



# Cyclic Block Coordinate Methods on a Finer Scale: Tighter Bounds and New Methods

Jelena Diakonikolas (UW-Madison)

Jason's Optimization Seminar



# Outline



- Recap of coordinate methods (cyclic vs randomized)
- A new cyclic method for variational inequalities
- Generalizations
- A fun problem I am looking at (if enough time!)



# A Quick Refresher on (Block) Coordinate Methods

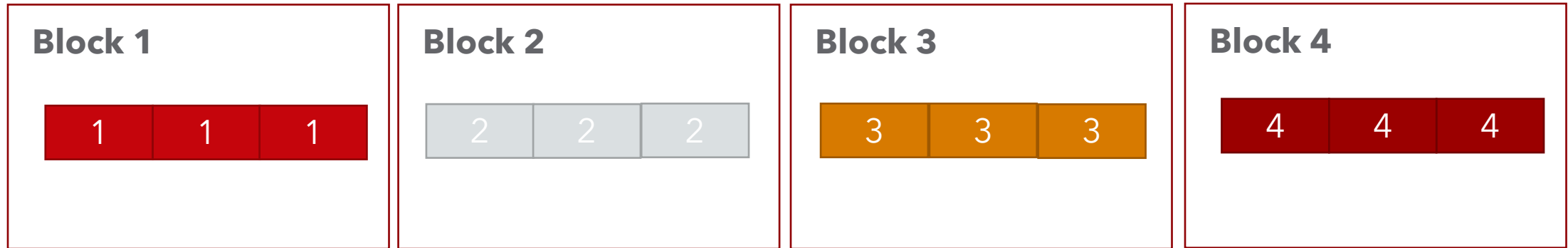
- Fix a partition of the vector of variables into  $m$  blocks:





# A Quick Refresher on (Block) Coordinate Methods

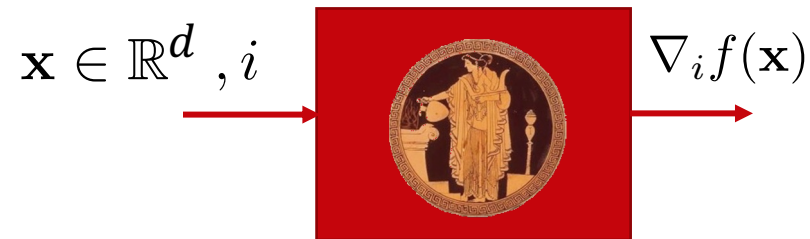
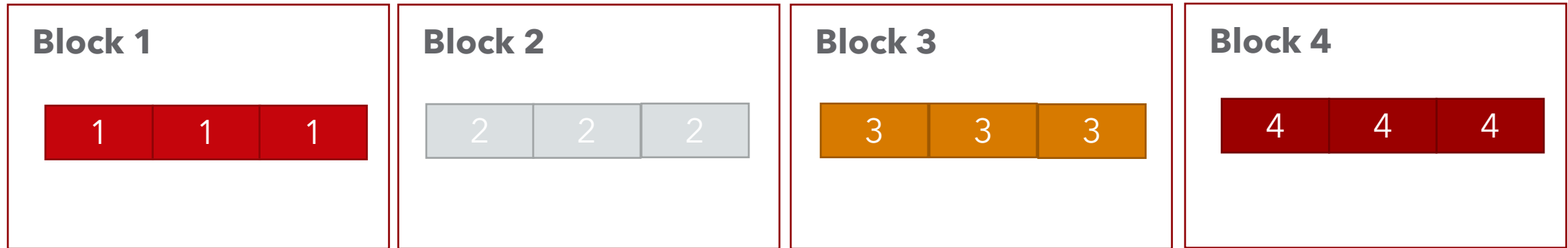
- Fix a partition of the vector of variables into  $m$  blocks:





# A Quick Refresher on (Block) Coordinate Methods

- Fix a partition of the vector of variables into  $m$  blocks:

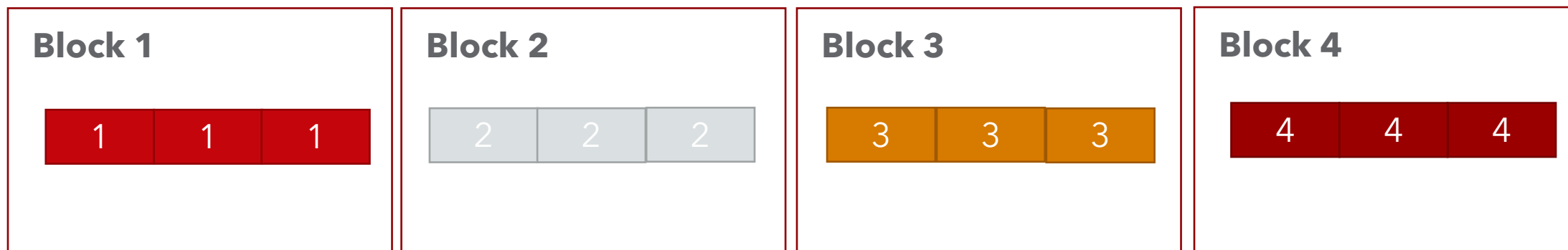


**first-order oracle**



# A Quick Refresher on (Block) Coordinate Methods

- Fix a partition of the vector of variables into  $m$  blocks:

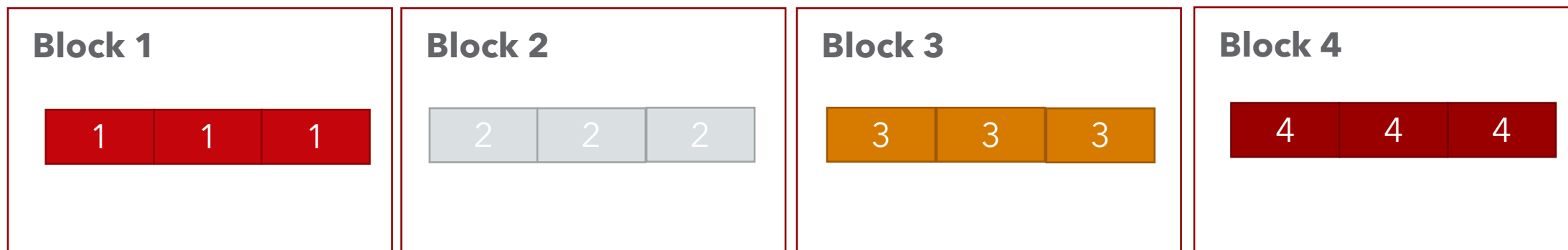


- Types of methods/orderings of updates:
  - cyclic**: fix an order of blocks, go through all of them in a cycle;
  - randomized**: pick blocks randomly, sample with replacement;
  - greedy**: pick the block that leads to the largest progress



# A Quick Refresher on (Block) Coordinate Methods

- Fix a partition of the vector of variables into  $m$  blocks:



- Types of methods/orderings of updates:

- **cyclic**: fix an order of blocks, go through all of them in a cycle;
- **randomized**: pick blocks randomly, sample with replacement;
- **greedy**: pick the block that leads to the largest progress



# Theory vs Practice?

- **Randomized methods** [Strohmer & Vershynin'08], [Nesterov'12]:
  - almost always faster than full gradient methods, assuming the problem is (block) coordinate friendly [Nesterov'12] + a lot of follow-up work;
  - generally the best theoretical guarantees among block coordinate methods;
  - **key property**: can relate the partial gradient to the full one by taking the expectation, helps much of the analysis carry over from full gradient methods
- **Cyclic methods** [Kaczmarz'37], [Ortega & Rheinboldt'70]:
  - often preferred in practice over randomized methods (e.g., in GLMNet, SparseNet);
  - much more challenging to analyze; hard to relate partial gradients to full ones;
  - complexity guarantees generally **worse than even for full gradient methods** and **smooth cvx opt**, by dimension-dependent factors [Beck-Tetruashvili'13]; this is tight for cyclic gradient descent-type method [Sun-Ye'19]



# State of the Art for Cyclic Methods (Prior to This Work)



- Essentially **no non-asymptotic results** for min-max opt/variational inequalities
  - Exception: [Chow-Wu-Yin'17], but requires cocoercivity and the rate is  $1/\sqrt{k}$
- For (non-accelerated) **smooth convex optimization**, number of full gradient queries (assuming coordinate-friendly) of the order  $O\left(\frac{mLD^2}{\epsilon}\right)$ , at best (**worse by a factor  $m = \#$  of blocks than gradient descent;  $m = d$  for coordinate descent**)

The only exception are convex **quadratic** problems [Gürbüzbalaban et al., 2017], [Lee & Wright, 2019]

- \*All analyses based on **relating the partial gradient to the full one**

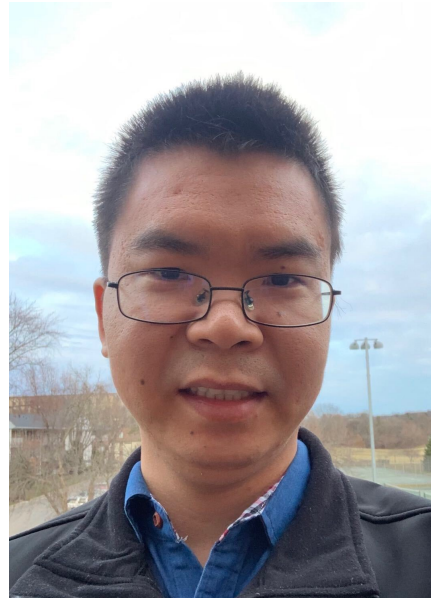
\*as far as I can tell, please correct me if wrong



# Cyclic Coordinate Dual Averaging with Extrapolation (CODER)

---

## Setup and Results



Joint Work With:

Chaobing Song (Huawei)

---

C. Song, J. Diakonikolas, "Cyclic Coordinate Dual Averaging with Extrapolation,"  
*SIAM Journal on Optimization*, vol. 33, no. 4, pp. 2935-2961, 2023.



# Problem Setup

(GMVI) Find  $x^* \in \mathbb{R}^d$  s.t.  $\forall x, \langle F(x), x - x^* \rangle + g(x) - g(x^*) \geq 0$

$$(P_{CO}) \min_{x \in \mathbb{R}^d} f(x) + g(x)$$

$$(P_{MM}) \min_x \max_y \Phi(x, y)$$

## Assumptions:

- There is a fixed partition of coordinates into  $m$  blocks;
- $F: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is
  - “(block) coordinate-friendly” according to that partition;
  - **monotone**:  $\forall x, y: \langle F(x) - F(y), x - y \rangle \geq 0$ ;
  - **Lipschitz**:  $\exists L < \infty$  s.t.  $\forall x, y: \|F(x) - F(y)\| \leq L\|x - y\|$
- $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is
  - **block-separable** over the given partition:  $g(x) = \sum_{j=1}^m g^j(x^j)$ ;
  - $\text{prox}_{\tau g^j}(x^j) = \arg \min_{y \in \mathbb{R}^{d_j}} \left\{ g^j(y) + \frac{1}{2\tau} \|y - x^j\|^2 \right\}$  is **easily computable**,  $\forall j \in \{1, \dots, m\}$ ;
  - possibly strongly **convex**, with modulus  $\gamma \geq 0$  and **lower semicontinuous**

I'll focus on the cyclic coordinate case ( $m = d$ ), for simplicity



# Coordinate Lipschitz Assumption?

- What is standard in convex optimization:

$$|\nabla^j f(x) - \nabla^j f(x + \underbrace{h}_{\text{scalar}} \underbrace{e_j}_{j^{\text{th}} \text{ standard basis vector}})| \leq L_j |h|, \quad \forall x \in \mathbb{R}^d, h \in \mathbb{R}$$

- For VIs unclear how to make useful. Take bilinear games  $\min_x \max_y x^T A y$ . Then

$$F \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) = \begin{bmatrix} 0 & A \\ -A^T & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} Ay \\ -A^T x \end{bmatrix},$$

so  $F^j \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) = A_{j:y}$  for  $j$  in the block belonging to the  $x$ -player. In particular,

$$F^j \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) - F^j \left( \begin{bmatrix} x + h e_j \\ y \end{bmatrix} \right) = 0$$



# A Different Coordinate Lipschitz Condition?

- Going back to the bilinear case, for any  $\begin{bmatrix} x \\ y \end{bmatrix}, \begin{bmatrix} x' \\ y' \end{bmatrix}$  and any  $j$  in the  $x$ -part,

$$F^j \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) - F^j \left( \begin{bmatrix} x' \\ y' \end{bmatrix} \right) = A_{j:} (y - y')$$

- From there, we can conclude:

$$\left| F^j \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) - F^j \left( \begin{bmatrix} x' \\ y' \end{bmatrix} \right) \right|^2 = \begin{bmatrix} x - x' \\ y - y' \end{bmatrix}^T \begin{bmatrix} 0 & 0 \\ 0 & A_{j:} A_{j:}^T \end{bmatrix} \begin{bmatrix} x - x' \\ y - y' \end{bmatrix}.$$

quadratic form

symmetric  
PSD matrix

- **Idea:** generalize to other (possibly nonlinear) operators  $F$



# New Lipschitz Condition

There exist symmetric positive semidefinite matrices  $Q^1, Q^2, \dots, Q^d$  such that for any  $z, z' \in \mathbb{R}^d$ ,

$$|F^j(z) - F^j(z')|^2 \leq (z - z')^T Q^j (z - z').$$

- Can be trivially satisfied with  $Q^j = LI$ , but can generally choose better
- Define:

$$\hat{Q}^j = \begin{bmatrix} \overset{1,2,\dots,j-1}{\mathbf{0}} & \mathbf{0} \\ \mathbf{0} & Q^j \end{bmatrix} \quad ]^{1,2,\dots,j-1}$$

- Summary Lipschitz constant:

$$\hat{L} = \sqrt{\left\| \sum_{j=1}^d \hat{Q}^j \right\|}$$

# Main Result



A method (CODER) that for any  $\epsilon > 0$  outputs a solution with primal-dual gap at most  $\epsilon$ , in

$$O\left(\min\left\{\frac{\hat{L}D^2}{\epsilon}, \frac{\hat{L}}{\gamma} \log\left(\frac{\hat{L}D}{\epsilon}\right)\right\}\right)$$

iterations.

- This is **the same** as for full-vector update methods, but with  $L$  replaced by  $\hat{L}$





# How Large is $\hat{L}$ ?

- Worst case:

$$\hat{L} = \sqrt{\left\| \sum_{j=1}^d \hat{Q}^j \right\|} \leq \sqrt{\sum_{j=1}^d \|\hat{Q}^j\|} \leq \sqrt{\sum_{j=1}^d \|Q^j\|} \leq L\sqrt{d}$$

triangle inequality       $\hat{Q}^j \preceq Q^j$        $Q^j \preceq LI$

- Even for the special case of smooth convex optimization, the resulting bound is better by a factor  $\sqrt{d}$  than what was known for any (unaccelerated) cyclic coordinate method.
- For variational inequalities, the obtained complexity result is state-of-the-art, even compared to randomized methods [Kotsalis et al., 2022] (it is actually better by a factor  $\sqrt{d}$  in the worst case).
- It is also the first cyclic method for general variational inequalities with monotone operators with provable convergence guarantees.

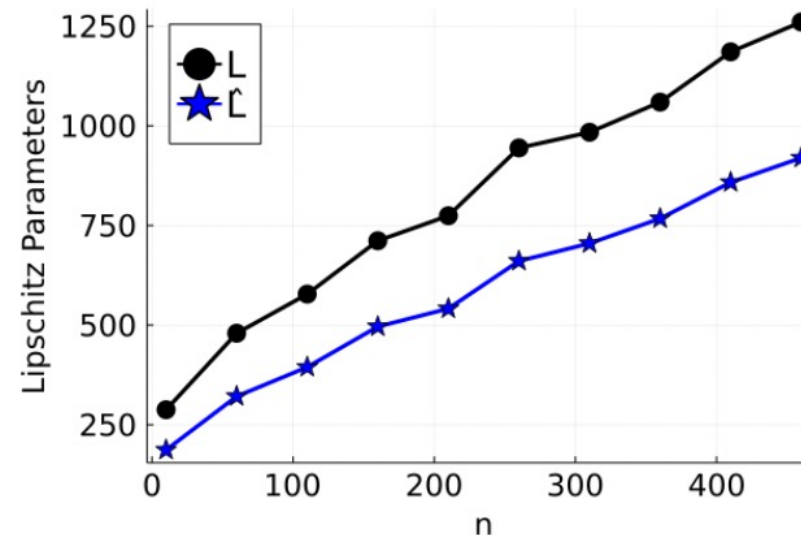
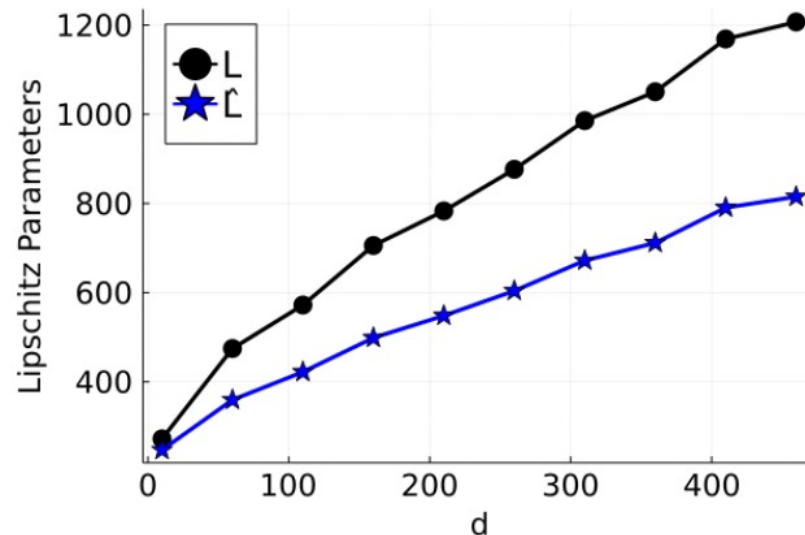
# How Large is $\hat{L}$ ?

- In practice,  $\hat{L}$  is no larger than  $L$ , usually smaller (constants below are for min-max SVM)

- Real data sets:

Dataset	a9a	australia	madelon	colon	mnist
$L$	15389.6	340.5	1992.4	9.0	24410.5
$\hat{L}$	10358.8	238.1	1269.7	5.7	15236.1

- Synthetic (standard Gaussian) data; fixed dataset size  $n$  or fixed dimension  $d$ , respectively



# Algorithm



---

**Algorithm 3.1** Cyclic cOordinate Dual avEraging with extRapolation (CODER)

---

1: **Input:**  $\mathbf{x}_{-1} = \mathbf{x}_0 \in \text{dom}(g)$ ,  $\gamma \geq 0$ ,  $\hat{L} > 0$ ,  $m$ ,  $\{\mathcal{S}^1, \dots, \mathcal{S}^m\}$

2: **Initialization:**  $\mathbf{p}_0 = \mathbf{F}(\mathbf{x}_0)$ ,  $\mathbf{z}_0 = \mathbf{0}$ ,  $a_0 = A_0 = 0$

3: **for**  $k = 1$  to  $K$  **do**

4:  $a_k = \frac{1+\gamma A_{k-1}}{2\hat{L}}$ ,  $A_k = A_{k-1} + a_k$

5: **for**  $j = 1$  to  $m$  **do**

6:  $\mathbf{p}_k^j = \mathbf{F}^j(\mathbf{x}_k^1, \dots, \mathbf{x}_k^{j-1}, \mathbf{x}_{k-1}^j, \dots, \mathbf{x}_{k-1}^m)$  ← partial “gradient” at intermediate point

7:  $\mathbf{q}_k^j = \mathbf{p}_k^j + \frac{a_{k-1}}{a_k} (\mathbf{F}^j(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}^j)$  ← partial “gradient” extrapolation

8:  $\mathbf{z}_k^j = \mathbf{z}_{k-1}^j + a_k \mathbf{q}_k^j$

9:  $\mathbf{x}_k^j = \text{prox}_{A_k g^j}(\mathbf{x}_0^j - \mathbf{z}_k^j)$  } dual averaging step with extrapolated gradient

10: **end for**

11: **end for**

12: **return**  $\tilde{\mathbf{x}}_K = \frac{1}{A_K} \sum_{k=1}^K a_k \mathbf{x}_k$ ,  $\mathbf{x}_K$

---

Similar gradient extrapolation in [Hamedani & Aybat., 2018] and [Kotsalis et al., 2022]



# How the Analysis Works

(GMVI) Find  $x^* \in \mathbb{R}^d$  s.t.  $\forall x, \langle F(x), x - x^* \rangle + g(x) - g(x^*) \geq 0$

- We can only solve this approximately to error  $\epsilon > 0$ , so equivalently want to

$$\text{find } x_\epsilon^* \in \mathbb{R}^d \text{ s.t. } \forall x, \langle F(x), x - x_\epsilon^* \rangle + g(x) - g(x_\epsilon^*) \geq -\epsilon$$

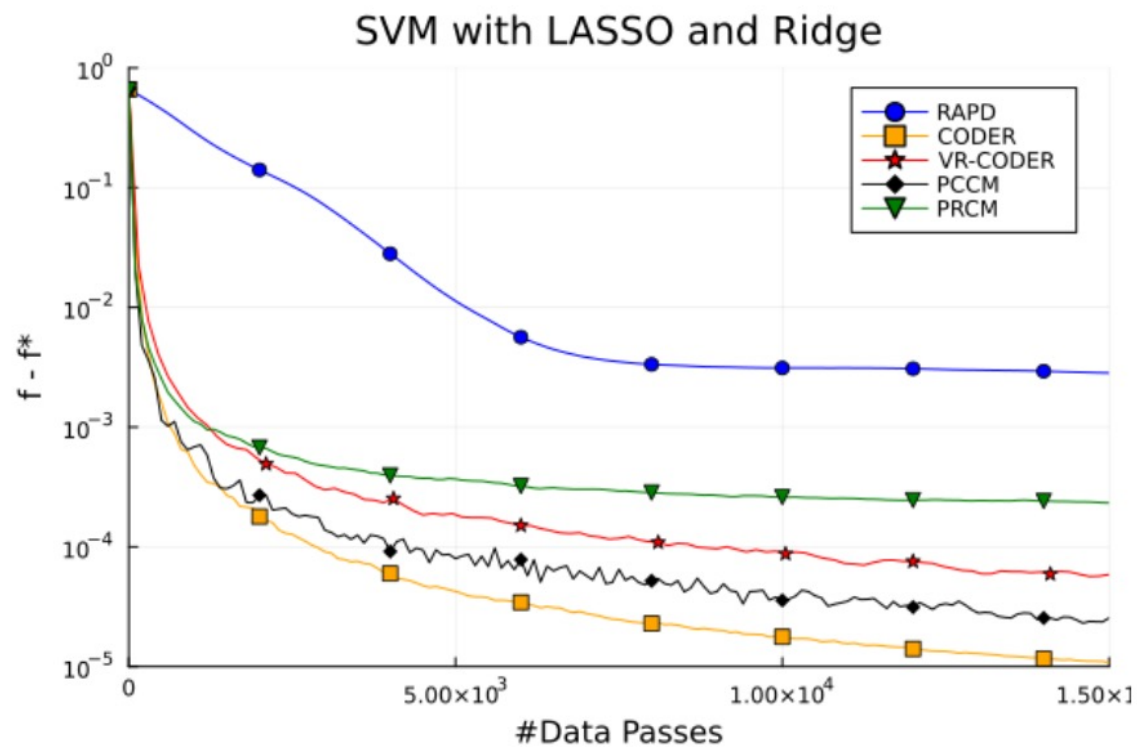
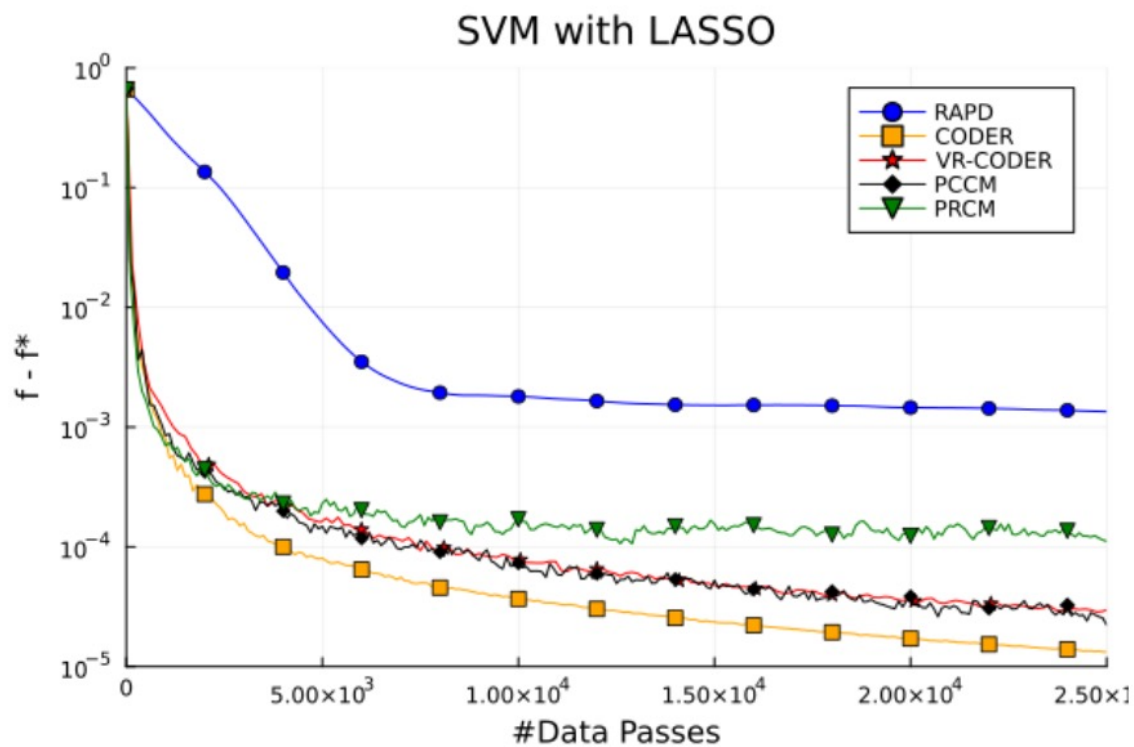
which is the same as find  $x_\epsilon^* \in \mathbb{R}^d$  s.t.  $\forall x, \text{Gap}(x, x_\epsilon^*) := \langle F(x), x_\epsilon^* - x \rangle - g(x) + g(x_\epsilon^*) \leq \epsilon$

- So we may try to bound  $\text{Gap}(x, x_k)$  for iterations  $k = 1, 2, \dots$
- The key step:

$$\begin{aligned} \text{Gap}(x, x_k) &= \langle F(x), x_k - x \rangle - g(x) + g(x_k) && \text{this part defines the} \\ &\leq \langle F(x_k), x_k - x \rangle - g(x) + g(x_k) && \text{step, by taking a max} \\ &= \langle q_k, x_k - x \rangle - g(x) + g(x_k) - \frac{1}{2} \|x - x_k\|_2^2 + \frac{1}{2} \|x - x_k\|_2^2 \\ &\quad + \langle F(x_k) - q_k, x_k - x \rangle && \text{the art of choosing } q_k \text{ is for} \\ &&& \text{controlling this term} \end{aligned}$$

# Numerical Experiments (Illustration)

SVM with LASSO or ridge,  
on a1a LibSVM dataset  
( $d = 123, n = 1605$ )





# Further Developments

---

Two Examples: Bilinear Extensive Form Games and Shuffled SGD



Q: Are cyclic methods always slower than full vector update methods in the worst case?

A: No, and they may have advantages. A specific example (the first of its kind!) to follow.



Darshan Chakrabarti  
(Columbia University)



Christian Kroer  
(Columbia University)

## Joint Work With:

---

D. Chakrabarti, **J. Diakonikolas**, C. Kroer "*Block Coordinate Methods and Restarting for Solving Extensive Form Games*," in *Proc. NeurIPS 2023*. ( $\alpha\beta$  ordering of authors)





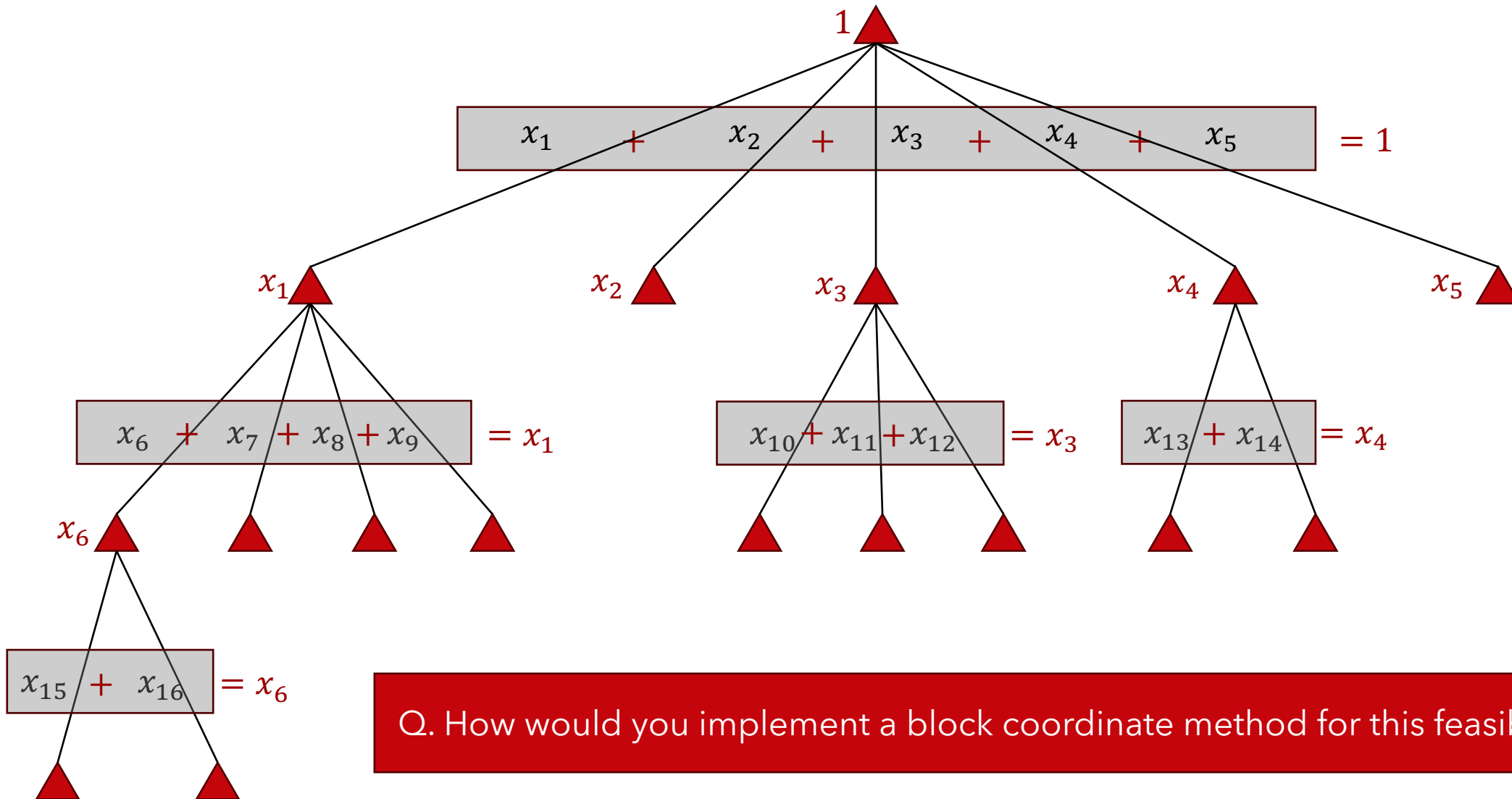
# Problem: Bilinear Extensive Form Games (EFGs)

- A general class of game-theoretic models that capture both simultaneous and sequential moved, private/imperfect information, and stochasticity
- In optimization language, we have a bilinear min-max problem:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} y^T M x$$

sequence form polytopes: treeplexes

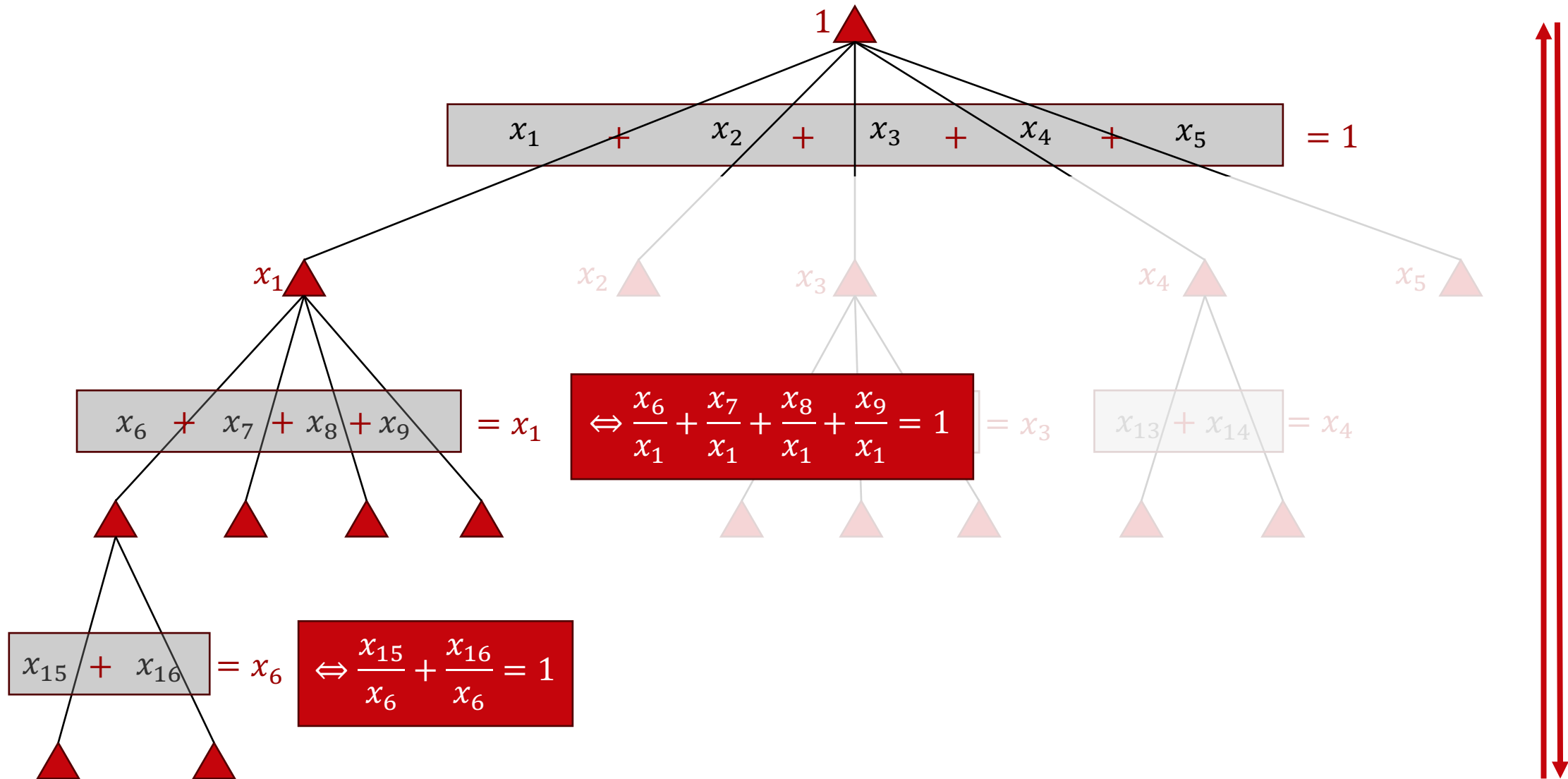
# Treeplex?



Q. How would you implement a block coordinate method for this feasible set?



# Why Cyclic Updates are OK



# Main Result



A method (ECyclicPDA) that for any  $\epsilon > 0$  outputs a solution with primal-dual gap at most  $\epsilon$  with iteration complexity that is **no worse** than the iteration complexity of full-vector update methods like Mirror-Prox.

This is the first example (that I know of) of a cyclic method with no scaling with the number of blocks in the worst case.



Q: Do these ideas extend beyond basic cyclic coordinate methods?

A: Yes. They extend to the incremental gradient method and shuffled SGD.



Xufeng Cai  
(UW-Madison)



Eric (Cheuk-Yin) Lin  
(UW-Madison)

## Joint Work With:

---

X. Cai, C-Y. Lin, [J. Diakonikolas](#), "Empirical Risk Minimization with Shuffled SGD: A Primal-Dual Perspective and Improved Bounds," arXiv preprint, arXiv:2306.12498, 2023.



# Problem Setup

Empirical Risk Minimization problems over data  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$

$$\min_{\mathbf{x}} \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{a}_i^\top \mathbf{x})$$

are (in practice) usually solved by **incremental/shuffled SGD** methods, which make updates

$$\mathbf{x}_{k,i+1} = \mathbf{x}_{k,i} - \eta \nabla \ell_i(\mathbf{a}_i^\top \mathbf{x}_{k,i})$$

going through all the data vectors in a cyclic manner, possibly permuting the order at the beginning of a cycle.

**Convergence guarantees:** not understood until recently [[Gürbüzbalaban et al., 2021](#)], [[Shamir, 2016](#)], [[Haochen & Sra, 2019](#)], [[Nagaraj et al., 2019](#)], [[Rajput et al., 2020](#)], [[Ahn et al., 2020](#)], [[Mishchenko et al., 2020](#)], [[Nguyen et al., 2021](#)], [[Cha et al., 2023](#)]



# Main Insight

- We can view incremental gradient/shuffled SGD as a **primal-dual method** with **cyclic updates on the dual side**

---

**Algorithm 1** Shuffled SGD (Primal-Dual View)

---

- 1: **Input:** Initial point  $\mathbf{x}_0 \in \mathbb{R}^d$ , batch size  $b > 0$ , step size  $\{\eta_k\} > 0$ , number of epochs  $K > 0$
  - 2: **for**  $k = 1$  to  $K$  **do**
  - 3:   Generate any permutation  $\pi^{(k)}$  of  $[n]$  (either deterministic or random)
  - 4:    $\mathbf{x}_{k-1,1} = \mathbf{x}_{k-1}$
  - 5:   **for**  $i = 1$  to  $m$  **do**
  - 6:      $\mathbf{y}_k^{(i)} = \arg \max_{\mathbf{y} \in \mathbb{R}^b} \left\{ \mathbf{y}^\top \mathbf{A}_k^{(i)} \mathbf{x}_{k-1,i} - \sum_{j=1}^b \ell_{\pi_{b(i-1)+j}^{(k)}}^*(\mathbf{y}^j) \right\}$
  - 7:      $\mathbf{x}_{k-1,i+1} = \arg \max_{\mathbf{x} \in \mathbb{R}^d} \left\{ \mathbf{y}_k^{(i)\top} \mathbf{A}_k^{(i)} \mathbf{x} + \frac{b}{2\eta_k} \|\mathbf{x} - \mathbf{x}_{k-1,i}\|^2 \right\}$
  - 8:   **end for**
  - 9:    $\mathbf{x}_k = \mathbf{x}_{k-1,m+1}$ ,  $\mathbf{y}_k = (\mathbf{y}_k^{(1)}, \mathbf{y}_k^{(2)}, \dots, \mathbf{y}_k^{(m)})^\top$
  - 10: **end for**
  - 11: **Return:**  $\hat{\mathbf{x}}_K = \sum_{k=1}^K \eta_k \mathbf{x}_k / \sum_{k=1}^K \eta_k$
-



# Main Results



Tighter, data-dependent, convergence bounds for all standard variants of incremental gradient/shuffled SGD, for smooth convex loss functions. The obtained bounds are tighter by a factor that can be as large as  $\sqrt{n}$ , both in theory and in practice.



# How Tighter in Practice?

DATASET	#FEATURES ( $d$ )	#DATAPOINTS ( $n$ )	$L/\hat{L}$	$\log_n L/\hat{L}$
A1A	123	1605	5.50	0.231
A9A	123	32561	5.49	0.164
BBBC005	361920	19201	18.3	0.295
BBBC010	361920	201	7.04	0.368
CIFAR10	3072	50000	10.0	0.213
<b>DUKE</b>	<b>7129</b>	<b>44</b>	<b>38.0</b>	<b>0.962</b>
E2006TRAIN	150360	16087	5.35	0.173
GISETTE	5000	6000	3.52	0.145
<b>LEU</b>	<b>7129</b>	<b>38</b>	<b>32.8</b>	<b>0.960</b>
MNIST	780	60000	19.1	0.268
NEWS20	1355191	19996	42.1	0.378
RCV1	47236	20242	111	0.475
<b>REAL-SIM</b>	<b>20958</b>	<b>72309</b>	<b>194</b>	<b>0.471</b>
SONAR	60	208	6.26	0.344
TMC2007	30438	21519	10.9	0.239



# Other Extensions I Did Not Talk About

- Variance reduction for cyclic methods [Song, D, 2021], [Lin, Song, D, ICML 2023], [Cai, Song, Wright, D, ICML 2023]
- Acceleration in smooth convex optimization [Lin, Song, D, ICML 2023]
- Nonconvex optimization [Cai, Song, Wright, D, ICML 2023]
- Incremental gradient methods and continual learning [Cai, D, forthcoming]

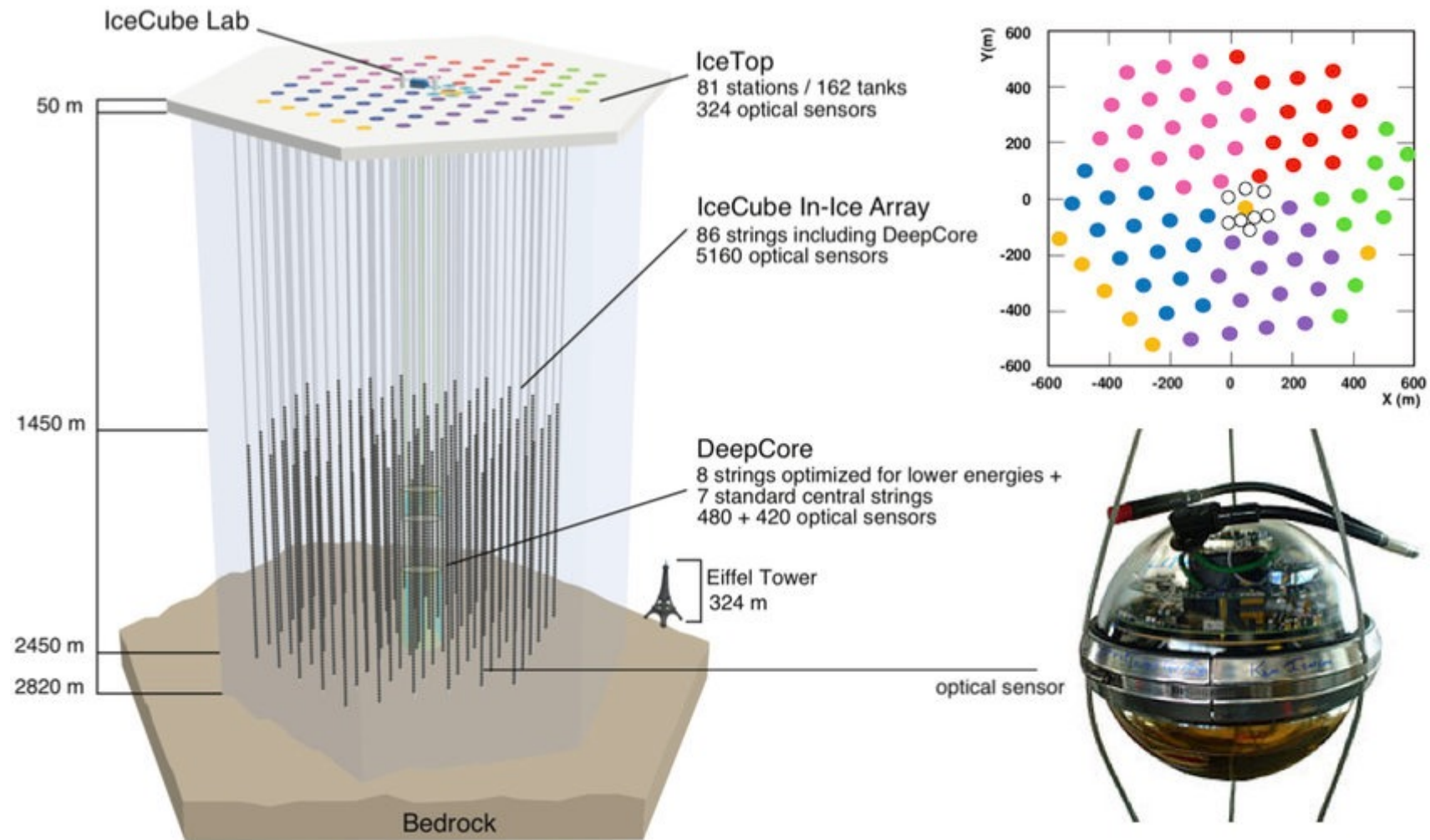


# A Toy Project I am Excited About

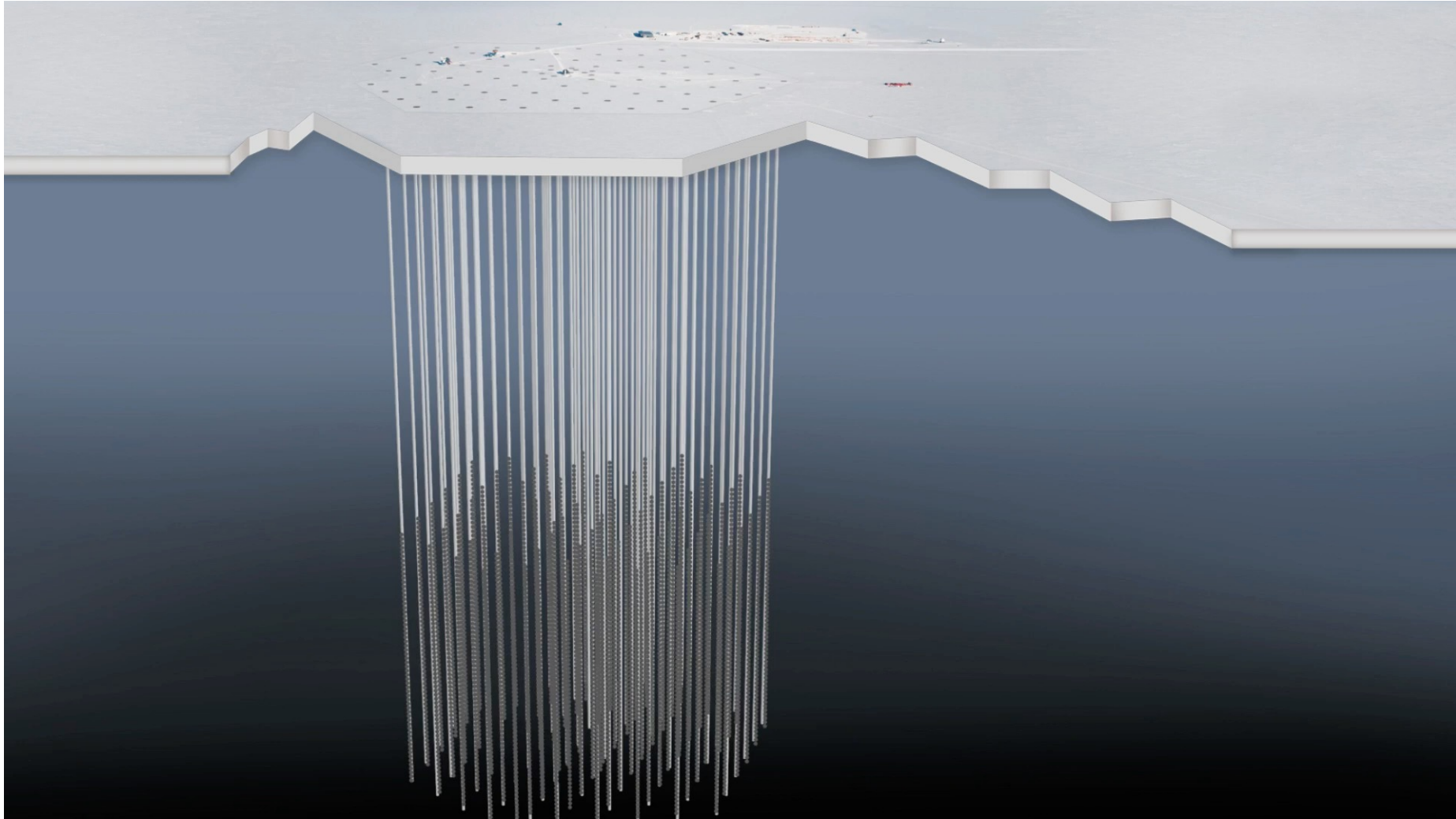
---

and it is for a “real” application!

# Ice Cube Neutrino Detector



# Ice Cube Neutrino Detector

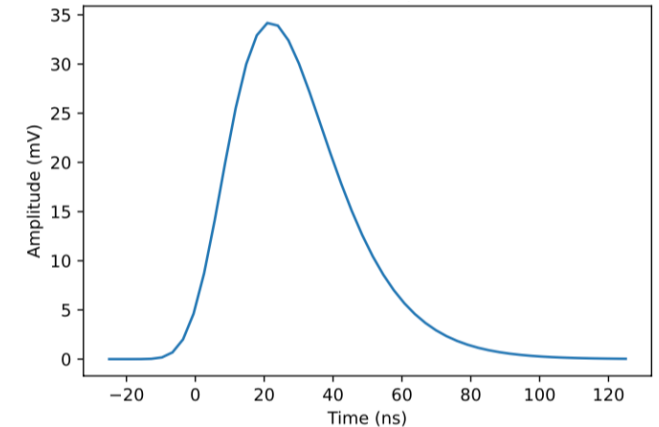




# Neutrino Detection & Non-negative Least Squares

Physics models for detected photon waveforms:

- $f(t) = A \left( e^{-\frac{t-x_0}{b_1}} + e^{\frac{t-x_0}{b_2}} \right)$
- $f(t) = A \left( C e^{-\frac{t-x_0}{b_1}} + e^{\frac{t-x_0}{b_2}} \right)$



Problem can be formulated as a regularized non-negative least-square problem:

$$\min_{\mathbf{x} \geq 0} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + r(\mathbf{x})$$

where  $\mathbf{A}$  is non-negative, highly sparse and structured, and the desired form of regularization  $r(\mathbf{x})$  is unclear.



# Summary

- We should still care about and study cyclic methods!
- We just need to be more careful about how we look at them
- **What next?**
  - For what classes of problems are cyclic methods particularly effective and why?
  - What kind of cyclic methods should we use in practice?

Questions?

[jelena@cs.wisc.edu](mailto:jelena@cs.wisc.edu)