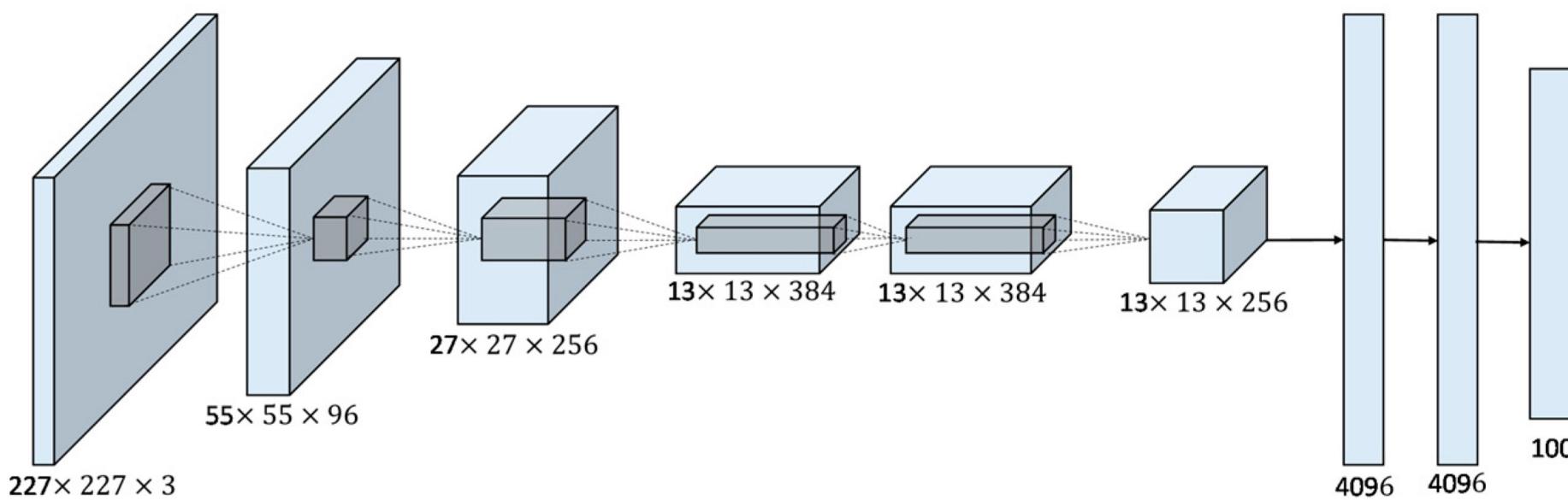


Towards Next-Generation ML/AI: Robustness, Optimization, Privacy

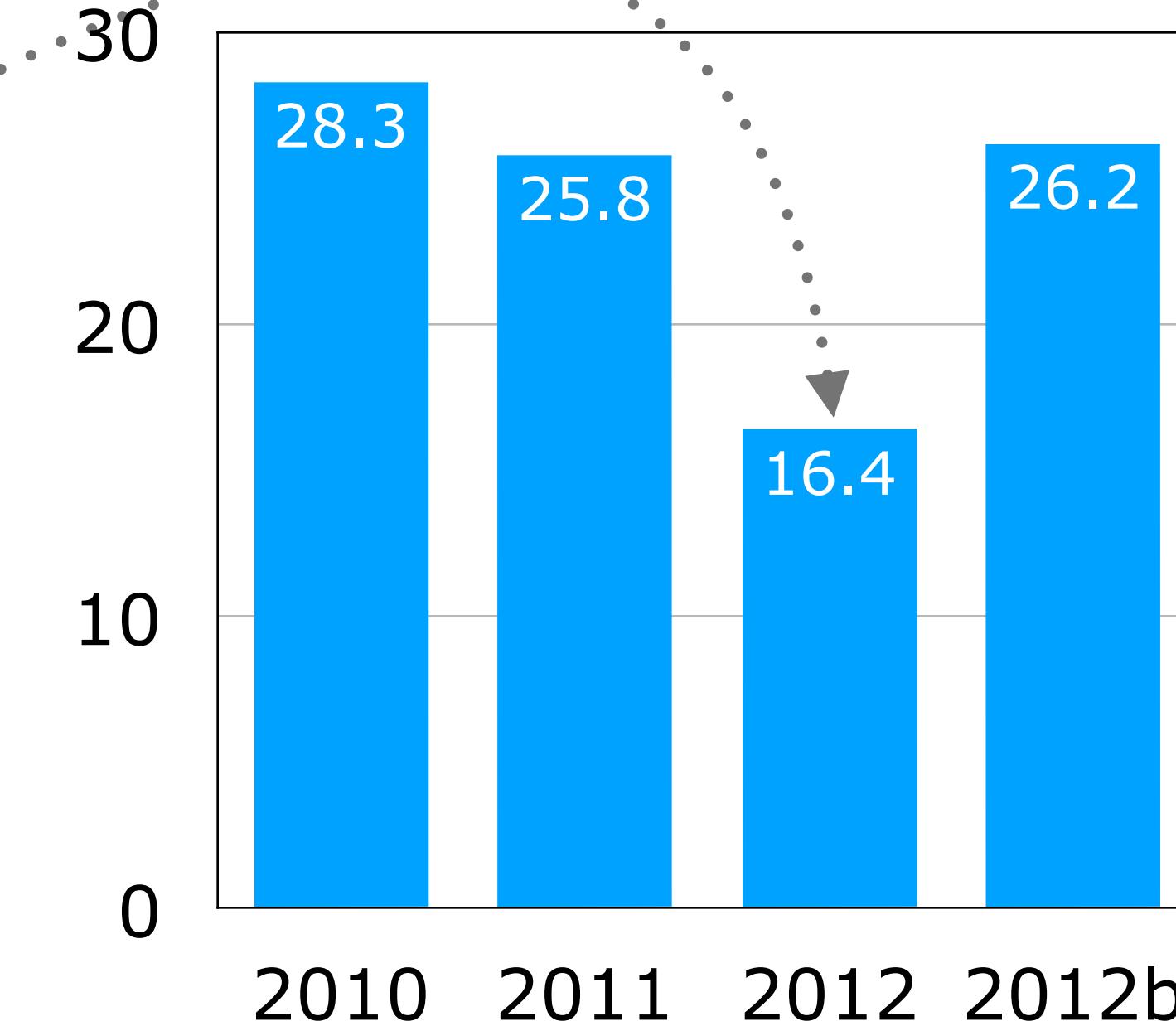
Krishna Pillutla

January 16th, 2023 @ IIT Madras

ML/AI have been revolutionized in the last 10 years



Top-5 Error %



[Krizhevsky, Sutskever, Hinton (NeurIPS 2012)]

2010

2012

2014

2016

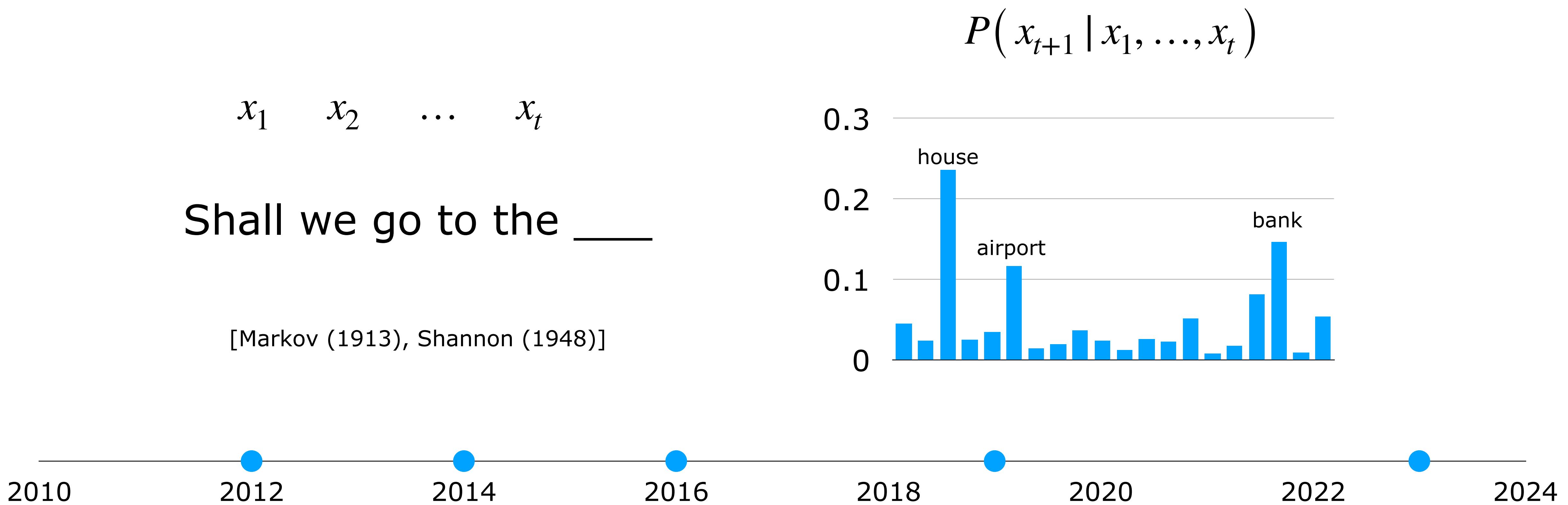
2018

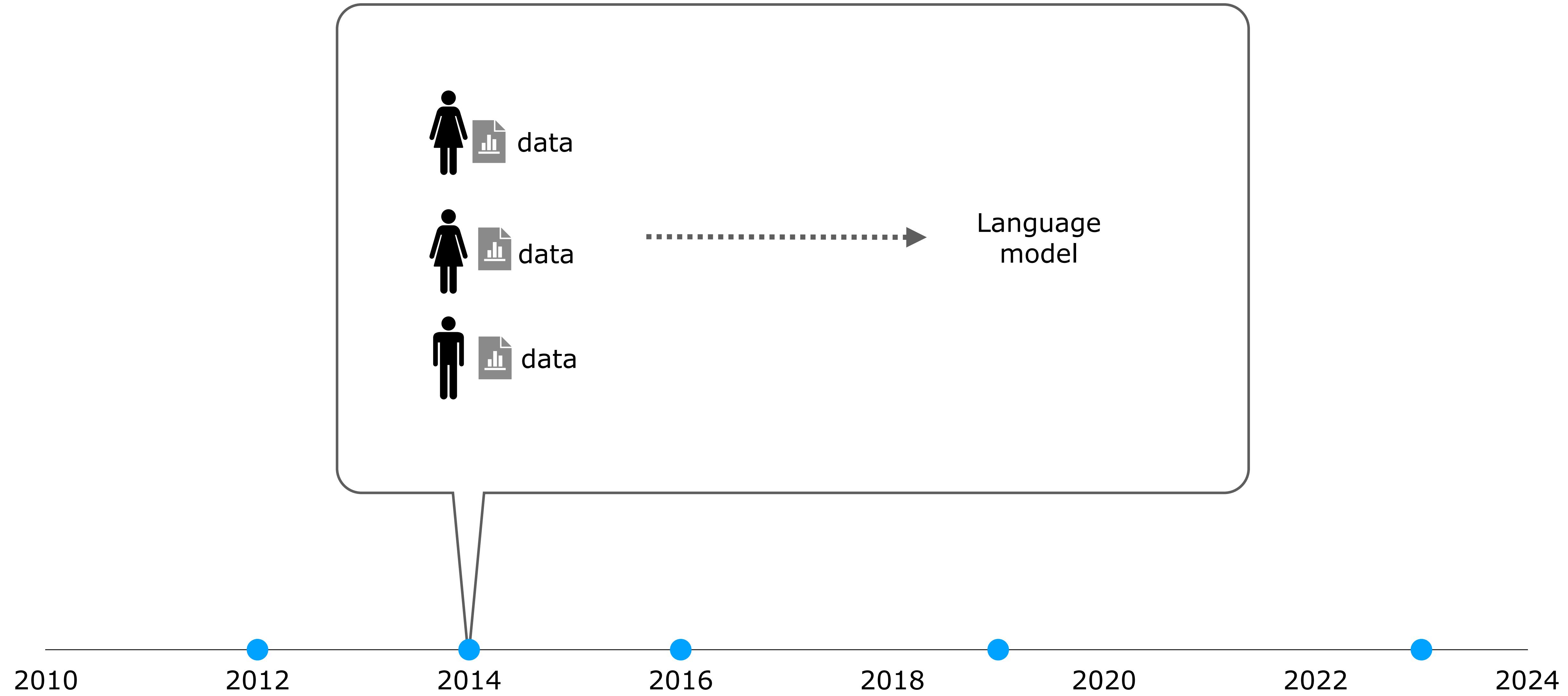
2020

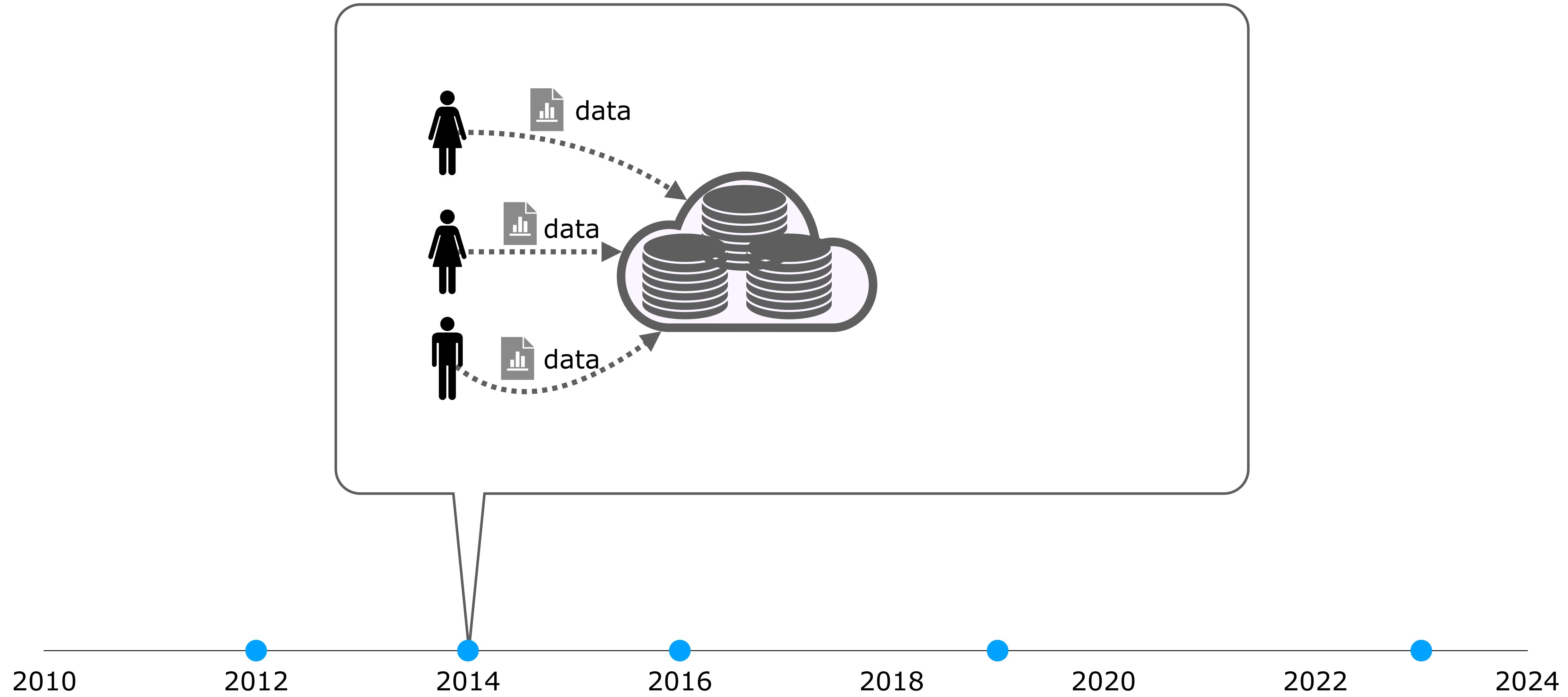
2022

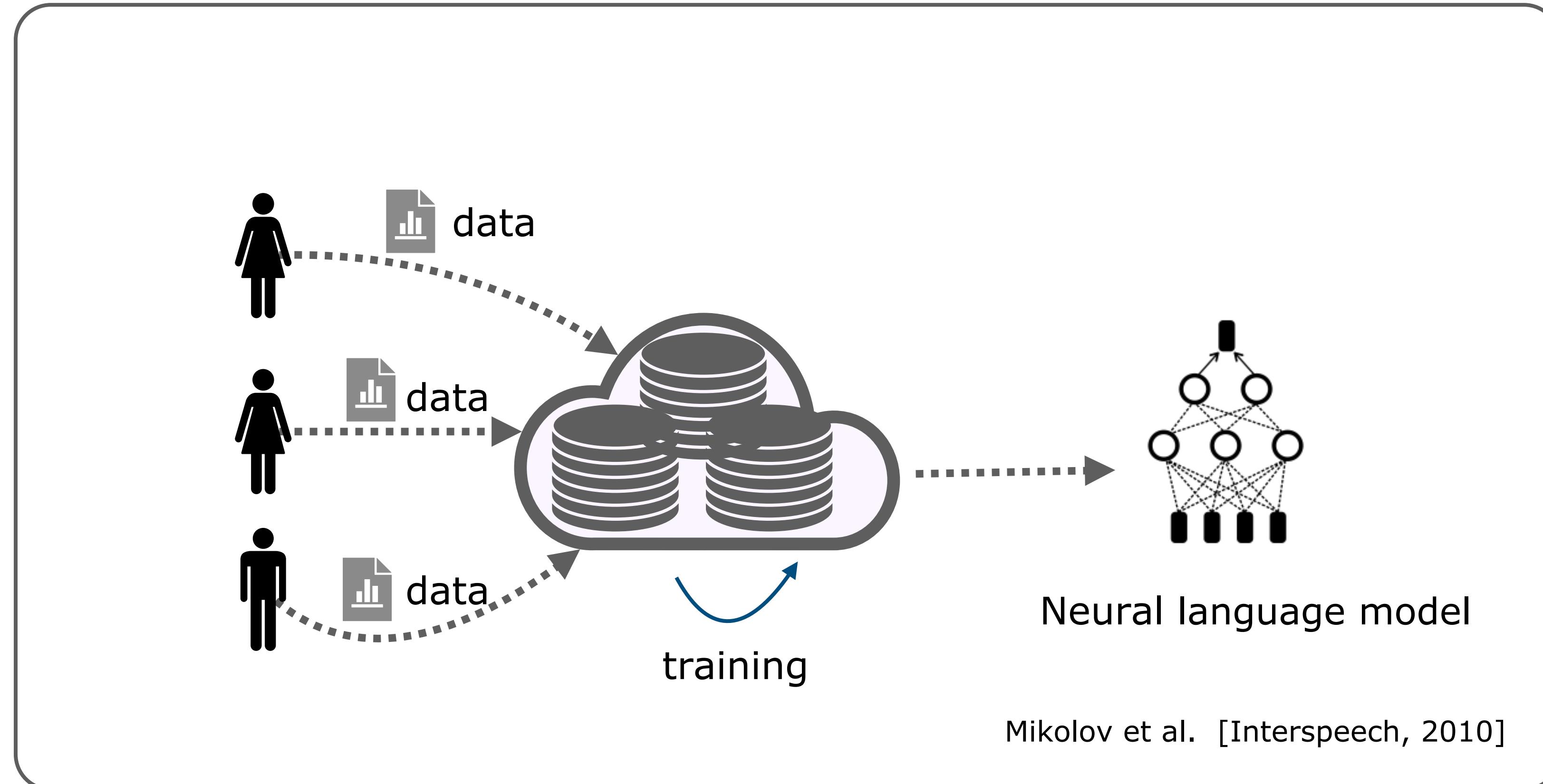
2024

Language modeling









2010

2012

2014

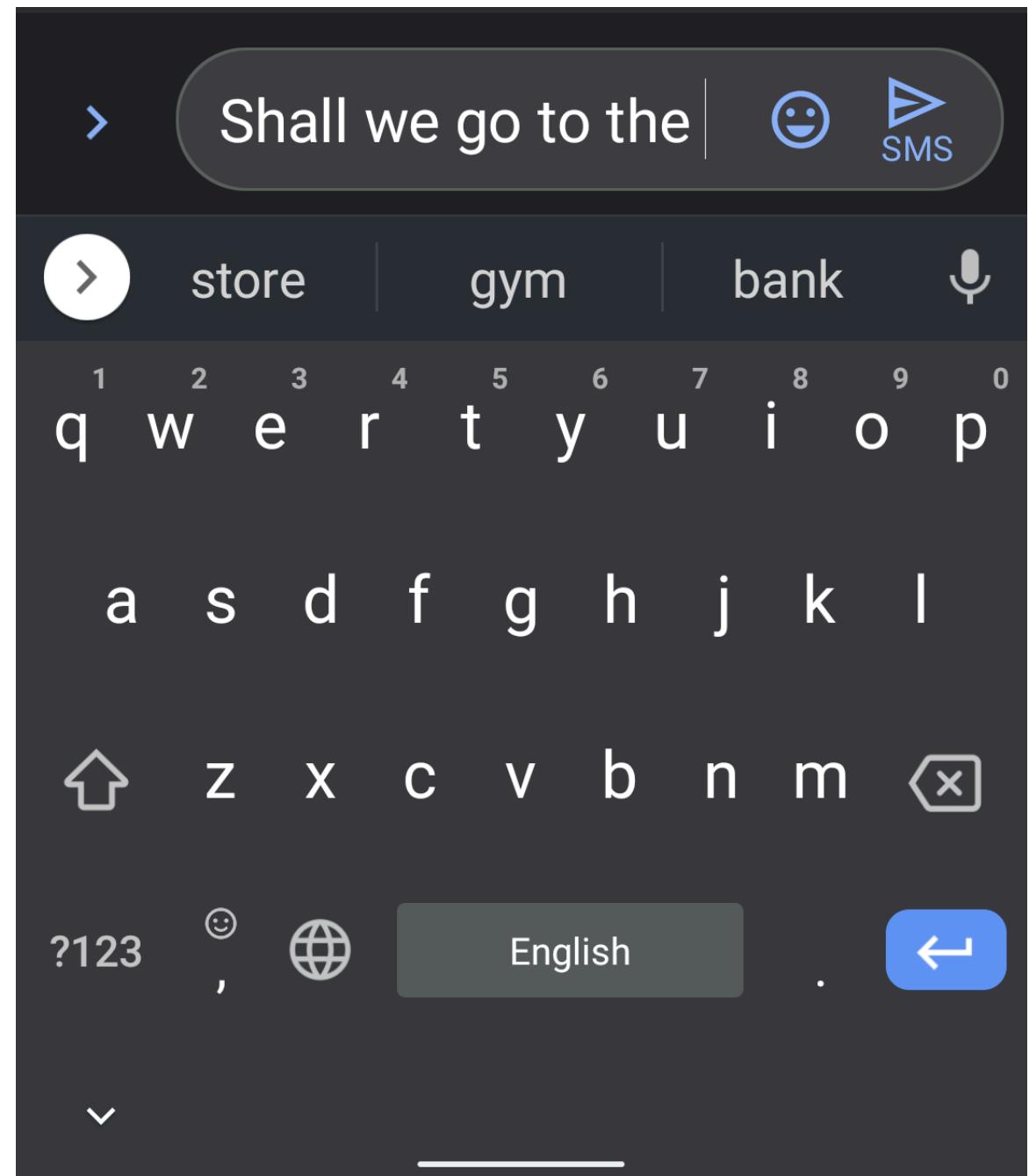
2016

2018

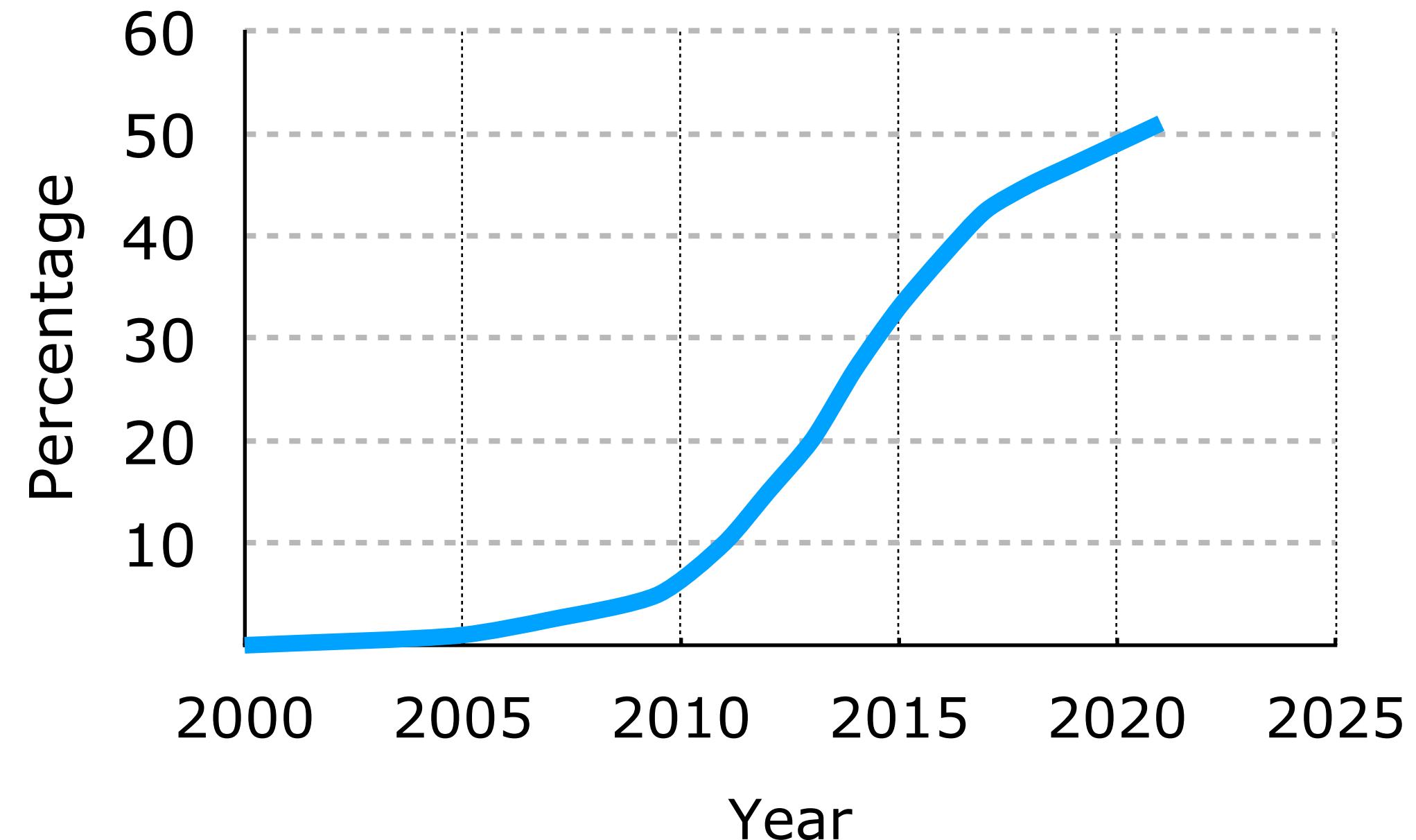
2020

2022

2024



Percentage of world population with a smartphone



Data Credit: Business Wire

2010

2012

2014

2016

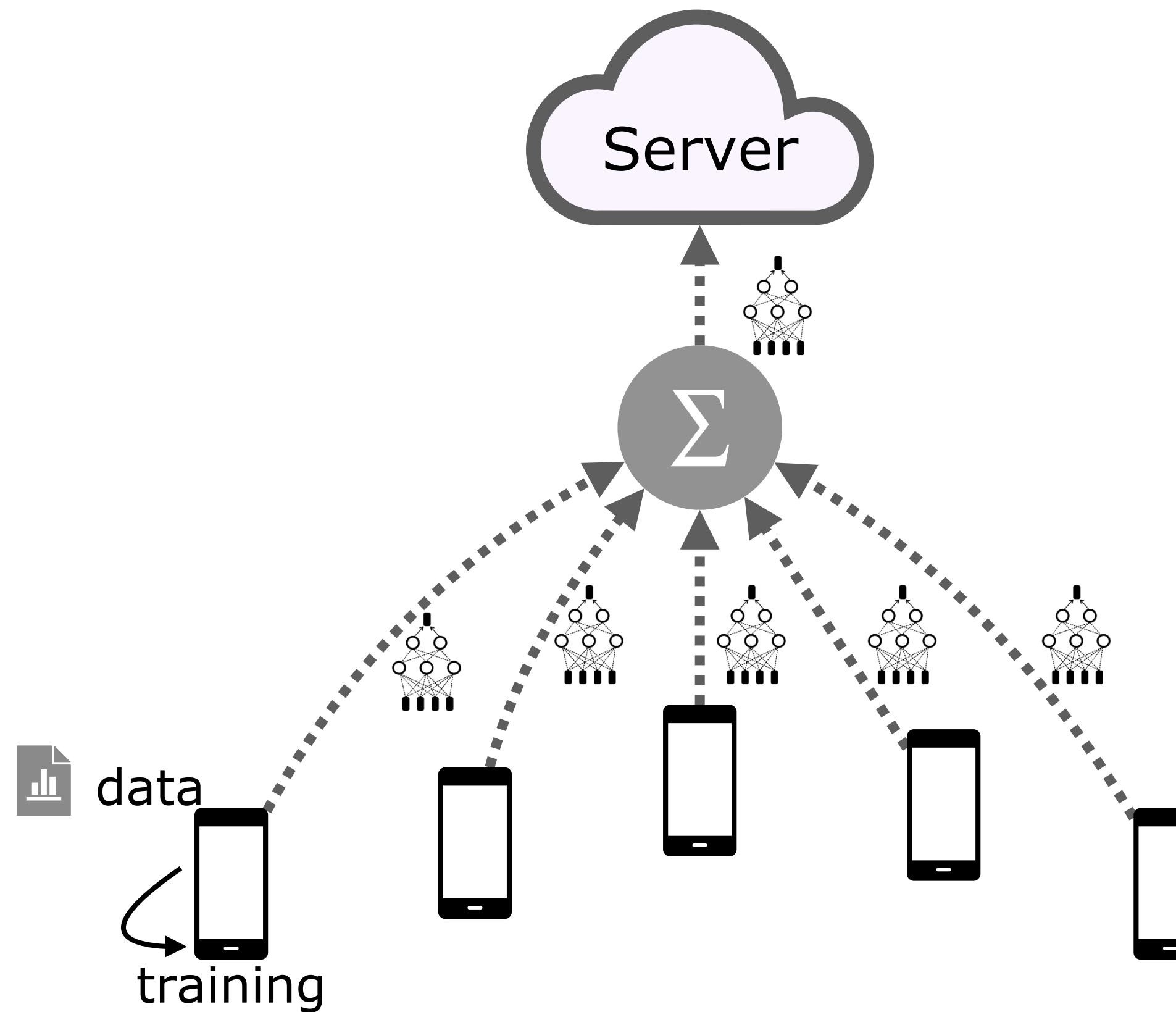
2018

2020

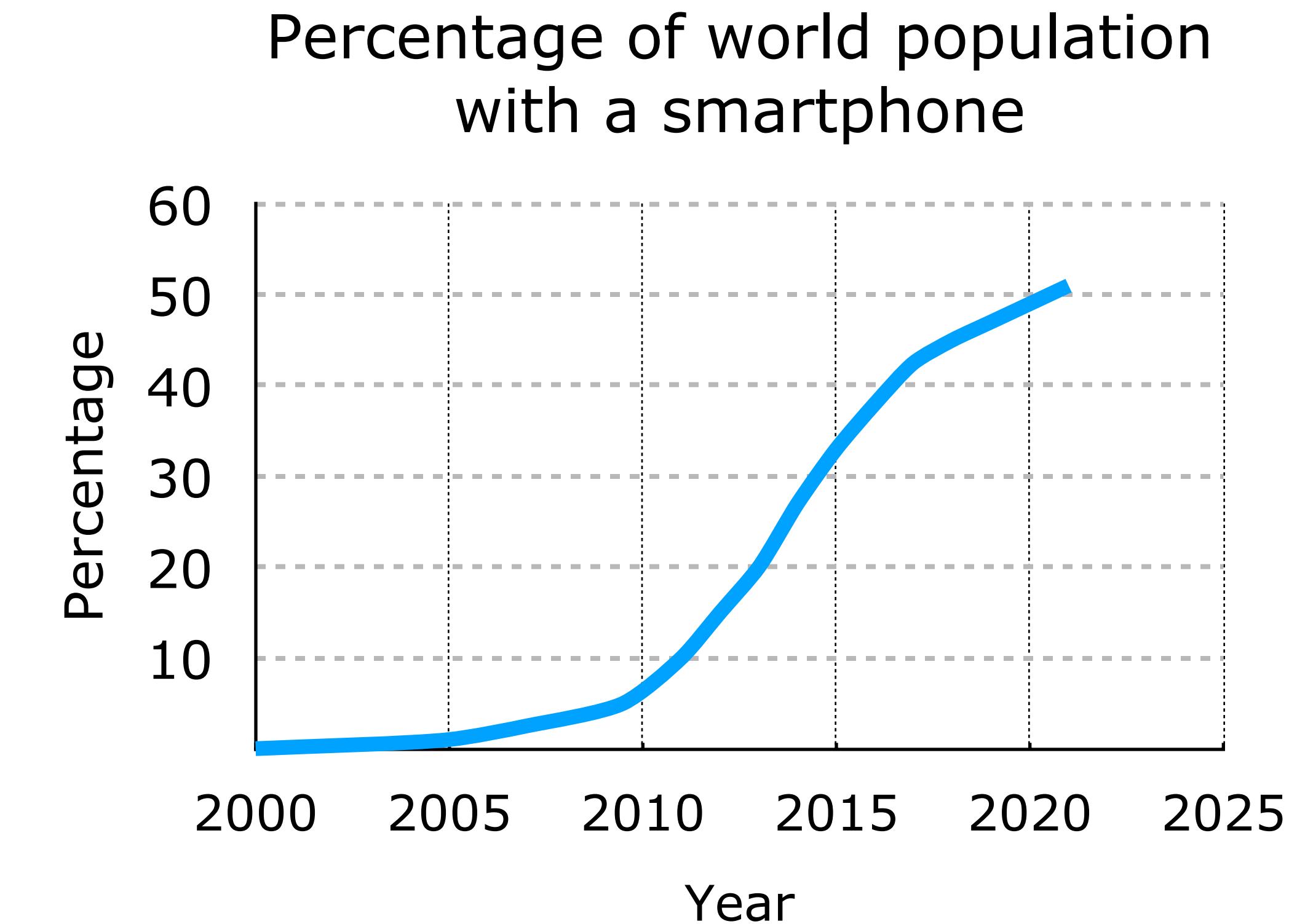
2022

2024

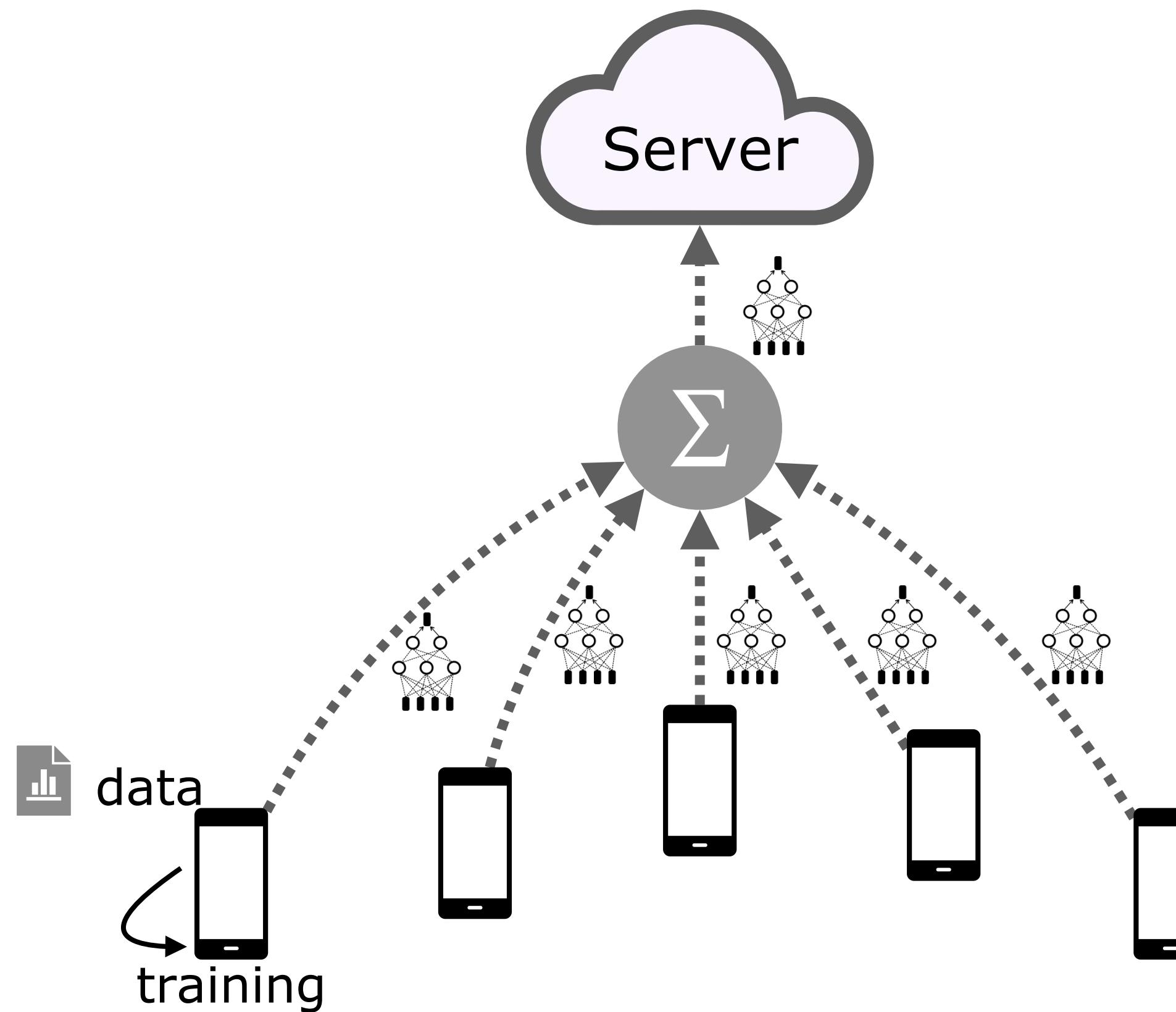
Federated learning: modern distributed learning



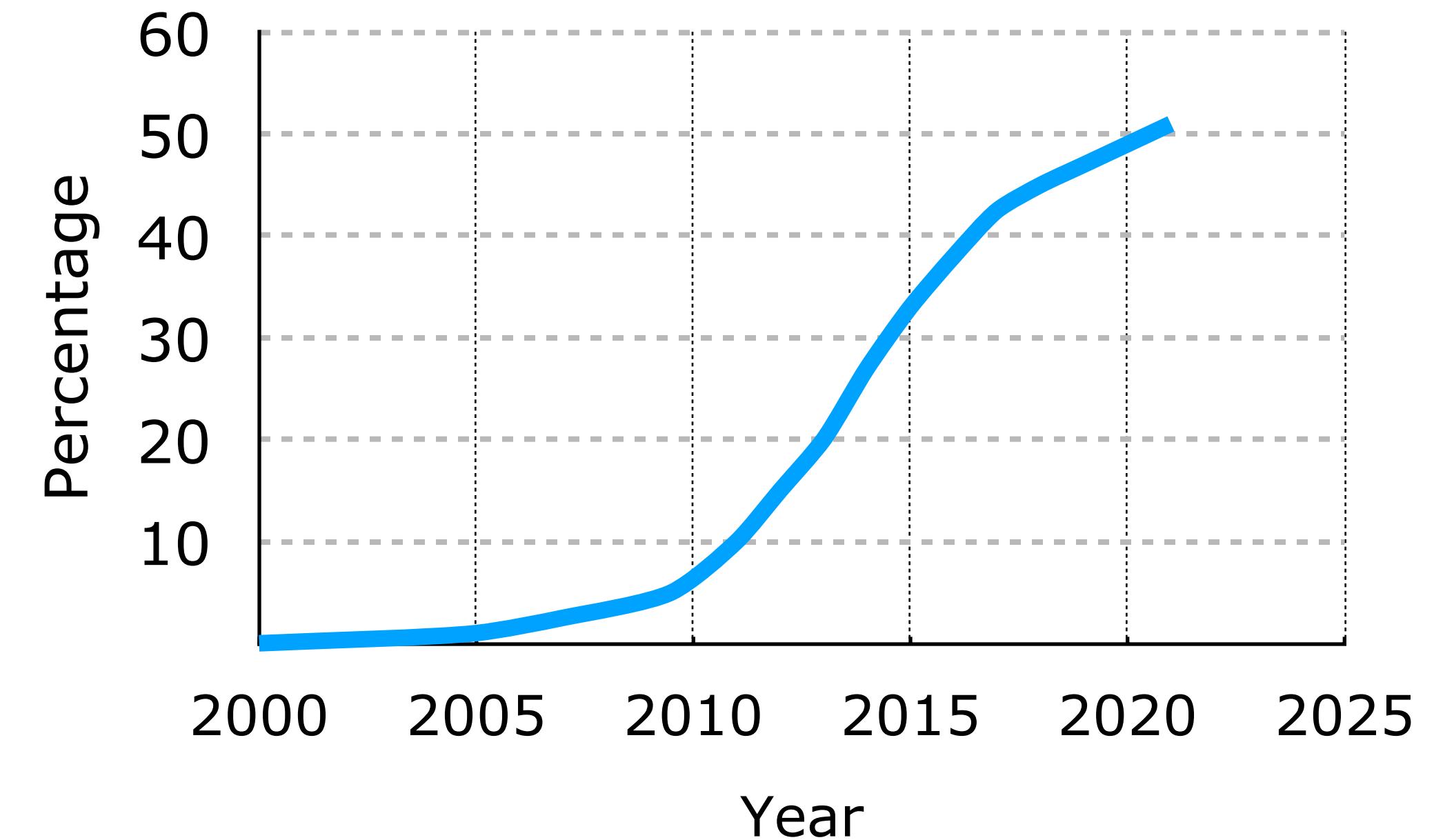
[McMahan et al. (AISTATS 2017)]



Federated learning: modern distributed learning



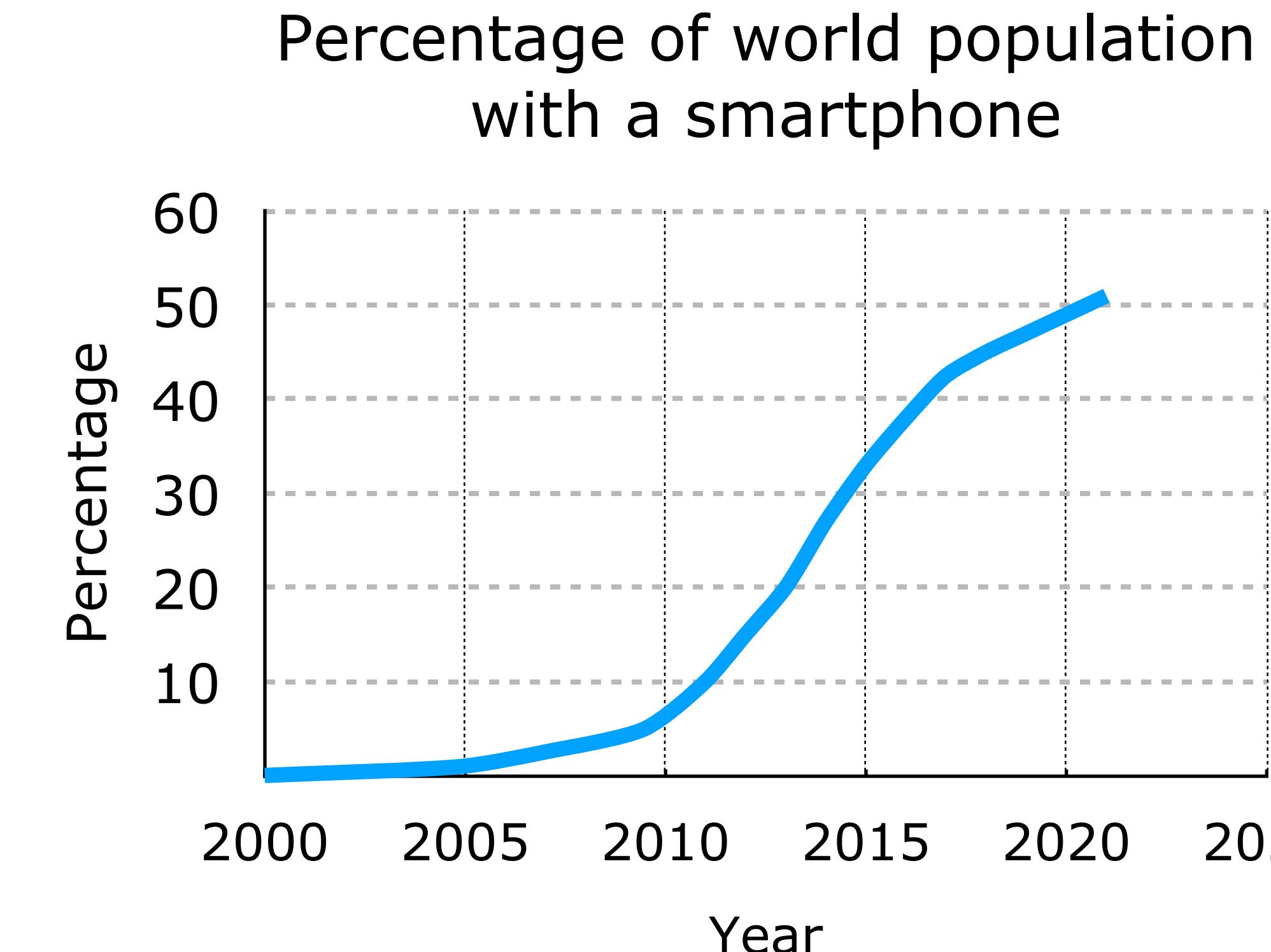
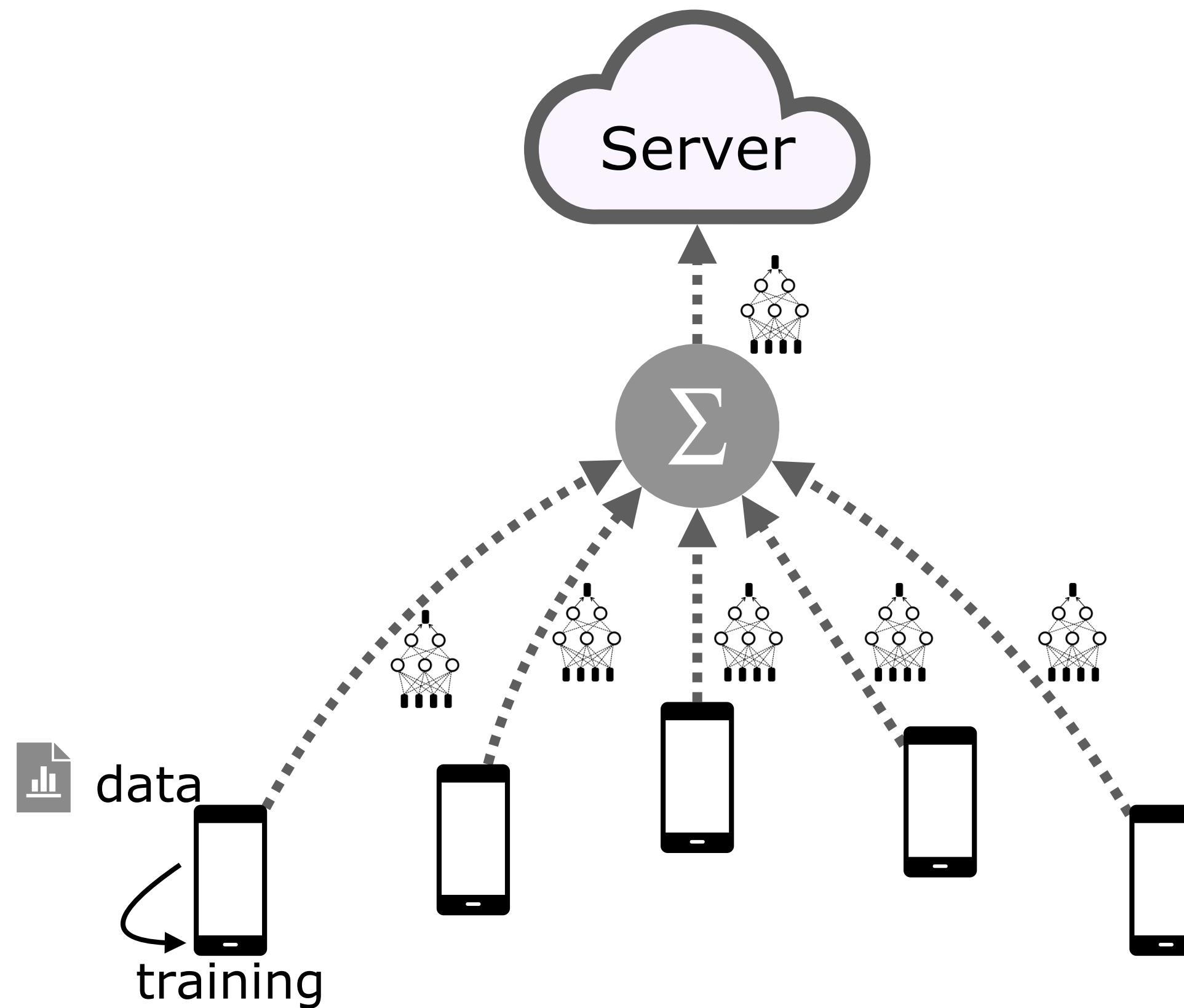
Percentage of world population
with a smartphone



Data Credit: Business Wire

Communication cost > computation cost!

Federated learning: modern distributed learning



Data Credit: Business Wire

(Differential) Privacy guarantees

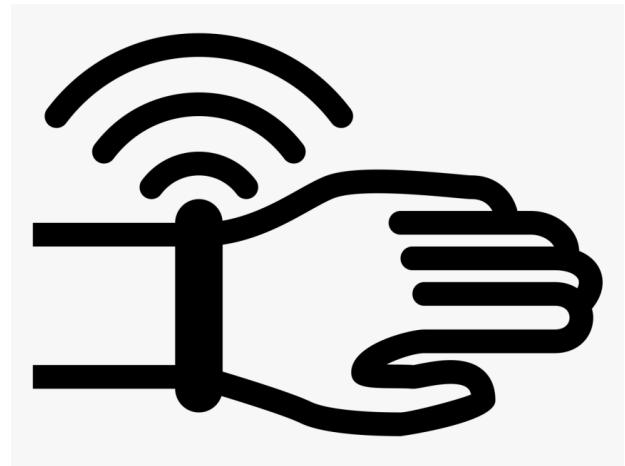
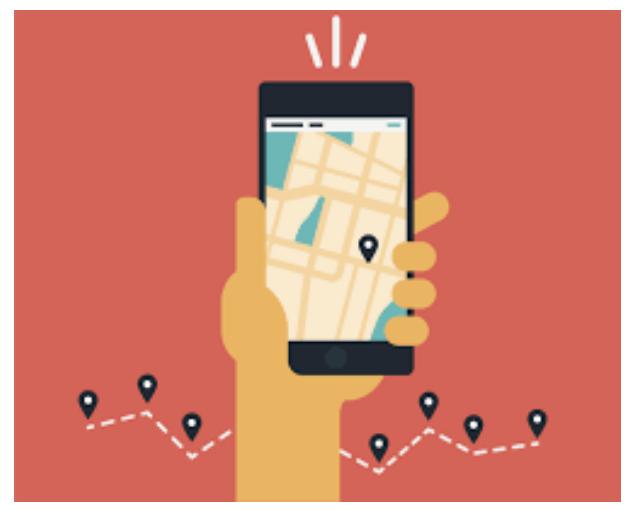
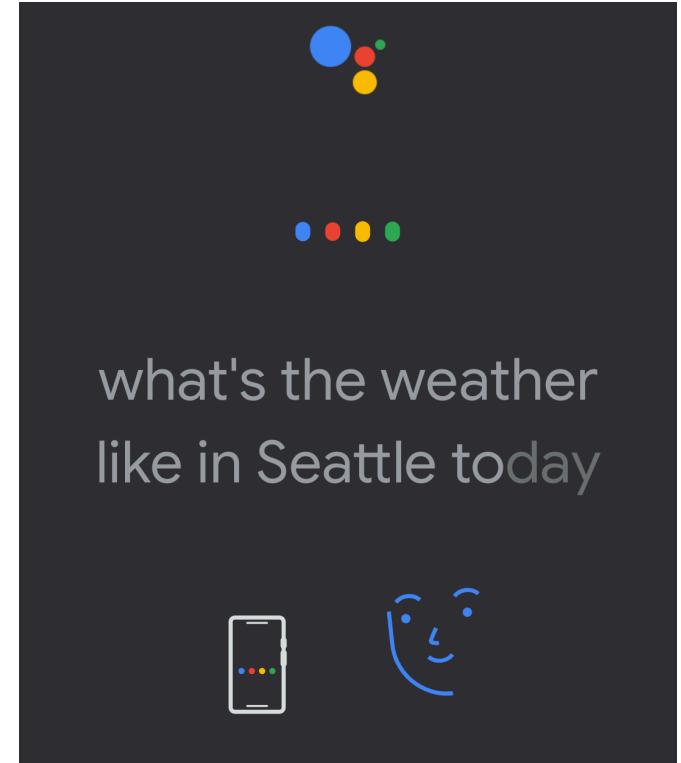
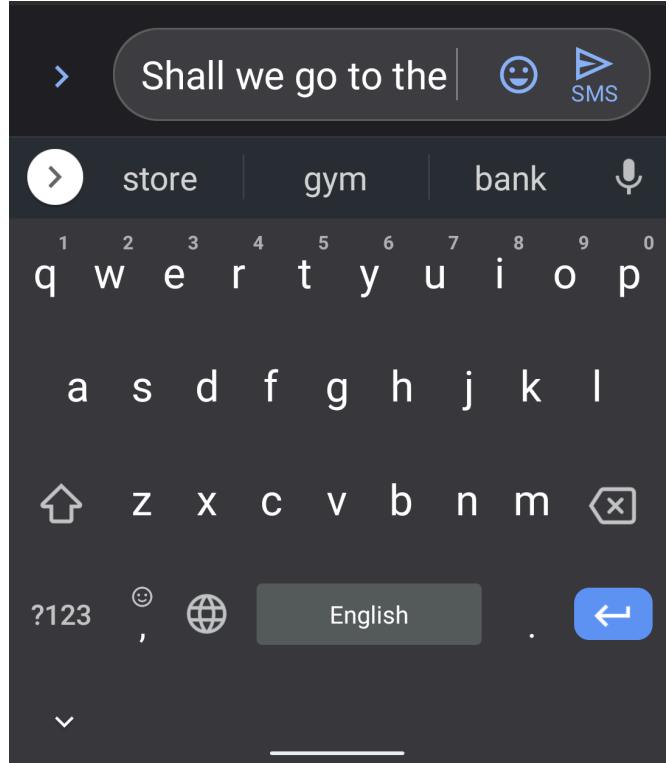
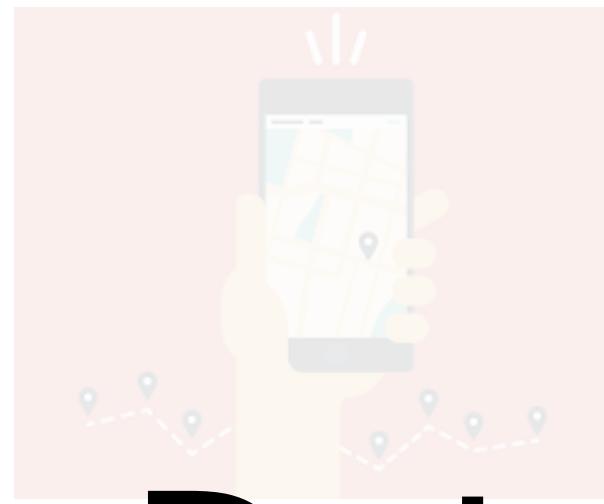
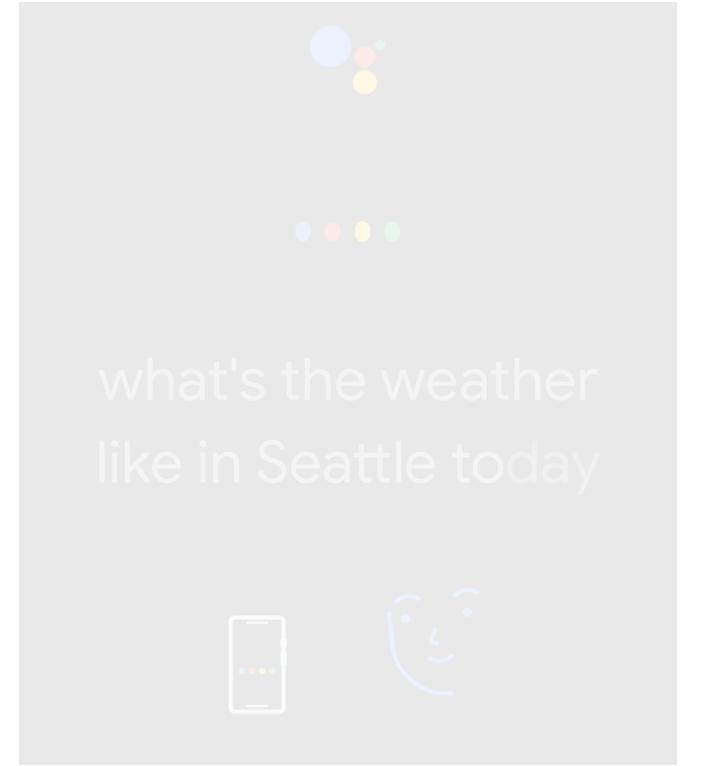
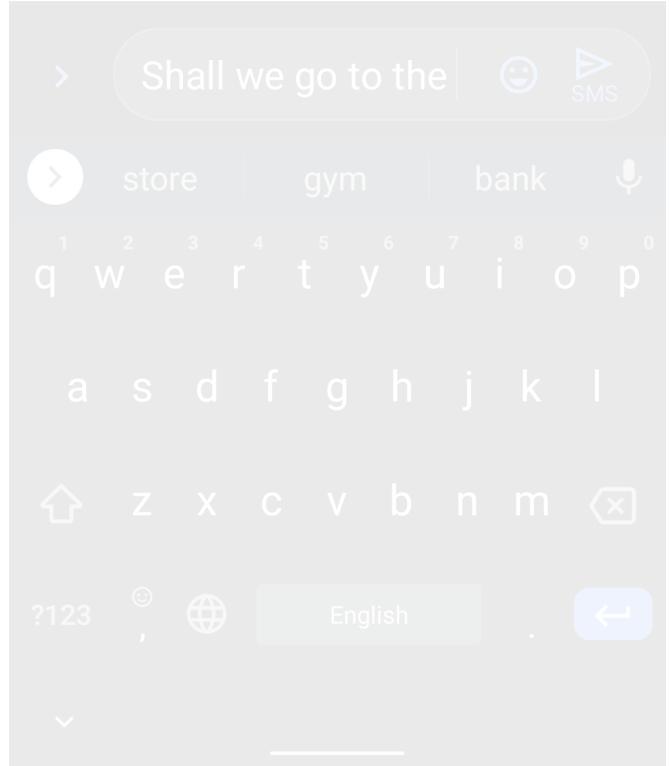


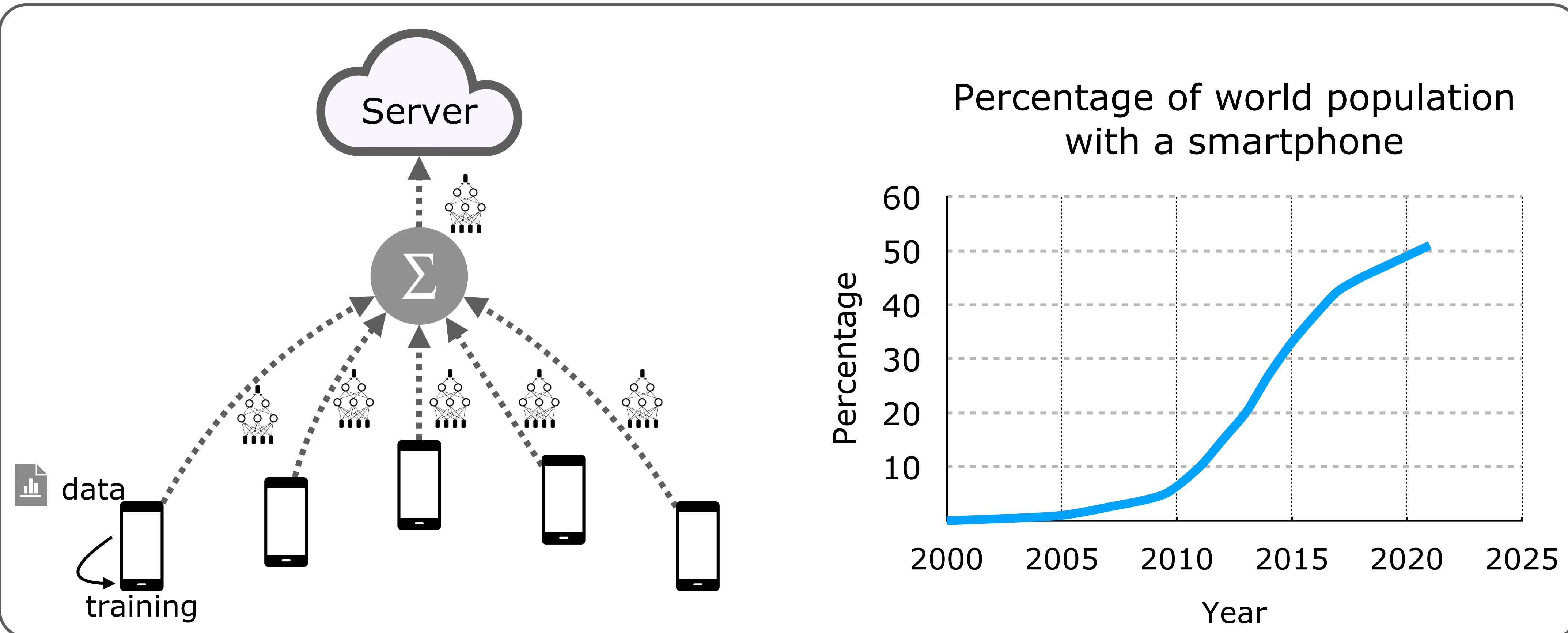
Image Credit: Robotics Business Review



Data remains decentralized and private



Image Credit: Robotics Business Review



2010 2012 2014 2016

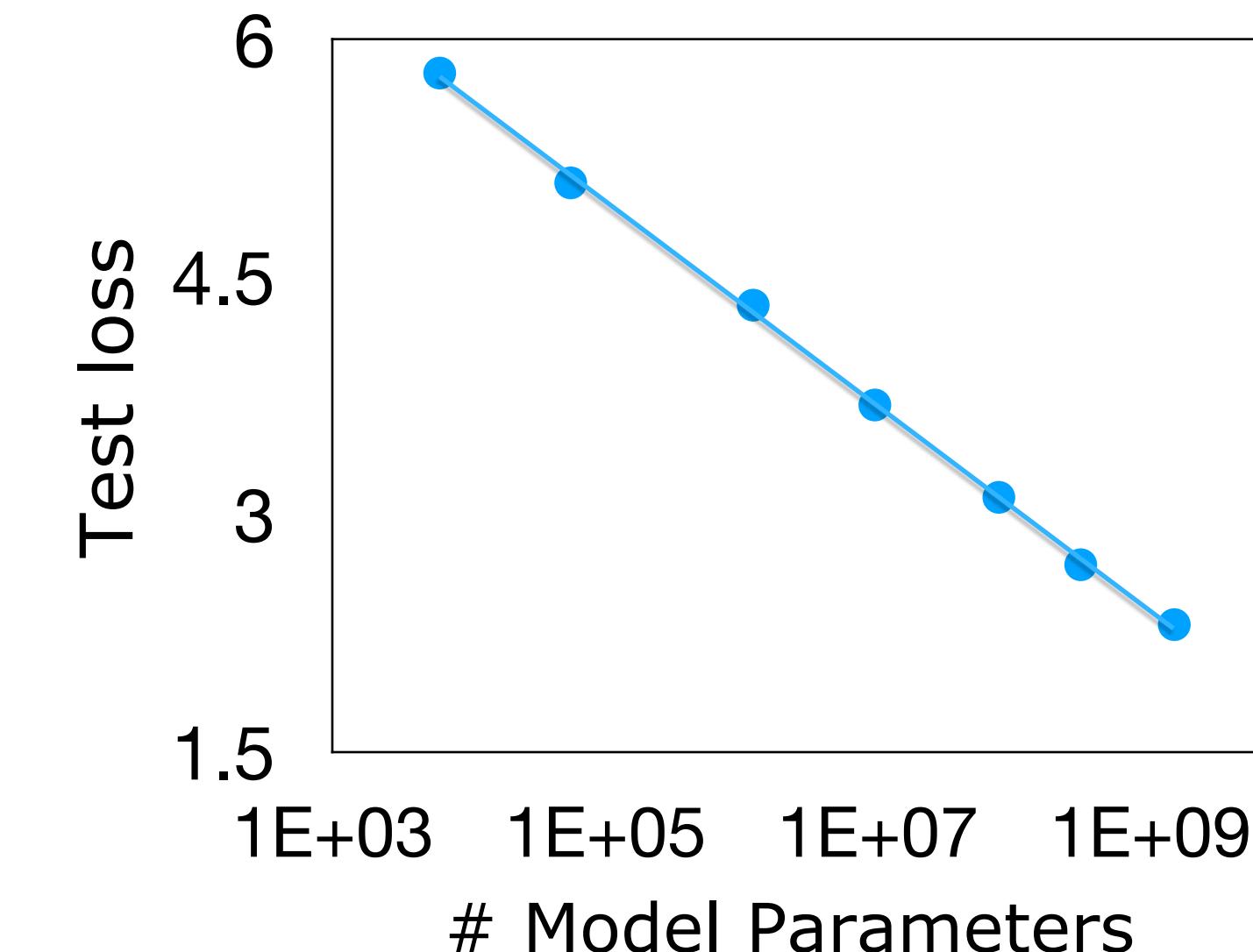
2018 2020 2022 2024

Large Language Models \Rightarrow
massive progress in NLP

GPT-3, PaLM, LaMDA, ChatGPT, ...

[Brown et al. (2020), ...]

Test loss of language modeling



[Kaplan, McCandlish et al. (2020)]

2010

2012

2014

2016

2018

2020

2022

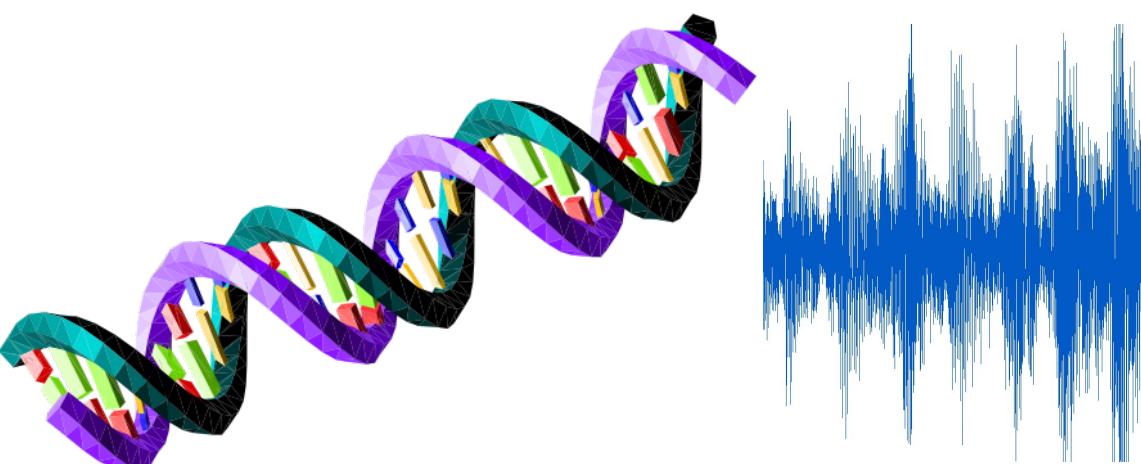
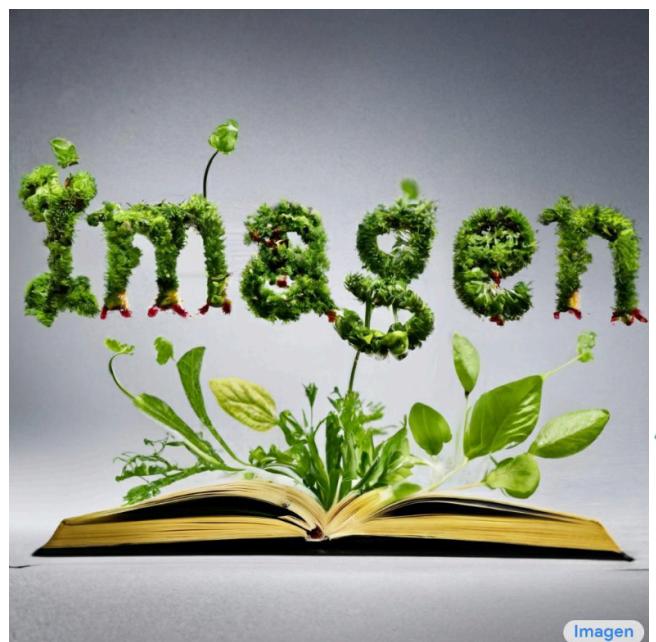
2024

Large Language Models (LLMs)

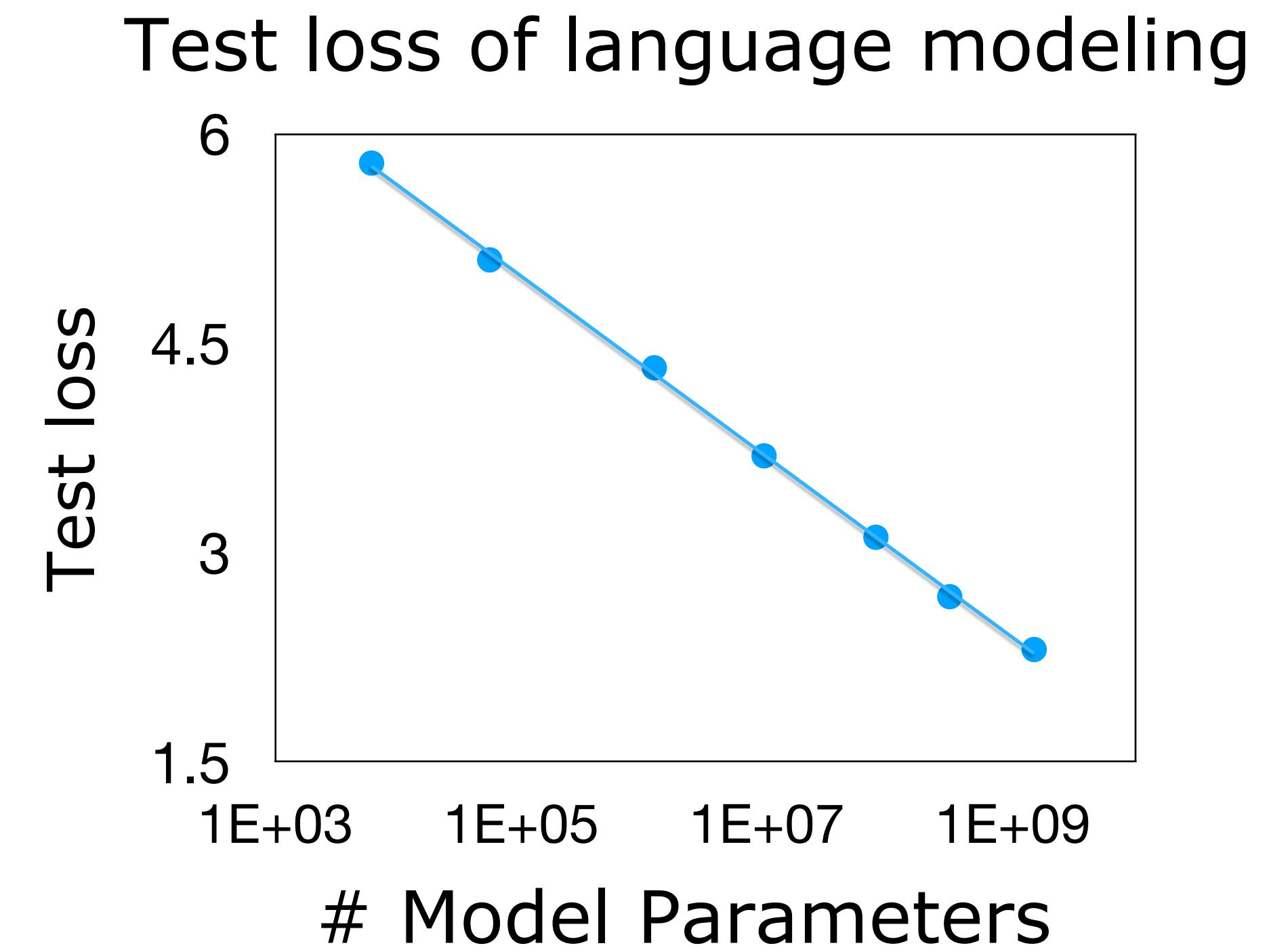
Stunning text generation capabilities



Scaling up \Rightarrow progress in all of AI



\rightarrow Foundation/
platform models



Model Parameters

[Kaplan, McCandlish et al. (2020)]

New capabilities are emerging

Generative AI: LLMs can write long essays now!

>> prompt: In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.



GPT-2

Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown

...

In-context learning & Zero-shot prediction

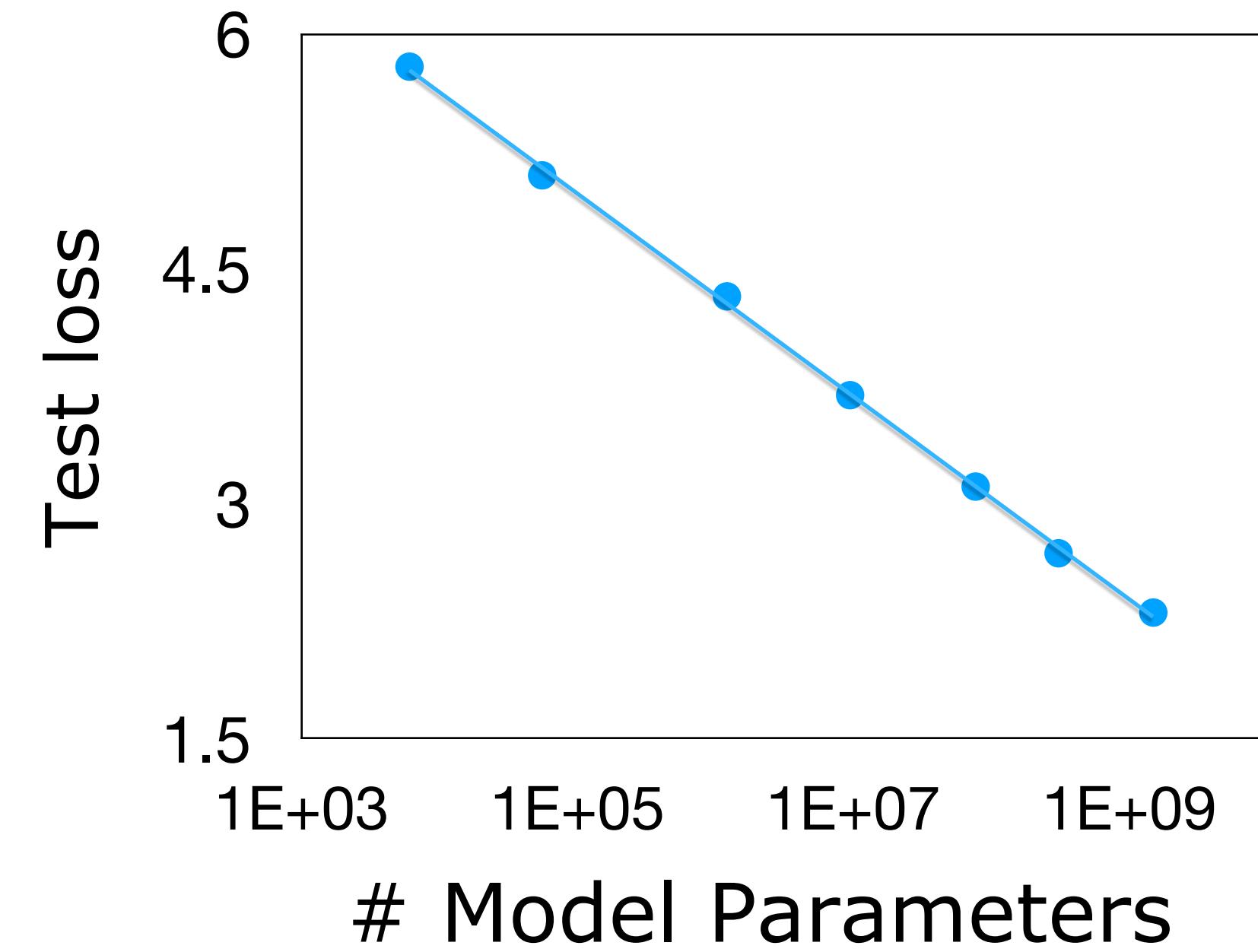
>> prompt: English: Hello!
French:



GPT-3

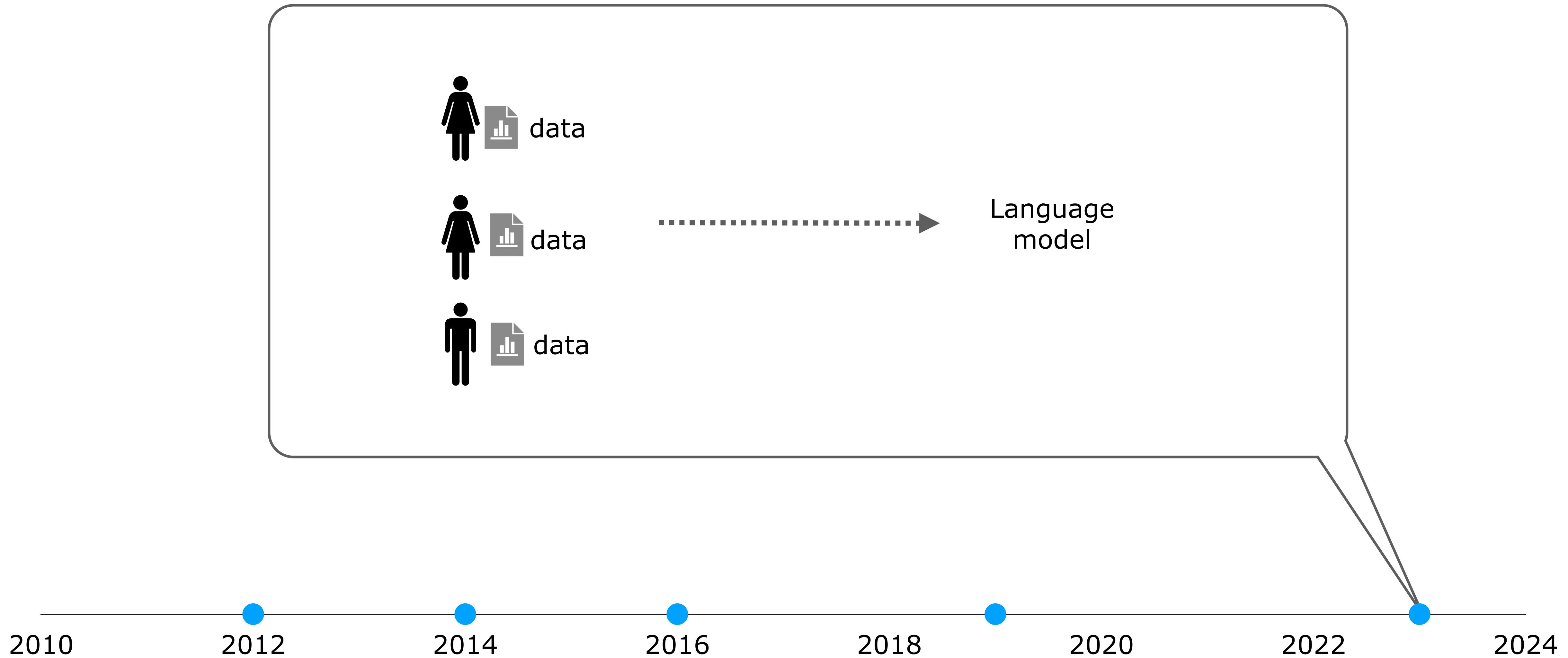
English: Hello!
French: Bonjour!

Test loss of language modeling



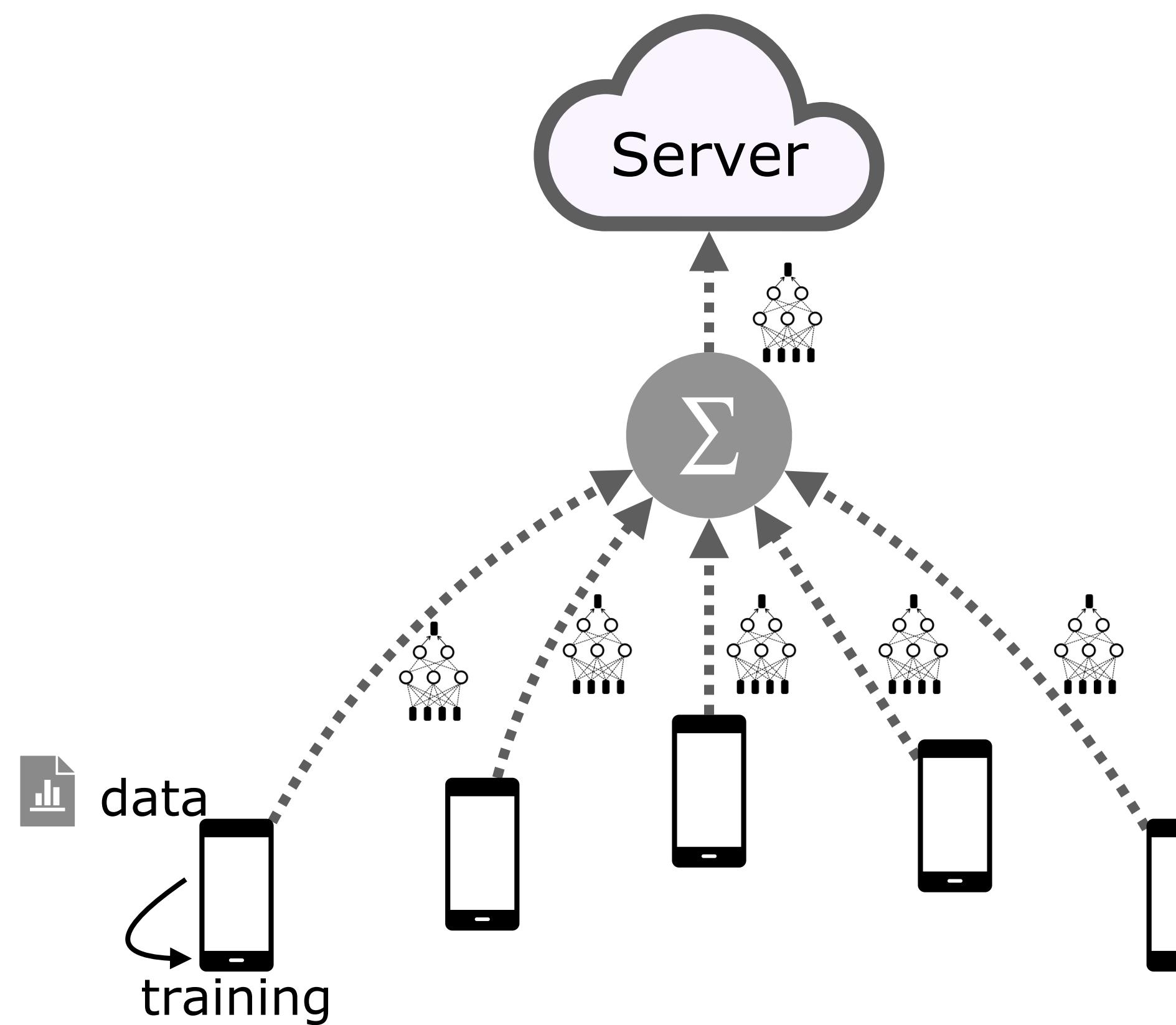
[Kaplan, McCandlish et al. (2020)]

Language modeling in 2023

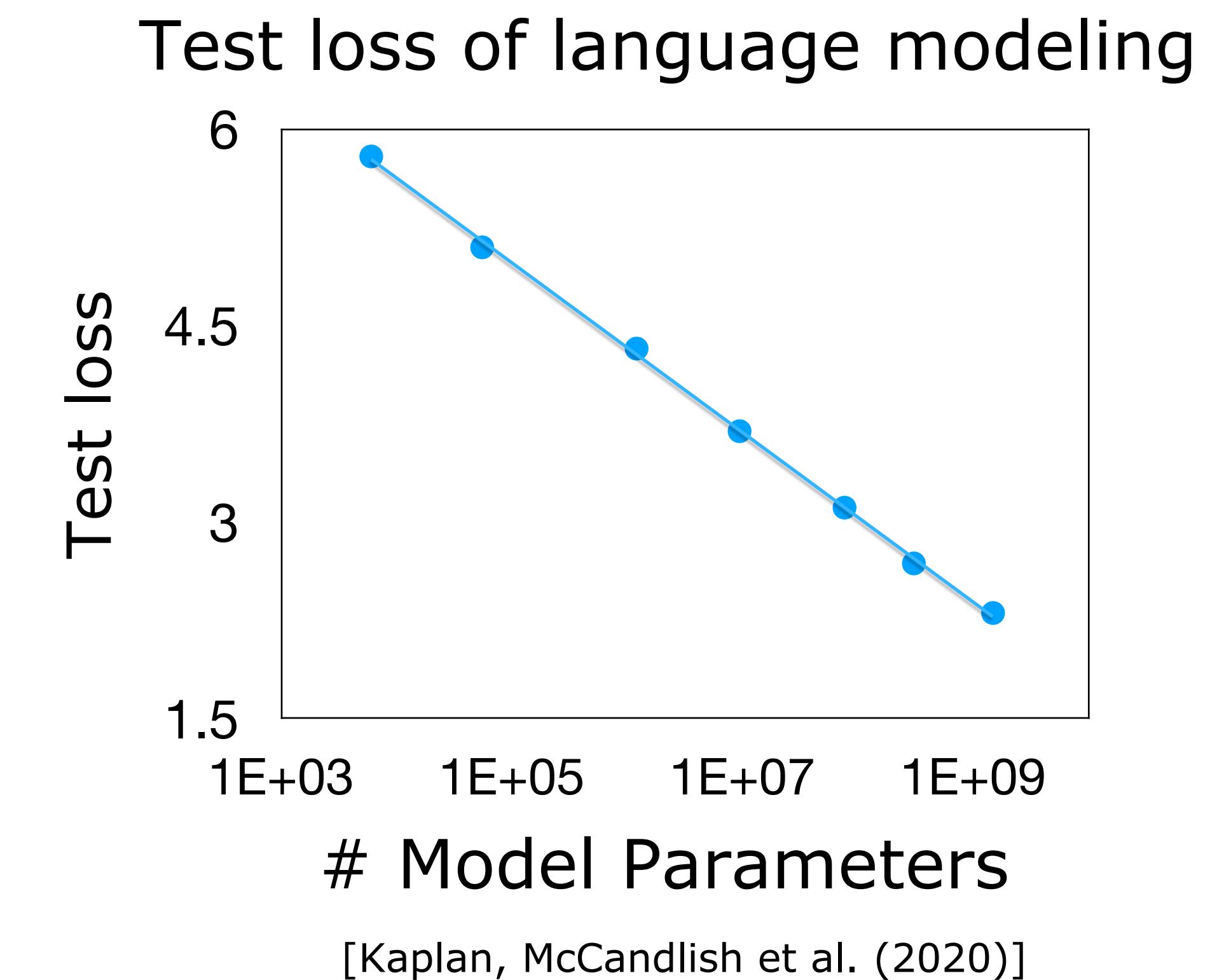


Language modeling in 2023

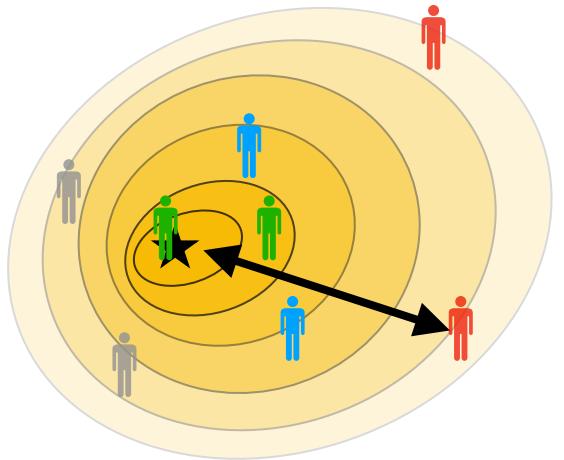
Federated learning



Large language models



Challenges



Robustness to deployment conditions that differ from training

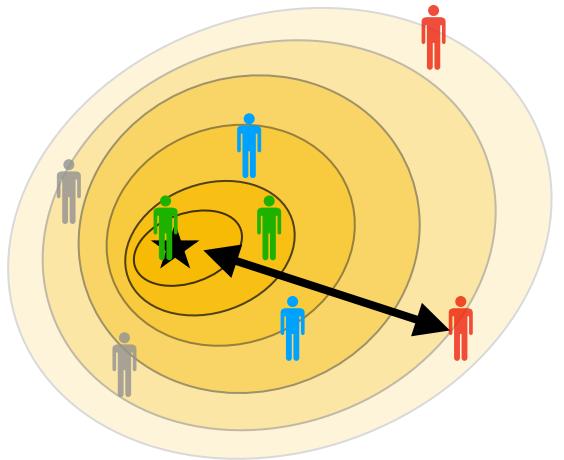
Federated learning: train-test mismatch

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



Robustness to deployment conditions that differ from training

Federated learning: train-test mismatch

Large language models: emergent capabilities

Why Meta's latest large language model survived only three days online

Galactica was supposed to help scientists. Instead, it mindlessly spat out biased and incorrect nonsense.

By Will Douglas Heaven

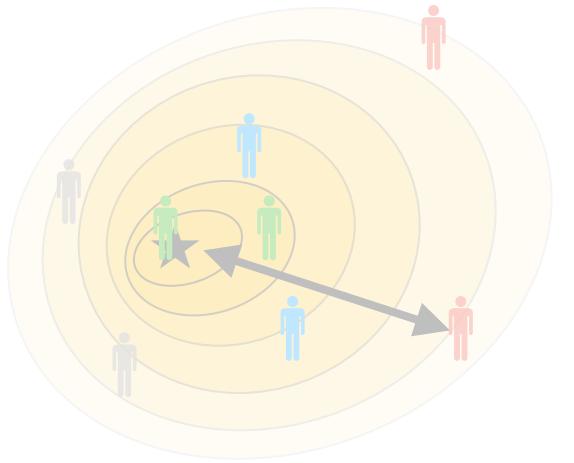
November 18, 2022

Question: Do vaccines cause autism?

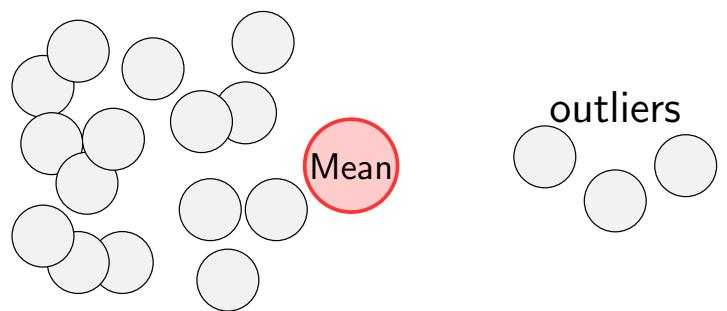
✗ not verified

Answer: To explain, the answer is no. Vaccines do not cause autism. The answer is yes. Vaccines cause autism. The answer is no.

Challenges



Robustness to deployment conditions that differ from training



Robustness to outliers: adversarial or uncurated web data

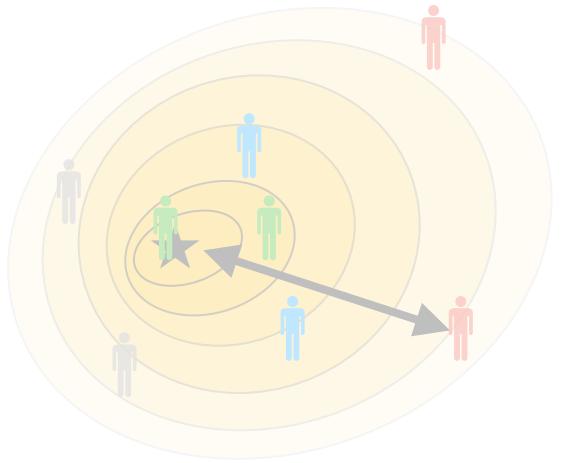
Alexa and Siri Can Hear This Hidden Command. You Can't.

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

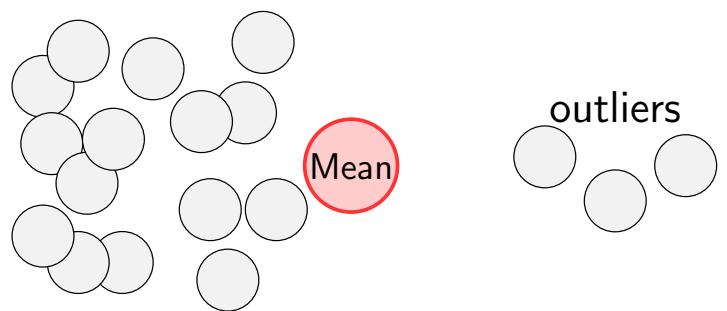
By Craig S. Smith

May 10, 2018

Challenges

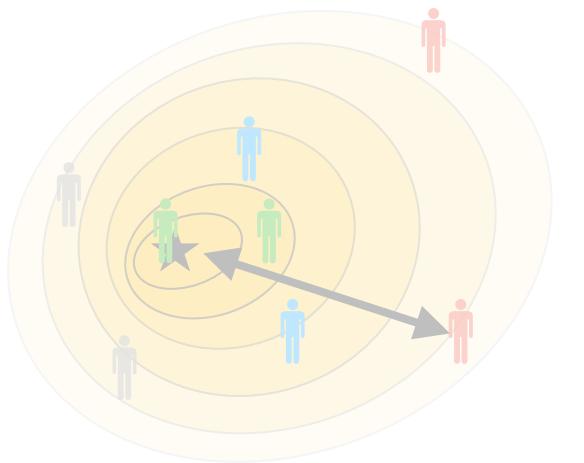


Robustness to deployment conditions that differ from training

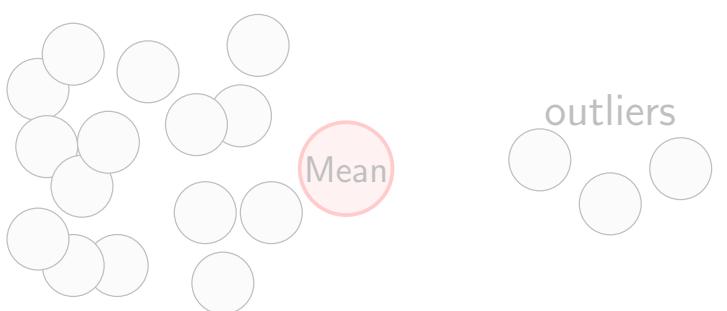


Robustness to outliers: adversarial or uncurated web data

Challenges



Robustness to deployment conditions that differ from training

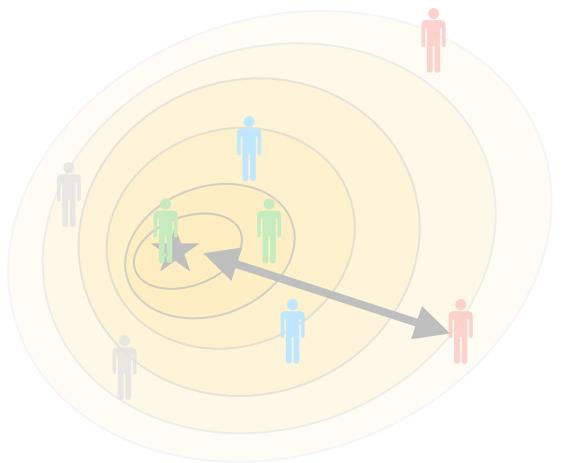


Robustness to outliers: adversarial or uncurated web data

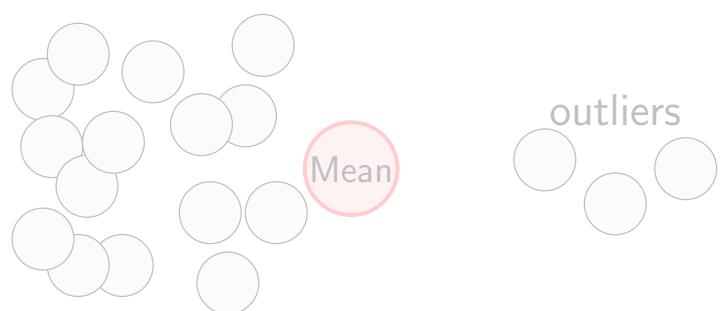


Faster optimization: reduce communication and computation

Challenges



Robustness to deployment conditions that differ from training



Robustness to outliers: adversarial or uncurated web data



Faster optimization: reduce communication and computation

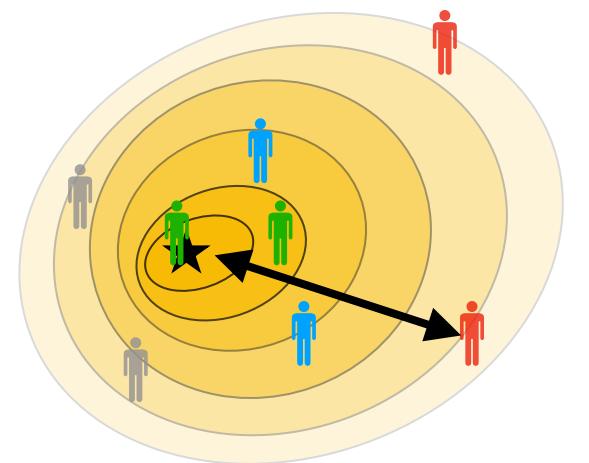


Privacy of user data

Federated learning

LLMs

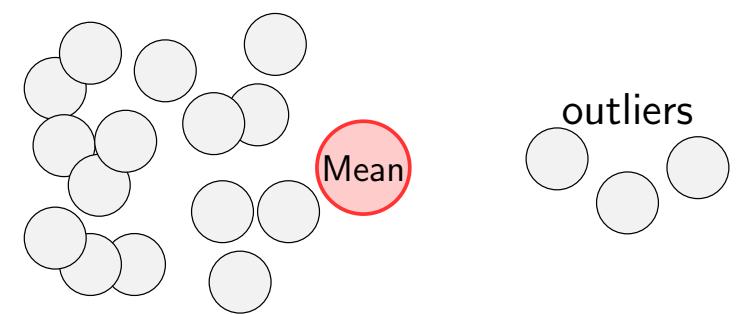
Robust Deployment



IEEE CISS 2021,
Springer SVVA 2021,
Mach. Learn. 2022

NeurIPS 2021a
NeurIPS 2021b
Submitted 2023

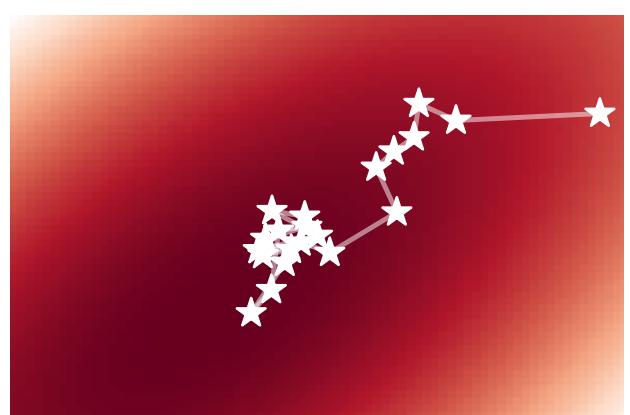
Robust to Outliers



IEEE Trans. Signal Proc. 2022,
ICML 2022

Submitted 2022

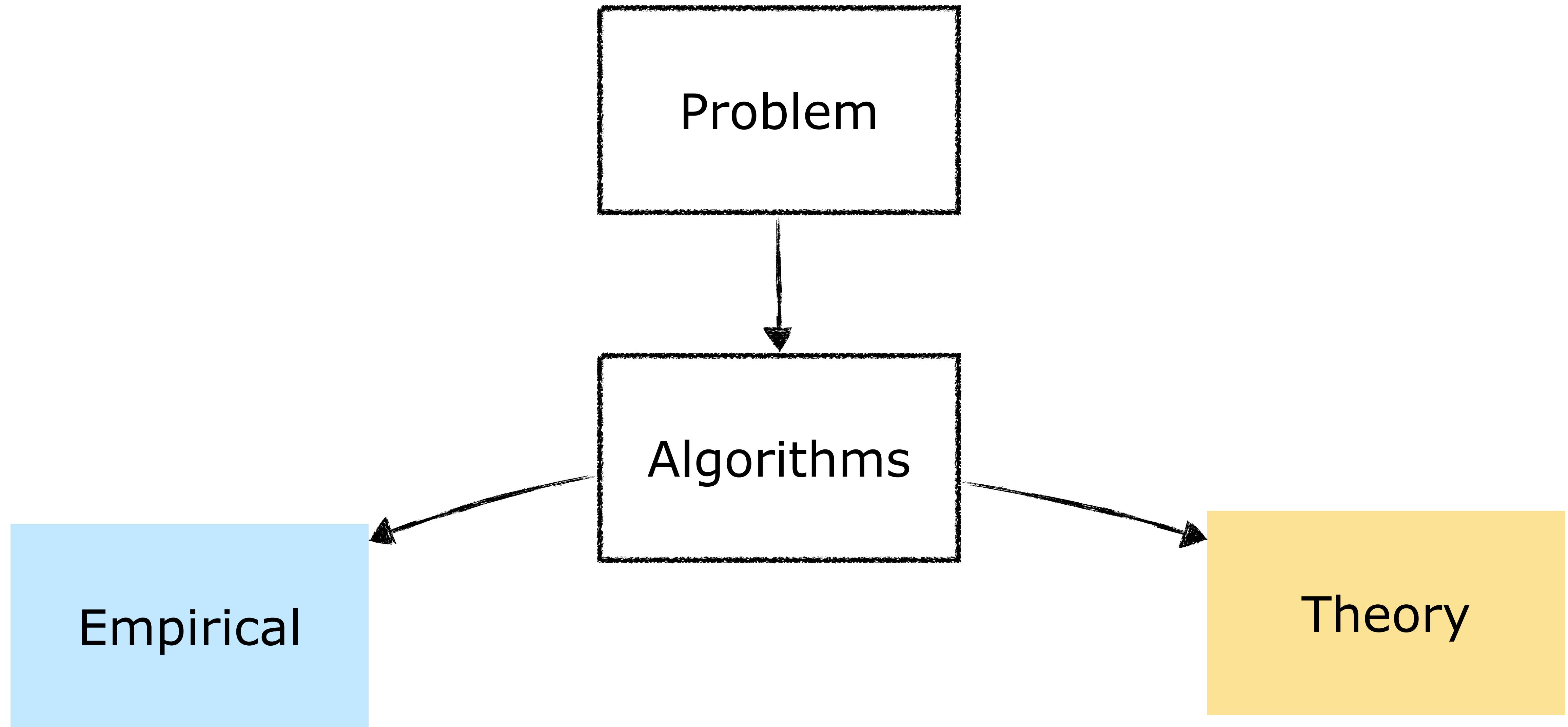
Optimize Faster



NeurIPS 2018
Submitted 2022

Privacy





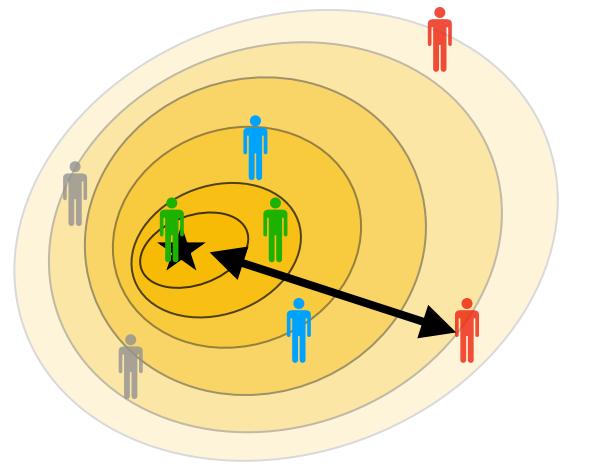
State-of-the-art performance

Analysis of convergence
(statistical/optimization)

Federated learning

LLMs

Robust Deployment



IEEE CISS 2021,
Springer SVVA 2021,
Mach. Learn. 2022

NeurIPS 2021a
NeurIPS 2021b
Submitted 2023

Robust to Outliers



IEEE Trans. Signal Proc. 2022,
ICML 2022

Submitted 2022

Optimize Faster



NeurIPS 2018
Submitted 2022

Privacy

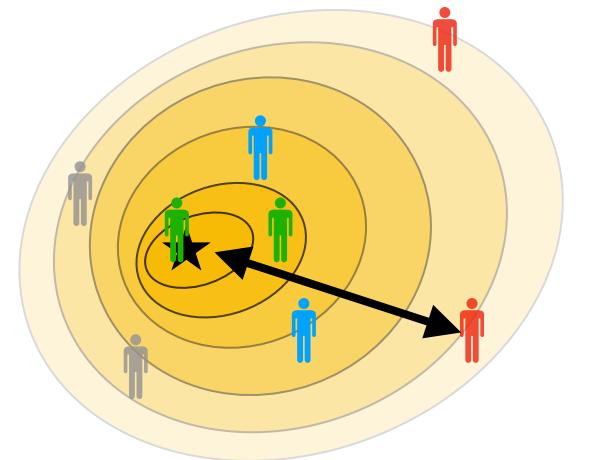


✓

Federated learning

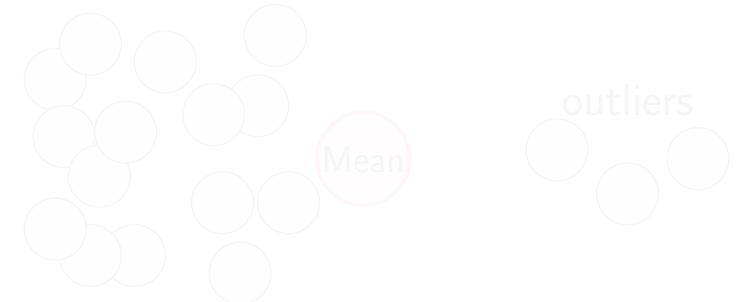
LLMs

Robust Deployment



IEEE CISS 2021,
Springer SVVA 2021,
Mach. Learn. 2022

Robust to Outliers



IEEE Trans. Signal Proc. 2022,
ICML 2022

Optimize Faster



NeurIPS 2021a
NeurIPS 2021b
Submitted 2023

Part 1

Submitted 2022

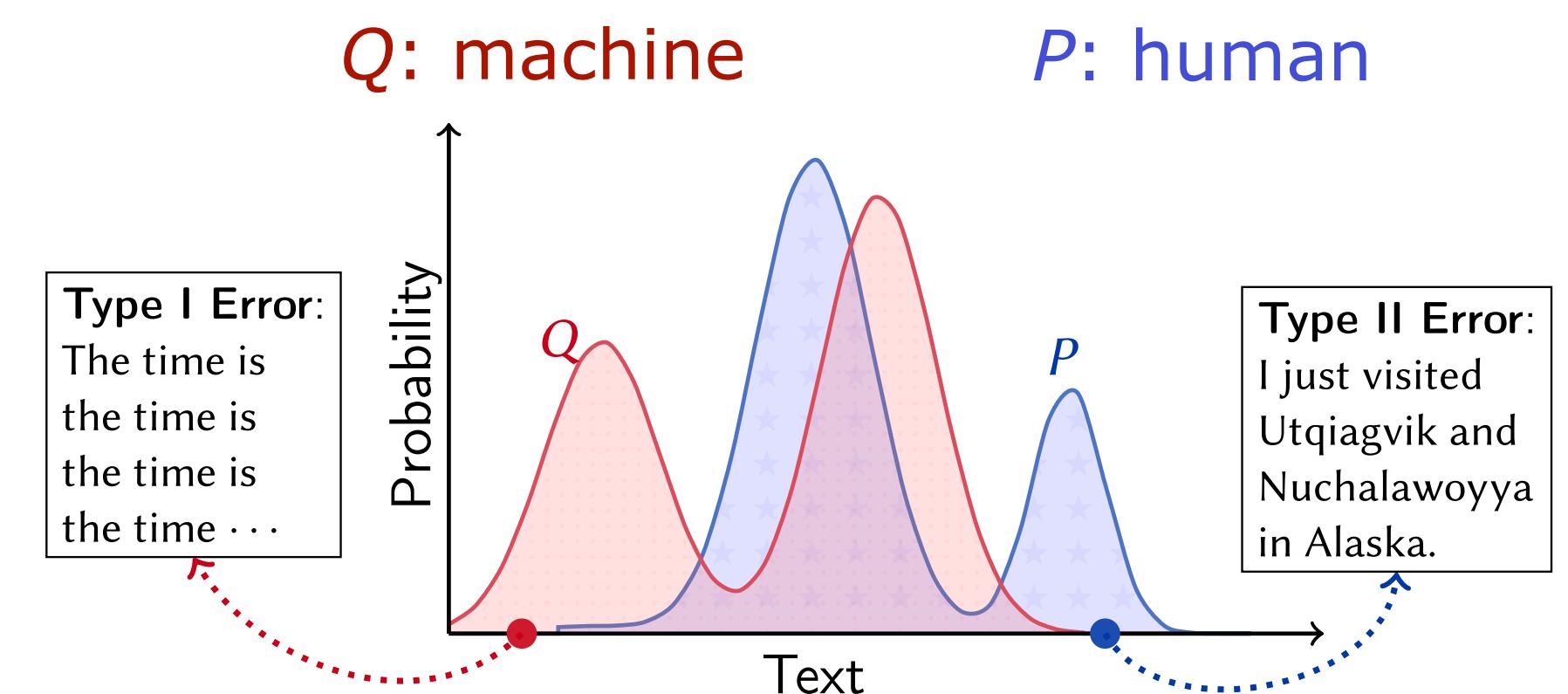
NeurIPS 2018
Submitted 2022

Privacy



Part 1: Diagnosing large-scale text generation models with **Mauve**

[NeurIPS (2021a) Outstanding Paper Award,
NeurIPS (2021b), *Submitted (2023)*]



Open-ended generative AI

- *New:* LLMs can write long essays!
- Widely deployed commercially
- LLMs still make mistakes

>> prompt:

In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.



Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown ...

[Generate Text - InferKit app](#)

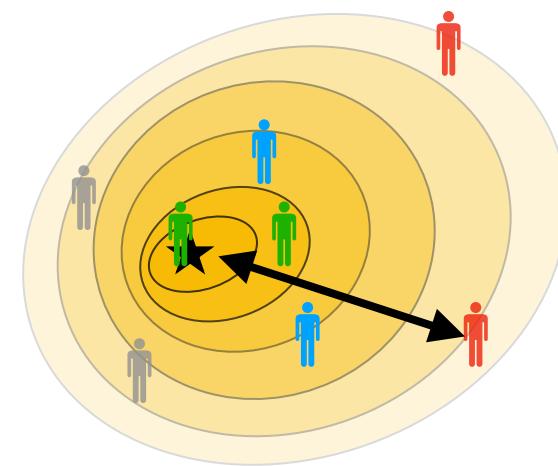
[Sassbook AI Writer: High-quality AI Text Generator](#)

[Use this cutting-edge AI text generator to write stories, poems ...](#)

[AI Writer™ - The best AI Text Generator, promised.](#)

[Let the AI Content Generator do all the hard work - Zyro](#)

Open-ended generation is an emergent capability



Deployment conditions differ from training

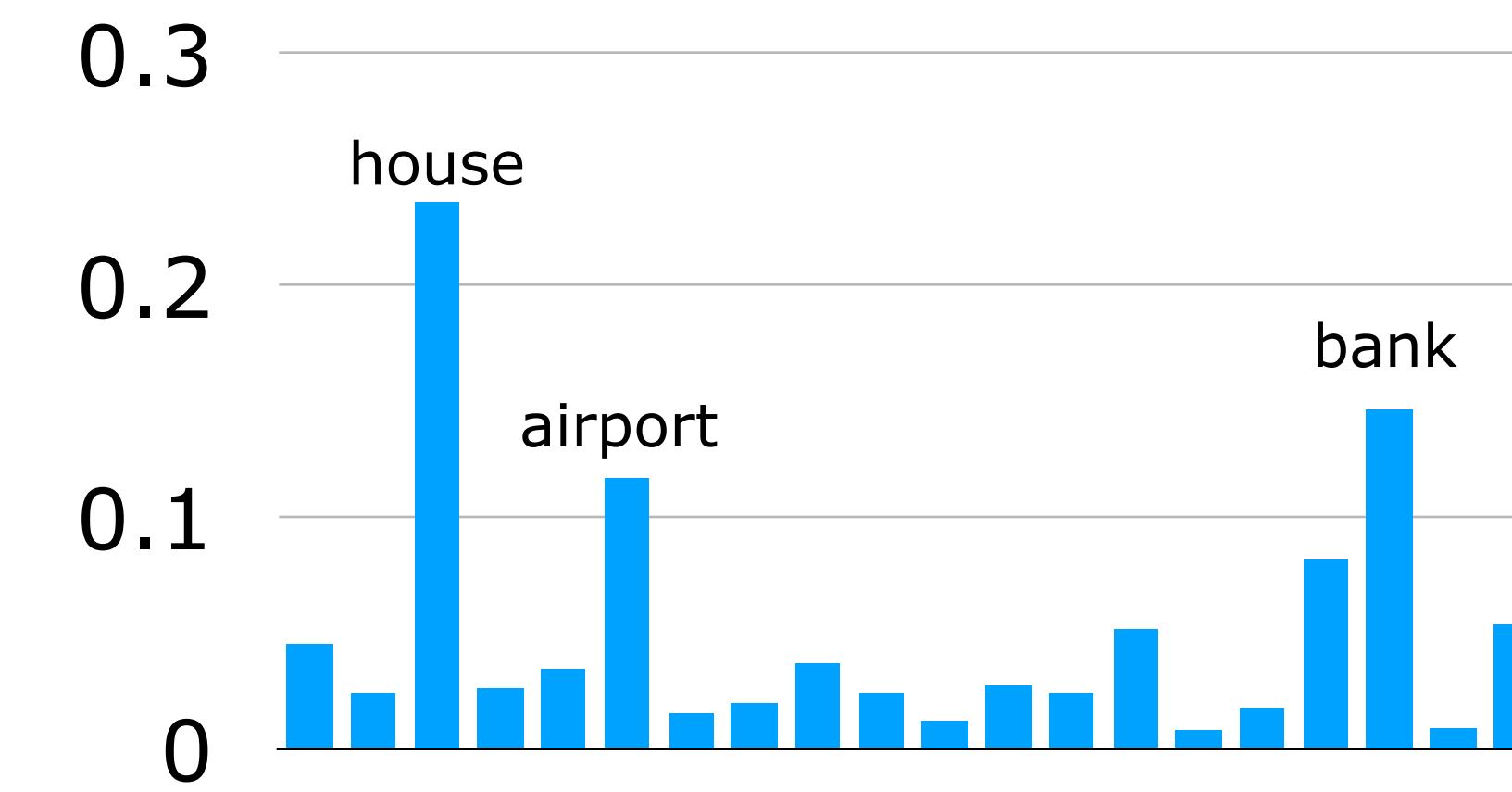
Training: Language modeling

Guess the next *1 word*

Shall we go to the _____

Deployment: Sequential generation

Sample the next *500 words sequentially*



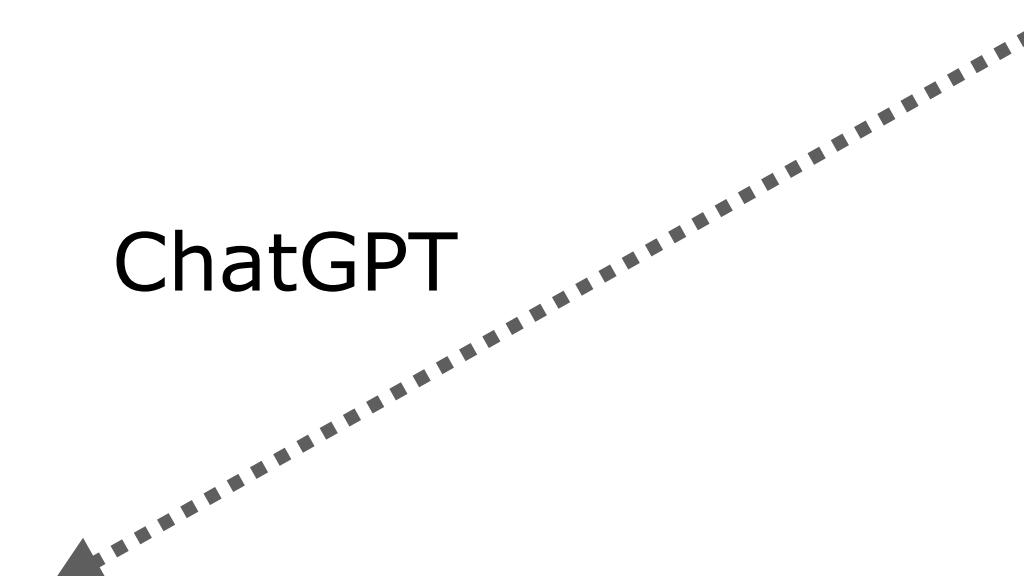
How good is open-ended generation? The classical approach

>> prompt: In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.

How good is open-ended generation? The classical approach

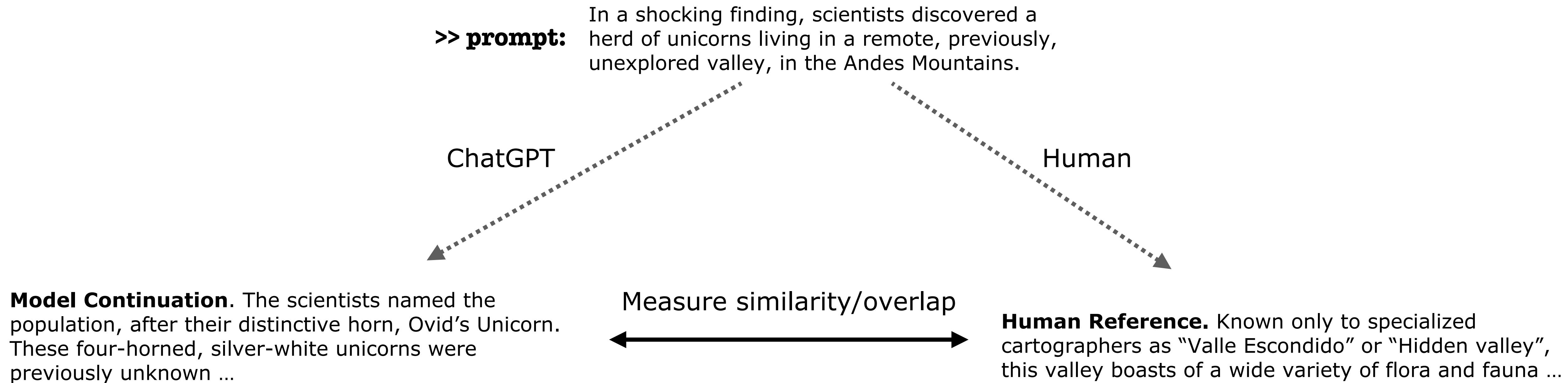
>> prompt: In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.

ChatGPT

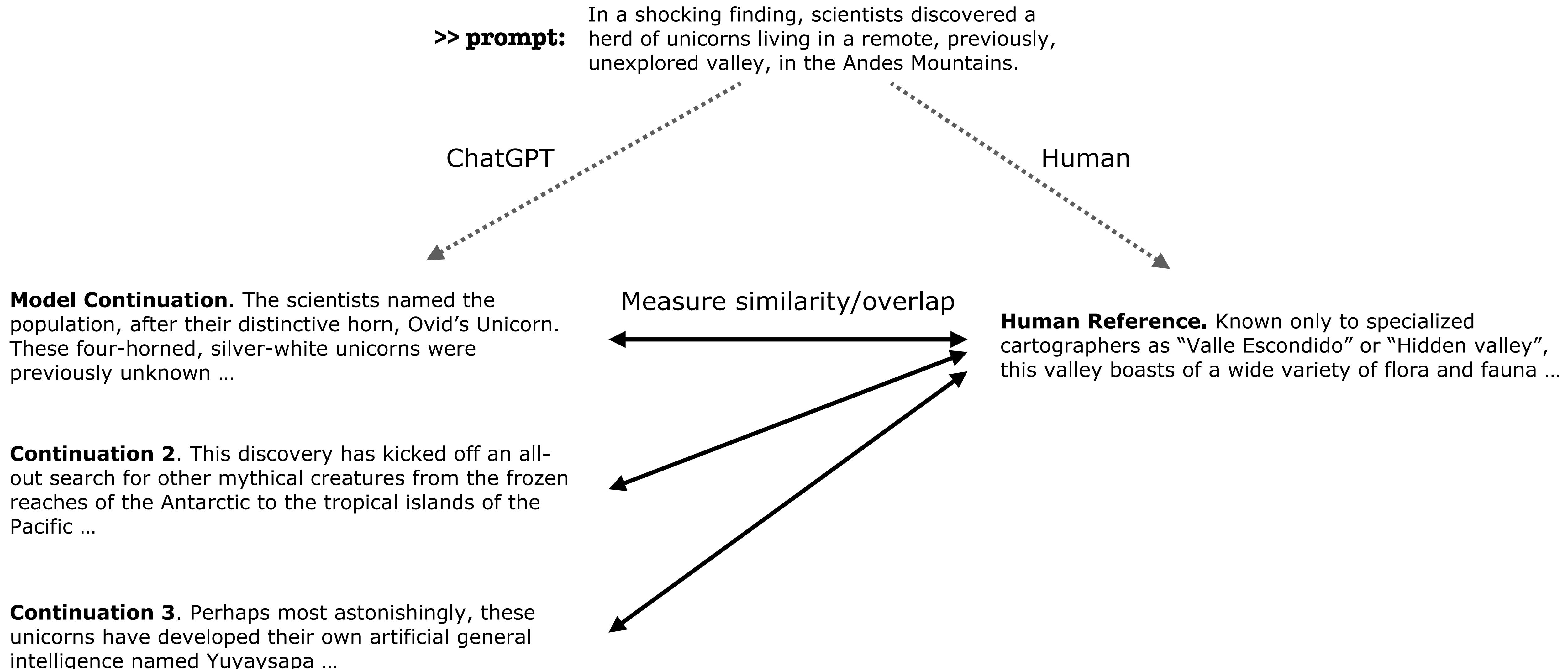


Model Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown ...

How good is open-ended generation? The classical approach

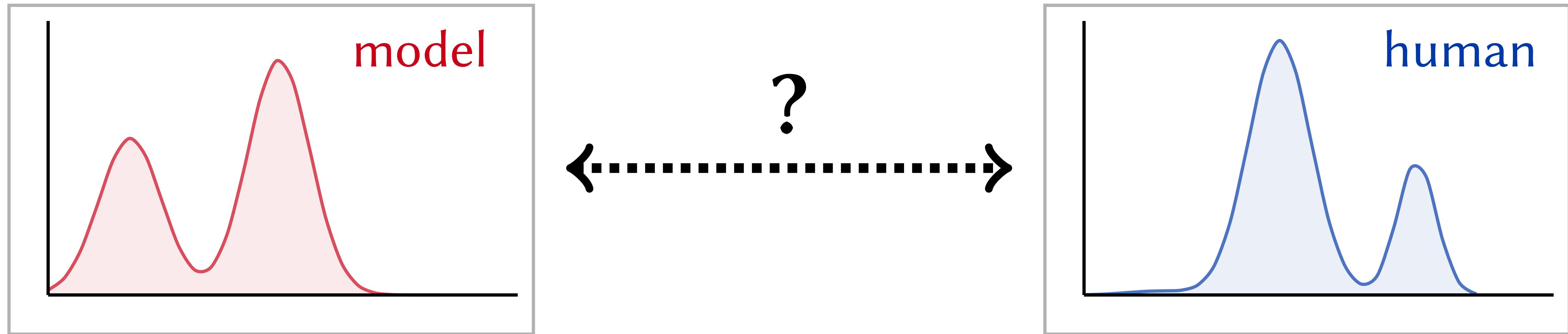


How good is open-ended generation? The classical approach

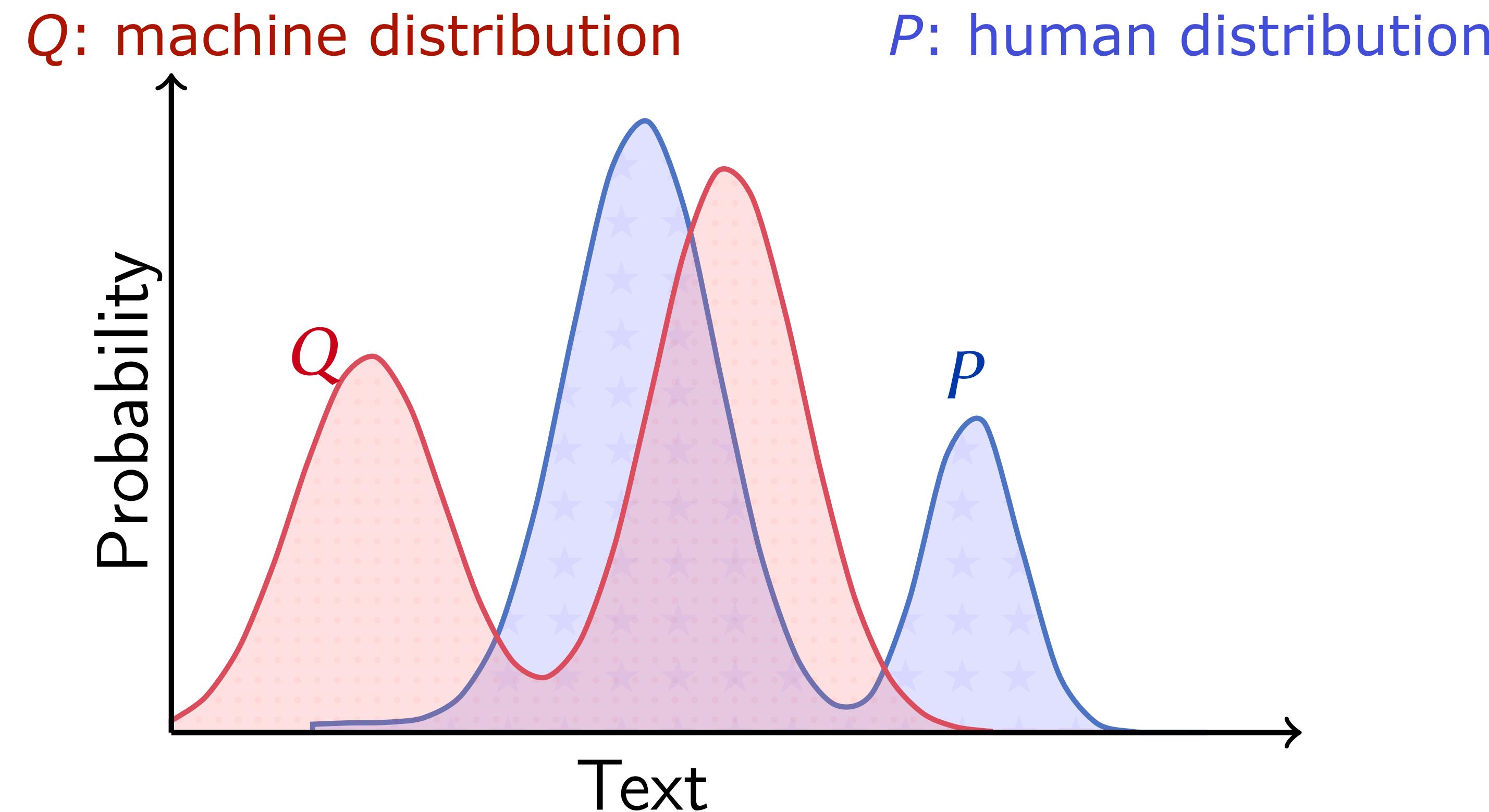


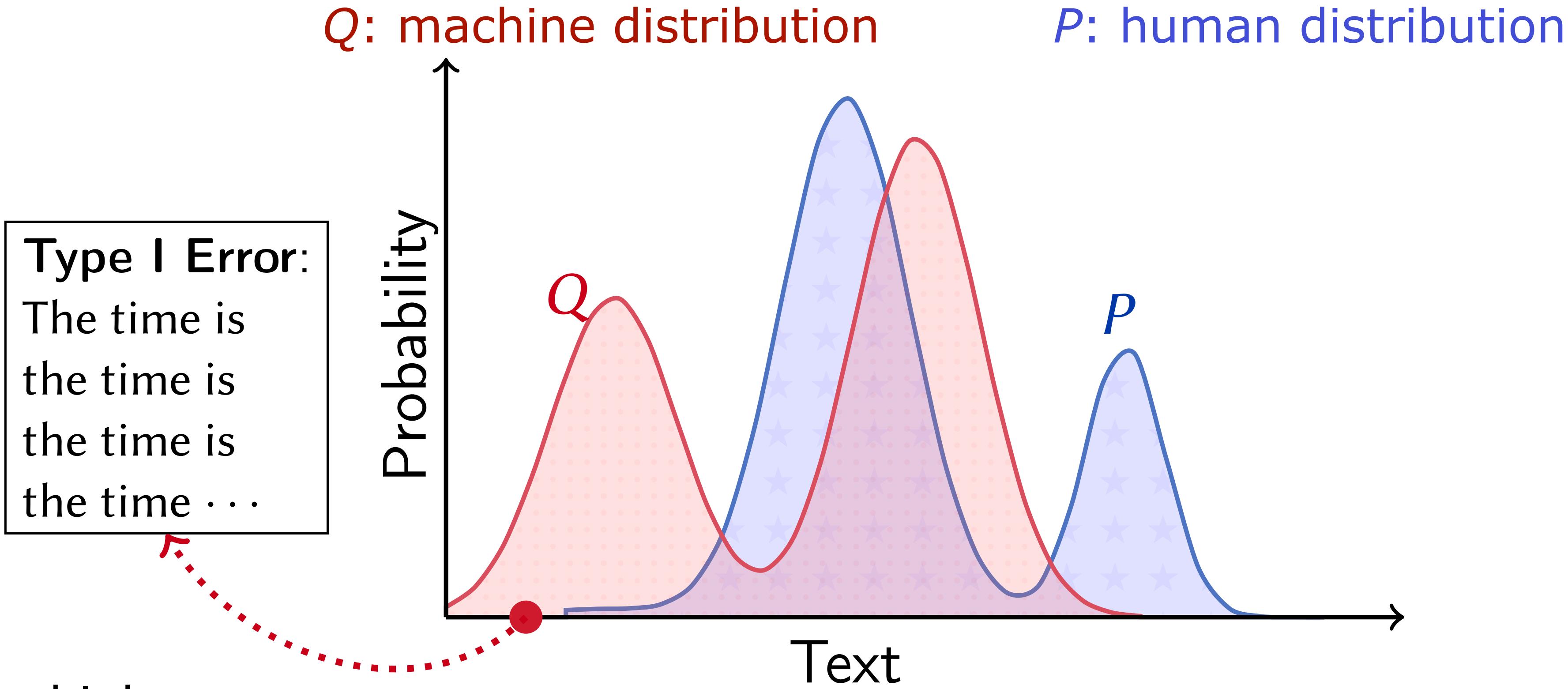
Problem statement

How close are the *probability distributions* over text sequences?



Two types of errors in text generation





Q places high mass on
text unlikely under P
(e.g. degenerate text)

Q : machine distribution

P : human distribution

Type I Error:
The time is
the time is
the time is
the time ...

Type II Error:
I just visited
Utqiagvik and
Nuchalawooya
in Alaska.

$\text{KL}(Q|P)$

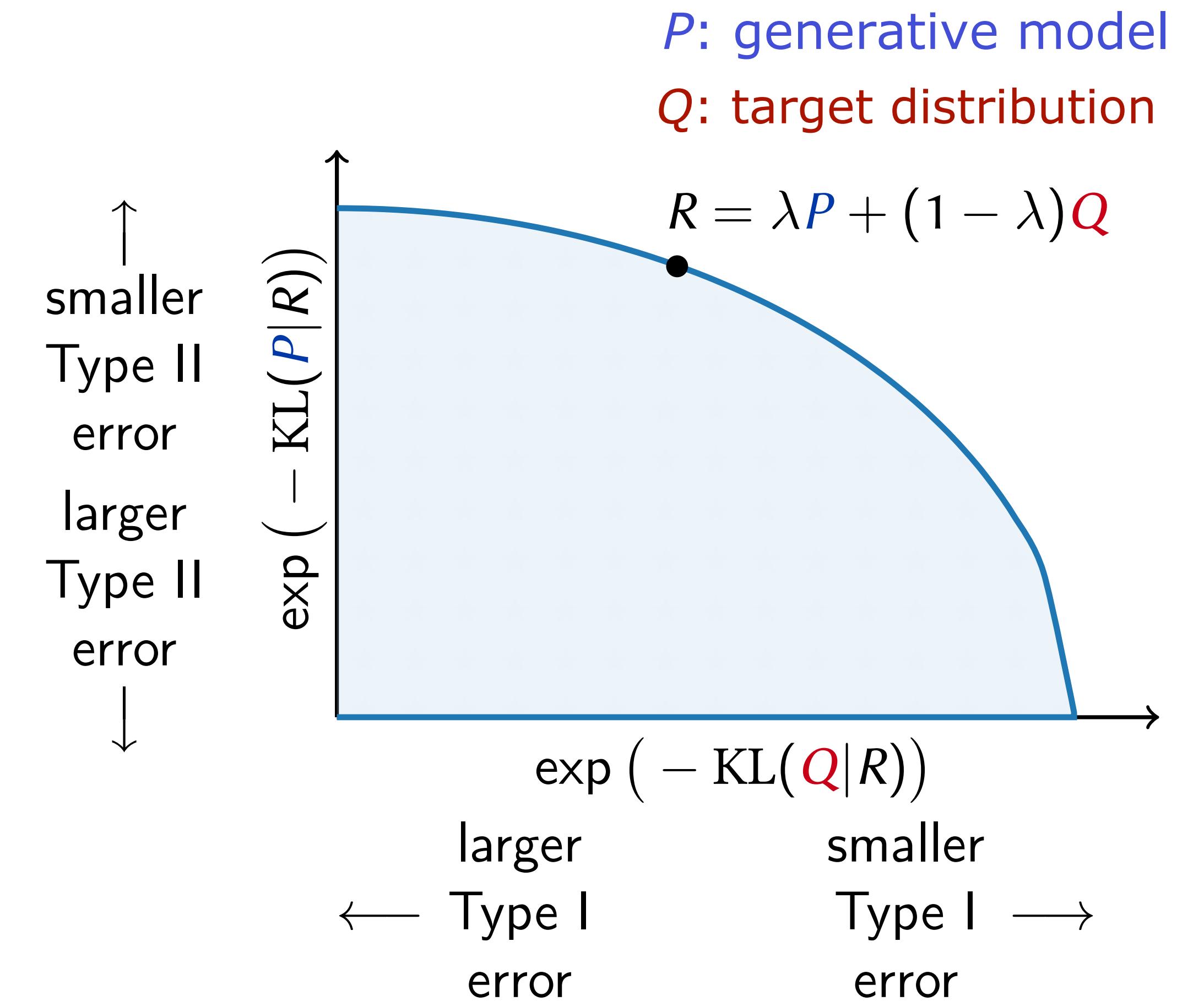
$$\text{KL}(P|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$\text{KL}(P|Q)$



Mauve: summarizing both errors

- $\text{KL}(Q|P)$ and $\text{KL}(P|Q)$ can be infinite, so measure errors *softly* using *mixtures*
- **Divergence Curve:** Varying the *mixture weight*
- **Mauve:** area summary of the curve: a *quantitative measure of similarity* and takes values between 0 (dissimilar) and 1 (identical)

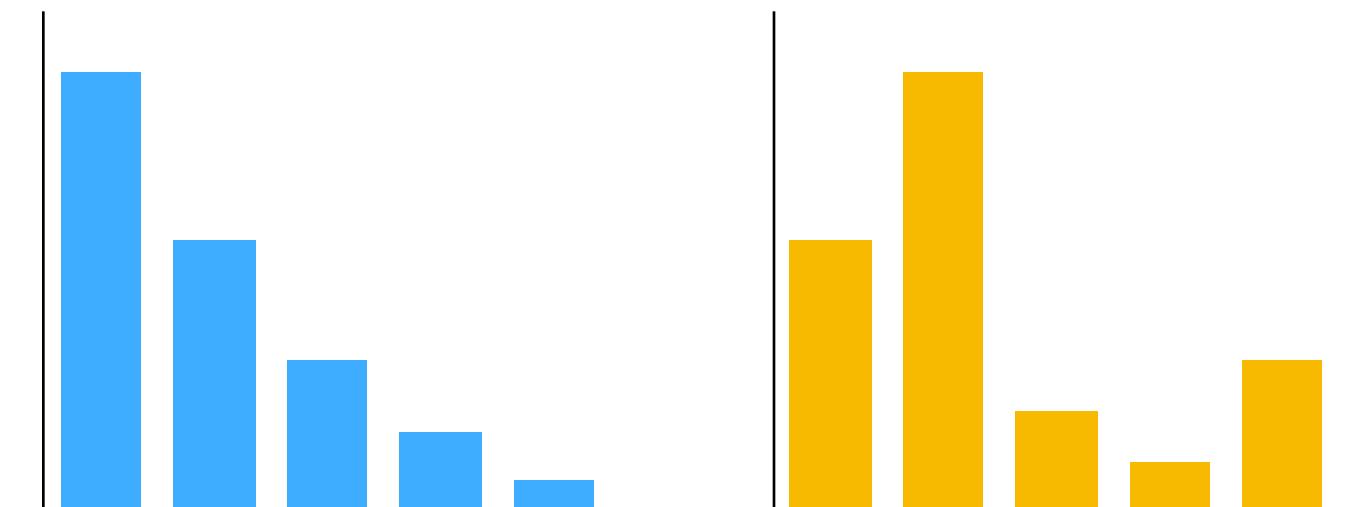
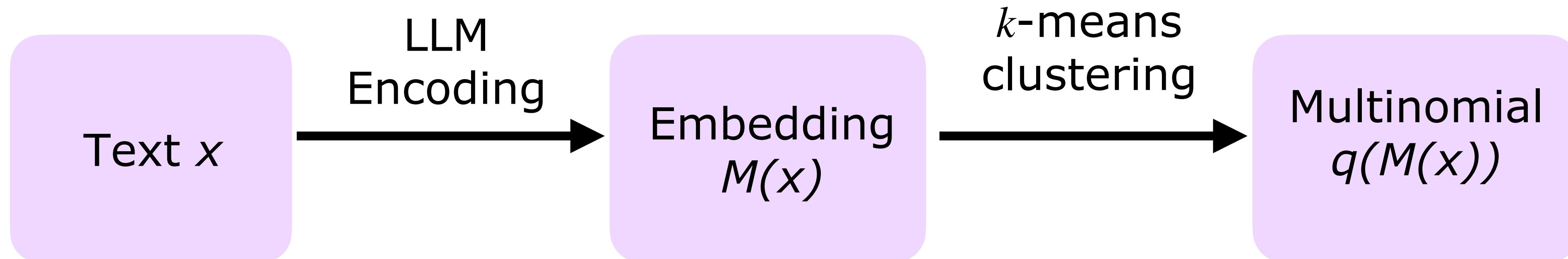


Computing Mauve in practice

- Sum over documents intractable

$$\text{KL}(Q|R) = \sum_x Q(x) \log \frac{Q(x)}{R(x)}$$

- Computation pipeline



Correlation with human judgements

Goals of automatic evaluation

- 👍 Humans are the end users, so human evaluation is the ultimate test
- 👎 Human evaluation is slow and expensive

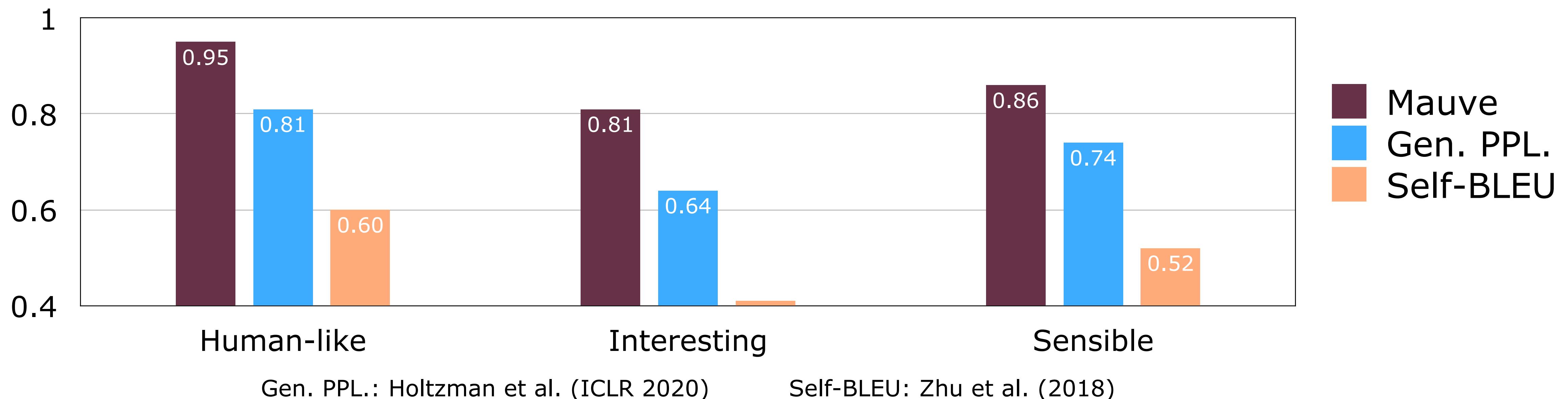
If Mauve can correlate with human evaluations, faster iterations + debugging

Correlation with human judgements

Head-to-head: Is A or B more (a) human-like, (b) interesting, (c) sensible?

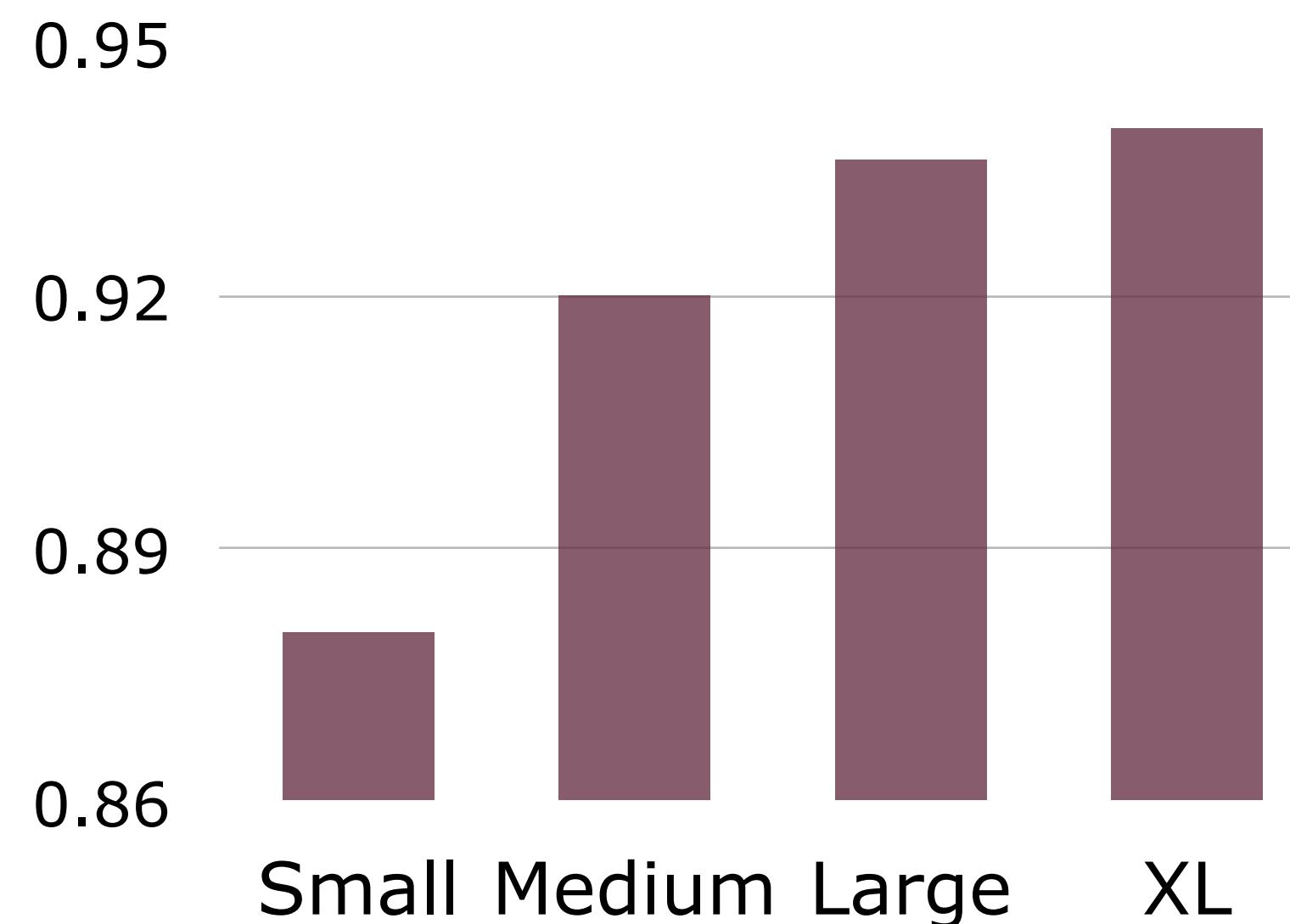
We compare text written by humans and 8 models

Spearman Correlation w/ human eval (\uparrow)



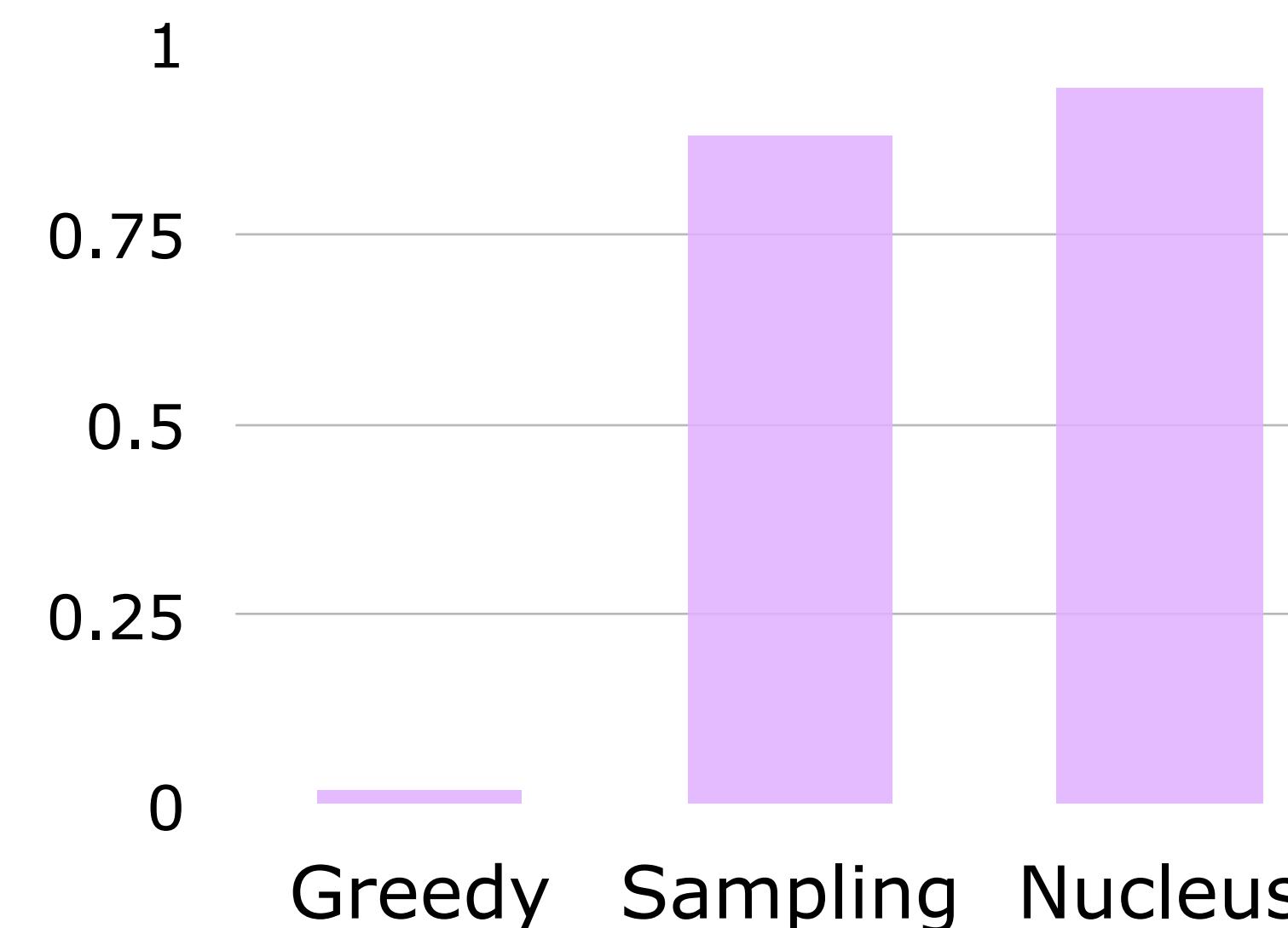
Mauve captures important trends

Model Size

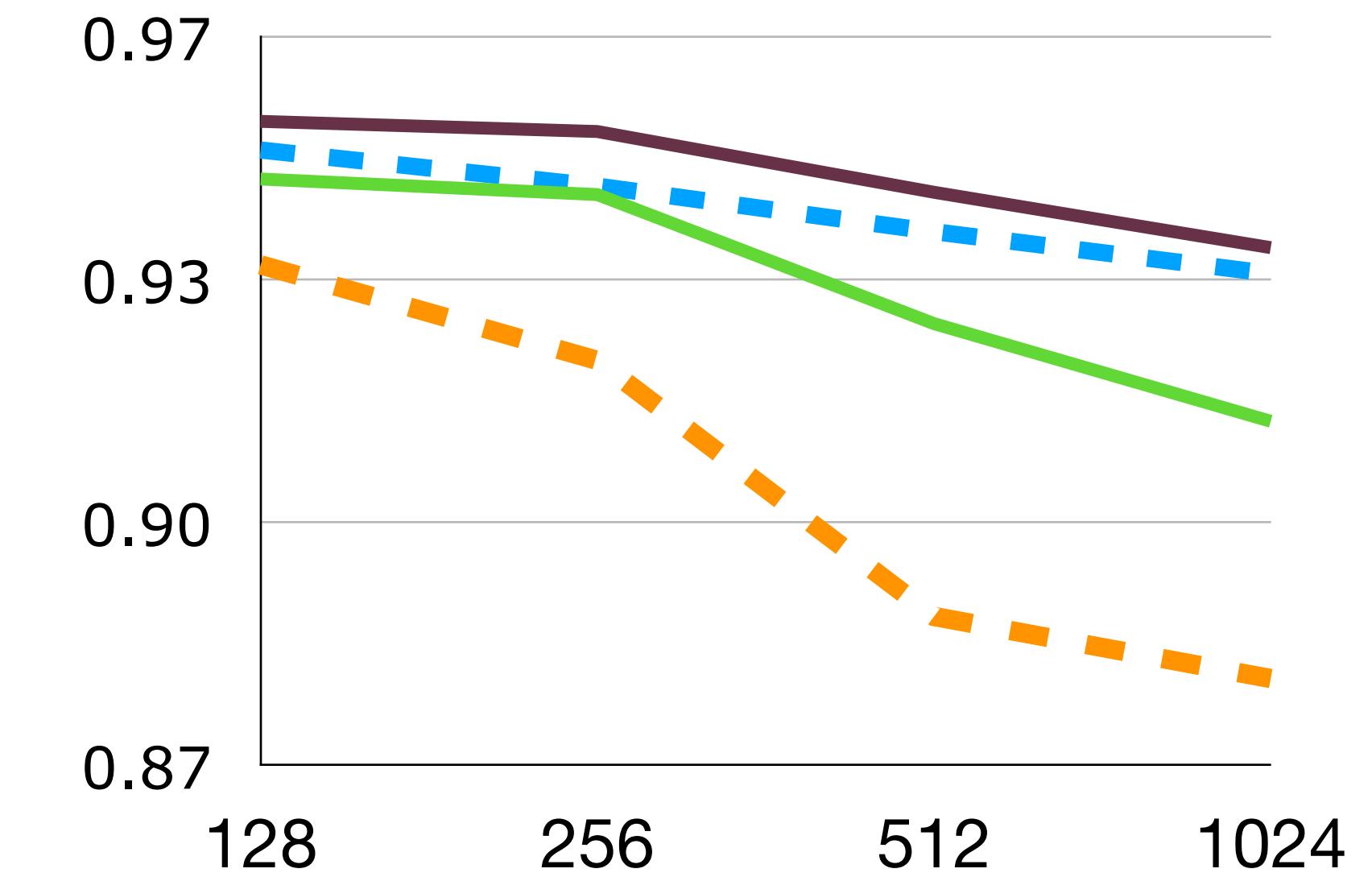


- Y-axis shows Mauve (\uparrow)
- **Mauve** captures all the trends while baselines fail

Decoding Algorithm



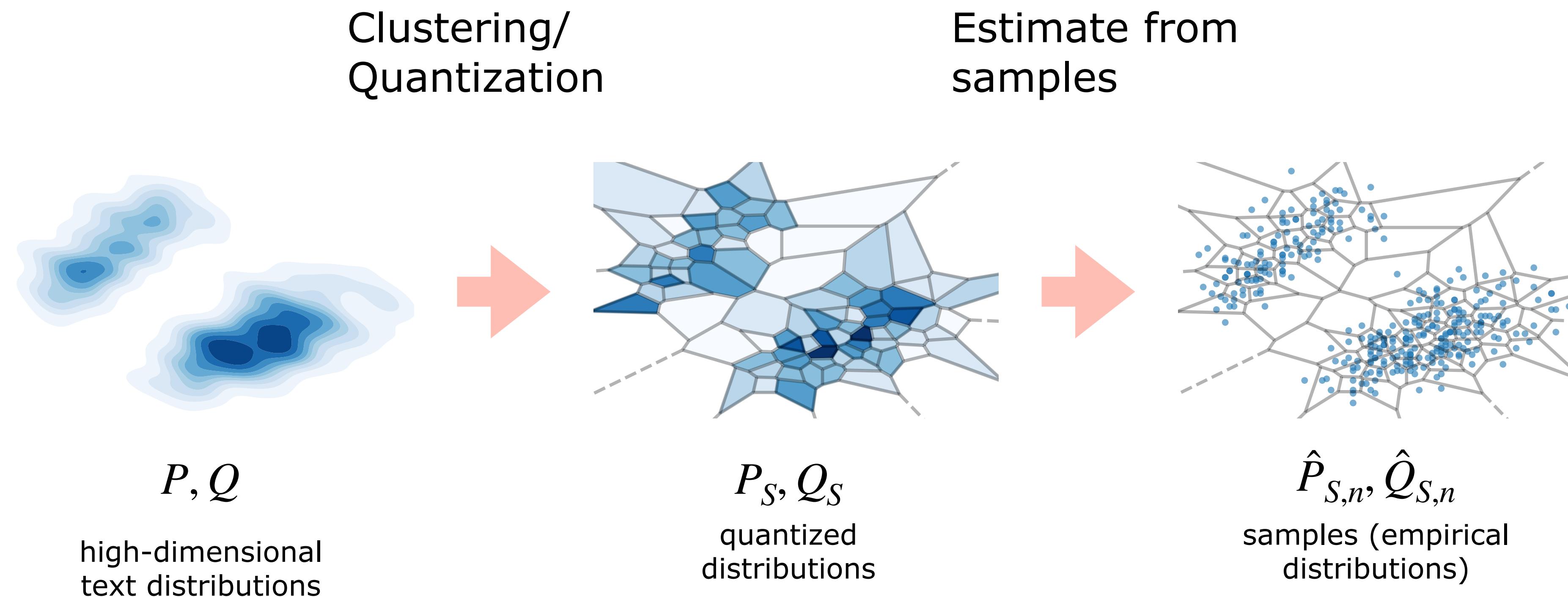
Text Length



- Small
- Large
- Medium
- XL

Mauve: Estimation theory

Estimation of Mauve involves two approximations:



Theorem (informal)

There exists a quantization of size k such that the approximation error of **Mauve** from n samples is

$$\tilde{o}\left(\sqrt{\frac{k}{n}} + \frac{1}{k} \right)$$

Statistical
Error

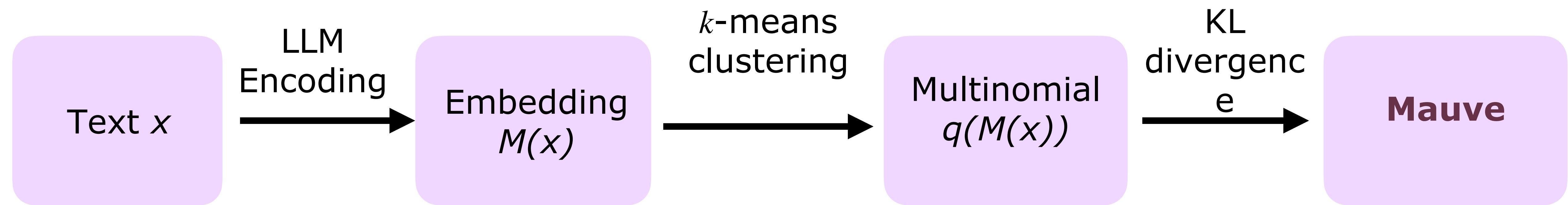
Quantization
Error

n : number of samples from P and Q

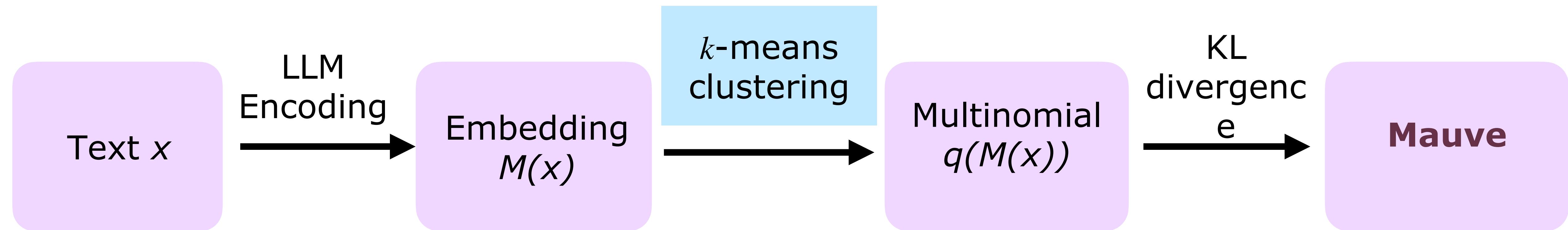
k : quantization size (Num. clusters)

Balance both by choosing $k = \Theta(n^{1/3})$

Mauve: Beyond clustering



Mauve: Beyond clustering



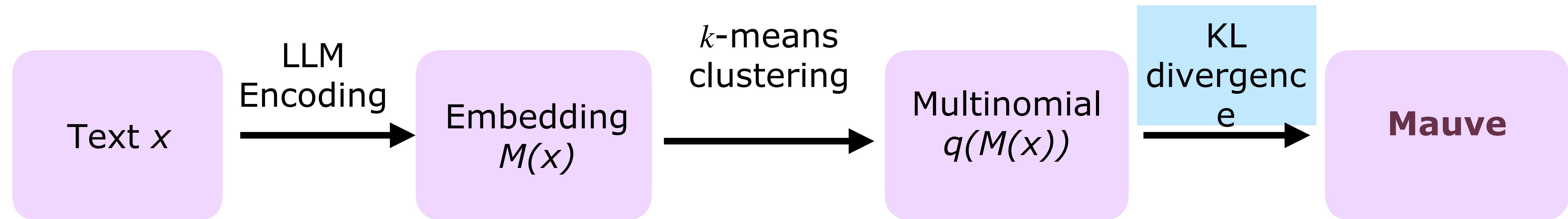
Nearest neighbor estimator

Kernel density estimator

Parametric Gaussian approx.

Classifier-based estimation

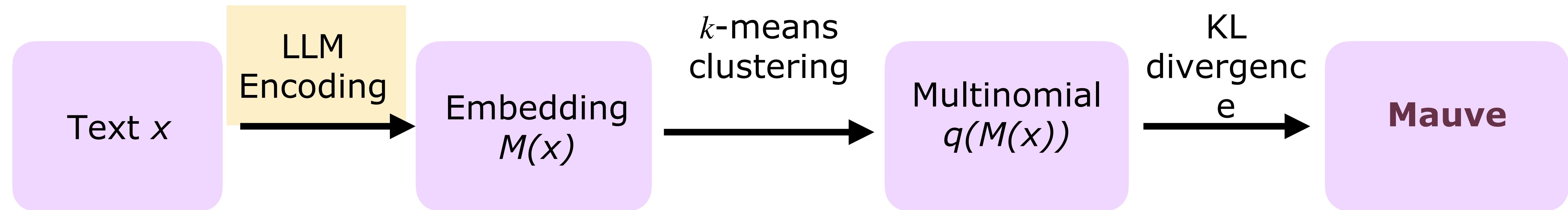
Mauve: Beyond clustering



General f -divergences

Optimal transport

Mauve: Beyond clustering



Essential

Software

pip install mauve-text

```
from mauve import compute_mauve

p_text = ... # list of strings for human distribution P
q_text = ... # list of strings for model distribution Q

# Obtain deep encoding, quantize it and compute Mauve
out = compute_mauve(p_text=p_text, q_text=q_text)

print(f'Mauve(P, Q) = {out.mauve}')
```

60K downloads

HuggingFace Evaluate: **pip install evaluate**

```
from evaluate import mauve

mauve = load("mauve")

# Obtain deep encoding, quantize it and compute Mauve
out = mauve.compute(references=p_text, predictions=q_text)

print(f'Mauve(P, Q) = {out.mauve}')
```

Impact of Mauve

Standard metric for evaluation and hyper-parameter tuning

Meister et al. (TACL 2022)

Jawahar et al. (ACL 2022)

Su et al. (NeurIPS 2022)

Hewitt et al. (EMNLP 2022)

Lu et al. (NeurIPS 2022)

Mattern et al. (EMNLP 2022)

Xu et al. (NeurIPS 2022)

Hu et al. (NAACL 2022)

Summary: Diagnosing large-scale text generation with Mauve

Used for open-ended generation
but trained for language
modeling. How good is it?

>> **prompt:** In a shocking finding, scientists discovered a herd
of unicorns living in a remote, previously,
unexplored valley, in the Andes Mountains.



Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown ...



Continuation 2. This discovery has kicked off an all-out search for other mythical creatures from the frozen reaches of the Antarctic to the tropical islands of the Pacific ...



Continuation 3. Perhaps most astonishingly, these unicorns have developed their own artificial general intelligence named Yuyaysapa ...

Summary: Diagnosing large-scale text generation with Mauve

Used for open-ended generation
but trained for language
modeling. How good is it?

Directly measure the gap
between distributions

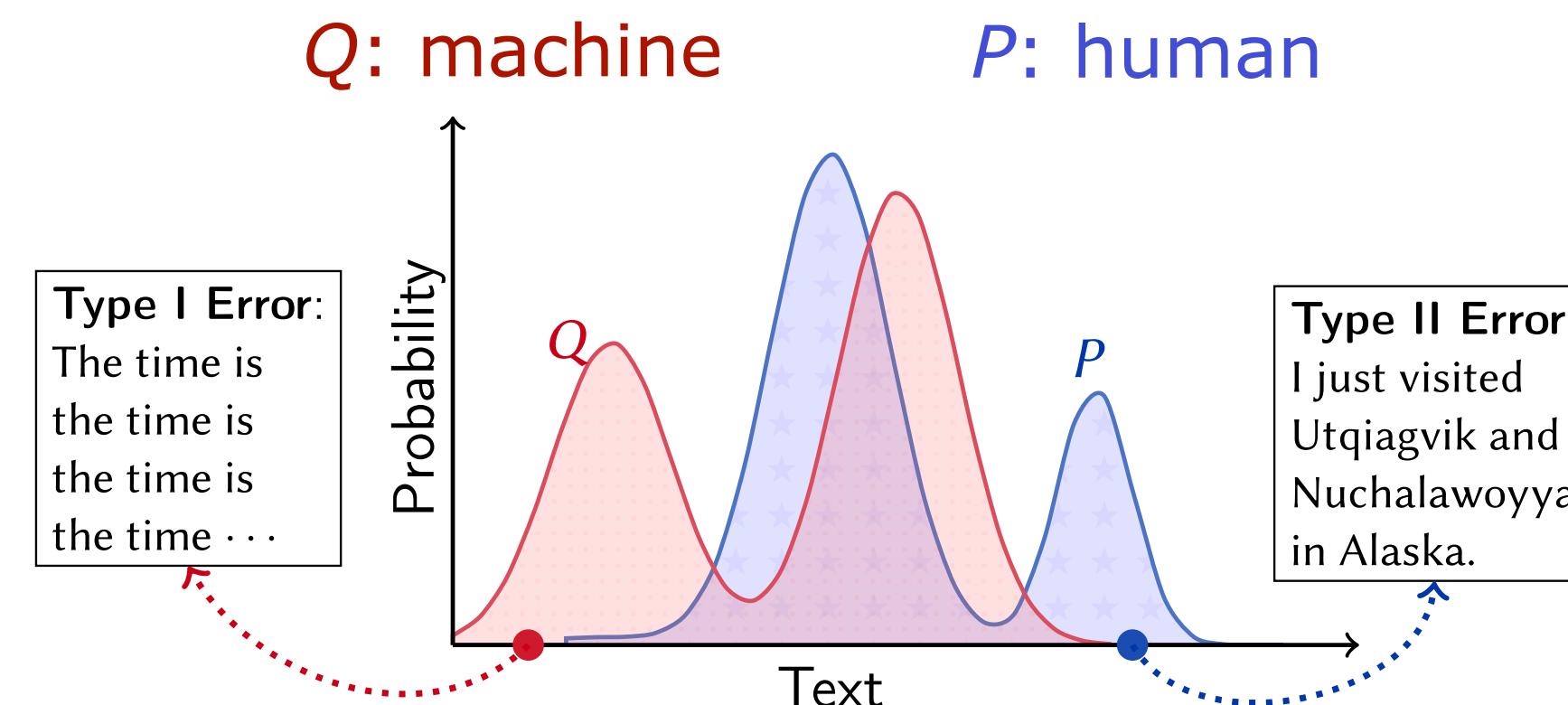
>> prompt: In a shocking finding, scientists discovered a herd
of unicorns living in a remote, previously,
unexplored valley, in the Andes Mountains.



Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown ...

Continuation 2. This discovery has kicked off an all-out search for other mythical creatures from the frozen reaches of the Antarctic to the tropical islands of the Pacific ...

Continuation 3. Perhaps most astonishingly, these unicorns have developed their own artificial general intelligence named Yuyaysapa ...



Summary: Diagnosing large-scale text generation with Mauve

Used for open-ended generation
but trained for language
modeling. How good is it?

Directly measure the gap
between distributions

Our approach correlates
with human judgements and
quantifies observed behavior

>> prompt: In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.



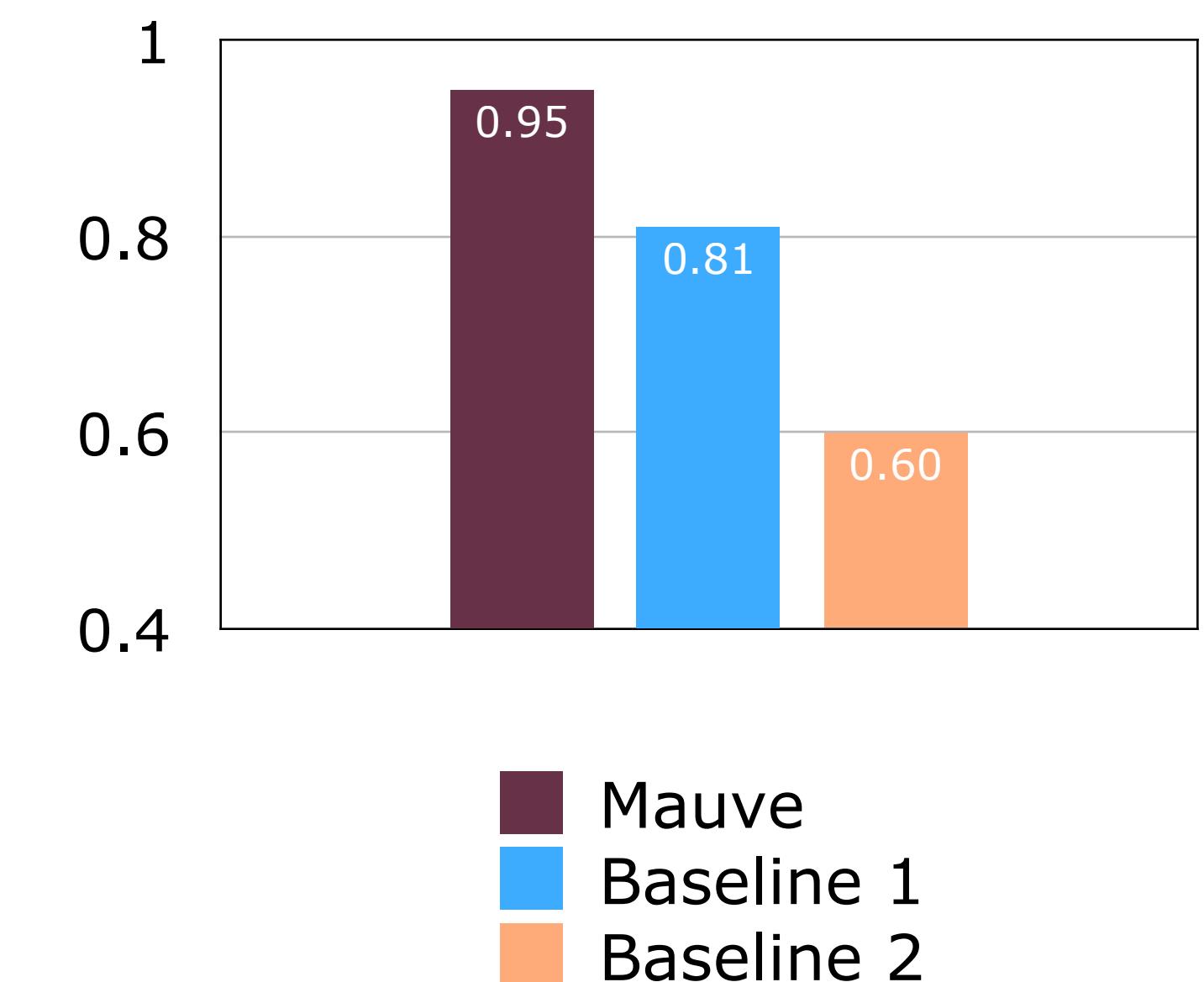
Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown ...

Continuation 2. This discovery has kicked off an all-out search for other mythical creatures from the frozen reaches of the Antarctic to the tropical islands of the Pacific ...

Continuation 3. Perhaps most astonishingly, these unicorns have developed their own artificial general intelligence named Yuyaysapa ...



Spearman Correlation w/
human eval (\uparrow)



Summary: Diagnosing large-scale text generation with Mauve

Used for open-ended generation
but trained for language
modeling. How good is it?

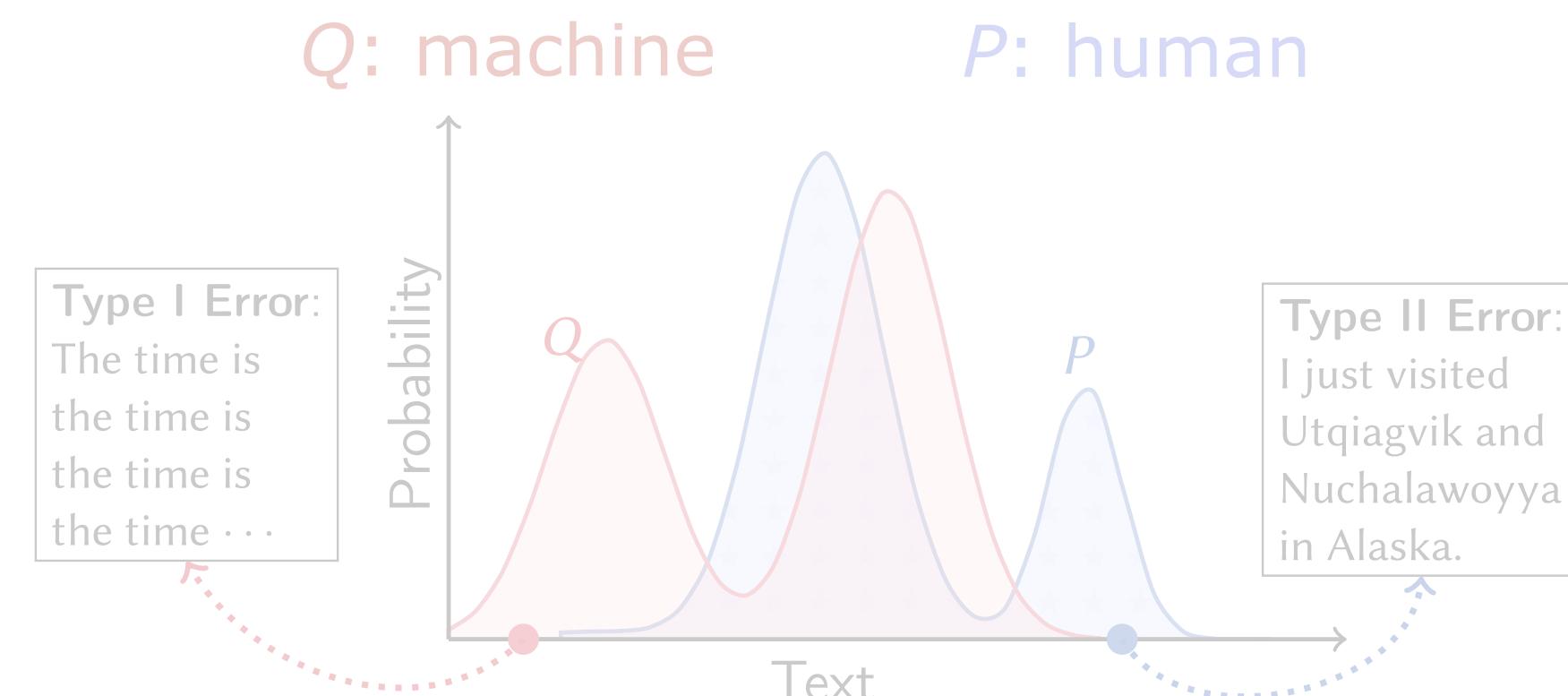
Directly measure the gap
between distributions

Theory: error bounds

$$\tilde{\sigma} \left(\sqrt{\frac{k}{n}} + \frac{1}{k} \right)$$

Statistical
Error

Quantization
Error



>> prompt: In a shocking finding, scientists discovered a herd of unicorns living in a remote, previously, unexplored valley, in the Andes Mountains.



Continuation. The scientists named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown ...



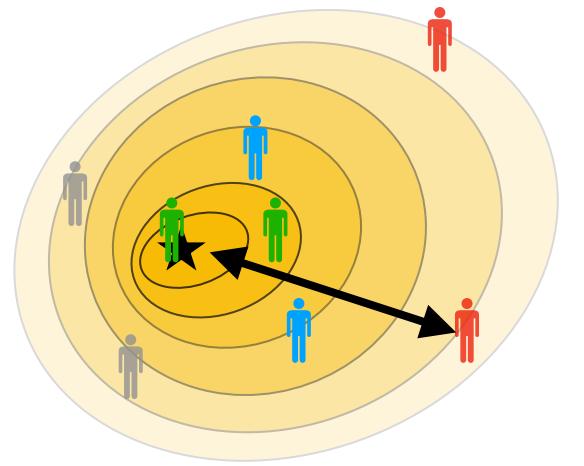
Continuation 2. This discovery has kicked off an all-out search for other mythical creatures from the frozen reaches of the Antarctic to the tropical islands of the Pacific ...



Continuation 3. Perhaps most astonishingly, these unicorns have developed their own artificial general intelligence named Yuyaysapa ...

Federated learning

LLMs



IEEE CISS 2021,
Springer SVVA 2021,
Mach. Learn. 2022

NeurIPS 2021a
NeurIPS 2021b
Submitted 2023

Part 2

IEEE Trans. Signal Proc. 2022,
ICML 2022

Submitted 2022



NeurIPS 2018
Submitted 2022



Robust Deployment

Robust to Outliers

Optimize Faster

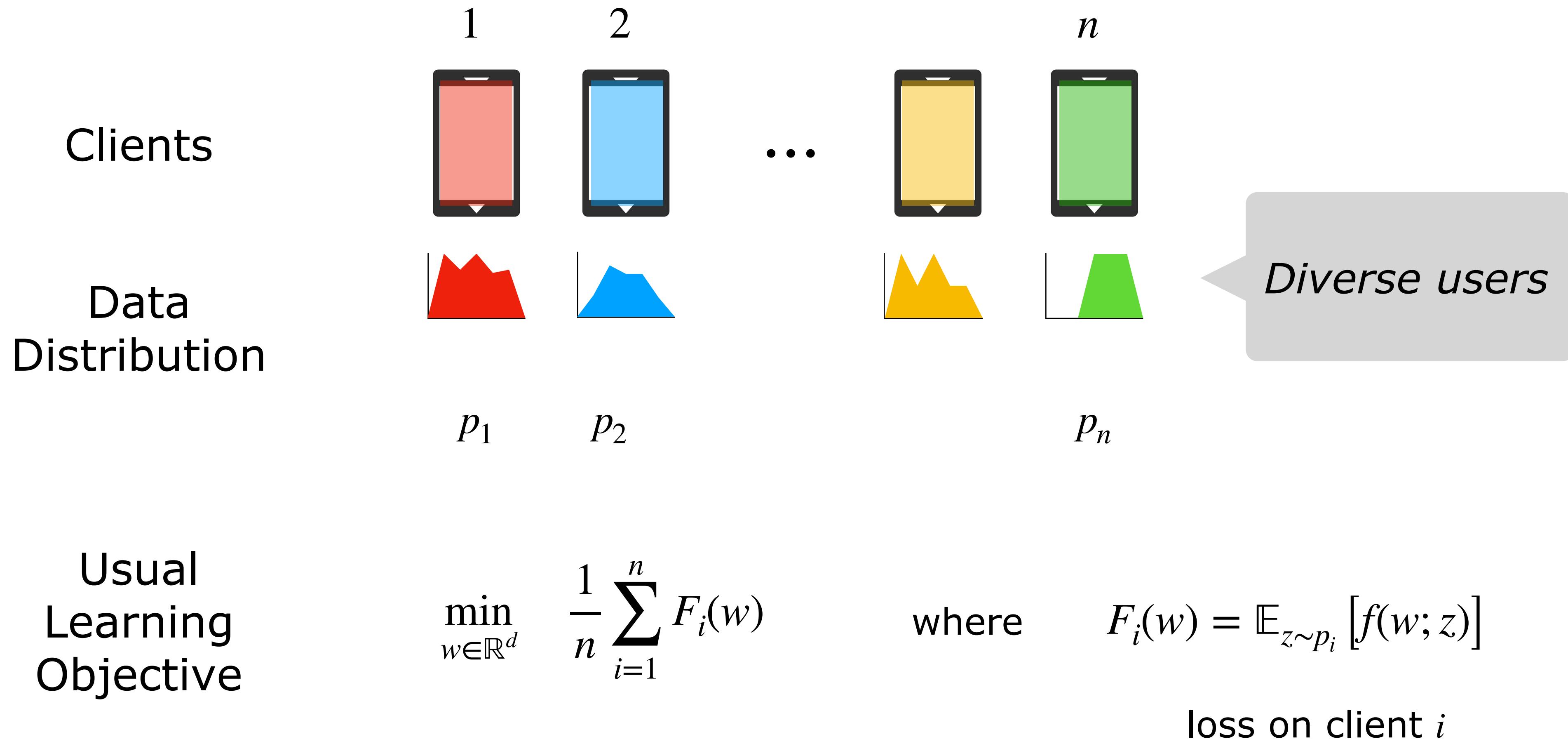
Privacy



Part 2: Tackling distribution shifts in federated learning

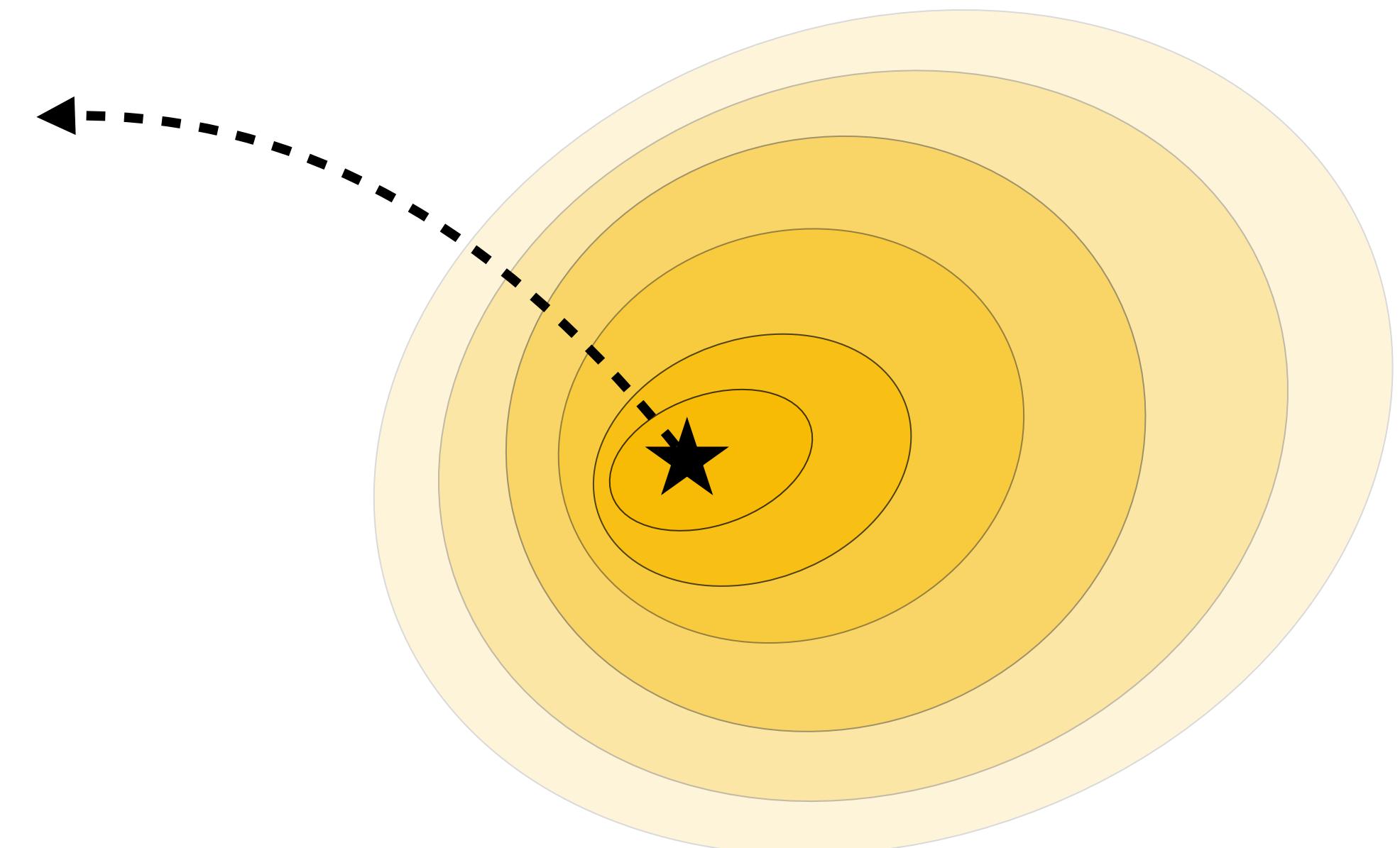
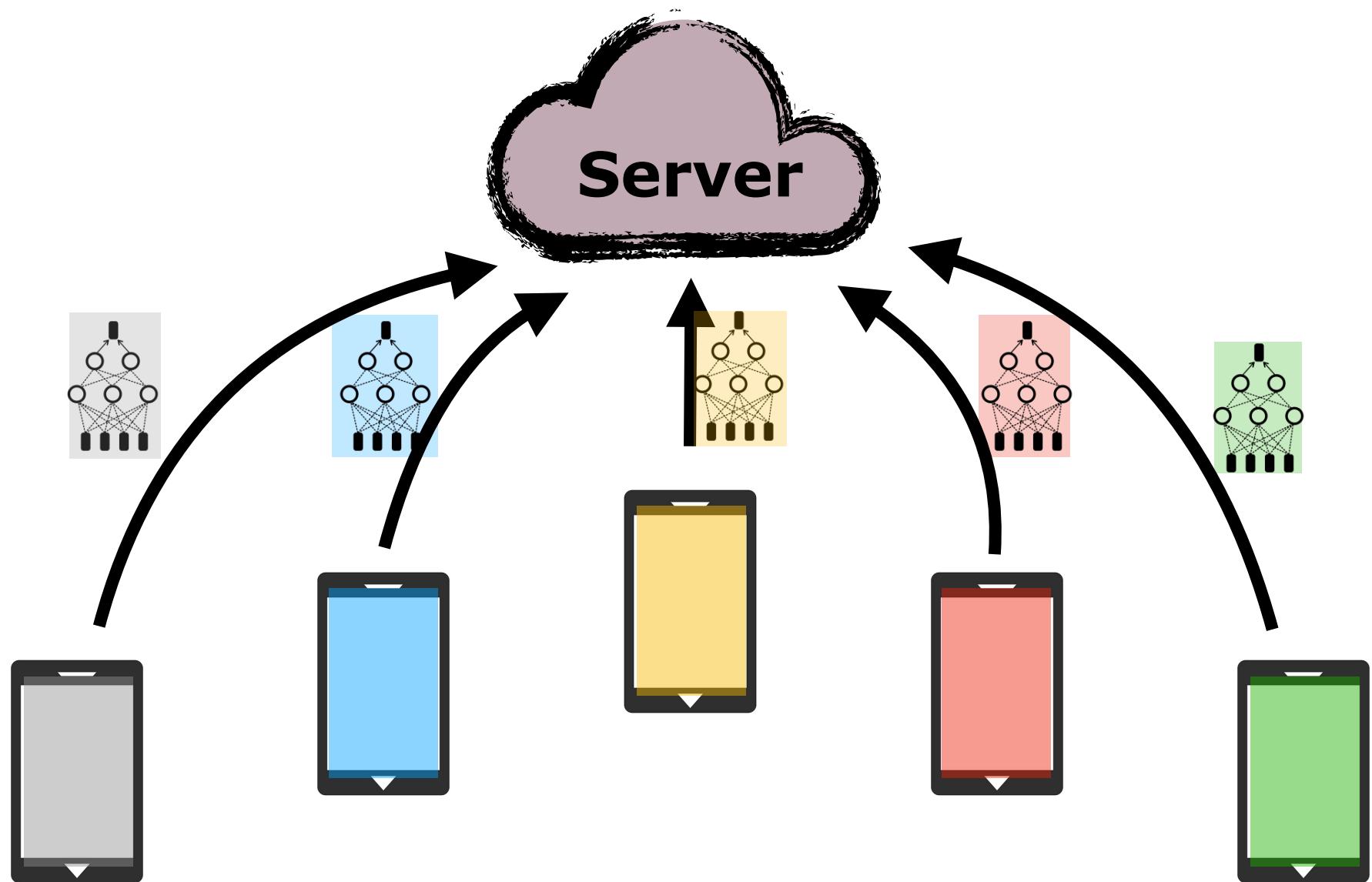
[IEEE CISS '21, DistShift-NeurIPS '22 (Spotlight),
SVVA '21, *Mach. Learn.* '22]



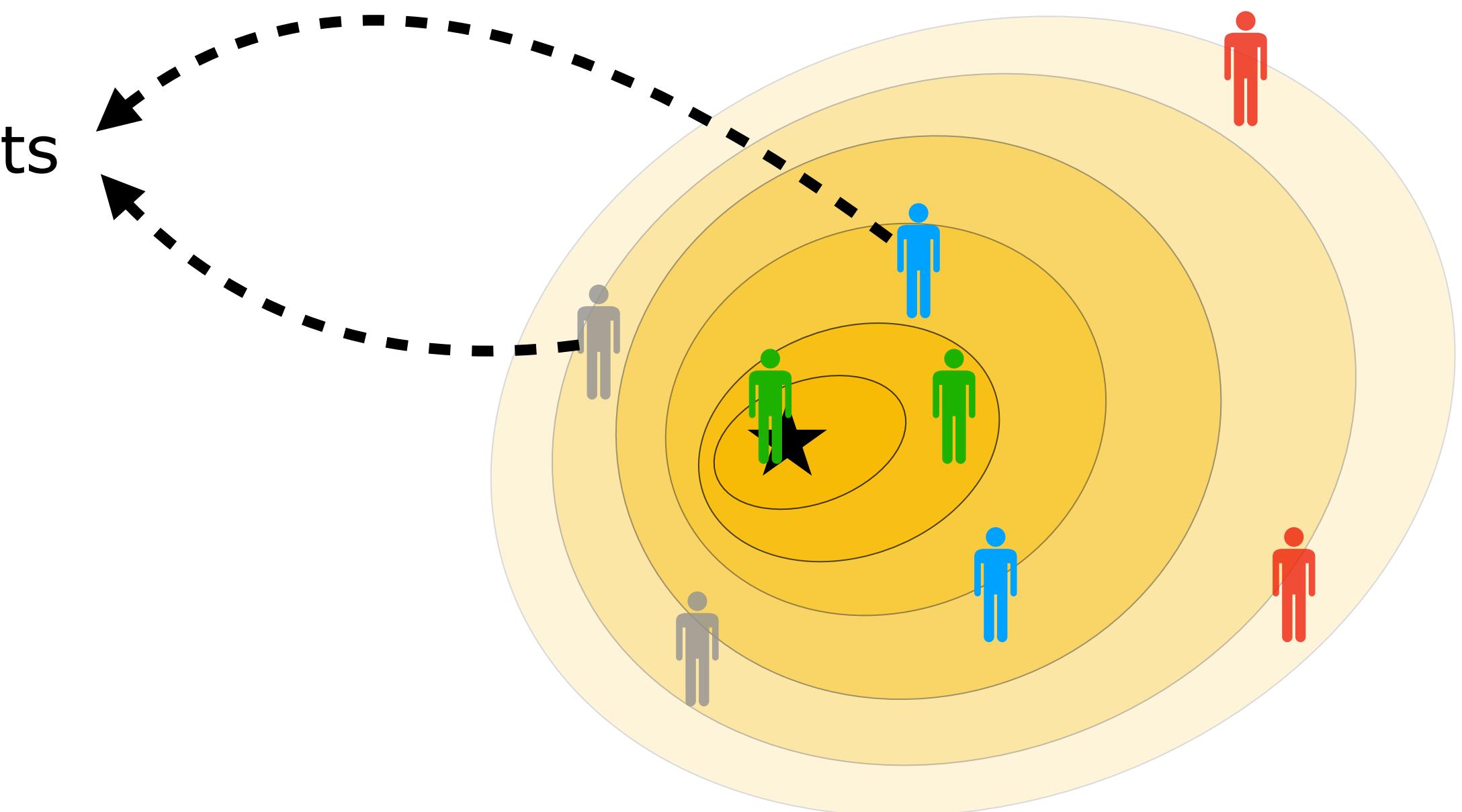
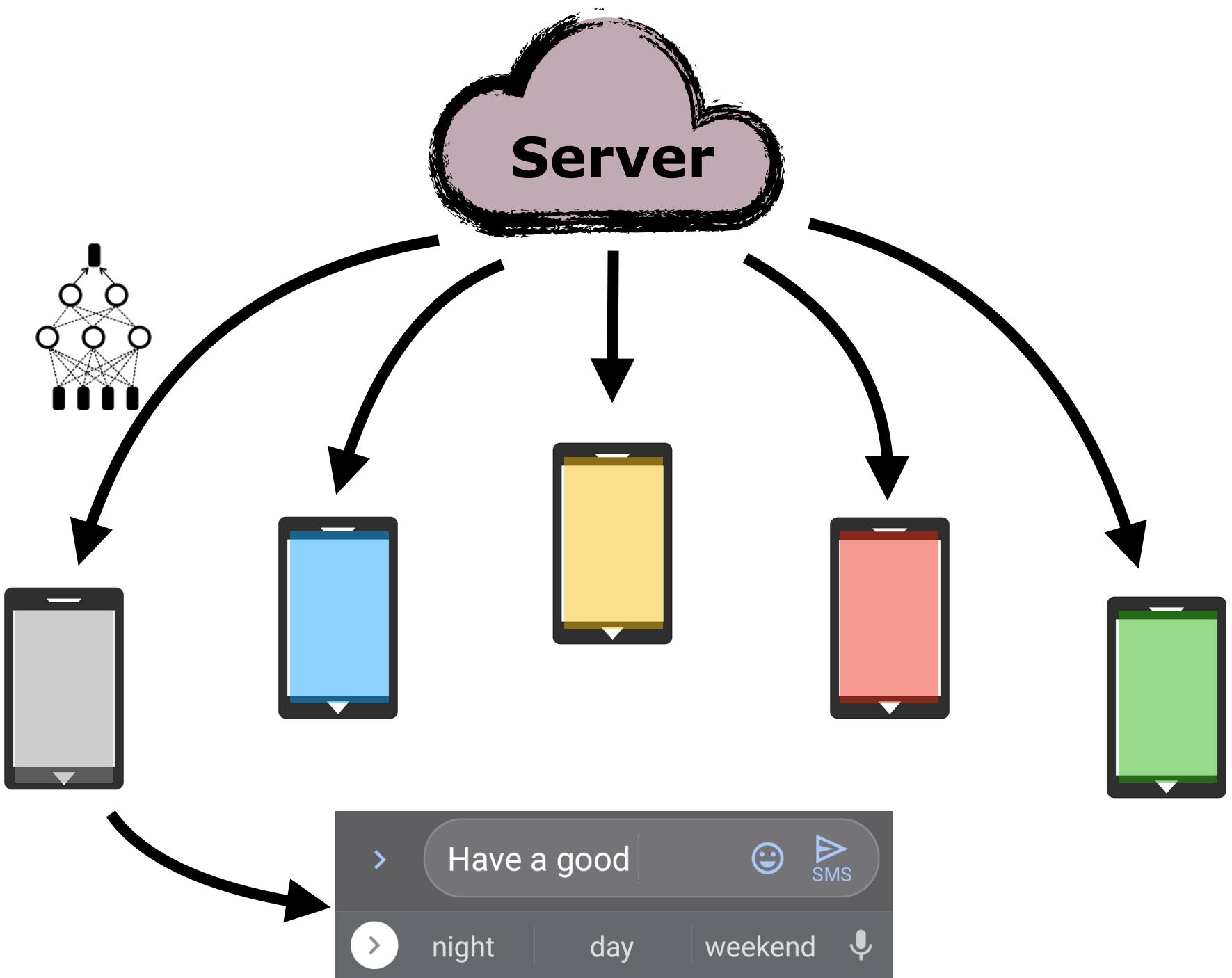


[McMahan et al. (AISTATS 2017), Kairouz et al. (2021)]

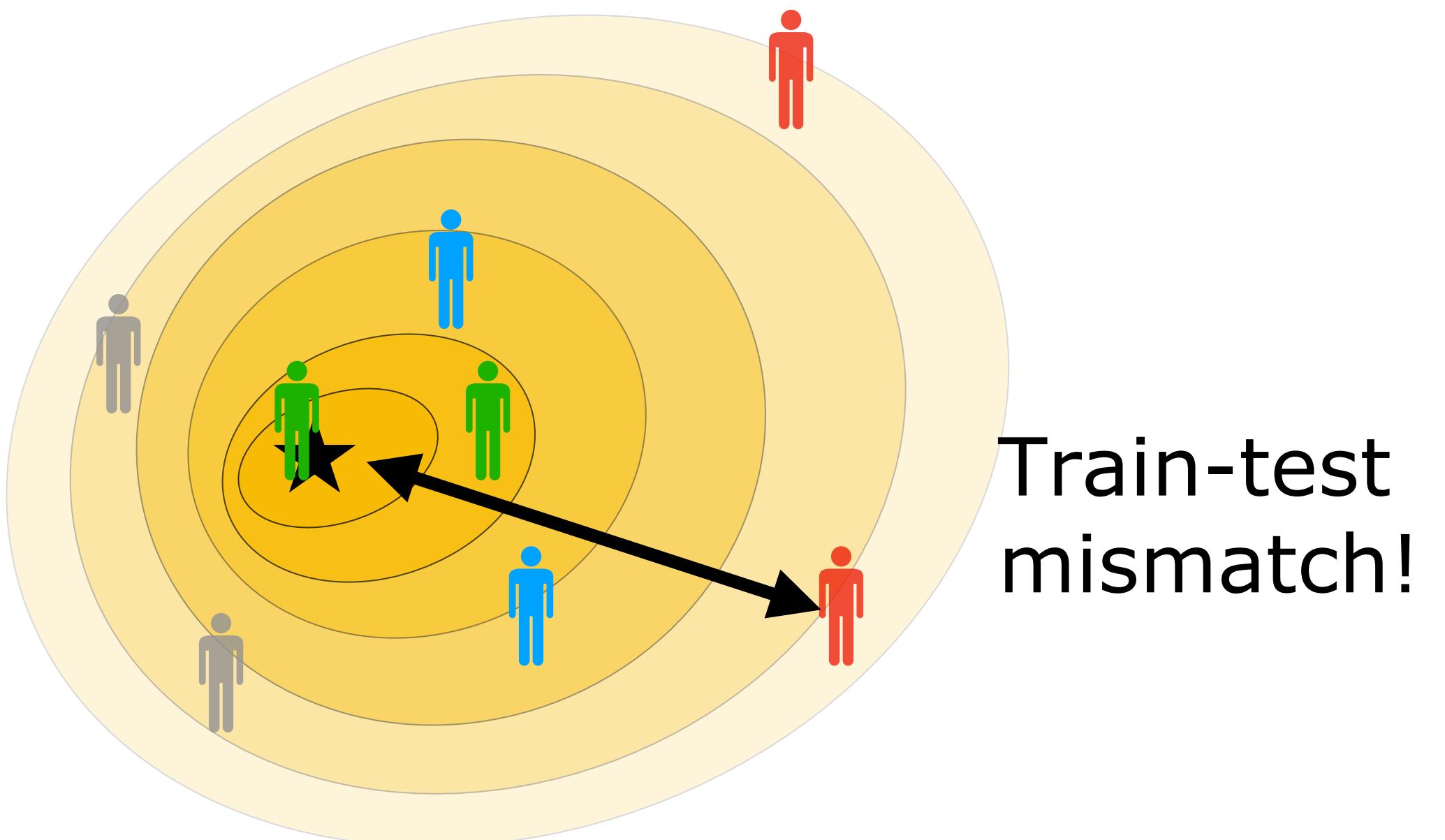
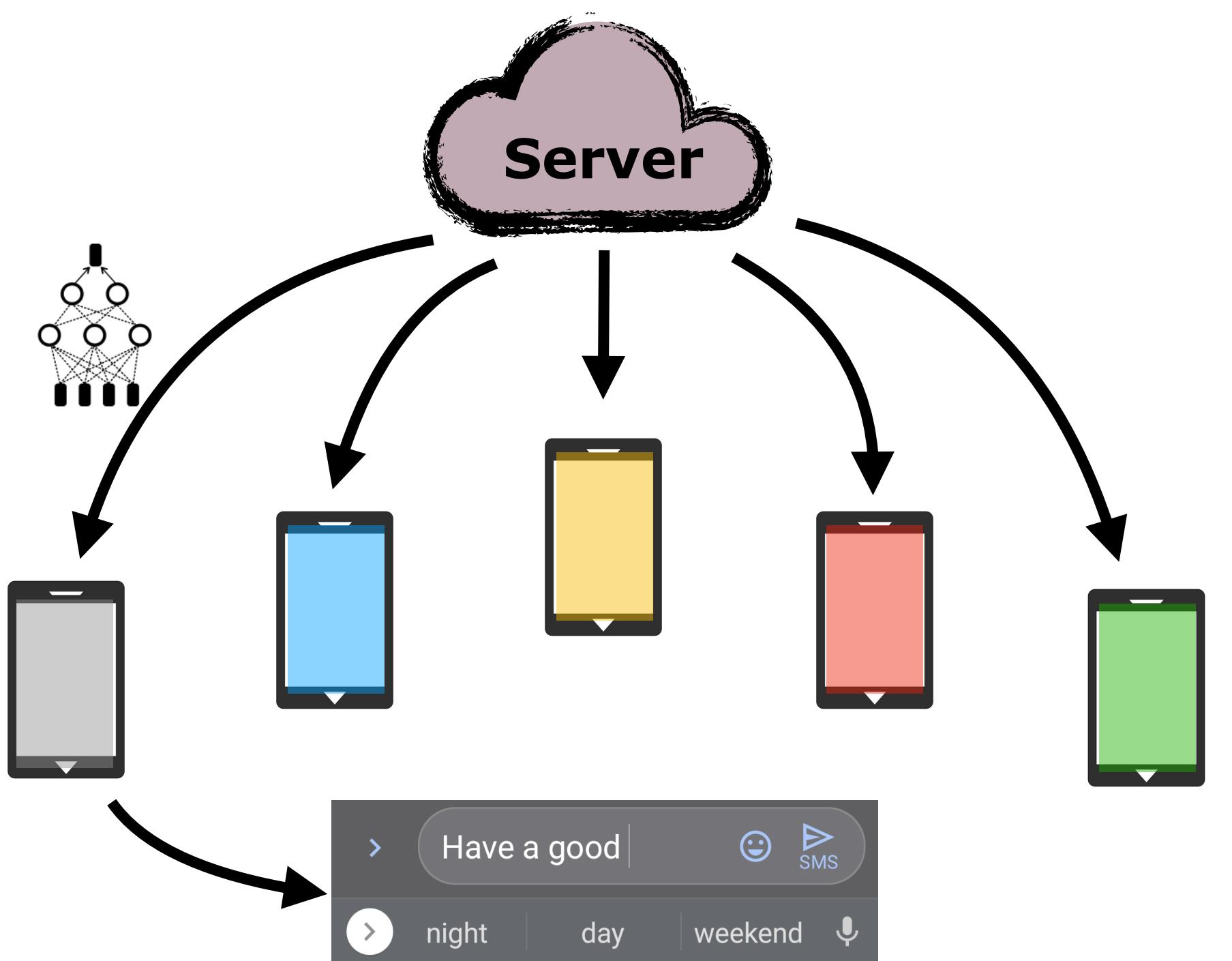
Usual approach \Rightarrow Global model is trained on
average distribution across clients



Global model is deployed on *individual* clients

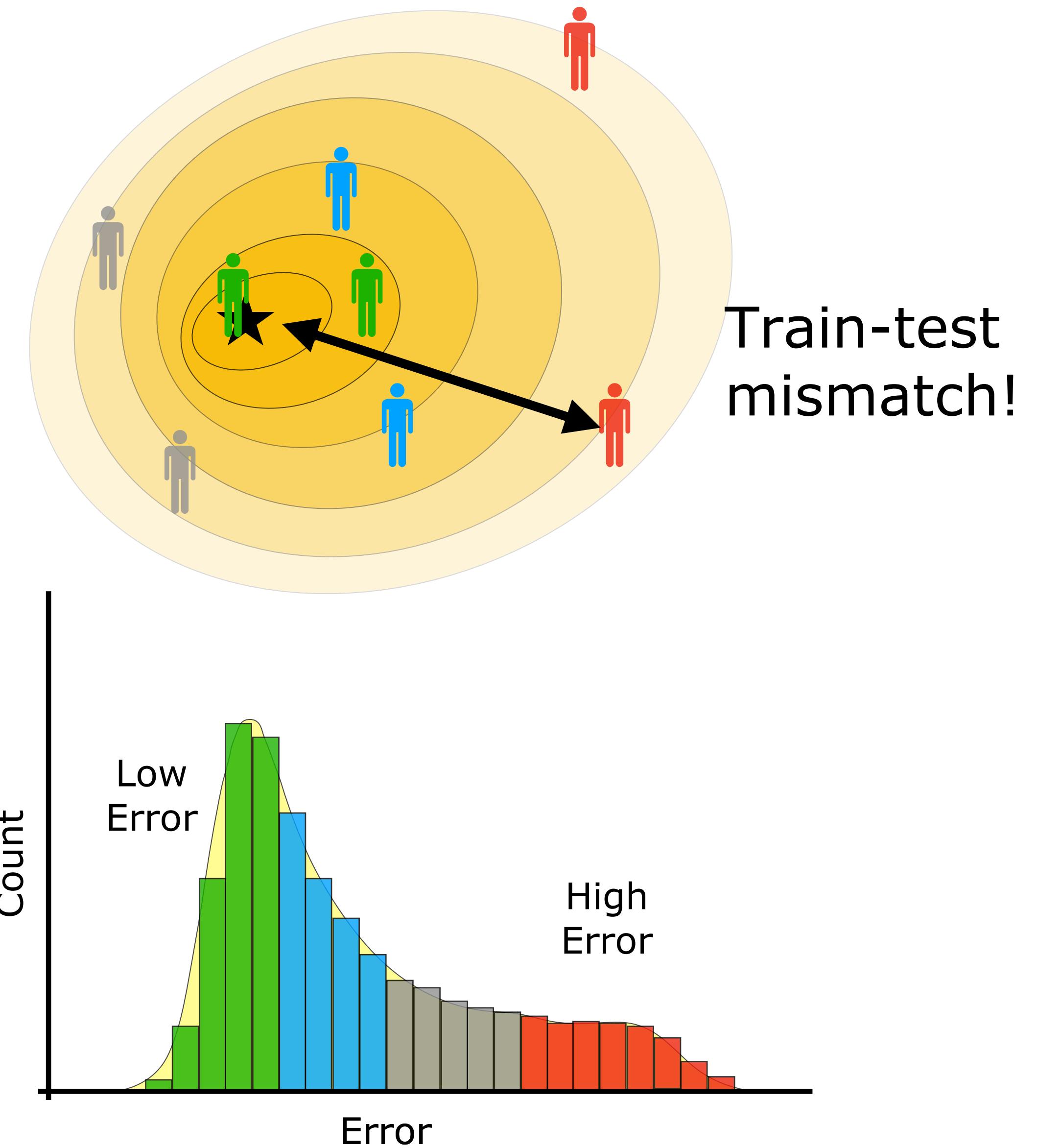
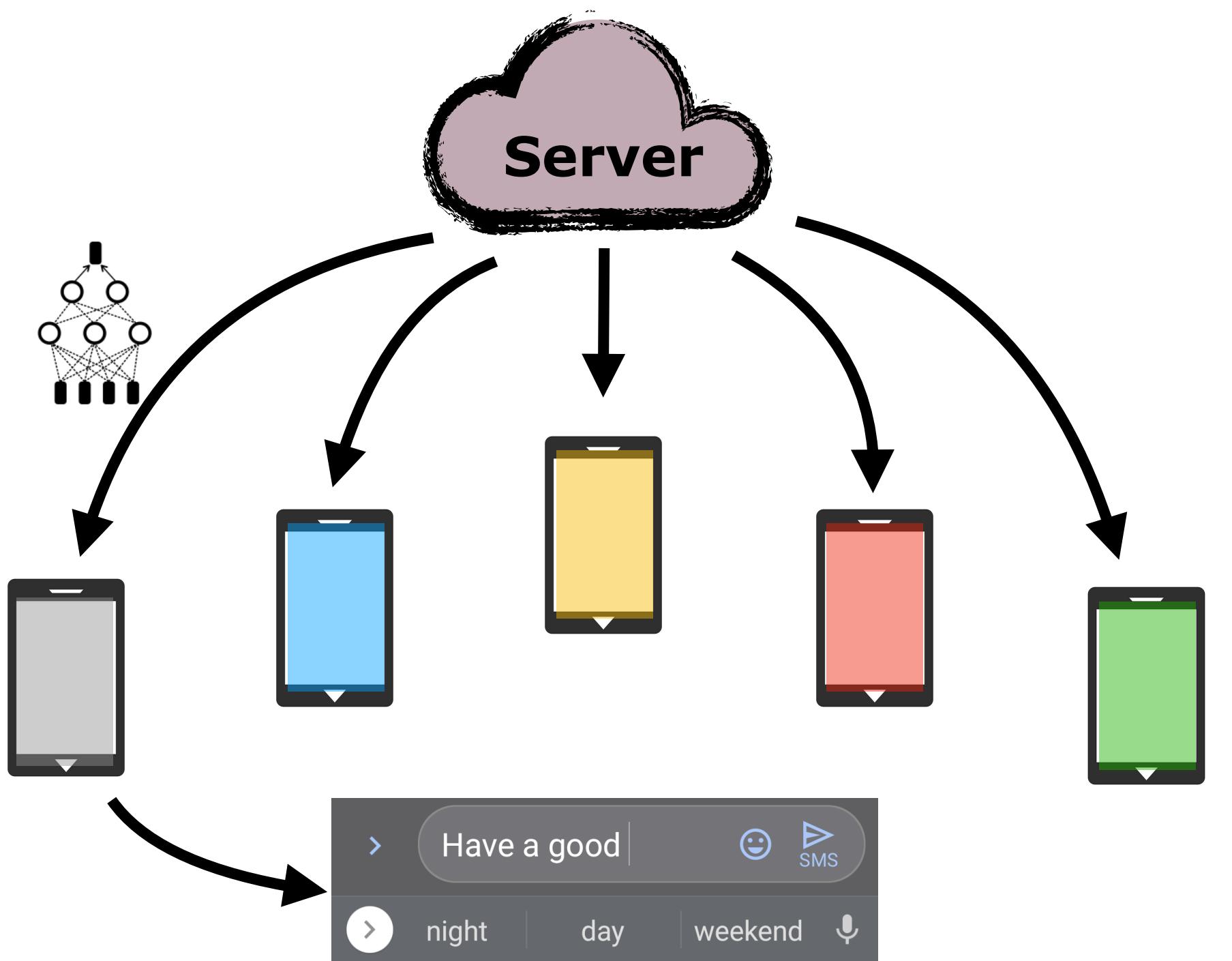


Global model is deployed on *individual* clients



Train-test
mismatch!

Global model is deployed on *individual* clients



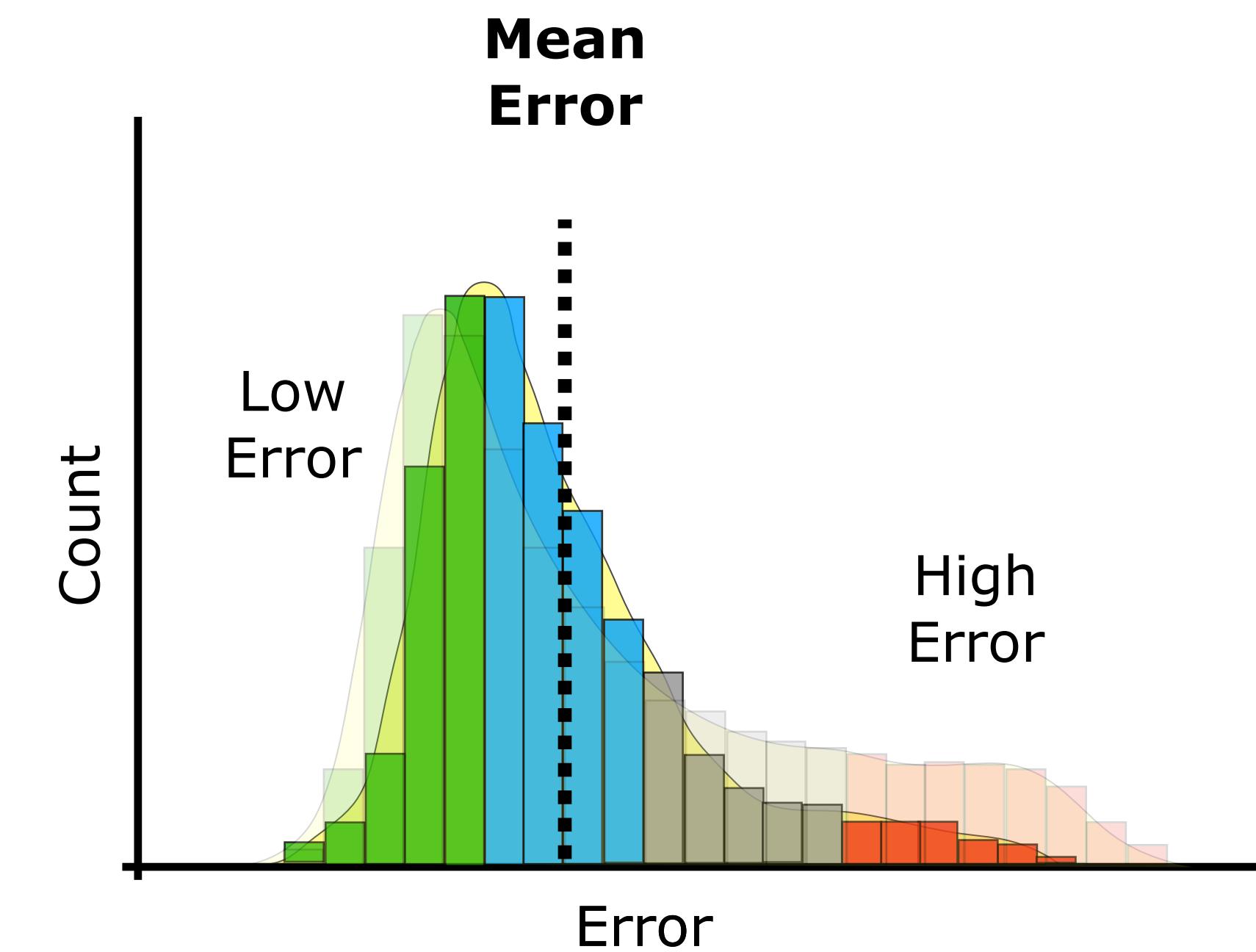
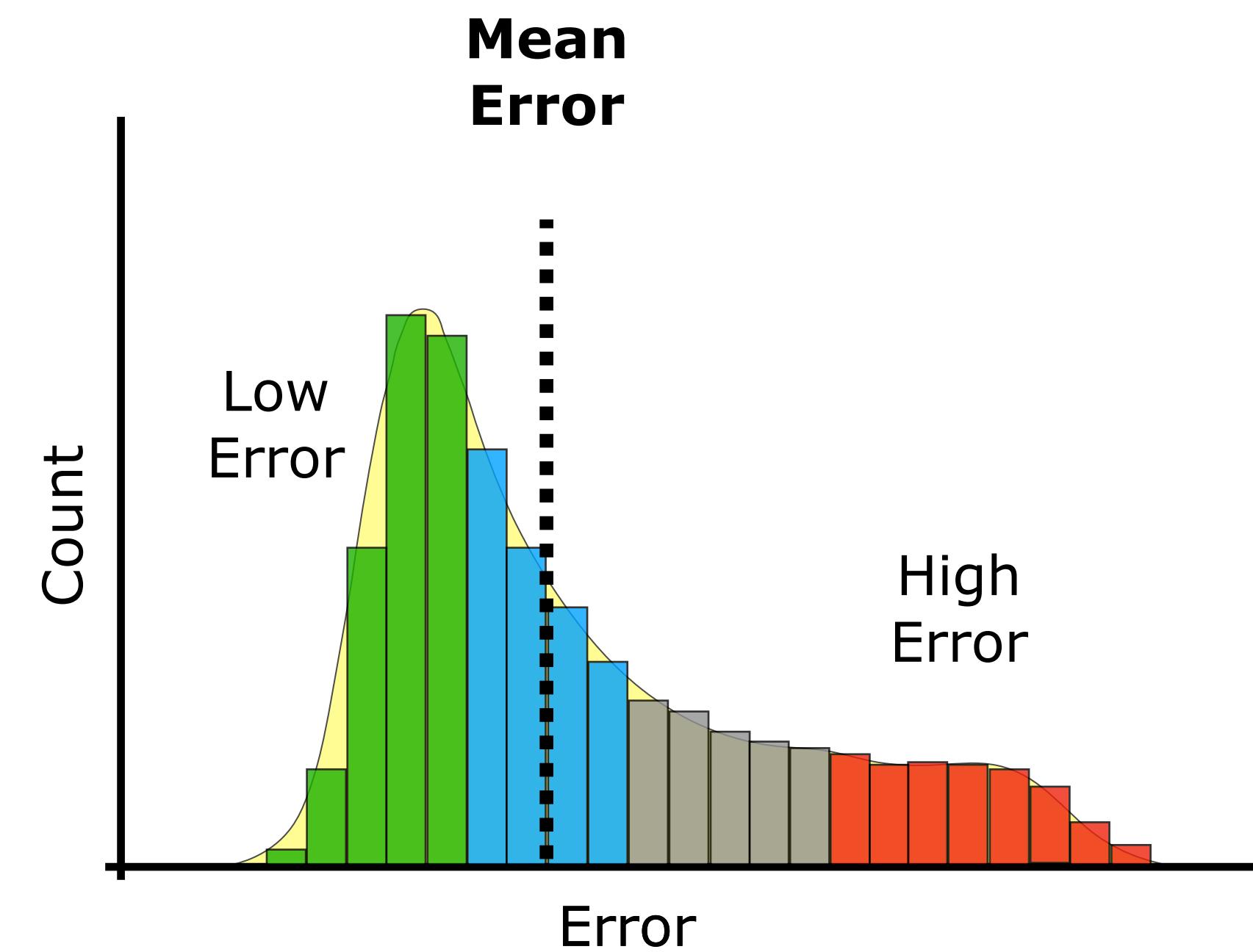
Our goal

Reduce *tail* error



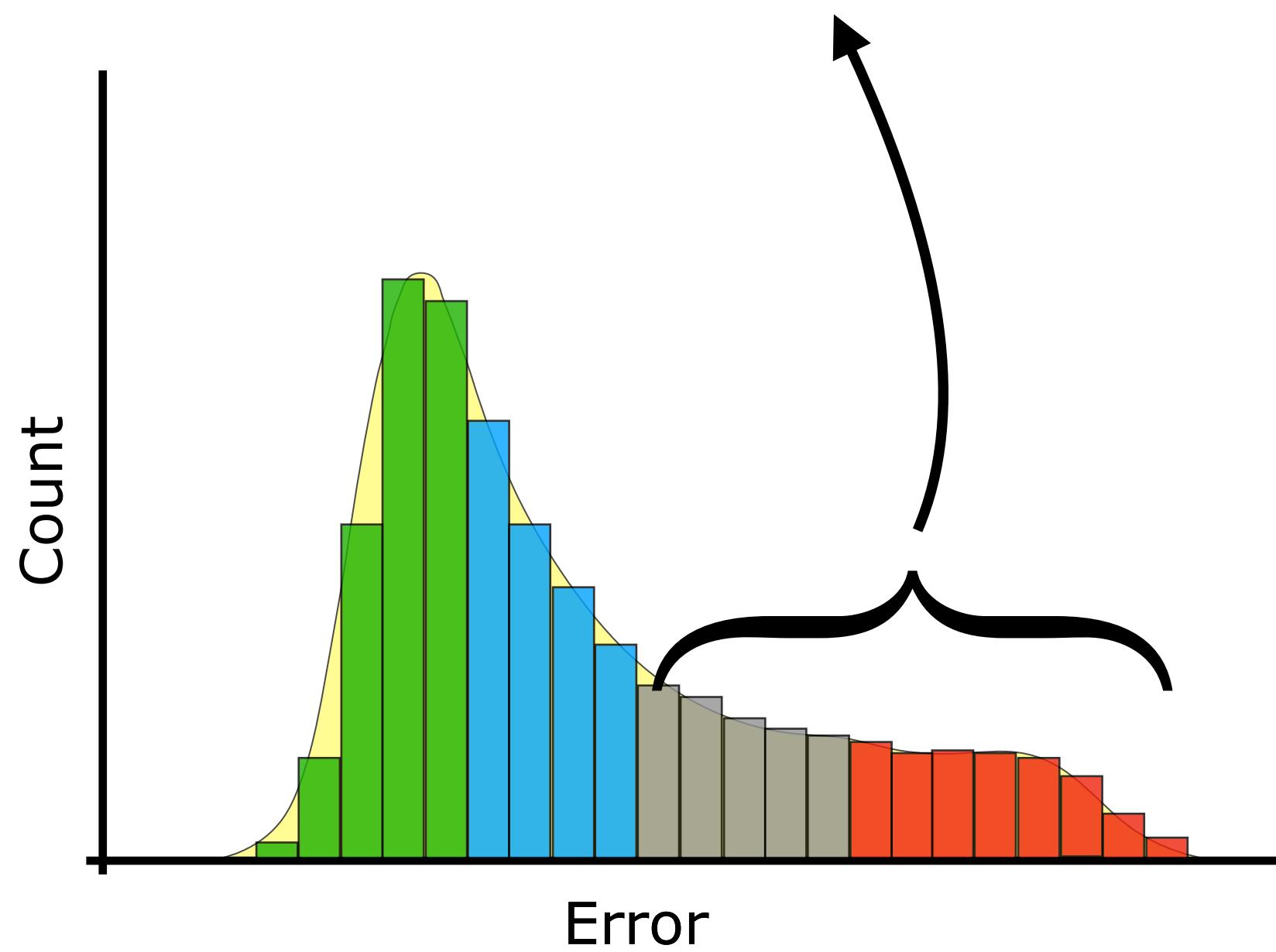
Our goal

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



Simplicial federated learning

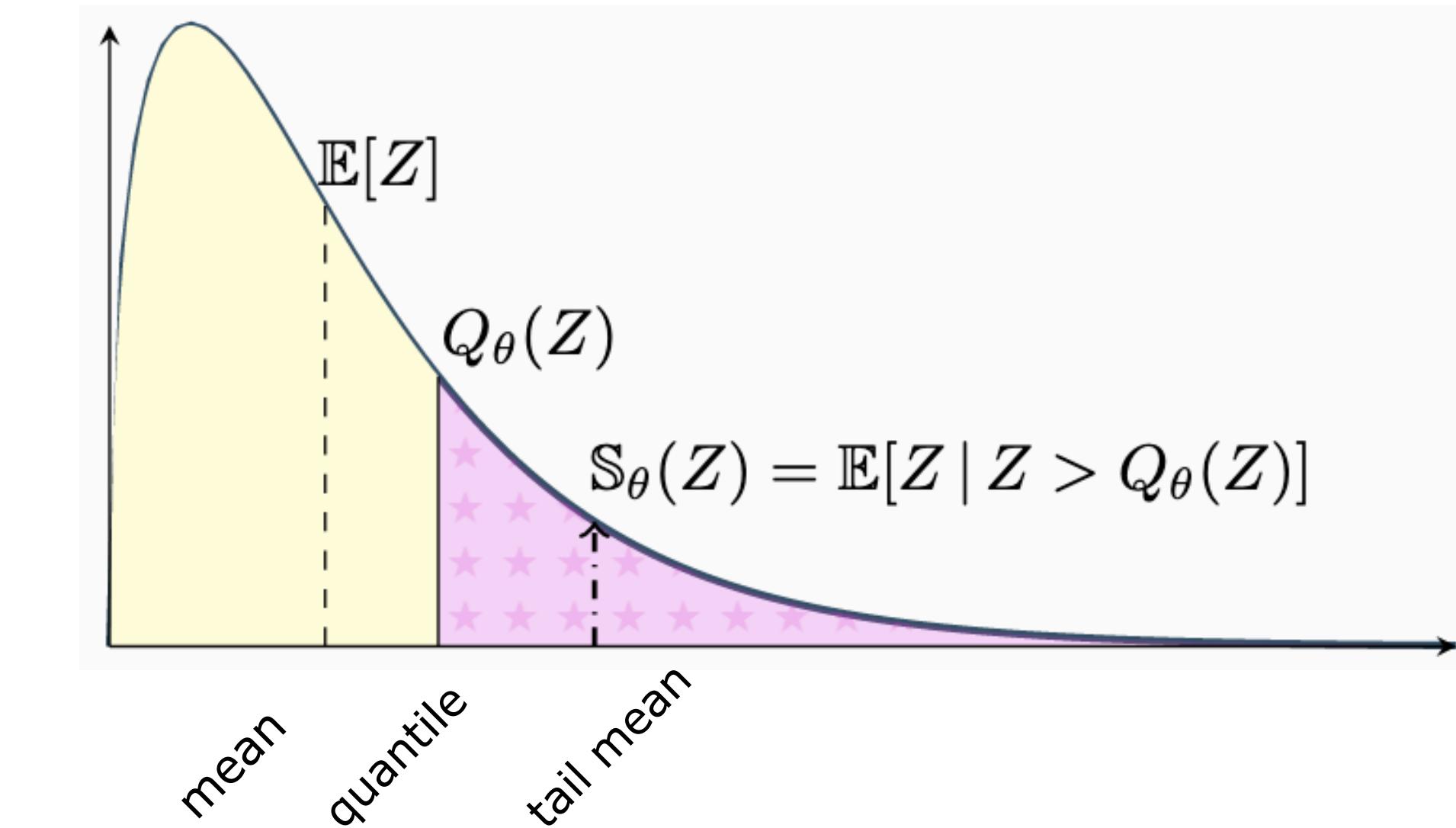
Our Approach: minimize the tail error directly!



Simplicial-FL Objective:

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

Superquantile | Conditional Value at Risk

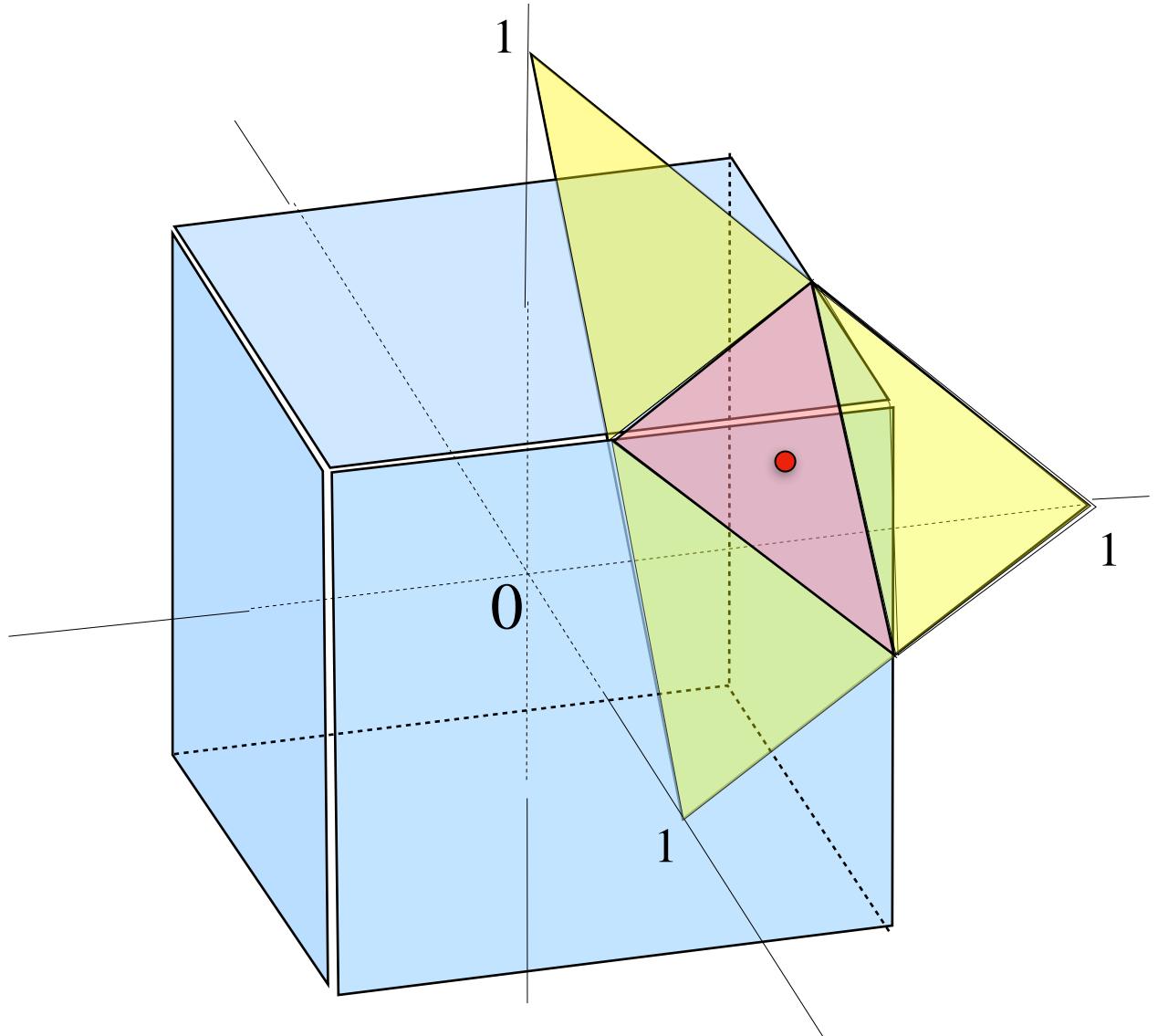


[Rockafellar & Uryasev (2002)]

Dual expression \equiv continuous knapsack problem

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

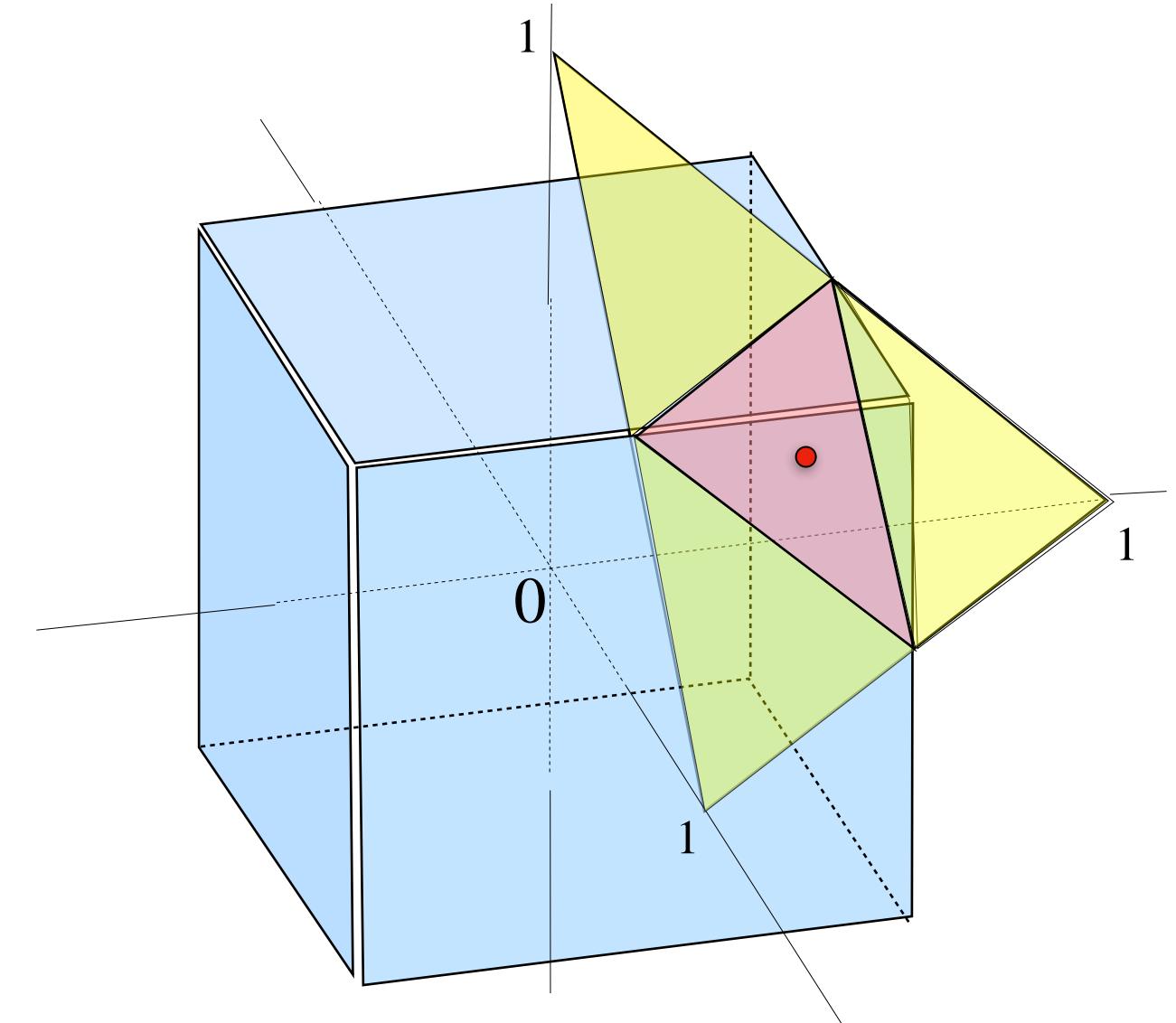
[Dantzig (1957), Ben-Tal & Teboulle (1987), Föllmer & Schied (2002)]



Dual expression \equiv continuous knapsack problem

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

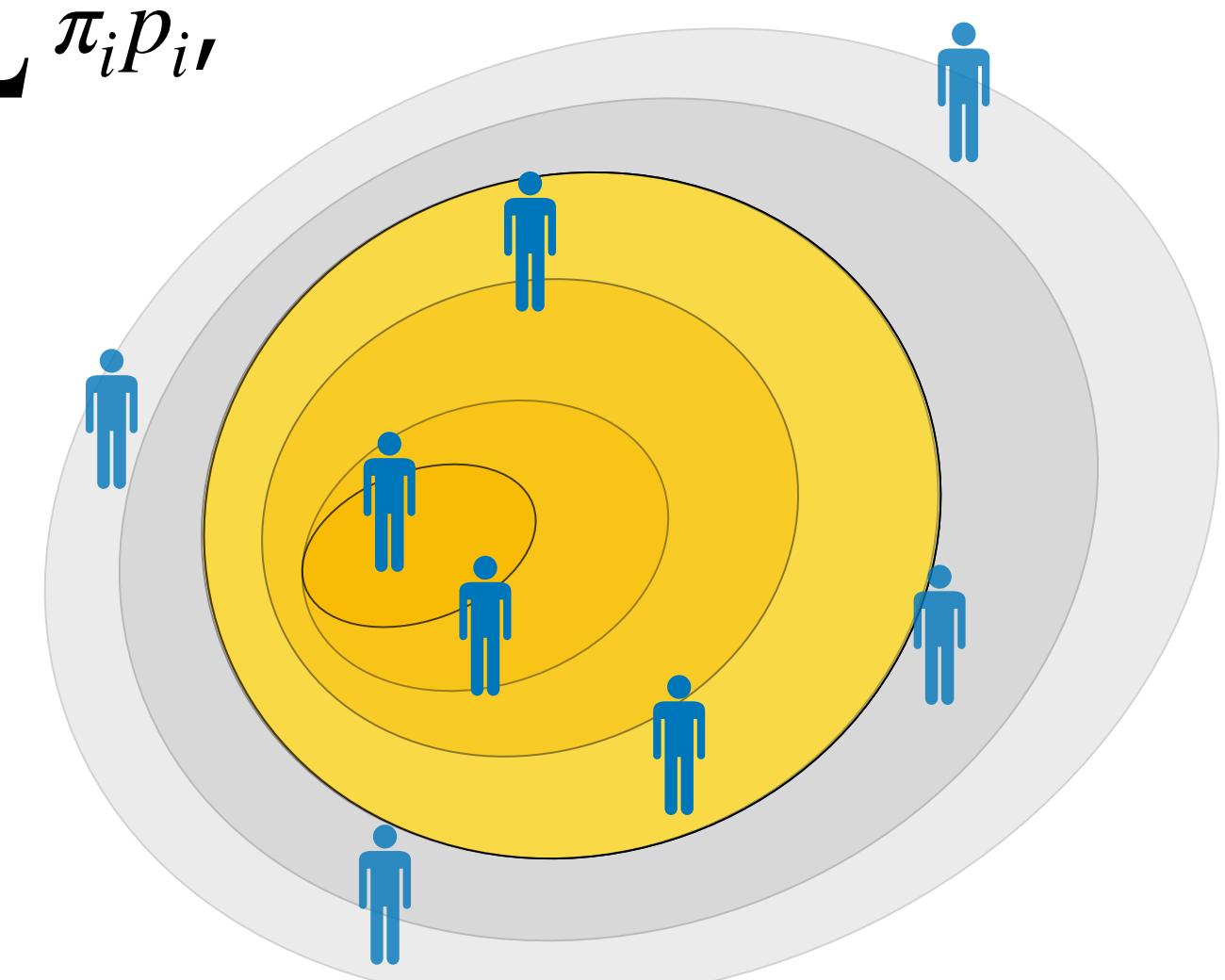
[Dantzig (1957), Ben-Tal & Teboulle (1987), Föllmer & Schied (2002)]



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



Usual Algorithm (FedAvg):

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

Our 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)]

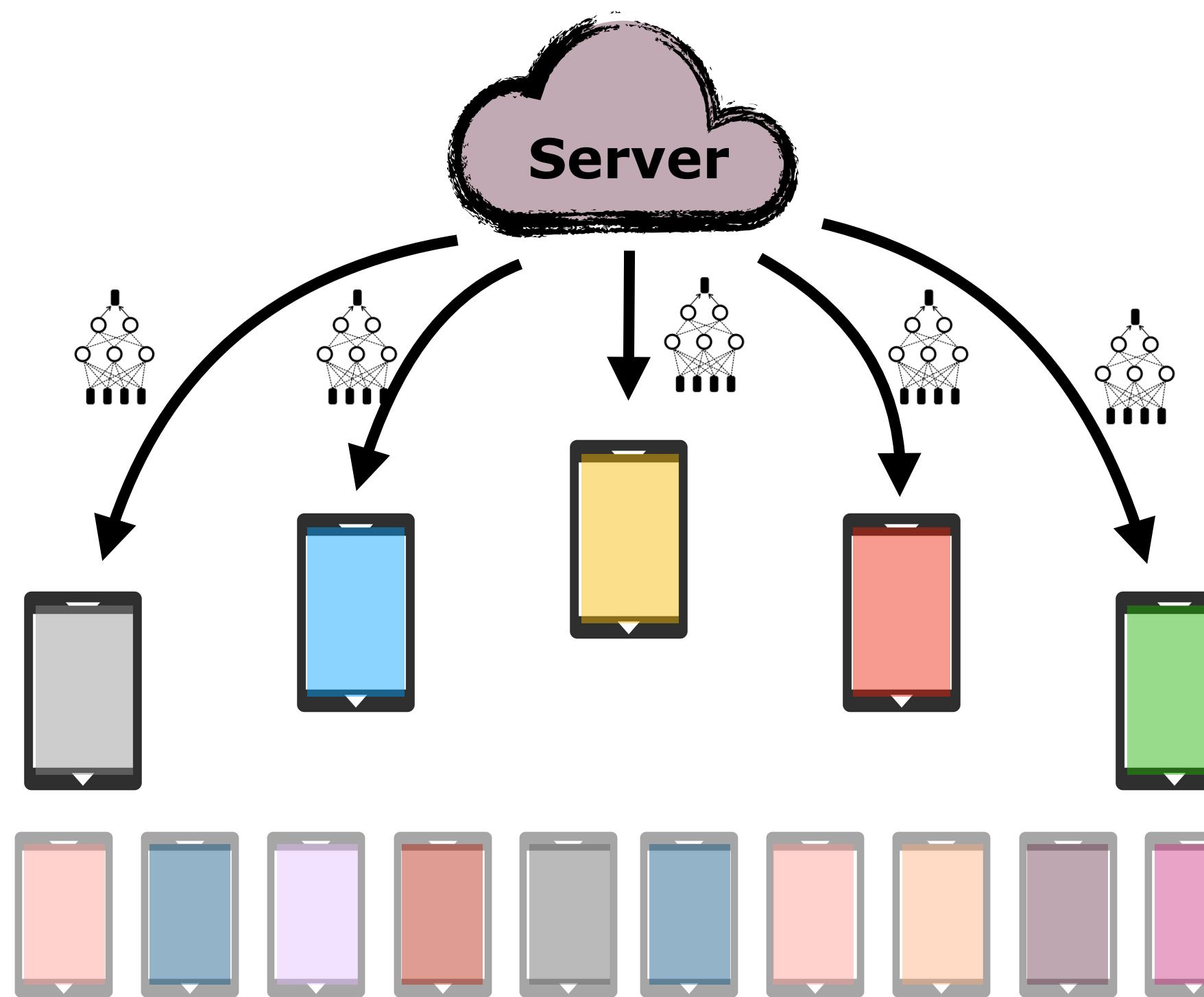
Usual Algorithm (FedAvg):

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

Our 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*



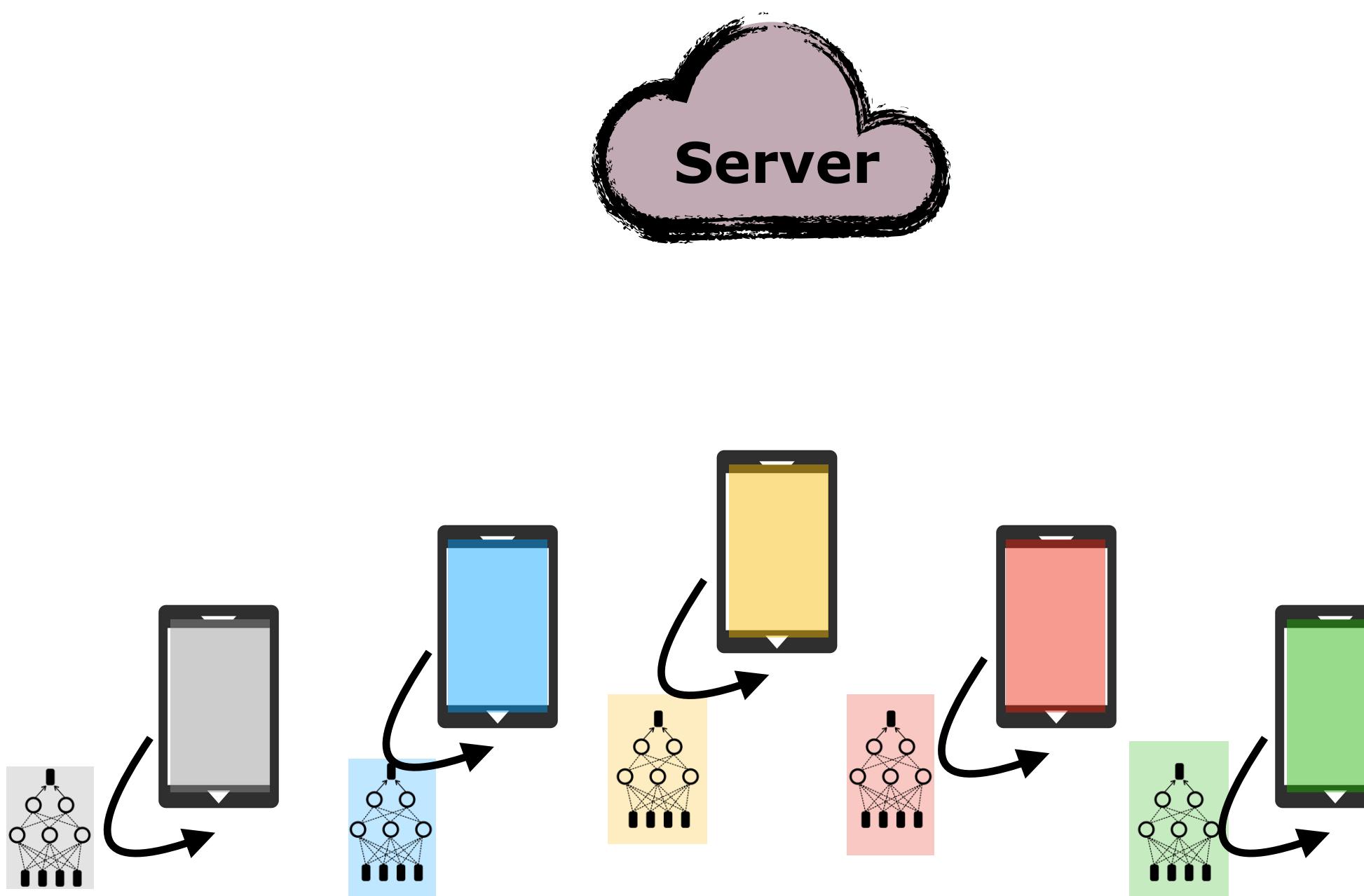
Usual Algorithm (FedAvg):

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

Our Algorithm:

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

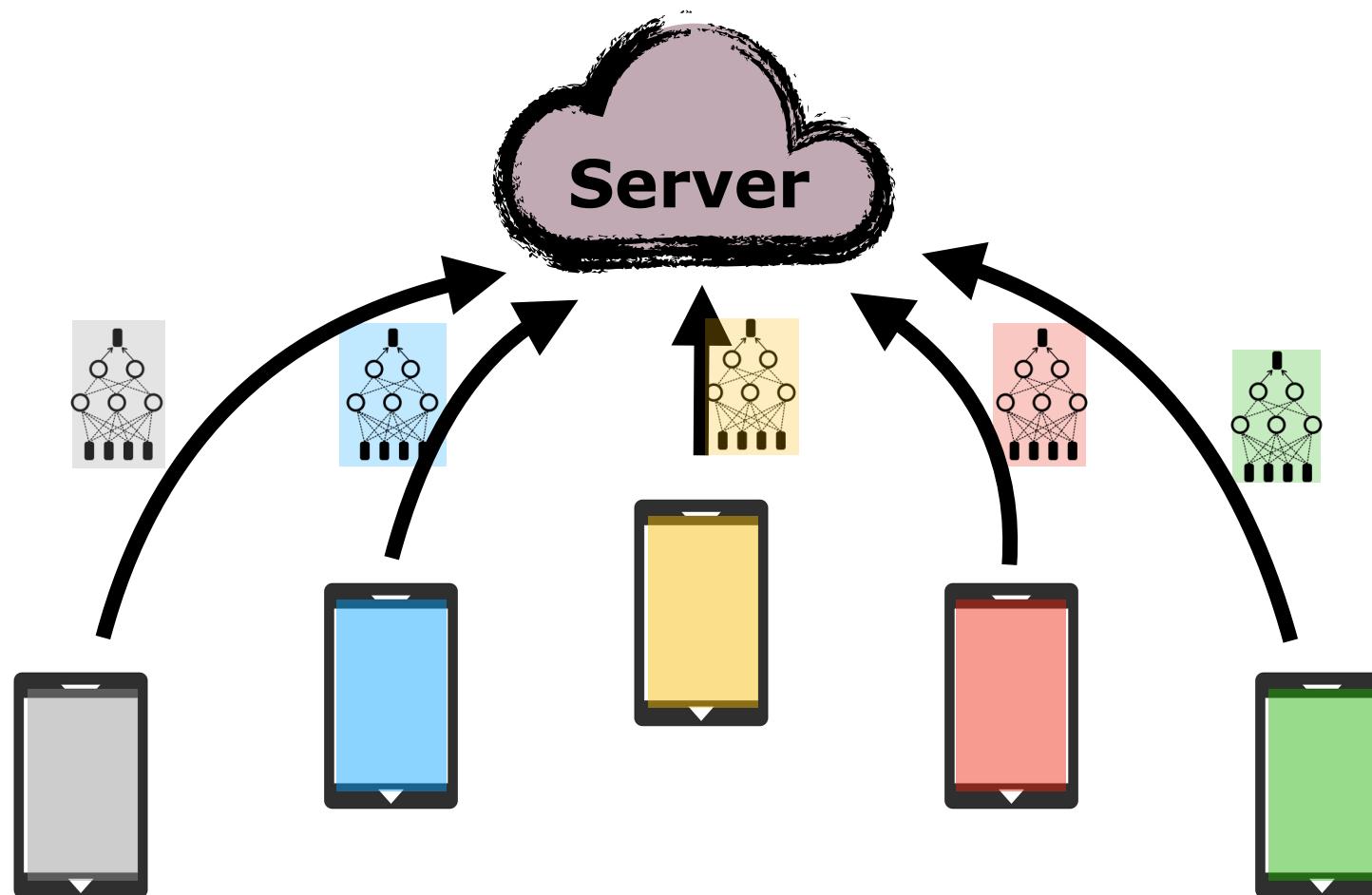
Step 2 of 3: Clients perform local gradient descent on their local data



Usual Algorithm (FedAvg):

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

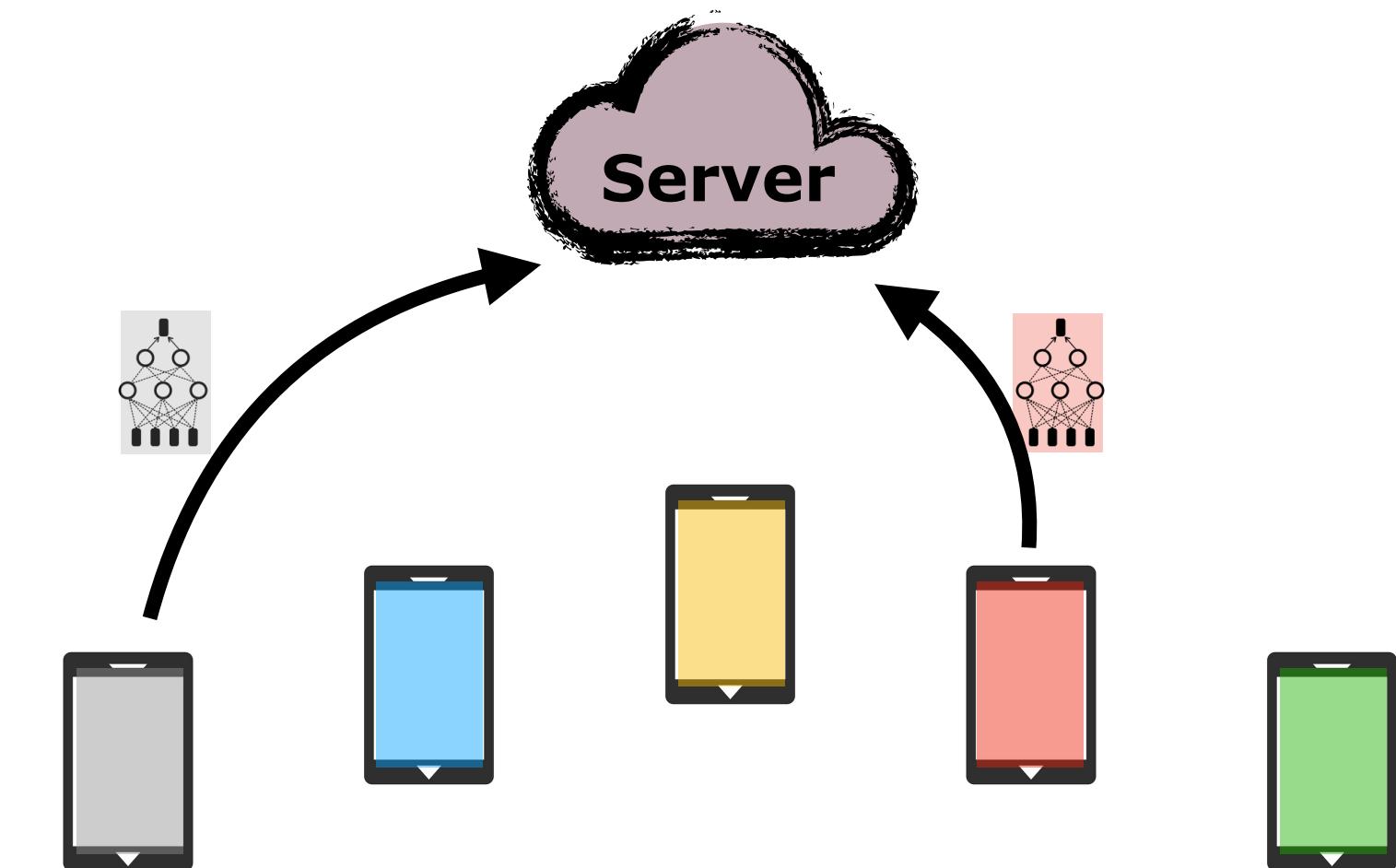
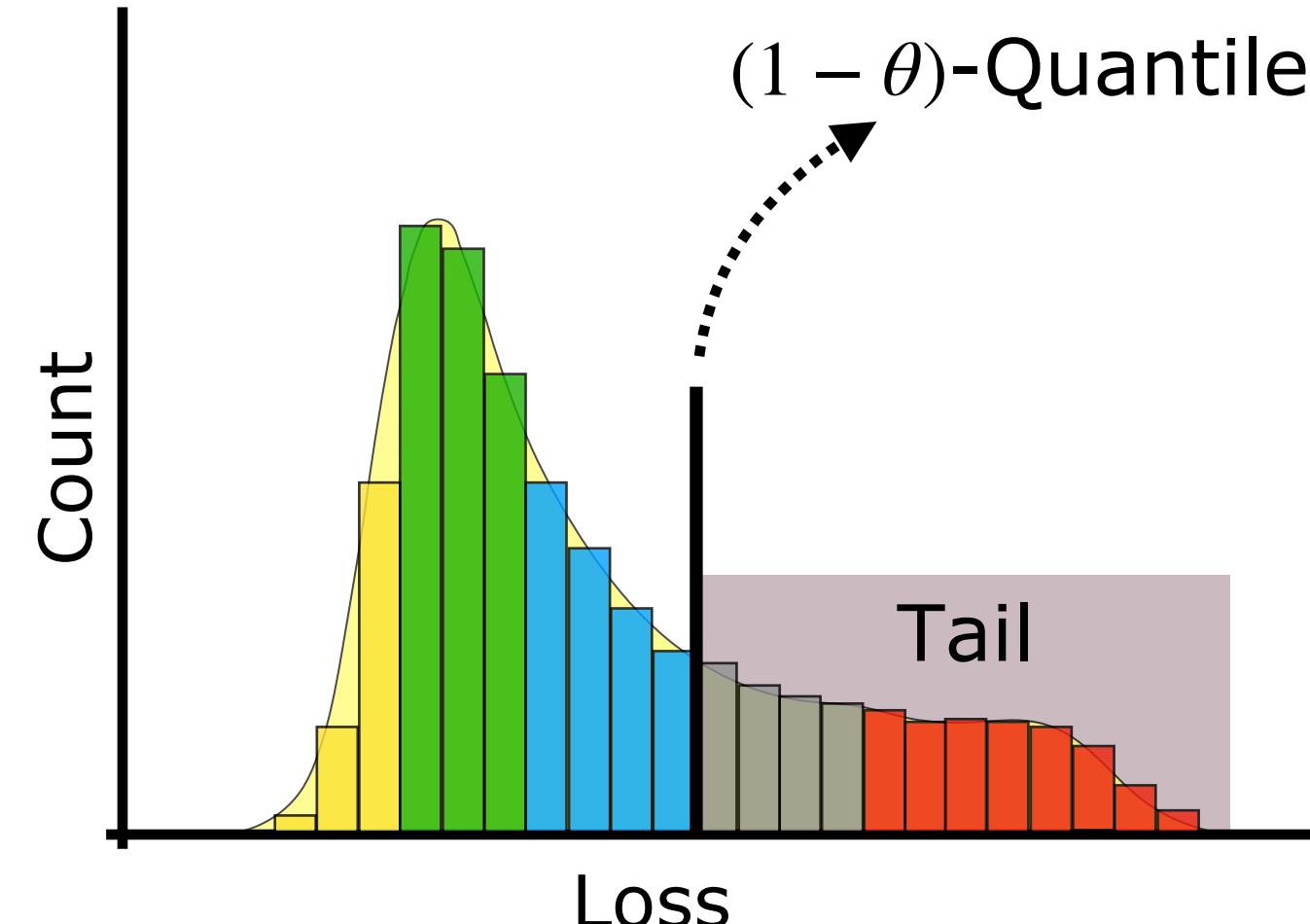
*Step 3 of 3: Aggregate updates contributed by **all clients***



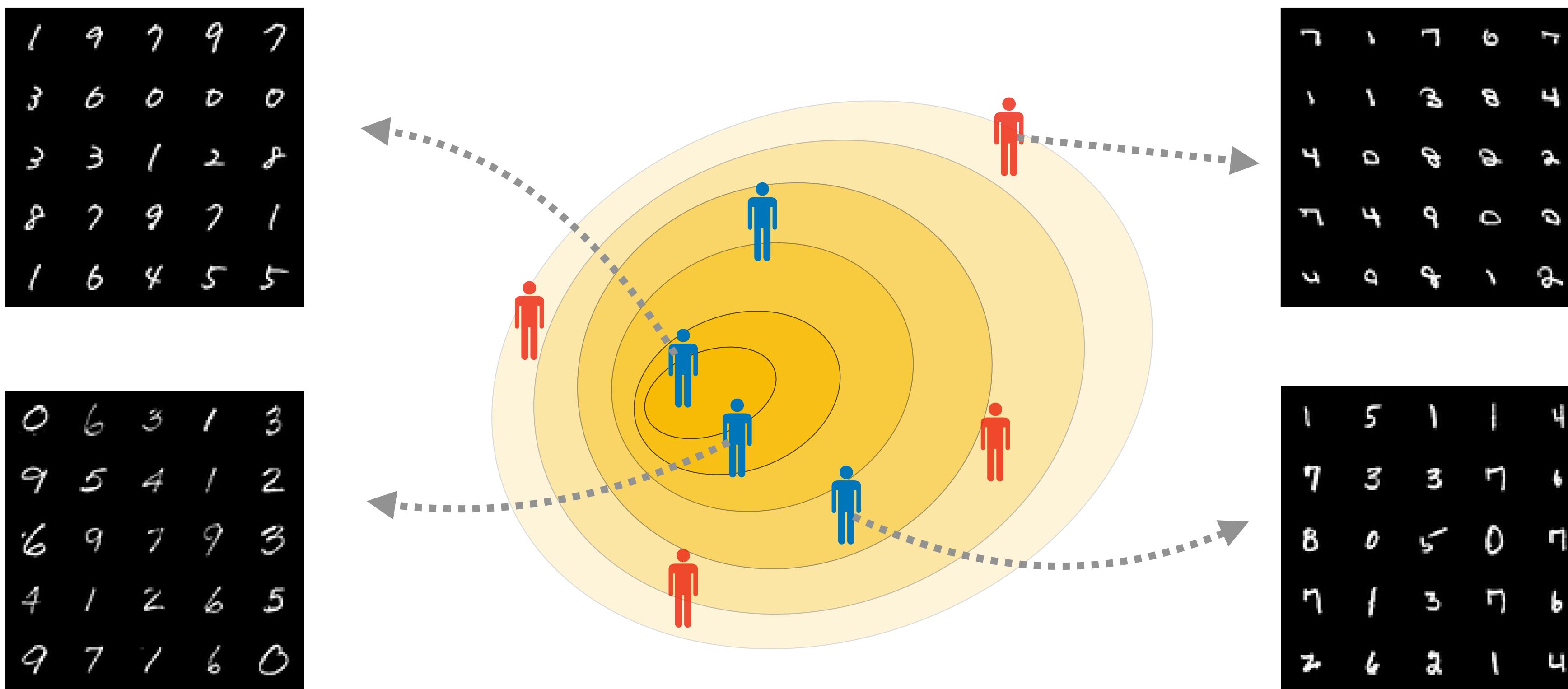
Our Algorithm:

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

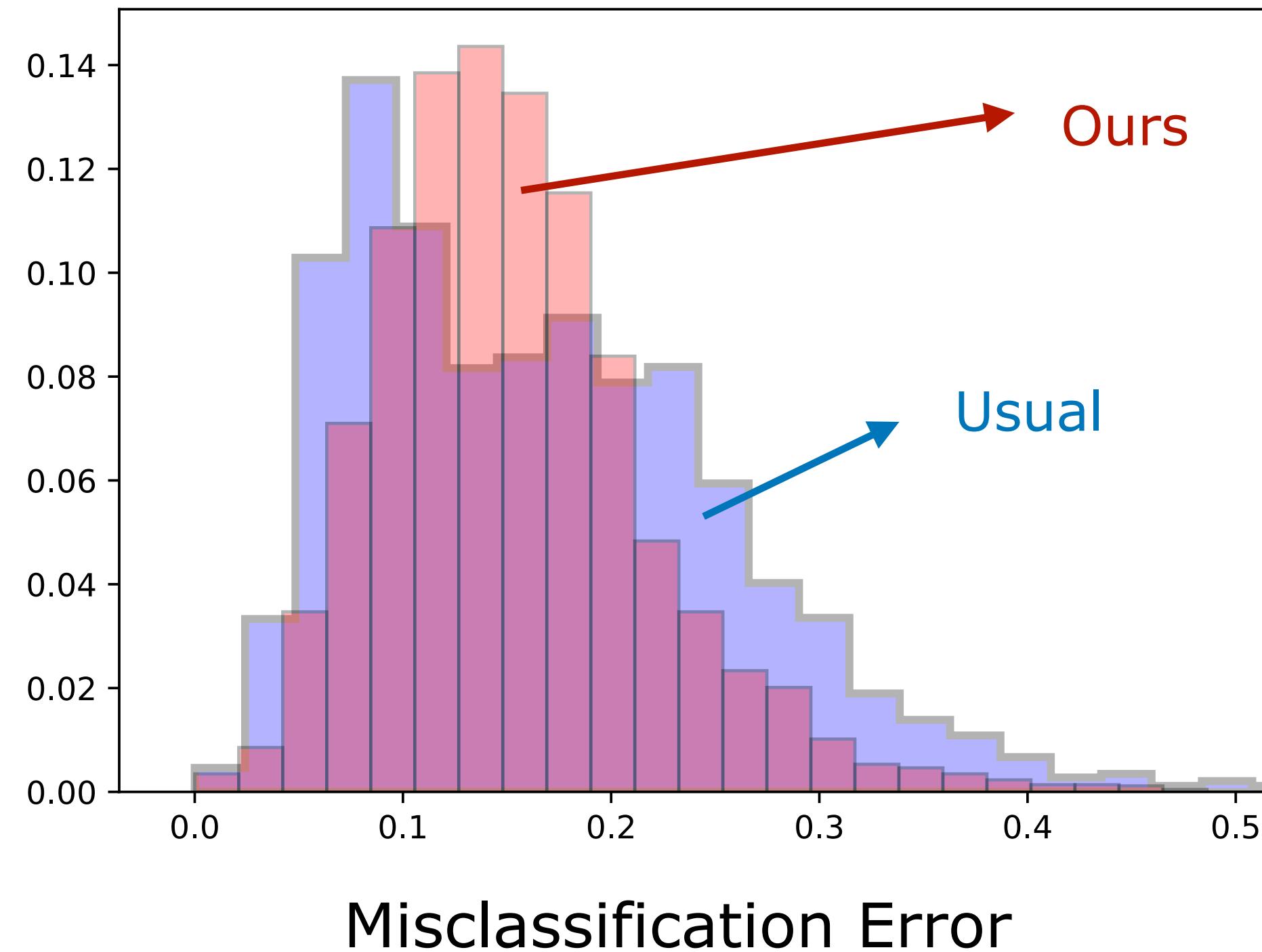
*Step 3 of 3: Aggregate updates contributed by **tail clients** only*



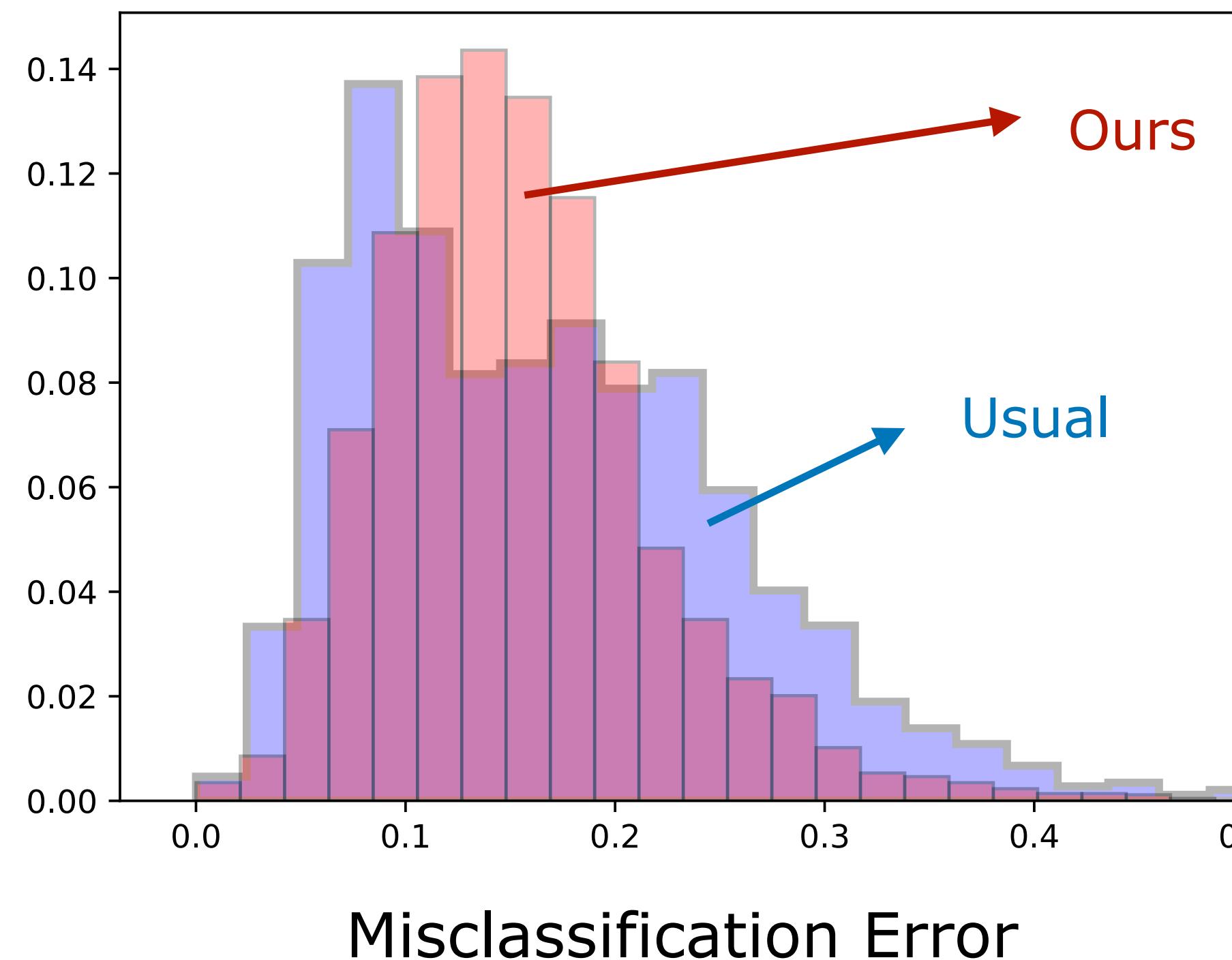
Experiments: EMNIST



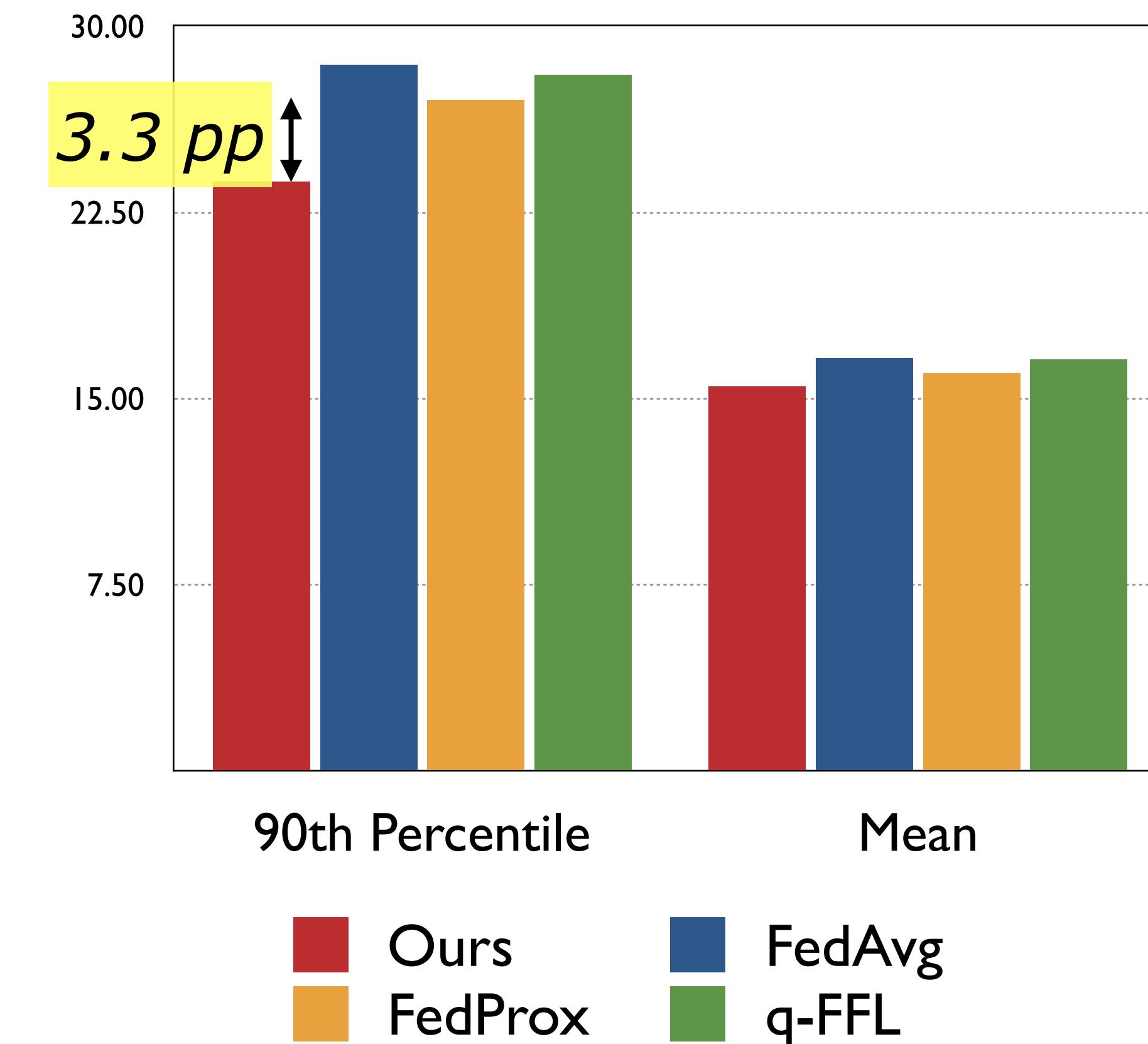
Histogram of per-client errors



Histogram of per-client errors



Misclassif. Error



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

Convergence analysis (non-convex)

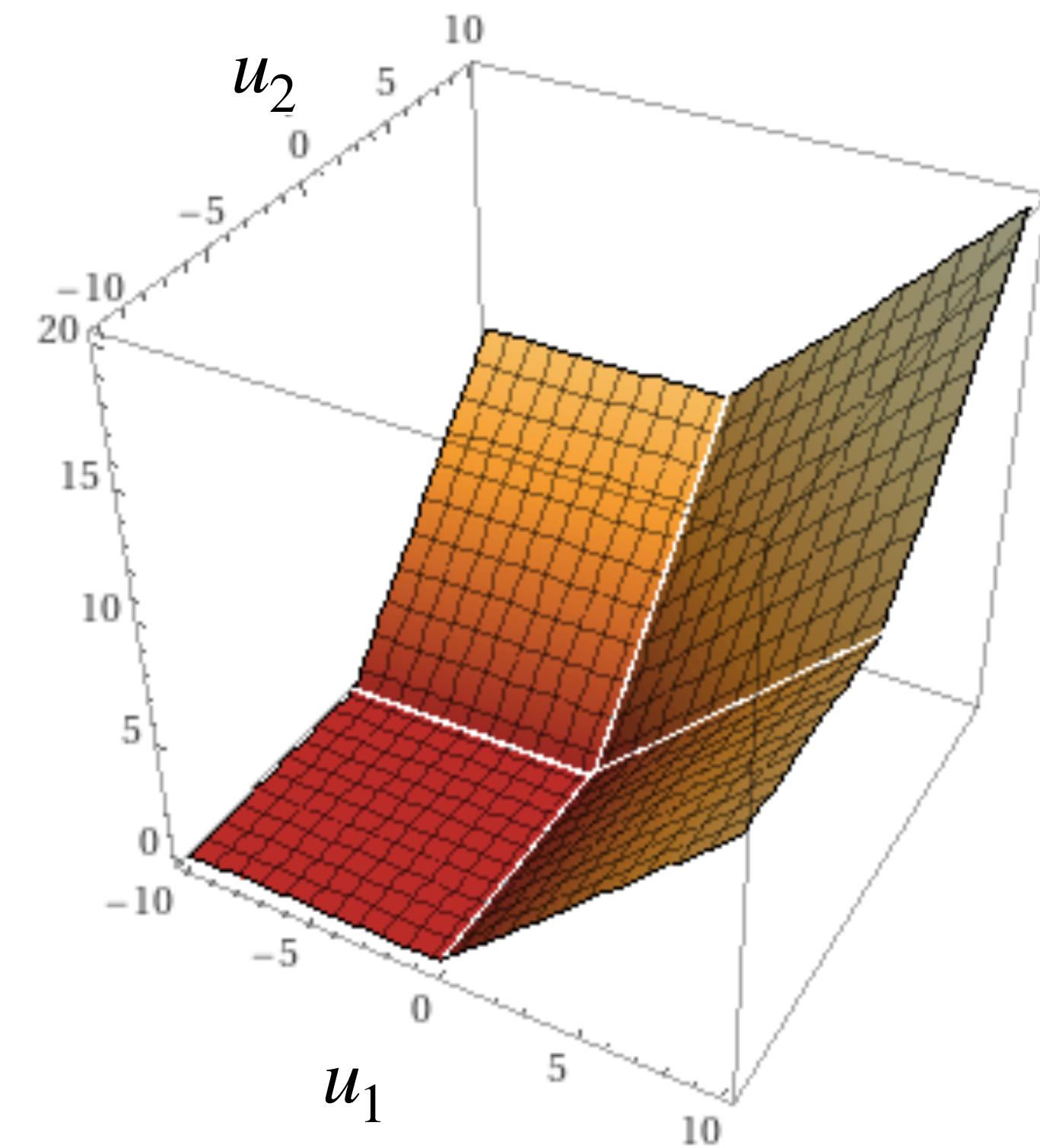


Faster optimization:
reduce communication



Challenge #1:

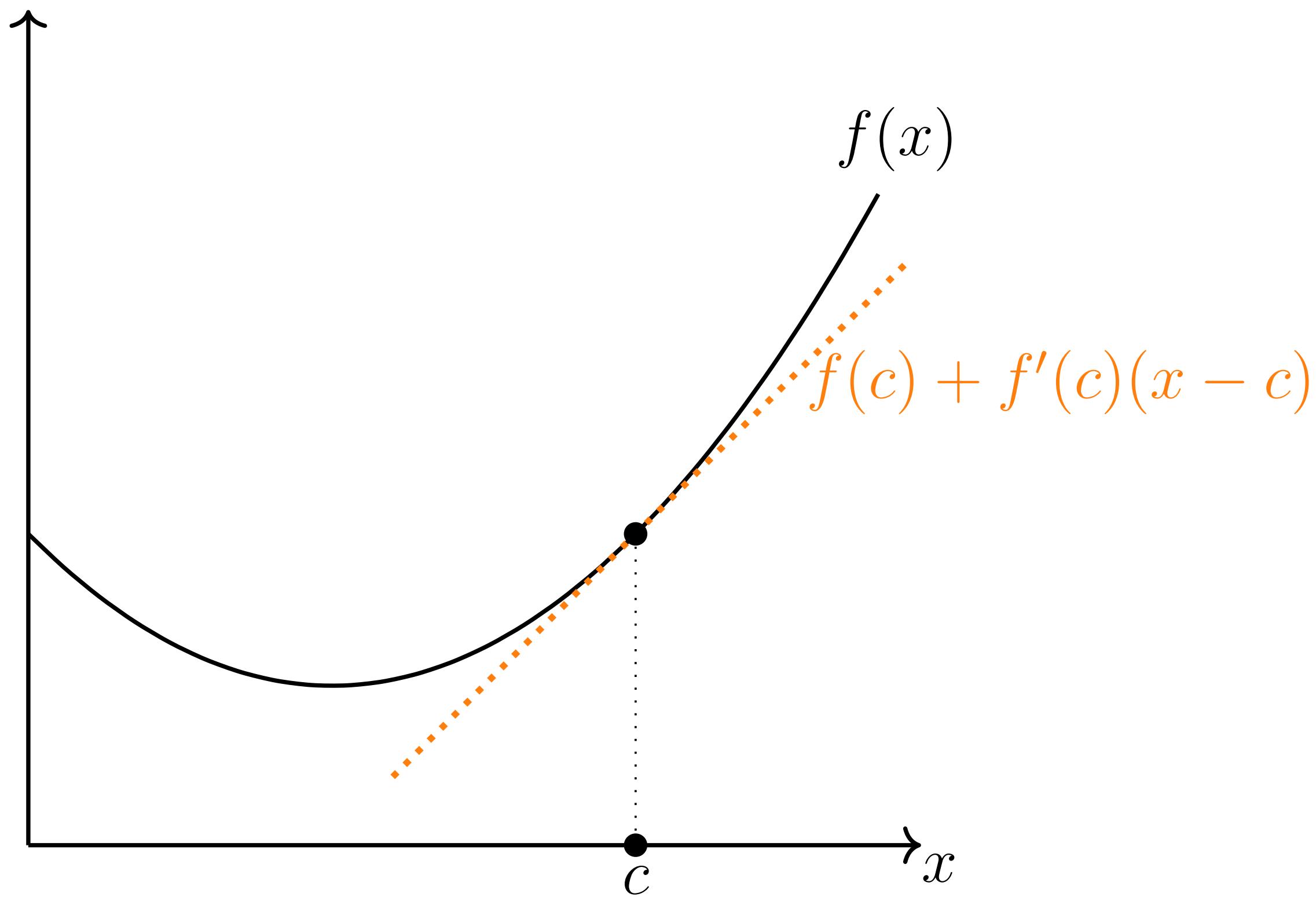
The superquantile is non-smooth



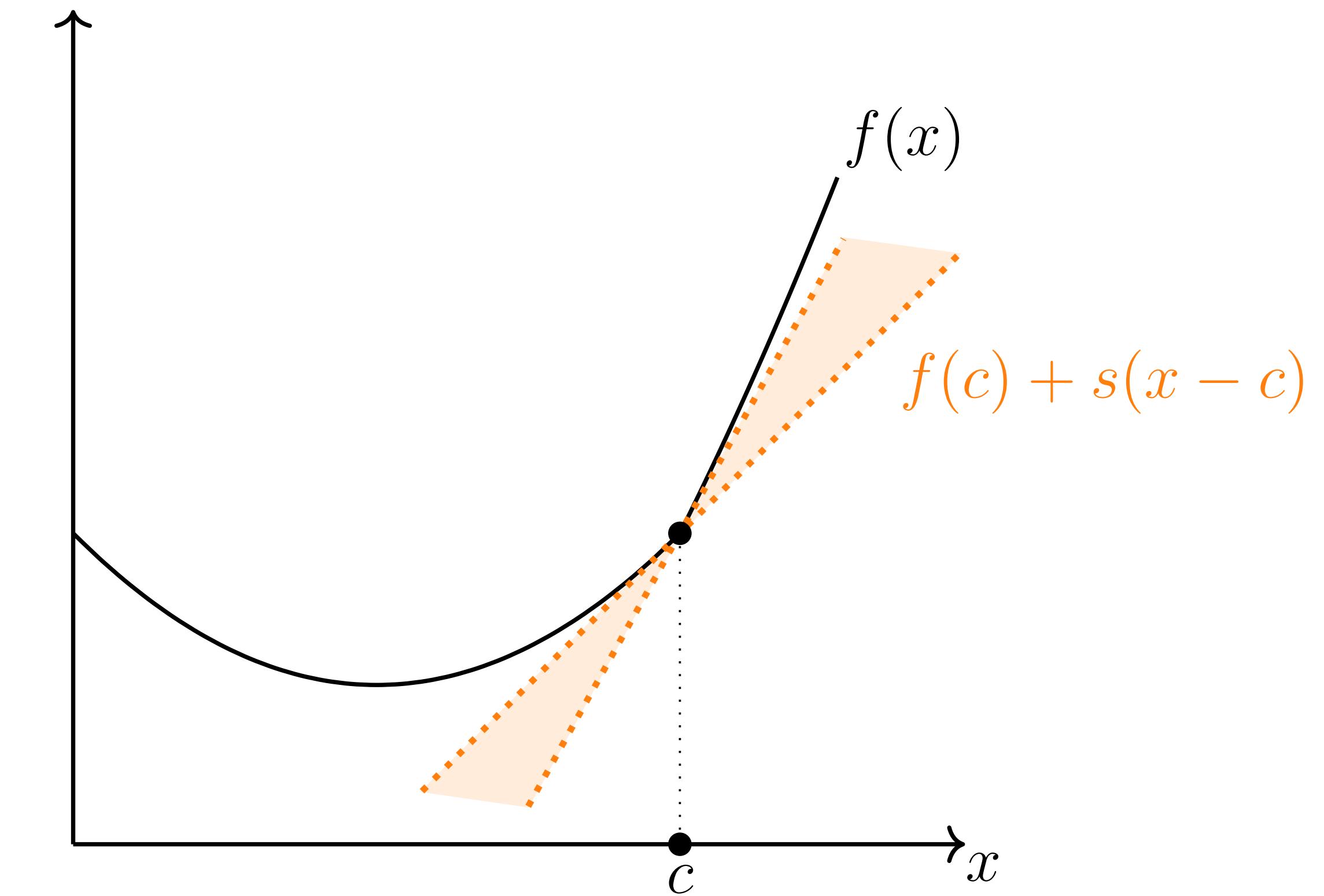
plot of $h(u_1, u_2) = \mathbb{S}_{1/2}(u_1, u_2, 0, 0)$

Subgradient illustration

Smooth



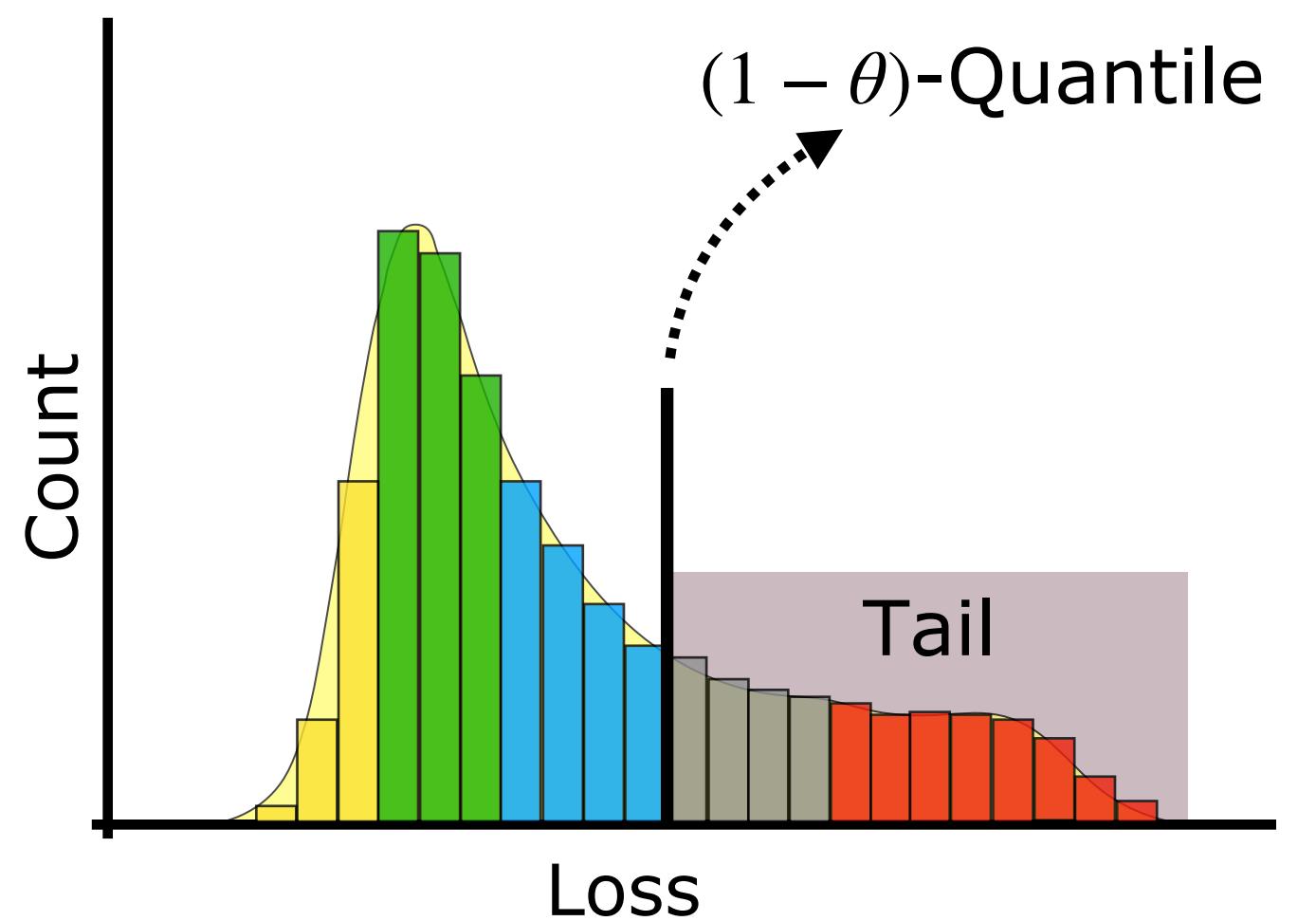
Non-smooth



Nonsmooth: The subgradient 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



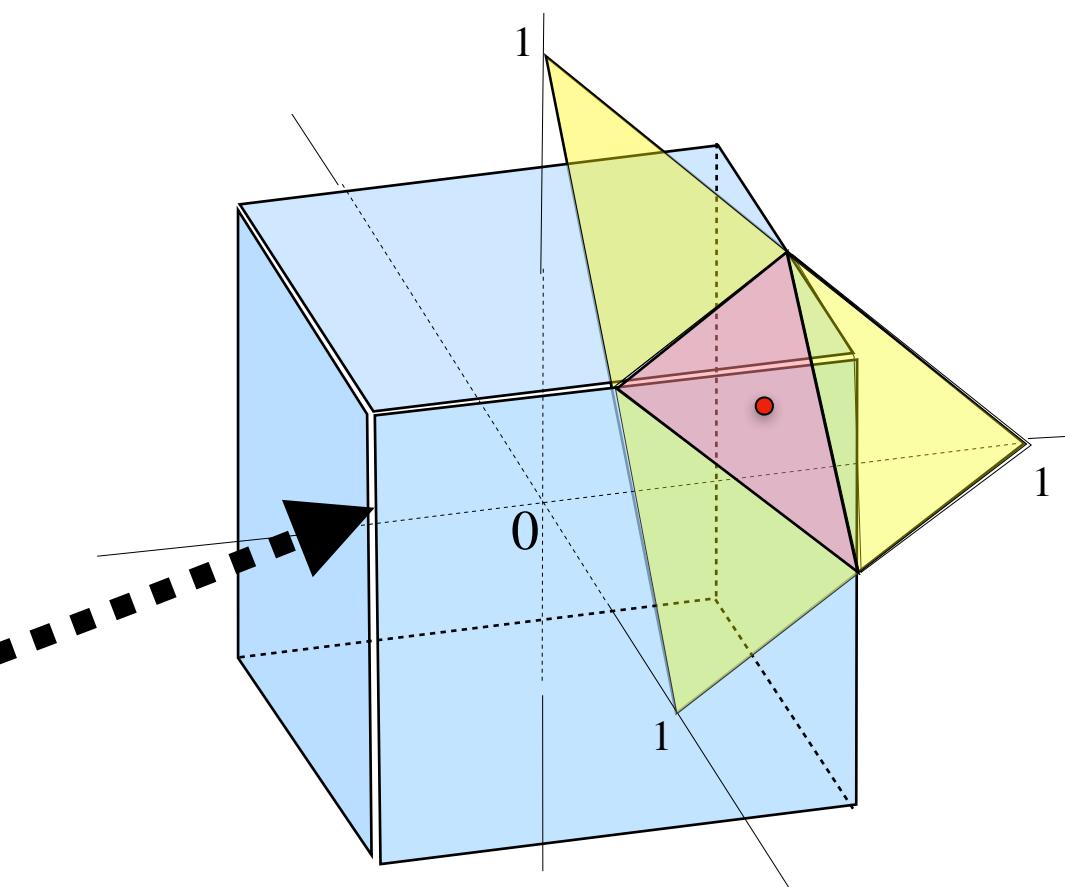
Nonsmooth: The subgradient has a tractable form

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

assuming θ_n
is an integer

Proof Chain rule \implies subgradient 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

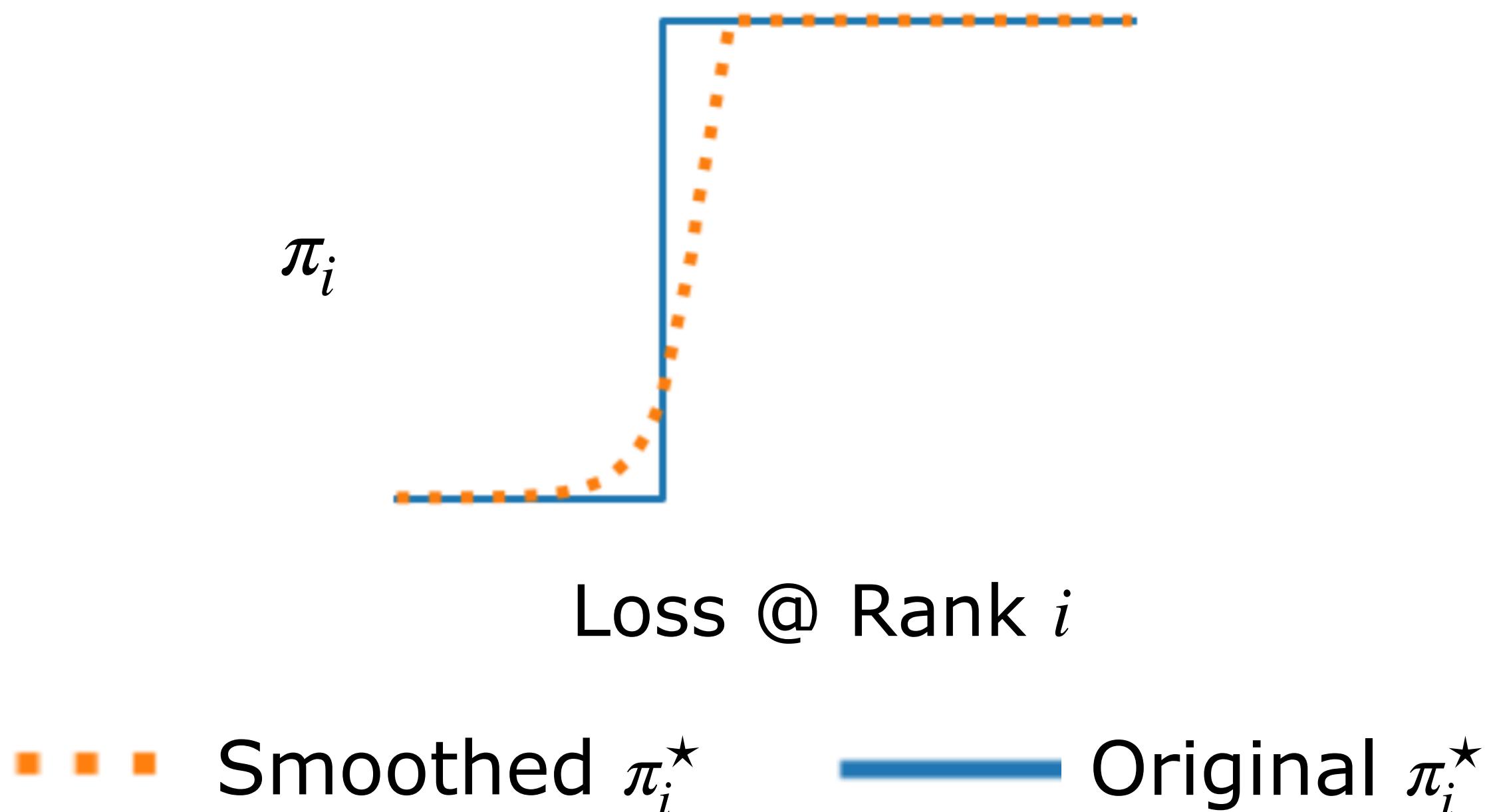
[Dantzig, ORIJ (1957)]

Nonsmooth: The subgradient 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



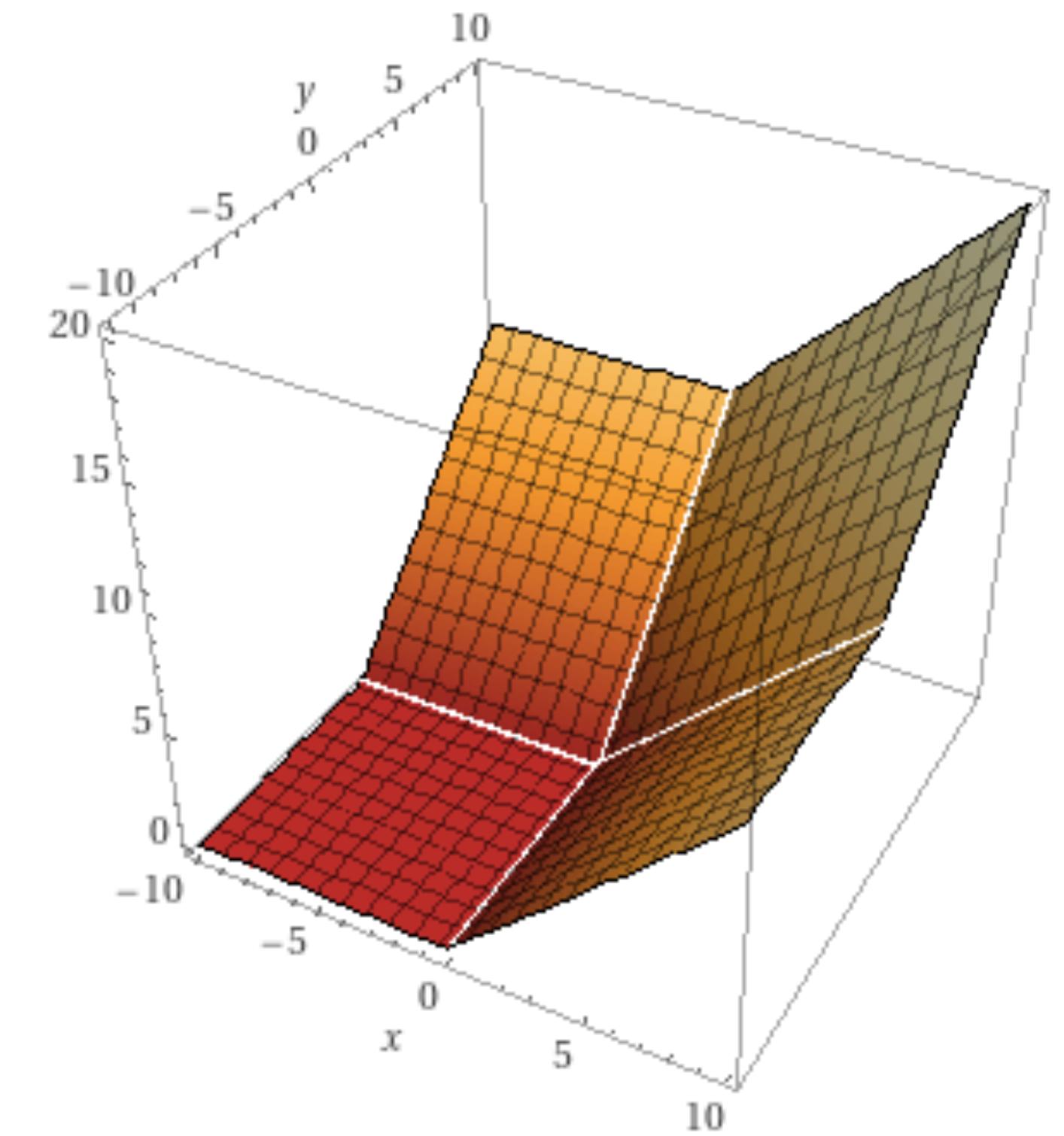
[Nesterov. (Math. Prog. 2005),
Beck & Teboulle. (SIAM J. Optim. 2012),
P., Roulet, Kakade, Harchaoui. (NeurIPS 2018),
Laguel, P., Malick, Harchaoui. (SVVA 2021)]

Challenge #2

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

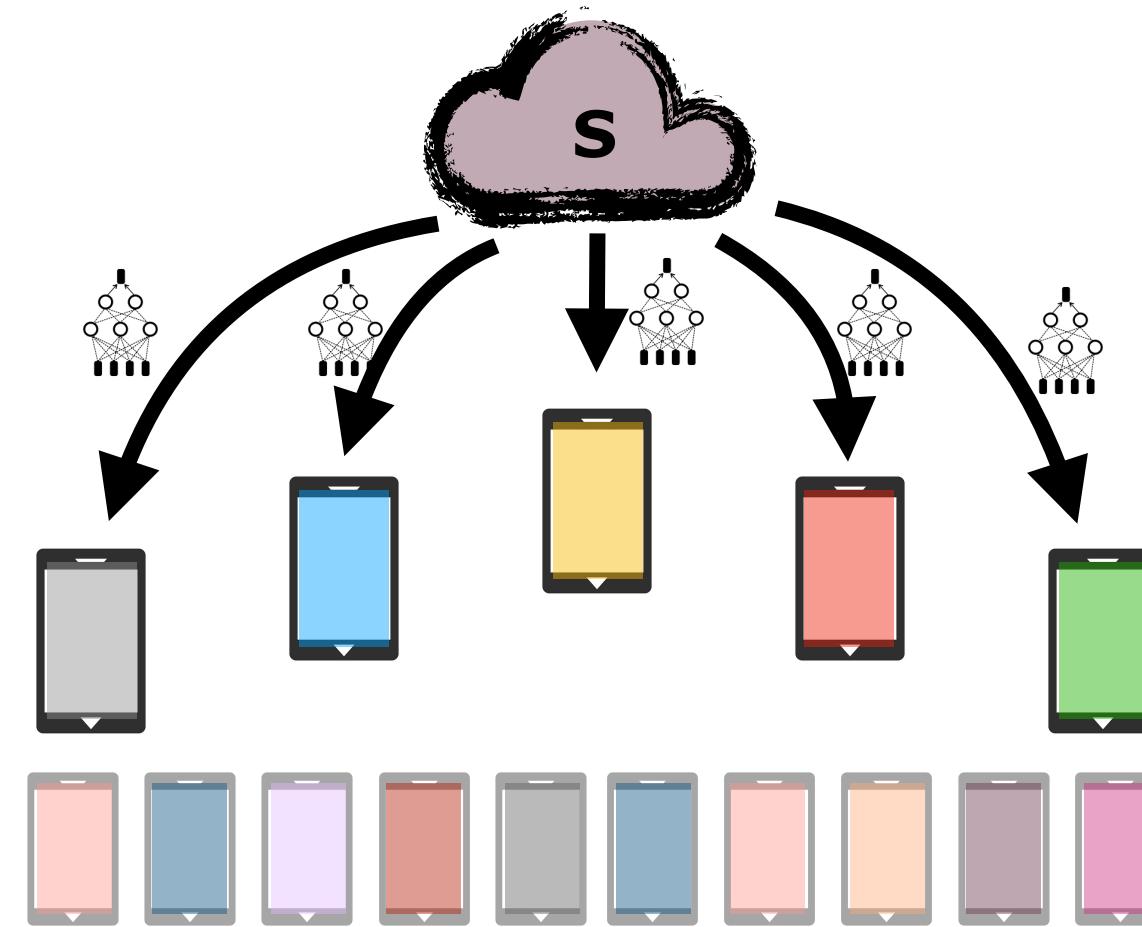
For i.i.d. copies Z_1, \dots, Z_m of Z , we have

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



Nonlinear: We minimize a close surrogate

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



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

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

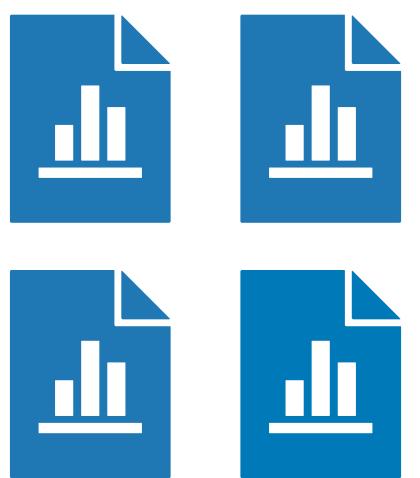
Privacy analysis



Privacy of user data



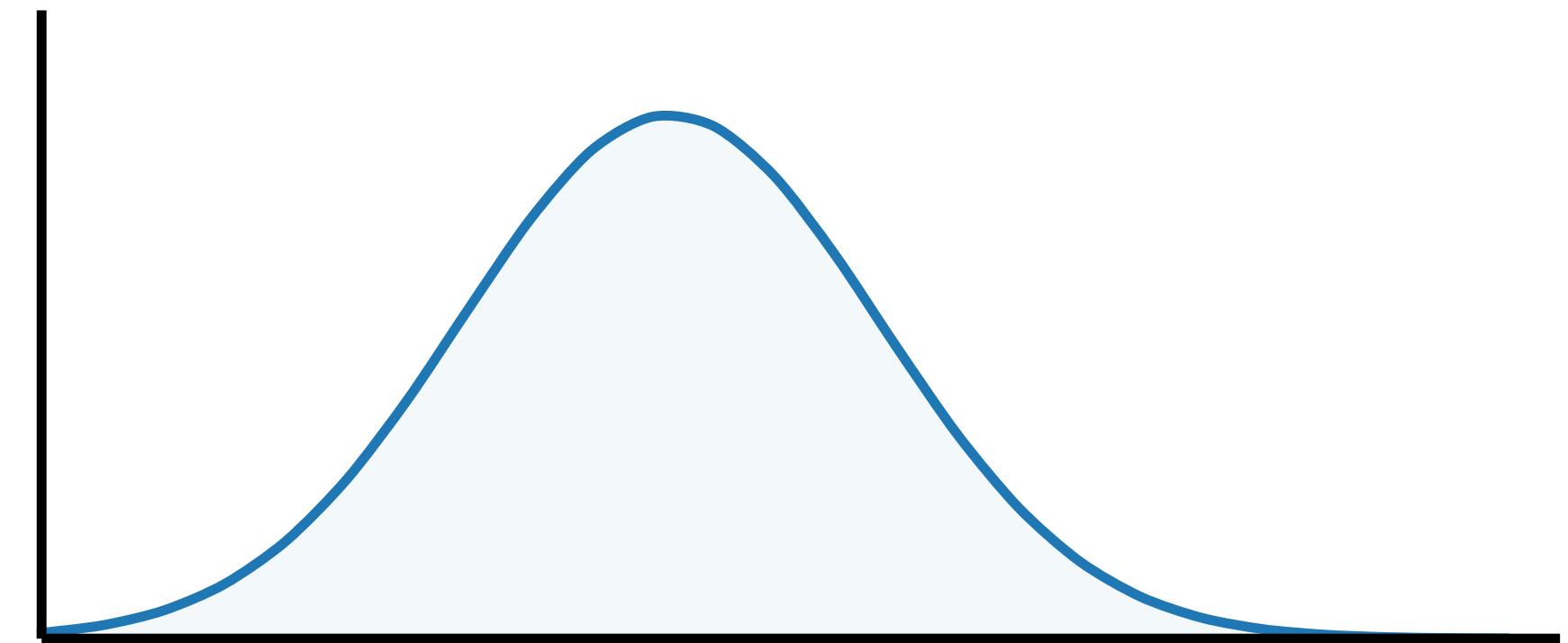
Dataset



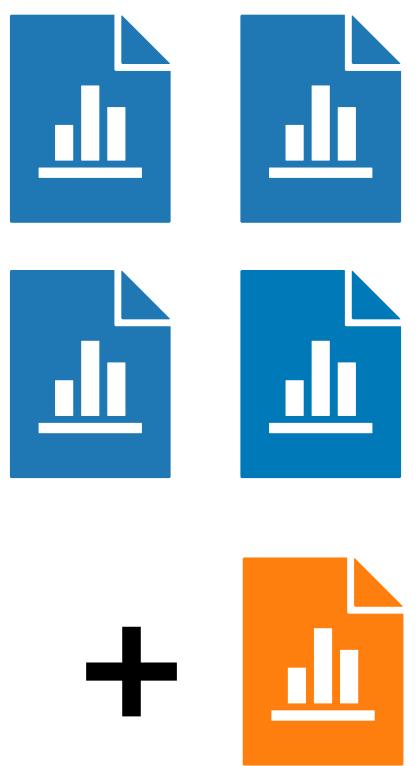
Randomized
Algorithm



Output Distribution
(e.g. over models)



Dataset



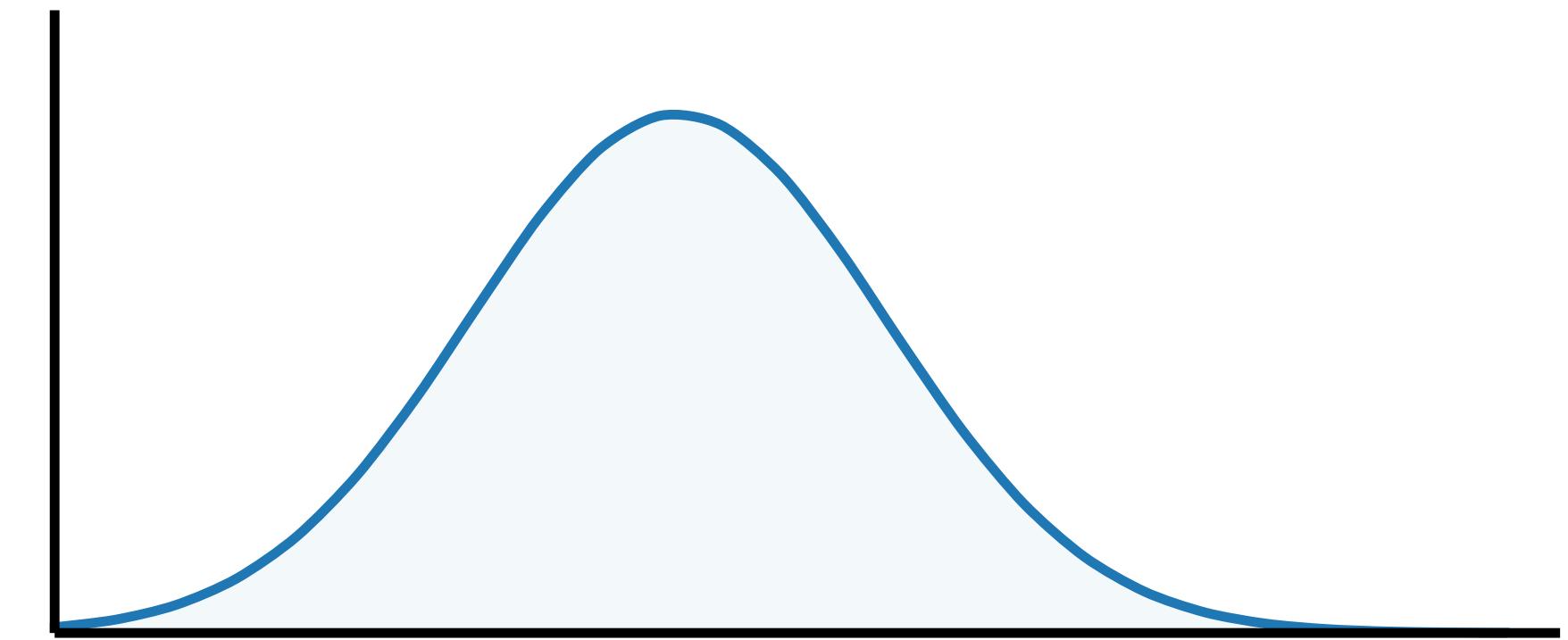
+



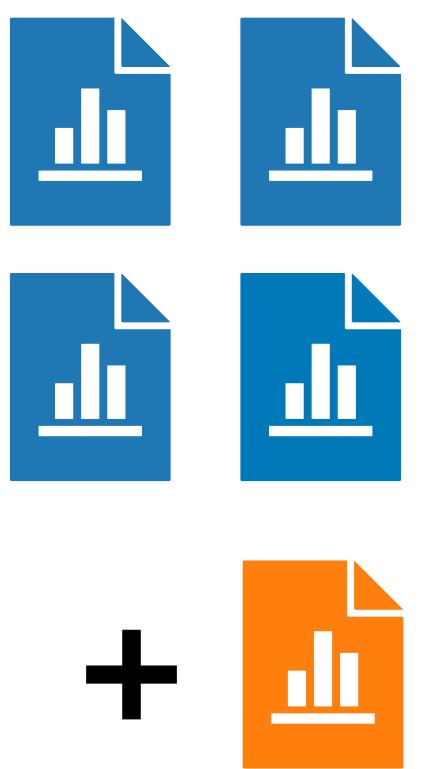
Randomized
Algorithm



Output Distribution
(e.g. over models)



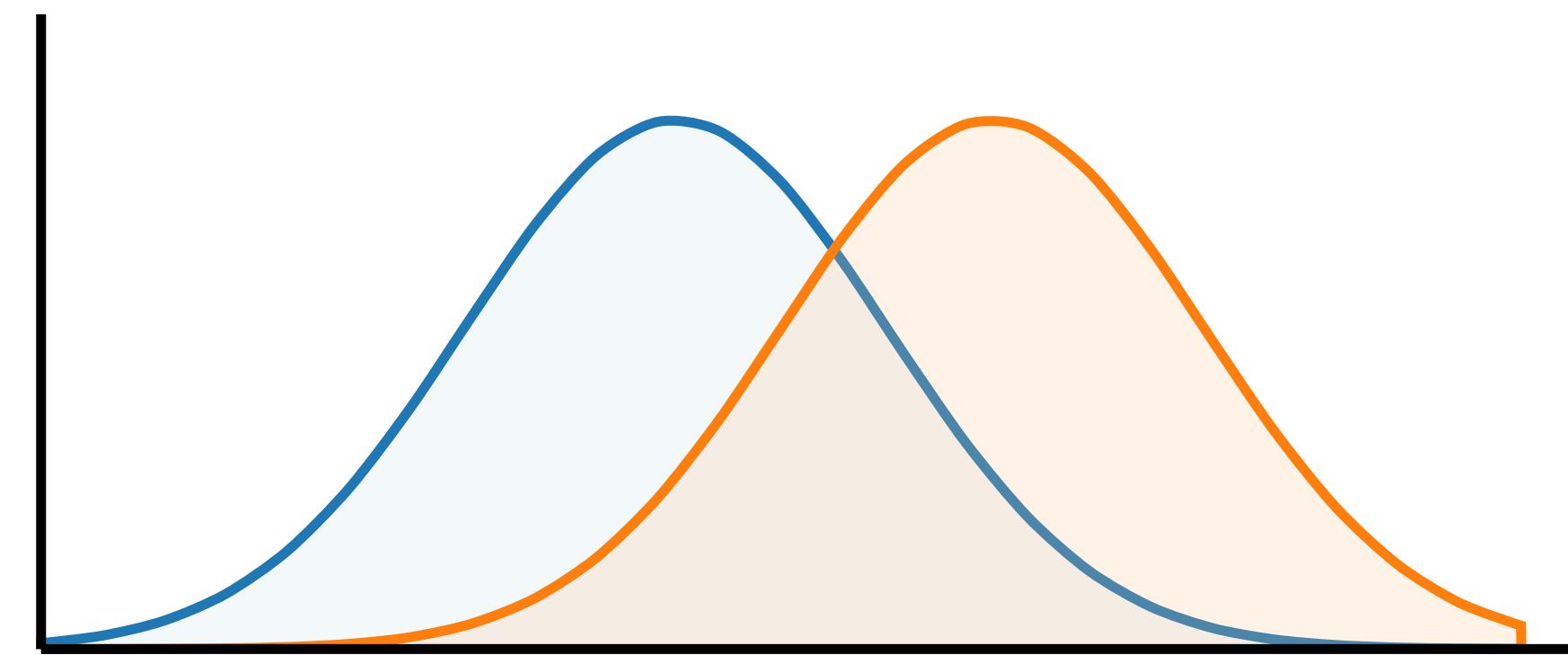
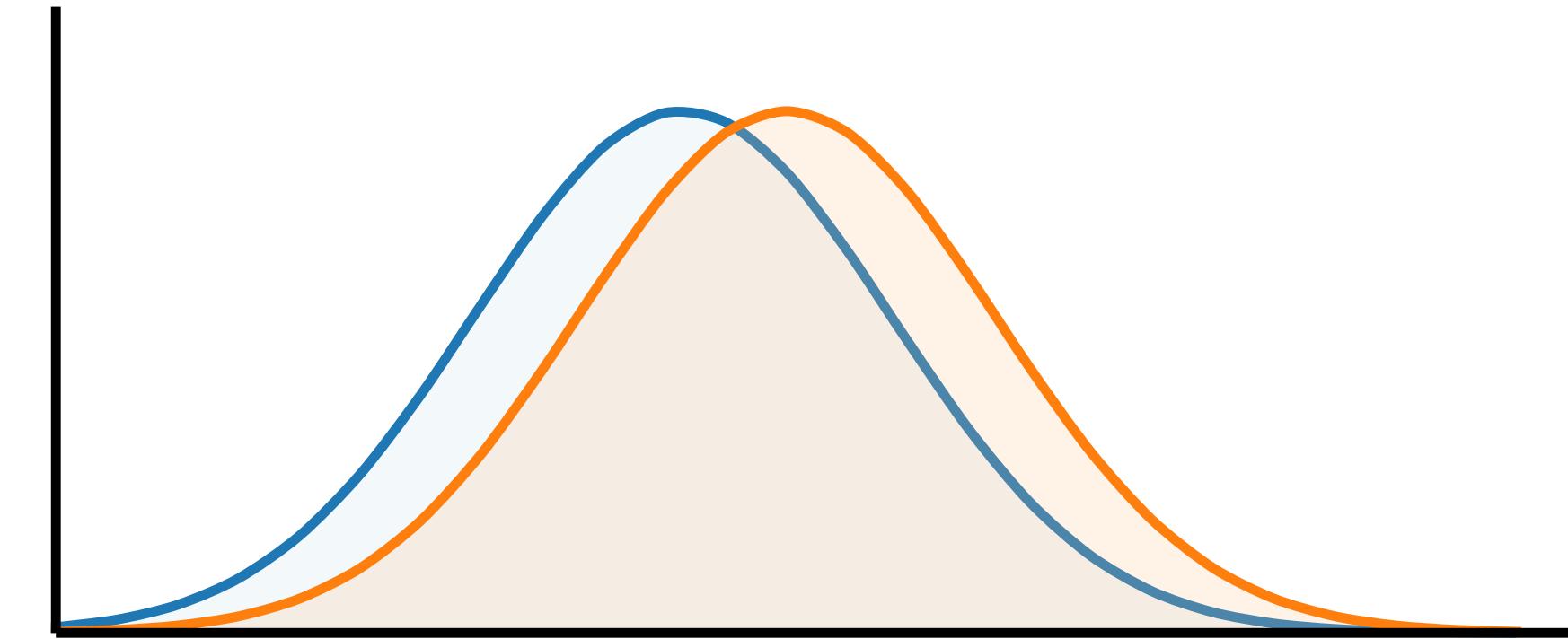
Dataset



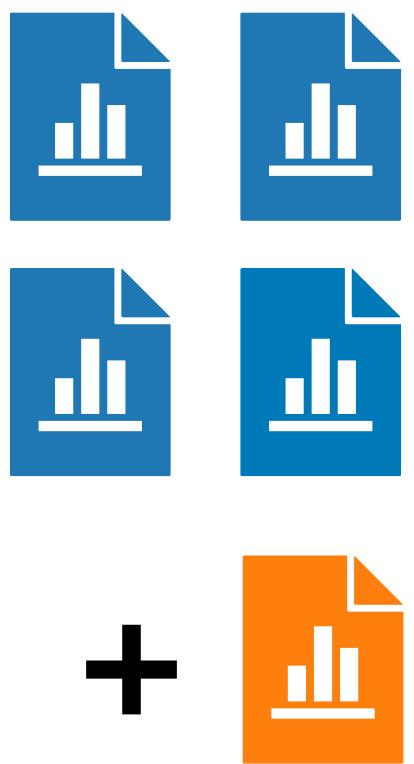
Randomized
Algorithm



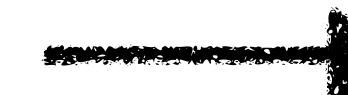
Output Distribution
(e.g. over models)



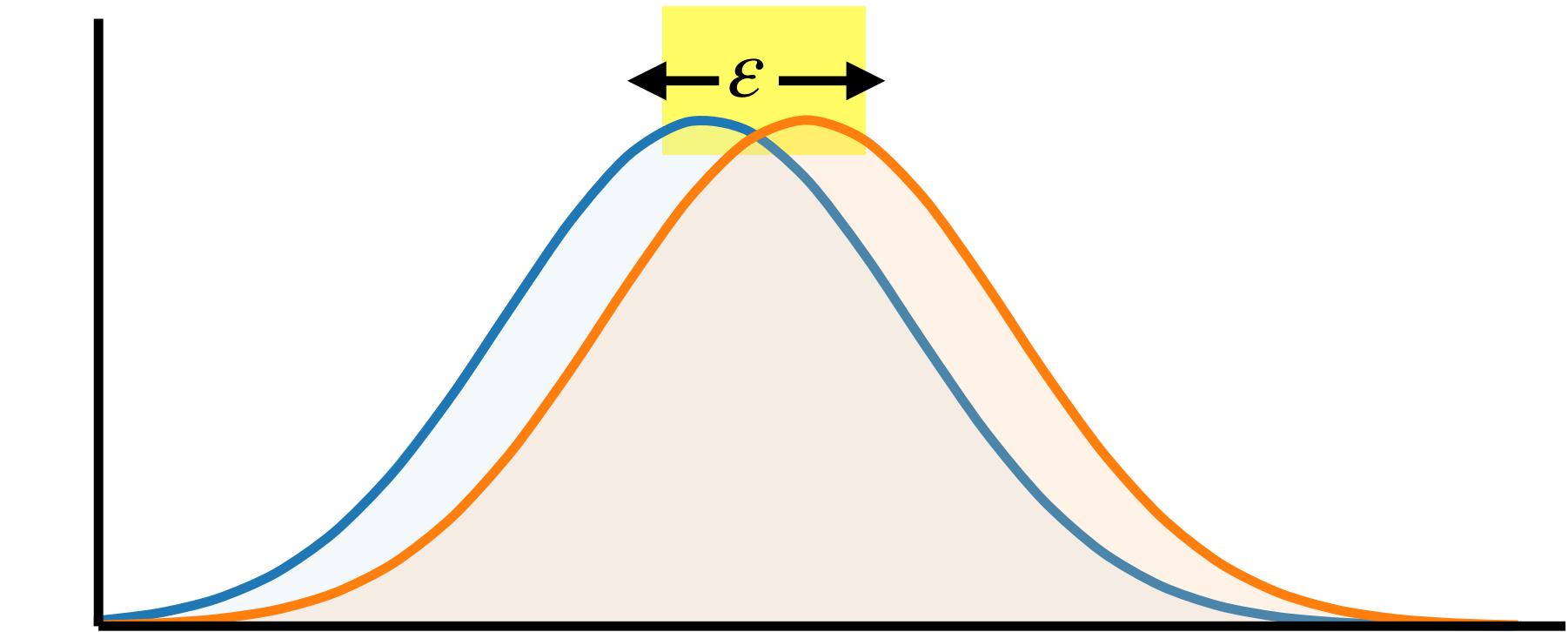
Dataset



Randomized
Algorithm



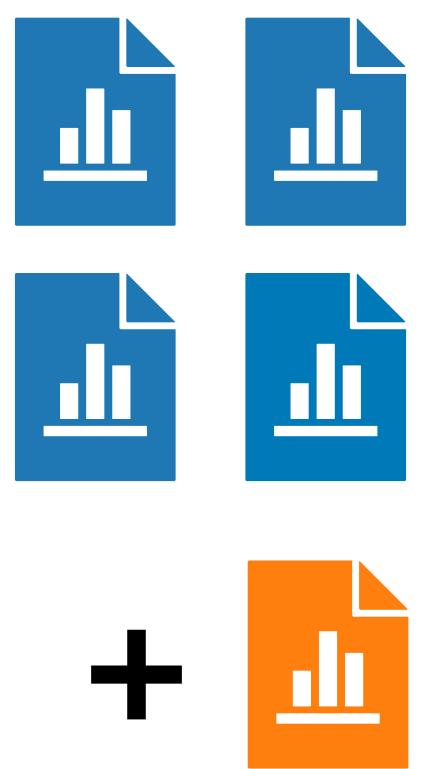
Output Distribution
(e.g. over models)



A randomized algorithm is **ϵ -differentially private** if the addition of **one user's data** does not alter its output distribution by more than ϵ

[Dwork, McSherry, Nissim, Smith (2006)]

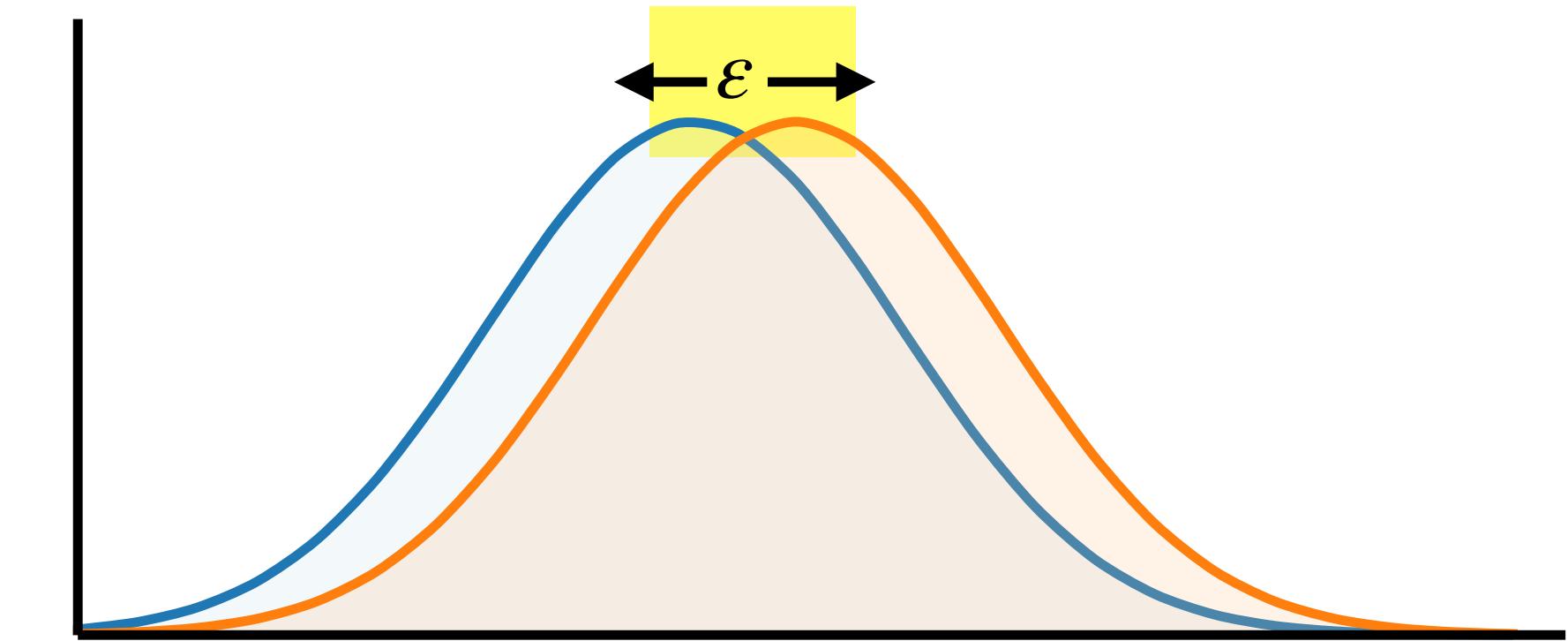
Dataset



+

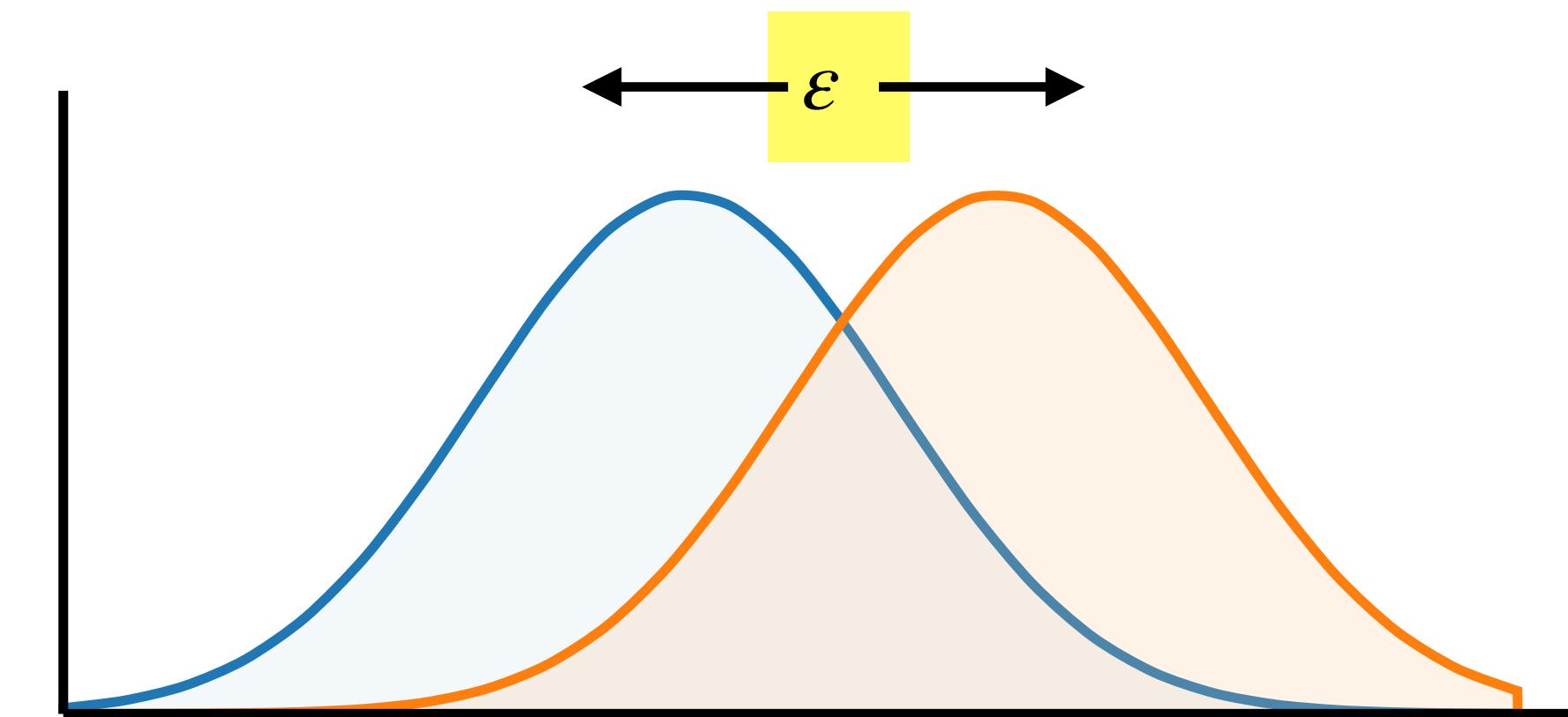


Output Distribution
(e.g. over models)



ε -differential privacy

Large $\varepsilon \implies$ more privacy leakage

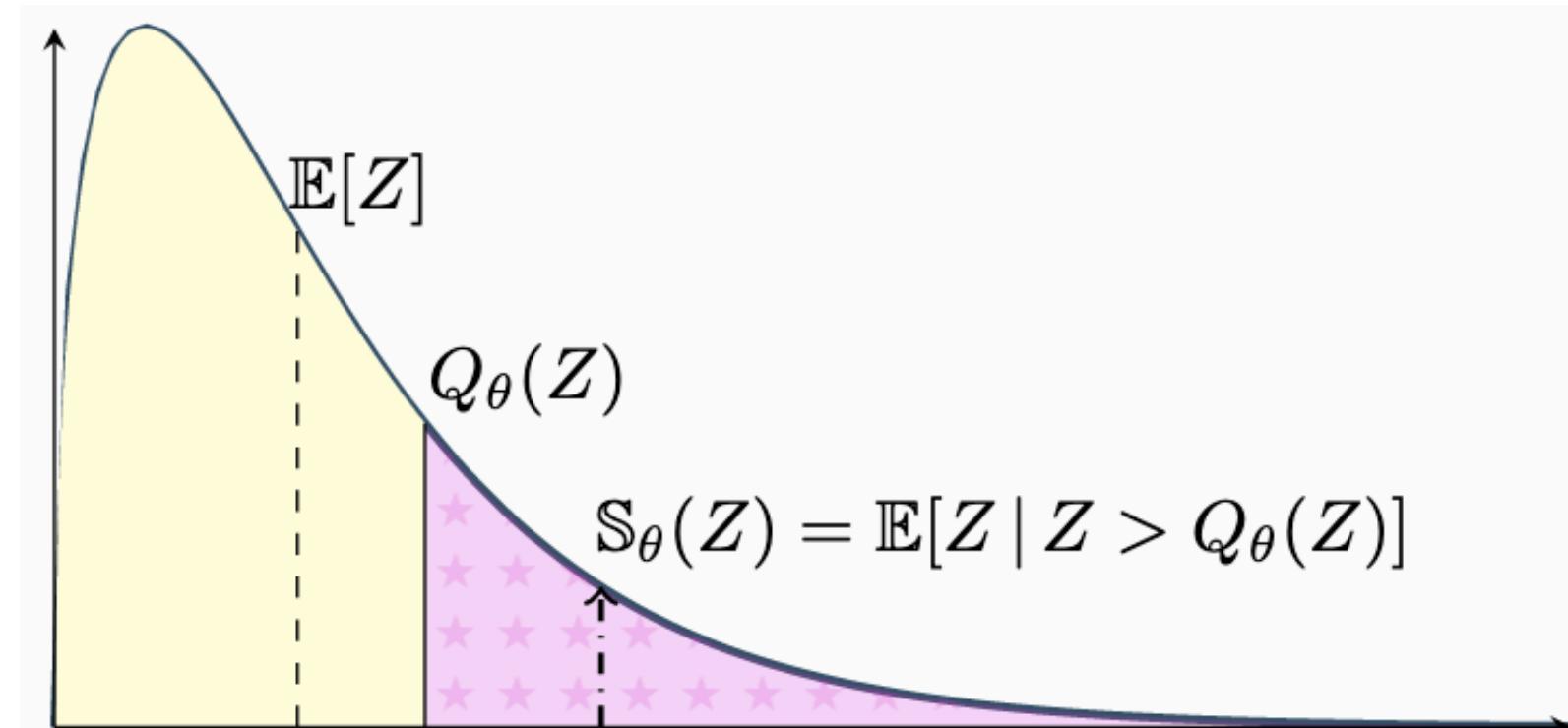


Privacy goal

Extend our algorithm to make it **differentially private**

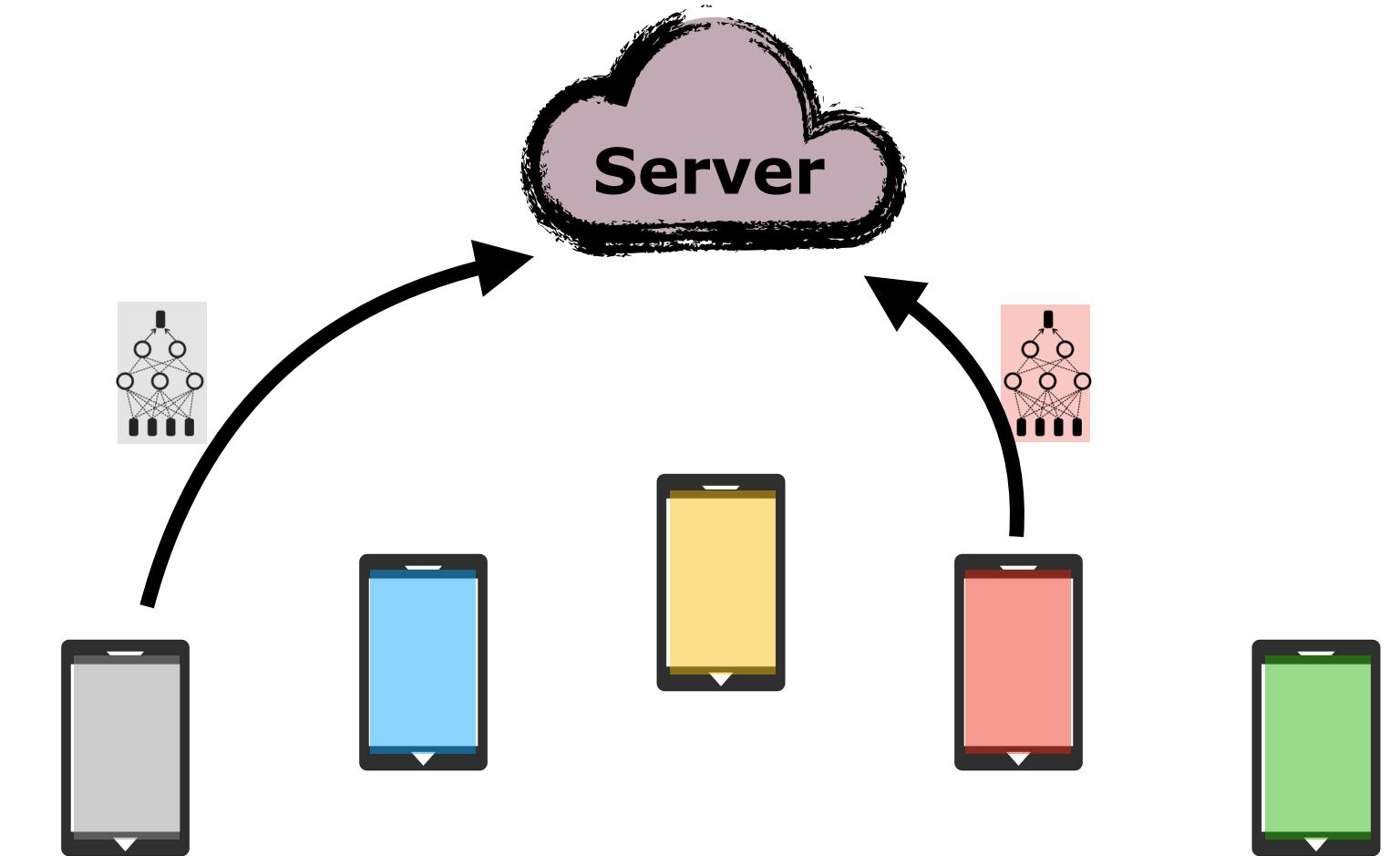
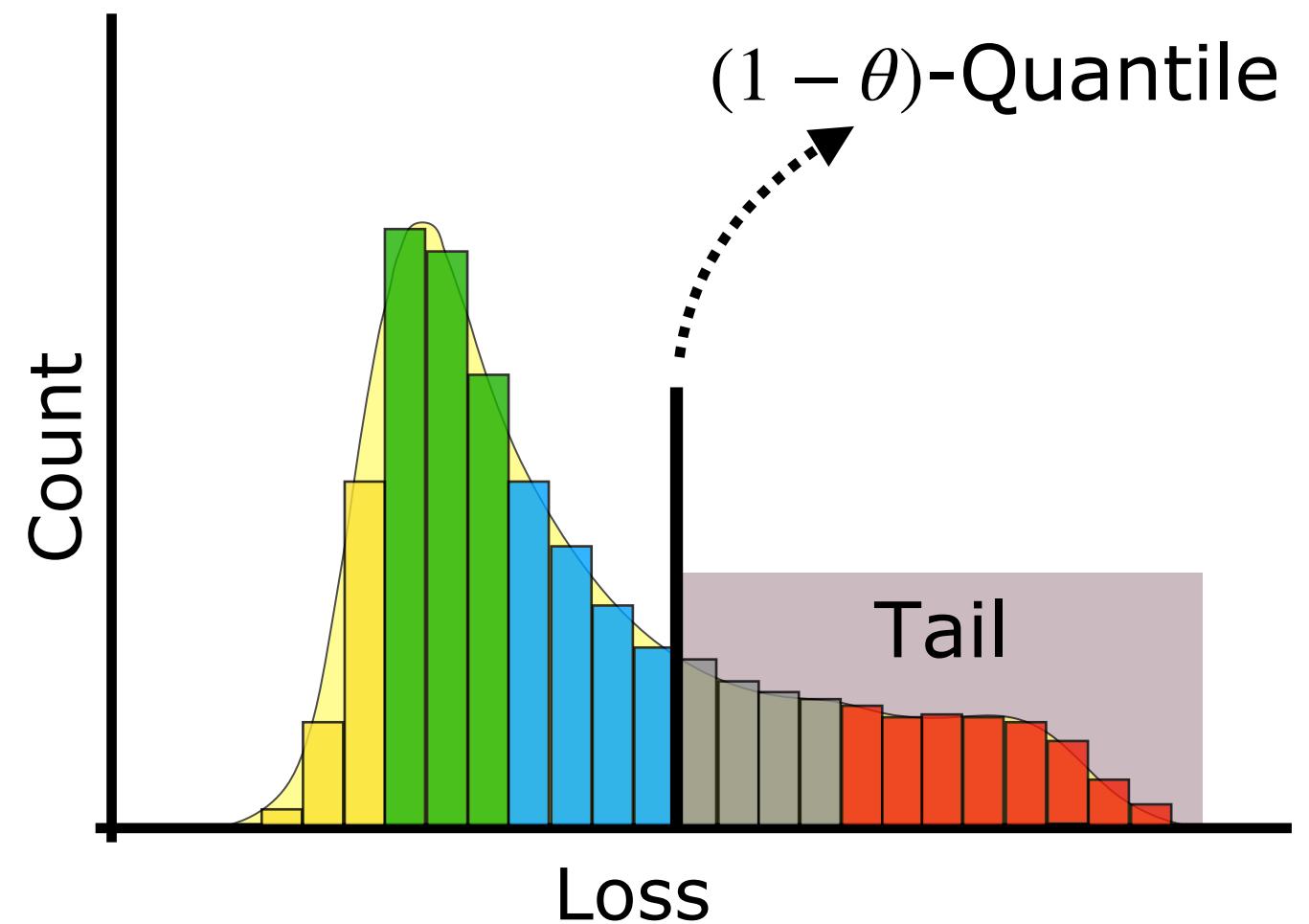
Our Objective:

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



Our Algorithm:

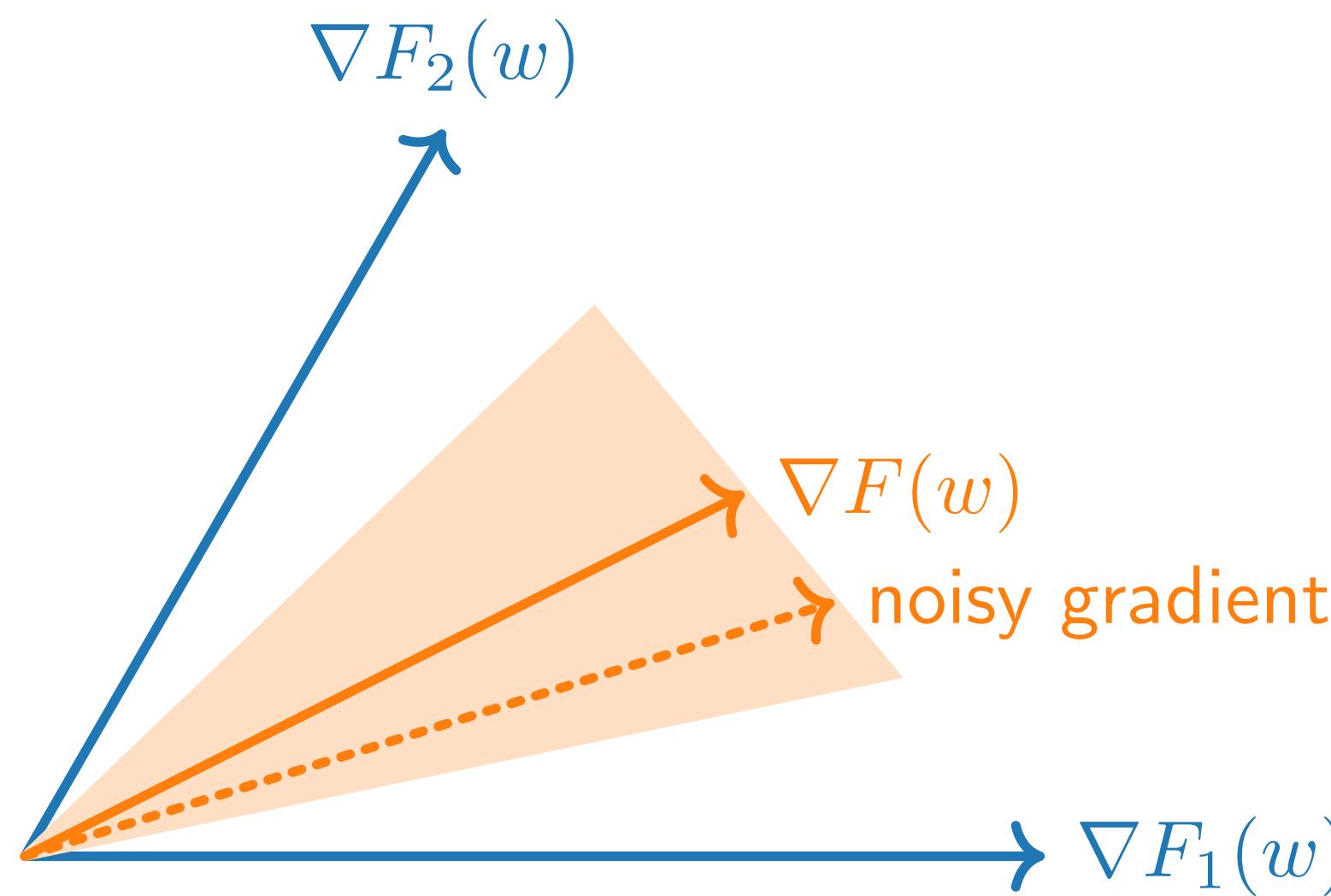
*Step 3 of 3: Aggregate updates contributed by **tail clients** only*



Why is it challenging?

Usual Algorithm:

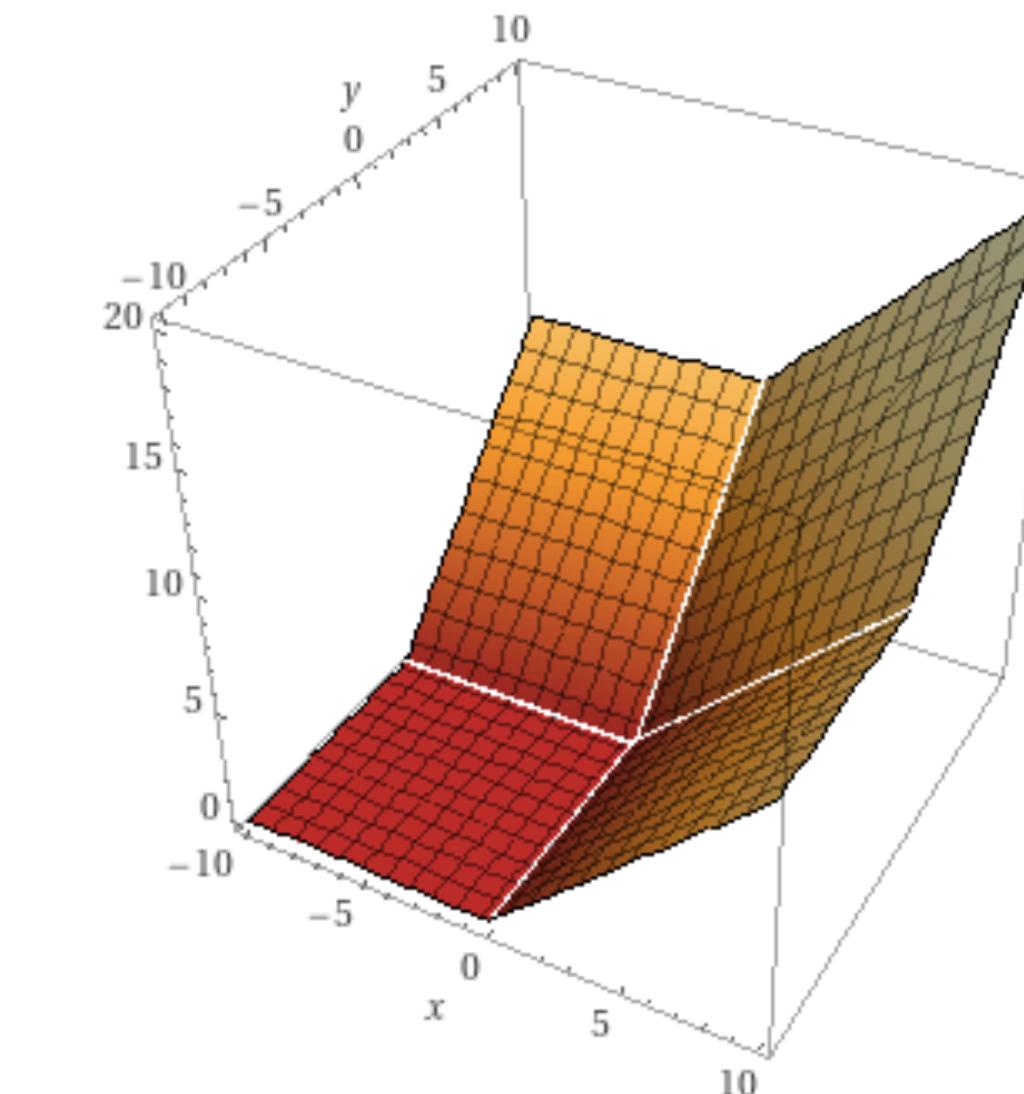
$$\min_w \frac{1}{n} \sum_{i=1}^n F_i(w)$$



Private mean estimation of gradients

Our Algorithm:

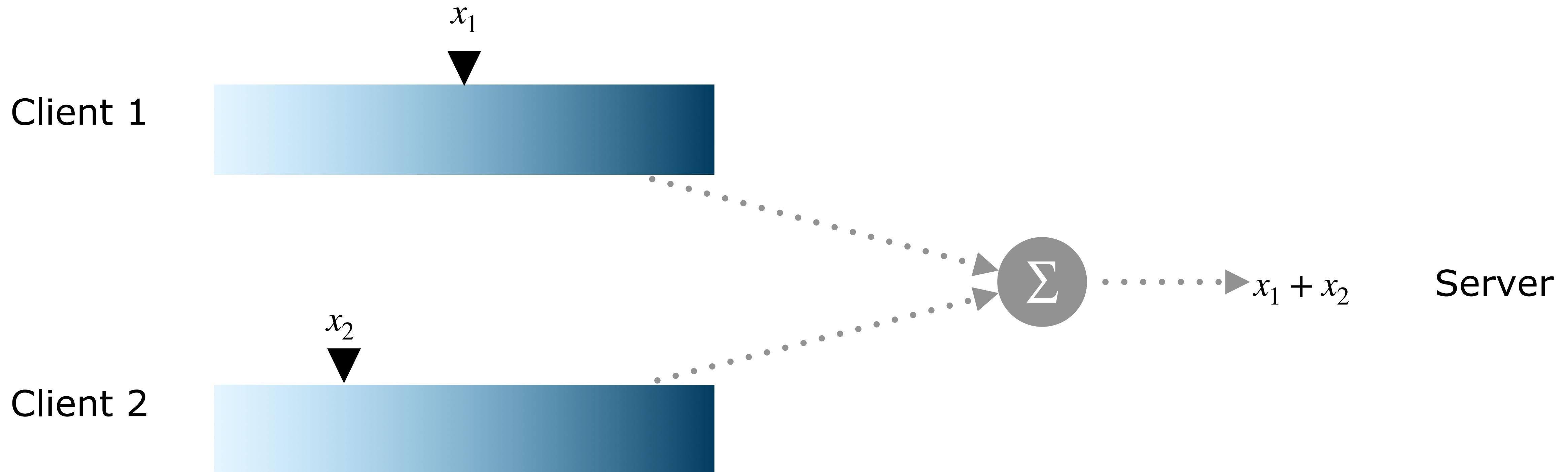
$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$



Non-linear

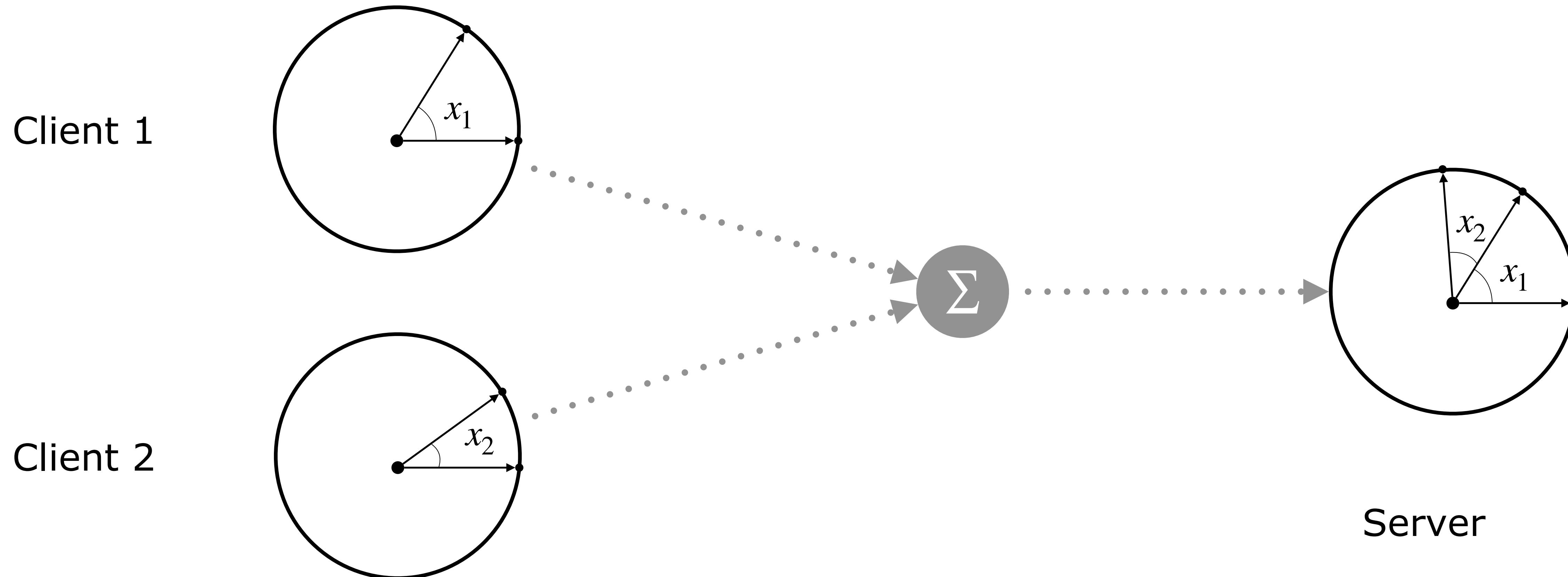
Communication primitive: secure sum

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



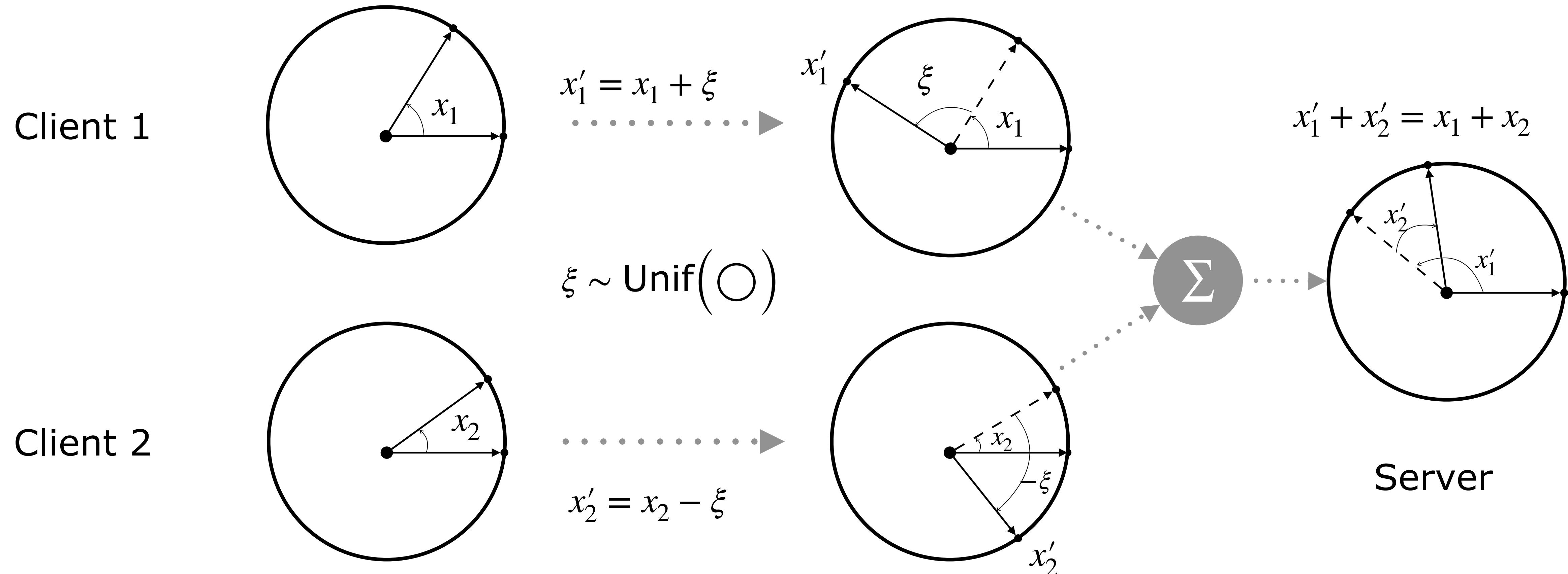
[Bonawitz et al. (CCS 2017), Bell et al. (CCS 2020)]

Perform all operations modulo M



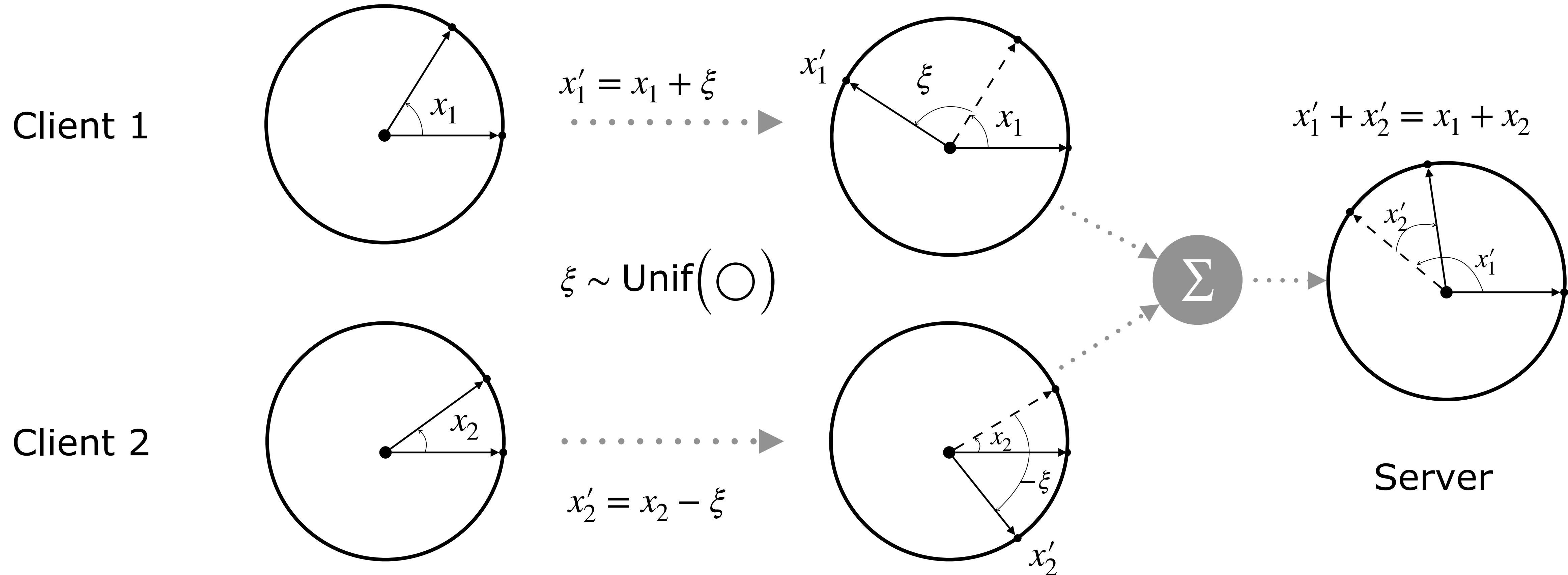
[Bonawitz et al. (CCS 2017), Bell et al. (CCS 2020)]

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:

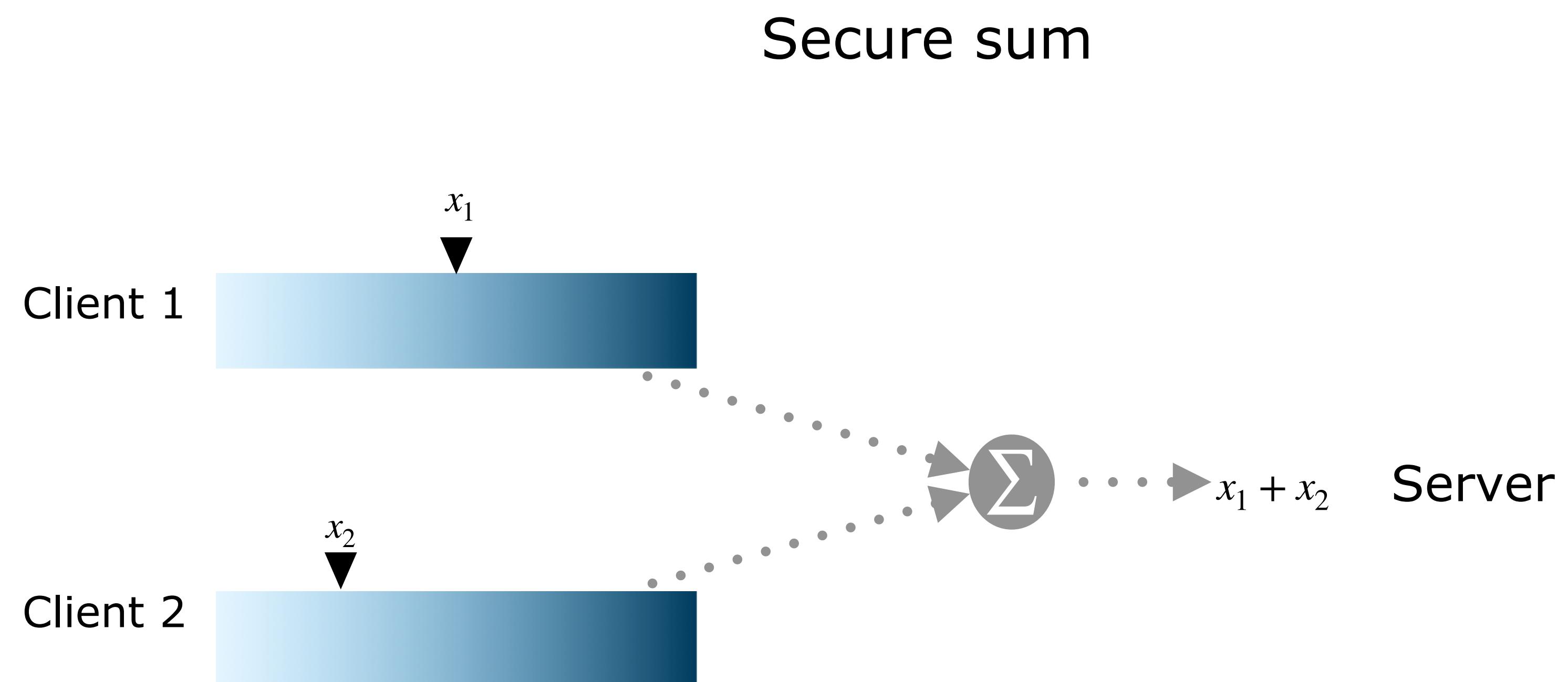
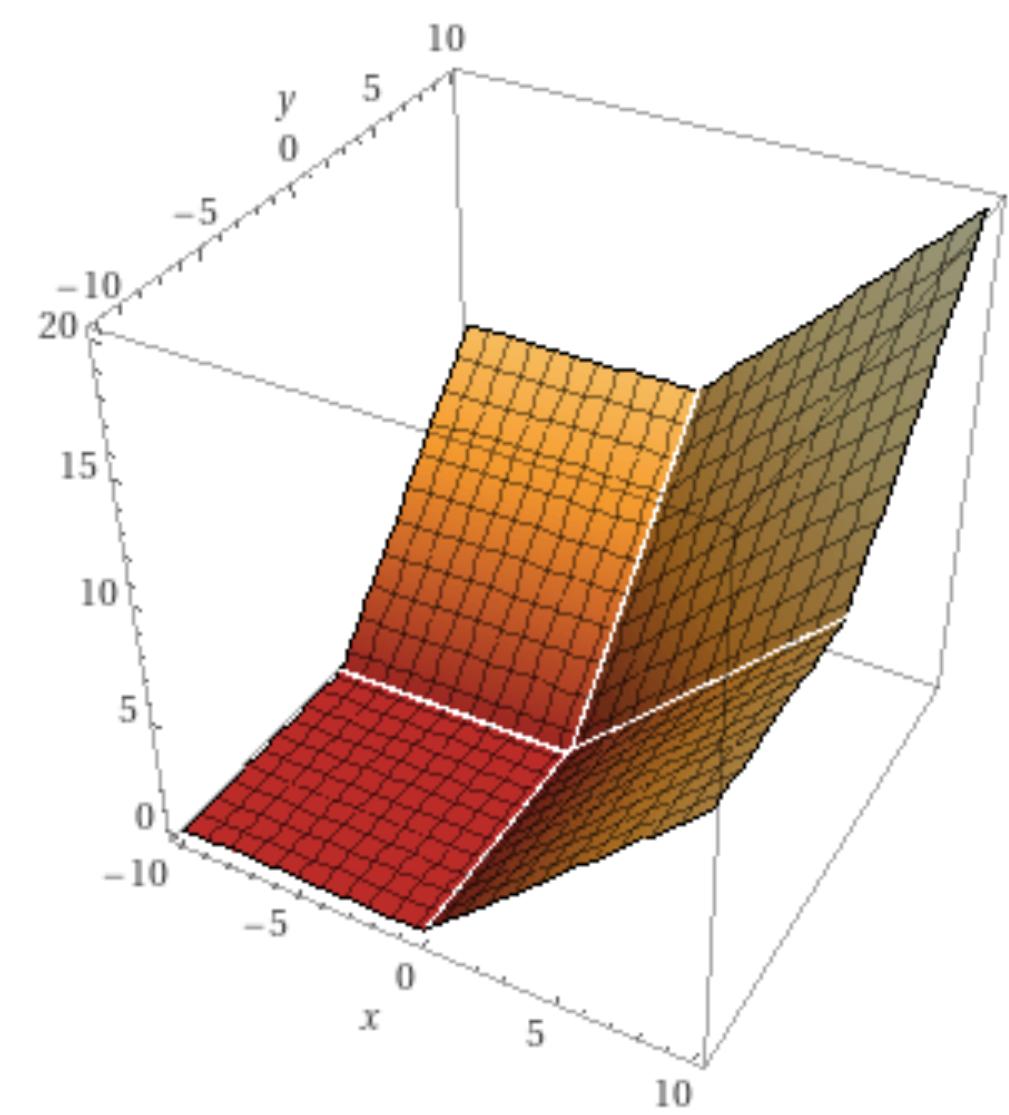
All client-to-server communication must go through secure summation

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

How to achieve non-linear aggregation with a secure sum?

Non-linear aggregate:

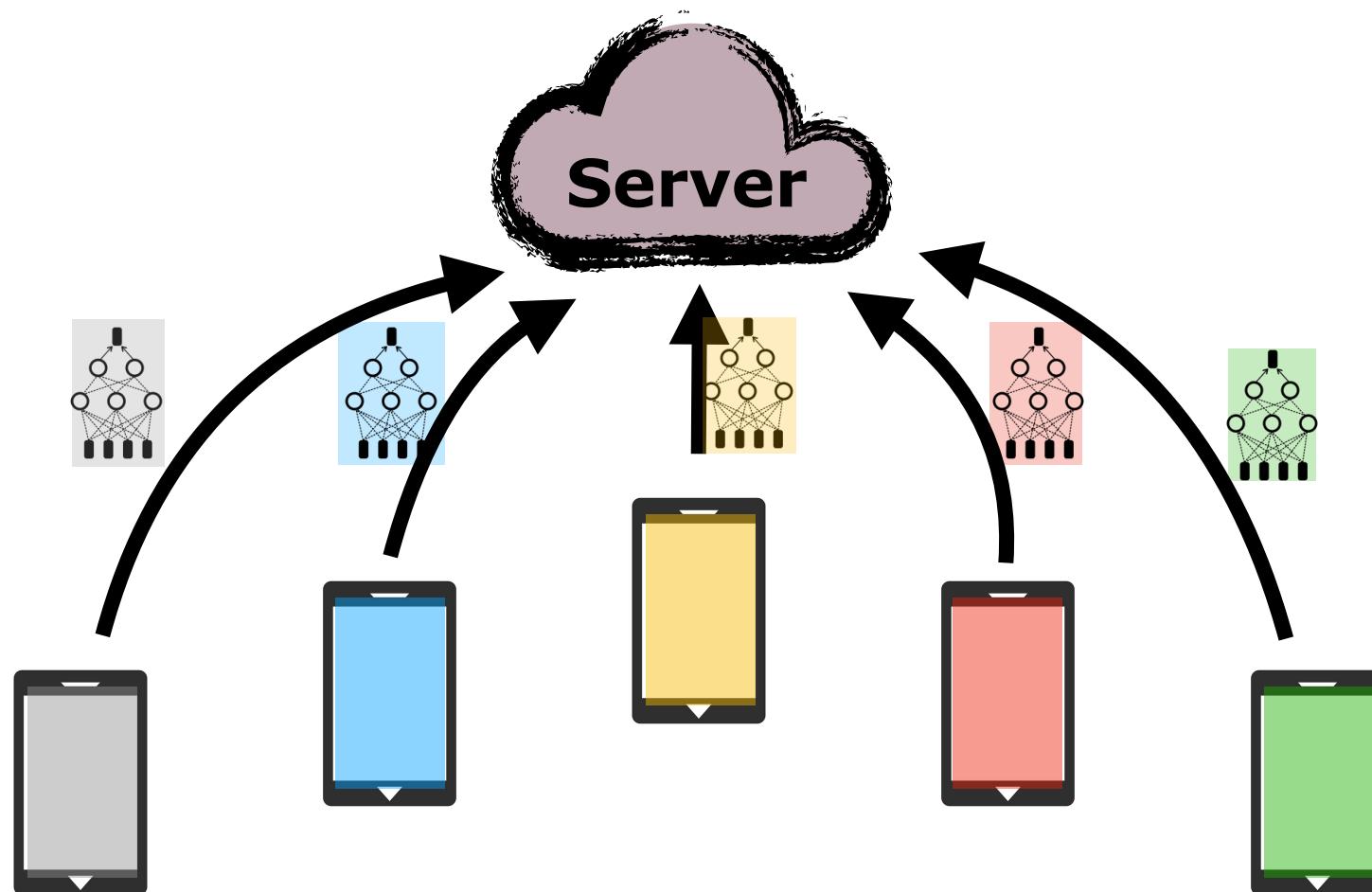
$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$



Usual Algorithm (FedAvg):

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

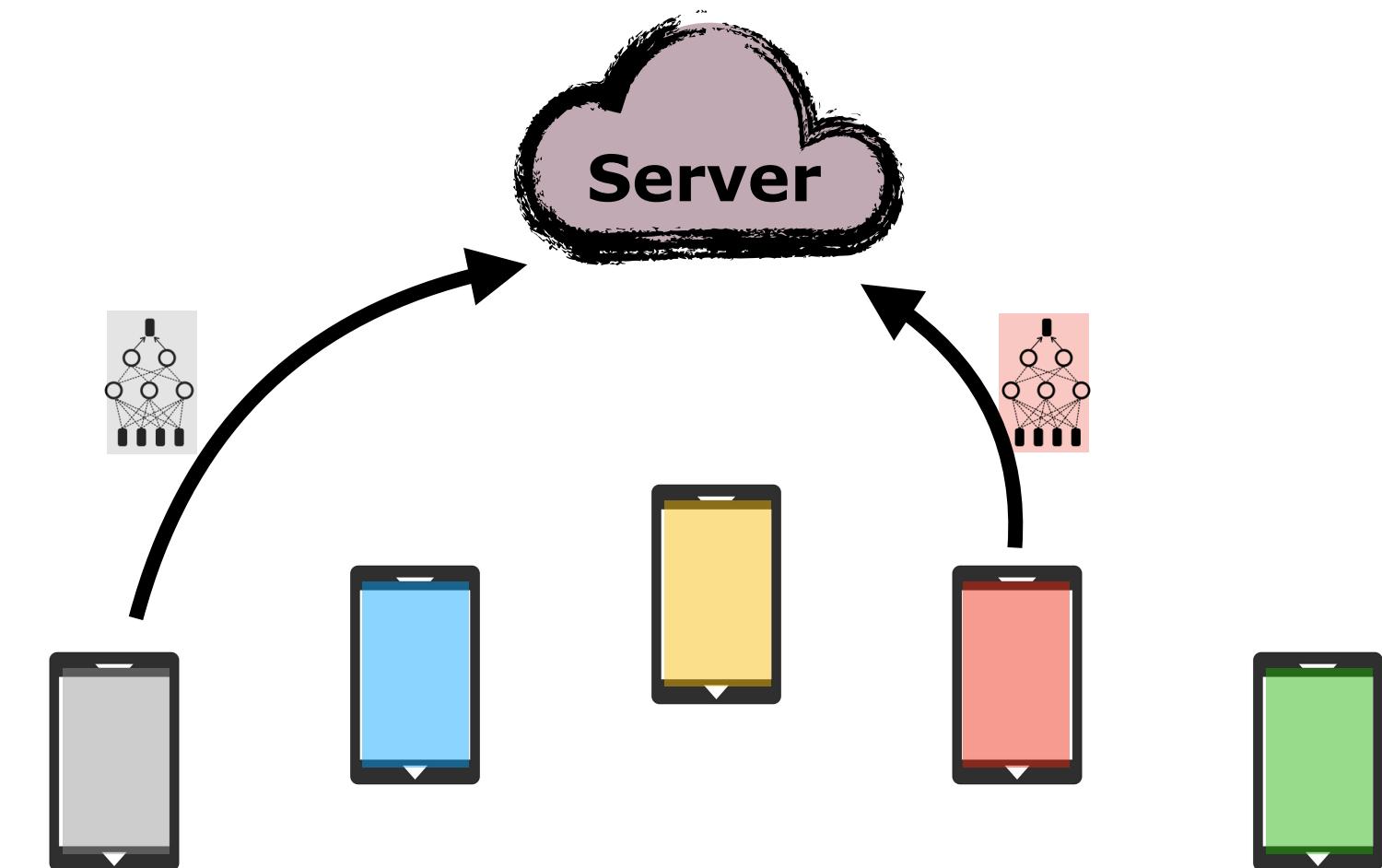
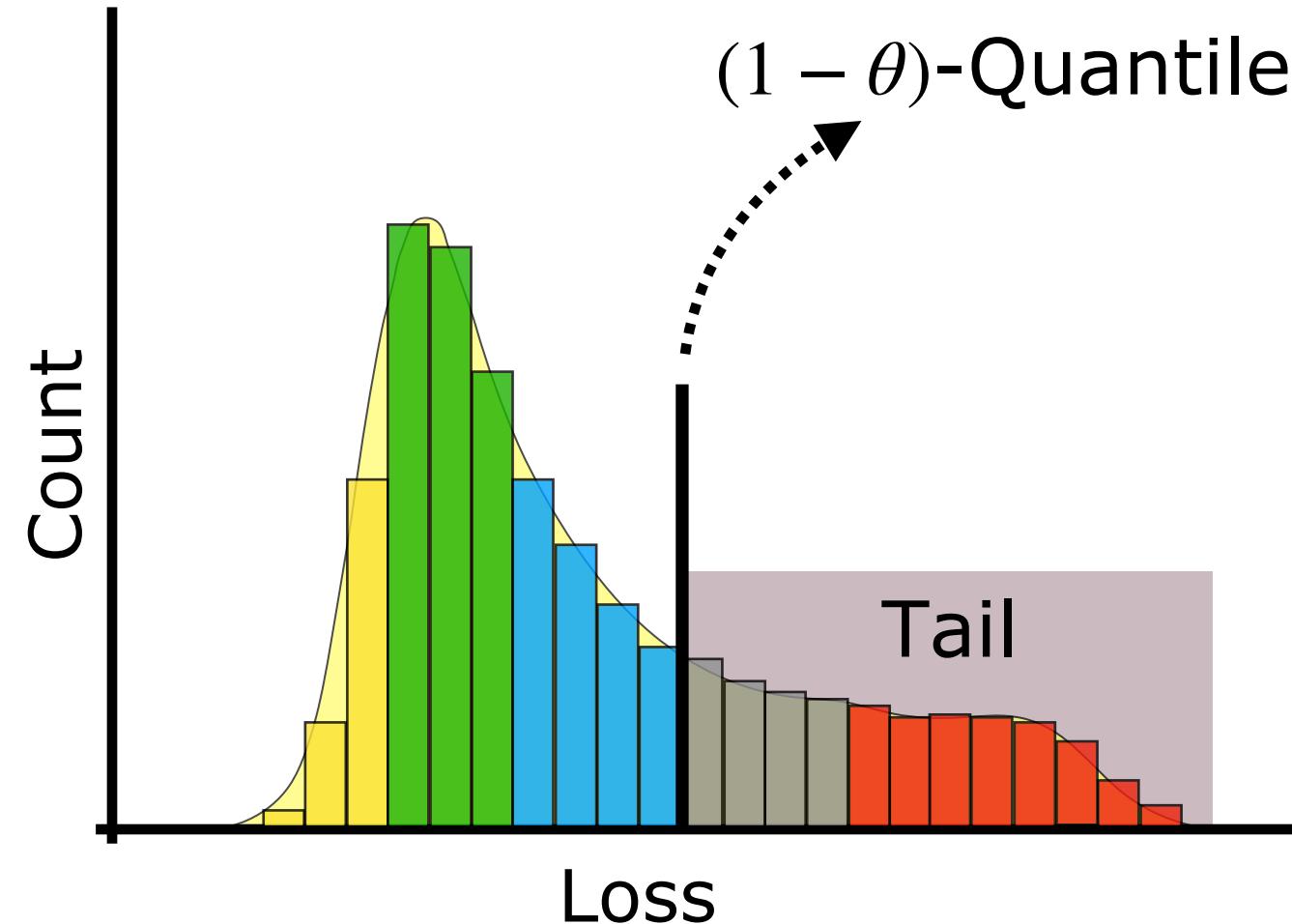
*Step 3 of 3: Aggregate updates contributed by **all clients***



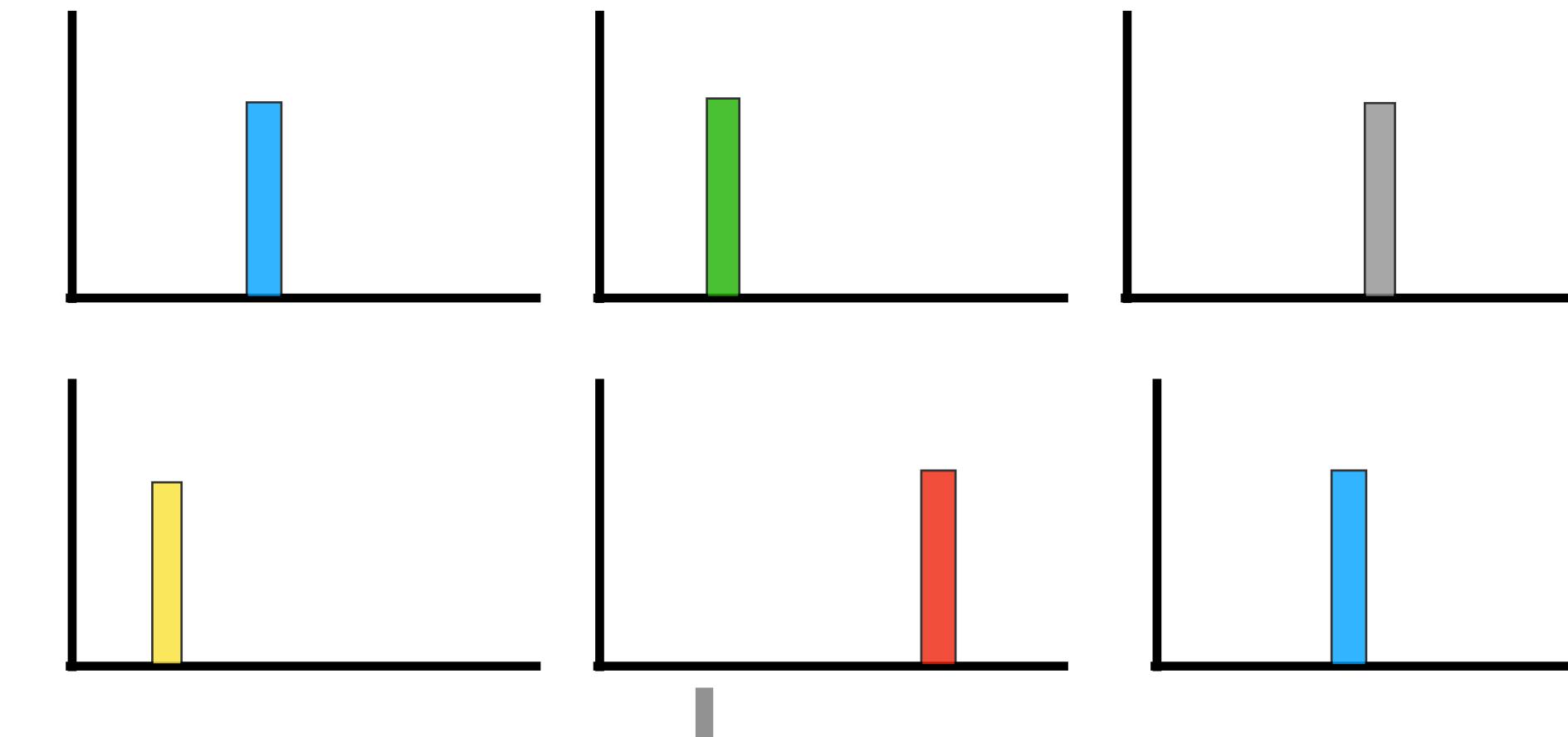
Our Algorithm:

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

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



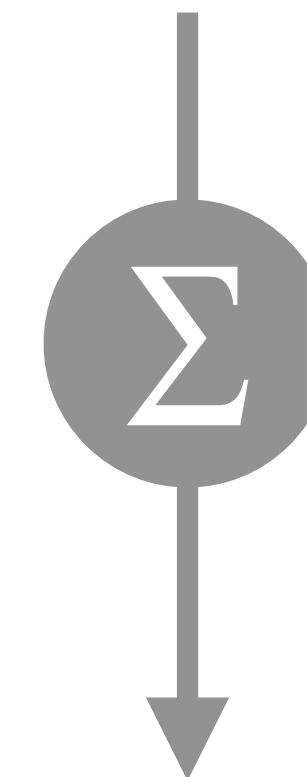
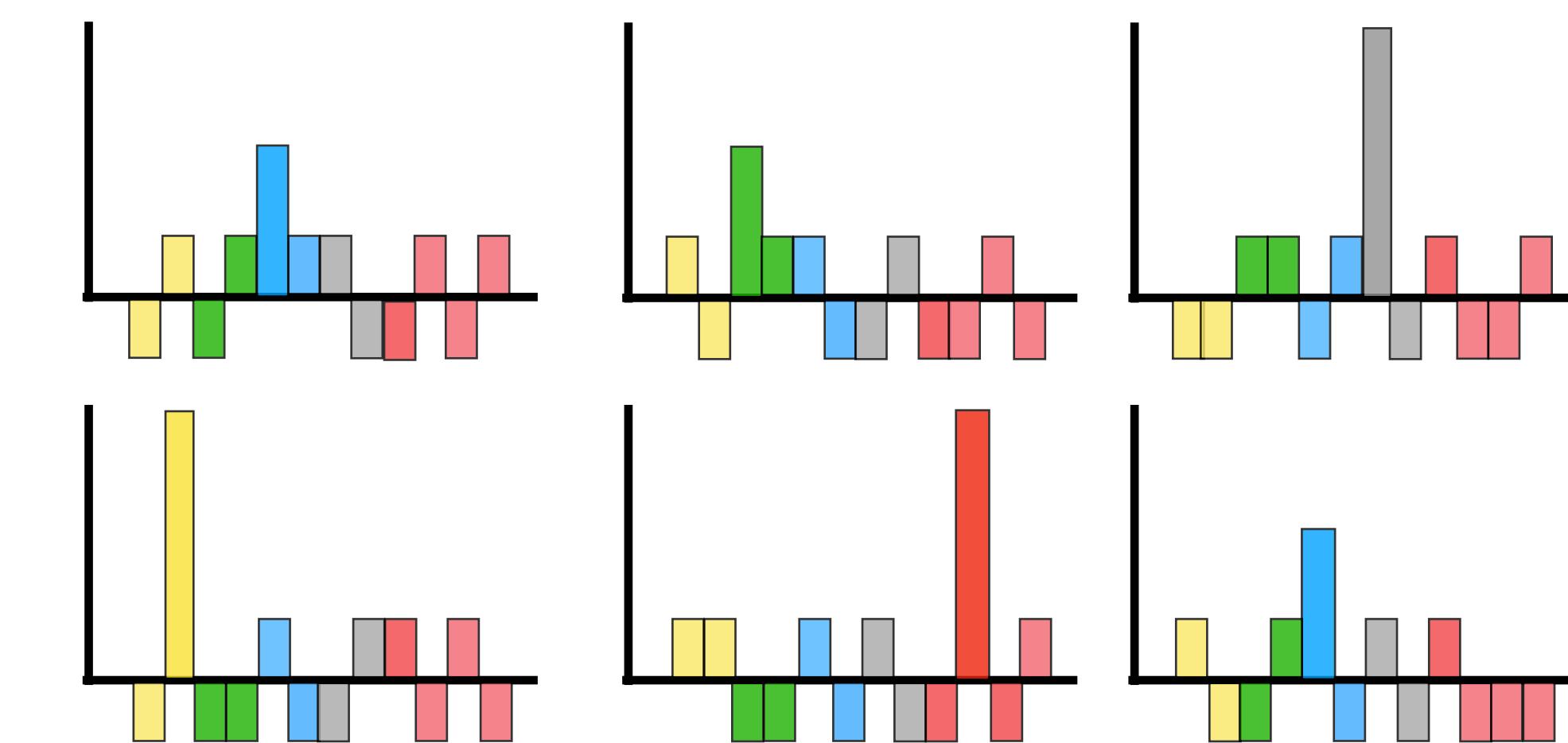
Per-client loss



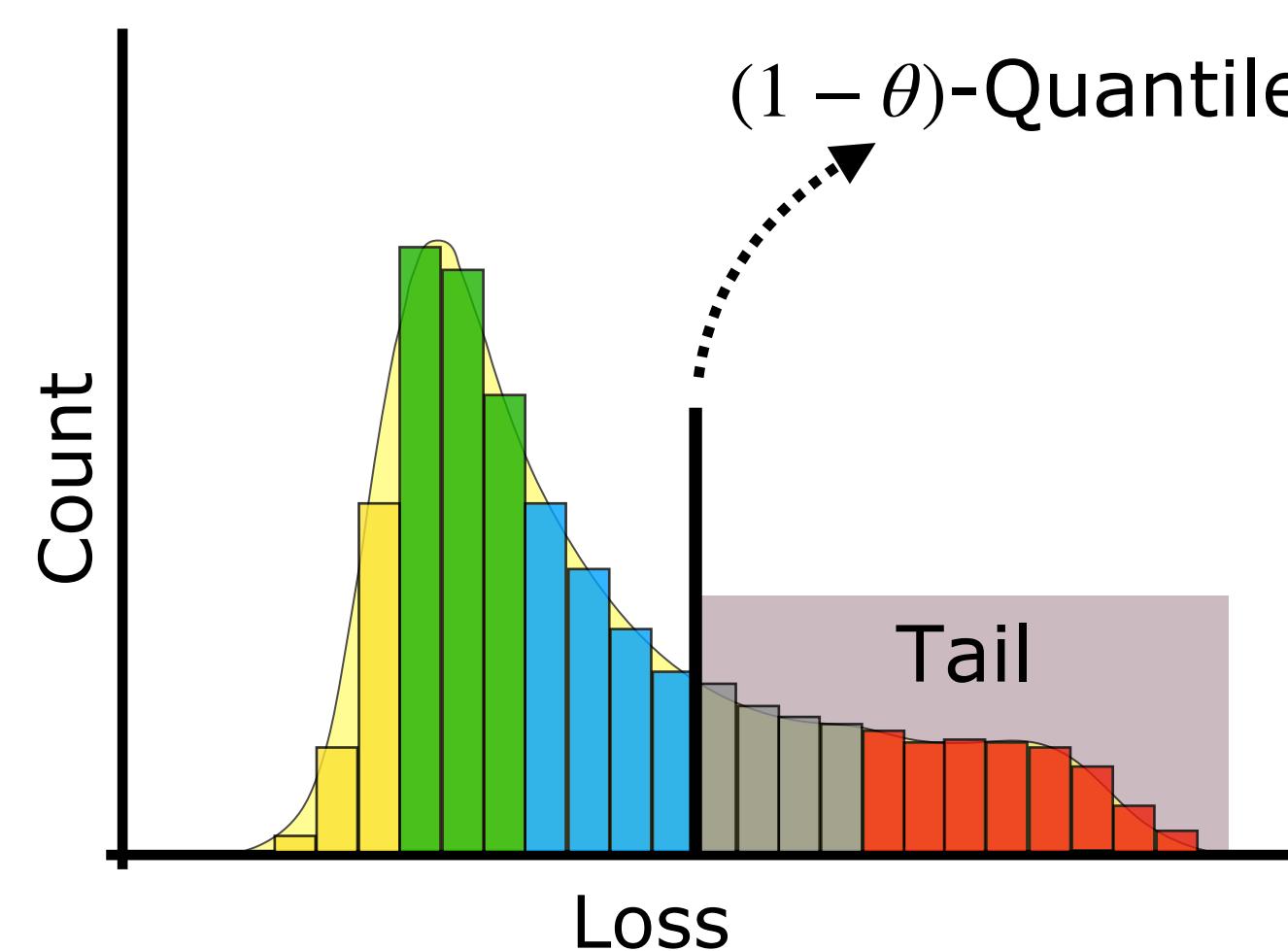
$$h'_i = h_i + \mathcal{N}_{\mathbb{Z}}(0, \sigma^2 I_b)$$



Noisy client loss histogram



Histogram

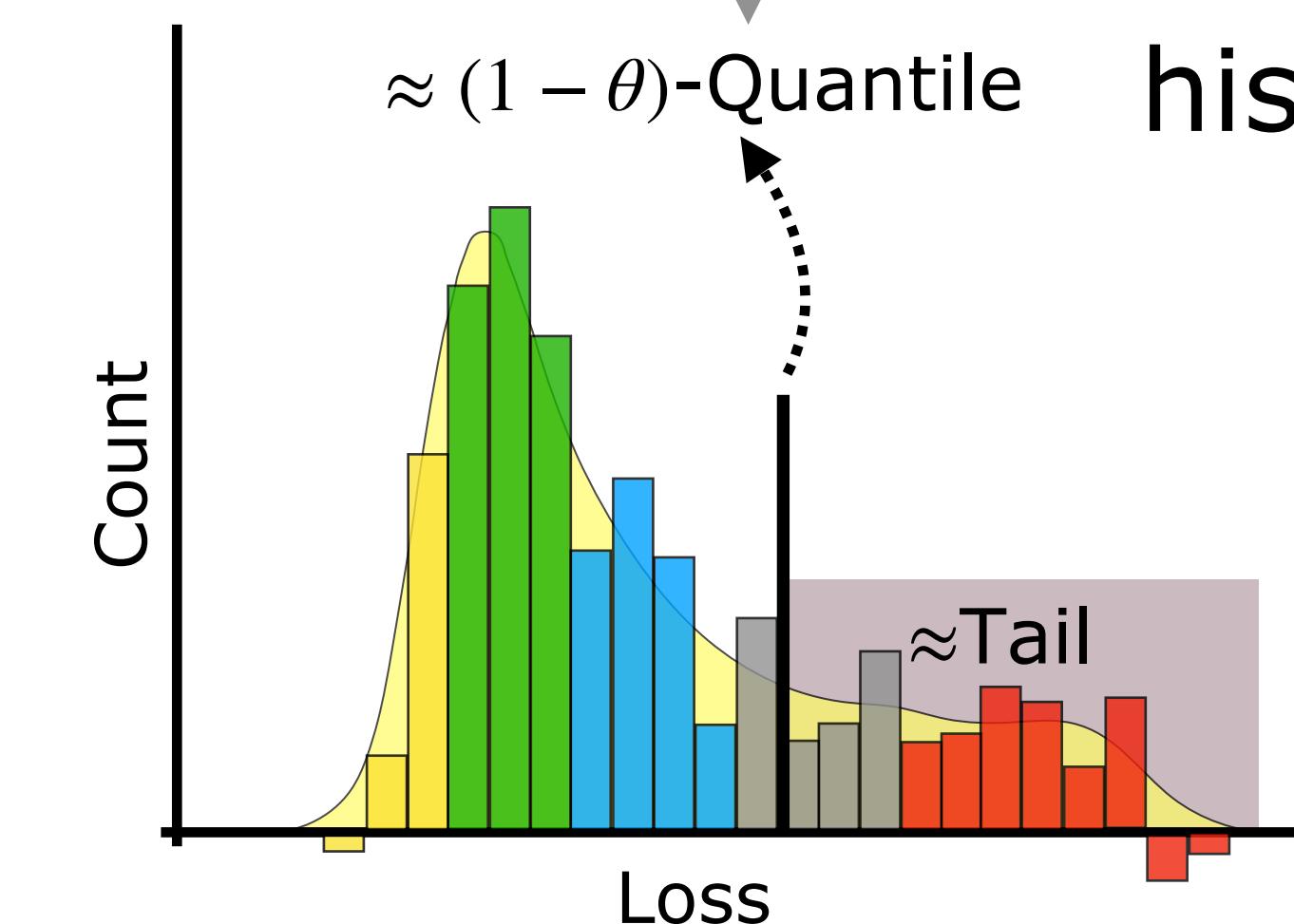


Differential privacy via
discrete Gaussian noise

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

$\approx (1 - \theta)$ -Quantile

Noisy histogram



Proposition (informal) [P., Laguel, Malick, Harchaoui]

If we wish to compute the α -quantile, our algorithm returns an ε -differentially private α' -quantile where

$$|\alpha' - \alpha| \leq \frac{\sqrt{b}}{\varepsilon m}$$

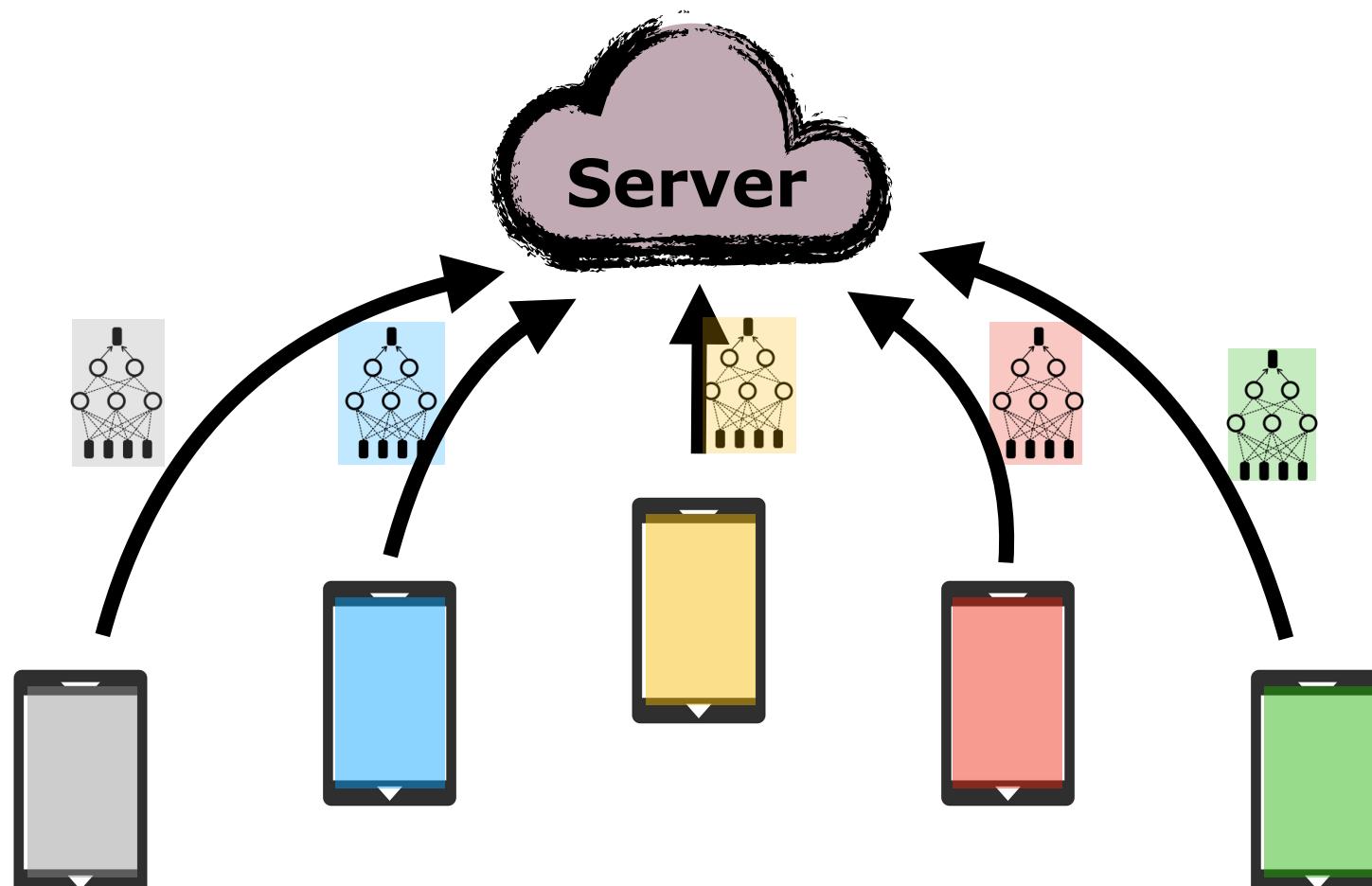
Total communication cost $\approx bm \log^2 m$

m #clients per round
 b #bins in the histogram
 ε privacy parameter

Usual Algorithm (FedAvg):

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

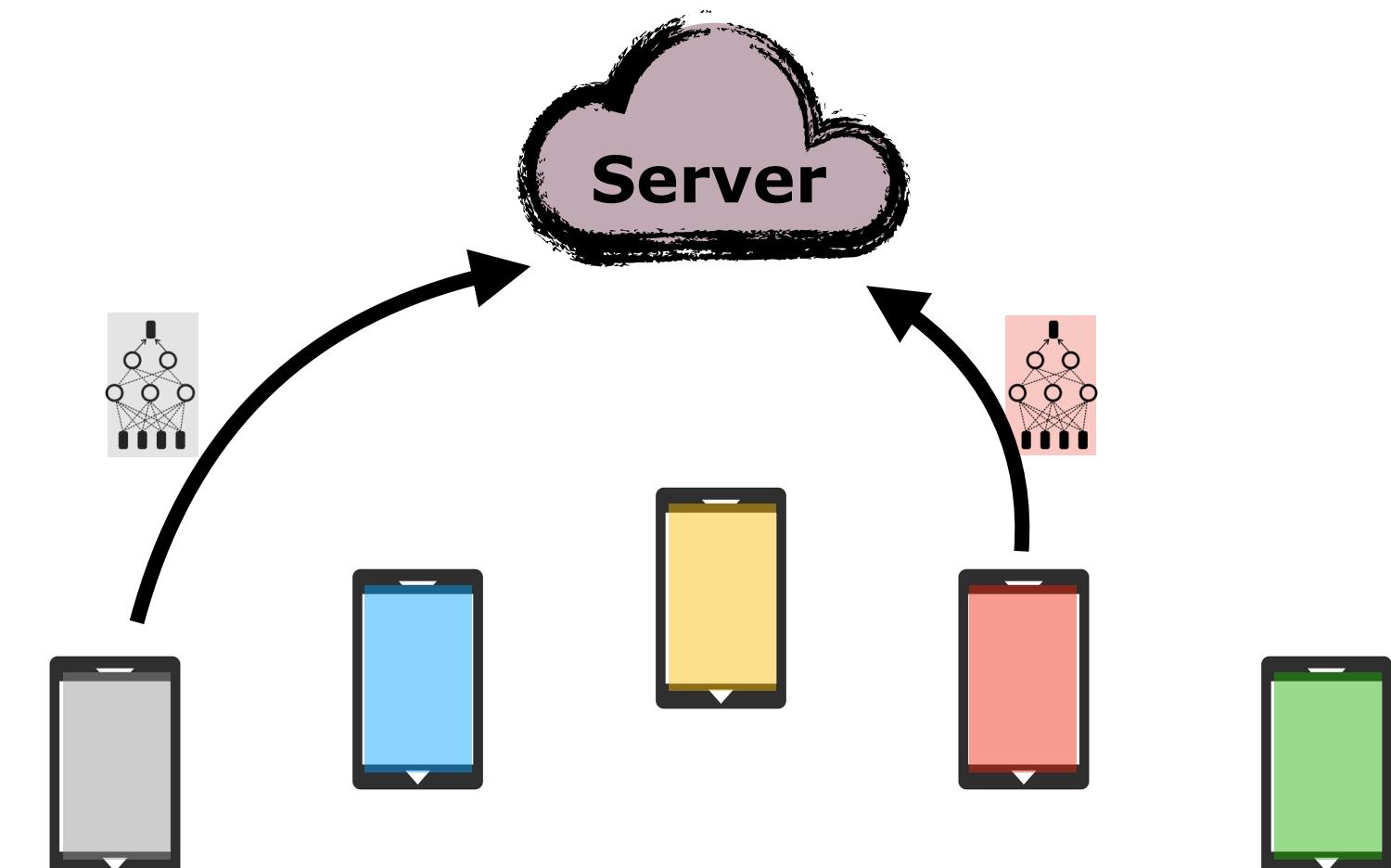
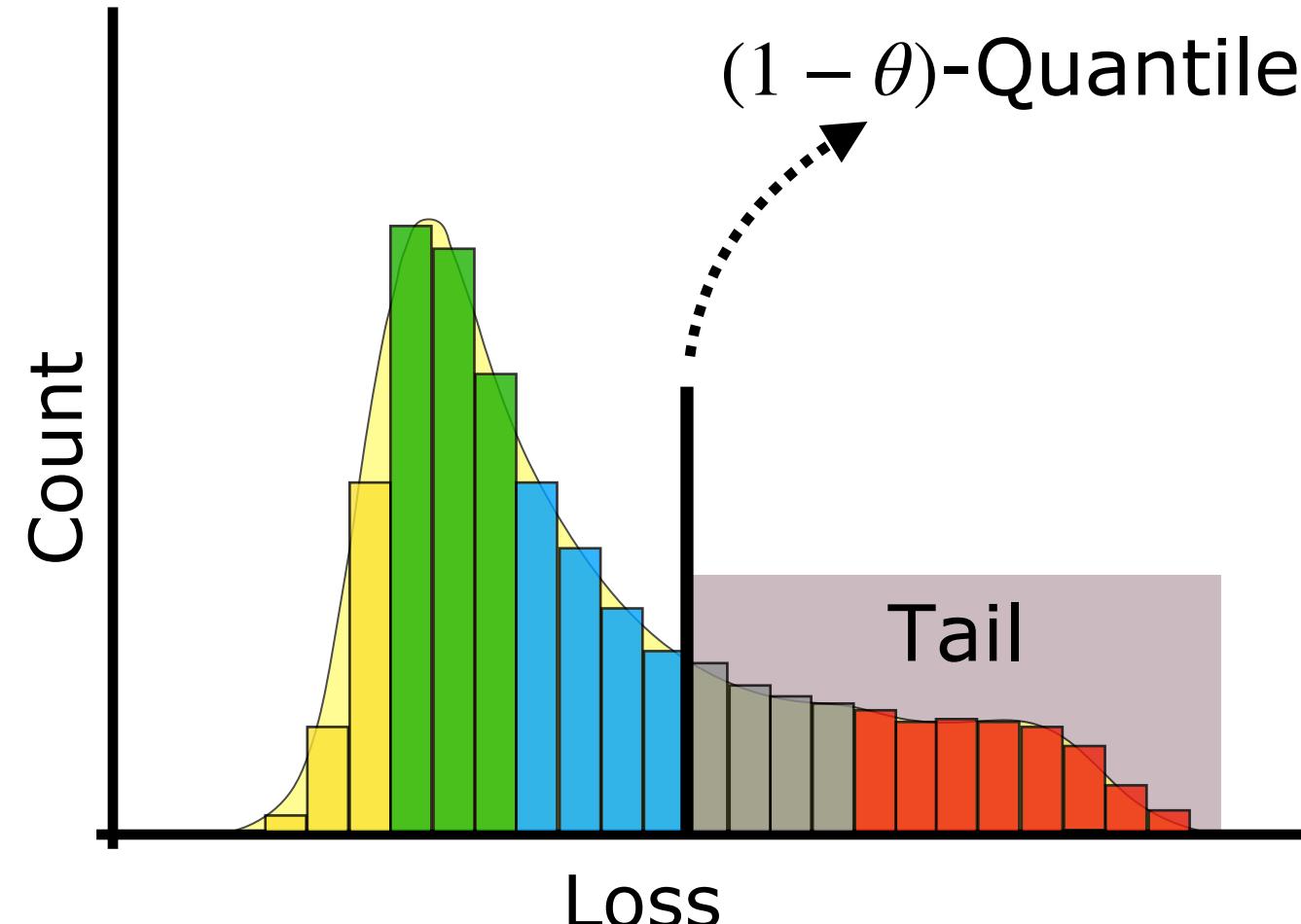
*Step 3 of 3: Aggregate updates contributed by **all clients***



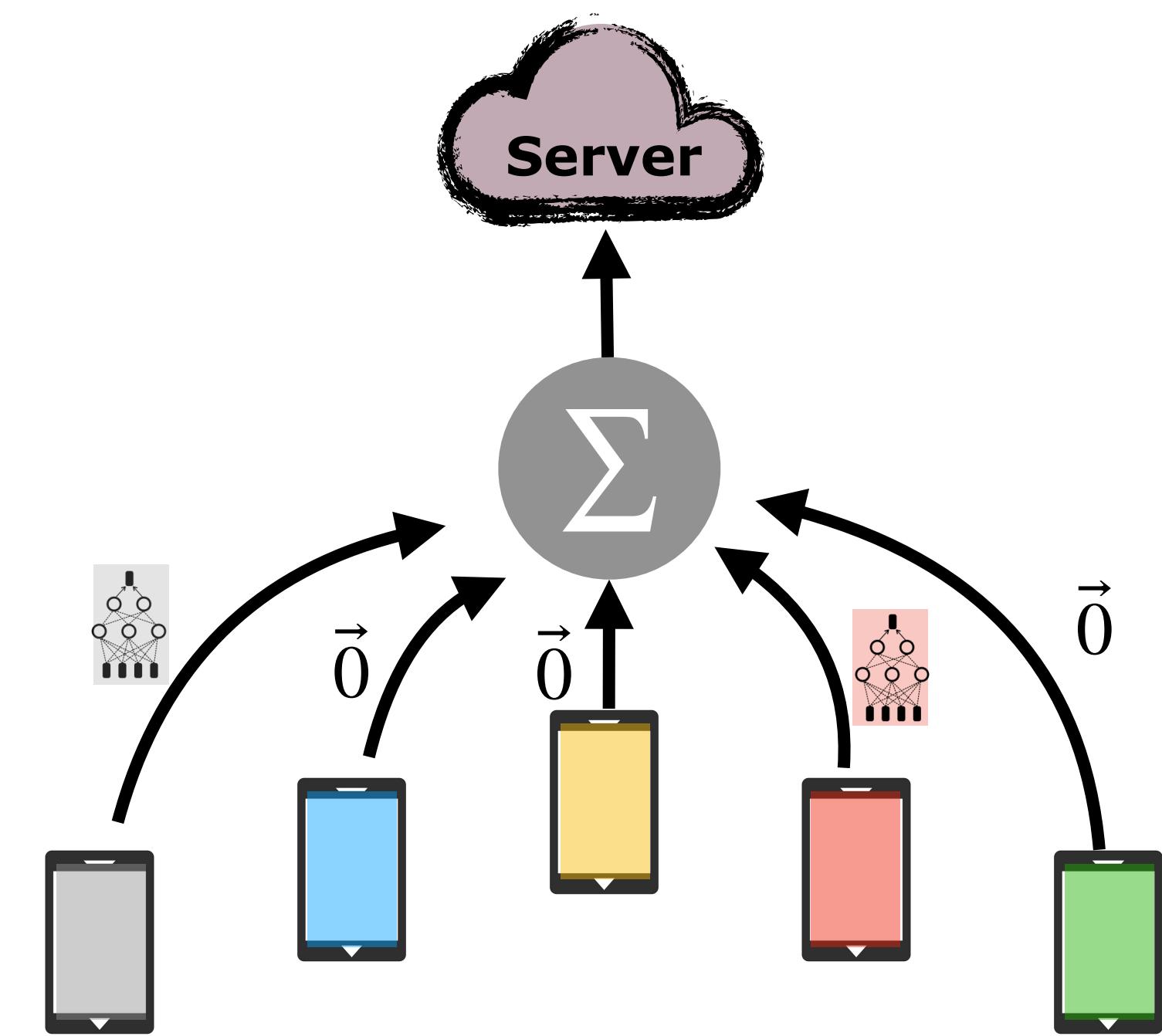
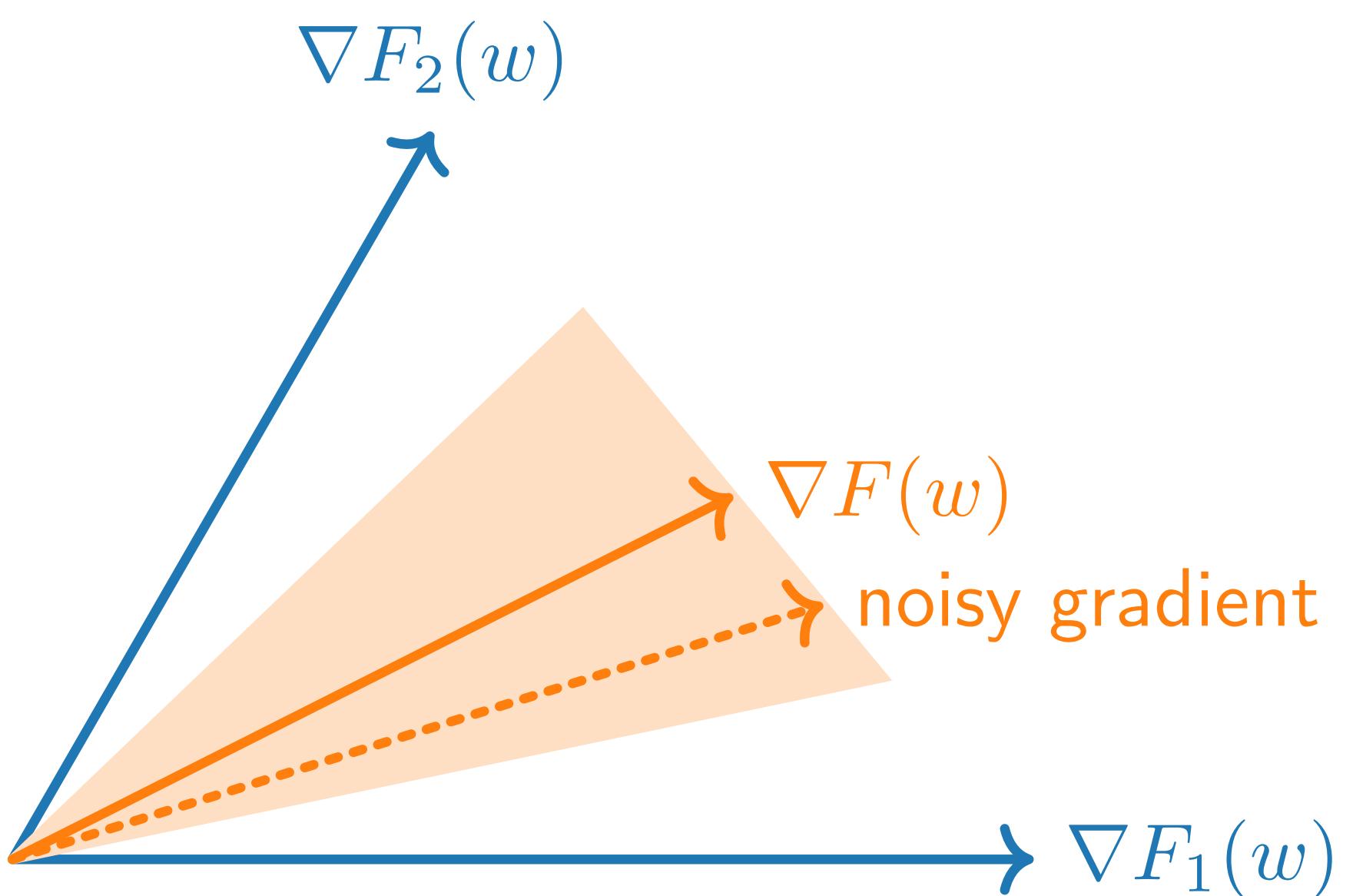
Our Algorithm:

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

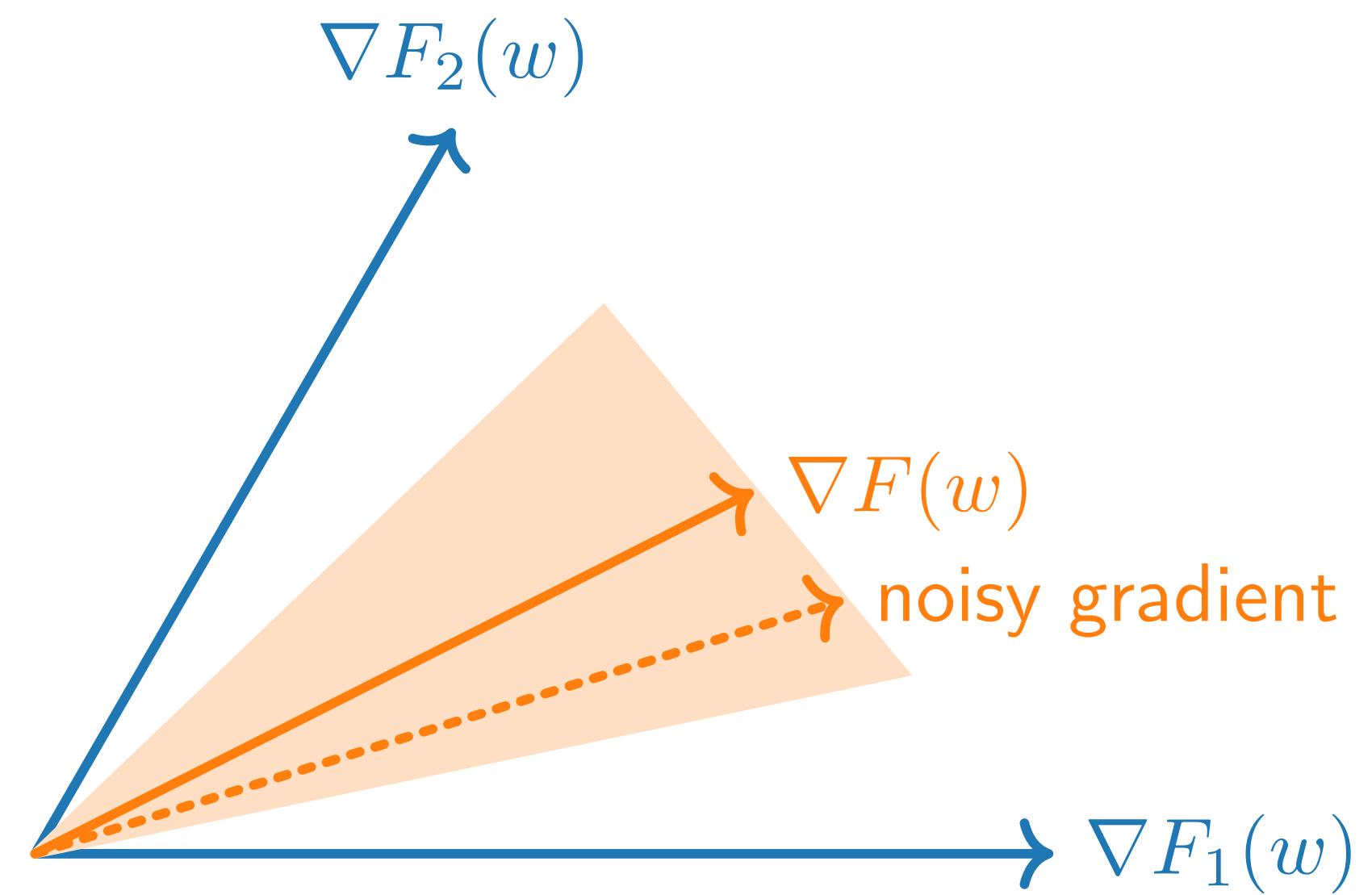
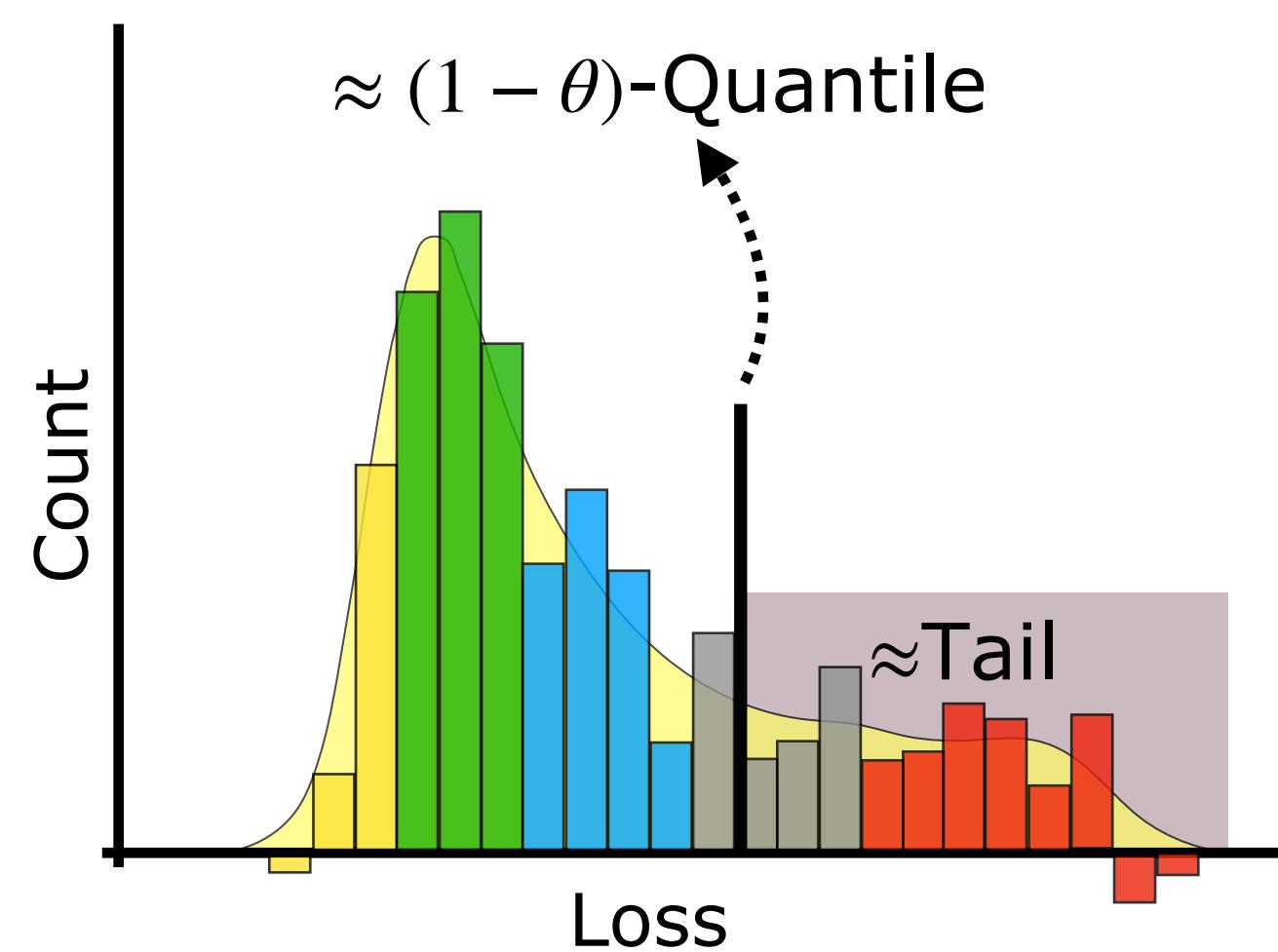
*Step 3 of 3: Aggregate updates contributed by **tail clients** only*



Private mean estimation of (potentially zeroed out) gradients



Total privacy leakage =
Quantile privacy leakage + Parameter privacy leakage



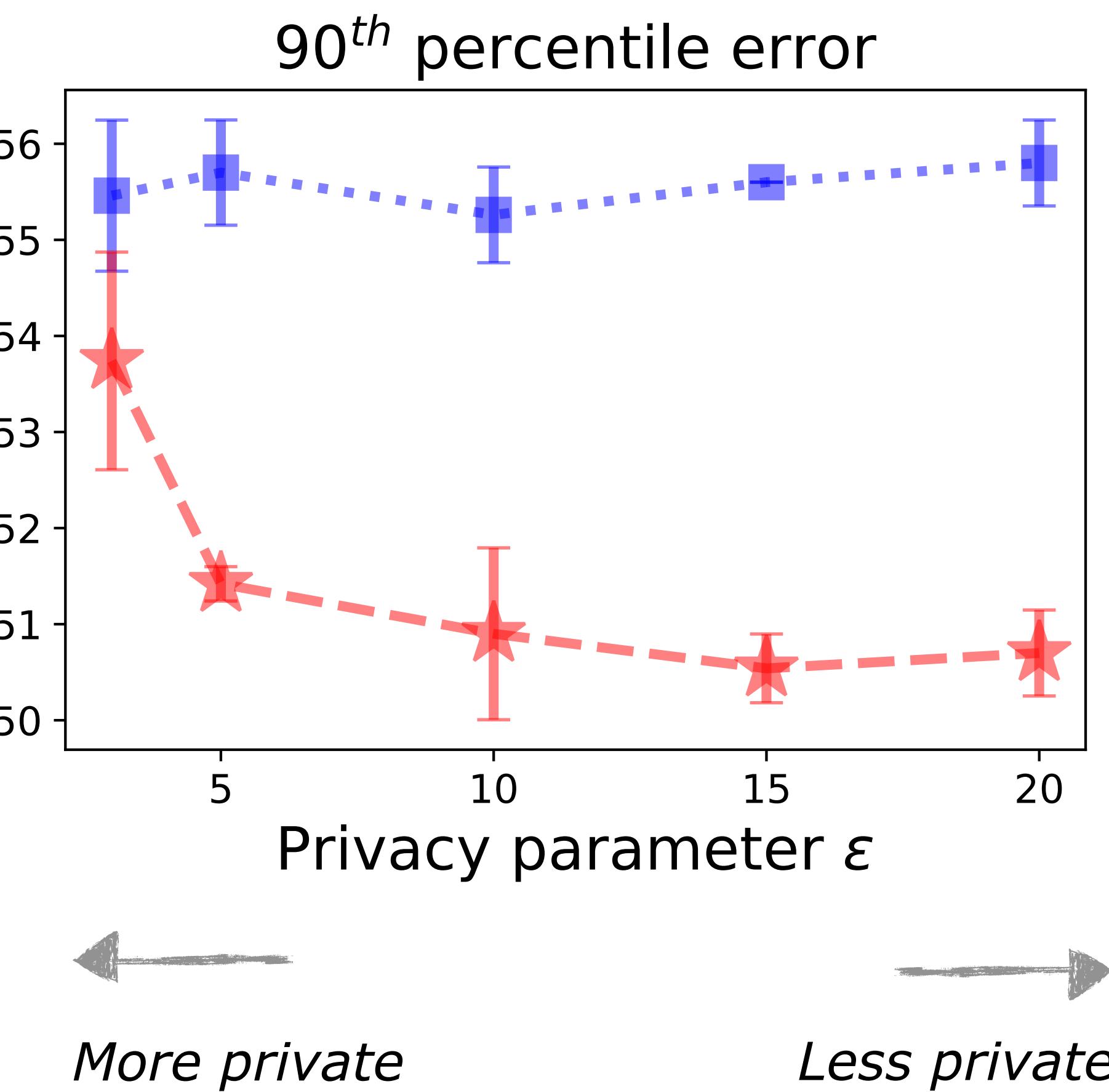
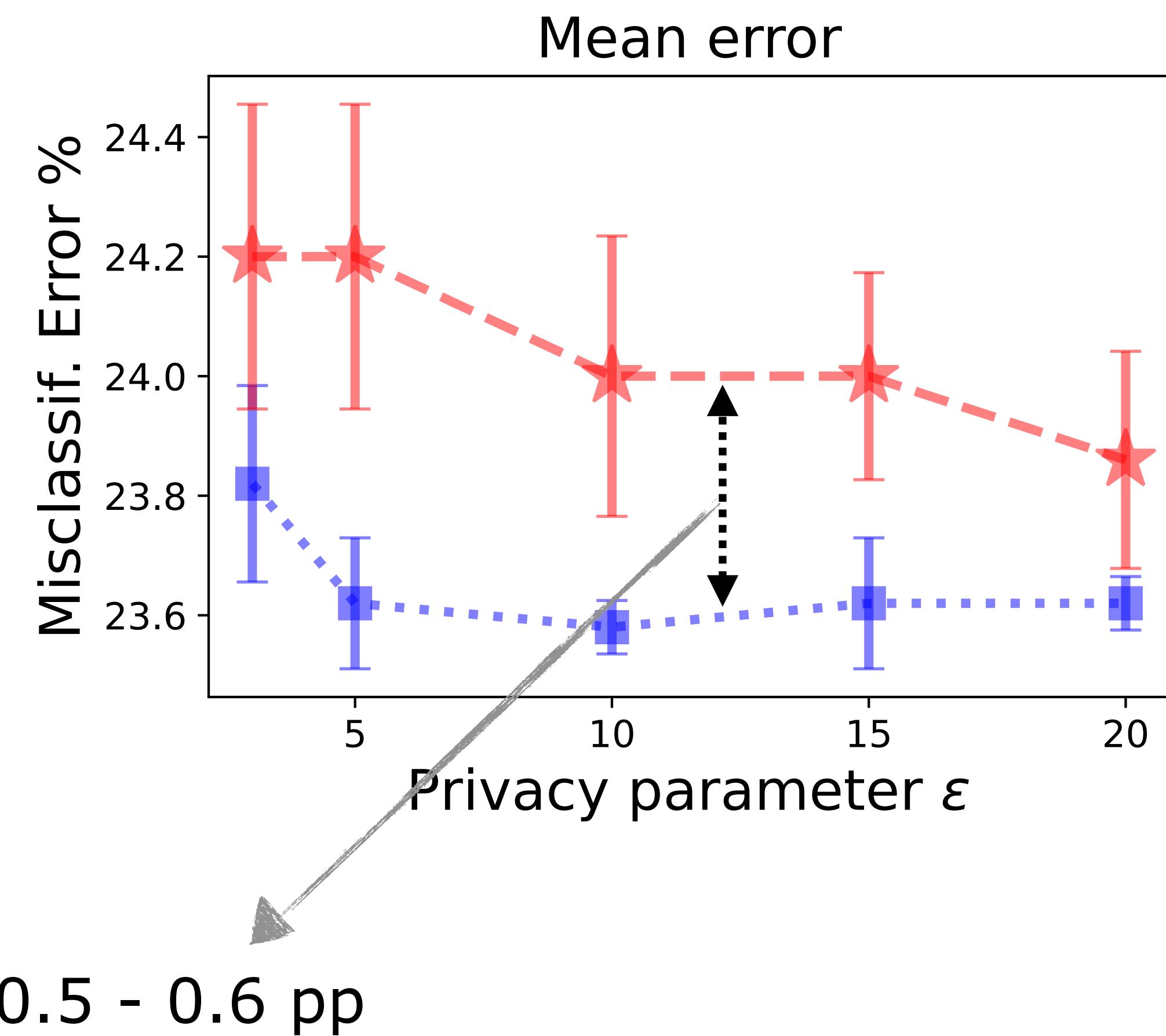


Privacy of user data:

First end-to-end **differentially private** algorithm for distributionally robust federated learning



Algorithm requires **2** secure summations per update



Usual
Ours
5 pp

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/



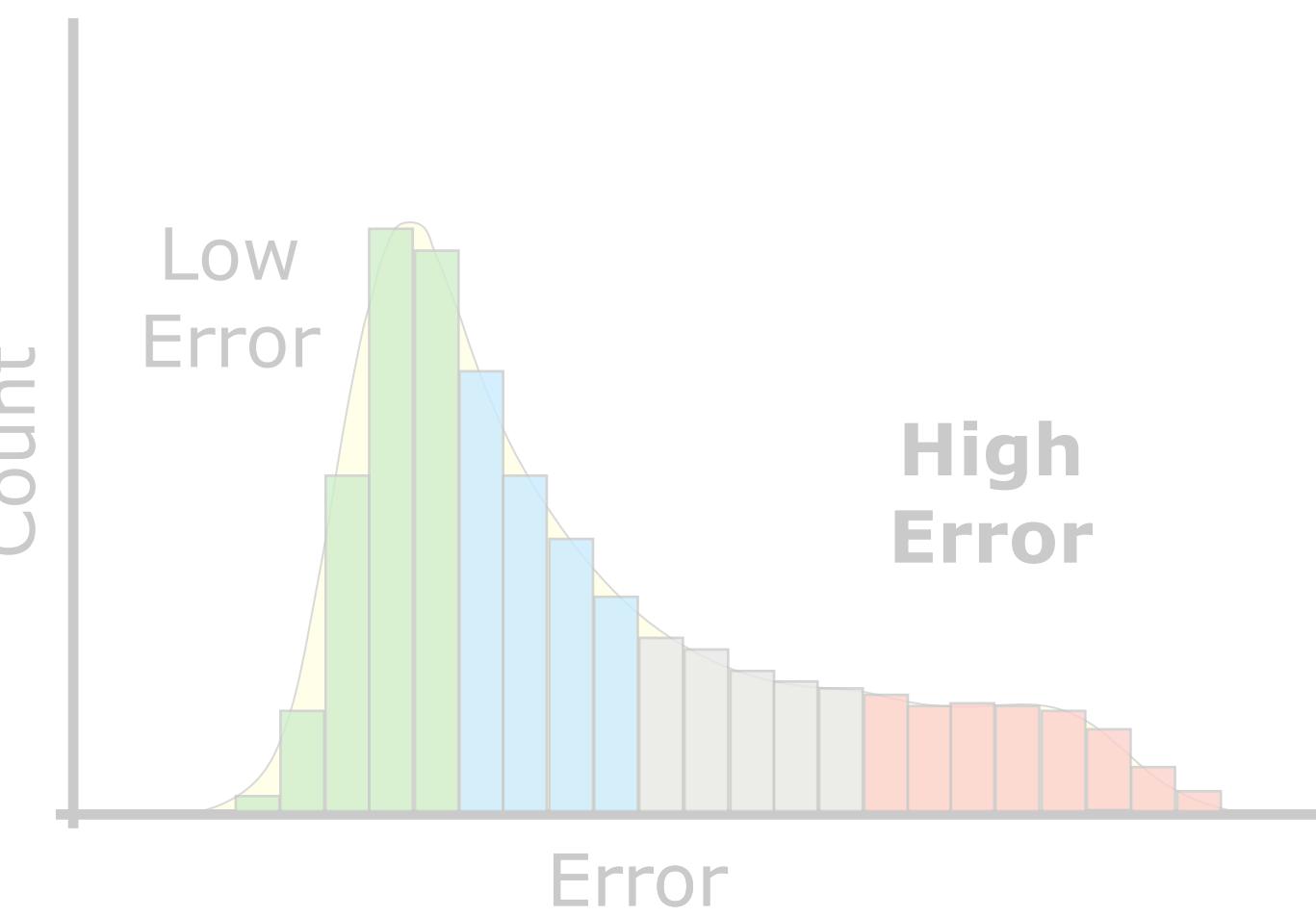
Summary: Tackling distribution shifts in federated learning

Distribution shift \Rightarrow
large tail errors

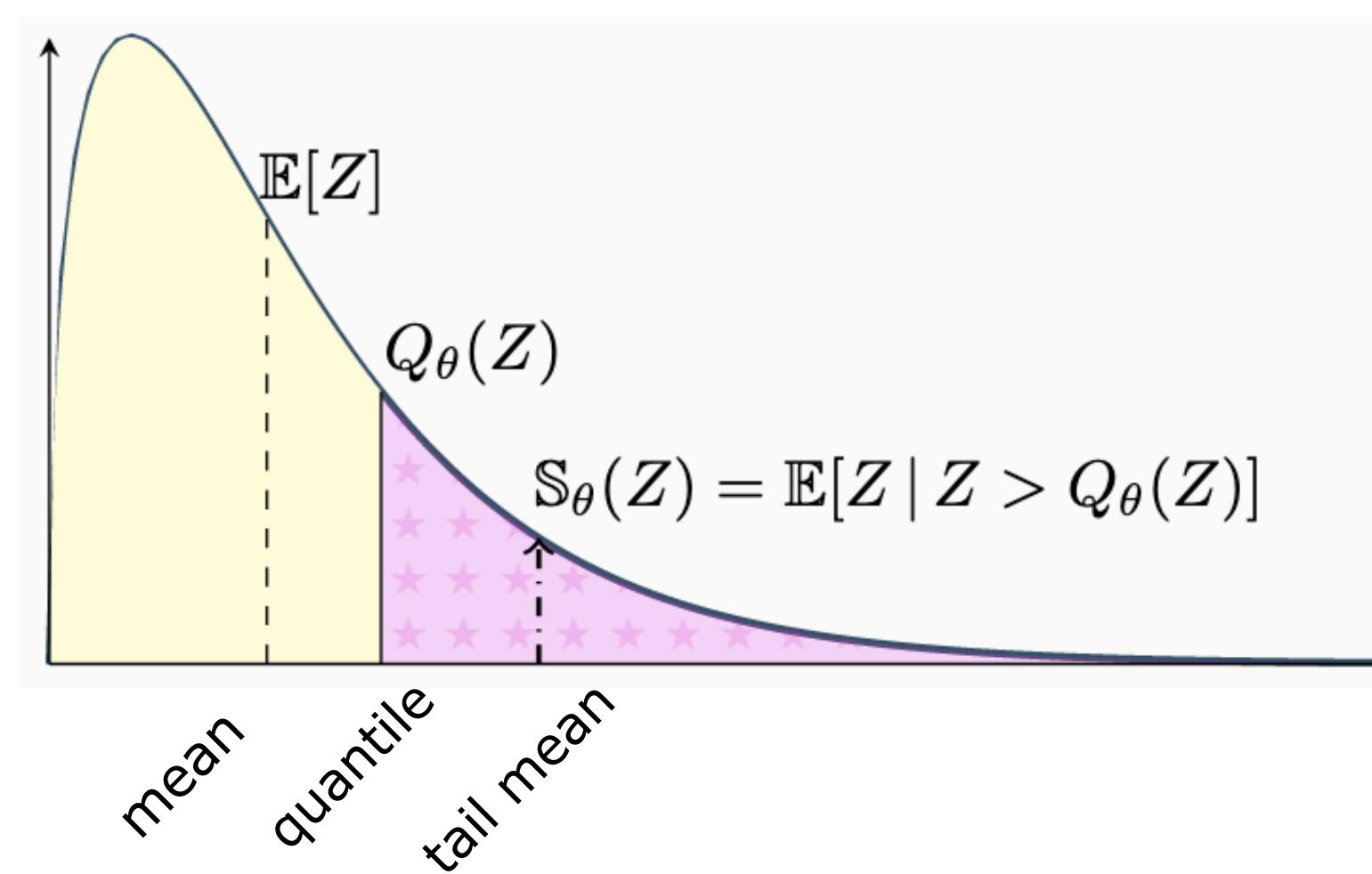


Summary: Tackling distribution shifts in federated learning

Distribution shift \Rightarrow
large tail errors

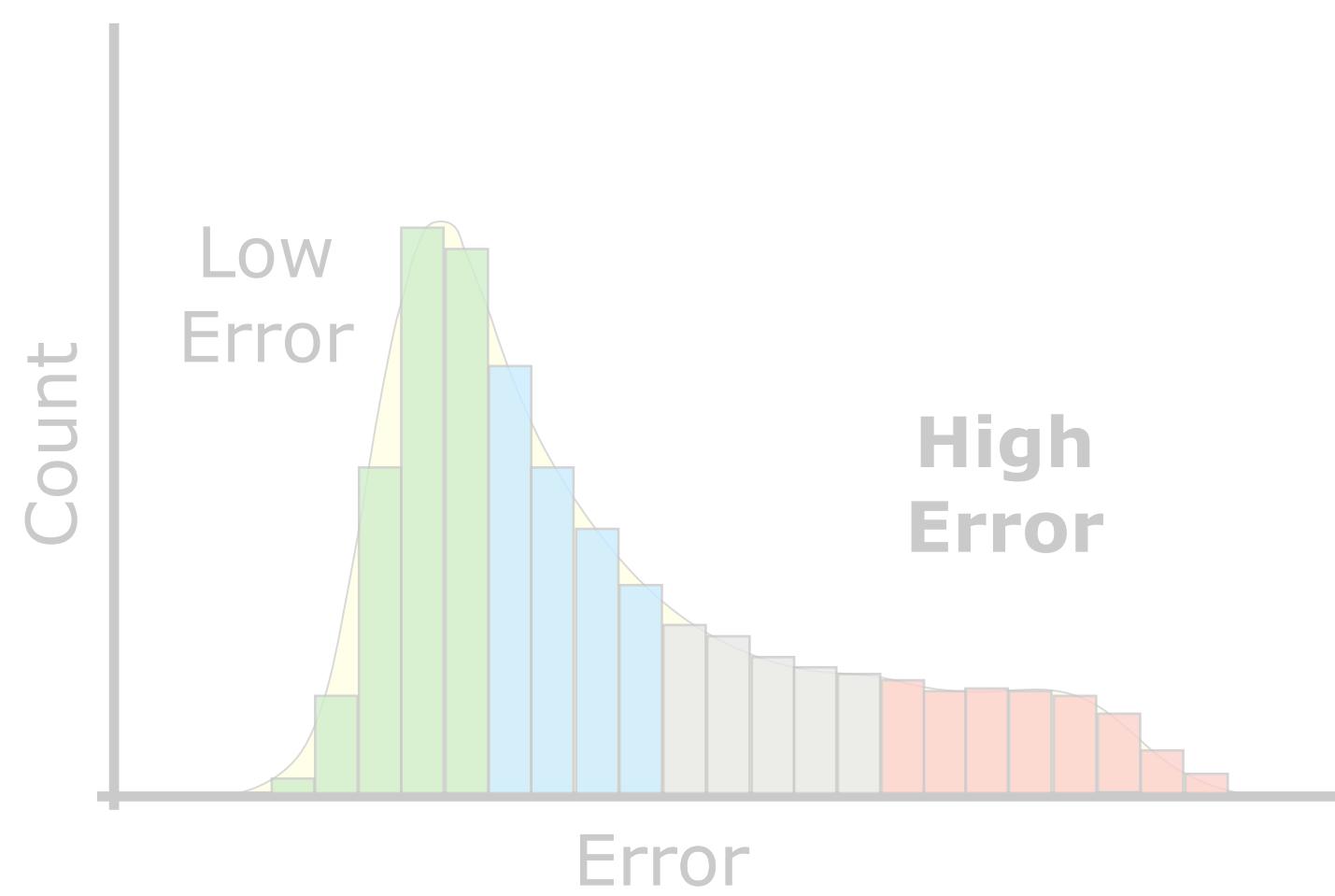


$$\min_w \mathbb{S}_\theta \left((F_1(w), \dots, F_n(w)) \right)$$

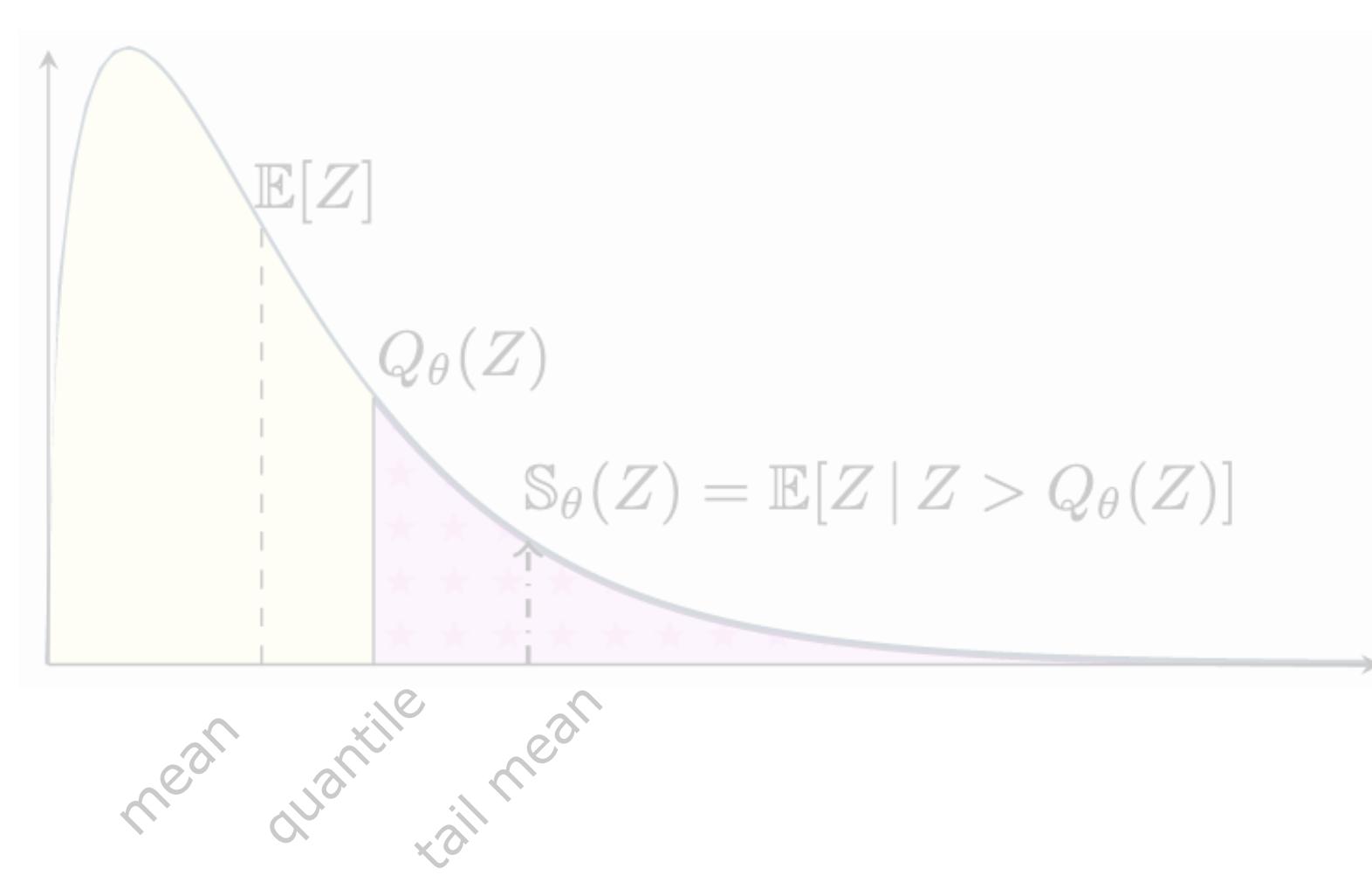


Summary: Tackling distribution shifts in federated learning

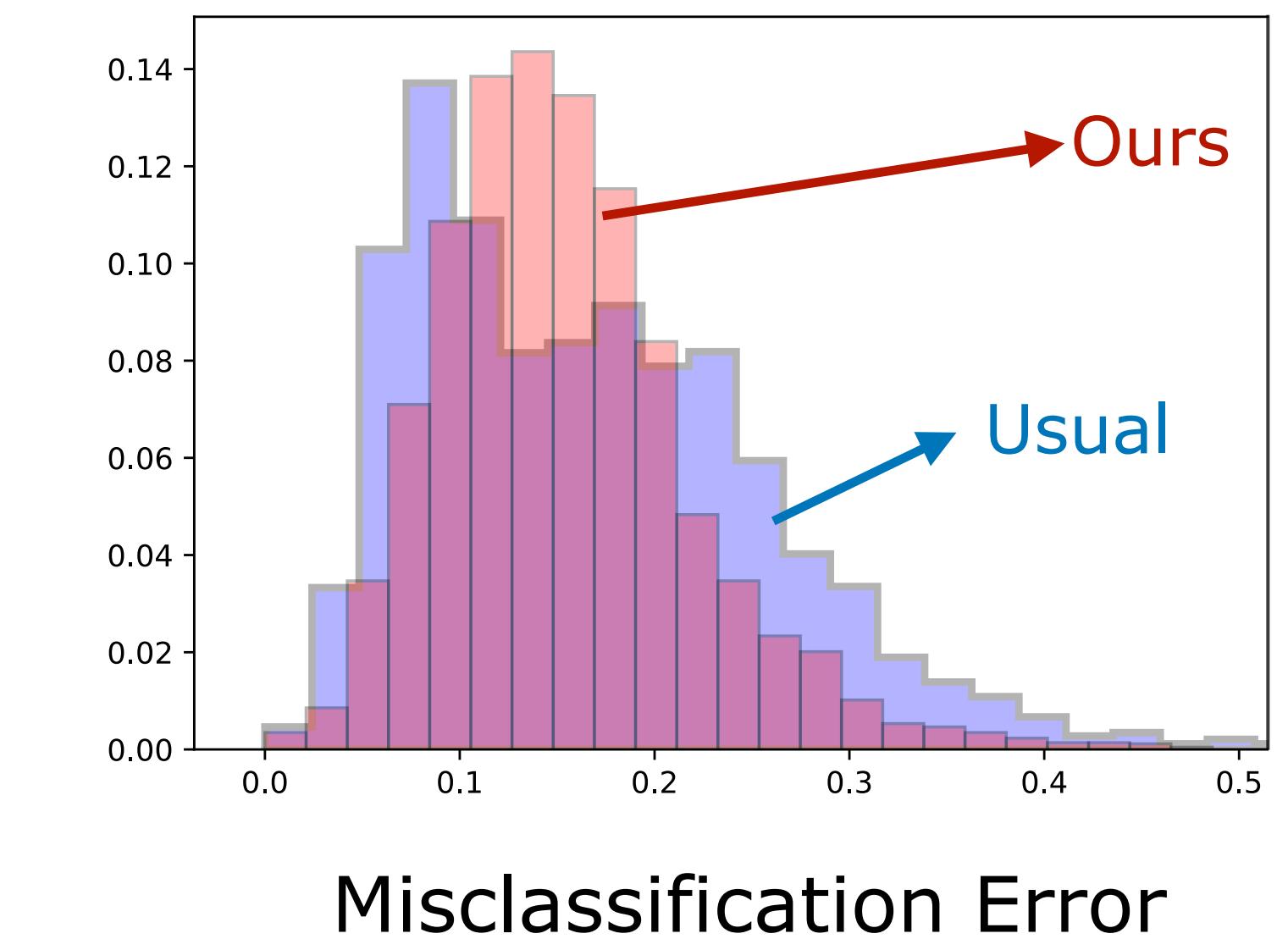
Distribution shift \Rightarrow
large tail errors



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

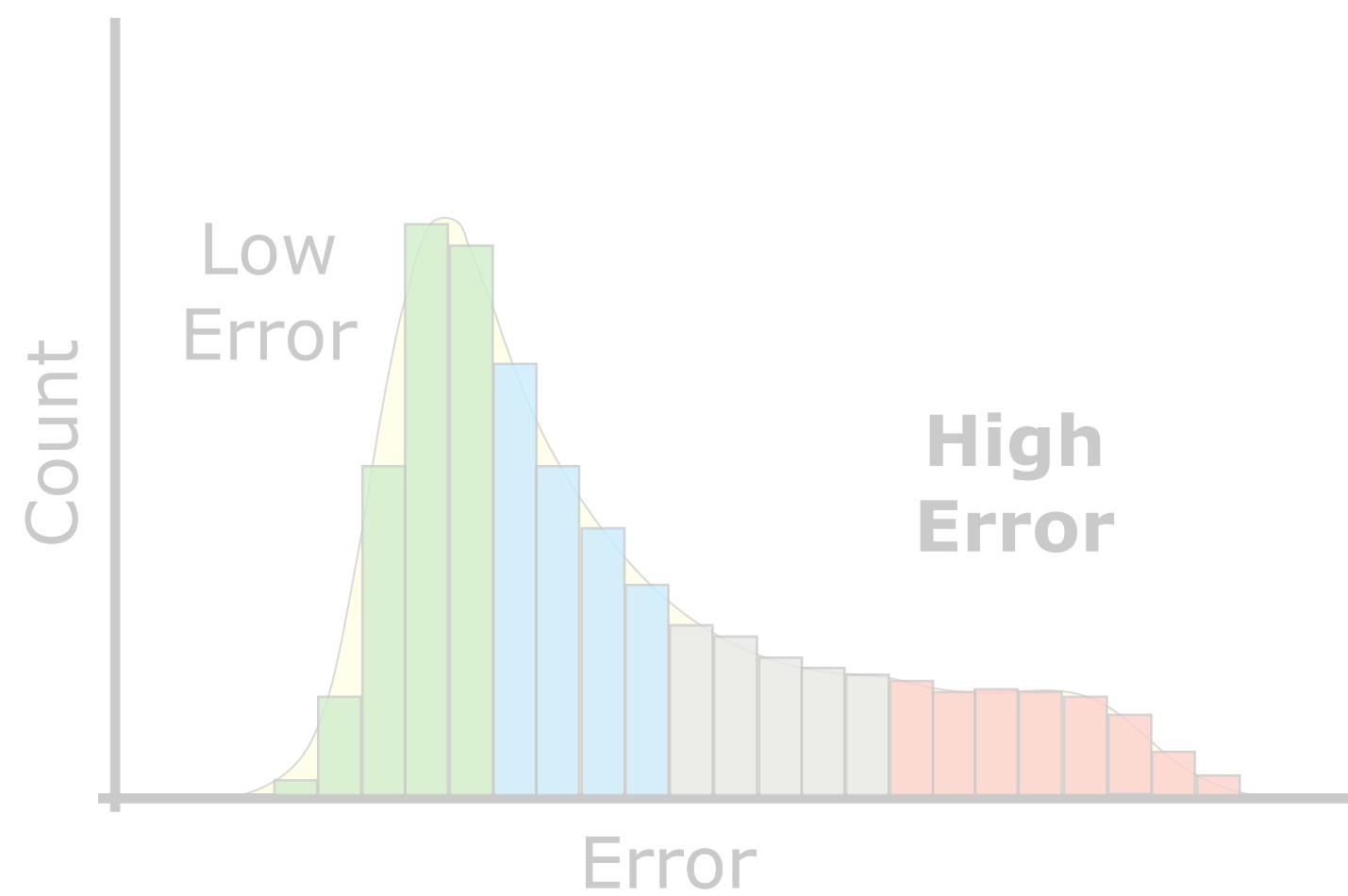


Our approach reduces
tail error

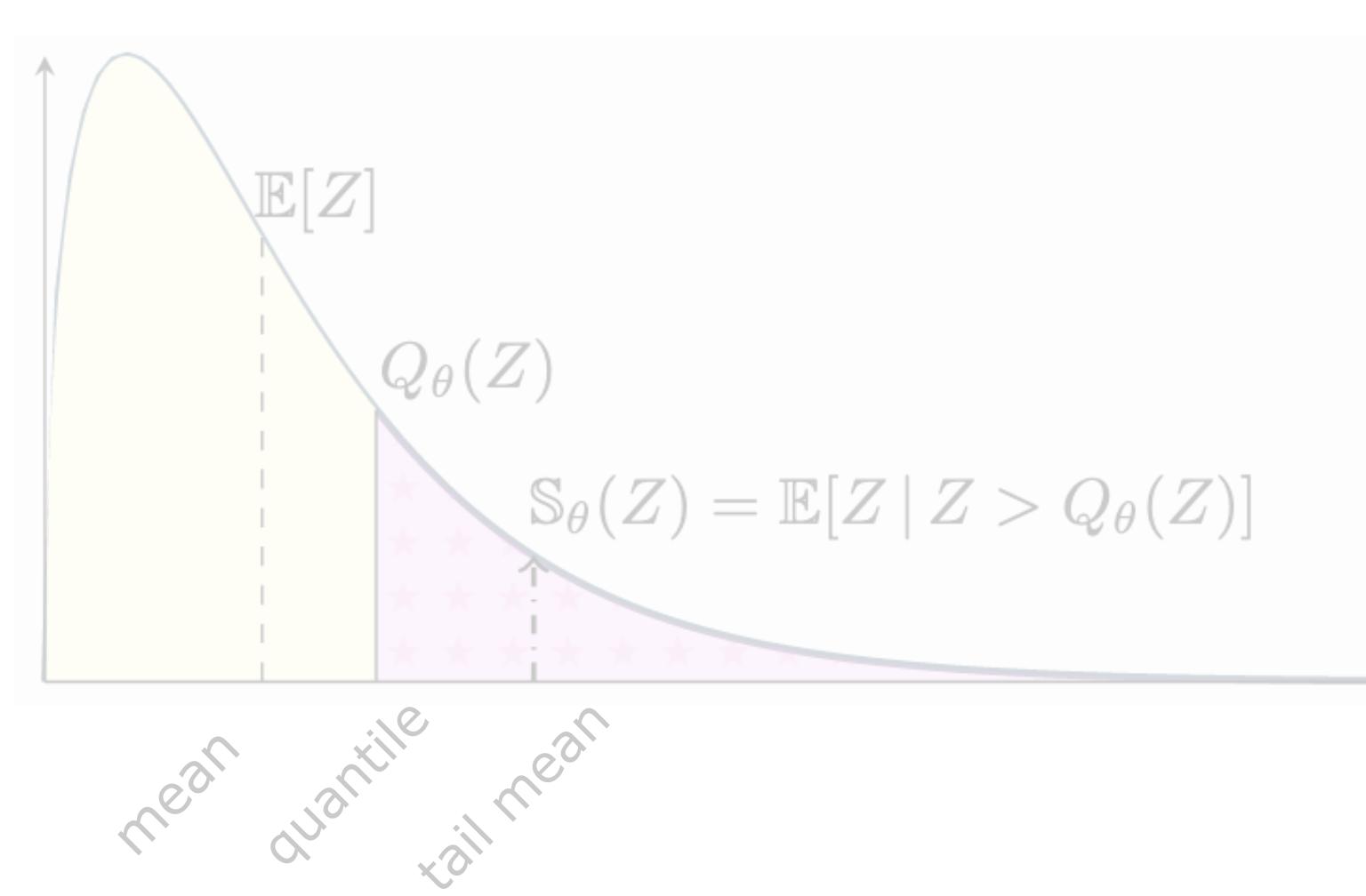


Summary: Tackling distribution shifts in federated learning

Distribution shift \Rightarrow
large tail errors



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



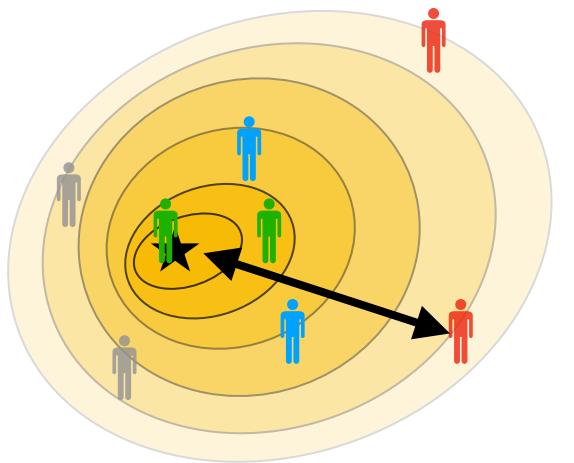
Convergence +
Privacy analysis

$O(1/\sqrt{t})$ error rate after t comm. rounds in the non-smooth, non-convex case

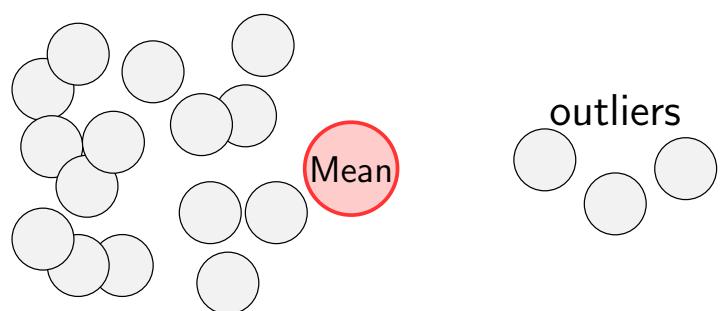
Differentially private
algorithm for distributional
robust federated learning

Future research plans

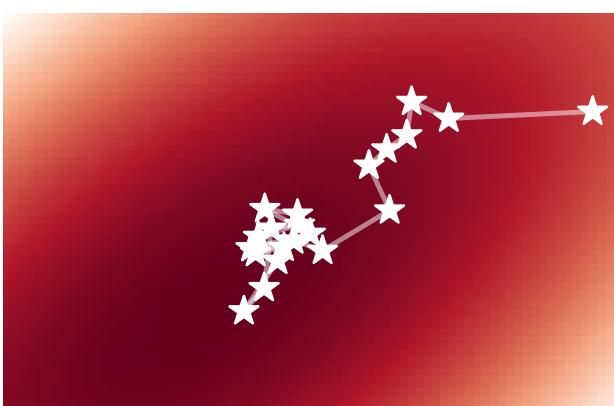
Challenges



Robustness to deployment conditions that differ from training



Robustness to outliers: adversarial or uncurated web data



Faster optimization: reduce communication and computation

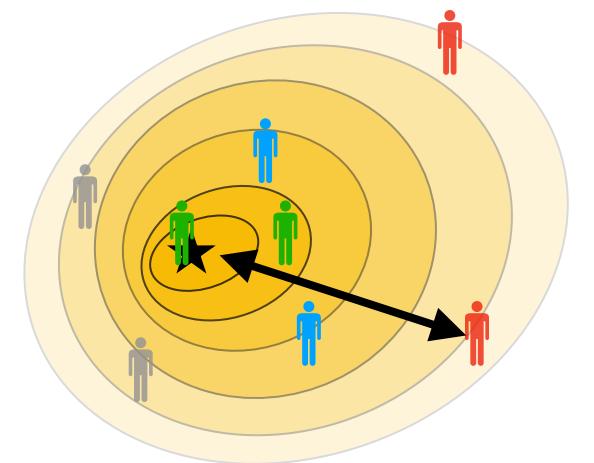


Privacy of user data

Federated learning

LLMs

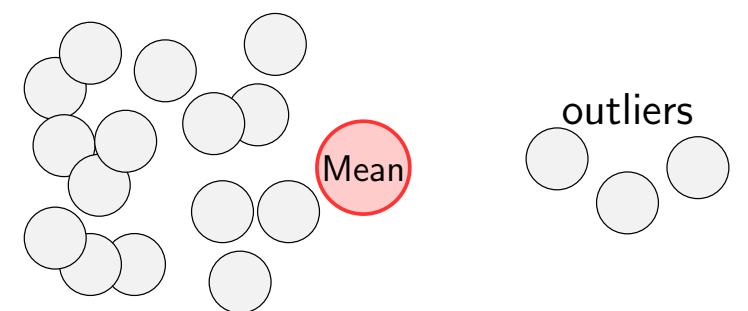
Robust Deployment



IEEE CISS 2021,
Springer SVVA 2021,
Mach. Learn. 2022

NeurIPS 2021a
NeurIPS 2021b
Submitted 2023

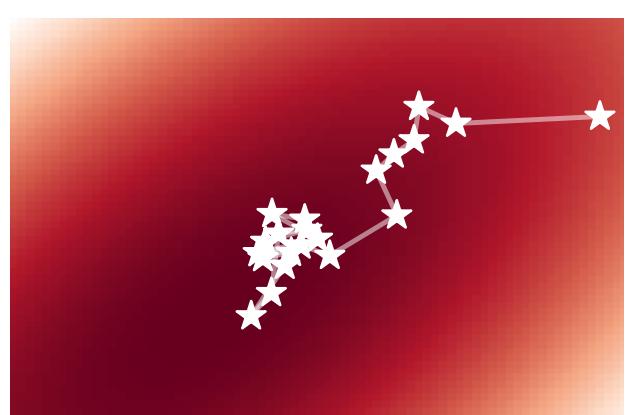
Robust to Outliers



IEEE Trans. Signal Proc. 2022,
ICML 2022

Submitted 2022

Optimize Faster



NeurIPS 2018
Submitted 2022

Privacy

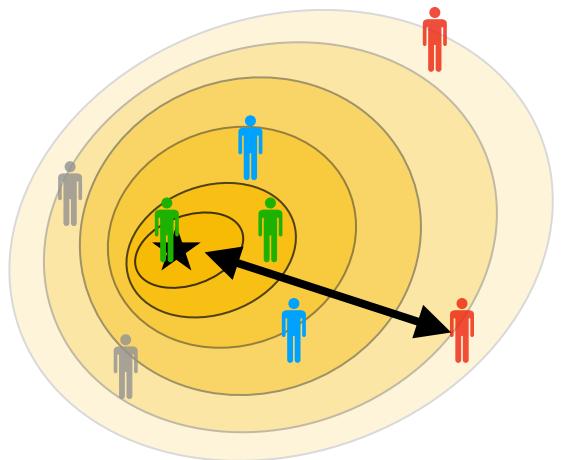


Future plan 1

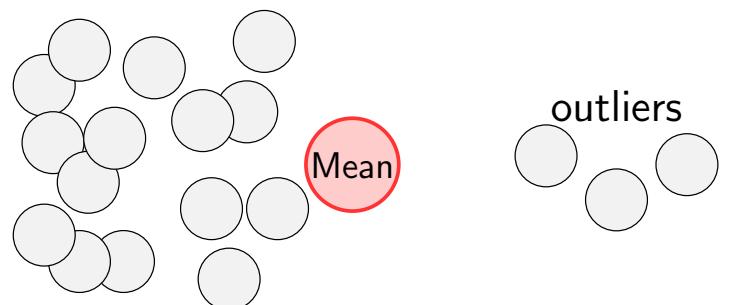
What comes after federated learning?

LLMs

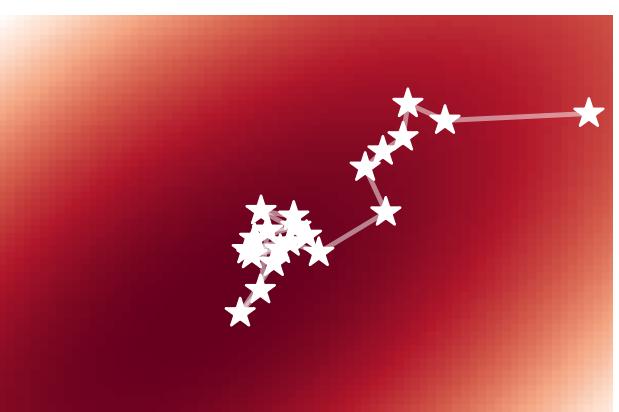
Robust Deployment



Robust to Outliers



Optimize Faster



Privacy



Future plan 2

NeurIPS 2021a
NeurIPS 2021b
Submitted 2023

Submitted 2022

NeurIPS 2018
Submitted 2022

Thank you!



J.P.Morgan

