

Optimization for Machine Learning in Practice

Martin Jaggi

EPFL Machine Learning and Optimization Laboratory

mlo.epfl.ch

Where are we?



Machine
Learning

Optimization

Systems



Applications



Trends - General

- ❖ **privacy in ML**
- ❖ **decentralized training**
- ❖ **new hardware & systems**
- ❖ **trust in ML (e.g. robust & secure against adversarial attacks)**

Adversarial Attacks

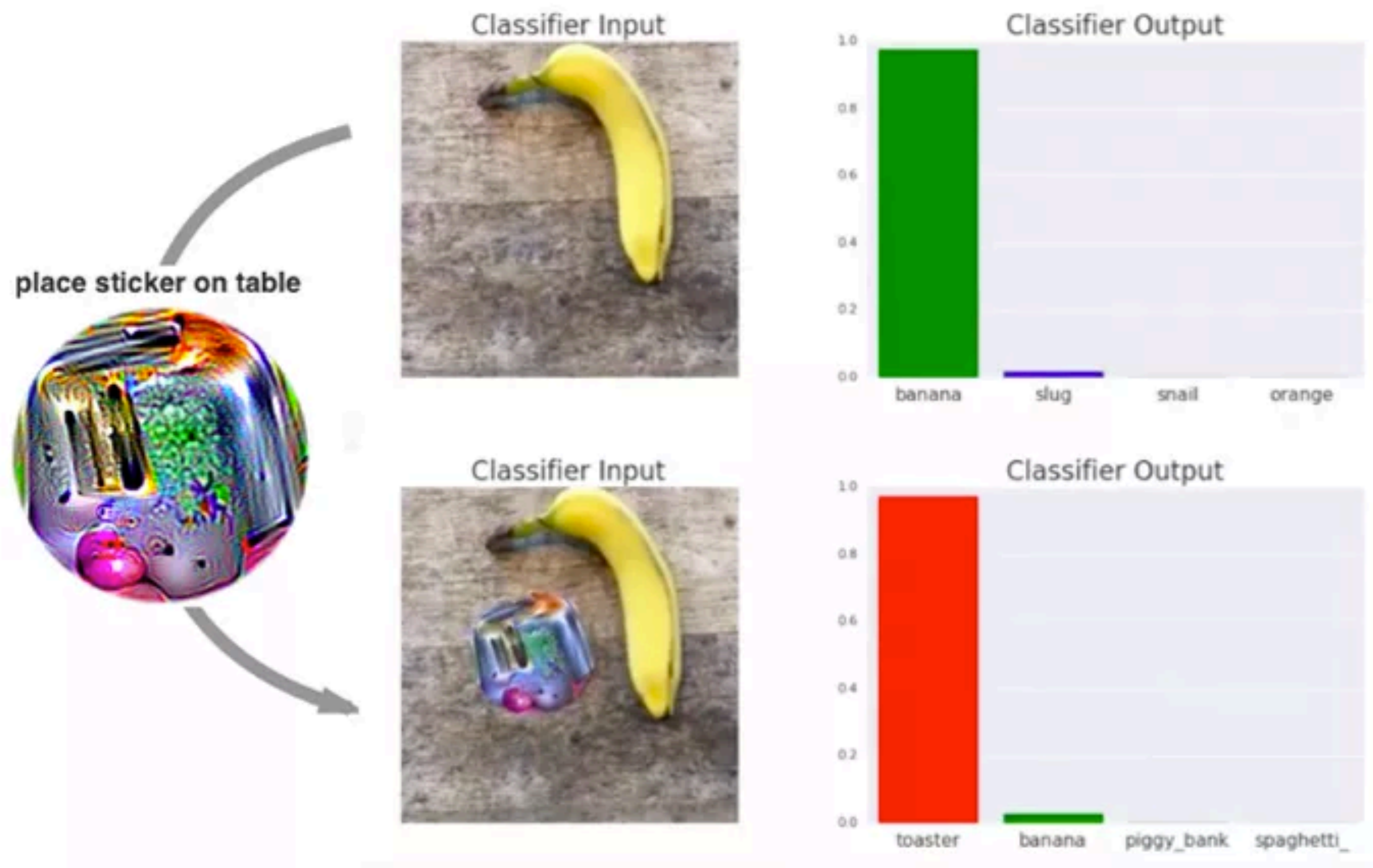


Image: Elsayed, Papernot et al 2018

Adversarial Attacks

- ❖ training

$$\min_w (f_w(\mathbf{x}_i) - y_i)^2$$

$$\nabla_w f$$

change model

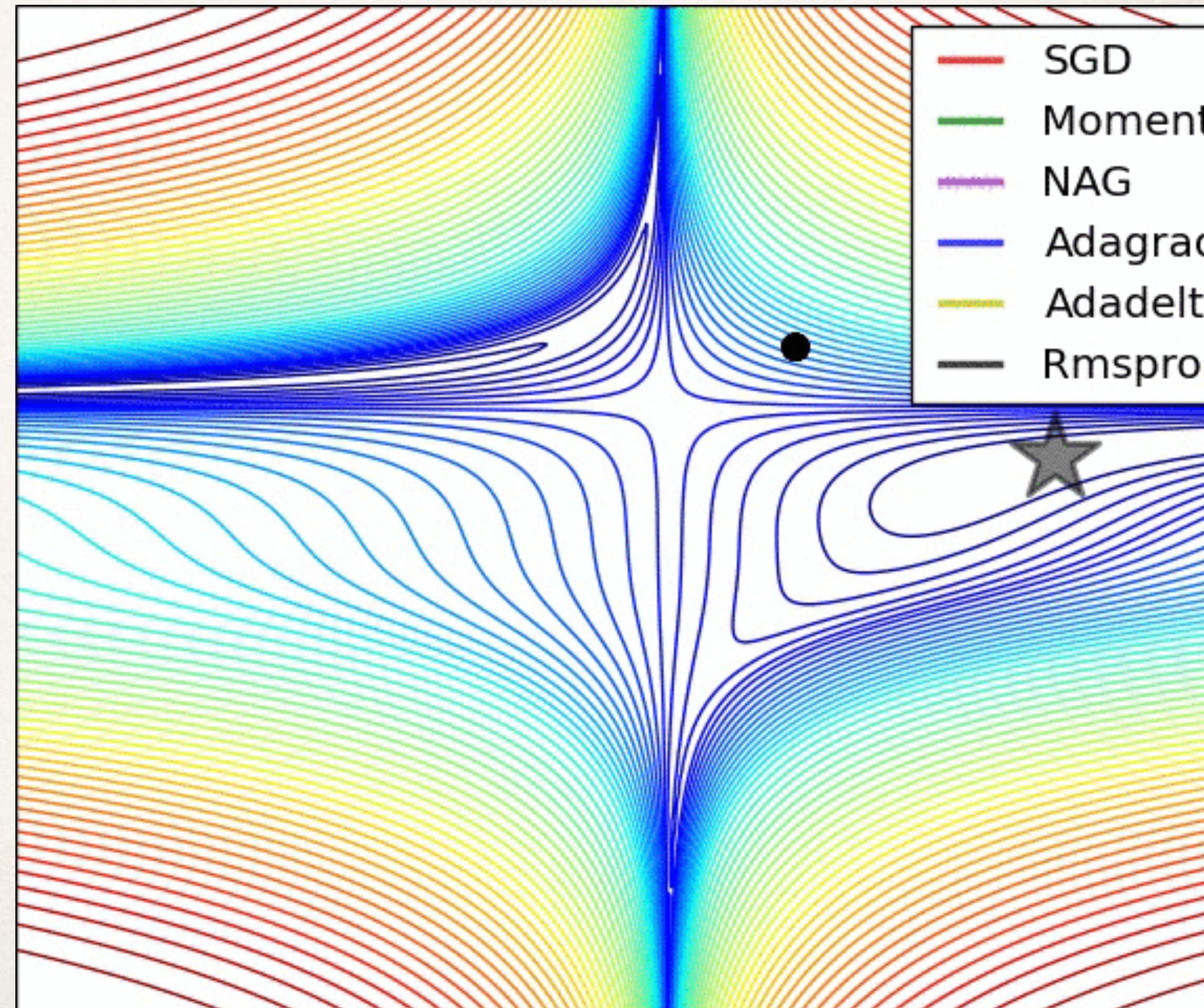
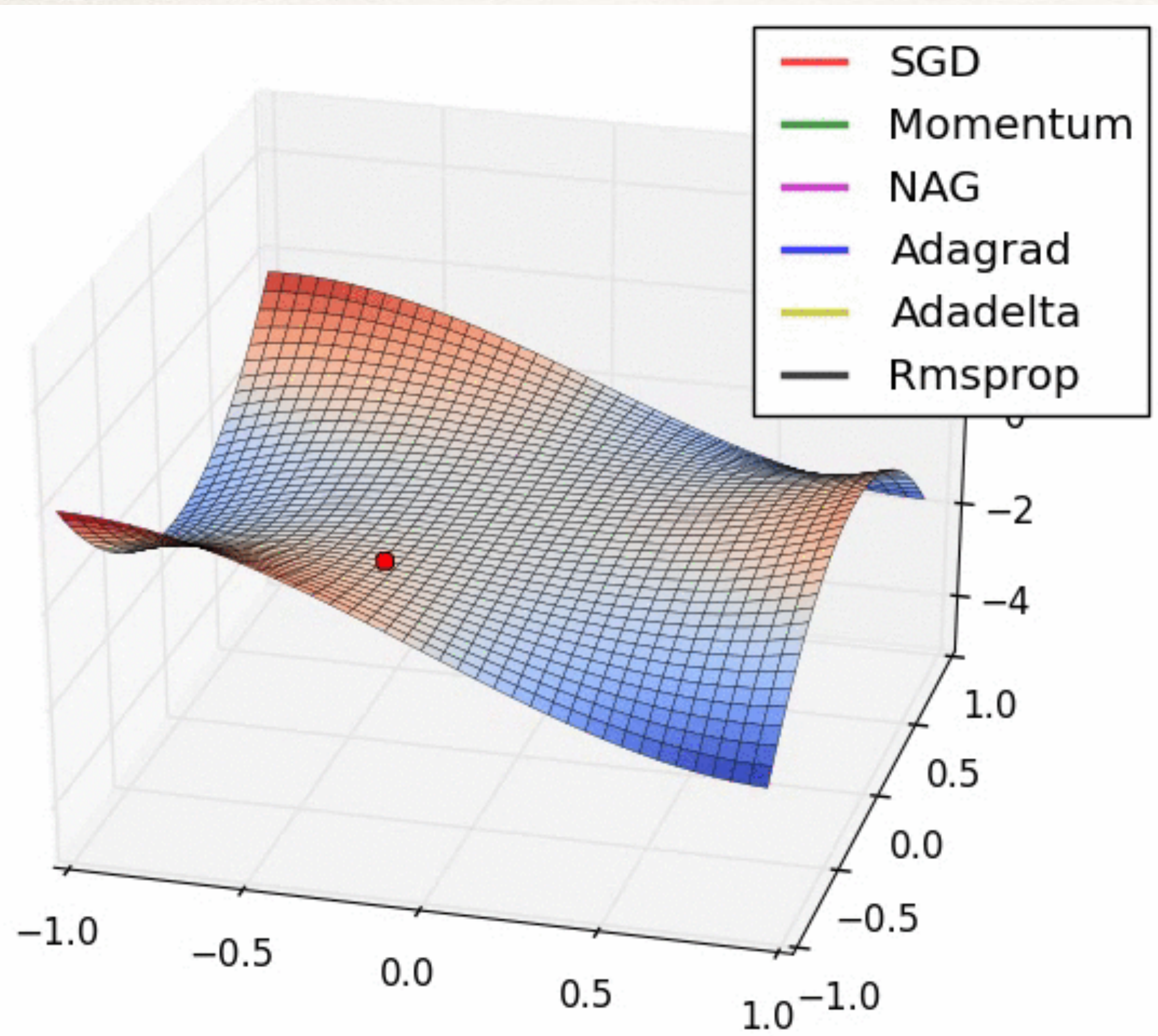
- ❖ attacking

$$\min_{\mathbf{x}_i} (f_w(\mathbf{x}_i) - y_i)^2$$

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

change data

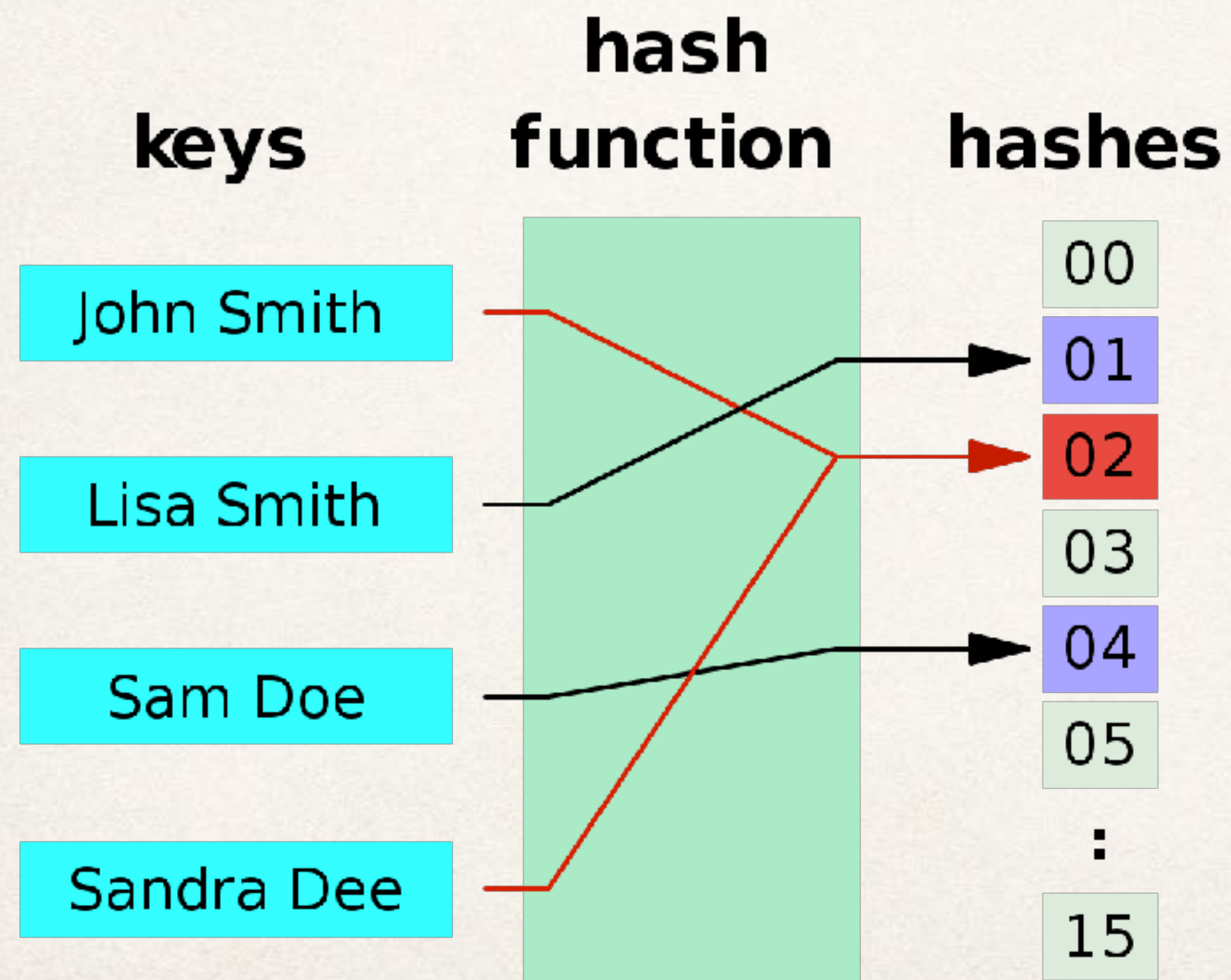
Practical comparison of algorithms



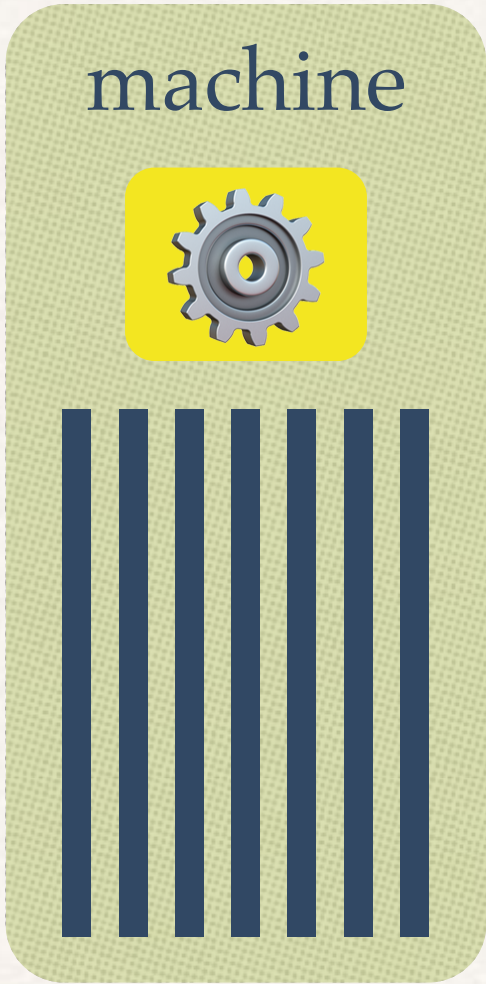
Practical tricks

❖ feature hashing

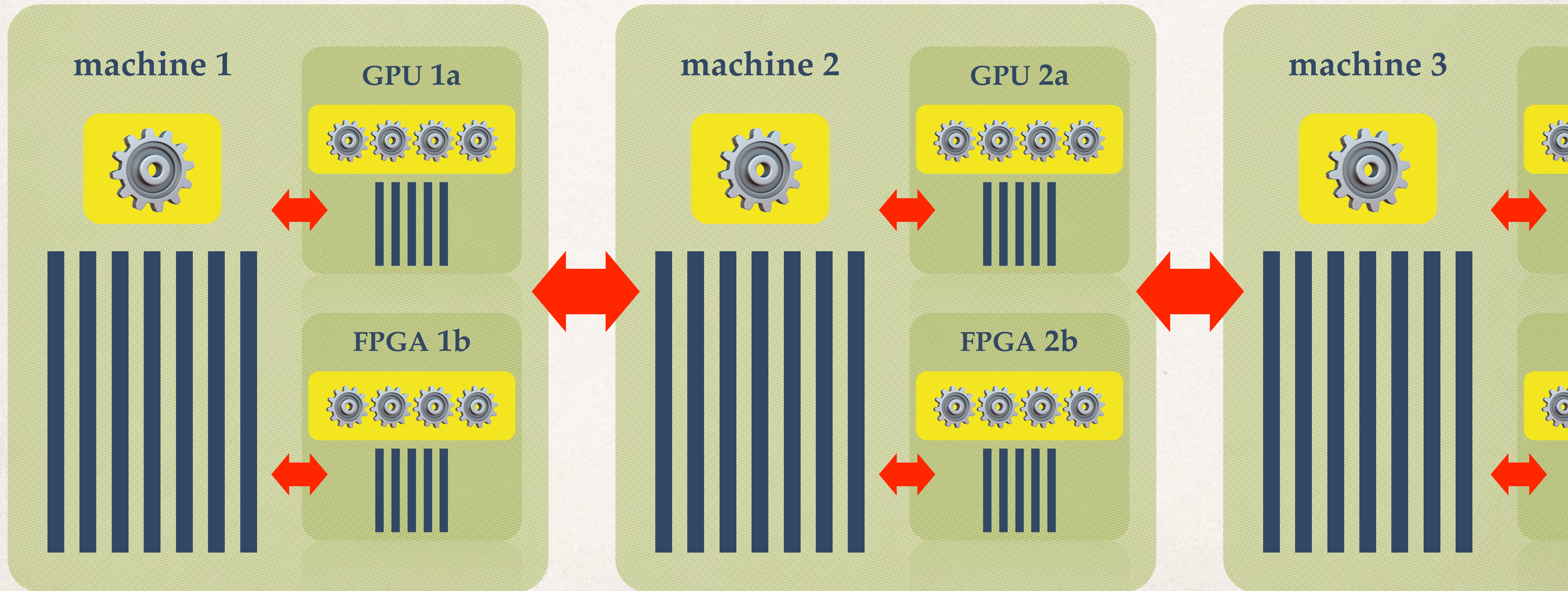
❖ limited precision operations



Systems ...then

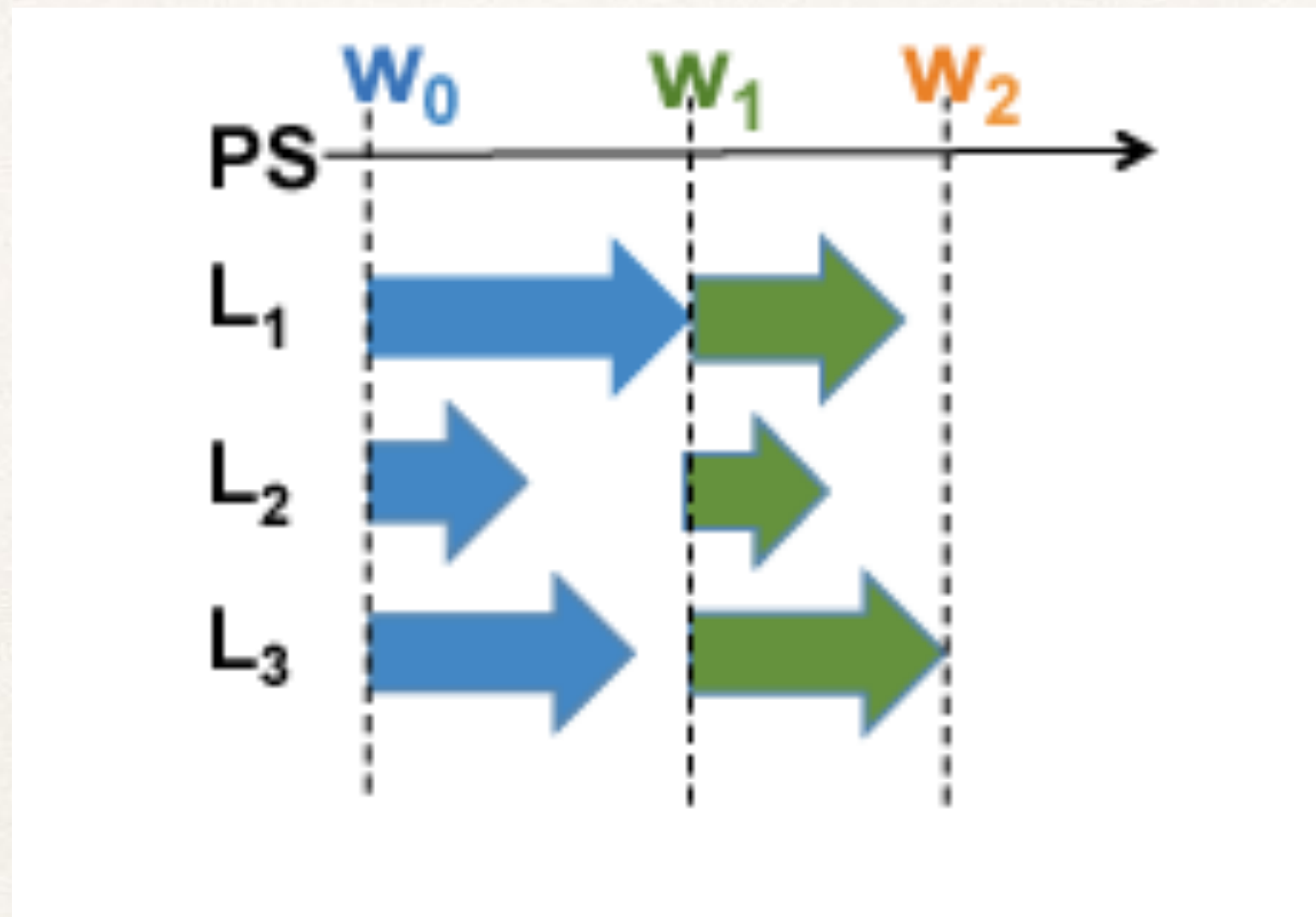


Systems ...now

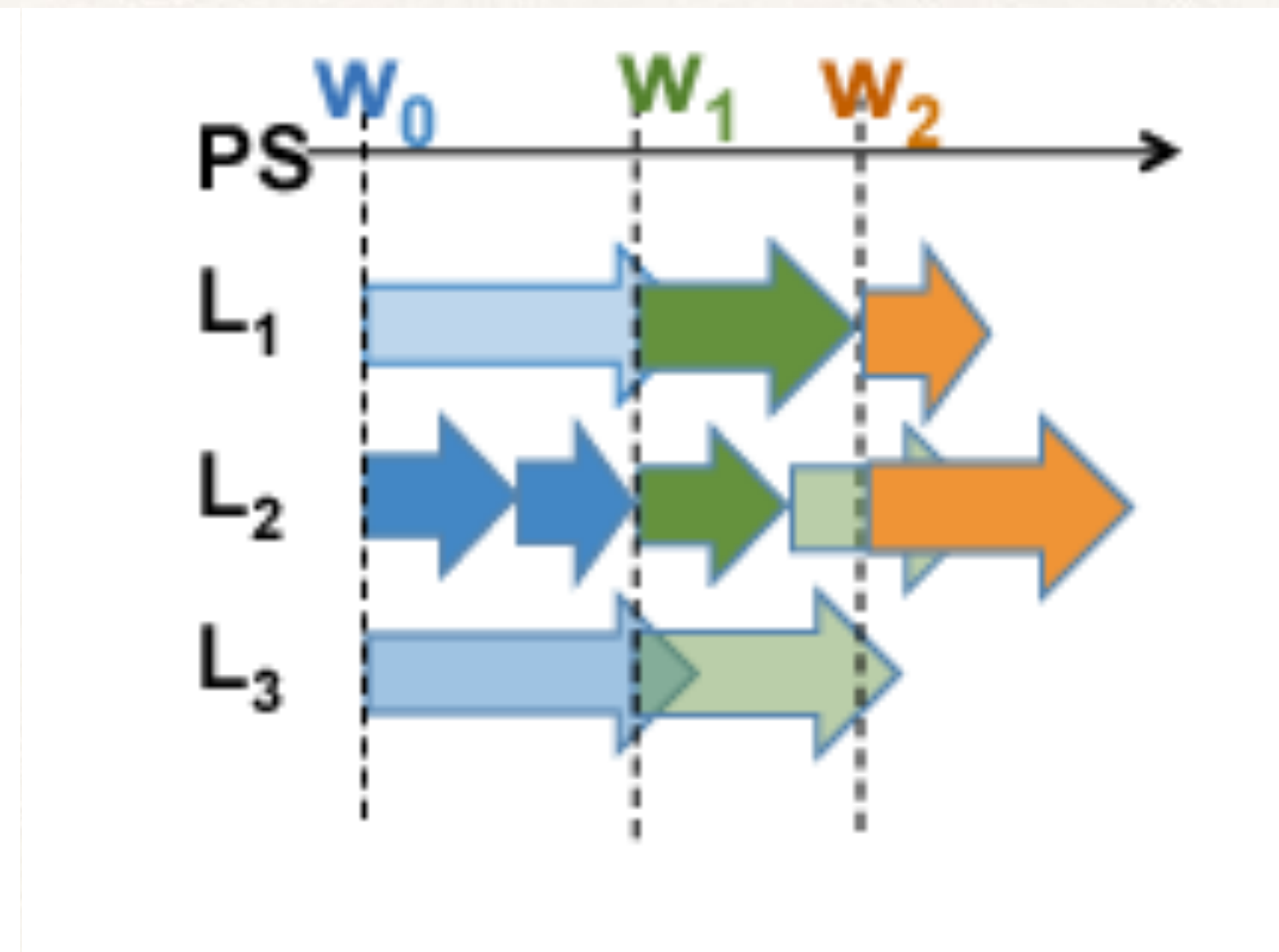


Parallel SGD

❖ Synchronous



❖ Asynchronous

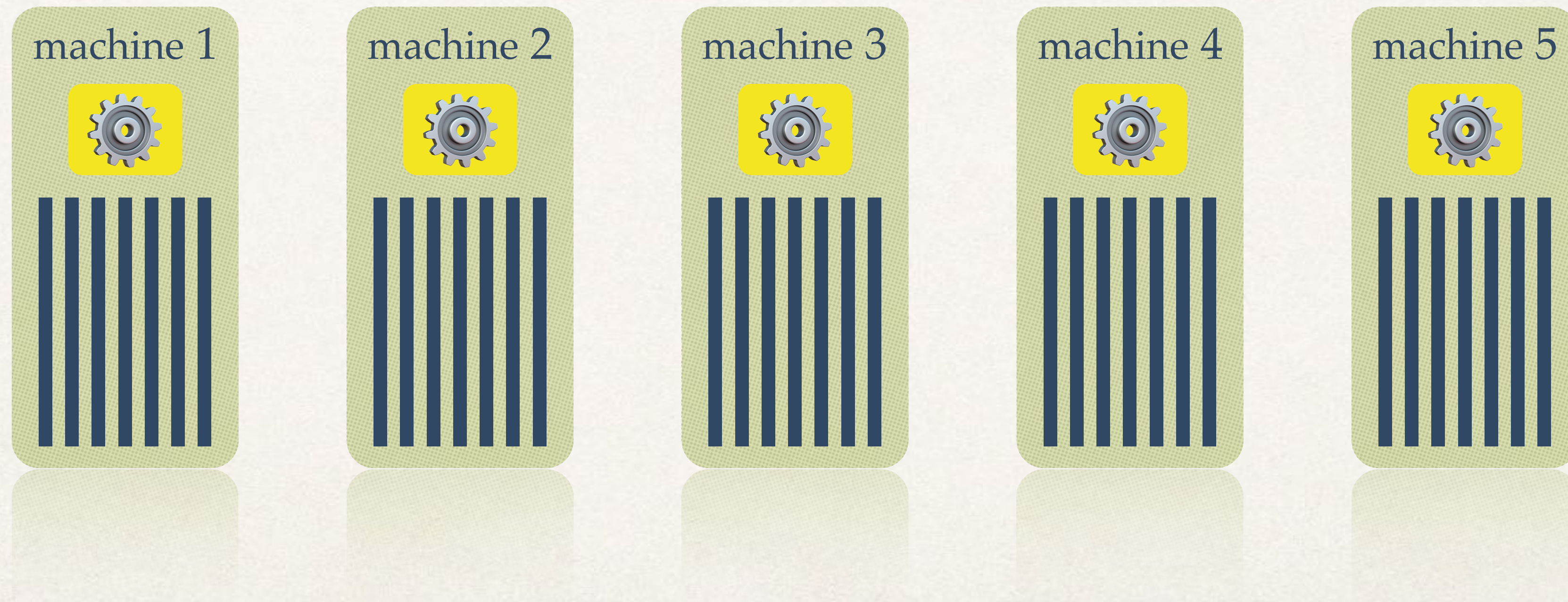


Mini-Batch!

1

Distributed

What if the data does not fit onto one device anymore?



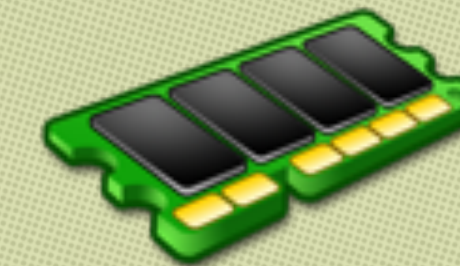
↔ Challenge

The Cost of Communication

$$v \in \mathbb{R}^{100}$$

- ❖ Reading v from memory (RAM)

100 ns



- ❖ Sending v to another machine

500'000 ns

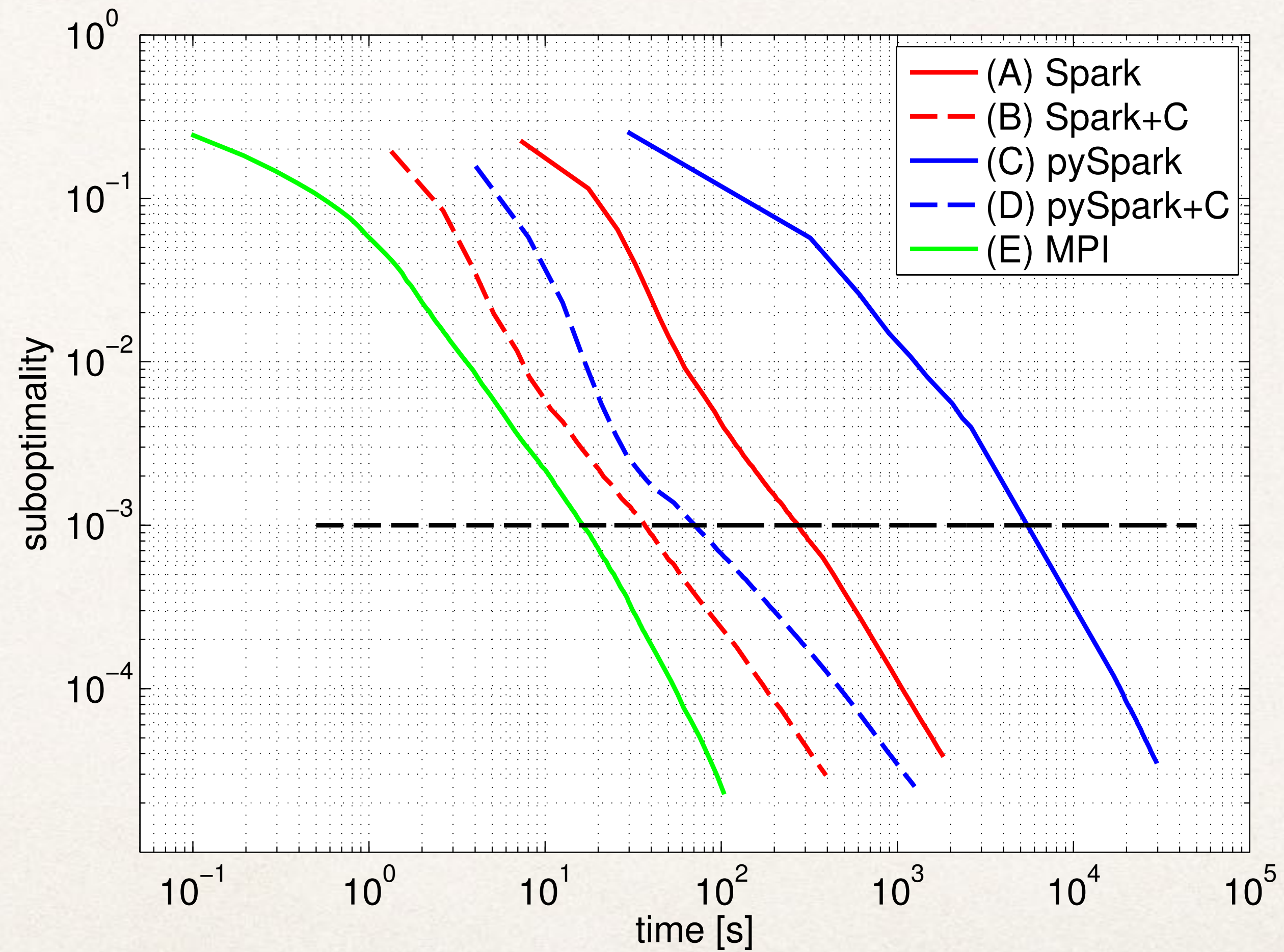
- ❖ Typical Map-Reduce iteration

10'000'000'000 ns

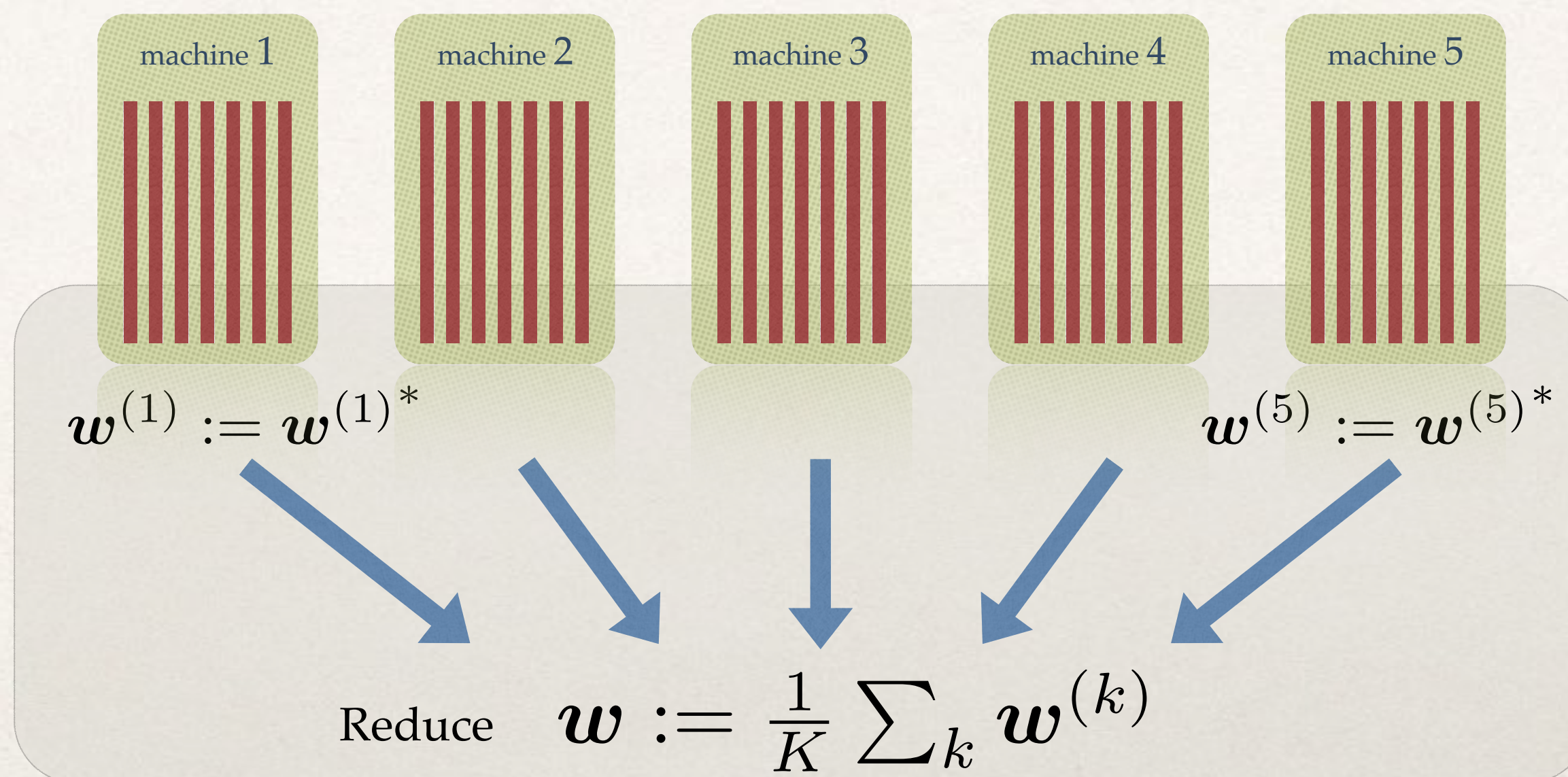
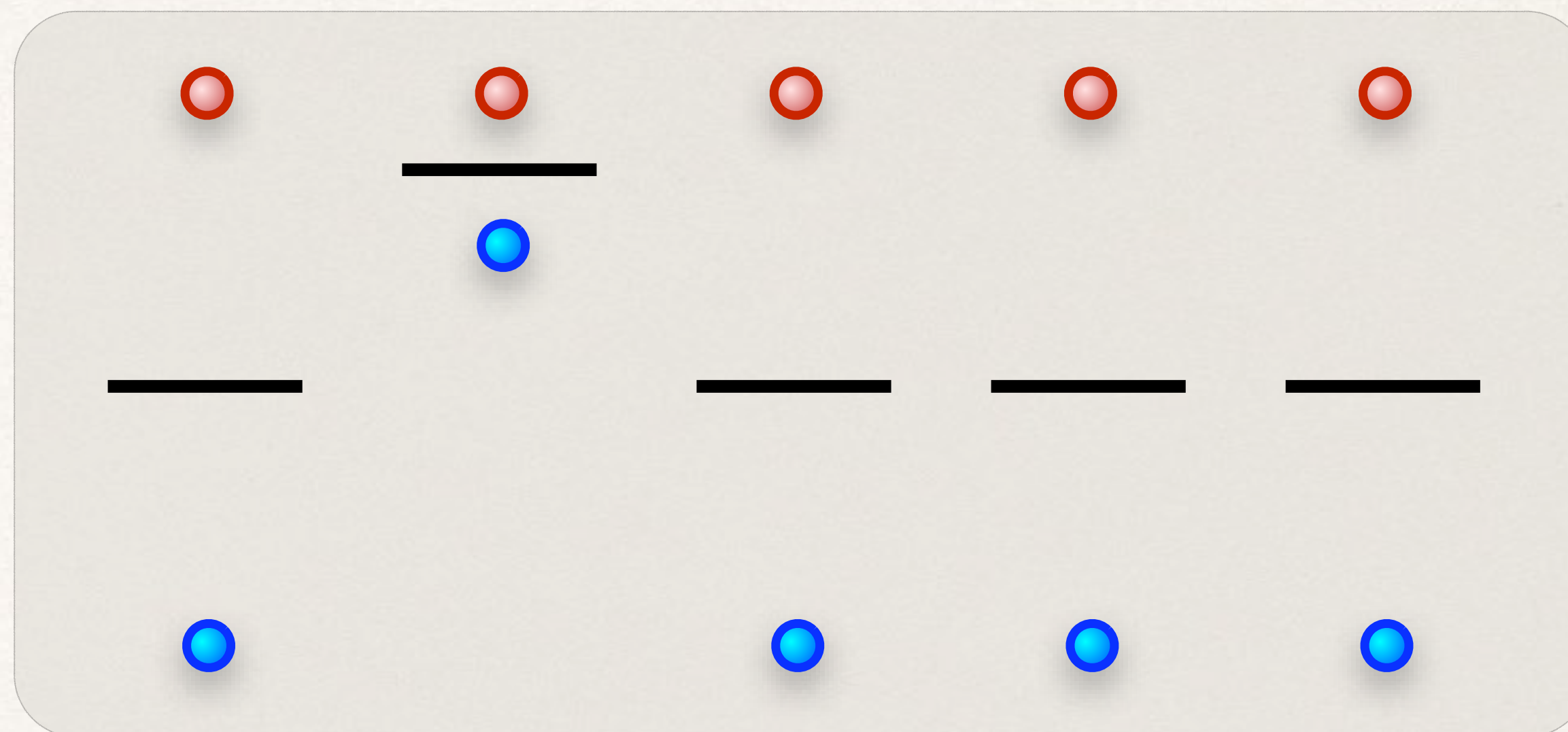


↔ Challenge

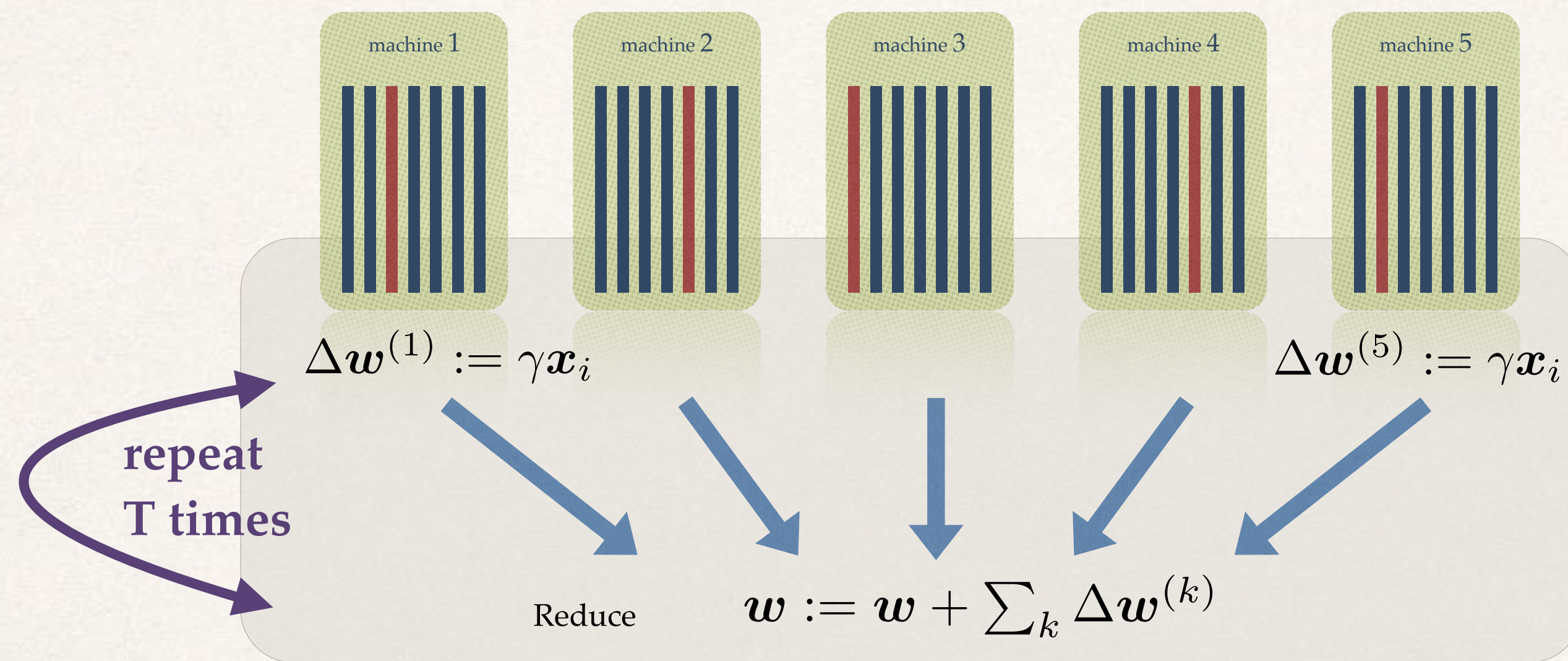
The Cost of Communication



One-Shot Averaging Does Not Work



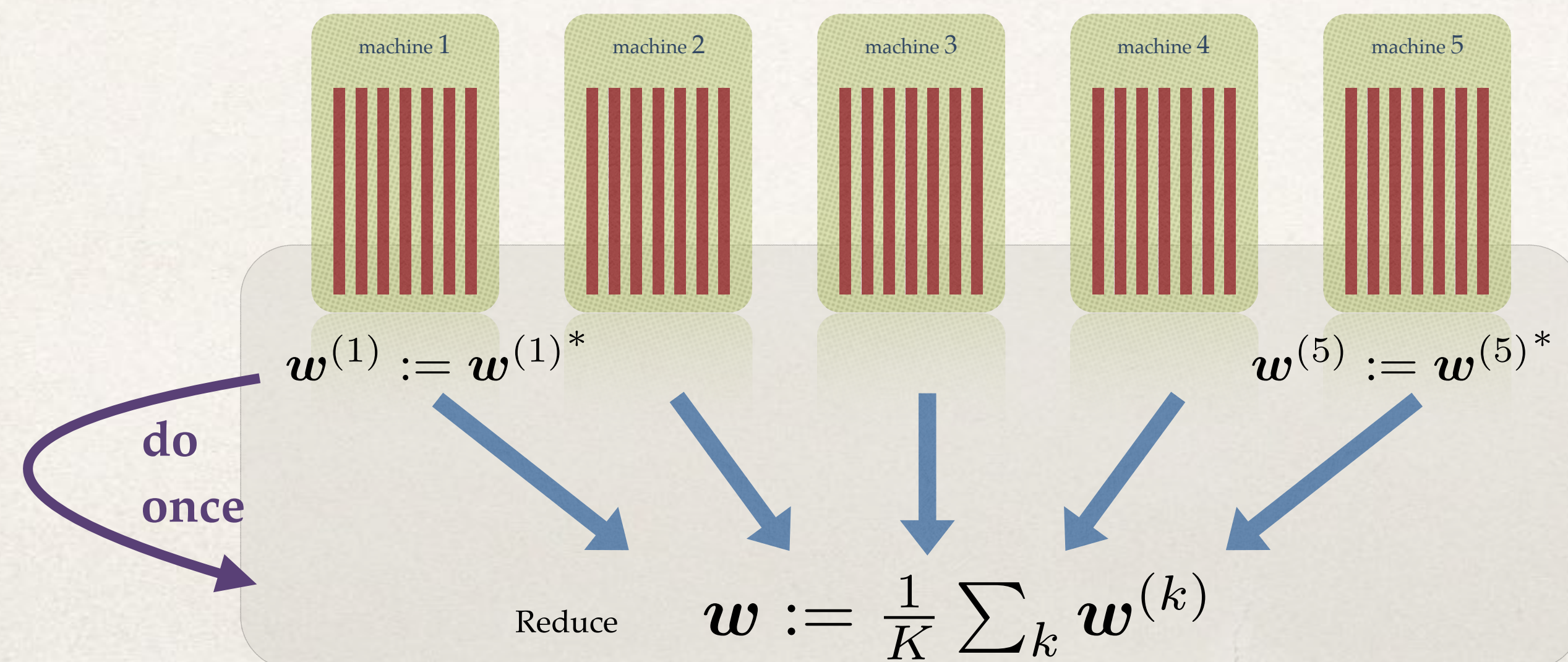
Communication: Always / Never



Naive Distributed SGD

local datapoints read: T
communications: T
convergence: ✓

“always communicate”



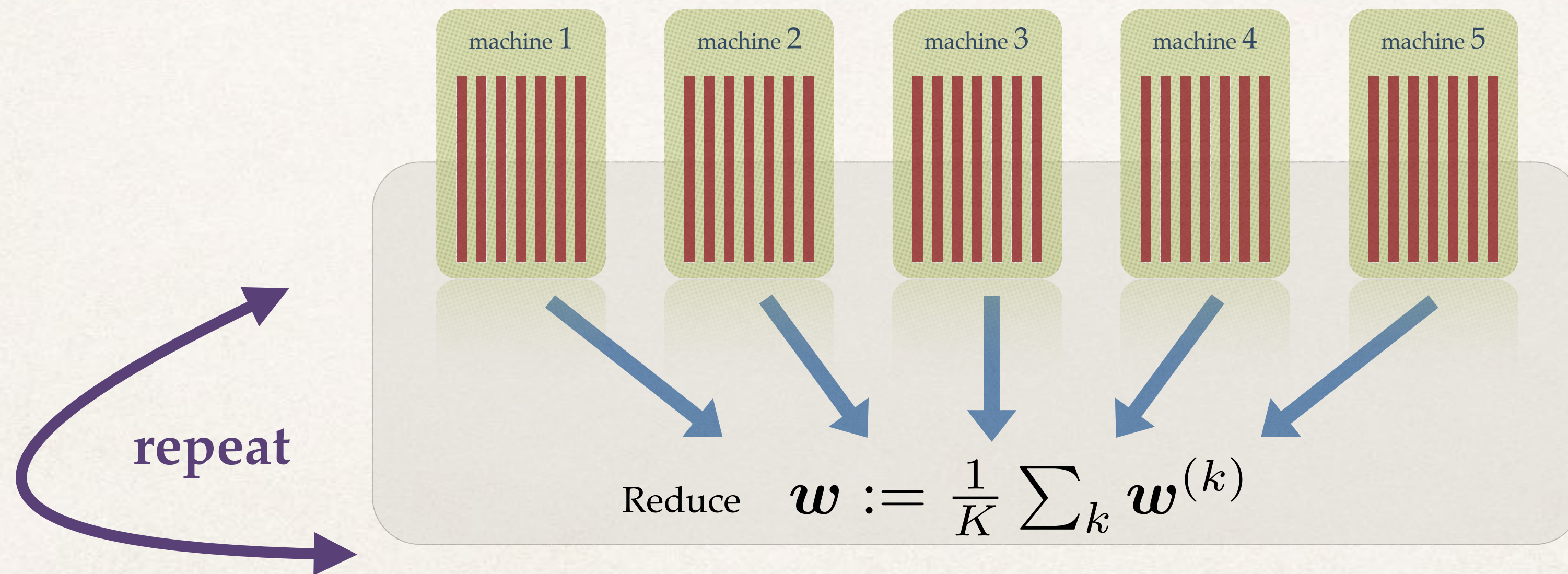
One-Shot Averaged Distributed Optimization

local datapoints read: T
communications: 1
convergence: ✗

“never communicate”

Distributed Full Gradient, L-BFGS

(just distribute the full gradient computation)



Problem class

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} f(A\boldsymbol{\alpha}) + g(\boldsymbol{\alpha})$$

Optimization: Primal-Dual Context

$$A_{\text{loc}} \Delta \alpha_{[k]} + w$$

$$\min_{\alpha \in \mathbb{R}^n} \left[\mathcal{O}_A(\alpha) := f(A\alpha) + g(\alpha) \right]$$

primal Lasso

dual L2-reg SVM/Log-Reg

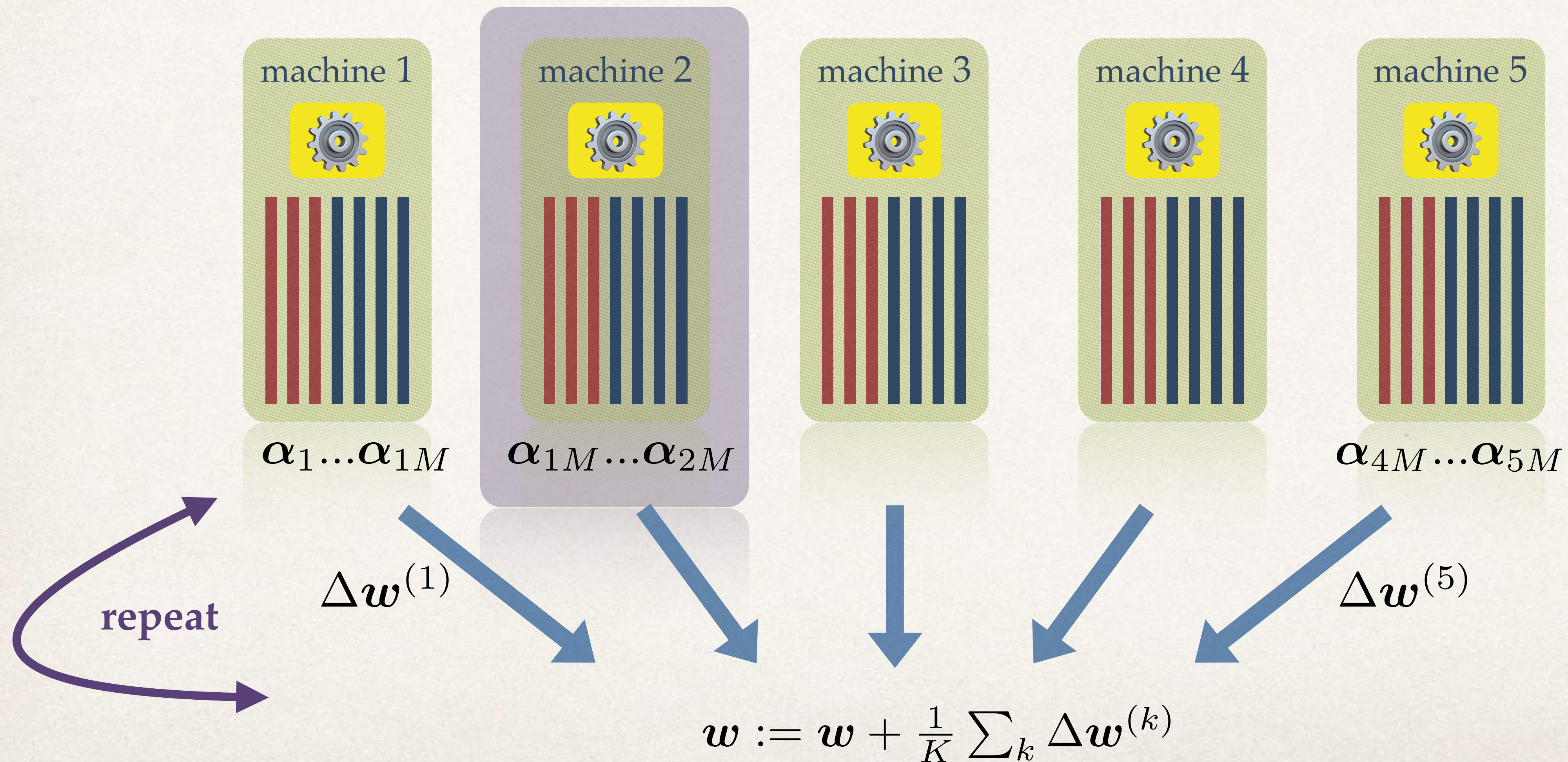
primal L1-reg SVM/Log-Reg

correspondence

$$w := \nabla f(A\alpha)$$

$$\min_{w \in \mathbb{R}^d} \left[\mathcal{O}_B(w) := g^*(-A^\top w) + f^*(w) \right]$$

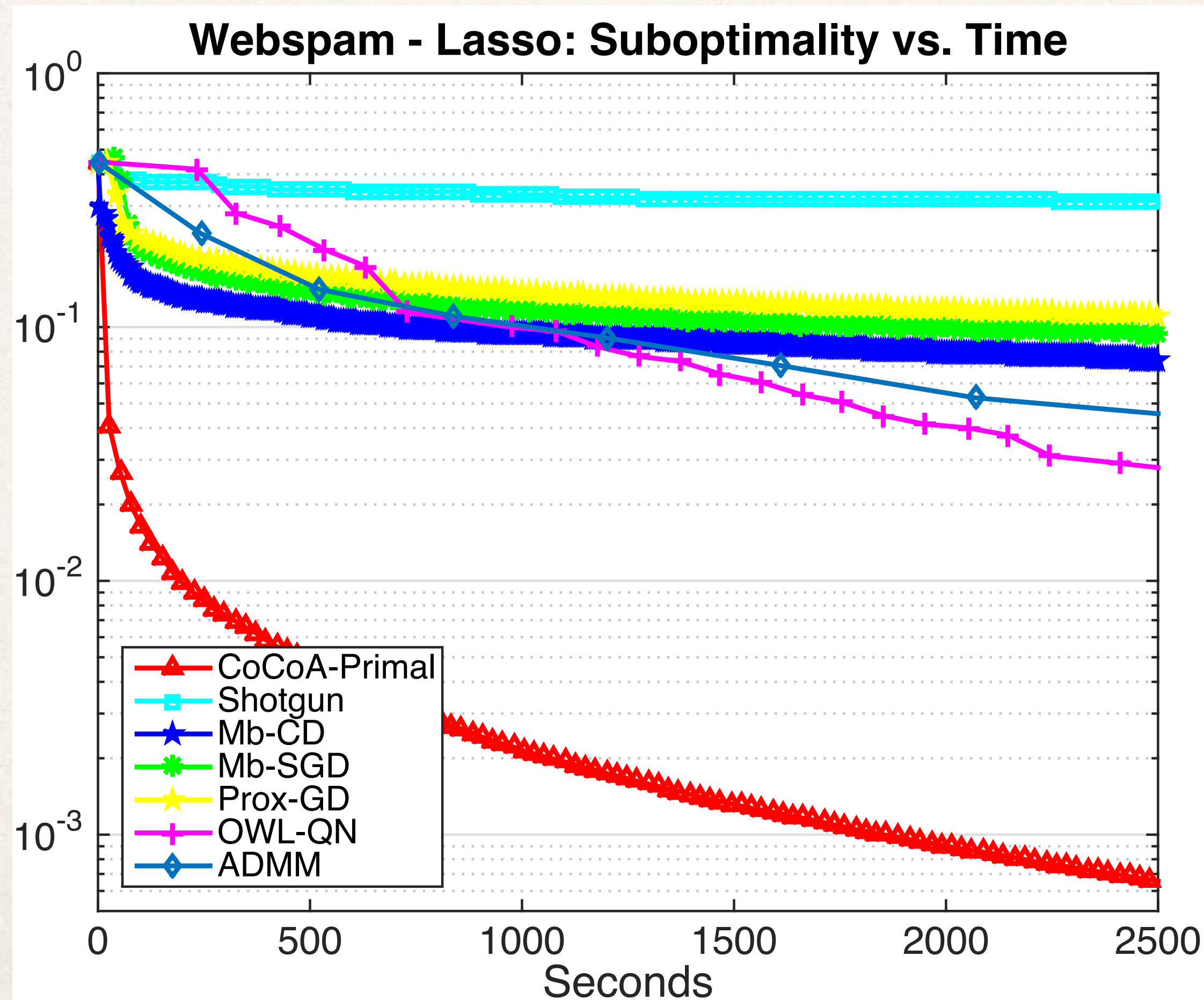
CoCoA - Communication Efficient Distributed Optimization



Distributed Experiments

L1-Regularized Linear Regression

Dataset	Training	Features	Sparsity
url	2,396,130	3,231,961	3.5e-3%
epsilon	400,000	2,000	100%
kddb	19,264,097	29,890,095	9.8e-5%
webspam	350,000	16,609,143	0.02%



NIPS 2014, ICML 2015, JMLR 2018
arxiv.org/abs/1611.02189

- part of TensorFlow core (L2)
- code in pytorch, TF, spark, C (also L1)

Summary

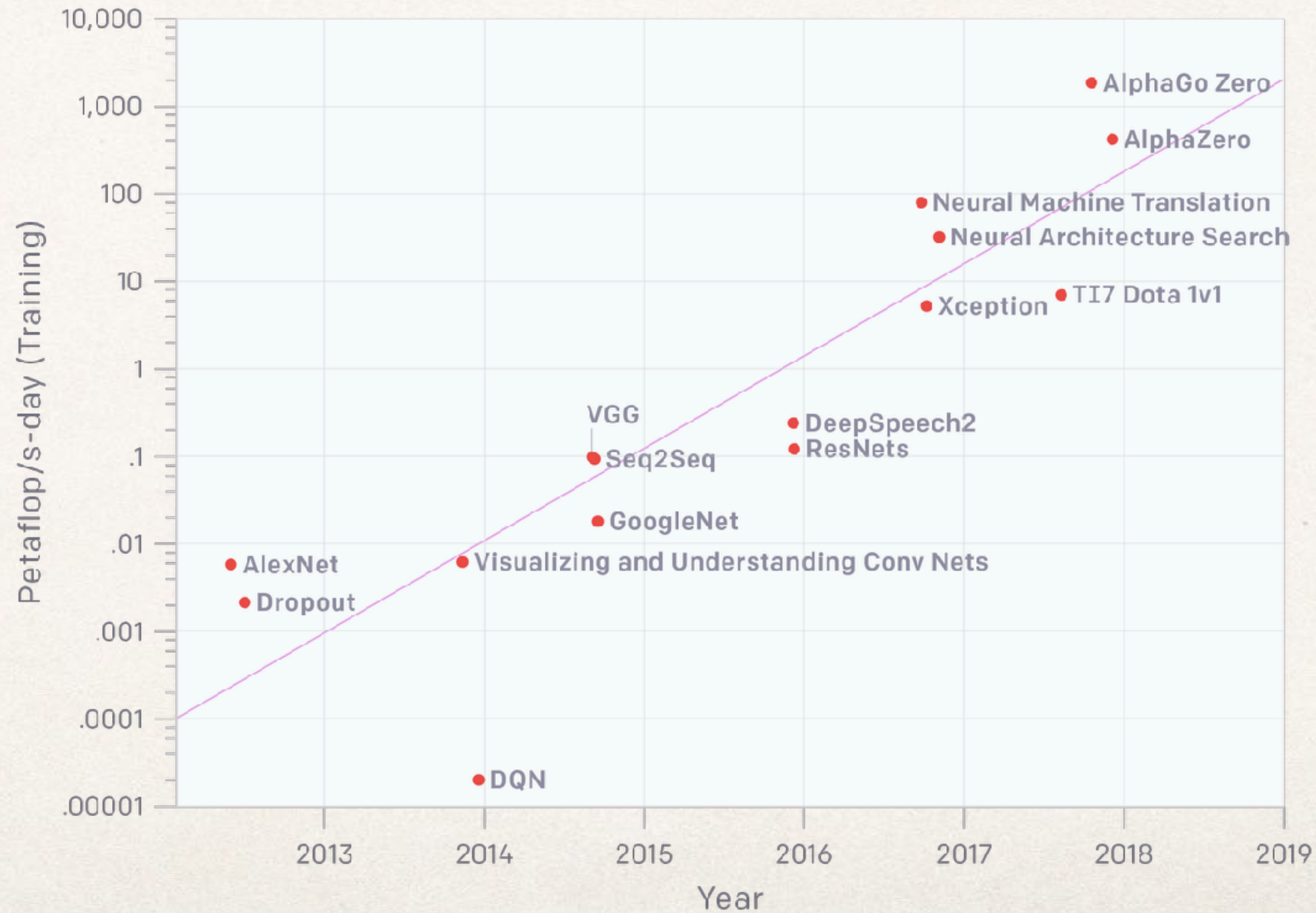
- ❖ **adaptivity** to the communication cost
- ❖ **re-usability** of good existing solvers
- ❖ **accuracy** certificates
- ❖ **second-order** and **trust-region** version (local Hessian)

Next Steps

- ❖ **adaptivity** to the degree of separability
- ❖ generalization to **deep learning, SGD**
- ❖ **decentralized version** (communication on graph)
- ❖ **benchmarking & code**

Deep Learning

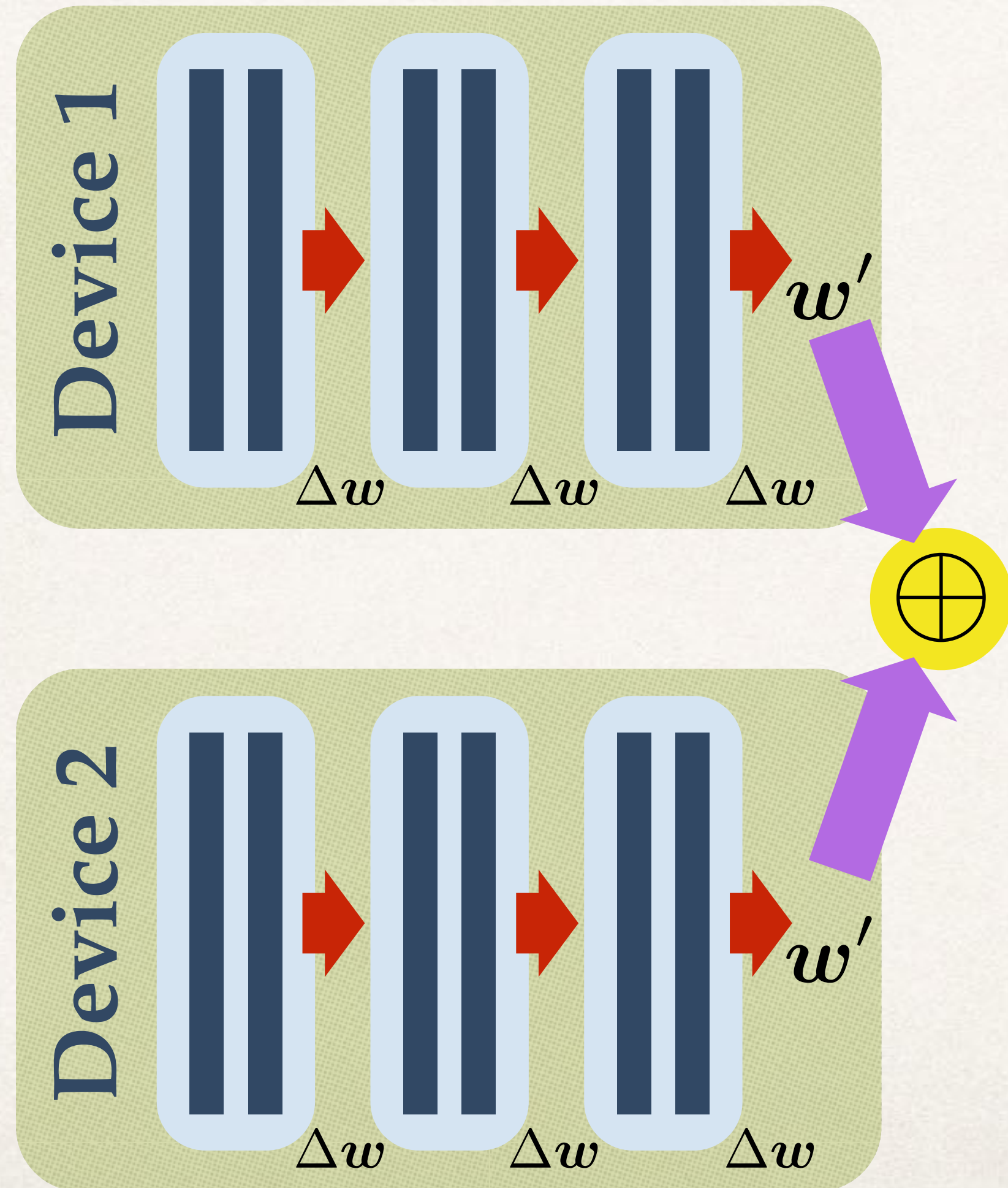
AlexNet to AlphaGo Zero: A 300,000x Increase in Compute



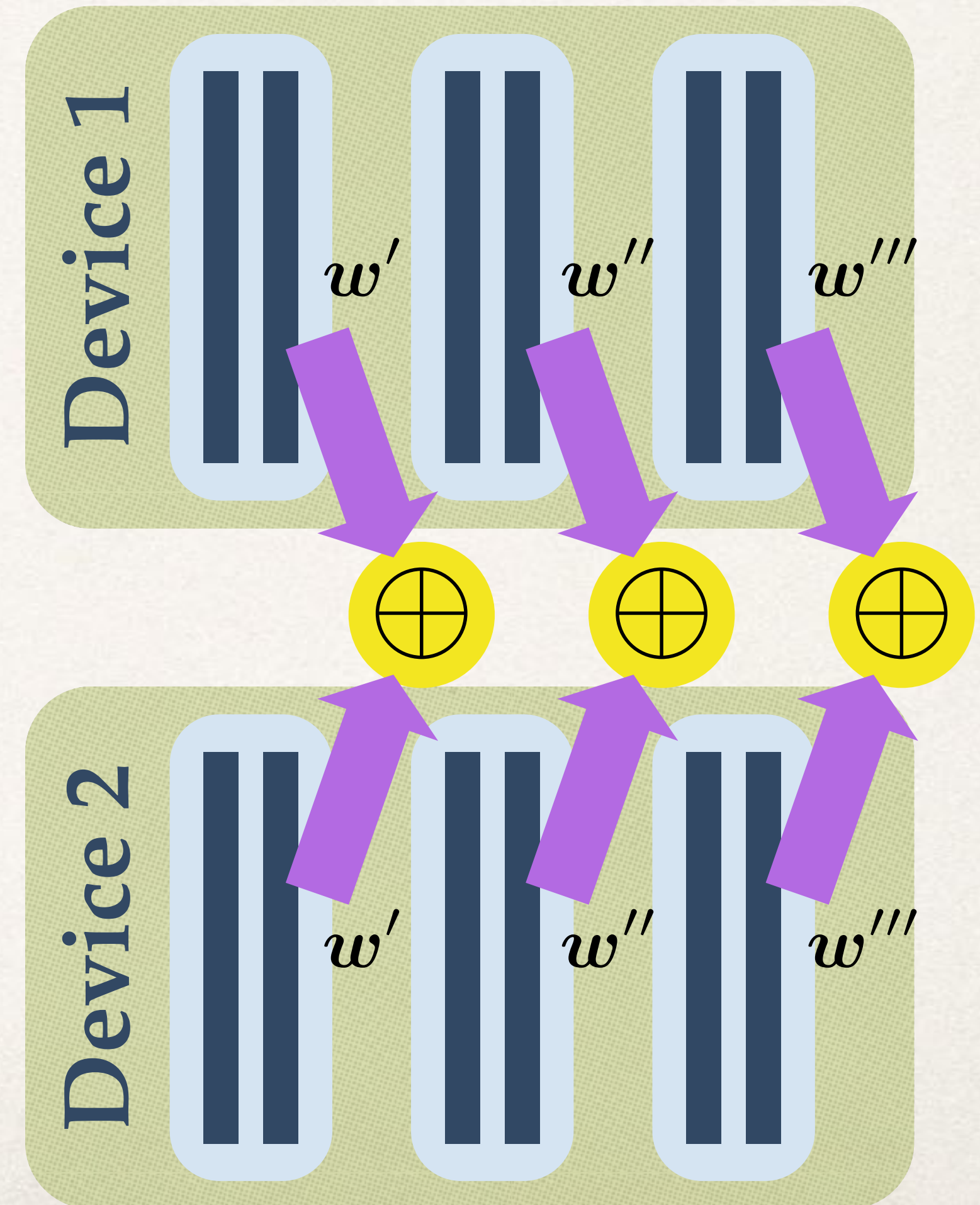
(Data Parallel)

Distributed DL

Local SGD

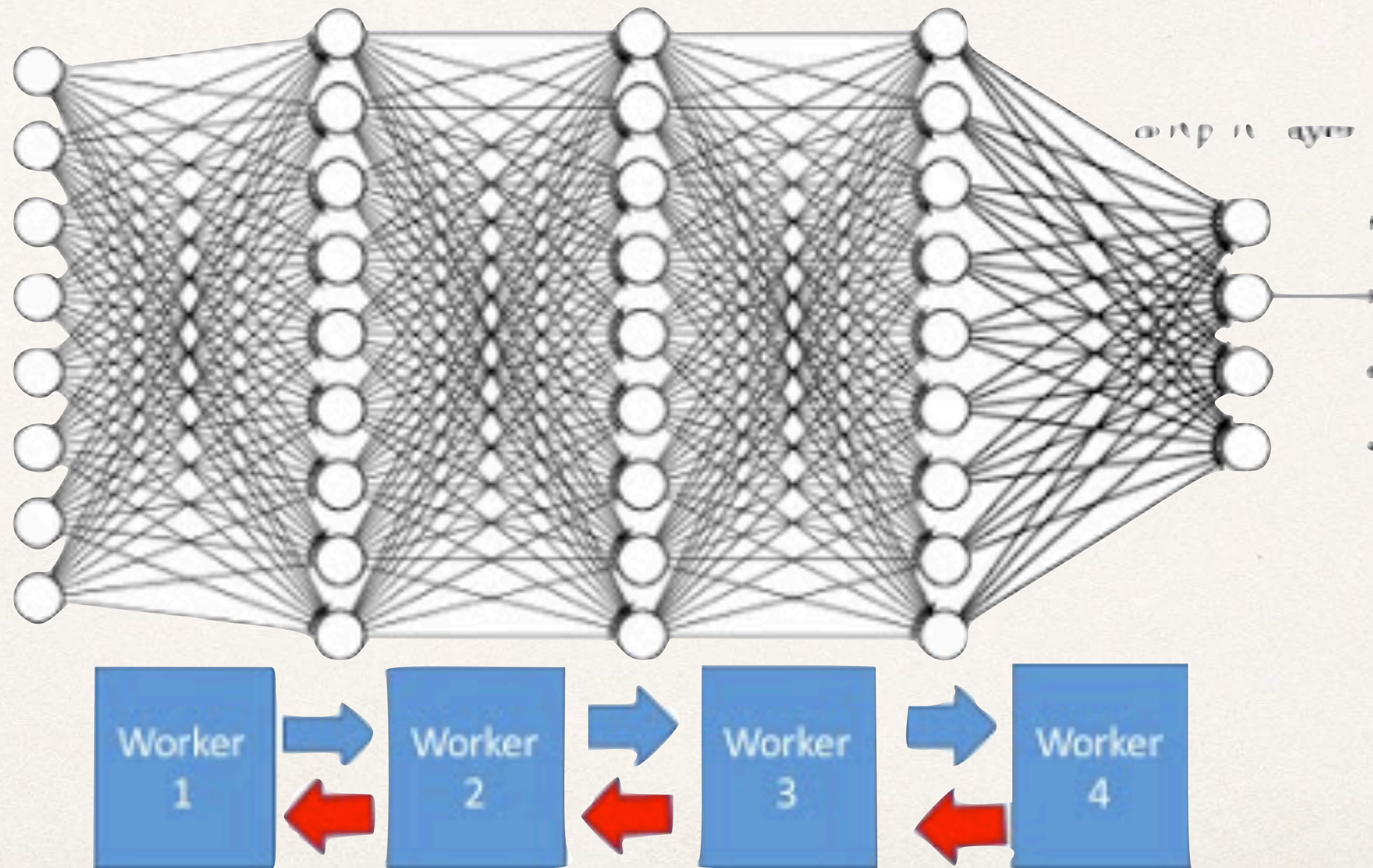


Mini-batch SGD



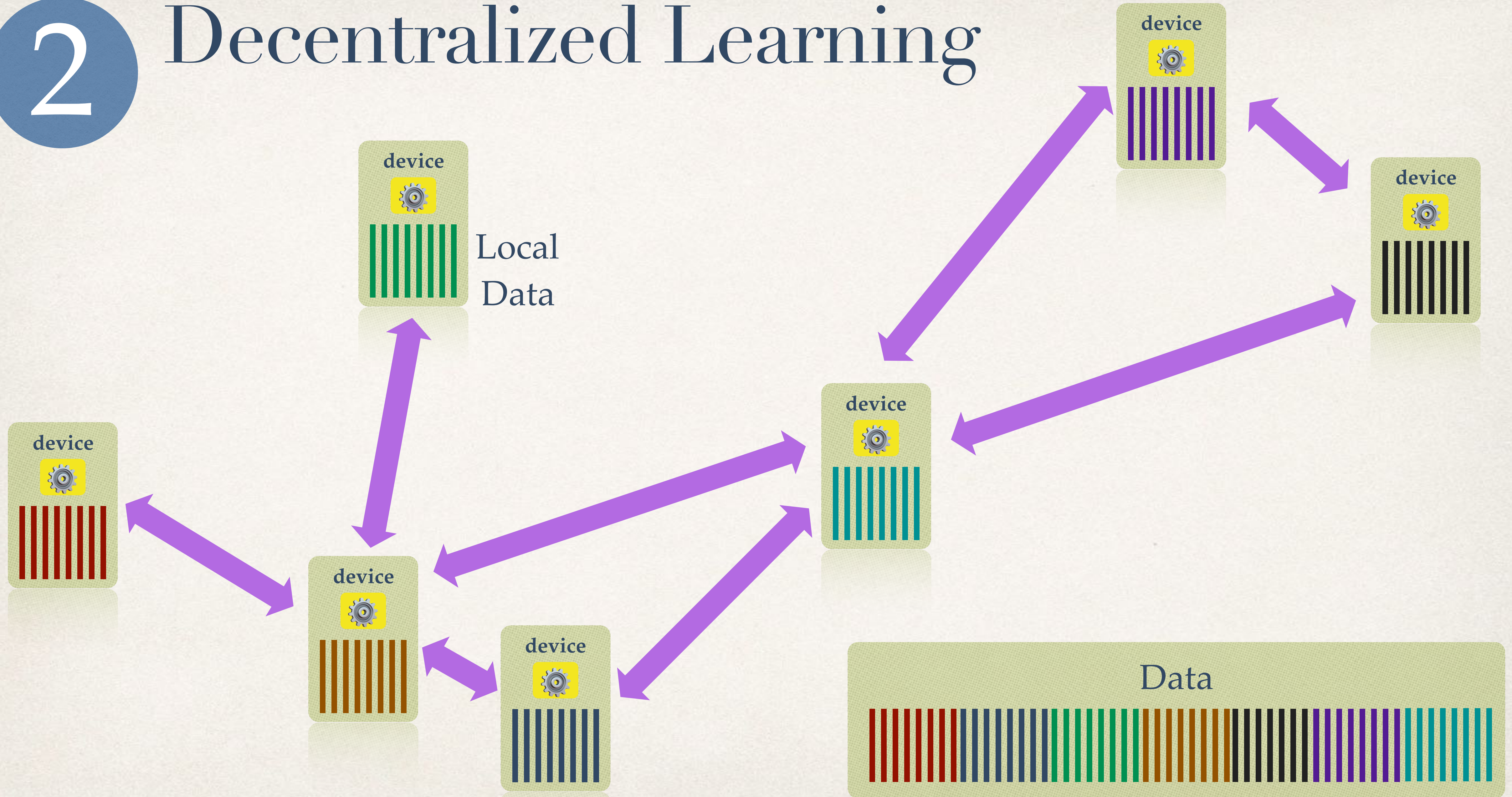
(Model Parallel)

Distributed DL



2

Decentralized Learning



Motivation



[image source](#)

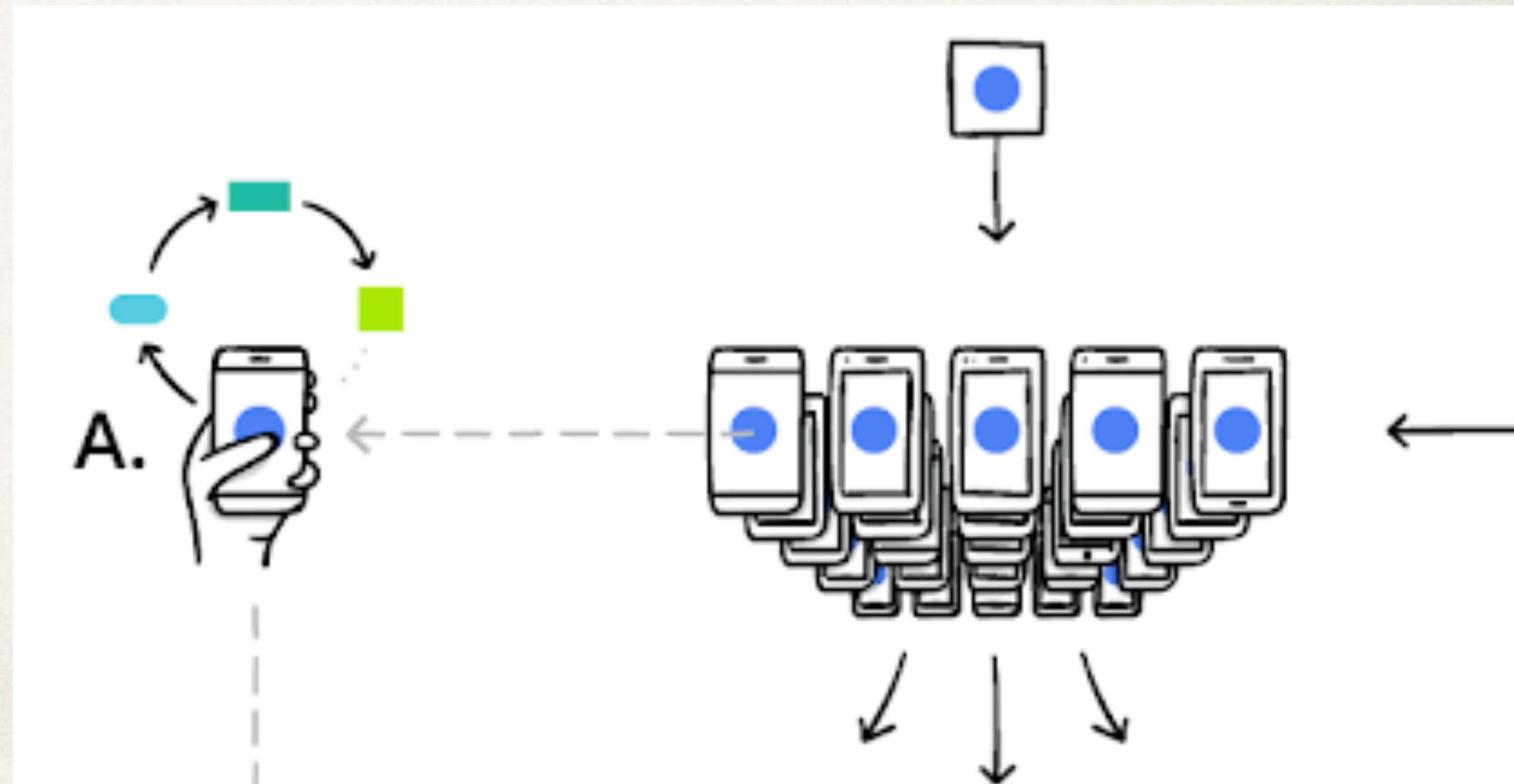
- ❖ Medical data is very sensitive
- ❖ Data cannot be sent outside of the hospital
- ❖ People could use the data to prevent diseases

Motivation



[image source](#)

- ❖ learn from users writing on smartphones
- ❖ data is very sensitive



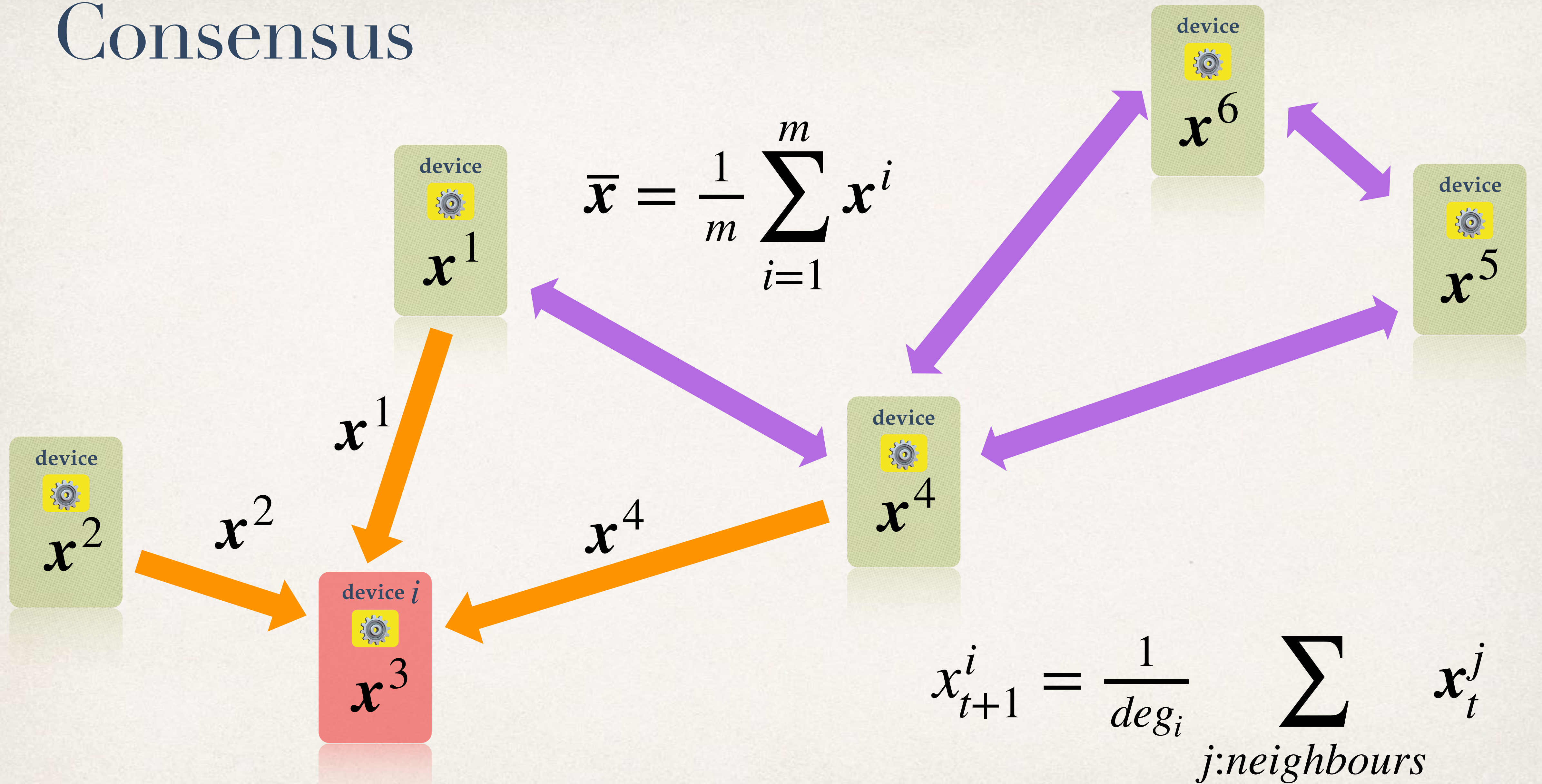
Motivation



[image source](#)

- ❖ Decentralized learning can speed up training time in datacenters

Consensus



Communication Compression

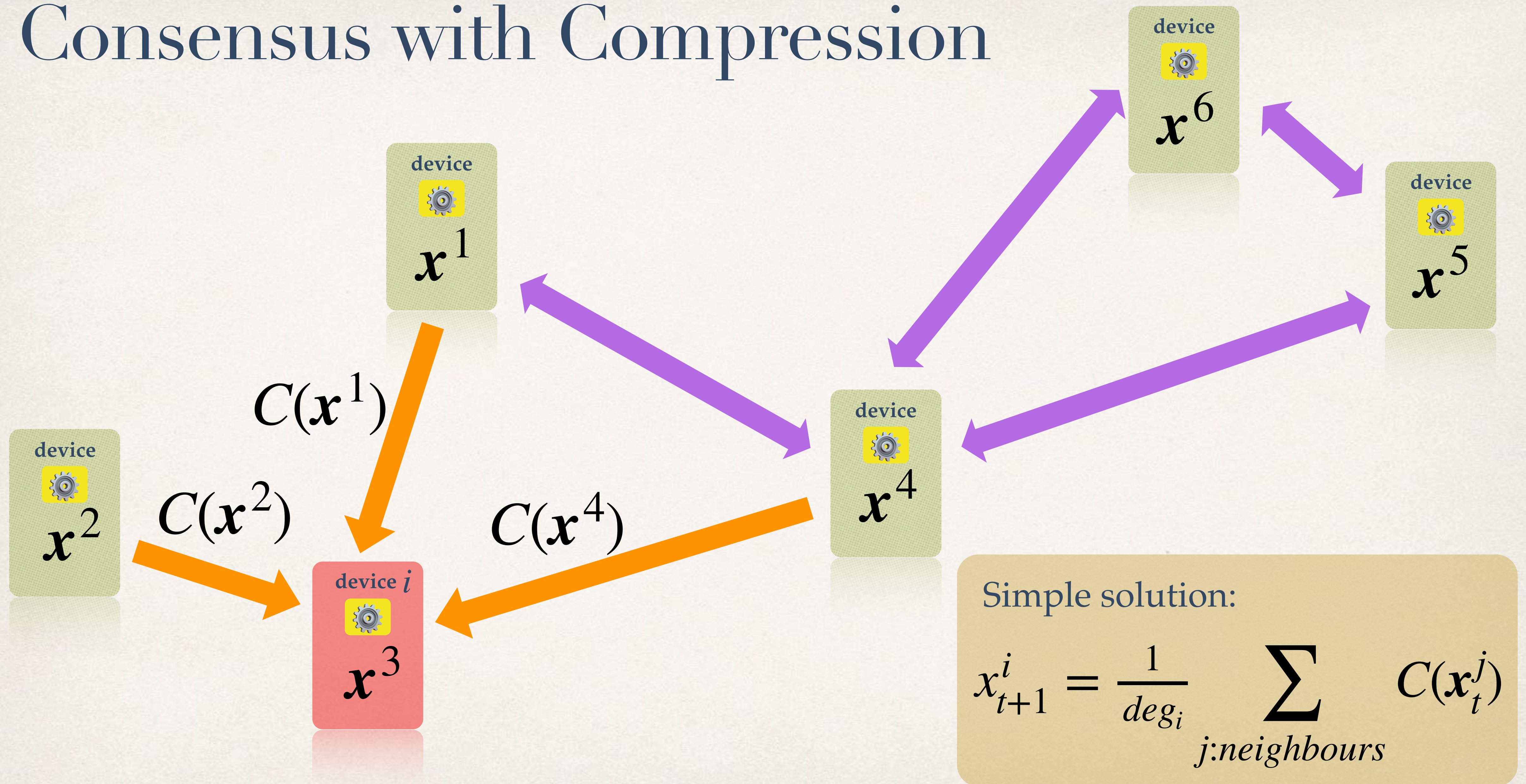
- ❖ limited-bit precision vector

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

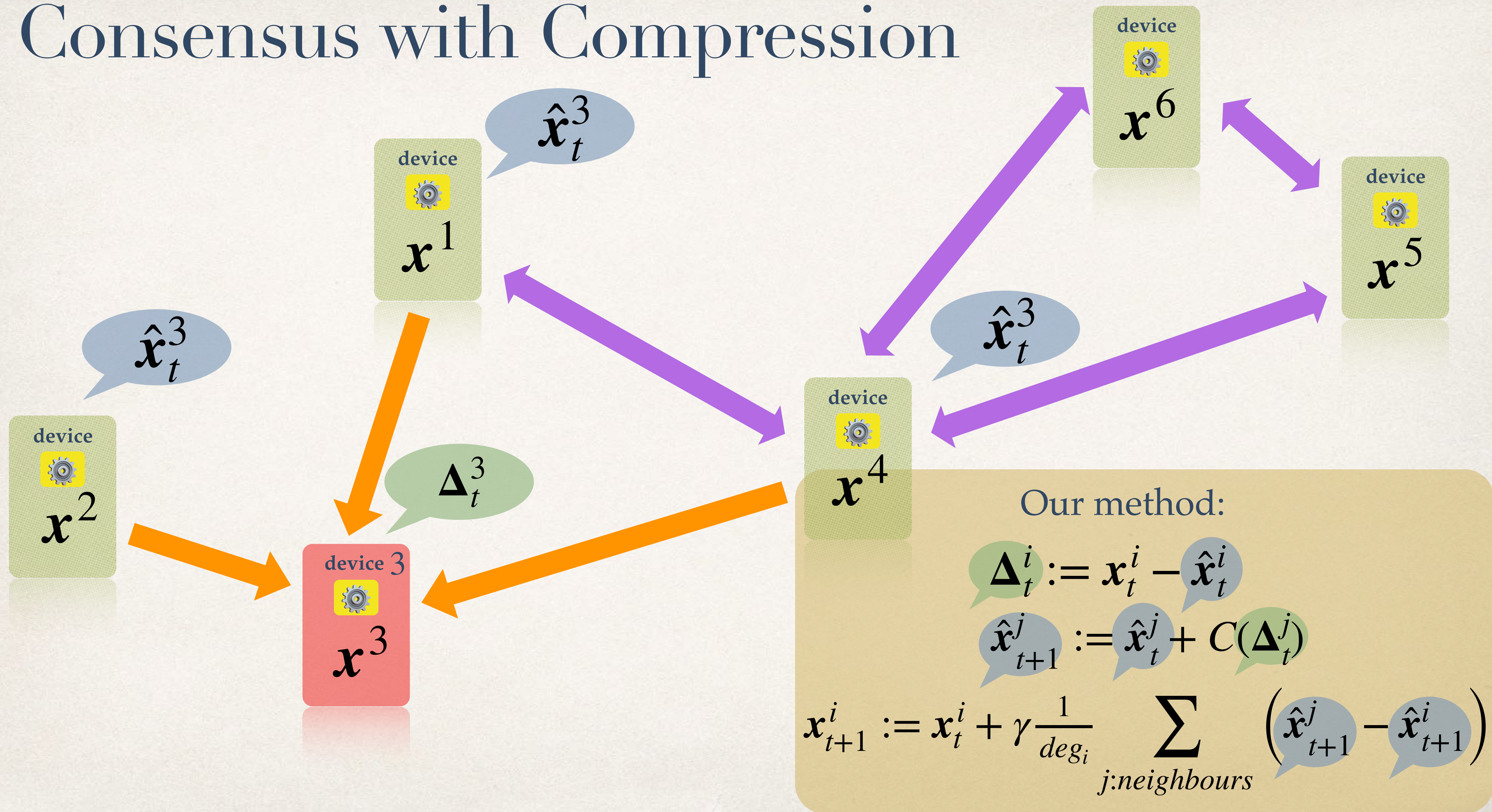
- ❖ random / top $k\%$ of all the entries

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

Consensus with Compression



Consensus with Compression

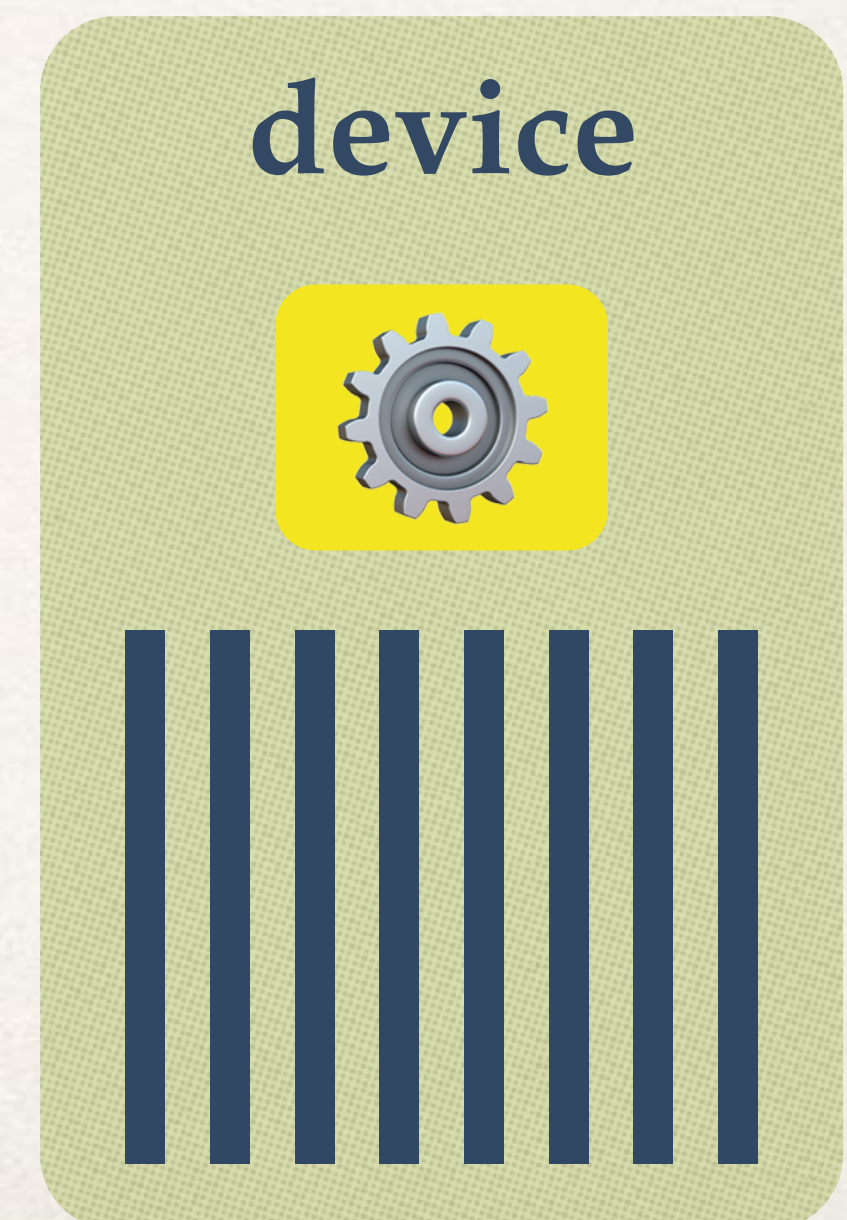


Stochastic Gradient Descent (SGD)

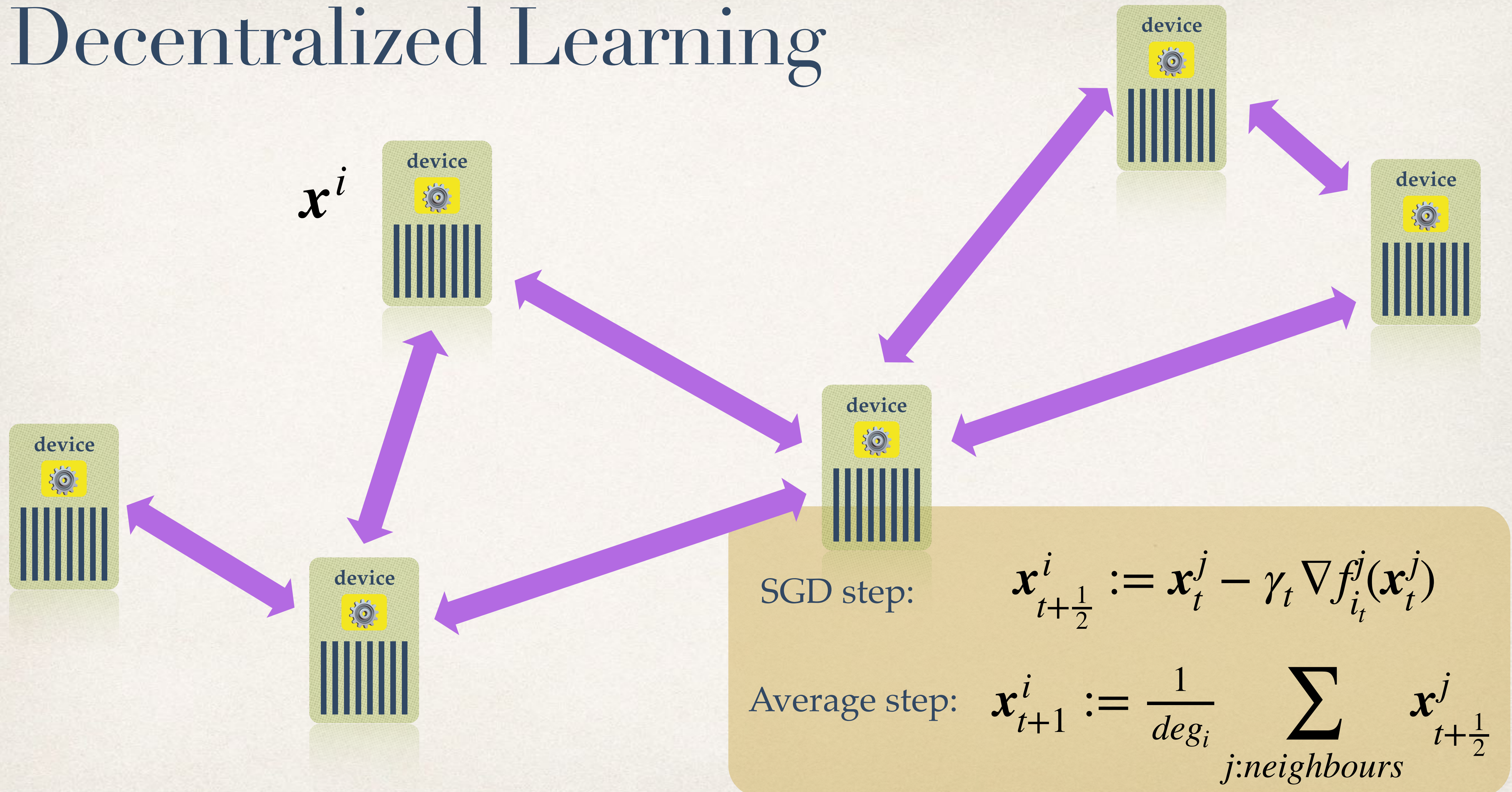
$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{|data|} \sum_{i \in data} f_i(\mathbf{x})$$

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

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t \nabla f_{i_t}(\mathbf{x}_t)$$



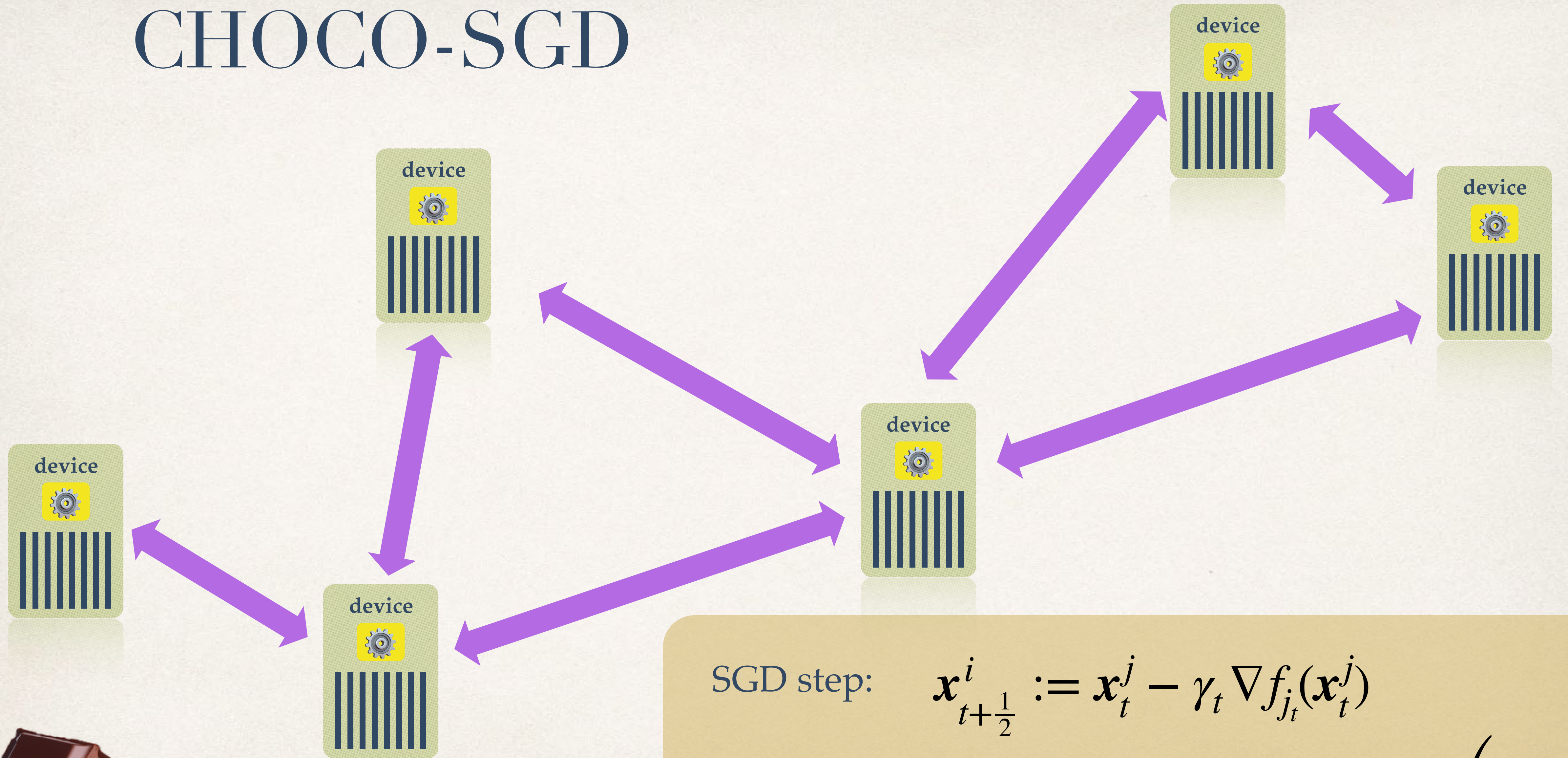
Decentralized Learning



CHOCO-SGD



CHOCO-SGD



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

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



Theoretical Result

$$f(\mathbf{x}_{avg}^{(T)}) - f^{\star} = \mathcal{O}\left(\frac{1}{nT}\right) + \mathcal{O}\left(\frac{1}{\delta^2 T^2}\right) + \mathcal{O}\left(\frac{1}{\delta^3 T^3}\right)$$

δ — compression ratio

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

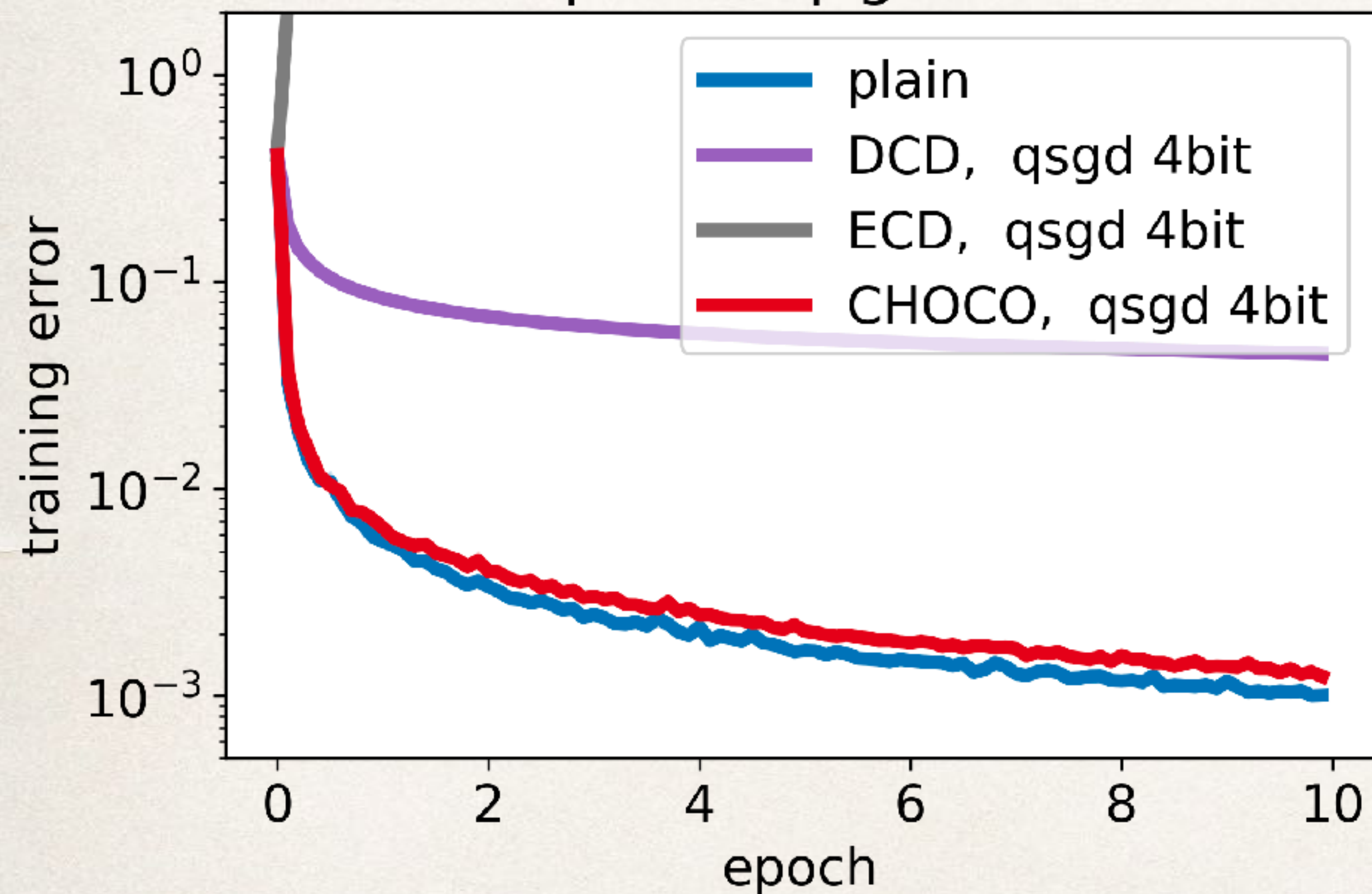
- ❖ (almost) same convergence rate as full communication



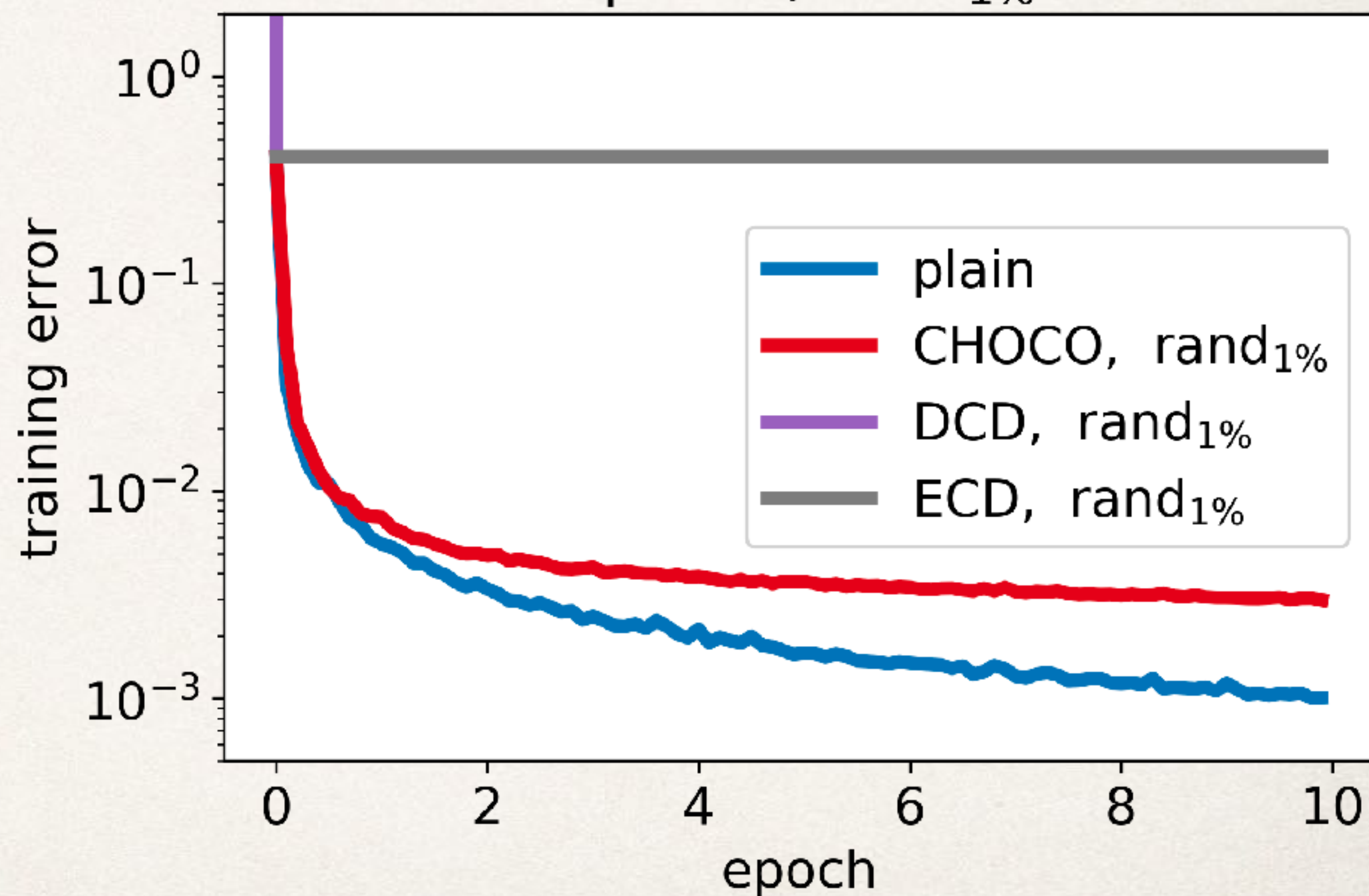
Experimental Results



Epsilon, qsgd 4bit



Epsilon, rand_{1%}



CHOCO SGD

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

Future work

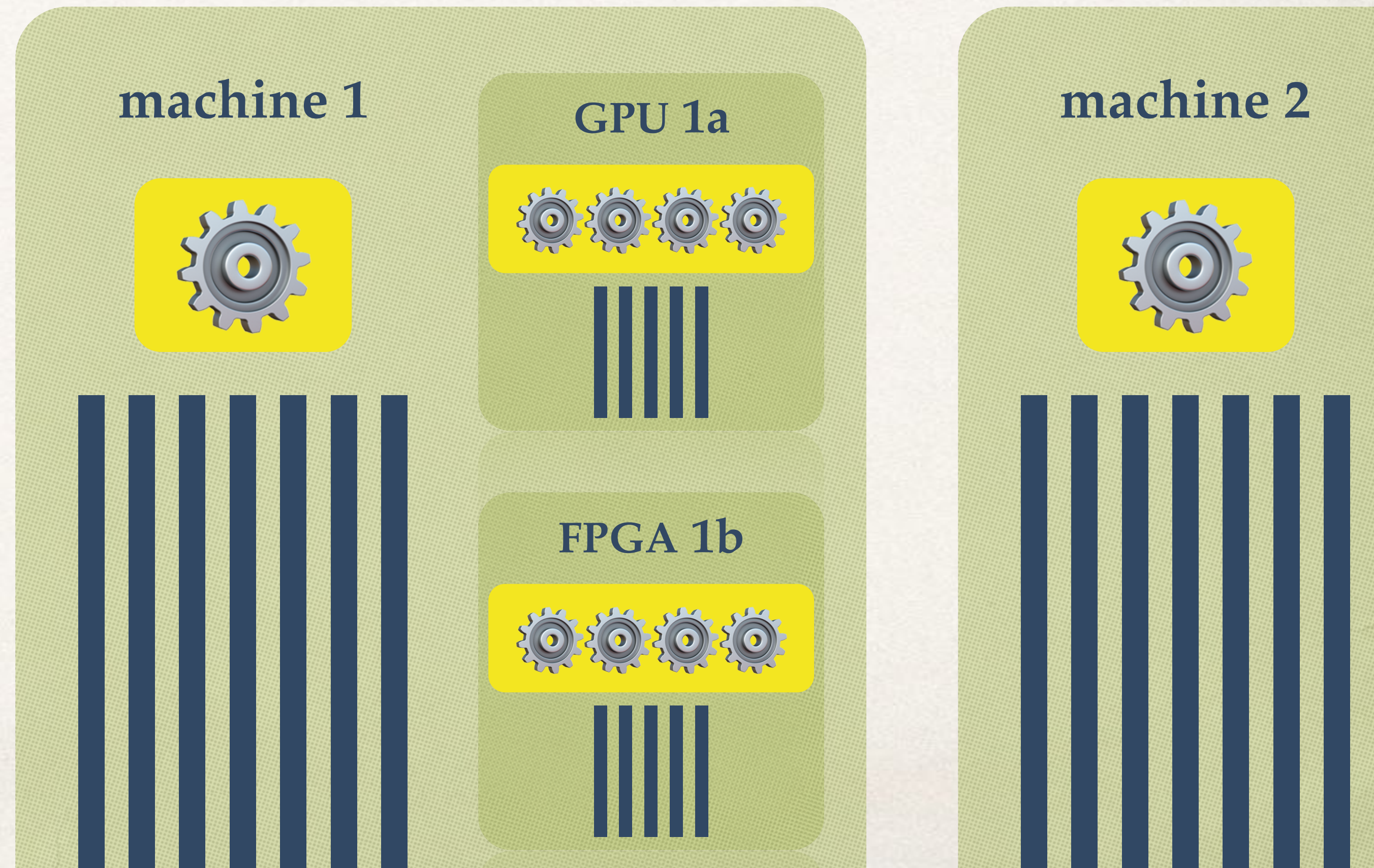
- ❖ Run it on deep learning experiments (multi GPU)
- ❖ Non-convex theory



3

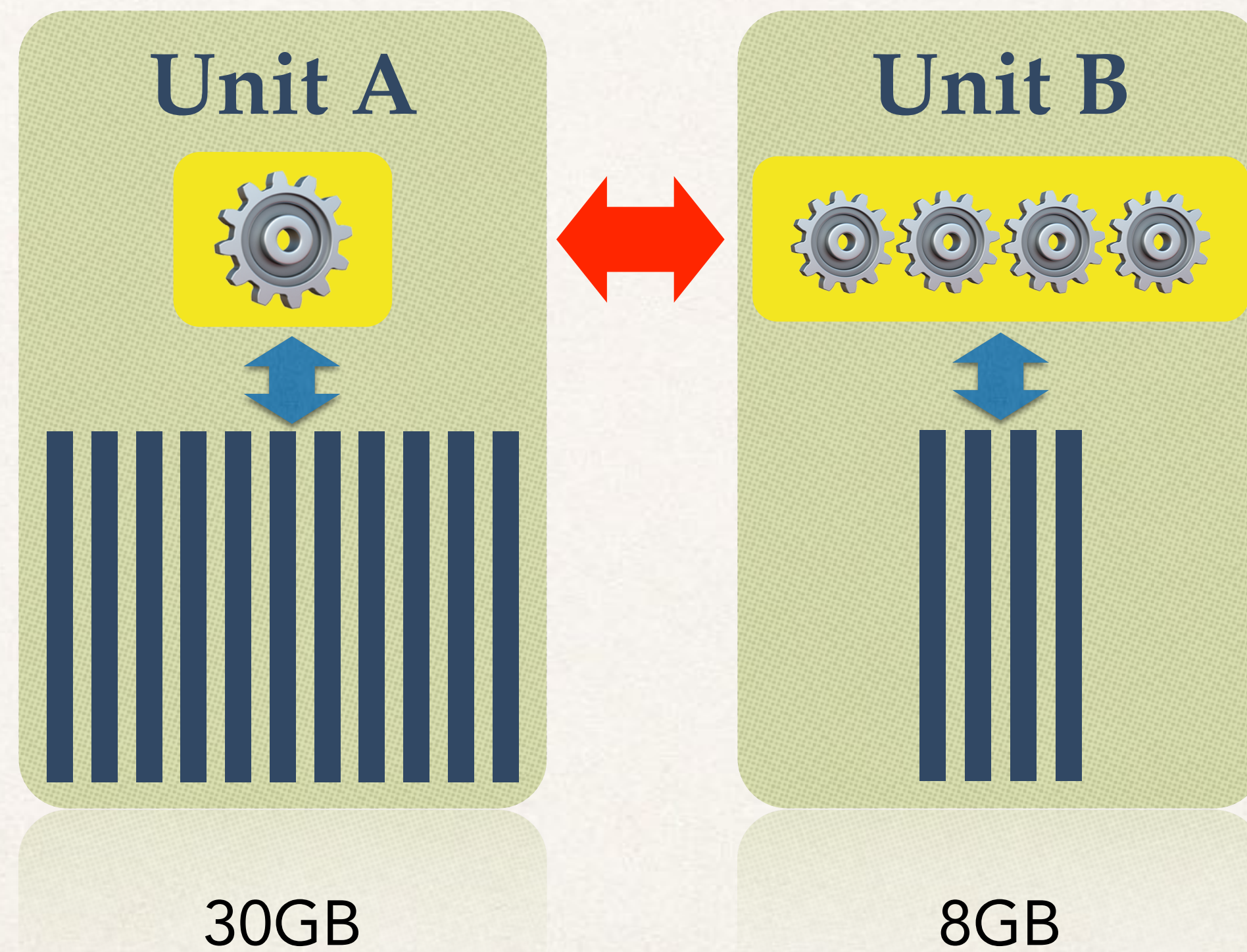
Leveraging Heterogenous Systems

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



Leveraging Heterogenous Systems

duality gap as selection criterion

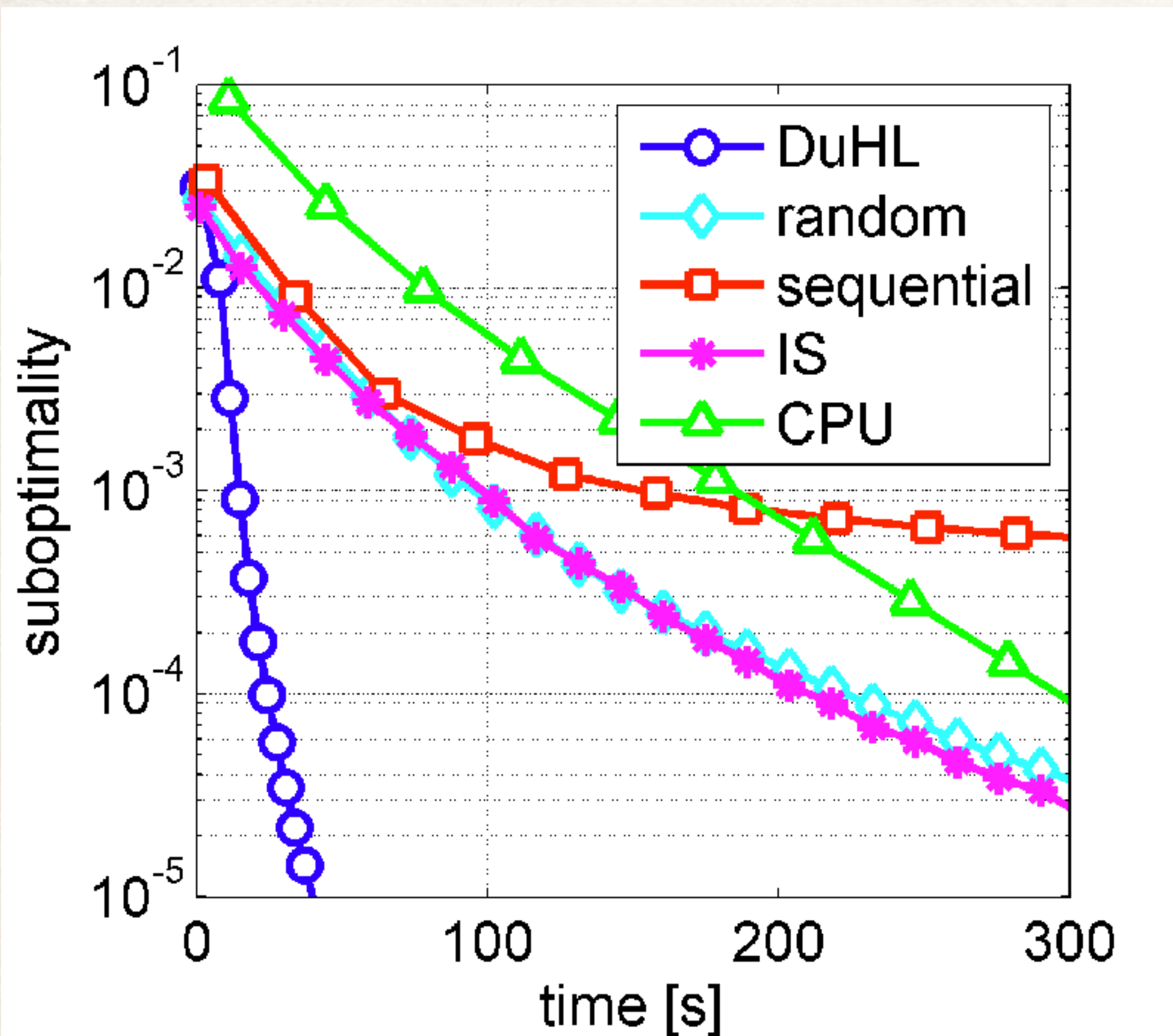


adaptive importance sampling

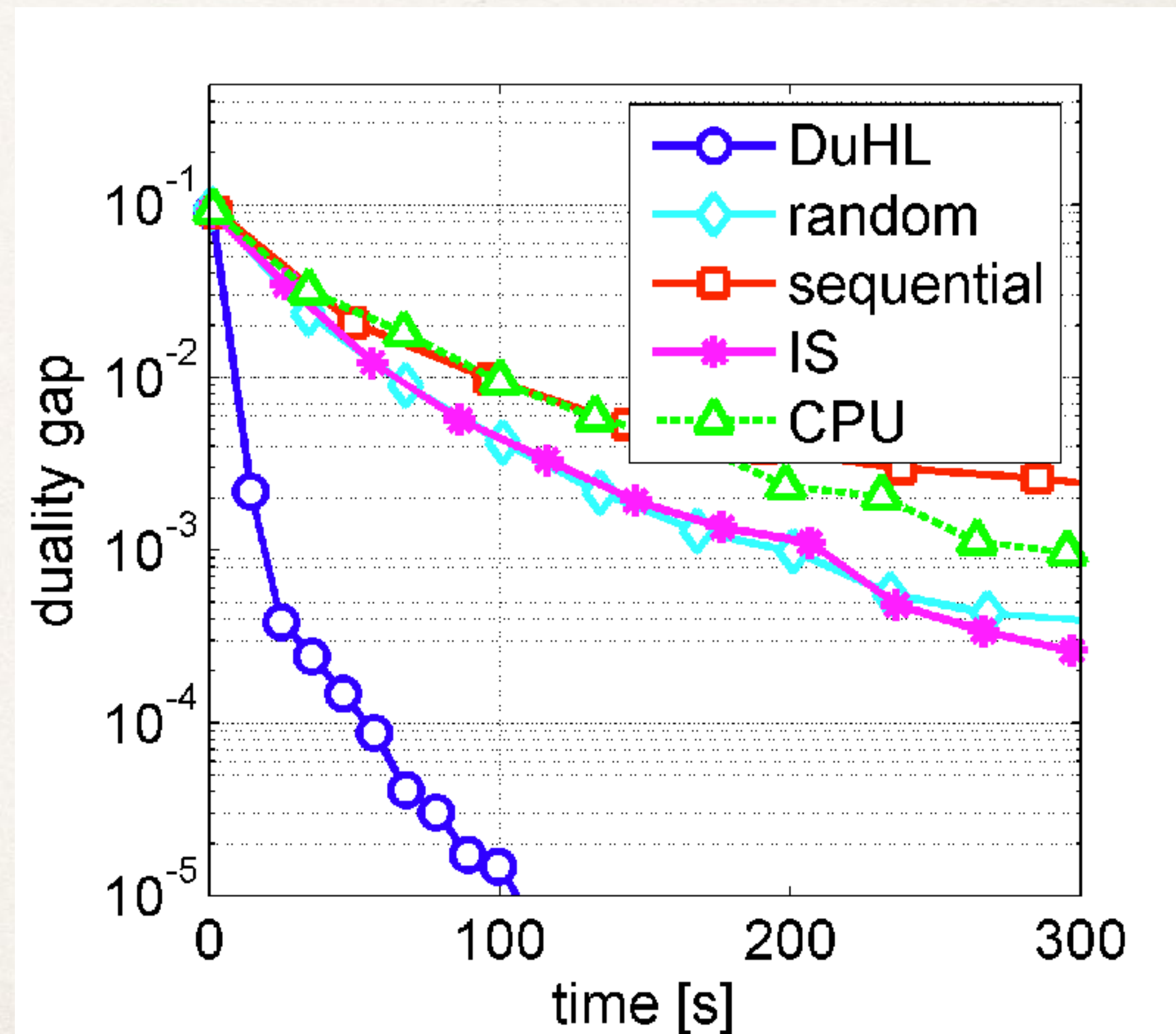
AISTATS 2017, 2018
NIPS 2017a,b

Experiments

RAM ↔ GPU, 30GB dataset

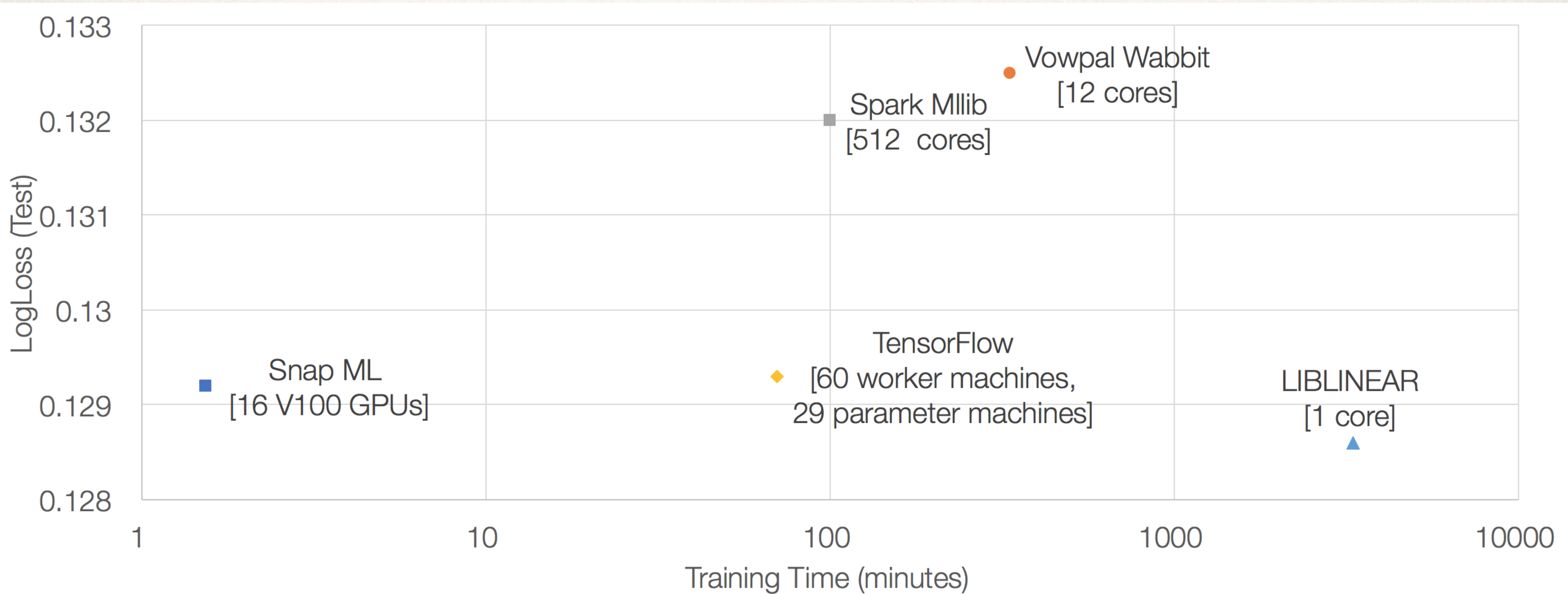


Lasso



SVM

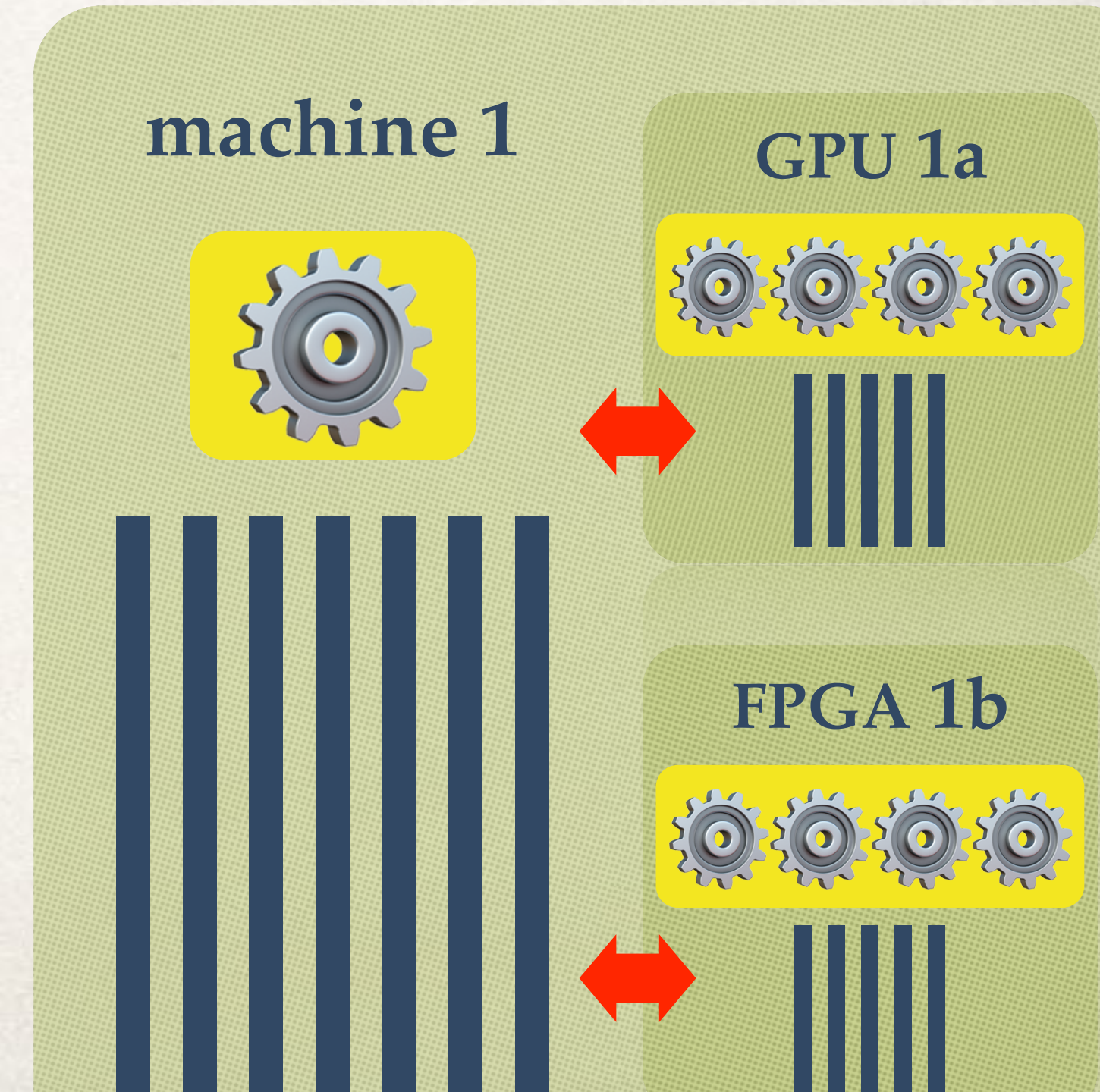
Experiments



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

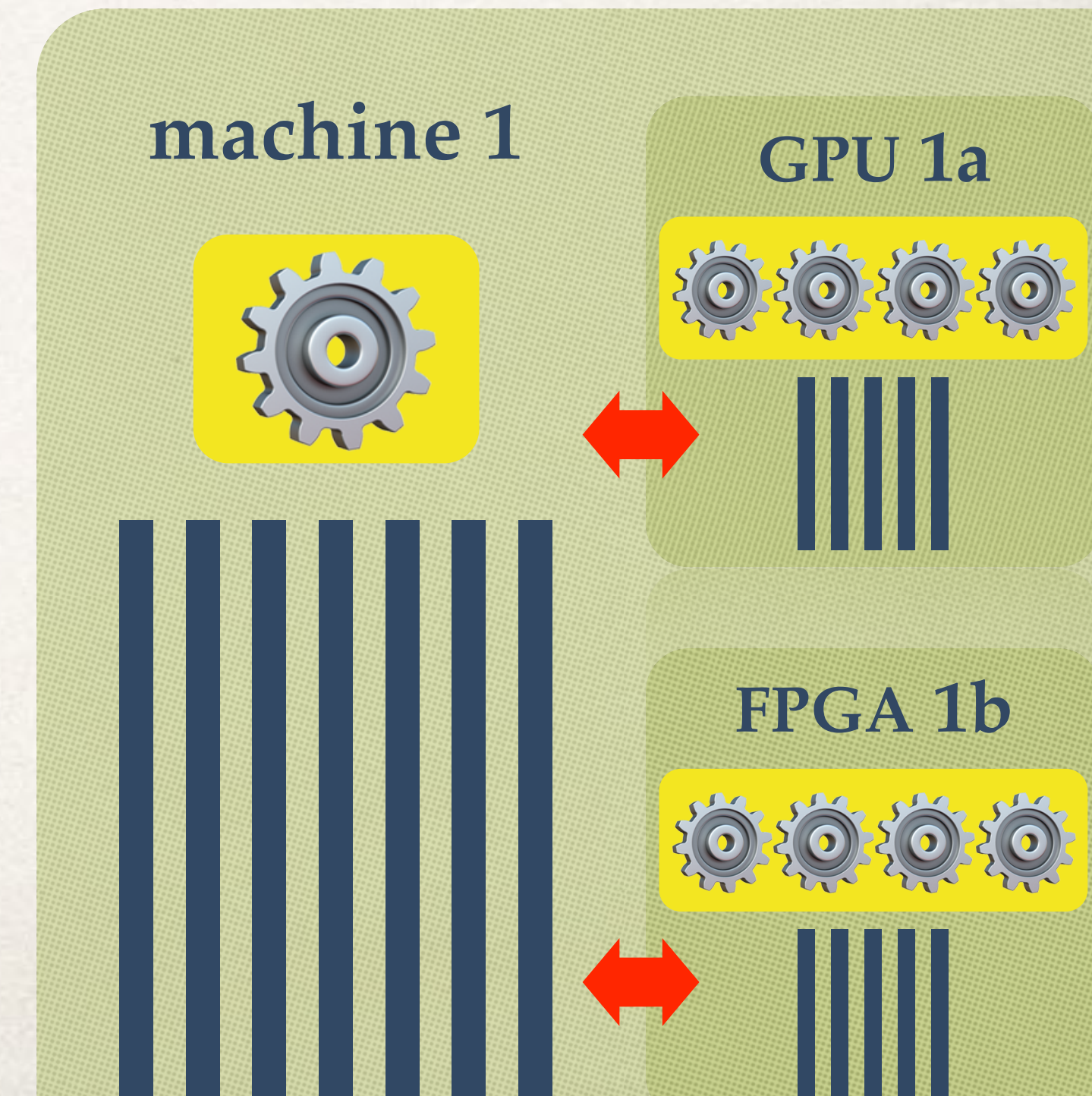
Open Research

- ❖ **limited precision operations** for efficiency of communication and computation
- ❖ **asynchronous and fault tolerant** algorithms
- ❖ heterogenous systems
- ❖ more **re-usable** algorithmic building blocks
 - for more systems and problems



Trends - Systems

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



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...

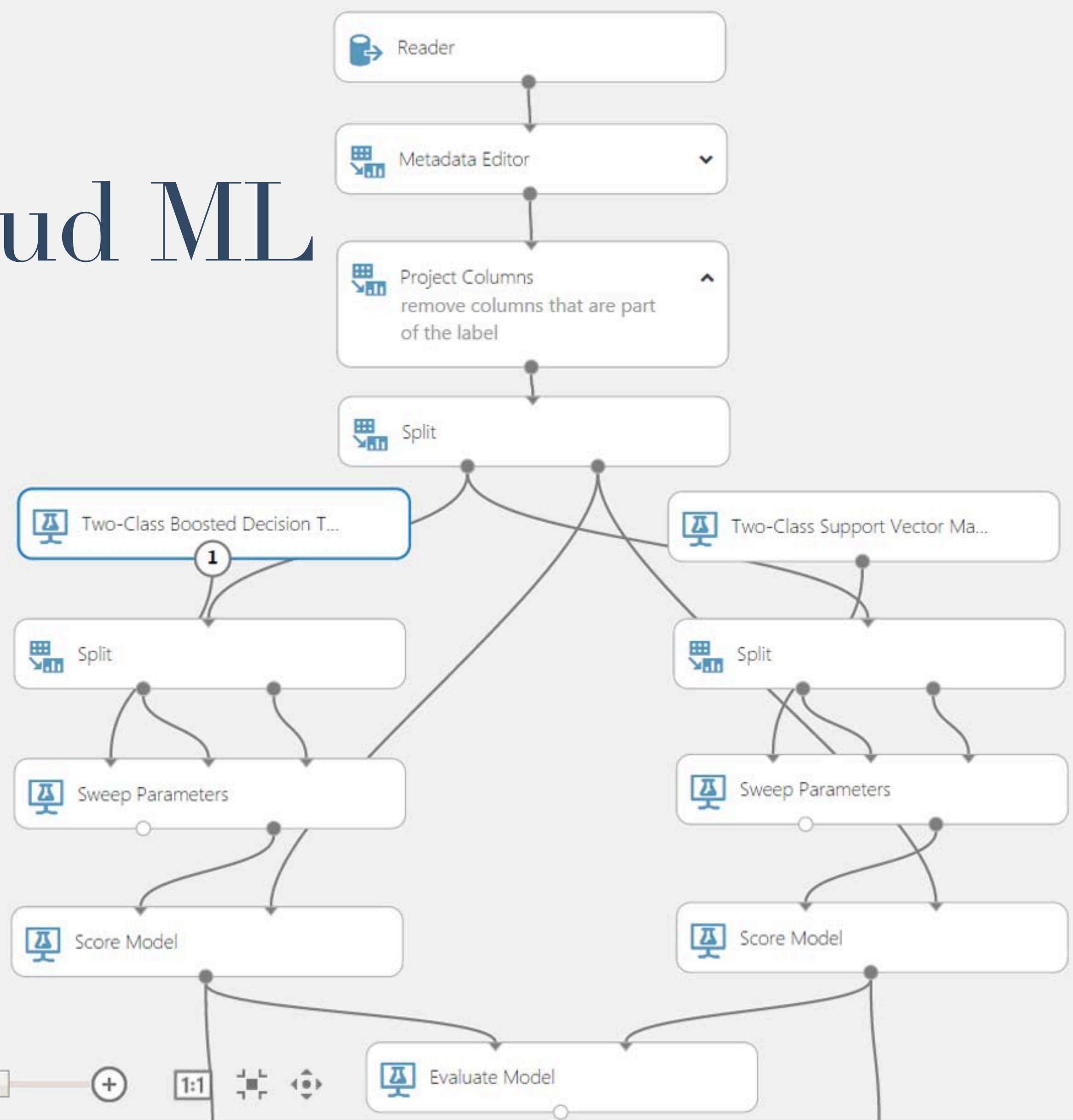
Quick Help

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

Cloud ML

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



Zoom and pan controls: minus, slider, plus, 1:1, pan, zoom

Project:

Distributed Machine Learning Benchmark

Goal:

Public and Reproducible
Comparison of Distributed Solvers

github.com/mlbench/mlbench

PyTorch



Google



Apache



HPC



Auto ML

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

Thanks!

mlo.epfl.ch