# Learning with Strange Gradients

Martin Jaggi
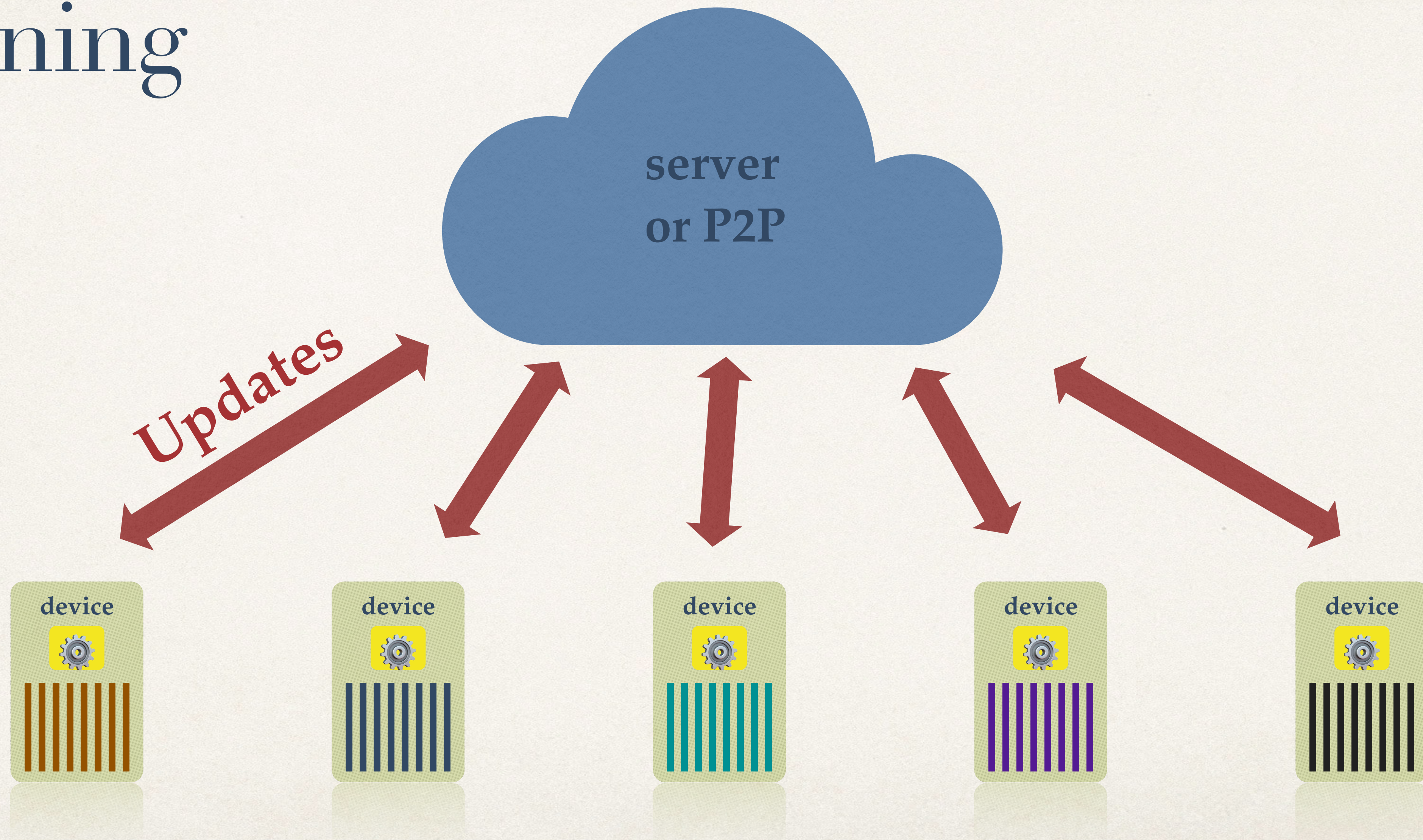
**EPFL**

*Machine Learning and Optimization Laboratory*

*mlo.epfl.ch*

# Collaborative & Federated Training

Data

server or P2P

Updates

device
device
device
device
device

**1** Gradients from strange collaborators:
   - **Federated Learning**

**2** Gradients from strange collaborators:
   - **Personalization**

**3** Gradients from strange **architectures**

**4** Gradients from faulty/**malicious collaborators**:
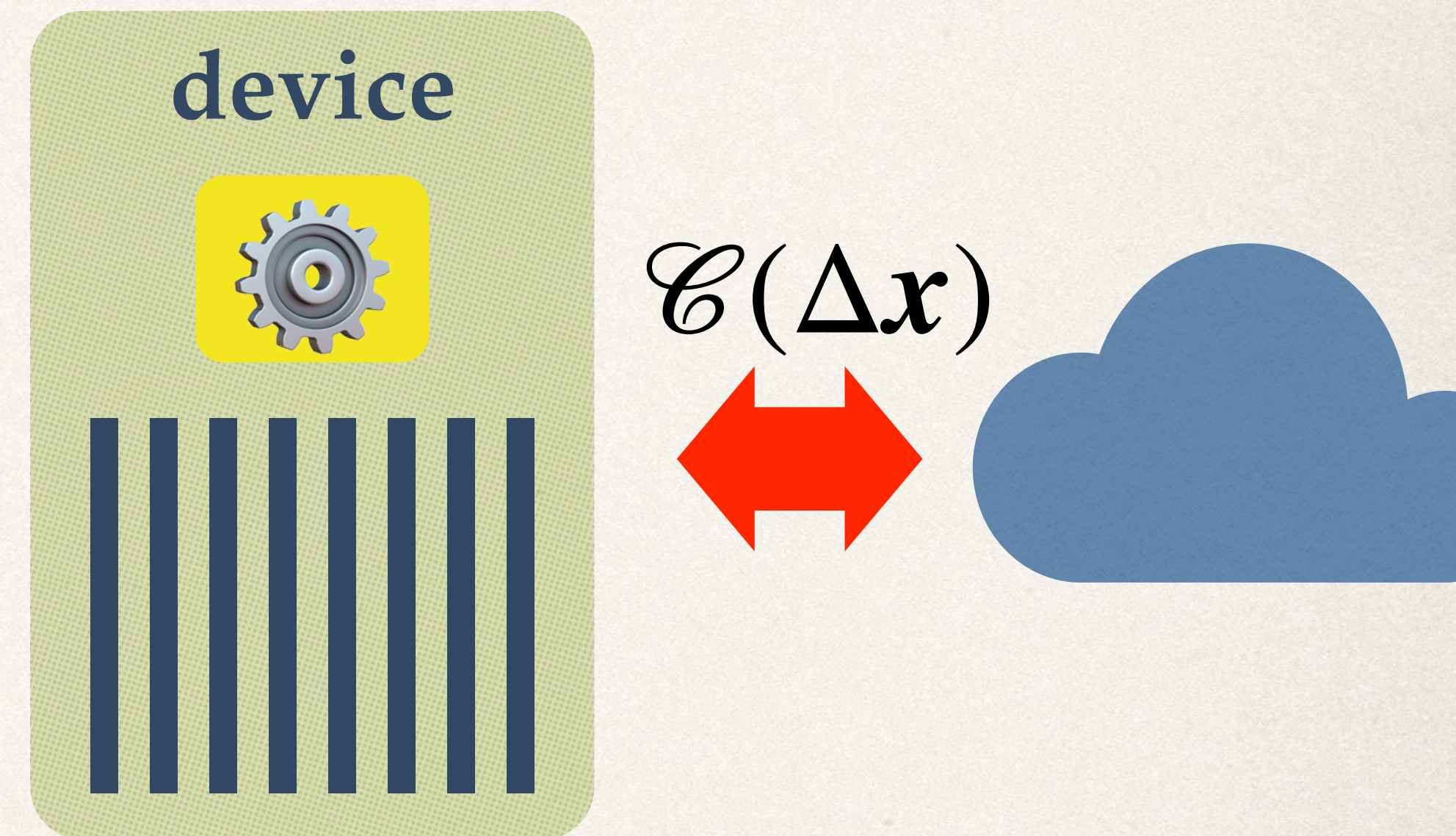   - Byzantine-robust Training

# Stochastic Gradient Descent (SGD)

$$\min_{x} \; f(x) = \frac{1}{|data|} \sum_{i \,\in\, data} f_i(x)$$

$$i_t \sim \mathrm{Uniform}(1, |data|)$$

$$x_{t+1} := x_t + \Delta x$$

**device**

$$\mathscr{C}(\Delta x)$$

$$\Delta x = -\gamma_t \nabla f_{i_t}(x_t) \quad \textbf{from backpropagation}$$

Gradients from strange
collaborators:
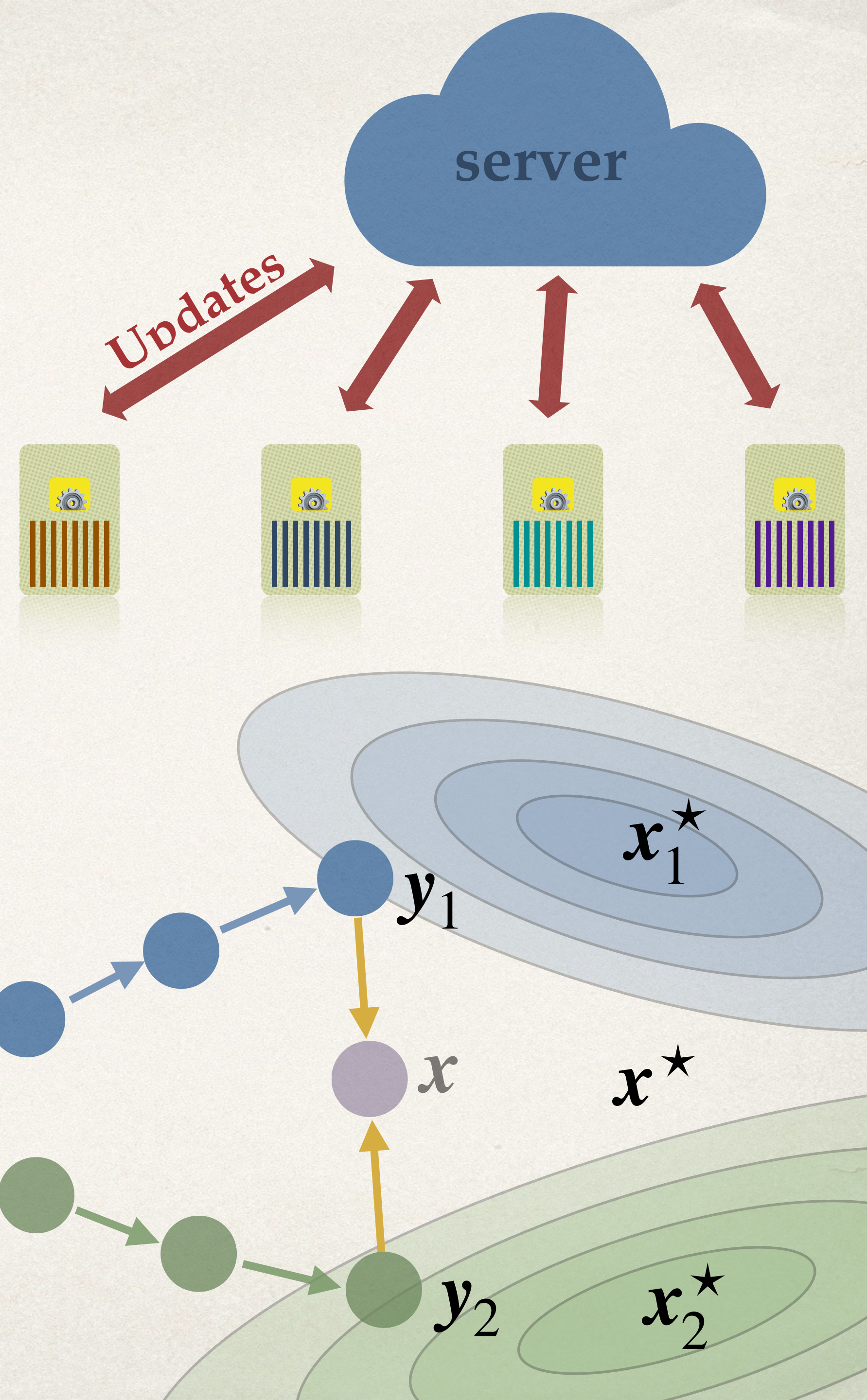        - Federated Learning

# Client drift



✤ **Federated Learning**

$$\min_{x} \; \frac{1}{n} \sum_{i}^{n} f_i(x)$$
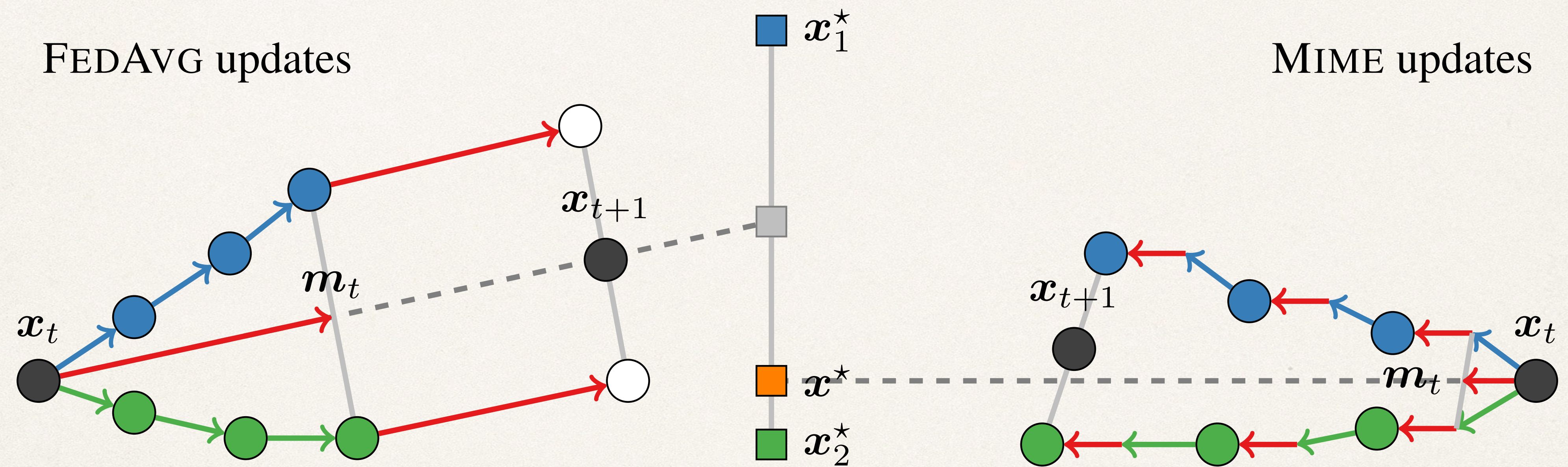
✤ **Fed Avg / Local SGD**

*for* some local steps

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

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

# Client drift

# Mime algorithm framework

*for* some local steps

$$y_i := y_i - \eta\big((1-\beta)\,\nabla f_i(y_i) + \beta m\big)$$

$$m := (1-\beta)\,\nabla f_i(x) + \beta m$$

*aggregated on server after each round*

# Mime convergence

**Number of rounds to reach**

$$\mathbb{E}\left[\|\nabla f(\boldsymbol{x}^{out})\|^2\right] \leq \varepsilon \ :$$

$$\mathscr{O}\left(\left(\frac{n}{S}\right)^{3/2}\frac{L}{\varepsilon}\right) \qquad \textbf{Scaffold}$$

$$\mathscr{O}\left(\frac{\delta(\zeta+\sigma)}{\varepsilon^{3/2}} + \frac{\zeta^2+\sigma^2}{\varepsilon} + \frac{\delta}{\varepsilon}\right) \qquad \textbf{MimeLiteMVR}$$

$$\mathscr{O}\left(\frac{\delta\zeta}{\sqrt{S}\varepsilon^{3/2}} + \frac{\zeta^2}{S\varepsilon} + \frac{\delta}{\varepsilon}\right) \qquad \textbf{MimeMVR}$$

$$\Omega\left(\frac{L\zeta}{\sqrt{S}\varepsilon^{3/2}} + \frac{\zeta^2}{S\varepsilon} + \frac{L}{\varepsilon}\right) \qquad \begin{array}{l}\textbf{Lower bound}\\ \textbf{(server-only)}\end{array}$$

**Data Heterogeneity:**

$\delta \ll L$ inter-cl. Hessian similarity

$\zeta$     inter-cl. gradient variance
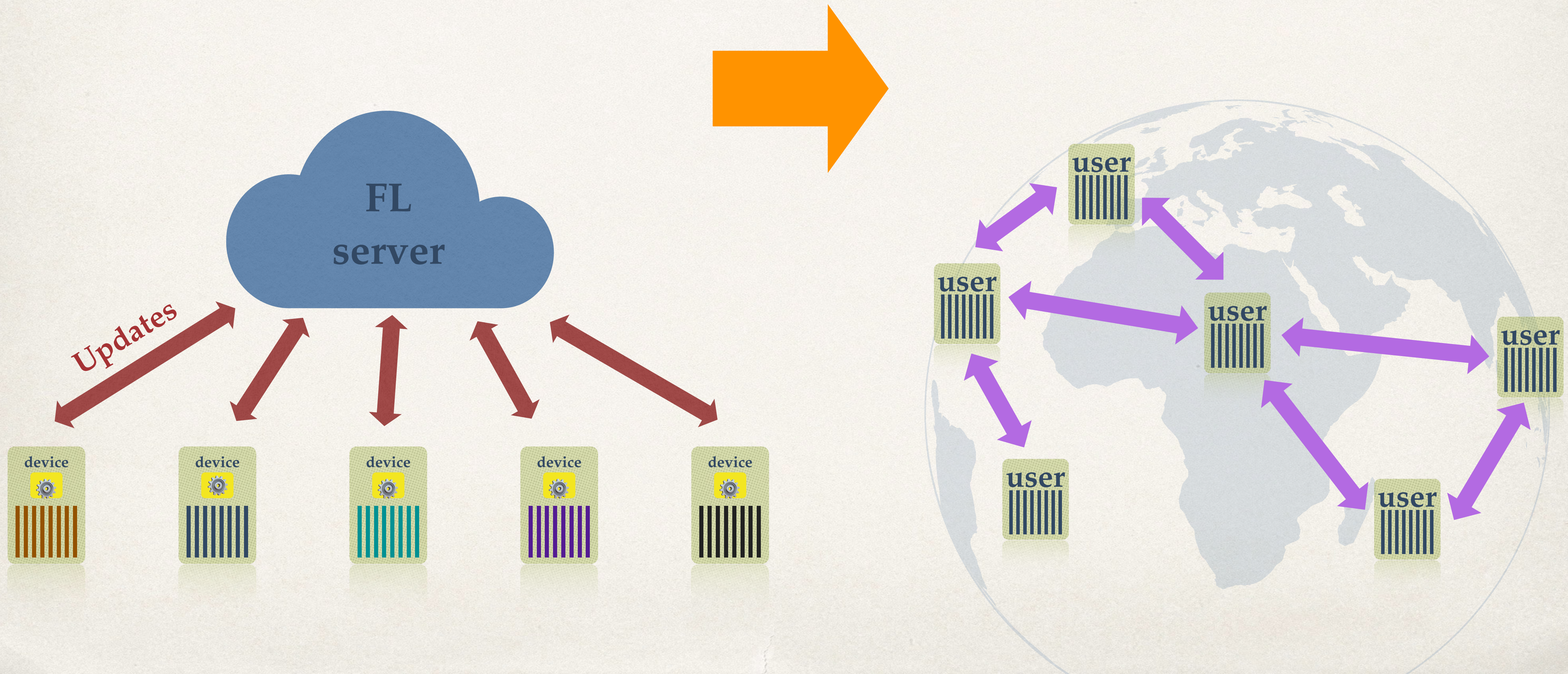
$\sigma$     intra-cl. gradient variance

Gradients from strange collaborators:
       - Personalization

# From **federated** towards **decentralized**
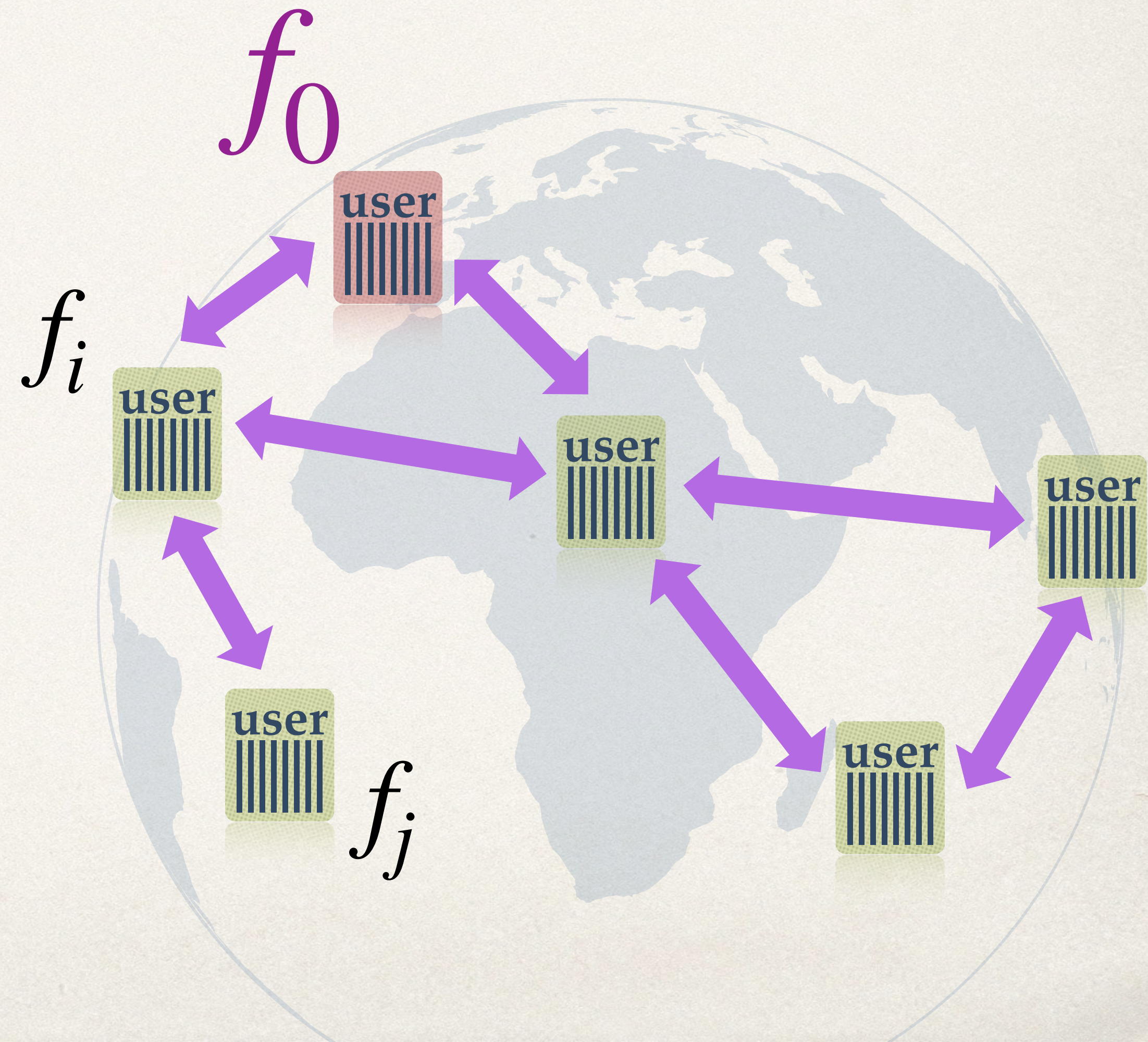
# Collaborative Learning

✤ **Federated**

$$\min_{\boldsymbol{x}} \; \frac{1}{n} \sum_{i}^{n} f_i(\boldsymbol{x})$$

✤ **Collaborative / Personalized**

$$\min_{\boldsymbol{x}} \; f_0(\boldsymbol{x}) \quad \Bigg| \quad \begin{array}{c} \min_{\boldsymbol{x}} \; f_1(\boldsymbol{x}) \\[1em] \min_{\boldsymbol{x}} \; f_n(\boldsymbol{x}) \end{array}$$

# Personalized learning / optimization

✤ **Weighted averaging**

$$x := x - \gamma \sum_{i=0}^{n} \alpha_i \nabla f_i(x)$$

✤ **Weighted averaging with bias correction**

$$x := x - \gamma \sum_{i=0}^{n} \left( \alpha_i \nabla f_i(x) + c_i \right)$$

*idea similar to Scaffold*

$f_0$

$f_i$

$f_j$

**Theorem:** Convergence on personal objective $f_0$
for non-convex smooth objectives,
using exponential moving average to learn $\boldsymbol{c}_i$

$$\mathbb{E}\left[\|\nabla f_0(\boldsymbol{x}^{out})\|^2\right] = \mathcal{O}\left(\sqrt{\frac{LF_0\,\sigma_0^2}{(\boldsymbol{n}+1)T}}\right)$$
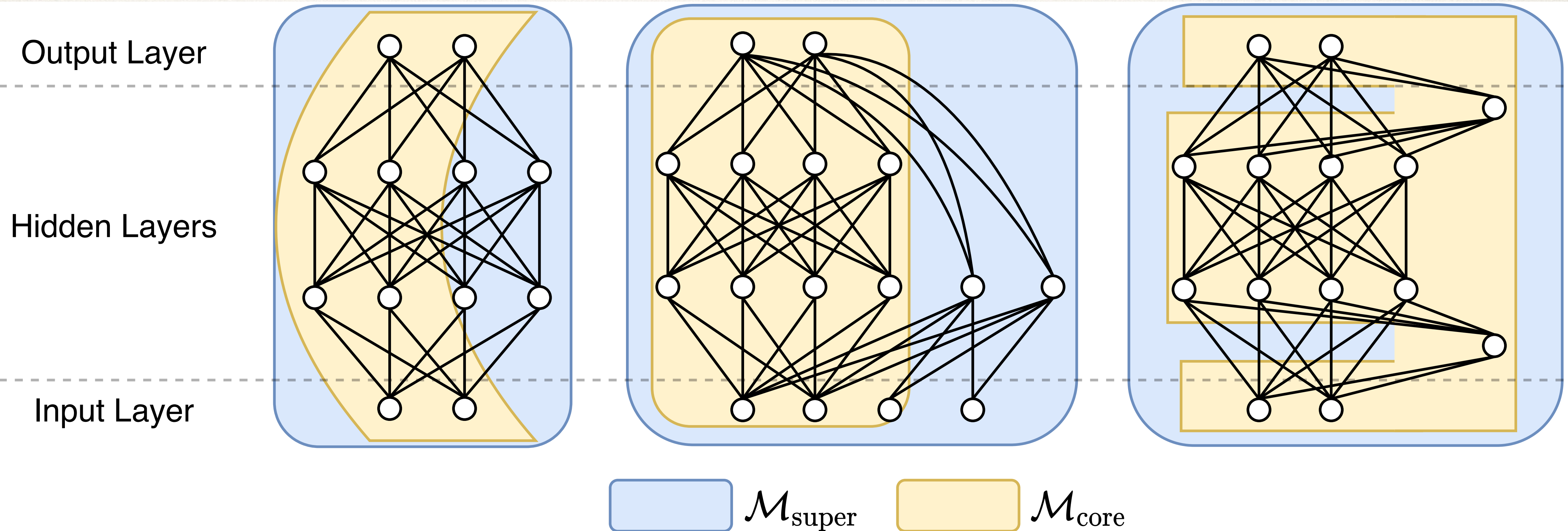
# Gradients from strange architectures

# Alternating Partial Training for Neural Nets



Output Layer

Hidden Layers

Input Layer

$\mathcal{M}_{\text{super}}$   $\mathcal{M}_{\text{core}}$

**Theorem:** Convergence on original network $f$
for non-convex smooth objectives,

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\left[\|\nabla f(\boldsymbol{x}_t)\|^2\right] = \mathcal{O}\left(\sqrt{\frac{q^4 L F_0 \sigma^2}{T}}\right)$$

and similarly for smaller core network

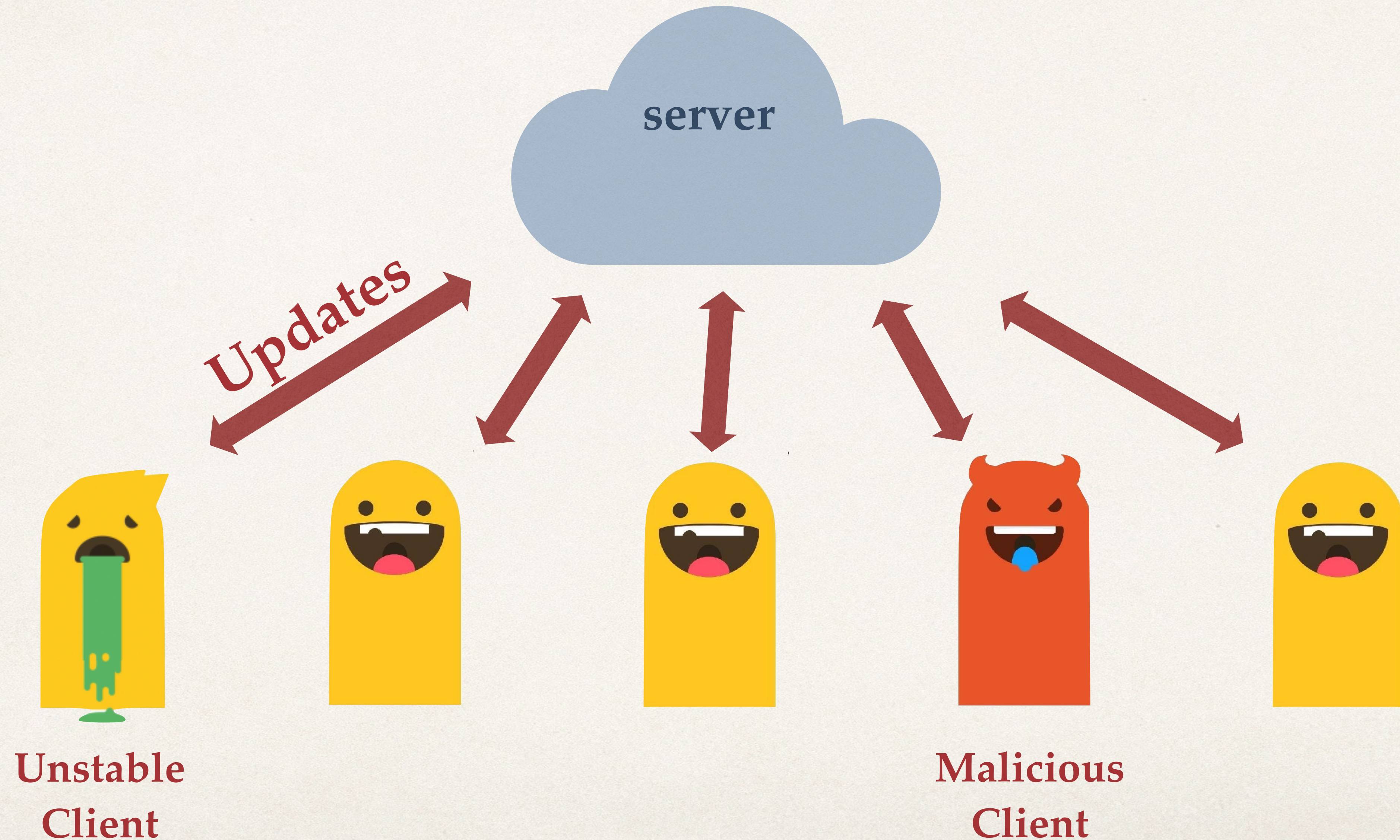$q$ :  "gradient alignment" between parent and core network

- Masked Training of Neural Networks with Partial Gradients, arXiv
- AC/DC: Alternating Training of Deep Neural Networks, Peste et al, NeurIPS 2021

**4**

Gradients from
faulty/malicious collaborators:
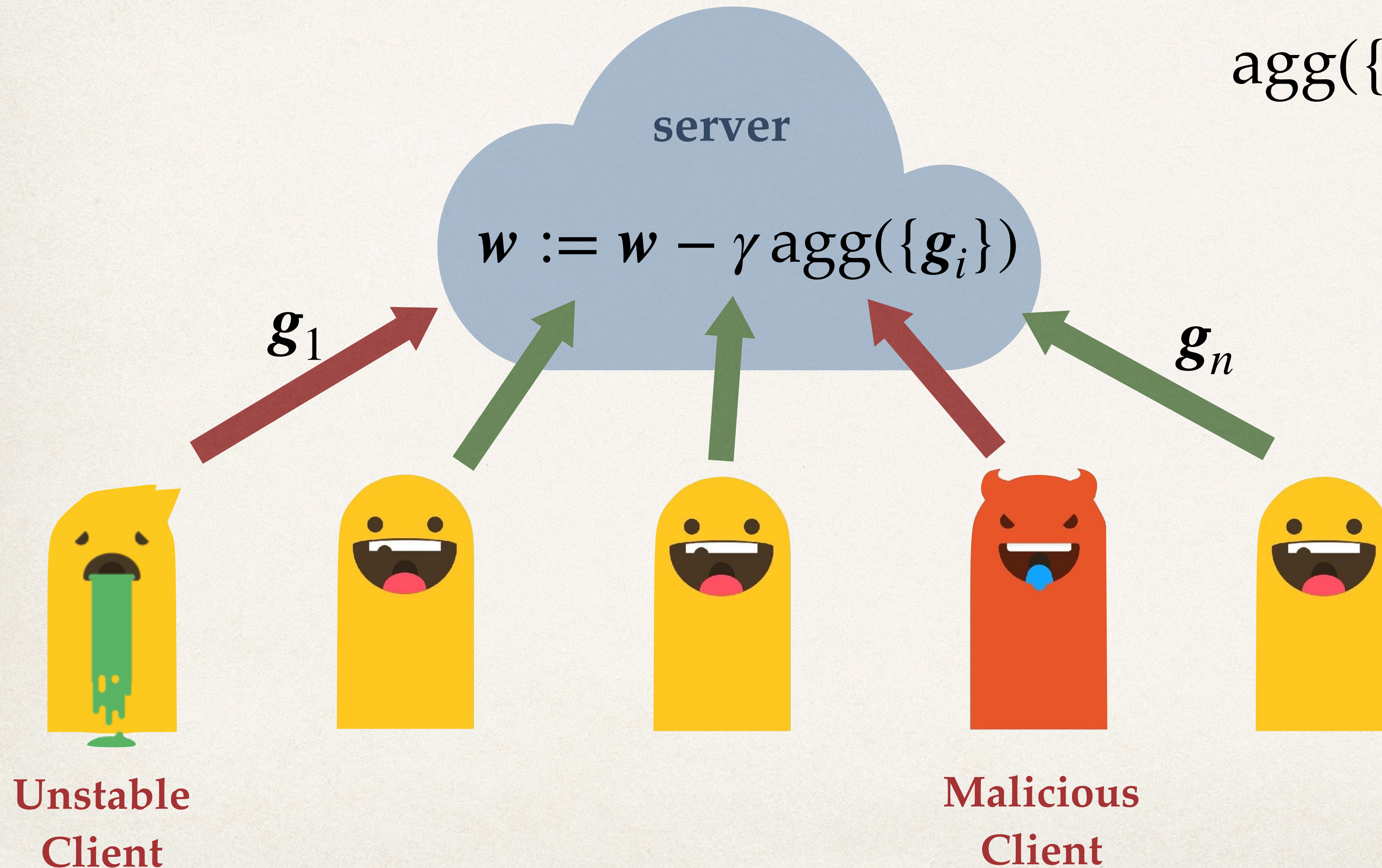 - Byzantine-robust Training

# Malicious actors in FL

# Byzantine Robust Training



server

$$w := w - \gamma \, \mathrm{agg}(\{g_i\})$$

$g_1$

$g_n$

$$\mathrm{agg}(\{g_i\}) := \mathrm{avg}(\{g_i\})$$

$$:= \mathrm{CM}(\{g_i\})$$

Examples:

- Coordinate-wise median [Yin et al. 2017]

- Krum [Blanchard et al. 2018]

- Geometric median / RFA [Pillutla et al. 2019]

**Unstable Client**

**Malicious Client**

# Fall of Empires

- Robustness of the aggregation rule agg($\{g_i\}$) does it imply robust training?

- **NO!**

- Time-coupled attacks:
Little is enough

# Strong negative result

✤ Any aggregation rule which does not use history will **<span style="color:red">fail</span>** training (convergence)

# Fix: Using history with momentum

✤ Simply use worker momentum

$$m_i := (1 - \beta)g_i + \beta m_i$$

✤ Effectively averages past gradients, reducing variance

✤ Aggregate worker momentum instead of gradients

$$w := w - \gamma \, \text{agg}(\{m_i\})$$

# Aggregation with Centered Clipping

✤ Norm-based clipping, before averaging

$$CC = v + \text{clip}_\tau(g_i - v)$$

✤ Removes outliers

✤ Center at previous aggregated update

# Robustness theorem

**Theorem:** Given any **($\delta_{\max}$, c)-robust aggregator**, under a **$\delta$**-fraction of attackers and **$\sigma^2$** variance, our algorithm outputs $\mathbf{x^{out}}$ s.t.

$$\mathbb{E}\|\nabla f(x^{\text{out}})\|^2 \leq \mathcal{O}\left(\sqrt{\frac{\sigma^2}{T}\left(\delta + \frac{1}{n}\right)}\right)$$

# References

1. ✤ **Mime: Mimicking Centralized Stochastic Algorithms in Federated Learning**
   - ✤ NeurIPS 2021  *paper link*

2. ✤ **Optimal Model Averaging: Towards Personalized Collaborative Learning**
   - ✤ FL workshop at ICML 2021  *paper link*
   ✤ **Linear Speedup in Personalized Collaborative Learning**
   - ✤ arXiv  *paper link*

3. ✤ **Masked Training of Neural Networks with Partial Gradients**
   - ✤ arXiv  *paper link*

4. ✤ **Learning from History for Byzantine Robust Optimization**
   - ✤ ICML 2021  *paper link*

# Thanks

Sai Praneeth Karimireddy, Sebastian U. Stich, Lie He, El Mahdi Chayti, Amirkeivan Mohtashami, Felix Grimberg, Nicolas Flammarion,  Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Ananda Theertha Suresh

**EPFL**

Machine Learning and Optimization Laboratory

website:    mlo.epfl.ch