

Optimization
for Machine Learning
in Practice II

Martin Jaggi

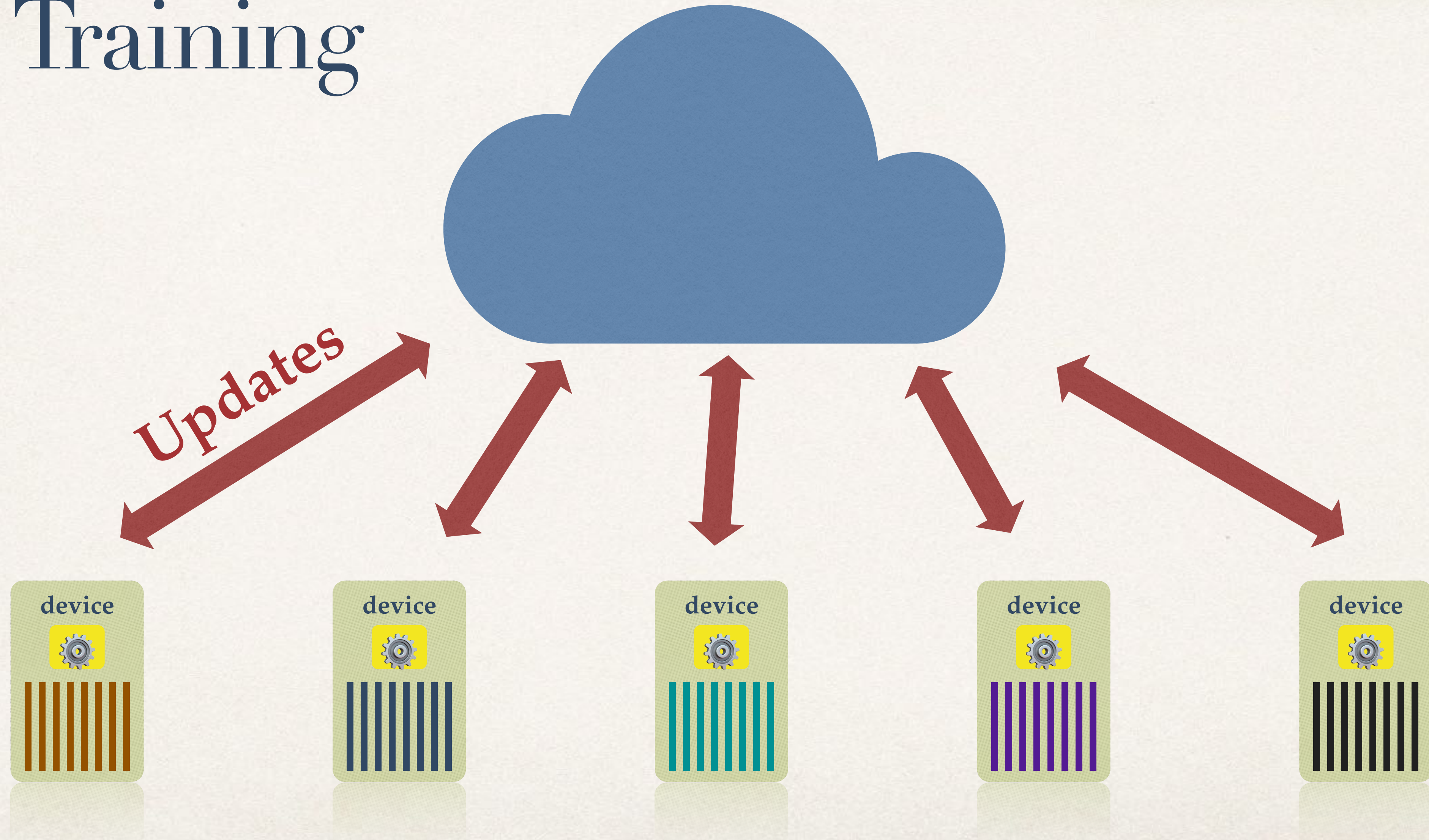
EPFL

Machine Learning and Optimization Laboratory

mlo.epfl.ch

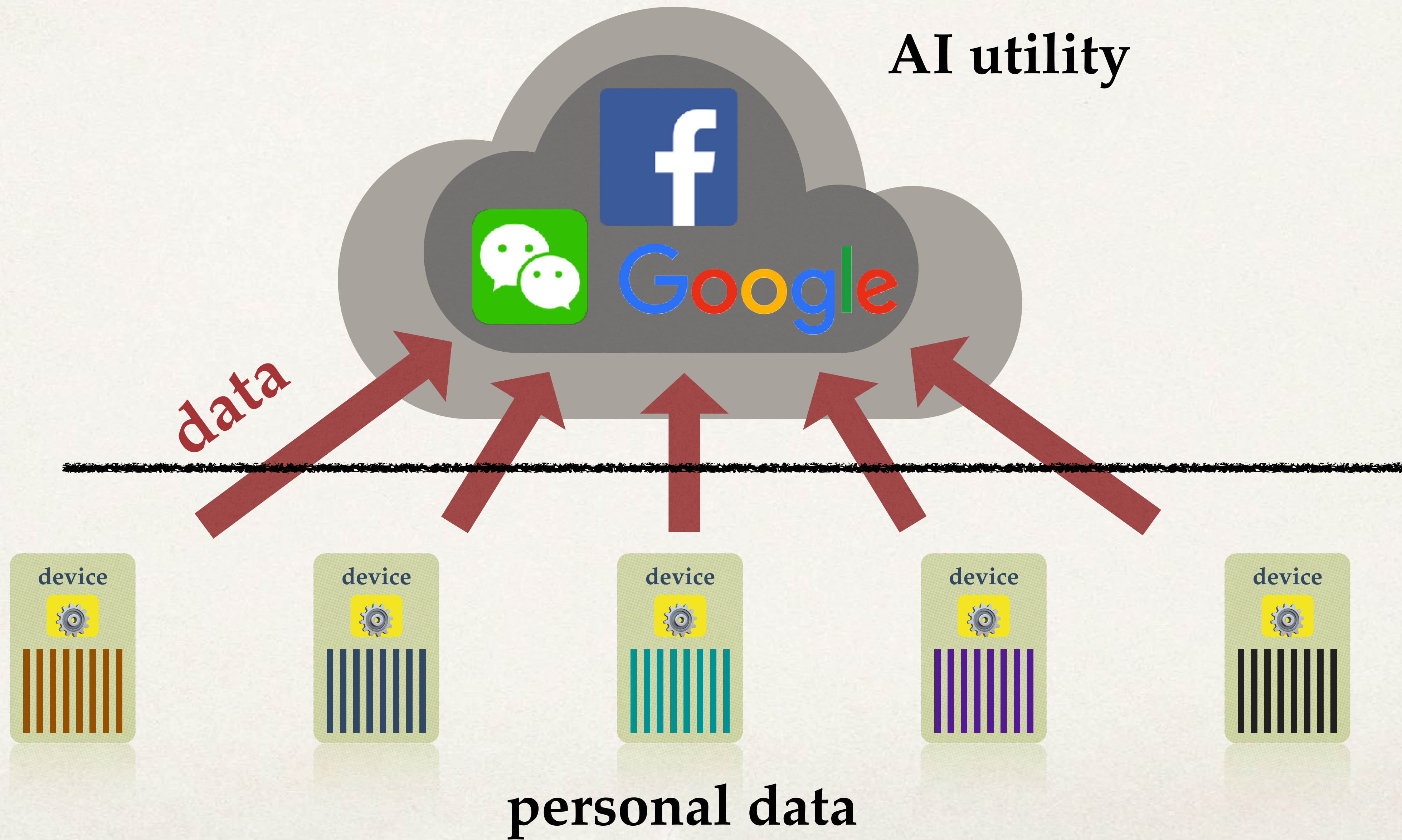
2

Collaborative Training



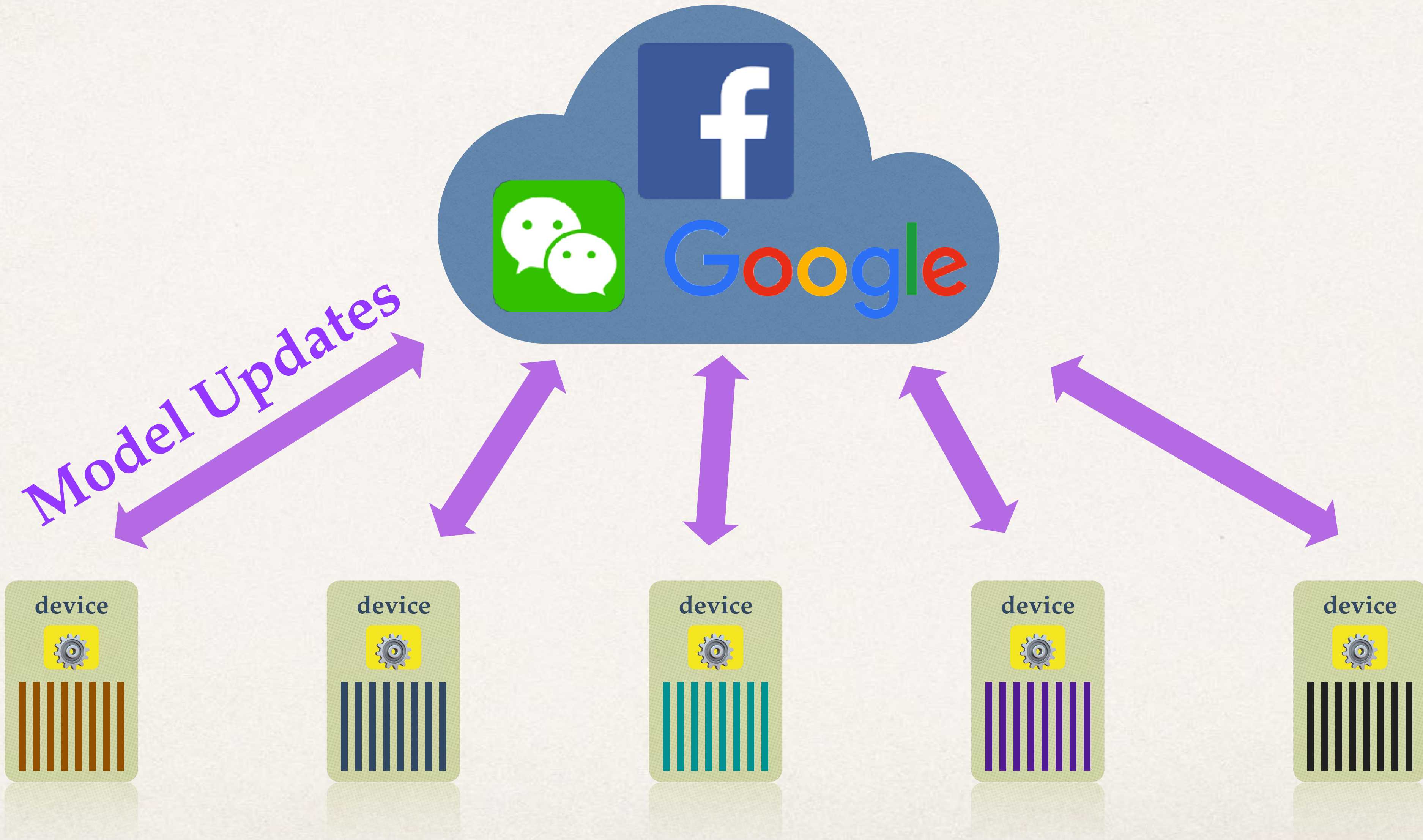
Big Picture

ONE WAY



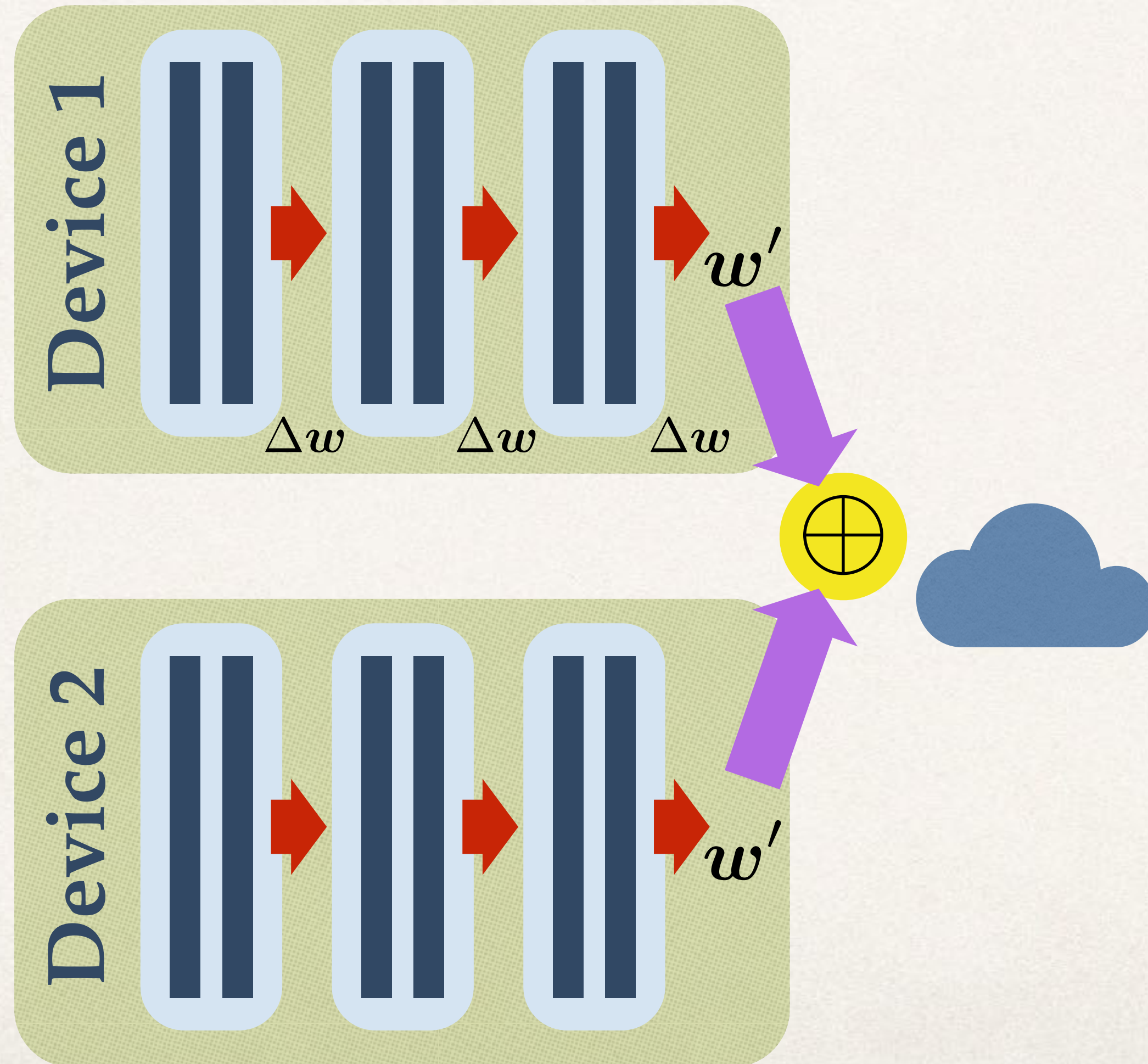
2a

Federated Learning



2a

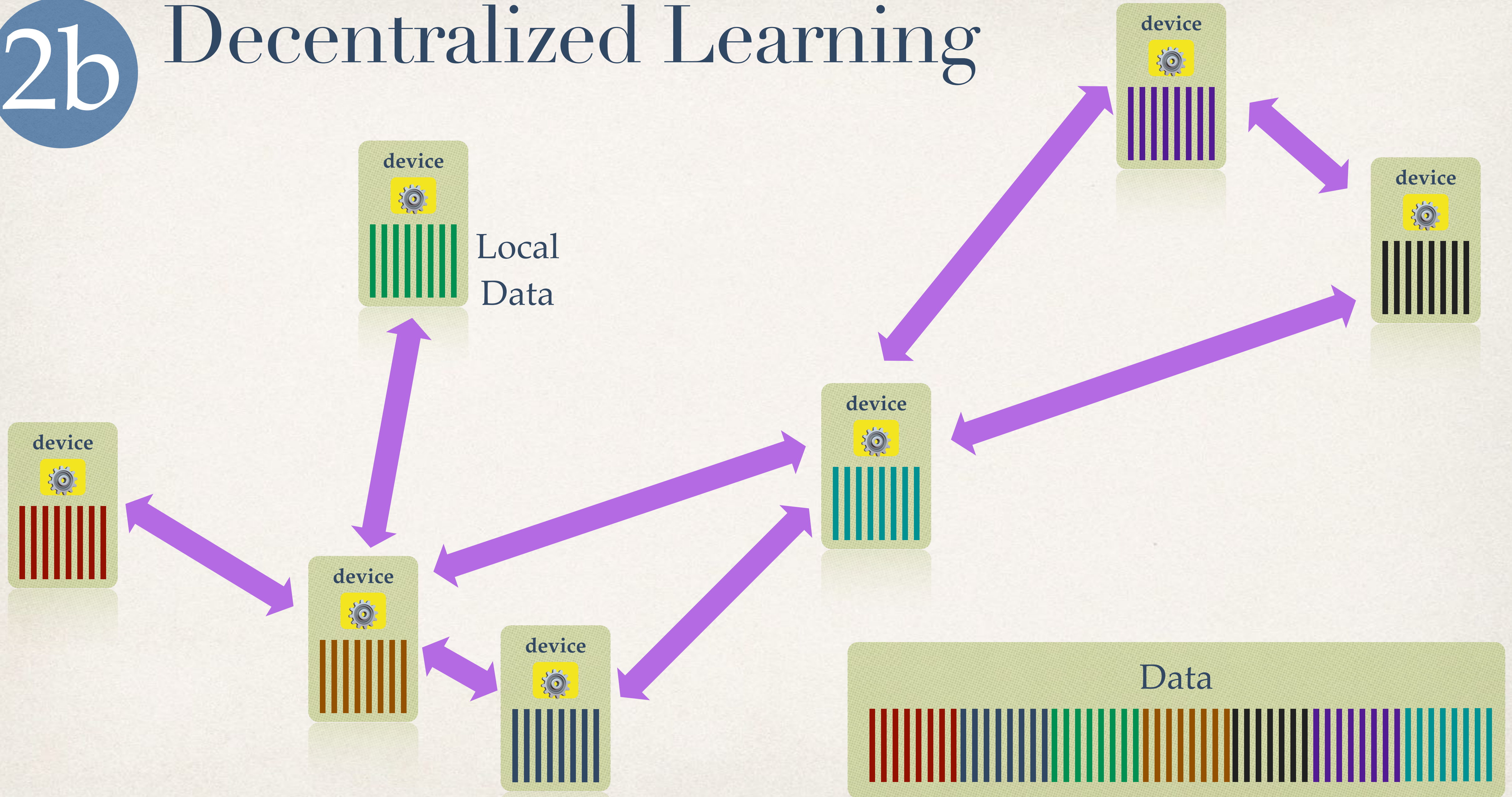
Federated Learning



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

2b

Decentralized Learning



Motivation

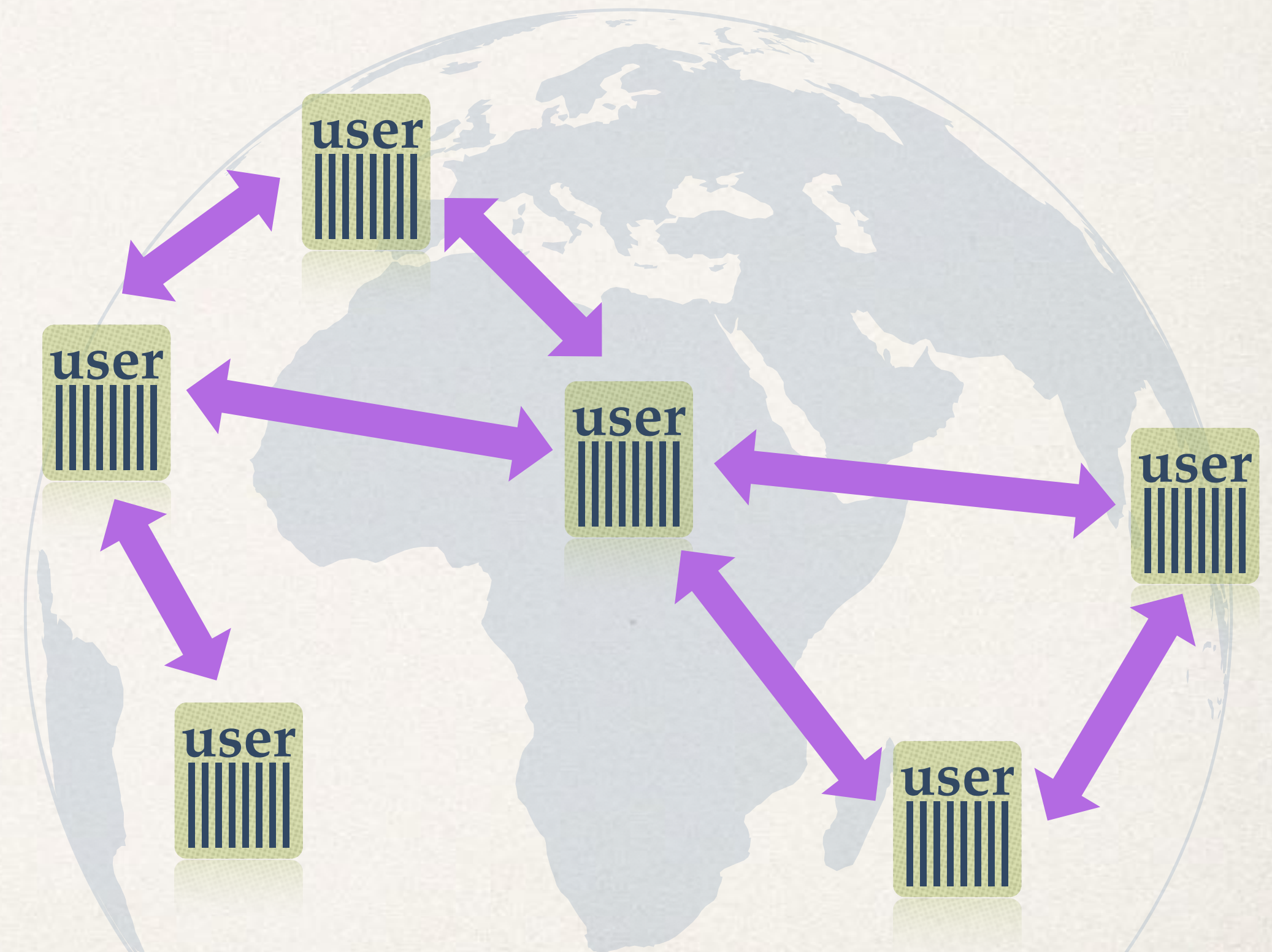
❖ Applications:

any ML system with user data
servers, devices, sensors, hospitals, ...



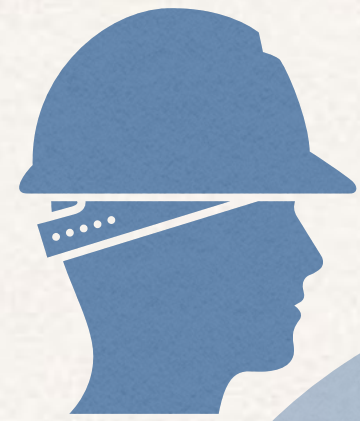
[image source](#)

❖ Advantages:



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

Required Building Blocks



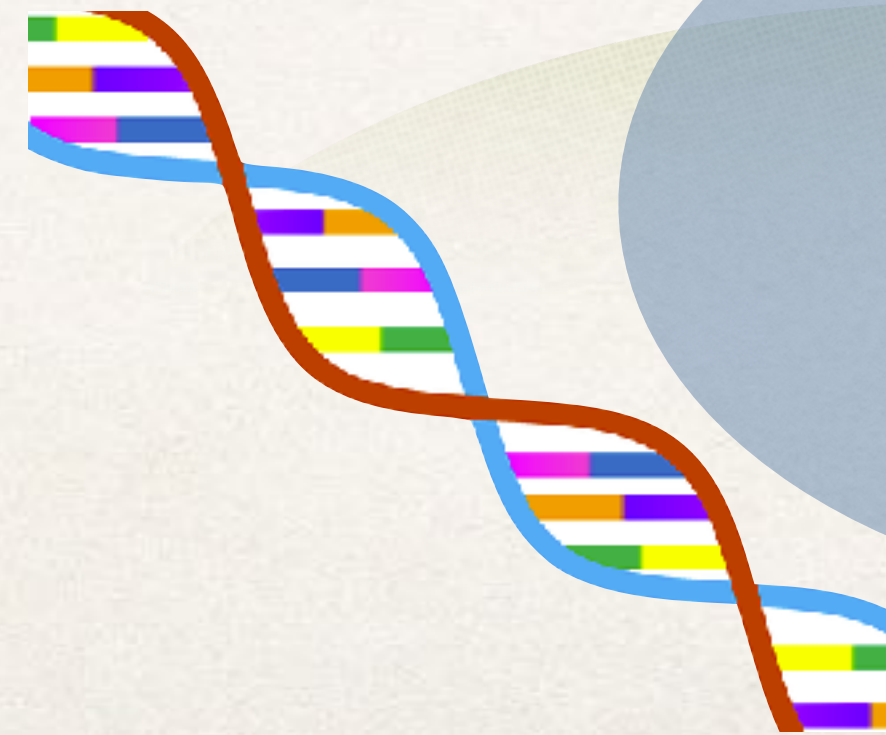
Robustness

Decentralized
ML

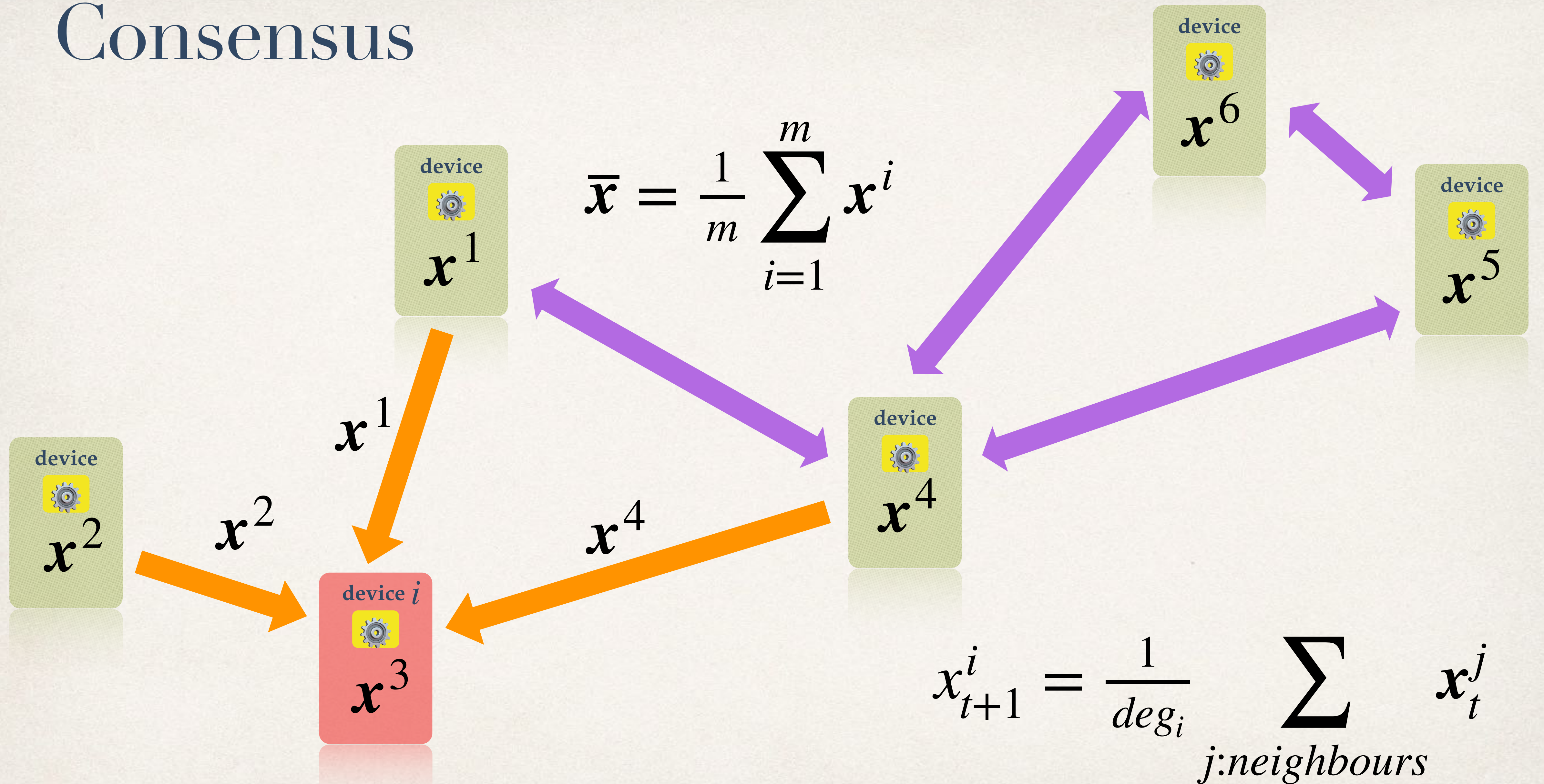
Efficiency



Privacy



Consensus



Communication Compression

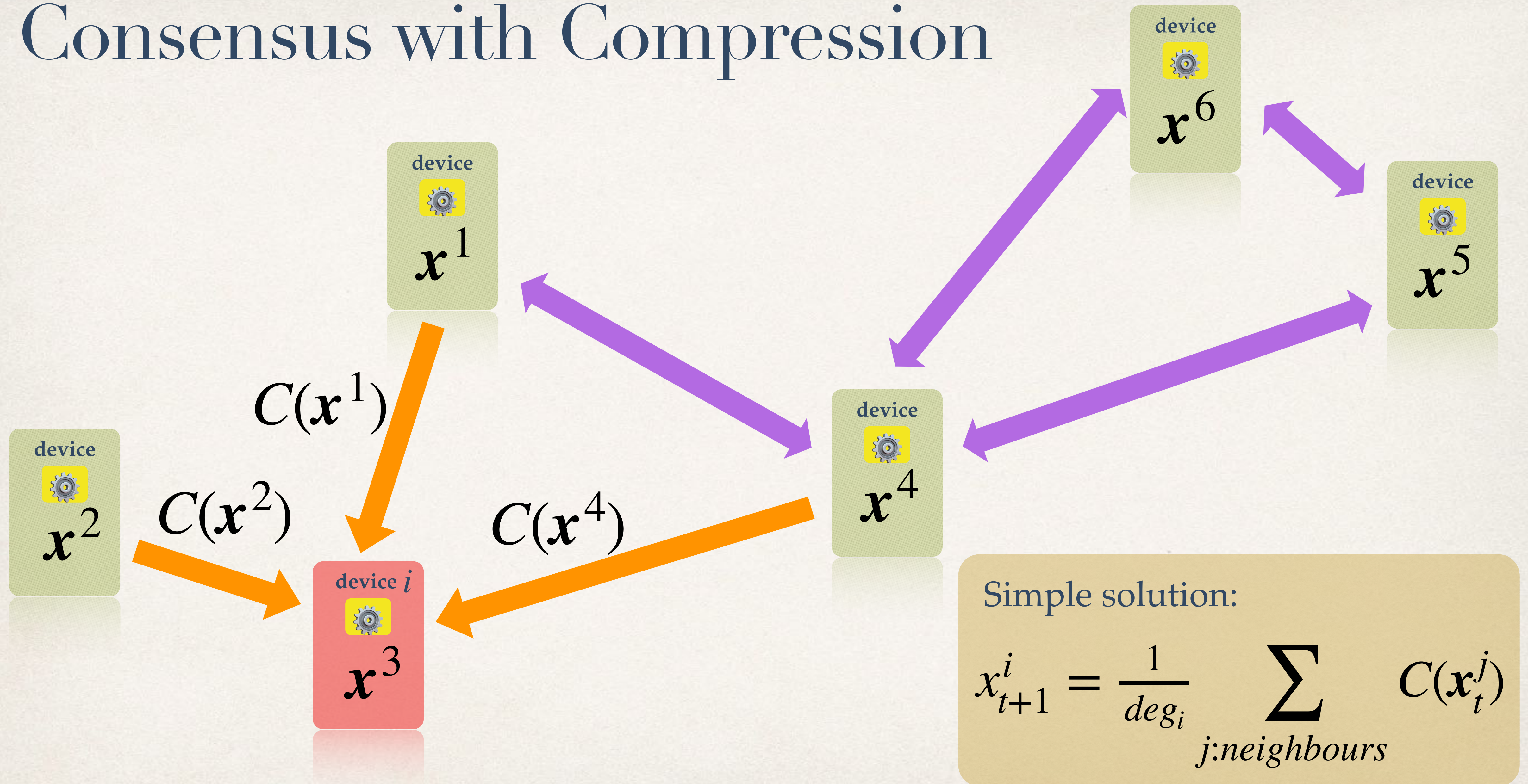
- ❖ limited-bit precision vector

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

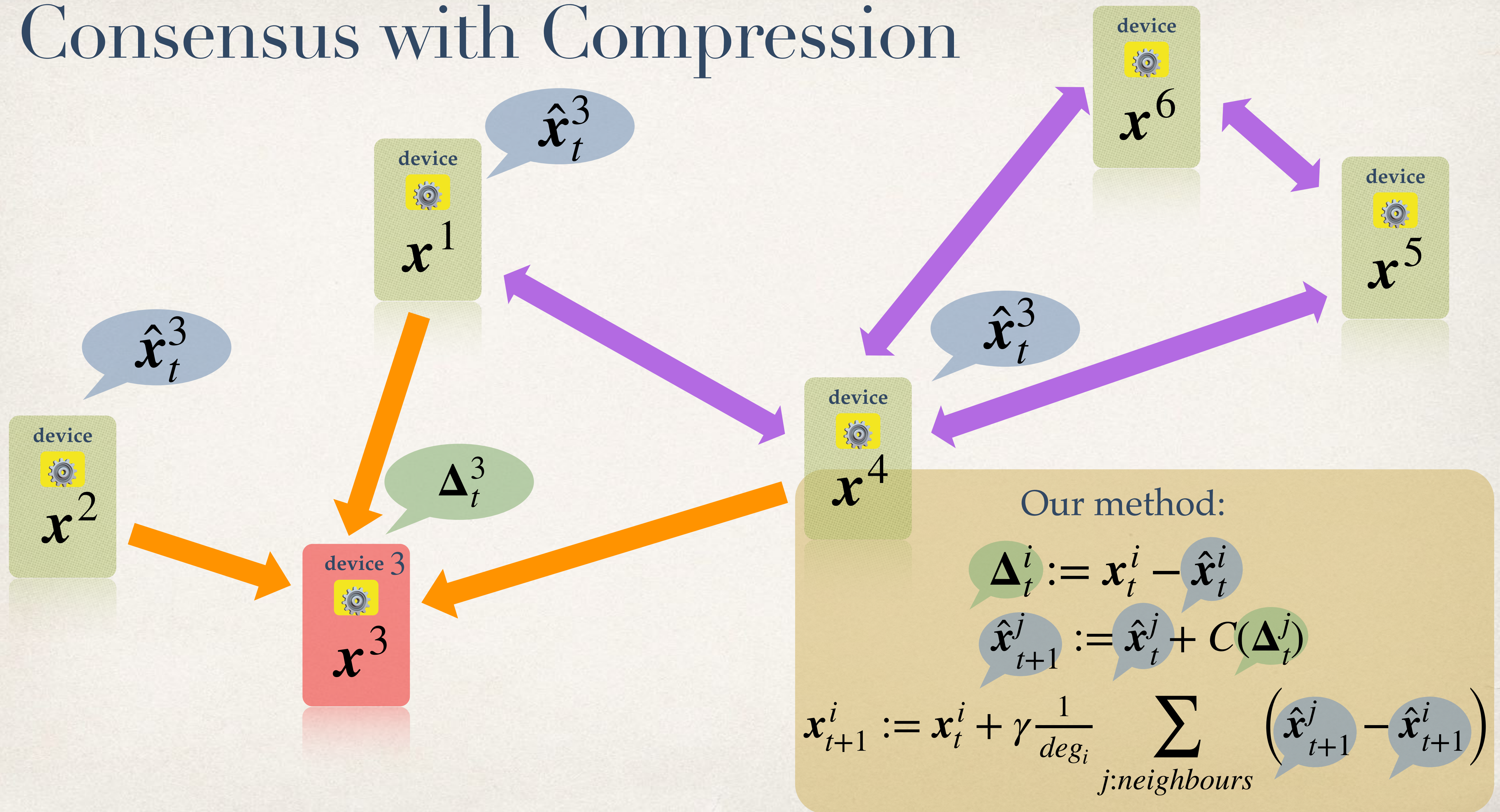
- ❖ random / top $k\%$ of all the entries

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

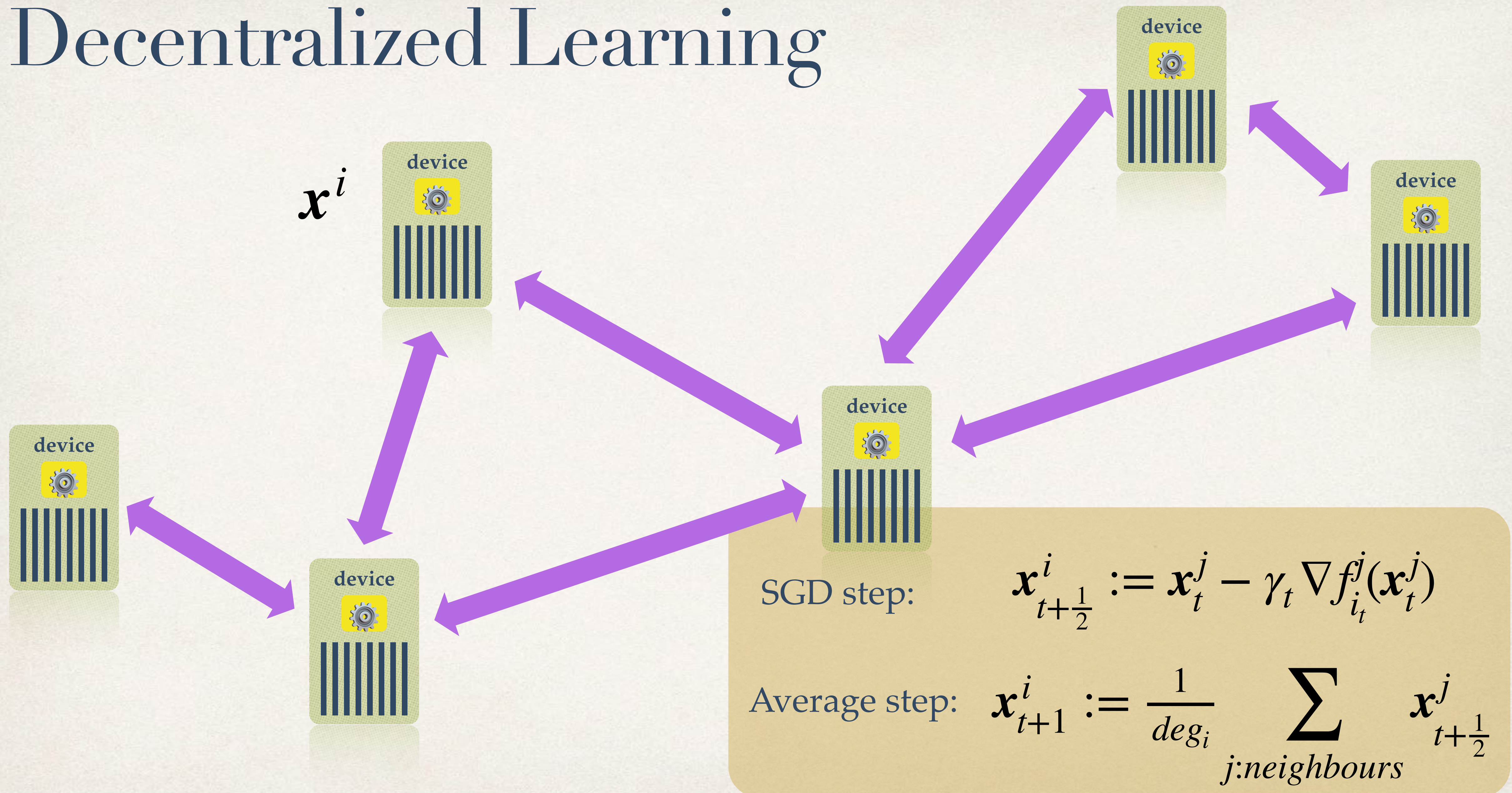
Consensus with Compression



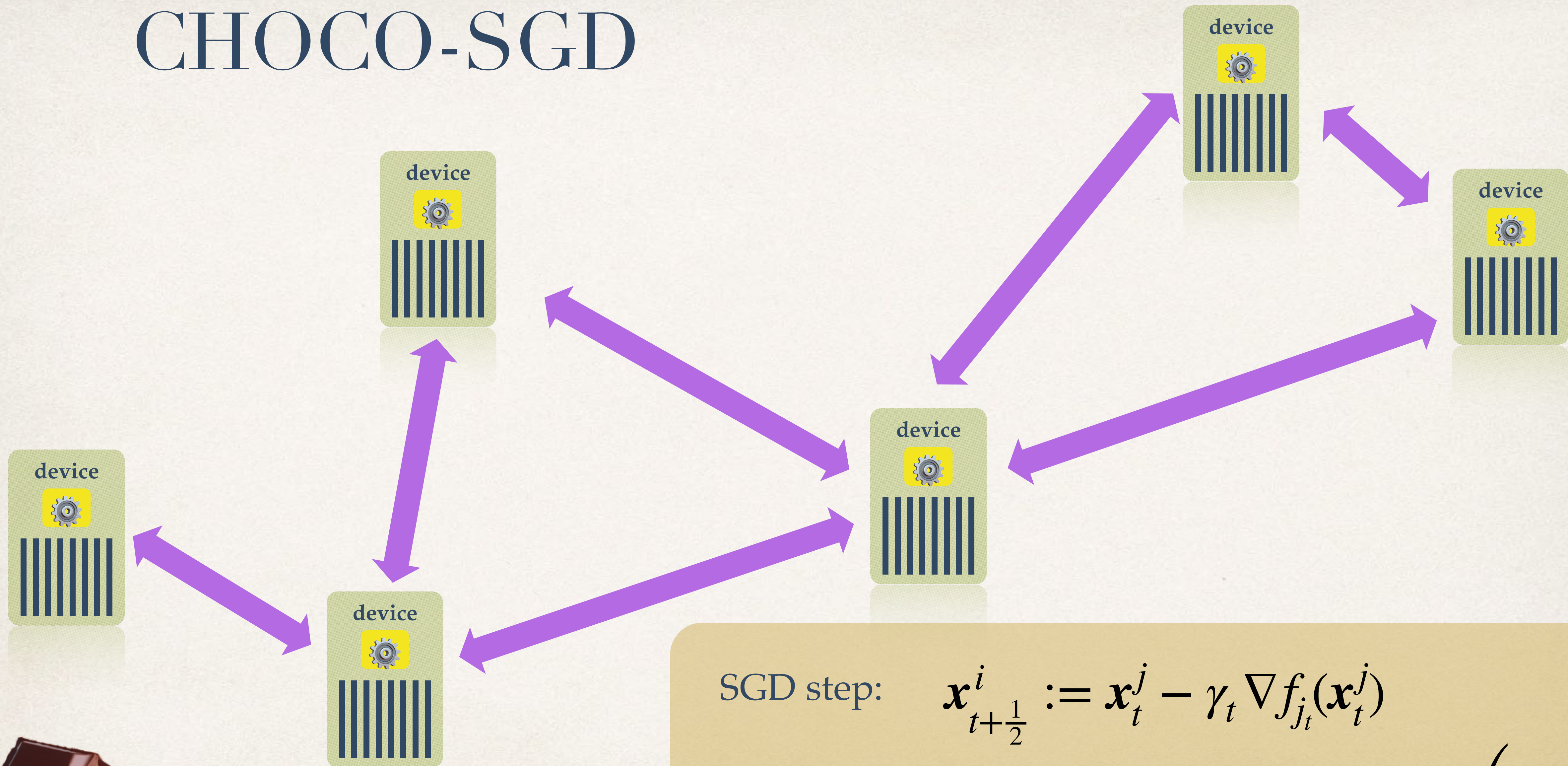
Consensus with Compression



Decentralized Learning



CHOCO-SGD



SGD step:
$$\mathbf{x}_{t+\frac{1}{2}}^i := \mathbf{x}_t^j - \gamma_t \nabla f_{j_t}(\mathbf{x}_t^j)$$

$$\mathbf{x}_{t+1}^i := \text{consensus_with_compression} \left(\mathbf{x}_{t+\frac{1}{2}}^j \right)$$



Convergence (Non-Convex Case)

$$\frac{1}{T+1} \sum_{t=0}^T \|\nabla f(\bar{x}_t)\|^2 = \mathcal{O}\left(\frac{1}{\sqrt{nT}} + \frac{n}{\delta^2 \rho^4 T}\right)$$

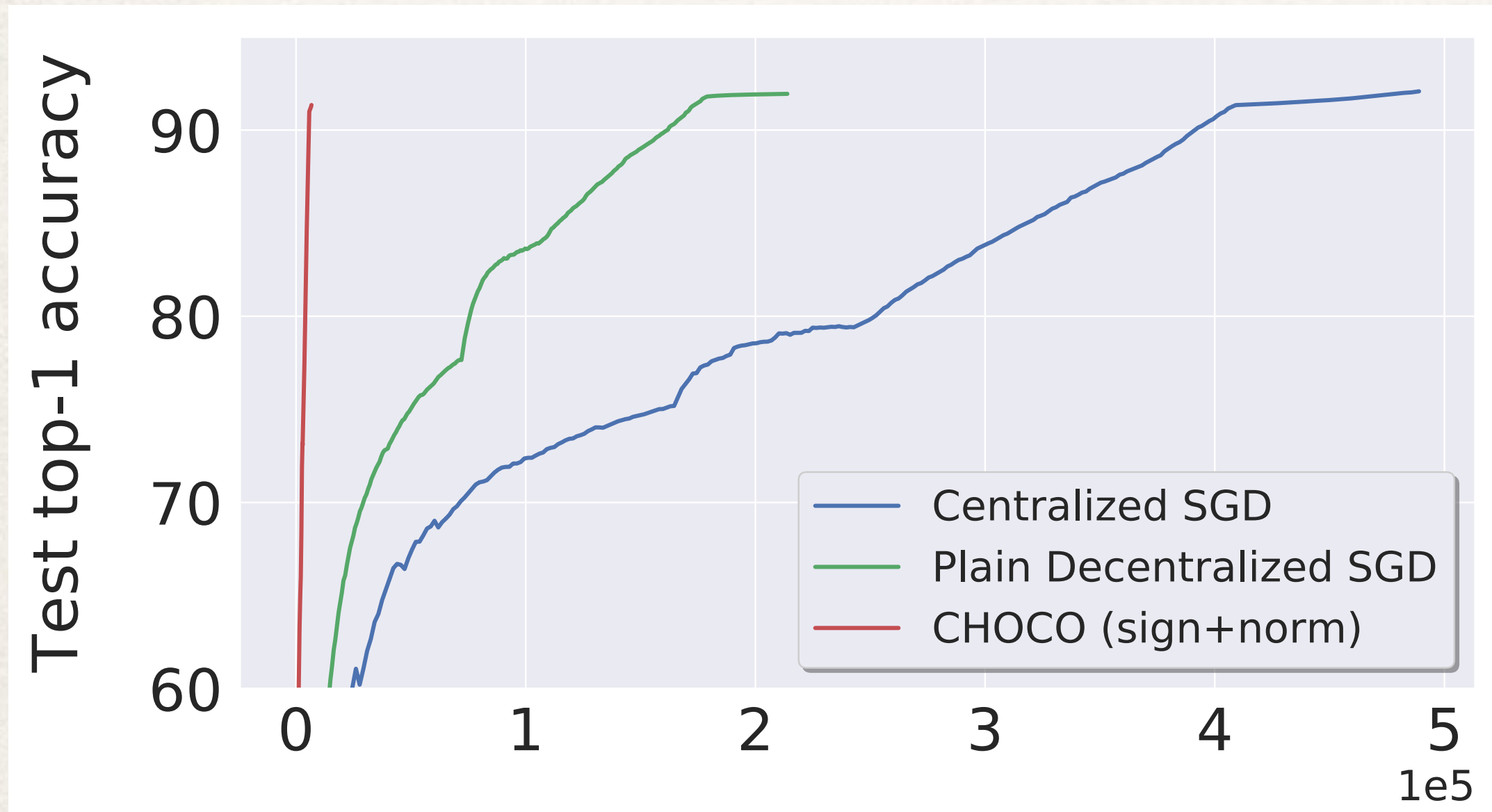
δ — compression ratio $\delta \in [0,1]$, $\delta = 1$ for no compression

ρ — spectral gap of the graph topology



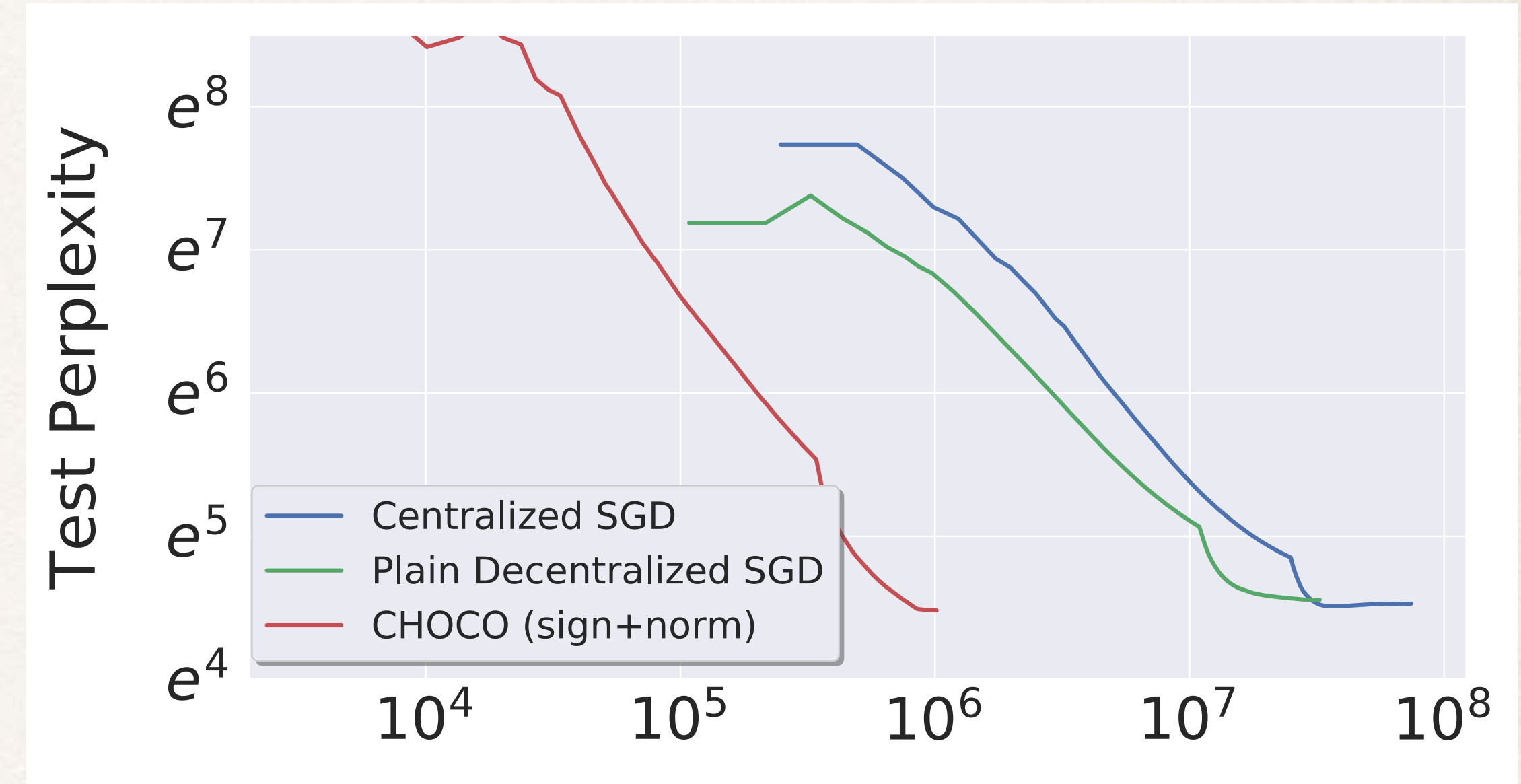
❖ linear speedup in the number of workers

Decentralized DL



data transmitted (MB)

Resnet20 on Cifar 10

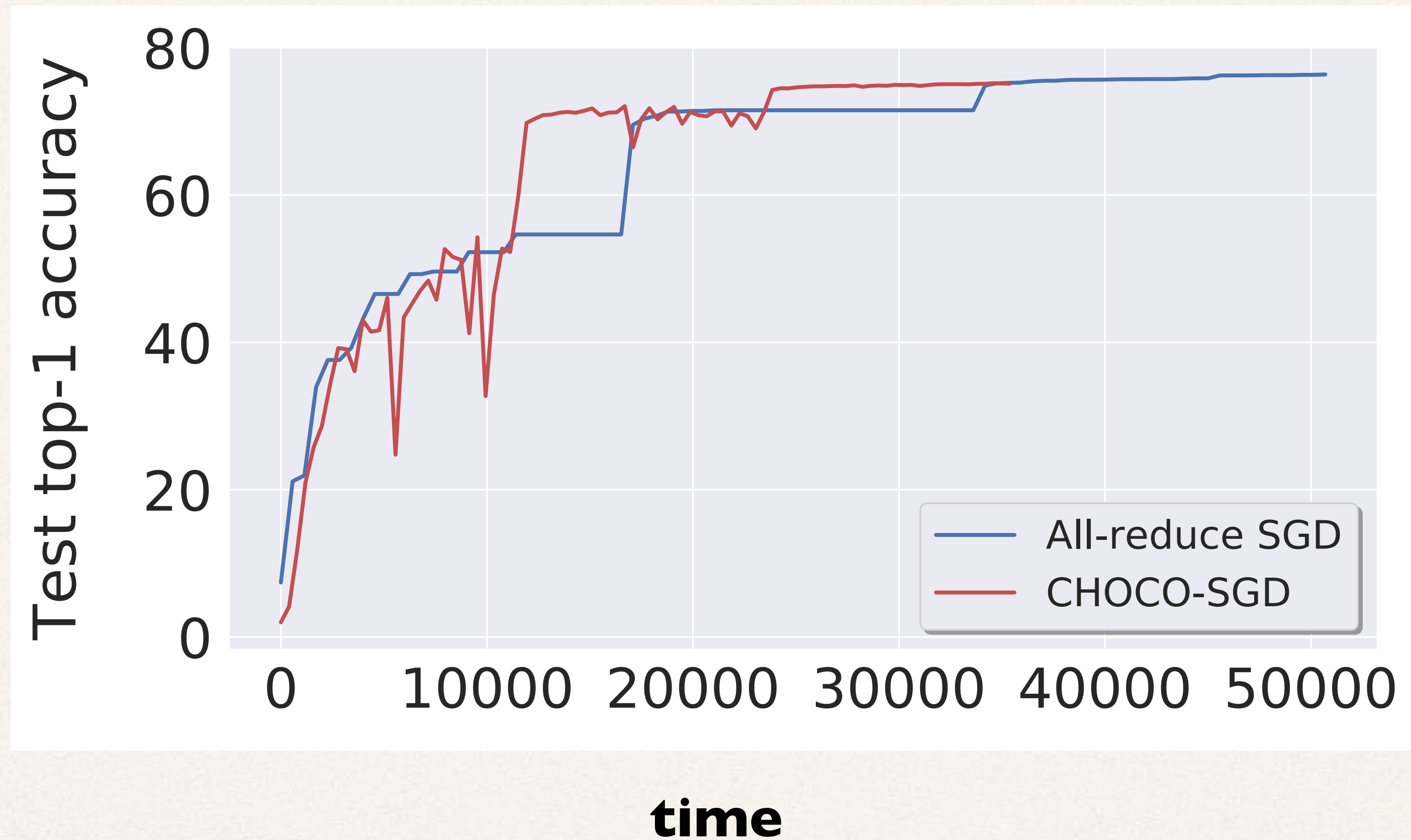


data transmitted (MB)

Language model (3-layer LSTM) on WikiText-2

Social Network Topology, 32 nodes of max deg 14
Sign quantization

DL in Datacenter



Resnet50 on ImageNet-1k
Ring of 8 nodes, each has 4 P100 GPUs

Conclusions - Choco

- ❖ First **consensus algorithm** that converges linearly with arbitrary compression
- ❖ First **decentralized SGD** algorithm that converges with arbitrary compression
- ❖ **Practical performance**



Building Blocks for Decentralized ML

- ❖ **Efficiency: Communication & Compute**

on-device learning, Edge AI
peer-to-peer communication

- ❖ **Privacy**

data locality, leakage?, attacks?

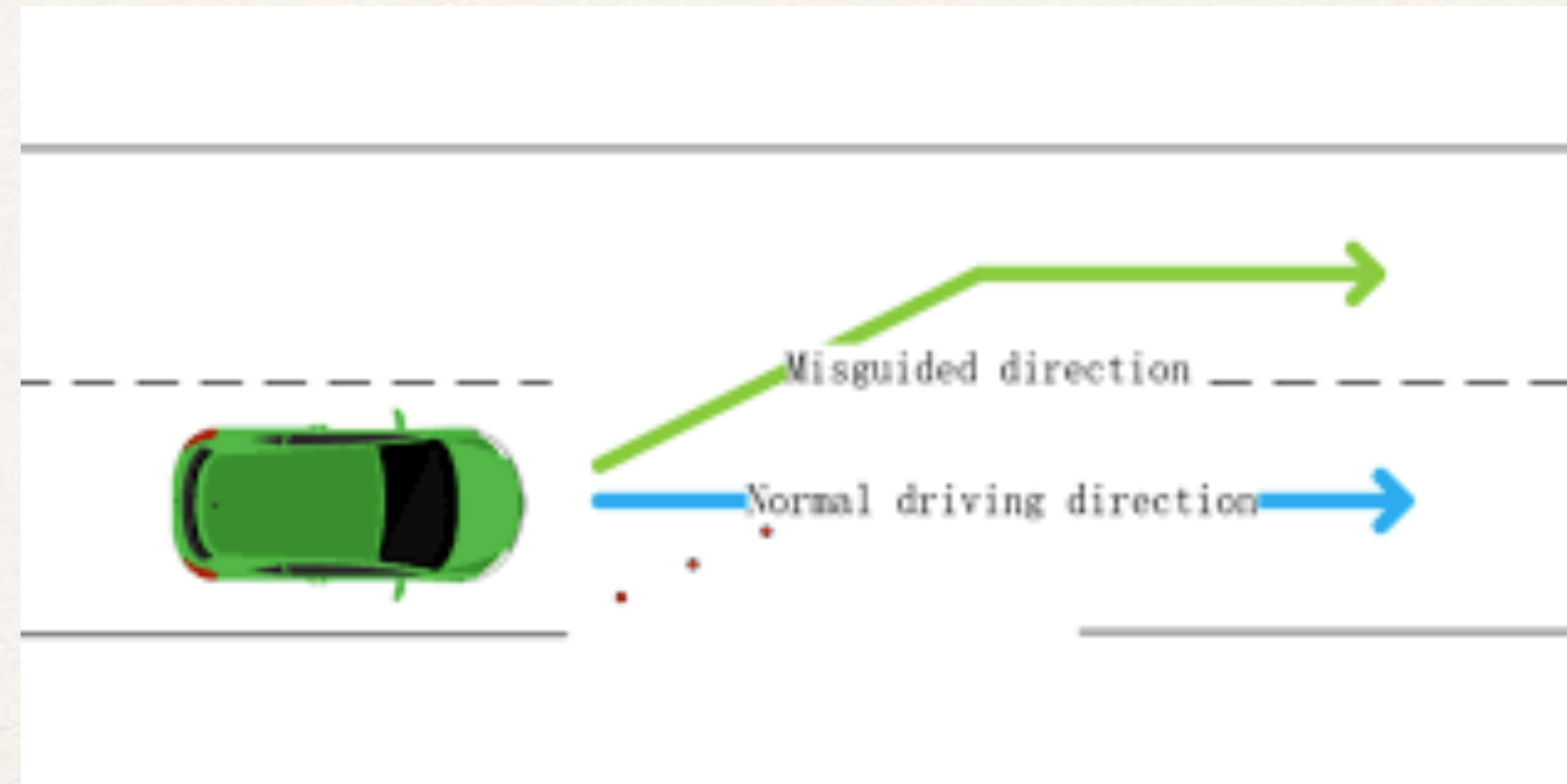
- ❖ **Robustness & Incentives**

tolerate bad players, reward collaboration

3

Robustness

During Training and Inference



Byzantine-robust training



❖ **Mean vs median**

Adversarial Attacks (at inference time)

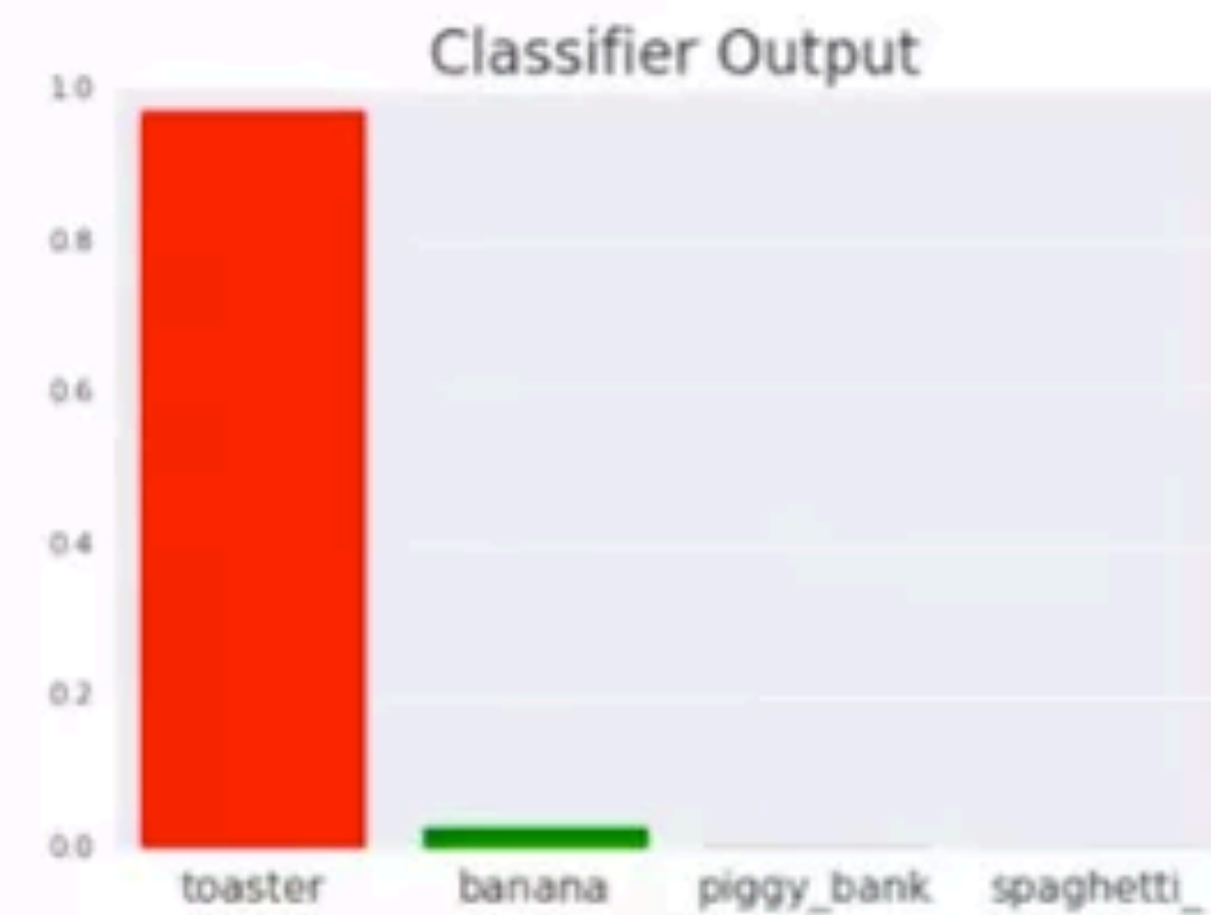
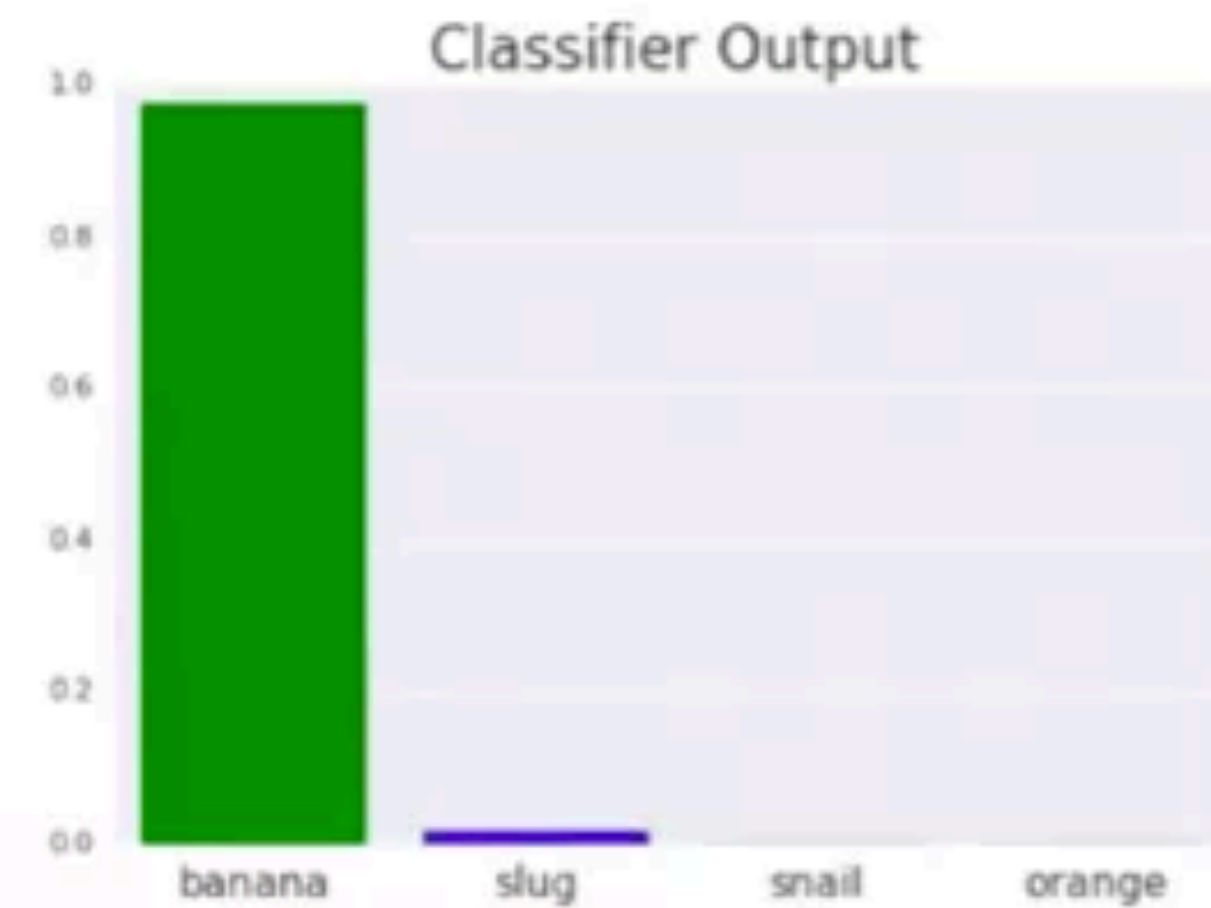


Image: Elsayed, Papernot et al 2018

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

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

$$\nabla_{\mathbf{x}_i} f$$

change data

- ❖ by Projected Gradient Descent!

4

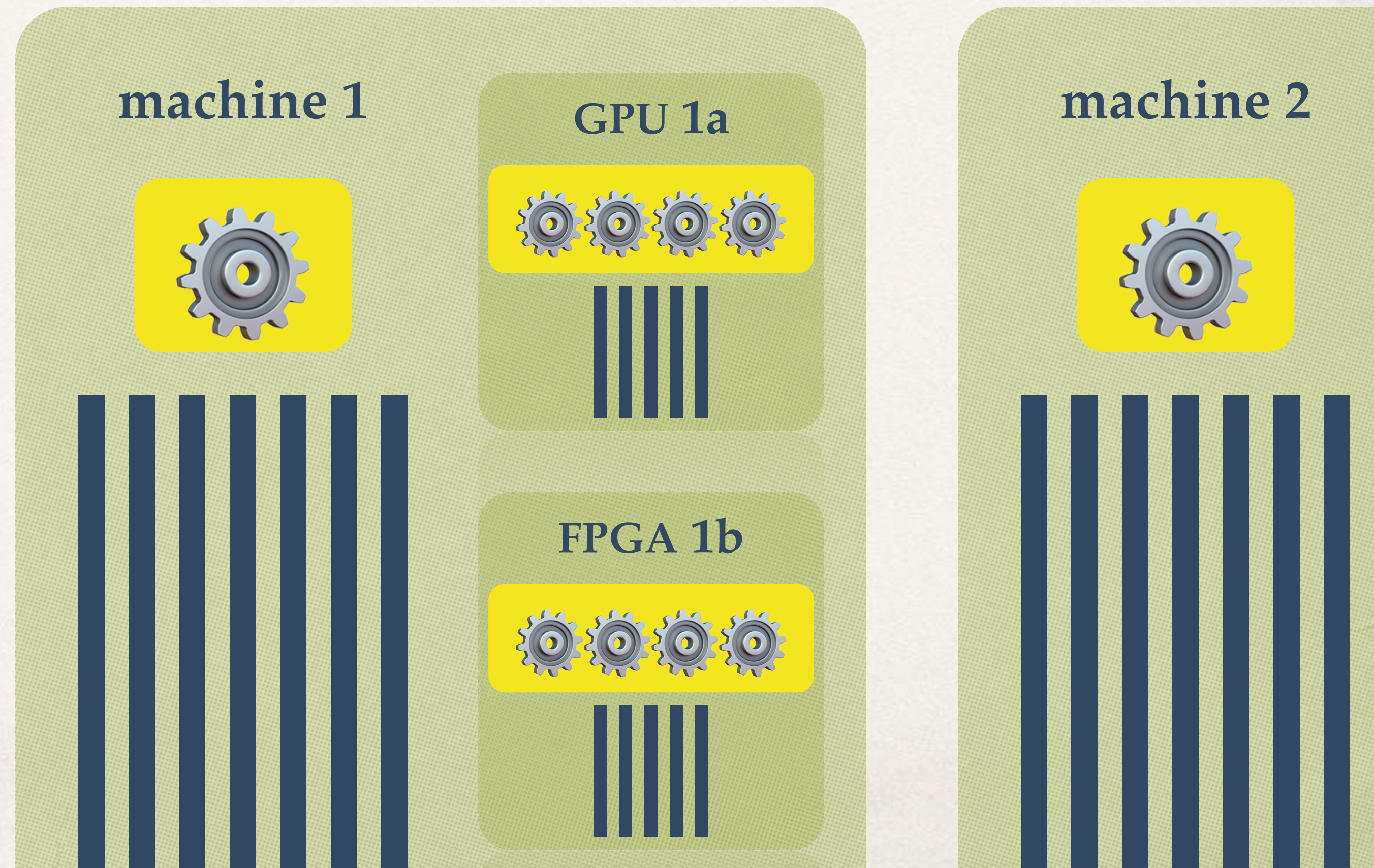
Privacy

- ❖ Secure Multiparty Computation
- ❖ Differential Privacy
- ❖ Privacy / inference Attacks

5

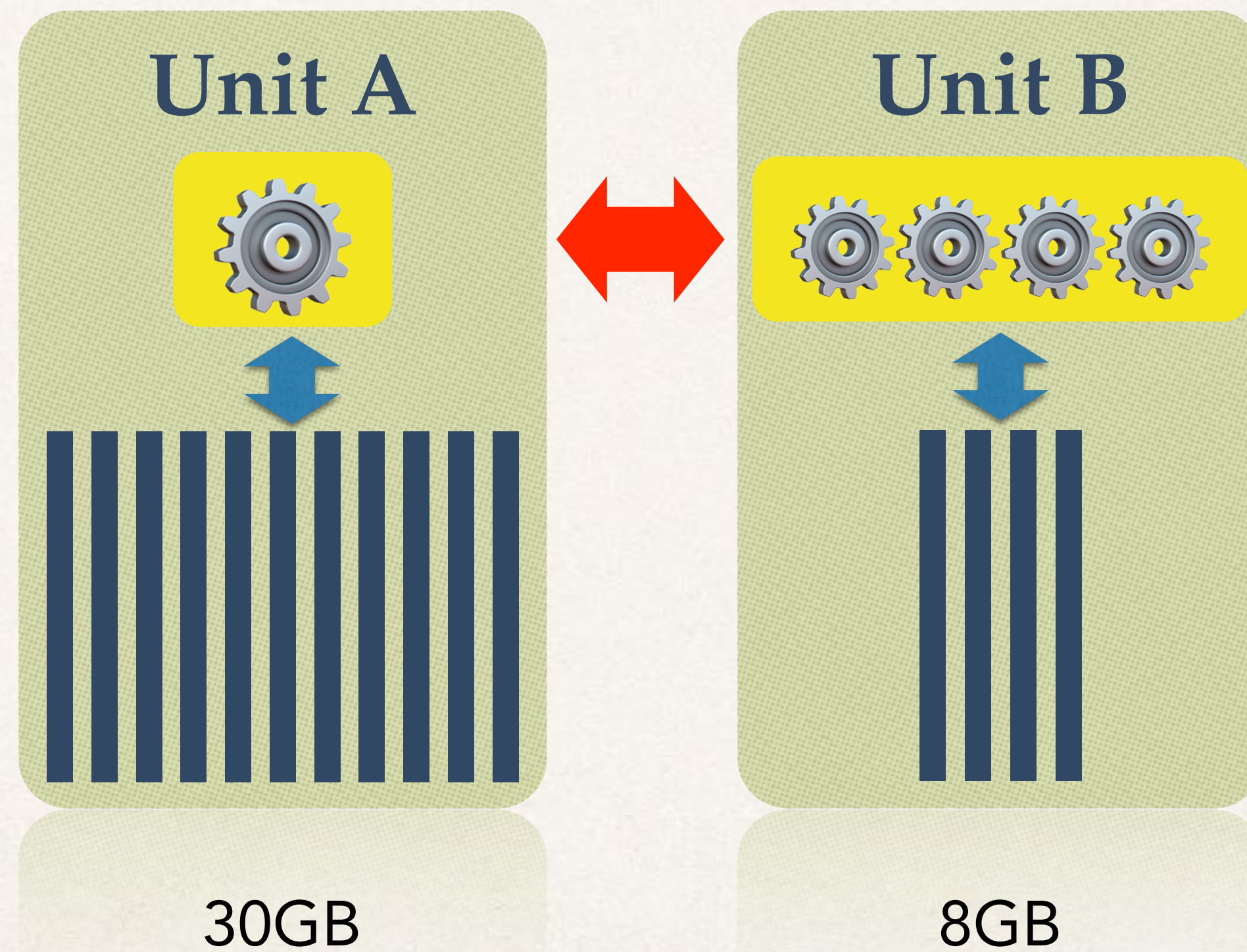
Leveraging Heterogenous Systems

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



Leveraging Heterogenous Systems

duality gap as selection criterion

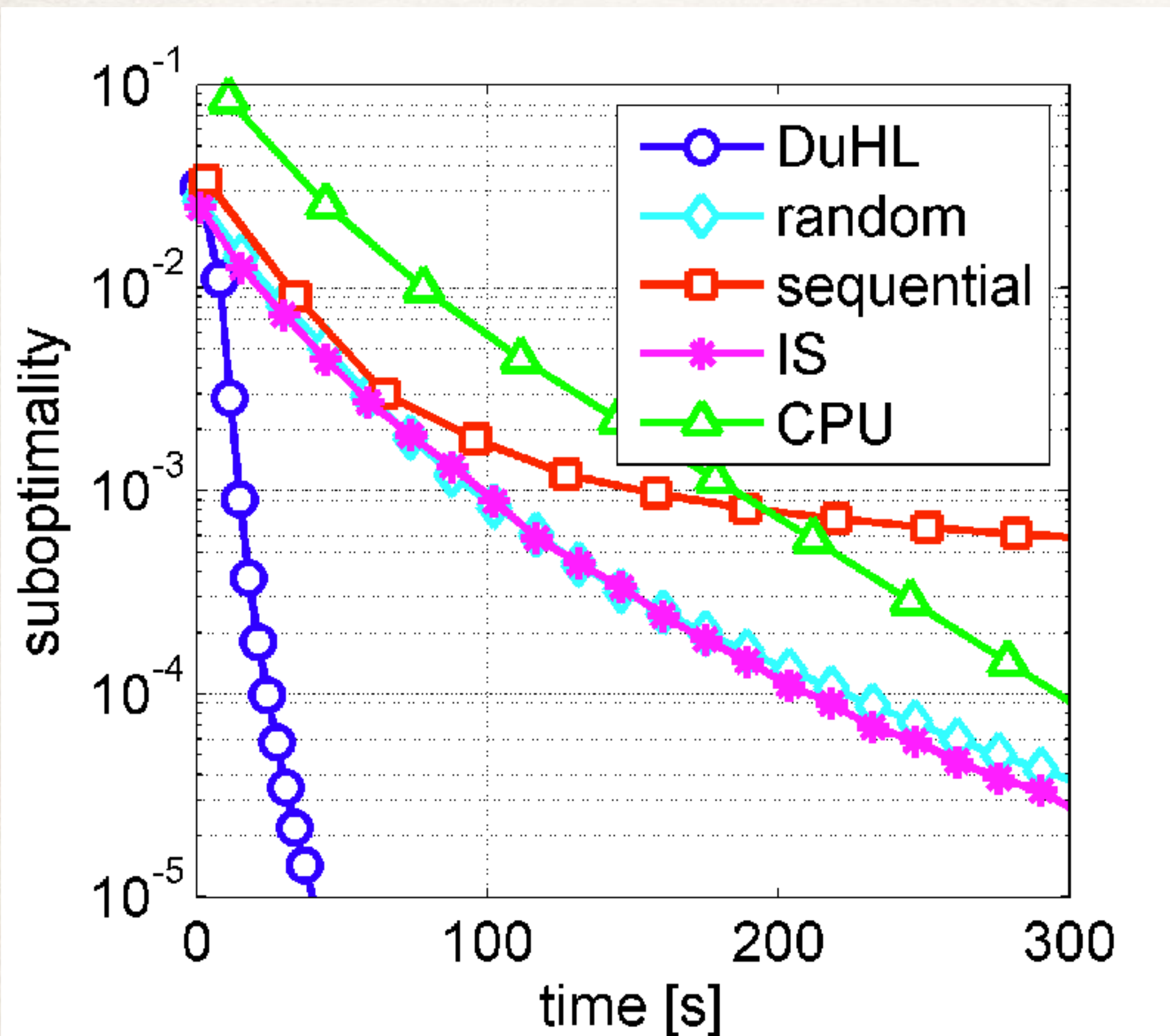


adaptive importance sampling

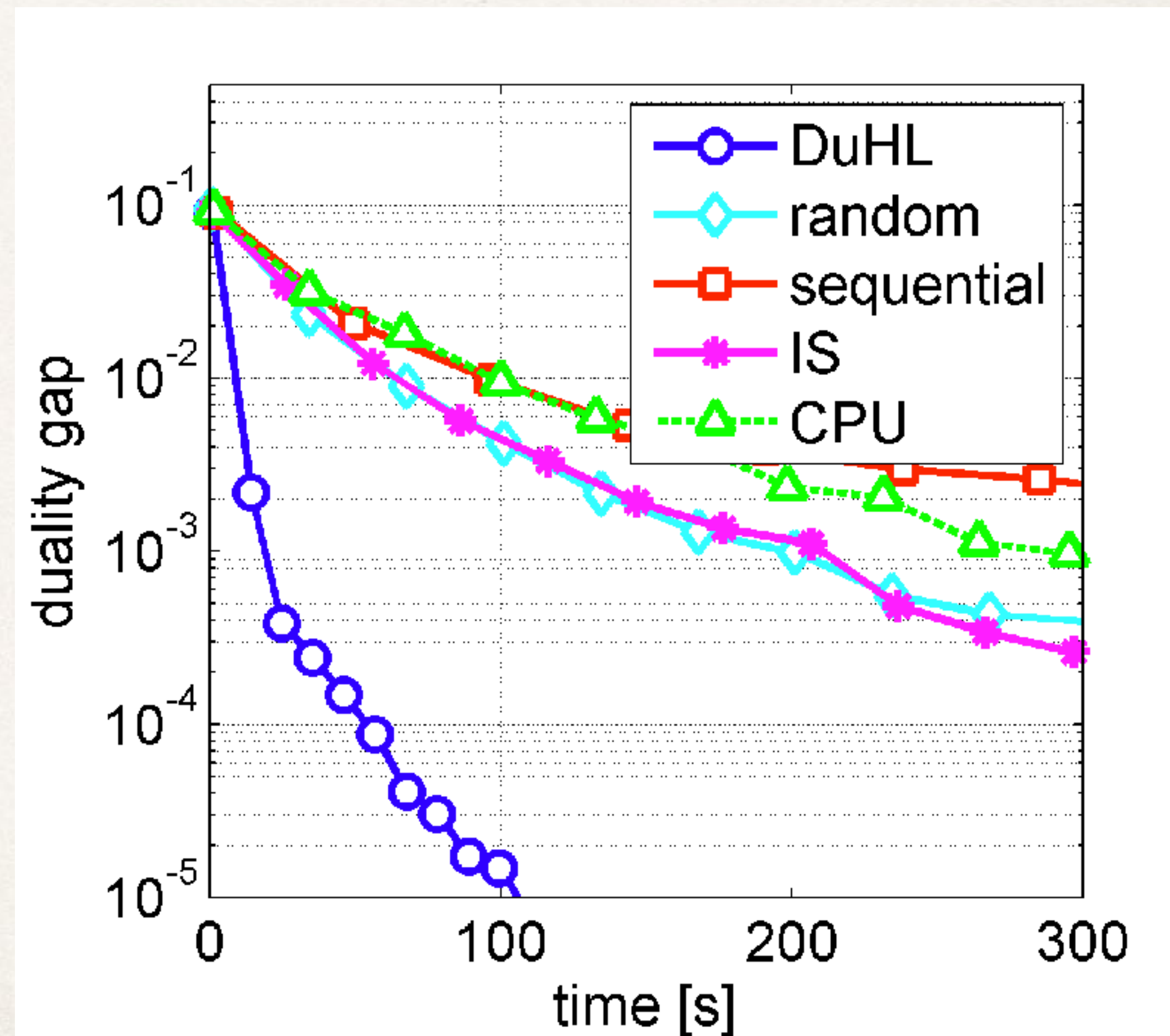
*AISTATS 2017, 2018
NIPS 2017a,b*

Experiments

RAM ↔ GPU, 30GB dataset

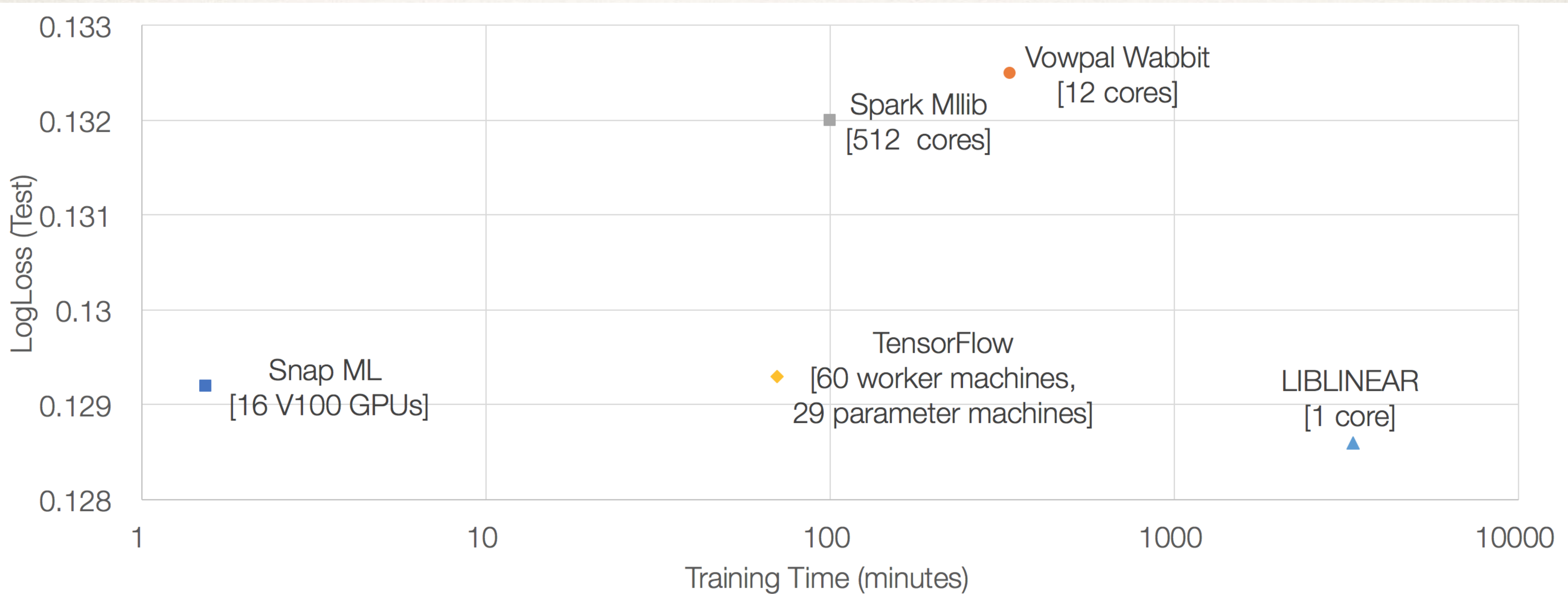


Lasso



SVM

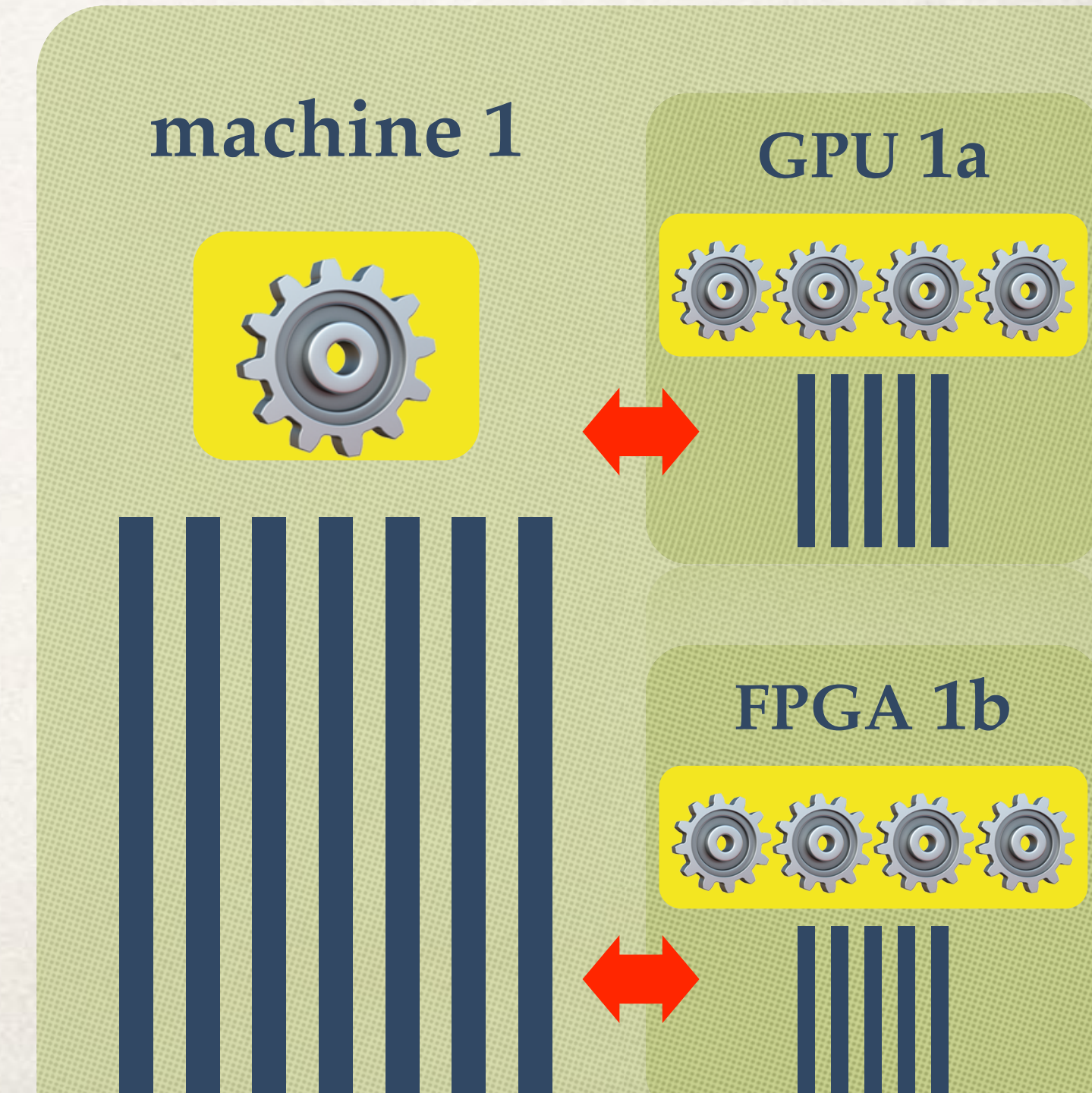
Experiments



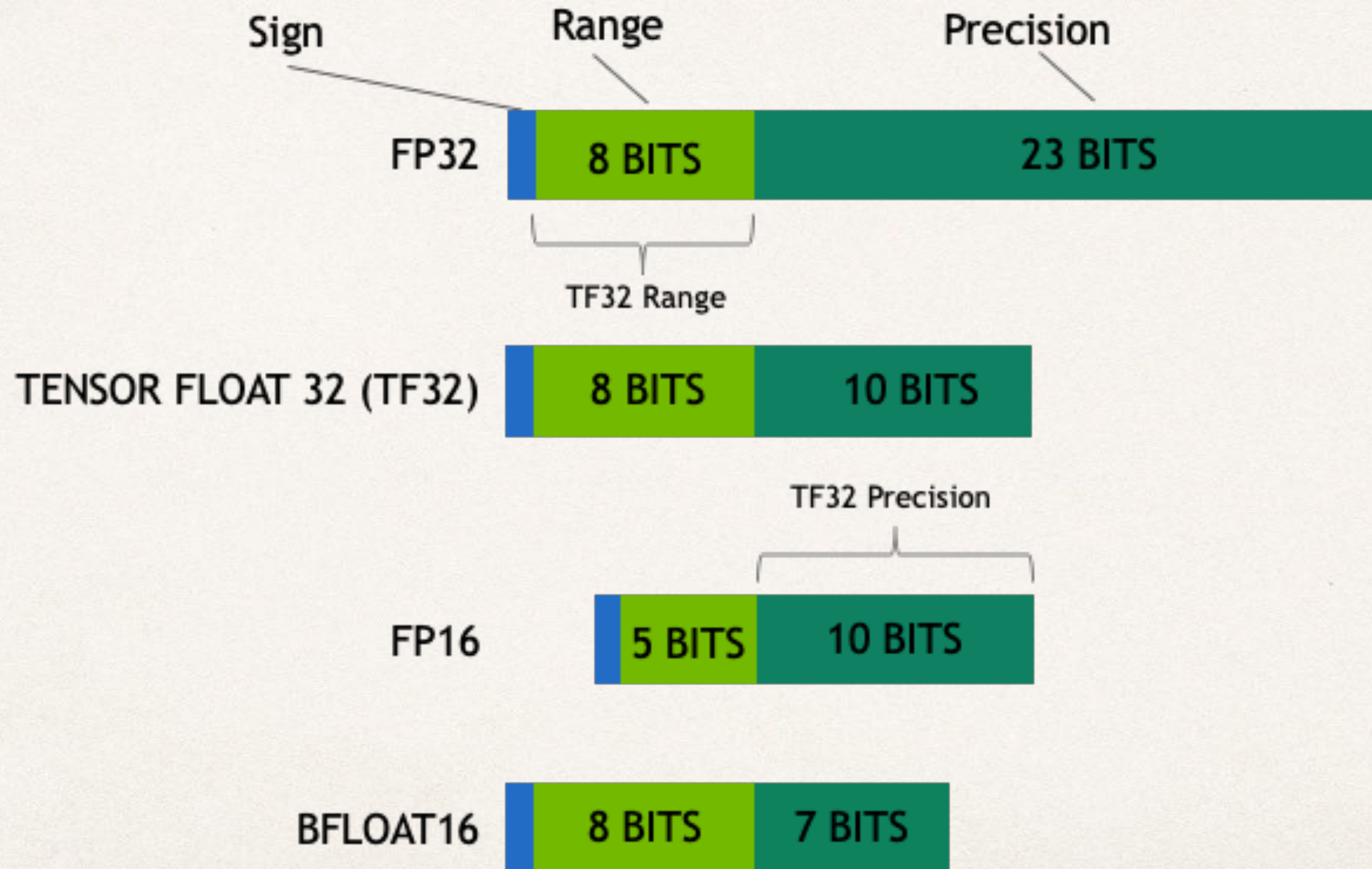
terabyte click log dataset, IBM cloud implementation [*\[arXiv\]*](#)

Trends - Systems

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



Number formats for DL



Open Source Project:

MLbench - Distributed Machine Learning Benchmark

Public and reproducible reference
implementations and benchmarks
for distributed machine learning
algorithms, frameworks and systems.

mlbench.github.io





Search experiment items

- ▶ Saved Datasets
- ▶ Data Format Conversions
- ▶ Data Input and Output
- ▶ Data Transformation
- ▶ Feature Selection
- ▶ Machine Learning
- ▶ OpenCV Library Modules
- ▶ Python Language Modules
- ▶ R Language Modules
- ▶ Statistical Functions
- ▶ Text Analytics
- ▶ Web Service
- ▶ Deprecated

Binary Classification: Direct marketing

In draft

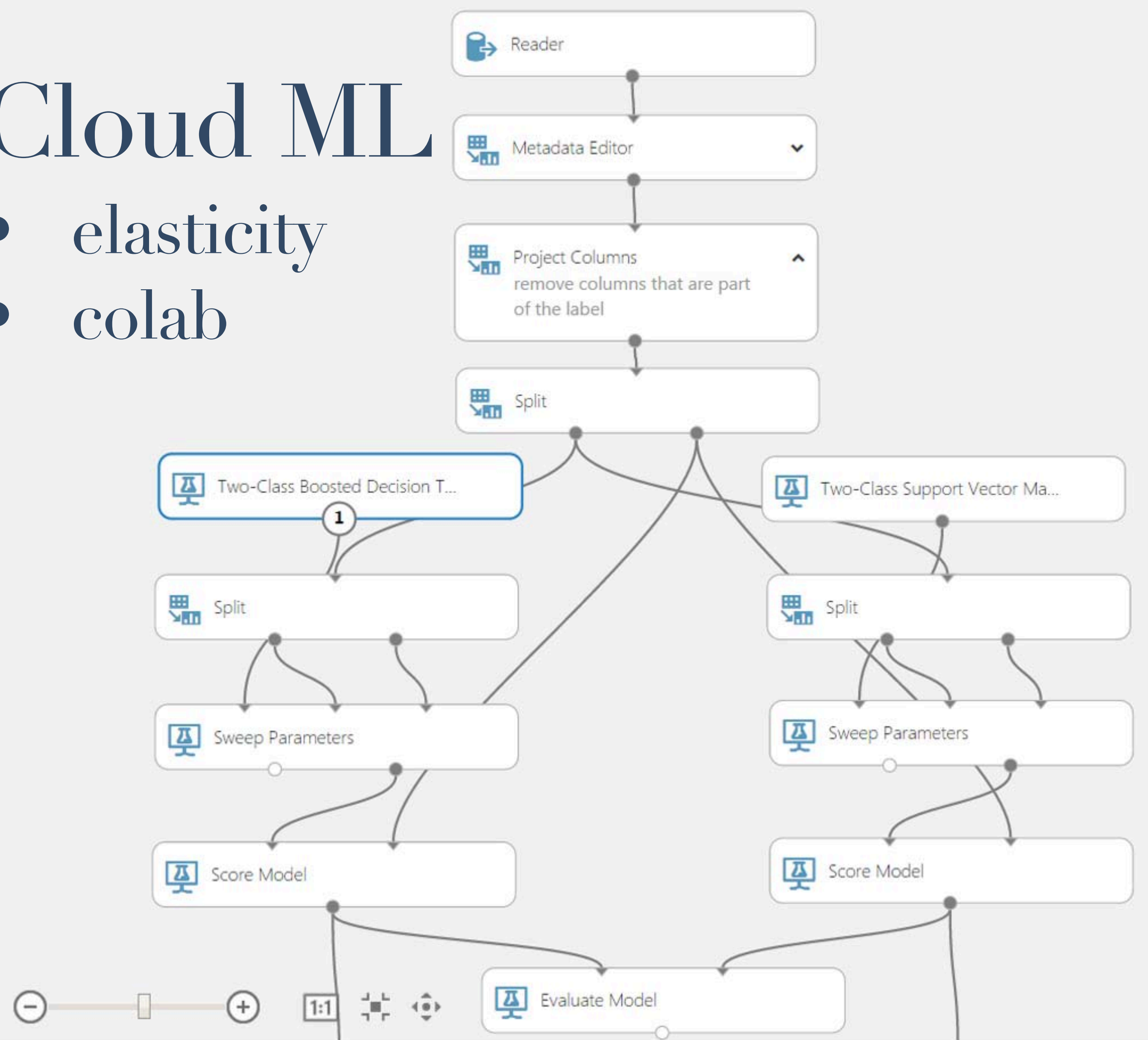
Properties

Two-Class Boosted Decision Tree

- Create trainer mode: Single Parameter
- Maximum number of leav...: 20
- Minimum number of sam...: 10
- Learning rate: 0.2
- Number of trees construct...: 100
- Random number seed: 0
- Allow unknown categ...

Cloud ML

- elasticity
- colab



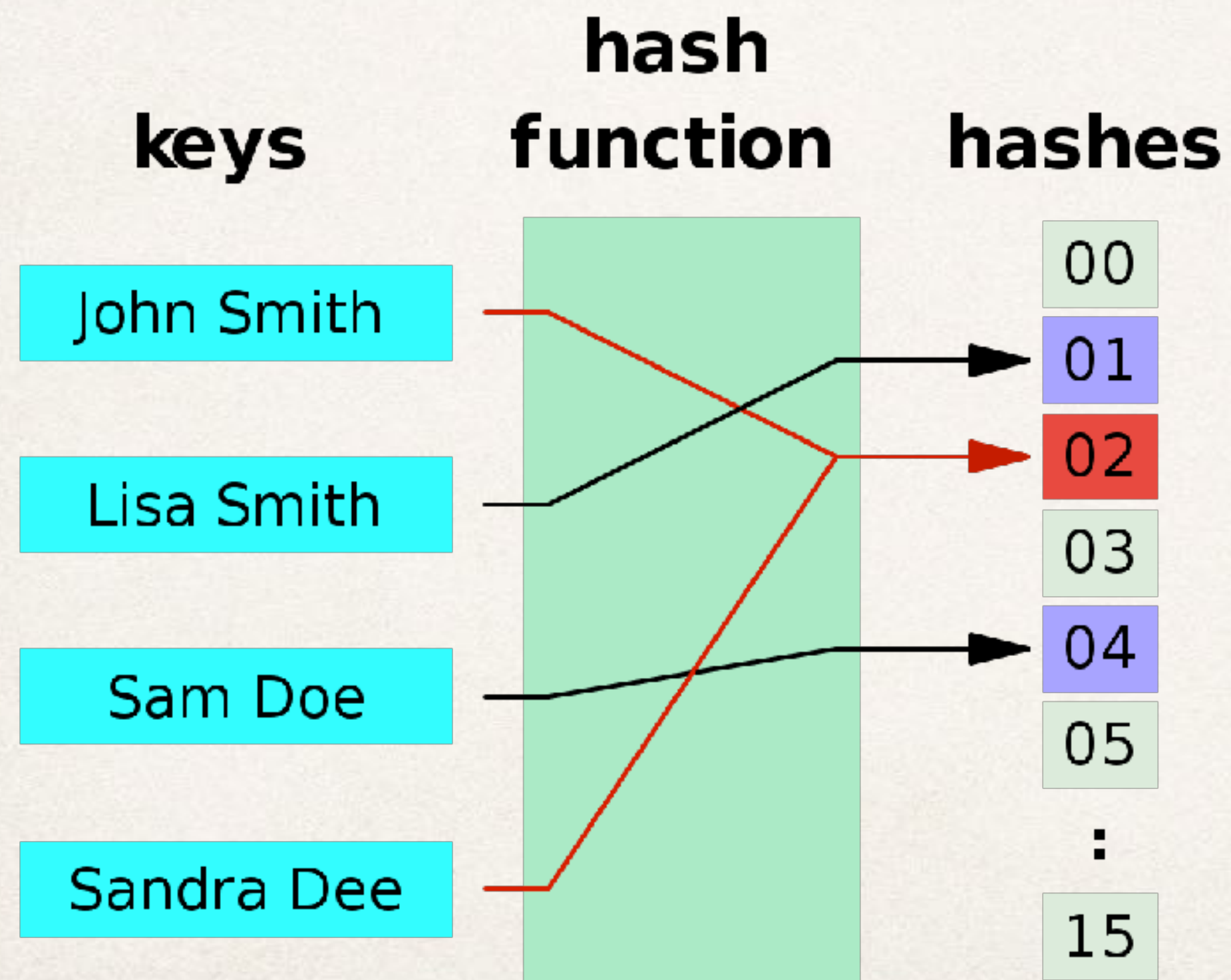
Quick Help

Creates a binary classifier using a boosted decision tree algorithm (more help...)

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