

# Distributed Optimization for Machine Learning

Martin Jaggi

*EPFL Machine Learning and Optimization Laboratory*

[mlo.epfl.ch](http://mlo.epfl.ch)



# Machine Learning Methods to Analyze Large-Scale Data



Machine  
Learning

Optimization

Systems



Applications





# Machine Learning?

software that can

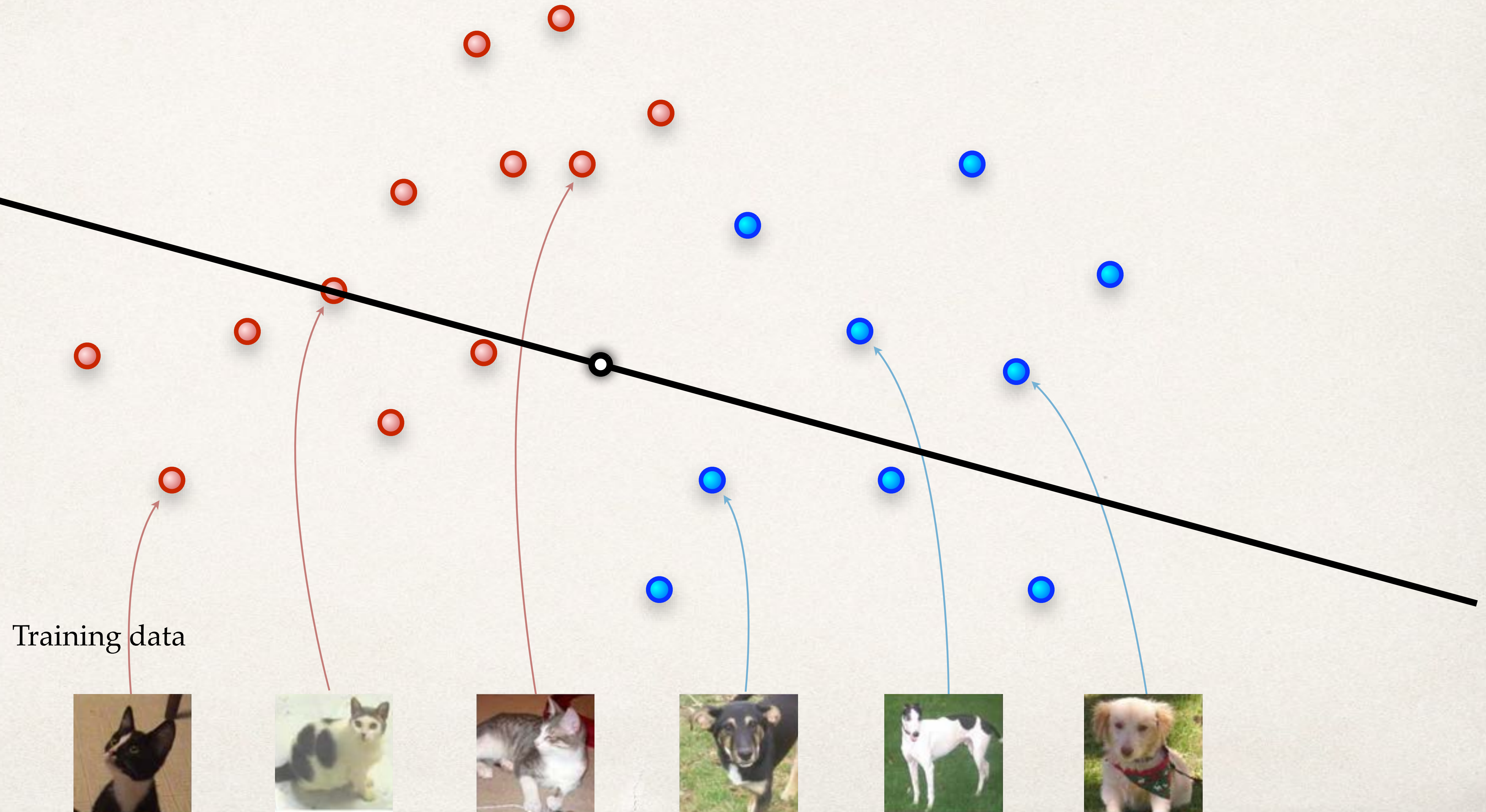
**learn from data**







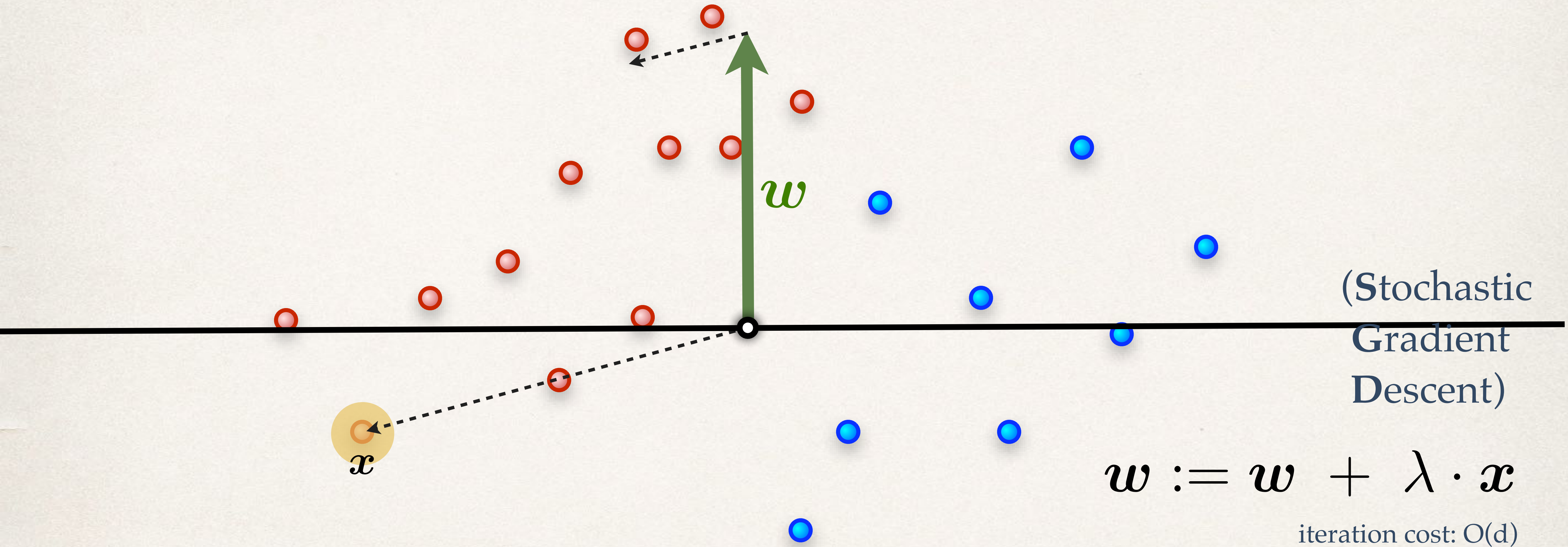
# Machine Learning Example





# The Learning Algorithm

$$\mathbf{x}_i \in \mathbb{R}^d$$



Perceptron

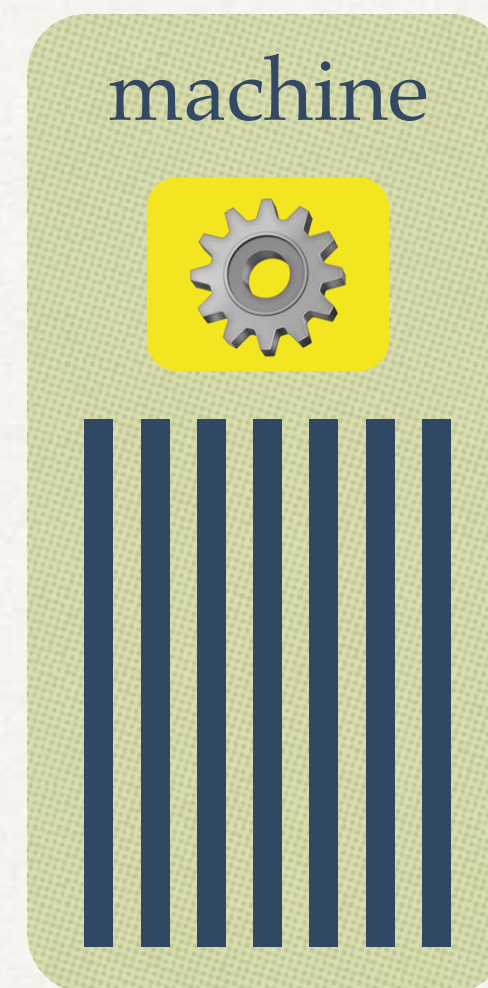
(Rosenblatt 1957)

Support-Vector-Machine

(Cortes & Vapnik 1995)



# Machine Learning Systems





# Machine Learning Systems

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

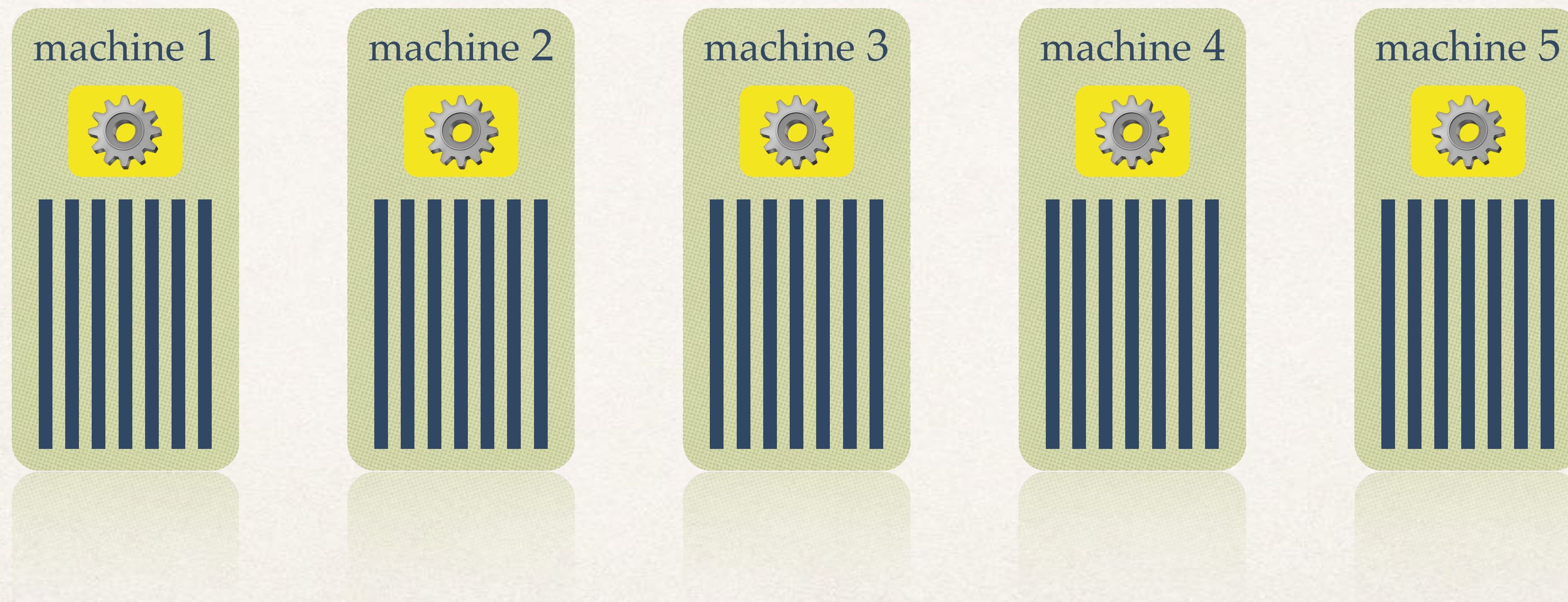


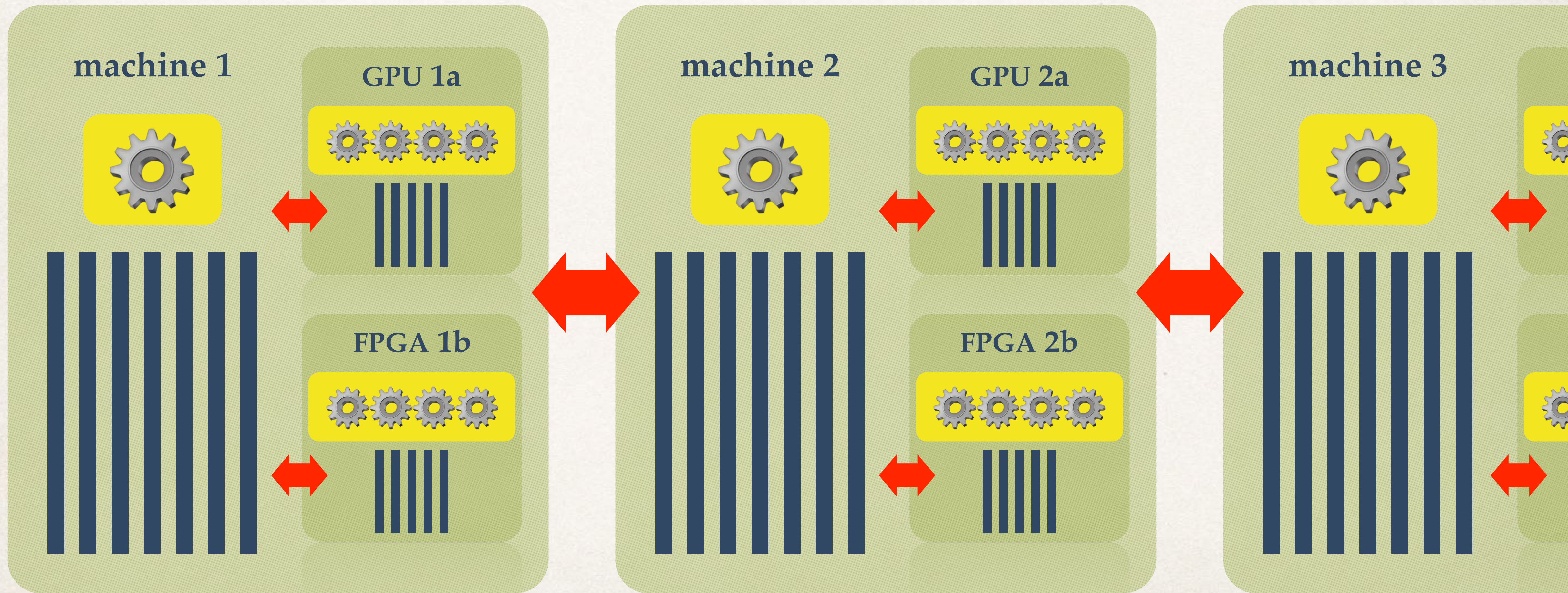




Foto: Florian Hirzinger



# Machine Learning Systems





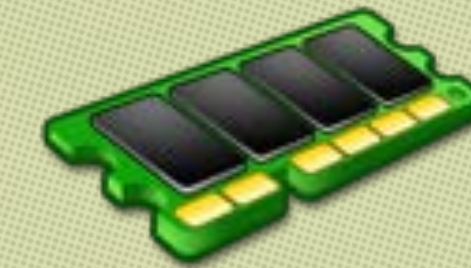
↔ Challenge

# The Cost of Communication

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

- ✦ Reading  $v$  from memory (RAM)

$100\text{ ns}$



- ✦ Sending  $v$  to another machine

$500'000\text{ ns}$

- ✦ Typical Map-Reduce iteration

$10'000'000'000\text{ ns}$





Challenge

# Usability

Parallel Programming is Hard

✦ no **reusability** of good  
single machine algorithms & code

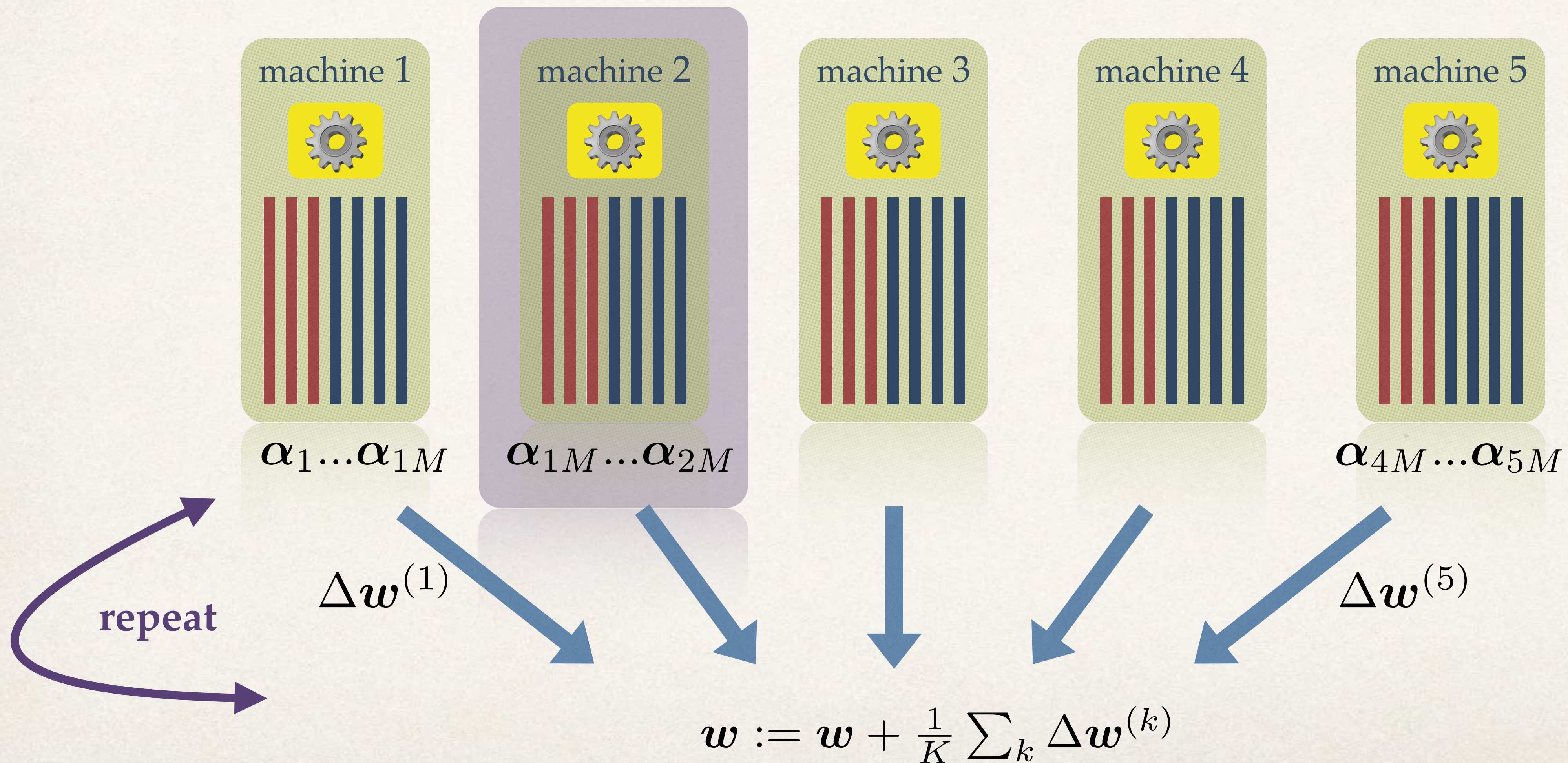


# Problem class

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



# CoCoA - Communication Efficient Distributed Optimization





# Optimization: Primal-Dual Context

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \left[ \mathcal{O}_A(\boldsymbol{\alpha}) := f(\overbrace{A\boldsymbol{\alpha}}^{A_{\text{loc}}\Delta\boldsymbol{\alpha}_{[k]} + \boldsymbol{w}}) + g(\boldsymbol{\alpha}) \right]$$

*primal Lasso*  
*dual L2-reg SVM/Log-Regr*  
*primal L1-reg SVM/Log-Reg*

correspondence

$$\boldsymbol{w} := \nabla f(A\boldsymbol{\alpha})$$

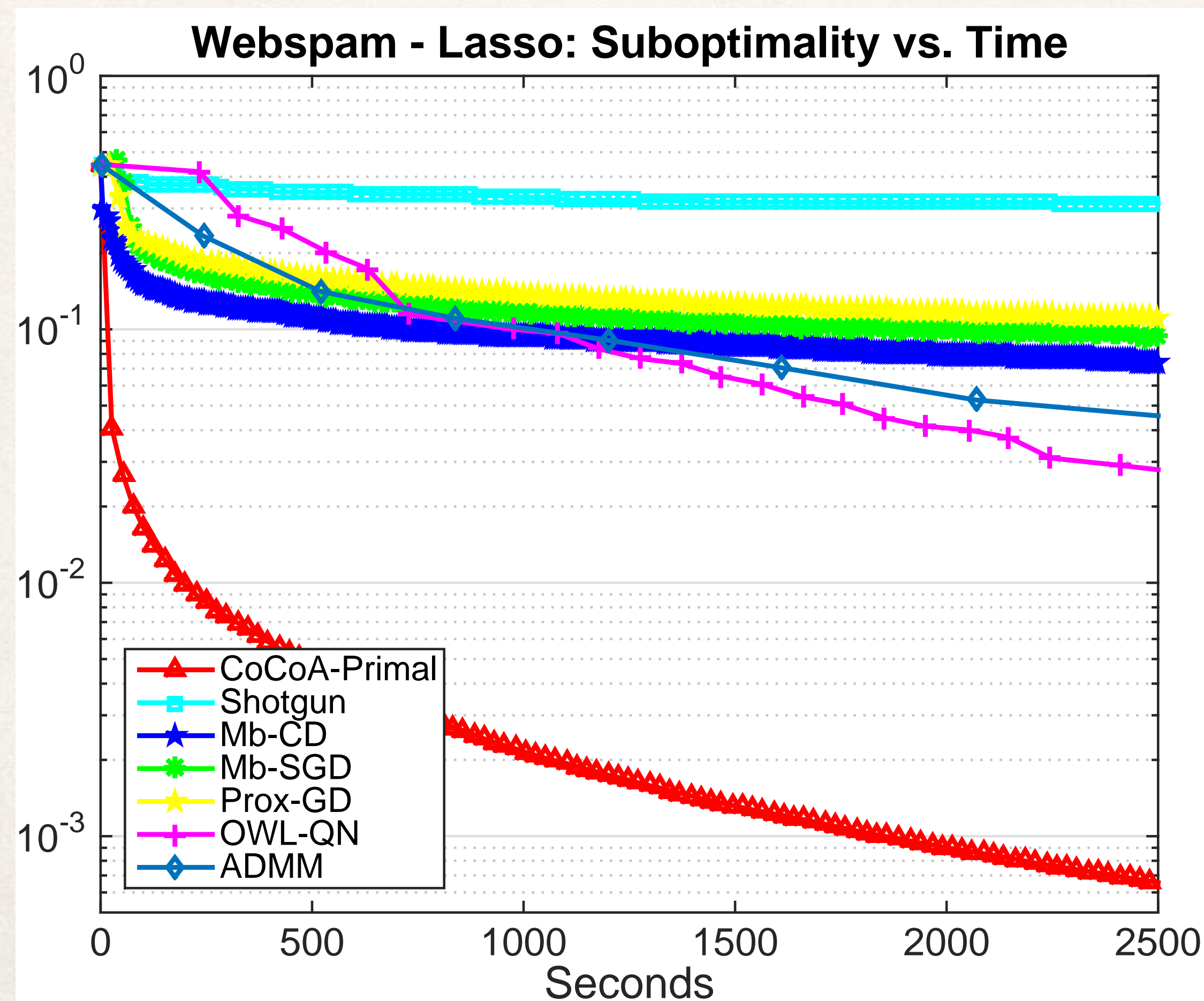
$$\min_{\boldsymbol{w} \in \mathbb{R}^d} \left[ \mathcal{O}_B(\boldsymbol{w}) := g^*(-A^\top \boldsymbol{w}) + f^*(\boldsymbol{w}) \right]$$



# Distributed Experiments

## Sparse 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,*  
[arxiv.org/abs/1611.02189](https://arxiv.org/abs/1611.02189)

*Spark Code:*  
[github.com/gingsmith/proxcocoa](https://github.com/gingsmith/proxcocoa)

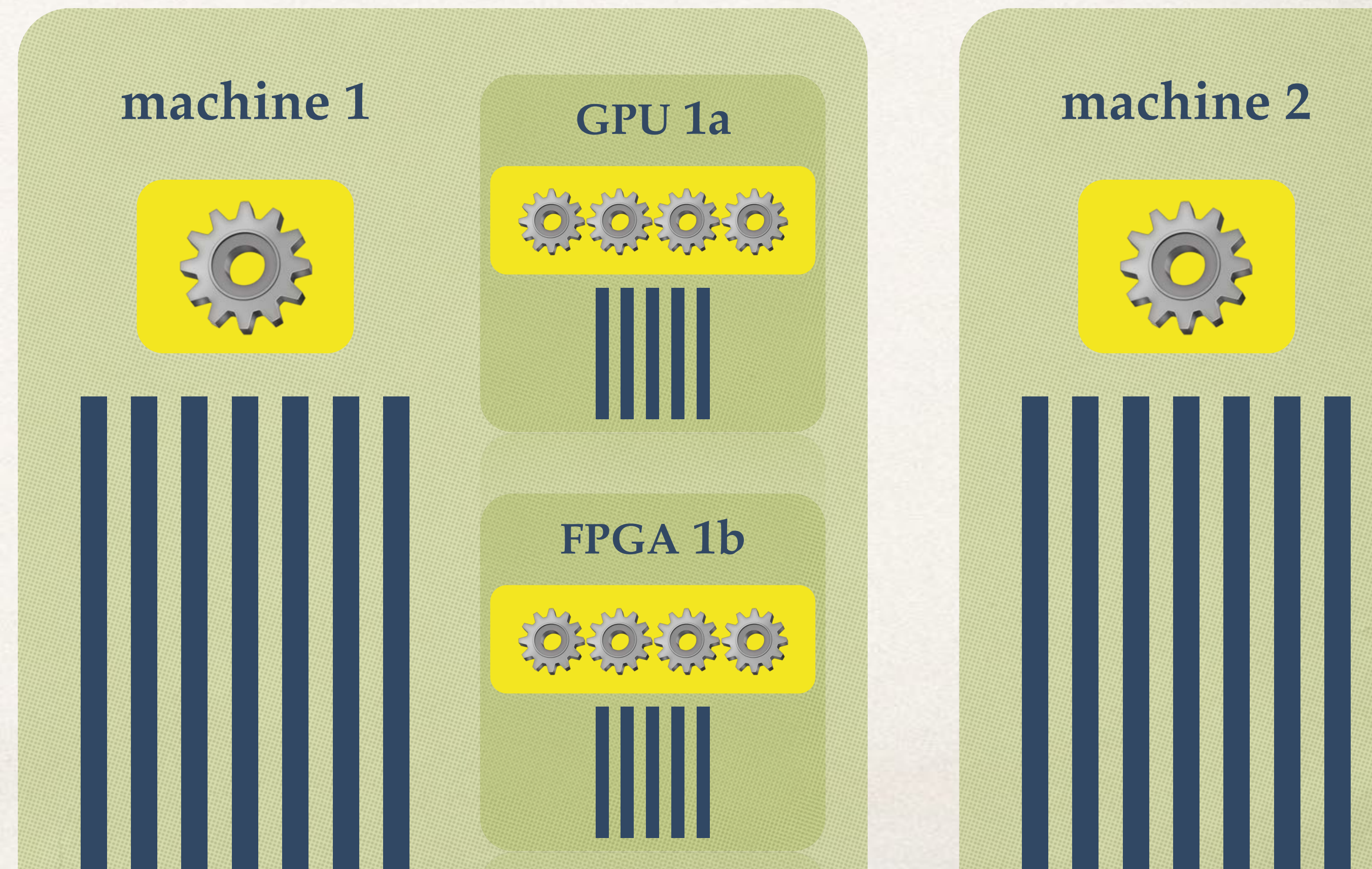
+ TensorFlow  
+ Apache Flink



Challenge

# Leveraging Memory Hierarchy

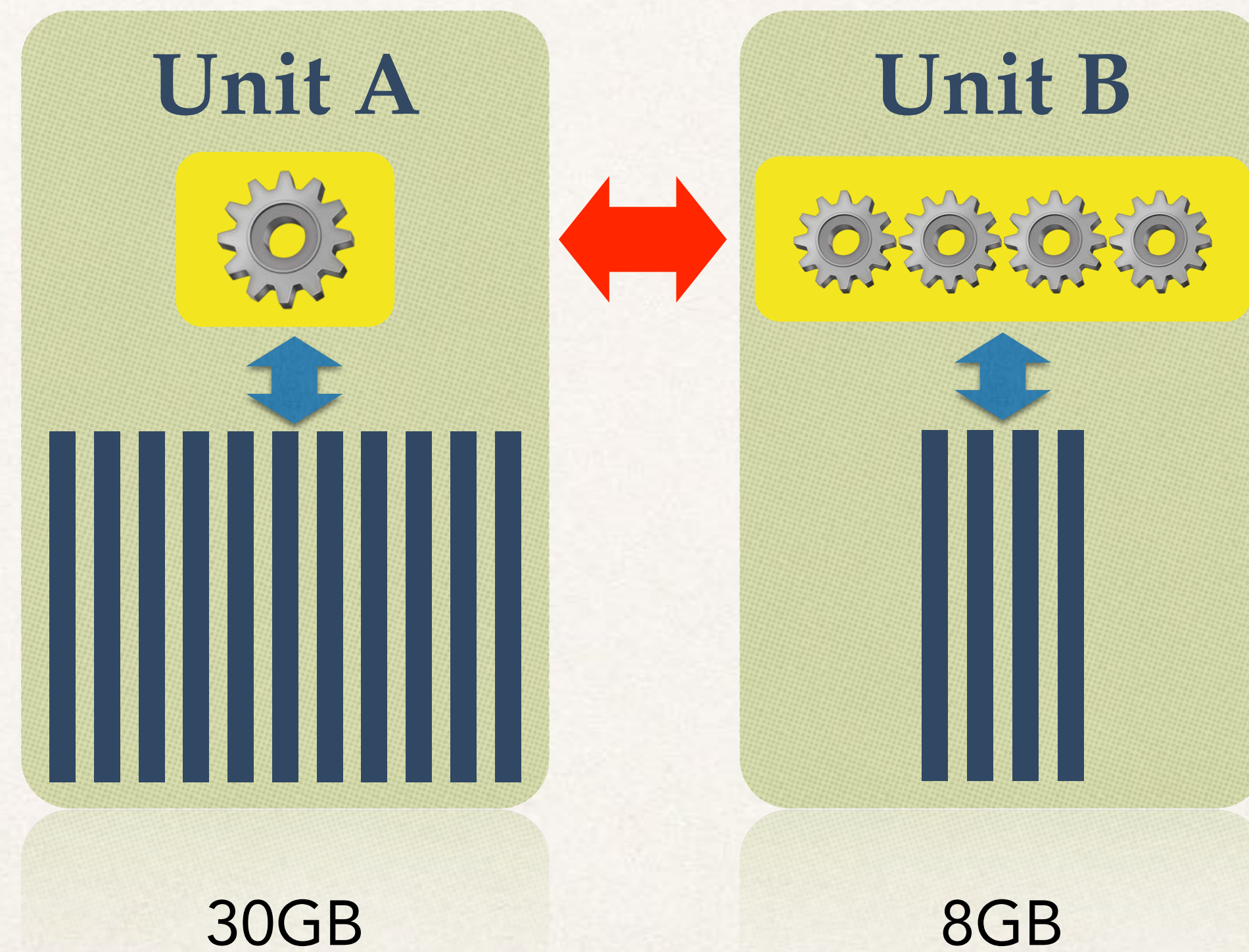
Which data to put in which memory?





# Leveraging Memory Hierarchy

duality gap as selection criterion



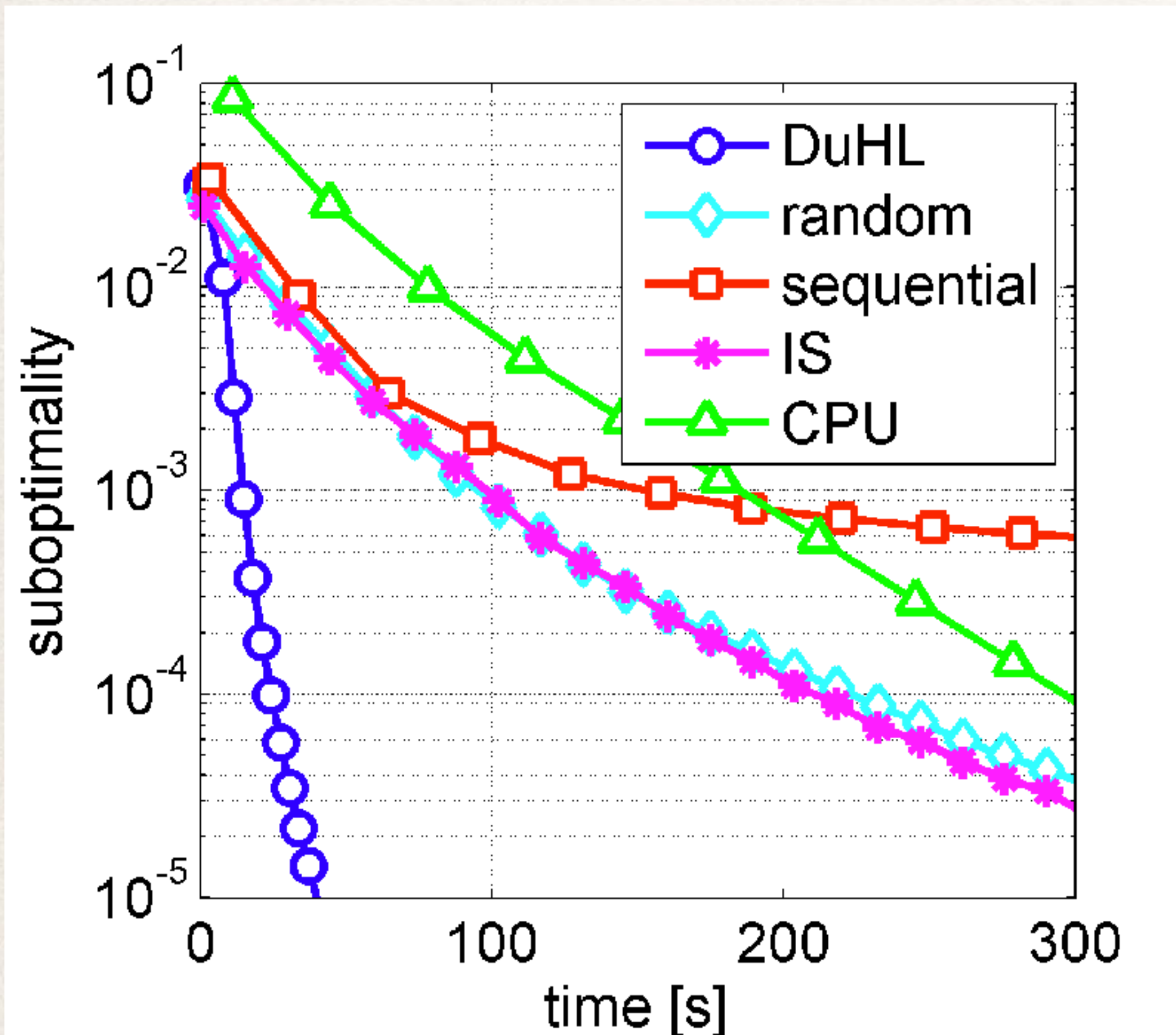
adaptive importance sampling

*AISTATS 2017*

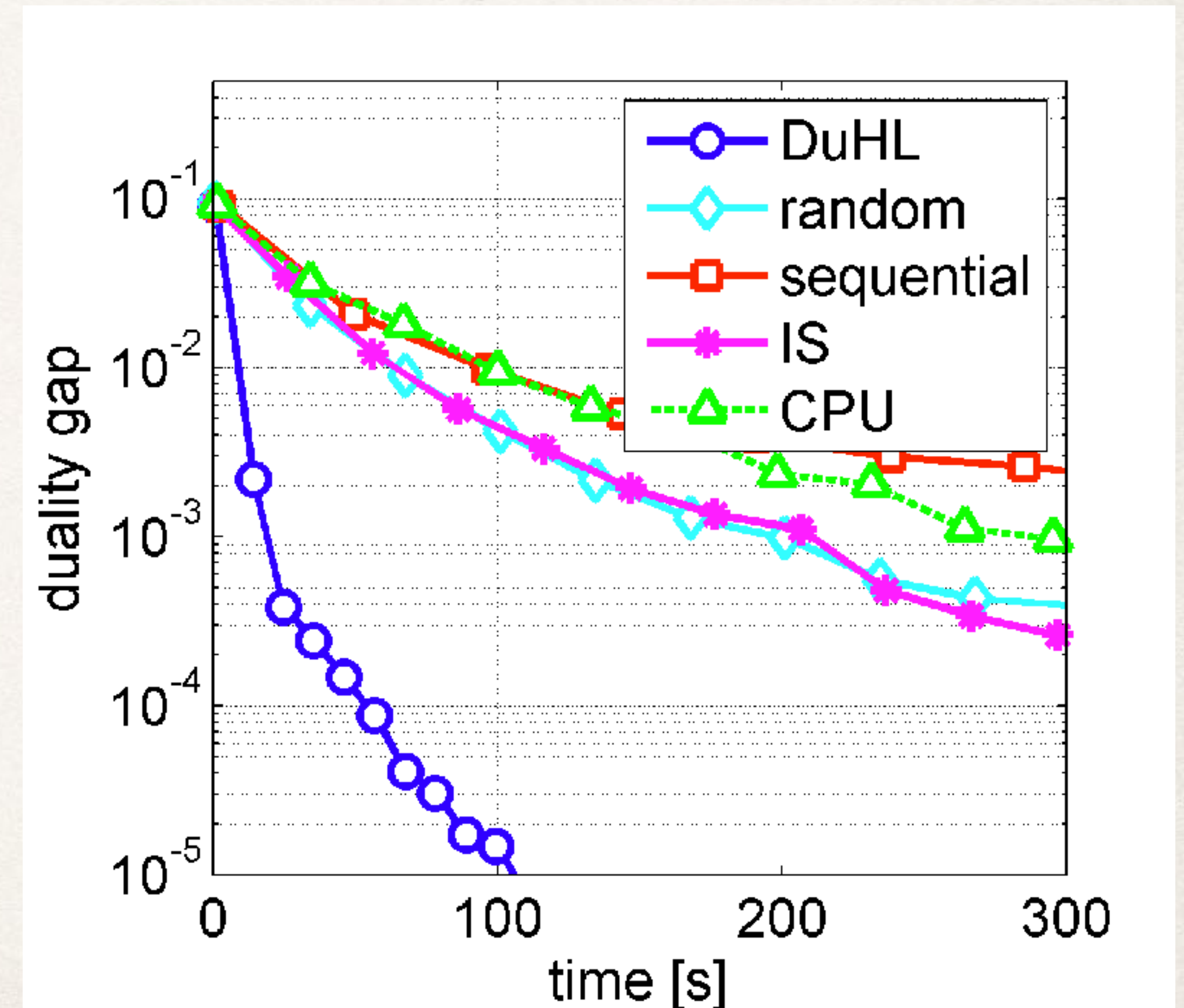


# Experiments

RAM  $\longleftrightarrow$  GPU, 30GB dataset



Lasso



SVM



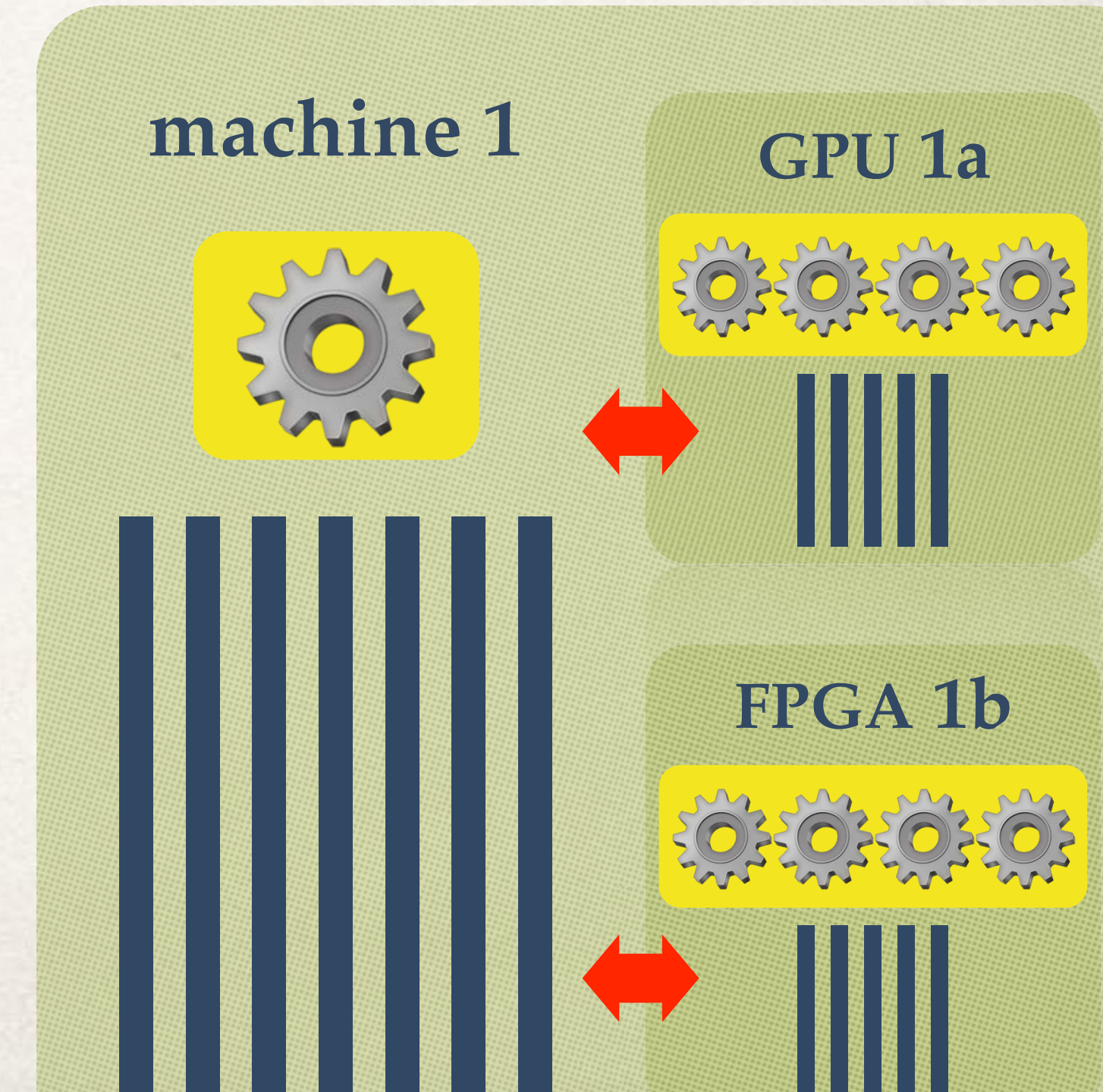
# Conclusion

- ❖ try to improve **usability** of large-scale ML
- ❖ full **adaptivity** to the communication cost, memory hierarchy and bandwidth
- ❖ **re-usability** of good single machine solvers
- ❖ **accuracy** certificates



# Open Research

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





*Project:*

# Distributed Machine Learning Benchmark

*Goal:*

Public and Reproducible  
Comparison of Distributed Solvers

[github.com/mlbench/mlbench](https://github.com/mlbench/mlbench)

Apache



Google



Apache



HPC





Thanks!

[mlo.epfl.ch](http://mlo.epfl.ch)

Celestine Dünnér, Virginia Smith, Simone Forte, Chenxin Ma, Martin Takac,  
Dmytro Perekrestenko, Volkan Cevher, Michael I. Jordan, Thomas Hofmann