

Optimization for Machine Learning in Practice II

Martin Jaggi

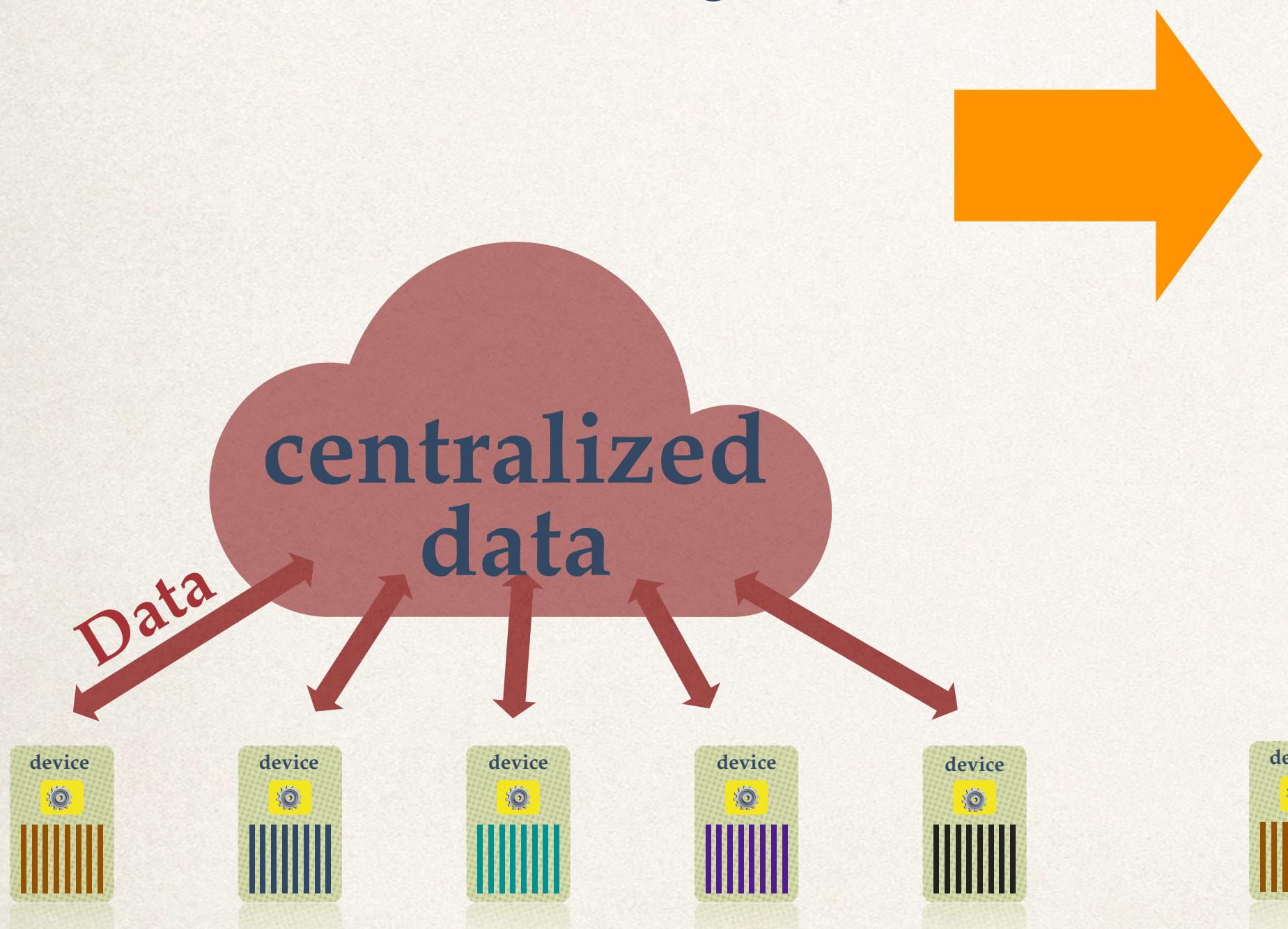


Machine Learning and Optimization Laboratory
mlo.epfl.ch

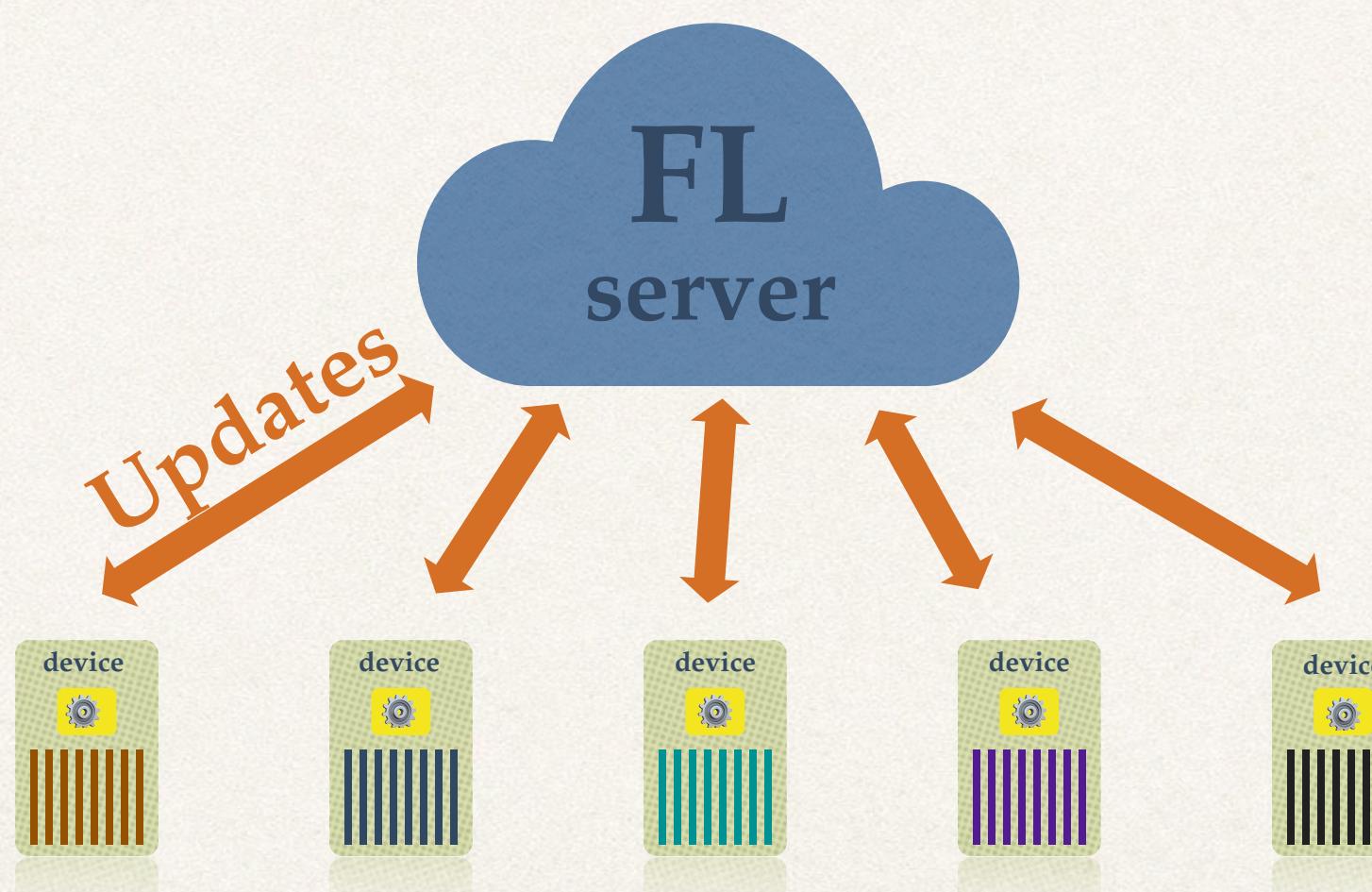
Collaborative Learning

Evolution

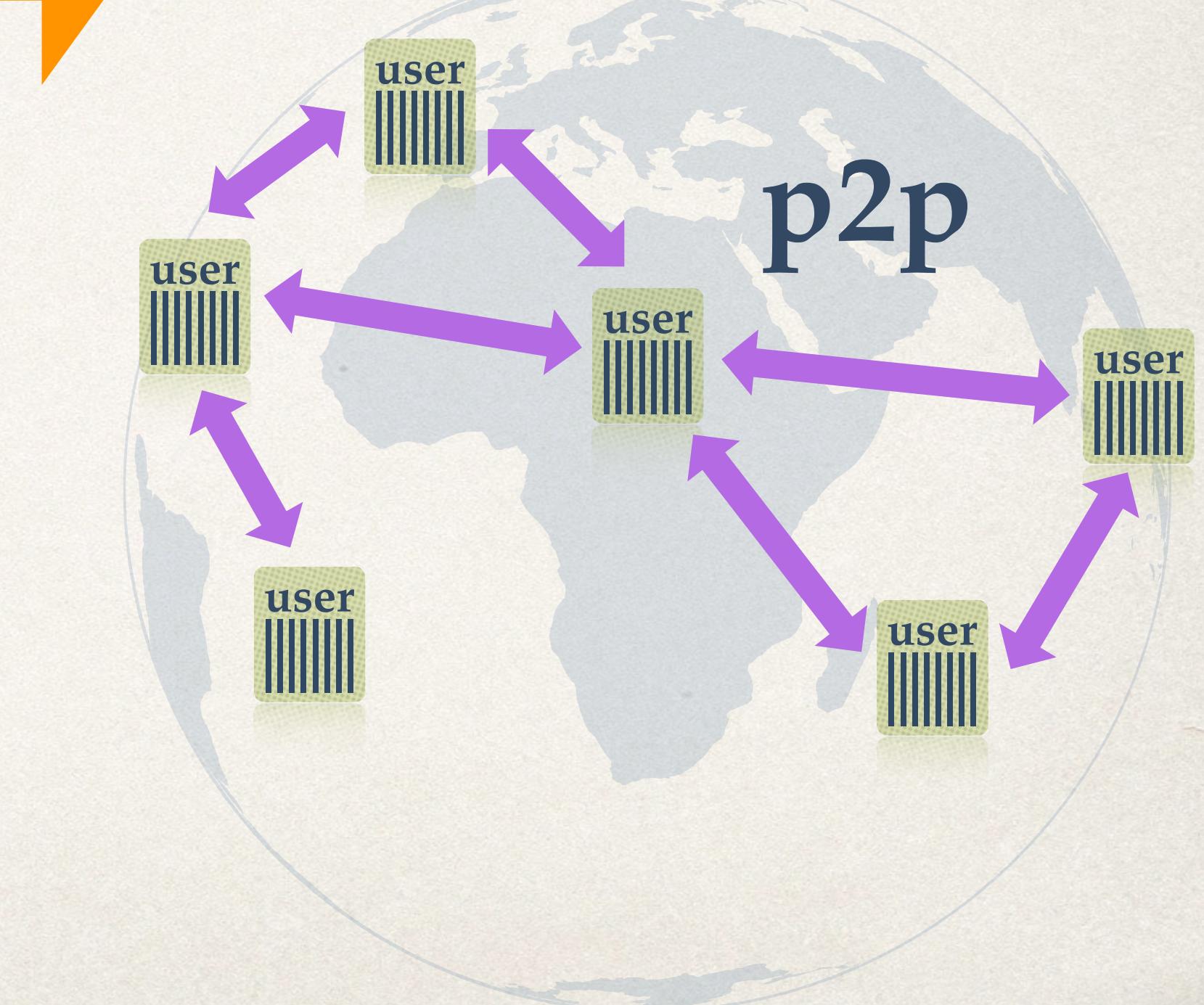
centralized
traditional, sharing data



federated
sharing model updates

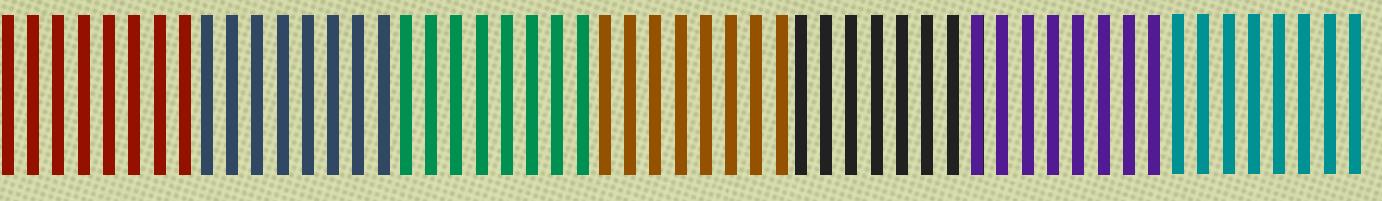


decentralized
collaborative learning

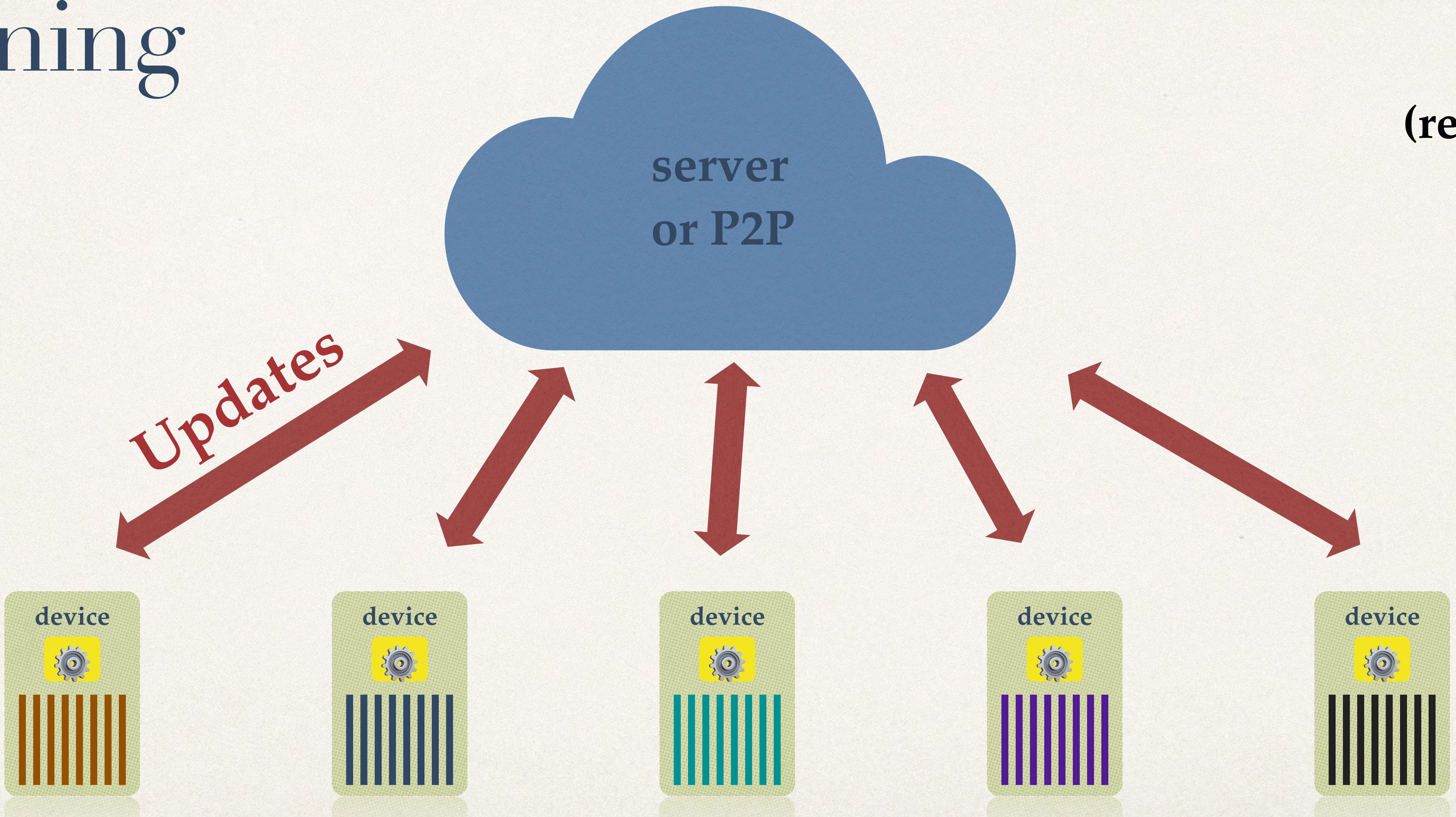


Collaborative & Federated Training

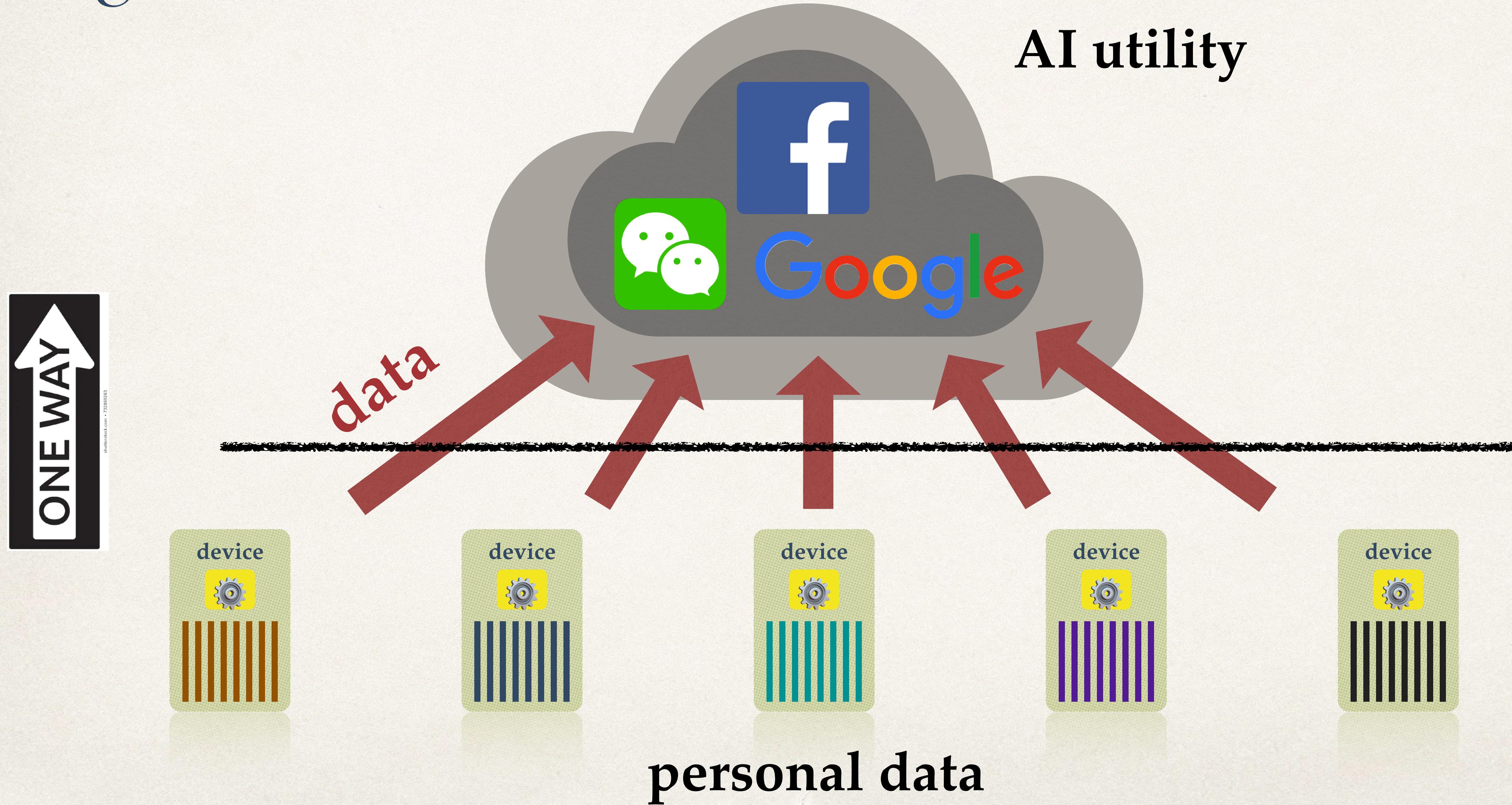
Data



(recap)

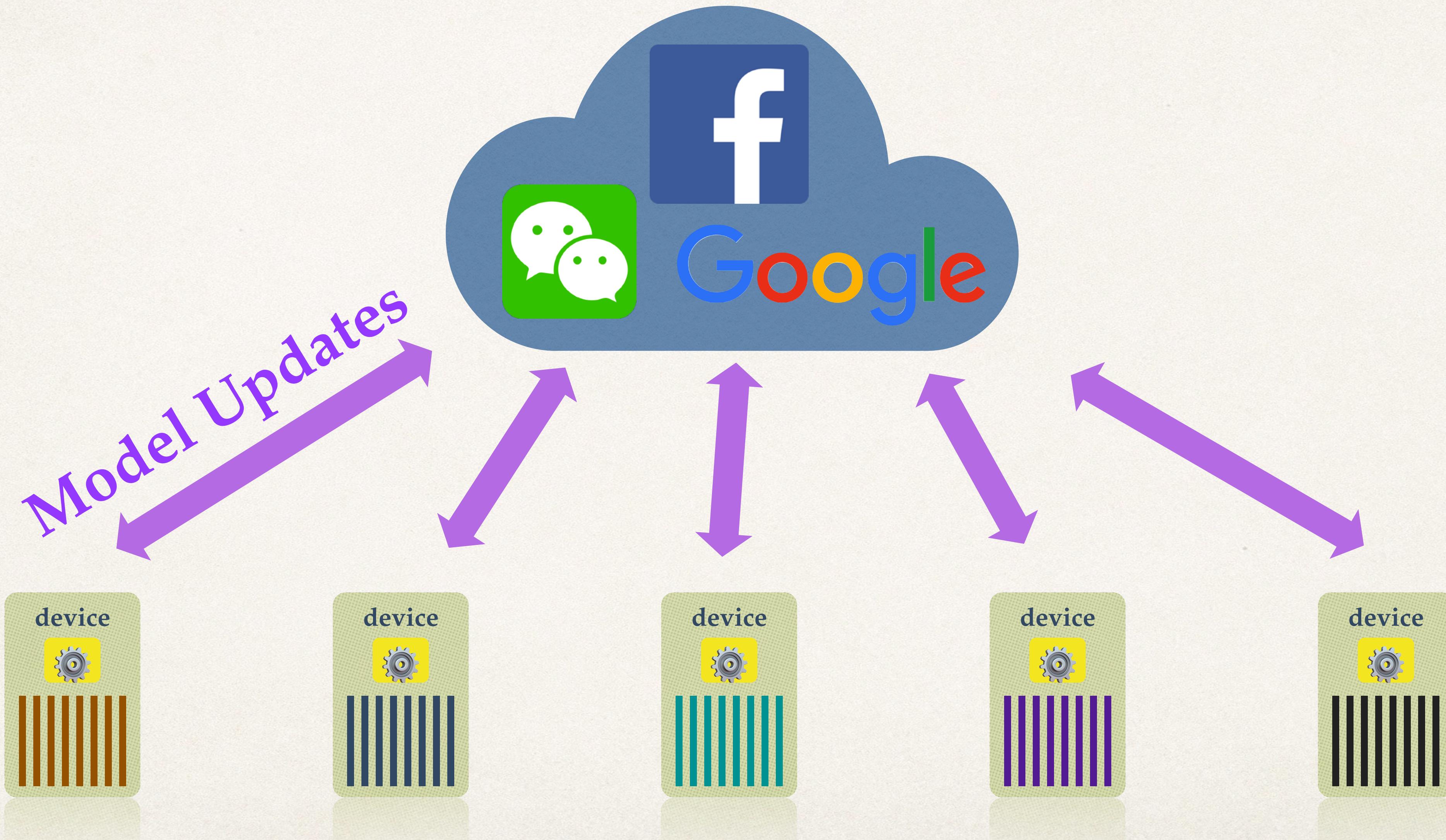


Big Picture



2a

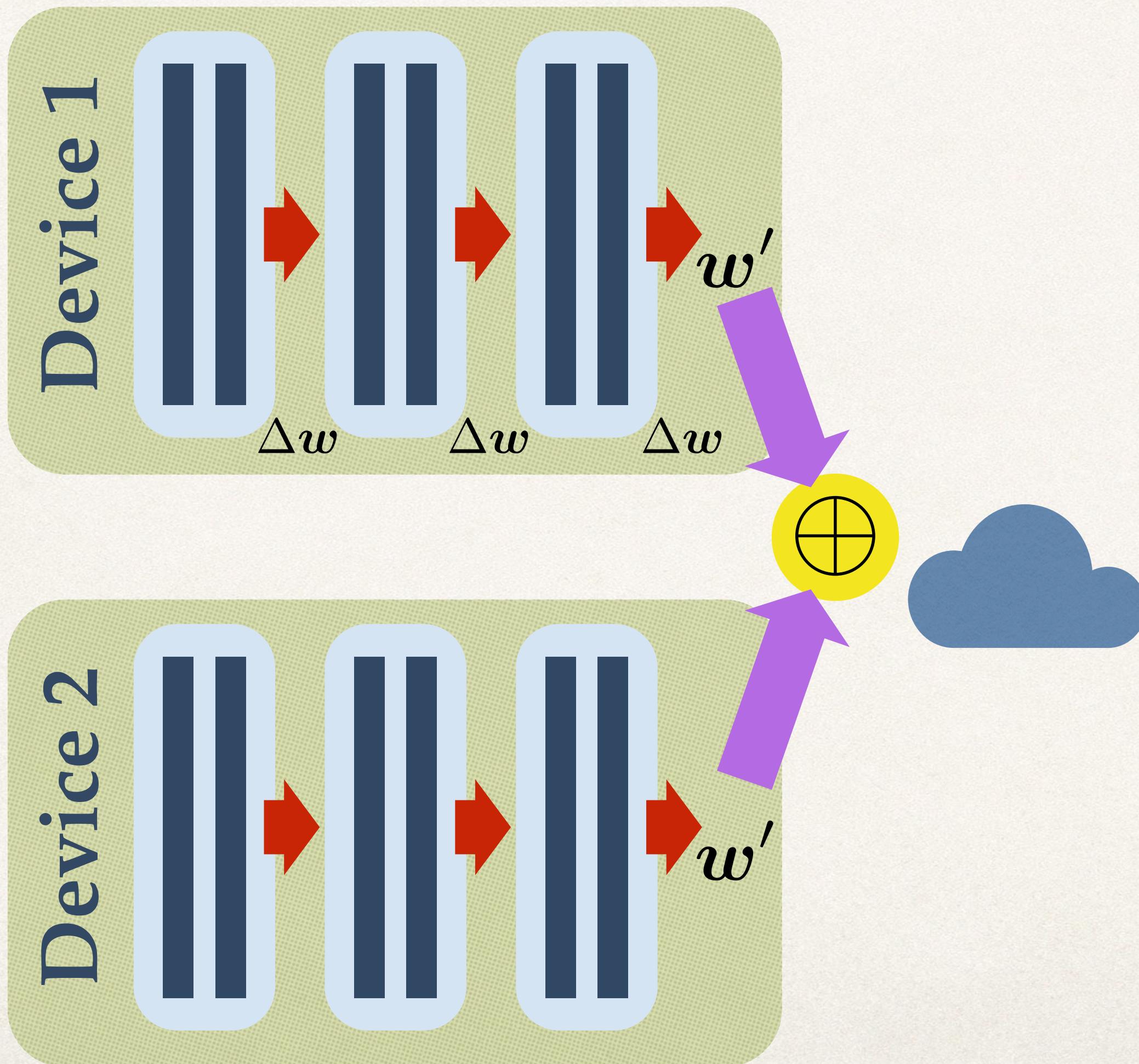
Federated Learning



data never leaves device

2a

Federated Learning



- ✿ Local SGD steps = “Federated averaging”
- ✿ Google Android Keyboard

Client drift

- * Federated Learning

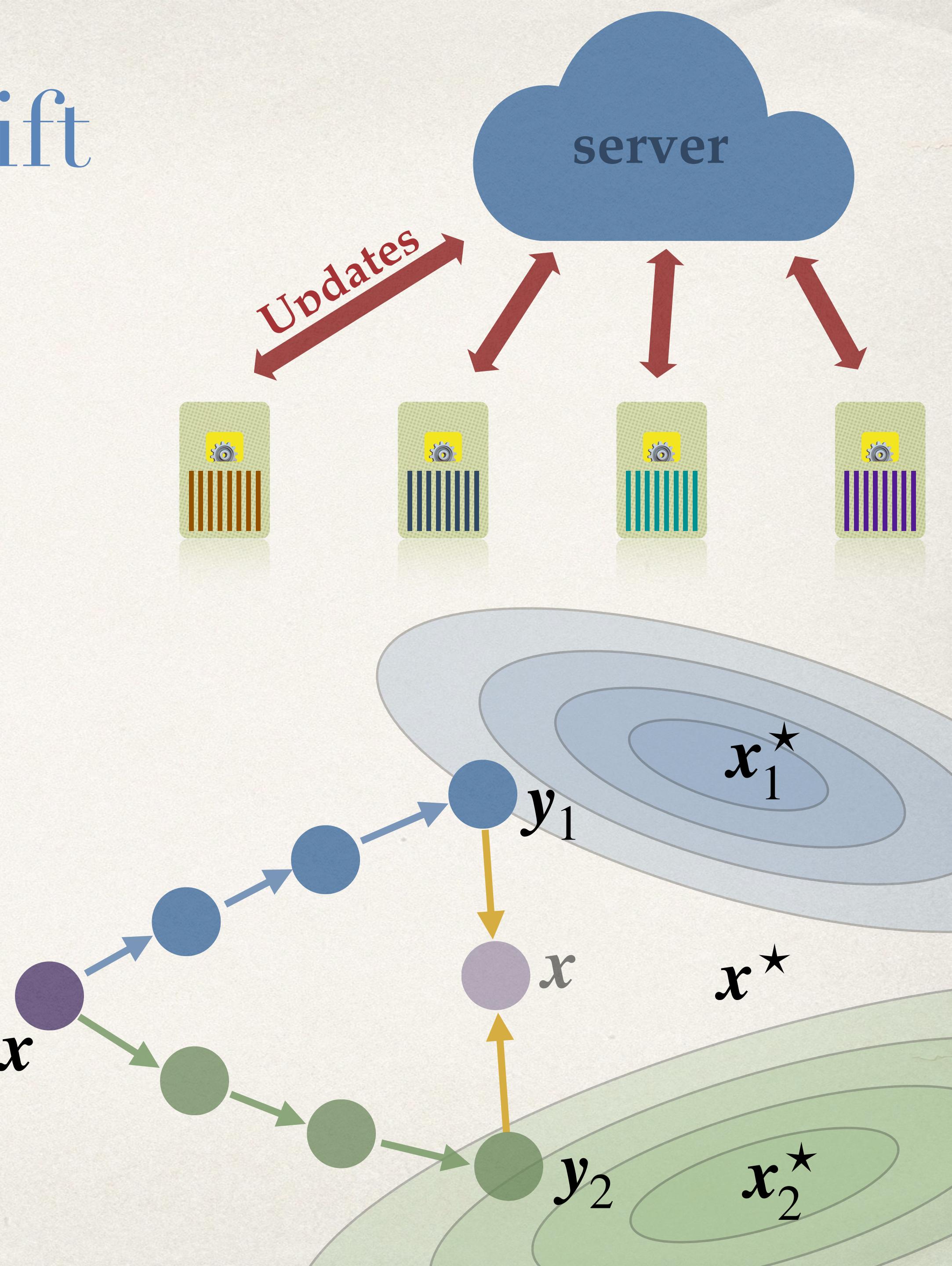
$$\min_{\mathbf{x}} \frac{1}{n} \sum_i^n f_i(\mathbf{x})$$

- * Fed Avg / Local SGD

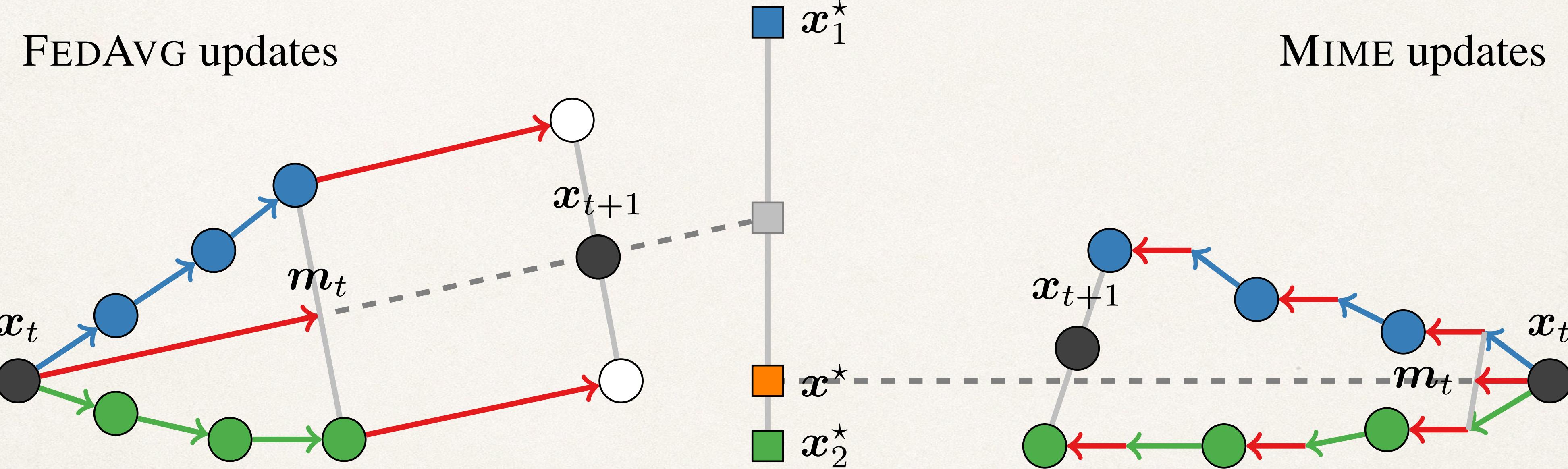
for some local steps

$$y_i := y_i - \eta \nabla f_i(y_i)$$

$$\mathbf{x} := \frac{1}{n} \sum_{i=1}^n y_i \quad (\text{aggregation})$$



Client drift



Mime algorithm framework

for some local steps

$$y_i := y_i - \eta \left((1 - \beta) \nabla f_i(y_i) + \beta \mathbf{m} \right)$$

$$\mathbf{m} := (1 - \beta) \nabla f_i(\mathbf{x}) + \beta \mathbf{m}$$



*aggregated on server
after each round*

Federated vs Personalized Learning

- **Federated**

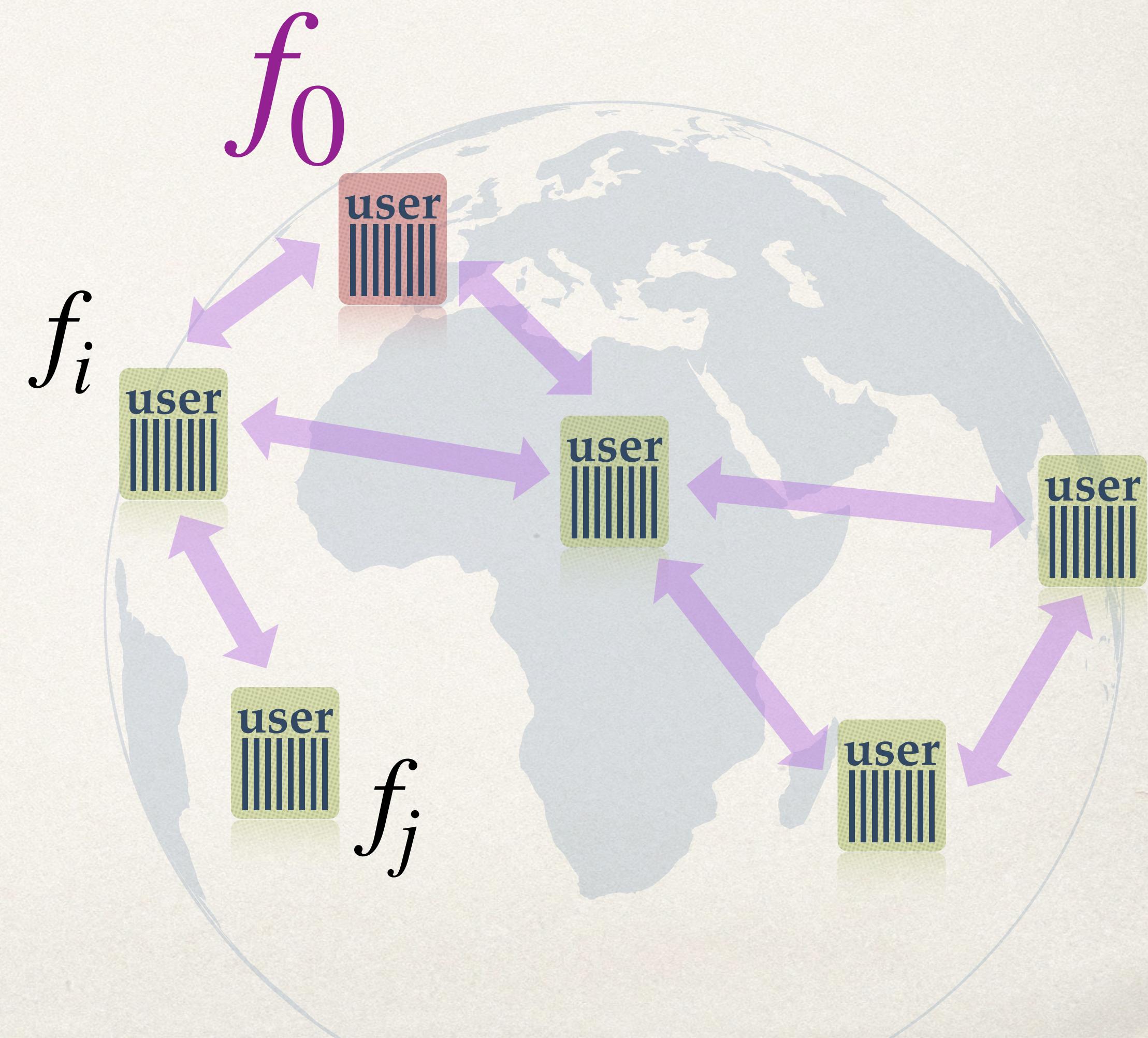
$$\min_x \frac{1}{n} \sum_i f_i(x)$$

- **Collaborative / Personalized**

$$\min_x f_0(x)$$

$$\min_x f_1(x)$$

$$\min_x f_n(x)$$



Federated vs Personalized Learning

- ❖ **Federated**

$$\min_{\mathbf{x}} \frac{1}{n} \sum_i^n f_i(\mathbf{x})$$

- ❖ **Collaborative / Personalized**

$$\min_{\mathbf{x}} f_0(\mathbf{x})$$

$$\min_{\mathbf{x}} f_1(\mathbf{x})$$

$$\min_{\mathbf{x}} f_n(\mathbf{x})$$

- ❖ **Ordering of training**

Set of active clients evolves (how?)

- ❖ **Clients = Tasks**

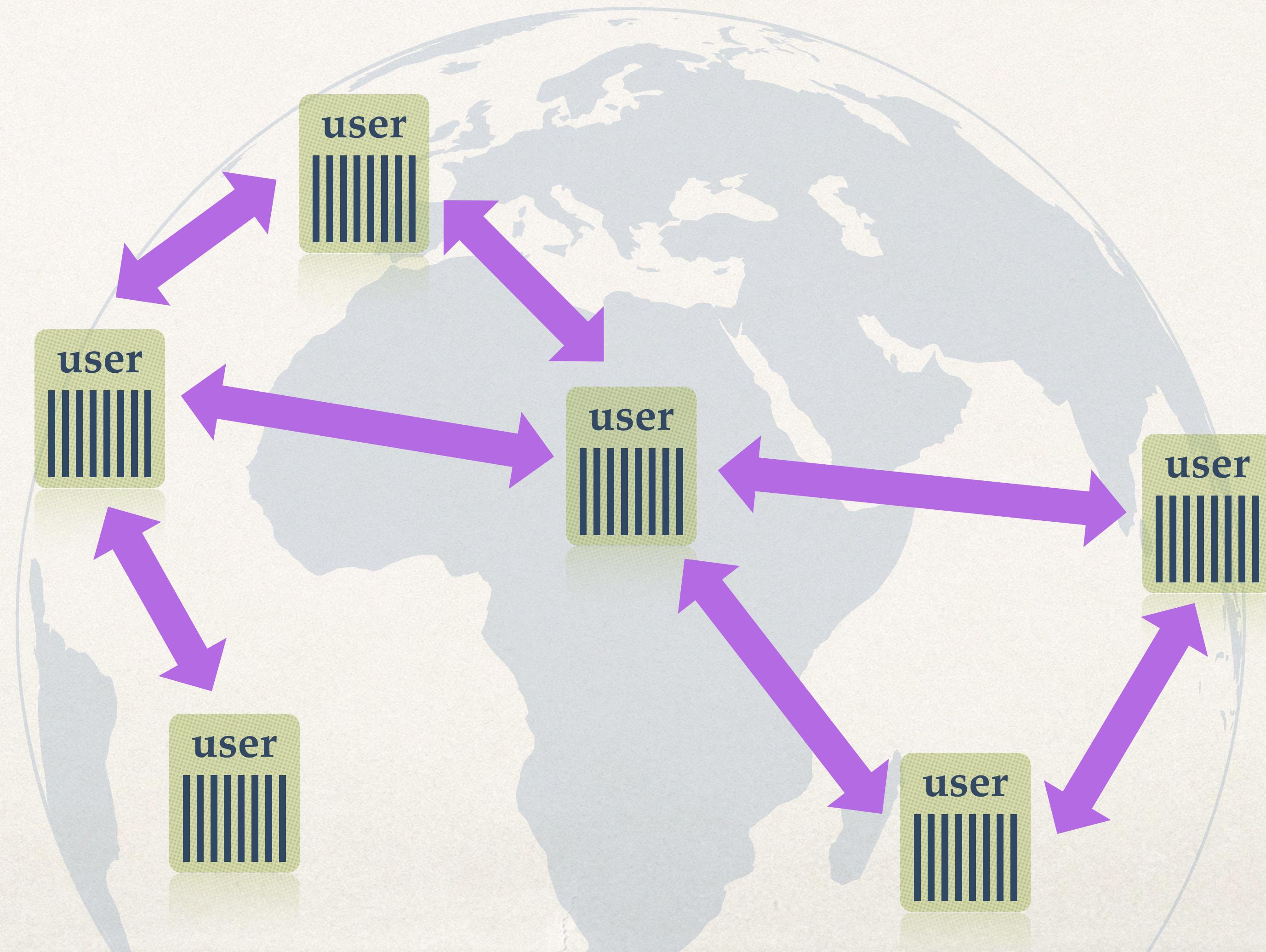
Sequential fine-tuning

Transfer learning,
overparameterized models?

- ❖ **Train alone or collaborate?**

2c

Decentralized Learning

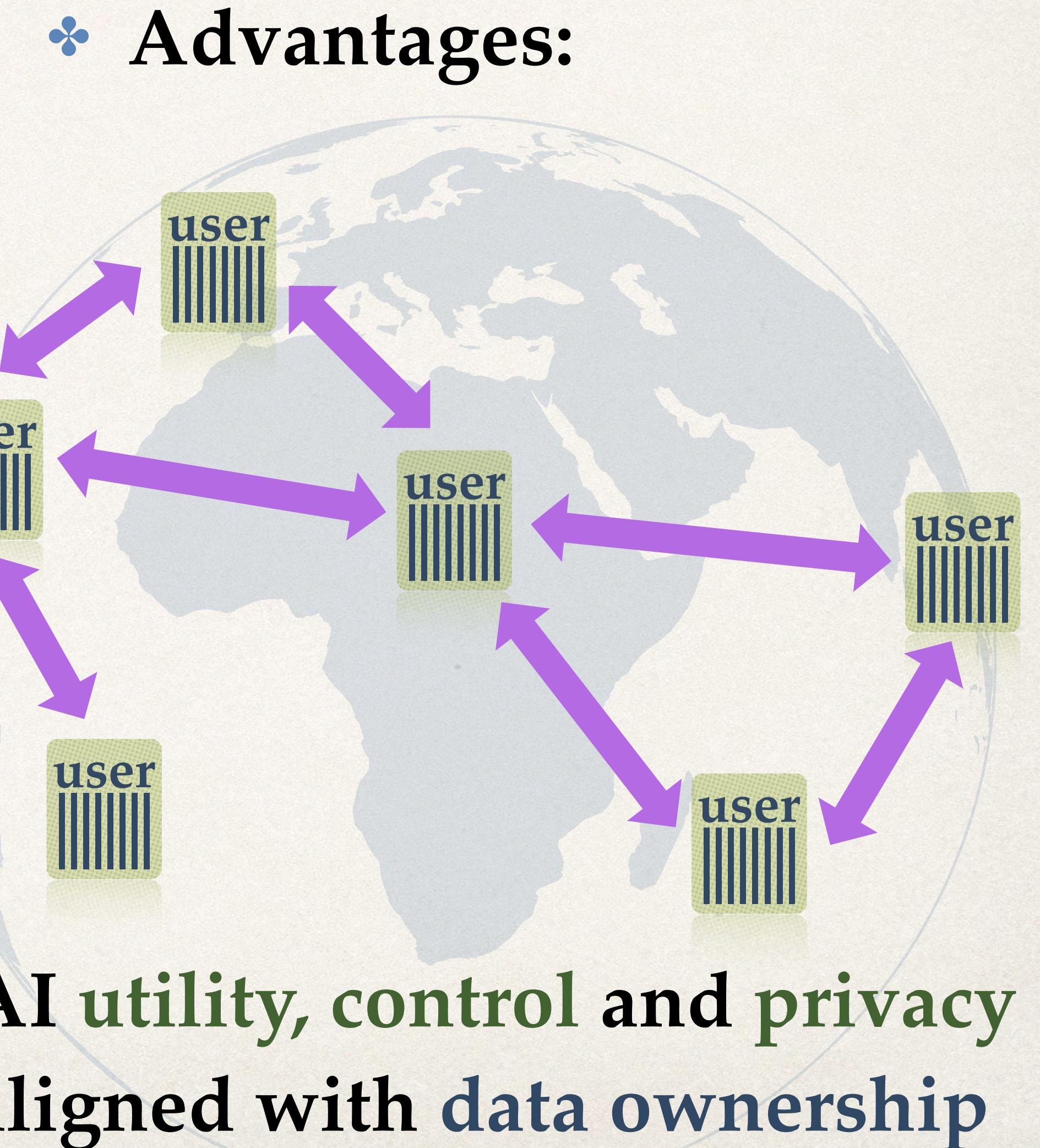


Motivation

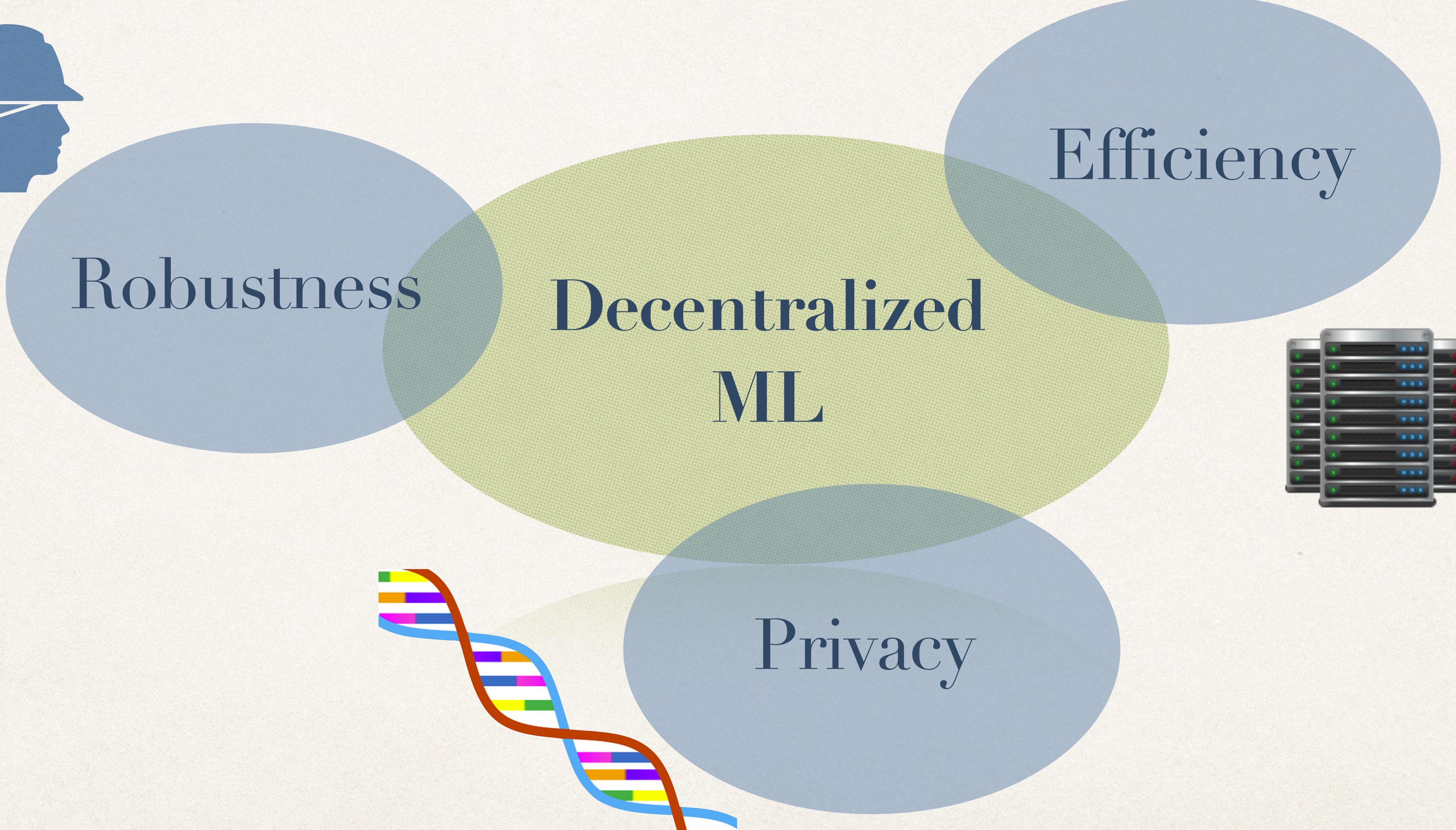
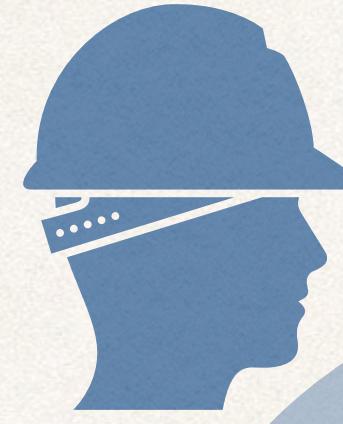
- ❖ **Applications:**
any ML system with user data
servers, devices, sensors, hospitals, ...



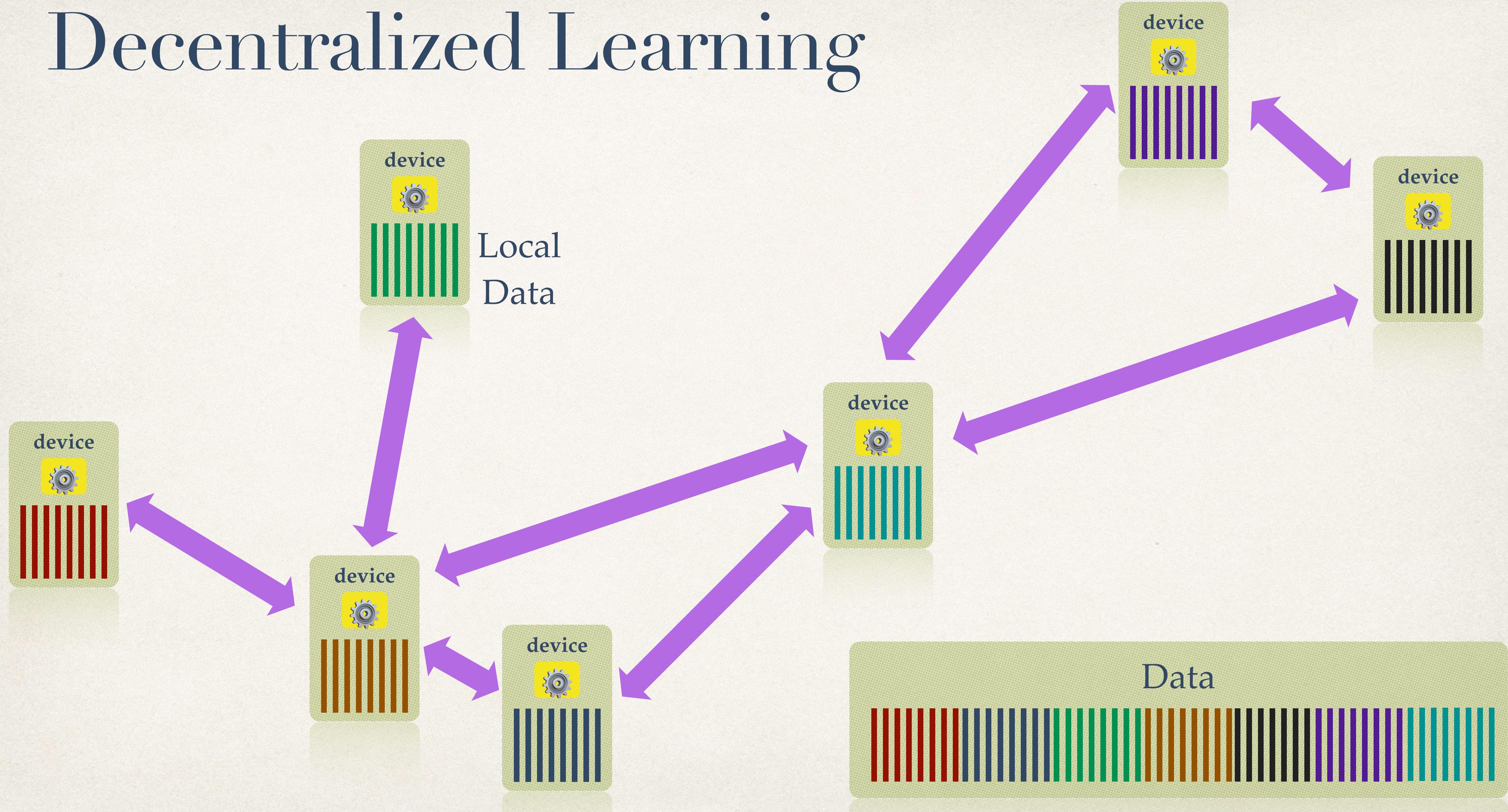
[image source](#)



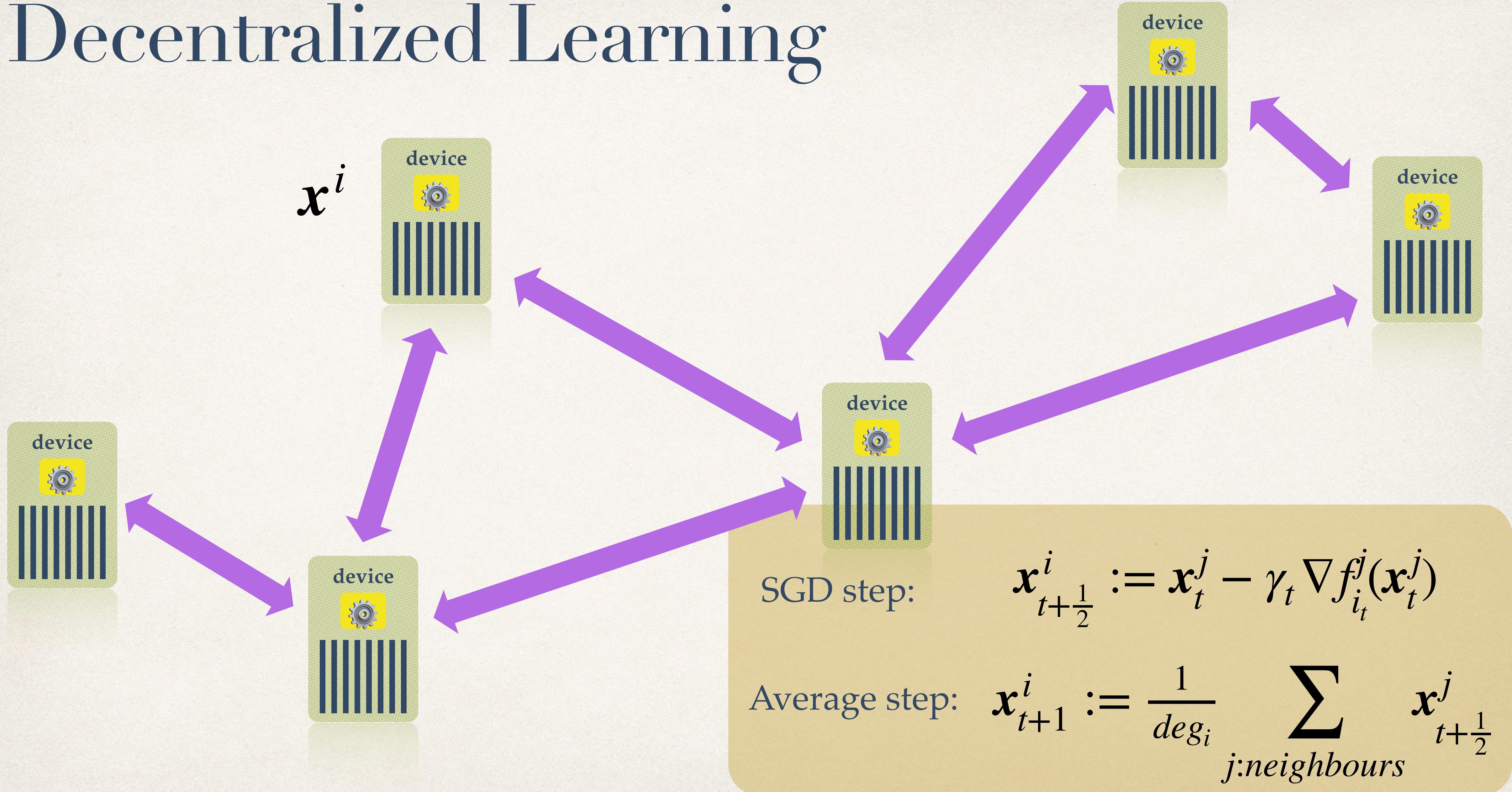
Required Building Blocks



Decentralized Learning



Decentralized Learning



Communication Compression

- ✿ limited-bit precision vector

e.g. 1-bit per entry reduces communication 32 times

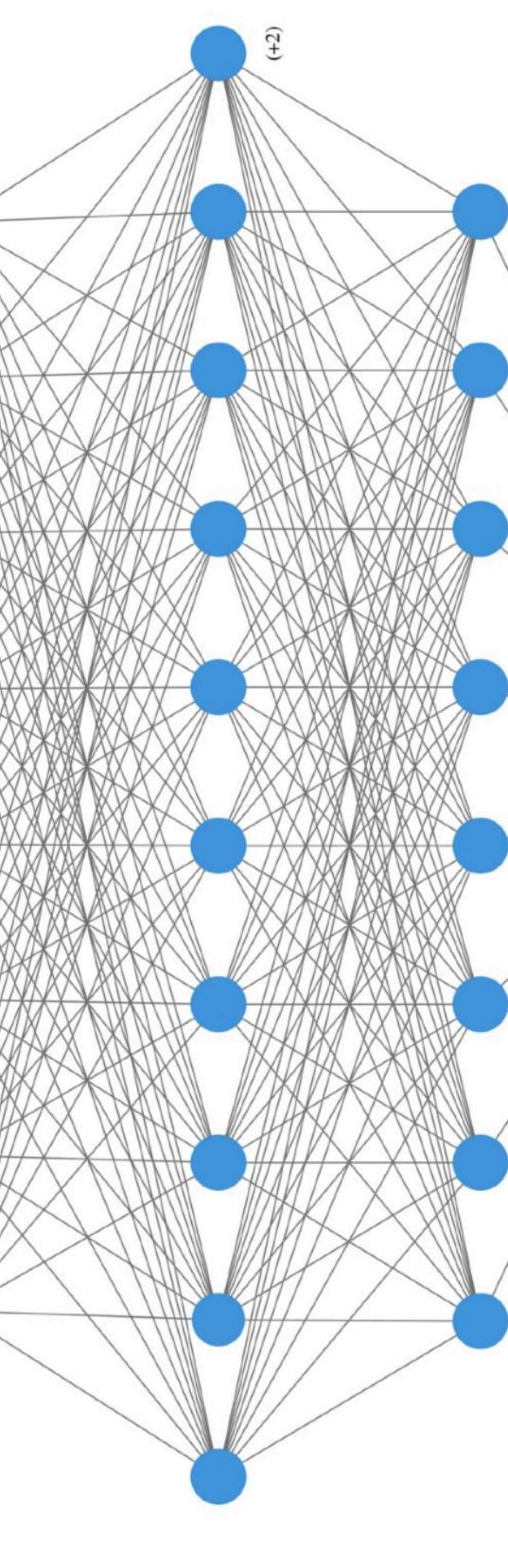
- ✿ random / top k% of all the entries

e.g. k=0.1% reduces communication 1000 times

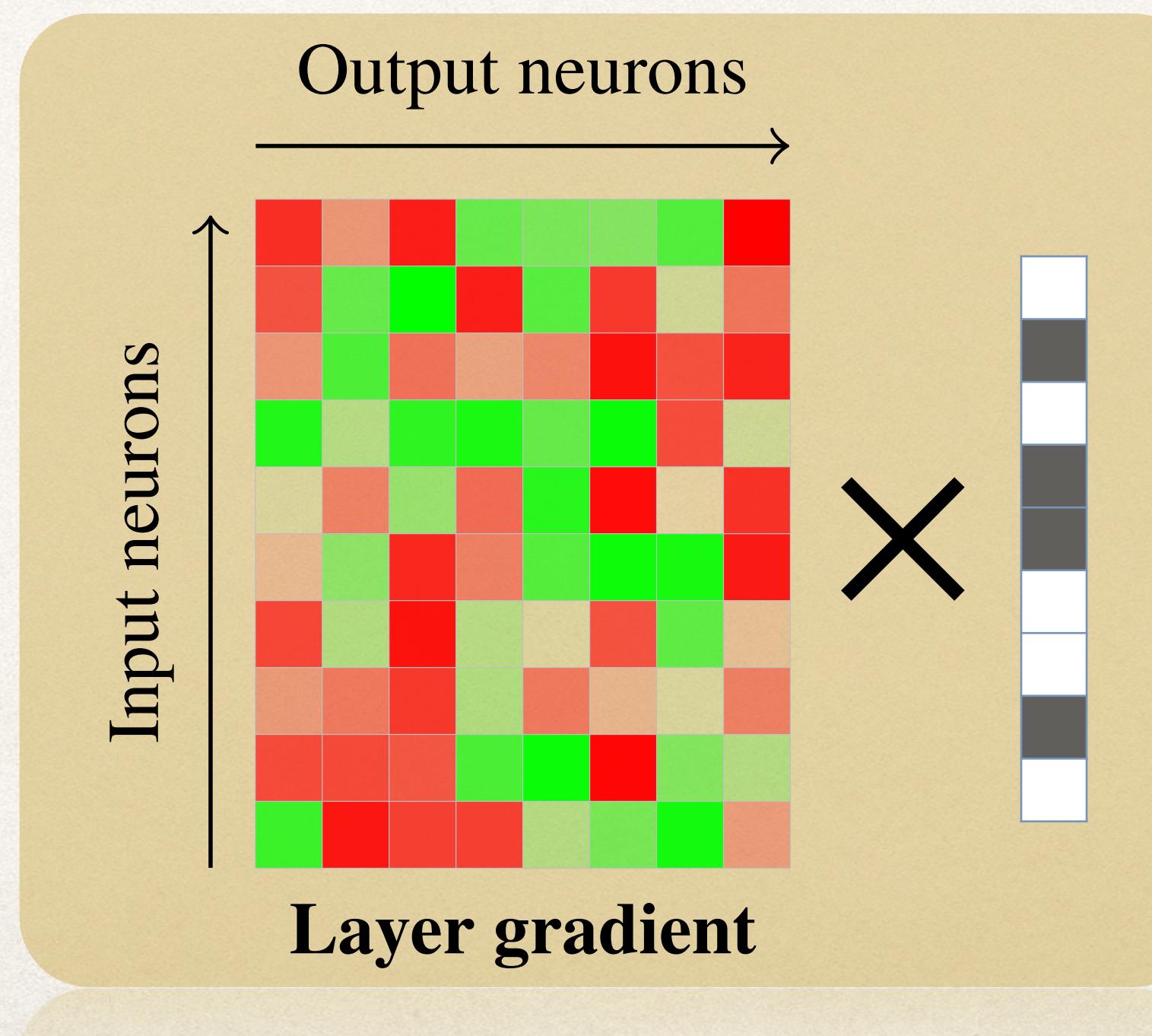
- ✿ low rank version of the gradient?

Low-Rank Communication Compression

- PowerSGD



backprop is fast:
linear time



fast compression?

Fast power iterations

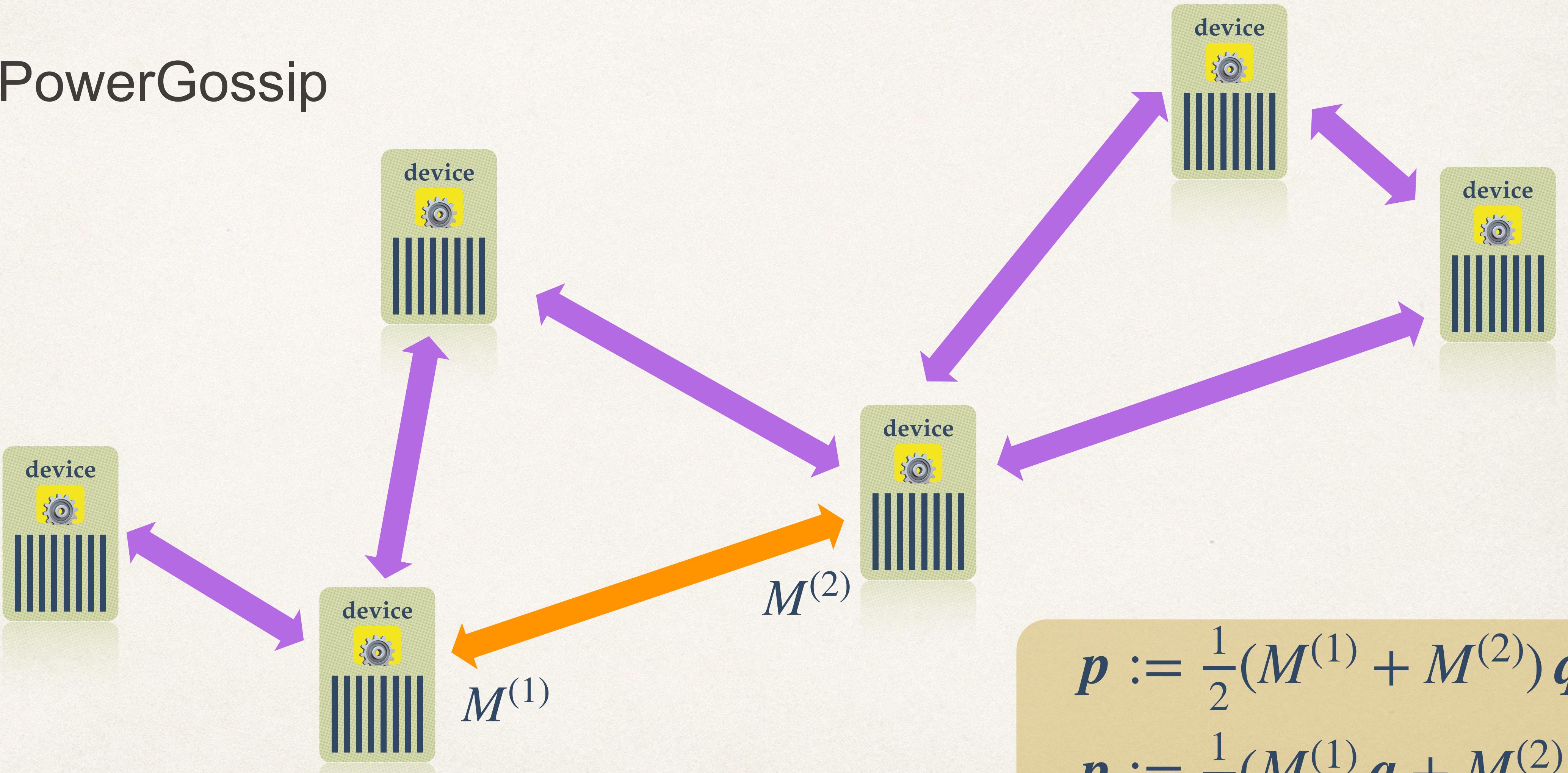
$$p := G q$$

$$q := G^\top p$$

$$\mathcal{C}(G) = pq^\top$$

Decentralized Learning with Compression

- PowerGossip



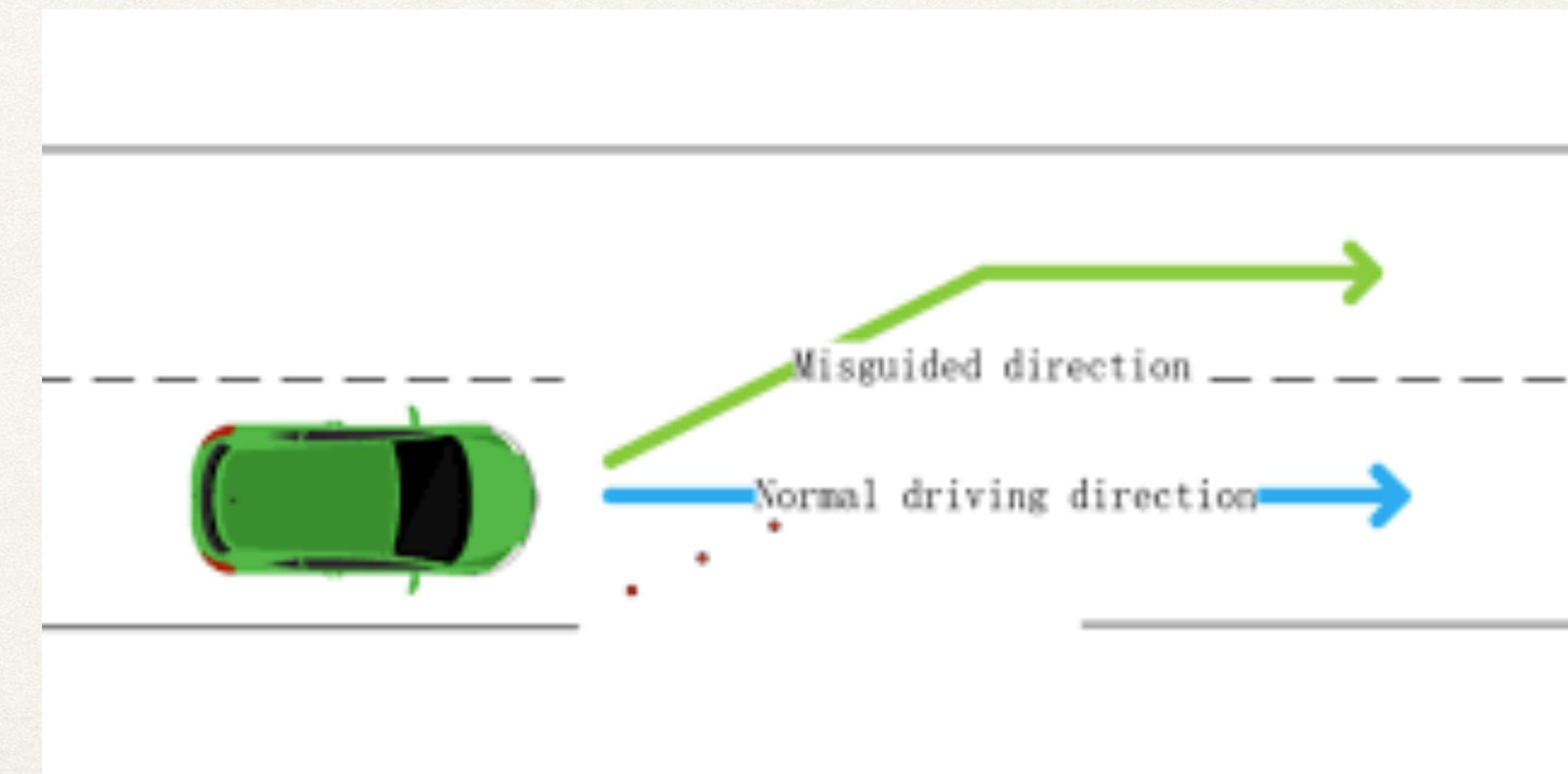
Building Blocks for Decentralized ML

- ✿ **Efficiency: Communication & Compute**
on-device learning, Edge AI
peer-to-peer communication
- ✿ **Privacy**
data locality, leakage?, attacks?
- ✿ **Robustness & Incentives**
tolerate bad players, reward collaboration

3

Robustness

During Training and Inference

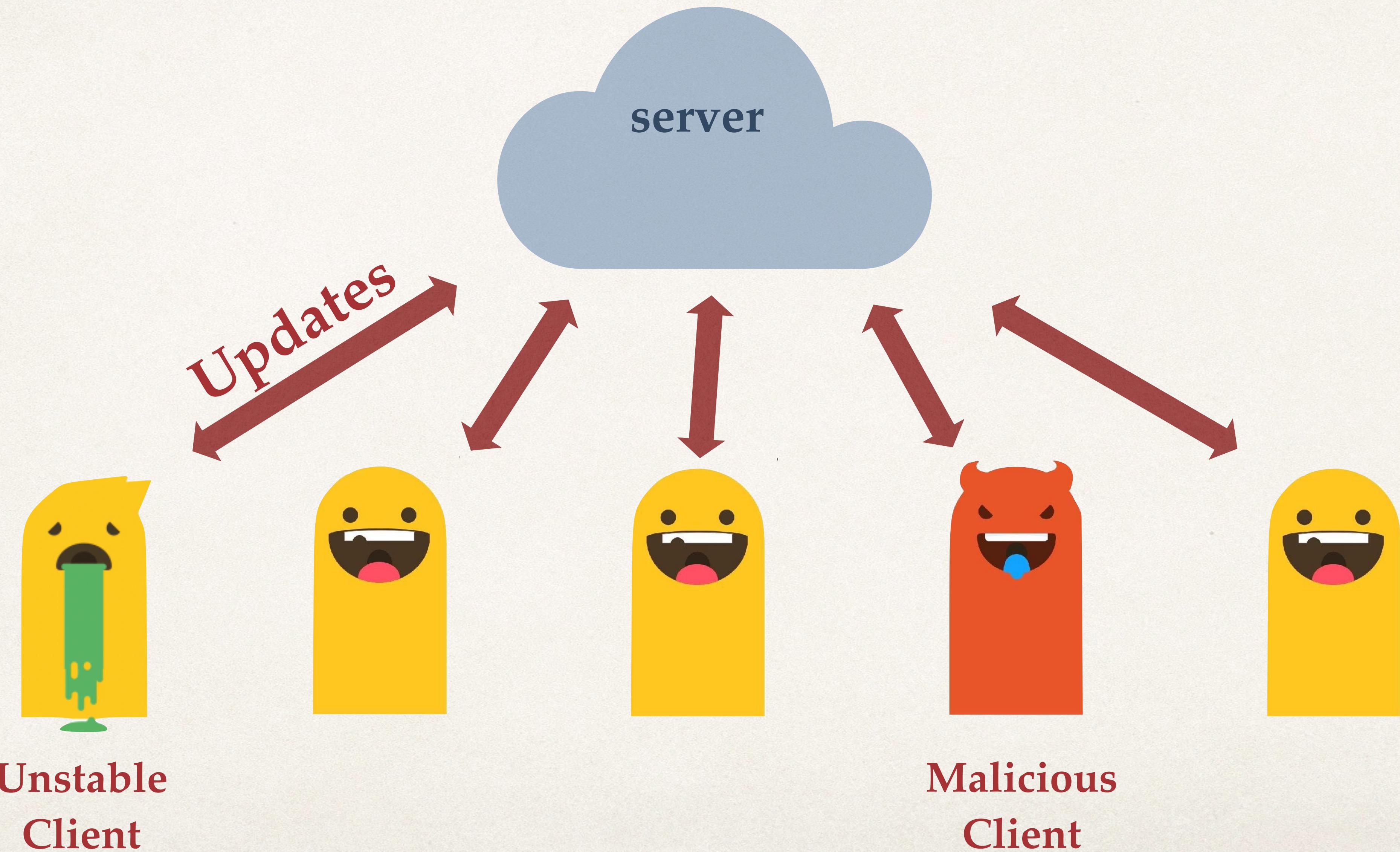


3a

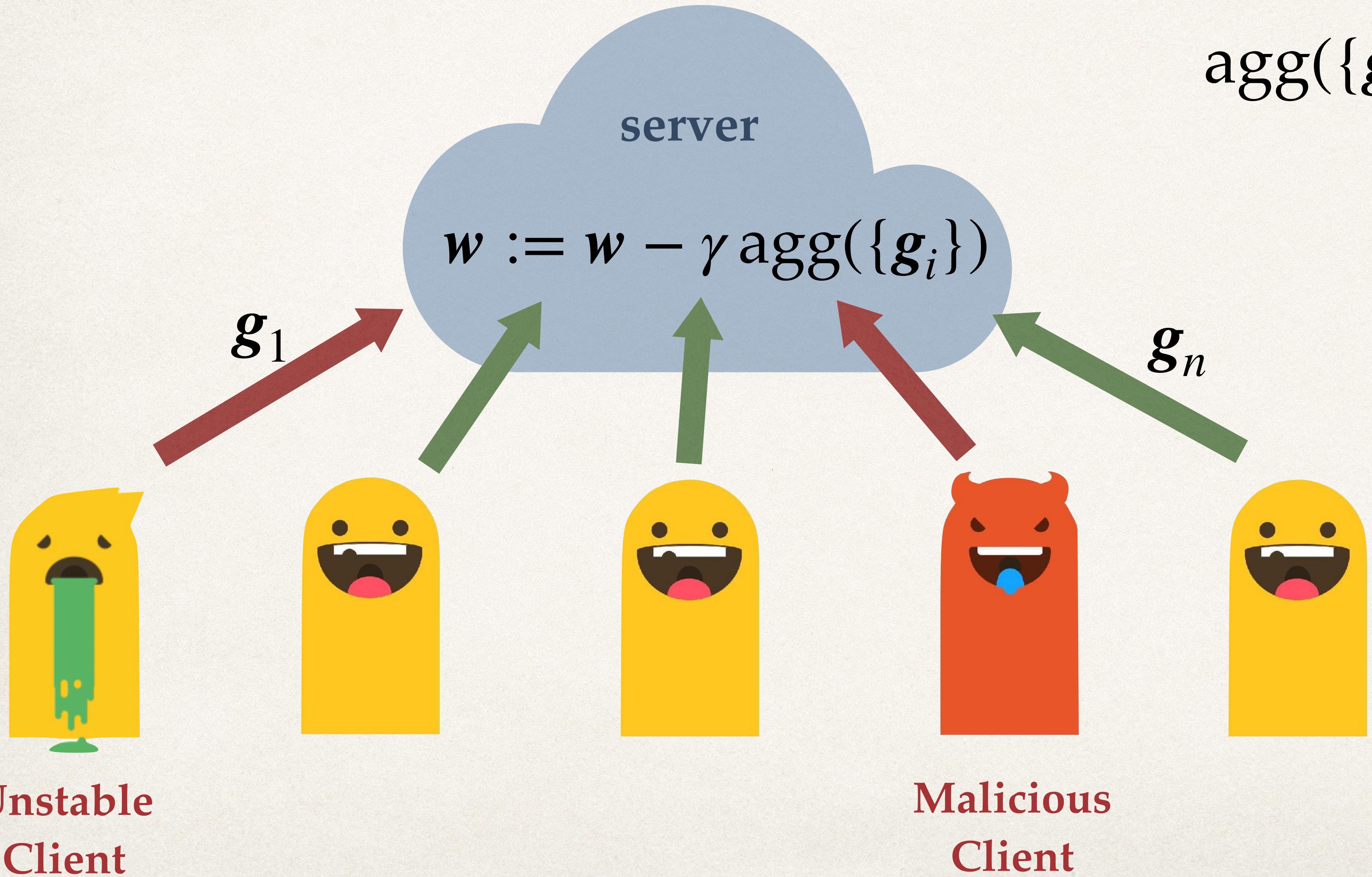
Gradients from
faulty/malicious collaborators:

- Byzantine-robust Training

Malicious actors in FL



Byzantine Robust Training



$$\begin{aligned}\text{agg}(\{g_i\}) &:= \text{avg}(\{g_i\}) \\ &:= \text{CM}(\{g_i\})\end{aligned}$$

Examples:

- Coordinate-wise median
[Yin et al. 2017]
- Krum
[Blanchard et al. 2018]
- Geometric median
/ RFA [Pillutla et al. 2019]

Byzantine-robust training



❖ Mean vs median

Negative result

- ❖ Robustness of the aggregation rule $\text{agg}(\{g_i\})$ does **not** imply robust training:
time-coupled attacks - “little is enough”
- ❖ Any aggregation rule which does not use history can **fail** for training (convergence)

Fix: Using history with momentum

- Simply use worker momentum

$$\mathbf{m}_i := (1 - \beta)\mathbf{g}_i + \beta\mathbf{m}_i$$

- Effectively averages past gradients, reducing variance
- (Robustly) aggregate worker momentum instead of gradients

$$\mathbf{w} := \mathbf{w} - \gamma \text{agg}(\{\mathbf{m}_i\})$$

Robustness vs Fairness

Objective

Robust mean

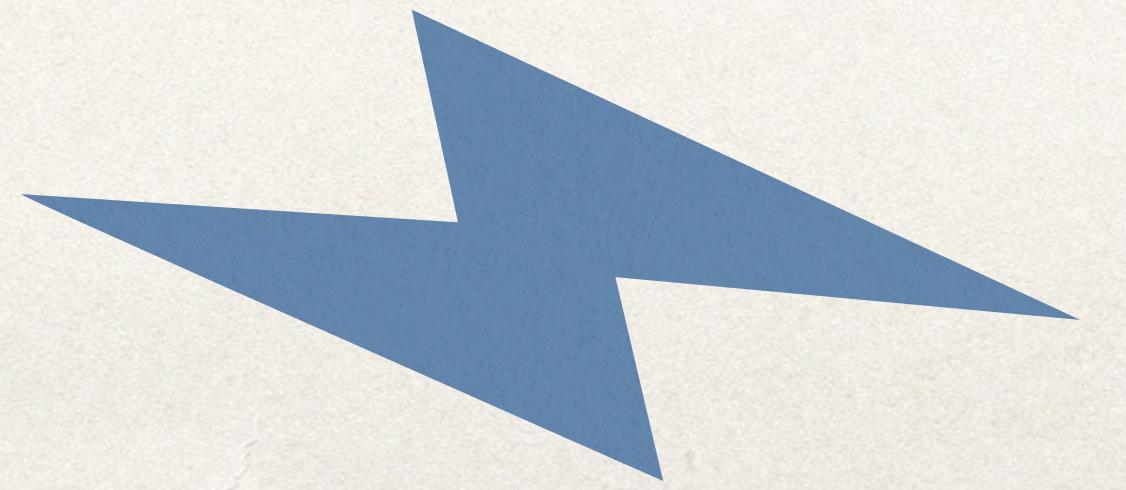
$$\text{robust-mean}_i f_i(\mathbf{x}) = \frac{1}{|good|} \sum_{i \in good} f_i(\mathbf{x})$$

Federated

$$\frac{1}{n} \sum_i^n f_i(\mathbf{x})$$

Fairness

$$\max_i f_i(\mathbf{x})$$



3b

Adversarial Attacks (at inference time)

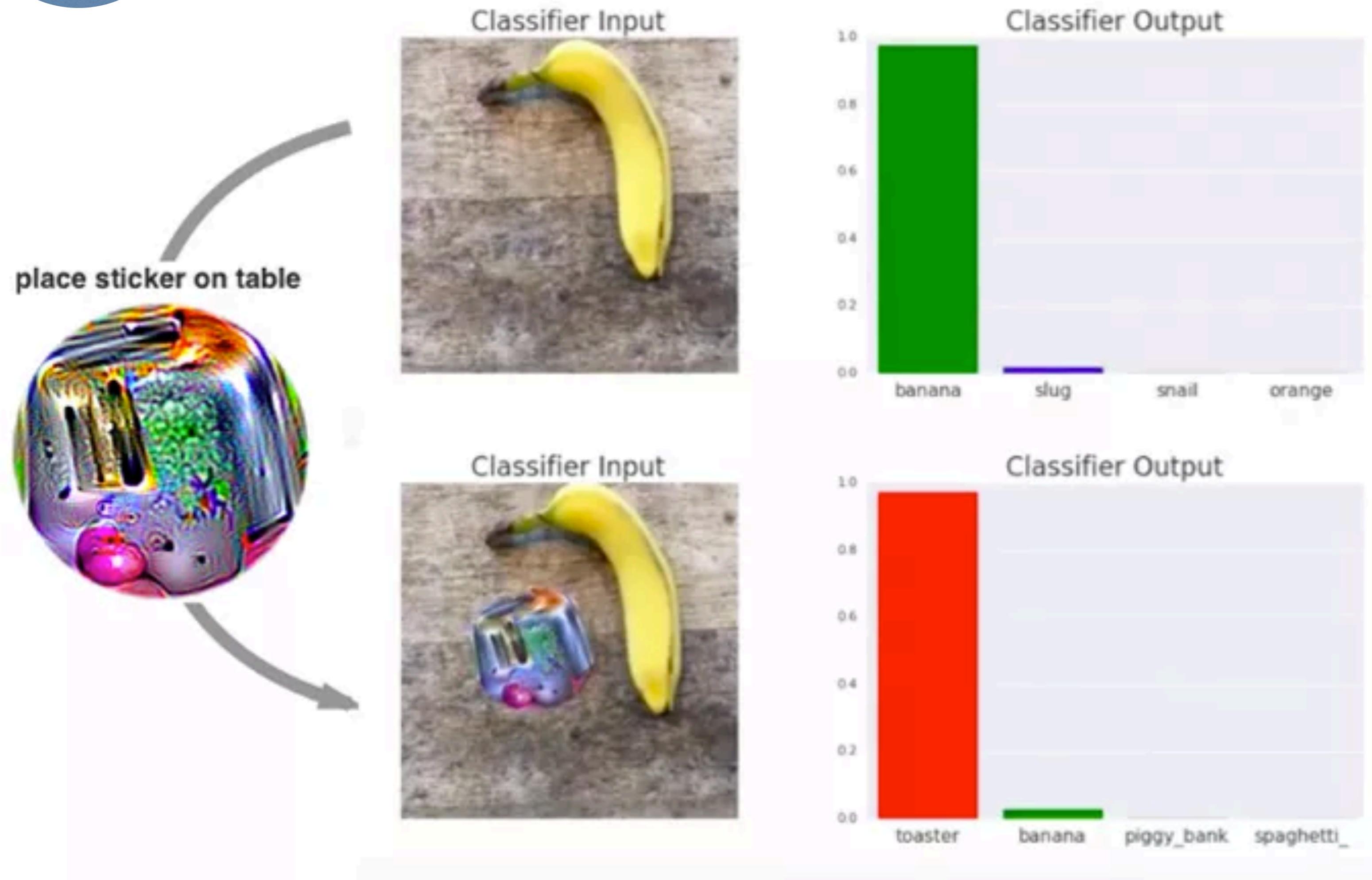


Image: [Tom B. Brown/Dandelion Mané](#)



Image: Elsayed ,Papernot et al 2018

Adversarial Attacks (at inference time)

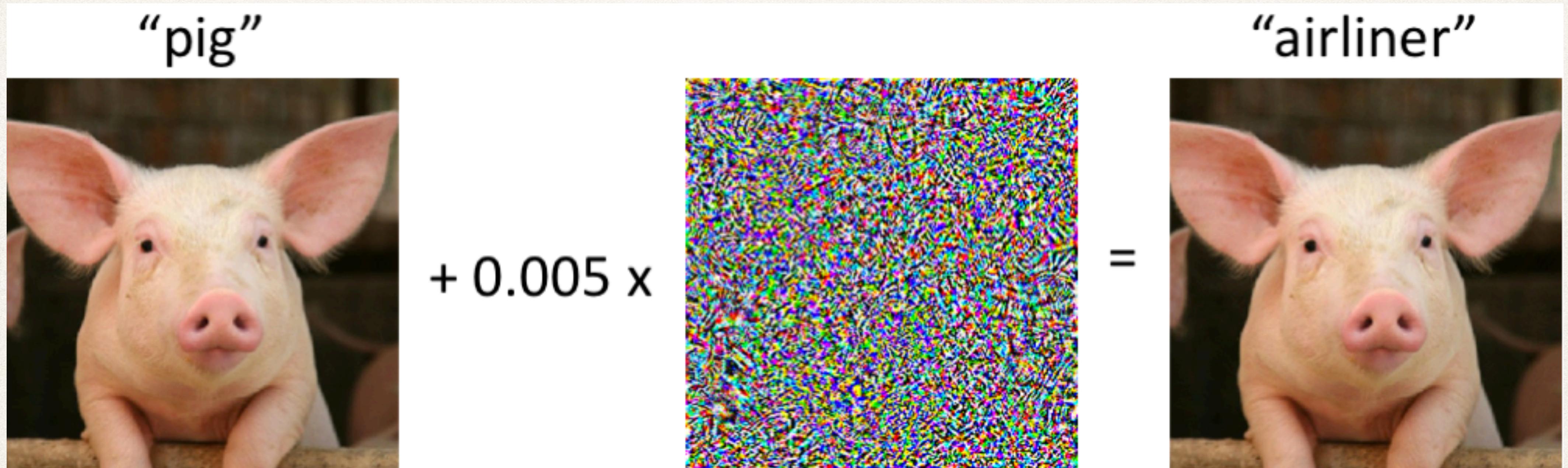


Image: [Mądry, Schmidt](#)

More info:
http://gradientscience.org/intro_adversarial/

Adversarial Attacks

- ✿ Standard **training**

$$\min_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x}_i)$$

$\nabla_{\mathbf{w}} f$
change **model**

- ✿ Attacking

$$\max_{\mathbf{x} \in R_{\infty}(\mathbf{x}_i, \varepsilon)} f_{\mathbf{w}}(\mathbf{x}_i)$$

$\nabla_{\mathbf{x}_i} f$
change **data**

- ✿ by **Projected Gradient Descent!**

4

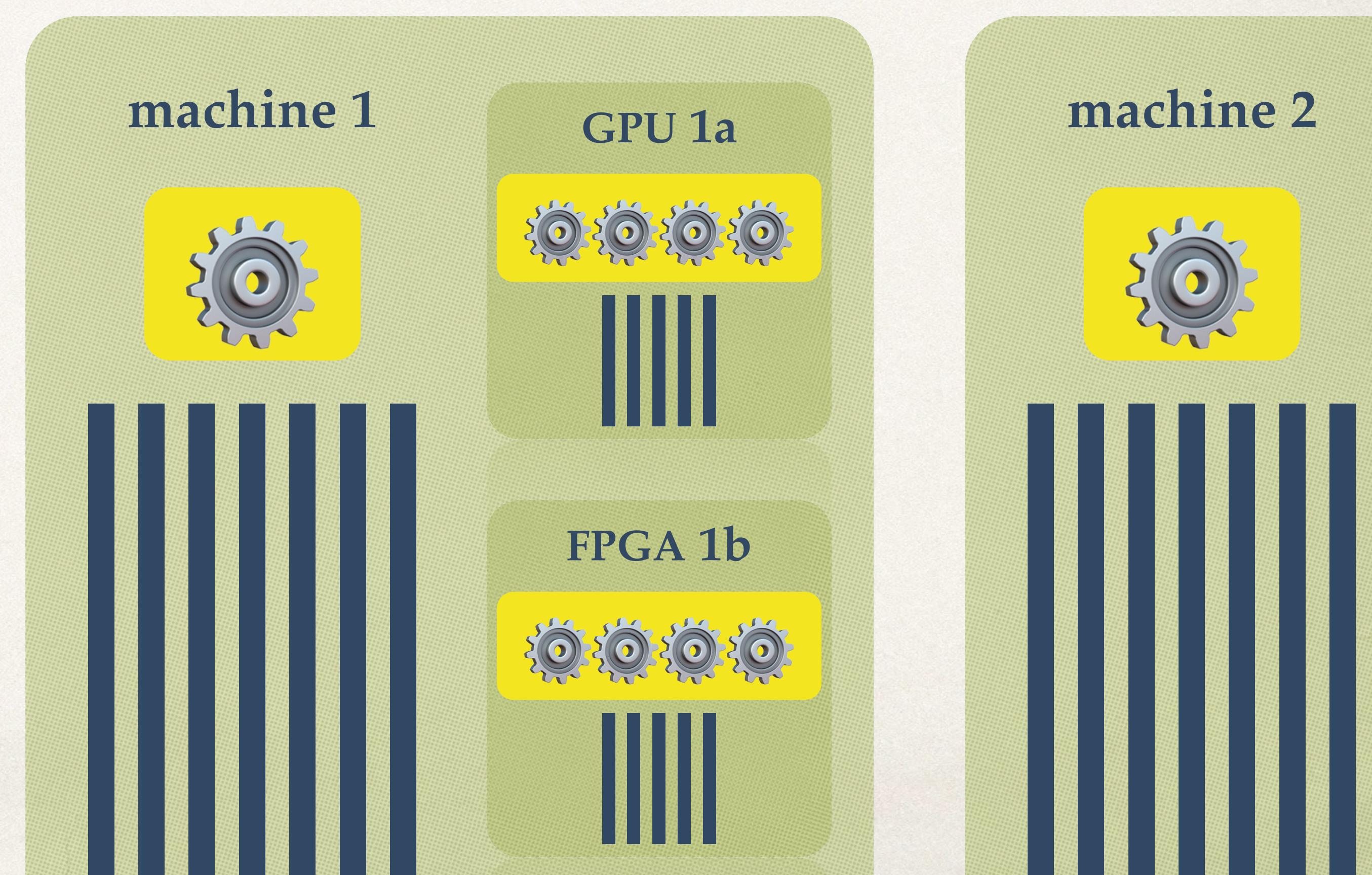
Privacy

- ❖ Secure Multiparty Computation
 - ❖ secure aggregation
(private gradients, public model)
- ❖ Differential Privacy
- ❖ Privacy/inference Attacks

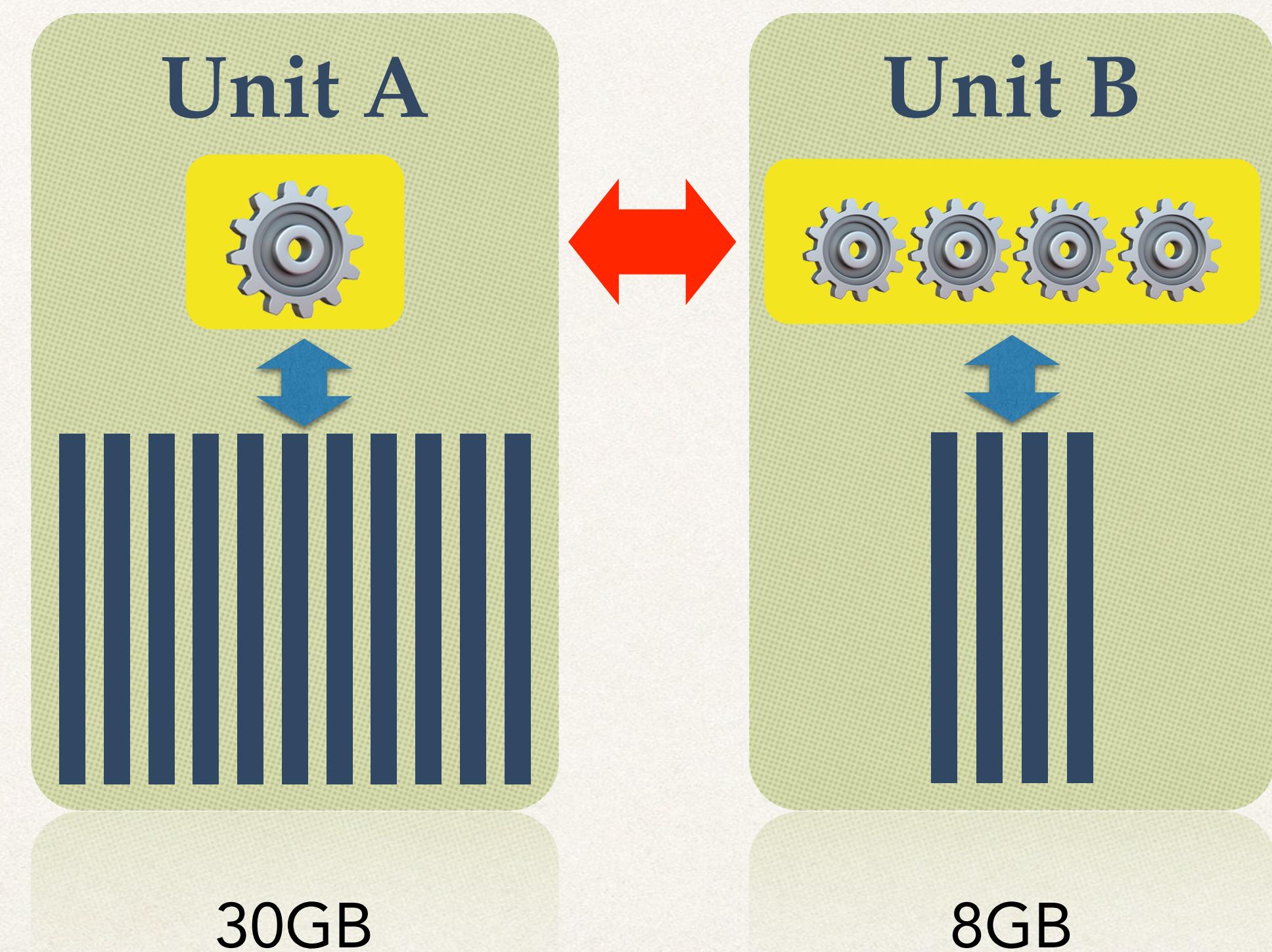
5

Leveraging Heterogenous Systems

Compute & Memory Hierarchy: Which data to put in which device?



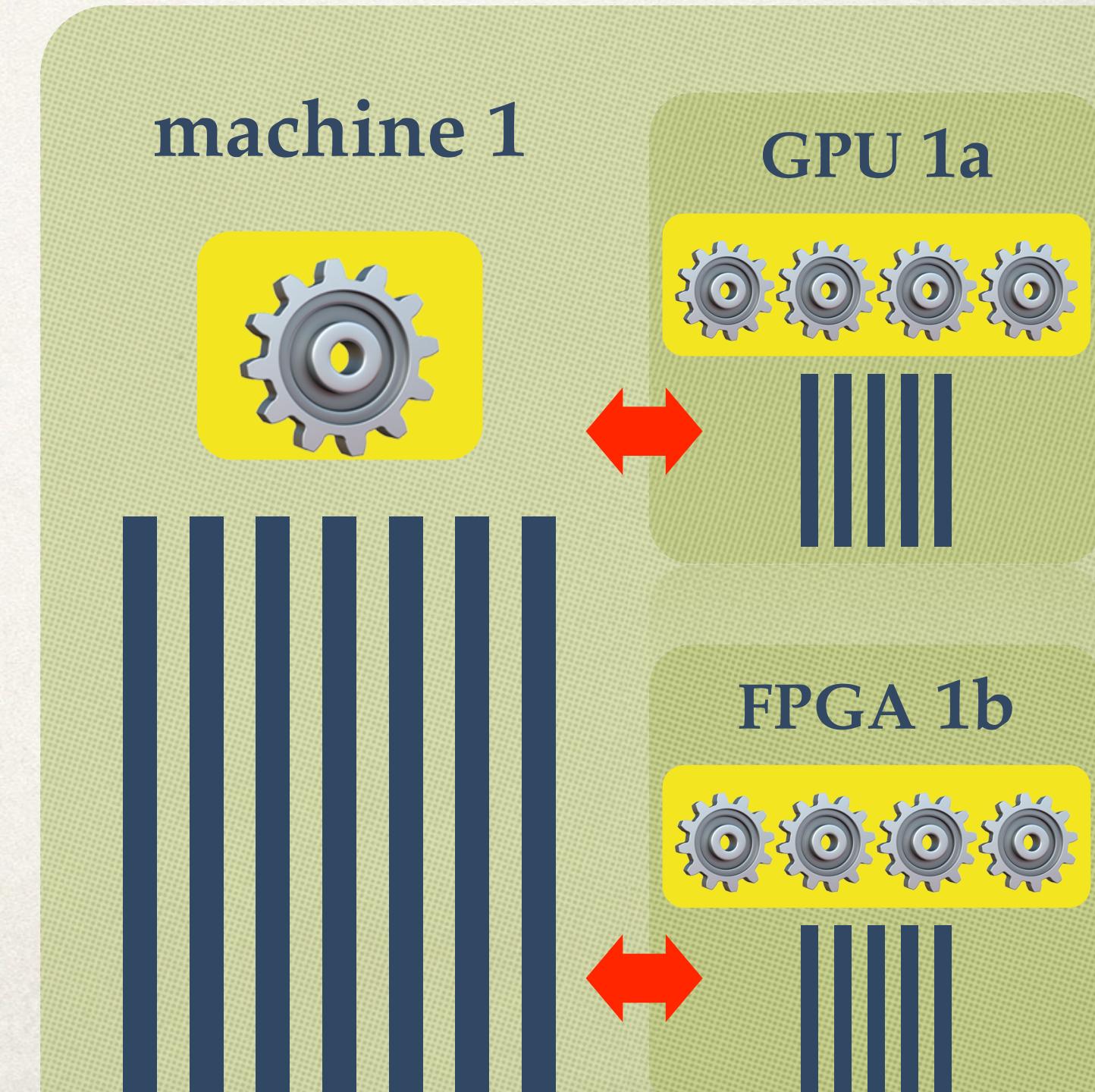
Leveraging Heterogenous Systems



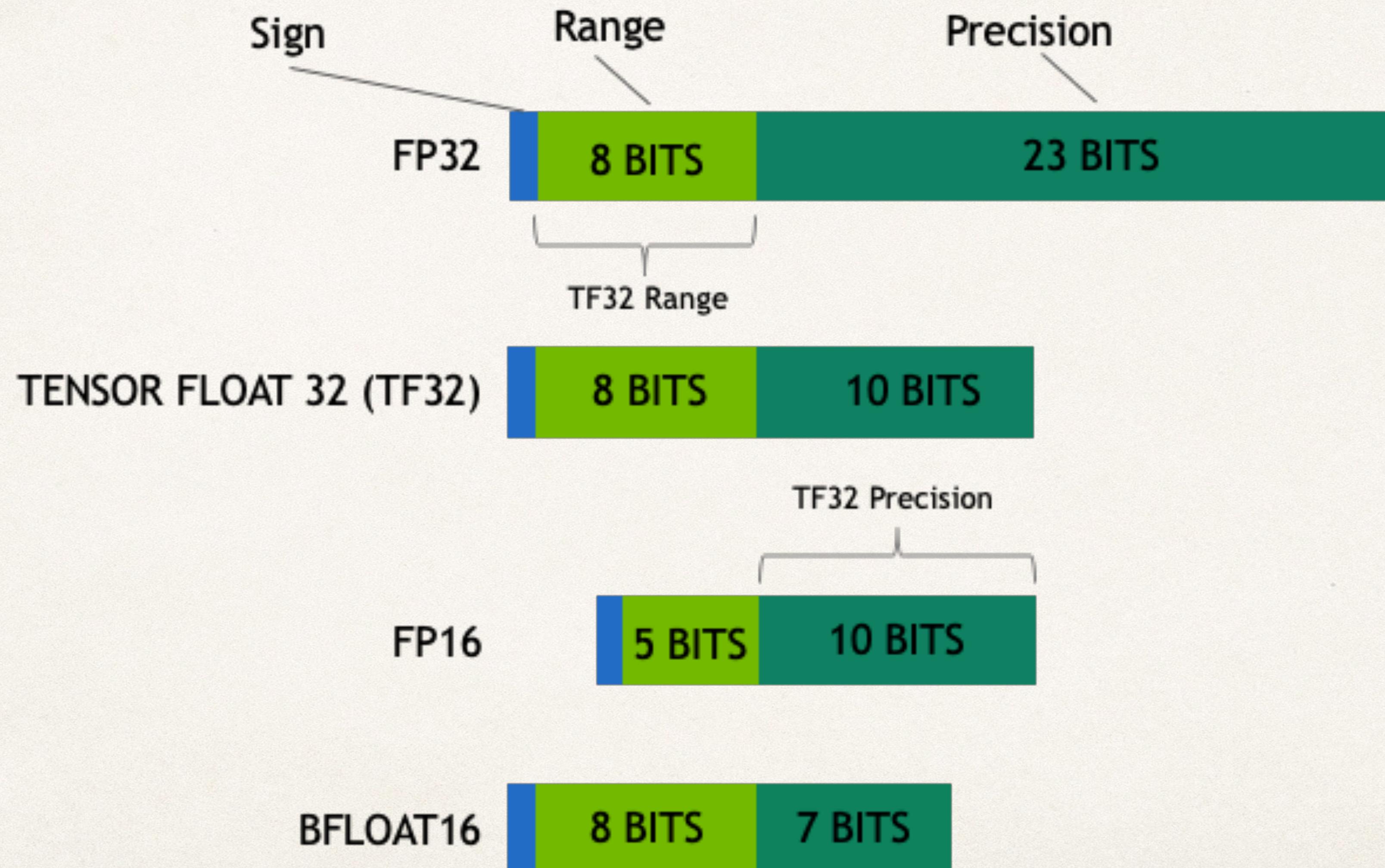
adaptive importance sampling of datapoint
e.g. for general linear models, or word2vec

Trends - Systems

- ❖ new hardware
 - ❖ TPU, GraphCore, Cerebras
- ❖ sparse ops
- ❖ efficient numerics (limited precision), model compression
- ❖ Software frameworks
 - ❖ AutoGrad (Jax, PyTorch, TensorFlow etc)
 - ❖ Backends for new hardware

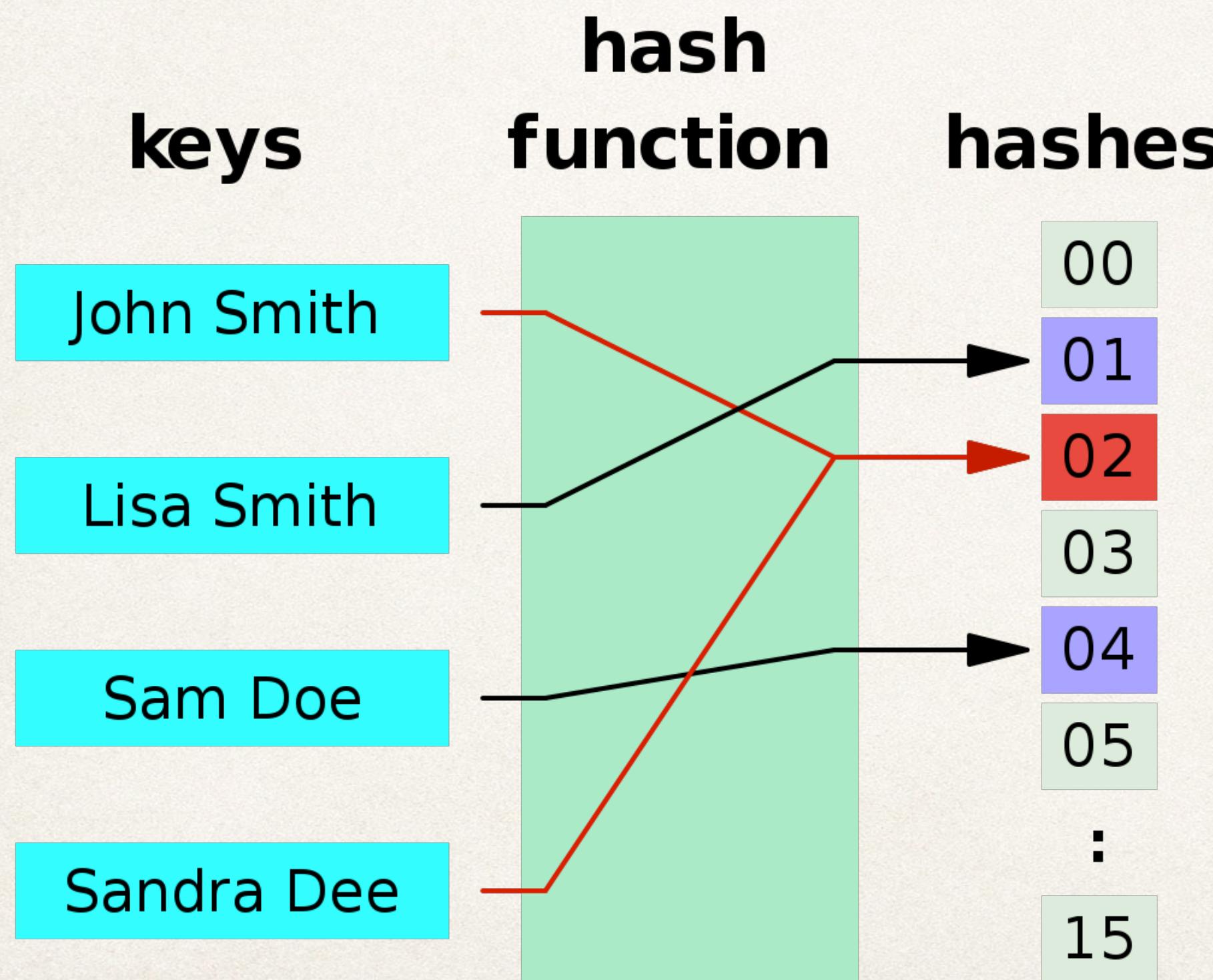


Number formats for DL



Practical tricks

- ❖ feature hashing



- ❖ limited precision operations

Auto ML

- ✿ **hyper-parameter optimization**
zero-order methods
- ✿ **learning to learn**
adaptive methods
- ✿ **neural architecture search**
zero-order, warm-start

Thanks!

mlo.epfl.ch
tml.epfl.ch