

On the **Unreasonable** Effectiveness Of Single Vector Krylov for Low-Rank Approximation

By: Raphael A. Meyer, Cameron Musco, Christopher Musco

NYU

UMass Amherst

NYU

 *Accepted at SODA 2024!*

Low-Rank Approximation

Given $A \in \mathbb{R}^{n \times n}$, K , $\varepsilon > 0$ find ortho $Q \in \mathbb{R}^{n \times K}$ with

$$\|A - QQ^T A\|_{F,2} \leq (1 + \varepsilon) \|A - A_K\|_{F,2}$$

Ideally, Q = top K eigenvectors of A , so use Krylov

[Rokhlin et al. '09], [Halko et al. '11], [Drineas Ipsen '19], [Tropp '22], ...

Block Krylov

1. Pick a start block

$$B \in \mathbb{R}^{n \times b}$$

Usually Gaussian

$b =$ **block size**

2. Build Krylov subspace

$$Z = \text{orth}(K) = \text{orth}([B \ AB \ \dots \ A^t B])$$

3. Return a solution

$$Q = Z^T U_K \quad \text{where} \quad U_K = \text{top } K \text{ eigvecs of } Z^T A A^T Z$$

But how should we pick b ?

How should we pick b ?

1. Large block size $b \geq K$

Rich line of work [Tropp, Halko, Martinson, Gu, Drineas, Ipsen, Woodruff, ...]

Strong theoretical results for L.R.A. specifically

Gap-Independent Convergence

[Musco Musco '15]

$b=K, [B]_{i,j} \sim \mathcal{N}(0,1) \implies t=O\left(\frac{1}{\sqrt{\epsilon}} \log\left(\frac{n}{\epsilon}\right)\right)$ suffices

Let $g_{K \rightarrow b} = \frac{\lambda_K - \lambda_{b+1}}{\lambda_K}$

Spectral Decay Convergence

[Musco Musco '15]

$$b \geq K, [B]_{ij} \sim \mathcal{N}(0,1) \implies t = O\left(\frac{1}{\sqrt{g_{K \rightarrow b}}} \log\left(\frac{n}{\epsilon}\right)\right) \text{ suffices}$$

Let $b = K+2, K+5, K+10$

How should we pick b ?

2. Small block size $b \ll K$

$b=1$ is called "Single Vector Krylov"

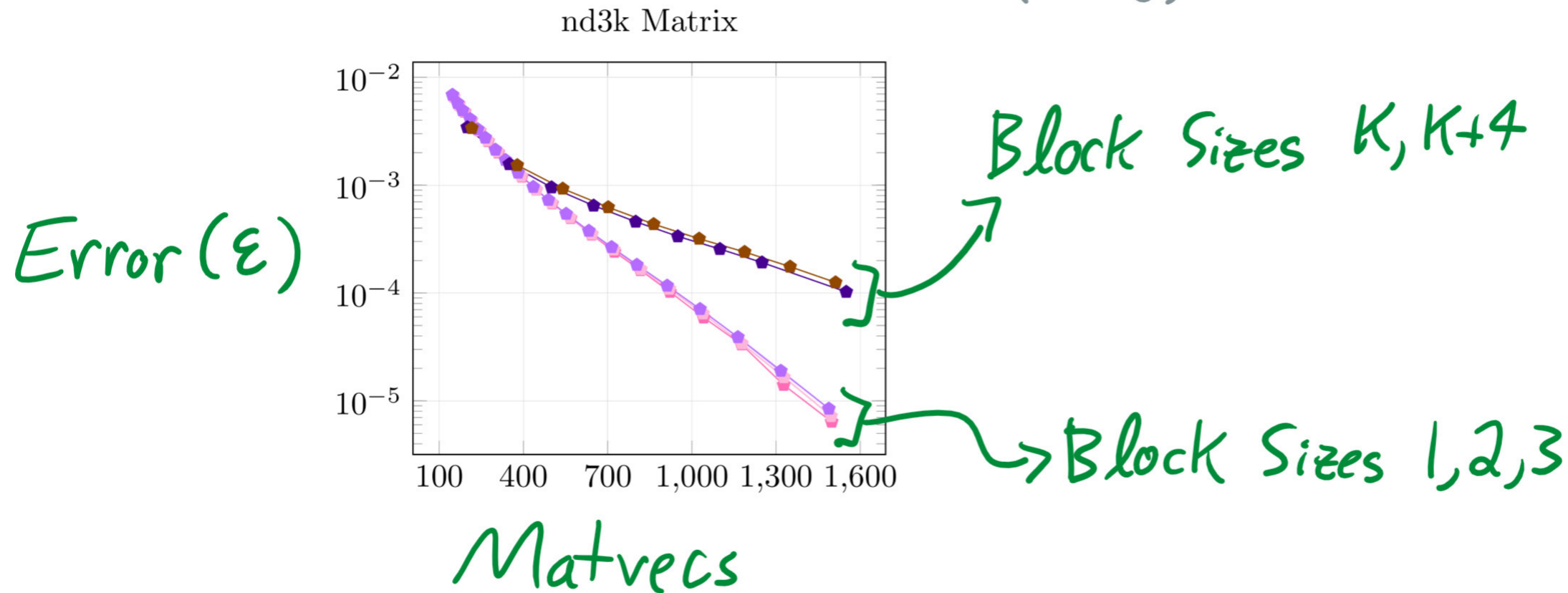
Classical NLA suggests $b \approx$ size of eigval clusters

Lack results* for Low-Rank Approximation

Cannot be gap-independent

In practice, $b=1$ often just works well

($K=50$)



A theory/practice gap!

When and why do small block methods match/outperform large block methods for low-rank approximation?

Caveats: Infinite Precision, Matvec Complexity

Main Result

For all $b \geq K$

Number of matvecs needed by Single Vector Krylov

is less than

Number of matvecs needed by block size b Krylov

If any $b \geq k$ is fast, then single vector is fast

Up to log dependence on eigengaps

Main Result (Rigorous)

Let g_{\min} = smallest gap between any of top K eigs

$$= \min_{i=1, \dots, K-1} \frac{\lambda_i - \lambda_{i+1}}{\lambda_{i+1}}$$

Then,

$$b=1 \text{ converges in } t = O\left(\frac{K}{\sqrt{\varepsilon}} \log\left(\frac{1}{g_{\min}}\right) + \frac{1}{\sqrt{\varepsilon}} \log\left(\frac{n}{\varepsilon}\right)\right)$$

$$\text{For all } l \geq K, \text{ converges in } t = O\left(\frac{l}{\sqrt{g_{K \rightarrow l}}} \log\left(\frac{1}{g_{\min}}\right) + \frac{1}{\sqrt{g_{K \rightarrow l}}} \log\left(\frac{n}{\varepsilon}\right)\right)$$

If some $b \asymp K$ gets linear convergence,

$$O\left(\frac{K}{\sqrt{g_{K \rightarrow b}}} \log\left(\frac{n}{\varepsilon}\right)\right) \quad \text{vs} \quad O\left(\frac{K}{\sqrt{g_{K \rightarrow b}}} \log\left(\frac{1}{g_{\min}}\right) + \frac{1}{\sqrt{g_{K \rightarrow b}}} \log\left(\frac{n}{\varepsilon}\right)\right)$$

Upside: separate K from ε

Downside: depends on g_{\min}

Key Observation: A Silly Manipulation

Suppose $b=1$, so for $\underline{x} \sim \mathcal{N}(0, I)$

$$\text{span}(\mathbb{R}^n) = \text{span}([\underline{x} \quad A\underline{x} \quad A^2\underline{x} \quad \dots \quad A^t\underline{x}])$$

Now, repeat some columns

$$\begin{aligned} &= \text{span}([\underbrace{\underline{x} \quad A\underline{x} \quad \dots \quad A^l\underline{x}}_{S_l} \quad \underbrace{A\underline{x} \quad A^2\underline{x} \quad \dots \quad A^{l+1}\underline{x}}_{AS_l} \quad \underbrace{A^2\underline{x} \quad A^3\underline{x} \quad \dots \quad A^{l+2}\underline{x}}_{A^2S_l} \quad \dots \quad \underbrace{A^{t-l}\underline{x} \quad \dots \quad A^t\underline{x}}_{A^{t-l}S_l}]) \\ &= \text{span}([S_l \quad AS_l \quad A^2S_l \quad \dots \quad A^{t-l}S_l]) \end{aligned}$$

Where $S_l = [\underline{x} \quad A\underline{x} \quad \dots \quad A^l\underline{x}]$ is our **Simulated Start Block**

$b=1$ Krylov Subspace
of degree t
starting from $\underline{x} \sim \mathcal{N}(0, I)$

$=$

$b=l$ Krylov Subspace
of degree $t-l$
starting from S_l

Upside: 1 matvec = 1 iteration of block krylov

Downside: S_l is a bad starting block
$$S_l = [\underline{x} \quad A\underline{x} \quad \dots \quad A^{l-1}\underline{x}]$$

Let $B \in \mathbb{R}^{n \times b}$ be an **L-good Starting Matrix**. Then,

[Musco Musco '15]

$b=K$ converges in $O\left(\frac{1}{\sqrt{\varepsilon}} \log\left(\frac{nL}{\varepsilon}\right)\right)$ iterations

$b \geq K$ converges in $O\left(\frac{1}{\sqrt{g_{K \rightarrow b}}} \log\left(\frac{nL}{\varepsilon}\right)\right)$ iterations

$[B]_{ij} \sim \mathcal{N}(0,1)$ has $L=O(nb)$

$$b=K \Rightarrow O\left(\frac{1}{\sqrt{\varepsilon}} \log\left(\frac{n}{\varepsilon}\right)\right)$$

[New Result]

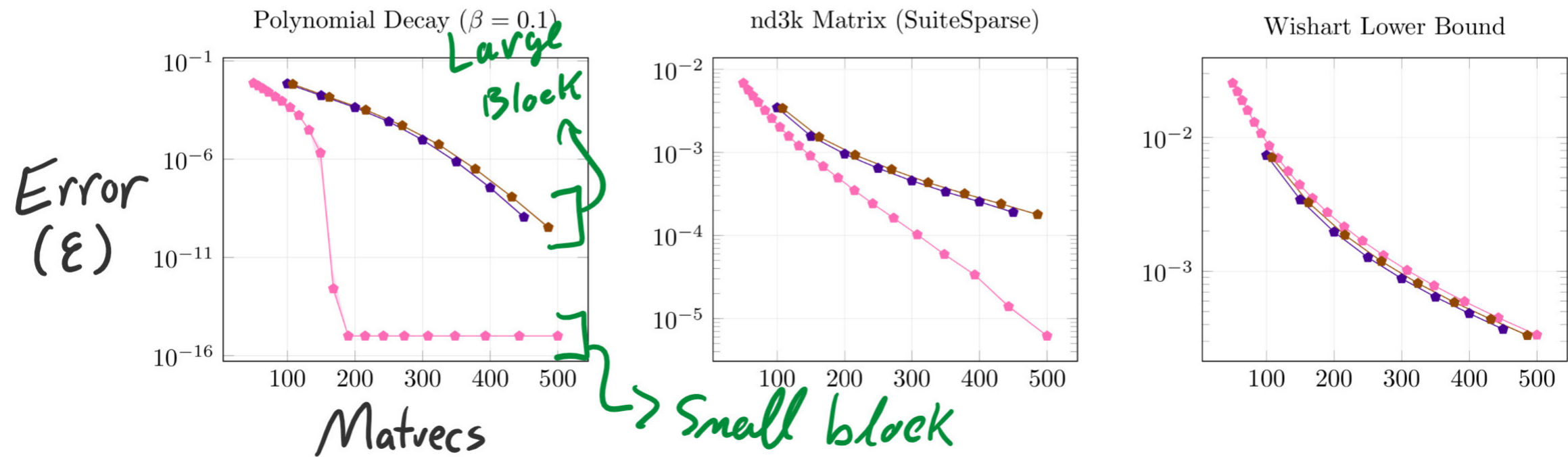
S_l has $L=O\left(\frac{nl^3}{g_{\min}^{4l}}\right)$

$$l=K \Rightarrow O\left(\frac{K}{\sqrt{\varepsilon}} \log\left(\frac{1}{g_{\min}}\right) + \frac{1}{\sqrt{\varepsilon}} \log\left(\frac{n}{\varepsilon}\right)\right)$$

$$l \geq K \Rightarrow O\left(\frac{l}{\sqrt{g_{K \rightarrow l}}} \log\left(\frac{1}{g_{\min}}\right) + \frac{1}{\sqrt{g_{K \rightarrow l}}} \log\left(\frac{n}{\varepsilon}\right)\right)$$

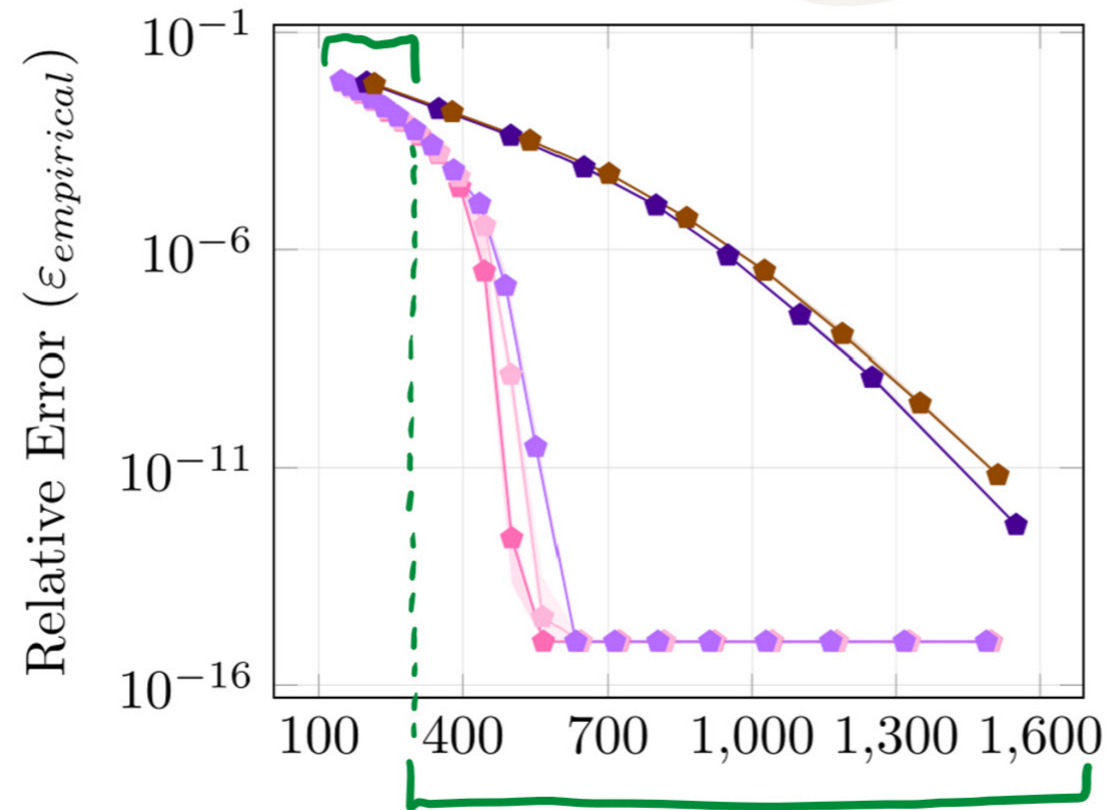
Q: *When are single vector methods faster?*

A: *When some block size achieves linear convergence*



Simulated blocks may explain slow-then-fast convergence

Can only simulate small blocks



Can now simulate large blocks

In the paper: Grab bag of more implications

- Beyond $b=1$
- Smoothed Analysis shatters g_{\min}
- Simplify Fast-Frobenius L.R.A. [Bakshi et al. '22]
- Faster-ish Schatten-norm L.R.A
- Single Vector Subspace Iteration
- Experiments

Any questions?