# Accelerating Linear Systems Solution by Exploiting Low-Rank Approximations to Factorization Error

Theo Mary, joint work with Nick Higham
University of Manchester, School of Mathematics

M\cr NA
Manchester Numerical Analysis

## Objective

- Compute solution to linear system $Ax = b$
- $A \in \mathbb{R}^{n \times n}$ is ill conditioned

## LU-based preconditioner

1. Compute approximate factorization $A = \widehat{L}\widehat{U} + \Delta A$
   - Half-precision factorization
   - Incomplete LU factorization
   - Structured matrix factorization: Block Low-Rank, $\mathcal{H}$, HSS,...
2. Solve $\Pi_{LU}Ax = \Pi_{LU}b$ with $\Pi_{LU} = \widehat{U}^{-1}\widehat{L}^{-1}$ via some iterative method

- Convergence to solution may be slow or fail
⇒ **Objective: accelerate convergence**

A New Preconditioner based on Low-Rank Error <span>Theo Mary</span>

Matrix lund_a ($n = 147$, $\kappa(A) = 2.8e+06$)



SVD of $A$

SVD of $A^{-1}$

- Often, $A$ is ill conditioned due to a small number of small singular values
- Then, $A^{-1}$ is numerically low-rank

A New Preconditioner based on Low-Rank Error Theo Mary

## Factorization error might be low-rank?

Let the error $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I = \widehat{U}^{-1}\widehat{L}^{-1}(\widehat{L}\widehat{U} + \Delta A) - I$

$$= \widehat{U}^{-1}\widehat{L}^{-1}\Delta A \approx A^{-1}\Delta A$$

Does $E$ retain the low-rank property of $A^{-1}$?

## A novel preconditioner

Consider the preconditioner

$$\Pi_{E_k} = (I + E_k)^{-1}\Pi_{LU}$$

with $E_k$ a rank-$k$ approximation to $E$.

- If $E = E_k$, $\Pi_{E_k} = A^{-1}$
- If $E \approx E_k$ for some small $k$, $\Pi_{E_k}$ can be computed cheaply

## Preprint

N. J. Higham and T. Mary, *A New Preconditioner that Exploits Low-Rank Approximations to Factorization Error*, MIMS EPrint 2018.10.

# Problem statement

## Low-rank gap

$$\varepsilon_k(A) = \min_{W_k} \left\{ \frac{\|A - W_k\|}{\|A\|} : \text{rank}(W_k) \le k \right\}$$

## Eckart-Young-Mirsky

$$\varepsilon_k(A) = \frac{\sigma_{k+1}(A)}{\sigma_1(A)}$$

## Problem statement

Quantify worst-case reduction of the low-rank gap from $A^{-1}$ to $E = \widehat{U}^{-1}\widehat{L}^{-1}\Delta A$, i.e find some $\beta$ such that

$$\varepsilon_k(E) \le \beta \varepsilon_k(A^{-1})$$

## Theorem

$$\varepsilon_k(E) \leq \beta_1 \beta_2 \varepsilon_k(A^{-1})$$

with

$$\varepsilon_k(\widehat{U}^{-1}\widehat{L}^{-1}) \leq \beta_1 \varepsilon_k(A^{-1})$$

$$\beta_1 = \left(1 + \|A^{-1}\Delta A\|\right)\left(1 + \|\widehat{U}^{-1}\widehat{L}^{-1}\Delta A\|\right)$$

$$\varepsilon_k(E) = \varepsilon_k(\widehat{U}^{-1}\widehat{L}^{-1}\Delta A) \leq \beta_2 \varepsilon_k(\widehat{U}^{-1}\widehat{L}^{-1})$$

$$\beta_2 = \frac{\|\widehat{U}^{-1}\widehat{L}^{-1}\| \|\Delta A\|}{\|\widehat{U}^{-1}\widehat{L}^{-1}\Delta A\|}$$

- $\beta_1$: maximal deviation of the sing. vals. by additive perturbation
- $\beta_2$: should be small for typical $\Delta A$

A New Preconditioner based on Low-Rank Error    Theo Mary

A New Preconditioner based on Low-Rank Error
Theo Mary

Matrix cz308

A New Preconditioner based on Low-Rank Error Theo Mary

Matrix steam1

Matrix rajat14

We did **not** specifically select matrices for which $A^{-1}$ is low-rank!

We need to build

$$\Pi_{E_k} = (I + E_k)^{-1}\Pi_{LU} = (I + E_k)^{-1}\widehat{U}^{-1}\widehat{L}^{-1}$$

where $E_k$ is a rank-$k$ approximation of $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I$

$E$ cannot be built explicitly! $\Rightarrow$ Use **randomized** method

---

**Algorithm 1** Randomized SVD via direct SVD of $V^T E$.

---

1: {Input: the error matrix $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I$, stored implicitly.}
2: Sample $E$: $S = E\Omega$, with $\Omega$ a $n \times (k + p)$ random matrix.
3: Orthonormalize $S$: $V = qr(S)$.
4: Compute SVD of $V^T E$: $X\Sigma Y^T = V^T E$.
5: Truncate $X$, $\Sigma$, $Y$ into $X_k$, $\Sigma_k$, $Y_k$.
6: The SVD of $E_k$ is given by $(VX_k)\Sigma_k Y_k^T$.

---

　　　　　　A New Preconditioner based on Low-Rank Error　　　　　　Theo Mary

# Computational cost analysis

**Algorithm 1** Randomized SVD via direct SVD of $V^T E$.

1: {Input: the error matrix $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I$, stored implicitly.}
2: Sample $E$: $S = E\Omega$, with $\Omega$ a $n \times (k + p)$ random matrix.
3: Orthonormalize $S$: $V = \text{qr}(S)$.
4: Compute SVD of $V^T E$: $X\Sigma Y^T = V^T E$.
5: Truncate $X$, $\Sigma$, $Y$ into $X_k$, $\Sigma_k$, $Y_k$.
6: The SVD of $E_k$ is given by $(VX_k)\Sigma_k Y_k^T$.

$$\ell = k + p$$

|  | setup | solve |
|---|---|---|
| $\Pi_{LU}$ | $\frac{2}{3}n^3$ | $2n^2$ |
| $\Pi_{E_k}^{(1)}$ | $\frac{2}{3}n^3 + 8n^2\ell + O(n\ell^2)$ | $2n^2 + O(nk)$ |

**Algorithm 1** Randomized SVD via direct SVD of $V^T E$.

1: {Input: the error matrix $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I$, stored implicitly.}
2: Sample $E$: $S = E\Omega$, with $\Omega$ a $n \times (k+p)$ random FFT matrix.
3: Orthonormalize $S$: $V = qr(S)$.
4: Compute SVD of $V^T E$: $X\Sigma Y^T = V^T E$.
5: Truncate $X$, $\Sigma$, $Y$ into $X_k$, $\Sigma_k$, $Y_k$.
6: The SVD of $E_k$ is given by $(VX_k)\Sigma_k Y_k^T$.

$$\ell = k + p$$

|  | setup | solve |
|---|---|---|
| $\Pi_{LU}$ | $\frac{2}{3}n^3$ | $2n^2$ |
| $\Pi_{E_k}^{(1)}$ | $\frac{2}{3}n^3 + 8n^2\ell + O(n\ell^2)$ | $2n^2 + O(nk)$ |
| $\Pi_{E_k}^{(2)}$ | $\frac{2}{3}n^3 + 6n^2\ell + 4n^2\log\ell + O(n\ell^2)$ | $2n^2 + O(nk)$ |

**Algorithm 1** Randomized SVD via row extraction.

1: {Input: the error matrix $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I$, stored implicitly.}
2: Sample $E$: $S = E\Omega$, with $\Omega$ a $n \times (k+p)$ random FFT matrix.
3: Orthonormalize $S$: $V = \mathrm{qr}(S)$.
4: Compute ID of $V$: $V = (I_k \ W)^T V_{(K,:)}$.
5: Extract $E_{(K,:)}$ and compute a QR factorization $E_{(K,:)}^T = QR$.
6: Compute SVD of $(I_k \ W)^T R^T$: $X\Sigma Y^T = (I_k \ W)^T R^T$.
7: Truncate $X$, $\Sigma$, $Y$ into $X_k$, $\Sigma_k$, $Y_k$.
8: The SVD of $E_k$ is given by $(VX_k)\Sigma_k Y_k^T$.

$$\ell = k + p$$

|  | setup | solve |
|---|---|---|
| $\Pi_{LU}$ | $\frac{2}{3}n^3$ | $2n^2$ |
| $\Pi_{E_k}^{(1)}$ | $\frac{2}{3}n^3 + 8n^2\ell + O(n\ell^2)$ | $2n^2 + O(nk)$ |
| $\Pi_{E_k}^{(2)}$ | $\frac{2}{3}n^3 + 6n^2\ell + 4n^2\log\ell + O(n\ell^2)$ | $2n^2 + O(nk)$ |
| $\Pi_{E_k}^{(3)}$ | $\frac{2}{3}n^3 + 2n^2\ell + 4n^2\log\ell + O(n\ell^2)$ | $2n^2 + O(nk)$ |

# Experimental setting

- Three types of approximate LU factorization:
  - Half-precision
  - Incomplete LU with drop tolerance $10^{-5} \leq \tau \leq 10^{-1}$
  - Block Low-Rank with low-rank threshold $10^{-9} \leq \tau \leq 10^{-1}$

- Three types of approximate LU factorization:
  - Half-precision
  - Incomplete LU with drop tolerance $10^{-5} \leq \tau \leq 10^{-1}$
  - Block Low-Rank with low-rank threshold $10^{-9} \leq \tau \leq 10^{-1}$

- Iterative solver is GMRES-based iterative refinement with three precisions
  - FP64 working precision and residual is computed in FP128
  - Max nb of GMRES iterations per IR step is 100
  - Max nb of IR steps is 10

- Three types of approximate LU factorization:
  - Half-precision
  - Incomplete LU with drop tolerance $10^{-5} \leq \tau \leq 10^{-1}$
  - Block Low-Rank with low-rank threshold $10^{-9} \leq \tau \leq 10^{-1}$

- Iterative solver is GMRES-based iterative refinement with three precisions
  - FP64 working precision and residual is computed in FP128
  - Max nb of GMRES iterations per IR step is 100
  - Max nb of IR steps is 10

- Large set of real-life matrices
  - $53 \leq n \leq 494$ and $10^3 \leq \kappa(A) \leq 10^{14}$
  - Most are sparse, but treated as dense
  - 149 tests on 40 different matrices

# Experimental setting

- Three types of approximate LU factorization:
  - Half-precision
  - Incomplete LU with drop tolerance $10^{-5} \le \tau \le 10^{-1}$
  - Block Low-Rank with low-rank threshold $10^{-9} \le \tau \le 10^{-1}$

- Iterative solver is GMRES-based iterative refinement with three precisions
  - FP64 working precision and residual is computed in FP128
  - Max nb of GMRES iterations per IR step is 100
  - Max nb of IR steps is 10

- Large set of real-life matrices
  - $53 \le n \le 494$ and $10^3 \le \kappa(A) \le 10^{14}$
  - Most are sparse, but treated as dense
  - 149 tests on 40 different matrices

- MATLAB code running on laptop
  - We measure nb of iterations and flops
  - Time is only estimated, not measured

Performance profile: $\rho$ is the percentage of problems solved for less than $\alpha \times$ the cost of the best choice $\Rightarrow$ **higher is better**

Performance profile: $\rho$ is the percentage of problems solved for less than $\alpha \times$ the cost of the best choice $\Rightarrow$ **higher is better**



We seek a compromise between number of iterations and flops to minimize time, which we estimate assuming BLAS-2 is $10\times$ slower than BLAS-3

A New Preconditioner based on Low-Rank Error
Theo Mary

Performance profile: $\rho$ is the percentage of problems solved for less than $\alpha \times$ the cost of the best choice $\Rightarrow$ **higher is better**



Need to set oversampling $p$ differently depending on preconditioner variant

Performance profile: $\rho$ is the percentage of problems solved for less than $\alpha \times$ the cost of the best choice $\Rightarrow$ **higher is better**



Similar trend for low-rank threshold $\varepsilon$

A New Preconditioner based on Low-Rank Error Theo Mary

Black-box setting: use $\Pi^{(3)}_{E_k}$ with $p = 10$ and $\varepsilon = 10^{-7}$

Black-box setting: use $\Pi^{(3)}_{E_k}$ with $p = 10$ and $\varepsilon = 10^{-7}$

Black-box setting: use $\Pi_{E_k}^{(3)}$ with $p = 10$ and $\varepsilon = 10^{-7}$

📄 P. R. Amestoy, A. Buttari, J.-Y. L'Excellent, and T. Mary, *Performance and Scalability of the Block Low-Rank Multifrontal Factorization on Multicore Architectures*.

Application to BLR-MUMPS sparse multifrontal solver
BLR threshold $= 10^{-2}$, iterate until converged to accuracy $10^{-9}$

| Matrix | $n$ | $\Pi_{LU}$ | | $\Pi_{E_k}$ | |
|--------|-----|------|------|------|------|
| | | Iter. | Time | Iter. | Time |
| audikw_1 | 1.0M | 691 | 1163 | 331 | 625 |
| Bump_2911 | 2.9M | — | — | 284 | 1708 |
| Emilia_923 | 0.9M | 174 | 304 | 136 | 267 |
| Fault_639 | 0.6M | — | — | 294 | 345 |
| Ga41As41H72 | 0.3M | — | — | 135 | 143 |
| Hook_1498 | 1.5M | 417 | 902 | 356 | 808 |
| Si87H76 | 0.2M | — | — | 131 | 116 |

**Good potential to improve low-precision, low-memory BLR solvers**

# Conclusion

## Summary

- Ill-conditioned matrices often have a numerically low-rank inverse
- Theoretical justification of why the error $E = \widehat{U}^{-1}\widehat{L}^{-1}A - I$ retains this property
- Novel preconditioner based on a low-rank approximation to the error to accelerate linear systems solution

## Future work

- High-performance implementation for FP16 and ILU
- Well suited for GPUs (FP16 $8\times$ faster than FP32!)

## Slides and paper available here

http://personalpages.manchester.ac.uk/staff/theo.mary/

# Backup slides

# Ingredient 1: $\widehat{U}^{-1}\widehat{L}^{-1}$ is low-rank if $A^{-1}$ is

## Lemma

$$\sigma_i(X + \Delta X) \leq \sigma_i(X)\left(1 + \|X^{-1}\Delta X\|\right)$$

Apply lemma twice:

Maximum growth

$$X = \widehat{L}\widehat{U} \text{ and } \Delta X = \Delta A \quad \Rightarrow \quad \sigma_i(A) \leq \sigma_i(\widehat{L}\widehat{U})\overbrace{\left(1 + \|\widehat{U}^{-1}\widehat{L}^{-1}\Delta A\|\right)}$$

$$X = A \text{ and } \Delta X = -\Delta A \quad \Rightarrow \quad \sigma_i(\widehat{L}\widehat{U}) \leq \sigma_i(A)\underbrace{\left(1 + \|A^{-1}\Delta A\|\right)}$$

Maximum shrinkage

## Theorem

$$\varepsilon_k(\widehat{U}^{-1}\widehat{L}^{-1}) \leq \beta_1 \varepsilon_k(A^{-1})$$

with

$$\beta_1 = \left(1 + \|A^{-1}\Delta A\|\right)\left(1 + \|\widehat{U}^{-1}\widehat{L}^{-1}\Delta A\|\right)$$

A New Preconditioner based on Low-Rank Error      Theo Mary

# Ingredient 2: $\widehat{U}^{-1}\widehat{L}^{-1}\Delta A$ is low-rank if $\widehat{U}^{-1}\widehat{L}^{-1}$ is

### Theorem

with

$$\varepsilon_k(\widehat{U}^{-1}\widehat{L}^{-1}\Delta A) \leq \beta_2 \varepsilon_k(\widehat{U}^{-1}\widehat{L}^{-1})$$

$$\beta_2 = \frac{\|\widehat{U}^{-1}\widehat{L}^{-1}\| \, \|\Delta A\|}{\|\widehat{U}^{-1}\widehat{L}^{-1}\Delta A\|}$$

### Corollary

$$\varepsilon_k(E) \leq \beta_1 \beta_2 \varepsilon_k(A^{-1})$$

### Theorem

$$\beta_2 \leq \bar{\beta}_2 = \frac{\sigma_{n+1-k}(\widehat{L}\widehat{U})}{\sigma_n(\widehat{L}\widehat{U})} \frac{\|\Delta A\|}{\|P_k \Delta A\|}$$

with $P_k = X_k X_k^T$ and $X_k$ the last $k$ left singular vectors of $\widehat{L}\widehat{U}$.

A New Preconditioner based on Low-Rank Error            Theo Mary

Typical $\Delta A$

Special $\Delta A$

A New Preconditioner based on Low-Rank Error

Theo Mary