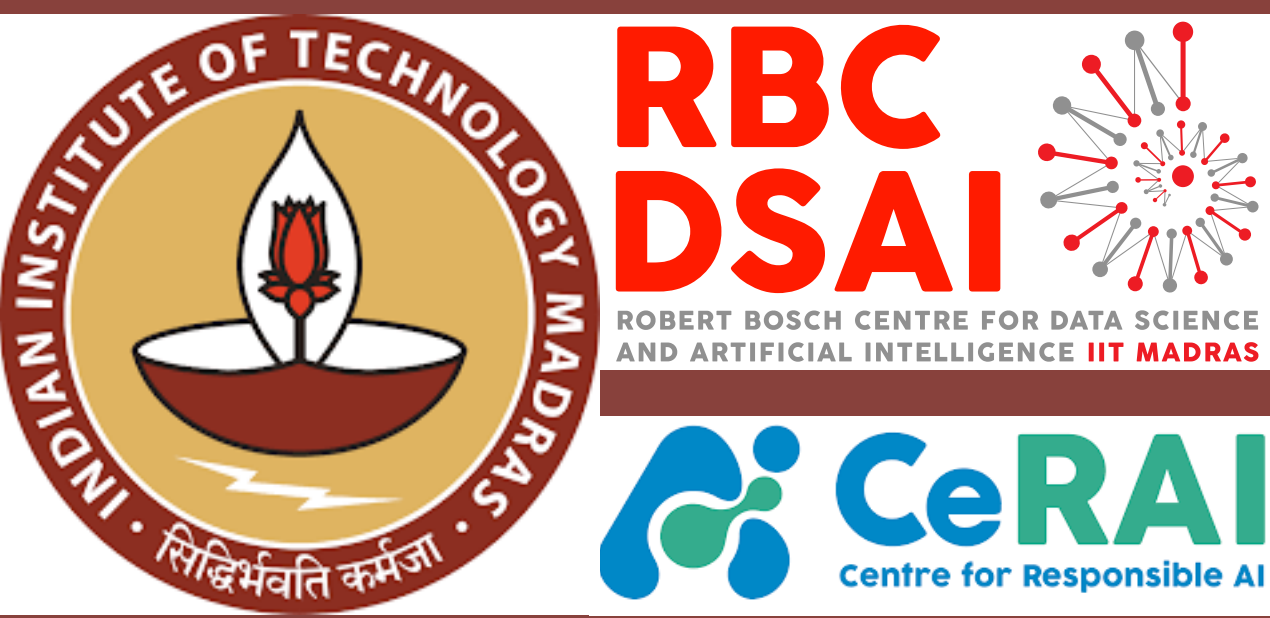# Robust Aggregation for Federated Learning

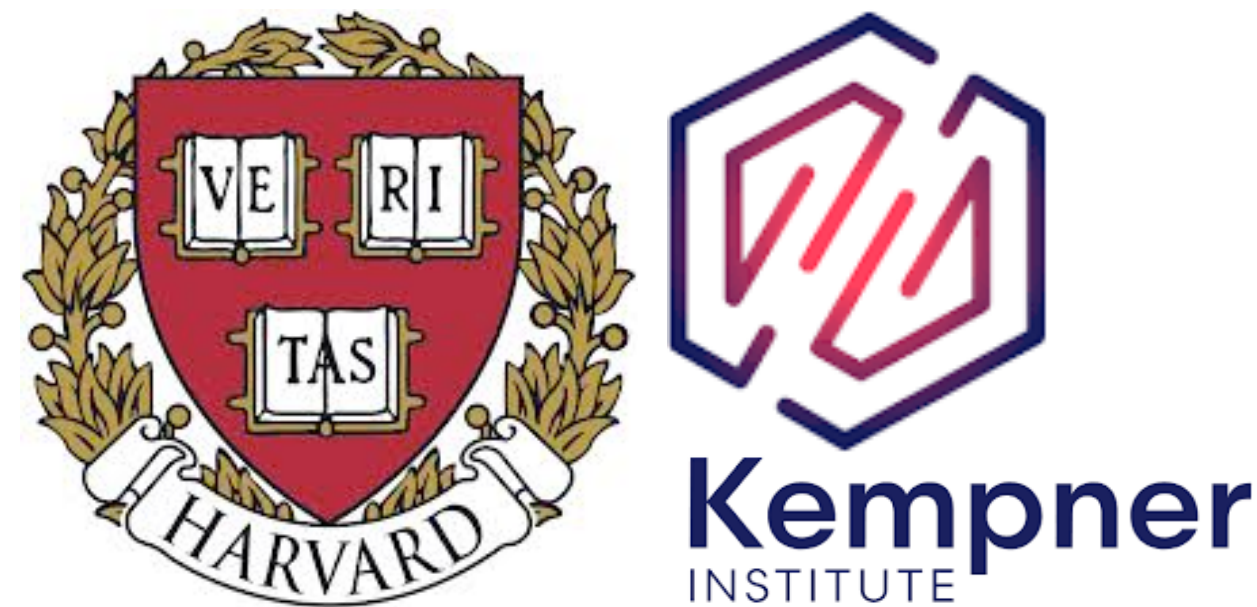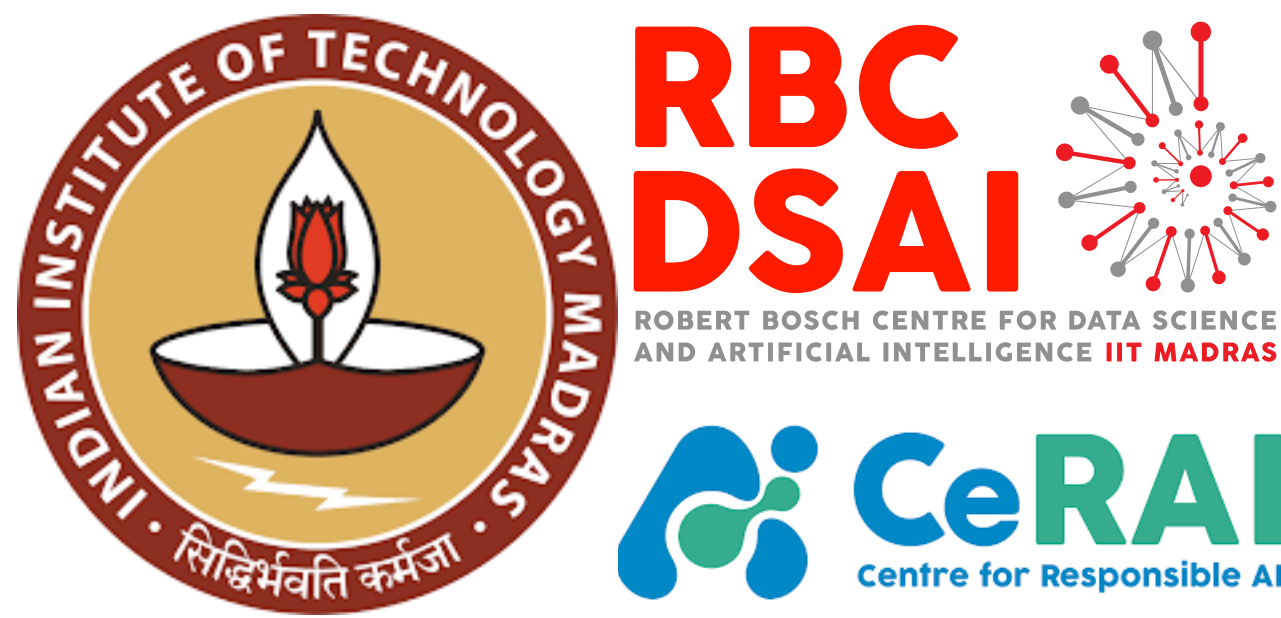IEEE Transactions on Signal Processing (2022)

**Krishna Pillutla**
IIT Madras

# Team

Krishna Pillutla      Sham Kakade      Zaid Harchaoui
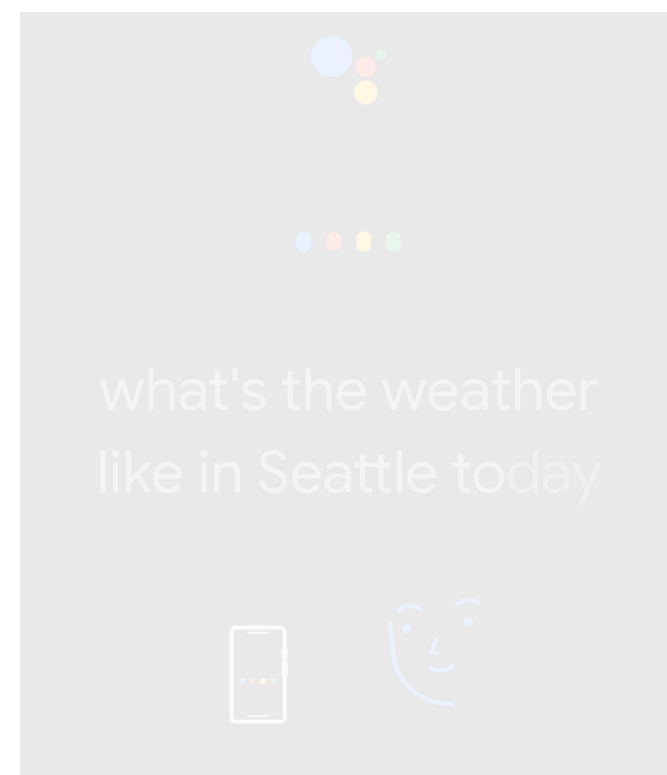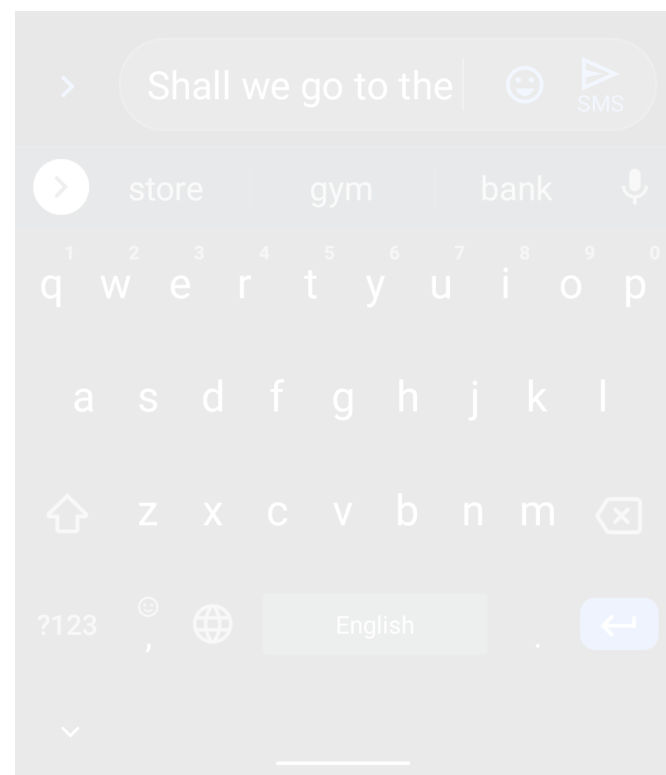
what's the weather like in Seattle today



Image Credit: Robotics Business Review



Rieke et al. NPJ Digit. Med. (2020)    Image Credit: Wellcome

# Data is *decentralized* and *private*

Image Credit: Robotics Business Review
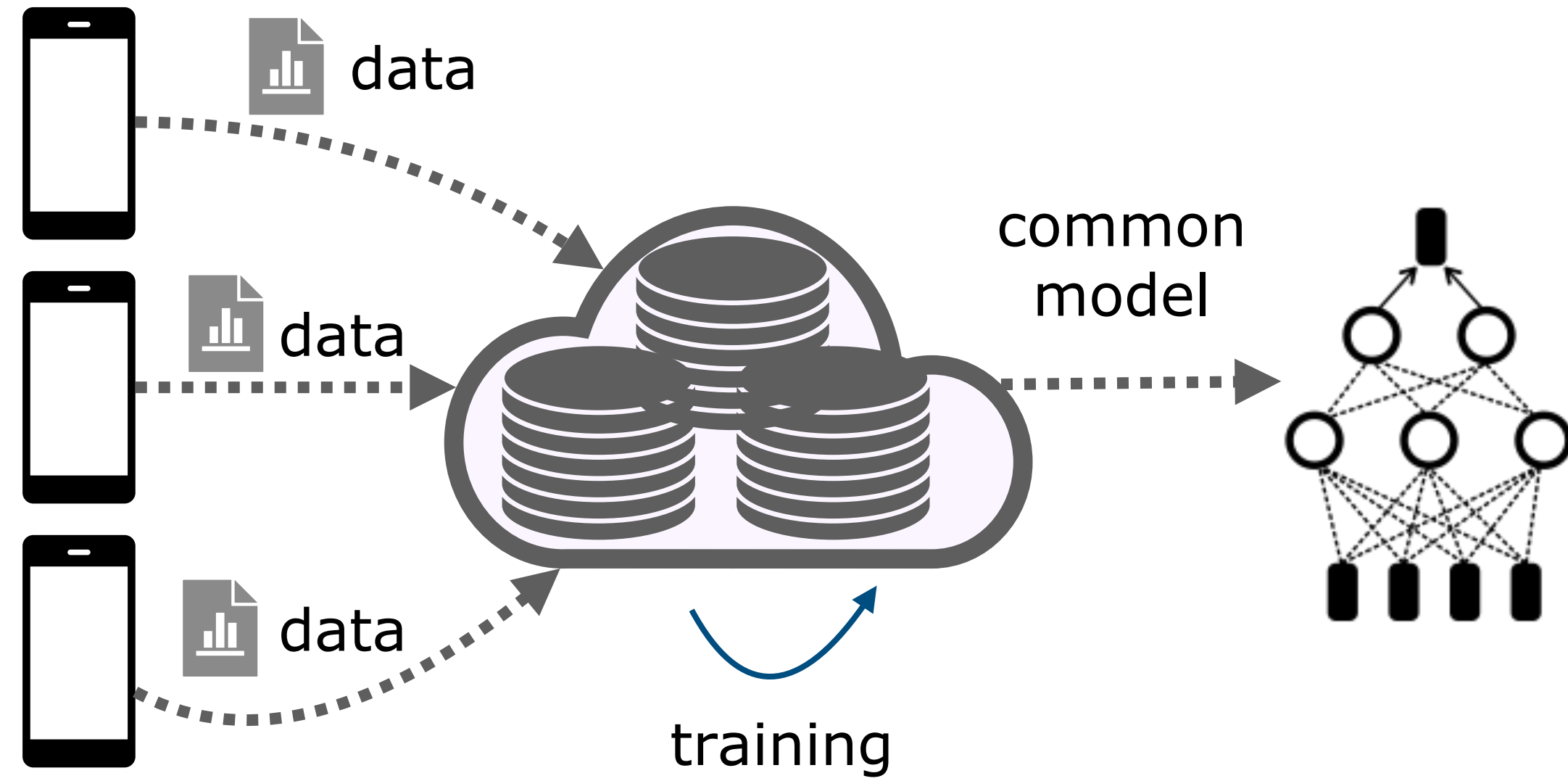
Rieke et al. NPJ Digit. Med. (2020)      Image Credit: Wellcome

4

# Datacenter



data

data

data

common
model

training

# Datacenter



data

data

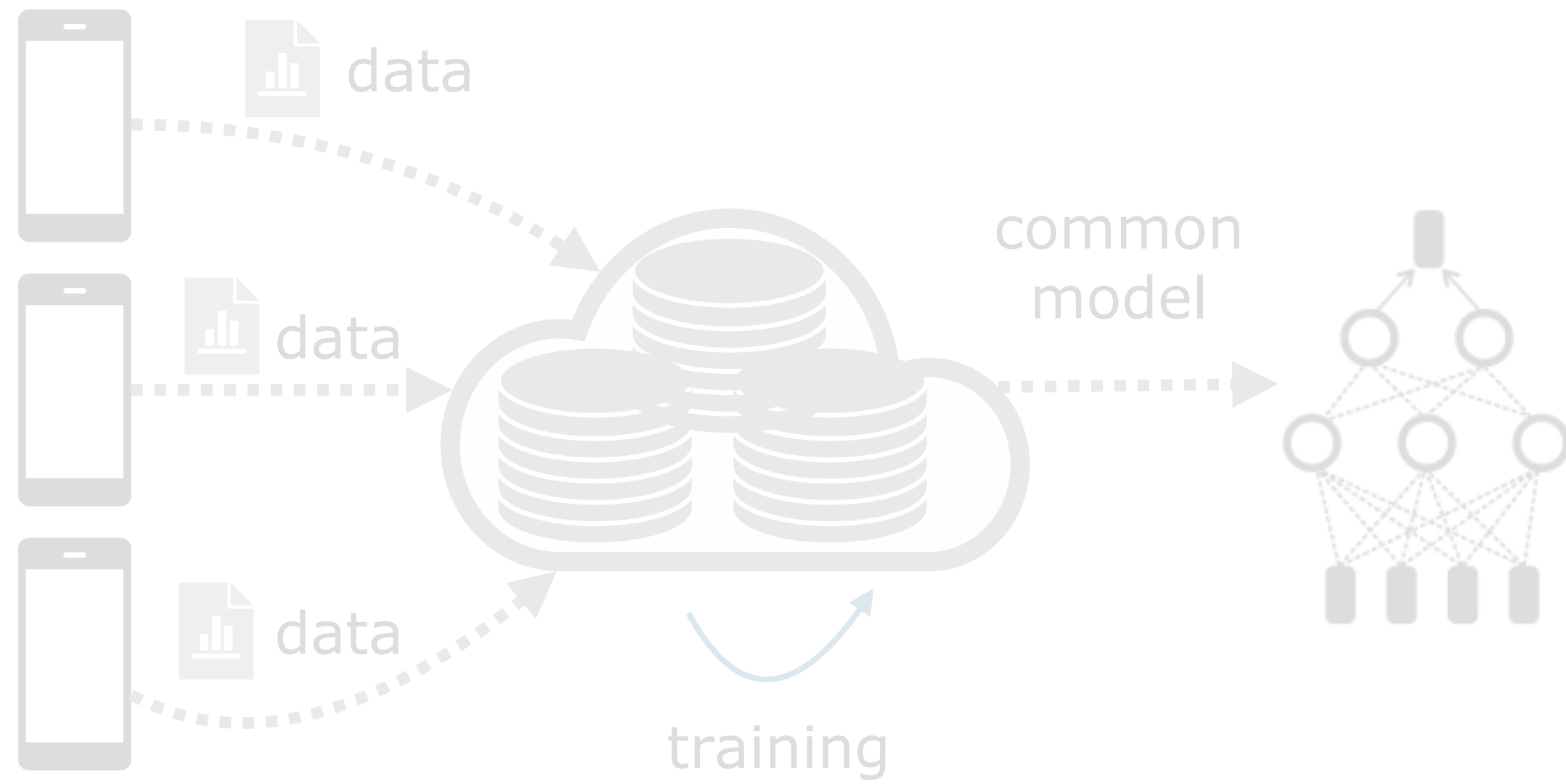data

common
model

training

# Non-collaborative



data

model 1

training

⋮

data

model $n$

training

# Datacenter

# Non-collaborative

data

data

common
model

data

training

data

model 1

data

model $n$

training

training

# Peer-to-peer

data

training

# Advances and Open Problems in Federated Learning

**Peter Kairouz**
Google Research
Kairouz@google.com

**H. Brendan McMahan**
Google Research

*et al.*

Percentage of world population with a smartphone (Data: Business Wire)

# Federated Learning



Percentage of world population with a smartphone (Data: Business Wire)

# Federated Learning



Percentage of world population with a smartphone (Data: Business Wire)

# Federated Learning



Percentage of world population with a smartphone (Data: Business Wire)

# Federated Learning



## Percentage of world population with a smartphone (Data: Business Wire)

# Federated Learning



Percentage of world population with a smartphone (Data: Business Wire)
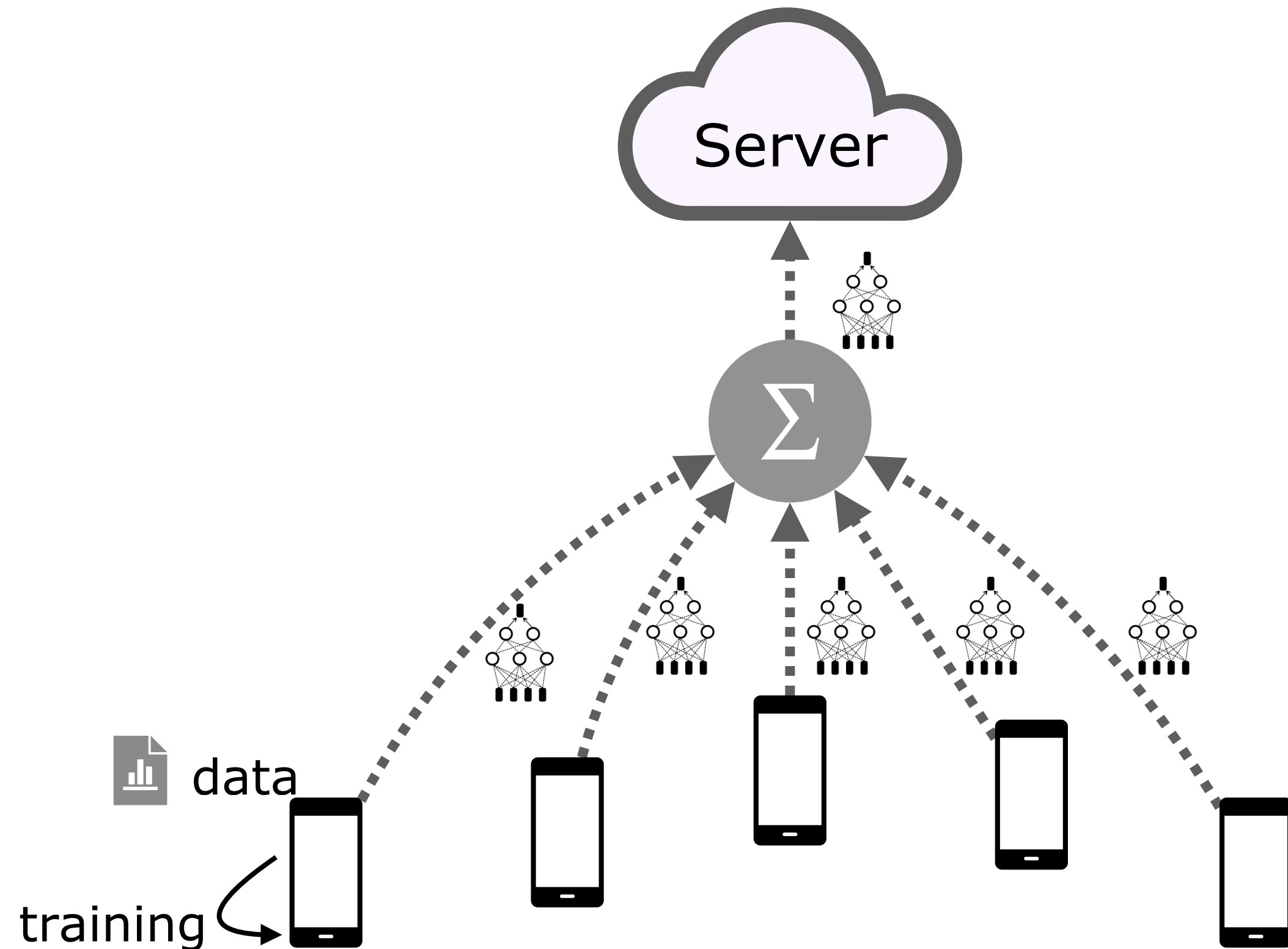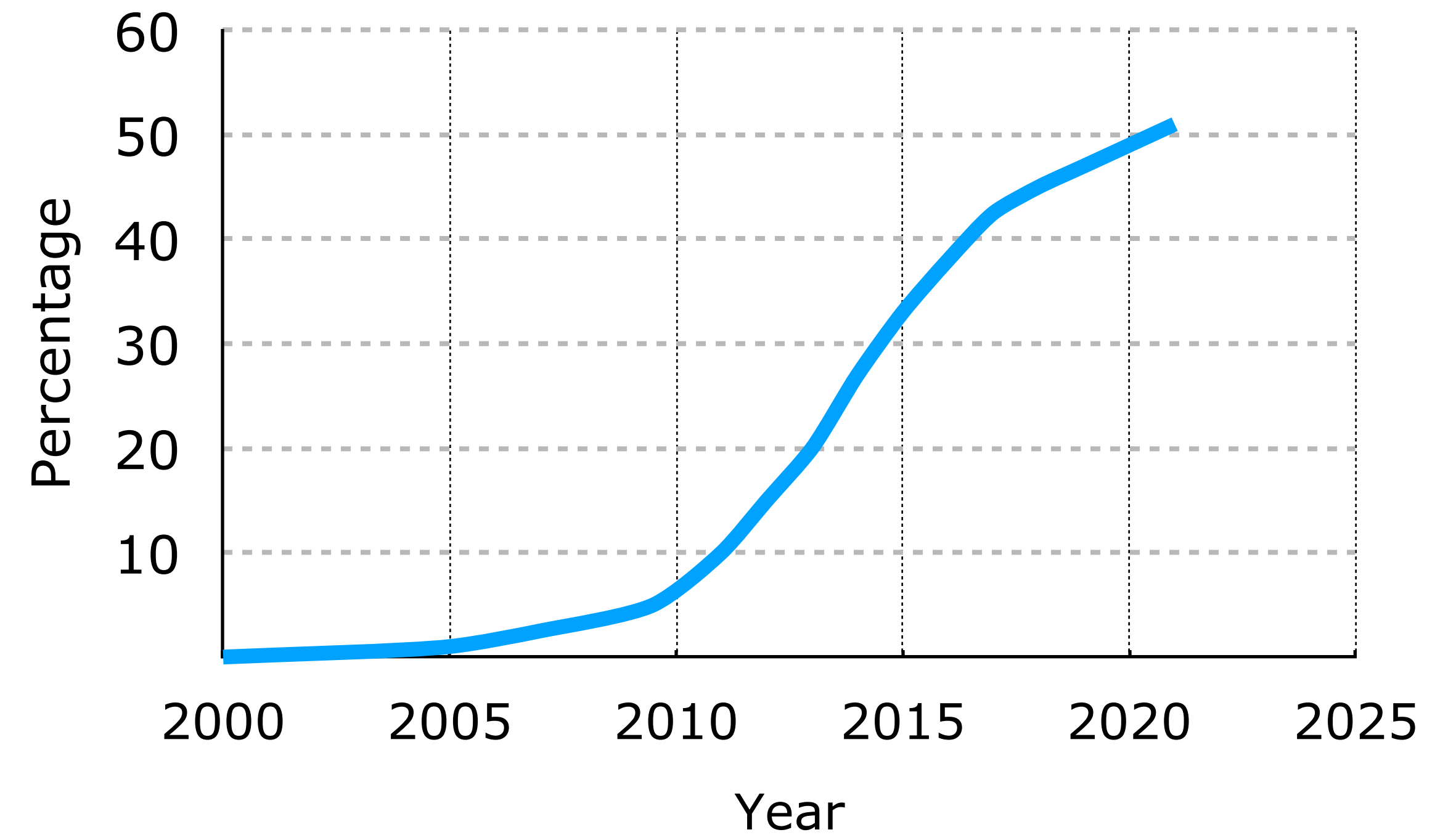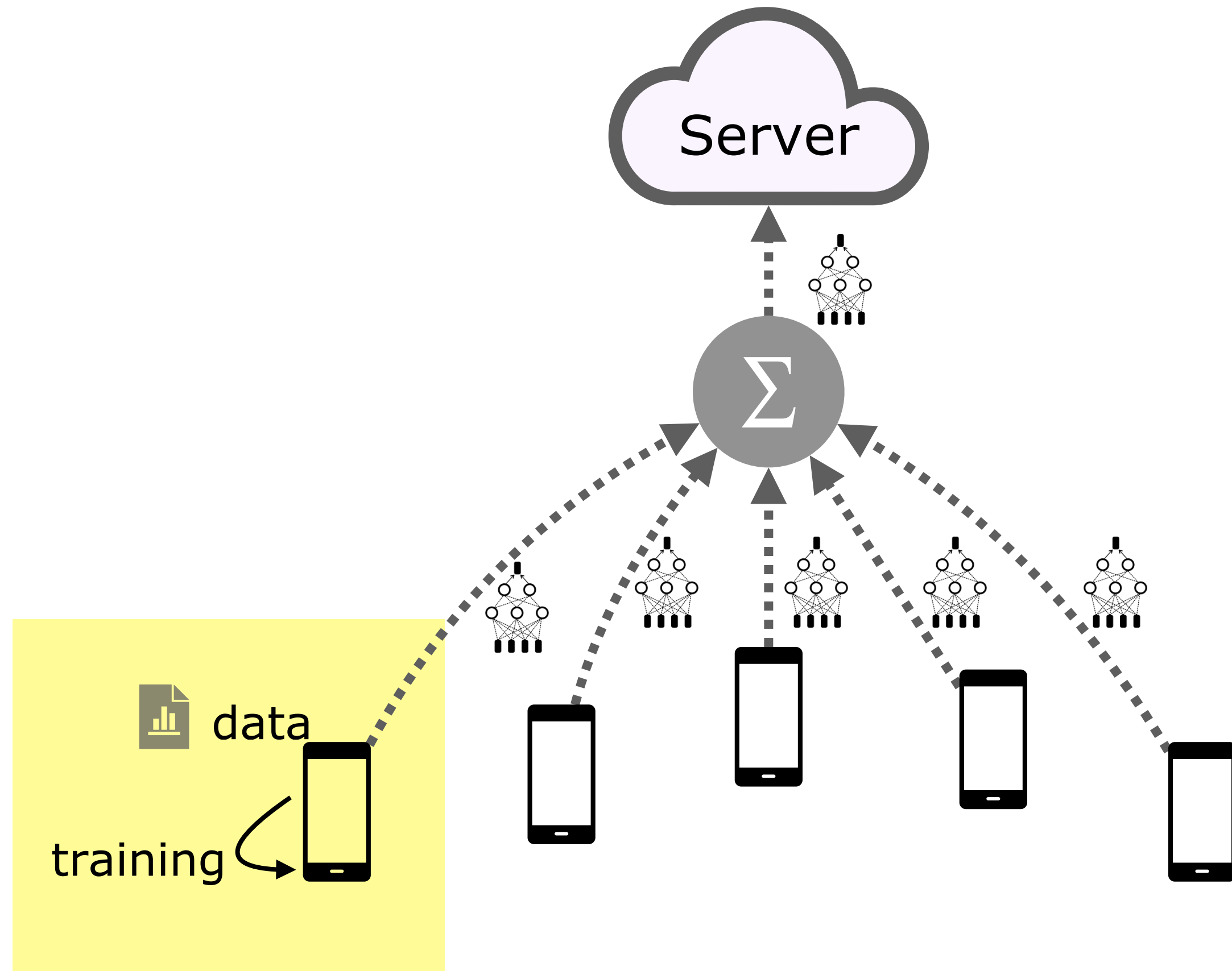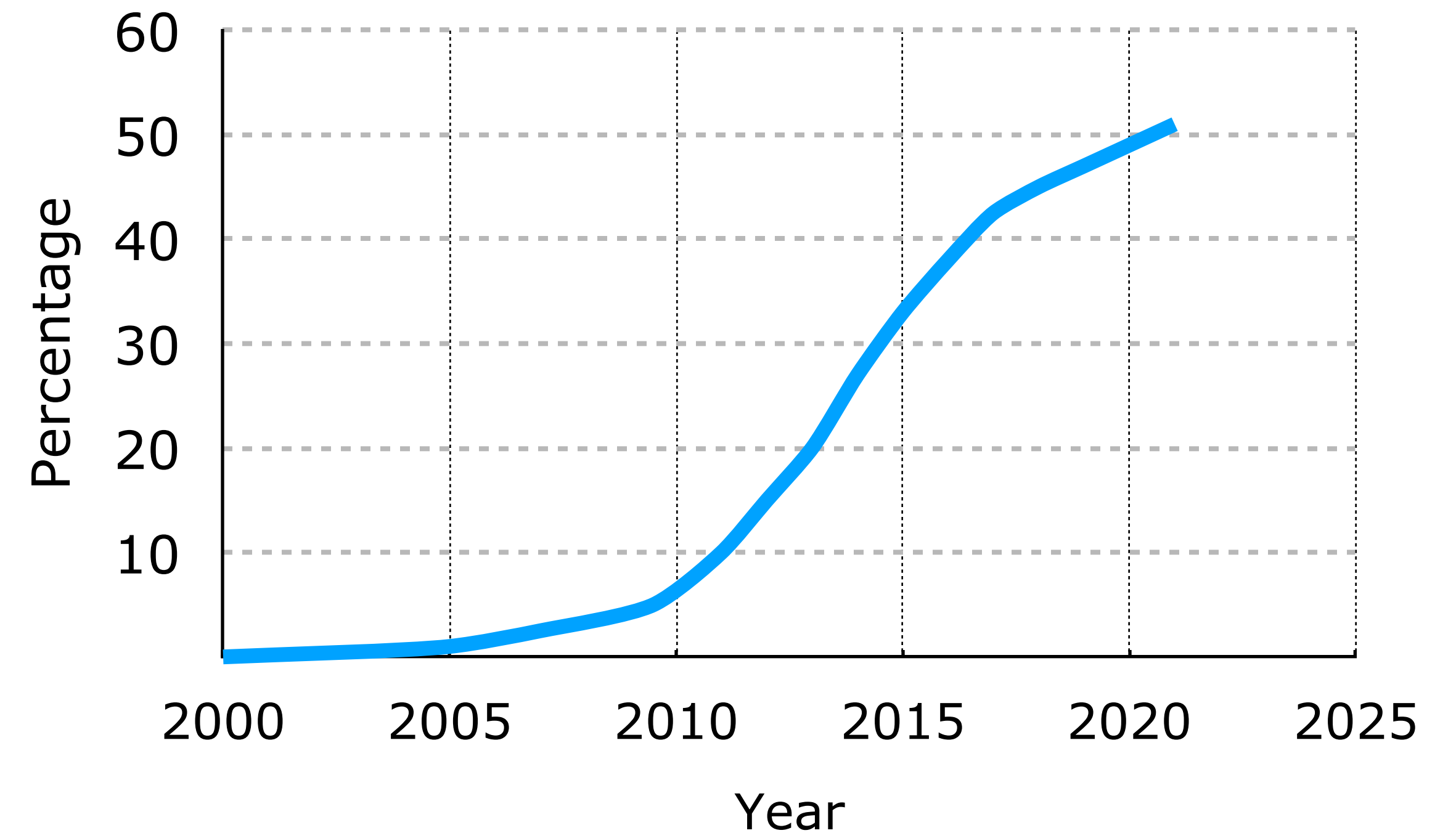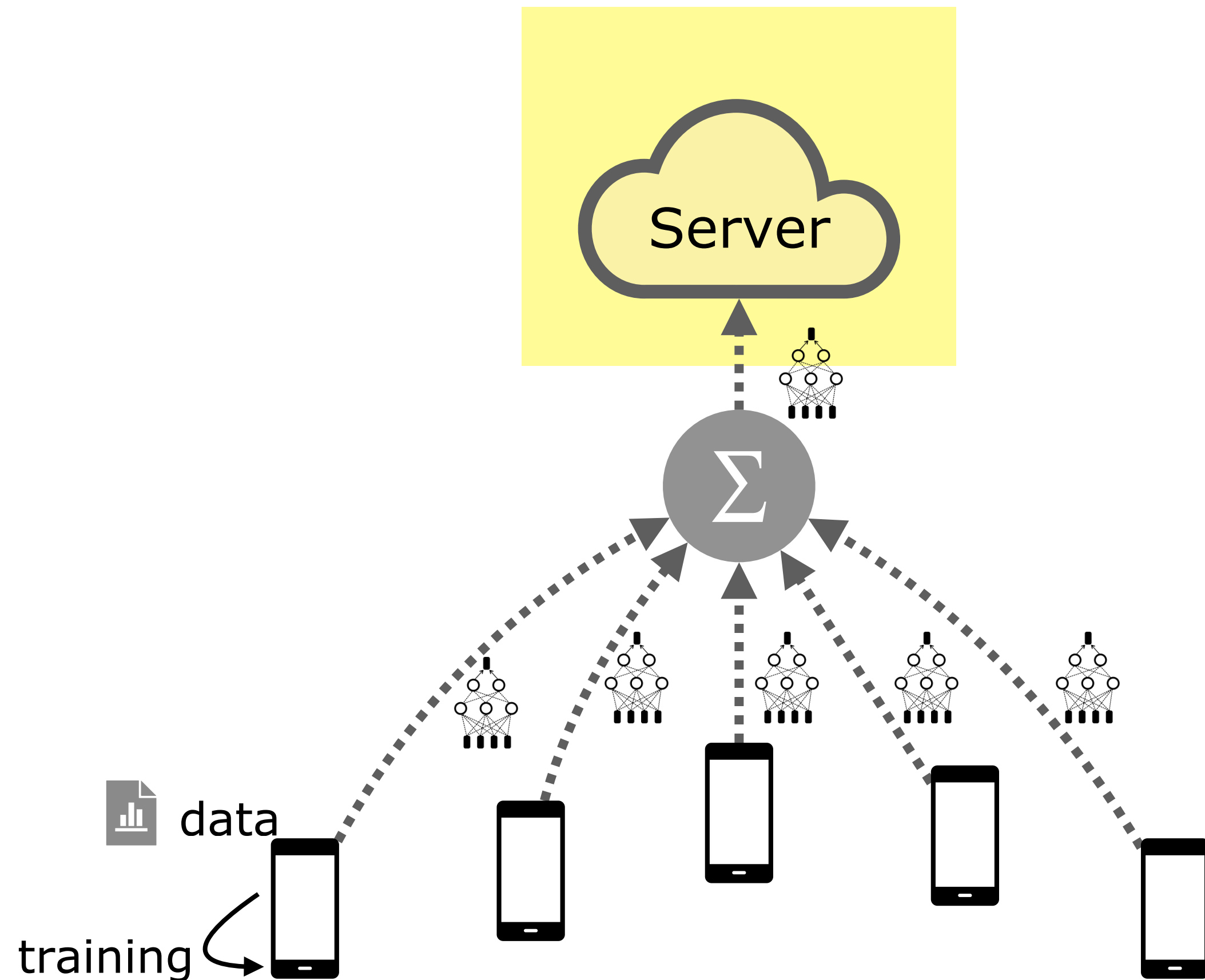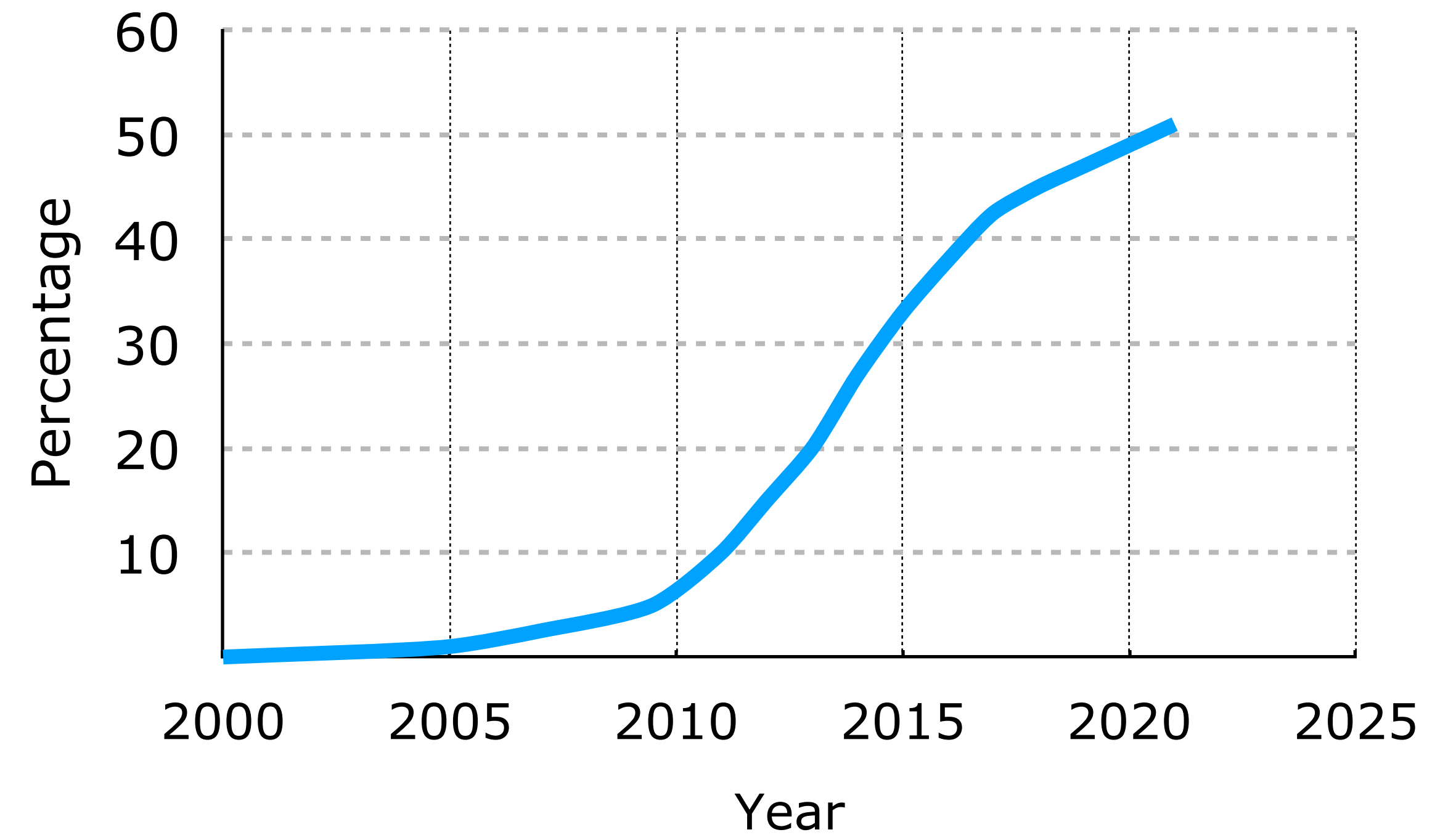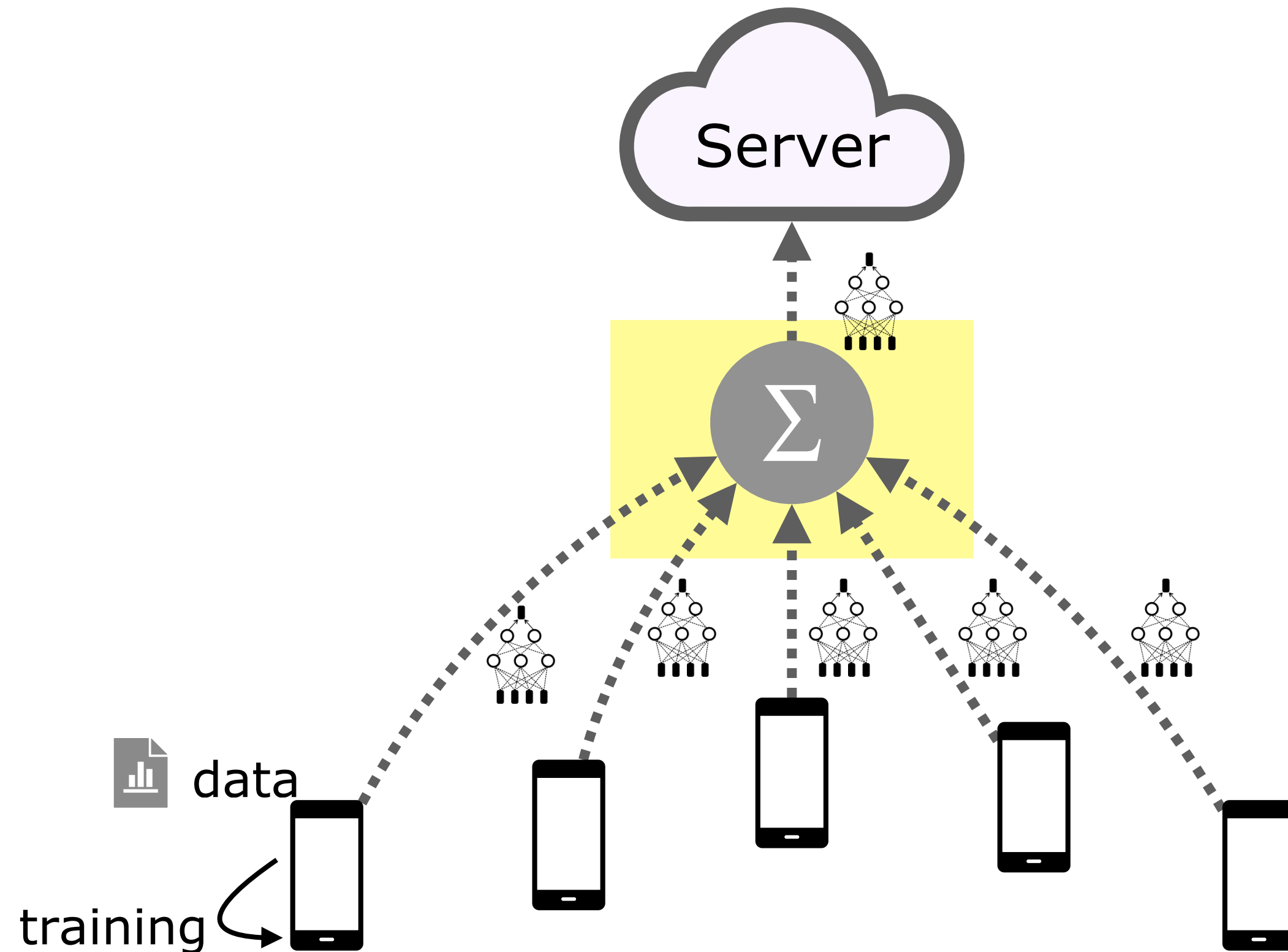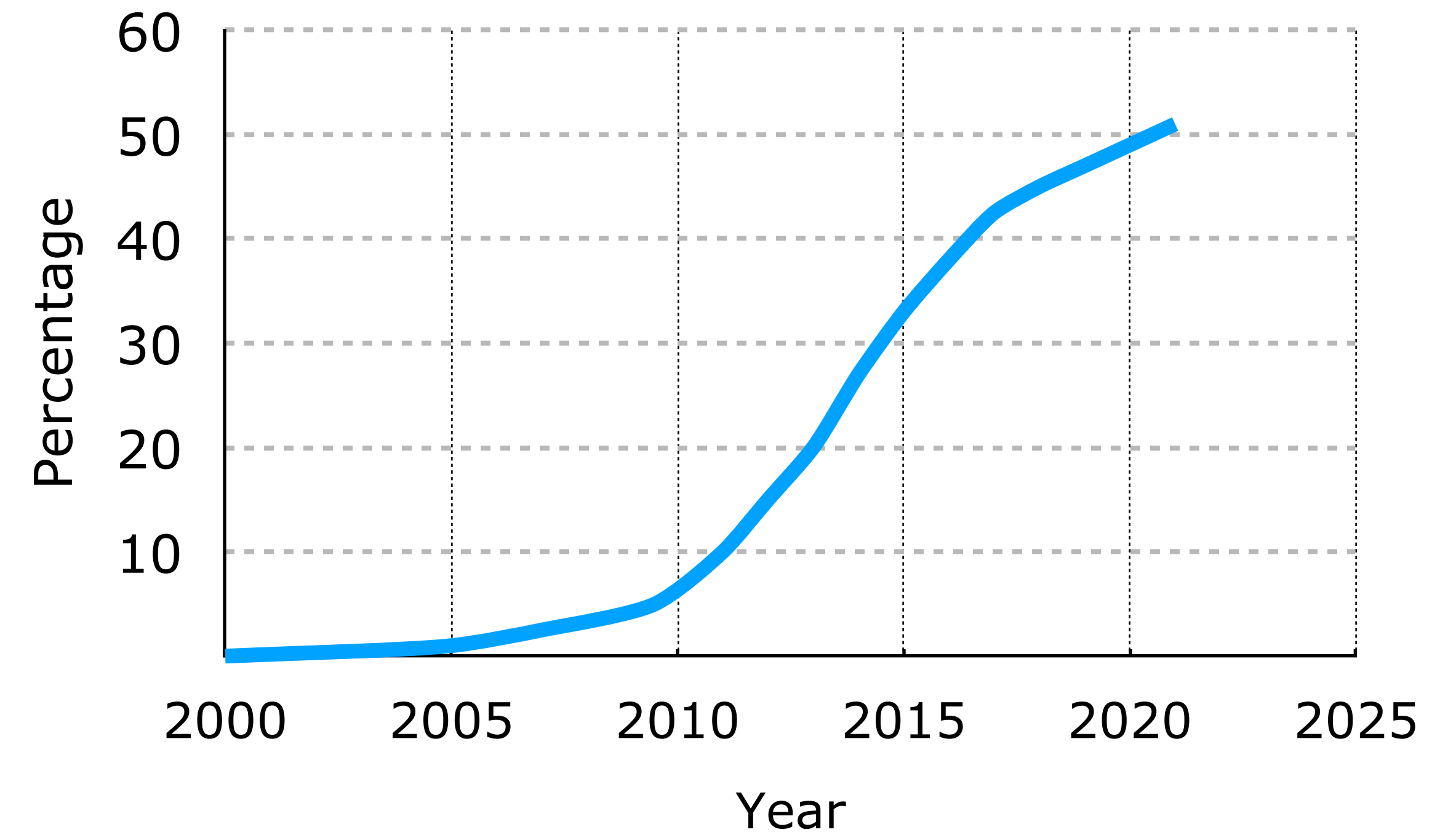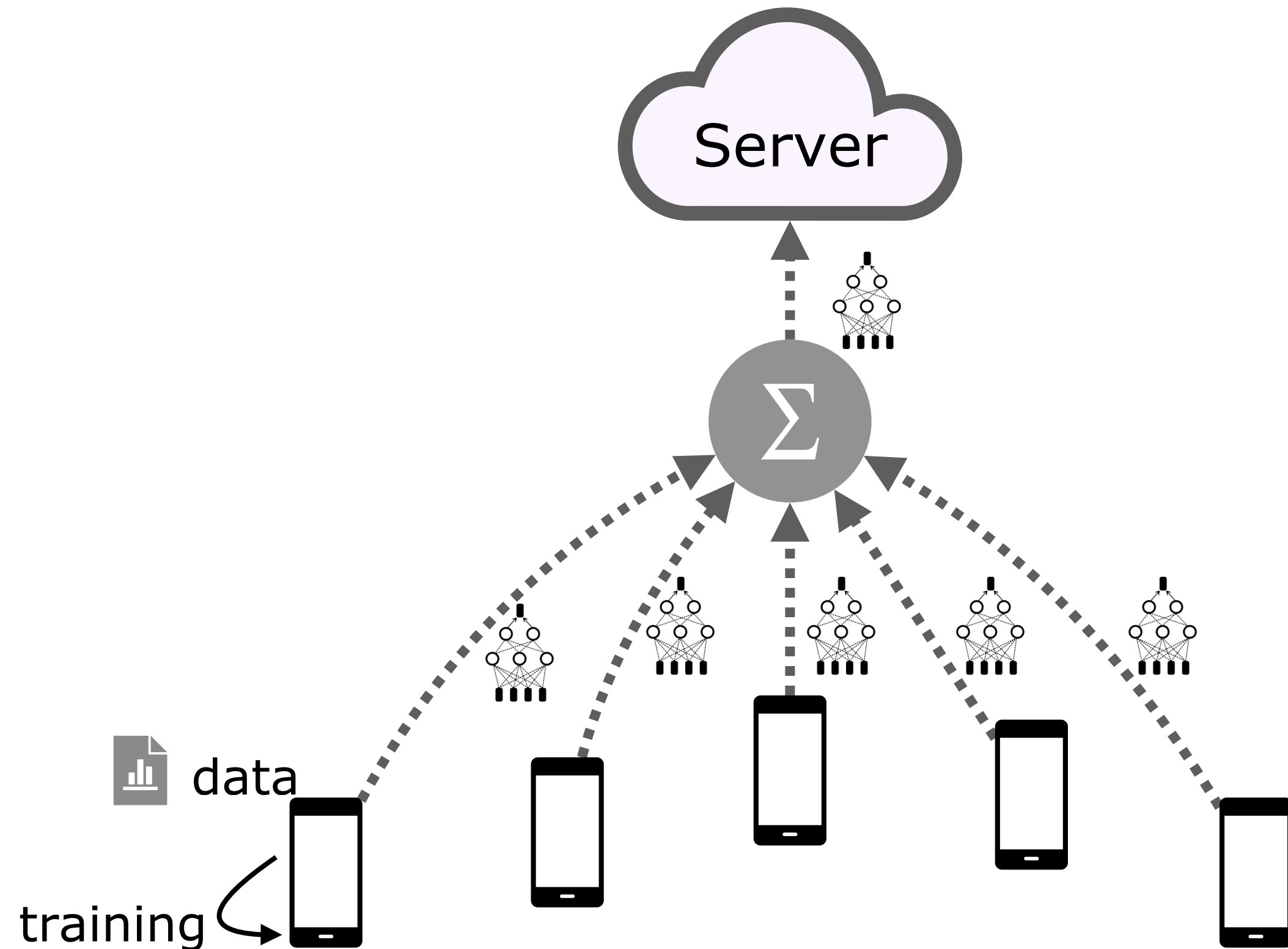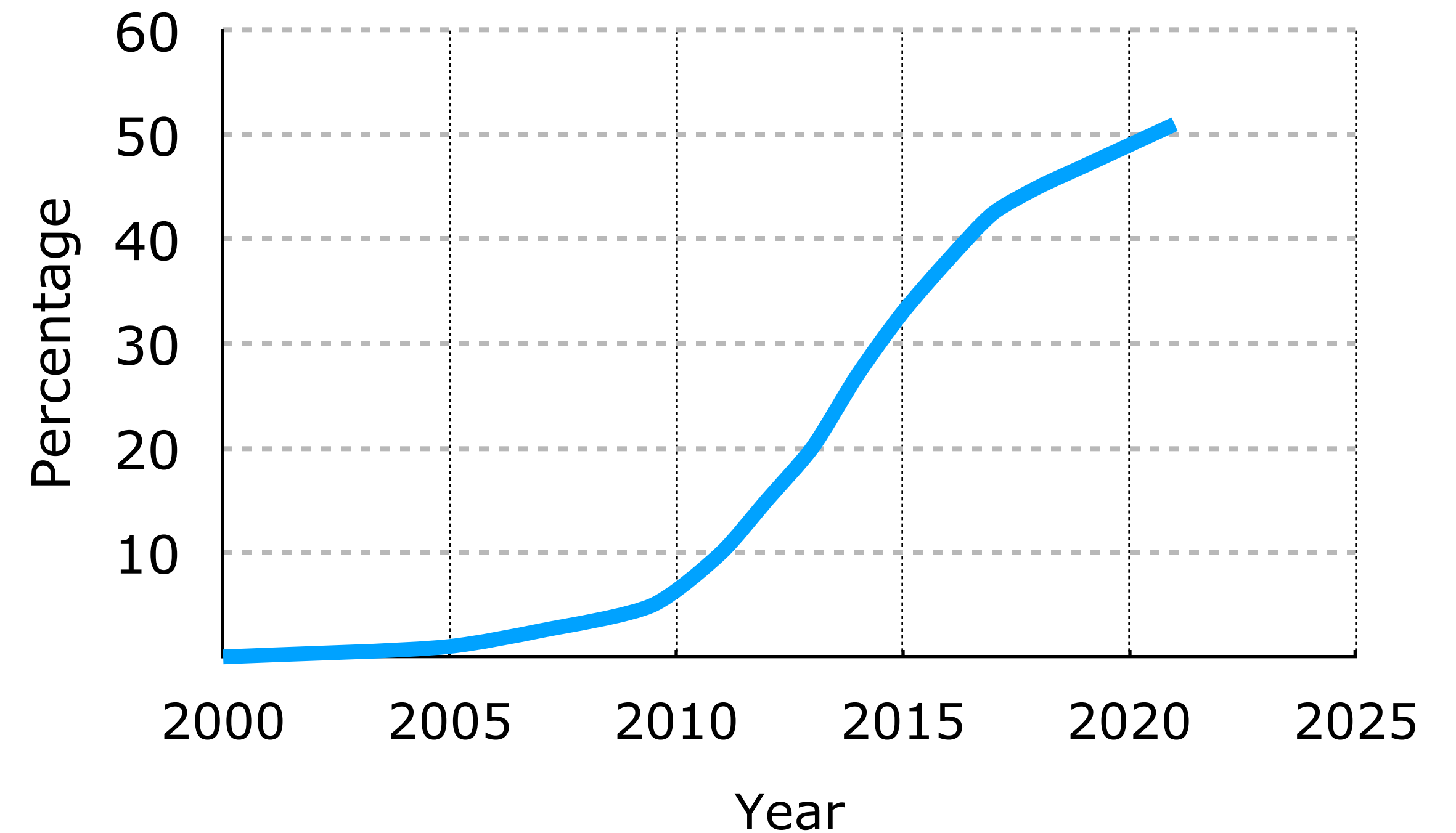


Communication cost > computation cost!

Google Research

Federated Learning: Collaborative Machine Learning without Centralized Training Data

April 6, 2017 · Posted by Brendan McMahan and Daniel Ramage, Research Scientists

How Apple personalizes Siri without hoovering up your data

The tech giant is using privacy-preserving machine learning to improve its voice assistant while keeping your data on your phone.

By Karen Hao

December 11, 2019

Engineering at Meta

POSTED ON JUNE 14, 2022 TO AI RESEARCH, ML APPLICATIONS, PRODUCTION ENGINEERING, SECURITY

Applying federated learning to protect data on mobile devices

Federated Learning for Postoperative Segmentation of Treated glioblastoma (FL-PoST)

IBM Federated Learning

Federated Learning

BANKING & PAYMENTS

Tencent's WeBank applying "federated learning" in A.I.

China's first mobile bank, Tencent's WeBank, is partnering with a H.K. startup to access decentralized sources of data.

Published 5 years ago on July 29, 2019

Federated learning in healthcare: the future of collaborative clinical and biomedical research

OWKIN

intel    PRODUCTS    SUPPORT    SOLUTIONS    DEVELOPERS    PARTNERS    FOUNDRY

Developers    /    Topics & Technologies ∨    /    Open Ecosystem ∨    /    Try Federated Learning with OpenFL

Try Federated Learning with OpenFL

14

# *Challenge:*
## 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

# Training

# Deployment



$$\bar{w} = \frac{1}{m}\sum_{i=1}^{m} w_i$$

Data poisoning

Model poisoning

Usual mean aggregation is ***not robust*** to corruptions $\implies$ Poor predictions!

**Usual approach**
(Direct)

**Our approach**
(Variational)

**Robust**

Robust to outliers/poisoning

**Usual approach**
(Direct)

**Our approach**
(Variational)

**Robust**

Robust to outliers/poisoning

**Communication efficient**

$O(1)$ times the communication cost as non-robust aggregation

**Secure aggregation**

Individual updates not revealed

**Usual approach** (Direct)   **Our approach** (Variational)

**Robust**

Robust to outliers/poisoning

**Communication efficient**

$O(1)$ times the communication cost as non-robust aggregation

**Secure aggregation**

Individual updates not revealed

??

??

# Robust aggregation approach

$w_1, \ldots, w_m$: updates sent by the clients



Server

*Robust Aggregation*

$w_1$    $w_2$                    $w_m$

Data poisoning          Model poisoning

**Facility location**

Fermat (~1600s)    Torricelli    Weber (1909)    Fréchet (~1940s)

**Facility location**

Fermat (~1600s)  Torricelli  Weber (1909)  Fréchet (~1940s)

**Facility location**

Fermat (~1600s)   Torricelli   Weber (1909)   Fréchet (~1940s)

**Facility location**

Fermat (~1600s)    Torricelli    Weber (1909)    Fréchet (~1940s)

**_Geometric Median_ /
Spatial Median /
$L_1$ Median /
Facility location**

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

| Fermat | Torricelli | Weber | Fréchet |
|--------|-----------|-------|---------|
| (~1600s) | | (1909) | (~1940s) |

# **Robustness**: Breakdown point = 1/2

(In **1D**, we have that *geometric median* ≡ *usual median*)



Inliers

Outliers

*Mean*

*Geometric Median*

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

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de $n$ points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$ & Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$



Weiszfeld a.k.a. Vázsonyi (1916-2003)

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de *n* points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$ & Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$



Weiszfeld a.k.a. Vázsonyi (1916-2003)

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de _n_ points donnes est minimum**. _Tohoku Mathematical Journal._

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$      &      Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$



Weiszfeld a.k.a. Vázsonyi (1916-2003)

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de _n_ points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$ & Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de _n_ points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$    &    Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de _n_ points donnes est minimum**. _Tohoku Mathematical Journal._

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2,\ \nu\}}$     &     Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$
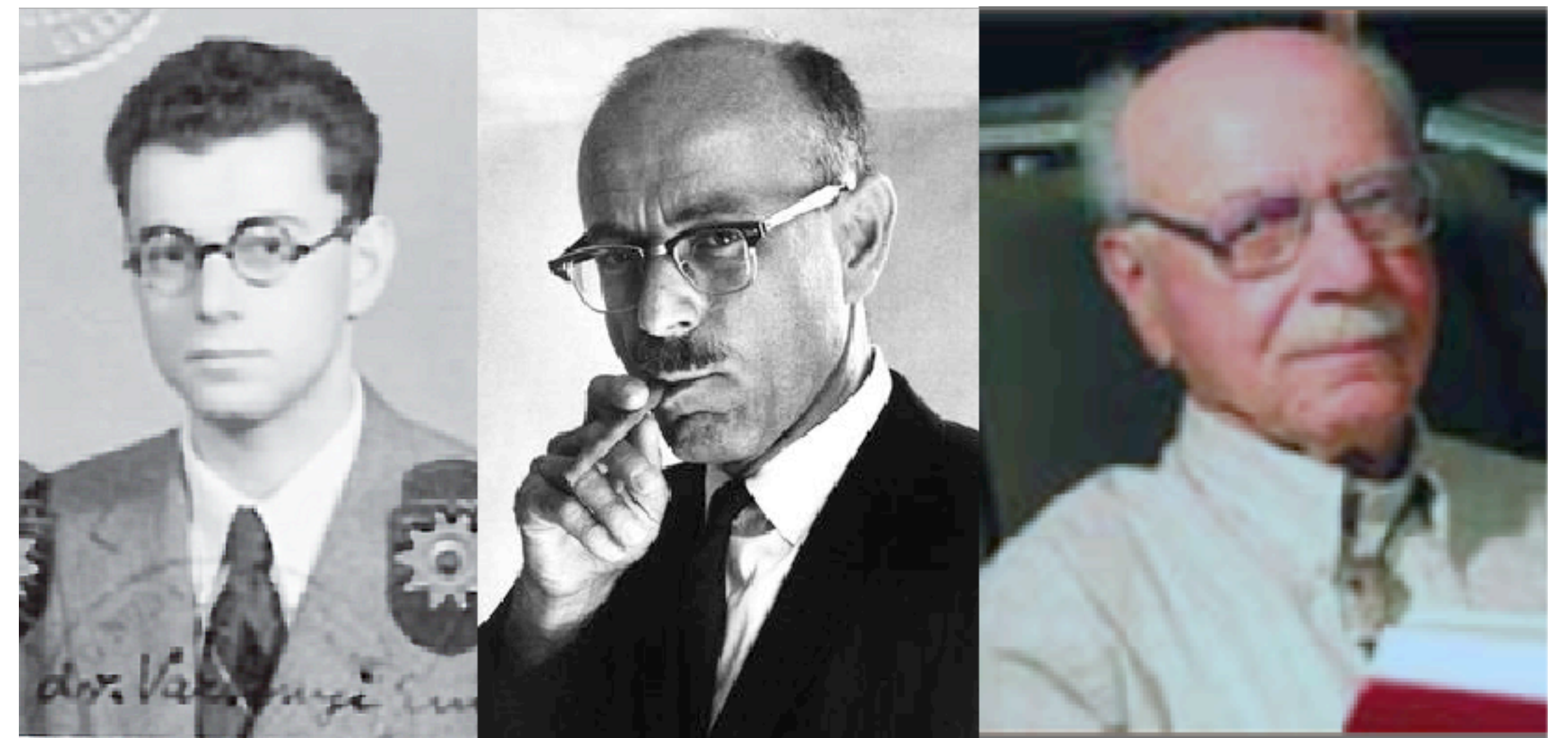
# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de n points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$  &  Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$
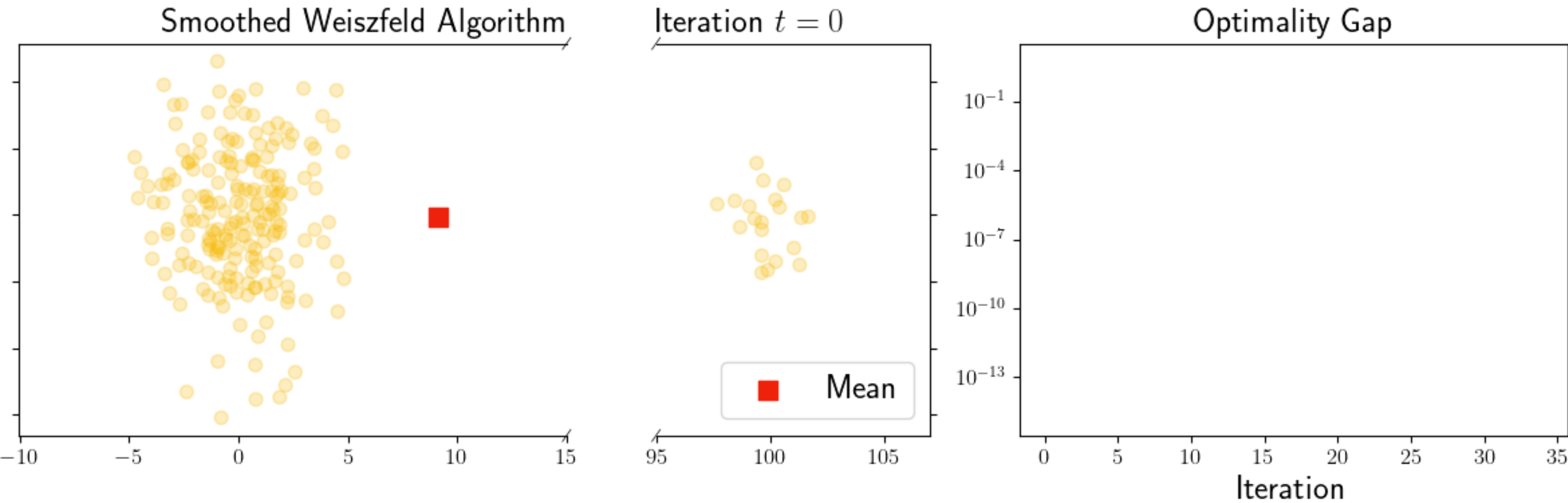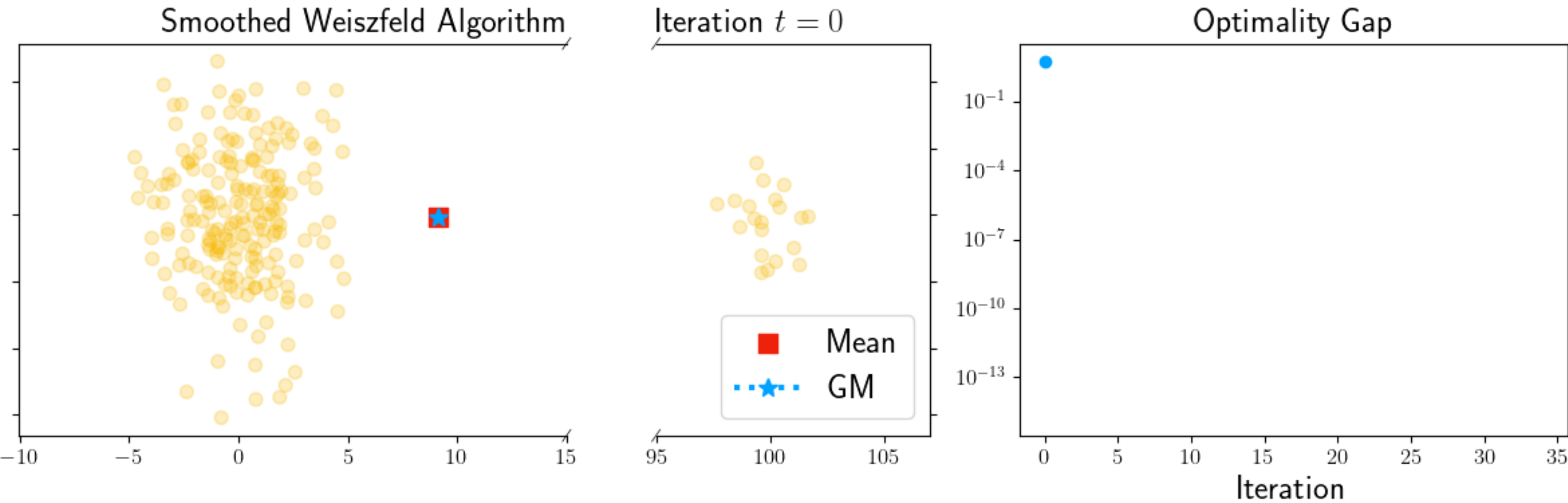
# Usual approach
(Direct)

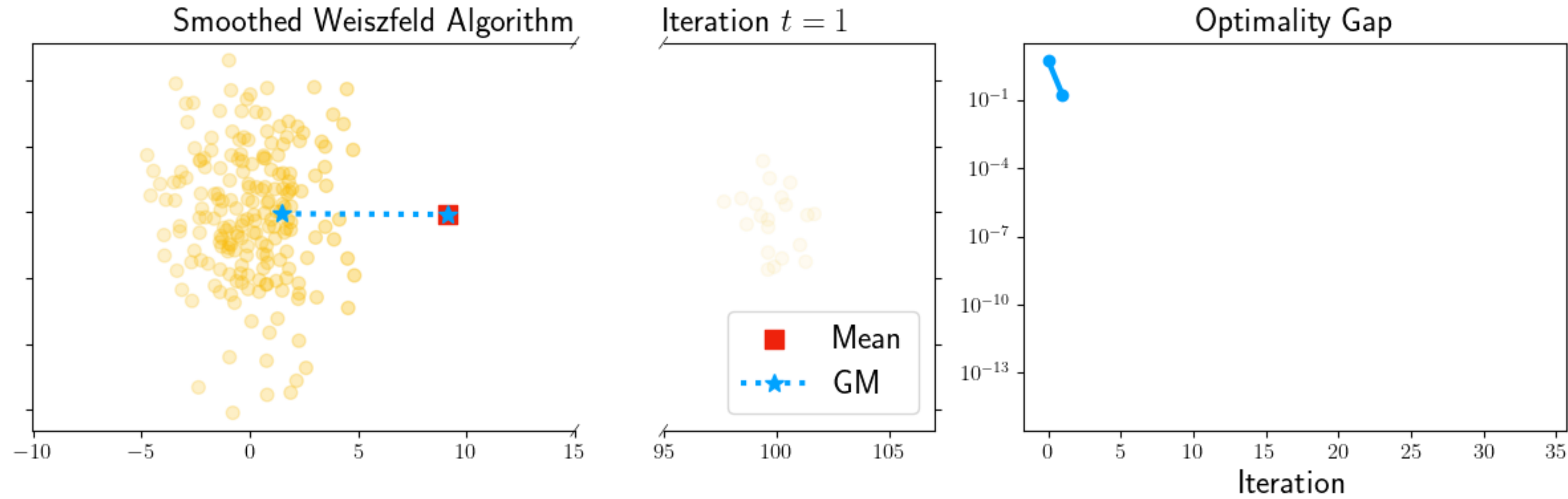# Our approach
(Variational)

**Robust**

Robust to outliers/poisoning

❌ ✔️

**Communication efficient**

$O(1)$ times the communication cost as non-robust aggregation

✔️ **??**

**Secure aggregation**

Individual updates not revealed

✔️ **??**

# Communication primitive: secure sum/average

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

$x_1$

Client 1

$x_2$

Client 2

$\Sigma$

$x_1 + x_2$ Server

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

# Perform all operations modulo $M$



Client 1

Client 2

$\Sigma$

Server

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

Client 1

$$x_1' = x_1 + \xi$$

$$\xi \sim \mathsf{Unif}\left(\bigcirc\right)$$

Client 2

$$x_2' = x_2 - \xi$$

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

38

Server only sees $x_1', x_2' \sim \text{Unif}(\bigcirc)$ but calculates the correct sum (and average)

Client 1

$x_1' = x_1 + \xi$

$\xi \sim \text{Unif}(\bigcirc)$

$x_1' + x_2' = x_1 + x_2$

Client 2

$x_2' = x_2 - \xi$

$\Sigma$

Server

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

38

Server only sees $x_1', x_2' \sim \text{Unif}\left(\bigcirc\right)$ but calculates the correct sum (and average)



Client 1

$x_1' = x_1 + \xi$

$\xi \sim \text{Unif}\left(\bigcirc\right)$

$x_1' + x_2' = x_1 + x_2$

Client 2

$x_2' = x_2 - \xi$

$\Sigma$

Server

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

$x_1' = x_1 + \xi$

$x_1'$    $\xi$

Client 1    $x_1$    $x_1$    $x_1' + x_2' = x_1 + x_2$

**Real-world communication constraint**:
All client-to-server communication must go through secure average

Client 2    $x_2$    $x_2$

Server

$x_2' = x_2 - \xi$    $-\xi$

**Extensions:** weighted sums/averages

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de $n$ points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2,\ \nu\}}$ & Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de $n$ points donnes est minimum**. *Tohoku Mathematical Journal.*

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$     &     Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$

1. Server **broadcasts** current estimate $z_t$ of the geometric median

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de $n$ points donnes est minimum**. *Tohoku Mathematical Journal.*
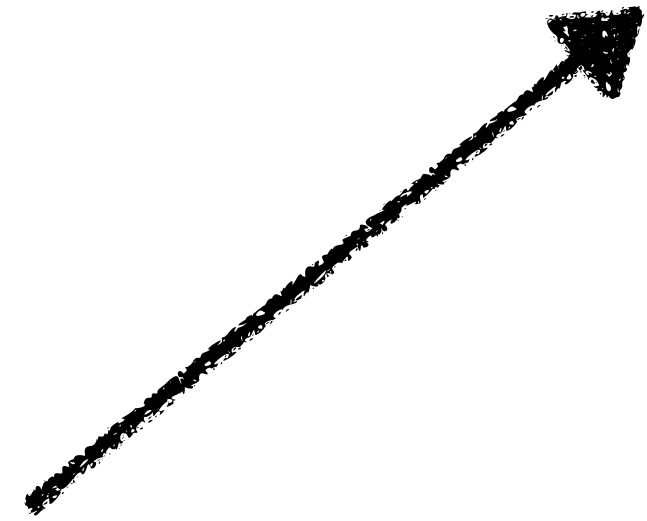
Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2,\ \nu\}}$      &      Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$

2. **Clients** compute new weights

# Smoothed Weiszfeld Algorithm

Weiszfeld (1937). **Sur le point par lequel la somme des distances de _n_ points donnes est minimum**. _Tohoku Mathematical Journal._

Compute new weights $\beta_{i,t} = \dfrac{1}{\max\{\|z_t - w_i\|_2, \nu\}}$    &    Reweighted average $z_{t+1} = \dfrac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$
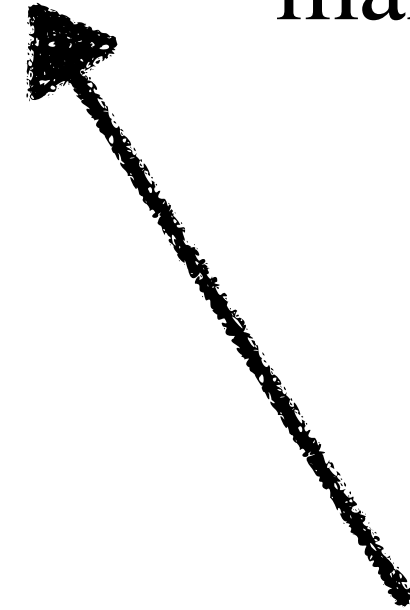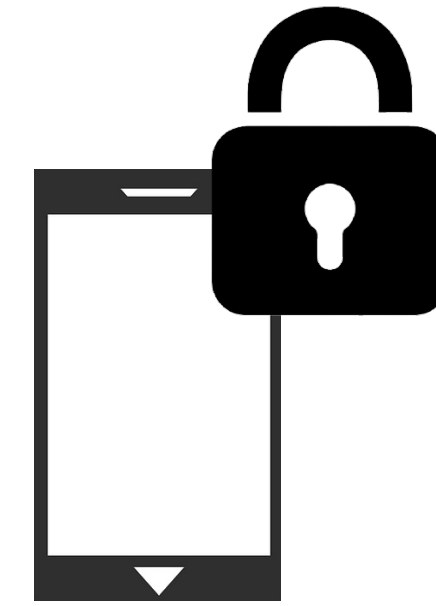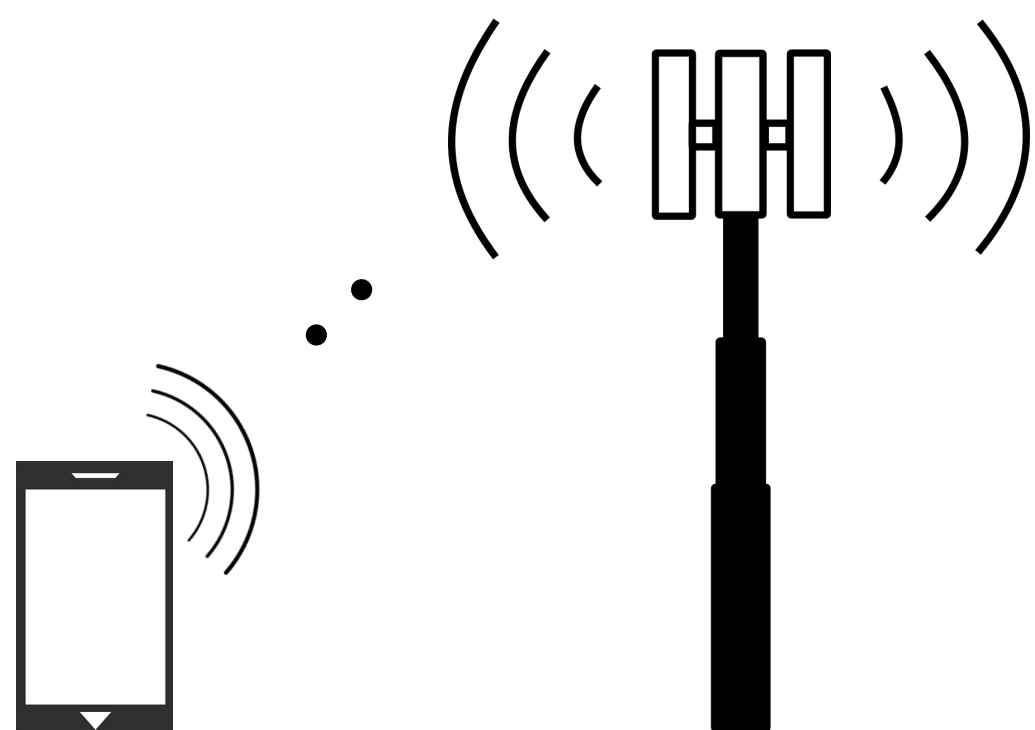
3. Obtain new estimate by **secure averaging**

# **Secure aggregation**

Only client-server communication is via **secure average** in the **Smoothed Weiszfeld Algorithm**

$$z_{t+1} = \frac{\sum_i \beta_{i,t} w_i}{\sum_i \beta_{i,t}}$$
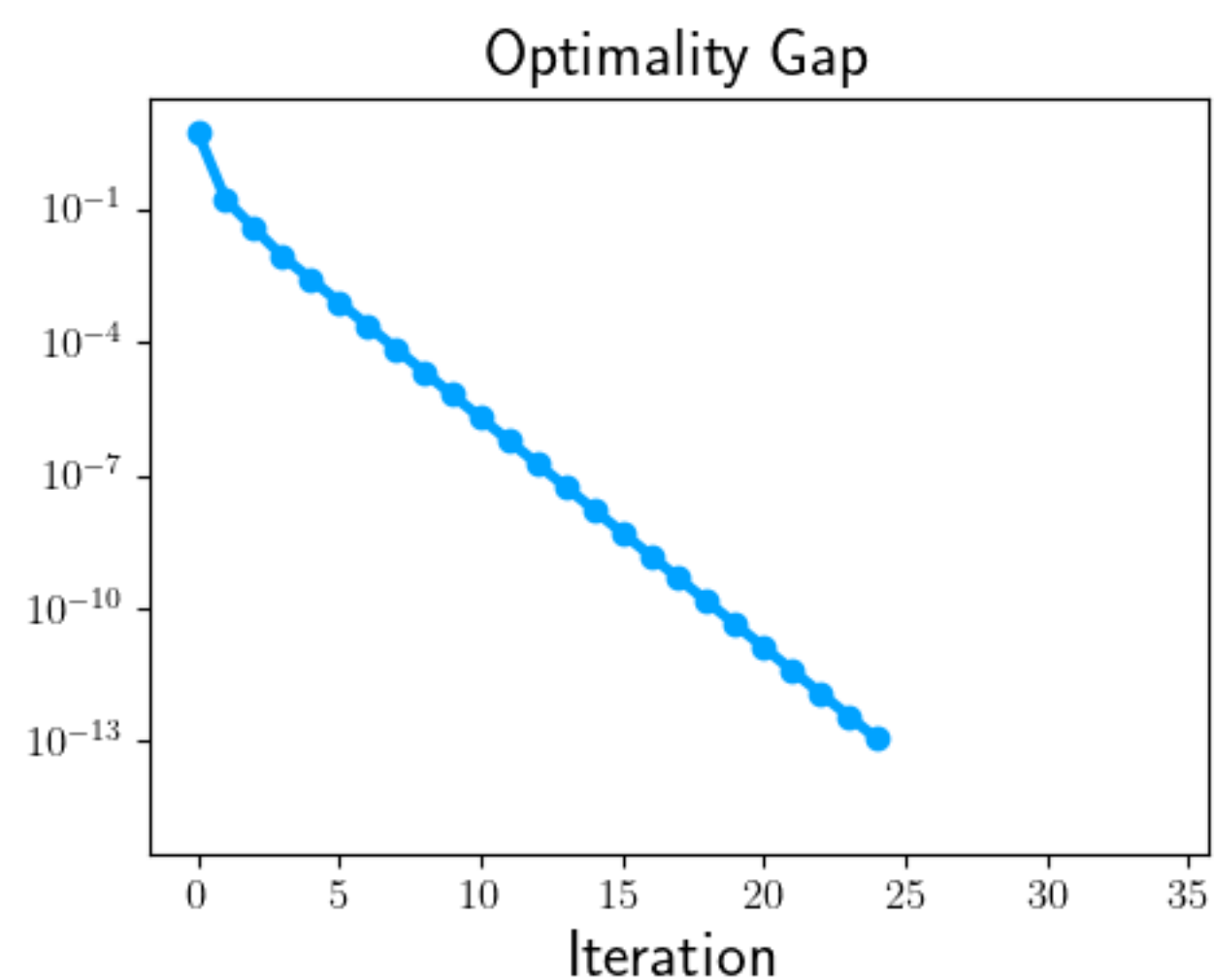
**Communication efficient!**

Empirically, **3-5** iterations suffice:

provably rapid convergence



Even **1** iteration improves robustness!

**Usual approach**
(Direct)

**Our approach**
(Variational)

**Robust**

Robust to outliers/poisoning

**Communication efficient**

$O(1)$ times the communication cost as non-robust aggregation

**Secure aggregation**

Individual updates not revealed

# Robust Federated Aggregation (RFA)

**More robust federated learning** =

Local SGD steps +

Geometric median + secure aggregation

*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 federated averaging

*Step 3 of 3: Aggregate with multiple rounds of secure average (weights $\beta_i$ from the Smoothed Weiszfeld Algorithm)*

Round 1 of Aggregation

Round 2 of Aggregation

Round 3 of Aggregation

Weights

# See the paper for:

Discussion on *heterogeneity*

*Convergence* analysis

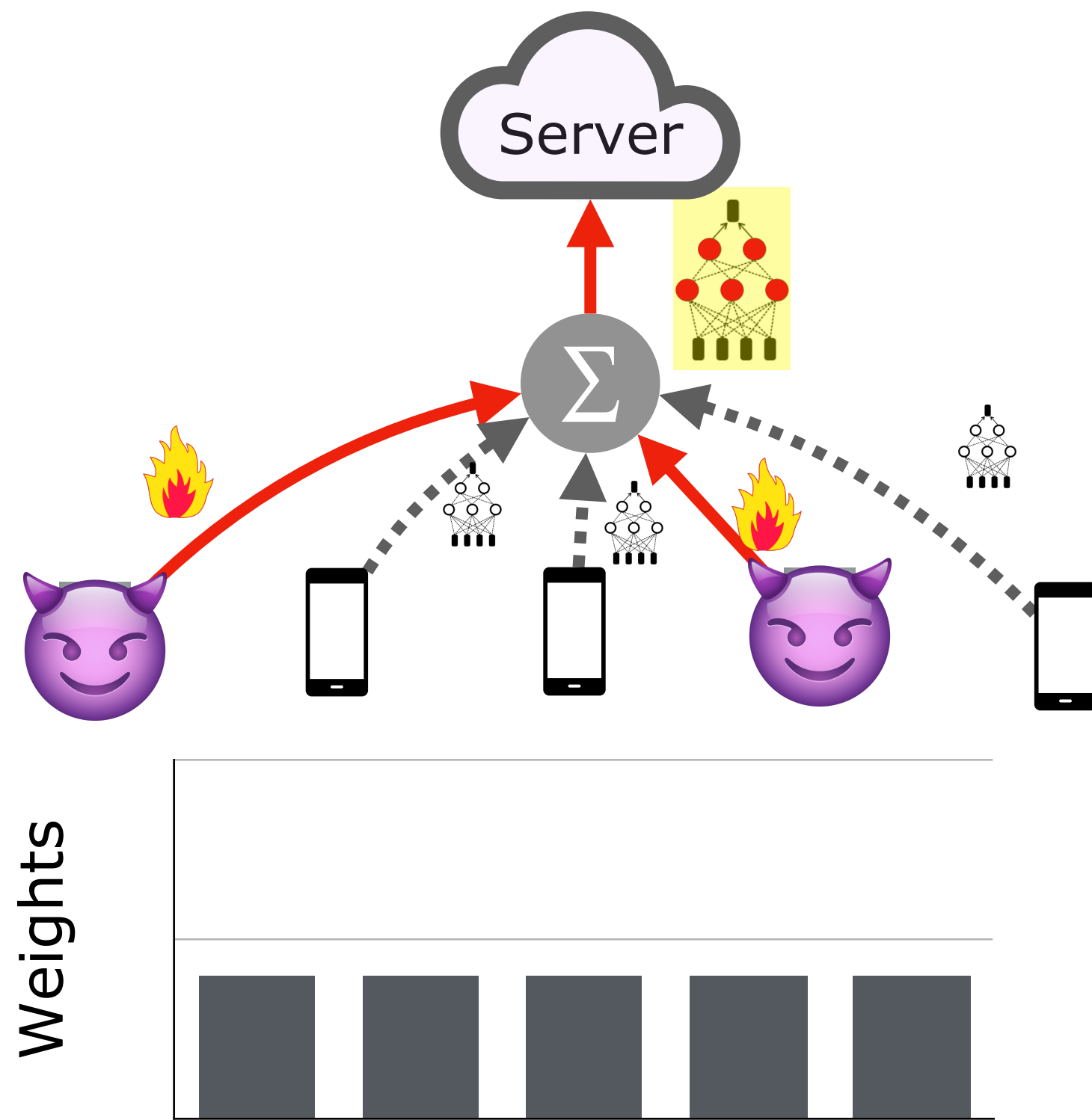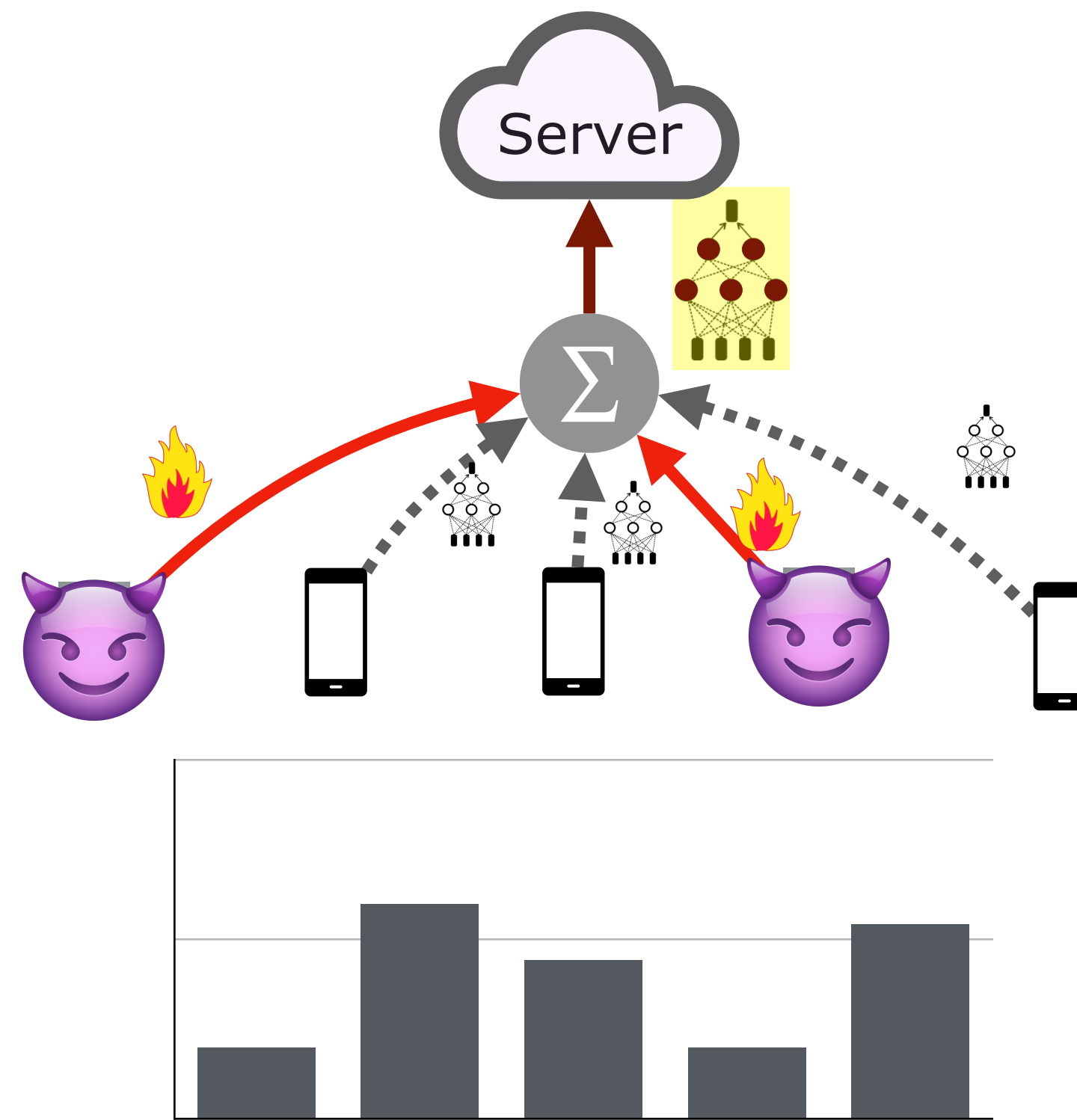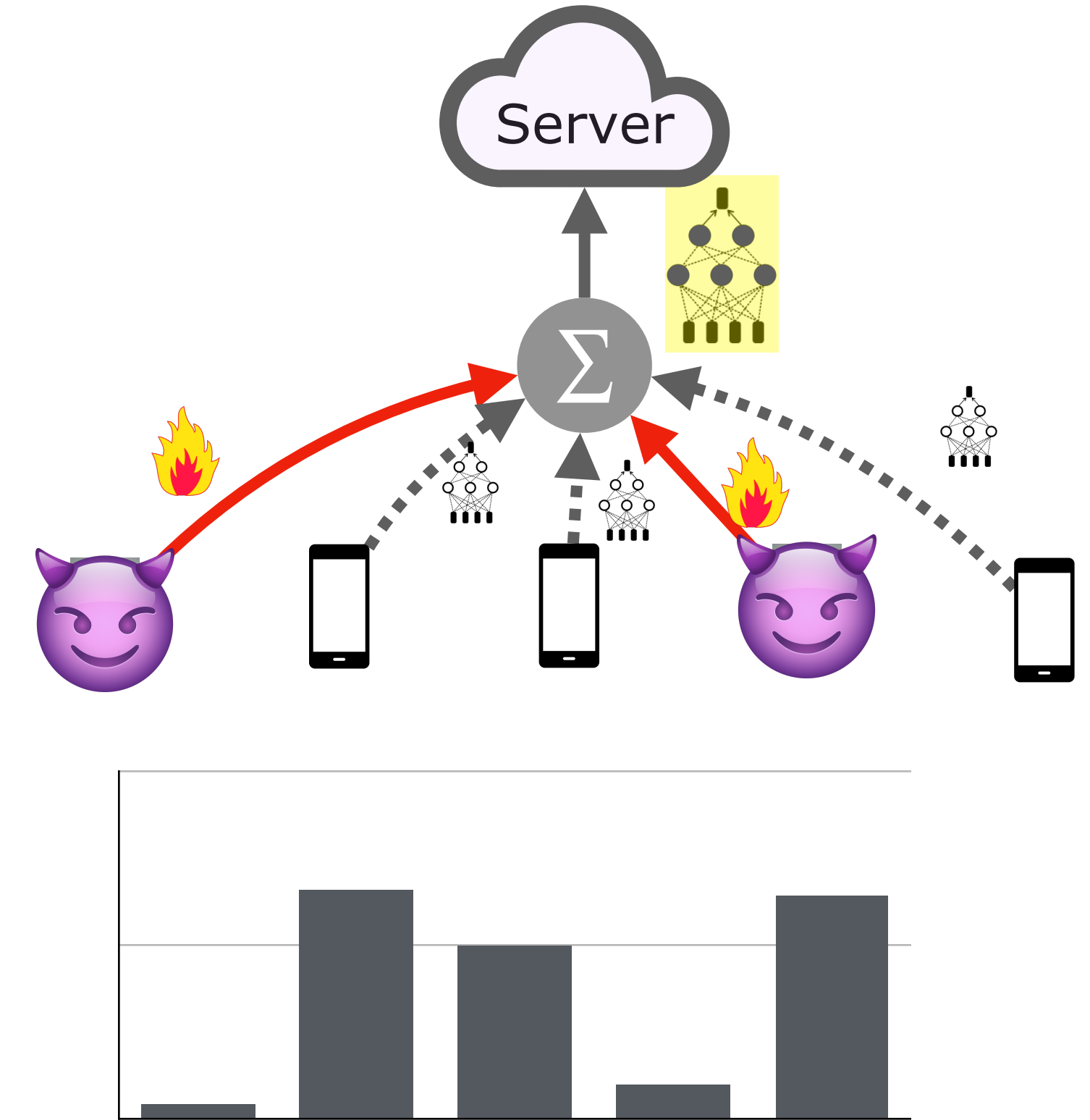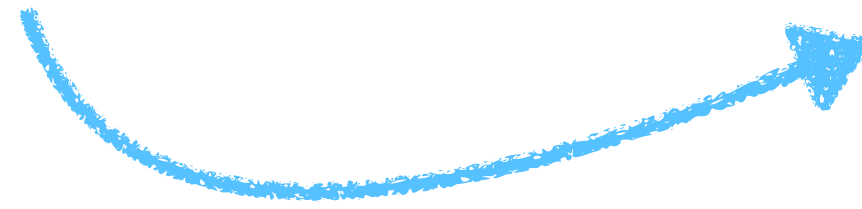**The Tension Between Robustness and Heterogeneity:** Heterogeneity is a key property of federated learning. The distribution $D_i$ of device $i$ can be quite different from the distribution $D_j$ of some other device $j$, reflecting the heterogeneous data generated by a diverse set of users.

To analyze the effect of heterogeneity on robustness, consider the simplified scenario of robust mean estimation in Huber's contamination model [34]. Here, we wish to estimate the mean $\mu \in \mathbb{R}^d$ given samples $w_1, \ldots, w_m \sim (1 - \rho)\mathcal{N}(\mu, \sigma^2 I) + \rho Q$, where $Q$ denotes some outlier distribution that $\rho$-fraction of the points (designated as outliers) are drawn from. Any aggregate $\bar{w}$ must satisfy the lower bound $\|\bar{w} - \mu\|^2 \geq \Omega(\sigma^2 \max\{\rho^2, d/m\})$ with constant probability [69, Theorem 2.2]. In the federated learning setting, more heterogeneity corresponds to a greater variance $\sigma^2$ among the inlier points, implying a larger error in mean estimation. This suggests a tension between robustness and heterogeneity, where increasing heterogeneity makes robust mean estimation harder in terms of $\ell_2$ error.

In this work, we strike a compromise between robustness and heterogeneity by considering a family $\mathcal{D}$ of allowed data

**Convergence:** We now analyze RFA where the local SGD updates are equipped with "tail-averaging" [73] so that $w_i^{(t+1)} = (2/\tau) \sum_{k=\tau/2}^{\tau} w_{i,k}^{(t)}$ is averaged over the latter half of the trajectory of iterates instead of line 9 of Algorithm 1. We show that this variant of RFA converges up to the dissimilarity level $\Omega = \Omega_X \Omega_{Y|X}$ when the corruption level $\rho < 1/2$.
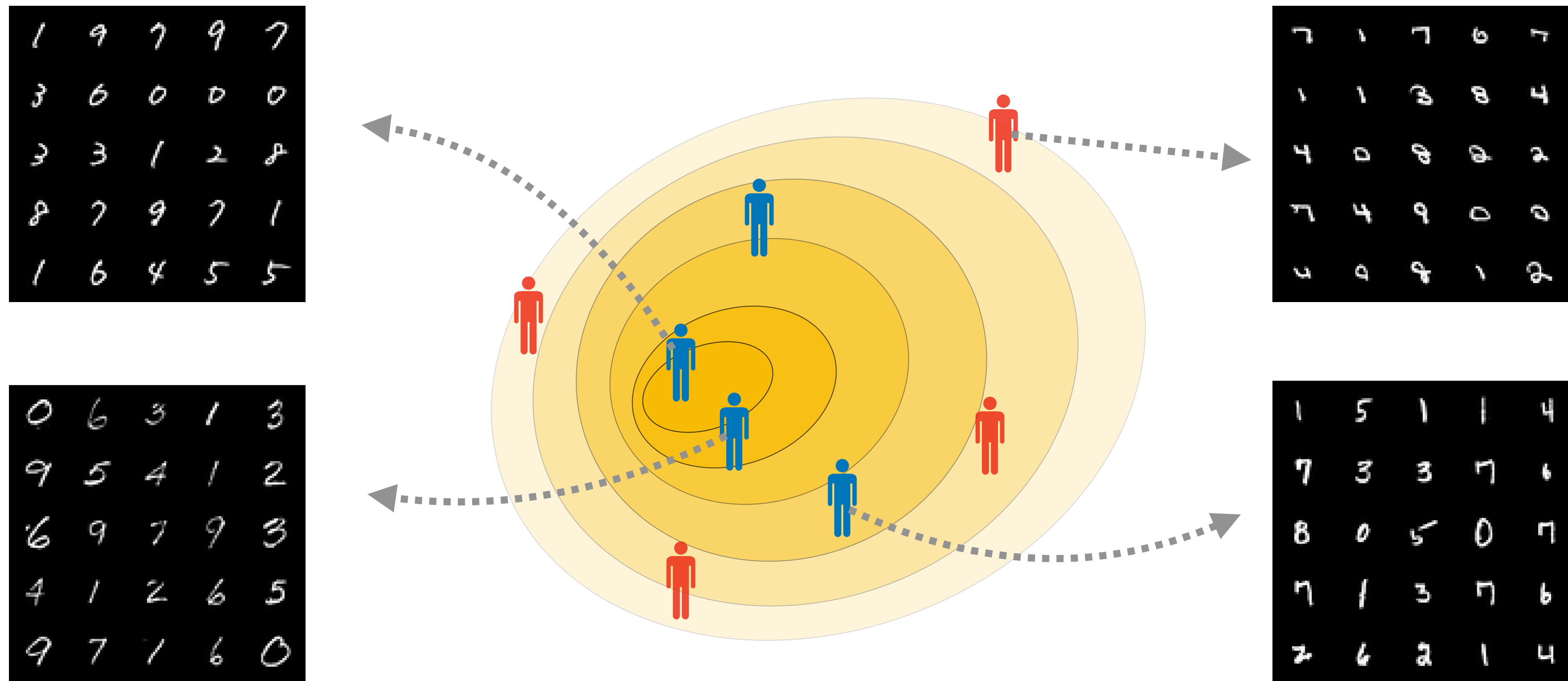
*Theorem 4:* Consider $F$ defined in (7) and suppose the corruption level satisfies $\rho < 1/2$. Consider Algorithm 1 run for $T$ outer iterations with a learning rate $\gamma = 1/(2R^2)$, and the local updates are run for $\tau_t$ steps in outer iteration $t$ with tail averaging. Fix $\delta > 0$ and $\theta \in (\rho, 1/2)$, and set the number of devices per iteration, $m$ as

$$m \geq \frac{\log(T/\delta)}{2(\theta - \rho)^2}. \tag{11}$$
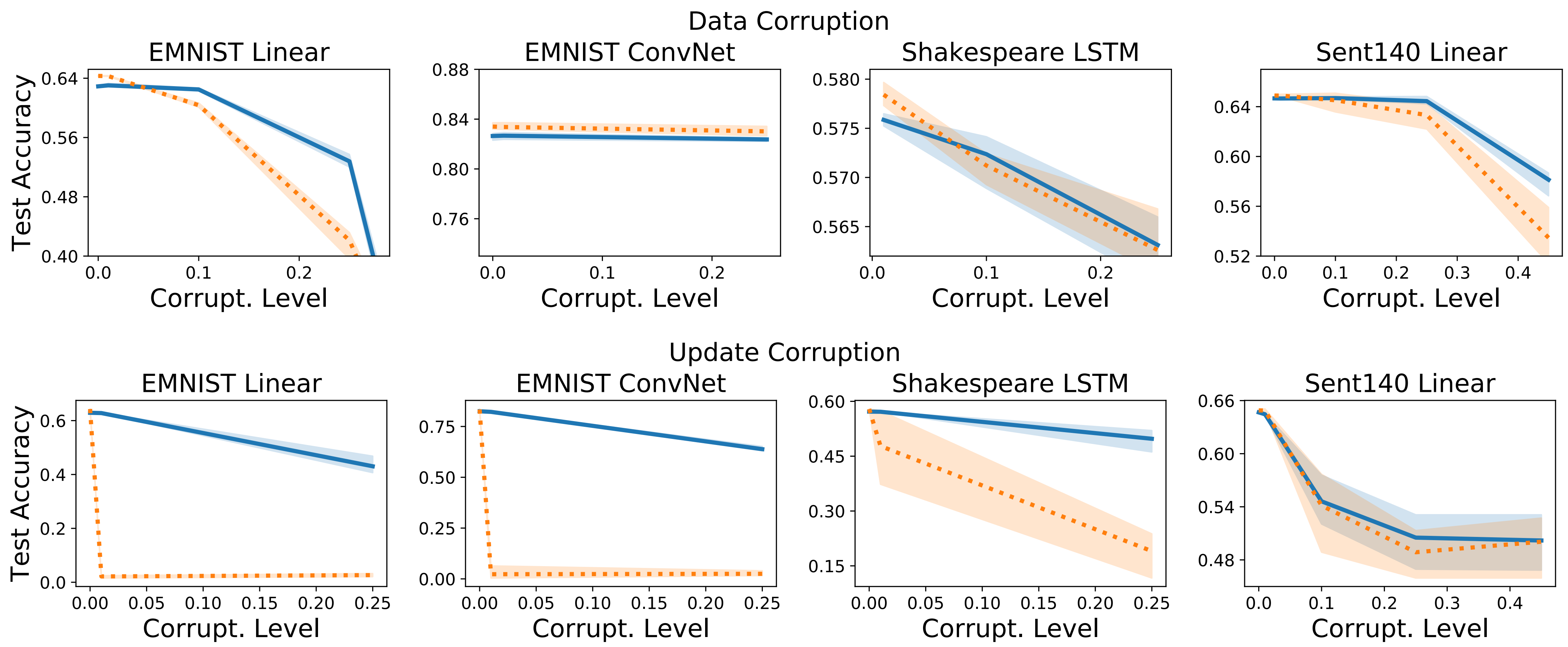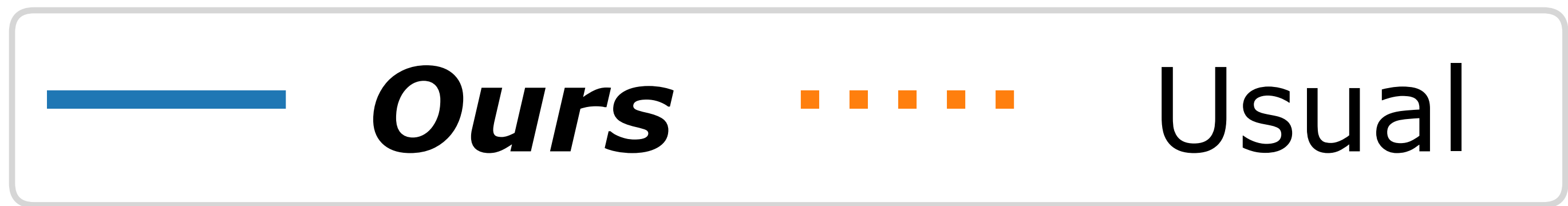
Define $C_\theta := (1 - 2\theta)^{-2}, w^\star = \arg\min F, F^\star = F(w^\star), \kappa := R^2/\mu$ and $\Delta_0 := \|w^{(0)} - w^\star\|^2$. Let $\tau \geq 4\kappa \log(128 C_\theta \kappa)$. We have that the event $\mathcal{E} = \bigcap_{t=0}^{T-1} \{|S_t \cap \mathcal{C}| \leq \theta m\}$ holds with probability at least $1 - \delta$. Further, if $\tau_t = 2^t \tau$ for each iteration $t$, then the output $w^{(T)}$ of Algorithm 1 satisfies,

$$\mathbb{E}\left[\|w^{(T)}) - w^\star\|^2 \,\middle|\, \mathcal{E}\right] \leq \frac{\Delta_0}{2^T} + CC_\theta \left(\frac{d\sigma^2 T}{\mu\tau 2^T} + \frac{\epsilon^2}{m^2} + \Omega^2\right)$$
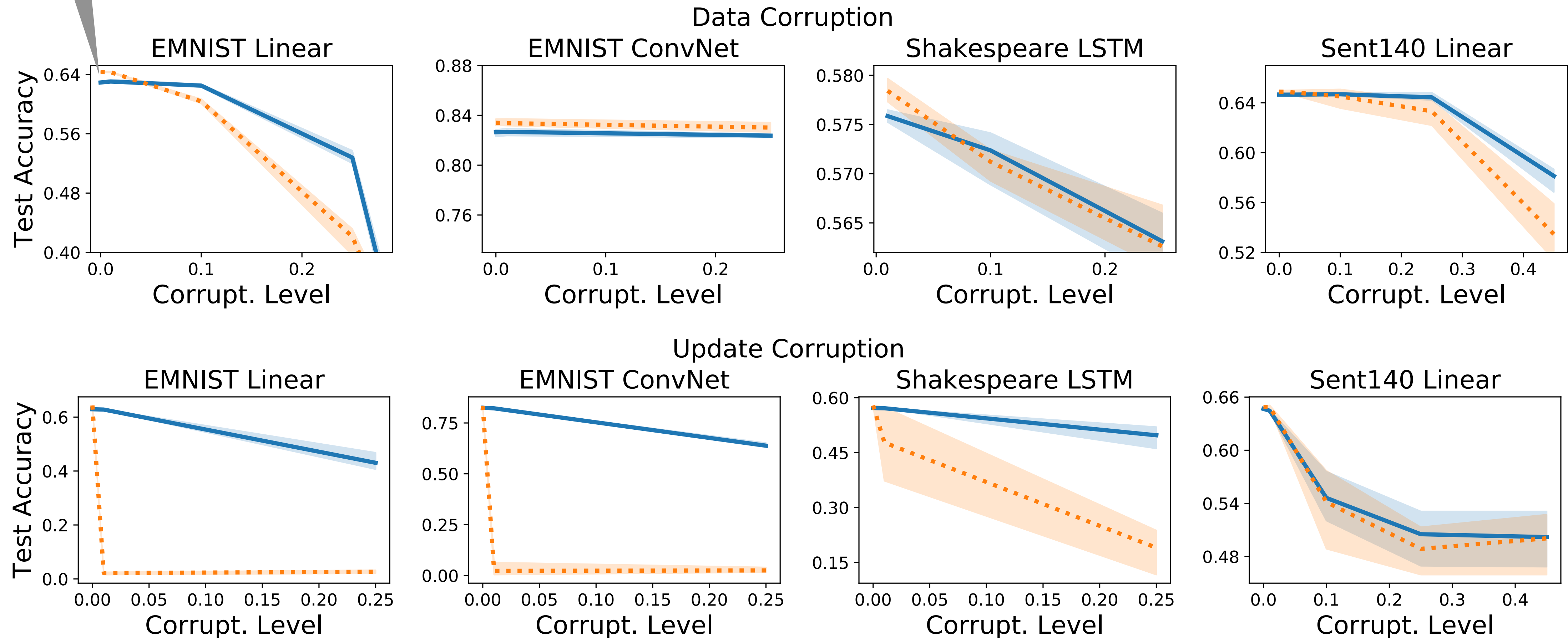
# Experiments and Improvements

Data Corruption

EMNIST Linear

EMNIST ConvNet

Shakespeare LSTM

Sent140 Linear

Update Corruption

EMNIST Linear

EMNIST ConvNet

Shakespeare LSTM

Sent140 Linear

Data Corruption

| EMNIST Linear | EMNIST ConvNet | Shakespeare LSTM | Sent140 Linear |

Update Corruption

| EMNIST Linear | EMNIST ConvNet | Shakespeare LSTM | Sent140 Linear |

1.4pp gap at zero corruption

*Ours*    Usual
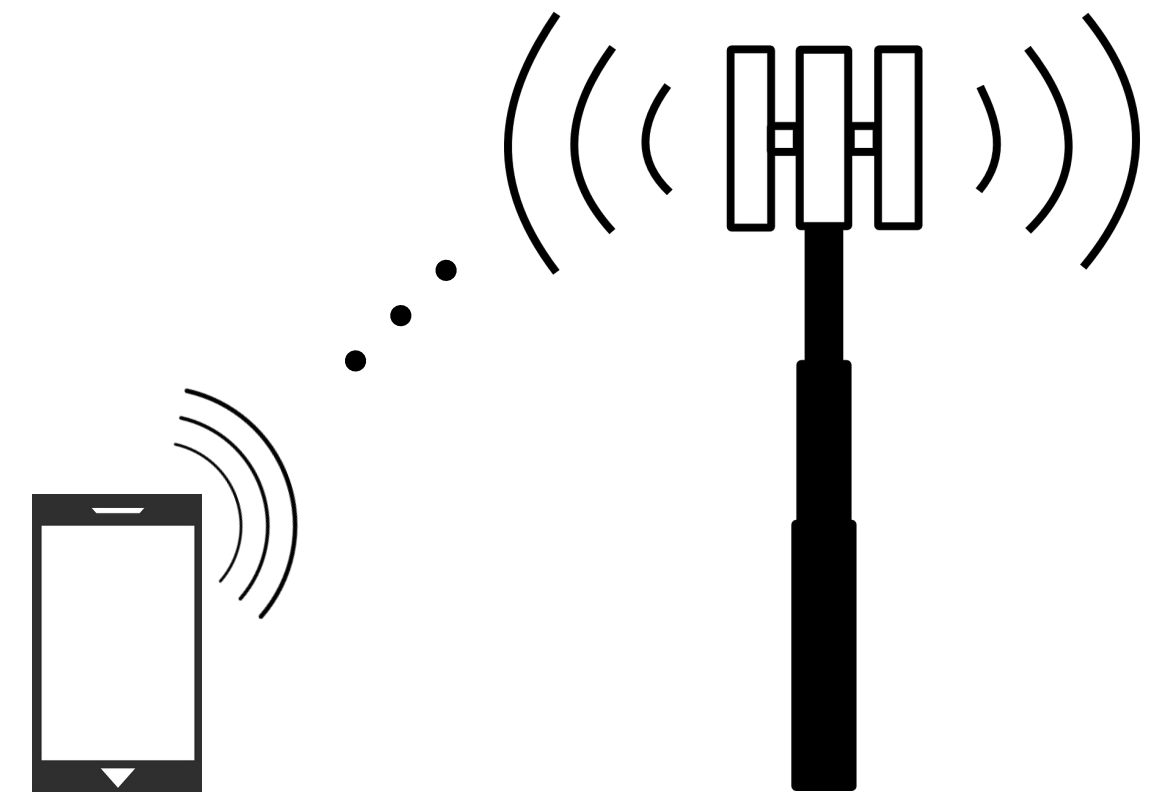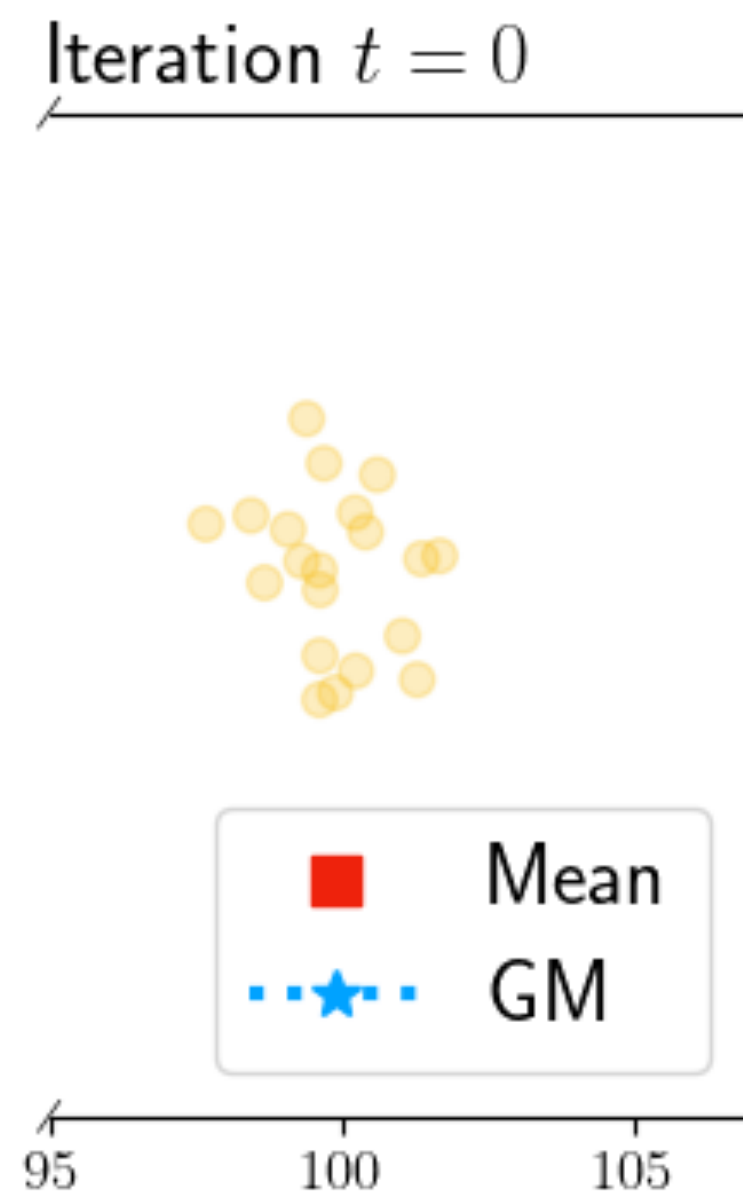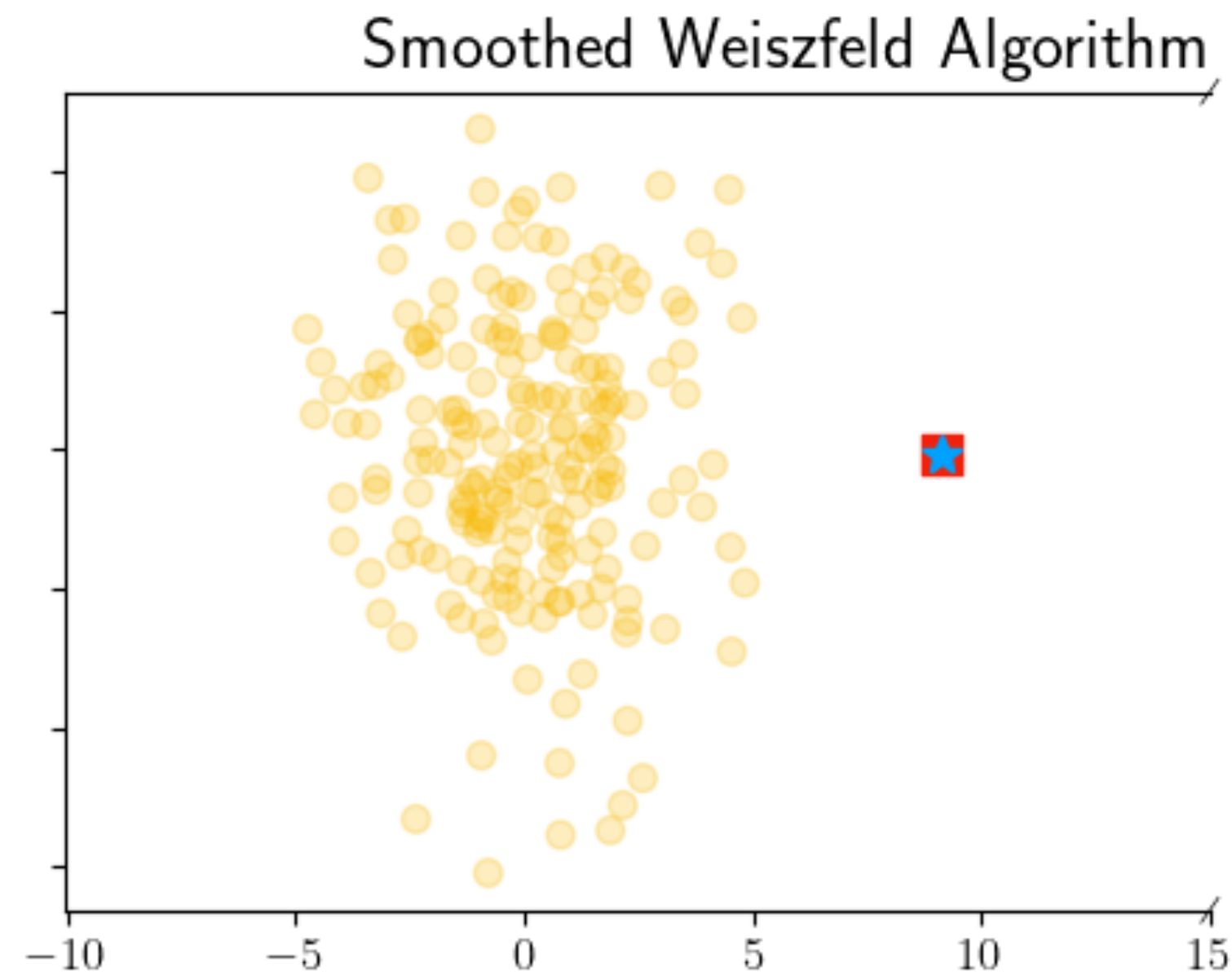
# Reducing the communication cost

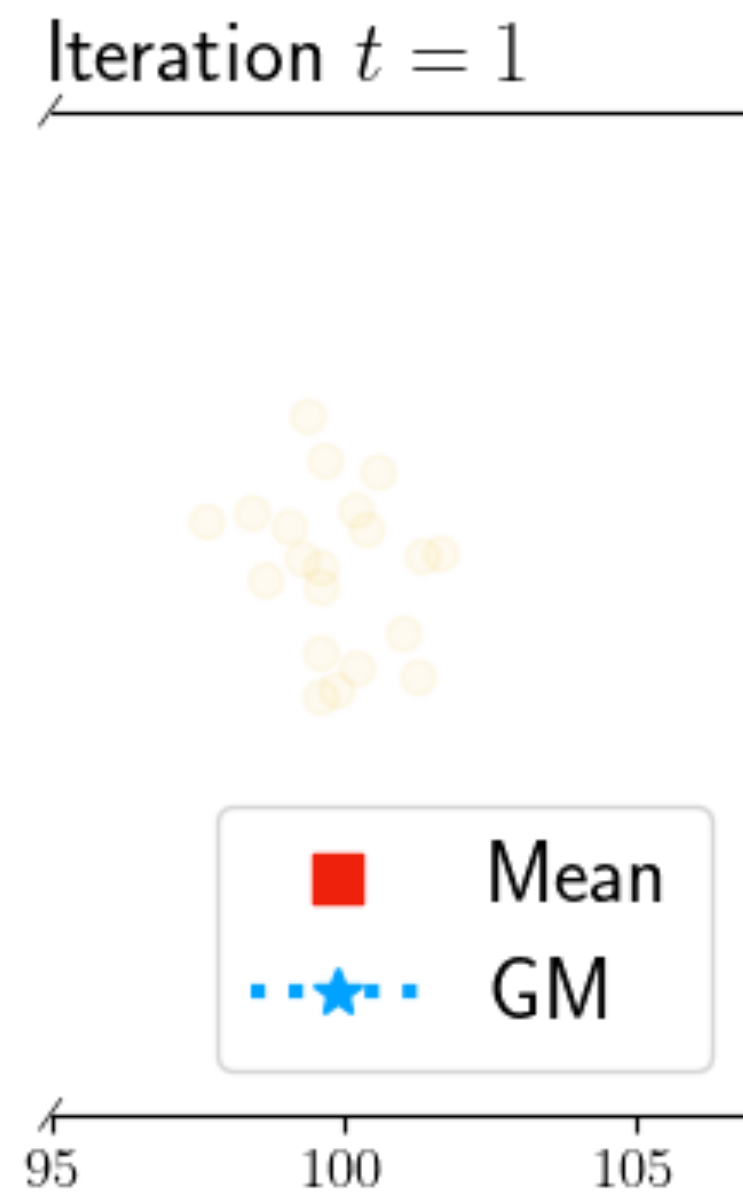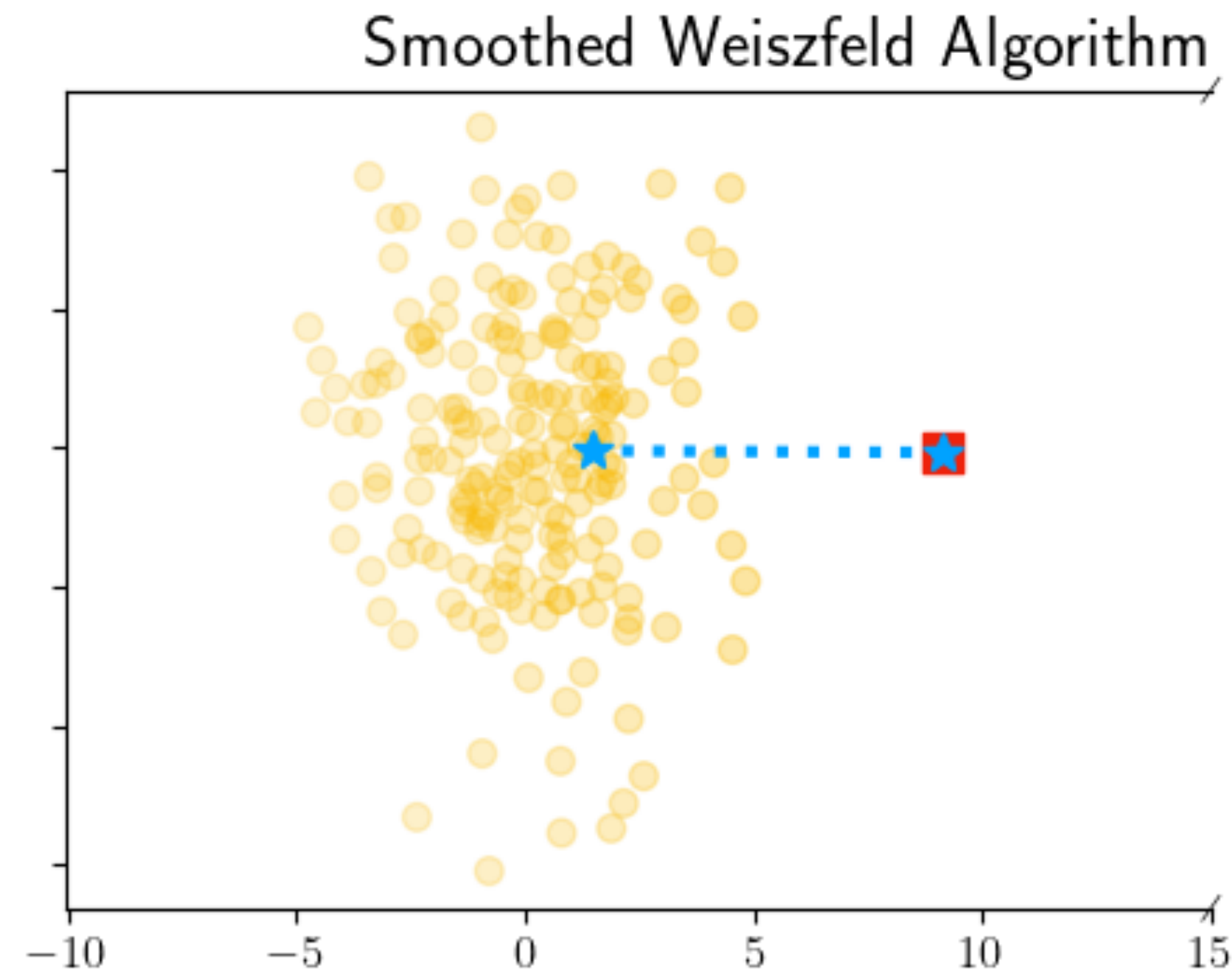One round of our algorithm $\implies$ 3-5 rounds of communication

Due to iterations of the smoothed Weiszfeld algorithm
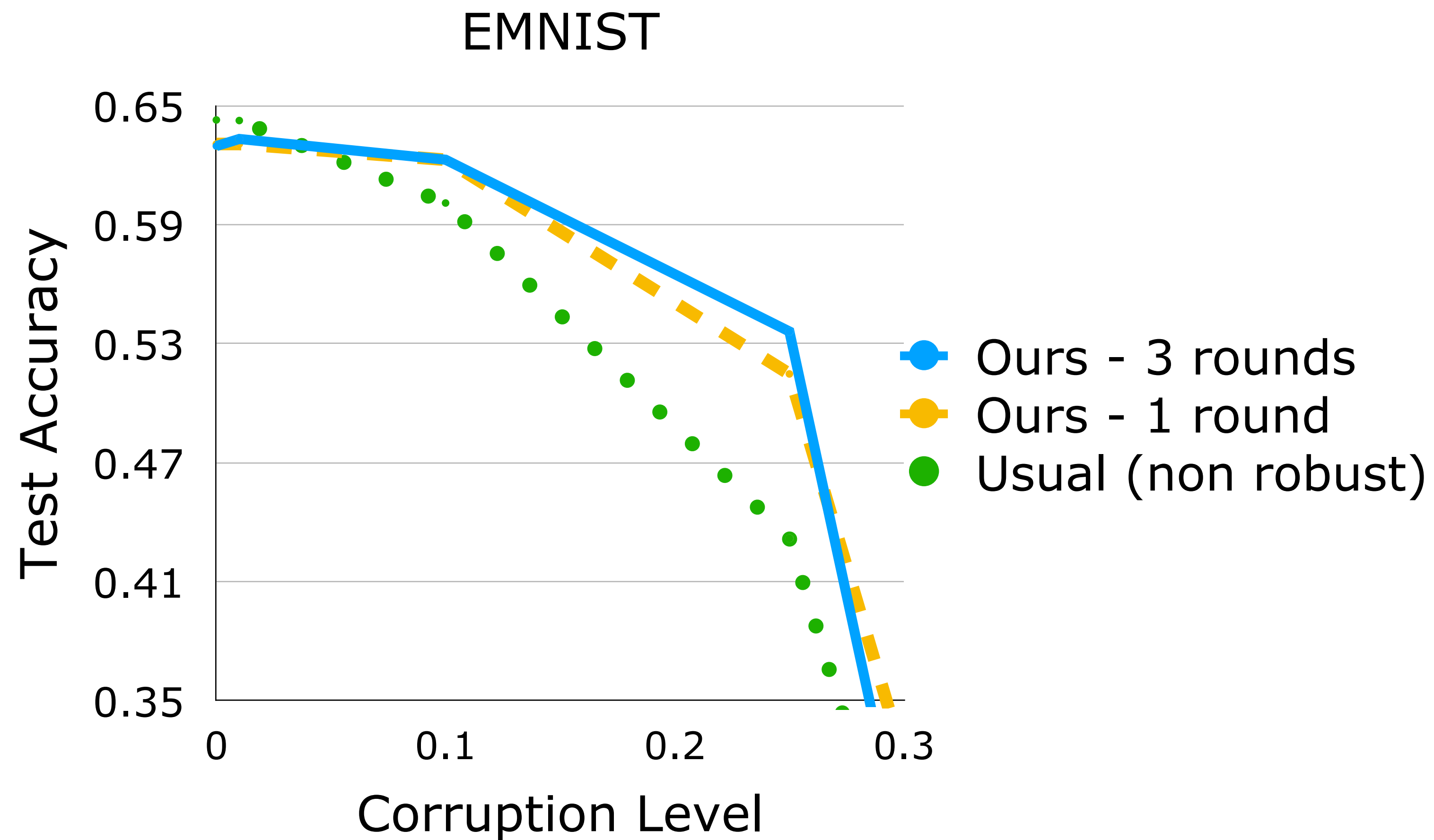
Does 1 round of communication improve robustness?

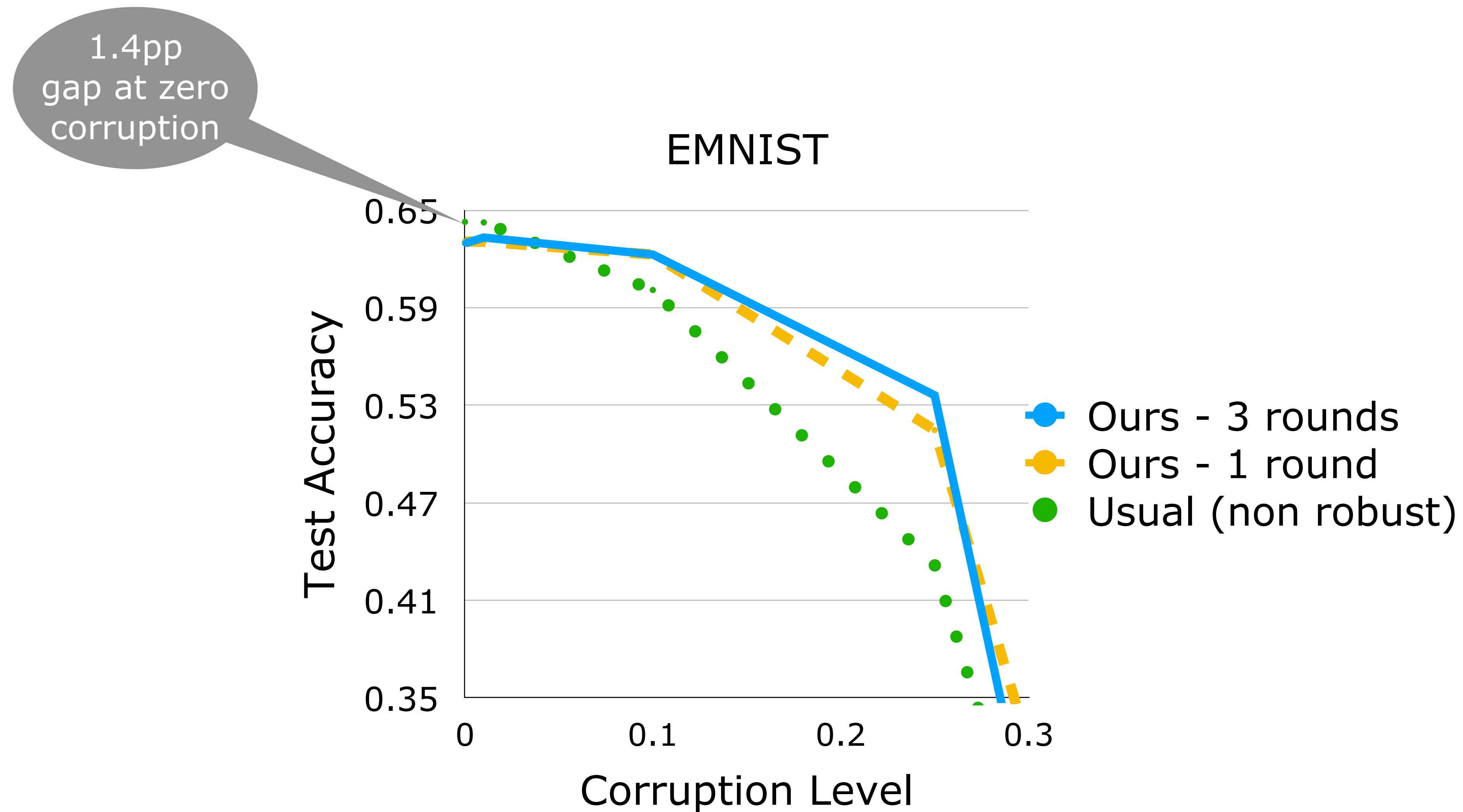$$\beta_i = \frac{1}{\max\{\|w_i\|_2, \nu\}}$$

$$z = \frac{\sum_i \beta_i w_i}{\sum_i \beta_i}$$

# *1 communication round* already improves robustness



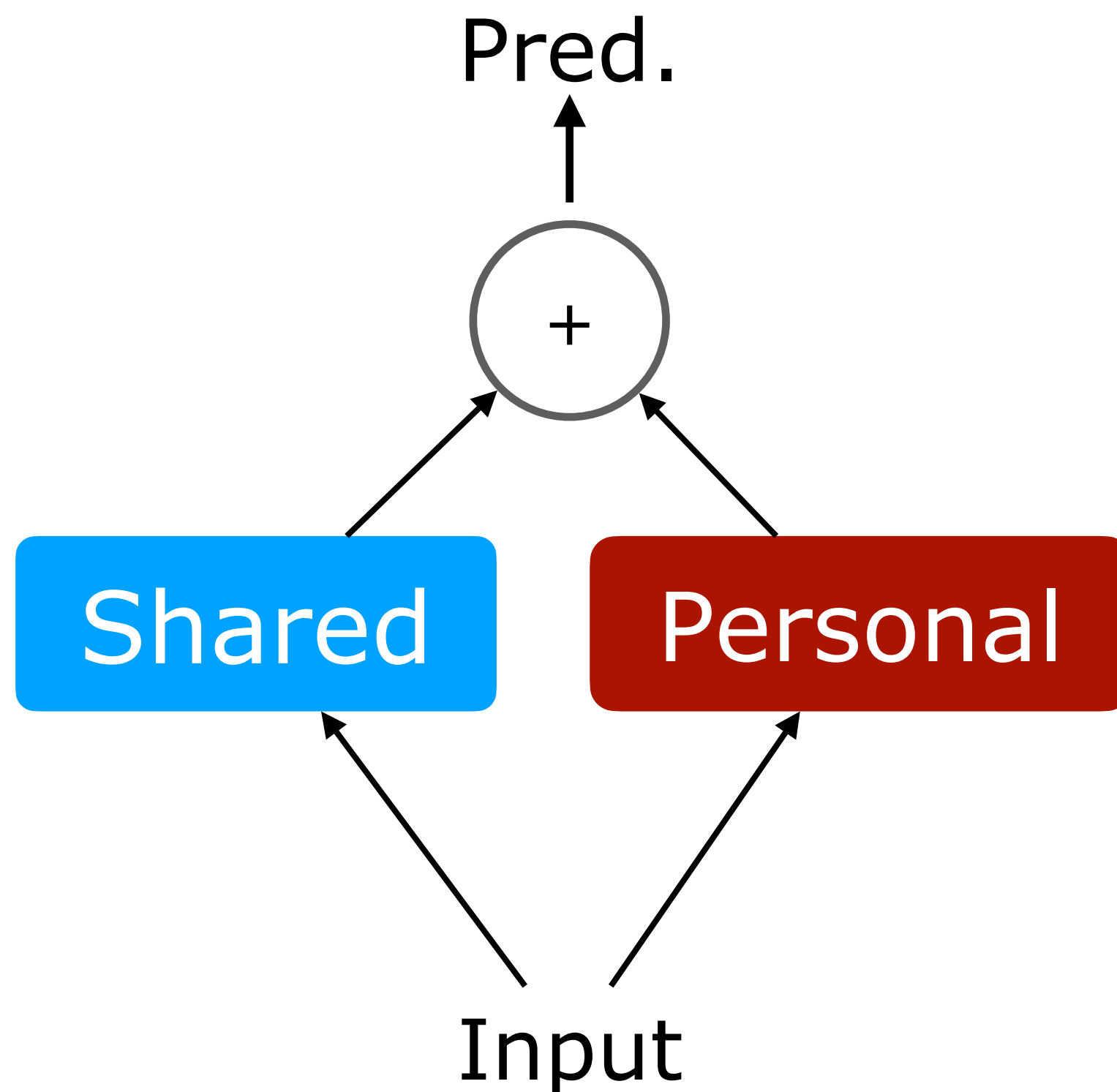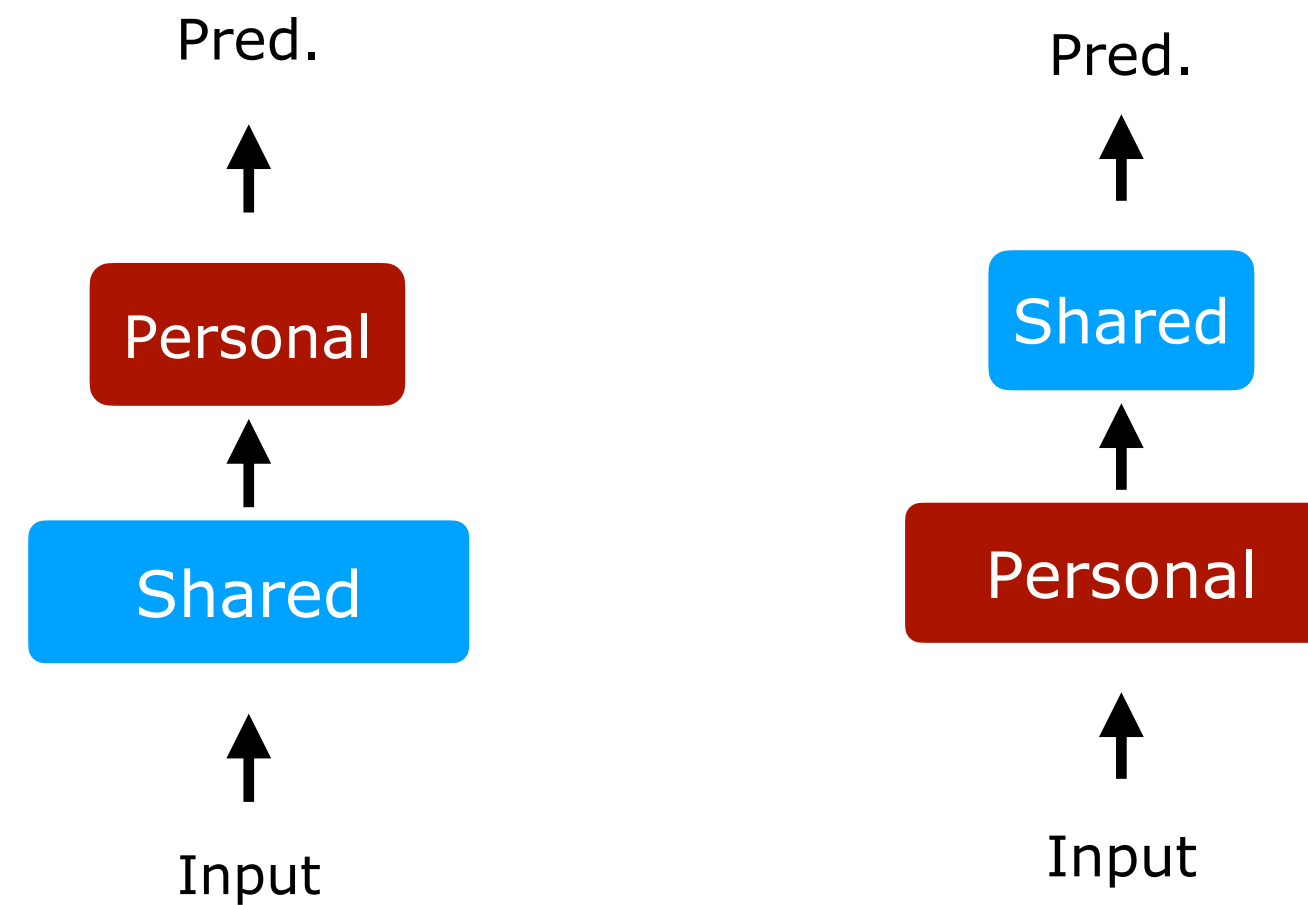EMNIST

# How do we get rid of this gap?

# Model personalization

The model has a global component
and a per-client component

Pred.

$+$

Shared    Personal

Input

Shared Params $u$

$+$    Personal Params $v_i$

$=$    Full model $w_i = (u, v_i)$

# Personalization Architectures



Pred.

Personal

Shared

Input

Arivazhagan et al. (2019)
Collins et al. (2021)

Pred.

Shared

Personal

Input

Liang et al. (2019)

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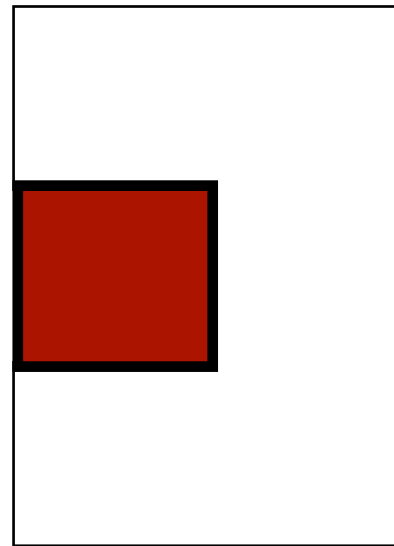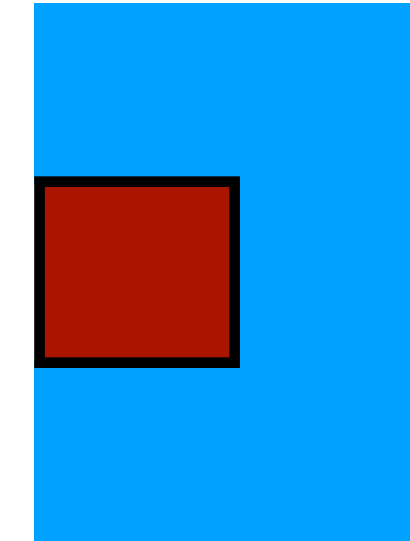
# Optimization

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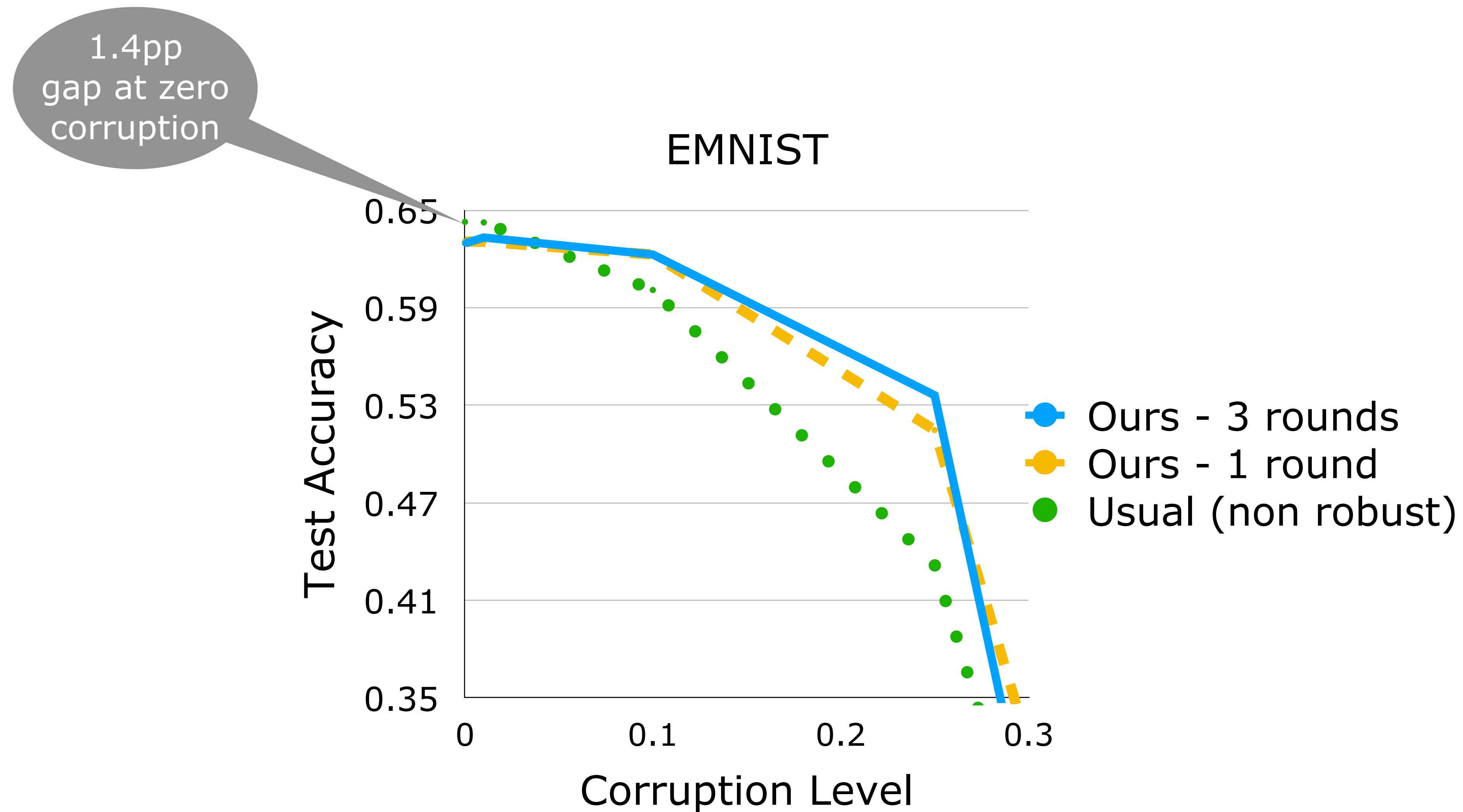


$+$ 　　　 $=$

***Shared part of the model*** *is updated with robust aggregation*　　　***Personal part of the model*** *stays with the client*

# Does personalization get rid of this gap?



1.4pp gap at zero corruption

EMNIST

Test Accuracy

Corruption Level

- Ours - 3 rounds
- Ours - 1 round
- Usual (non robust)

# Yes, we can improve robust aggregation with personalization!

# *In the literature*:

# Robust Federated Aggregation (RFA)

# RFA is certifiably more robust to backdoor attacks



[Xie et al. (ICML 2021)]
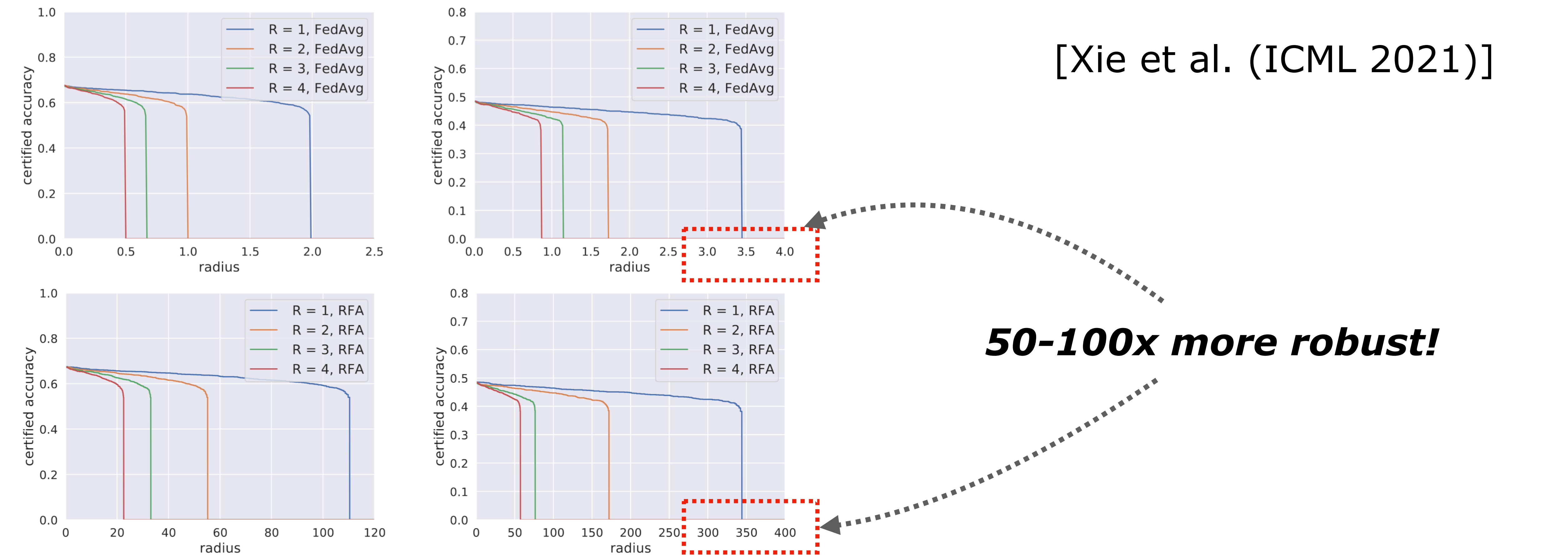
***50-100x more robust!***

*Figure 6.* Certified accuracy on MNIST (left) and EMNIST (right) with different $R$ when FL is trained under the robust aggregation RFA (Pillutla et al., 2019).

Xie, Chen, Chen, Li. **CRFL: Certifiably Robust Federated Learning against Backdoor Attacks.** *ICML 2021*.

64

# RFA is asymptotically strategy-proof

**Strategy-proof**: Can a device lie to bring the aggregate to a desired point?

With a large number of independent devices, RFA is approximately strategy-proof

## On the Strategyproofness of the Geometric Median

**El-Mahdi El-Mhamdi**
Calicarpa, École Polytechnique
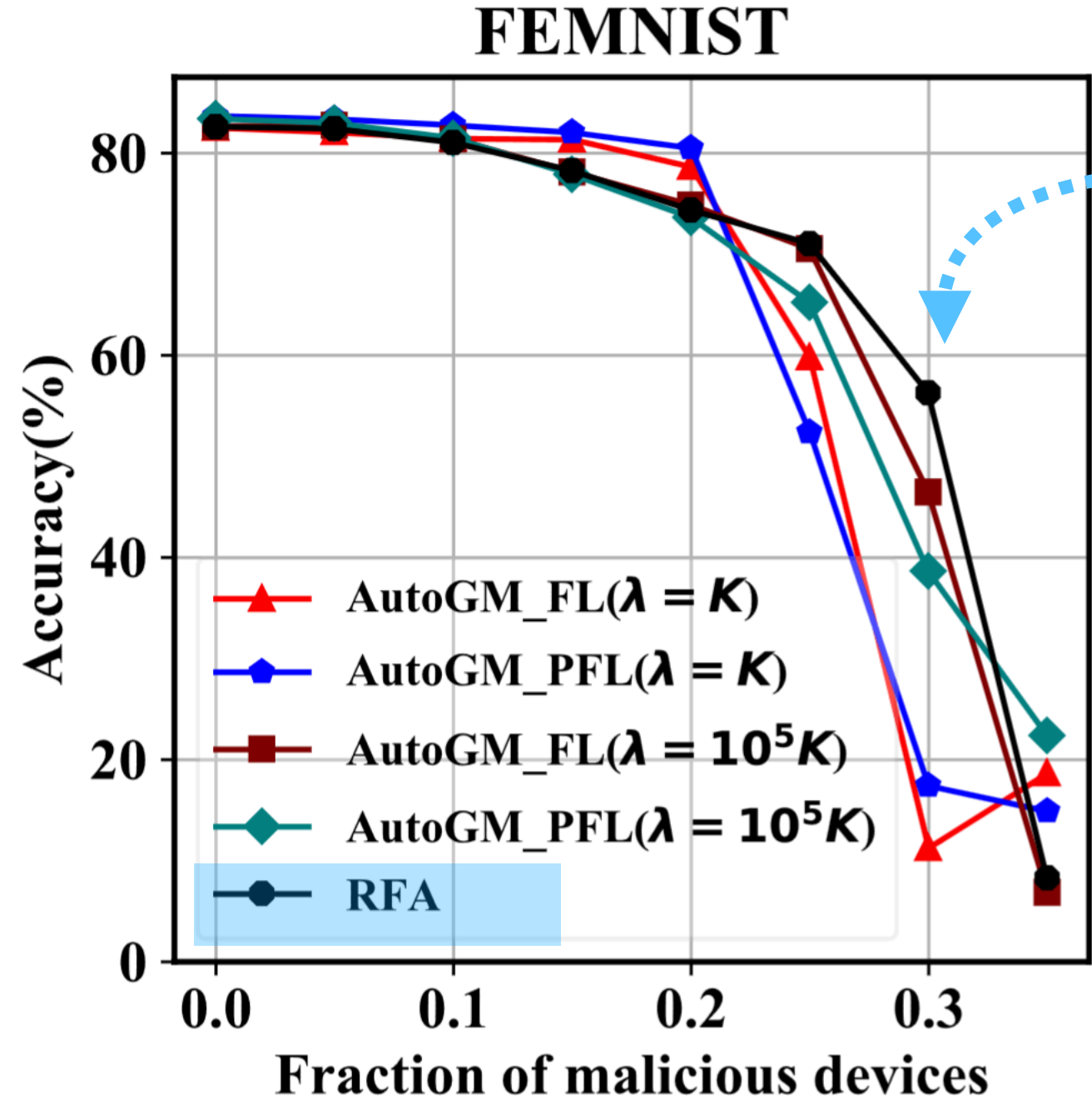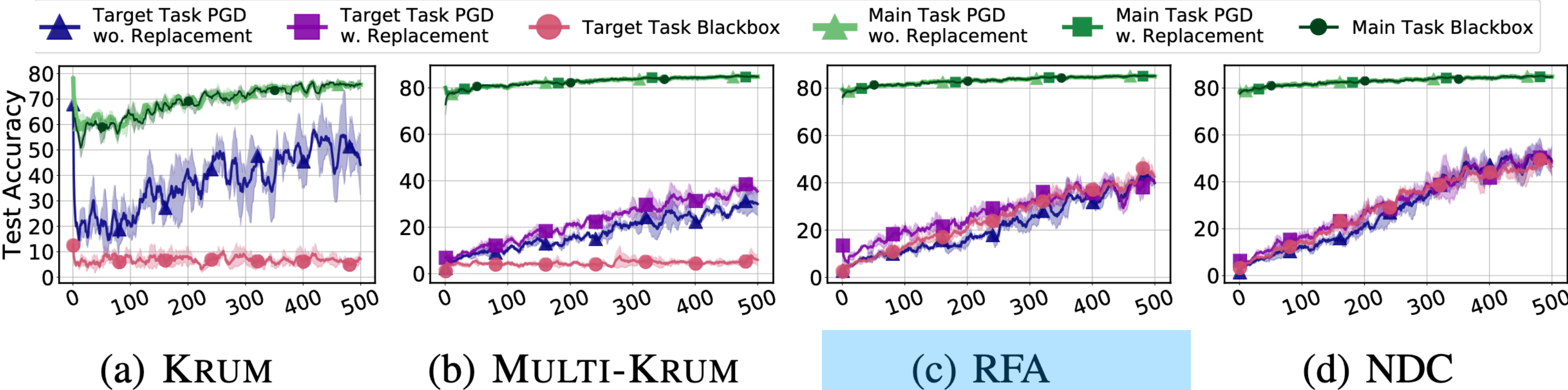
**Sadegh Farhadkhani**[*]
EPFL

**Rachid Guerraoui**
EPFL

**Lê-Nguyên Hoang**[*]
Calicarpa, Tournesol

**[AISTATS 2023]**

# RFA is a strong baseline

Wang et al.
(NeurIPS 2020)



Target Task PGD wo. Replacement | Target Task PGD w. Replacement | Target Task Blackbox | Main Task PGD wo. Replacement | Main Task PGD w. Replacement | Main Task Blackbox

(a) KRUM    (b) MULTI-KRUM    (c) RFA    (d) NDC



FEMNIST

AutoGM_FL($\lambda = K$)
AutoGM_PFL($\lambda = K$)
AutoGM_FL($\lambda = 10^5 K$)
AutoGM_PFL($\lambda = 10^5 K$)
RFA

Li et al. (IEEE Trans. Industrial Informatics 2023)

See also Sejwalkar et al. (IEEE Security & Privacy 2022), Jin & Li (Medical Image Analysis 2023), Li et al. (IEEE Trans. Big Data 2023), …

66

# Algorithmic advances based on RFA

Park et al. (NeurIPS 2021): RFA + Entropy-based reweighting

Karimireddy et al. (ICLR 2022): RFA + Bucketing

Li et al. (IEEE Trans. Ind. Inform. 2023): RFA + adaptive weighting

Allouah et al. (AISTATS 2023): RFA + nearest neighbhors

$\vdots$

# Fast and differentiable geometric median

```python
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,)
```

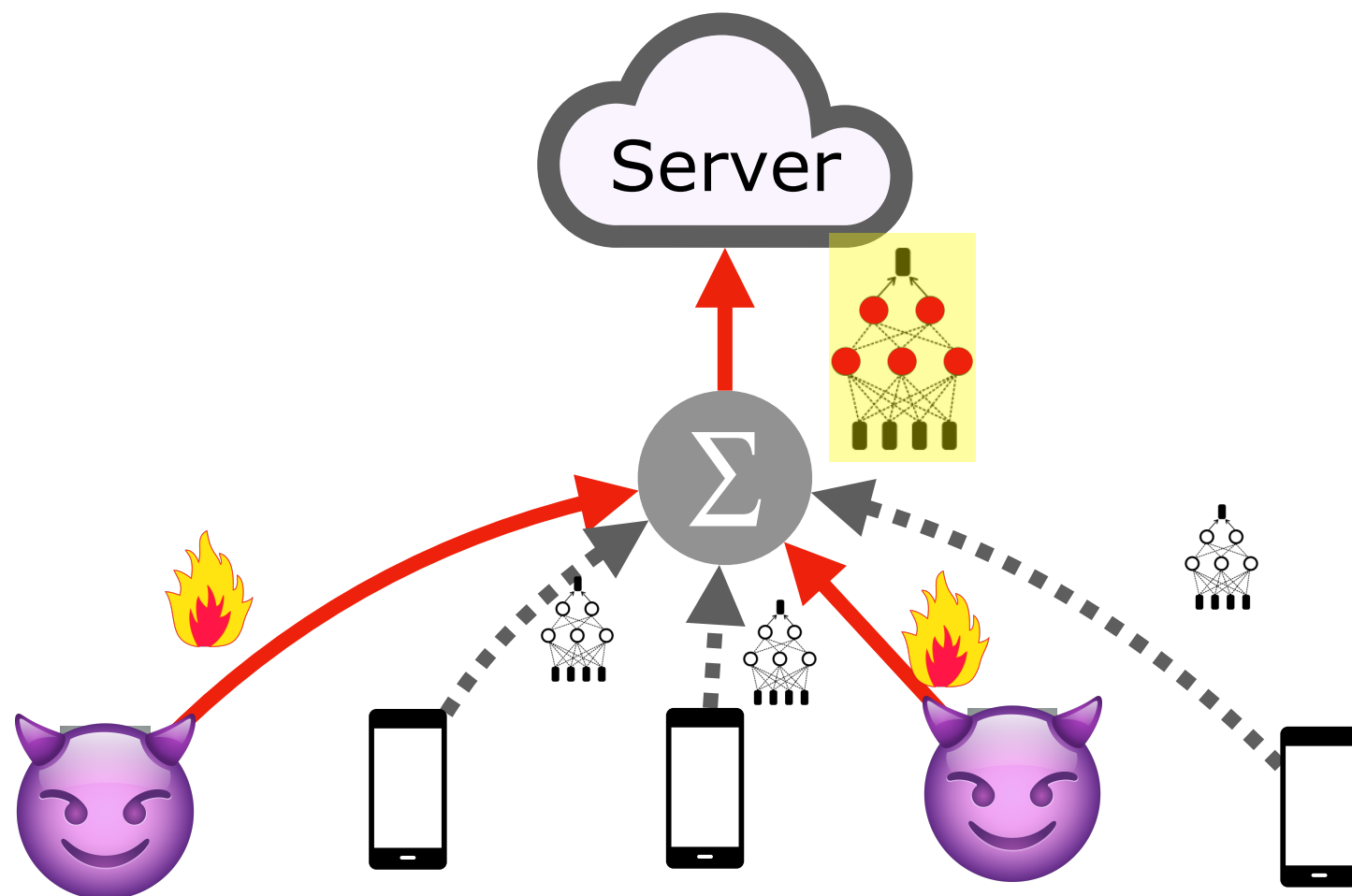Install: **pip install geom-median**

Documentation: github.com/krishnap25/geom-median

GitHub Link
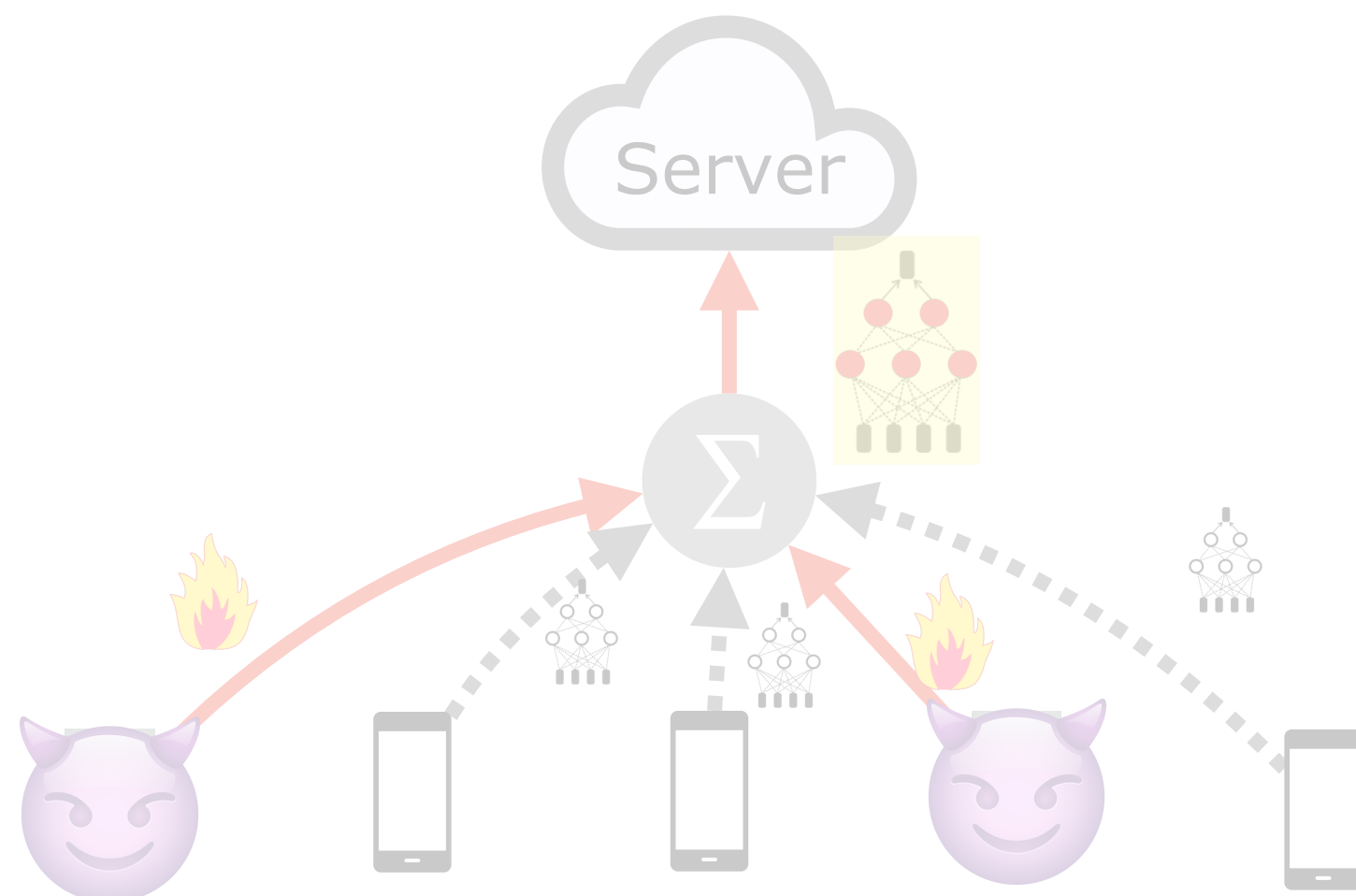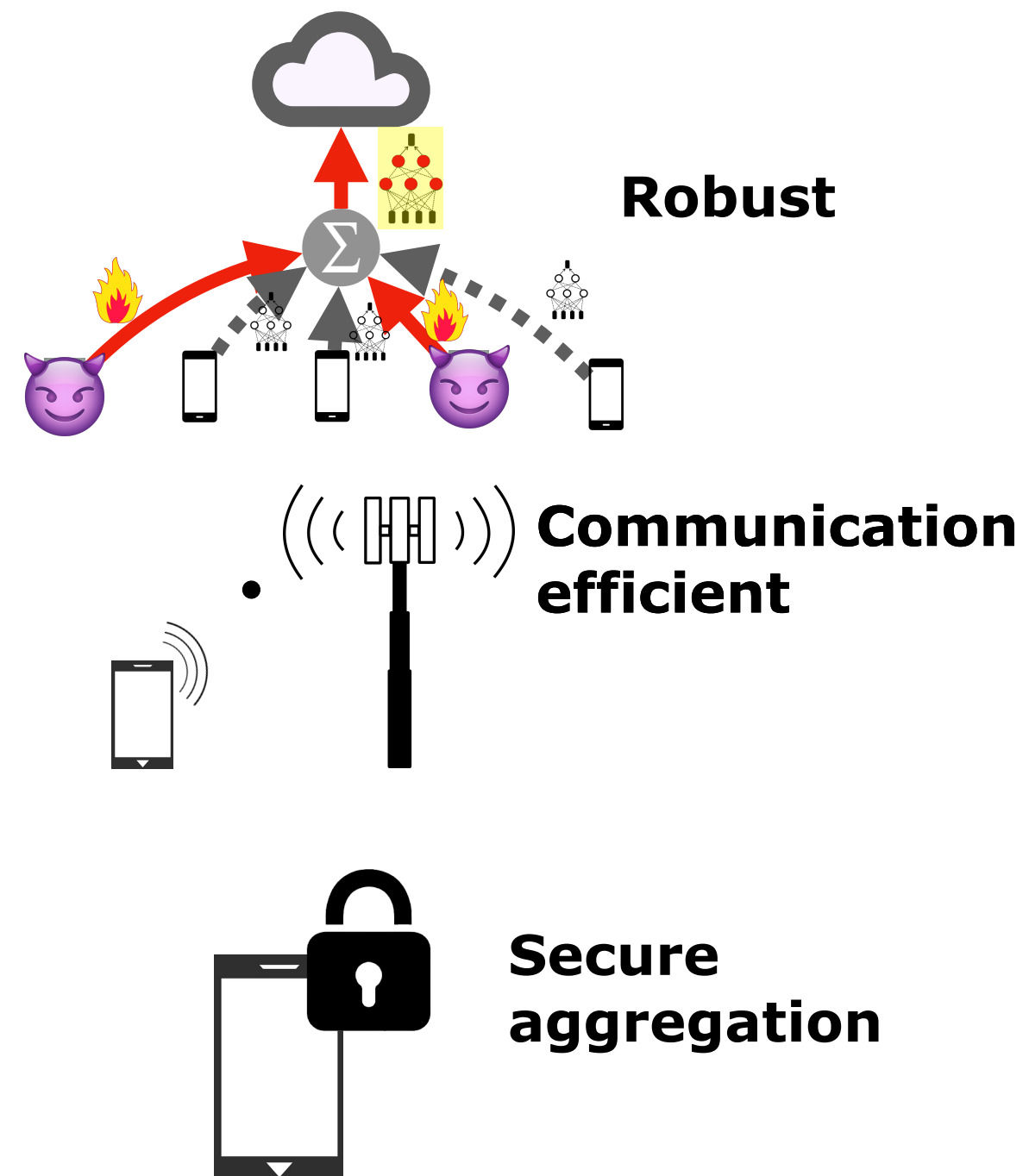
# Summary

Federated learning is
**not robust** to
poisoned updates

# Summary

Federated learning is **not robust** to poisoned updates

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

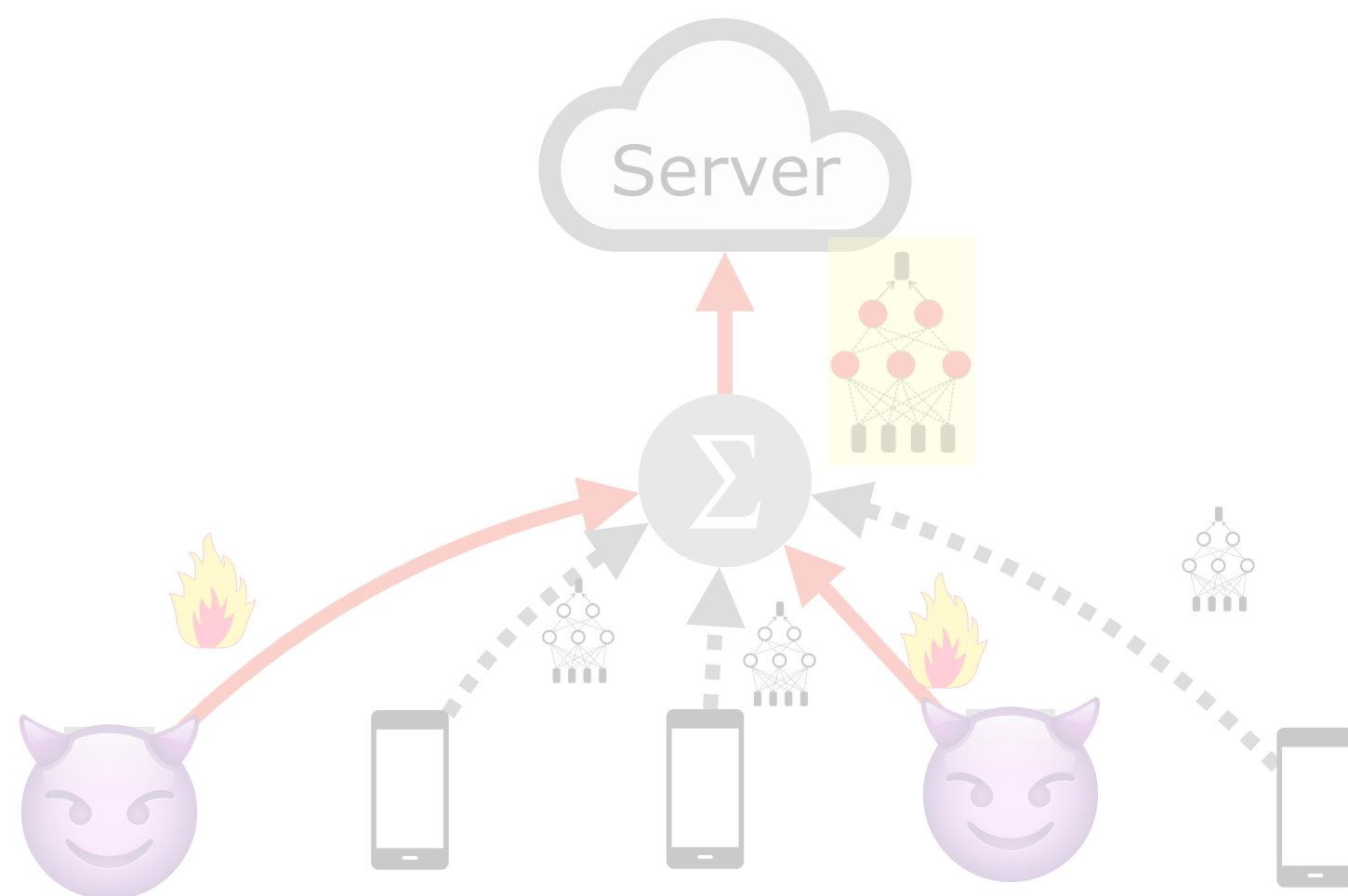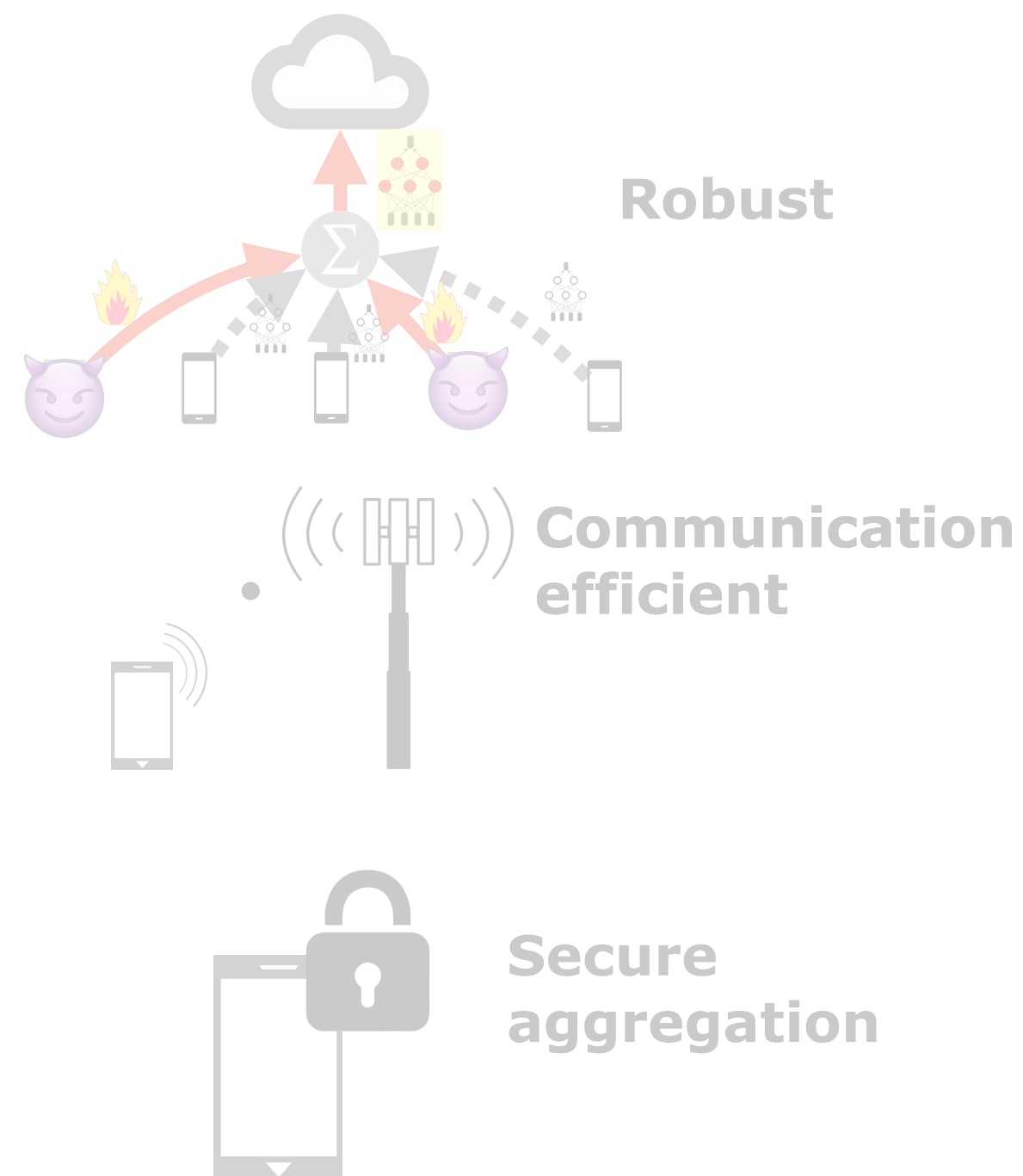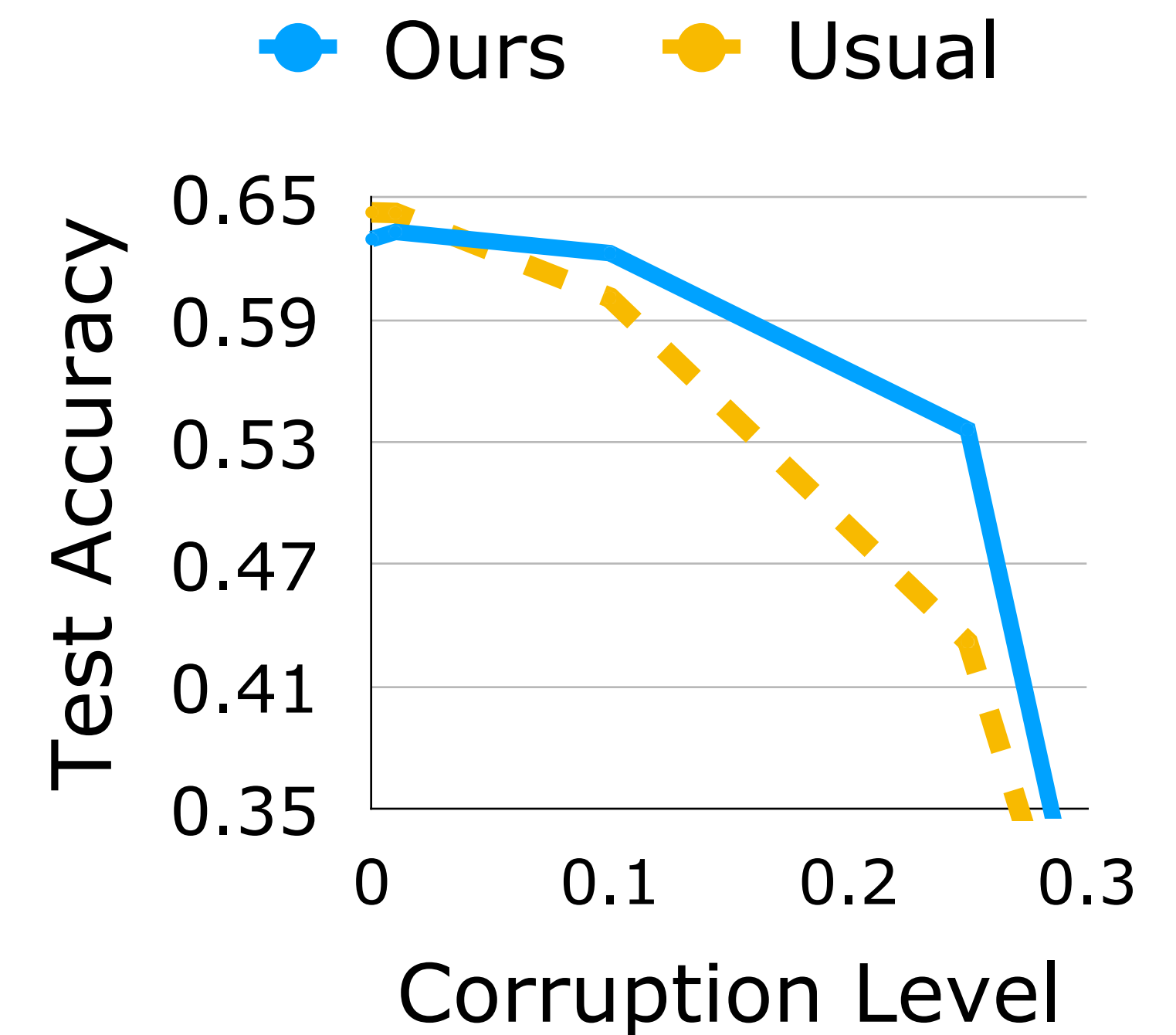**Robust**

**Communication efficient**

**Secure aggregation**

# Summary

Federated learning is *not robust* to poisoned updates

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

**Robust**

**Communication efficient**

**Secure aggregation**

Our approach gives **greater robustness**

Ours      Usual

Test Accuracy

0.65
0.59
0.53
0.47
0.41
0.35

0      0.1      0.2      0.3

Corruption Level

71

# Heterogeneity, fairness, equity with differential privacy in federated learning
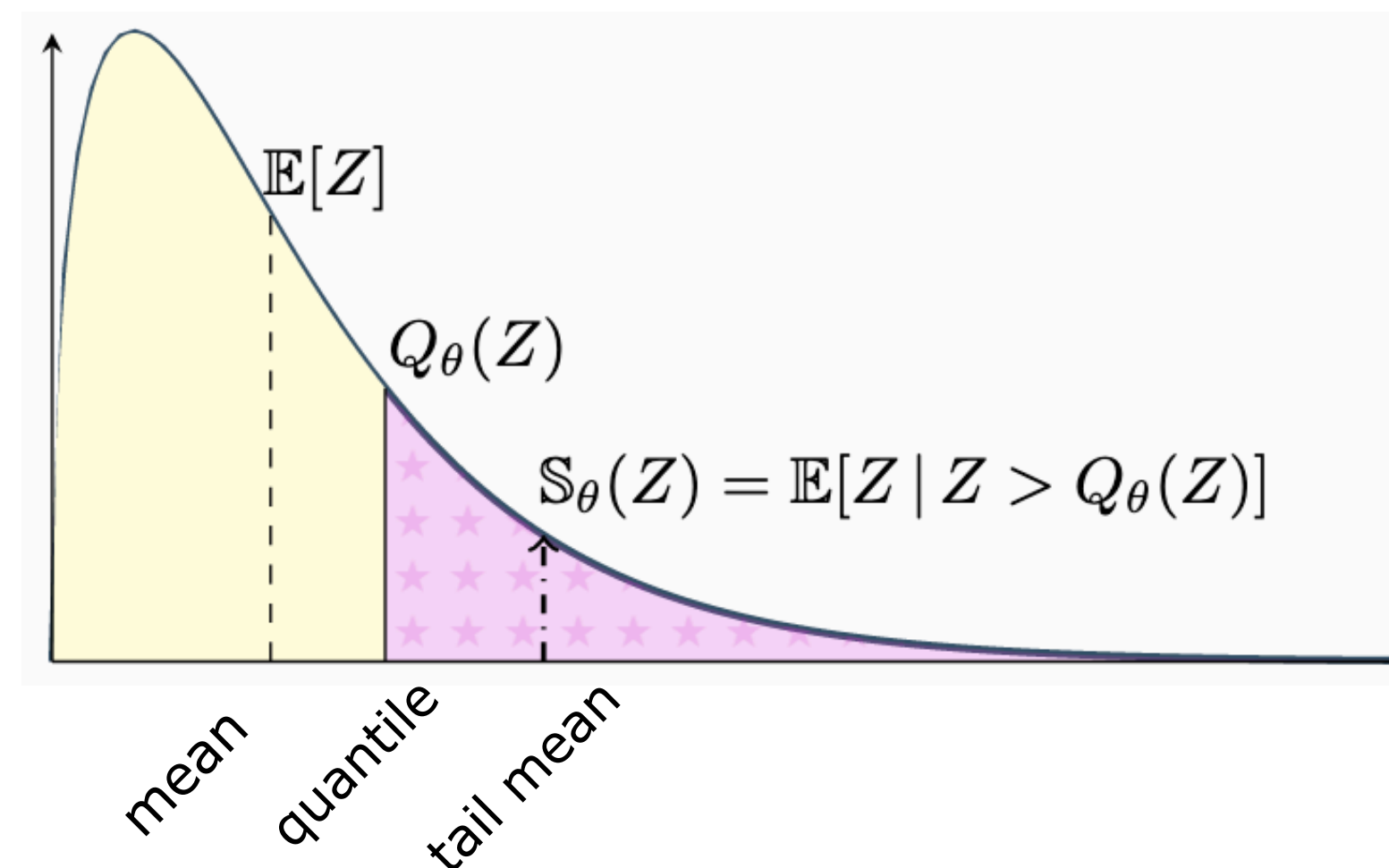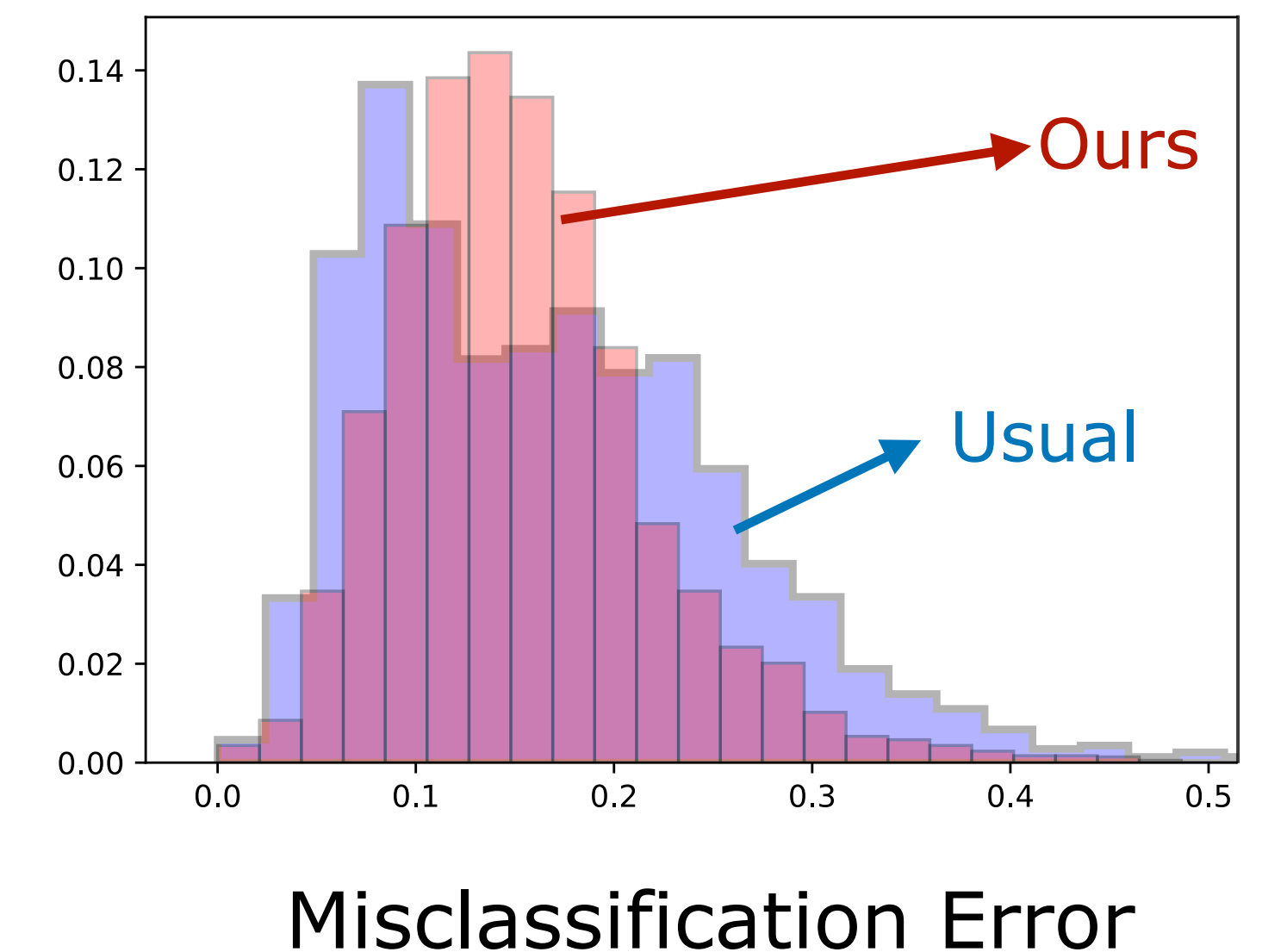
*Paper:*

Distribution shift $\implies$ large tail errors

Minimize the tail error directly

We reduce tail error + support differential privacy





$$\mathbb{E}[Z]$$

$$Q_\theta(Z)$$
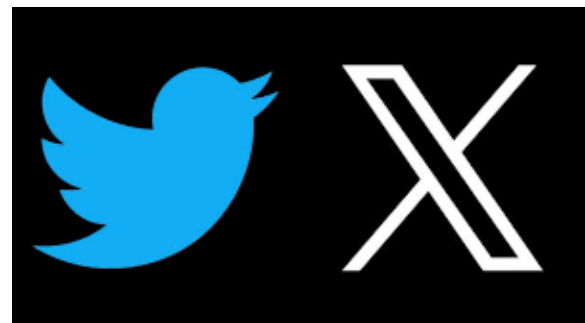
$$\mathbb{S}_\theta(Z) = \mathbb{E}[Z \mid Z > Q_\theta(Z)]$$

# Thank you!

Software

```
pip install geom-median
```

Code    https://github.com/krishnap25/tRFA

@KrishnaPillutla