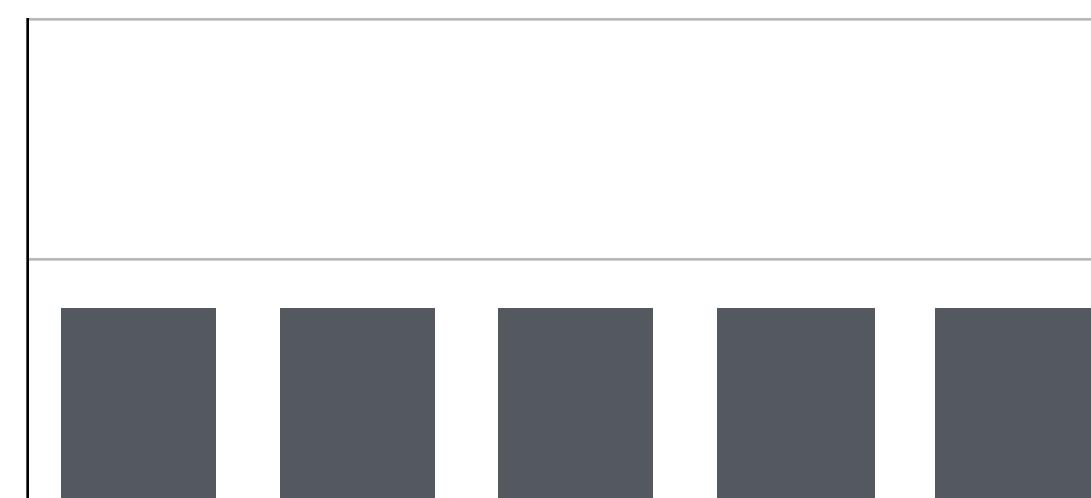
Round 2 of Aggregation



Round 3 of Aggregation



Weights



Convergence analysis (least squares)



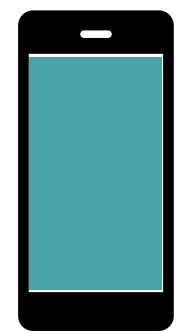
Data
poisoning

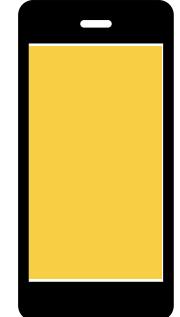
Model
poisoning

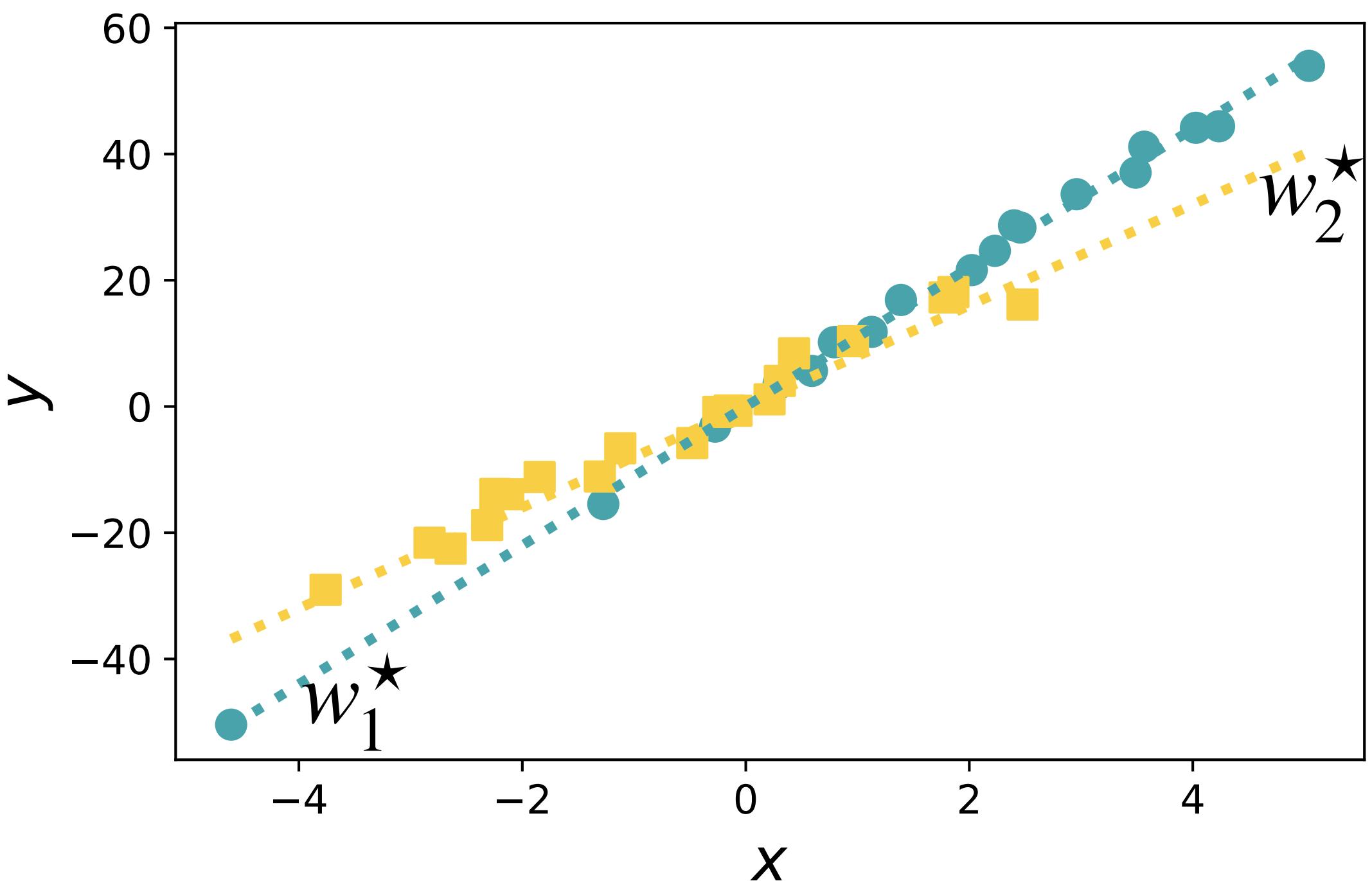
Data on client i : $(X_i, Y_i) \sim p_i$ satisfies

$$Y_i = X_i^\top w_i^* + \xi_i \quad \text{where} \quad \xi_i \sim \mathcal{N}(0, \sigma^2)$$

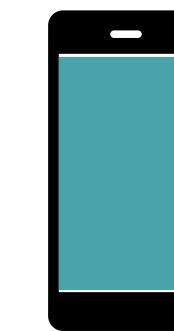
$$\begin{array}{c} \curvearrowleft \\ \mathbb{R} \end{array} \quad \begin{array}{c} \curvearrowright \\ \mathbb{R}^d \end{array}$$


$$X_1 \sim \mathcal{N}(0, H_1)$$

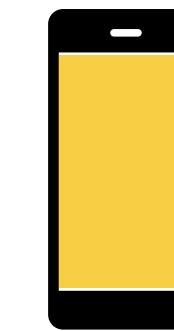

$$X_2 \sim \mathcal{N}(0, H_2)$$



$$X_1 \sim \mathcal{N}(0, H_1)$$



$$X_2 \sim \mathcal{N}(0, H_2)$$



Data on client i : $(X_i, Y_i) \sim p_i$ satisfies

$$Y_i = X_i^\top w_i^* + \xi_i \quad \text{where} \quad \xi_i \sim \mathcal{N}(0, \sigma^2)$$

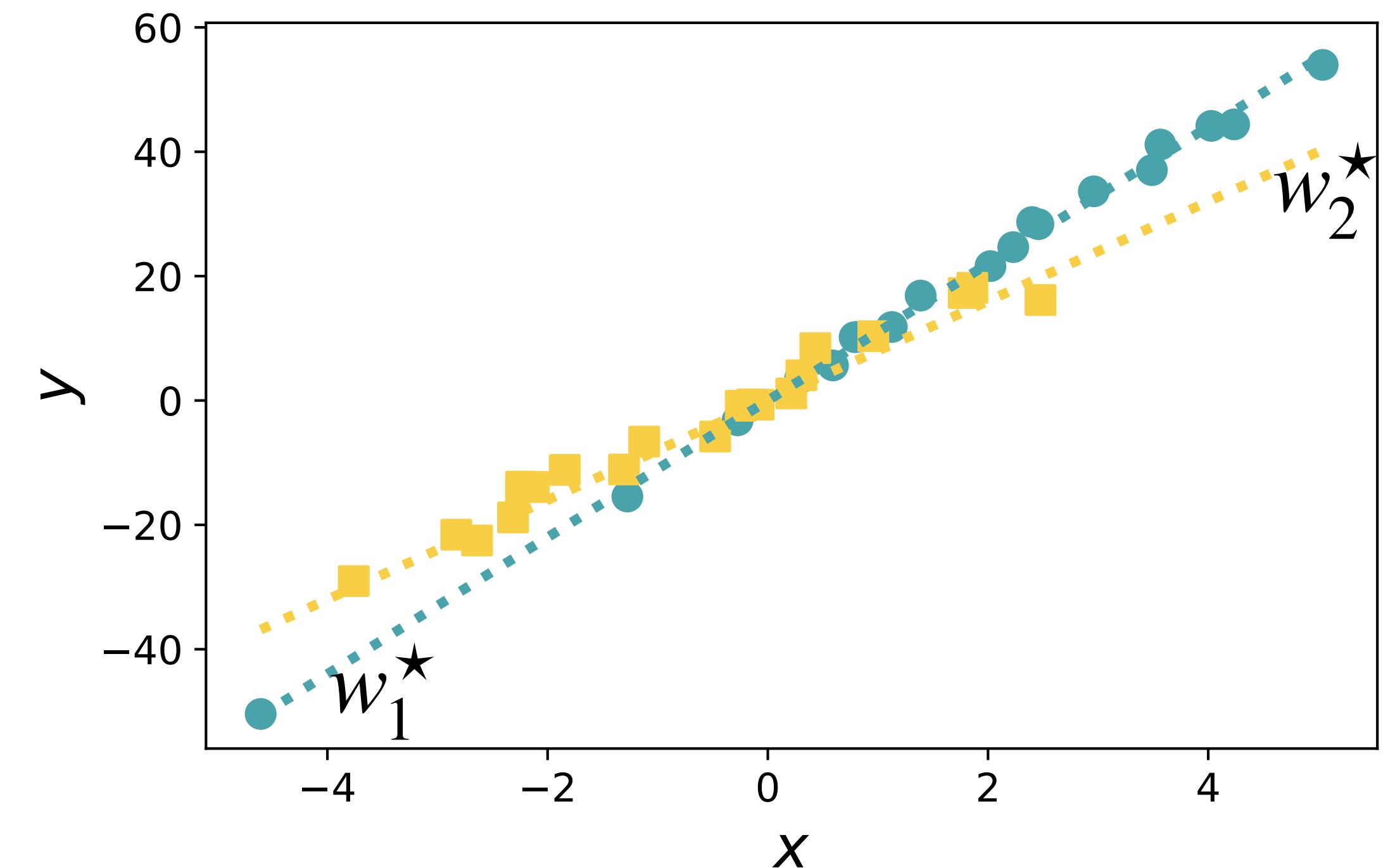
$$\begin{array}{ccc} & \swarrow & \searrow \\ \mathbb{R} & & \mathbb{R}^d \end{array}$$

Measure of heterogeneity

$$\Omega_X = \max_{\substack{i \\ \text{inlier}}} \lambda_{\max}(H^{-1/2} H_i H^{-1/2}) \geq 1$$

$$H_i = \mathbb{E}[X_i X_i^\top] \quad \text{marginal covariance of } i$$

$$H = \frac{1}{n} \sum_i H_i \quad \text{marginal covariance of mixture}$$



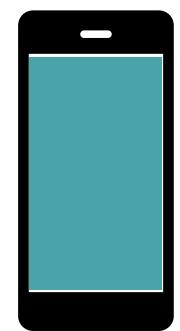
Data on client i : $(X_i, Y_i) \sim p_i$ satisfies

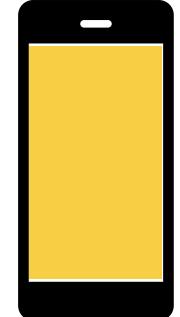
$$Y_i = X_i^\top w_i^* + \xi_i \quad \text{where} \quad \xi_i \sim \mathcal{N}(0, \sigma^2)$$

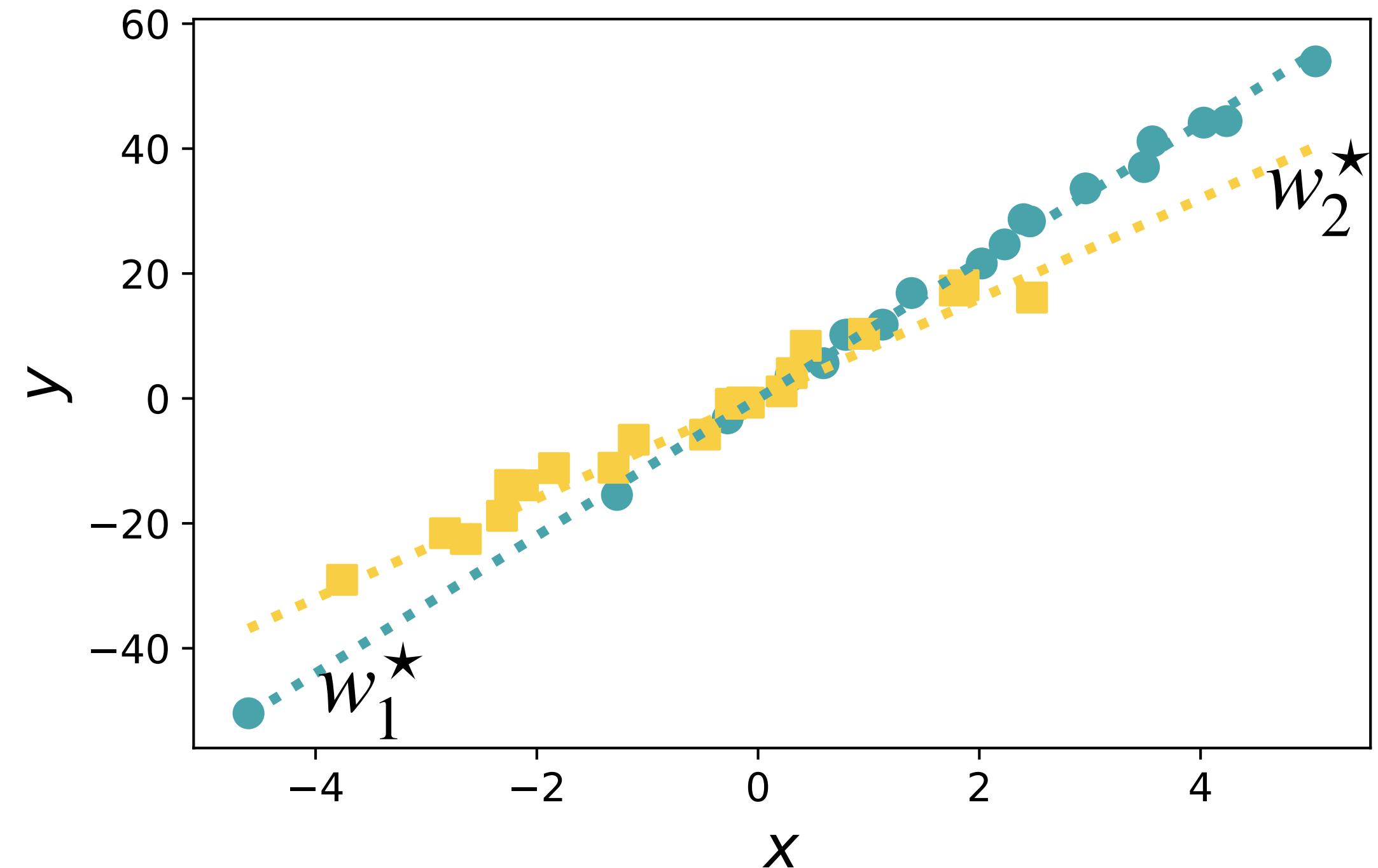
$$\begin{array}{ccc} & \swarrow & \searrow \\ \mathbb{R} & & \mathbb{R}^d \end{array}$$

Measure of heterogeneity

$$\Omega_Y = \max_{\substack{i,j \\ \text{inlier}}} \|w_i^* - w_j^*\|_2 \geq 0$$


$$X_1 \sim \mathcal{N}(0, H_1)$$


$$X_2 \sim \mathcal{N}(0, H_2)$$



$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

$$\text{Number of clients per round} = m$$

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that at least $1/2 + c/2$ non-corrupted devices are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

Optimization error

Statistical error

Heterogeneity Error

GM Approx. Error

$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

Number of clients per round = m

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that at least $1/2 + c/2$ non-corrupted devices are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

Optimization error

Statistical error

Heterogeneity Error

GM Approx. Error

$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

$$\text{Number of clients per round} = m$$

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that **at least $1/2 + c/2$ non-corrupted devices** are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

Optimization error

Statistical error

Heterogeneity Error

GM Approx. Error

$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

$$\text{Number of clients per round} = m$$

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that at least $1/2 + c/2$ non-corrupted devices are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

Optimization error

Statistical error

Heterogeneity Error

GM Approx. Error

$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

$$\text{Number of clients per round} = m$$

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that at least $1/2 + c/2$ non-corrupted devices are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

Optimization error

Statistical error

Heterogeneity Error

GM Approx. Error

$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

$$\text{Number of clients per round} = m$$

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that at least $1/2 + c/2$ non-corrupted devices are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

Optimization error

Statistical error

Heterogeneity Error

GM Approx. Error

$$\text{Fraction of non-corrupted clients} = \frac{1}{2} + c$$

$$\text{Number of clients per round} = m$$

\mathcal{E} holds w.h.p. if $m \gtrsim \frac{1}{c^2}$

Theorem [P., Kakade, Harchaoui]

Assume that $F(w)$ is strongly convex, $\|X_i\| \leq 1$ and number of local steps $\propto 2^t$. Let \mathcal{E} denote the event that at least $1/2 + c/2$ non-corrupted devices are chosen in each round.

Then, RFA with ε -approximate GM satisfies

$$\mathbb{E} [\|w_t - w^\star\|^2 | \mathcal{E}] \lesssim \frac{\|w_0 - w^\star\|^2}{2^t} + \frac{1}{c^2} \left(d\sigma^2 \frac{t}{2^t} + \frac{\varepsilon^2}{m^2} + \Omega_X^2 \Omega_Y^2 \right)$$

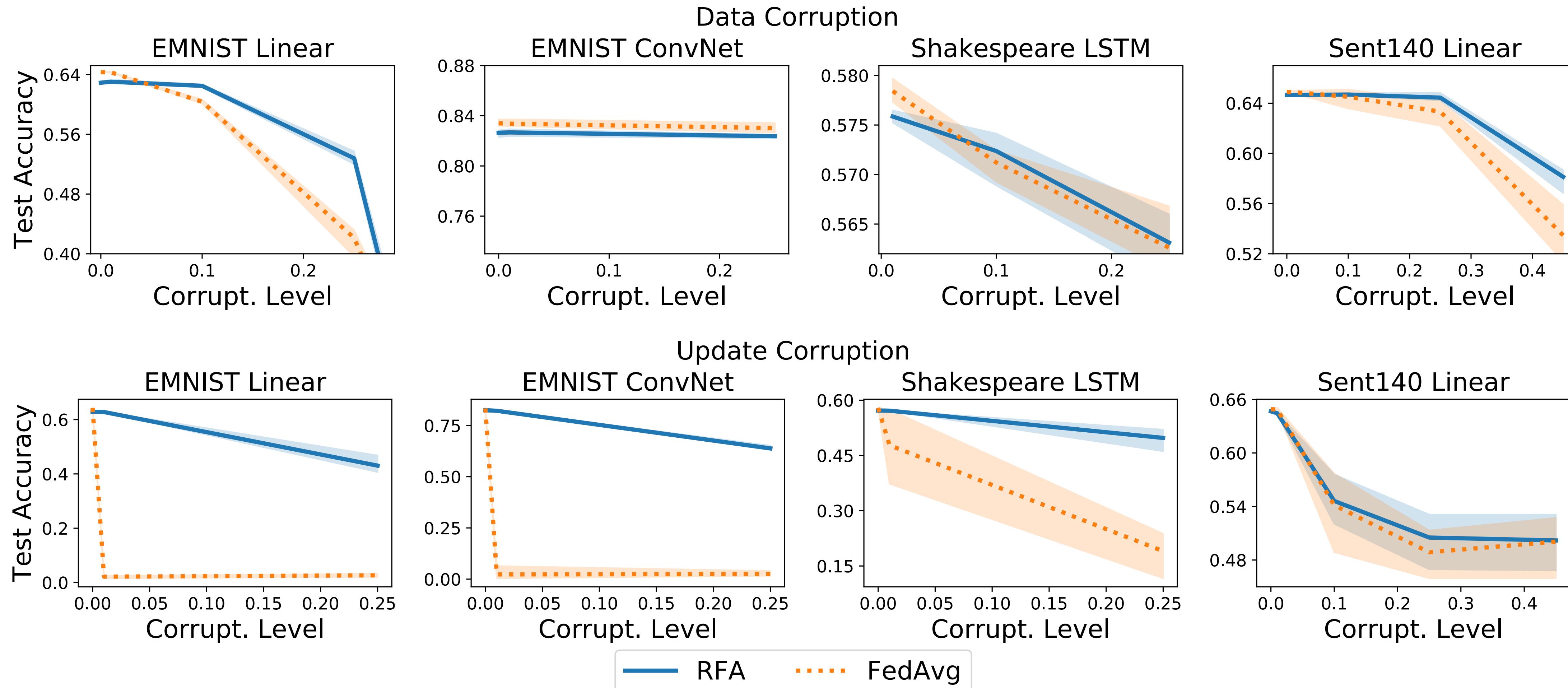
Optimization error

Statistical error

Heterogeneity Error

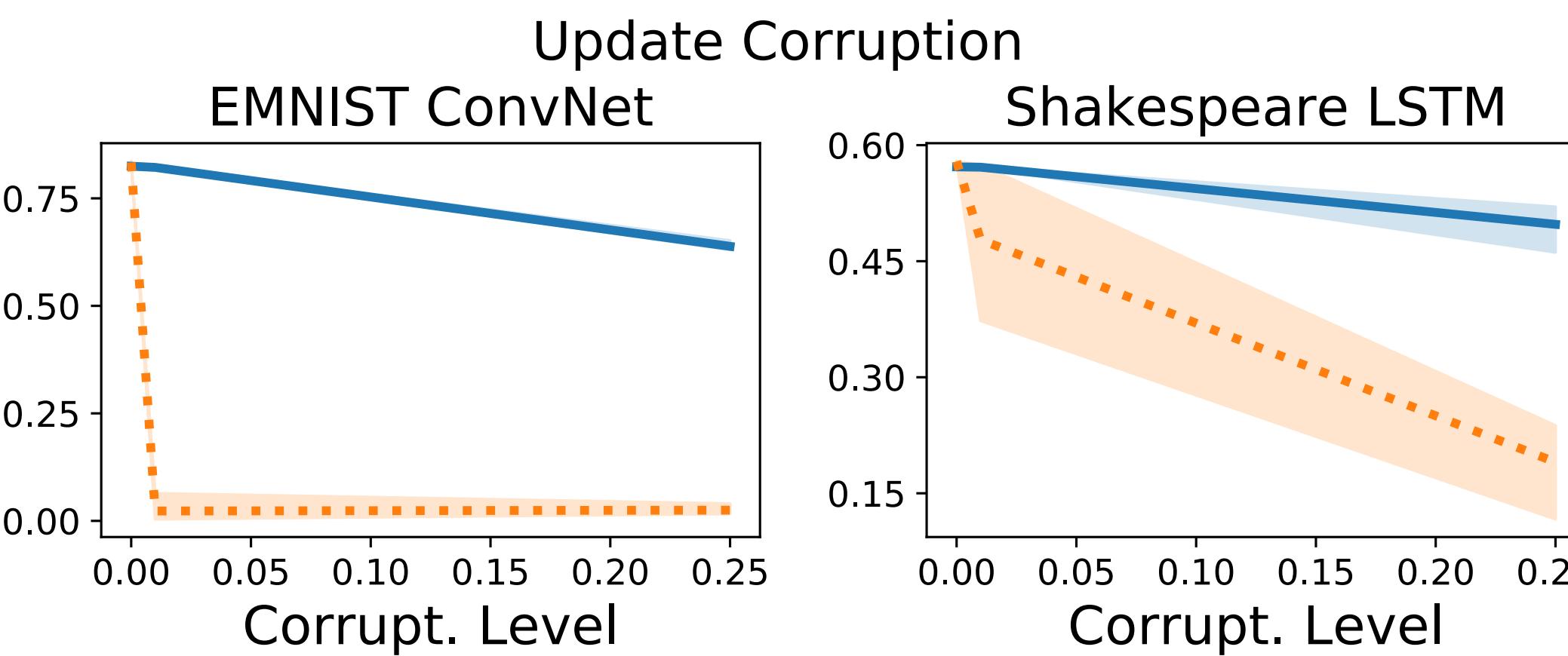
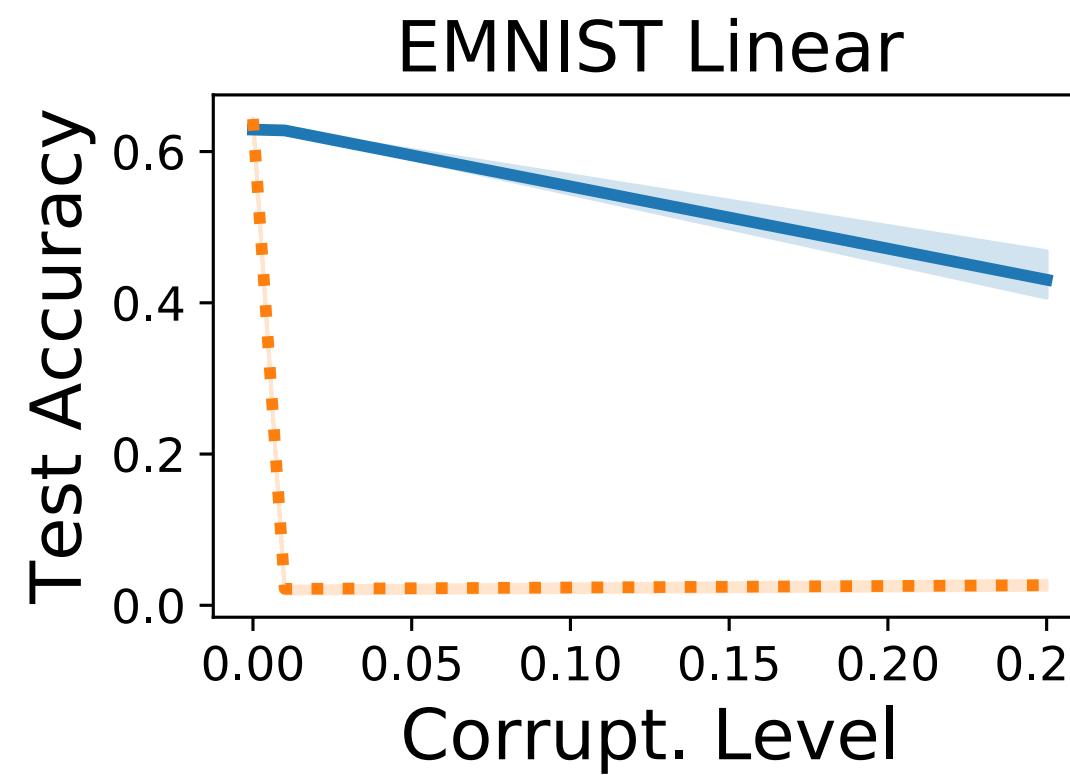
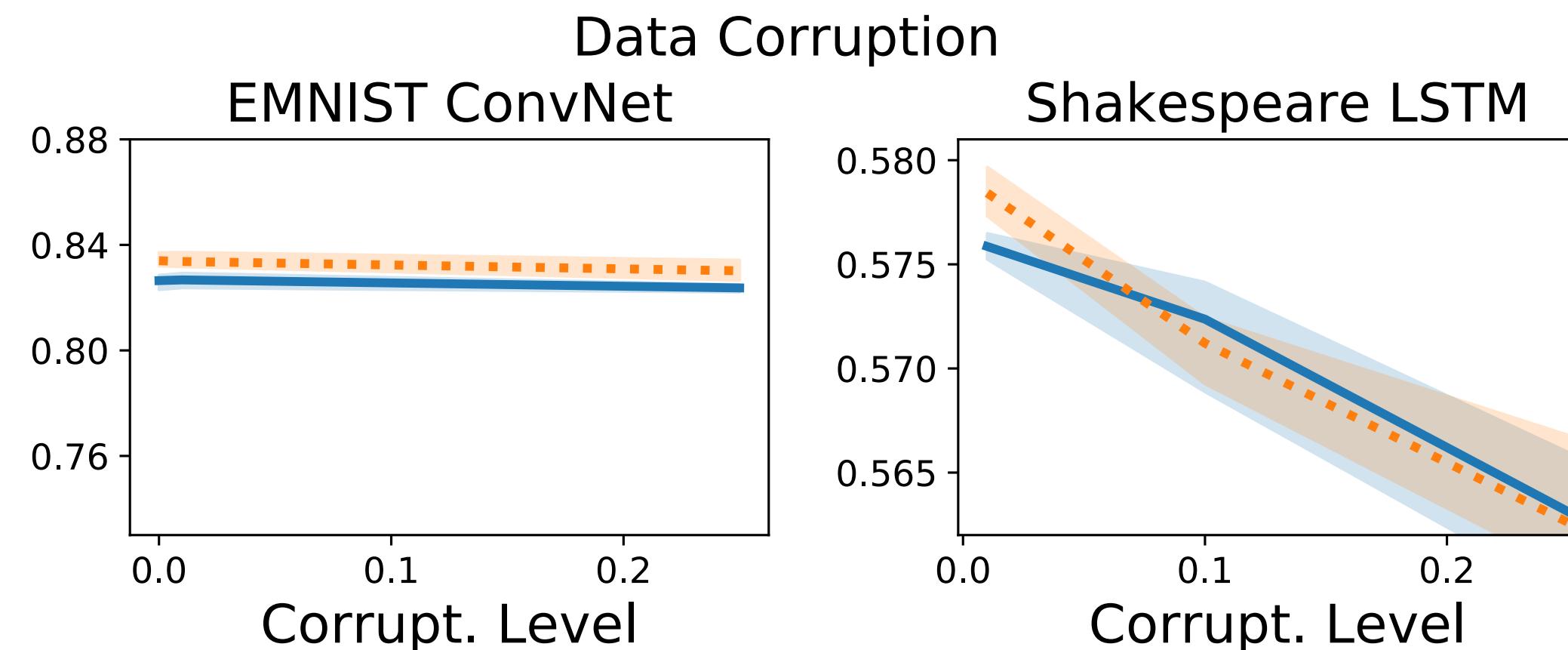
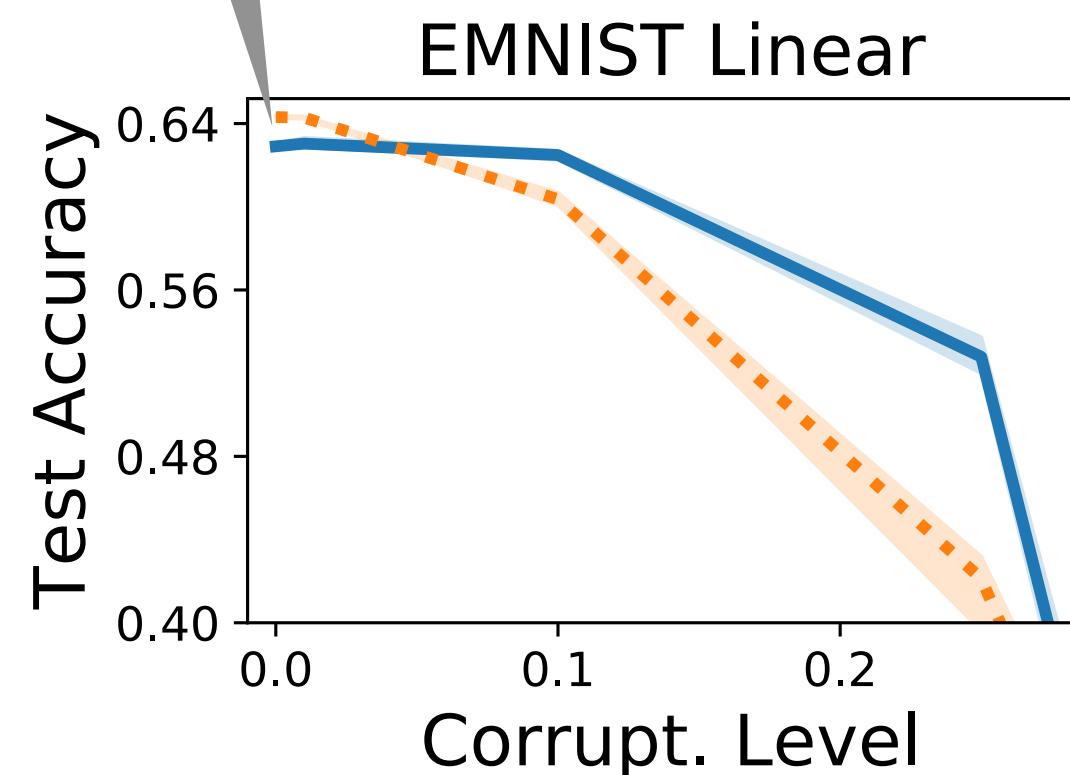
GM Approx. Error

Experiments



Experiments

1.4pp
gap at zero
corruption



— RFA - - - FedAvg

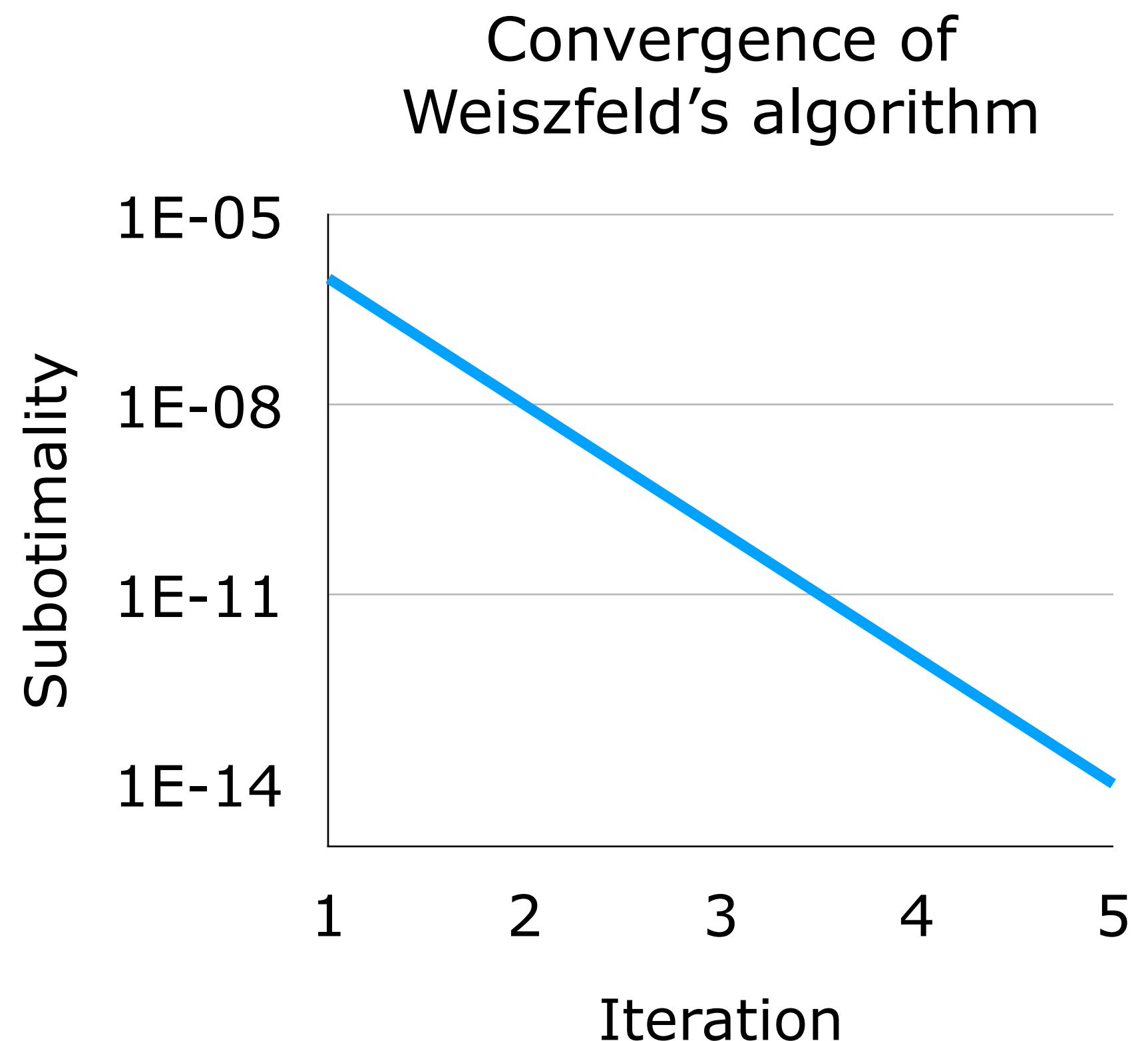
Reducing RFA's communication

One round of RFA \implies 3-5 rounds of communication (Weiszfeld)

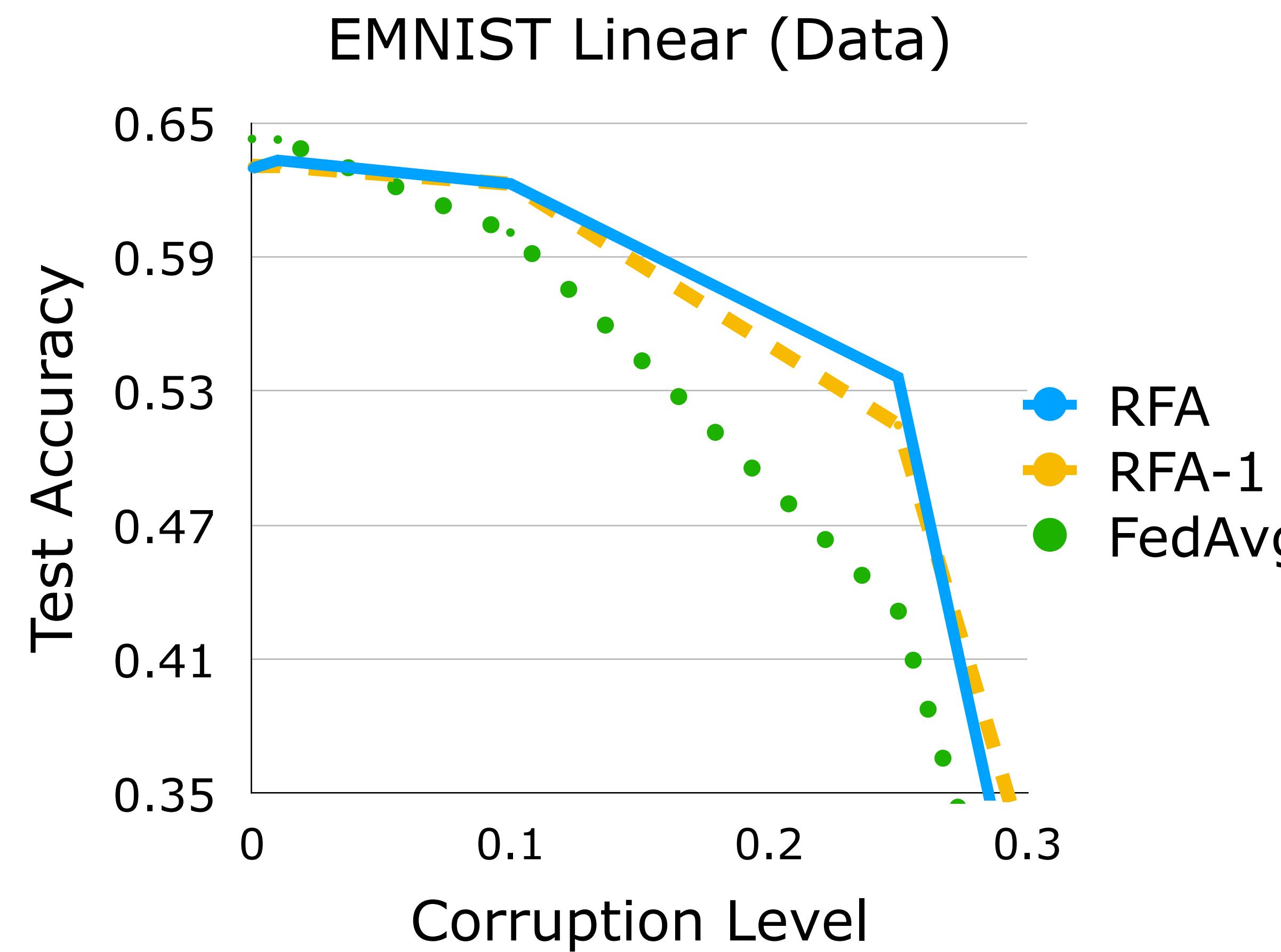
Does 1 round of communication give robustness?

$$\beta_i = \frac{1}{\max\{\|w_i\|_2, \nu\}}$$

$$z = \frac{\sum_i \beta_i w_i}{\sum_i \beta_{i,t}}$$



Does 1 round of communication give robustness? **Yes!**



Fast and differentiable geometric median

```
import torch
from geom_median.torch import compute_geometric_median # PyTorch API
# from geom_median.numpy import compute_geometric_median # NumPy API

points = [torch.rand(d) for _ in range(n)] # list of n tensors of shape (d,)
# The shape of each tensor is the same and can be arbitrary (not necessarily 1-dimensional)
weights = torch.rand(n) # non-negative weights of shape (n,)
out = compute_geometric_median(points, weights)
# Access the median via `out.median`, which has the same shape as the points, i.e., (d,)
```

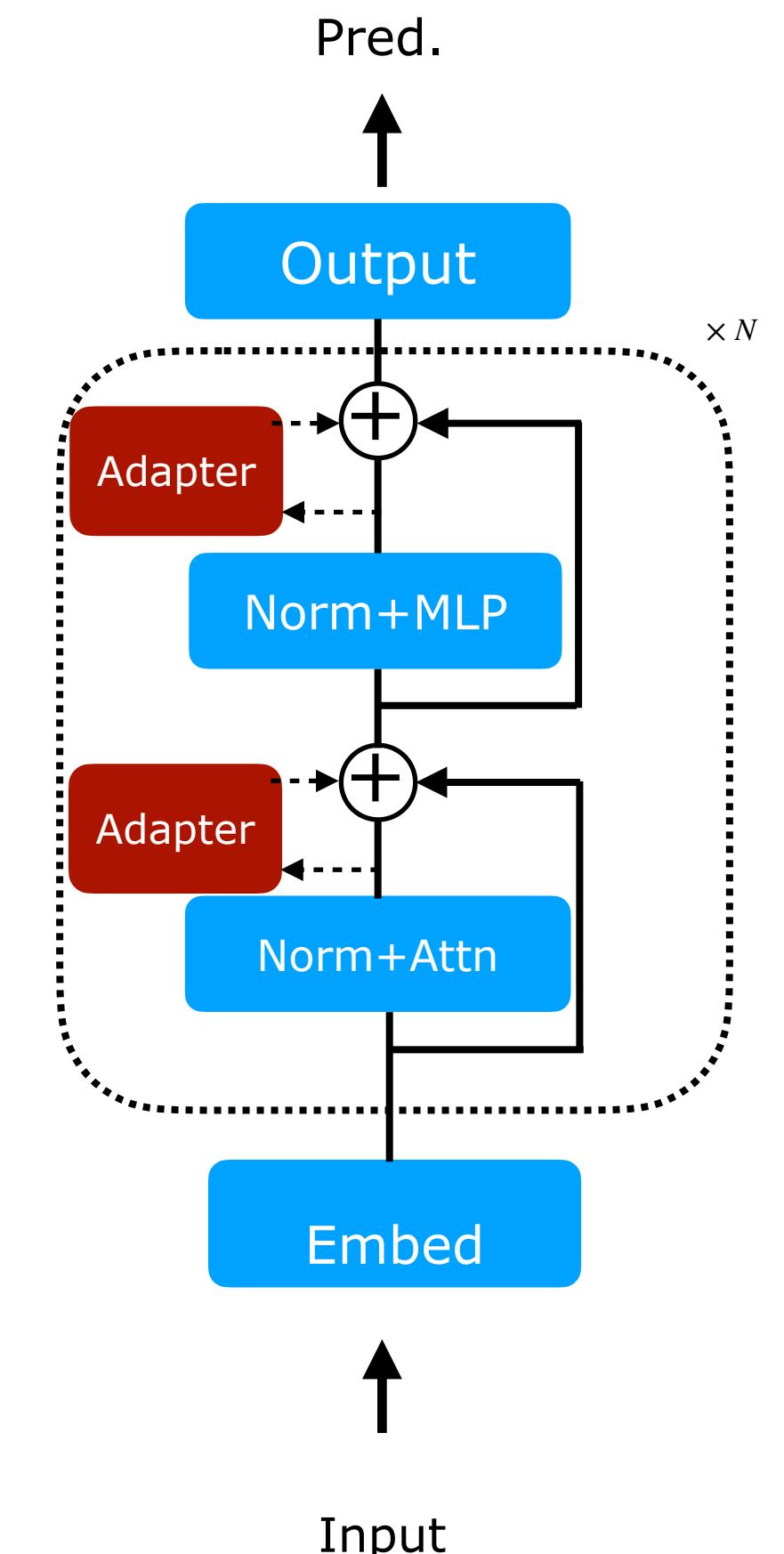
Install: [pip install geom-median](#)

Documentation: github.com/krishnap25/geom-median

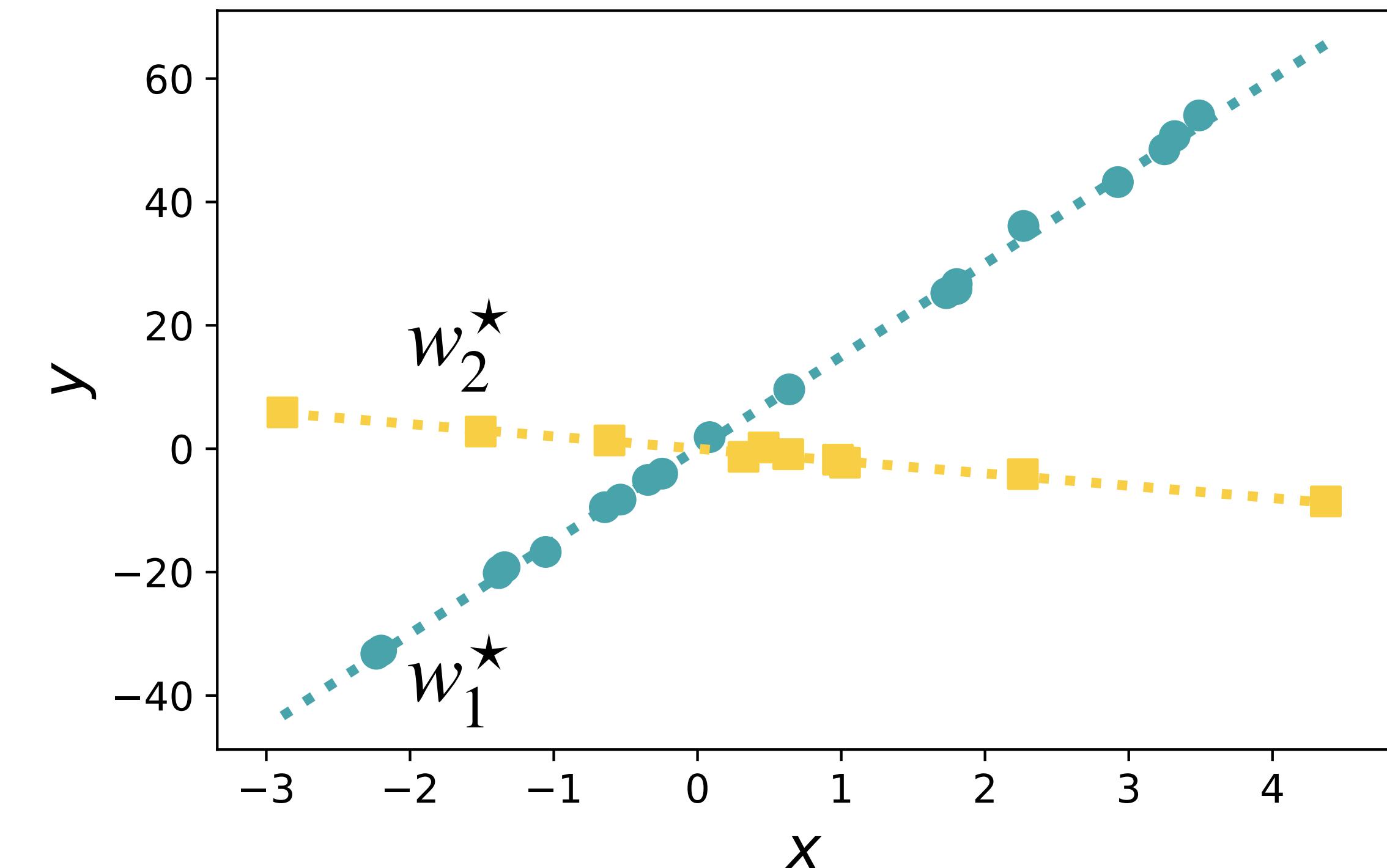
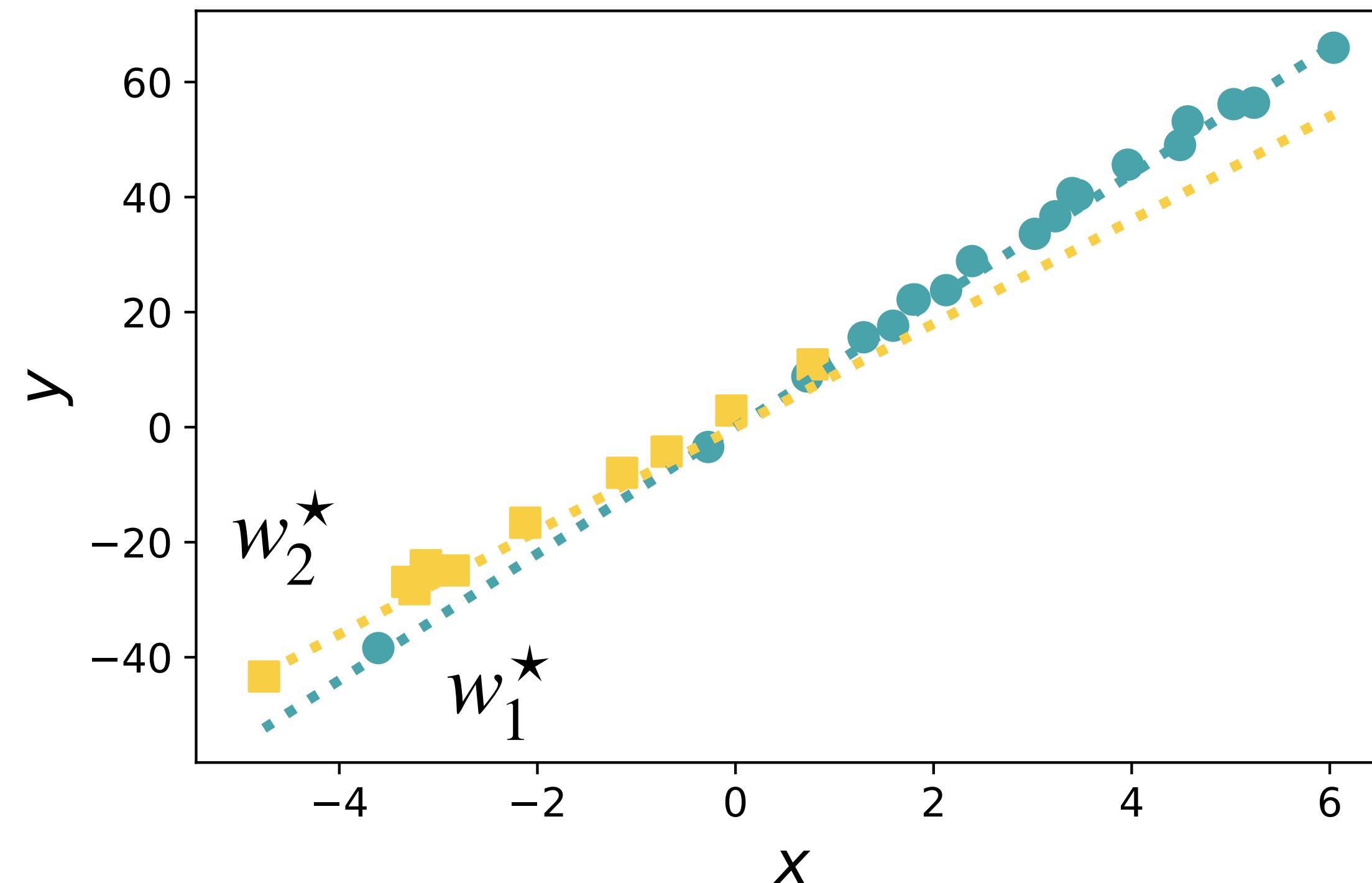


Part 3: Model personalization for federated learning

[TSP '22, ICML '22]



Two regimes



Objective

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n F_i(w)$$

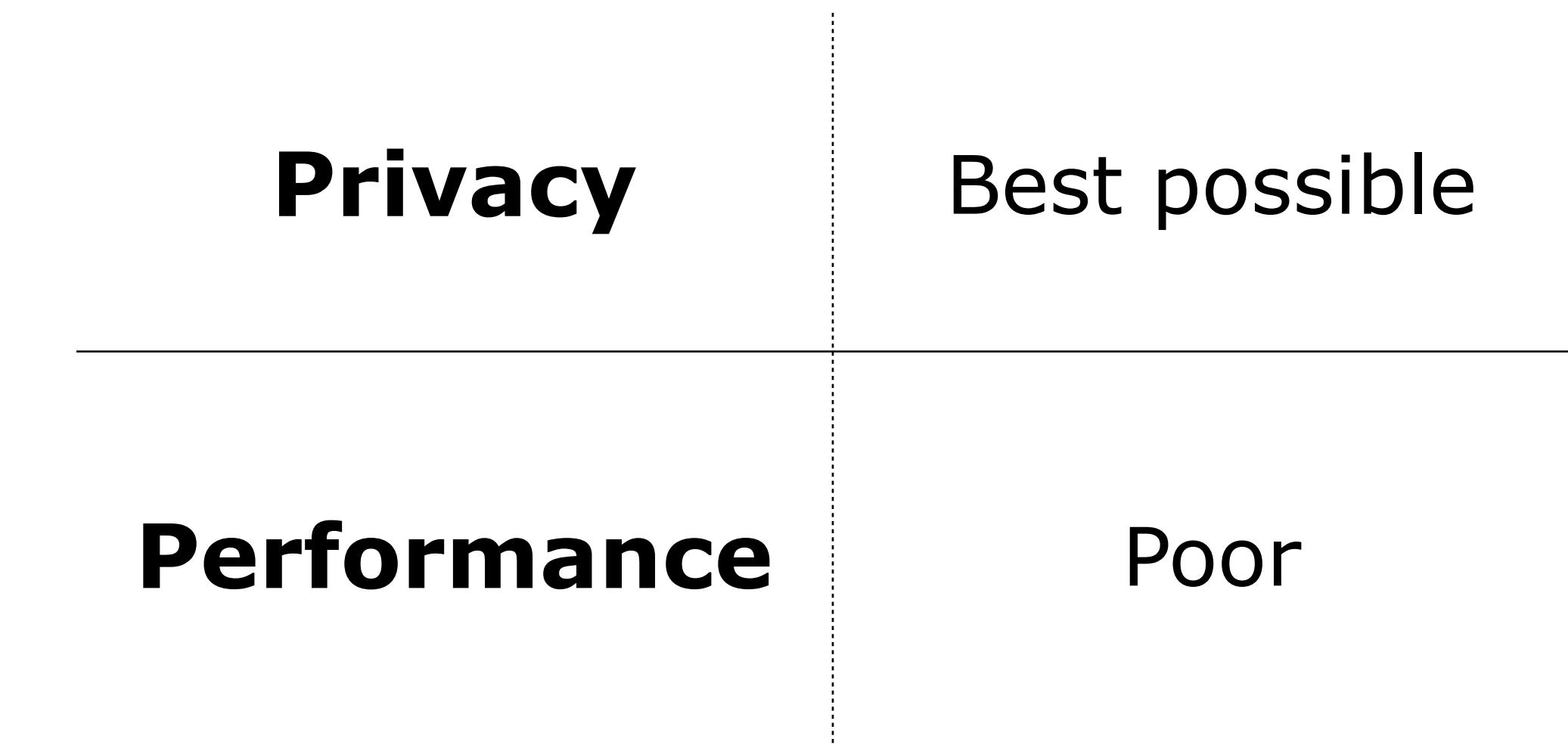
where

$$F_i(w) = \mathbb{E}_{z \sim p_i} [f(w; z)]$$

loss on client i

Option 1: Train a separate model per client (no collaboration)

Objective: $\min_{w_i} F_i(w_i)$



Option 2: The model has a global component and a per-client component

Shared Params u + Personal Params v_i = Full model $w_i = (u, v_i)$

Objective: $\min_{u, v_1, \dots, v_n} \frac{1}{n} \sum_{i=1}^n F_i(u, v_i)$

Example: $F_i(u, v_i) = \mathbb{E}_{(X, Y) \sim p_i} \left(\phi_g(X; u) + \phi_l(X; v_i) - Y \right)^2$

Option 2: The model has a global component and a per-client component

$$\text{Shared Params } u + \text{Personal Params } v_i = \text{Full model } w_i = (u, v_i)$$

Objective: $\min_{u, v_1, \dots, v_n} \frac{1}{n} \sum_{i=1}^n F_i(u, v_i)$

Privacy	Yes
---------	-----

Performance	
-------------	--

*data and personal
params on client*

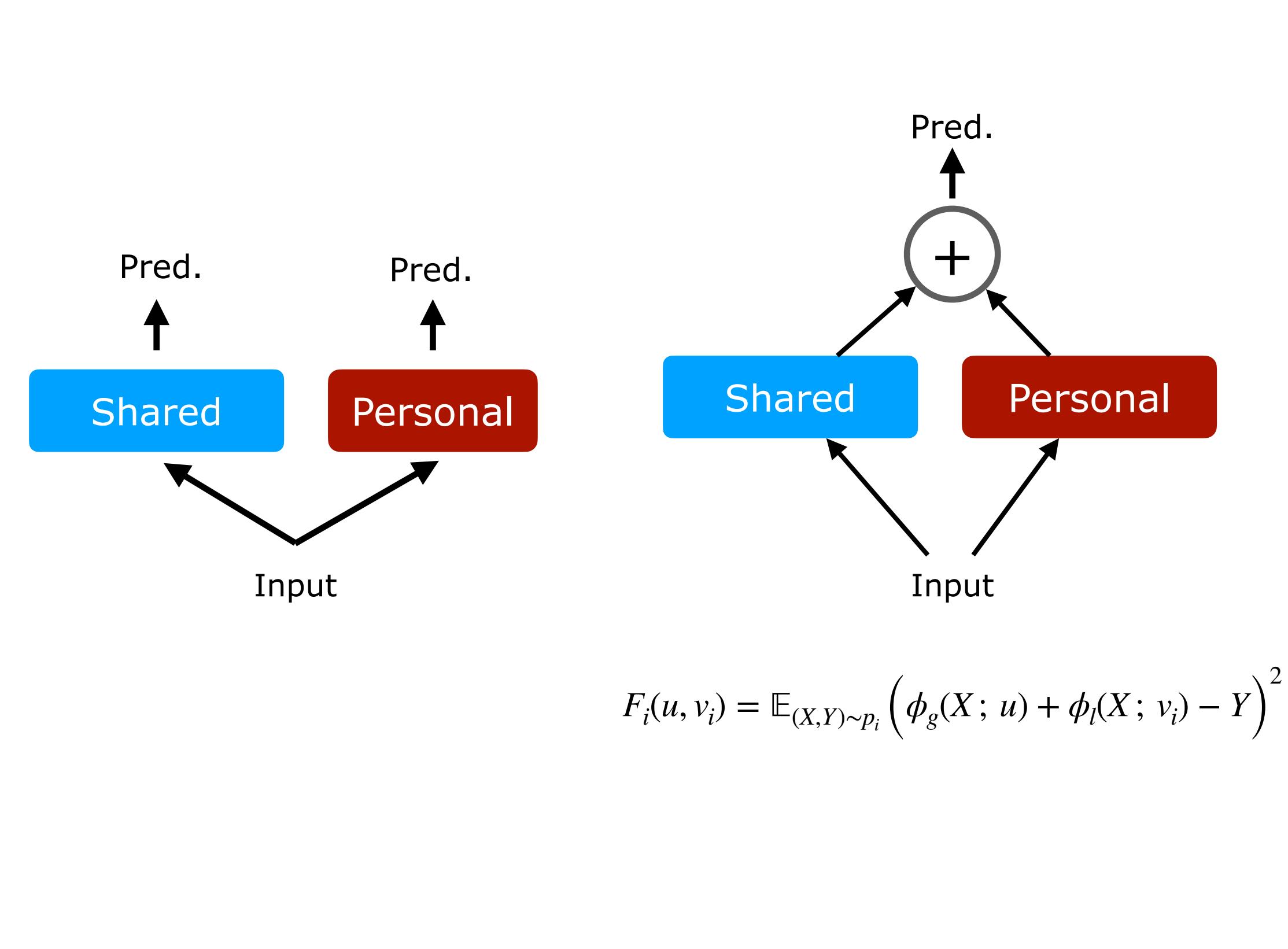
Option 2: The model has a global component and a per-client component

$$\text{Shared Params } u + \text{Personal Params } v_i = \text{Full model } w_i = (u, v_i)$$

Objective: $\min_{u, v_1, \dots, v_n} \frac{1}{n} \sum_{i=1}^n F_i(u, v_i)$

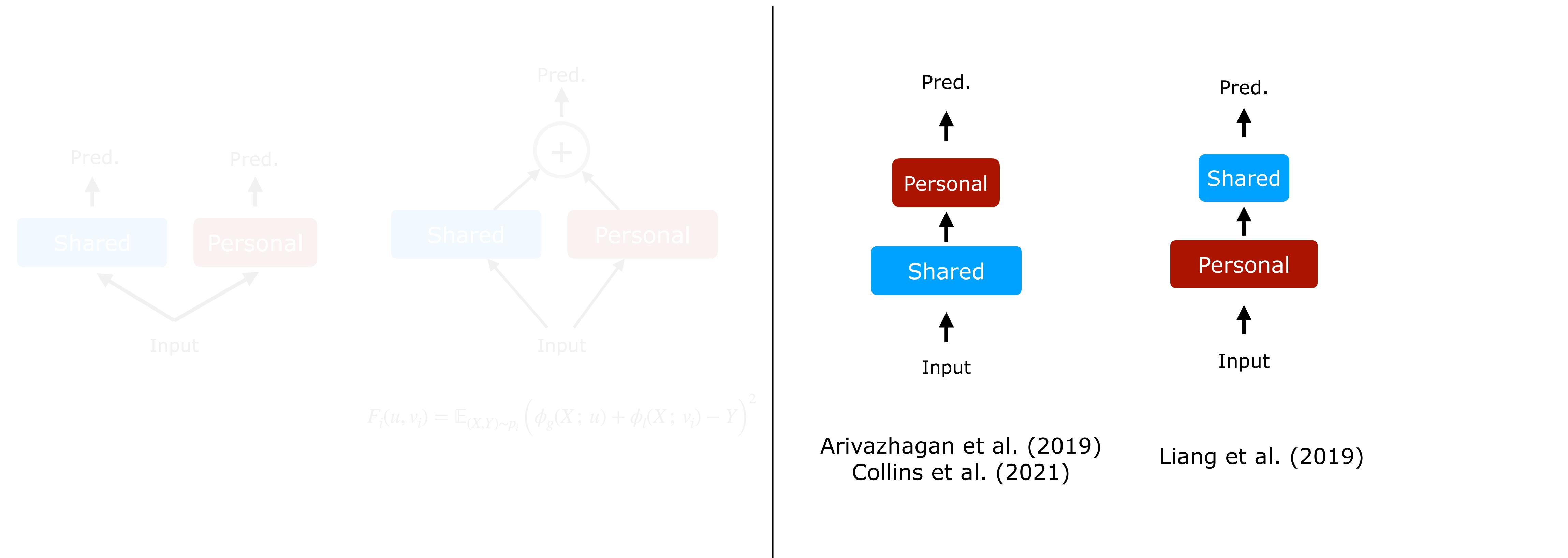
Privacy	Yes	<i>data and personal params on client</i>
Performance	Yes	

Personalization Architectures



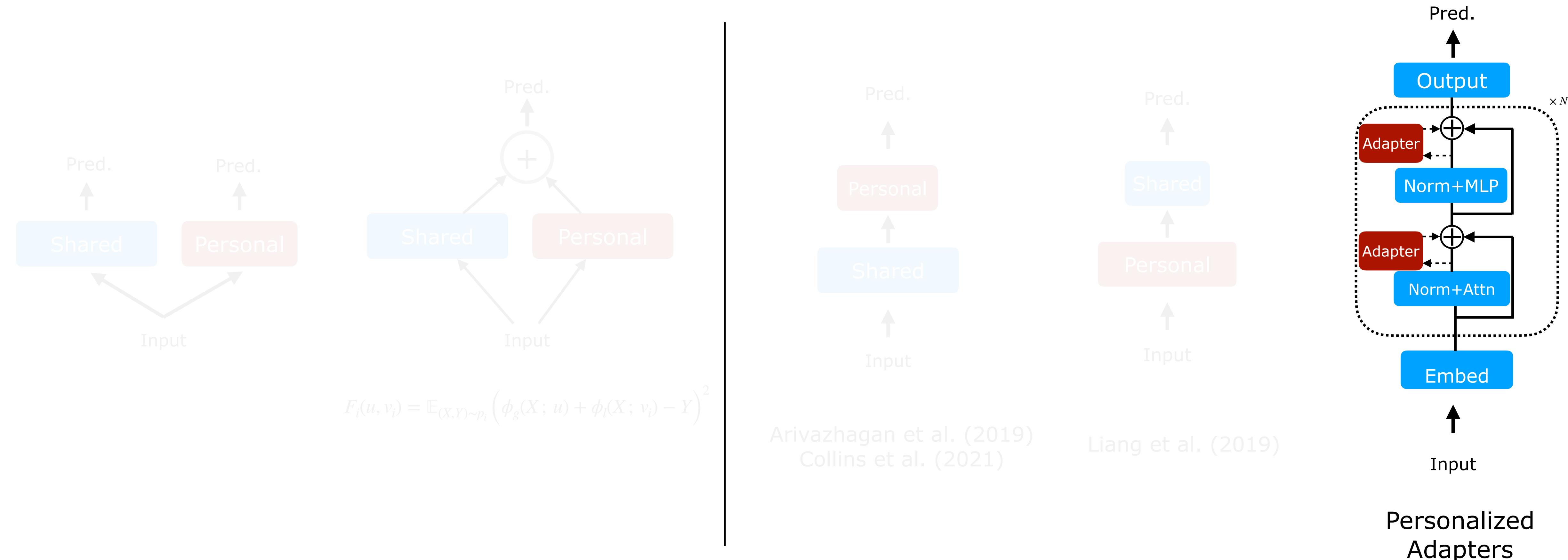
Multi-task learning: Caruana (1997), Baxter (2000), Evgeniou & Pontil (2004), Collobert & Weston (2005), Argyriou et al. (2008), ...

Personalization Architectures



Multi-task learning: Caruana (1997), Baxter (2000), Evgeniou & Pontil (2004), Collobert & Weston (2005), Argyriou et al. (2008), ...

Personalization Architectures



Multi-task learning: Caruana (1997), Baxter (2000), Evgeniou & Pontil (2004), Collobert & Weston (2005), Argyriou et al. (2008), ...

Optimization

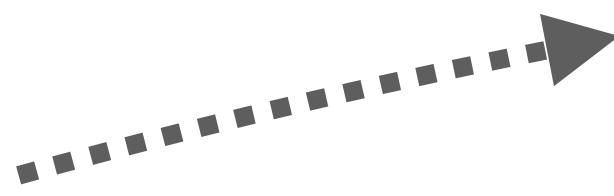
Alternating update

- Server samples m clients and broadcast global model u

$$v_i^+ = v_i - \gamma \nabla_v F_i(u, v_i)$$

- **Local updates** on client i :

$$(u_i^+, v_i^+) = \text{LocalUpdate}_i(u, v_i)$$



$$u_i^+ = u - \gamma \nabla_u F_i(u, v_i^+)$$

- Aggregate updates to global part of the model:

$$u^+ = \frac{1}{m} \sum_i u_i^+$$



$$v_i^+ = v_i - \gamma \nabla_v F_i(u, v_i)$$

$$u_i^+ = u - \gamma \nabla_u F_i(u, v_i)$$

Simultaneous update

Theorem [P., Malik, Mohamed, Rabbat, Sanjabi, Xiao]

For smooth, nonconvex functions, we have the rates:

Alternating update: $\frac{\sigma_1^2}{\sqrt{t}}$

Simultaneous update: $\frac{\sigma_2^2}{\sqrt{t}}$

where $\sigma_1^2 < \sigma_2^2$ under typical scenarios

Experimentally, small but consistent trend of alternating > simultaneous

Alternating update

$$\textcolor{red}{v_i^+} = v_i - \gamma \nabla_v F_i(u, v_i)$$

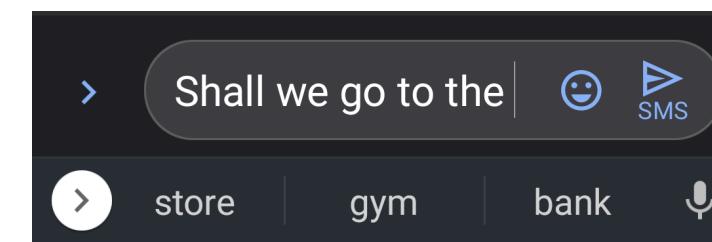
$$\textcolor{red}{u_i^+} = u - \gamma \nabla_u F_i(u, \textcolor{red}{v_i^+})$$

Simultaneous update

$$\textcolor{red}{v_i^+} = v_i - \gamma \nabla_v F_i(u, v_i)$$

$$\textcolor{red}{u_i^+} = u - \gamma \nabla_u F_i(u, v_i)$$

Experiments



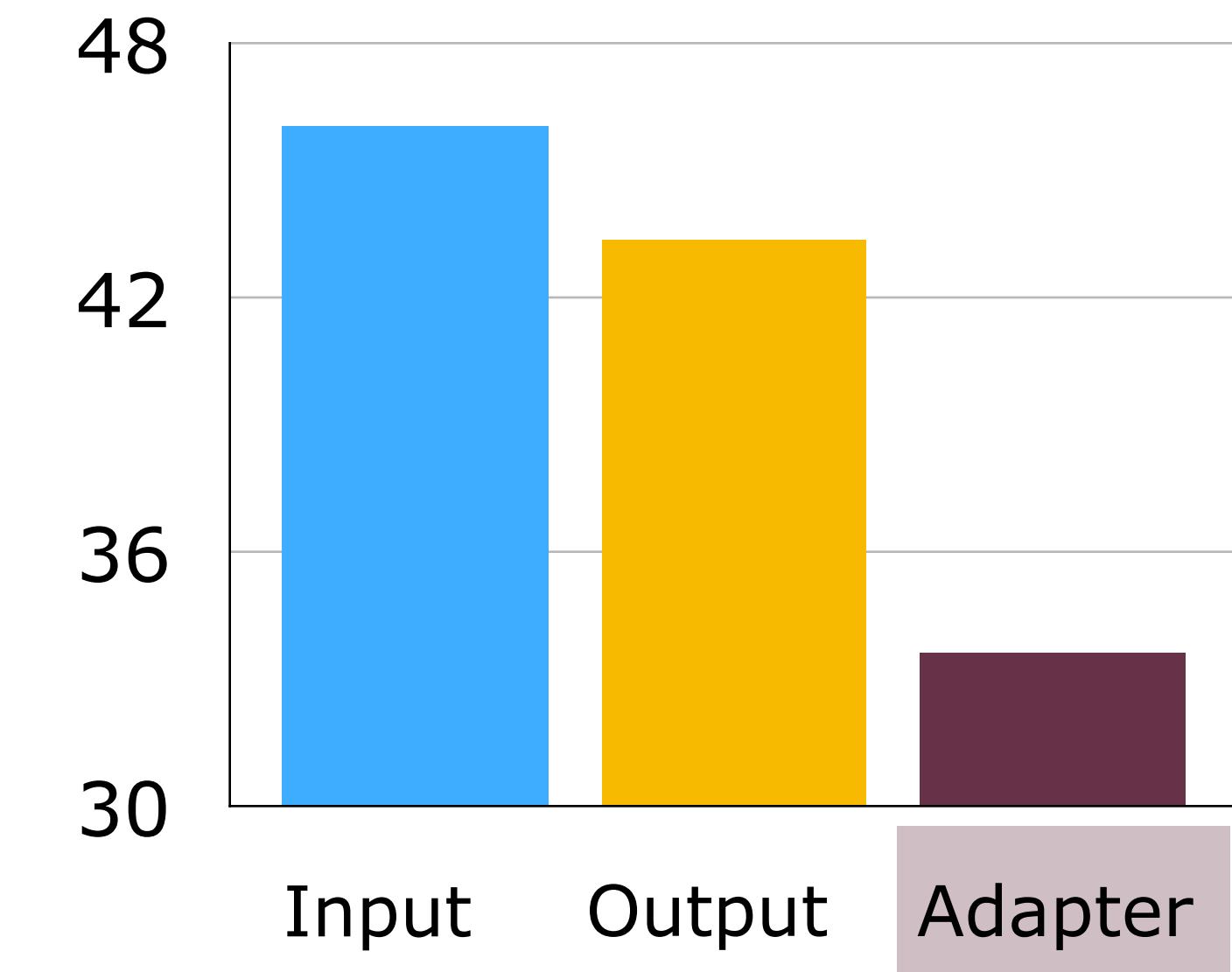
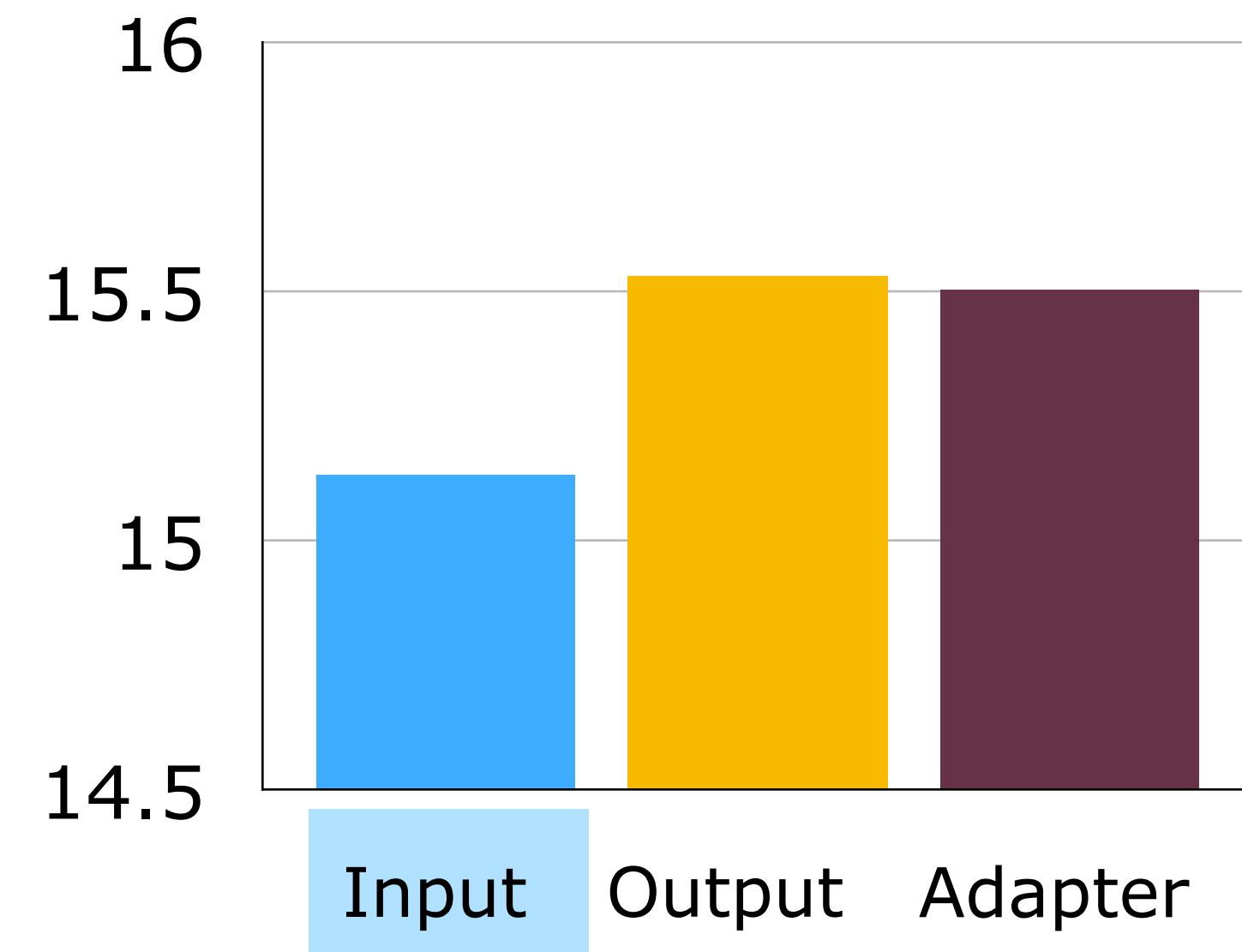
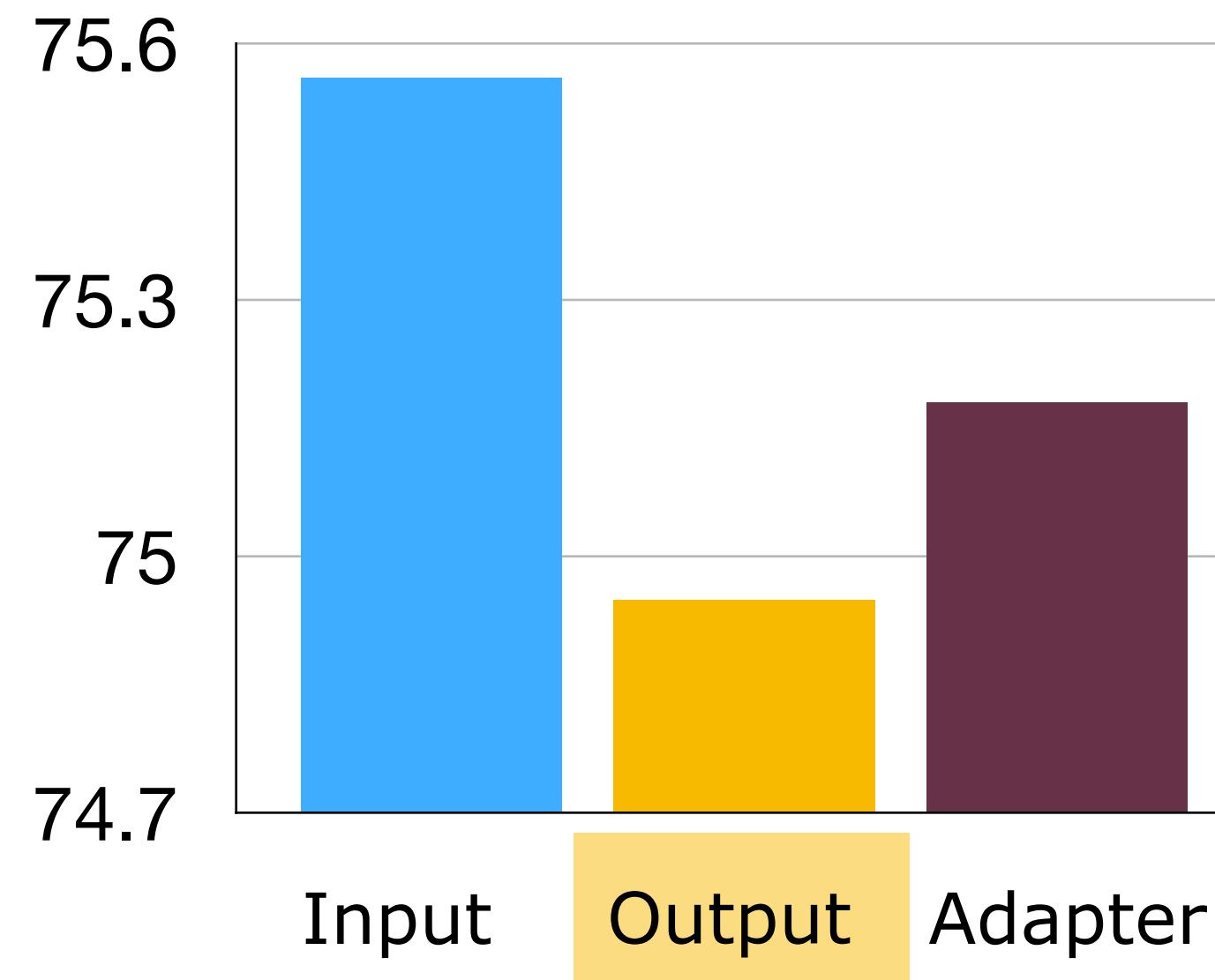
Next word prediction



Speech recognition

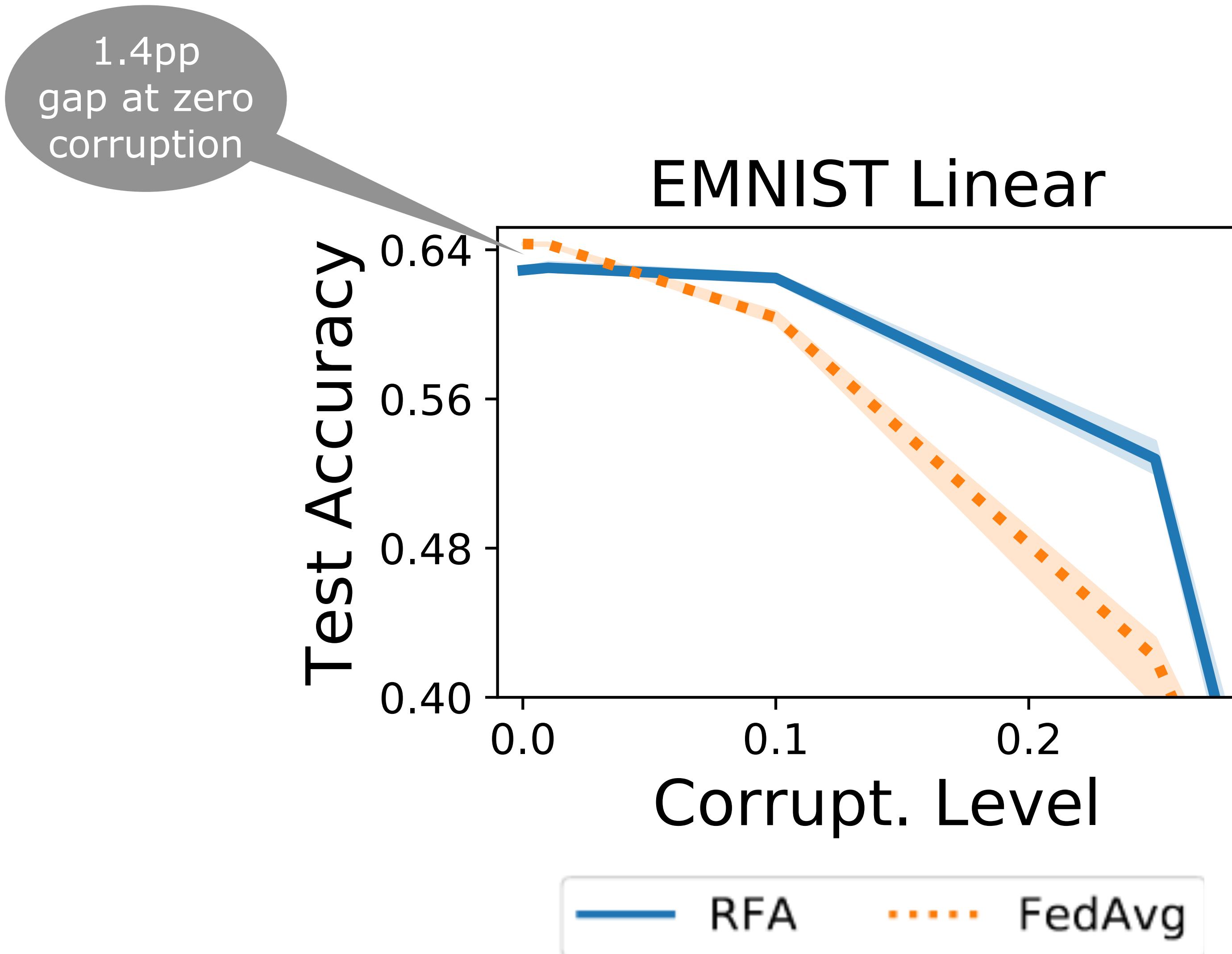


Landmark detection

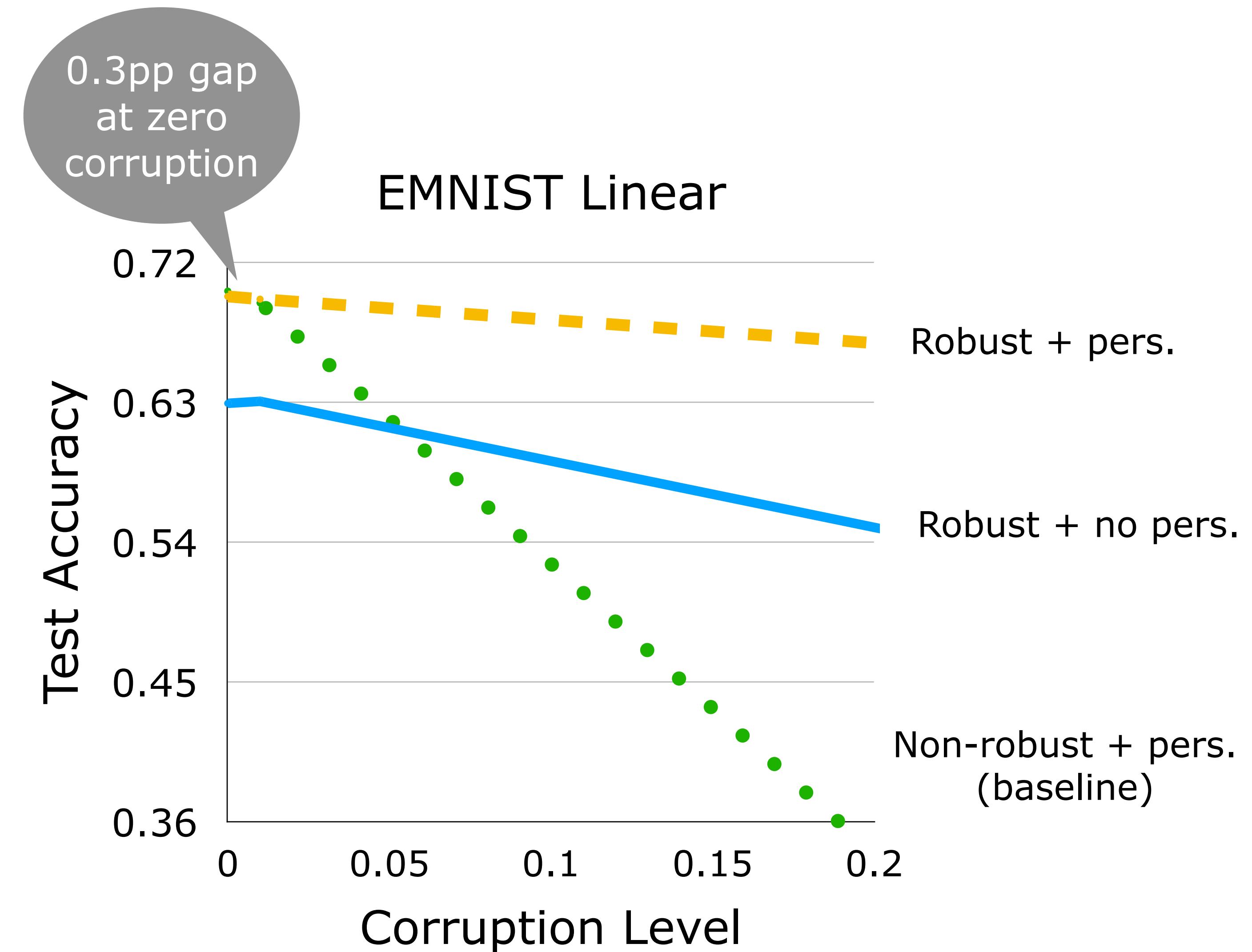
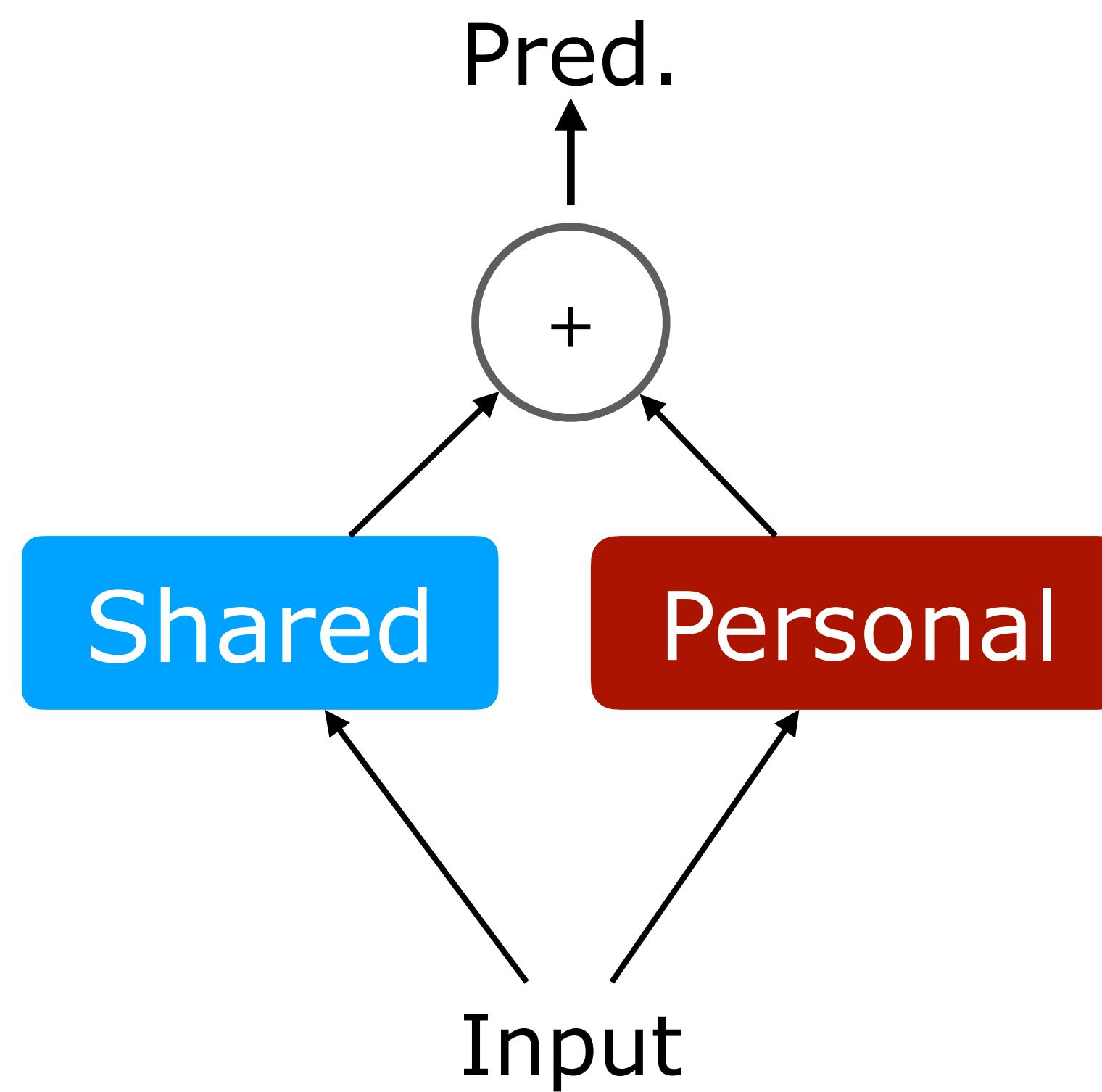


y-axis shows error: lower is better

Recall: robust aggregation experiments



Improving robust aggregation with personalization



Summary

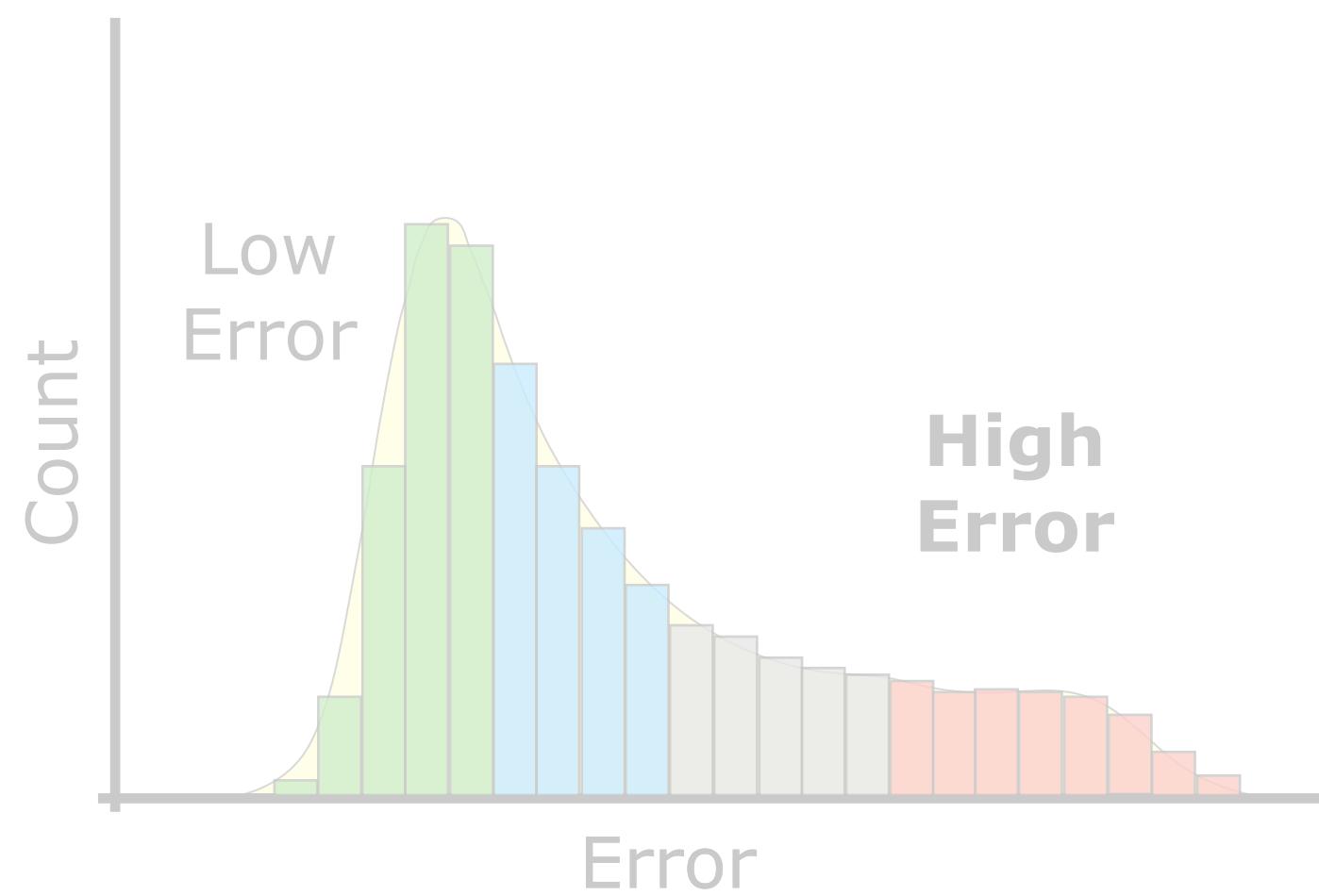
Part 1: Heterogeneity-aware objectives for federated learning

Heterogeneity \implies
large tail errors

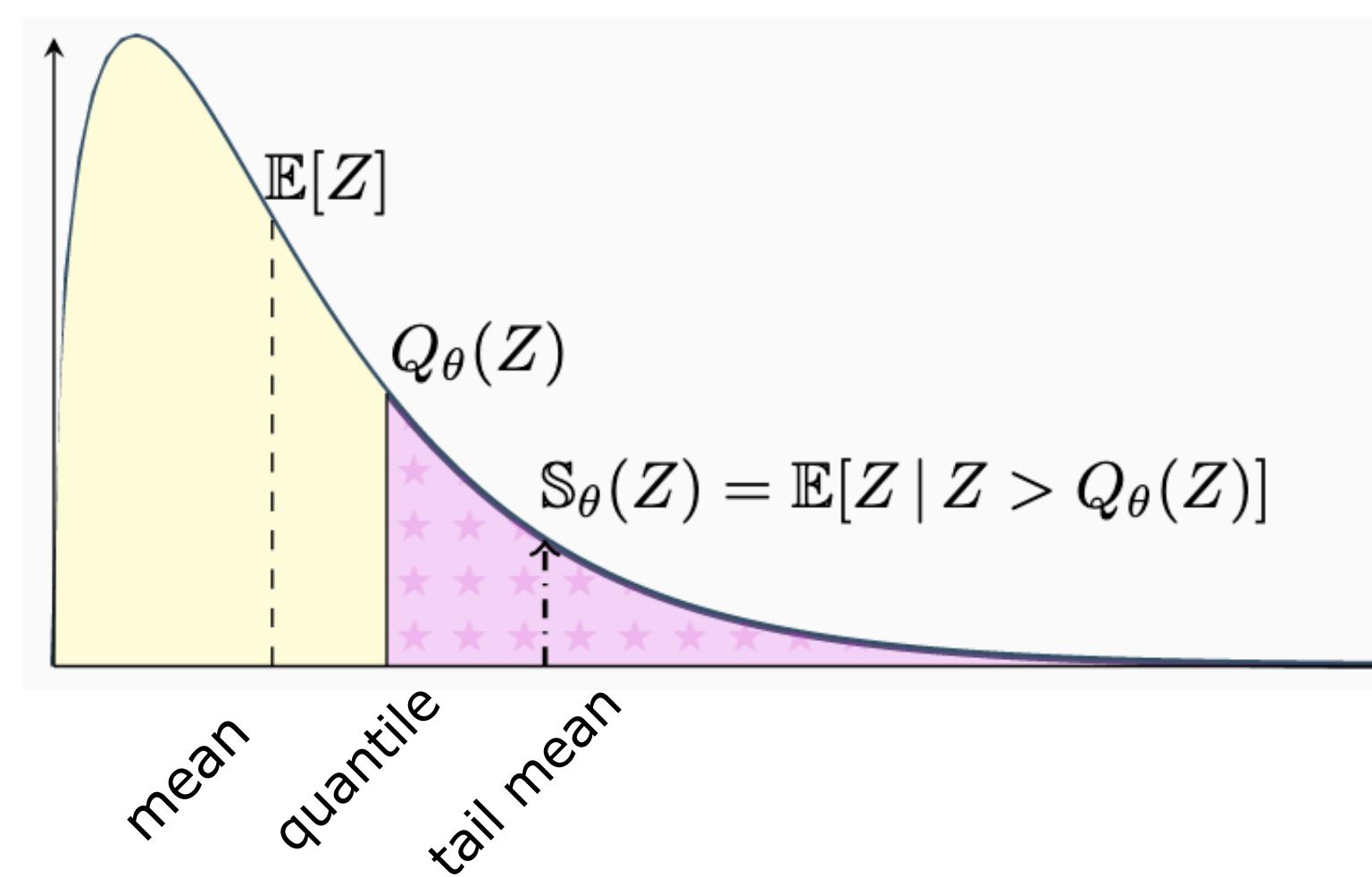


Part 1: Heterogeneity-aware objectives for federated learning

Heterogeneity \implies
large tail errors

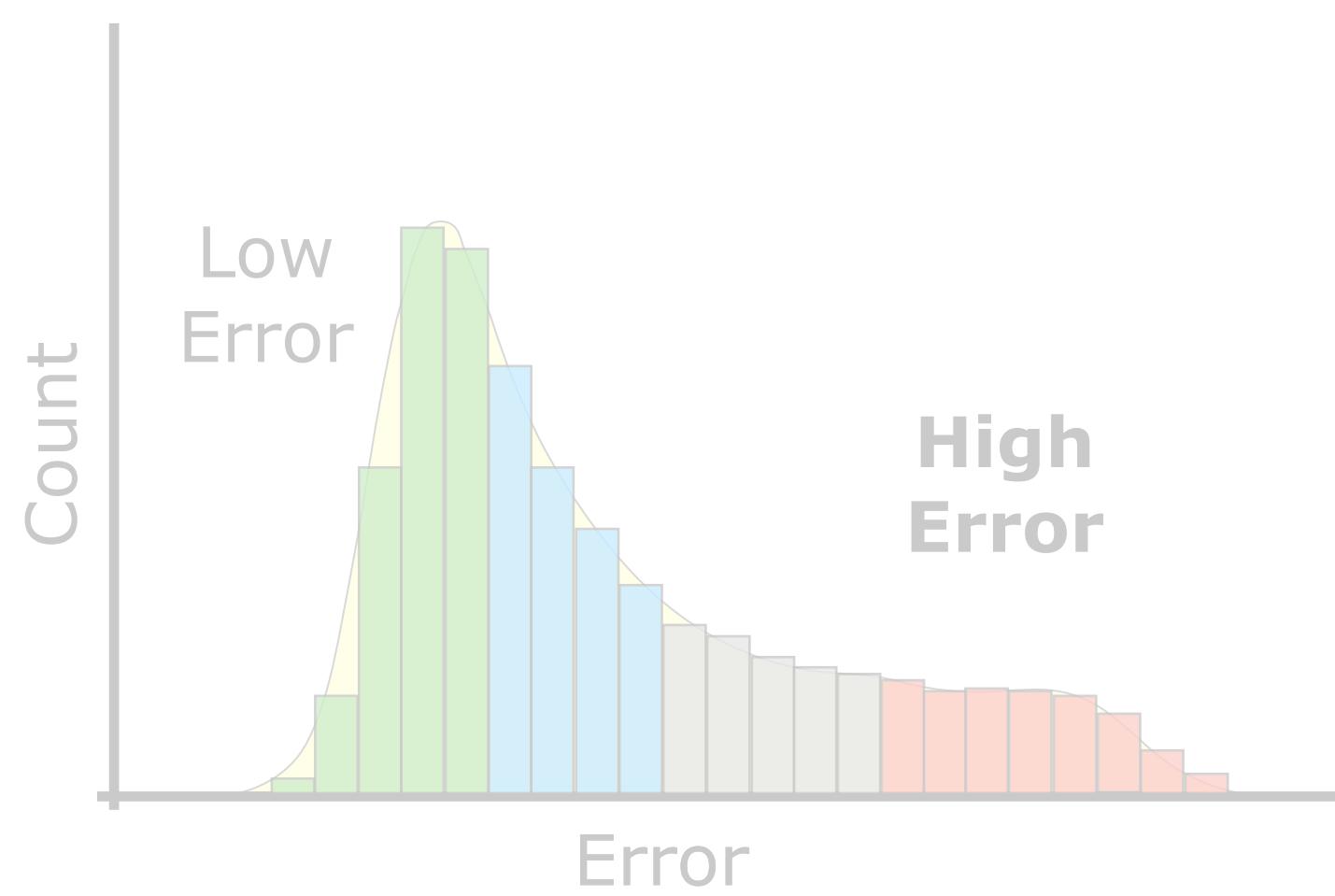


$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

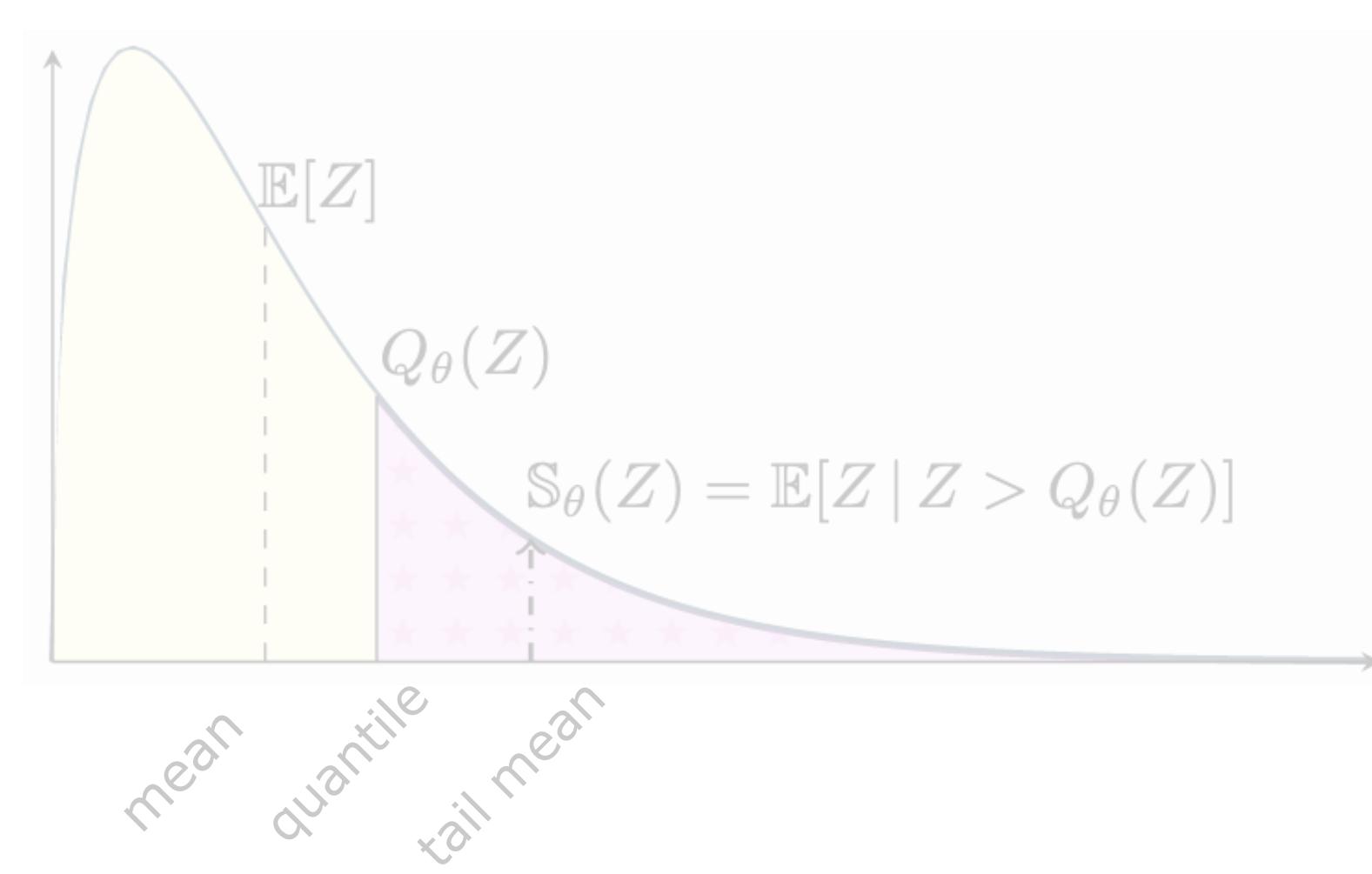


Part 1: Heterogeneity-aware objectives for federated learning

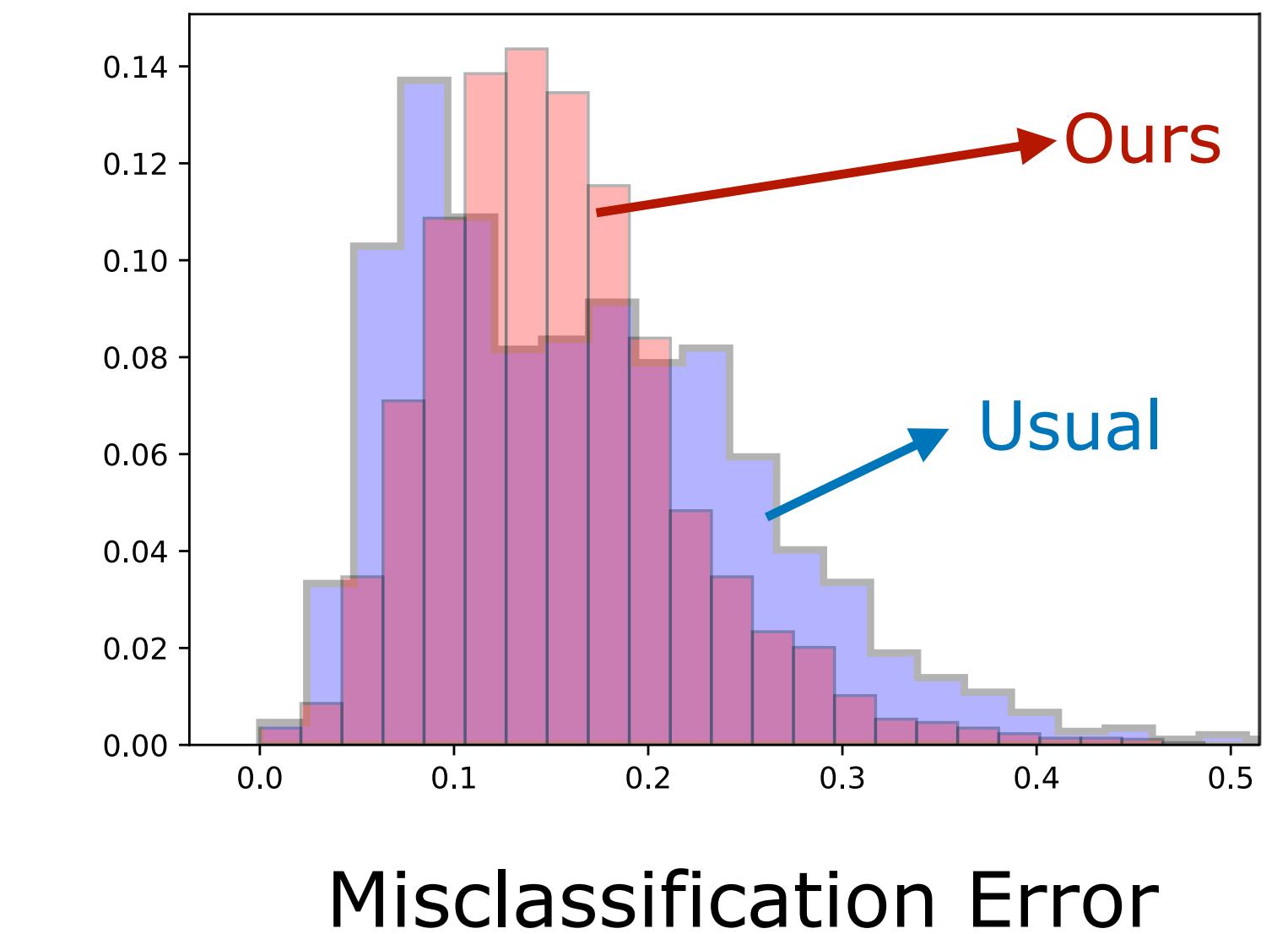
Heterogeneity \implies
large tail errors



$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

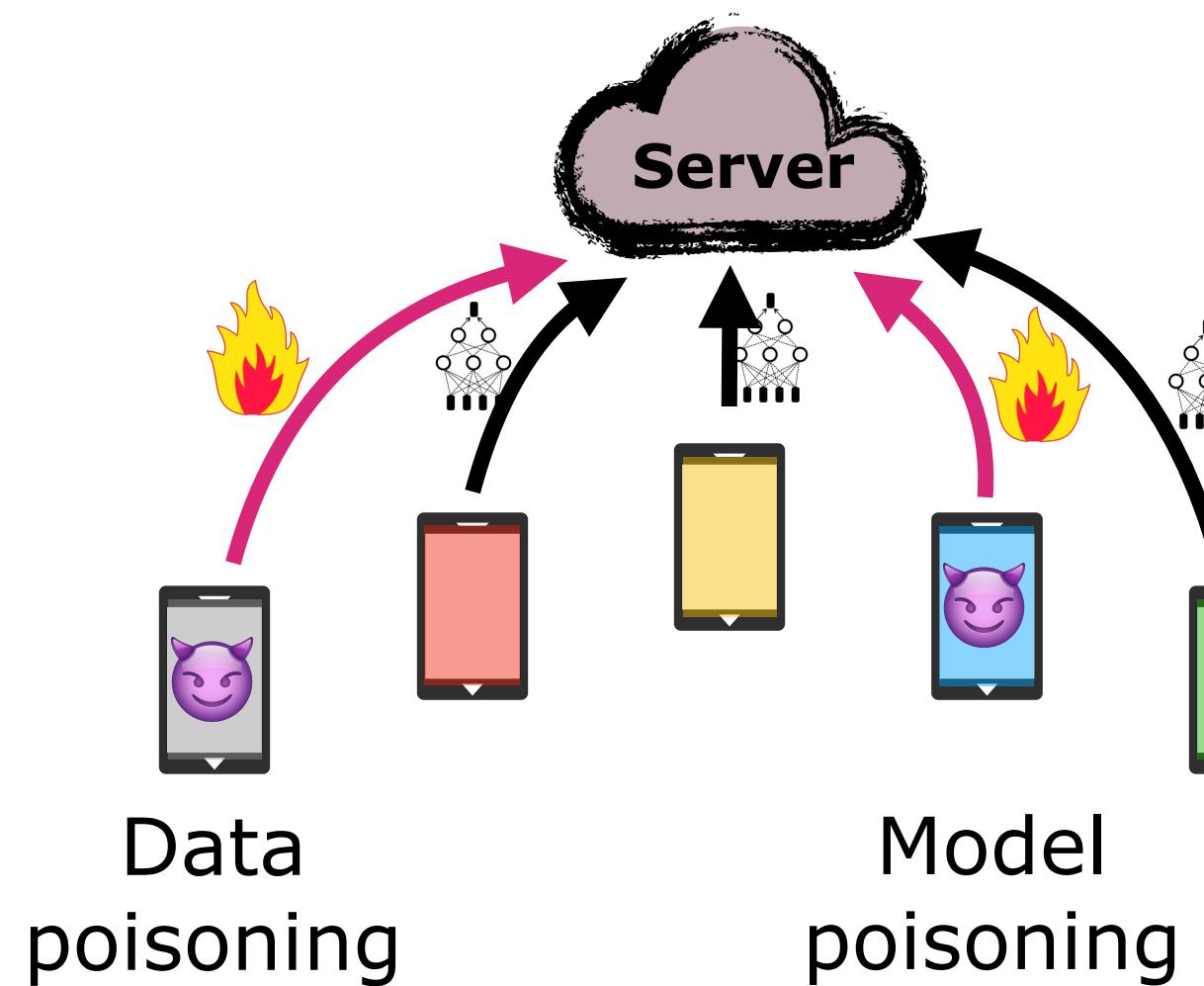


Our approach reduces
tail error



Part 2: Robust aggregation for federated learning

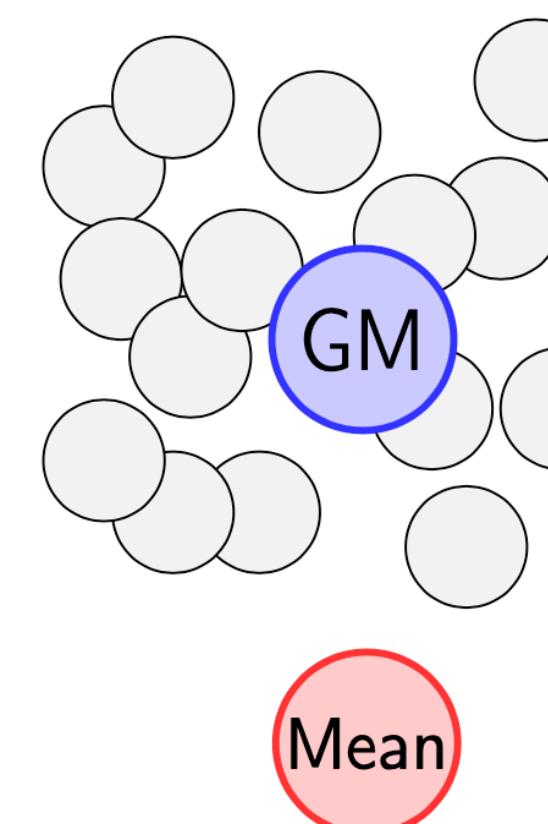
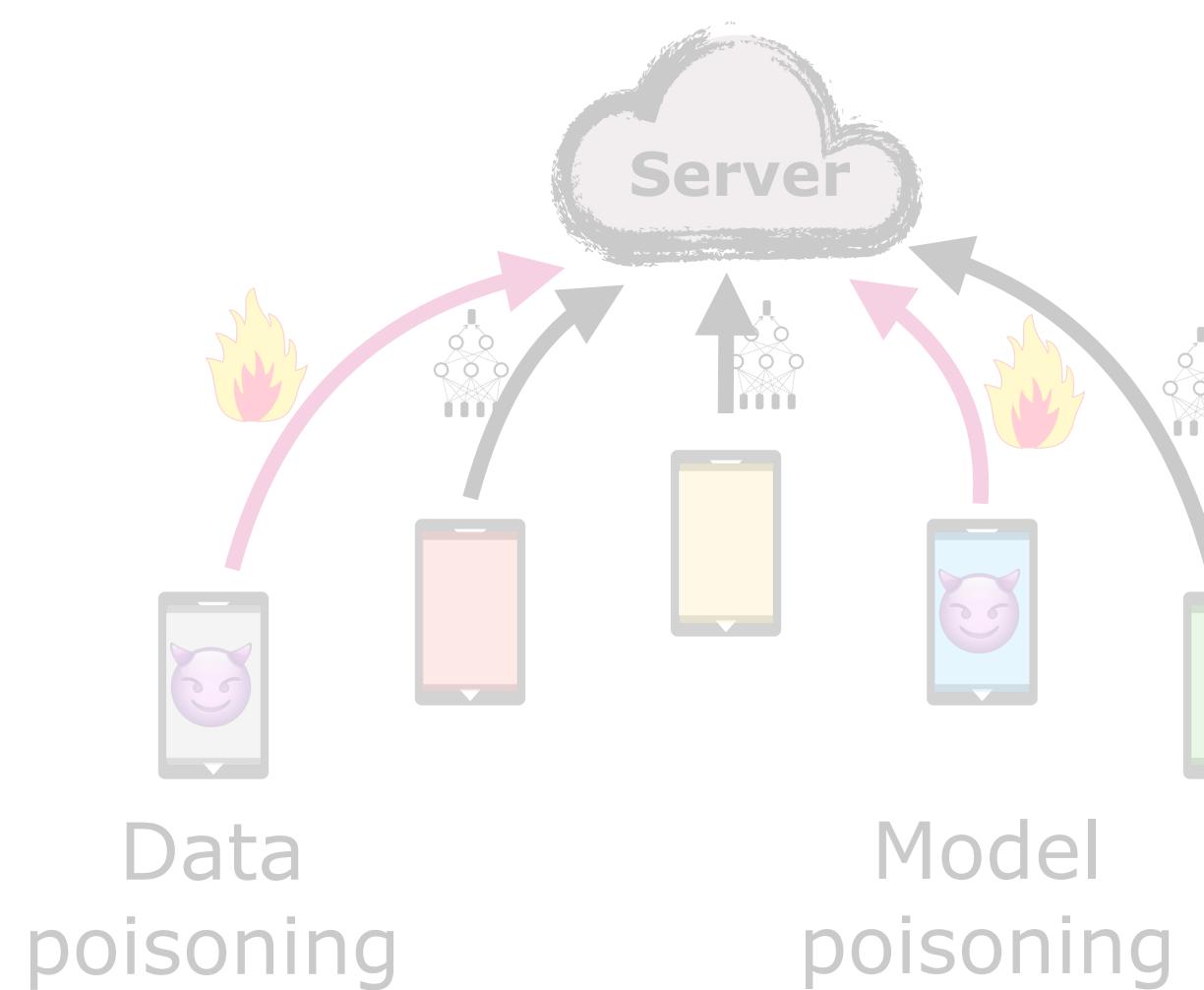
Arithmetic mean \Rightarrow
not robust to
poisoned updates



Part 2: Robust aggregation for federated learning

Arithmetic mean \Rightarrow
not robust to
poisoned updates

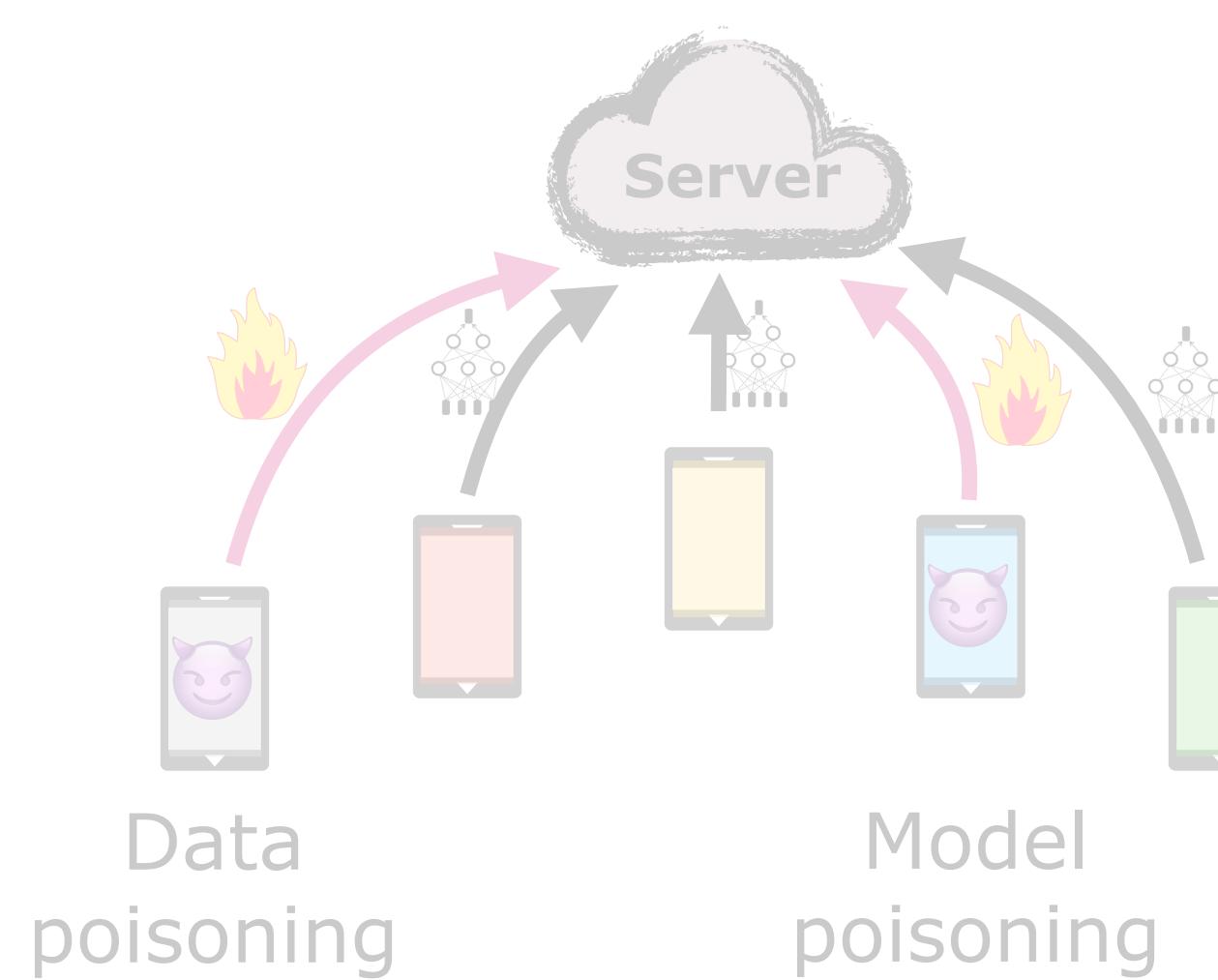
$$GM = \arg \min_z \sum_{i=1}^m \|z - w_i\|_2$$



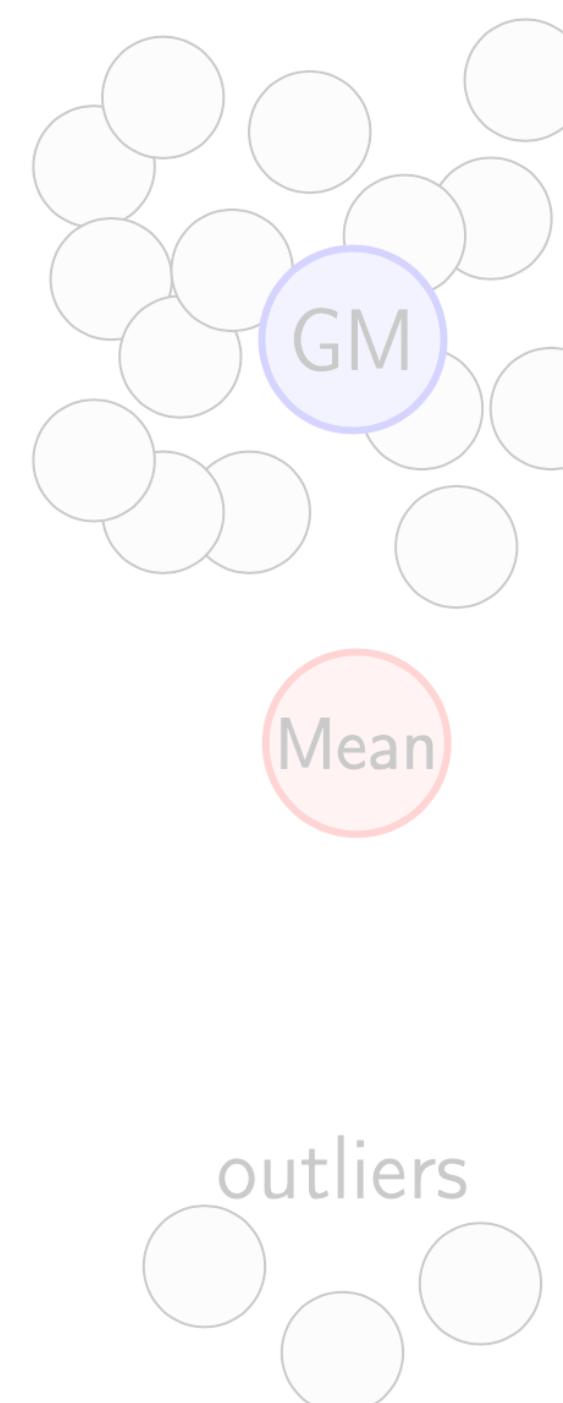
outliers

Part 2: Robust aggregation for federated learning

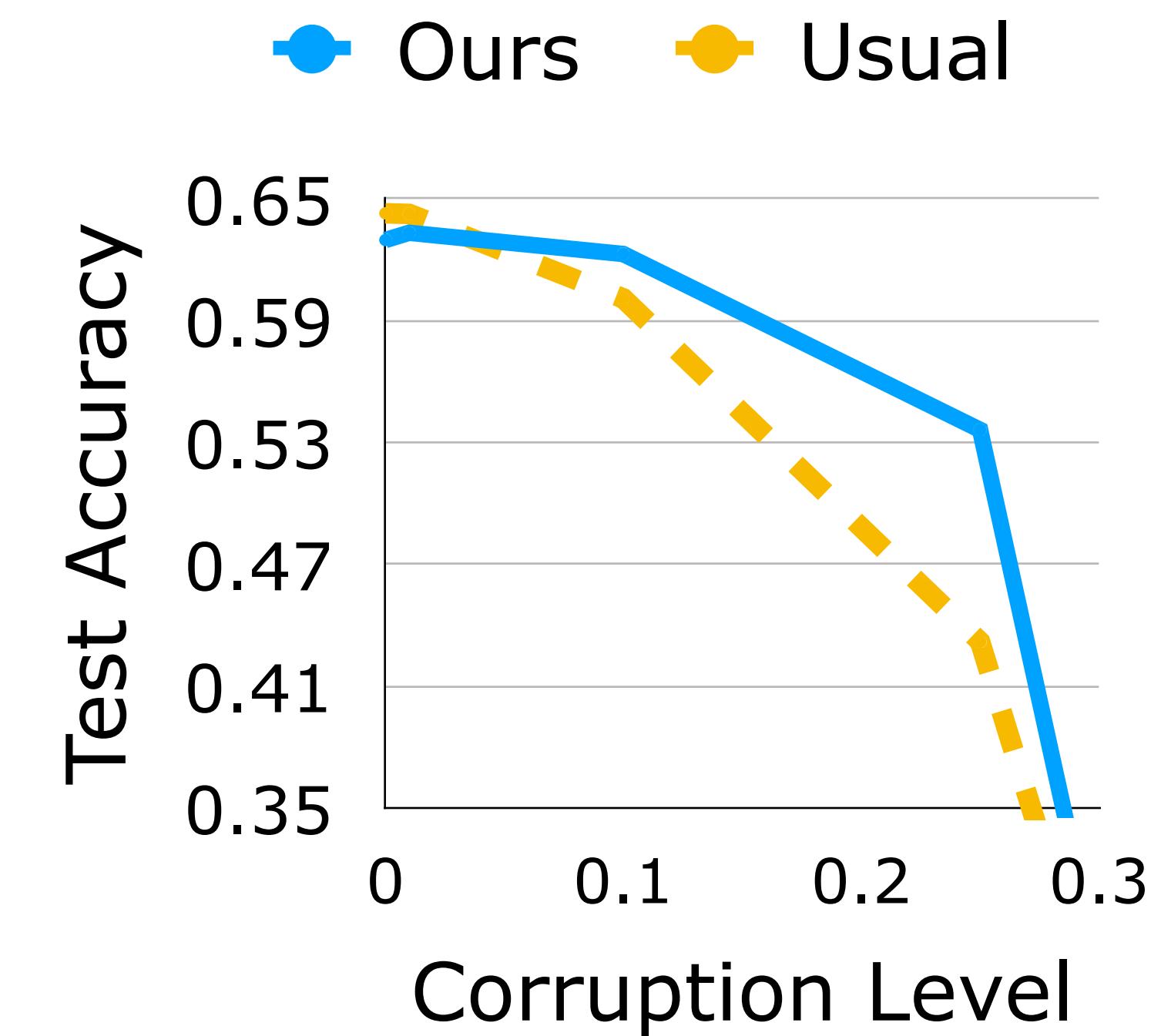
Arithmetic mean \Rightarrow
not robust to
poisoned updates



$$GM = \arg \min_z \sum_{i=1}^m \|z - w_i\|_2$$

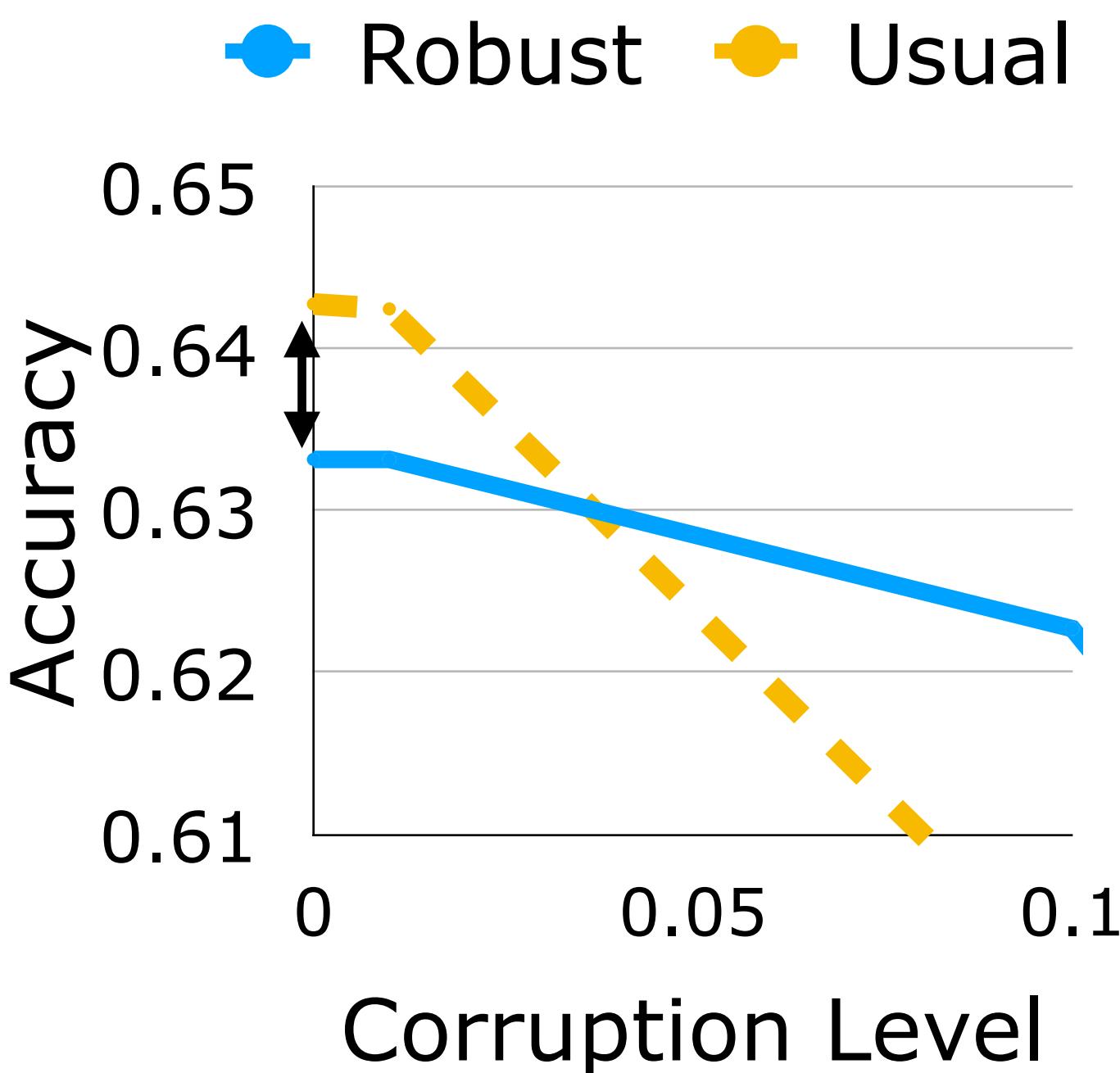


Our approach gives
greater robustness



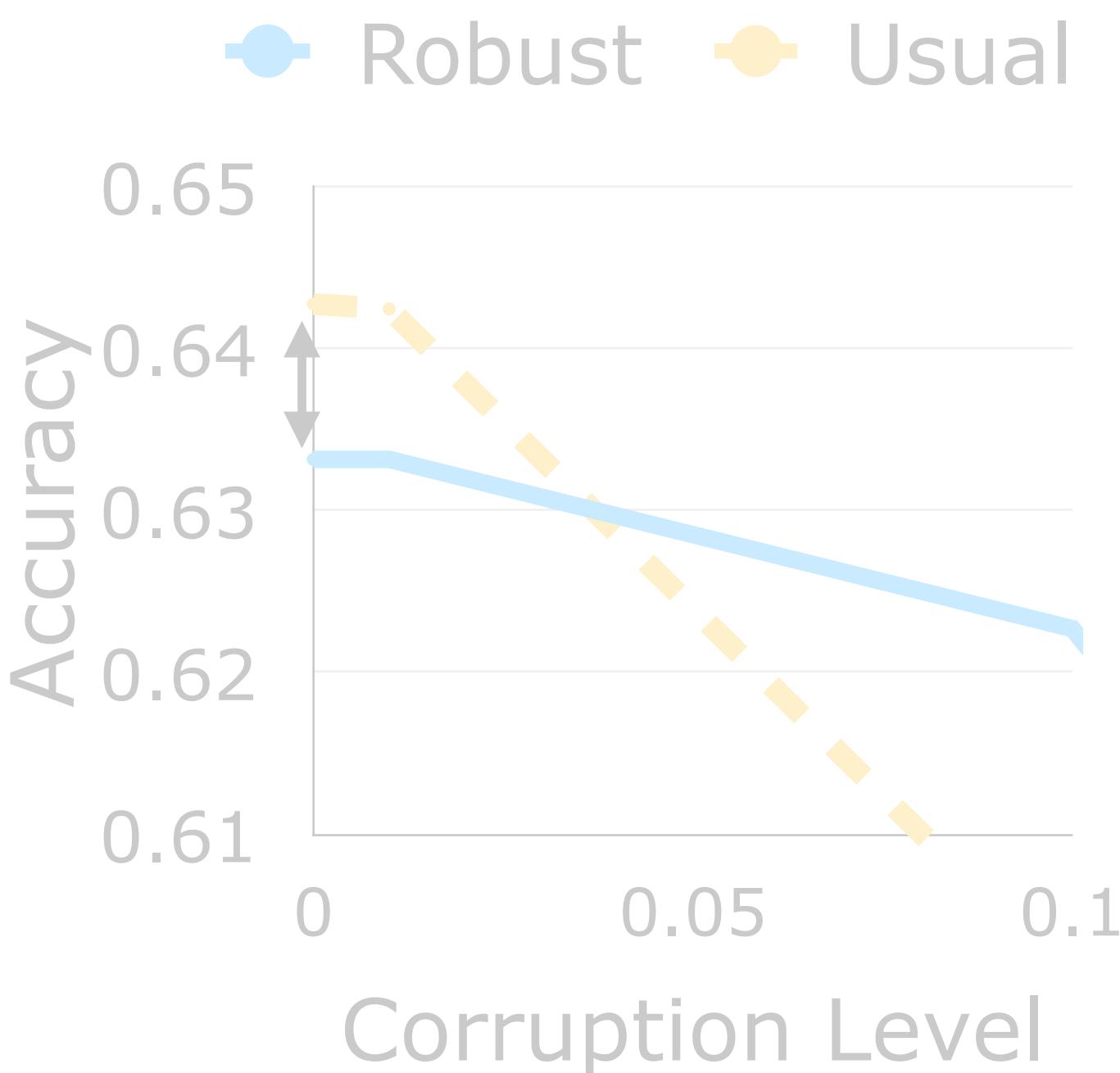
Part 3: Model personalization for federated learning

Heterogeneity &
robustness at odds

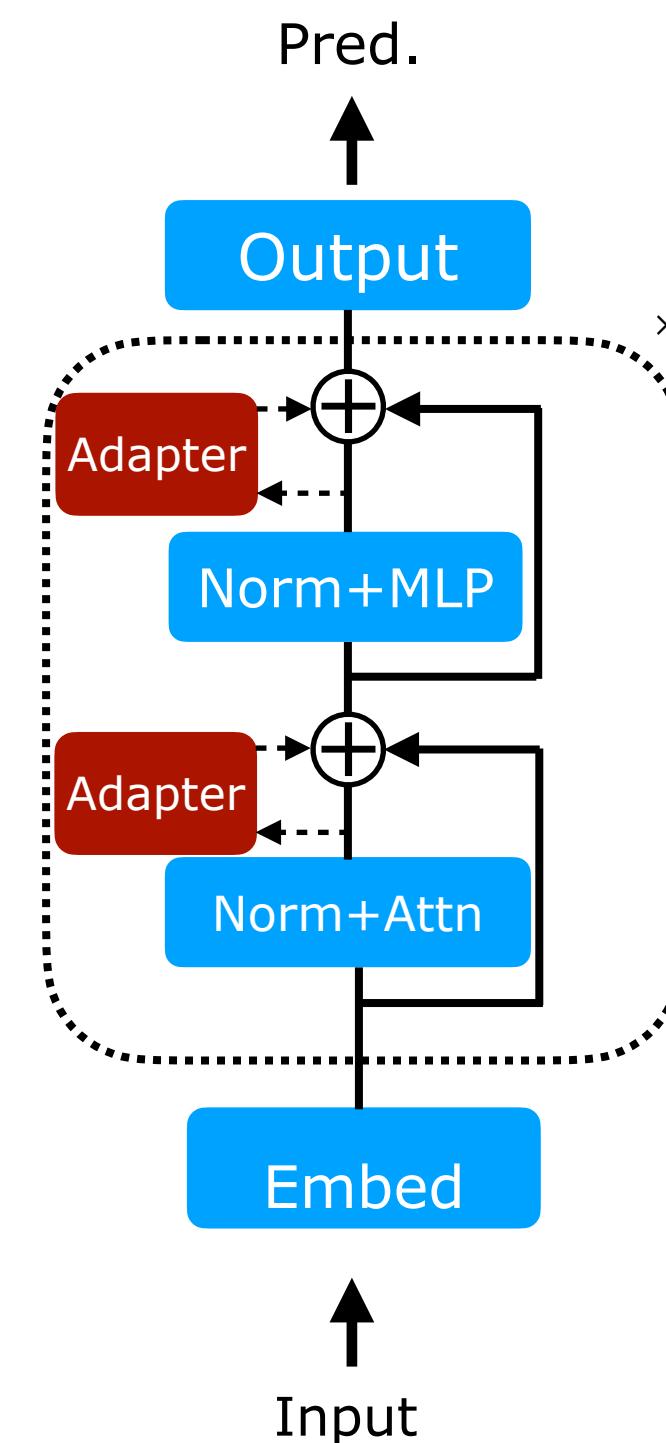


Part 3: Model personalization for federated learning

Heterogeneity &
robustness at odds

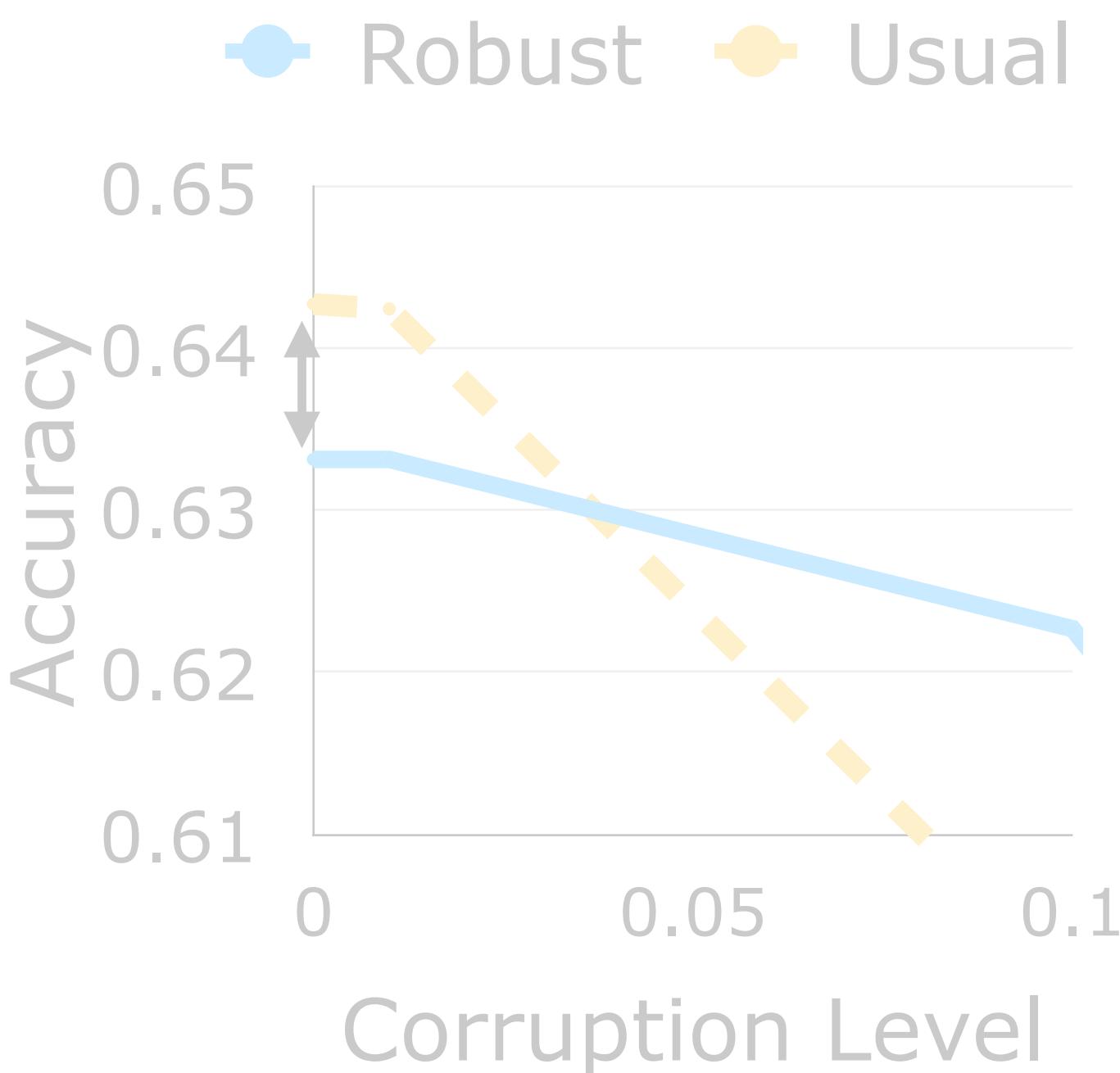


$$\min_{u, v_1, \dots, v_n} \frac{1}{n} \sum_{i=1}^n F_i(u, v_i)$$

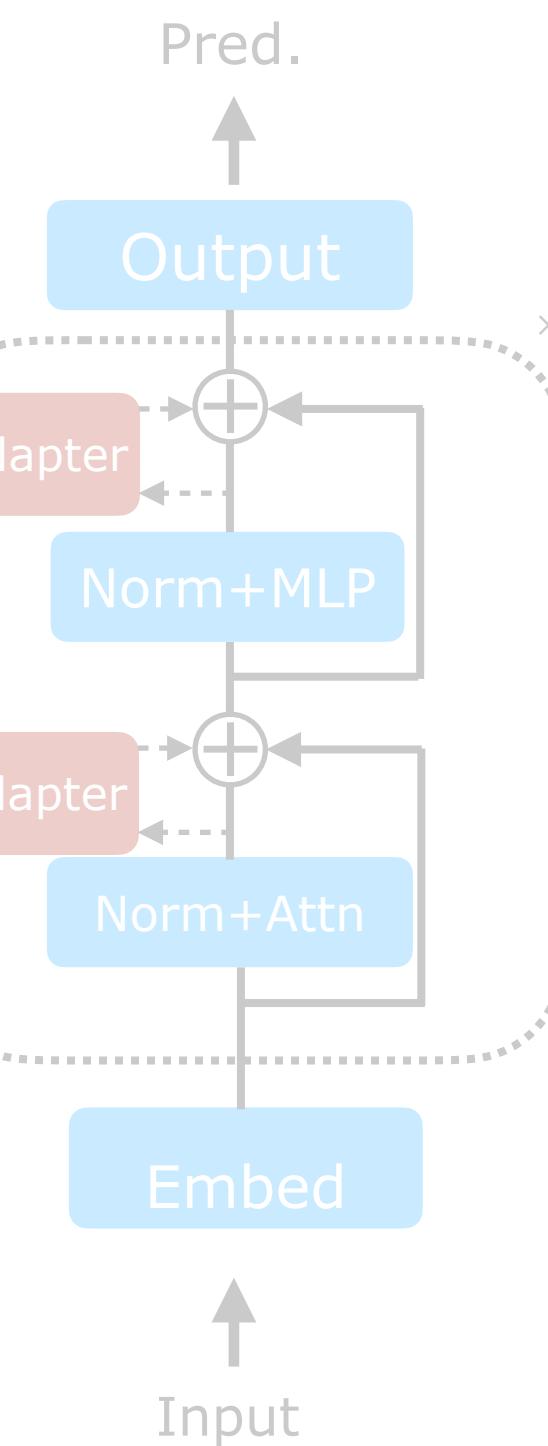


Part 3: Model personalization for federated learning

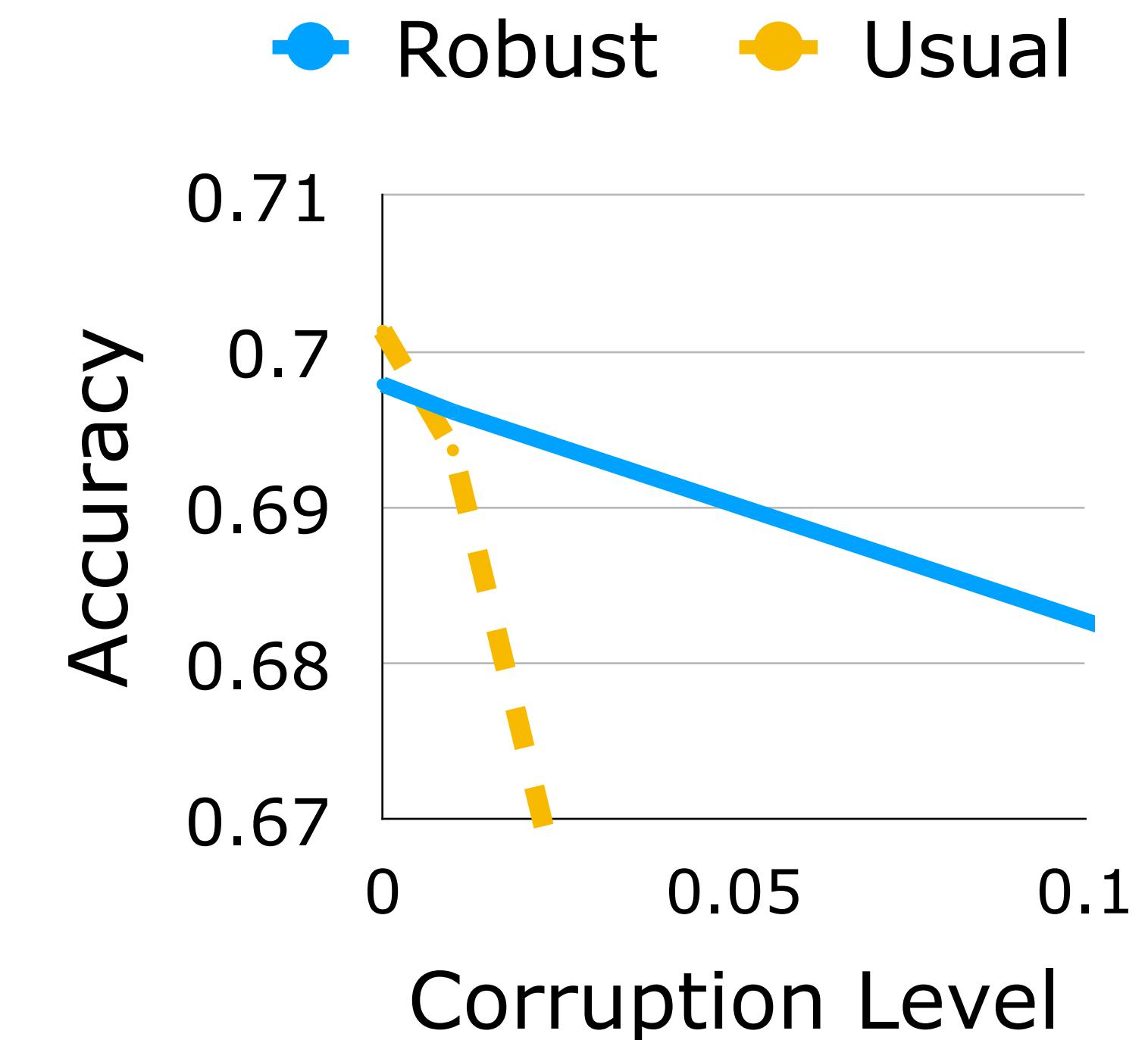
Heterogeneity & robustness at odds



$$\min_{u, v_1, \dots, v_n} \frac{1}{n} \sum_{i=1}^n F_i(u, v_i)$$



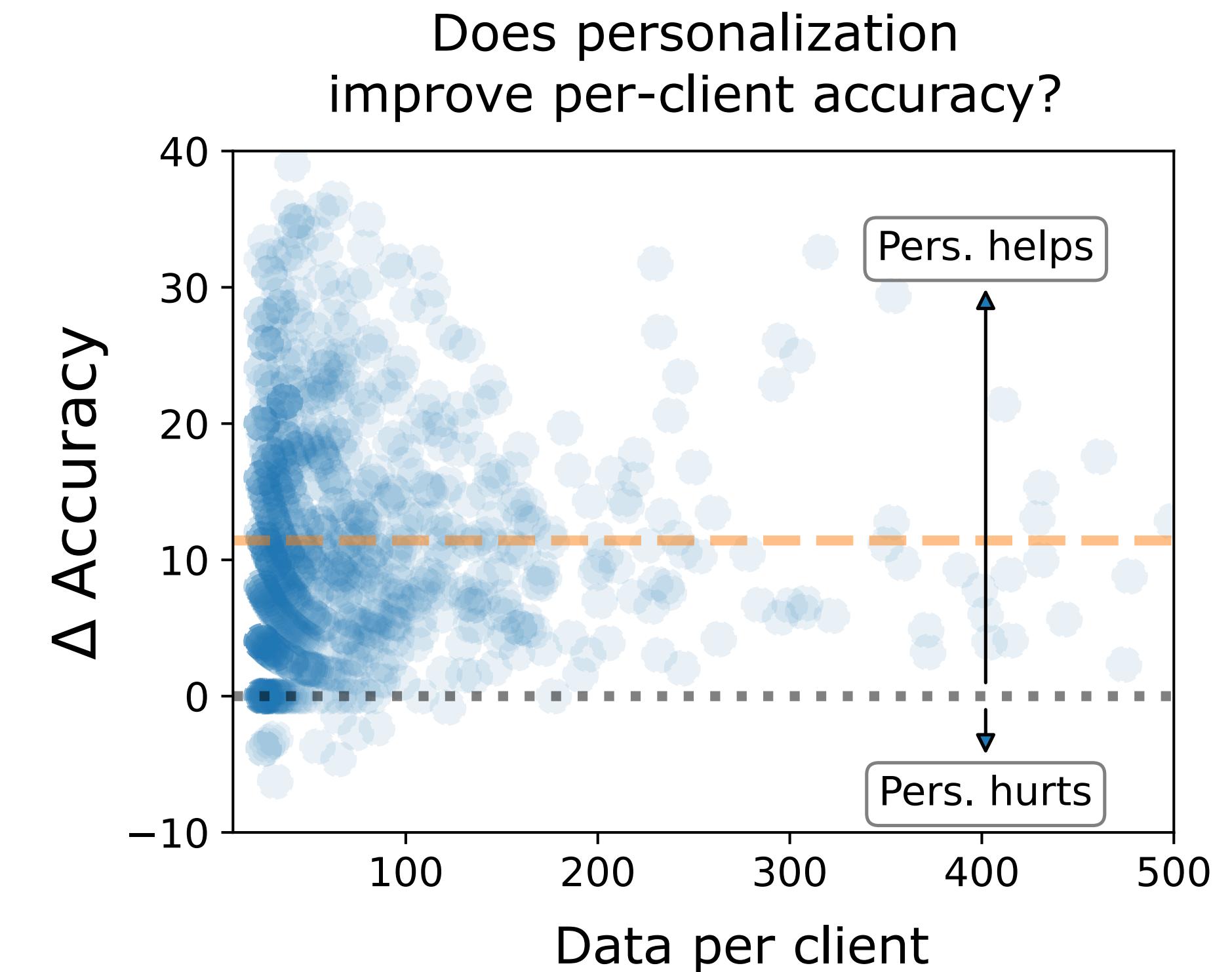
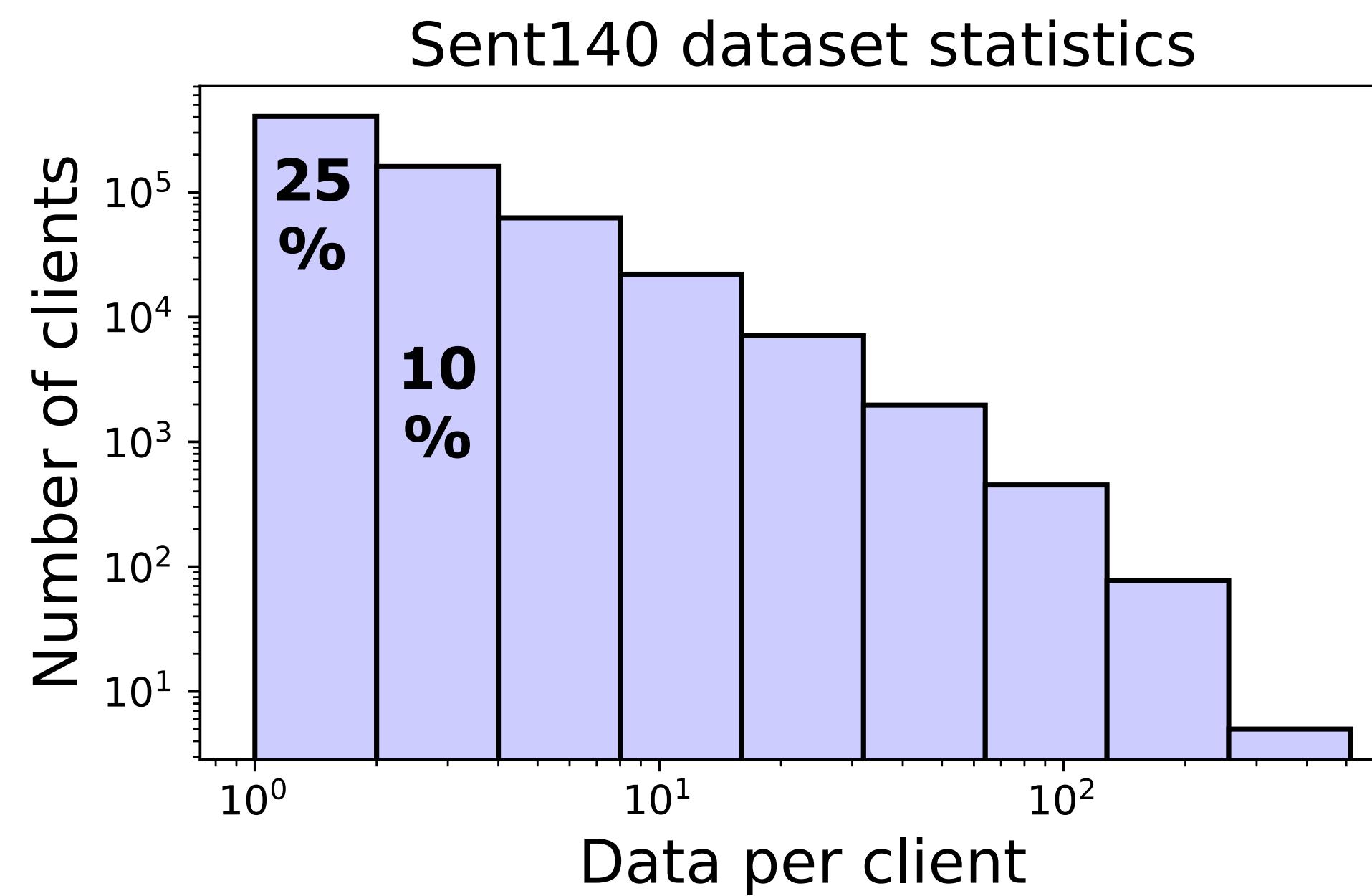
Can tailor to heterogeneity & retain robustness



Perspectives and conclusion

The small data problem

- Per-client evaluations are not reliable
- Personalization can overfit



Understanding heterogeneity

Many negative results: optimization can slow down, makes robustness harder, ...

Yet, federated learning works.

Understanding heterogeneity

Many negative results: optimization can slow down, makes robustness harder, ...

Yet, federated learning works.

Quantify heterogeneity:

Measure gaps between
distributions: **MAUVE**

[**P.**, Swayamdipta, Zellers, Thickstun, Welleck, Choi, Harchaoui. NeurIPS (2021),
Liu, **P.**, Welleck, Oh, Choi, Harchaoui. NeurIPS (2021)]

Understanding heterogeneity

Many negative results: optimization can slow down, makes robustness harder, ...

Yet, federated learning works.

Quantify heterogeneity:

Measure gaps between distributions: **MAUVE**

[P., Swayamdipta, Zellers, Thickstun, Welleck, Choi, Harchaoui. NeurIPS (2021),
Liu, P., Welleck, Oh, Choi, Harchaoui. NeurIPS (2021)]

Statistical assumptions under which heterogeneity is benign?

What measures of heterogeneity impact optimization?

Tension between heterogeneity and privacy

Thank you!



J.P.Morgan

