

# On Momentum Acceleration for Randomized Coordinate Descent in Matrix Completion



Matthew Callahan<sup>1</sup>, Trung Vu<sup>2</sup>, and Raviv Raich<sup>1</sup>

<sup>1</sup> School of EECS, Oregon State University, Corvallis, OR 97331-5501, USA

<sup>2</sup> Department of CSEE, University of Maryland, Baltimore County, MD 21250-0002, USA  
callamat@orgonstate.edu



## Matrix Completion

**Goal:** given  $i$ .  $\Omega$ , a set of known entries of an  $m \times n$  matrix  $\mathbf{X}$  and  $ii$ . its rank  $r$ , fill-in the remaining entries:

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \frac{1}{2} \|\mathcal{P}_{\Omega}(\mathbf{X} - \mathbf{M})\|_F^2$$

subject to:  $\text{rank}(\mathbf{X}) \leq r$

$$[\mathcal{P}_{\Omega}(\mathbf{X})]_{ij} = \begin{cases} X_{ij} & (i, j) \in \Omega, \\ 0 & (i, j) \notin \Omega. \end{cases}$$

Movies				
	4	?	?	
	?	?	?	4
	?	?	2	?
	4	?	?	4

## Applications

- Recommender Systems [1]
- Image Inpainting [2]
- Device Localization [3]
- Challenge:** The matrix can be high-dimensional

Node	1	2	3	4	5
1	0	31	?	?	?
2	31	0	27	45	?
3	?	27	0	23	26
4	?	45	23	0	29
5	?	?	26	29	0

## Randomized Coordinate Descent (RCD) for Matrix Completion

- Unconstrained reformulation:

$$\min_{\mathbf{A} \in \mathbb{R}^{m \times r}, \mathbf{B} \in \mathbb{R}^{n \times r}} \frac{1}{2} \|\mathcal{P}_{\Omega}(\mathbf{A}\mathbf{B}^T - \mathbf{M})\|_F^2$$

- Algorithm [11]  $\mathcal{O}(|\Omega|r)$ :

1. Minimize with respect to a coordinate of  $\mathbf{A}^{(k)}$  or of  $\mathbf{B}^{(k)}$

2. After  $(n+m)r$  repetitions of 1., refactor  $\mathbf{A}^{(k)}\mathbf{B}^{(k)T} = \mathbf{X}^{(k)}$

Movies				
	4	?	?	
	?	?	?	4
	?	?	2	?
	4	?	?	4

- Refactor as

$$\mathbf{A}^{(k)} = \mathbf{U}^{(k)}\sqrt{\Sigma^{(k)}}, \quad \mathbf{B}^{(k)} = \mathbf{V}^{(k)}\sqrt{\Sigma^{(k)}}$$

where the SVD of  $\mathbf{X}^{(k)}$  is

$$\mathbf{U}^{(k)}\Sigma^{(k)}\mathbf{V}^{(k)T}$$

## RCD Error Analysis [11]

The expected value of the projected error  $\delta^{(k)} = (\mathbf{Z}^T \mathbf{S} \mathbf{S}^T \mathbf{Z})^{1/2} \mathbf{Z}^T (\text{vec}(\mathbf{A}^{(k)} \mathbf{B}^{(k)T} - \mathbf{M}))$  follows:

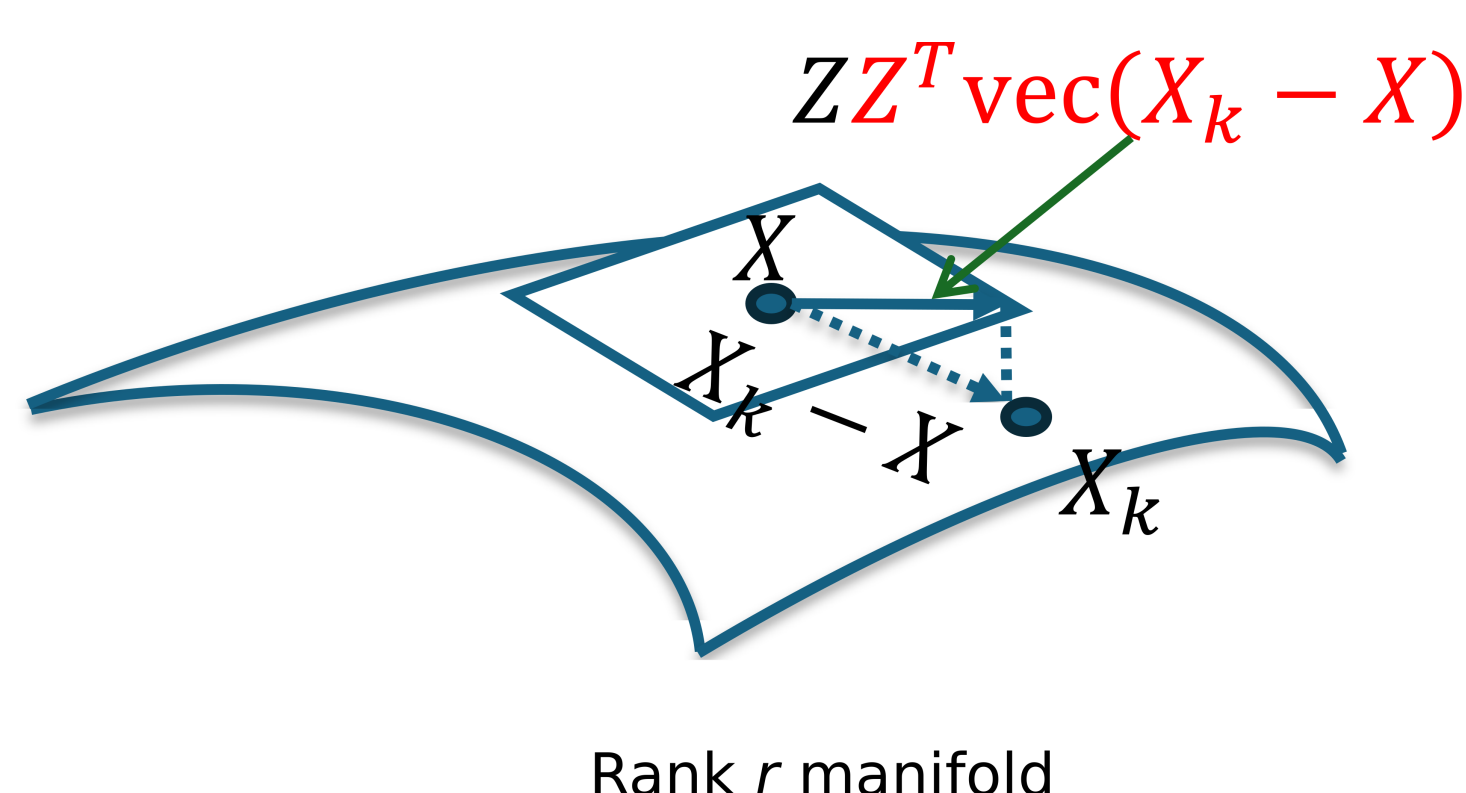
$$\mathbb{E}[\delta^{(k+1)}] = \mathbb{E} \left[ \underbrace{\mathbf{I} - \frac{\mathbf{q}^*(\mathbf{q}^*)^T}{(\mathbf{q}^*)^T \mathbf{q}^*}}_{\hat{\mathbf{T}}} \right] \mathbb{E}[\delta^{(k)}] + o(\|\delta^{(k)}\|^2)$$

$$\mathbf{q}^* = \begin{cases} \sqrt{\sigma_j} (\mathbf{Z}^T \mathbf{S} \mathbf{S}^T \mathbf{Z})^{1/2} \mathbf{Z}^T (\mathbf{v}_j \otimes \mathbf{e}_i^{(m)}) & \text{for } \mathbf{A}_{ij}^{(k)}, \\ \sqrt{\sigma_j} (\mathbf{Z}^T \mathbf{S} \mathbf{S}^T \mathbf{Z})^{1/2} \mathbf{Z}^T (\mathbf{e}_i^{(n)} \otimes \mathbf{u}_j) & \text{for } \mathbf{B}_{ij}^{(k)} \end{cases}$$

**Key result:**

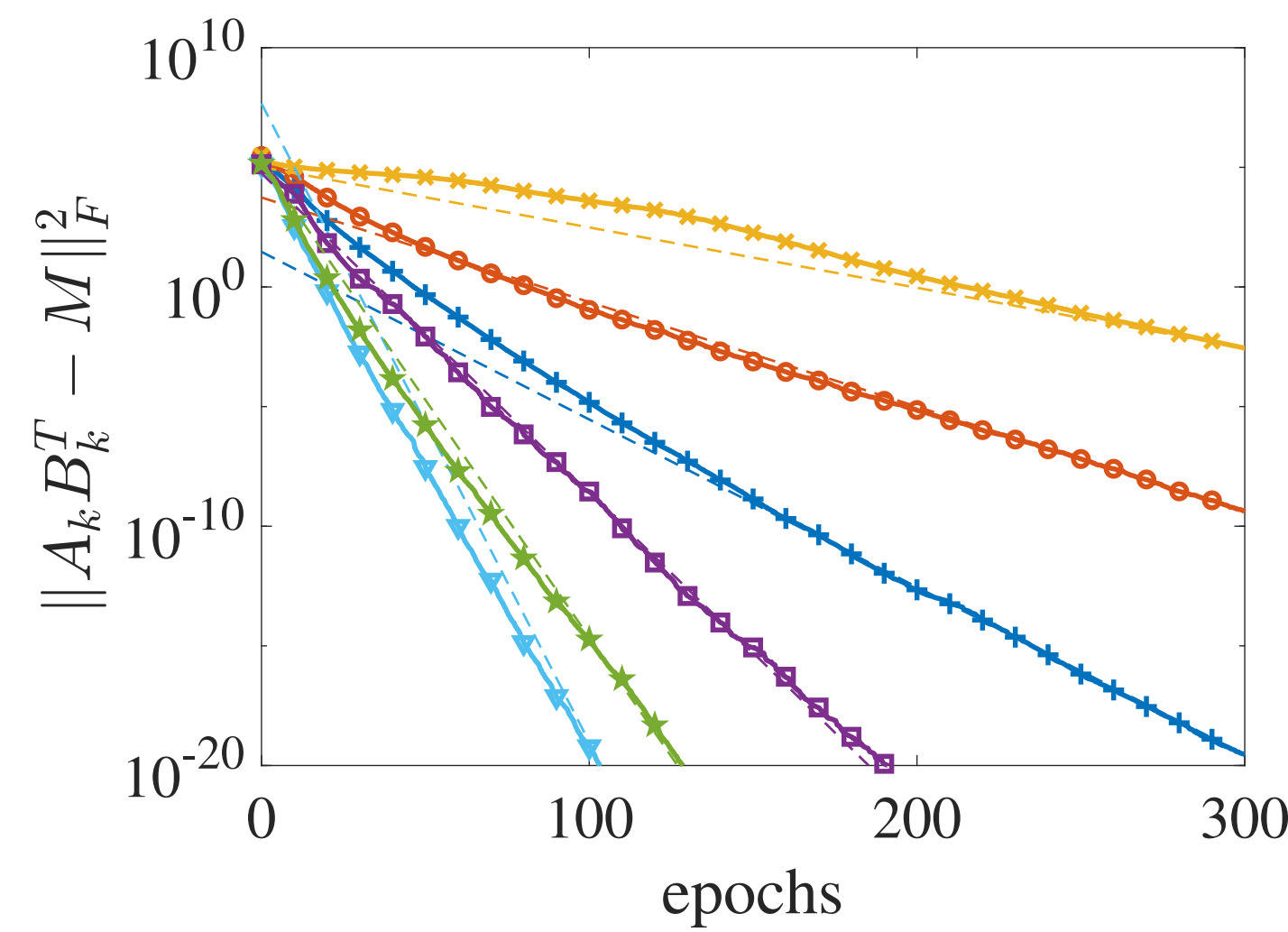
Linear convergence rate  $\rho(\mathbf{T})$  when  $\rho(\mathbf{T}) < 1$

**Goal:** Develop efficient tuning-free accelerated RCD



Rank  $r$  manifold

## RCD Error Analysis Experiments



Comparison of asymptotic rate (dashed) to empirical rate (solid) on  $120 \times 100$  matrices

## Polyak's Momentum Acceleration

Fixed point iteration:

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k), \quad \text{for } k = 0, 1, 2, \dots,$$

Error:

$$\boldsymbol{\epsilon}_{k+1} = \mathbf{T} \boldsymbol{\epsilon}_k + \mathbf{q}(\boldsymbol{\epsilon}_k), \quad \|\mathbf{q}(\boldsymbol{\epsilon})\| \leq q \|\boldsymbol{\epsilon}\|^2$$

Convergence rate:

$$\rho = \rho(\mathbf{T})$$

Accelerated fixed point:

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k) + \beta(\mathbf{x}_k - \mathbf{x}_{k-1}).$$

Accelerated error:

$$\begin{bmatrix} \boldsymbol{\epsilon}_{k+1} \\ \boldsymbol{\epsilon}_k \end{bmatrix} = \mathbf{H}(\beta) \begin{bmatrix} \boldsymbol{\epsilon}_k \\ \boldsymbol{\epsilon}_{k-1} \end{bmatrix} + \begin{bmatrix} \mathbf{q}(\boldsymbol{\epsilon}_k) \\ \mathbf{q}(\boldsymbol{\epsilon}_{k-1}) \end{bmatrix}$$

$$\mathbf{H}(\beta) = \begin{bmatrix} \mathbf{T} + \beta \mathbf{I} & -\beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$$

Convergence rate with optimal  $\beta^*$ :

$$\rho = 1 - \sqrt{1 - \rho(\mathbf{T})}$$

## Momentum Accelerated RCD

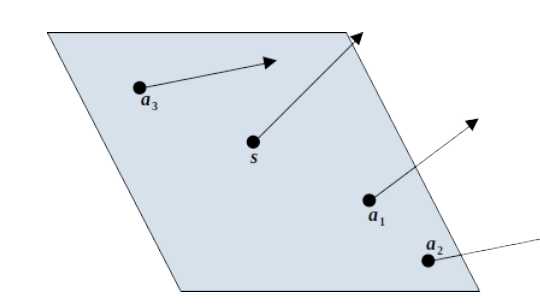
- Initialize  $\mathbf{A}_0, \mathbf{B}_0$
- Set  $\mathbf{A}_1 = \mathbf{A}_0, \mathbf{B}_1 = \mathbf{B}_0$
- for  $k = 1, 2, \dots$ 
  - Generate  $(\tilde{\mathbf{A}}_k, \tilde{\mathbf{B}}_k)$  from iterated RCD:  
 $(\tilde{\mathbf{A}}_k, \tilde{\mathbf{B}}_k) = \underbrace{\text{RCD}(\text{RCD}(\dots \text{RCD}(\mathbf{A}_k, \mathbf{B}_k)))}_{t \text{ nested functions}} \quad \mathcal{O}(t|\Omega|r)$
  - Update the factors as follows:

$$\begin{aligned} \mathbf{A}_{k+1} &= \tilde{\mathbf{A}}_k + \beta(\mathbf{A}_k - \mathbf{A}_{k-1}) \quad \mathcal{O}(mr), \\ \mathbf{B}_{k+1} &= \tilde{\mathbf{B}}_k + \beta(\mathbf{B}_k - \mathbf{B}_{k-1}) \quad \mathcal{O}(nr). \end{aligned}$$

## Modified Refactorization

- RCD applies refactorization for unique solution

$$(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \text{Refactor}(\mathbf{A}, \mathbf{B})$$



- Requirement: Maintain sign consistency
- Solution: Pick direction vector  $\mathbf{s}$  and ensure alignment of  $\mathbf{A}$  on same side of  $\mathbf{s}$  hyperplane

$$\tilde{\mathbf{A}} = \hat{\mathbf{A}} \text{diag}(\text{sign}(\hat{\mathbf{A}}^T \mathbf{s}))$$

$$\tilde{\mathbf{B}} = \hat{\mathbf{B}} \text{diag}(\text{sign}(\hat{\mathbf{A}}^T \mathbf{s}))$$

## Analysis and Stepsize Selection

- Unaccelerated  $t$ -epoch convergence rate:

$$\rho_t = \lambda_{\max}(\mathbf{I} - \mathbf{Q})^{(n+m)rt}$$

- Optimal momentum selection:

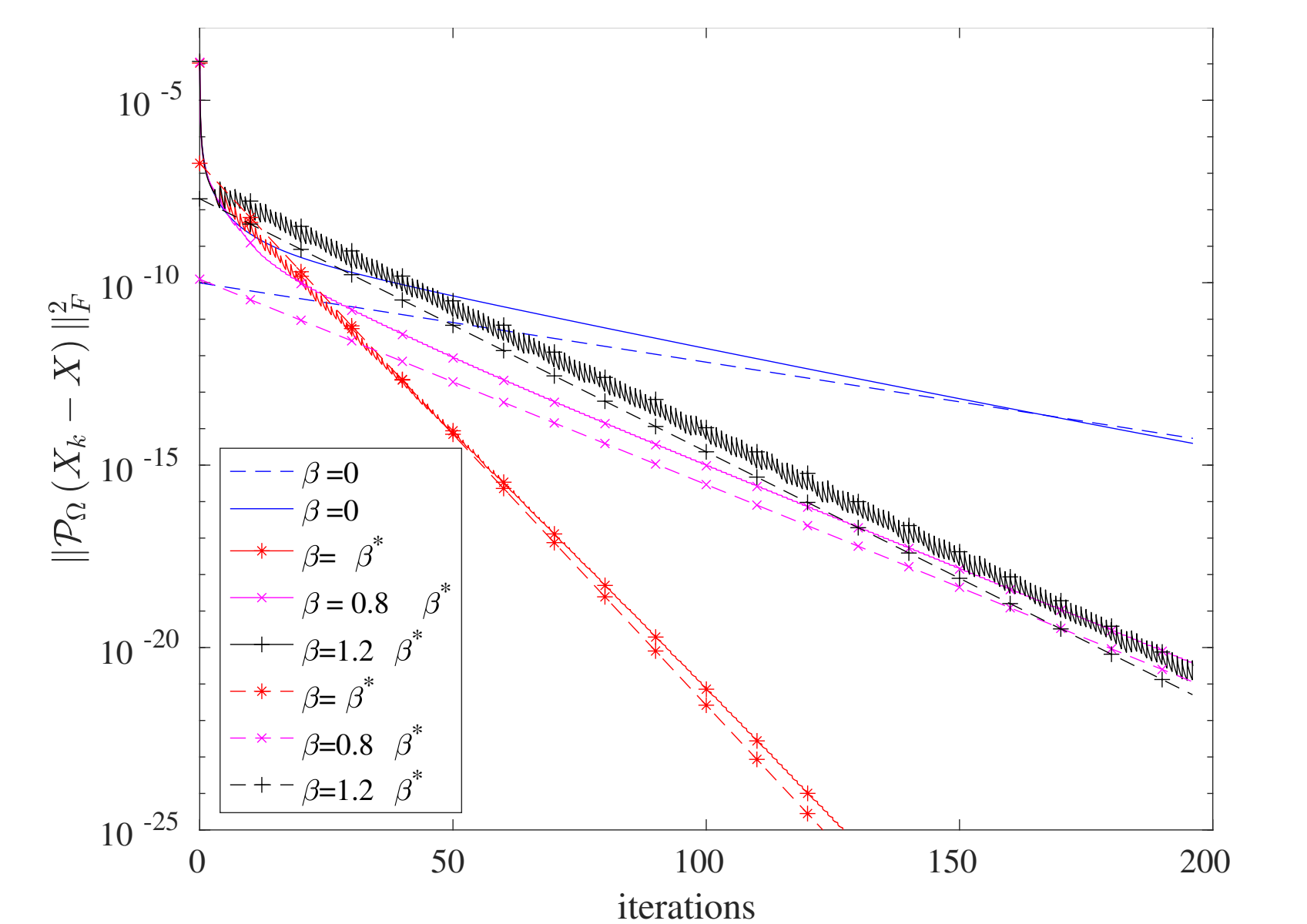
$$\beta^* = (1 - \sqrt{1 - \rho_t})^2$$

- Improved convergence results:

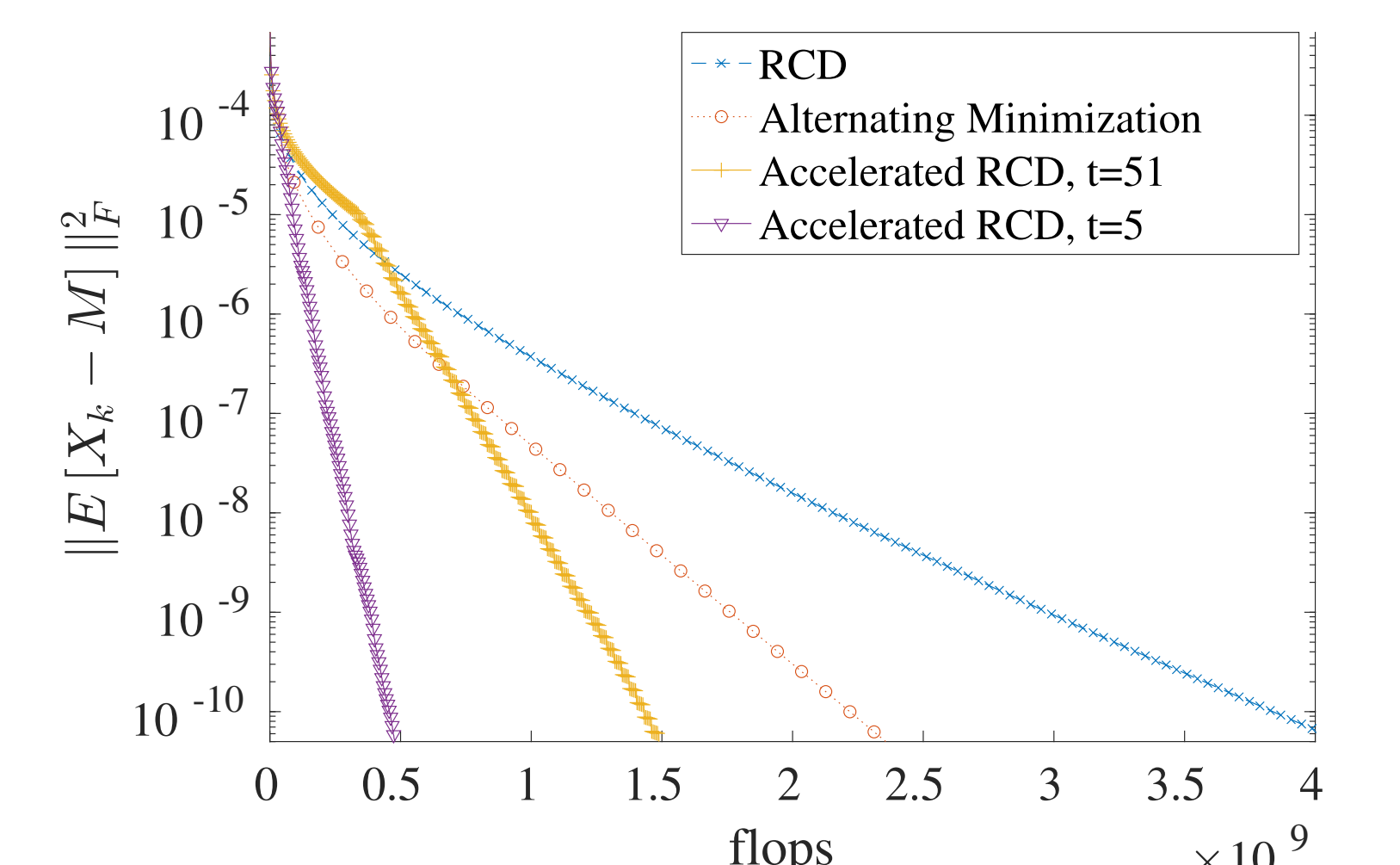
$$\|\mathbb{E}[\mathbf{A}_k \mathbf{B}_k^T - \mathbf{M}]\|_F \leq C(1 - \sqrt{1 - \rho_t})^k$$

- Per-epoch computation complexity remains  $\mathcal{O}(|\Omega|r)$

## Results



Twice objective value vs. epochs for accelerated randomized coordinate descent with various values of momentum parameter. Number of epochs per acceleration step is  $t = 51$ .



Squared Frobenius error of mean distance for an  $80 \times 80$  matrix vs. flops for different algorithms. The algorithms shown are the unaccelerated RCD (\*), optimally accelerated RCD (+ and  $\nabla$ ), and alternating minimization (o).

## Conclusion

- Polyak's momentum acceleration to RCD provides an efficient method
- Epoch-level acceleration preserves RCD complexity
- Tight analysis of unaccelerated RCD allows for optimal hyper-parameter selection