

# Optimization for Machine Learning in Practice II

Martin Jaggi

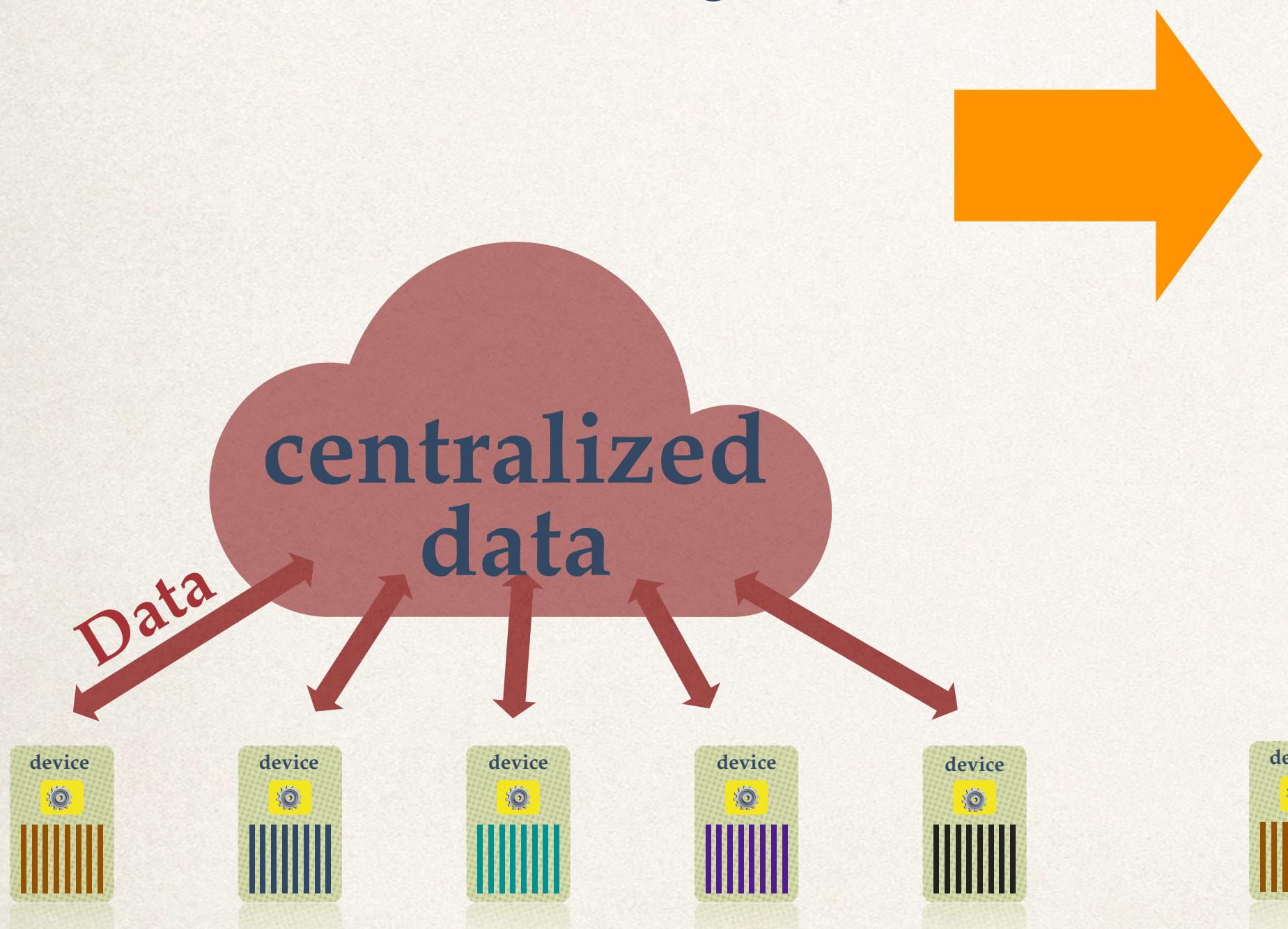


*Machine Learning and Optimization Laboratory*  
[mlo.epfl.ch](http://mlo.epfl.ch)

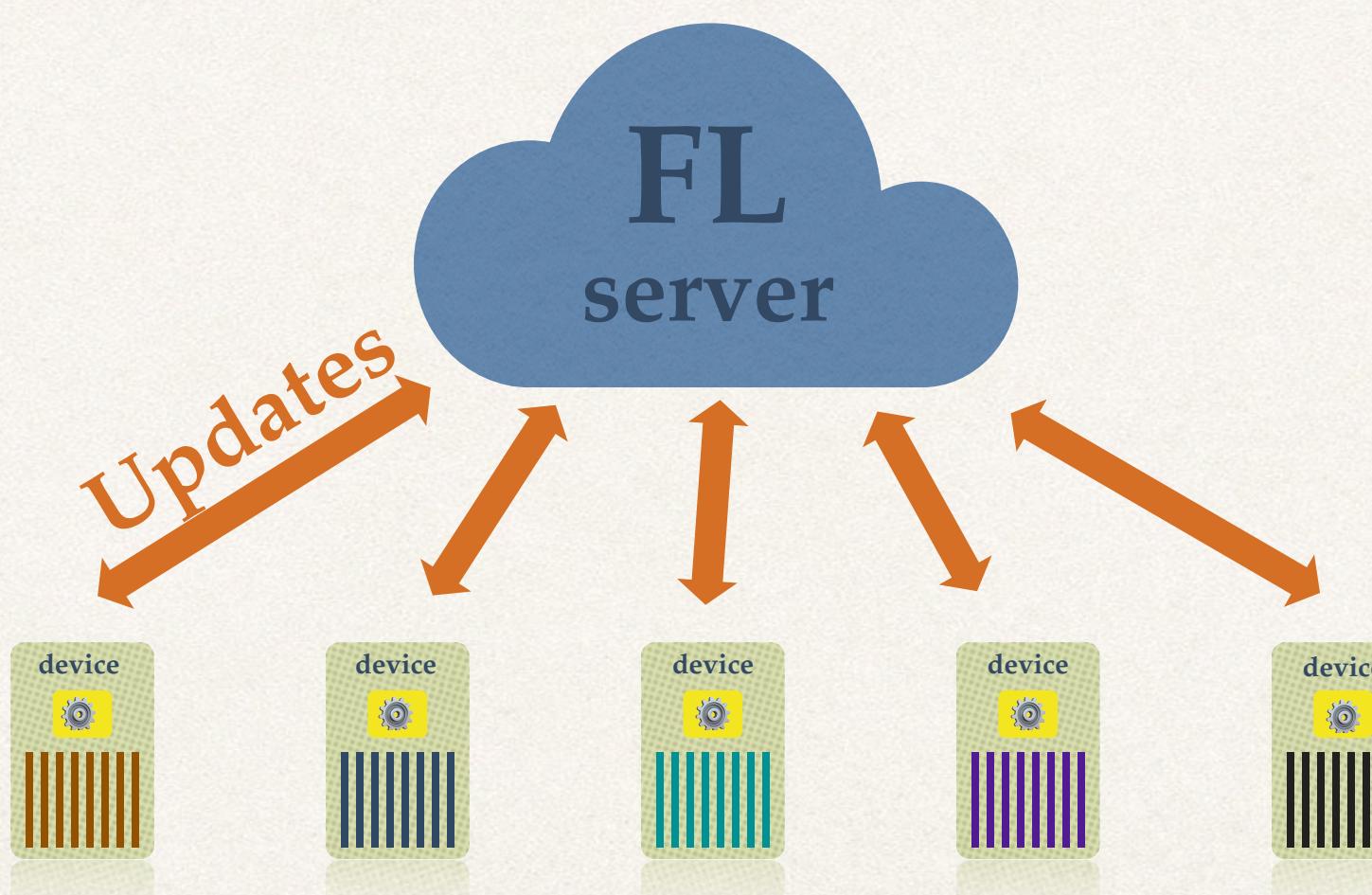
# Collaborative Learning

# Evolution

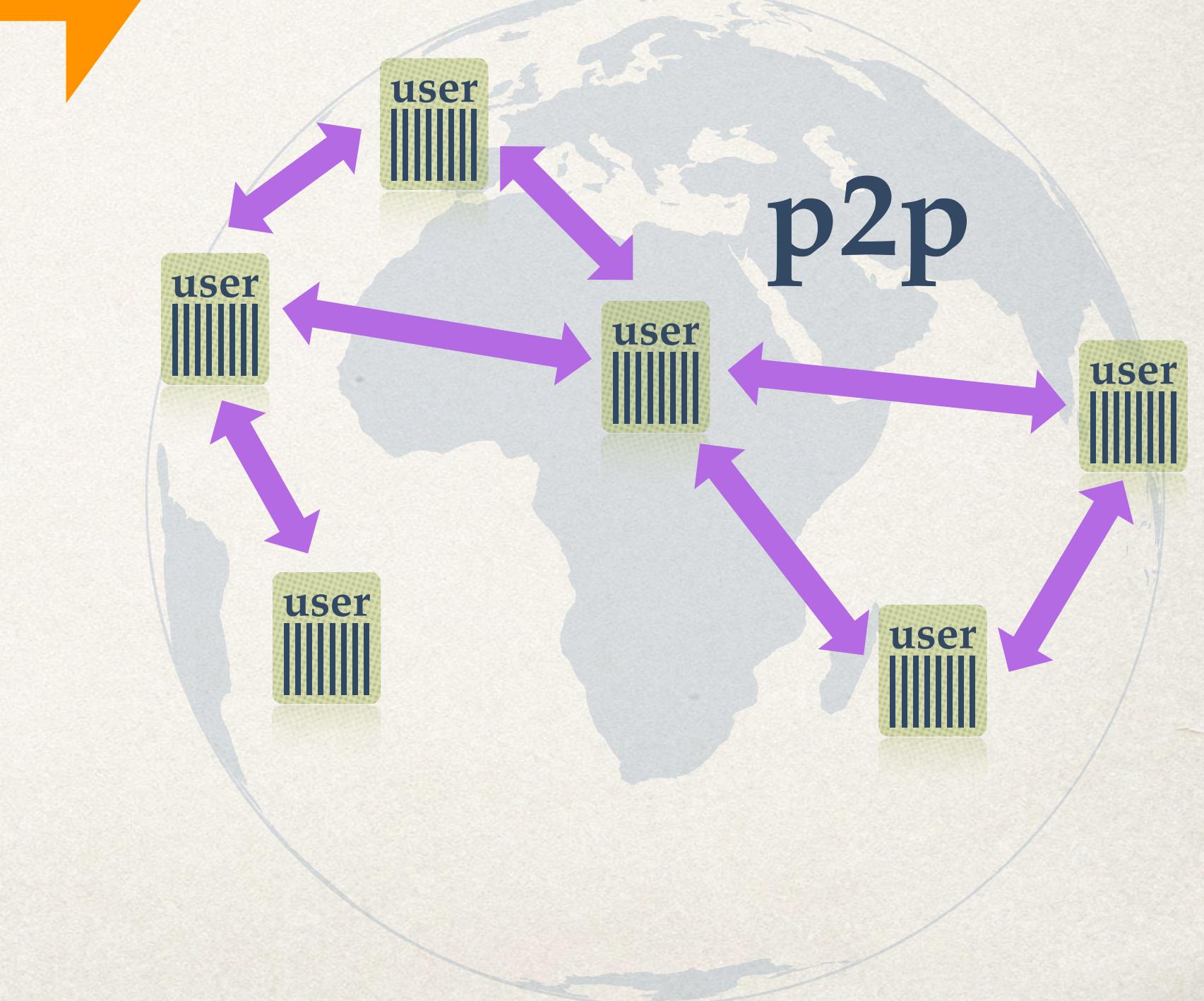
**centralized**  
traditional, sharing data



**federated**  
sharing model updates

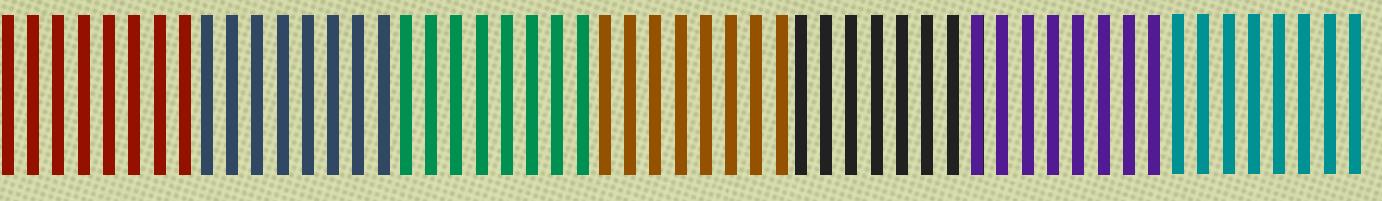


**decentralized**  
collaborative learning

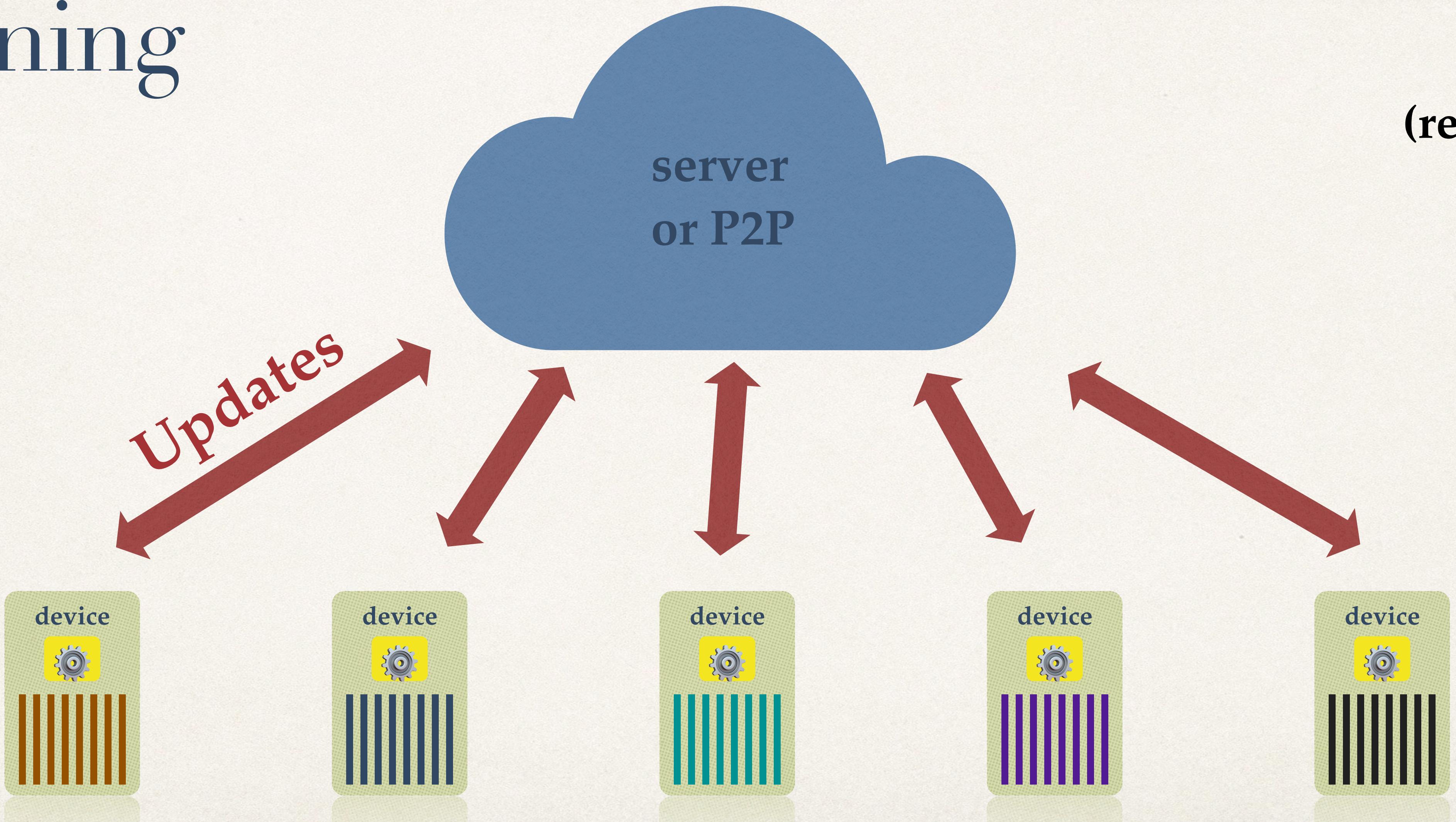


# Collaborative & Federated Training

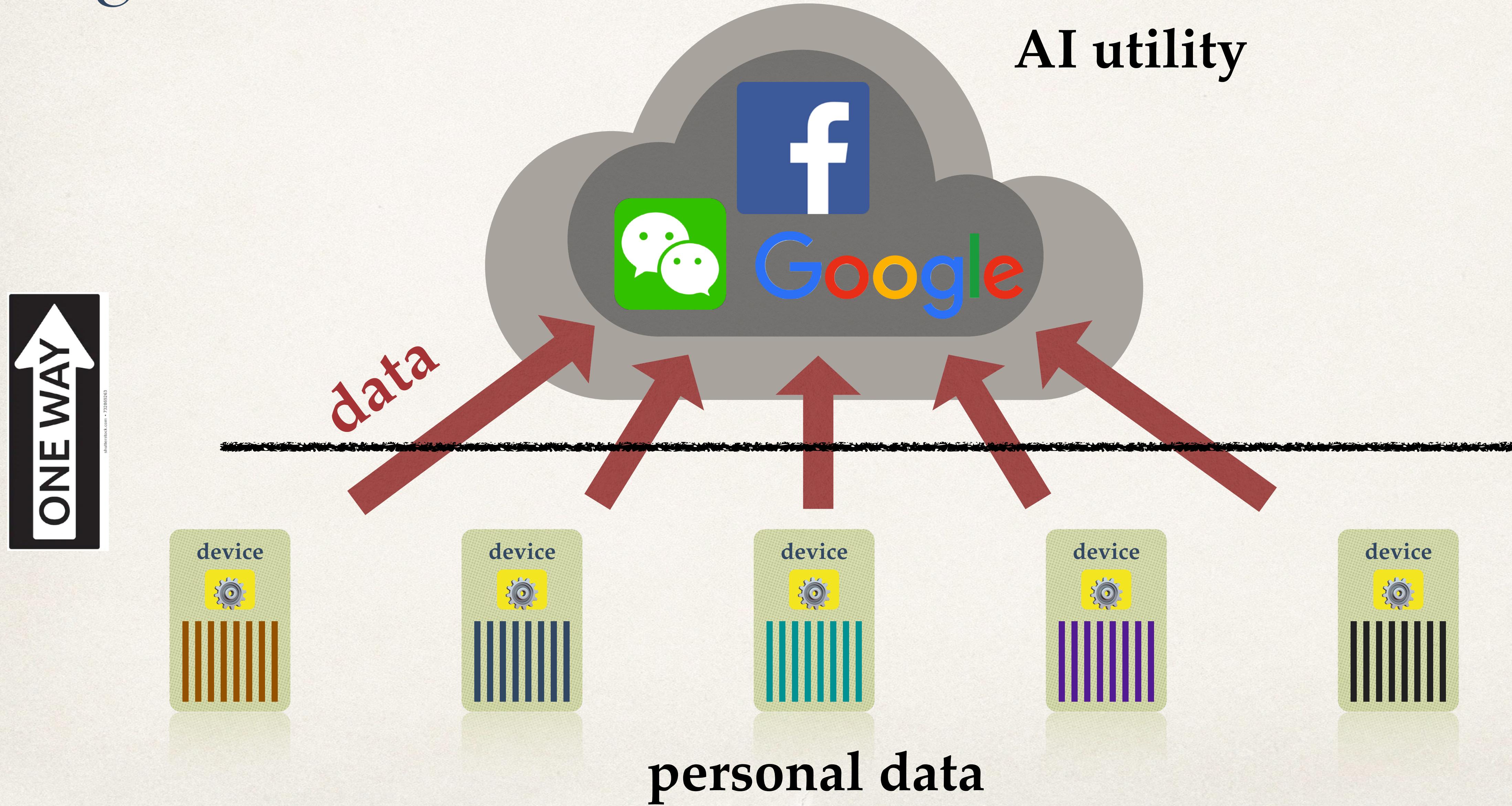
Data



(recap)

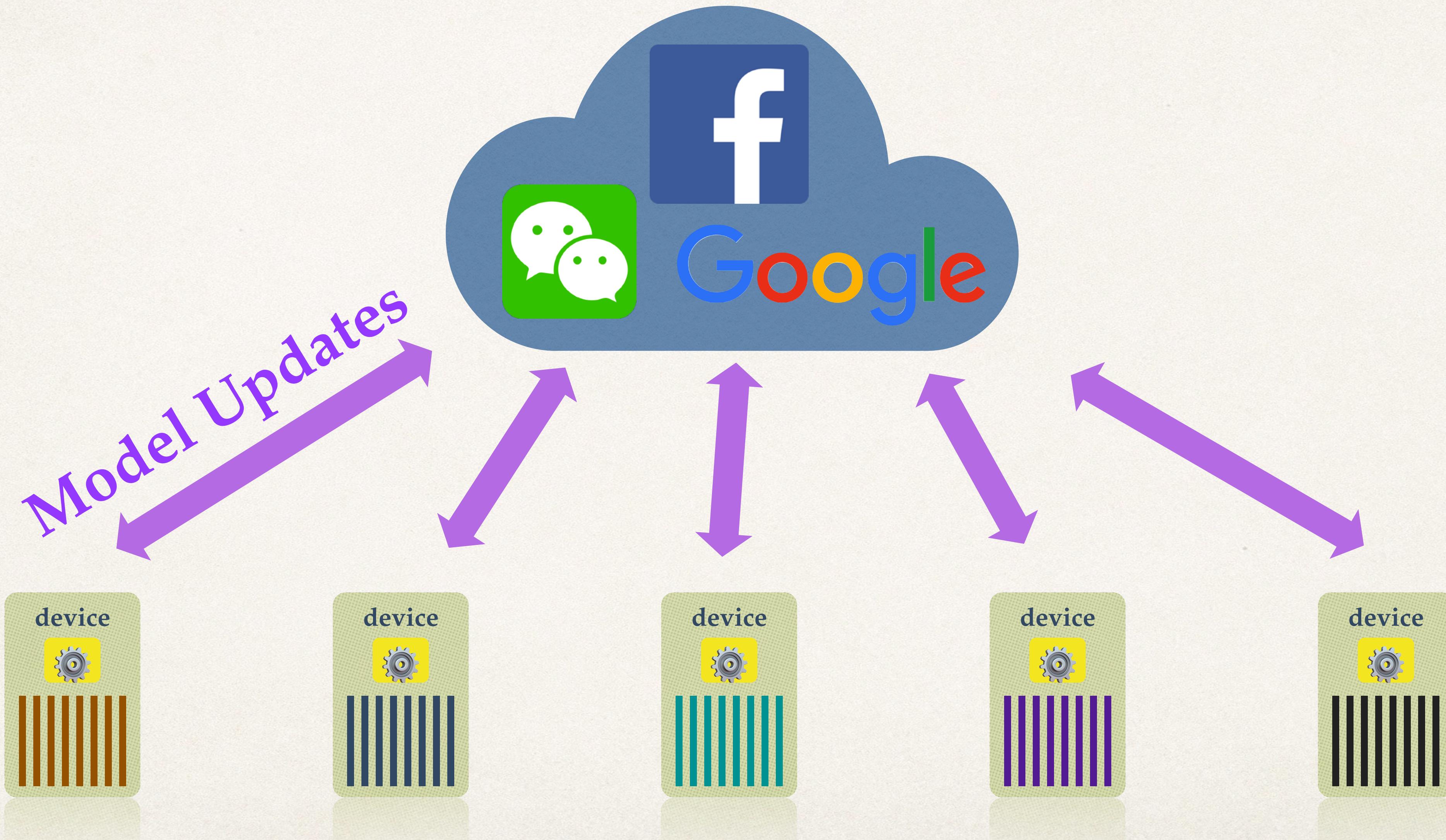


# Big Picture



2a

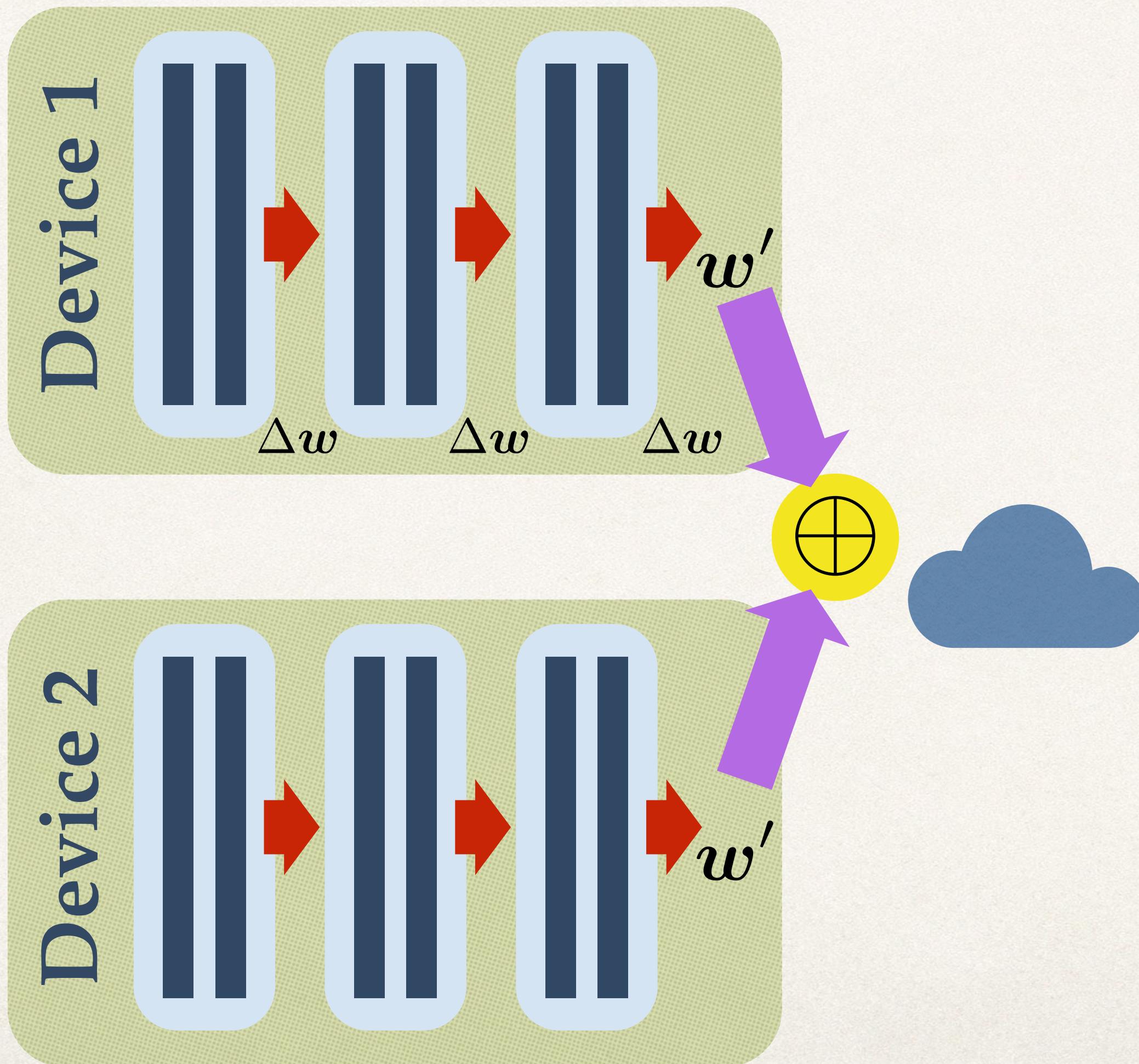
# Federated Learning



data never leaves device

2a

# Federated Learning



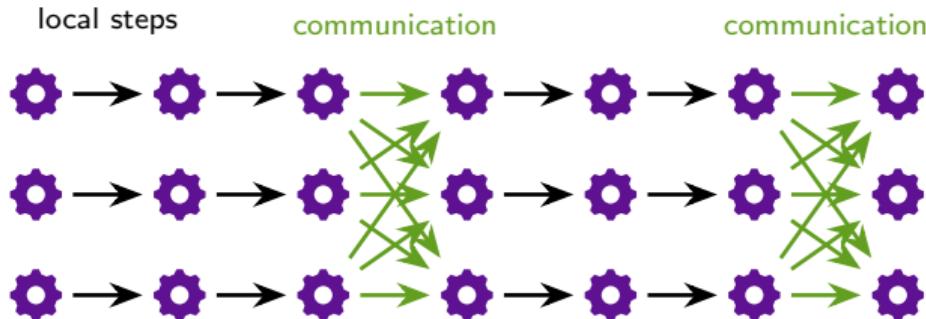
- ✿ Local SGD steps = “Federated averaging”
- ✿ Google Android Keyboard

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left[ f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n \underbrace{f_i(\mathbf{x})}_{\text{data } \mathcal{D}_i \text{ on client } i} \right] \quad f_i(\mathbf{x}) = \begin{cases} \mathbb{E}_{\xi \sim \mathcal{D}_i} F(\mathbf{x}, \xi) \\ \frac{1}{m} \sum_{j=1}^m f_{ij}(\mathbf{x}) \end{cases}$$

- ▶ Collaboratively solve **a (joint)** machine learning problem
- ▶ **efficiently**, in terms of:
  - ▶ computation (stochastic gradients, mini-batches),
  - ▶ communication (server  $\leftrightarrow$  client).

## Other very relevant scenarios:

- personalization • heterogeneity • privacy • robustness



## Notation/Setting

- ▶  $n$  machines
- ▶  $f_i(\mathbf{x})$  denote the function (data) available locally at node  $i$
- ▶ local gradient oracle  $\mathbb{E}[\mathbf{g}^{(i)}(\mathbf{x})] = \nabla f_i(\mathbf{x})$ ,  $\forall i \in [n]$ , with bounded variance:

$$\mathbb{E} \left\| \mathbf{g}^{(i)}(\mathbf{x}) - \nabla f_i(\mathbf{x}) \right\|^2 \leq \sigma^2$$

- ▶ let  $\mathbf{x}_t^{(i)} \in \mathbb{R}^d$  denote the local iterate at node  $i$ ,  $\forall i \in [n]$

# Local SGD

Input:  $\mathbf{x}_0 \in \mathbb{R}^d$ ,  $\mathbf{x}_0^{(i)} = \mathbf{x}_0$ ,  $\forall i \in [n]$ , stepsize  $\gamma$ ,  $\tau \geq 1$  (number of local steps)

At iteration  $t$  (in parallel on all nodes  $i \in [n]$ ):

$$\mathbf{g}_t^i = \mathbf{g}^{(i)}(\mathbf{x}_t^{(i)}) \quad (\text{stochastic gradient locally on each node})$$

if  $t + 1$  is a multiple of  $\tau$ :

$$\mathbf{x}_{t+1}^{(i)} = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{x}_t^{(i)} - \gamma \mathbf{g}_t^i \right) \quad (\text{global averaging})$$

otherwise:

$$\mathbf{x}_{t+1}^{(i)} = \mathbf{x}_t^{(i)} - \gamma \mathbf{g}_t^i \quad (\text{local step})$$

# Homogeneous Data

For simplicity, assume

$$f_1(\mathbf{x}) = f_2(\mathbf{x}) = \dots = f_n(\mathbf{x})$$

(note that in general  $\mathbf{g}_t^i \neq \mathbf{g}_t^j$  for  $i \neq j$ ).

- ▶ This means stochastic gradients are uniformly sampled from the whole dataset (similar as for mini-batch SGD).
- ▶ For  $\tau = 1$  Local SGD is identical to mini-batch SGD with batch size  $b = n$ .

## Theorem (Homogeneous Case, [Sti18, KLB<sup>+</sup>20])

Let  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  be  $L$ -smooth,  $\forall i \in [n]$  and  $f_i = f_j$ ,  $\forall i, j \in [n]$ , with  $\Delta = f(\mathbf{x}_0) - f^\star$ . Then there exists a stepsize  $\gamma \leq \gamma_{\text{crit}} := \frac{1}{20L\tau}$  such that after  $T$  steps (that is,  $T/\tau$  communication rounds) of Local SGD it holds

$$\min_{t \leq T} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 = \mathcal{O} \left( \frac{\Delta L \tau}{T} + \frac{(\Delta L \sigma)^{2/3} \tau^{1/3}}{T^{2/3}} + \frac{\sqrt{L \Delta \sigma^2}}{\sqrt{Tn}} \right),$$

with  $\bar{\mathbf{x}}_t := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_t^{(i)}$ .

## Discussion

- ▶ Linear speedup if  $\sigma^2 > 0$ : the variance decreases linearly in the number of oracle calls ( $Tn$ ). This is optimal.
- ▶ The deterministic optimization term (the term not depending on  $\sigma^2$ ) is impacted by  $\tau$  (similarly as with mini-batch SGD).
  - ▶ ideally, we would have hoped to see there  $\mathcal{O}\left(\frac{\Delta L}{T}\right)$  (= progress in every iteration)  
vs.  $\mathcal{O}\left(\frac{\Delta L \tau}{T}\right)$  (= progress every communication round)
- ▶ The theorem shows almost the same convergence as for mini-batch SGD, up to the higher order  $\mathcal{O}(T^{-2/3})$  term.
  - ▶ There is no clear winner.

# Performance in Practice [LSPJ19]

ResNet-20 on CIFAR-10 (IID data)

	Top-1 acc.	local gradients	communication
Mini-batch SGD ( $n = 16, \tau = 128$ )	92.5%	2048	-
Mini-batch SGD ( $n = 16, \tau = 1024$ )	76.3%	16384	$\div 8$
Local-SGD ( $n = 16, \tau = 8 \times 128$ )	92.0%	16384	$\div 8$

# Proof I

We will use the **virtual sequence** technique. As virtual sequence we consider  $\bar{\mathbf{x}}_t$  (note that the average is not computed in every iteration).

## Lemma (Decrease)

For  $\gamma \leq \frac{1}{4L}$  it holds

$$\mathbb{E}f(\bar{\mathbf{x}}_{t+1}) \leq Ef(\bar{\mathbf{x}}_t) - \frac{\gamma}{4} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \gamma^2 L \frac{\sigma^2}{\textcolor{red}{n}} + \frac{\gamma L^2}{n} \sum_{i=1}^n \mathbb{E}\|\mathbf{x}_t^{(i)} - \bar{\mathbf{x}}_t\|^2$$

## Lemma (Difference)

For  $\gamma \leq \gamma_{\text{crit}} = \frac{1}{20L\tau}$ , with the notation for  $R_t = \frac{1}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}_t - \mathbf{x}_t^{(i)}\|^2$ , it holds

$$\mathbb{E}R_t \leq \frac{1}{20L^2\tau} \sum_{j=(t-1)-k}^{t-1} \mathbb{E}\|\nabla f(\bar{\mathbf{x}}_j)\|^2 + 5\gamma^2\tau\sigma^2$$

where  $(t-1) - k$  denotes the index of the last communication round ( $k \leq \tau - 1$ ).

## Proof II

Plug (Difference) into (Decrease), re-arrange and divide by  $\gamma$ :

$$\frac{1}{4} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 \leq \frac{1}{\gamma} (\mathbb{E} f(\bar{\mathbf{x}}_t) - \mathbb{E} f(\bar{\mathbf{x}}_{t+1})) + \gamma L \frac{\sigma^2}{n} + \frac{1}{20\tau} \sum_{j=(t-1)-k}^{t-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_j)\|^2 + 5\gamma^2 L^2 \tau \sigma^2$$

Now we divide by  $T$  and sum over  $t = 0, \dots, T-1$ :

$$\frac{1}{4T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 \leq \frac{1}{\gamma T} \sum_{t=0}^{T-1} \left[ (\mathbb{E} f(\bar{\mathbf{x}}_t) - \mathbb{E} f(\bar{\mathbf{x}}_{t+1})) + \frac{1}{20T} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 \right] + \gamma^2 L \frac{\sigma^2}{n} + 5\gamma L^2 \tau \sigma^2$$

Note that  $\sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2$  appears on both sides, with  $\frac{1}{4T} - \frac{1}{20T} = \frac{1}{5T}$ .

$$\frac{1}{5T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 \leq \frac{\Delta}{\gamma T} + \gamma L \frac{\sigma^2}{n} + 5\gamma^2 L^2 \tau \sigma^2.$$

Now the result follows by tuning  $\gamma$  (Exercise).

## Proof of Lemma (Decrease)

By  $L$ -smoothness:

$$\begin{aligned}\mathbb{E}f(\bar{\mathbf{x}}_{t+1}) &\leq \mathbb{E}f(\bar{\mathbf{x}}_t) - \frac{\gamma}{n} \sum_{i=1}^n \nabla f(\bar{\mathbf{x}}_t)^\top \mathbf{g}_t^i + \frac{\gamma^2 L}{2} \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{g}_t^i \right\|^2 \\ &\leq \mathbb{E}f(\bar{\mathbf{x}}_t) - \frac{\gamma}{n} \sum_{i=1}^n \nabla f(\bar{\mathbf{x}}_t)^\top \nabla f_i(\mathbf{x}_t^{(i)}) + \frac{\gamma^2 L \sigma^2}{2n} + \frac{\gamma^2 L}{2} \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t^{(i)}) \right\|^2\end{aligned}$$

Similarly as we have seen before, by adding and subtracting  $\nabla f(\bar{\mathbf{x}}_t)$ :

$$\begin{aligned}-\frac{1}{n} \sum_{i=1}^n \nabla f(\bar{\mathbf{x}}_t)^\top \nabla f_i(\mathbf{x}_t^{(i)}) &= -\nabla f(\bar{\mathbf{x}}_t)^\top \nabla f(\bar{\mathbf{x}}_t) + \frac{1}{n} \sum_{i=1}^n \nabla f(\bar{\mathbf{x}}_t)^\top (\nabla f(\bar{\mathbf{x}}_t) - \nabla f_i(\mathbf{x}_t^{(i)})) \\ &\leq -\frac{1}{2} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{2} \left\| \nabla f(\bar{\mathbf{x}}_t) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t^{(i)}) \right\|^2 \\ &\leq -\frac{1}{2} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \frac{L^2}{2n} \sum_{i=1}^n \left\| \bar{\mathbf{x}}_t - \mathbf{x}_t^{(i)} \right\|^2\end{aligned}$$

Note:  $-\mathbf{a}^\top \mathbf{b} \leq \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2$ .

## Continued

And by adding and subtracting  $\nabla f(\bar{\mathbf{x}}_t)$  in the last term:  $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$

$$\begin{aligned} \frac{1}{2} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t^{(i)}) \right\|^2 &\leq \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t^{(i)}) - \nabla f(\bar{\mathbf{x}}_t) \right\|^2 + \|\nabla f(\bar{\mathbf{x}}_t)\|^2 \\ &\leq \frac{L^2}{n} \sum_{i=1}^n \left\| \bar{\mathbf{x}}_t - \mathbf{x}_t^{(i)} \right\|^2 + \|\nabla f(\bar{\mathbf{x}}_t)\|^2 \end{aligned}$$

Now we plug everything together, and use  $\gamma \leq \frac{1}{4L}$ .

$$\mathbb{E}f(\bar{\mathbf{x}}_{t+1}) \leq \mathbb{E}f(\bar{\mathbf{x}}_t) + \left( \gamma^2 L - \frac{\gamma}{2} \right) \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \frac{\gamma^2 L \sigma^2}{2n} + \left( \frac{\gamma L^2}{2} + \gamma^2 L^3 \right) \frac{1}{n} \sum_{i=1}^n \left\| \bar{\mathbf{x}}_t - \mathbf{x}_t^{(i)} \right\|^2$$

## Proof of Lemma (Difference)

Note that if  $t$  is a multiple of  $\tau$ , then  $R_t = 0$  and there is nothing to prove. Otherwise note that

$$\begin{aligned}\mathbb{E}R_{t+1} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\| \bar{\mathbf{x}}_{t+1} - \mathbf{x}_{t+1}^i \right\|^2 = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\| \bar{\mathbf{x}}_t - \mathbf{x}_t^i + \gamma \mathbf{g}_t^i - \gamma \bar{\mathbf{g}}_t \right\|^2 \\ &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\| \bar{\mathbf{x}}_t - \mathbf{x}_t^i + \gamma \nabla f_i(\mathbf{x}_t^i) - \gamma \bar{\mathbf{v}}_t \right\|^2 + \gamma^2 \sigma^2,\end{aligned}$$

where  $\bar{\mathbf{g}}_t := \frac{1}{n} \sum_{i=1}^n \mathbf{g}_t^i$  and  $\bar{\mathbf{v}}_t := \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t^{(i)})$  denote the average of the client gradients. With the inequality  $\|\mathbf{a} + \mathbf{b}\|^2 \leq (1 + \tau^{-1}) \|\mathbf{a}\|^2 + 2\tau \|\mathbf{b}\|^2$  for  $\tau \geq 1$ ,  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ , we continue:

$$\begin{aligned}\mathbb{E}R_{t+1} &\leq \left(1 + \frac{1}{\tau}\right) \mathbb{E}R_t + \frac{2\tau\gamma^2}{n} \sum_{i=1}^n \mathbb{E} \left\| \nabla f(\mathbf{x}_t^{(i)}) - \bar{\mathbf{v}}_t \right\|^2 + \gamma^2 \sigma^2 \\ &\leq \left(1 + \frac{1}{\tau}\right) \mathbb{E}R_t + \frac{2\tau\gamma^2}{n} \sum_{i=1}^n \mathbb{E} \left\| \nabla f_i(\mathbf{x}_t^{(i)}) \right\|^2 + \gamma^2 \sigma^2 \\ &\leq \left(1 + \frac{1}{\tau}\right) \mathbb{E}R_t + \frac{2\tau\gamma^2}{n} \sum_{i=1}^n \mathbb{E} \left( 2 \|\nabla f_i(\bar{\mathbf{x}}_t)\|^2 + 2 \left\| \nabla f_i(\mathbf{x}_t^{(i)}) - \nabla f_i(\bar{\mathbf{x}}_t) \right\|^2 \right) + \gamma^2 \sigma^2\end{aligned}$$

(\*)<sub>11/22</sub>

## Continued

We now use  $f_1 = f_2 = \dots = f_n$  and  $\gamma \leq \frac{1}{20L\tau}$  to simplify:

$$\begin{aligned}\mathbb{E}R_{t+1} &\leq \left(1 + \frac{1}{\tau}\right) \mathbb{E}R_t + \frac{1}{100L^2\tau} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{100\tau} \mathbb{E}R_t + \gamma^2 \sigma^2 \\ &\leq \left(1 + \frac{3}{2\tau}\right) \mathbb{E}R_t + \frac{1}{100L^2\tau} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \gamma^2 \sigma^2\end{aligned}$$

The lemma now follows by unrolling, and noting that  $\left(1 + \frac{3}{2\tau}\right)^j \leq 5$  for all  $0 \leq j \leq \tau$ .

# Federated Learning on Heterogeneous Data

- We now want/need to drop the assumption that  $f_i(\mathbf{x}) = f_j(\mathbf{x})$ ,  $i \neq j$ .

Data-heterogeneity causes **drift** when doing local steps.

# Client drift

- \* Federated Learning

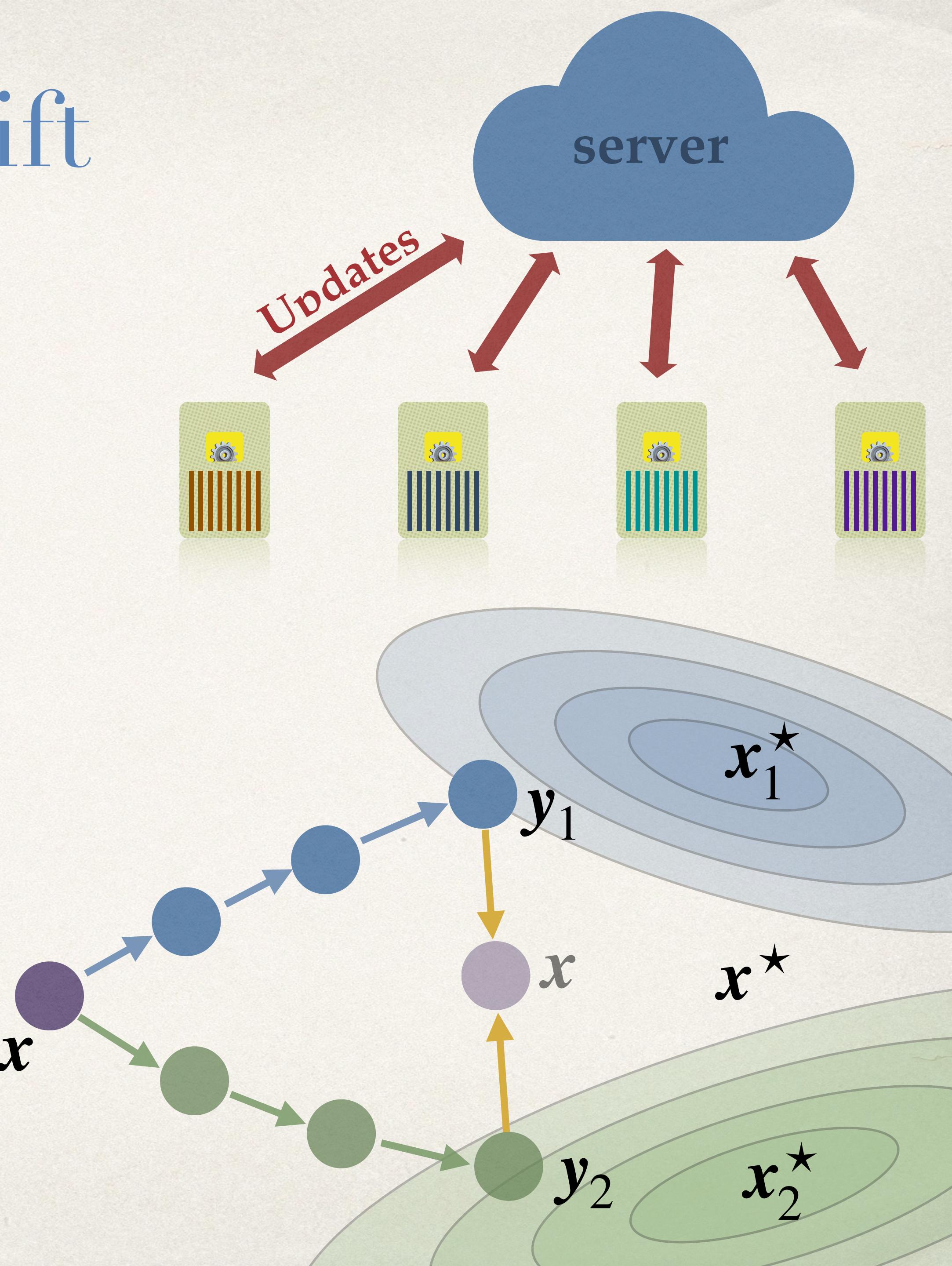
$$\min_{\mathbf{x}} \frac{1}{n} \sum_i^n f_i(\mathbf{x})$$

- \* Fed Avg / Local SGD

*for some local steps*

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

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



# Measuring Data Heterogeneity

## Definition (Heterogeneity $\zeta^2$ )

(The smallest) parameter  $\zeta^2$  such that for all  $\mathbf{x} \in \mathbb{R}^d$ :

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \zeta^2.$$

Similar to the definition of the variance, but now we measure the inter-worker variance.

# Convergence

Theorem (Heterogeneous Case, [KLB<sup>+</sup>20])

Let  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  be  $L$ -smooth,  $\forall i \in [n]$  and the function heterogeneity bounded by  $\zeta^2$ , with  $\Delta = f(\mathbf{x}_0) - f^*$ . Then for there exists a stepsize  $\gamma \leq \gamma_{\text{crit}} := \frac{1}{10L\tau}$  such that after  $T$  steps (that is,  $T/\tau$  communication rounds) of Local SGD it holds

$$\min_{t \leq T} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 = \mathcal{O} \left( \frac{\Delta L \tau}{T} + \left( \frac{L \Delta (\tau \zeta + \sqrt{\tau} \sigma)}{T} \right)^{2/3} + \frac{\sqrt{L \Delta \sigma^2}}{\sqrt{Tn}} \right),$$

with  $\bar{\mathbf{x}}_t := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_t^{(i)}$ .

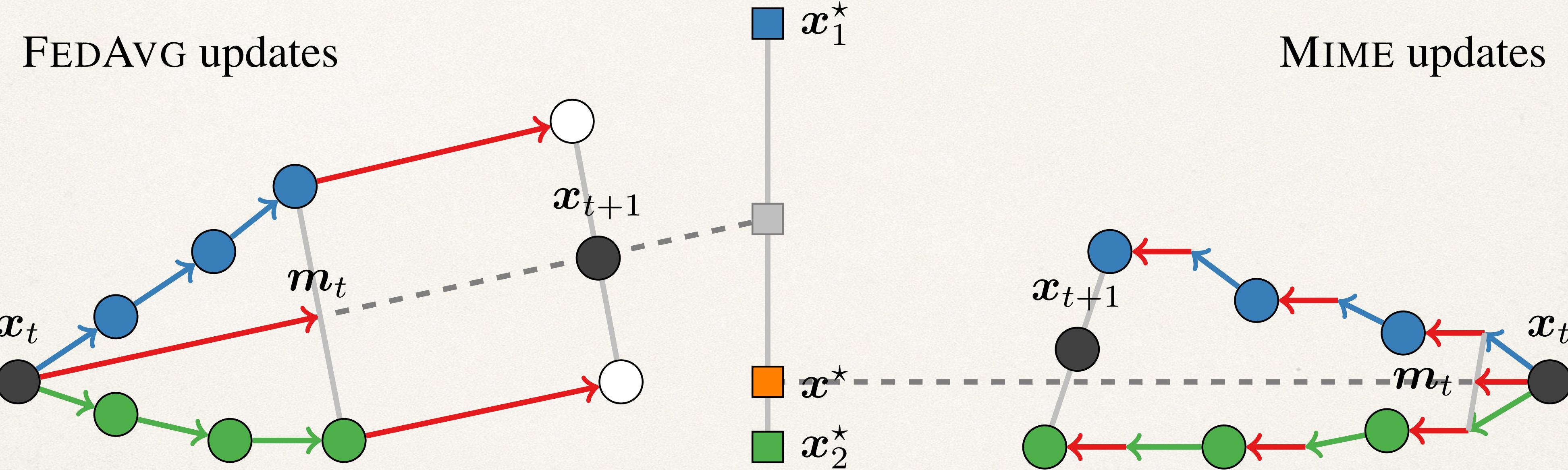
## Discussion

- ▶  $\zeta^2 > 0$  slows down the convergence.
- ▶ The dependency on  $\zeta$  is optimal.
- ▶ In general, the convergence is slower than for mini-batch SGD.
- ▶ Momentum (applied on the server) speeds up convergence, and solves the drift problem [KJK<sup>+</sup>21, CHY23]

**Why does then Local SGD behave well on many practical problems?**

- ▶ Speedup can be proven under different similarity assumptions.

# 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*

# Proof

We will prove a new [\(Difference\) Lemma](#), the rest of the proof will follow the usual pattern.

## Lemma (Difference)

For  $\gamma \leq \gamma_{\text{crit}} = \frac{1}{20L\tau}$ , with the notation for  $R_t = \frac{1}{n} \sum_{i=1}^n \left\| \bar{\mathbf{x}}_t - \mathbf{x}_t^{(i)} \right\|^2$ , it holds

$$\mathbb{E}R_t \leq \frac{1}{10L^2\tau} \sum_{j=(t-1)-k}^{t-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_j)\|^2 + 5\gamma^2\sigma^2\tau + 40\gamma^2\tau^2\zeta^2$$

where  $(t-1) - k$  denotes the index of the last communication round ( $k \leq \tau - 1$ ).

## Proof II

Plug the new (Difference) bound into (Decrease), re-arrange, divide by  $\gamma$ , divide by  $T$ , and sum over  $t = 0, \dots, T - 1 \dots$  (left as exercise!)

We will end up with:

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 = \mathcal{O} \left( \frac{\Delta}{\gamma T} + \gamma L \frac{\sigma^2}{n} + \gamma^2 L^2 (\tau^2 \zeta^2 + \tau \sigma^2) \right)$$

Now use previous exercise.

## Proof of Lemma (Difference) I

For the proof of the difference lemma, we need to be more careful at inequality (\*). Recall the inequality (\*):

$$\mathbb{E}R_{t+1} \leq \left(1 + \frac{1}{\tau}\right) \mathbb{E}R_t + \frac{2\tau\gamma^2}{n} \sum_{i=1}^n \mathbb{E} \left( 2 \|\nabla f_i(\bar{\mathbf{x}}_t)\|^2 + 2 \left\| \nabla f_i(\mathbf{x}_t^{(i)}) - \nabla f_i(\bar{\mathbf{x}}_t) \right\|^2 \right) + \gamma^2 \sigma^2$$

Previously we used  $f_i = f_j$ , to simplify. Instead, note that:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\bar{\mathbf{x}}_j)\|^2 &= \frac{1}{n} \sum_{i=1}^n \left( \|\nabla f_i(\bar{\mathbf{x}}_j) - \nabla f(\bar{\mathbf{x}}_j) + \nabla f(\bar{\mathbf{x}}_j)\|^2 \right) \\ &\leq \frac{2}{n} \sum_{i=1}^n \left( \|\nabla f_i(\bar{\mathbf{x}}_j) - \nabla f(\bar{\mathbf{x}}_j)\|^2 + \|\nabla f(\bar{\mathbf{x}}_j)\|^2 \right) \\ &\leq 2\zeta^2 + 2 \|\nabla f(\bar{\mathbf{x}}_j)\|^2 \end{aligned}$$

## Proof of Lemma (Difference) II

Therefore (and with  $\gamma \leq \frac{1}{20L\tau}$ )

$$\begin{aligned}\mathbb{E}R_{t+1} &\leq \left(1 + \frac{1}{\tau}\right) \mathbb{E}R_t + \frac{1}{50L^2\tau} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + \frac{1}{100\tau} \mathbb{E}R_t + 8\tau\gamma^2 \zeta^2 + \gamma^2\sigma^2 \\ &\leq \left(1 + \frac{3}{2\tau}\right) \mathbb{E}R_t + \frac{1}{50L^2\tau} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 + 8\tau\gamma^2 \zeta^2 + \gamma^2\sigma^2\end{aligned}$$

The lemma now follows by unrolling, and noting that  $\left(1 + \frac{3}{2\tau}\right)^j \leq 5$  for all  $0 \leq j \leq \tau$ .

# Bibliography I

-  Ziheng Cheng, Xinmeng Huang, and Kun Yuan.  
Momentum benefits non-iid federated learning simply and provably.  
*arXiv preprint arXiv:2306.16504*, 2023.
-  Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U Stich, and Ananda Theertha Suresh.  
Mime: Mimicking centralized stochastic algorithms in federated learning.  
In *NeurIPS*, 2021.
-  Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian U. Stich.  
A unified theory of decentralized SGD with changing topology and local updates.  
In *ICML 2020 - International Conference on Machine Learning*, pages 5381–5393, 2020.
-  Tao Lin, Sebastian U Stich, Kumar Kshitij Patel, and Martin Jaggi.  
Don't Use Large Mini-batches, Use Local SGD.  
In *ICLR - International Conference on Learning Representations*, 2019.

## Bibliography II



Sebastian U Stich.

Local SGD Converges Fast and Communicates Little.

In *ICLR - International Conference on Learning Representations*, 2018.

# Federated vs Personalized Learning

- **Federated**

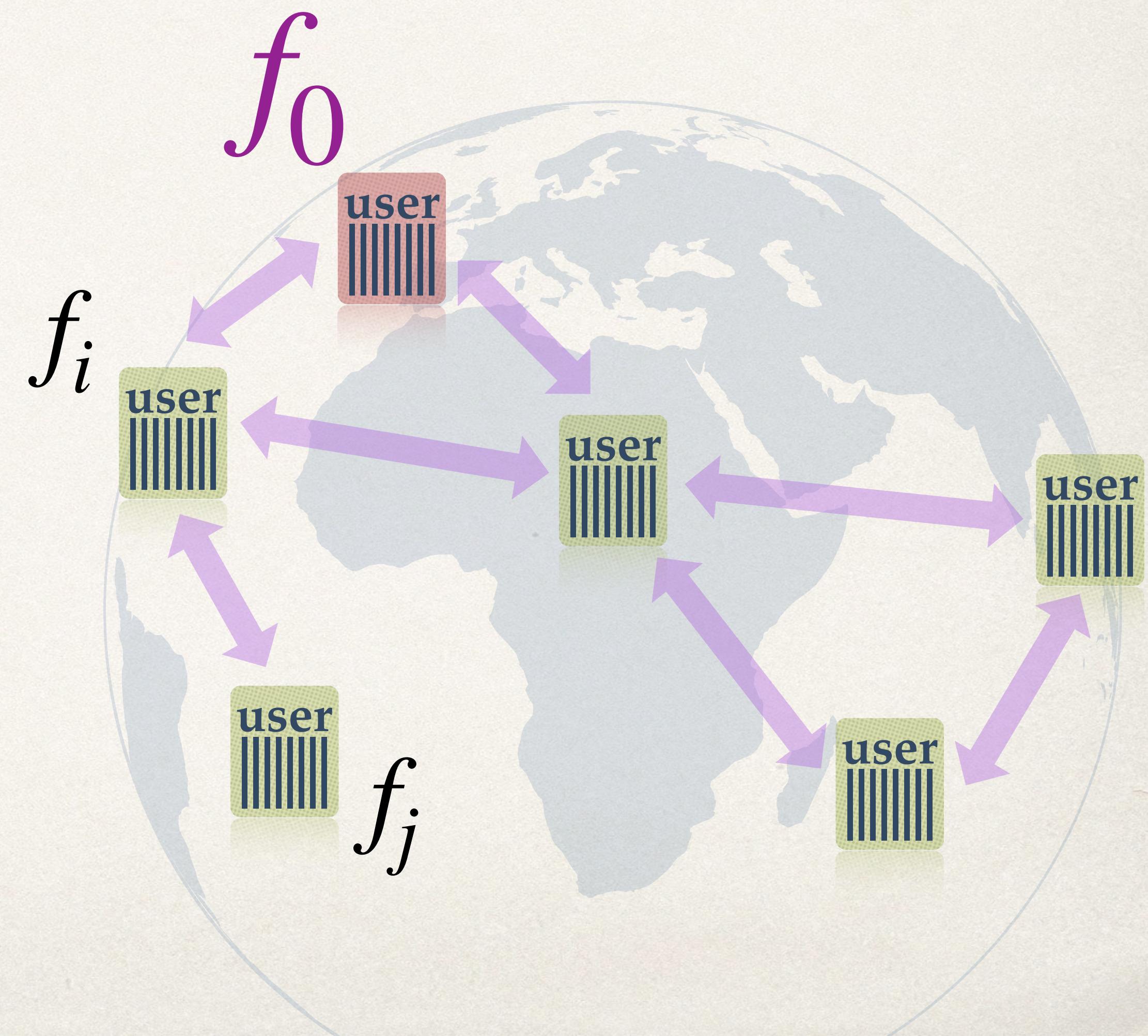
$$\min_x \frac{1}{n} \sum_i f_i(x)$$

- **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**

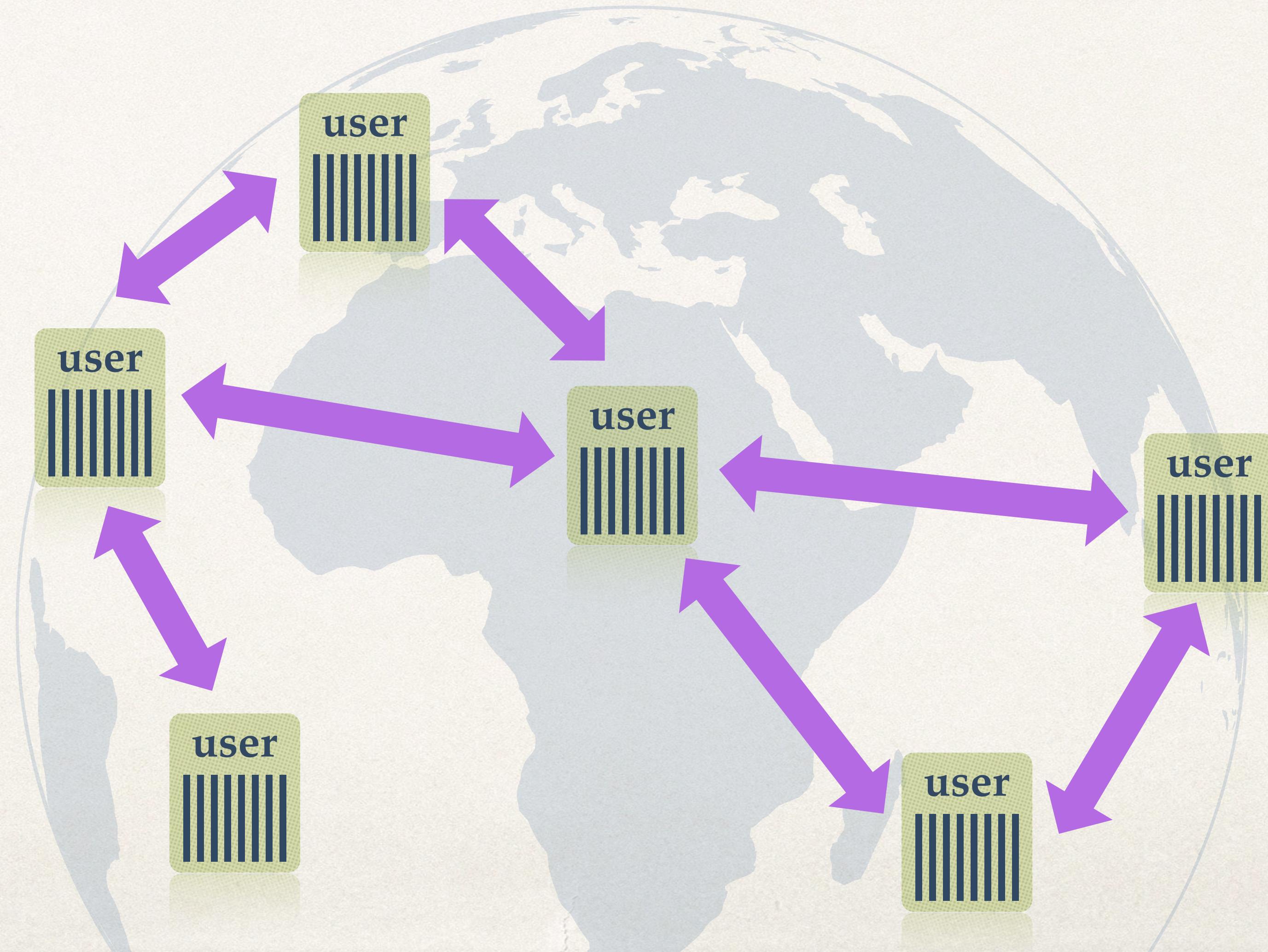
Sequential fine-tuning

Transfer learning,  
overparameterized models?

- ❖ **Train alone or collaborate?**

2c

# Decentralized Learning



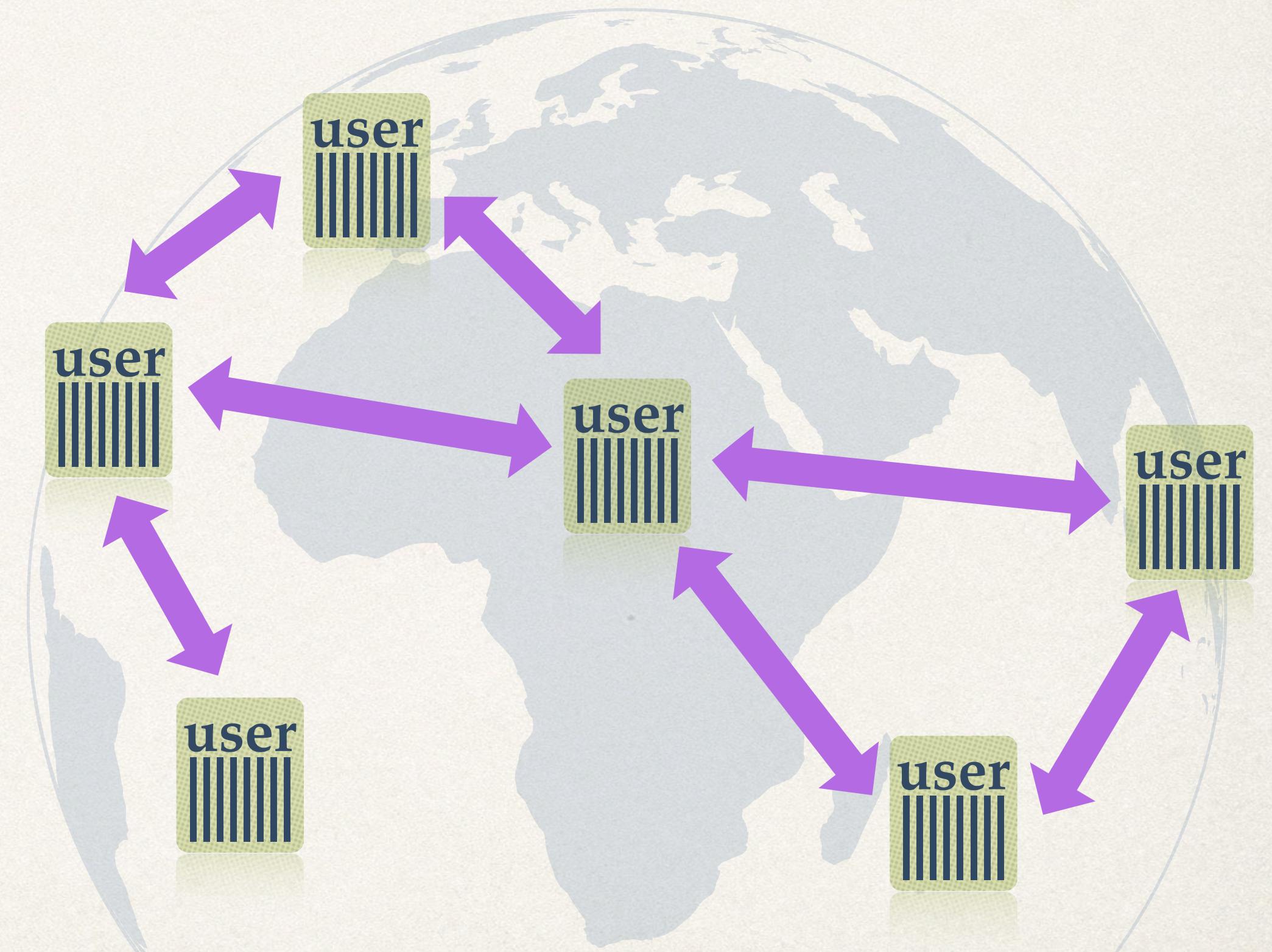
# Motivation

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



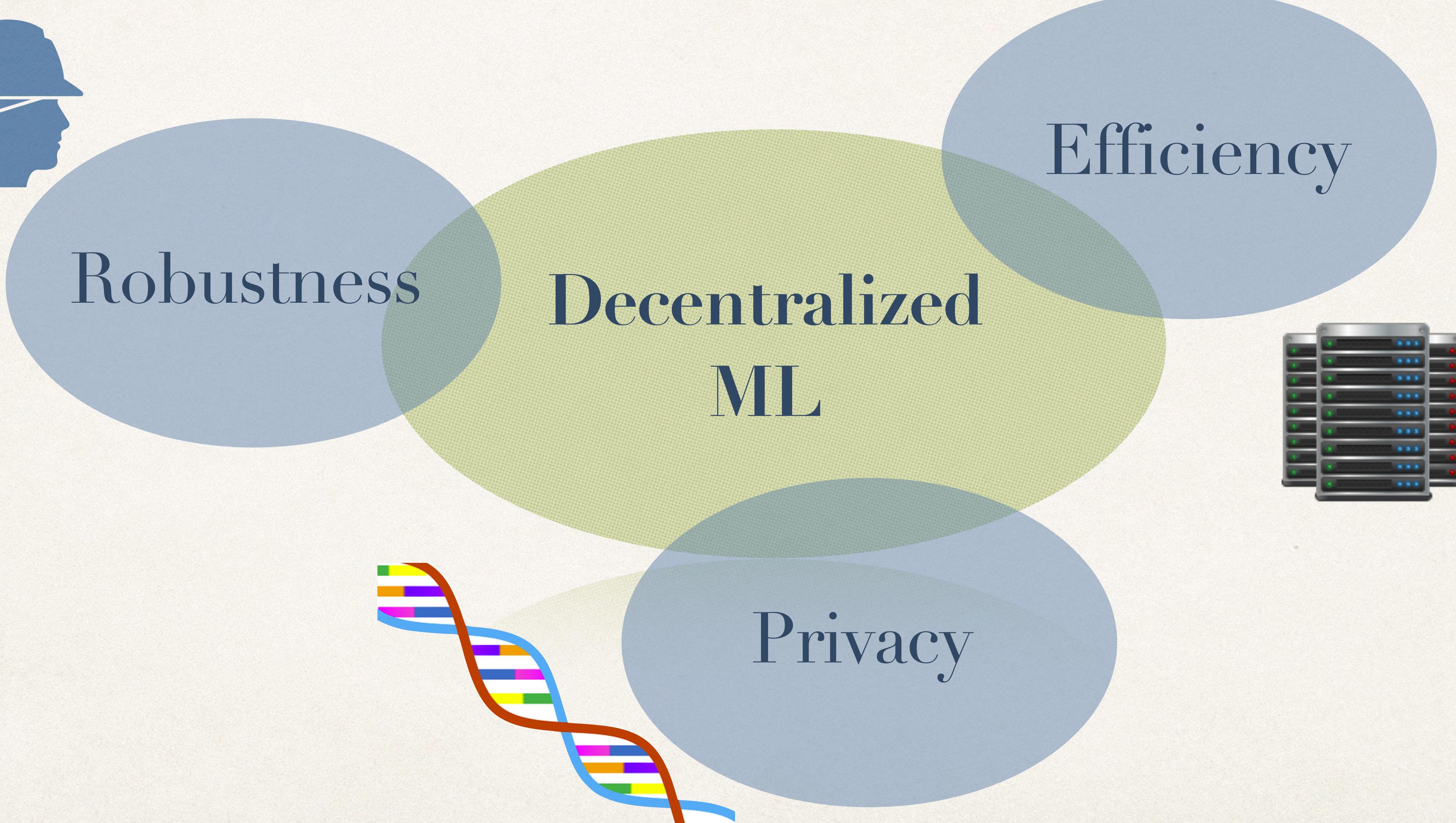
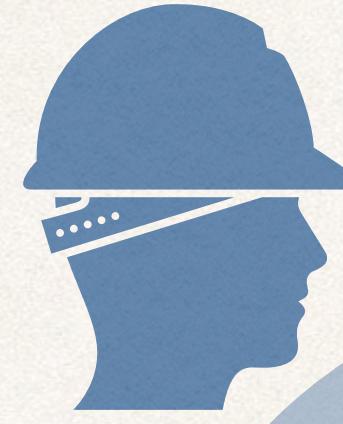
[image source](#)

- ❖ **Advantages:**

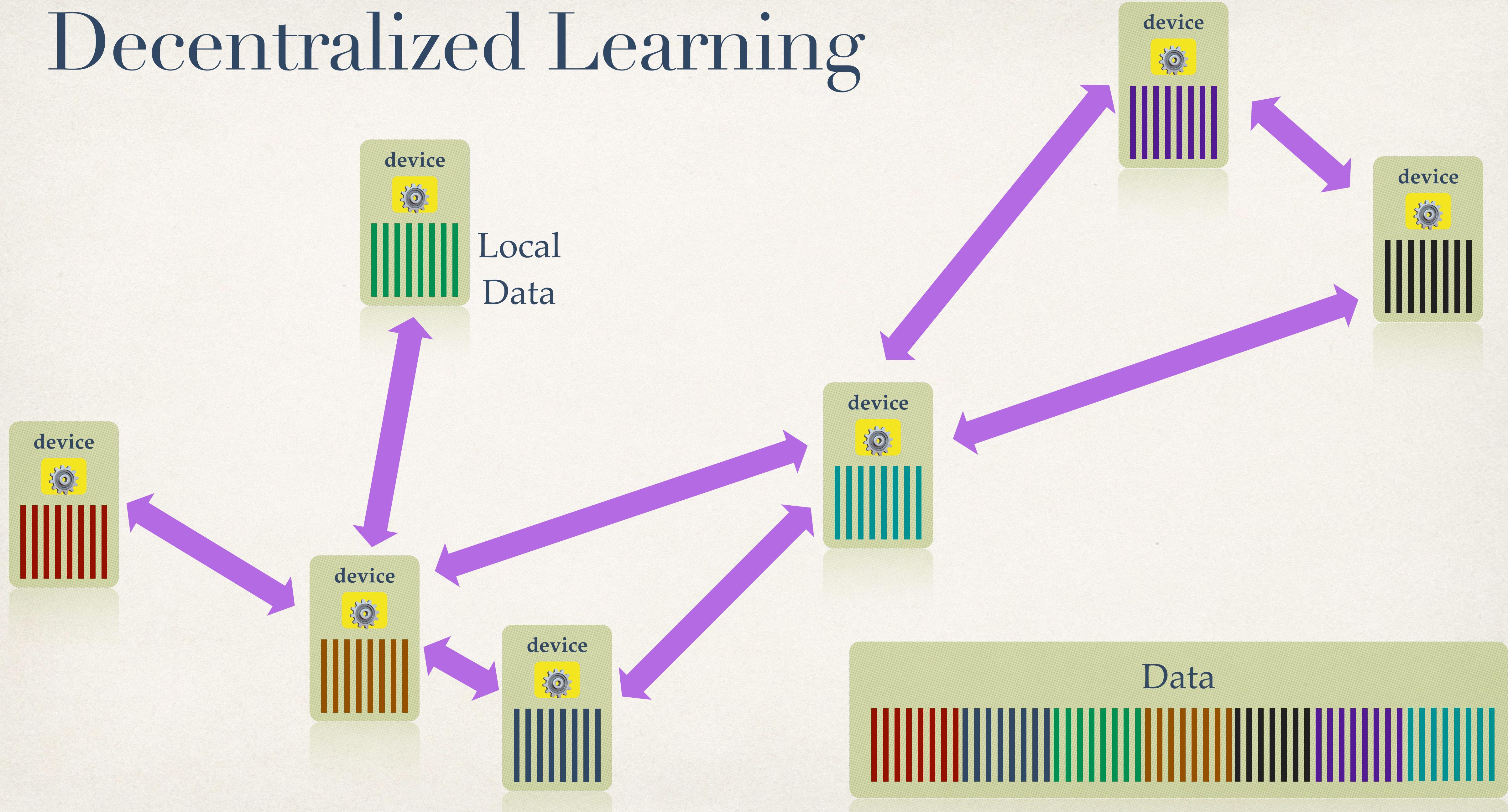


**AI utility, control and privacy aligned with data ownership**

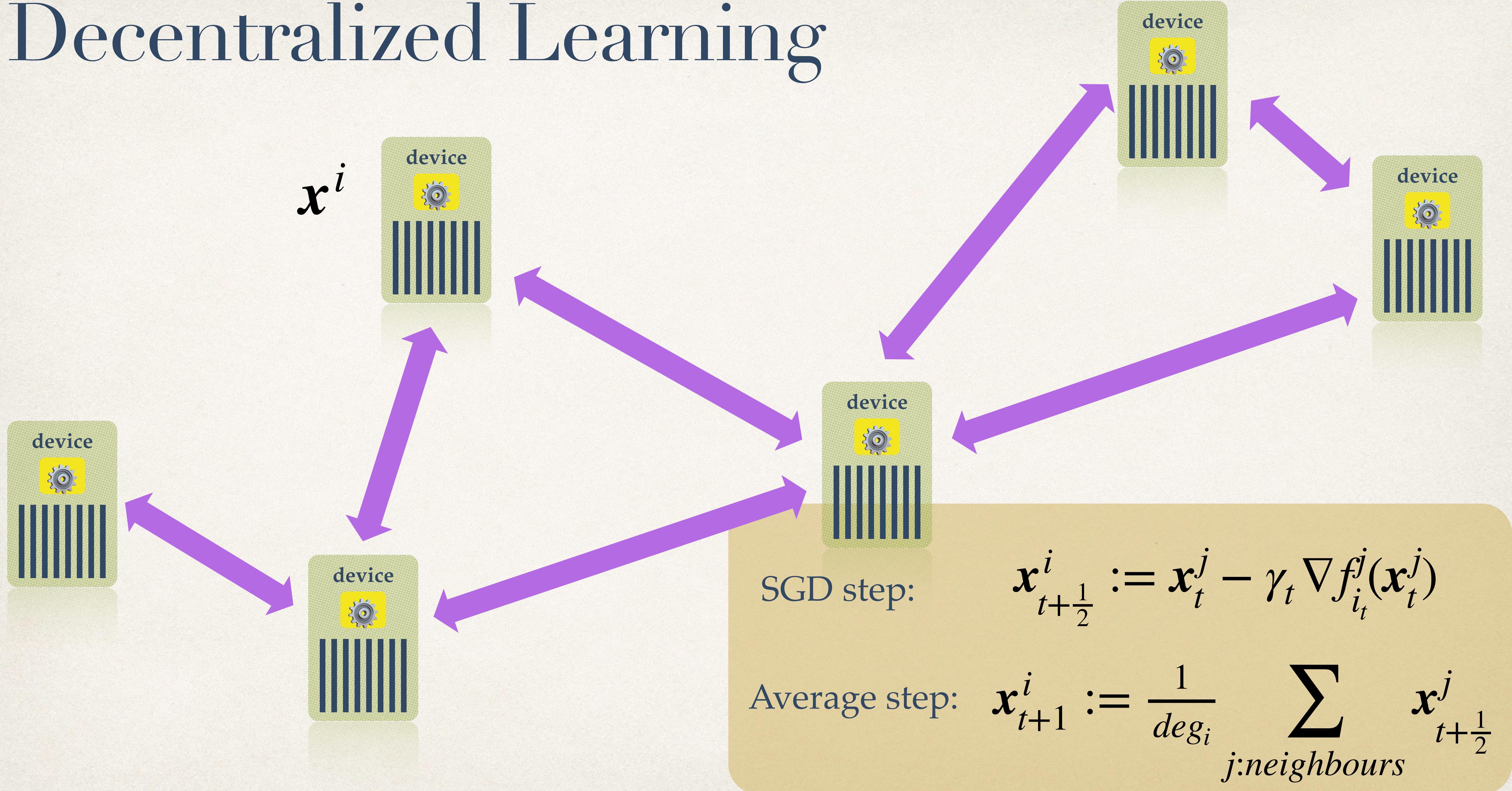
# Required Building Blocks



# Decentralized Learning



# Decentralized Learning



# Communication Compression

- ✿ limited-bit precision vector

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

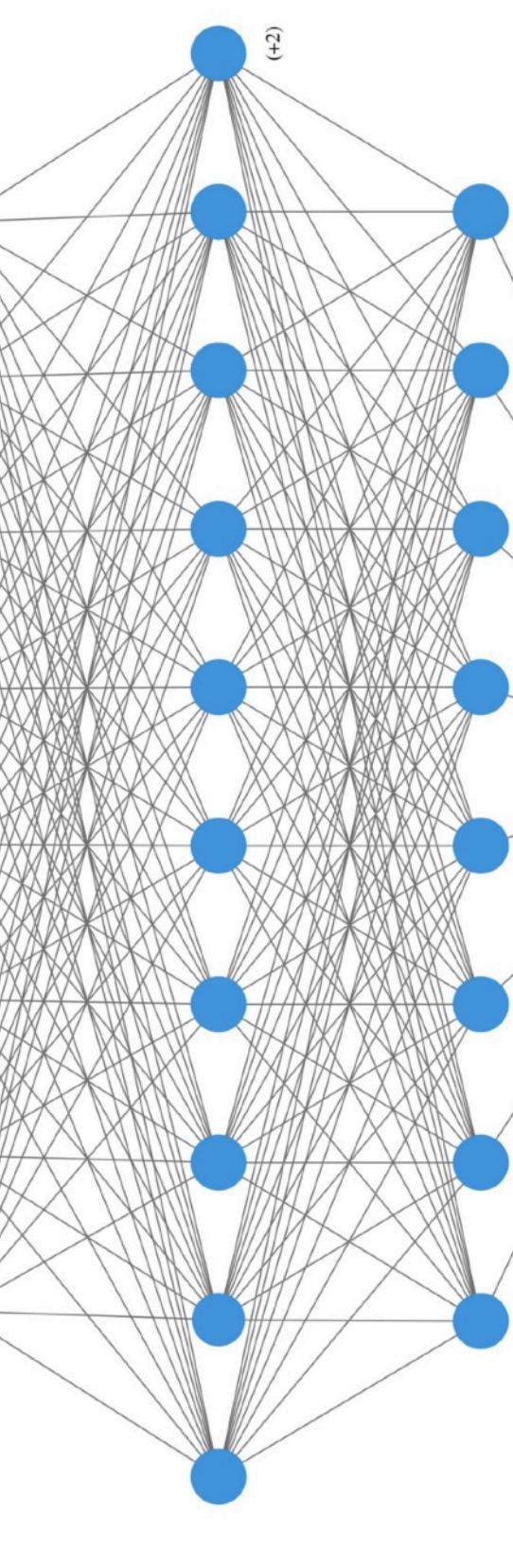
- ✿ random / top k% of all the entries

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

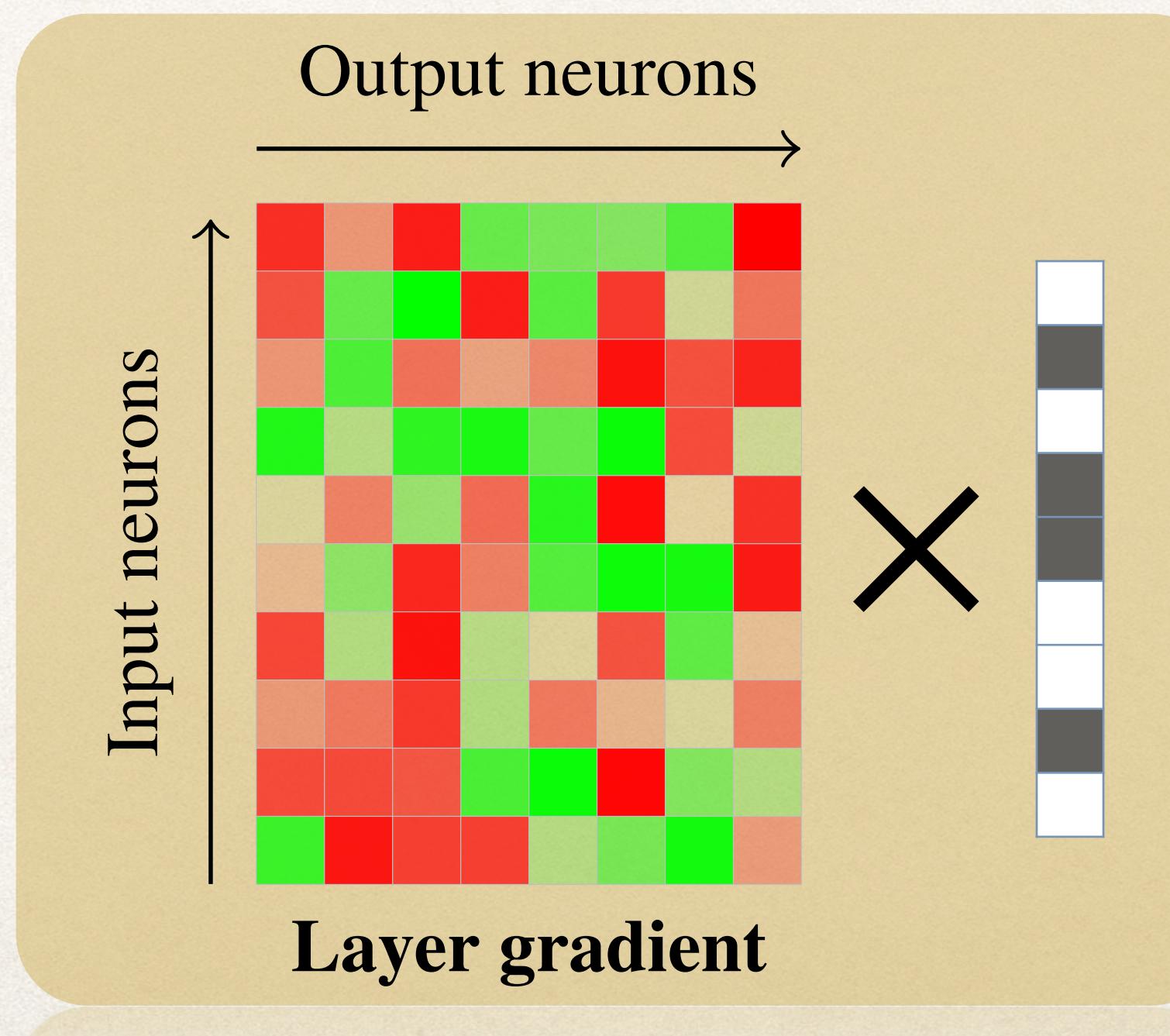
- ✿ low rank version of the gradient?

# Low-Rank Communication Compression

- PowerSGD



backprop is fast:  
linear time



fast compression?

**Fast power iterations**

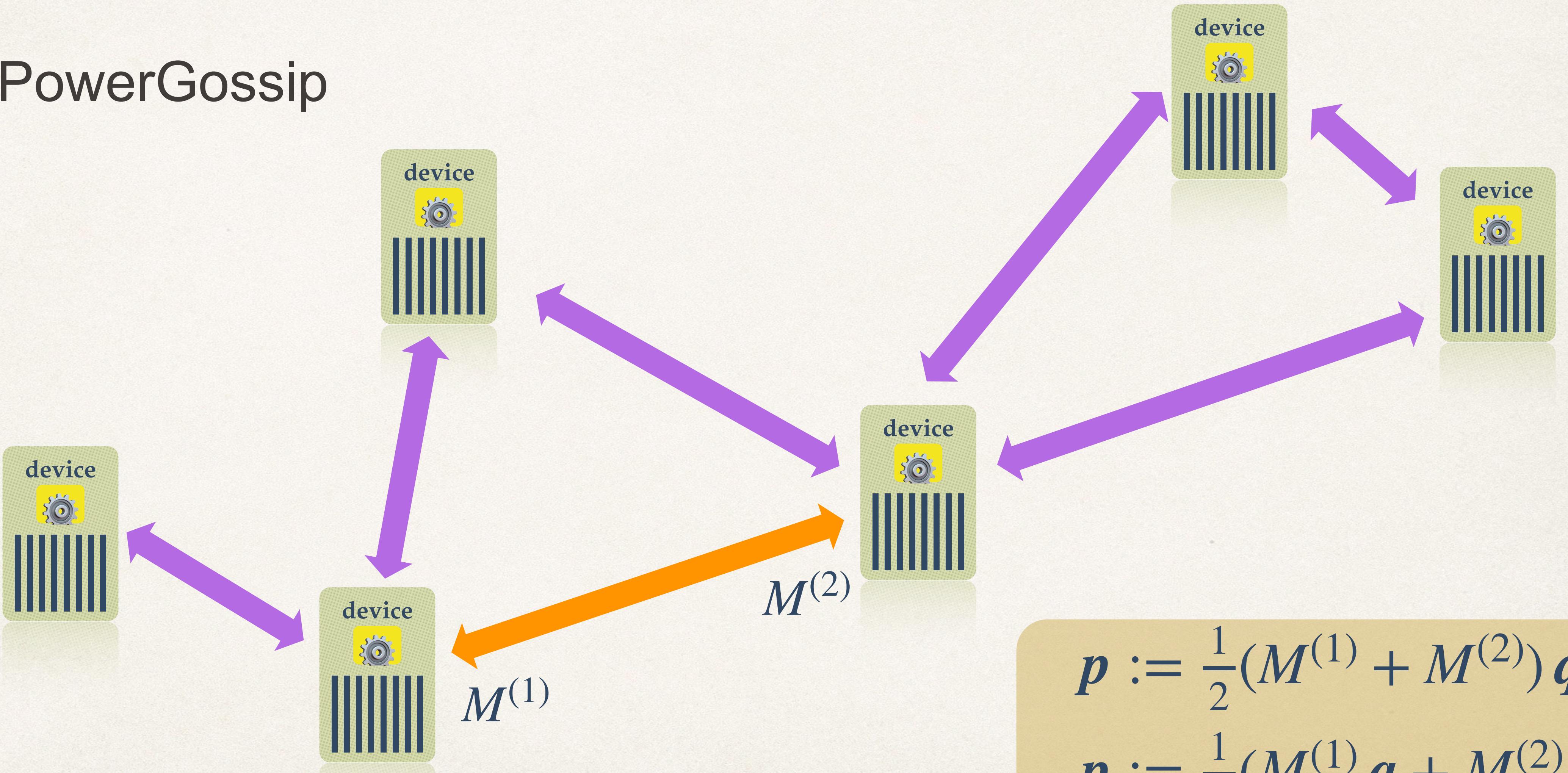
$$p := G q$$

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

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

# Decentralized Learning with Compression

- PowerGossip



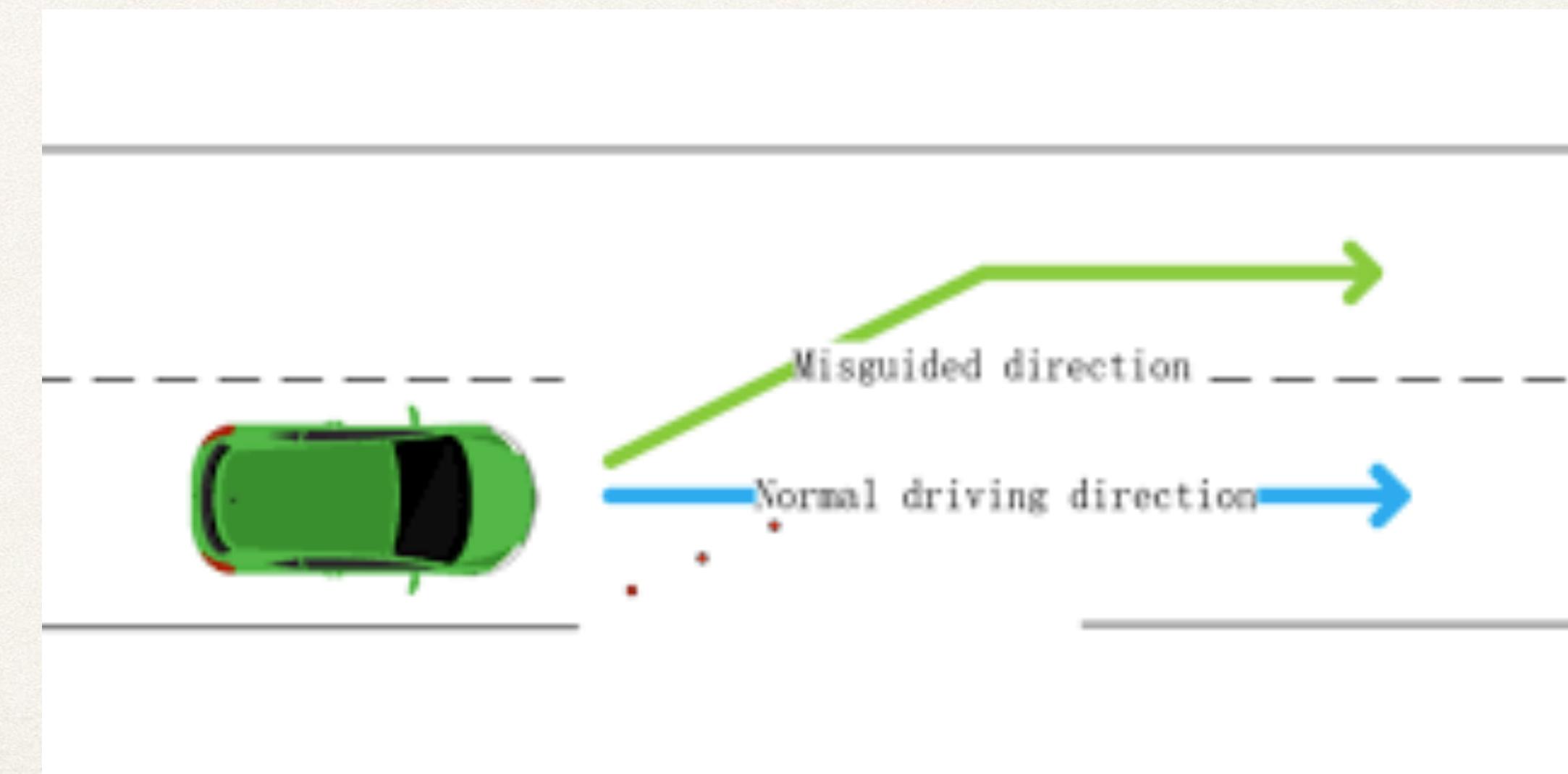
# Building Blocks for Decentralized ML

- ✿ **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

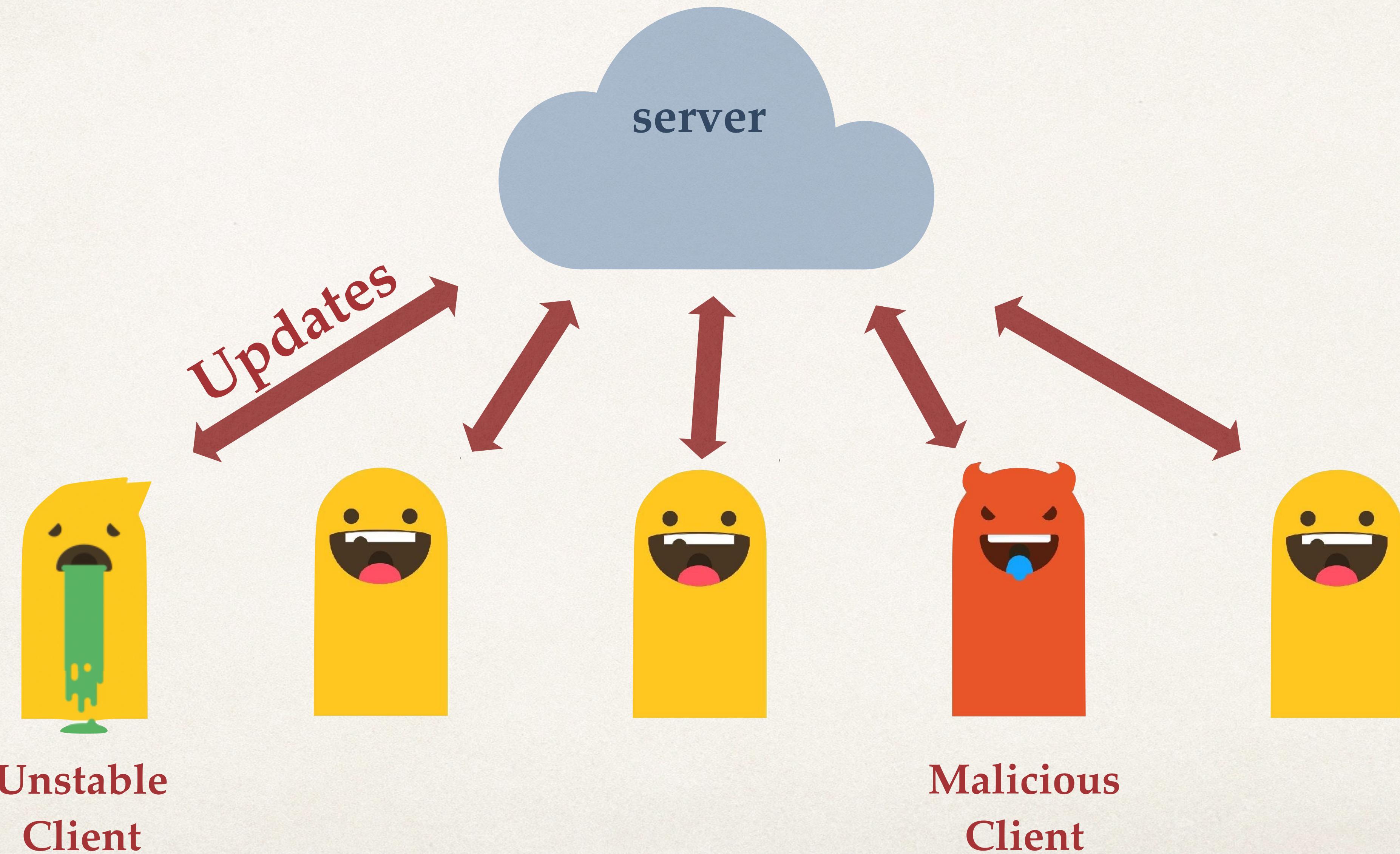


3a

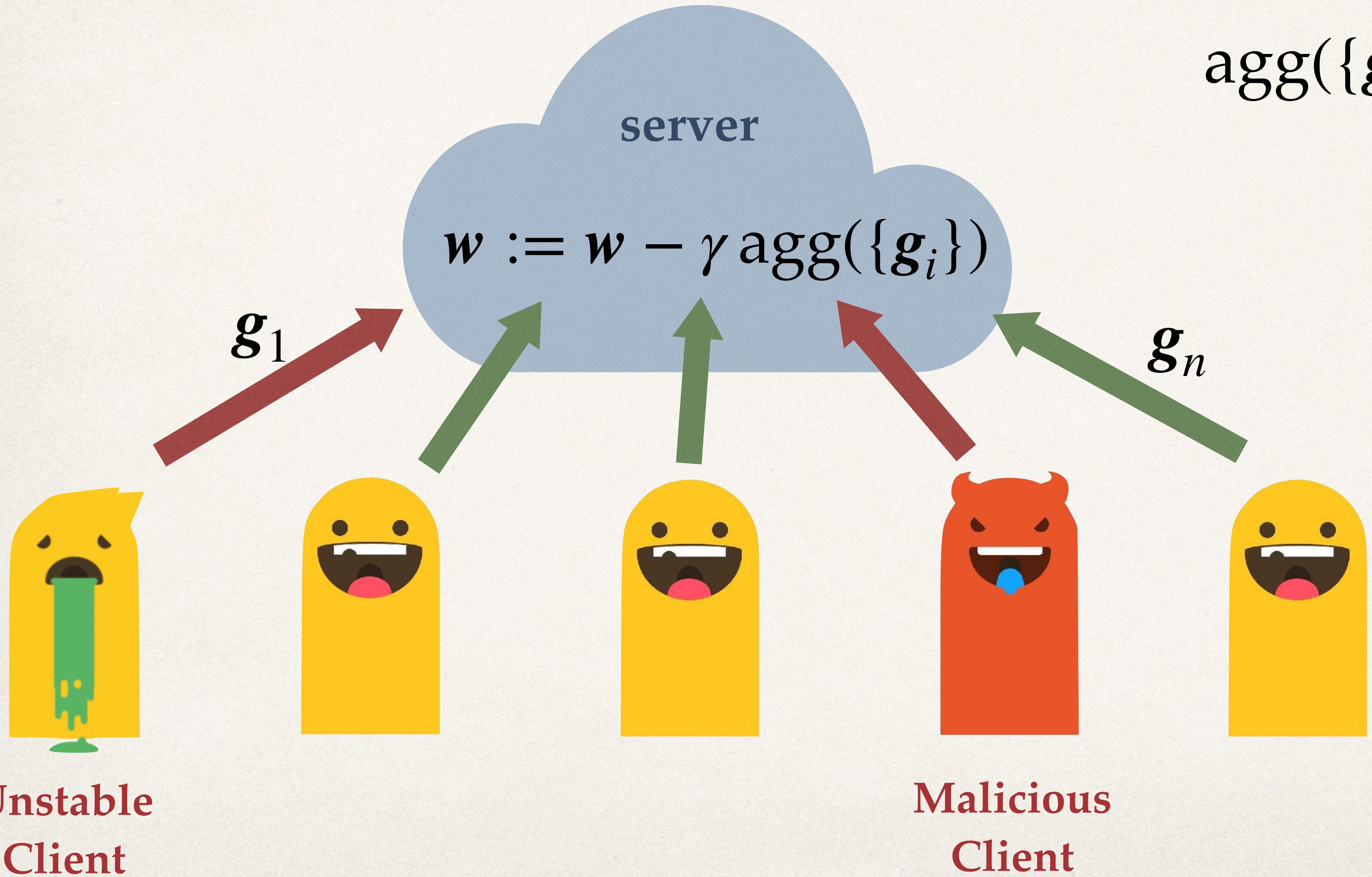
Gradients from  
faulty/malicious collaborators:

- 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

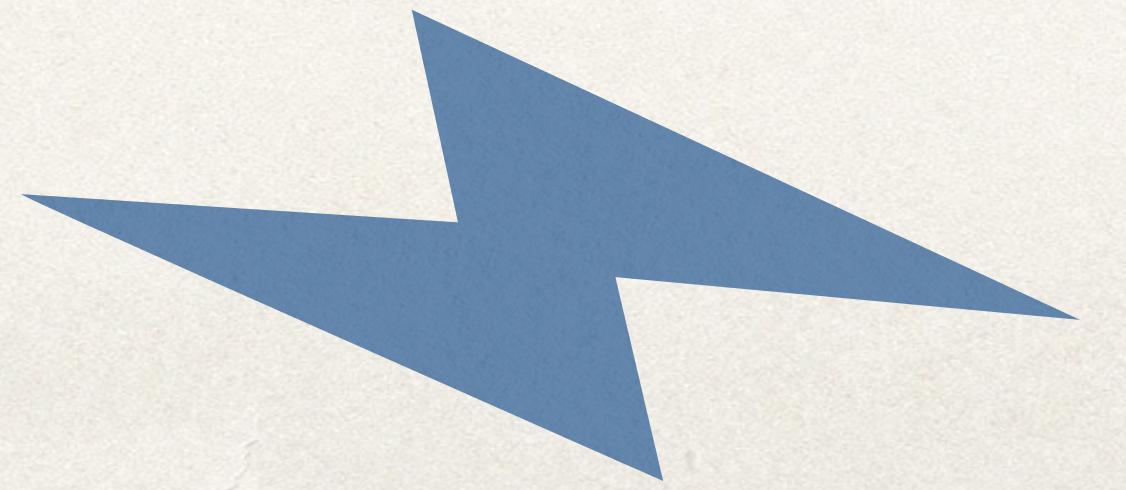
$$\text{robust-mean}_i f_i(\mathbf{x}) = \frac{1}{|good|} \sum_{i \in good} f_i(\mathbf{x})$$

Federated

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

Fairness

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



# 3b

## Adversarial Attacks (at inference time)

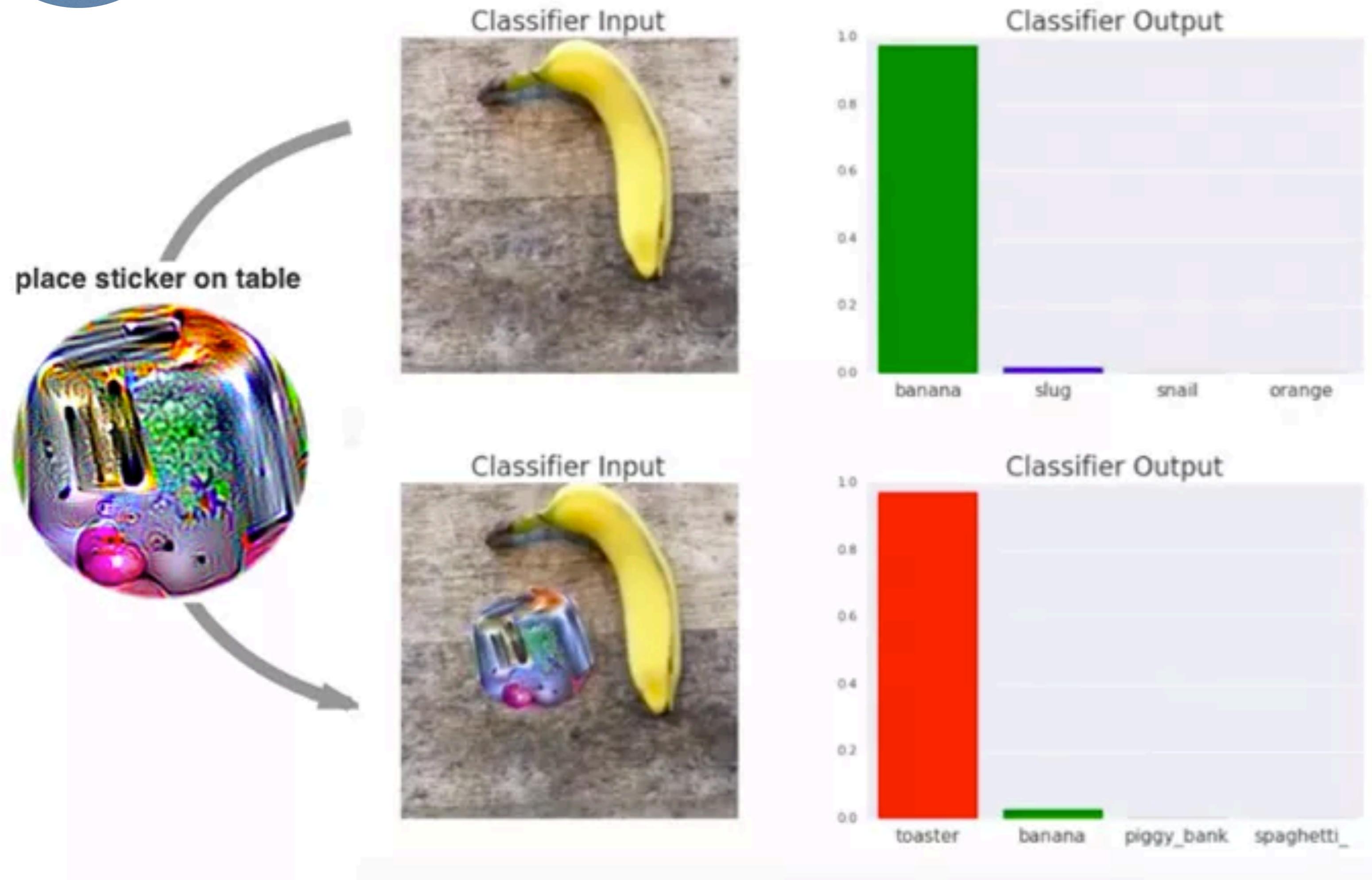


Image: [Tom B. Brown/Dandelion Mané](#)



Image: Elsayed ,Papernot et al 2018

# Adversarial Attacks (at inference time)

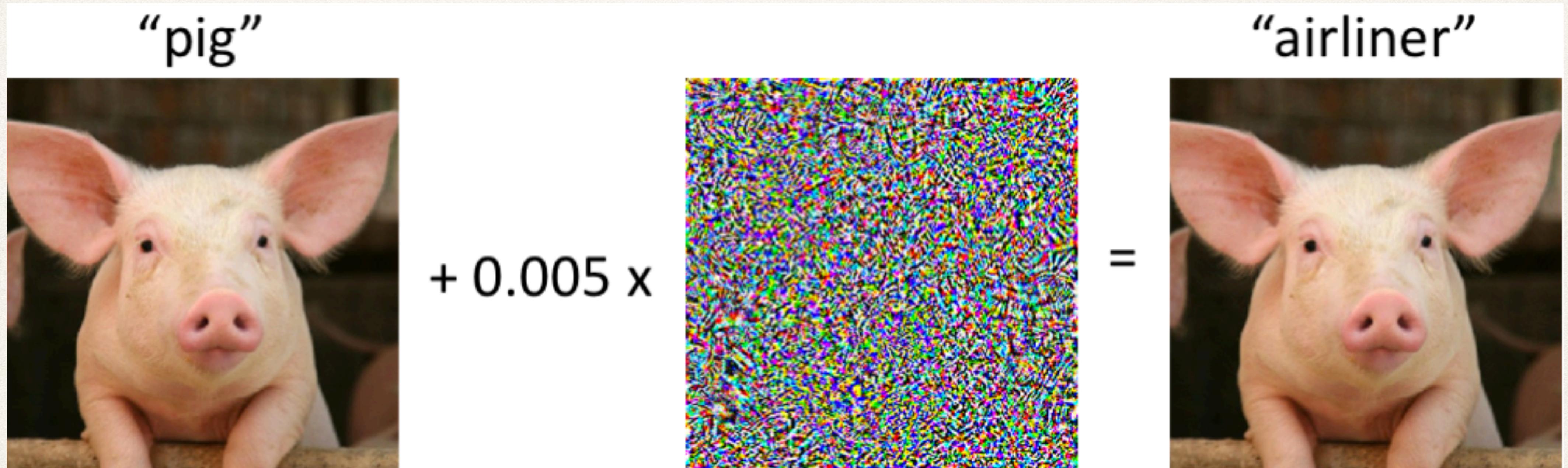


Image: [Mądry, Schmidt](#)

More info:  
[http://gradientscience.org/intro\\_adversarial/](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

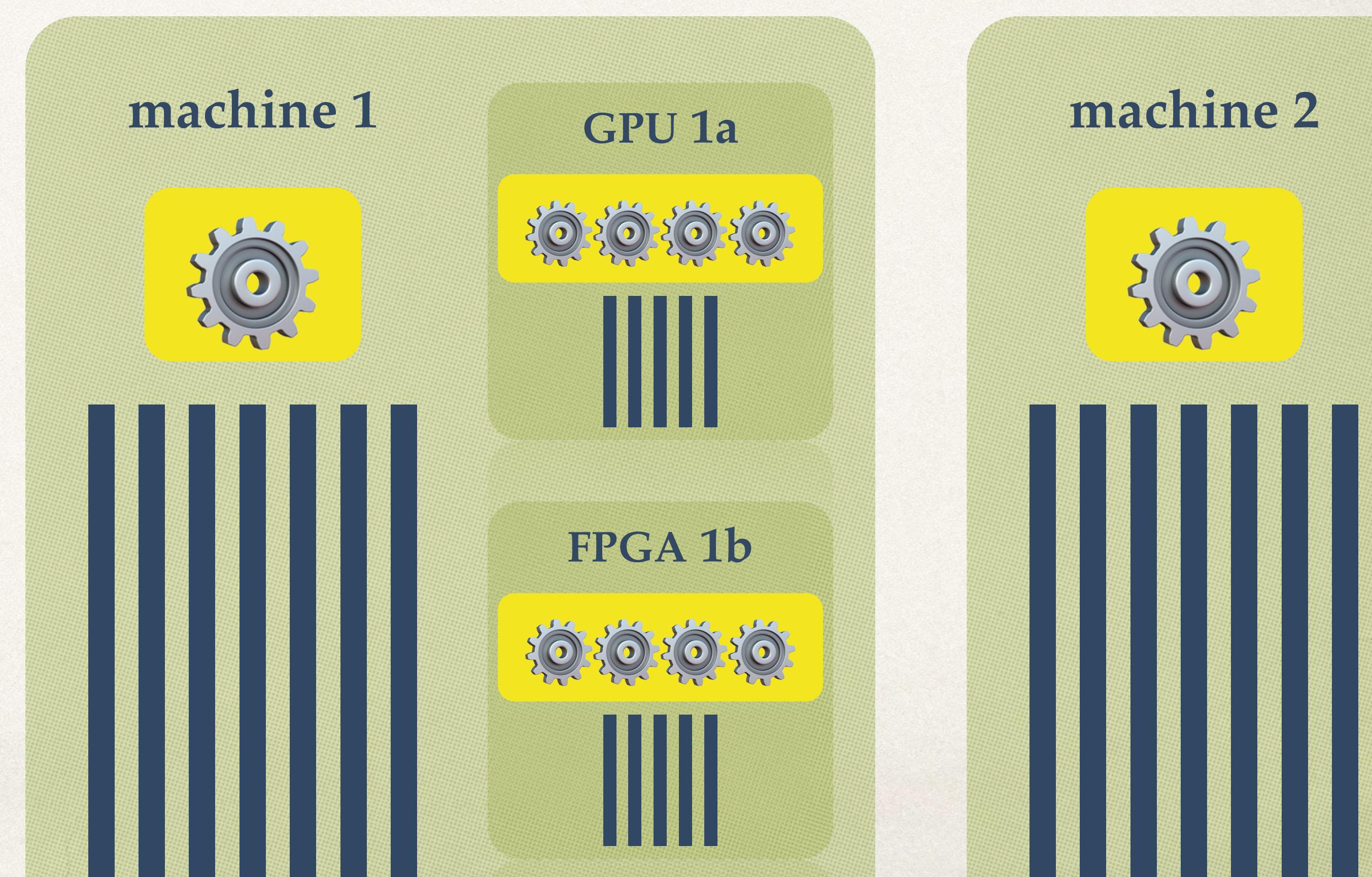
## Privacy

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

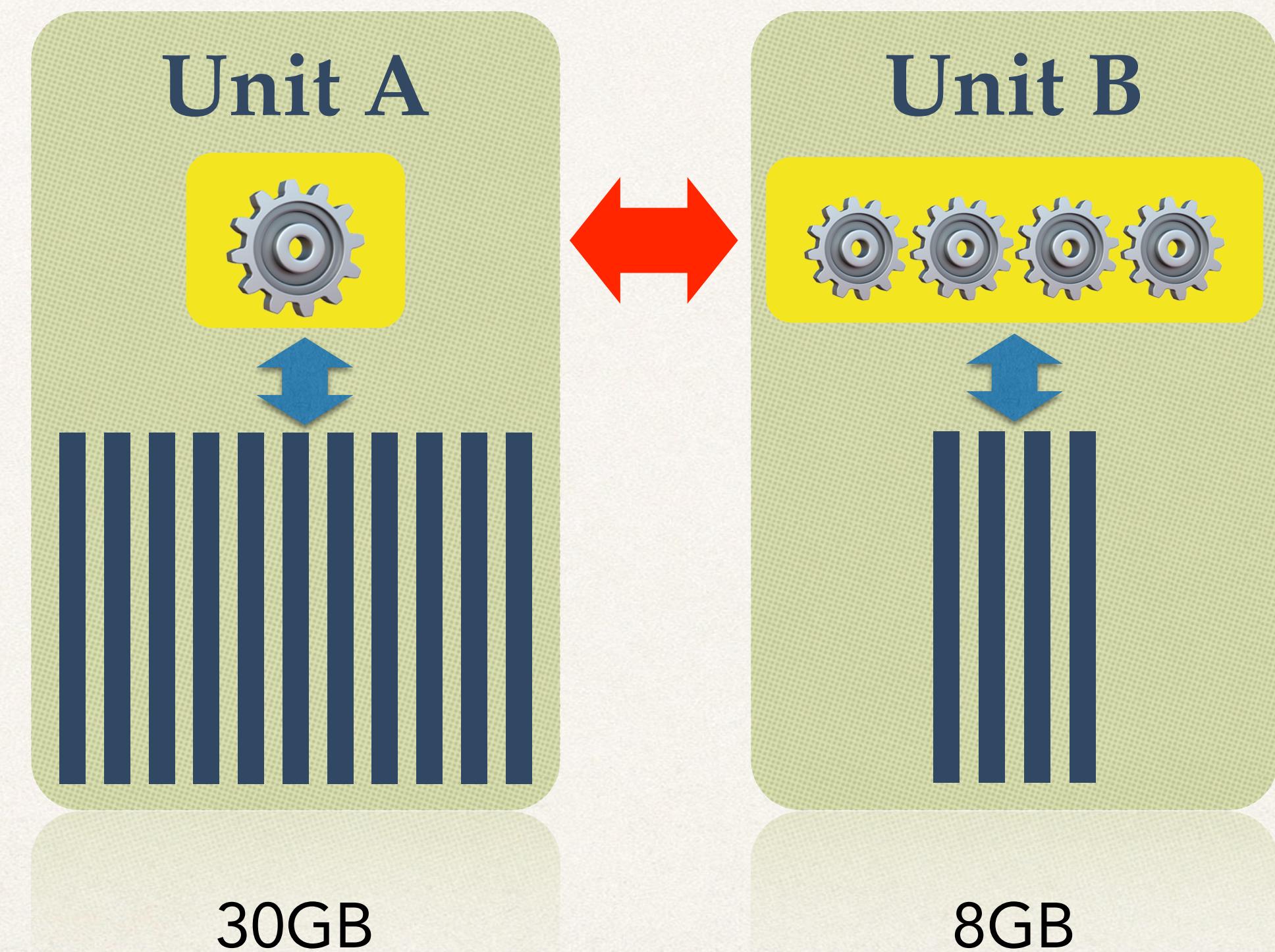
## 5

# Leveraging Heterogenous Systems

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



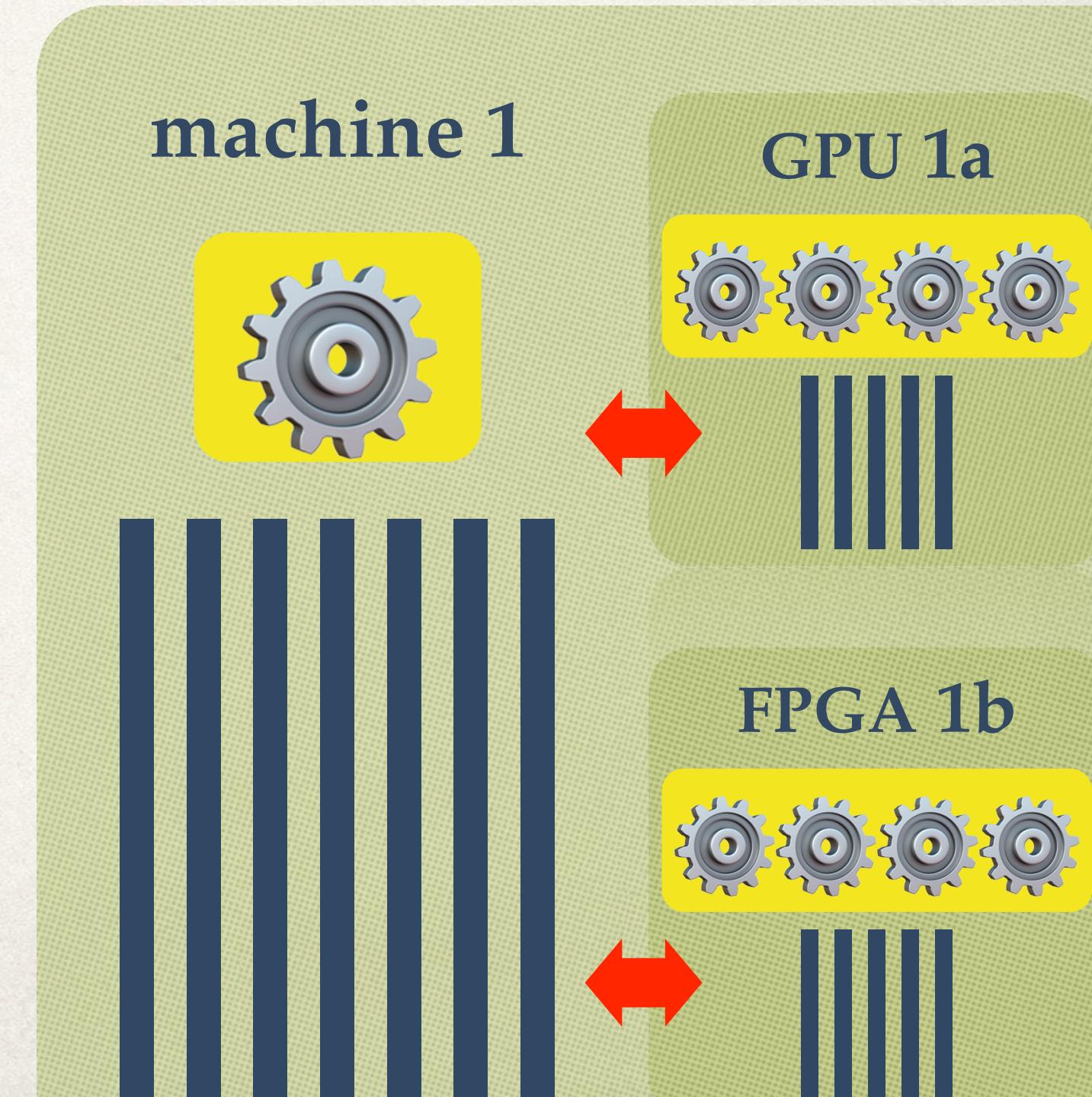
# Leveraging Heterogenous Systems



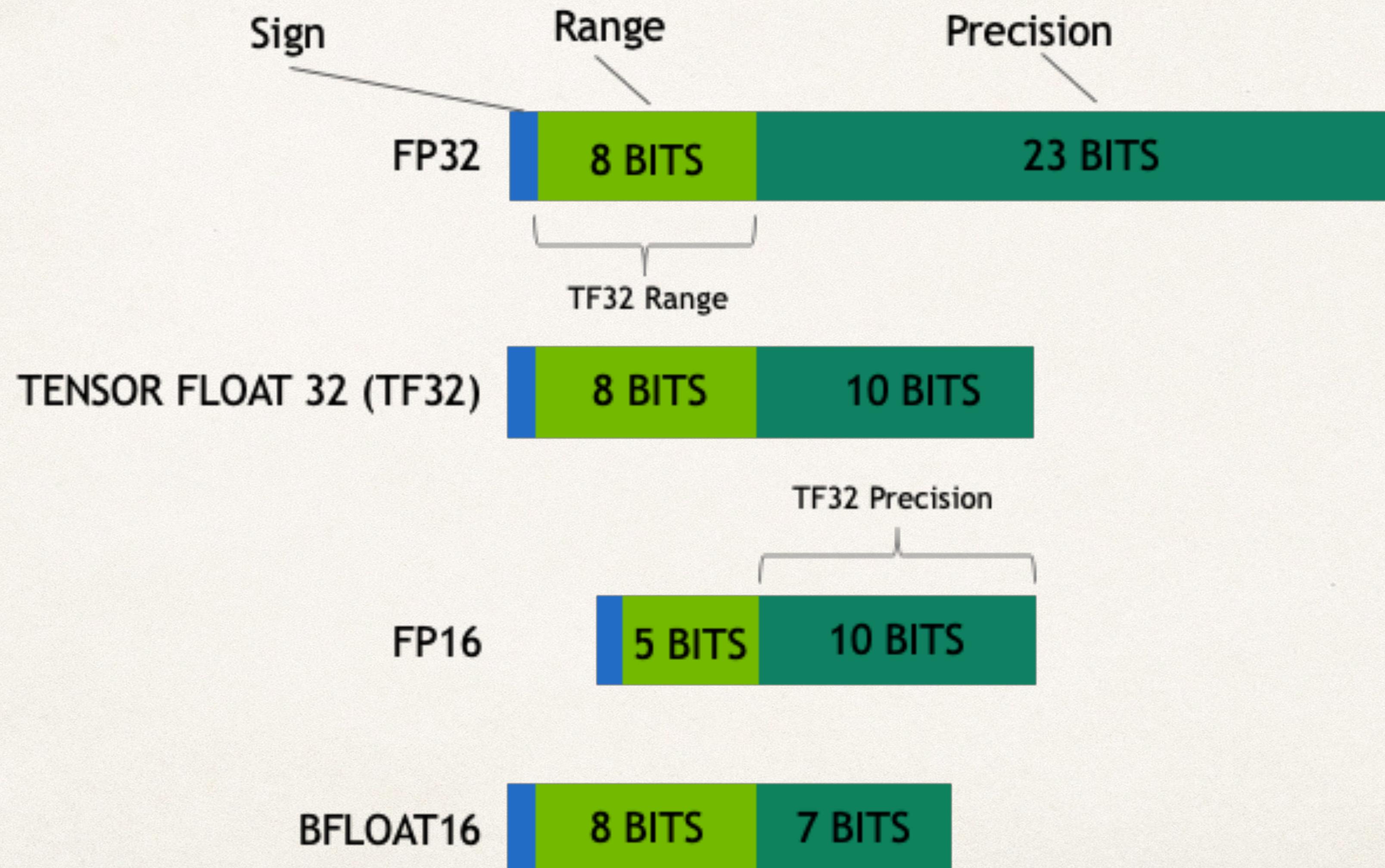
adaptive importance sampling of datapoint  
e.g. for general linear models, or word2vec

# 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



Thanks!

mlo.epfl.ch