

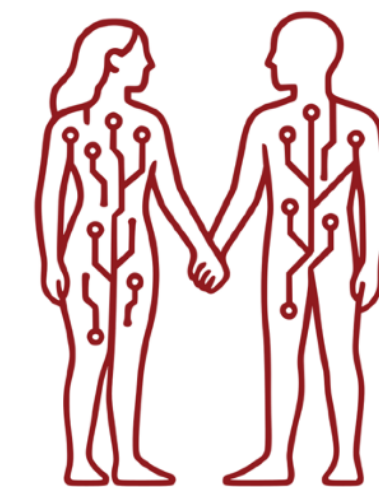
Robust Machine Learning

Towards Scalable Learning from Untrusted Data

Nirupam Gupta

Assistant Professor
Department of Computer Science

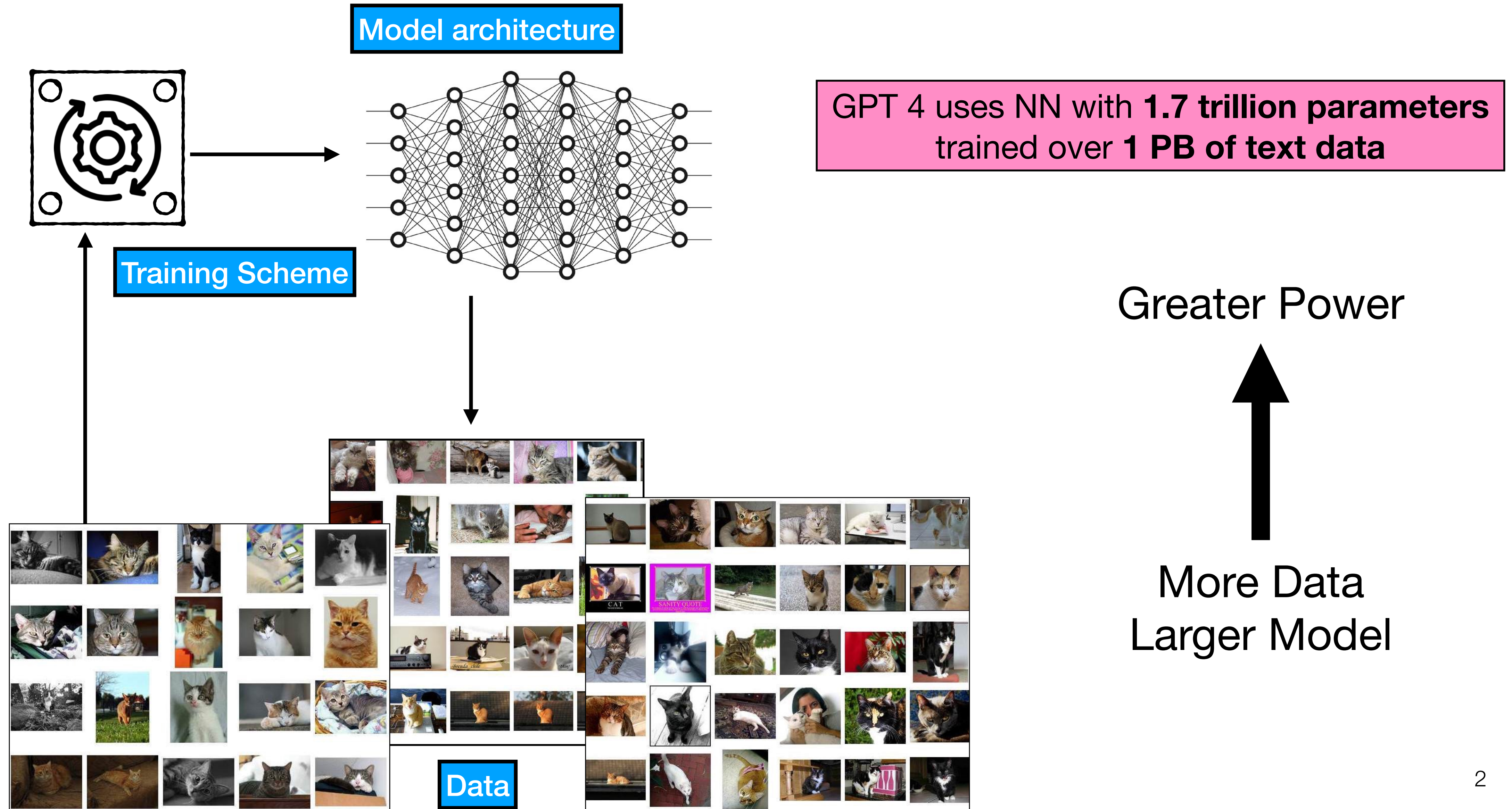
UNIVERSITY OF
COPENHAGEN



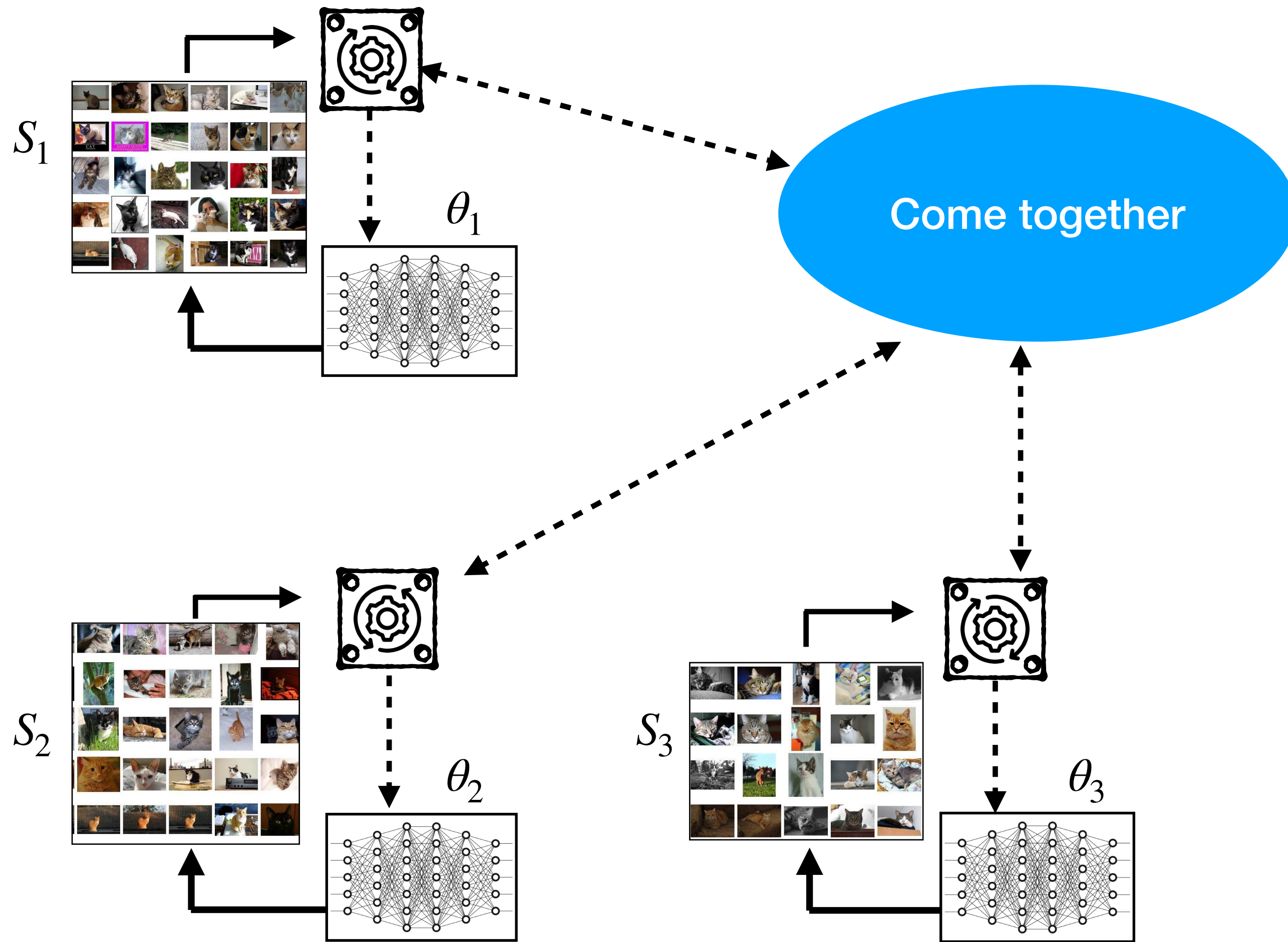
MACHINE LEARNING
UNIVERSITY OF COPENHAGEN

Learning Theory Summer School (LTSS), Copenhagen 2026

Ever-growing Scale of Machine Learning



Distributed Learning: Model Training on Decentralized Data



$$\text{Loss}^{(i)}(\theta) := \frac{1}{m} \sum_{z \in \mathcal{S}_i} \text{loss}(\theta, z)$$

Local loss function
Local samples

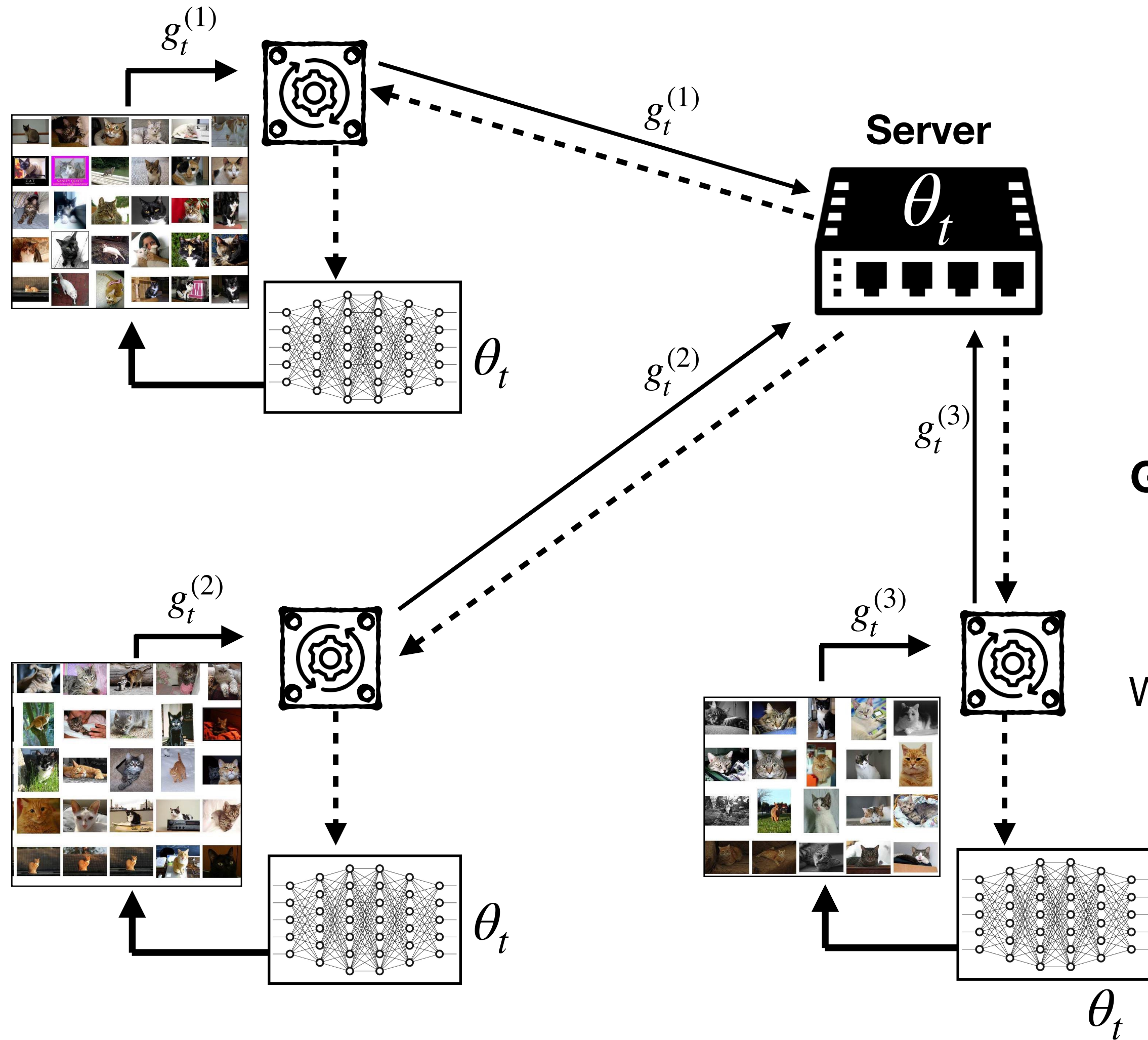
Goal: $\min_{\theta \in \Theta} \text{Loss}(\theta) \triangleq \frac{1}{n} \sum \text{Loss}^{(i)}(\theta)$

- Benefits**

 - Scalability
 - Data ownership control
 - Data privacy*
 - Democratization of ML
 - ...

Distributed Stochastic Gradient Descent (D-SGD)

Special Case of Federated Averaging



Local Phase: Each *node* i computes $g_t^{(i)} := \nabla \text{loss}(\theta_t, z_t^{(i)})$
 $z_t^{(i)} \sim U(S_i)$

Global Phase: The server updates the model:

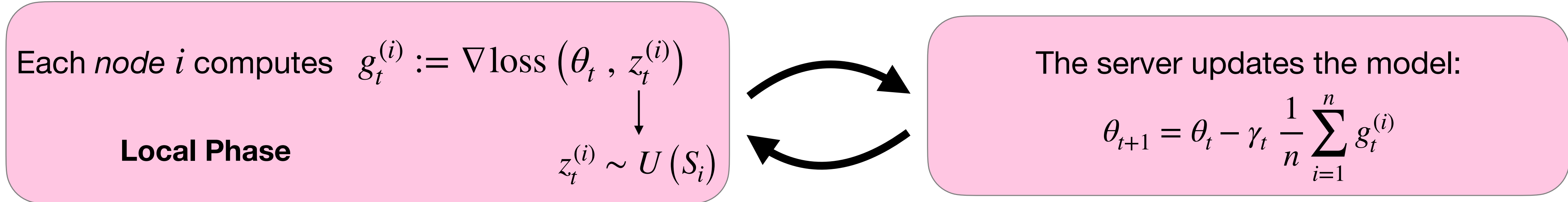
$$\theta_{t+1} = \theta_t - \gamma_t \text{Avg}(g_t^{(1)}, \dots, g_t^{(n)})$$

When nodes hold identical data, this reduces to mini-batch GD.

$$\theta_t \longrightarrow \arg \min_{\theta \in \Theta} \text{Loss}(\theta)$$

Learning Error Rate of D-SGD

Global Phase



Assumption I: Loss function is Lipschitz smooth, i.e., $\exists \lambda$ s.t. $\| \nabla \text{Loss}(\theta) - \nabla \text{Loss}(\theta') \| \leq \lambda \| \theta - \theta' \|$

$$\text{Loss}(\theta) - \text{Loss}(\theta') \leq \langle \nabla \text{Loss}(\theta'), \theta - \theta' \rangle + \frac{\lambda}{2} \| \theta - \theta' \|^2$$

Assumption II: Local gradients have bounded noise, $\exists \sigma$ s.t. $\| g_t^{(i)} - \mathbb{E}[g_t^{(i)}] \|^2 \leq \sigma^2$ where $\mathbb{E}[g_t^{(i)}] = \nabla \text{Loss}^{(i)}(\theta_t)$

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \| \nabla \text{Loss}(\theta_t) \|^2 \right] \leq \mathcal{O} \left(\sqrt{\frac{\sigma^2}{nT}} \right)$$

Learning Error Rate of D-SGD (quick proof)

$$\mathbb{E}_t \left[\text{Loss}(\theta_{t+1}) - \text{Loss}(\theta_t) \right] \leq -\gamma_t \left(1 - \frac{\lambda\gamma_t}{2} \right) \left\| \nabla \text{Loss}(\theta_t) \right\|^2 + \gamma_t^2 \frac{\lambda\sigma^2}{2n}$$

If $\gamma_t = \gamma \leq \frac{1}{\lambda}$ then
$$\sum_{t=1}^T \mathbb{E} \left[\left\| \nabla \text{Loss}(\theta_t) \right\|^2 \right] \leq \frac{2}{\gamma} \mathbb{E} \left[\text{Loss}(\theta_1) - \text{Loss}(\theta_{T+1}) \right] + \gamma \frac{\lambda\sigma^2}{n} T$$

Since $\text{Loss}(\theta_{T+1}) - \text{Loss}^*$ where $\text{Loss}^* = \min \text{Loss}(\theta)$, substituting $\gamma = 1/(\lambda\sqrt{T})$,

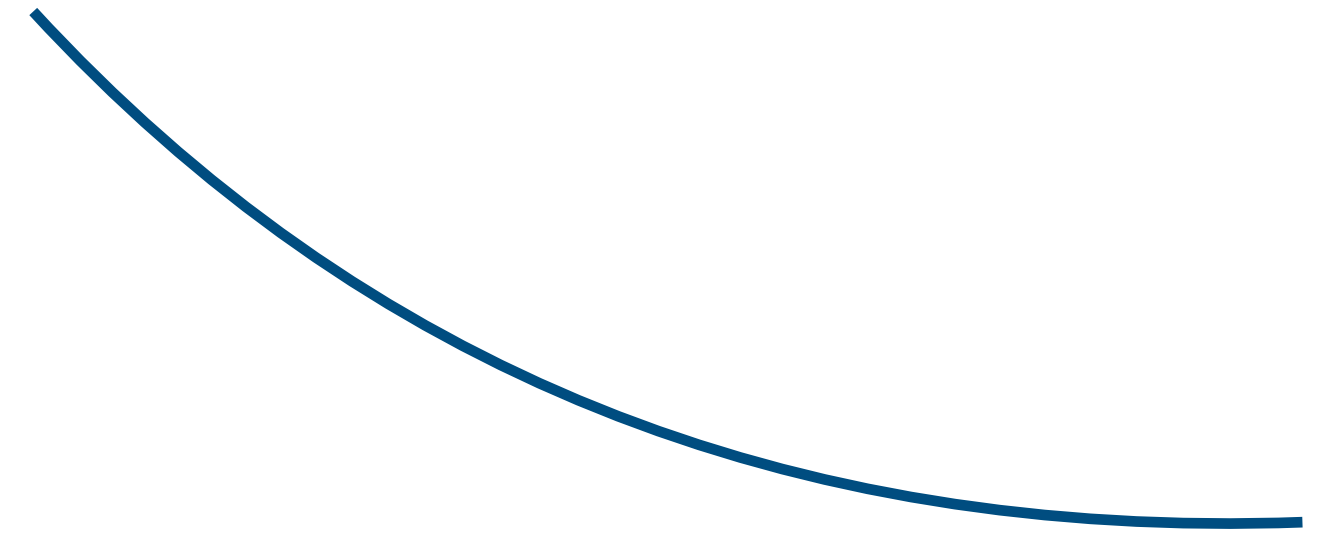
$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla \text{Loss}(\theta_t) \right\|^2 \right] \leq \left(2\lambda \left(\mathbb{E} \left[\text{Loss}(\theta_1) \right] - \text{Loss}^* \right) + \frac{\sigma^2}{n} \right) \frac{1}{\sqrt{T}}$$

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \left\| \nabla \text{Loss}(\theta_t) \right\|^2 \right] \leq \mathcal{O} \left(\frac{1}{n\sqrt{T}} \right)$$

Learning Error Rate of D-SGD (cont'd)

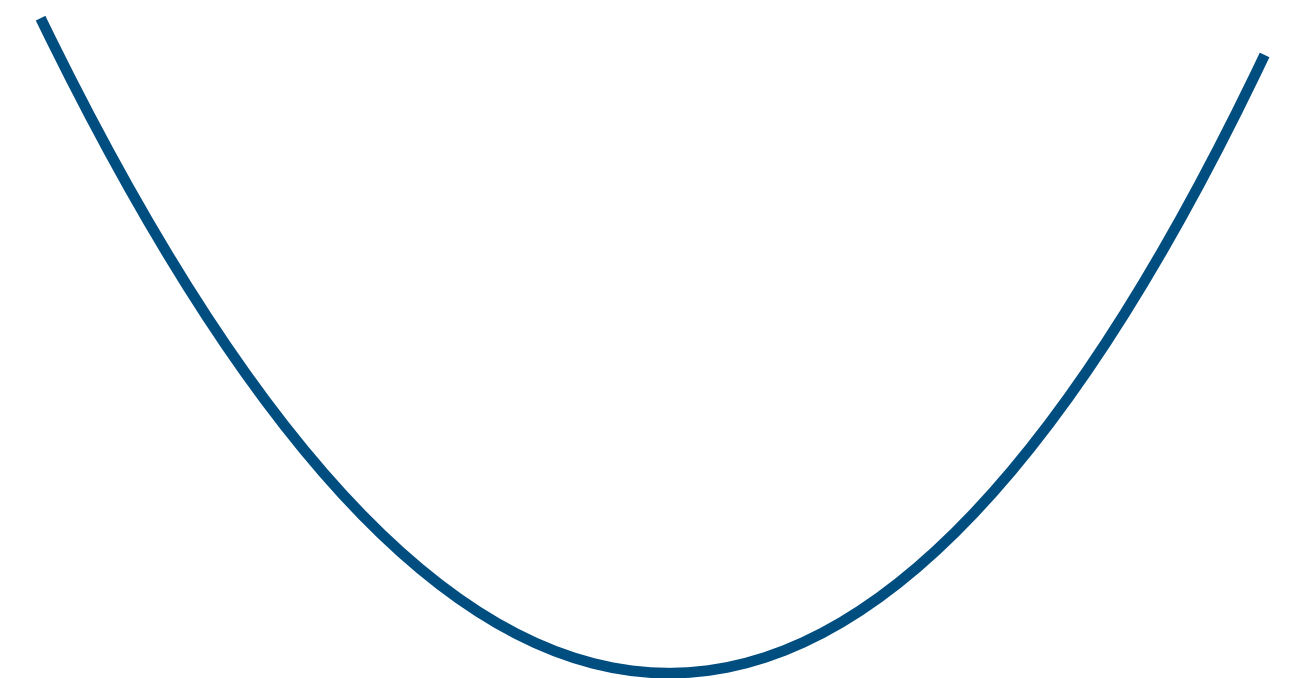
Convexity: $\text{Loss}(\theta') \geq \text{Loss}(\theta) + \langle \nabla \text{Loss}(\theta), \theta' - \theta \rangle$

$$\mathbb{E} \left[\text{Loss}(\theta_{T+1}) - \text{Loss}^* \right] \leq \mathcal{O} \left(\sqrt{\frac{\sigma^2}{nT}} \right)$$

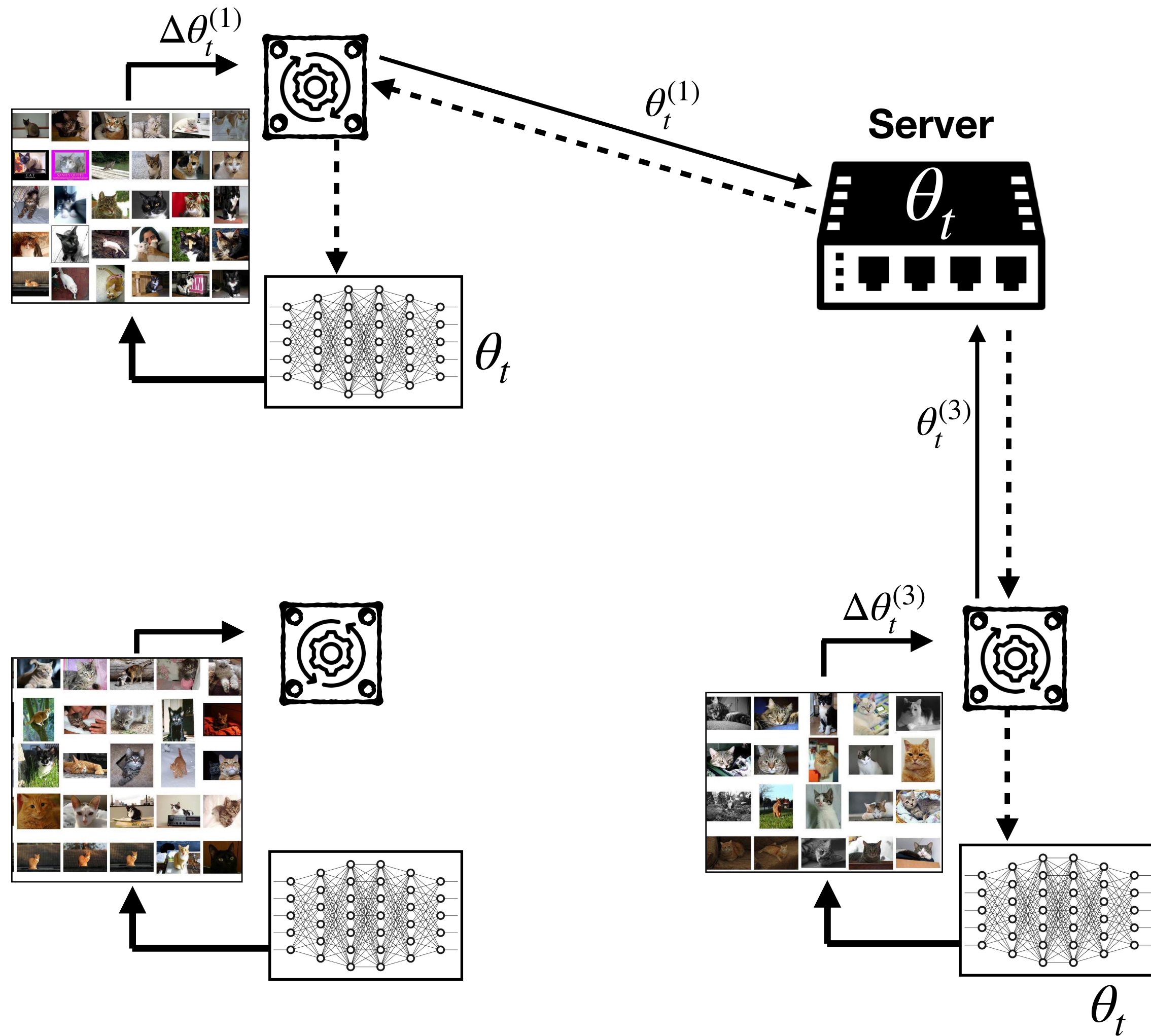


Strong convexity: $\text{Loss}(\theta') \geq \text{Loss}(\theta) + \langle \nabla \text{Loss}(\theta), \theta' - \theta \rangle + \mu \|\theta' - \theta\|^2$

$$\mathbb{E} \left[\text{Loss}(\theta_{T+1}) - \text{Loss}^* \right] \leq \mathcal{O} \left(\frac{\sigma^2}{nT} \right)$$



Federated Averaging: Local SGD



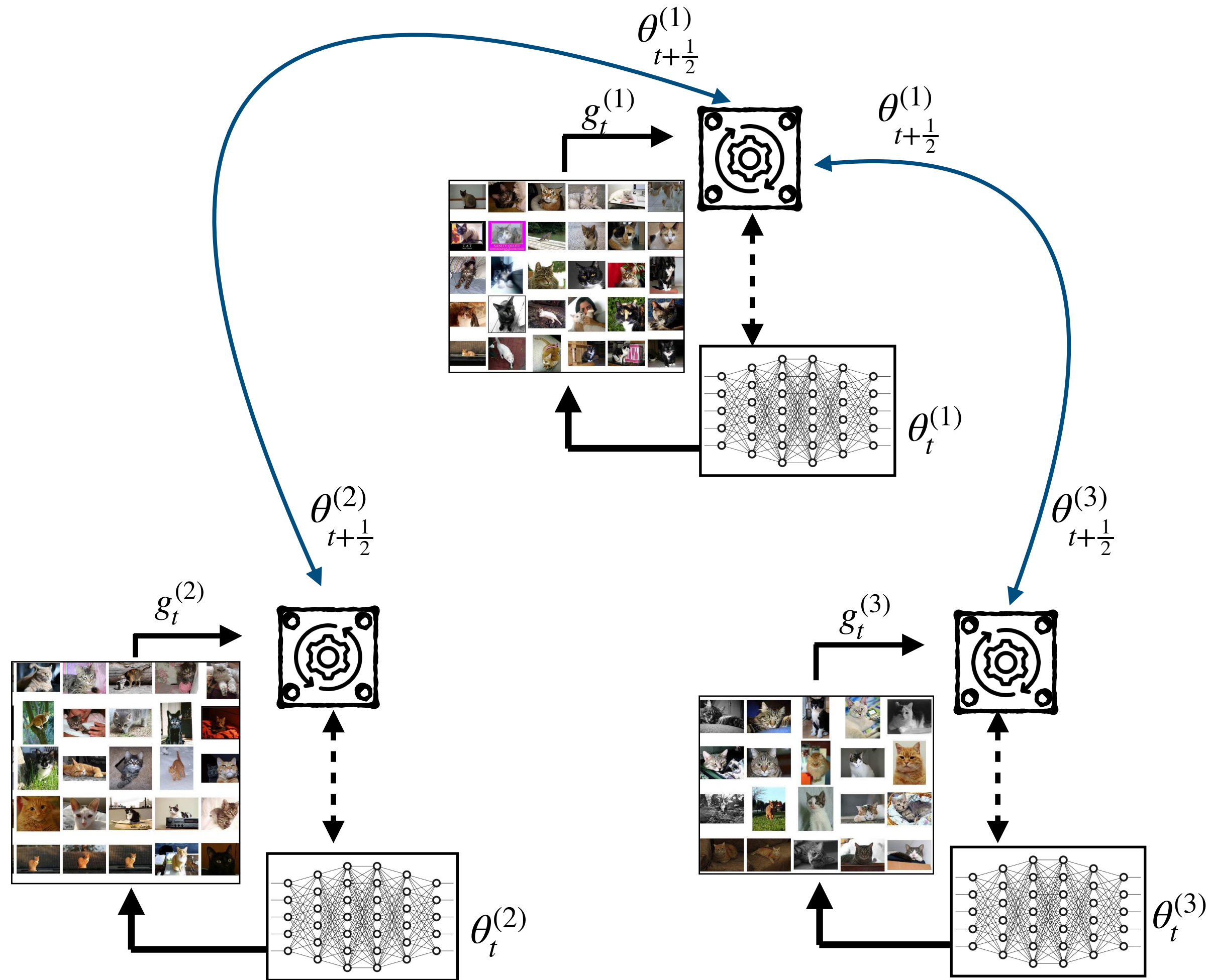
Partial participation: Server may only ask a few available clients to compute the updates.

Local Phase: Each selected client i computes an updated model $\theta_t^{(i)}$, upon running several local SGD steps.

Global Phase: The server updates the model:

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{|S_t|} \sum_{i \in S_t} (\theta_t - \theta_t^{(i)})$$

Distributed SGD in Peer-to-Peer Architecture (Serverless)



$$\min_{\theta^{(i)}, \forall i} \frac{1}{n} \sum \text{Loss}^{(i)}(\theta^{(i)})$$

$$\text{Subject to : } \theta^{(i)} = \theta^{(j)}, \forall i, j$$

Phase I: Each *node* i updates the local model:

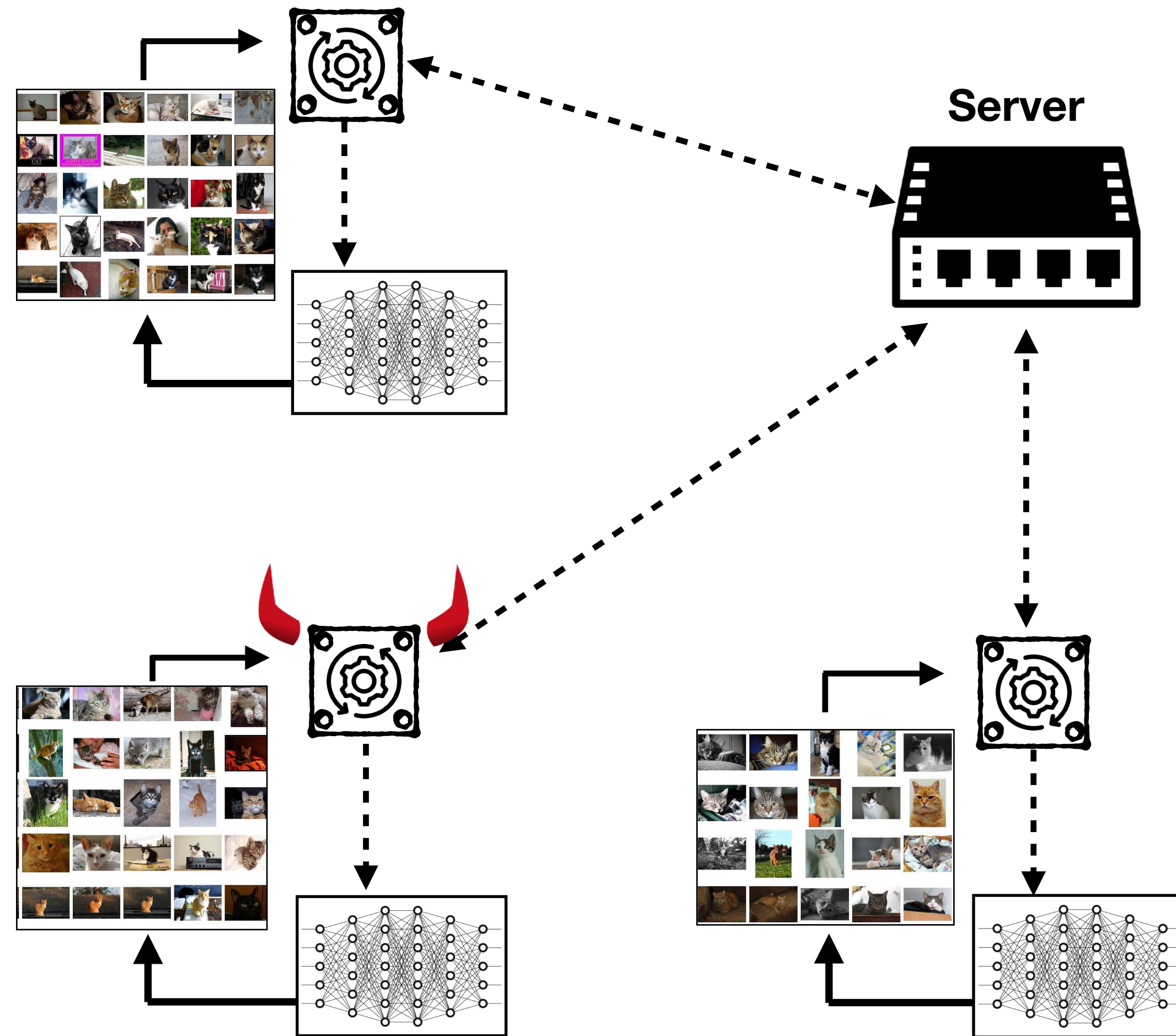
$$\theta_{t+\frac{1}{2}}^{(i)} = \theta_t^{(i)} - \gamma_t g_t^{(i)}$$

Phase II: Nodes gossip:

$$\theta_{t+1}^{(i)} = \theta_{t+\frac{1}{2}}^{(i)} - \alpha_t \sum_{j \in \mathcal{N}_t^{(i)}} \left(\theta_{t+\frac{1}{2}}^{(i)} - \theta_{t+\frac{1}{2}}^{(j)} \right)$$

Gossip rate

Vulnerability to Data & Model Poisonings

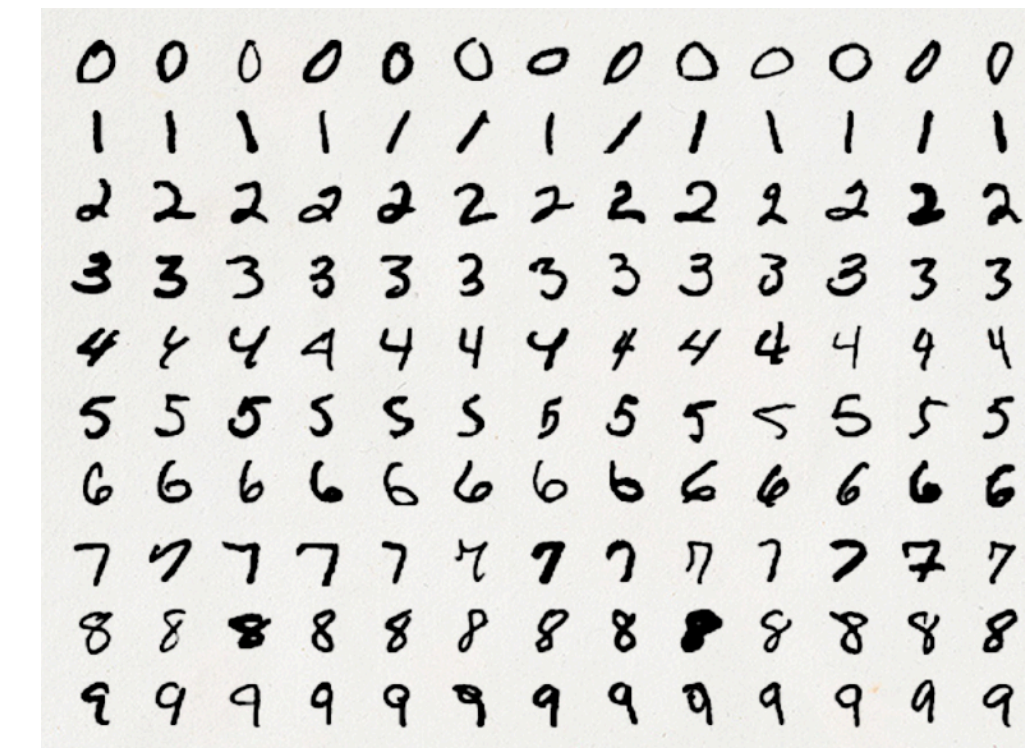
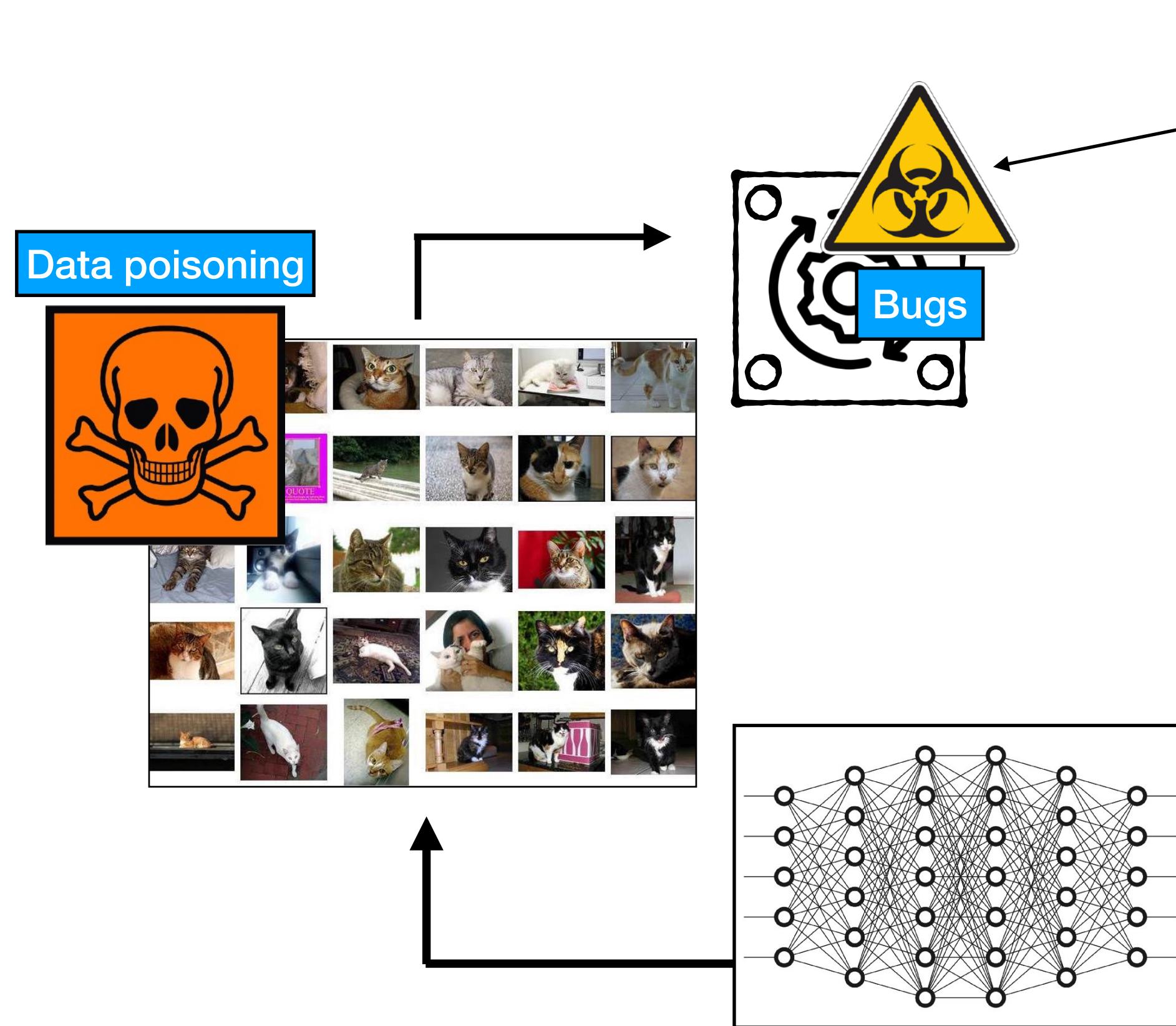


Not all data/machines can be certified (trusted)

$$\theta_{t+1} = \theta_t - \gamma_t \text{Avg} (g_t^{(1)}, \dots, g_t^{(n)})$$

Breaks down when some machines “misbehave”

Common sources of poisoning

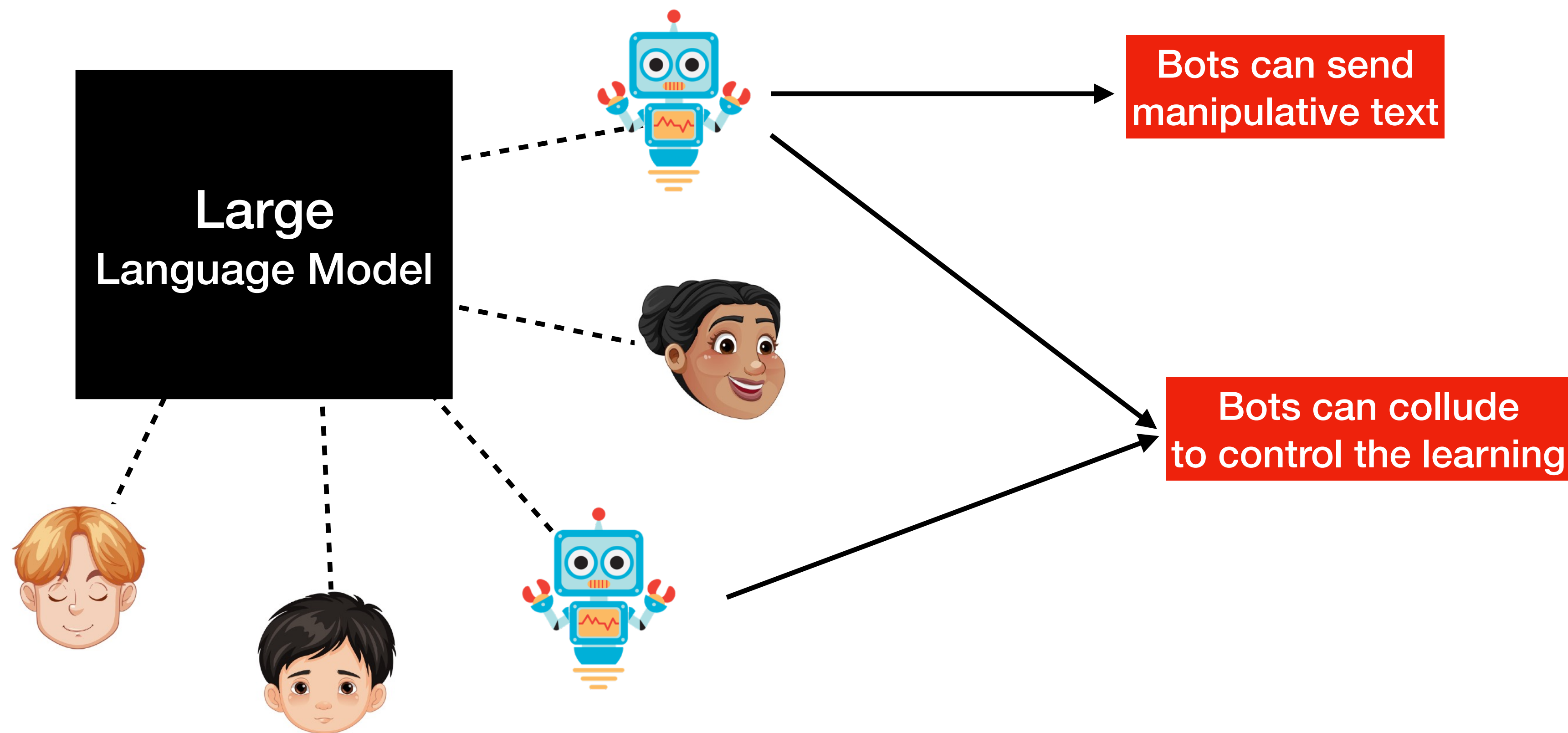


Some adversarial behaviors:

- *Label-flipping (data poisoning)*: swap the labels, e.g., $(4, 9) \rightarrow (9, 4)$
- *Sign-flipping (model poisoning)*: toggle gradient's sign $g_t^{(i)} \rightarrow -g_t^{(i)}$

Data Poisoning is Common in the Age of Bots

Some fake news is created not for humans, but for misleading language models



Bots are Spreading Faster & Getting Smarter

Communications of the ACM, 2016

B18717

Today's social bots are sophisticated and sometimes menacing. Indeed, their presence can endanger online ecosystems as well as our society.

BY EMILIO FERRARA, ONUR VAROL, CLAYTON DAVIS, FILIPPO MENCZER, AND ALESSANDRO FLAMMINI

The Rise of Social Bots

“exhibit human-like behavior”

CNBC, 2017

(~ 300 million total)

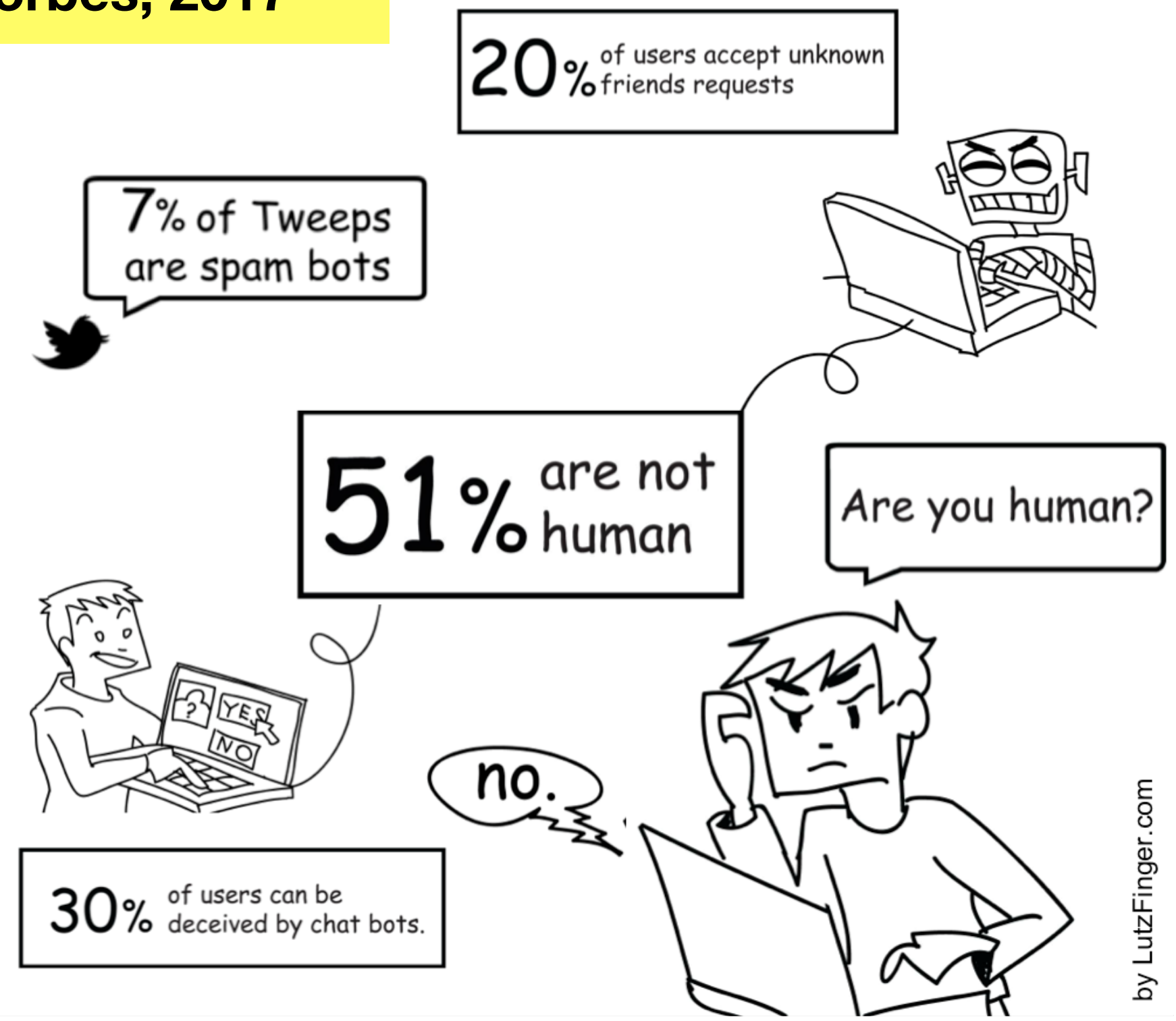
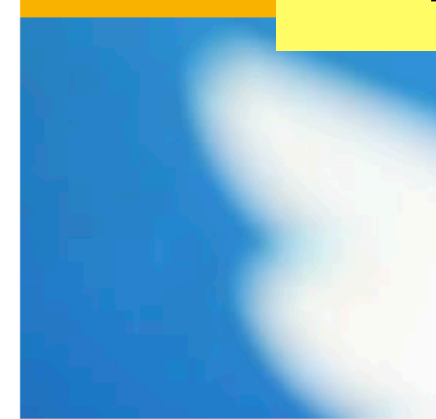
As many as 48 million Twitter accounts aren't people, says study

PUBLISHED FRI, MAR 10 2017·1:09 PM EST | UPDATED FRI, MAR 10 2017·7:56 PM EST

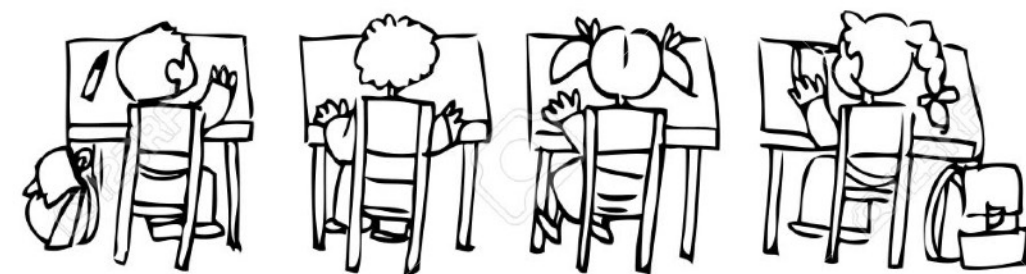
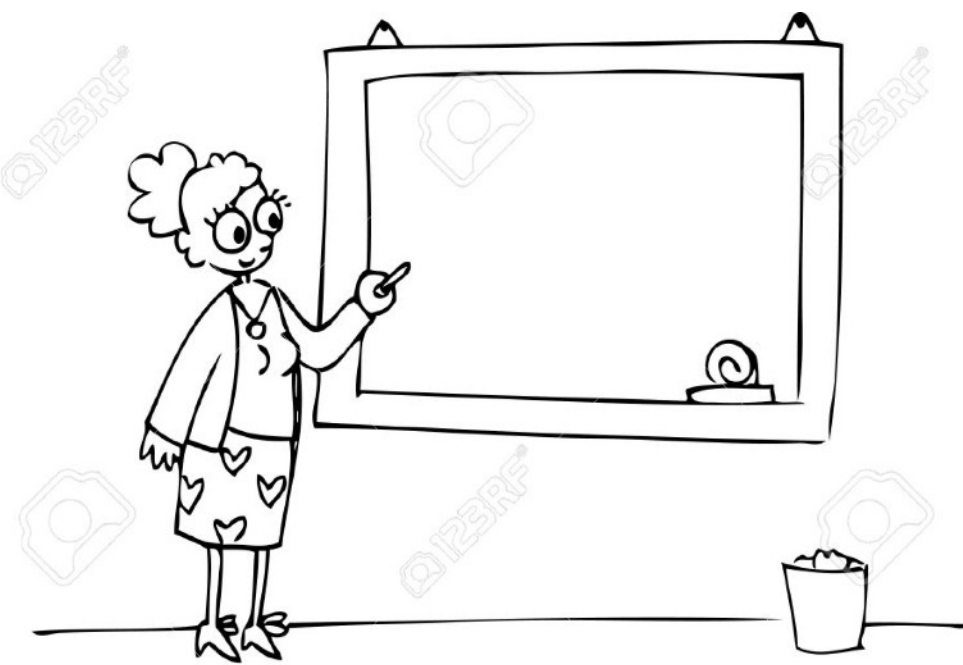
Michael Newberg
@MIKENEWBERG

SHARE [f](#) [t](#) [in](#) [✉](#)

Forbes, 2017

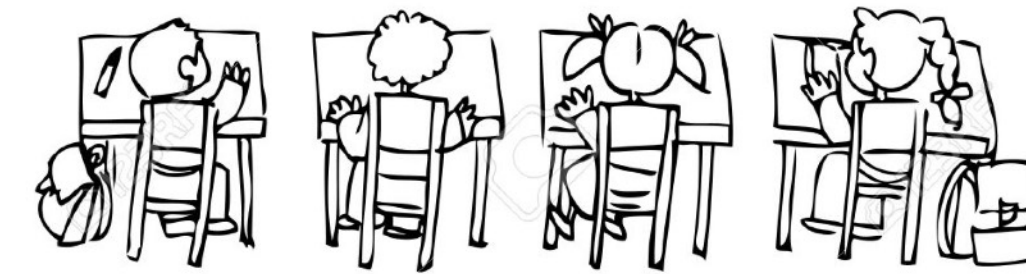


Can we Learn Without Trusting the Source?



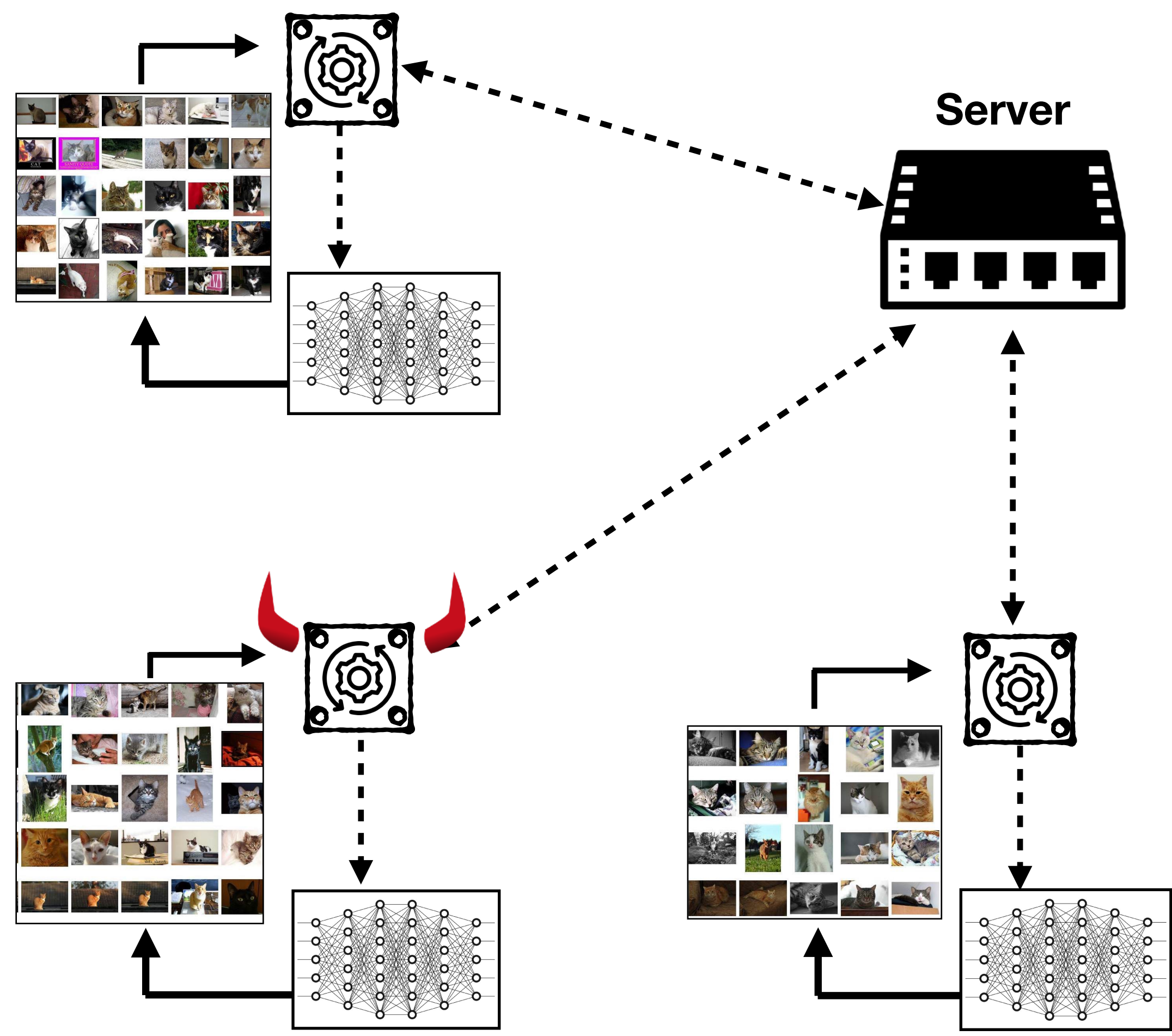
Copyright: Carla Francesca Castagno

Vs



Copyright: Carla Francesca Castagno

(Byzantine-)Robust Machine Learning



f nodes (silos) of *unknown identity* are **adversarial**

$< n/2$

Need not follow the algorithm

Robust DL Goal:

$$\min_{\theta \in \Theta} \frac{1}{n - f} \sum_{i \in H} \text{Loss}^{(i)}(\theta)$$

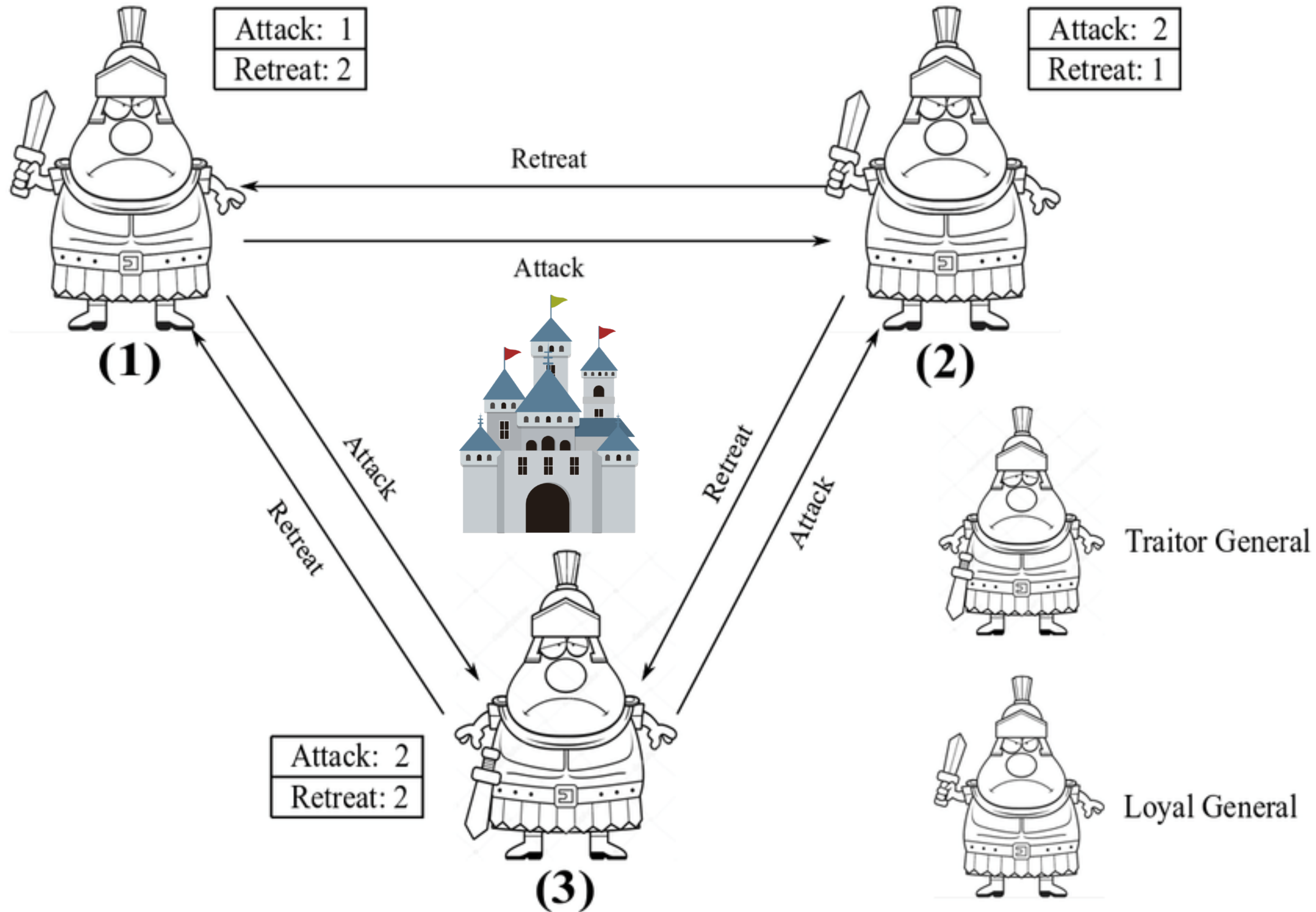
$$\text{Loss}^{(i)}(\theta) := \frac{1}{m} \sum_{z \in S_i} \text{loss}(\theta, z)$$

Set of *honest nodes*

Unknown! ?

Byzantine Generals' Problem

- Leslie Lamport (*Turing Award - 2013*)



Three Challenges in Robust FL

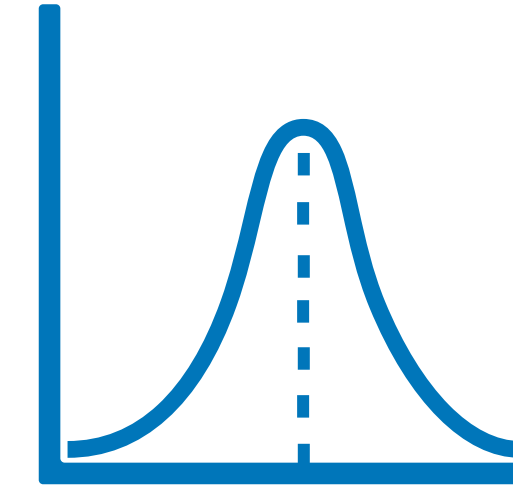
State machine replication
is not efficient

Adversarial nodes can
camouflage as honest ones

Local randomness gets
amplified

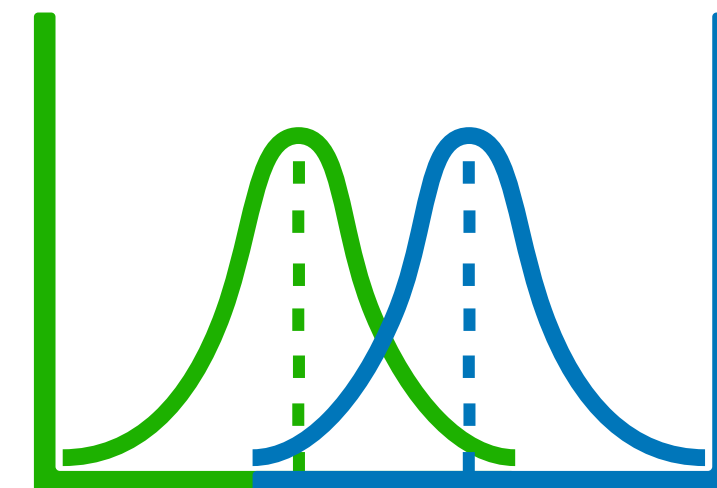
Local Randomness

Honest nodes/clients *independently*
compute noisy updates



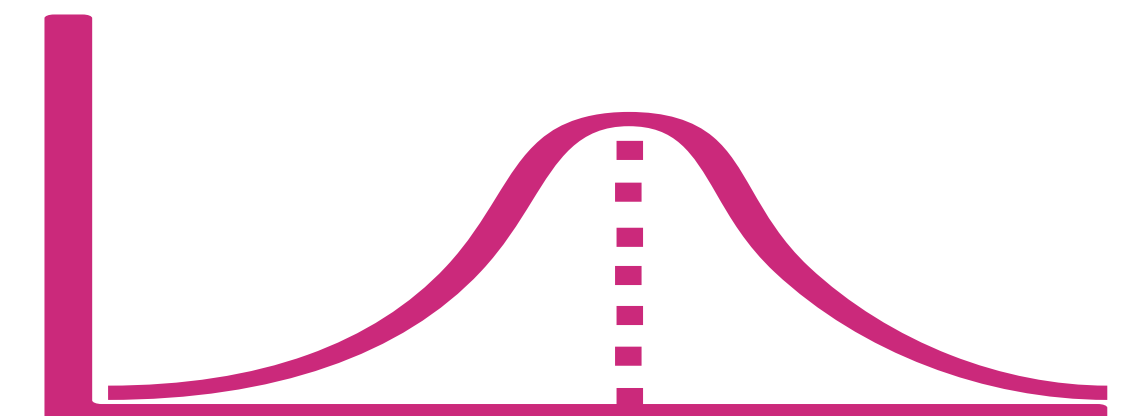
Data Heterogeneity

Different nodes sample data
from different distributions



Data Privacy

Nodes do not share exact
information of their data

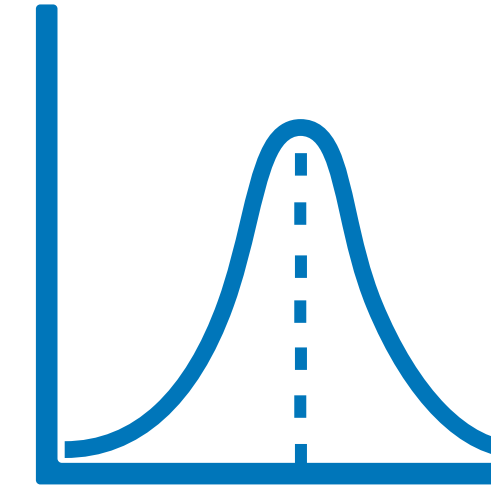


Challenge I: Tackling Local Randomness

State machine replication
is not efficient

Local Randomness

Honest nodes independently
compute noisy updates



Mini-batch GD

$$S_i = S \quad \forall i \in H$$

$$\text{Loss}^i(\theta) = \text{Loss}(\theta) = \frac{1}{m} \sum_{z \in S} \text{loss}(\theta, z)$$

Can be exploited by adversarial nodes

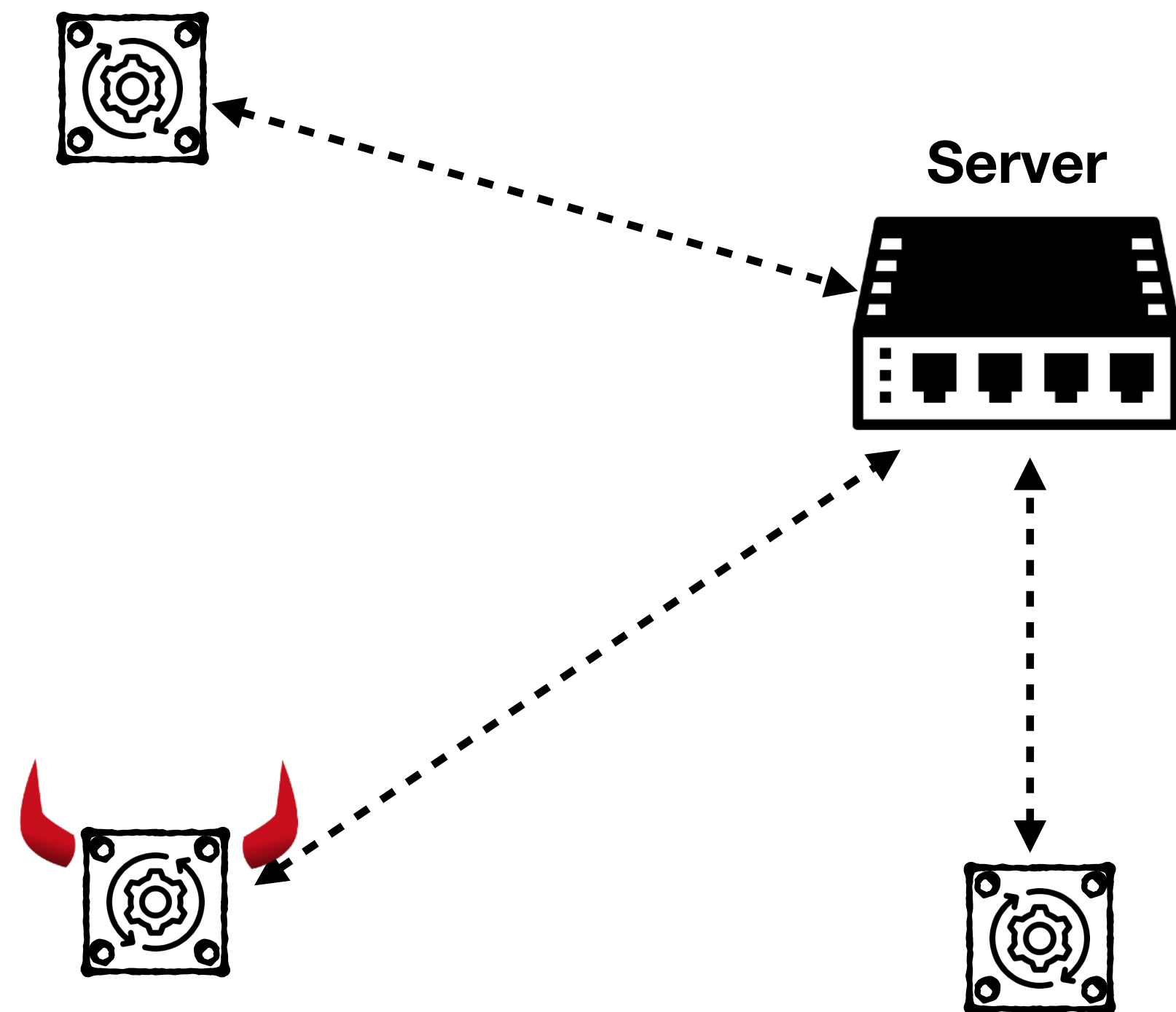
$$g_t^{(i)} := \nabla \text{loss}(\theta_t, z_t^{(i)})$$

where $z_t^{(i)} \sim \mathcal{U}(S)$

Noisy estimate

$$\nabla \text{Loss}(\theta_t)$$

Robust Mini-batch Gradient Descent



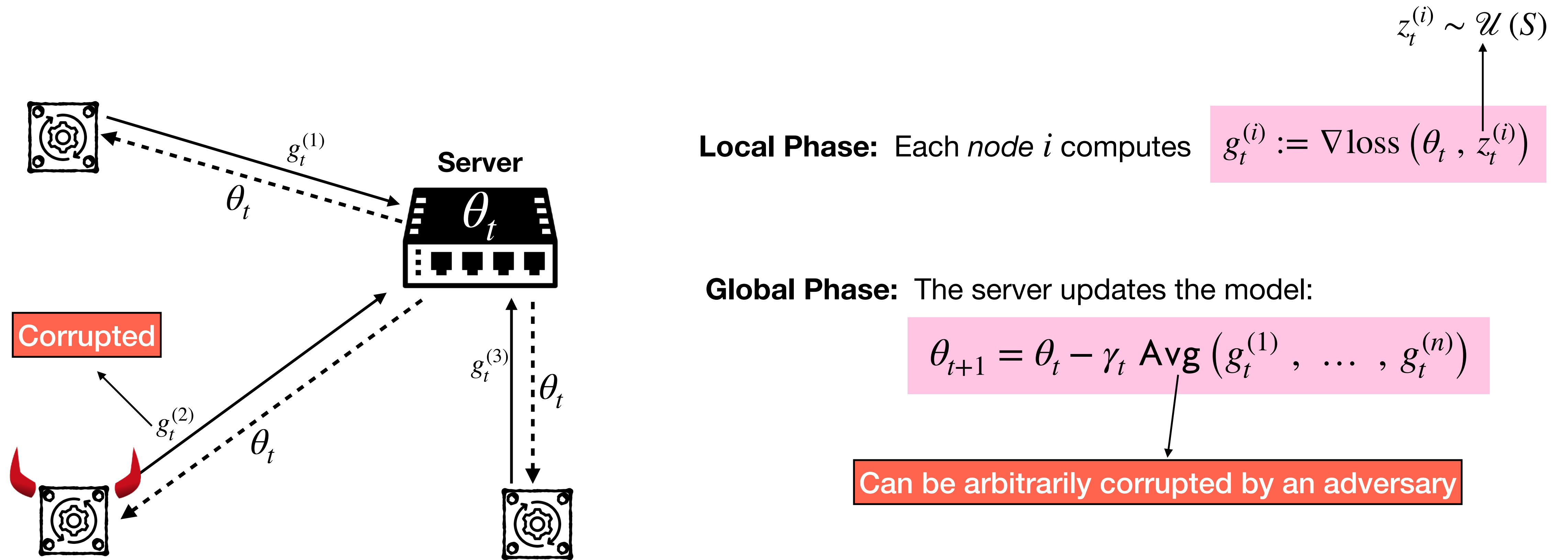
$f < \frac{n}{2}$ nodes are adversarial

$$S_i = S \quad \forall i \in H \quad \text{Loss}^i(\theta) = \text{Loss}(\theta) = \frac{1}{m} \sum_{z \in S} \text{loss}(\theta, z)$$

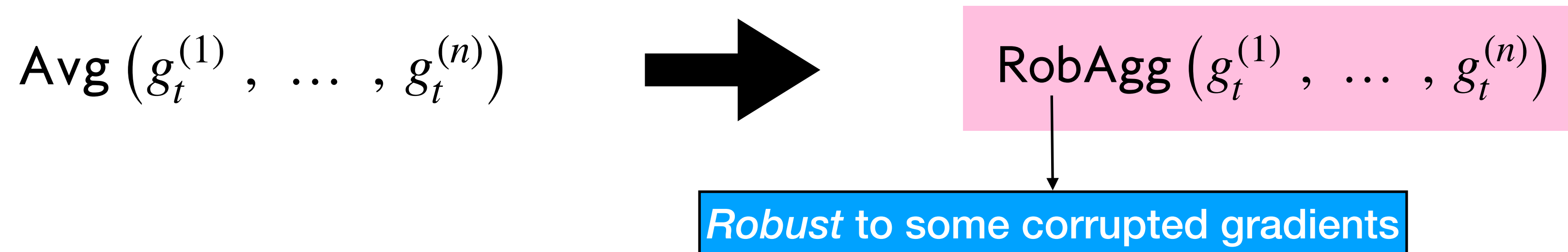
Minimize $\text{Loss}(\theta) := \frac{1}{m} \sum_{z \in S} \text{loss}(\theta, z)$

Problem reduces to estimating $\nabla \text{Loss}(\theta)$

Shortcoming of Classic Mini-batch Gradient Descent



From Classic to Robust Mini-Batch GD



Examples: $d = 2$, $n = 5$ and $f = 1$. Let $(g_t^{(1)}, \dots, g_t^{(5)}) = \begin{pmatrix} 2 & 3 & 2 & 1 & 12 \\ 5 & 6 & 4 & 5 & 0 \end{pmatrix}$

$\begin{matrix} & & & & & \begin{matrix} (2) \\ (5) \end{matrix} \\ & & & & \uparrow & \\ & & & & \hline & & & & & \begin{matrix} (12) \\ (0) \end{matrix} \end{matrix}$

Coordinate-wise Median (**CWMed**)

$$\text{CWMed} \begin{pmatrix} 2 & 3 & \textcircled{2} & 1 & 12 \\ 5 & 6 & 4 & \textcircled{5} & 0 \end{pmatrix} = \begin{pmatrix} 2 \\ 5 \end{pmatrix}$$

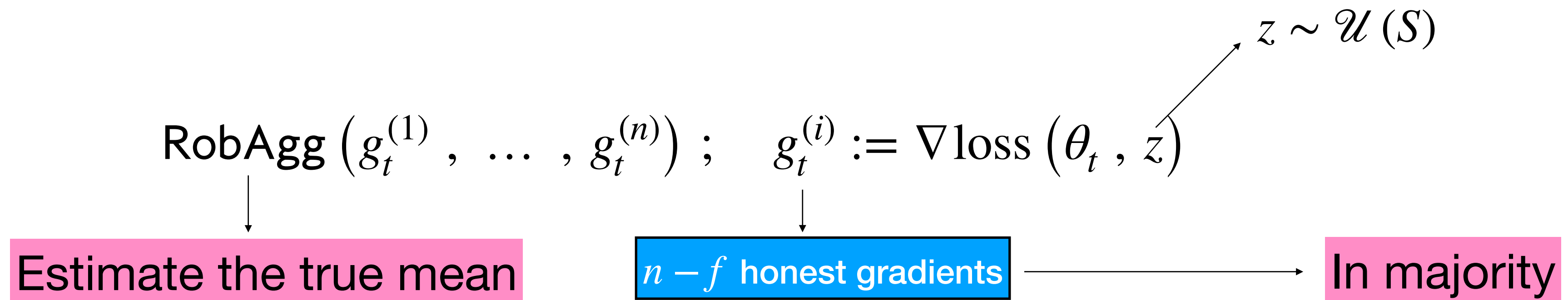
Averaging

$$\text{Avg} \begin{pmatrix} 2 & 3 & 2 & 1 & 12 \\ 5 & 6 & 4 & 5 & 0 \end{pmatrix} = \begin{pmatrix} 4 \\ 4 \end{pmatrix}$$

Coordinate-wise Trimmed Mean (**CWTM**)

$$\text{CWTM} \begin{pmatrix} 2 & 3 & 2 & \times & \times \\ 5 & \times & 4 & 5 & \times \end{pmatrix} = \begin{pmatrix} 2.3 \\ 4.6 \end{pmatrix}$$

Robust Mini-batch GD – Reducible to Robust Mean

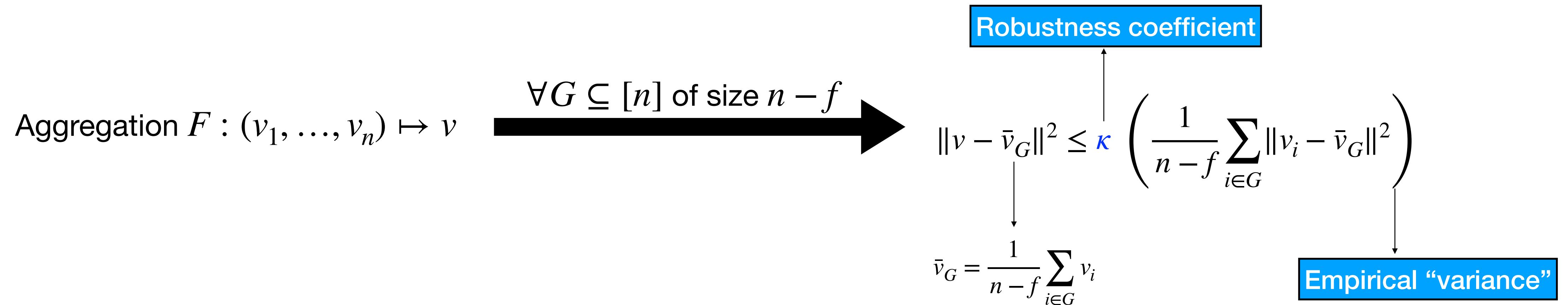


Minimize $\mathbb{E} \left[\left\| \text{RobAgg} (g_t^{(1)}, \dots, g_t^{(n)}) - \nabla \text{Loss} (\theta_t) \right\|^2 \right]$

$\nabla \text{Loss}(\theta_t) = \mathbb{E}_{z \sim \mathcal{U}(S)} \left[\nabla \text{loss} (\theta_t, z) \right]$
Mean of honest gradients

Bounded “variance”: $\mathbb{E}_{z \sim \mathcal{U}(S)} \left[\left\| \nabla \text{loss} (\theta_t, z) - \nabla \text{Loss}(\theta_t) \right\|^2 \right] \leq \sigma^2$

(f, κ) -Robust Averaging: A Solution to Robust Estimation



Sanity check: If $n - f$ vectors are identical the output is that vector.
 If honest vectors are separated by distance δ , the error is in $\mathcal{O}(\kappa \delta^2)$

When $\text{RobAgg}(g_t^{(1)}, \dots, g_t^{(n)})$ is (f, κ) -robust, we have

$$\mathbb{E} \left[\|\text{RobAgg}(g_t^{(1)}, \dots, g_t^{(n)}) - \nabla \text{Loss}(\theta_t)\|^2 \right] \leq \mathcal{O} \left(\left(\kappa + \frac{1}{n} \right) \sigma^2 \right)$$

(f, κ) -Robust Averaging (cont'd)

Aggregation $F : (v_1, \dots, v_n) \mapsto v \xrightarrow{\forall G \subseteq [n] \text{ of size } n-f} \|v - \bar{v}_G\|^2 \leq \kappa \left(\frac{1}{n-f} \sum_{i \in G} \|v_i - \bar{v}_G\|^2 \right)$

$$\mathbb{E} \left[\left\| \text{RobAgg} (g_t^{(1)}, \dots, g_t^{(n)}) - \frac{1}{n-f} \sum_{i \in H} g_t^{(i)} \right\|^2 \right] \leq 4\kappa \sigma^2$$

$$f < \frac{n}{2} \quad \kappa \in \Omega \left(\frac{f}{n} \right)$$

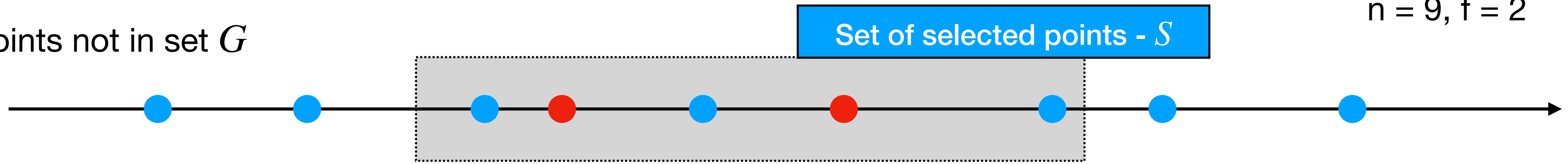
$$\mathbb{E} \left[\left\| \text{CWTM} (g_t^{(1)}, \dots, g_t^{(n)}) - \frac{1}{n-f} \sum_{i \in H} g_t^{(i)} \right\|^2 \right] \in \mathcal{O} \left(\frac{f}{n} \sigma^2 \right)$$

Agg.	κ	When $f \ll n$
CWMed	$4 \left(1 + \frac{f}{n-2f} \right)^2$	$\mathcal{O}(1)$
CWTM	$\frac{6f}{n-2f} \left(1 + \frac{f}{n-2f} \right)$	$\mathcal{O} \left(\frac{f}{n} \right)$

(f, κ) -Robustness of Trimmed Mean

● Points in set $G \subset [n]$, $|G| = n - f$

● Points not in set G



Let $v_1 \leq \dots \leq v_n$, then $\hat{v} = \text{TM}(v_1, \dots, v_n) = \frac{1}{n - 2f} \sum_{i=f}^{n-f-1} v_i$

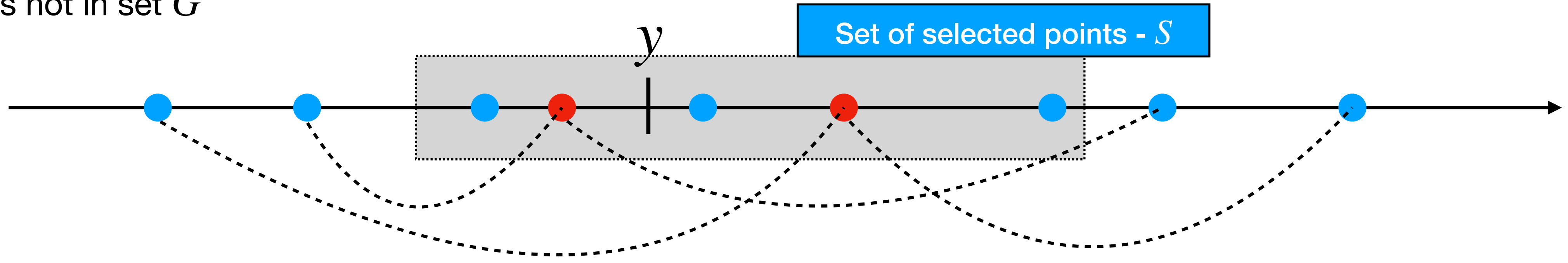
$$(\hat{v} - \bar{v}_G)^2 = \frac{1}{(n - 2f)^2} \left(\sum_{i \in S} v_i - \bar{v}_G \right)^2 = \frac{1}{(n - 2f)^2} \left(\sum_{i \in S \setminus G} v_i - \bar{v}_G + \sum_{i \in S \cap G} v_i - \bar{v}_G \right)^2$$

$$(\hat{v} - \bar{v}_G)^2 = \frac{1}{(n - 2f)^2} \left(\sum_{i \in S \setminus G} v_i - \bar{v}_G - \sum_{i \in G \setminus S} v_i - \bar{v}_G \right)^2 \leq \frac{|S \setminus G| + |G \setminus S|}{(n - 2f)^2} \left(\sum_{i \in S \setminus G} (v_i - \bar{v}_G)^2 + \sum_{i \in G \setminus S} (v_i - \bar{v}_G)^2 \right)$$

(f, κ) -Robustness of Trimmed Mean (Cont'd)

● Points in set $G \subset [n]$, $|G| = n - f$

● Points not in set G



$$(\hat{v} - \bar{v}_G)^2 \leq \frac{|S \setminus G| + |G \setminus S| \leq 3f}{(n - 2f)^2} \left(\sum_{i \in S \setminus G} (v_i - \bar{v}_G)^2 + \sum_{i \in G \setminus S} (v_i - \bar{v}_G)^2 \right)$$

$$\forall j \in S \setminus G, \exists i_j \in G \setminus S, \text{ such that } (v_j - y)^2 \leq (v_{i_j} - y)^2$$

$$(\hat{v} - \bar{v}_G)^2 \leq \frac{3f}{(n - 2f)^2} \left(\sum_{i \in G \setminus S} (v_i - \bar{v}_G)^2 + \sum_{i \in G \setminus S} (v_i - \bar{v}_G)^2 \right) \leq \frac{6f(n - f)}{(n - 2f)^2} \left(\frac{1}{n - f} \sum_{i \in G} (v_i - \bar{v}_G)^2 \right)$$

$\searrow \kappa$

(f, κ) -Robustness of Geometric Median

$$\hat{v} = \text{GM}(v_1, \dots, v_n) = \min_v \sum_{i=1}^n \|v_i - v\|$$

For all $i \in G$, use: $\|\hat{v} - \bar{v}_G\| \leq \|\hat{v} - v_i\| + \|v_i - \bar{v}_G\|$

For all $i \in [n] \setminus G$, use: $\|\hat{v} - \bar{v}_G\| \geq \|\bar{v}_G - v_i\| - \|v_i - \hat{v}\|$

Niceness of (f, κ) -Robust Averaging

Aggregation	Breakdown	Robustness Coeff. κ	When $f \ll n$
CWMed	$\frac{1}{2}$	$4 \left(1 + \frac{f}{n-2f}\right)^2$	$\mathcal{O}(1)$
CWTM	$\frac{f+1}{n}$	$\frac{6f}{n-2f} \left(1 + \frac{f}{n-2f}\right)$	$\mathcal{O}\left(\frac{f}{n}\right)$
GeoMed	$\frac{1}{2}$	$4 \left(1 + \frac{f}{n-2f}\right)^2$	$\mathcal{O}(1)$
Krum	$\frac{f+1}{n}$	$6 \left(1 + \frac{f}{n-2f}\right)^2$	$\mathcal{O}(1)$

Robustness analysis is agnostic to input distribution.

Incorporate input pre-processing (e.g., clipping) easily.

Extends easily to heterogeneous settings.

Learning Error Rate of Robust Mini-batch GD

$$\theta_{t+1} = \theta_t - \gamma_t \text{RobAgg} \left(g_t^{(1)}, \dots, g_t^{(n)} \right)$$

(f, κ) -Robust

Loss(θ) $\begin{cases} \text{Lipschitz smooth} \\ \text{Strongly convex} \end{cases}$

$$\mathbb{E}_{z \sim \mathcal{U}(S)} \left[\|\nabla \text{loss}(\theta, z) - \nabla \text{Loss}(\theta)\|^2 \right] \leq \sigma^2$$

Minimum value

After T iterations:

$$\mathbb{E} \left[\text{Loss}(\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{\sigma^2}{(n-f)T} + \kappa \sigma^2 \right)$$

DSGD with $n - f$ honest nodes

Gradient noise

Can We Do Better?

The squared-error in estimating $\nabla \text{Loss}(\theta_t)$ is in $\Omega\left(\frac{f}{n} \sigma^2\right)$

Trimmed-Mean
is optimal

Honest gradients

Adversarial gradients

$\nabla \text{Loss}(\theta_t)$

θ_t

Introduces **unavoidable** bias in
estimation of true gradient $\nabla \text{Loss}(\theta_t)$

Why Can't We Do Better?

n points sampled from $P' = \left(1 - \frac{f}{n}\right)P + \frac{f}{n}Q$, where $P = 0$ w.p. 1, and $Q = \sigma\sqrt{\frac{n}{f}}$ w.p. 1.

An algorithm \mathcal{A} cannot tell if the honest distribution is either P or P' ,

Huber's contamination model

Both p and p' are valid honest distributions: $\sigma_p^2 = 0$ and $\sigma_{p'}^2 = \sigma^2 \left(1 - \frac{f}{n}\right)$.

We have $\mu_p = 0$ and $\mu_{p'} = \sigma\sqrt{\frac{f}{n}}$. Thus, $|\mu_p - \mu_{p'}|^2 = \sigma^2 \frac{f}{n}$.

Therefore, \mathcal{A} incurs an error overhead in estimating the true mean $\in \Omega\left(\sigma^2 \frac{f}{n}\right)$.

Overall error $\in \Omega\left(\sigma^2 \left(\frac{1}{n} + \frac{f}{n}\right)\right)$.

Non-Vanishing Robustness Error is Problematic in Practice

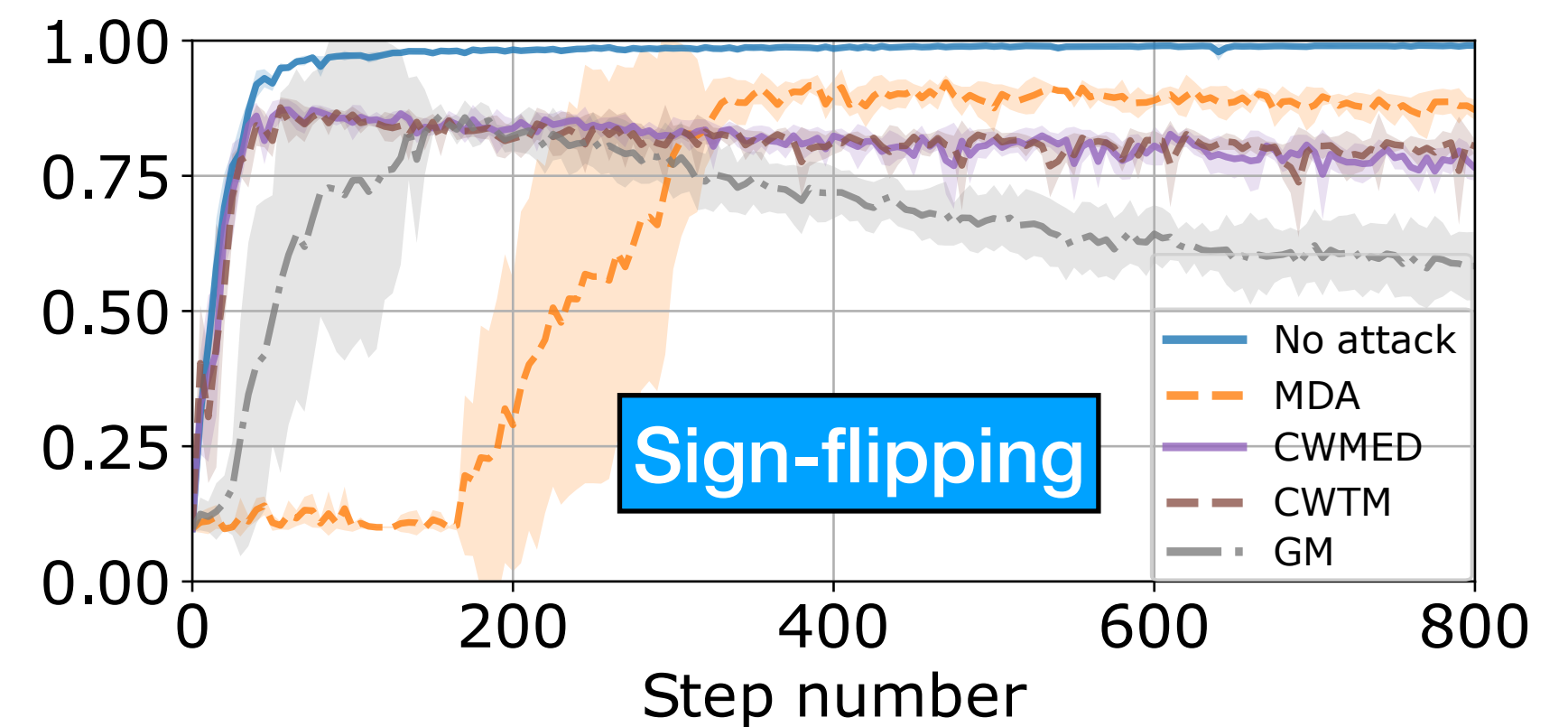
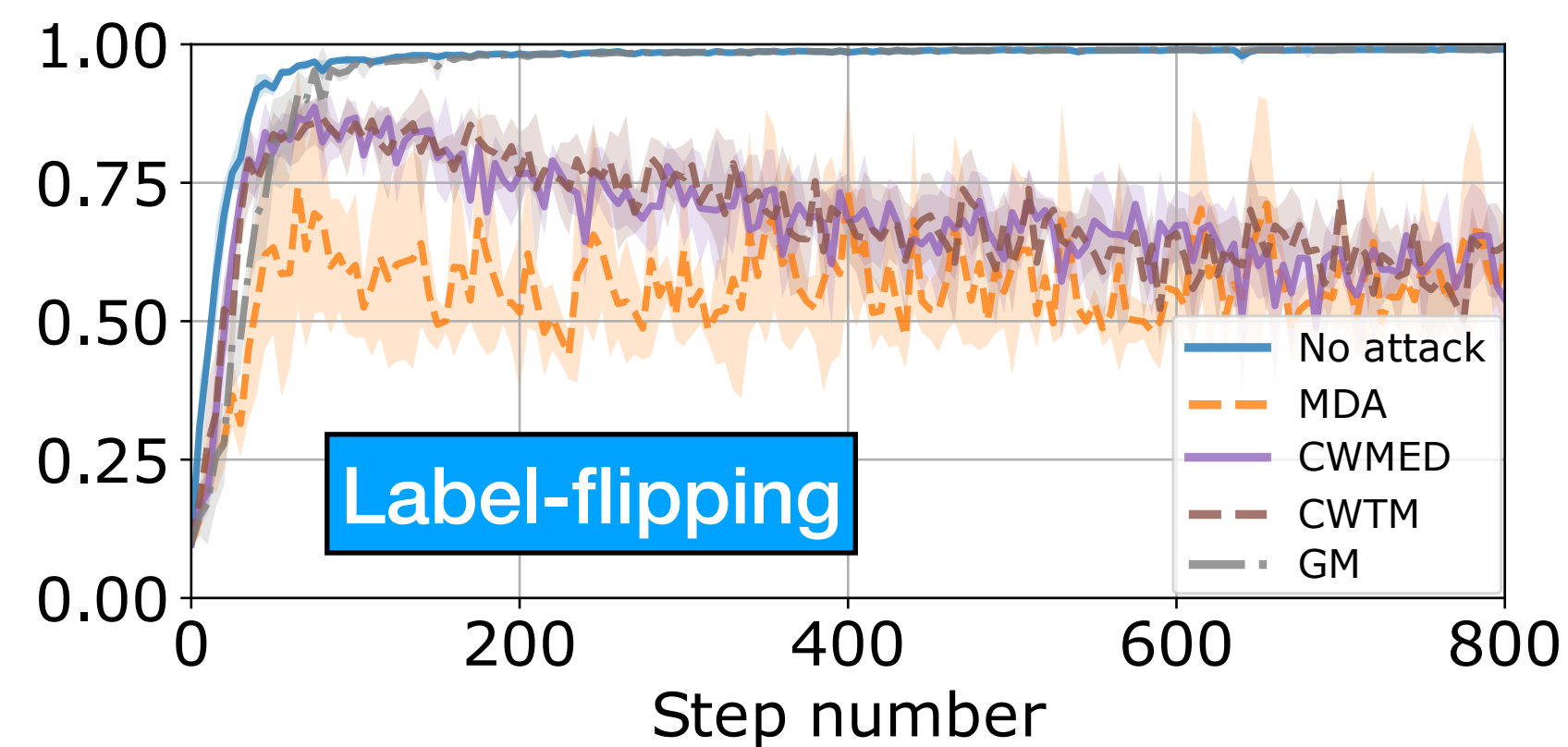
$$\mathbb{E} \left[\text{Loss} (\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{\sigma^2}{(n-f)T} + \kappa \sigma^2 \right)$$

Non-negligible in practice



Training **CNN** using $n = 15$ nodes where $f = 5$ nodes are adversarial

MNIST dataset is equally divided amongst the honest nodes

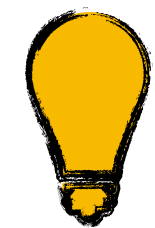


Mitigating the Impact of Noise

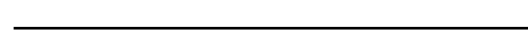
$$\mathbb{E} \left[\text{Loss} (\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{\sigma^2}{(n-f)T} + \kappa \sigma^2 \right)$$

↓
Reduce the noise

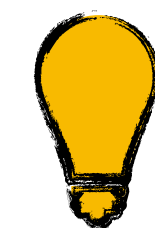
Few Ideas -



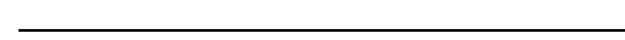
Increase local batch-sizes



Imposes high computation cost



Stochastic variance reduction



Requires computing full-batch gradients w.p. > 0



Incorporate Local Gradient Momentum

Assumption: Identity of corrupted nodes is fixed throughout the training procedure

Local Phase: Each honest node i computes $g_t^{(i)} := \nabla \text{loss}(\theta_t, z_t^{(i)})$ and computes (Polyak's) momentum:

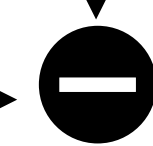
$$m_t^{(i)} = \beta m_{t-1}^{(i)} + (1 - \beta) g_t^{(i)}$$

Past gradients

Global Phase: $\theta_{t+1} = \theta_t - \gamma_t \text{RobAvg}(m_t^{(1)}, \dots, m_t^{(n)})$

Introduces bias $\approx \beta$

Honest variance reduces by factor $(1 - \beta)$



Controlled momentum rate: $\beta = 1 - c \gamma_t$

Byzantine Machine Learning Made Easy
by Resilient Averaging of Momentums
FGGPS. ICML, 2022.

$$\mathbb{E} \left[\text{Loss}(\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\left(\frac{1}{n-f} + \kappa \right) \frac{\sigma^2}{T} \right)$$

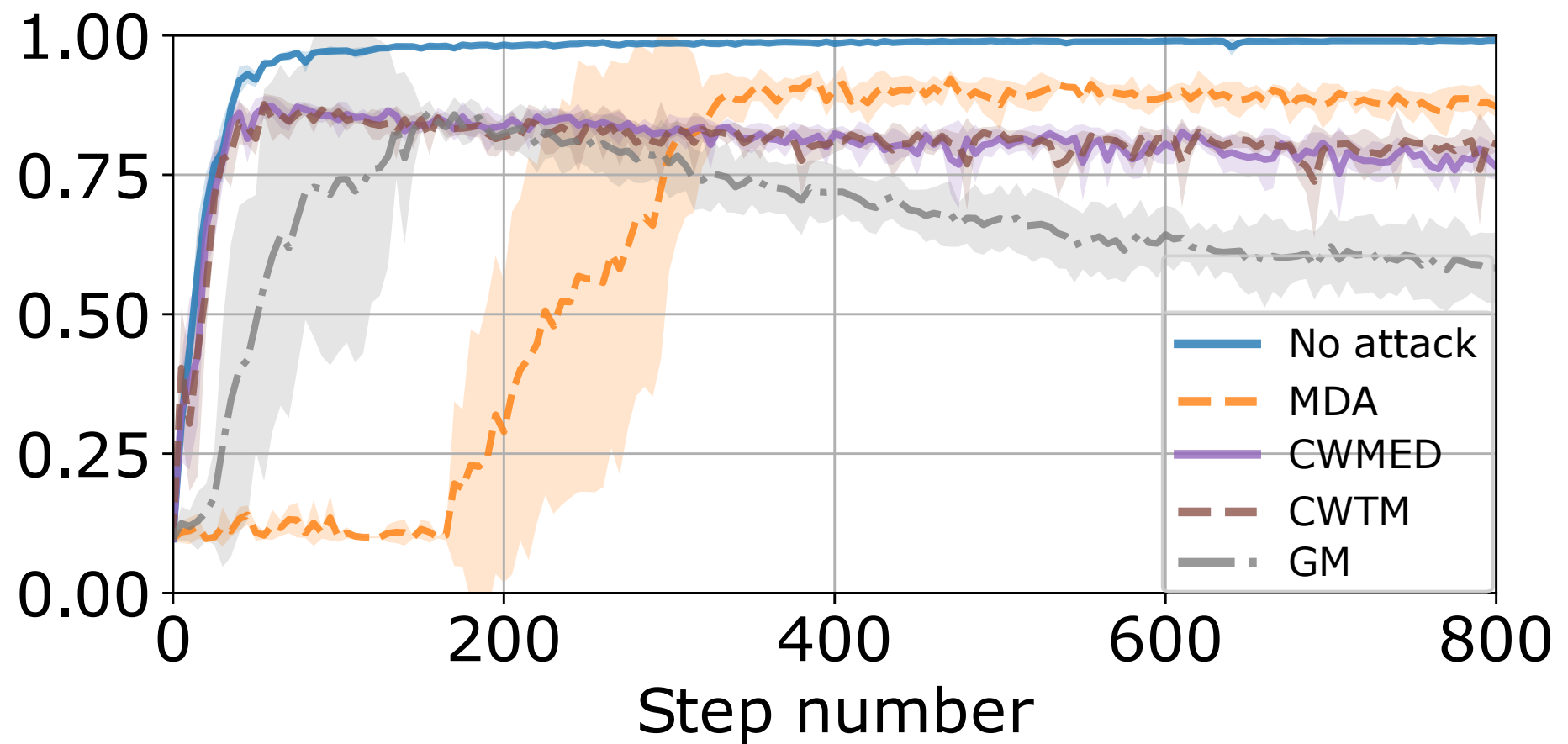
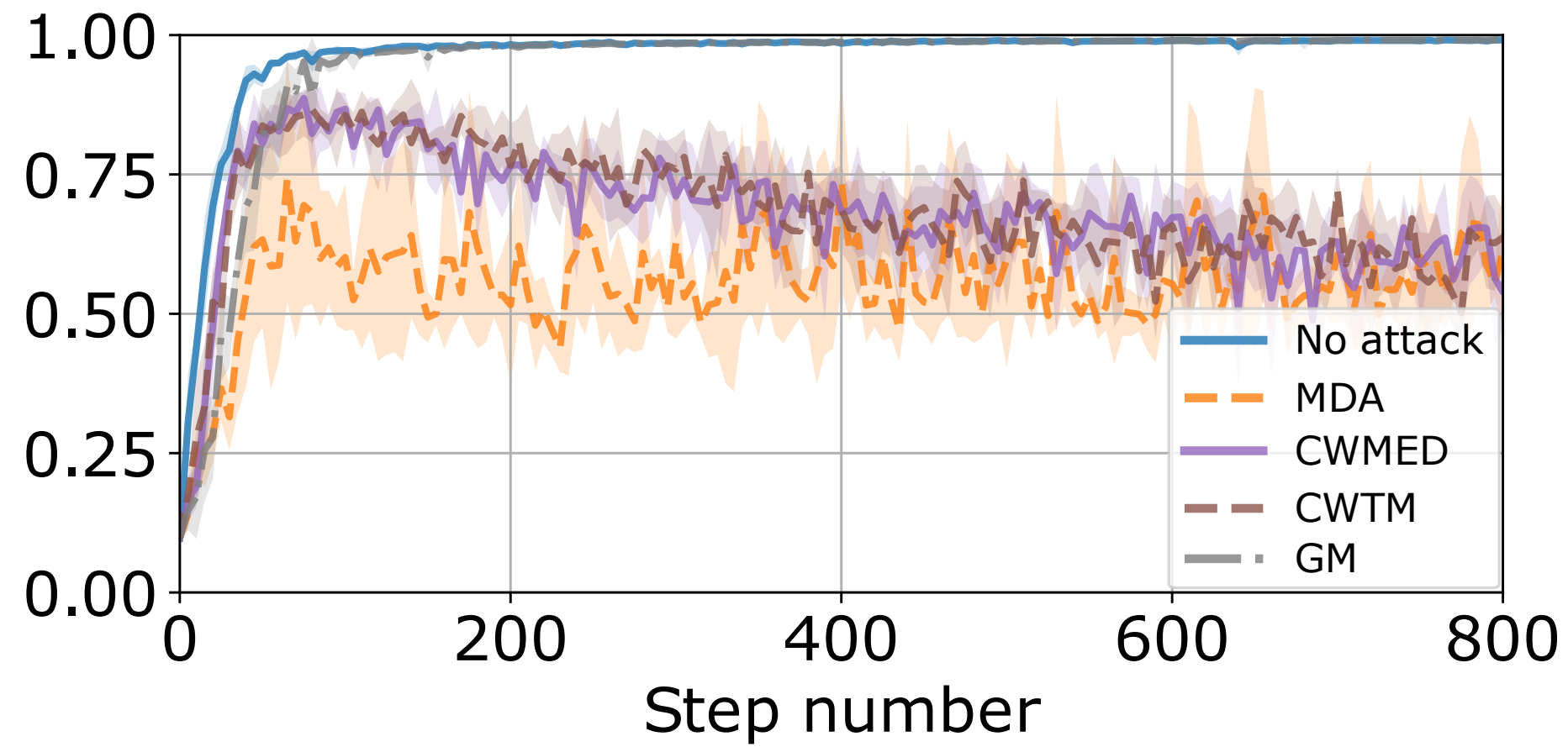
Trimmed-mean

Optimal when
 $\kappa \in \mathcal{O}(f/n)$

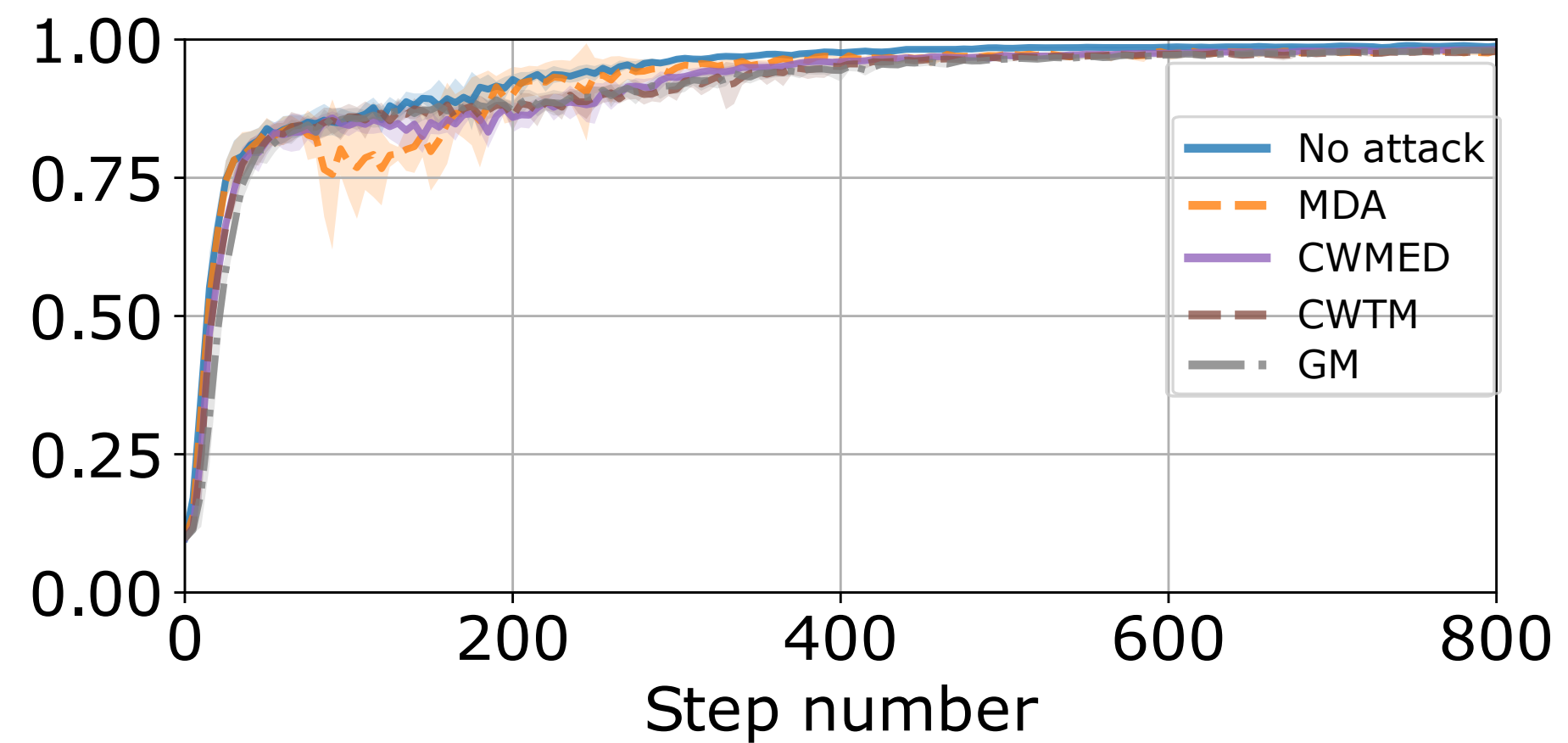
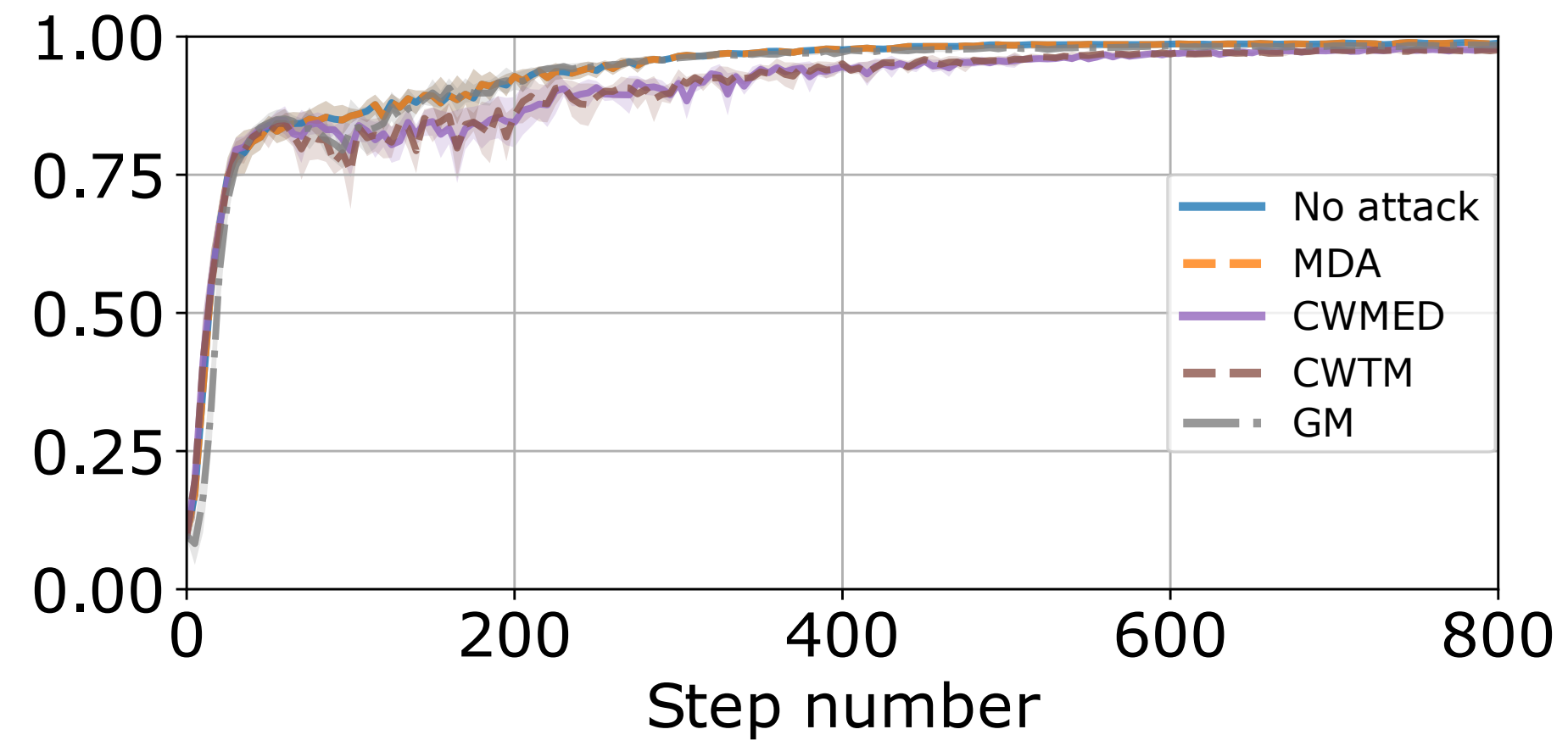
Augments gradient complexity

Empirical Benefits of Momentum

Without Momentum



With Momentum $\beta = 0.9$



Label-flipping

Sign-flipping

Training **CNN** on $n = 15$ nodes where $f = 5$ nodes are adversarial.
MNIST dataset is equally divided amongst the honest nodes

Data Poisoning \neq Model (gradient) Poisoning

Data poisoning



$$\mathbb{E} \left[\text{Loss}(\theta_T) - \text{Loss}^* \right] \in \Omega \left(\frac{1+f}{n} \cdot \frac{1}{T} \right)$$

Gradient poisoning



$$\in \mathcal{O} \left(\frac{1+f}{n} \cdot \frac{K}{T} \right)$$

Condition number of $\text{Loss}(\theta)$

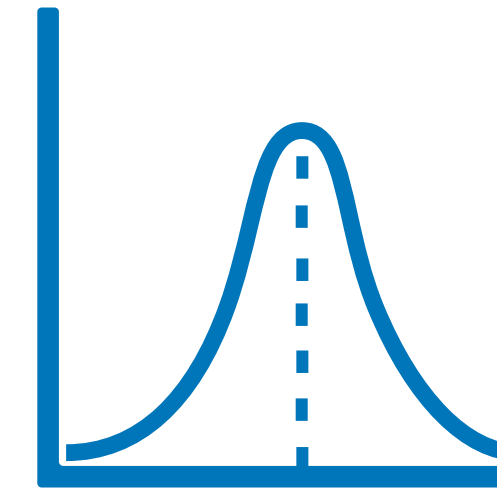
Adapting to the learning procedure does not help an adversary

Take Away: Robustness under Local Randomness

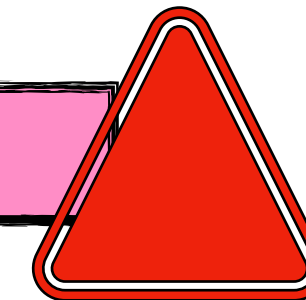
Robust Averaging
Local Momentum

Local Randomness

Even honest machines
compute noisy updates



Adapting to the learning procedure **does not help** an adversary



Byzantine Machine Learning Made Easy by Resilient Averaging of Momentums

S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *International Conference on Machine Learning (ICML)*, 2022.

Fixing by Mixing: A Recipe for Optimal Byzantine ML under Heterogeneity

Y. Allouah, S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023.

Robust Collaborative Learning with Linear Gradient Overhead

S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *ICML*, 2023.

**Serverless
architecture**

Byzantine Failures Hurt Generalization More Than Data Poisoning

T. Boudou, B. Le Bars, N. Gupta, and A. Bellet. In *International Conference on Machine Learning (ICML)*, 2026.

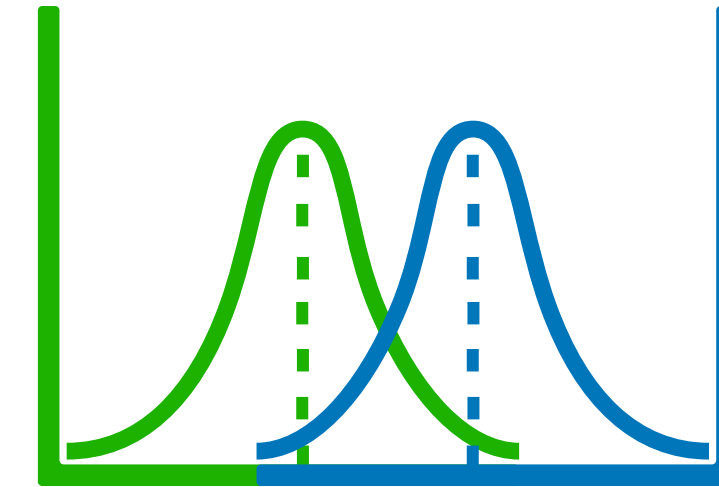


Challenge II: Tackling Heterogeneity

Adversarial nodes can camouflage as honest ones

Data Heterogeneity

Different nodes/clients sample data from different distributions



$S_i \neq S_j$ for any $i, j \in H$

Adversarial clients can exploit heterogeneity to inflict more damage

Impossibility for Robust FL under Heterogeneity

Compute parameters θ^* minimizing $\text{Loss}_H(\theta) := \frac{1}{n-f} \sum_{i \in H} \text{Loss}^i(\theta)$

This goal is generally impossible to achieve,
unless the heterogeneity is bounded

Unless $(2f, \varepsilon)$ – Redundancy

Let $\hat{\theta}_S = \arg \min \text{Loss}_S(\theta)$
We guarantee $\hat{\theta}$ s.t. $\|\hat{\theta} - \theta^*\|^2 \leq \varepsilon$ only if $\forall S' \subseteq S \subseteq [n]$,
with $|S'| = n - 2f$, $|S| = n - f$, we have $\|\hat{\theta}_{S'} - \hat{\theta}_S\|^2 \leq \varepsilon$.

Robust DSGD (with Momentum) under Heterogeneity

Local Phase: Each honest node i computes $g_t^{(i)} := \nabla \text{loss}(\theta_t, z_t^{(i)})$ and updates local Polyak's momentum: $z_t^{(i)} \sim \mathcal{U}(S_i)$

$$m_t^{(i)} = \beta_t m_{t-1}^{(i)} + (1 - \beta_t) g_t^{(i)}$$

Global Phase: The coordinator updates the model: $\theta_{t+1} = \theta_t - \gamma_t \text{RobAvg}(m_t^{(1)}, \dots, m_t^{(n)})$

(f, κ) -robust averaging

Heterogeneity bound

$$\frac{1}{n-f} \sum_{i \in H} \|\nabla \text{Loss}^{(i)}(\theta) - \nabla \text{Loss}(\theta)\|^2 \leq \zeta$$

After T iterations:

$$\mathbb{E} \left[\text{Loss}(\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\left(\frac{1}{n-f} + \kappa \right) \frac{\sigma^2}{T} + \kappa \zeta \right)$$

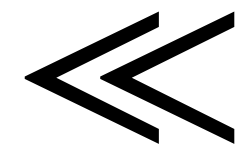
Optimal when
 $\kappa \in \mathcal{O}(f/n)$

Unavoidable in general

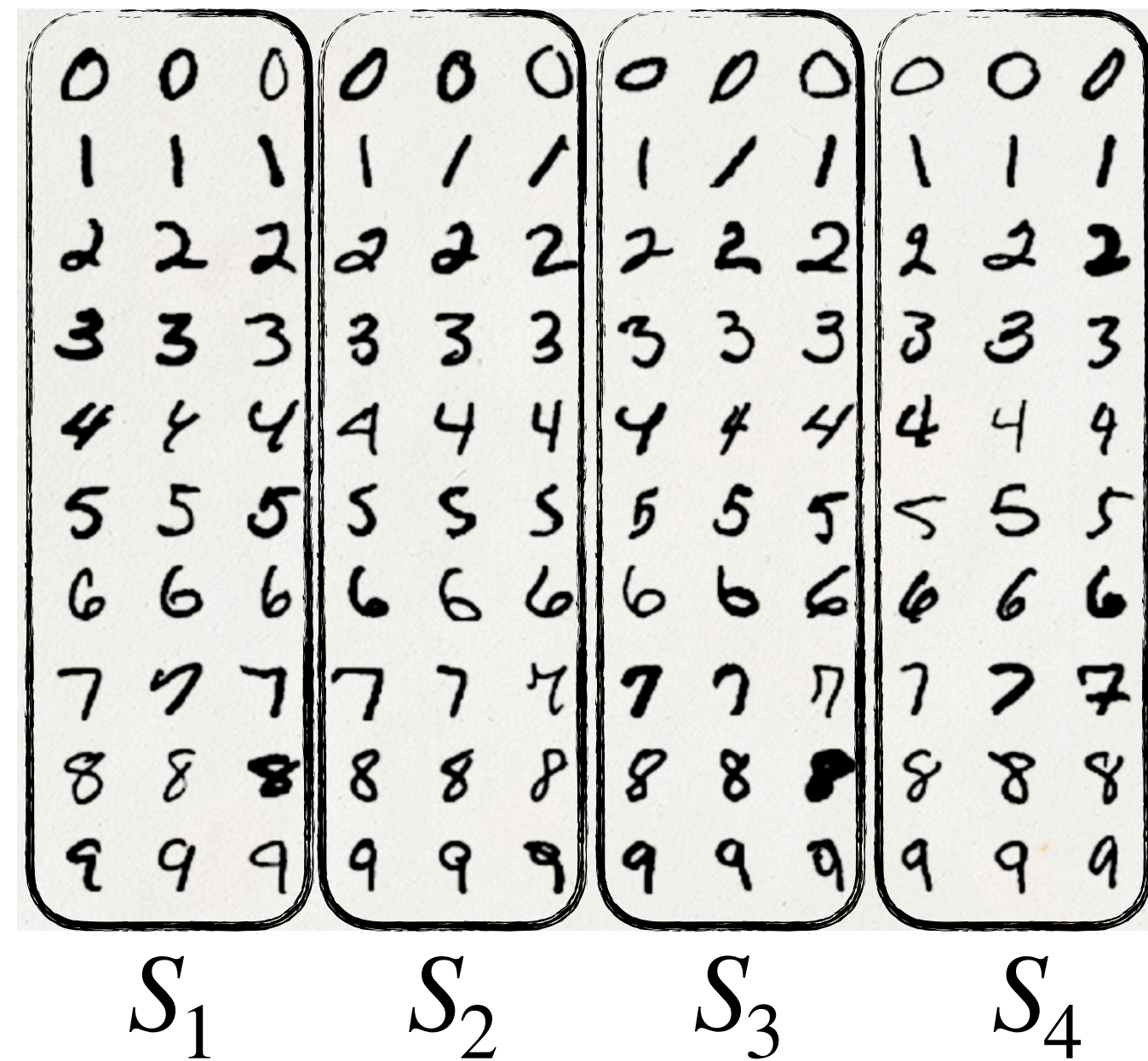
Different Training Samples \iff Heterogeneity

Label-wise heterogeneity

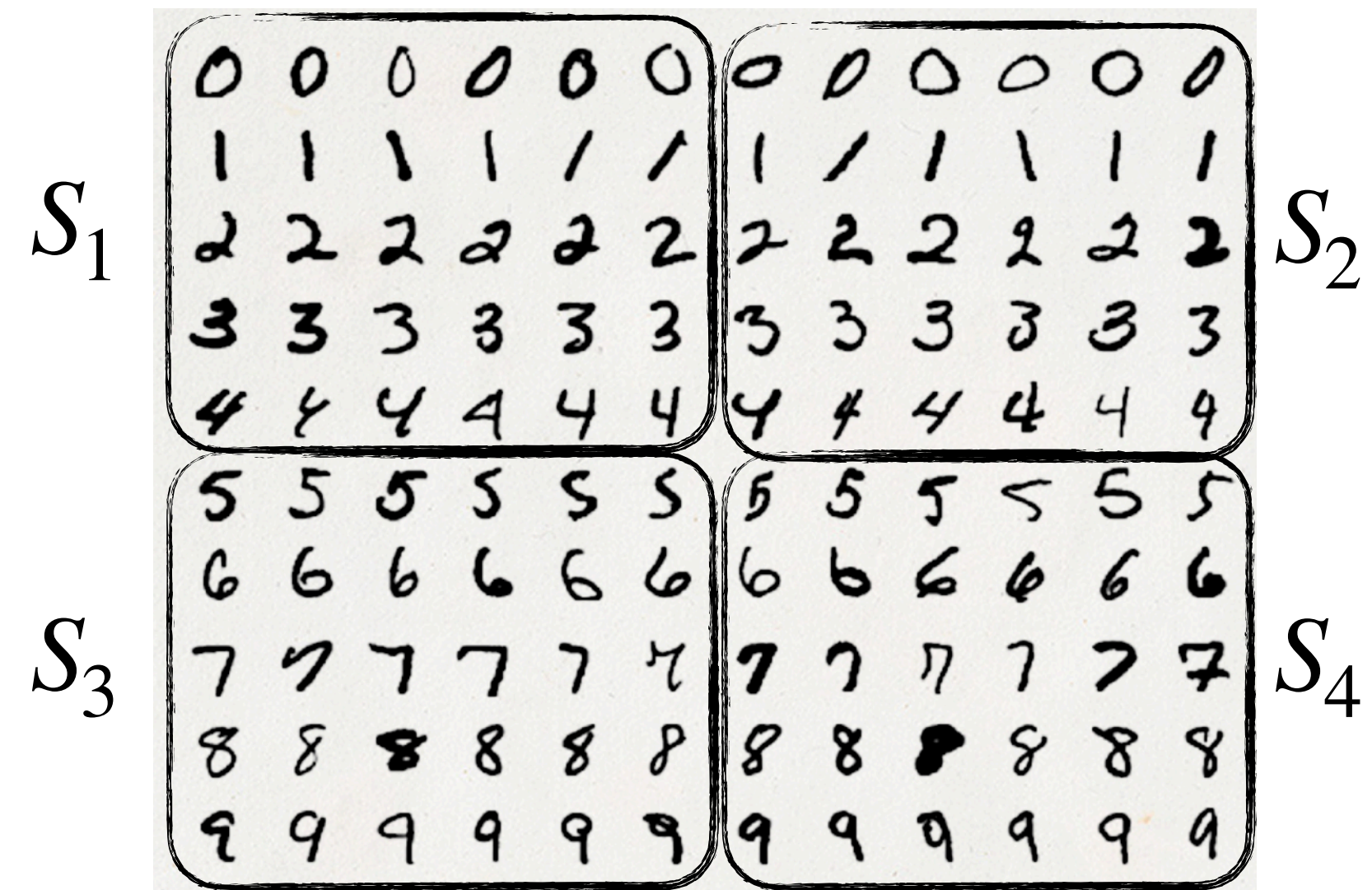
$$\frac{1}{n-f} \sum_{i \in H} \|\nabla \text{Loss}^{(i)}(\theta) - \nabla \text{Loss}(\theta)\|^2$$



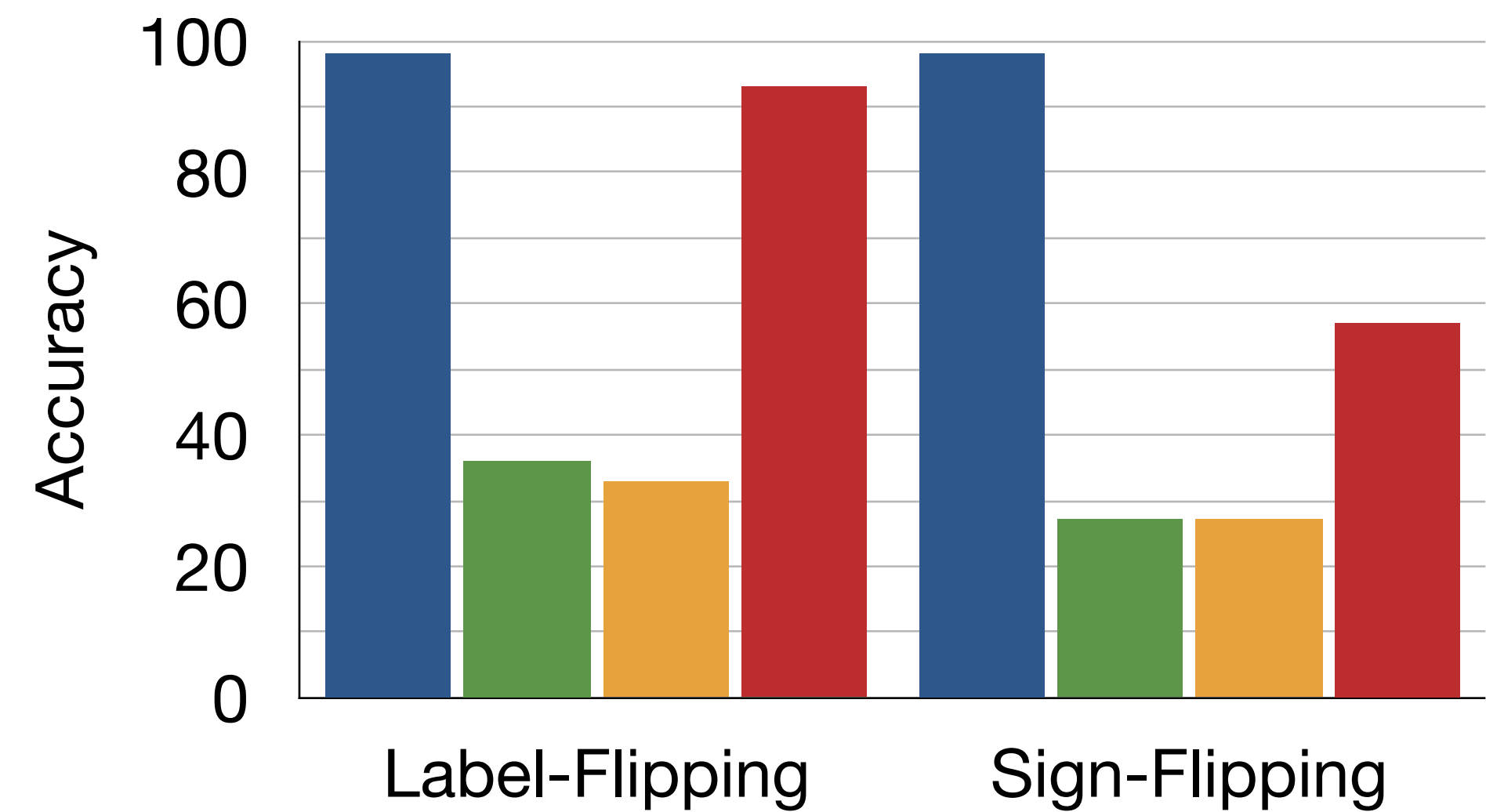
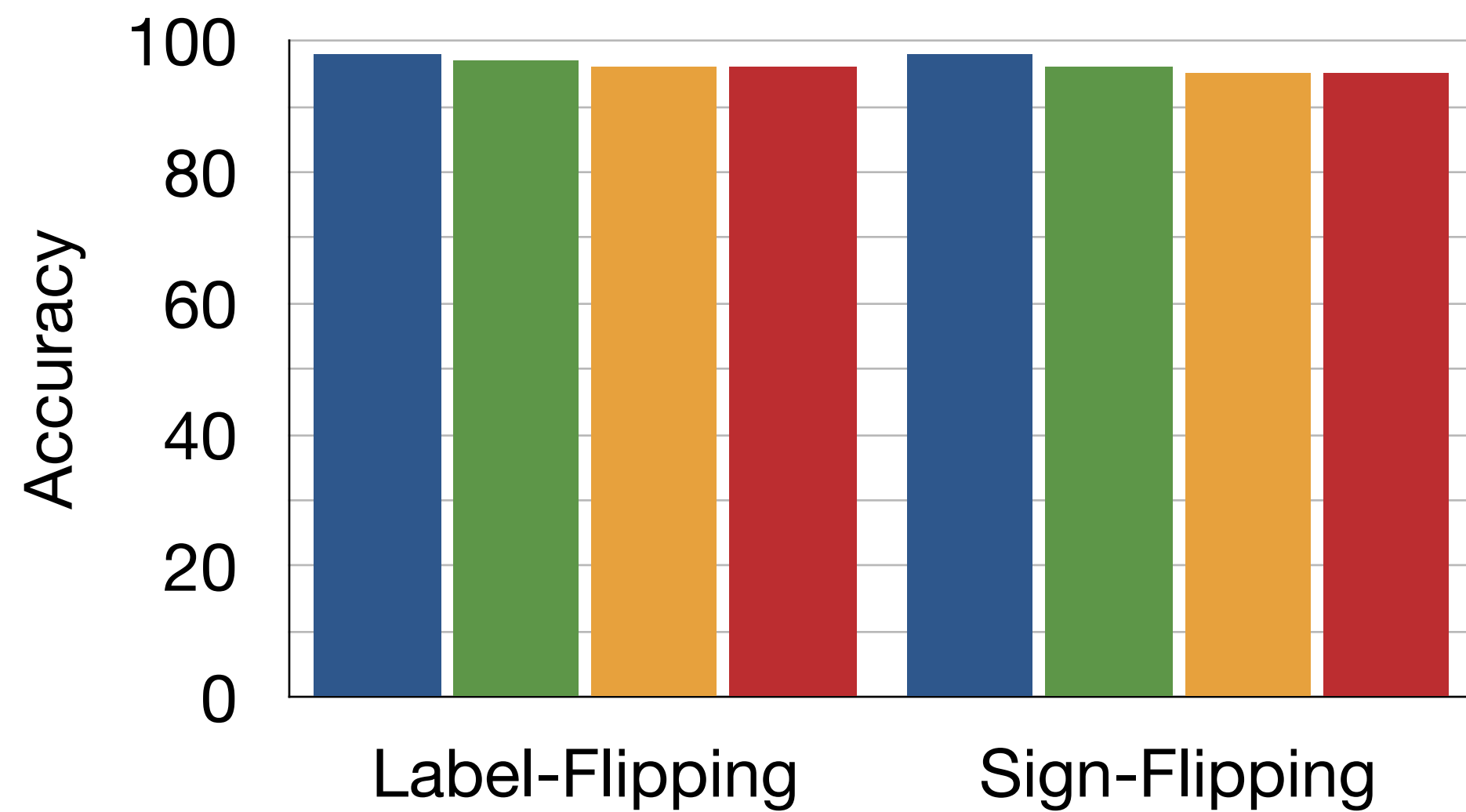
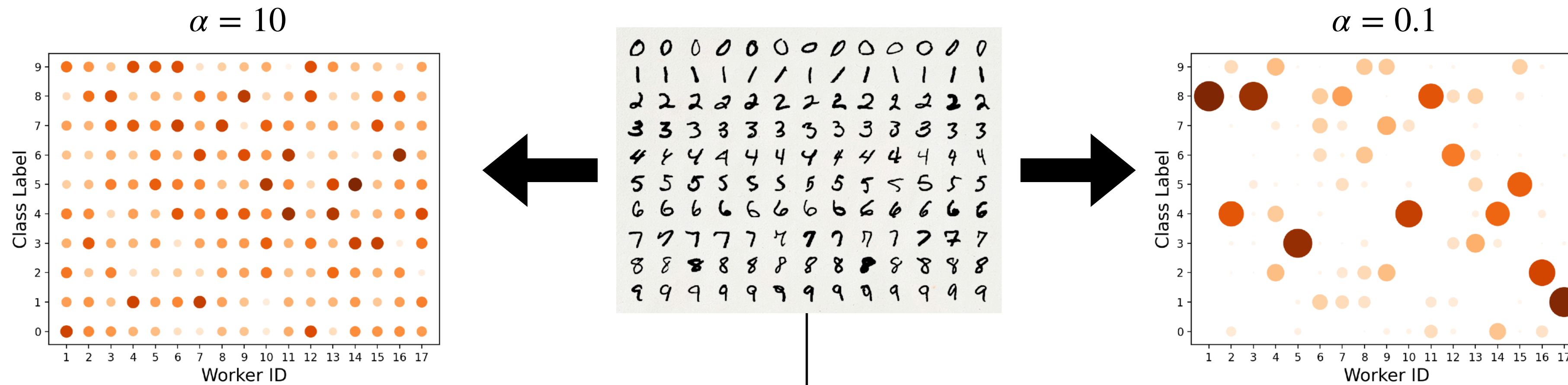
$$\frac{1}{n-f} \sum_{i \in H} \|\nabla \text{Loss}^{(i)}(\theta) - \nabla \text{Loss}(\theta)\|^2$$



$$\nabla \text{Loss}^i(\theta) \approx \nabla \text{Loss}^j(\theta)$$



Numerical Observations under Heterogeneity



Training **CNN** on $n = 17$ nodes where $f = 4$ nodes are adversarial. MNIST dataset is divided amongst nodes using Dirichlet distribution with parameters $\alpha = 10$ and $\alpha = 0.1$



Fixing by Mixing: Reducing Effective Heterogeneity

Pre-Aggregation - Nearest Neighbor Mixing (NNM)

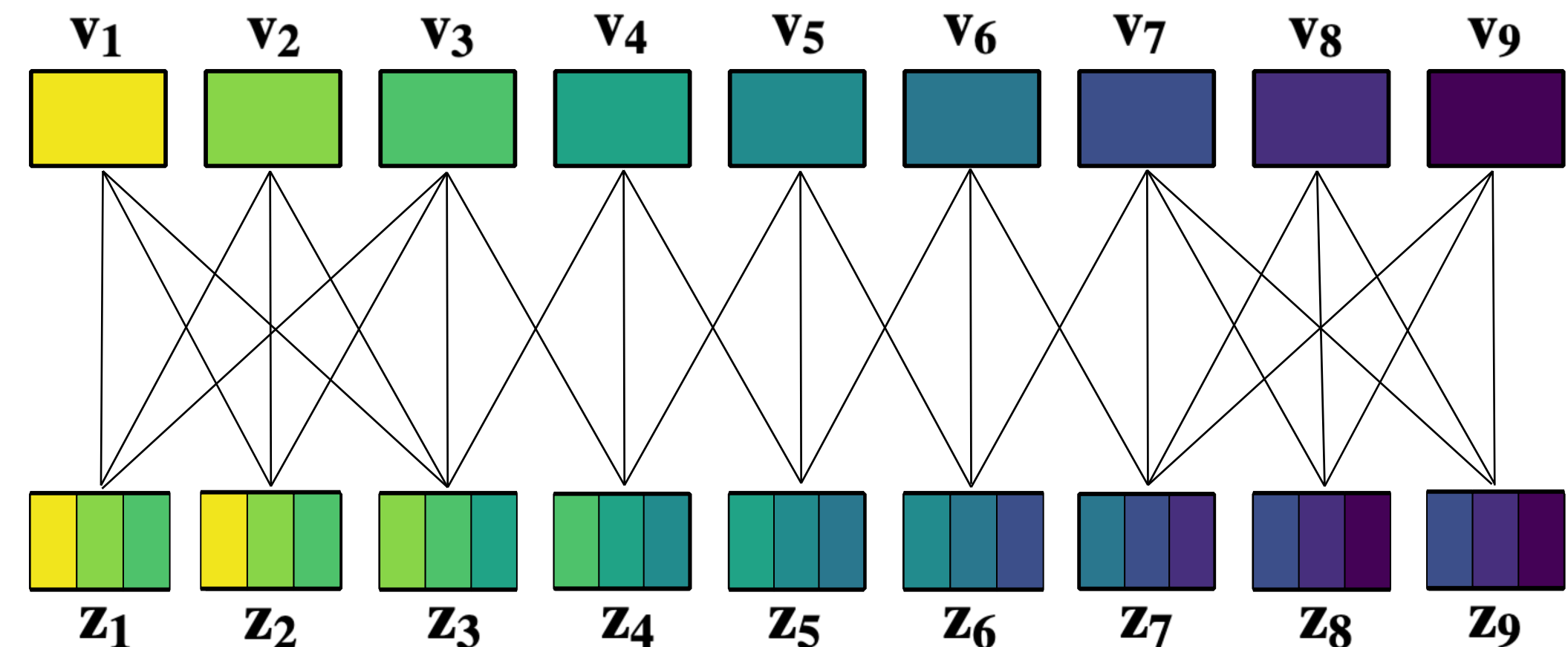
$$\text{RobAvg} \left(\text{NNM} \left(m_t^{(1)}, \dots, m_t^{(n)} \right) \right)$$

Any median-based aggregation

If F is (f, κ) -robust with $\kappa \in \mathcal{O}(1)$ then
 $F \circ \text{NNM}$ is (f, κ) -robust with $\kappa \in \mathcal{O} \left(\frac{f}{n} \right)$

Best possible robustness coefficient

Yields optimal learning guarantee

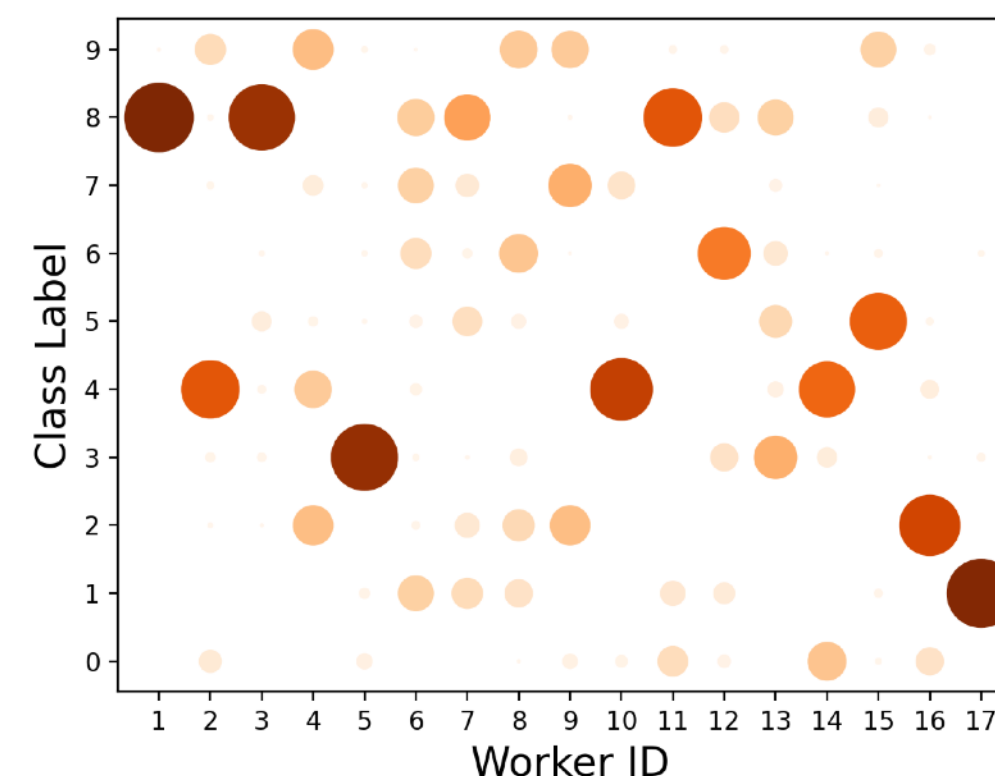
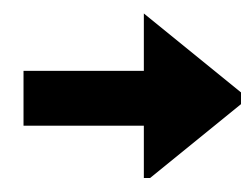


For each v_i choose $n - f$
 nearest neighbors in $\{v_1, \dots, v_n\}$

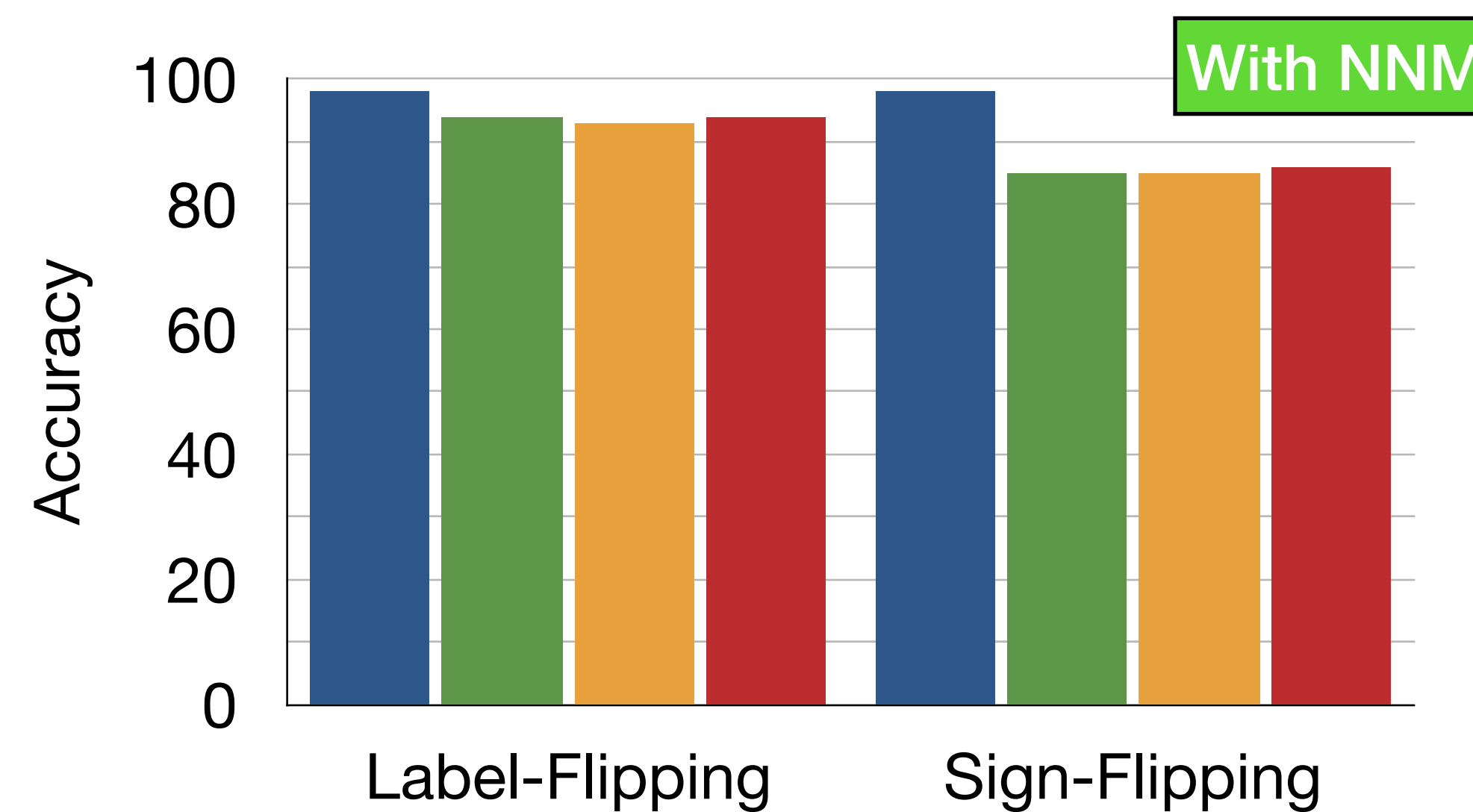
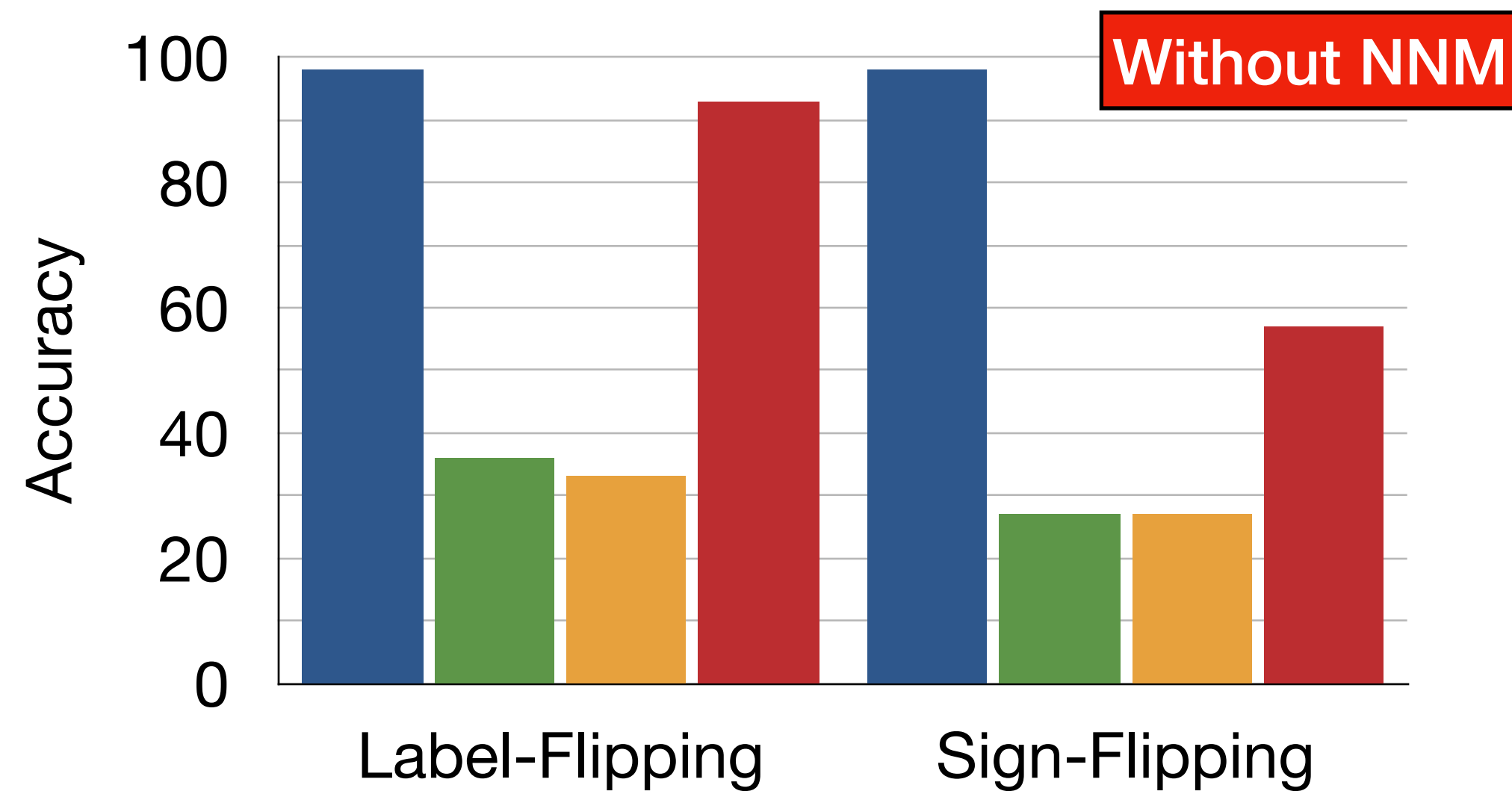
Let N_i be the set of $n - f$ vectors nearest to v_i

$$\text{Map } v_i \text{ to } z_i := \frac{1}{n - f} \sum_{v \in N_i} v$$

Practical Benefits of Nearest Neighbor Mixing (NNM)



$\alpha = 0.1$



Training **CNN** on $n = 17$ nodes where $f = 4$ nodes are adversarial. MNIST dataset is divided amongst nodes using Dirichlet distribution with parameters $\alpha = 0.1$



(G, B) -Dissimilarity & Adaptive Robust Clipping

(G, B) -Gradient Dissimilarity

Refined
heterogeneity
model

$$\frac{1}{n-f} \sum_{i \in H} \|\nabla \text{Loss}^{(i)}(\theta) - \nabla \text{Loss}_H(\theta)\|^2 \leq G + B \|\nabla \text{Loss}_H(\theta)\|^2$$

$$\mathbb{E} \left[\text{Loss}(\theta_T) - \text{Loss}^* \right] \in \Omega \left(\frac{fG}{n - (2+B)f} \right)$$

This lower bound can be remedied under “good” model initialization

ARC

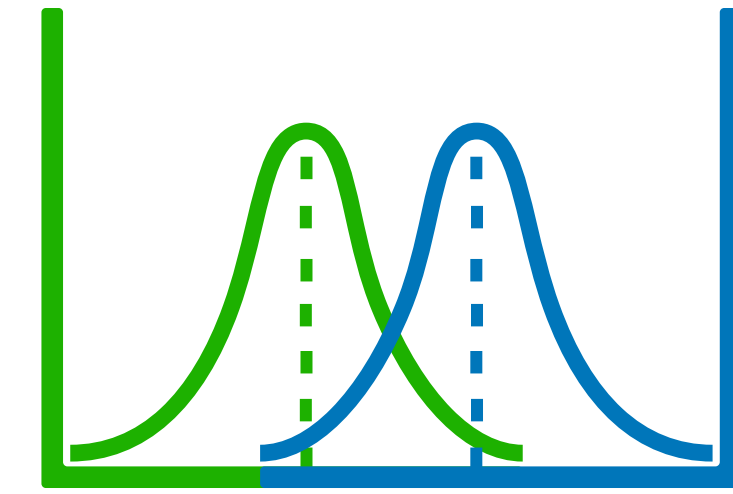
Sort the inputs: $\|v_1\| \leq \dots \leq \|v_n\|$ and clip them
with clipping threshold $C = \|v_k\|$ where $k = \frac{2f}{n}(n-f) + 1$

Take Away: Robustness under Heterogeneity

Nearest Neighbor Mixing &
Adaptive Robust Clipping

Data Heterogeneity

Different clients sample data
from different distributions



Fixing by Mixing: A Recipe for Optimal Byzantine ML under Heterogeneity

Y. Allouah, S. Farhadkhani, R. Guerraoui, [N. Gupta](#), R. Pinot, and J. Stephan. *In AISTATS, 2023.*

Robust distributed learning: Tight error bounds and breakdown point under data heterogeneity

Y. Allouah, R. Guerraoui, [N. Gupta](#), R. Pinot, and G. Rizk. *In NeurIPS, 2023.*

Adaptive Gradient Clipping for Robust Federated Learning

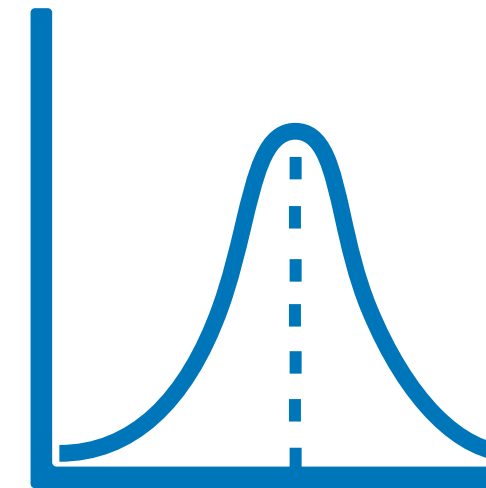
Y. Allouah, R. Guerraoui, [N. Gupta](#), A. Jellouli, G. Rizk, and J. Stephan. *In ICLR, 2025.*

Take Aways: Two Challenges in Robust ML

State machine replication
is not efficient

Local Randomness

Even honest machines
compute noisy updates

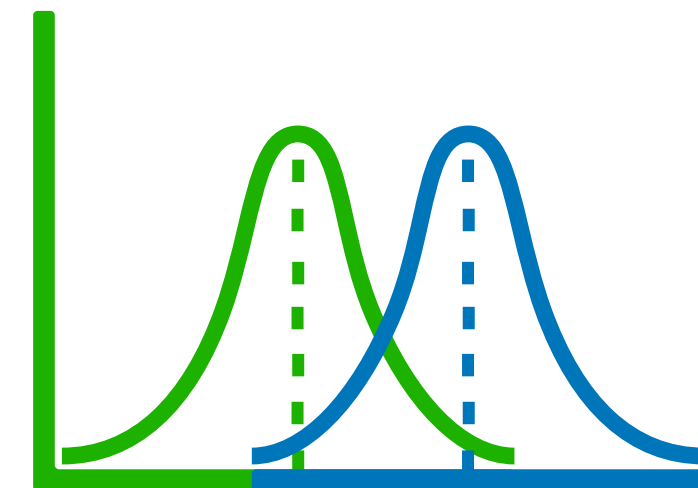


Robust Averaging
Local Momentum

Adversarial nodes can
camouflage as honest nodes

Data Heterogeneity

Different nodes sample data
from different distributions

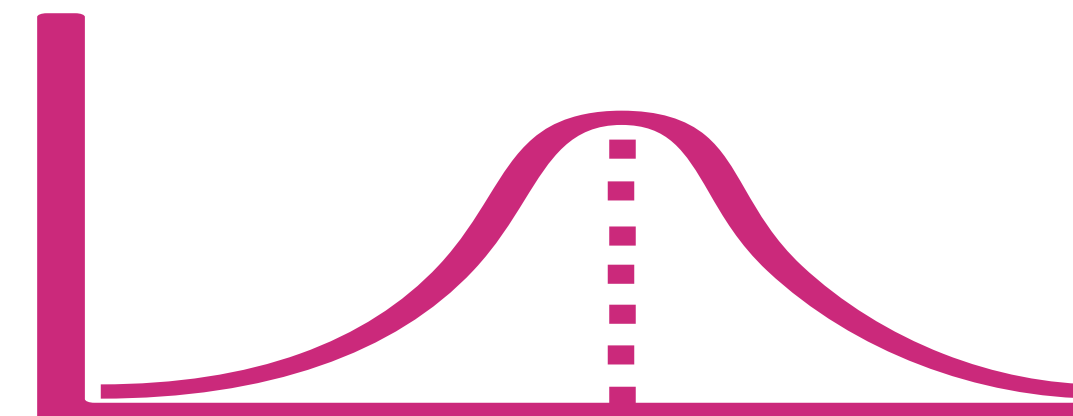


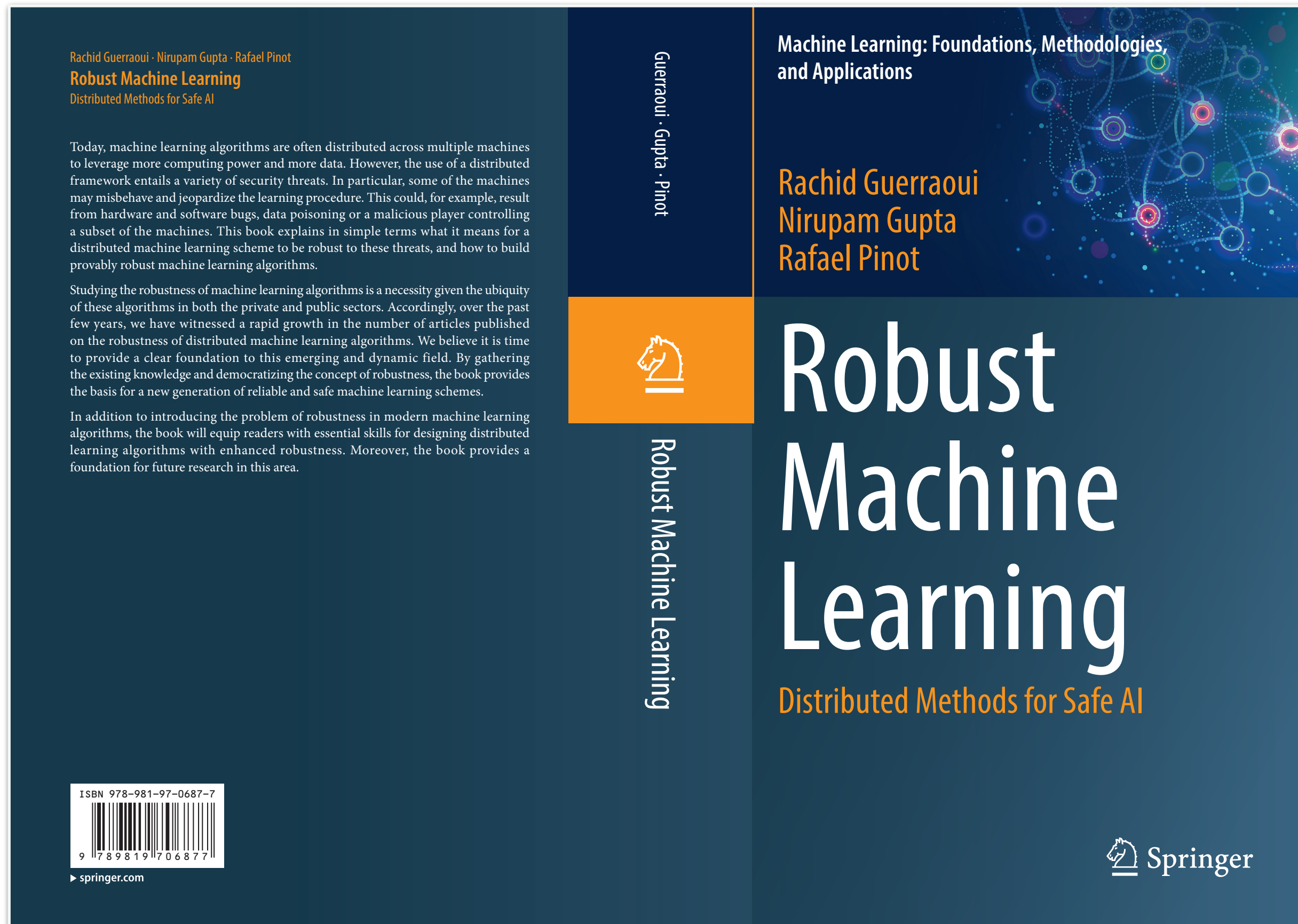
Nearest Neighbor Mixing &
Adaptive Robust Clipping

Augments local
randomness isotropically

Data Privacy

Nodes do not share exact
information on their data





Reading material:

1. Chapter 4 - Robust averaging
2. Chapter 5 - Nearest neighbor mixing
3. Chapter 6 - Local gradient momentum

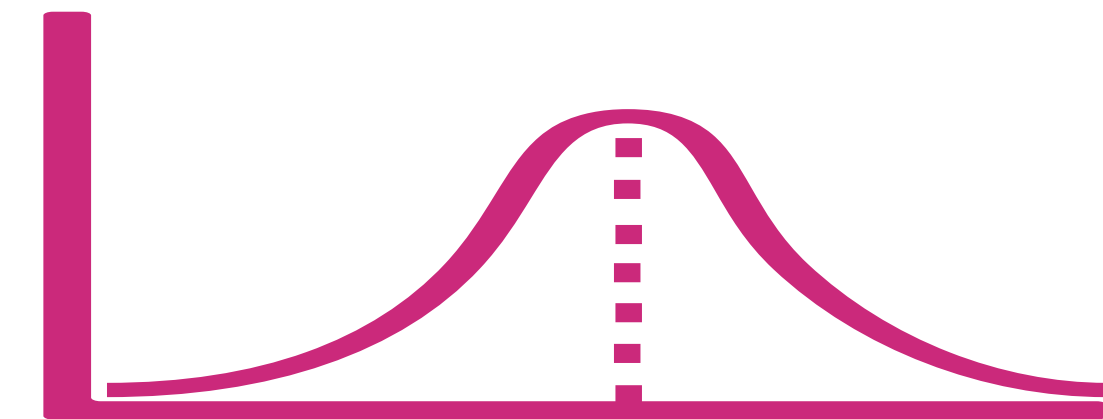
nigu@di.ku.dk

Challenge III: Robustness with Privacy

Augments local randomness isotropically

Data Privacy

Clients do not share exact information on their data

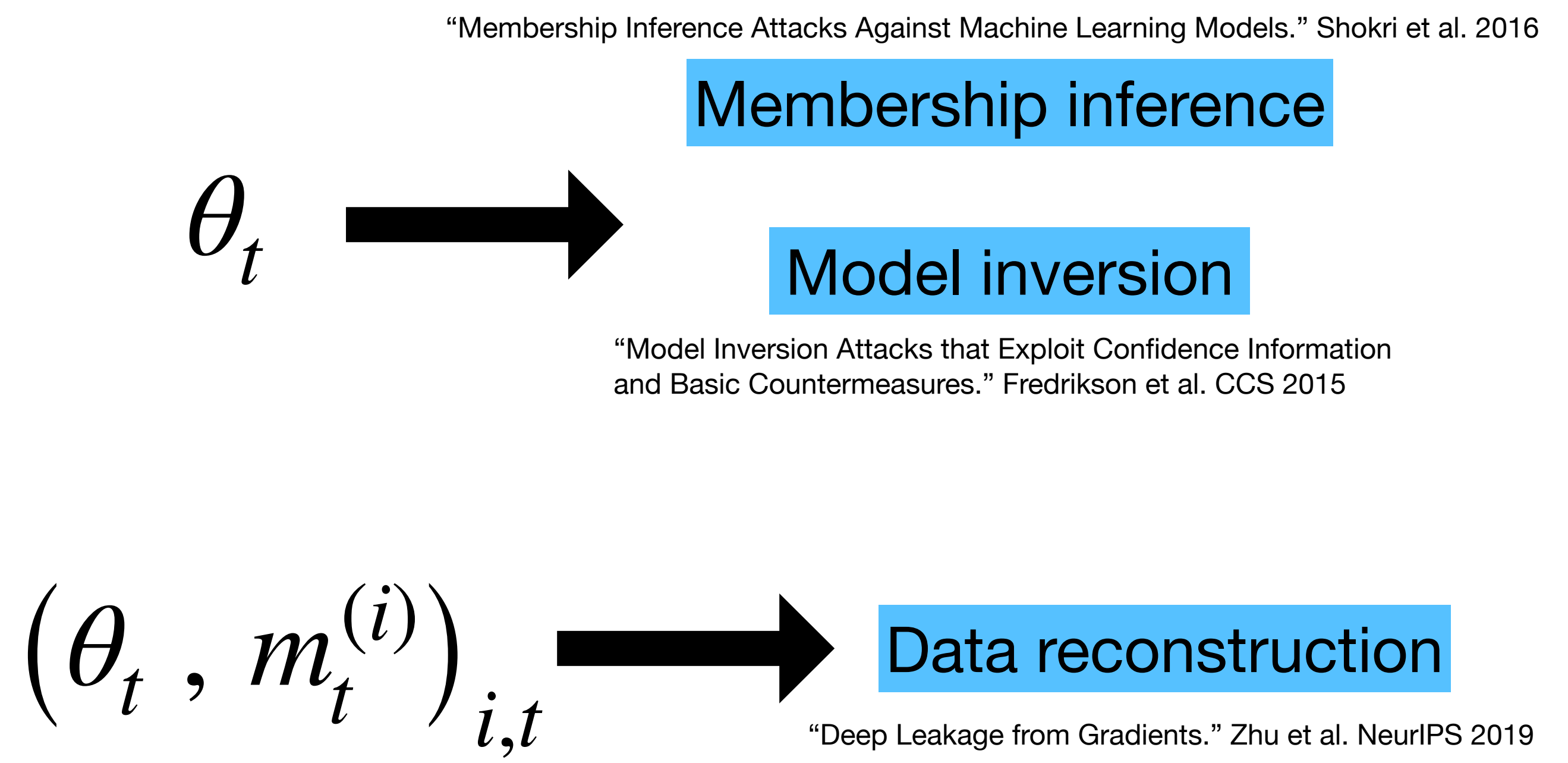
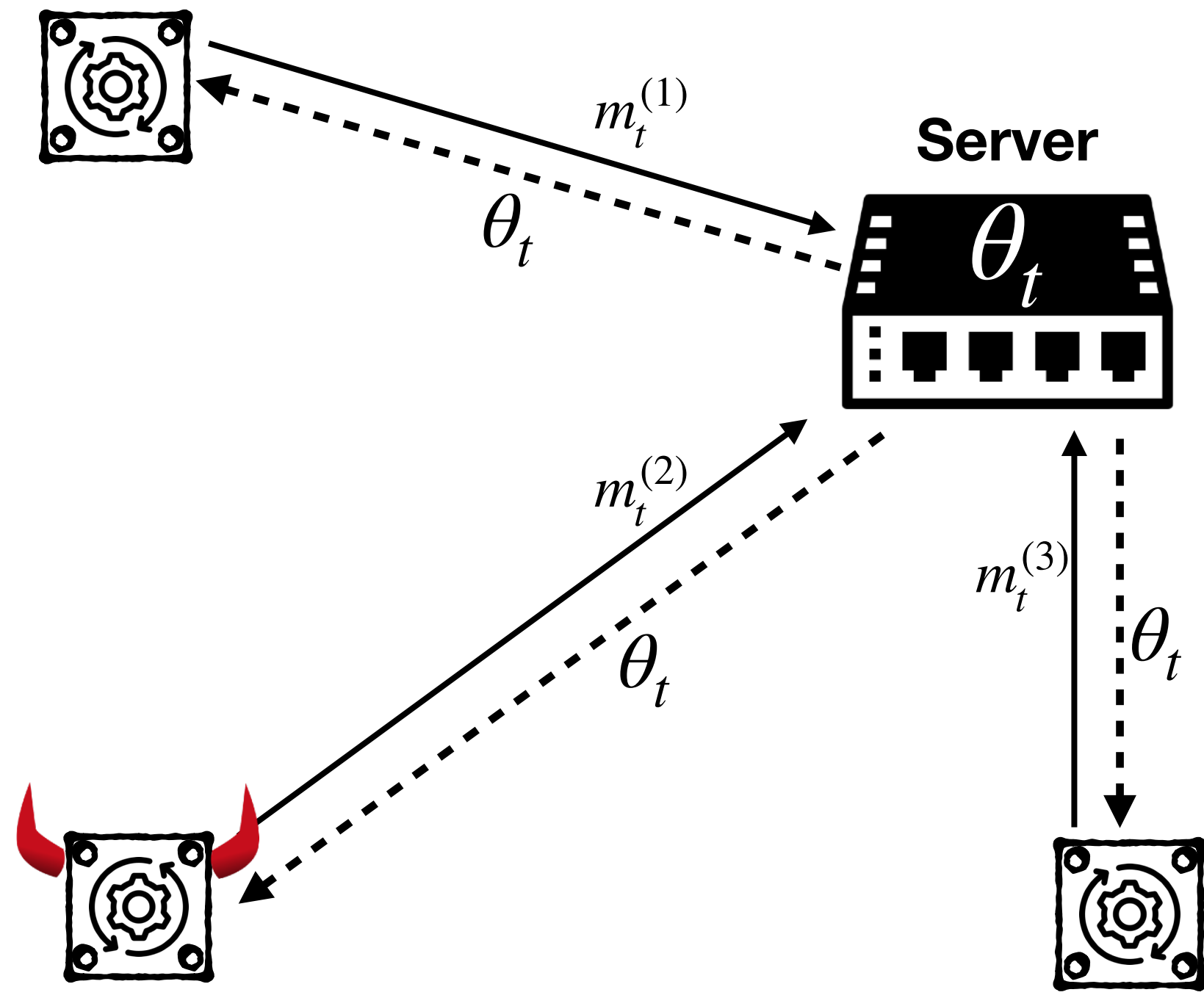


(ϵ, δ) -Distributed DP

Communication between any node i and server is (ϵ, δ) -DP w.r.t. the data held by node i

“Is Interaction Necessary Distributed Private Learning?” A. Smith et al. IEEE S&P 2017.

Robust-DSGD (Similar to DSGD) Leaks Privacy of Data



"On the Privacy-Robustness-Utility Trilemma in Distributed Learning". AGGPS. ICML 2023

Quick Recap of Differentially Private DSGD

Local Phase: Each *honest* node i computes noisy local gradient:

$$g_t^i := \text{Clip} \left(\nabla \text{loss}(\theta_t, z), C \right) + \eta_t$$

with $\eta_t \sim \mathcal{N} \left(0, \sigma_{\text{DP}}^2 I_d \right)$, where $\text{Clip}(v, C) = \min \left\{ 1, \frac{C}{\|v\|} \right\} v$

Gaussian ← η_t

↓ **Clipping**

Global Phase: The server averages the gradients $g_t^{(1)}, \dots, g_t^{(n)}$:

$$\hat{g}_t := \text{Avg} \left(g_t^{(1)}, \dots, g_t^{(n)} \right),$$

updates the current model: $\theta_{t+1} = \theta_t - \gamma_t \hat{g}_t$.

Distributed Differential Privacy (DP)

(ϵ, δ) -Distributed DP

Communication between any node i and server is
 (ϵ, δ) -DP w.r.t. the data held by node i

“Is Interaction Necessary Distributed Private Learning?” A. Smith et al. IEEE S&P 2017.

Algorithm \mathcal{A} is (ϵ, δ) -DP if for $D \sim D'$

$$\Pr(\mathcal{A}(D)) \leq e^\epsilon \Pr(\mathcal{A}(D')) + \delta$$

“Our Data, Ourselves: Privacy via Distributed Noise Generation” C. Dwork et al. Eurocrypt 2006.

Privacy by DP-DSGD

Running T iterations of DP-DSGD

By RDP composition and subsampling amplification theorems, we get

"Rényi Differential Privacy." *Mironov, Ilya*. IEEE CSF, 2017.

$$\text{If } \sigma_{\text{DP}} \approx C \max \left\{ 1, \frac{\sqrt{T \log(1/\delta)}}{m \epsilon} \right\}$$

Grows with T !!

then DP-DSGD ensures (ϵ, δ) -Distributed DP

Learning Error Rate of DP-DSGD

At the completion of T iterations: $\mathbb{E} \left[\text{Loss} (\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{\sigma^2 + d\sigma_{\text{DP}}^2}{nT} \right)$

Some constant

For providing (ϵ, δ) -distributed DP: $\sigma_{\text{DP}} = kC \max \left\{ 1, \frac{\sqrt{T \log(1/\delta)}}{m \epsilon} \right\}$

$\mathbb{E} \left[\text{Loss} (\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{\sigma^2}{nT} + \frac{d \log(1/\delta)}{nm^2\epsilon^2} \right)$ Privacy-Accuracy trade-off

What About Robust DP-DSGD?

Global Phase: The server robustly aggregates averages the gradients $g_t^{(1)}, \dots, g_t^{(n)}$:

$$\hat{g}_t := \text{RobAvg} \left(g_t^{(1)}, \dots, g_t^{(n)} \right)$$

Without distributed momentum:

$$\mathbb{E} \left[\text{Loss} \left(\theta_T \right) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{d \log(1/\delta)}{nm^2\epsilon^2} + \kappa T \frac{d \log(1/\delta)}{m^2\epsilon^2} + \kappa \zeta \right)$$

Overhead due to robustness grows with T

Robust Averaging with Distributed Momentum

$$\mathbb{E} \left[\text{Loss} (\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{d \log(1/\delta)}{m^2 \epsilon^2} \left(\frac{1}{n} + \kappa \right) + \kappa \zeta \right)$$

At best $\kappa \in \mathcal{O} (f/n)$

No dimension

Lower Bound:

$$\text{Loss} (\theta_T) - \text{Loss}^* \in \Omega \left(\frac{d \log(1/\delta)}{nm^2 \epsilon^2} + \frac{f}{n} \cdot \frac{\log(1/\delta)}{m^2 \epsilon^2} + \frac{f}{n} \zeta \right)$$

Dimension Can Be Big!

ML Model	Common Applications	Number of Parameters (d)
Resnet-50	Vision, regression	23×10^6
Resnet-18	Vision, regression	18×10^6
GPT-2	NLP, vision	1.5×10^9
GPT-3	NLP, vision	175×10^9

(f, κ) –Spectral Robust Averaging

Aggregation $F : (v_1, \dots, v_n) \mapsto \hat{v}$ for all $G \subseteq [n]$ with $|G| = n - f$,

(f, κ) -Robust averaging

$$\frac{1}{n-f} \sum_{i \in G} \|v_i - \bar{v}_G\|^2$$

$$\|\hat{v} - \bar{v}_G\|^2 \leq \kappa \lambda_{\max} \left(\frac{1}{n-f} \sum_{i \in G} (v_i - \bar{v}_S) (v_i - \bar{v}_G)^\top \right)$$

Spectral norm

No dimension!

$$\mathbb{E} \left[\text{Loss}(\theta_T) - \text{Loss}^* \right] \in \mathcal{O} \left(\frac{d \log(1/\delta)}{nm^2\epsilon^2} + \kappa \cdot \frac{\log(1/\delta)}{m^2\epsilon^2} + \kappa\zeta \right)$$

Matches LB if
 $\kappa \in \mathcal{O} \left(\frac{f}{n} \right)$

SMEA: Optimal Spectral Robust Averaging

Smallest Maximum Eigenvalue Averaging

$$G^* \in \arg \min_{G \subseteq [n], |G|=n-f} \lambda_{\max} \left(\frac{1}{n-f} \sum_{i \in G} (v_i - \bar{v}_G) (v_i - \bar{v}_G)^\top \right)$$

$$F(v_1, \dots, v_n) := \bar{v}_{G^*}$$

SMEA is (f, κ) -Spectral Robust with $\kappa \in \mathcal{O}\left(\frac{f}{n}\right)$

“On the Privacy-Robustness-Utility Trilemma in Distributed Learning”. AGGPS. ICML 2023

Take Away III: Robustness with Privacy

Augments local
randomness isotropically

Data Privacy

Nodes do not share exact
information on their data



Spectral Averaging
Local Momentum

“Differential Privacy and Byzantine Resilience in SGD: Do They Add Up?”. **GGPRS**. PODC 2022

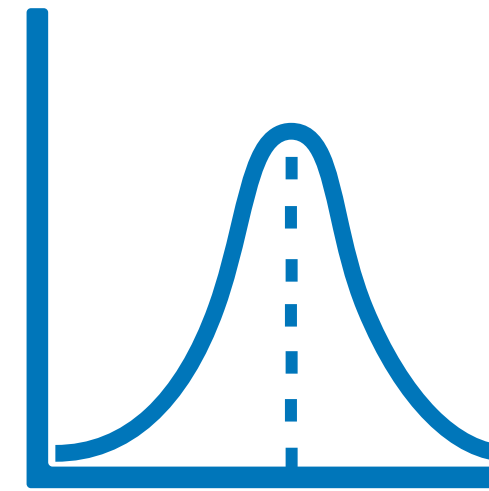
“On the Privacy-Robustness-Utility Trilemma in Distributed Learning”. **AGGPS**. ICML 2023

Take Aways: Three Challenges in Robust DL

State machine replication
is not efficient

Local Randomness

Even honest machines
compute noisy updates

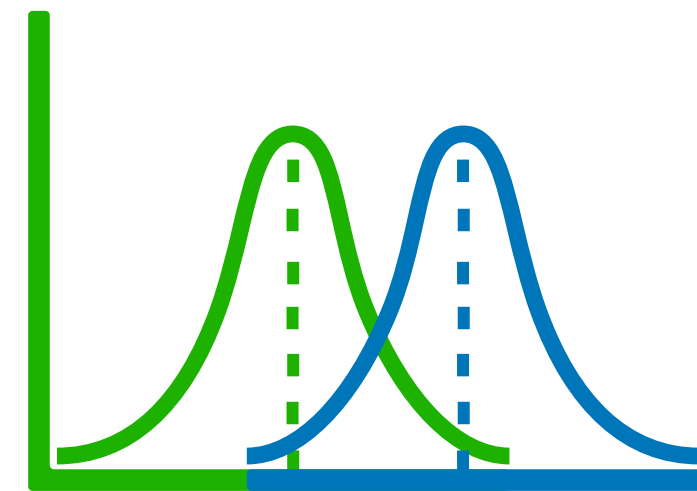


Robust Averaging
Local Momentum

Adversarial nodes can
camouflage as honest nodes

Data Heterogeneity

Different nodes sample data
from different distributions

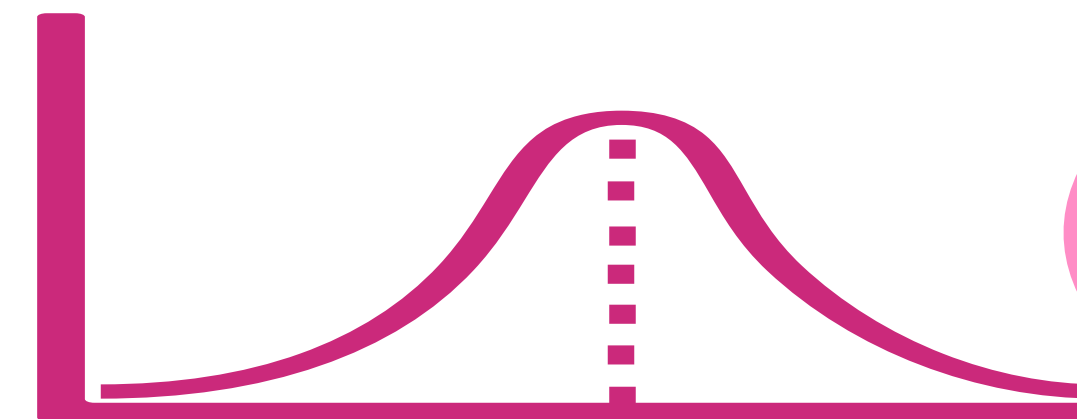


Nearest Neighbor Mixing &
Adaptive Robust Clipping

Augments local
randomness isotropically

Data Privacy

Nodes do not share exact
information on their data



Spectral Averaging
Local Momentum

Bibliography

Fault-Tolerance in Distributed Optimization: The Case of Redundancy

N. Gupta and N. H. Vaidya. In *Principles of Distributed Computing*. ACM, 2020.

Differential Privacy and Byzantine Resilience in SGD: Do They Add Up?

R. Guerraoui, N. Gupta, R. Pinot, S. Rouault and J. Stephan. In *Principles of Distributed Computing (PODC)*. ACM, 2021.

Presented in the talk

Approximate Byzantine Fault-Tolerance in Distributed Optimization

S. Liu, N. Gupta and N. H. Vaidya. In *Principles of Distributed Computing*. ACM, 2021.

Byzantine Machine Learning Made Easy by Resilient Averaging of Momentums

S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *International Conference on Machine Learning (ICML)*. PMLR, 2022.

Fixing by Mixing: A Recipe for Optimal Byzantine ML under Heterogeneity

Y. Allouah, S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*. PMLR, 2023.

Robust Collaborative Learning with Linear Gradient Overhead

S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *International Conference on Machine Learning (ICML)*. PMLR, 2023.

On the Privacy-Robustness-Utility Trilemma in Distributed Learning

Y. Allouah, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. In *International Conference on Machine Learning (ICML)*. PMLR, 2023.

Byzantine Machine Learning: A Primer

R. Guerraoui, N. Gupta and R. Pinot. *ACM Computing Surveys*, 2023.

Robust Distributed Learning: Tight Error Bounds and Breakdown Point under Data Heterogeneity

Y. Allouah, R. Guerraoui, N. Gupta, R. Pinot, and G.Rizk. In *Advances in Neural Information Processing Systems (NeurIPS) (Spotlight)*, 2023.

Robust Machine Learning: Distributed Methods for Safe AI

R. Guerraoui, N. Gupta and R. Pinot. *Springer Nature Publishing Company*, 2023. **[BOOK]**

Open Problems

***Generalizability* of Robust DSGD**

Influence of robustness on accuracy for unseen points?

Boudou, Thomas, Le Bars, Batiste, Gupta, N., & Bellet, Aurélien. (2025).

Generalization under Byzantine & Poisoning Attacks: Tight Stability Bounds in Robust Distributed Learning. *arXiv e-prints*, arXiv-2506.

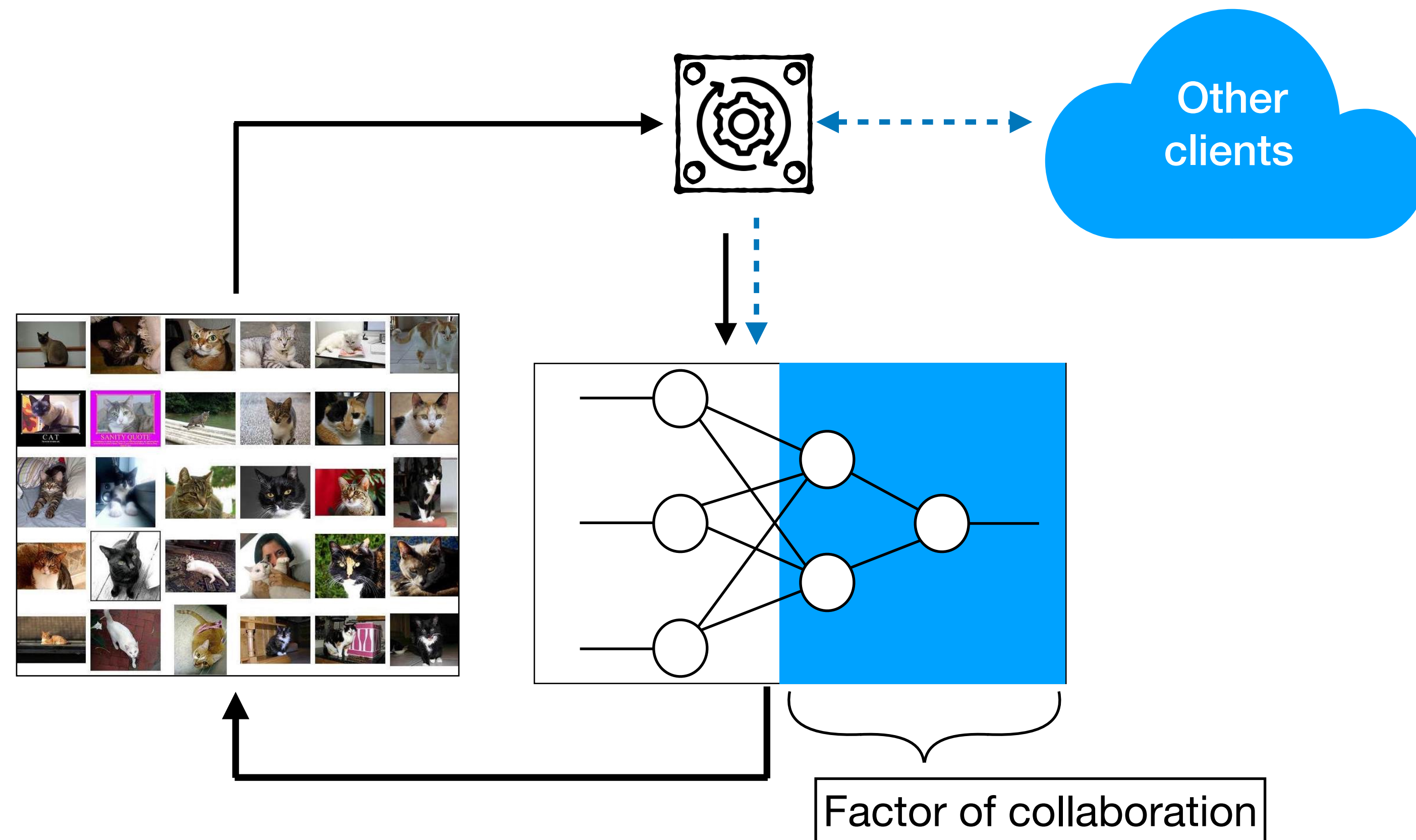
Connection to *adversarial examples*

Can Robust SGD defend against evasion attacks?

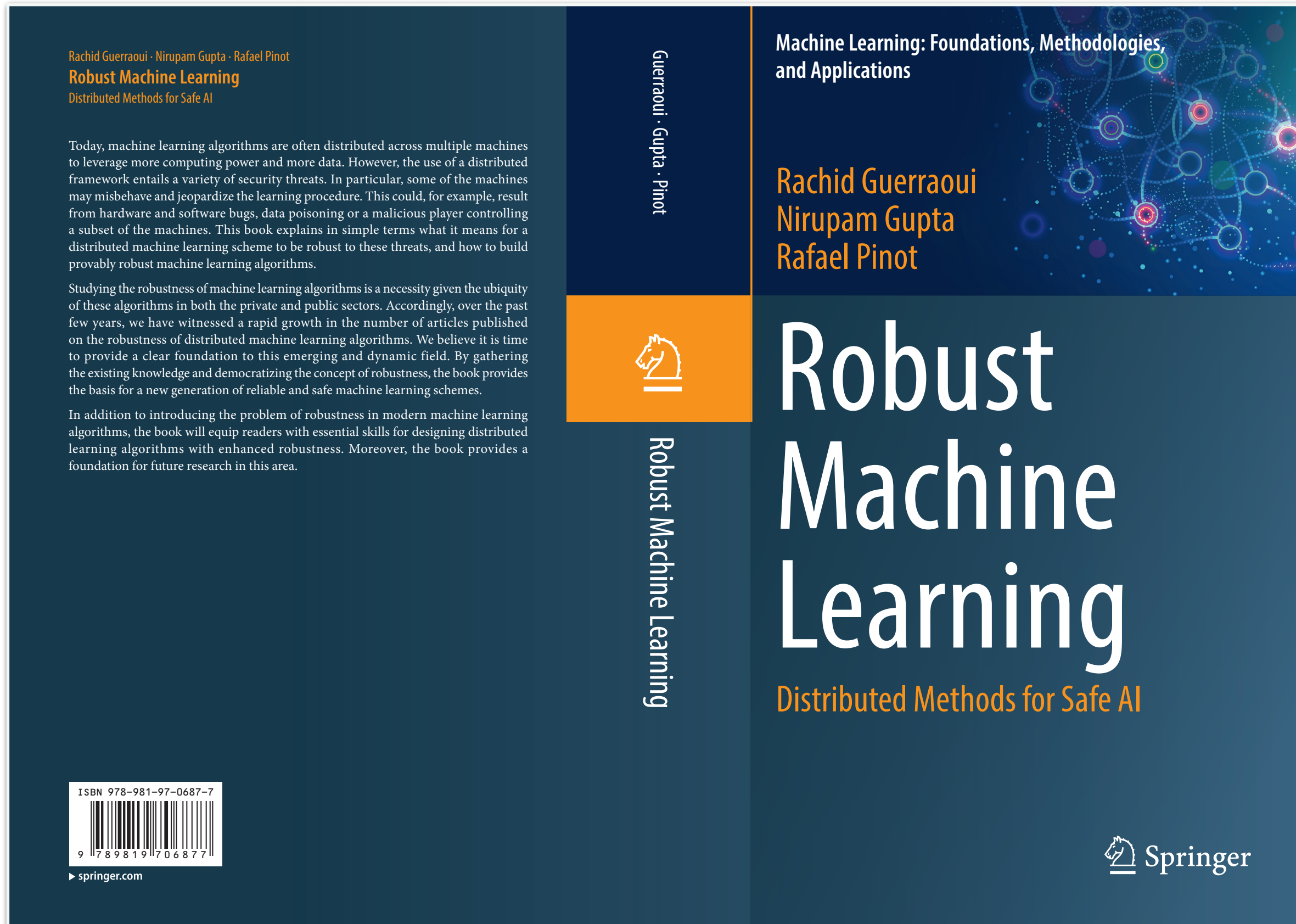
***Game-theoretic* approach**

Tight adversary-aware robustness guarantees

Future Avenue: Robustness with Personalization



Allouah, Youssef, **Abdellah El Mrini**, Rachid Guerraoui, N. Gupta, and Rafael Pinot.
"Fine-tuning personalization in federated learning to mitigate adversarial clients." *NeurIPS*, 2024.



Thank you!

nigu@di.ku.dk

BYZFL