# The implicit Hari–Zimmermann algorithm for the generalized SVD

Sanja Singer[1]     Vedran Novaković[2]     Saša Singer[3]

[1]University of Zagreb, Faculty of Mechanical Engineering and Naval Architecture, Croatia

[2]STFC, Daresbury Laboratory, United Kingdom

[3]University of Zagreb, Faculty of Science, Department of Mathematics, Croatia

Seminar for Numerical Analysis and Scientific Computing, November 27[th], 2015, Manchester, UK

Outline of the talk:

- the Jacobi type methods—two-sided vs. one sided,
- brief description of the original Falk–Langemeyer algorithm, and the Hari–Zimmermann (HZ) algorithm for the GEP,
- description of the HZ algorithm for the GSVD computation,
- accuracy of the pointwise HZ GSVD algorithm,
- some implementation details,
- results of numerical testing.

# Classical Jacobi method for eigensystem computation

Choose a pivot pair $(i, j)$ and then

- apply a plane rotation $R$ in $(i, j)$ plane, from both sides, to annihilate the element $a_{ij}$ of the Hermitian $A$

- since $R$ is $I$, except at the crossings of $i$th and $j$th rows and columns, where

$$\widehat{R} = \begin{bmatrix} \cos\varphi & \sin\varphi \\ -\sin\varphi & \cos\varphi \end{bmatrix},$$

- this transformation affects rows $i$ and $j$, and columns $i$ and $j$, and diagonalizes the $2 \times 2$ submatrix

$$\widehat{A} = \begin{bmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{bmatrix}, \quad \widehat{R}^* \widehat{A} \widehat{R} = \mathrm{diag}(a'_{ii}, a'_{jj}).$$

## Advantages of the one-sided over the two-sided methods

- One-sided methods: if matrix $A$ is, for example, (in the case of positive definite $A$) factored by the Cholesky factorization

$$A = F^T F,$$

- $2 \times 2$ restriction $\widehat{A}$ is obtained by three dot products

$$\widehat{A} = \begin{bmatrix} f_i^T f_i & f_i^T f_j \\ f_i^T f_j & f_j^T f_j \end{bmatrix}$$

- after computing the rotation angle, the transformation is applied on the columns $i$ and $j$ of $F$
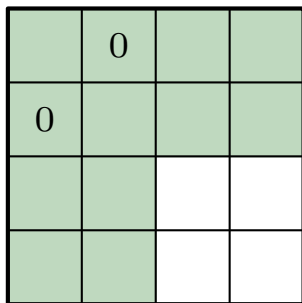
$$[f_i', f_j'] = [f_i, f_j]\widehat{R}.$$

## Advantages of the one-sided over the two-sided methods

- one-sided methods are faster than the two-sided methods — columnwise action
- and more accurate — $\kappa(A) = \kappa^2(F)$, no actual annihilation
- one-sided methods can be parallelized easily, if $i \neq j \neq k \neq \ell$, then columns $(i, j)$ and $(k, \ell)$ can be simultaneously orthogonalized
- if $A$ is positive definite, the one-sided algorithm computes the SVD of $F$ — the algorithm stops when $A_N$ is nearly diagonal, i.e., when $F_N$ has orthogonal columns.
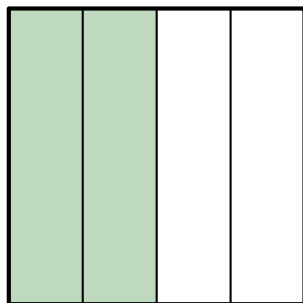
One sweep of the Jacobi-type algorithm under the row–cyclic strategy:
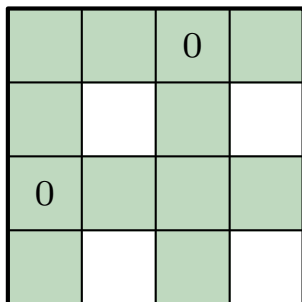
two-sided alg. on $F^*F$:



one-sided alg. on $F$:

# Jacobi type methods — two-sided vs. one-sided

One sweep of the Jacobi-type algorithm under the row–cyclic strategy:
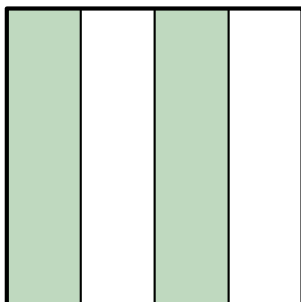
two-sided alg. on $F^*F$:

one-sided alg. on $F$:

One sweep of the Jacobi-type algorithm under the row–cyclic strategy:
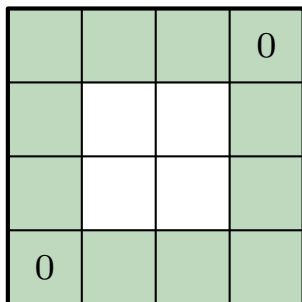
two-sided alg. on $F^*F$:



one-sided alg. on $F$:

# Jacobi type methods — two-sided vs. one-sided

One sweep of the Jacobi-type algorithm under the row–cyclic strategy:
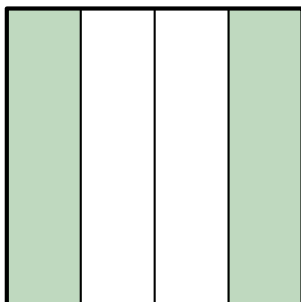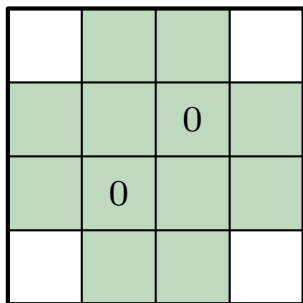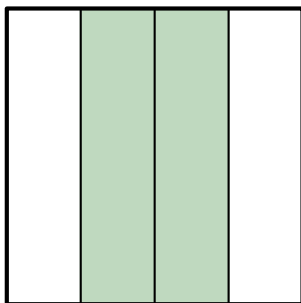
two-sided alg. on $F^*F$:

one-sided alg. on $F$:

One sweep of the Jacobi-type algorithm under the row–cyclic strategy:

two-sided alg. on $F^*F$:                    one-sided alg. on $F$:

One sweep of the Jacobi-type algorithm under the row–cyclic strategy:
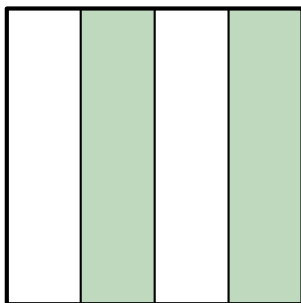
two-sided alg. on $F^*F$:                      one-sided alg. on $F$:

The main differences vs. the methods for the EP

- everything is doubled!
- we have two symmetric matrices $A$ and $B$ (at least one of them is positive definite) that should be simultaneously diagonalized
- in each step, two submatrices $\widehat{A}$ and $\widehat{B}$ are simultaneously diagonalized — therefore, we need two independent parameters in the "rotation" matrix
- if either $A = F^T F$ or $B = G^T G$ is positive definite, the one-sided application of the transformations on $F$ and $G$ leads to the GSVD of the pair $(F, G)$.

The Falk–Langemeyer method

- invented in 1960, paper published in two parts, in the collection Elektronische Datenverarbeitung,
- quadratic convergence of the cyclic method is proved in M.Sc. thesis of Slapničar (1989, supervised by Hari),
- the method solves the Generalized Eigenvalue Problem (GEP) for a symmetric and definite matrix pair $(A, B)$,
- it constructs a sequence of congruent pairs,

$$A^{(\ell+1)} = C_\ell^T A^{(\ell)} C_\ell, \quad B^{(\ell+1)} = C_\ell^T B^{(\ell)} C_\ell,$$

where $(A^{(1)}, B^{(1)}) := (A, B)$.

# The Falk–Langemeyer method for the GEP

The transformation matrix $C_\ell$

- ▶ resembles a scaled plane rotation: it is the identity matrix, except for its $(i,j)$-restriction $\widehat{C}_\ell$, where

$$\widehat{C}_\ell = \begin{bmatrix} 1 & \alpha_\ell \\ -\beta_\ell & 1 \end{bmatrix}.$$

- ▶ $\alpha_\ell$ and $\beta_\ell$ are determined so that the transformations diagonalize the pivot submatrices

$$\widehat{A}^{(\ell)} = \begin{bmatrix} a_{ii}^{(\ell)} & a_{ij}^{(\ell)} \\ a_{ij}^{(\ell)} & a_{jj}^{(\ell)} \end{bmatrix}, \qquad \widehat{B}^{(\ell)} = \begin{bmatrix} b_{ii}^{(\ell)} & b_{ij}^{(\ell)} \\ b_{ij}^{(\ell)} & b_{jj}^{(\ell)} \end{bmatrix},$$

- ▶ pivot indices $(i,j)$ are selected according to some pivot strategy.

# The Hari–Zimmermann method for the GEP

The Hari–Zimmermann method

- ▶ Zimmermann in her Ph.D. thesis (1969) briefly sketched a method for the GEP when $B$ is positive definite,
- ▶ Hari in his Ph.D. thesis (1984) filled in the missing details, proved global and quadratic convergence (cyclic strategies)
- ▶ before the iterative part, the pair is scaled so that the diagonal elements of $B$ are all equal to one,

$$A^{(1)} := DAD, \quad B^{(1)} := DBD,$$
$$D = \operatorname{diag}\left((b_{11})^{-1/2}, (b_{22})^{-1/2}, \ldots, (b_{kk})^{-1/2}\right),$$

- ▶ the method constructs a sequence of congruent pairs

$$A^{(\ell+1)} = Z_\ell^T A^{(\ell)} Z_\ell, \quad B^{(\ell+1)} = Z_\ell^T B^{(\ell)} Z_\ell.$$

## The transformation matrix $Z_\ell$

- resembles an ordinary plane rotation: it is the identity matrix, except for its $(i, j)$-restriction $\widehat{Z}_\ell$, where

$$\widehat{Z}_\ell = \frac{1}{\sqrt{1 - \left(b_{ij}^{(\ell)}\right)^2}} \begin{bmatrix} \cos\varphi_\ell & \sin\varphi_\ell \\ -\sin\psi_\ell & \cos\psi_\ell \end{bmatrix},$$

- $\varphi_\ell$ and $\psi_\ell$ are determined so that the transformations diagonalize the pivot submatrices $\widehat{A}^{(\ell)}$ and $\widehat{B}^{(\ell)}$
- the transformation keeps the diagonal elements of $B$ intact
- if $B = I$ then $Z_\ell$ is the ordinary rotation, the method is the ordinary Jacobi method for a single matrix.

## Computation of the elements of $\widehat{Z}_\ell$

- for simplicity, the transformation index $\ell$ is omitted $\boxed{\text{▸ Skip}}$

$$\tan(2\vartheta) = \frac{2a_{ij} - (a_{ii} + a_{jj})b_{ij}}{(a_{jj} - a_{ii})\sqrt{1 - (b_{ij})^2}}, \qquad -\frac{\pi}{4} < \vartheta \le \frac{\pi}{4}$$

$$\xi = \frac{b_{ij}}{\sqrt{1 + b_{ij}} + \sqrt{1 - b_{ij}}}$$

$$\eta = \frac{b_{ij}}{\left(1 + \sqrt{1 + b_{ij}}\right)\left(1 + \sqrt{1 - b_{ij}}\right)}$$

$$\cos\varphi = \cos\vartheta + \xi(\sin\vartheta - \eta\cos\vartheta)$$
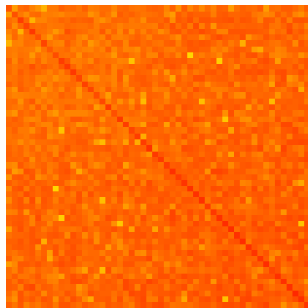
$$\cos\psi = \cos\vartheta - \xi(\sin\vartheta + \eta\cos\vartheta)$$

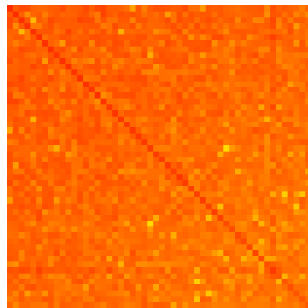$$\sin\varphi = \sin\vartheta - \xi(\cos\vartheta + \eta\sin\vartheta)$$

$$\sin\psi = \sin\vartheta + \xi(\cos\vartheta - \eta\sin\vartheta)$$

An example — $A$ and $B$ positive definite of order 52



$A$                                                    $B$

the starting pair

An example — $A$ and $B$ positive definite of order 52



$A$          $B$

end of sweep 1
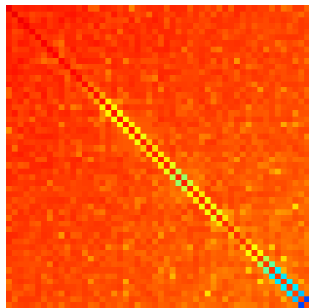
# The pointwise algorithm for the GEP

An example — $A$ and $B$ positive definite of order 52



$A$                    $B$

end of sweep 2

An example — $A$ and $B$ positive definite of order 52



$A$                                    $B$

end of sweep 3

An example — $A$ and $B$ positive definite of order 52



$A$                                        $B$

end of sweep 4

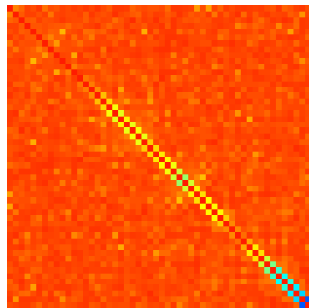An example — $A$ and $B$ positive definite of order 52
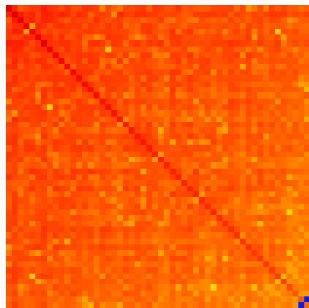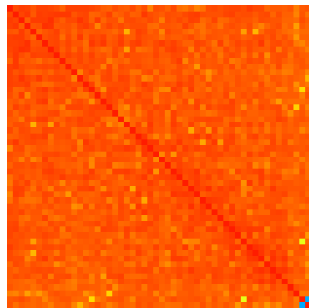


$A$                  $B$

end of sweep 5

# The pointwise algorithm for the GEP

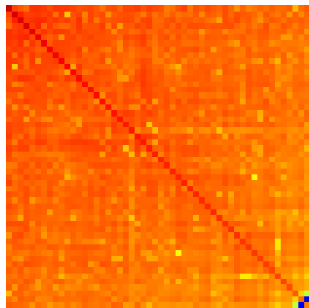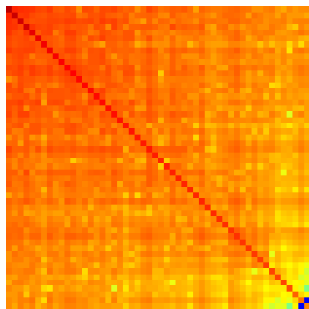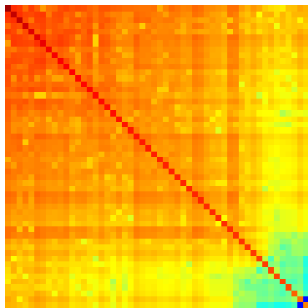An example — $A$ and $B$ positive definite of order 52



$A$

$B$

end of sweep 6

An example — $A$ and $B$ positive definite of order 52



$A$                    $B$

end of sweep 7

An example — $A$ and $B$ positive definite of order 52



$A$        $B$

end of sweep 8

An example — $A$ and $B$ positive definite of order 52



$A$             $B$

end of sweep 9

# The generalized SVD

## Definition

- For given matrices $F \in \mathbb{C}^{m \times n}$ and $G \in \mathbb{C}^{p \times n}$, where

$$K = \begin{bmatrix} F \\ G \end{bmatrix}, \quad k = \text{rank}(K),$$

  there exist unitary matrices $U \in \mathbb{C}^{m \times m}$, $V \in \mathbb{C}^{p \times p}$, and a matrix $X \in \mathbb{C}^{k \times n}$, such that

$$F = U\Sigma_F X, \quad G = V\Sigma_G X, \quad \Sigma_F \in \mathbb{R}^{m \times k}, \quad \Sigma_G \in \mathbb{R}^{p \times k}.$$

- $\Sigma_F$ and $\Sigma_G$ are real, "diagonal", and nonnegative.
- Furthermore, $\Sigma_F$ and $\Sigma_G$ satisfy

$$\Sigma_F^T \Sigma_F + \Sigma_G^T \Sigma_G = I.$$

- The ratios $(\Sigma_F)_{ii}/(\Sigma_G)_{ii}$ are called the generalized singular values of the pair $(F, G)$.

# The GEP and the GSVD

Connection between the GEP and the GSVD

- Given matrices: $F_0 \in \mathbb{R}^{m \times n}$ and $G_0 \in \mathbb{R}^{p \times n}$.
- If $G_0$ is not of full column rank, then use, for example, LAPACK preprocessing to obtain square matrices $(F, G)$, with $G$ of full rank $k$.
- For such $F$ and $G$, since $G^T G$ is a positive definite matrix, the pair $(F^T F, G^T G)$ in the corresponding GEP is symmetric and definite.
- There exist many nonsingular matrices $Z$ that simultaneously diagonalize $(F^T F, G^T G)$ by congruences,

$$Z^T F^T F Z = \Lambda_F, \quad Z^T G^T G Z = \Lambda_G,$$

where $\Lambda_F$ and $\Lambda_G$ are diagonal, $(\Lambda_F)_{ii} \geq 0$ and $(\Lambda_G)_{ii} > 0$, for $i = 1, \ldots, k$.

# The GEP and the GSVD

Connection between the GEP and the GSVD

- Since $\Lambda_F$ and $\Lambda_G$ are diagonal, the columns of $FZ$ and $GZ$ are orthogonal (not orthonormal),

$$FZ = U\Lambda_F^{1/2}, \quad GZ = V\Lambda_G^{1/2},$$

  where $U$ and $V$ are orthogonal matrices.

- If $\Lambda_F + \Lambda_G \neq I$, then the matrices in the GSVD are

$$X := SZ^{-1}, \qquad \Sigma_F := \Lambda_F^{1/2}S^{-1}, \qquad \Sigma_G := \Lambda_G^{1/2}S^{-1}.$$

  where $S = (\Lambda_F + \Lambda_G)^{1/2}$ is the diagonal scaling.

- If only the generalized singular values are needed, rescaling is not necessary, and $\sigma_i = (\Lambda_G^{-1/2}\Lambda_F^{1/2})_{ii}$, for $i = 1, \ldots, k$.

# The pointwise algorithm for the GSVD

The implicit HZ algorithm for the GSVD

$Z = I;$ $\qquad it = 0$

`repeat` // sweep loop

$\quad it = it + 1$

$\quad$ `for` all pairs $(i, j)$, $1 \leq i < j \leq k$

$\qquad$ compute

$$\widehat{A} = \begin{bmatrix} f_i^T f_i & f_i^T f_j \\ f_i^T f_j & f_j^T f_j \end{bmatrix}; \qquad \widehat{B} = \begin{bmatrix} g_i^T g_i & g_i^T g_j \\ g_i^T g_j & g_j^T g_j \end{bmatrix}$$

$\qquad$ compute the elements of $\widehat{Z}$

$\qquad\quad$ // transform $F$, $G$ and $Z$

$\qquad [f_i, f_j] = [f_i, f_j] \cdot \widehat{Z}$

$\qquad [g_i, g_j] = [g_i, g_j] \cdot \widehat{Z}$

$\qquad [z_i, z_j] = [z_i, z_j] \cdot \widehat{Z}$

`until` (no transf. in this sweep) or $(it \geq maxcyc))$

# Accuracy of the implicit HZ algorithm

Standard assumptions on $f\ell$ arithmetic

- $f\ell(x \circ y) = (1 + \varepsilon_\circ)(x \circ y), \quad |\varepsilon_\circ| \leq \varepsilon, \quad \circ = +, -, *, /,$
- $f\ell(\sqrt{x}) = (1 + \varepsilon_{\sqrt{}})\sqrt{x}, \quad |\varepsilon_{\sqrt{}}| \leq \varepsilon,$
- $f\ell(x + (y \cdot z)) = (1 + \varepsilon_{\mathrm{fma}})(x + (y \cdot z)), \quad |\varepsilon_{\mathrm{fma}}| \leq \varepsilon.$

Observations

- transformation $\widehat{Z}_\ell$ is determined by $\sin\varphi$, $\sin\psi$, and $b_{ij}$ (transformation indices omitted)
- both cosines are positive, and uniquely determined

$$\cos\varphi = \sqrt{1 - \sin^2\varphi}, \quad \cos\psi = \sqrt{1 - \sin^2\psi}.$$

# Accuracy of the implicit HZ algorithm

## Analysis

- Let $W$ be a certain HZ transformation in step $\ell$
- its submatrix of order $2$ in the $(i, j)$ plane is

$$\widehat{W} = \begin{bmatrix} \widehat{w}_{11} & \widehat{w}_{12} \\ \widehat{w}_{21} & \widehat{w}_{22} \end{bmatrix} = \frac{1}{\sqrt{1 - b^2}} \begin{bmatrix} \cos \tilde{\varphi} & \sin \tilde{\varphi} \\ -\sin \tilde{\psi} & \cos \tilde{\psi} \end{bmatrix}.$$

- $\widehat{W}$ is used to transform the pivot columns $i$ and $j$, to obtain the transformed matrices $FW$ and $GW$
- in $f\ell$ arithmetic, each computation involves rounding errors, therefore $W' = f\ell(W)$ is the actually computed transformation matrix
- $F' = f\ell(FW')$ and $G' = f\ell(GW')$ are the computed matrices after the transformation.

Forward bounds

- The computed matrices $F'$ and $G'$ can be written as

$$F' = FW + \delta F', \quad G' = GW + \delta G',$$

  $\delta F'$ and $\delta G'$ are the forward perturbations

- only the columns $i$ and $j$ are changed

$$[f_i', f_j'] = [f_i, f_j] \cdot \widehat{W} + [\delta f_i', \delta f_j'].$$

- Normwise bounds for the columns of $F$ are

$$\|\delta f_i'\|_2 \leq \frac{\varepsilon}{\sqrt{1-b^2}} \big(5\cos\tilde{\varphi} \cdot \|f_i\|_2 + 4.5\,|\sin\tilde{\psi}| \cdot \|f_j\|_2\big),$$

$$\|\delta f_j'\|_2 \leq \frac{\varepsilon}{\sqrt{1-b^2}} \big(4.5\,|\sin\tilde{\varphi}| \cdot \|f_i\|_2 + 5\cos\tilde{\psi} \cdot \|f_j\|_2\big).$$

- The same holds for $G$, with $\|g_i\|_2 = \|g_j\|_2 = 1$.

# Accuracy of the implicit HZ algorithm

Backward bounds

- The computed matrices $F'$ and $G'$ can be viewed as

$$F' = (F + \delta F)W, \quad G' = (G + \delta G)W,$$

where $\delta F$ and $\delta G$ denote the backward perturbations

- only the columns $i$ and $j$ are changed

$$[f_i', f_j'] = \big([f_i, f_j] + [\delta f_i, \delta f_j]\big)\widehat{W}.$$

- Normwise bounds for the columns of $F$ are

$$\|\delta f_i\|_2 \le \varepsilon c_{ij}\big(5\|f_i\|_2 + 4.25\|f_j\|_2\big),$$
$$\|\delta f_j\|_2 \le \varepsilon c_{ij}\big(4.25\|f_i\|_2 + 5\|f_j\|_2\big),$$

where $c_{ij} = 1/|\cos(\tilde{\varphi} - \tilde{\psi})|$.

- The same holds for $G$, with $\|g_i\|_2 = \|g_j\|_2 = 1$.

The main result

- Assumption: each pivot pair $(i, j)$ is ordered such that $\|f_i\|_2 \geq \|f_j\|_2$
- if $F$ is of full column rank, it holds

$$\|f_j\|_2 = r_{ij}\|f_i\|_2, \quad 0 < r_{ij} \leq 1.$$

- Let $r_s := \min r_{ij}$ and $c_s := \max c_{ij}$ over all pairs of pivot indices $(i, j)$ at this stage of the algorithm,
- and let

$$\epsilon := \varepsilon c_s \left( \frac{4.25}{r_s} + 5 \right).$$

## Theorem (Drmač)

Let $F$ and $G$ be of full column rank, and let the columns of perturbation matrices satisfy

$$\|\delta f_p\|_2 \leq \epsilon \|f_p\|_2, \qquad \|\delta g_p\|_2 \leq \epsilon \|g_p\|_2$$

for some constant $\epsilon$, such that $0 \leq \epsilon < 1$. Then, the relative errors in the perturbed generalized singular values $\tilde{\sigma}_p$ of the pair $(F + \delta F, G + \delta G)$ are bounded by

$$\frac{|\tilde{\sigma}_p - \sigma_p|}{\sigma_p} \leq \left(1 + \frac{\sigma_{\min}(G_S)}{\sigma_{\min}(F_S)}\right) \frac{\epsilon \sqrt{q}}{\sigma_{\min}(G_S) - \epsilon \sqrt{q}},$$

where $F_S = F \operatorname{diag}\left(\|f_p\|_2^{-1}\right)$, $G_S = G \operatorname{diag}\left(\|g_p\|_2^{-1}\right)$, and $q$ is the maximal number of nonzero elements in any row of $\delta F$ and $\delta G$.

Sequential algorithms

- blocking – each block has $k_i \approx k/nb$ columns

$$F = [F_1, F_2, \ldots, F_{nb}], \quad G = [G_1, G_2, \ldots, G_{nb}].$$

- each pivot block can either be fully orthogonalized – full-block algorithm, or,

- each pair of columns in each block is orthogonalized once in a sweep – block oriented algorithm

- pivoting – transformations are applied in such way that after each transformation it holds

$$\frac{\|f_i'\|_2}{\|g_i'\|_2} \geq \frac{\|f_j'\|_2}{\|g_j'\|_2}, \quad i < j.$$

# Numerical testing of the sequential algorithms

▶ **Implementation**: Fortran routines with MKL.

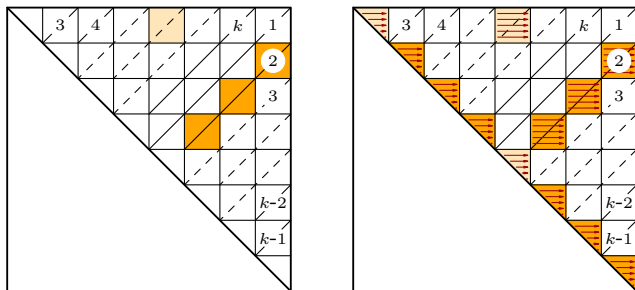| $k$ | with threaded MKL (12 cores) | | | |
|------|------------|--------------|----------|----------|
| | DTGSJA | pointwise HZ | HZ-FB-32 | HZ-BO-32 |
| 500 | 16.16 | 3.17 | 4.36 | 2.03 |
| 1000 | 128.56 | 26.89 | 18.50 | 7.65 |
| 1500 | 466.11 | 105.31 | 42.38 | 19.31 |
| 2000 | 1092.39 | 273.48 | 86.01 | 41.60 |
| 2500 | 2186.39 | 547.84 | 139.53 | 73.07 |
| 3000 | 3726.76 | 1652.14 | 203.00 | 109.46 |
| 3500 | 6062.03 | 2480.14 | 294.58 | 186.40 |
| 4000 | 8976.99 | 3568.00 | 411.71 | 239.89 |
| 4500 | 12805.27 | 4910.09 | 553.67 | 343.58 |
| 5000 | 20110.39 | 6599.68 | 711.86 | 426.76 |

Times (in seconds).

## Parallel pivoting strategy

- Choose pivot blocks independently in each step, for example, by using (block)-modulus strategy (not optimal!)



- stopping criterion
  - skip a transformation if cosines are 1
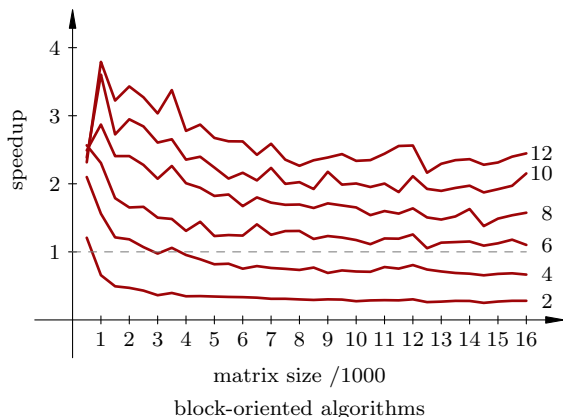  - final stop — all transformations are skipped.

# Shared memory algorithms

- Implementation: OpenMP in Fortran routines.

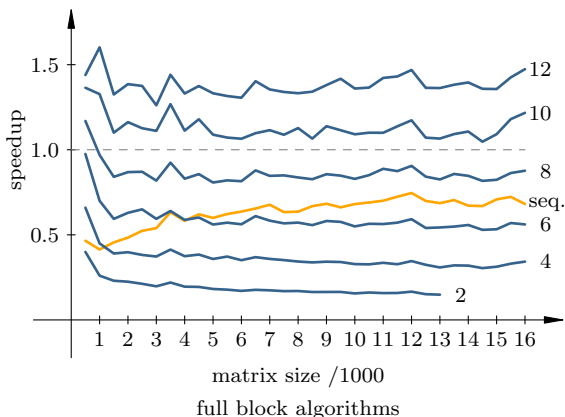| | with sequential MKL | |
| $k$ | P-HZ-FB-32 | P-HZ-BO-32 |
|---|---|---|
| 500 | 1.41 | 0.88 |
| 1000 | 4.78 | 2.02 |
| 1500 | 14.57 | 5.99 |
| 2000 | 30.02 | 12.13 |
| 2500 | 53.13 | 22.34 |
| 3000 | 86.78 | 36.08 |
| 3500 | 129.37 | 55.20 |
| 4000 | 180.32 | 86.36 |
| 4500 | 249.92 | 119.74 |
| 5000 | 320.39 | 159.59 |

Times (in seconds).

# Shared memory algorithms



Speedup of the shared memory block-oriented algorithms on 2–12 cores vs. the sequential block-oriented Hari–Zimmermann algorithm (threaded MKL on 12 cores).
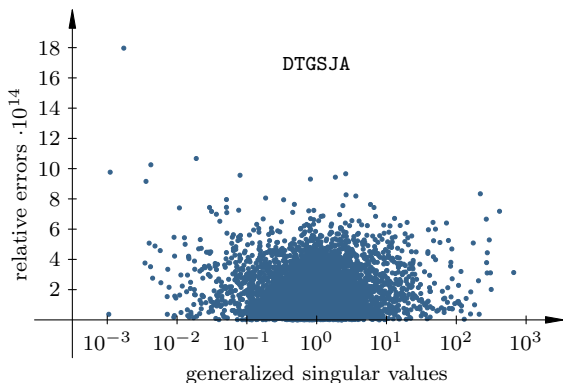
# Shared memory algorithms



Speedup of the shared memory full block algorithms on 2–12 cores vs. the sequential block-oriented Hari–Zimmermann algorithm (threaded MKL on 12 cores).

# Accuracy (matrix of order 5000)

Test matrix condition number $\max \sigma_i / \min \sigma_i \approx 6.32 \cdot 10^5$

Test matrix condition number $\max \sigma_i / \min \sigma_i \approx 6.32 \cdot 10^5$
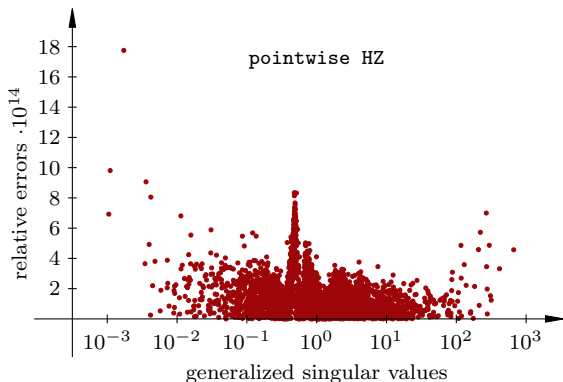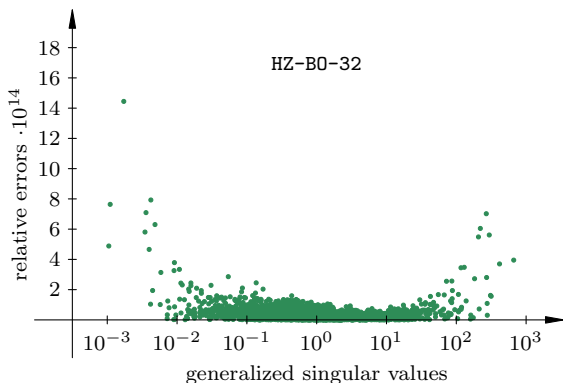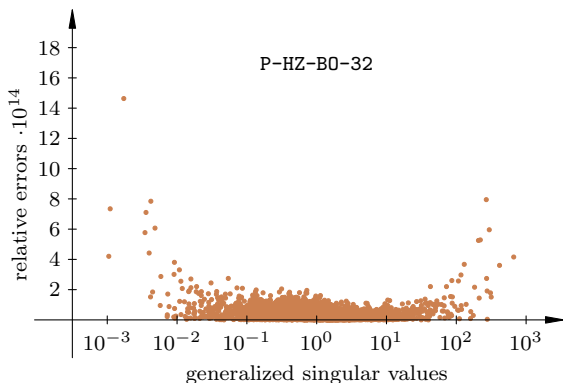
# Accuracy (matrix of order 5000)

Test matrix condition number $\max \sigma_i / \min \sigma_i \approx 6.32 \cdot 10^5$
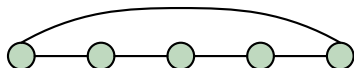
# Accuracy (matrix of order 5000)

Test matrix condition number $\max \sigma_i / \min \sigma_i \approx 6.32 \cdot 10^5$

# Distributed memory algorithms

Distributed algorithms = another level of hierarchy added

- shared-memory algorithm — a building block for the distributed memory algorithm (hybrid MPI/OpenMP)
- only conceptual difference between the distributed memory and the shared memory HZ algorithm — exchange updated block-columns among the MPI processes
- Cartesian topology — one dimensional torus of processes.



- each MPI process in each step sends only one block-column and receives only one block column.

## Distributed vs. shared memory algorithms

| number of | | time |
| MPI processes | cores | MPI-HZ-BO-32 |
| --- | --- | --- |
| 2 | 24 | 15323.72 |
| 4 | 48 | 8229.32 |
| 6 | 72 | 6049.77 |
| 8 | 96 | 4276.65 |
| 10 | 120 | 3448.90 |
| 12 | 144 | 3003.39 |
| 14 | 168 | 2565.29 |
| 16 | 192 | 2231.71 |

The running times of the hybrid MPI/OpenMP version HZ,
matrix pair of order 16000.

# Distributed vs. shared memory algorithms

| number of cores | time | |
|:---:|:---:|:---:|
| | P-HZ-FB-32 | P-HZ-BO-32 |
| 2 | – | 42906.93 |
| 4 | 35168.73 | 18096.72 |
| 6 | 21473.00 | 10936.10 |
| 8 | 13745.17 | 7651.86 |
| 10 | 9901.96 | 5599.25 |
| 12 | 8177.90 | 4925.56 |

The running times for the full block and block-oriented shared memory algorithms for the same matrix.

# Conclusion

On a particular hardware (with threaded MKL on 12 cores)

- Pointwise HZ method is 3 times faster than `DTGSJA` on matrices of order 5000.
- Sequential block-oriented `HZ-BO-32` algorithm is 15 times faster than the pointwise algorithm, i.e., more than 47 times faster than `DTGSJA`.
- For the fastest, explicitly parallel, shared memory algorithm `P-HZ-BO-32`, the speedup factor is 126!
- `DTGSJA` is unable to handle large matrices in any reasonable time.
- Triangularization is mandatory for `DTGSJA`, but not necessary for the Hari–Zimmermann method, when $G$ is of full column rank.