

# The implicit Hari–Zimmermann algorithm for the generalized SVD

Sanja Singer<sup>1</sup>   Vedran Novaković<sup>2</sup>   Saša Singer<sup>3</sup>

<sup>1</sup>University of Zagreb, Faculty of Mechanical Engineering and Naval Architecture, Croatia

<sup>2</sup>STFC, Daresbury Laboratory, United Kingdom

<sup>3</sup>University of Zagreb, Faculty of Science, Department of Mathematics, Croatia

SIAM Conference on Applied Linear Algebra (LA15),  
October 26–30, 2015, Atlanta, Georgia, USA

This work has been fully supported by Croatian Science Foundation  
under the project [IP-2014-09-3670](#).

## Outline of the talk:

- ▶ brief description of the original Falk–Langemeyer algorithm, and the **H**ari–**Z**immermann (**HZ**) algorithm for the GEP,
- ▶ description how to use the HZ algorithm for the GSVD computation,
- ▶ accuracy of the pointwise HZ GSVD algorithm,
- ▶ some implementation details,
- ▶ results of numerical testing.

# The Falk–Langemeyer method for the GEP

## The Falk–Langemeyer method

- ▶ invented in 1960, paper published in two parts in 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 **G**eneralized **E**igenvalue **P**roblem (**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  $(A^{(1)}, B^{(1)}) := (A, B)$ ,

- ▶ pairs are selected according to some pivot order.

# 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

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

# Hari–Zimmermann method for the GEP

## The Hari–Zimmermann method

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

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

- ▶ the method constructs a sequence of congruent pairs,

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

# Hari–Zimmermann method for the GEP

The transformation matrix  $Z_\ell$

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

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

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

# Hari–Zimmermann method for the GEP

Computation of the elements of  $\tilde{Z}_\ell$

- ▶ for simplicity, index of the transformation  $\ell$  is omitted

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

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

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

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

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

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

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

# 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  $(F, G)$ .



# Hari–Zimmermann method for 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  $i = 1, \dots, k$ .

# Hari–Zimmermann method for 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),

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

$U$  and  $V$  are orthogonal matrices.

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

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

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  $i = 1, \dots, k$ .

# Pointwise algorithm for the GSVD

Implicit HZ algorithm for the GSVD

$Z = I; \quad it = 0$

repeat // sweep loop

$it = it + 1$

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

compute

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

compute the elements of  $\hat{Z}$

// transform  $F, G$  and  $Z$

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

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

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

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

# Accuracy of the implicit HZ algorithm

## Standard assumptions on $fl$ arithmetic

- ▶  $fl(x \circ y) = (1 + \varepsilon_\circ)(x \circ y)$ ,  $|\varepsilon_\circ| \leq \varepsilon$ ,  $\circ = +, -, *, /$ ,
- ▶  $fl(\sqrt{x}) = (1 + \varepsilon_{\sqrt{\cdot}})\sqrt{x}$ ,  $|\varepsilon_{\sqrt{\cdot}}| \leq \varepsilon$ ,
- ▶  $fl(x + (y \cdot z)) = (1 + \varepsilon_{\text{fma}})(x + (y \cdot z))$ ,  $|\varepsilon_{\text{fma}}| \leq \varepsilon$ .

## Assumptions

- ▶ 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  $fl$  arithmetic, each computation involves **rounding errors**, therefore  $W' = fl(W)$  is the actually computed transformation matrix
- ▶  $F' = fl(FW')$  and  $G' = fl(GW')$  are the computed matrices after the transformation.

# Accuracy of the implicit HZ algorithm

## 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

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

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

- ▶ 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

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

- ▶ Normwise bounds for the columns of  $F'$  are

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

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

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

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

# Accuracy of the implicit HZ algorithm

## The main result

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

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

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

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



# Accuracy of the implicit HZ algorithm

## Theorem (Drmač)

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

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

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

# How to make the algorithm fast and accurate

## Sequential algorithms

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

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

- ▶ each pivot block can either be **fully orthogonalized** – **full-block algorithm**, or,
- ▶ in each pair of columns in each block are **orthogonalized once** – **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

---

with threaded MKL (12 cores)

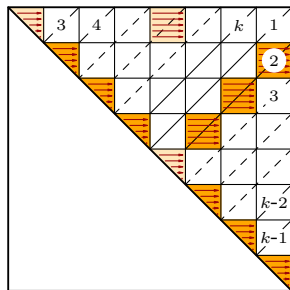
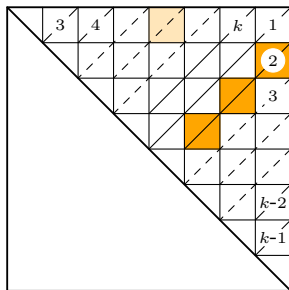
$k$	DTGSJA	pointwise HZ	HZ-FB-32	HZ-B0-32
500	16.16	3.17	4.36	2.03
1000	128.56	26.89	18.50	7.65
1500	466.11	105.31	42.38	19.31
2000	1092.39	273.48	86.01	41.60
2500	2186.39	547.84	139.53	73.07
3000	3726.76	1652.14	203.00	109.46
3500	6062.03	2480.14	294.58	186.40
4000	8976.99	3568.00	411.71	239.89
4500	12805.27	4910.09	553.67	343.58
5000	20110.39	6599.68	711.86	426.76

---

# How to make the algorithm fast and accurate

## Parallel algorithms

- ▶ Choose pivot blocks independently in each step, for example by using **(block)-modulus strategy**
- ▶ **shared-memory algorithm** – a building block for distributed memory algorithm



# How to make the algorithm fast and accurate

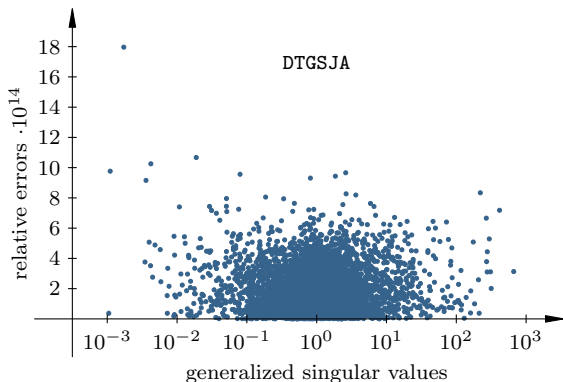
---

	with sequential MKL	
$k$	P-HZ-FB-32	P-HZ-B0-32
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

---

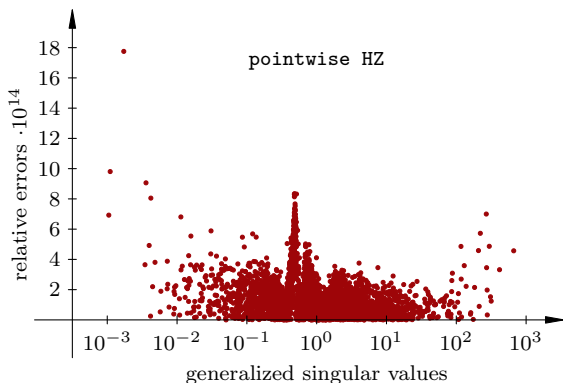
# Accuracy for matrix of order 5000

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



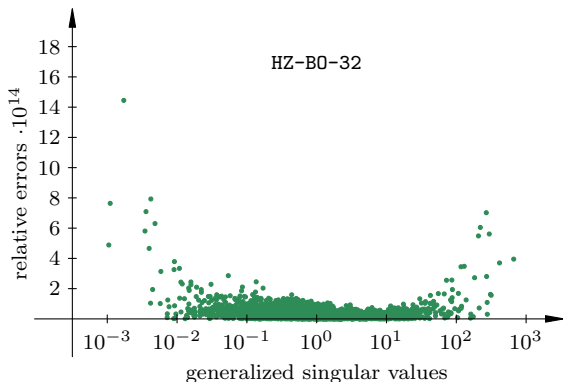
# Accuracy for matrix of order 5000

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



# Accuracy for matrix of order 5000

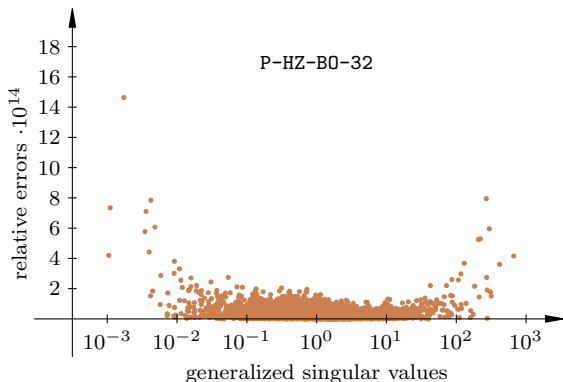
Test matrix condition number  $\max \sigma_i / \min \sigma_i \approx 6.32 \cdot 10^5$





# Accuracy for matrix of order 5000

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



# 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-B0-32 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-B0-32, 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.