

# Adding vs. Averaging in Distributed Primal-Dual Optimization

Chenxin Ma<sup>\*a</sup> · Virginia Smith<sup>\*b</sup> · Martin Jaggi<sup>c</sup> · Michael I. Jordan<sup>b</sup> · Peter Richtárik<sup>d</sup> · Martin Takáč<sup>a</sup>



## overview

**Goal:** Fixing the **communication bottleneck** for **distributed optimization** in supervised ML

### Contributions

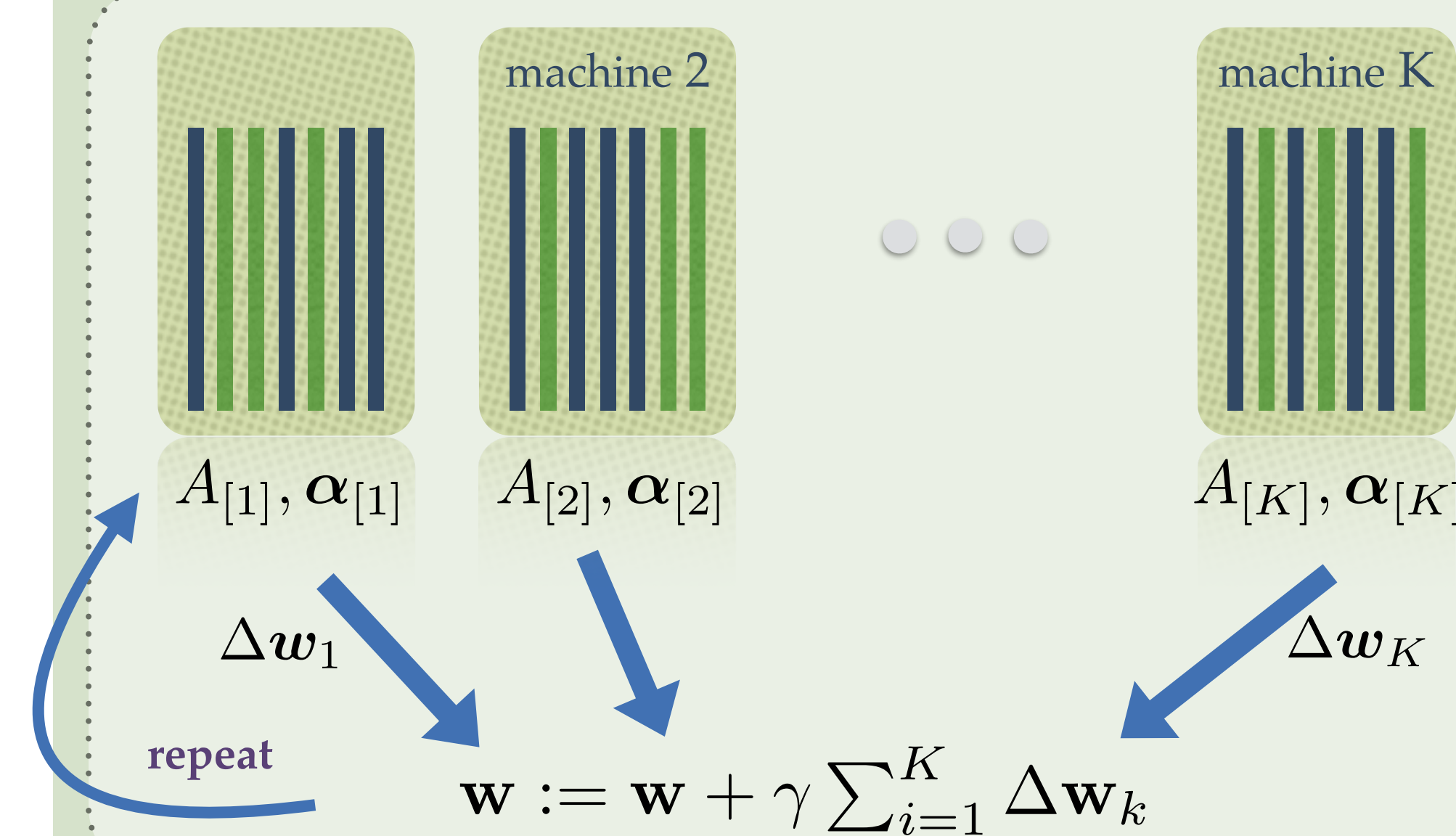
COCO<sup>+</sup>A is a *primal-dual* framework for distributed optimization

- ✓ efficient additive aggregation of local updates
- ✓ strong convergence guarantees
- ✓ framework: guarantees for *arbitrary* local solvers
- ✓ significant practical speedup

## COCO<sup>+</sup>A

### Main Idea

Propose a *local subproblem* to allow additive, data dependent aggregation



►  $\gamma = \frac{1}{K} \Rightarrow$  averaging,  $\gamma = 1 \Rightarrow$  adding

### Local Subproblem

$$\max_{\Delta \alpha_{[k]} \in \mathbb{R}^n} \mathcal{G}_k^{\sigma'}(\Delta \alpha_{[k]}; \mathbf{w}, \alpha_{[k]})$$

$$\mathcal{G}_k^{\sigma'}(\Delta \alpha_{[k]}; \mathbf{w}, \alpha_{[k]}) := -\frac{1}{n} \sum_{i \in \mathcal{P}_k} \ell_i^*(-\alpha_i - (\Delta \alpha_{[k]})_i)$$

$$- \frac{1}{K} \frac{\lambda}{2} \|\mathbf{w}\|^2 - \frac{1}{n} \mathbf{w}^T A \Delta \alpha_{[k]} - \frac{\lambda}{2} \sigma' \left\| \frac{1}{\lambda n} A \Delta \alpha_{[k]} \right\|^2$$

- $\sigma'$  measure of difficulty of data partitioning
- $\sigma' := \gamma K \in [1, K]$  safe in practice

### CoCoA<sup>+</sup> Framework

**Input:** Aggregation parameter  $\gamma \in (0, 1]$ , subproblem parameter  $\sigma'$   
**Initialize:**  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^d$   
**for**  $t = 0, 1, 2, \dots$  **do**  
    **for**  $k \in \{1, 2, \dots, K\}$  **in parallel over workers do**  
        call *local solver*: compute a  $\Theta$ -approximate solution  $\Delta \alpha_{[k]}$  of the local subproblem  $\mathcal{G}_k^{\sigma'}(\Delta \alpha_{[k]}; \mathbf{w}, \alpha_{[k]})$   
        update  $\alpha_{[k]}^{(t+1)} := \alpha_{[k]}^{(t)} + \gamma \Delta \alpha_{[k]}$   
        return  $\Delta \mathbf{w}_k := \frac{1}{\lambda n} A \Delta \alpha_{[k]}$   
    **end for**  
    reduce  $\mathbf{w}^{(t+1)} := \mathbf{w}^{(t)} + \gamma \sum_{k=1}^K \Delta \mathbf{w}_k$   
**end for**

► flexible: can use *arbitrary* local solver

## convergence results

### Local $\Theta$ -Approximation

For  $\Theta \in [0, 1)$ , we assume the local solver finds a (possibly) randomized approximate solution satisfying:

$$\mathbb{E}[\mathcal{G}_k^{\sigma'}(\Delta \alpha_{[k]}^*) - \mathcal{G}_k^{\sigma'}(\Delta \alpha_{[k]})] \leq \Theta \left( \mathcal{G}_k^{\sigma'}(\Delta \alpha_{[k]}^*) - \mathcal{G}_k^{\sigma'}(\mathbf{0}) \right)$$

**Theorem.** Let  $\ell_i(\cdot)$  be  $L$ -**Lipschitz**. Obtain suboptimality  $\epsilon$ , after  $T$  iterations, with:

- CoCoA, averaging  $\gamma = 1/K$   

$$T \geq \tilde{\mathcal{O}} \left( \frac{K}{1-\Theta} \left( \frac{8L^2}{\lambda K \epsilon} + \tilde{c} \right) \right)$$
- CoCoA<sup>+</sup>, adding  $\gamma = 1$   

$$T \geq \tilde{\mathcal{O}} \left( \frac{1}{1-\Theta} \left( \frac{8L^2}{\lambda \epsilon} + \tilde{c} \right) \right)$$

**Theorem.** Let  $\ell_i(\cdot)$  be  $1/\mu$ -**smooth**. Obtain suboptimality  $\epsilon$ , after  $T$  iterations, with:

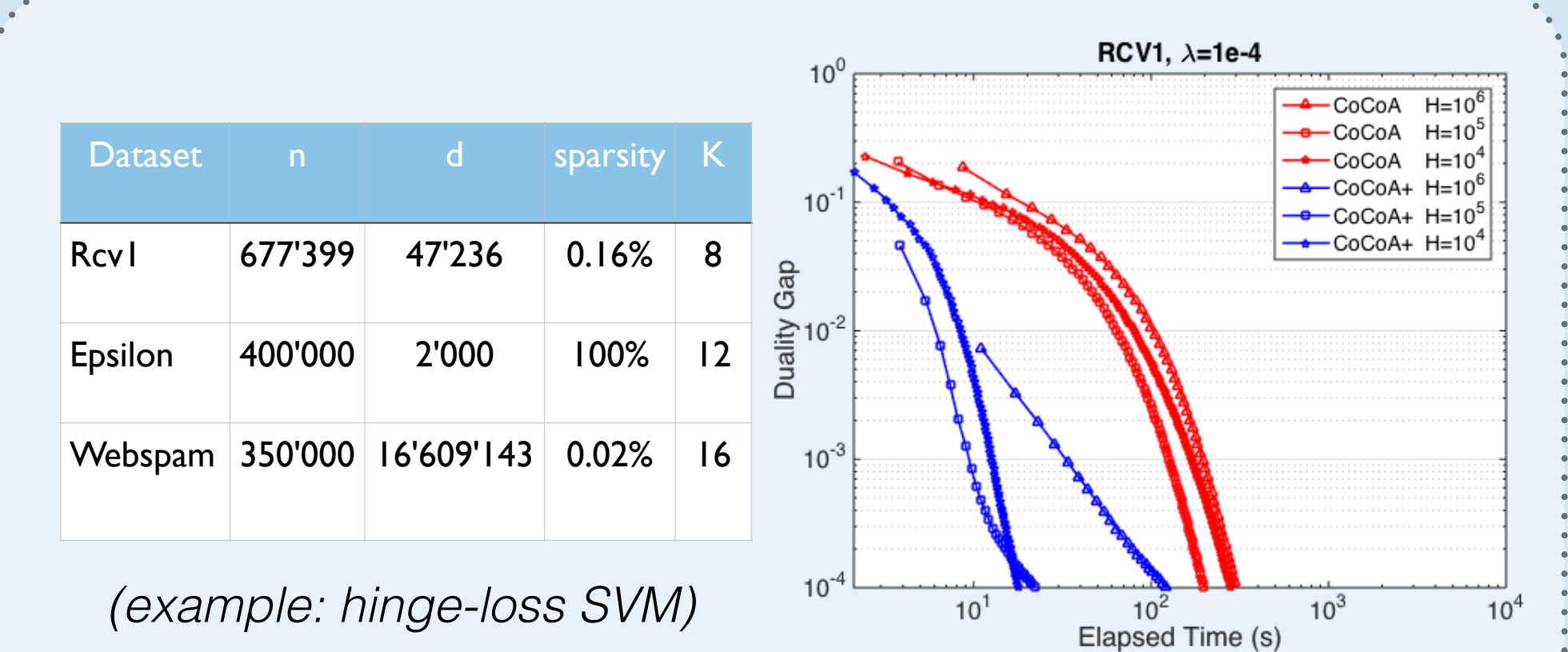
- CoCoA, averaging  $\gamma = 1/K$   

$$T \geq \frac{1}{1-\Theta} \frac{\lambda \mu K + 1}{\lambda \mu} \log \frac{1}{\epsilon}$$
- CoCoA<sup>+</sup>, adding  $\gamma = 1$   

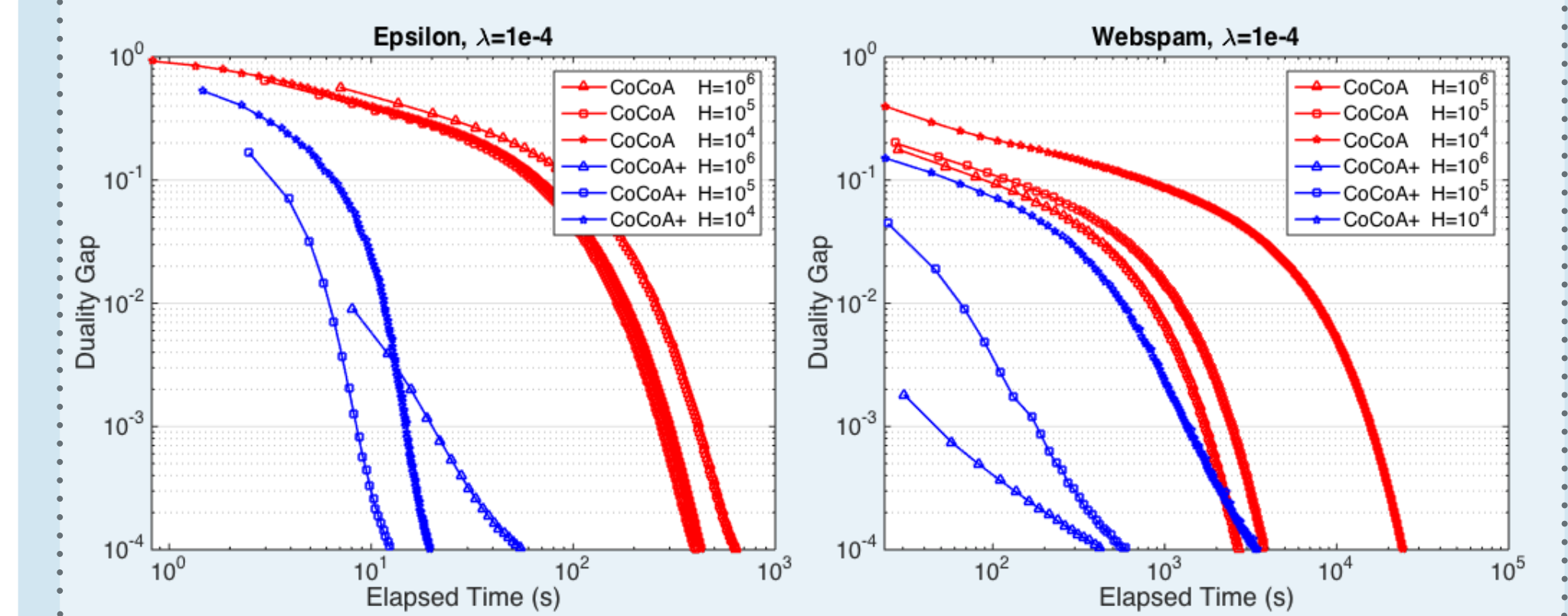
$$T \geq \frac{1}{1-\Theta} \frac{\lambda \mu + 1}{\lambda \mu} \log \frac{1}{\epsilon}$$

- CoCoA<sup>+</sup> rate *independent* of  $K$
- applies also to *duality gap*

## experiments in *Spark*



(example: hinge-loss SVM)



\*Experiments with SDCA as a local solver -> reduces to DisDCA-p [1]

$H$  = number of local updates per round

code at: [github.com/gingsmith/cocoa](https://github.com/gingsmith/cocoa)

## setup

Primal problem formulation

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

data partitioned by examples  $A_i := \mathbf{x}_i$

*primal-dual correspondence*  
 $\mathbf{w} = \frac{1}{\lambda n} A \alpha$

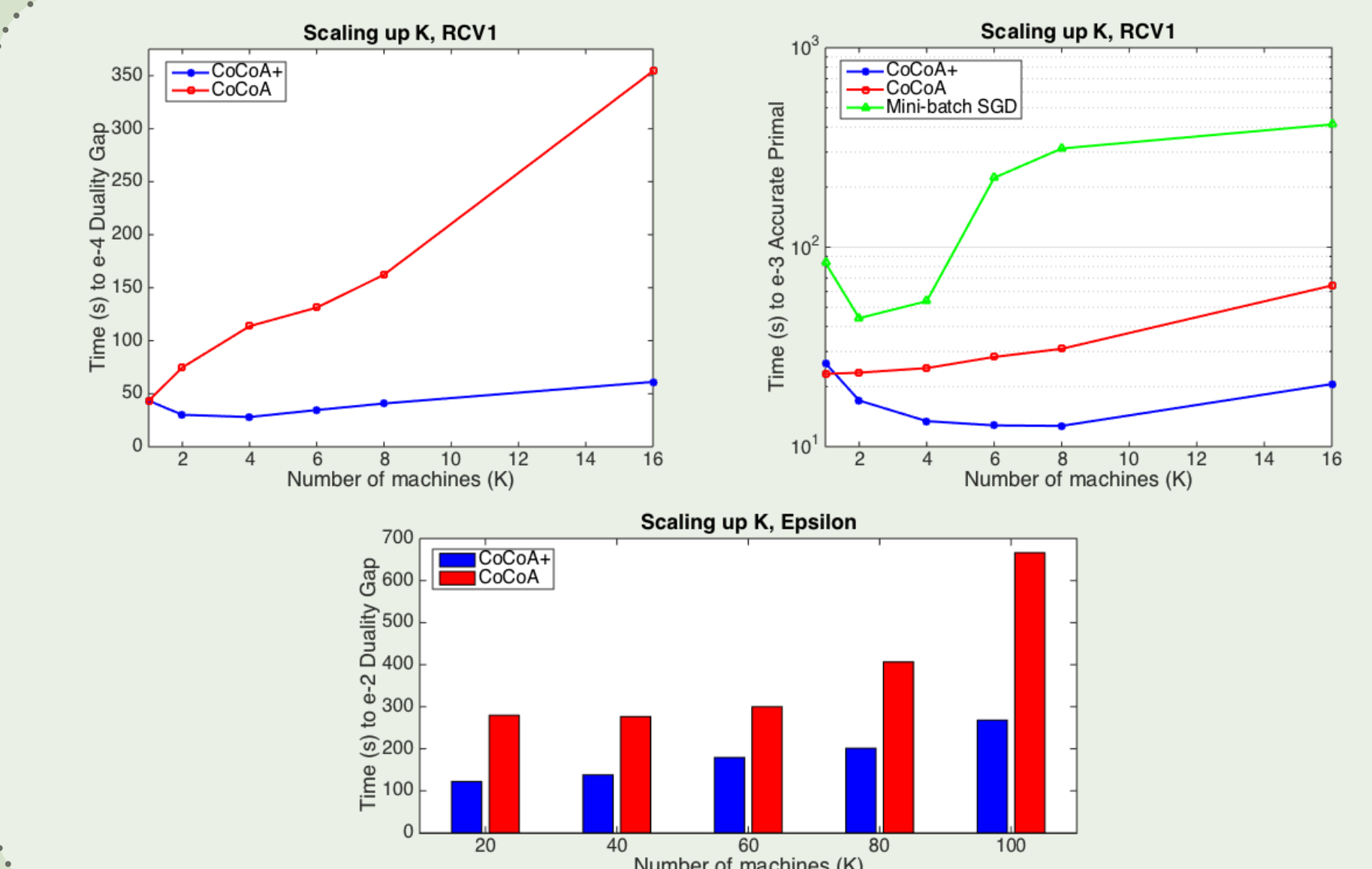
Dual problem

$$\max_{\alpha \in \mathbb{R}^n} -\frac{\lambda}{2} \left\| \frac{A \alpha}{\lambda n} \right\|^2 - \frac{1}{n} \sum_{i=1}^n \ell_i^*(-\alpha_i)$$

Information: local      shared

$$\ell_i^*(s) := \sup_{t \in \mathbb{R}} \{st - \ell_i(t)\}$$

## scaling



## references

- [1] Yang. *Trading computation for communication: Distributed stochastic coordinate ascent*. NIPS, 2013.
- [2] Shalev-Shwartz and Zhang. *Stochastic dual coordinate ascent methods for regularized loss minimization*. JMLR, 14:567–599, 2013.
- [3] Jaggi et al., *Communication-Efficient Distributed Dual Coordinate Ascent*. NIPS, 2014.