See discussions, stats, and author profiles for this publication at: https://www.researchgate.net/publication/362791003

# A block minimum residual norm subspace solver with partial convergence management for sequences of linear systems -IB-BGCRO-DR

Presentation · July 2022



reads 49



A block minimum residual norm subspace solver with partial convergence management for sequences of linear systems — IB-BGCRO-DR

DD27, Minisymposium MS7\_1 July 28, 2022, Prague, Czechia

Yanfei Xiang (☞: yanfei.xiang@inria.fr) Cerfacs, Inria Bordeaux Sud-Ouest, France

# Summary

Overview
 Scientific background
 IB-BGCRO-DR
 Experimental results
 Conclusion





### Overview



3 -IB-BGCRO-DR- Yanfei Xiang

### Overview

Minimum residual norm block "Krylov" solver: first two years of my  $\mathsf{PhD}$ 

• Joint work with

Luc Giraud (Inria Bordeaux - Sud-Ouest, my PhD supervisor) Yan-Fei Jing (UESTC)

- Outcome
  - SIMAX (SIAM J. Matrix Anal. Appl.) paper, 2022, pp.30 https://doi.org/10.1137/21M1401127
  - Research report in HAL-Inria, 2022, pp.60 https://hal.inria.fr/hal-03146213v3
  - C++ implementation in Fabulous https://gitlab.inria.fr/solverstack/fabulous/





## Scientific background



5 -IB-BGCRO-DR- Yanfei Xiang

• Generate Arnoldi relation (Arnoldi algorithm):

 $AV_j = V_{j+1}\underline{H}_j$ , span $(V_{j+1}) = \text{span}(v_1, v_2, \cdots, v_{j+1})$  (1) • Generate nested and **orthonormal basis**:

 $V_{j+1}^{H}V_{j+1} = I_{j+1}$ 

 $\mathcal{K}_{j+1}(A, v_1) = \operatorname{span}(r_0, Ar_0, \cdots, A^j r_0) = \operatorname{span}(V_{j+1}) \text{ with } v_1 = \frac{r_0}{\|r_0\|_2},$ 

and  $\underline{H}_i$  is upper Hessenberg matrix



• Generate Arnoldi relation (Arnoldi algorithm):

 $AV_j = V_{j+1}\underline{H}_j$ , span $(V_{j+1}) = \text{span}(v_1, v_2, \cdots, v_{j+1})$  (1) • Generate nested and **orthonormal basis**:

 $V_{j+1}^H V_{j+1} = I_{j+1}$ 

• Solve a small dimensional **least-squares problem** (Minimal Residual approach):

 $x_{j} = \underset{x \in x_{0} + \mathcal{K}_{j}(A, v_{1})}{\operatorname{argmin}} \|b - Ax\|_{2} \text{ with } x_{j} = x_{0} + V_{j}y_{j}, \text{ thus}$  $\underbrace{\|b - Ax_{j}\|_{2}}_{x_{j} \in \mathbb{C}^{n}} = \left\|V_{j+1}(V_{j+1}^{H}r_{0} - \underline{H}_{j}y_{j})\right\|_{2} = \underbrace{\|\beta e_{1} - \underline{H}_{j}y_{j}\|_{2}}_{y_{j} \in \mathbb{C}^{j}},$ 

where  $j \ll n$  and  $y_j = \underset{y \in \mathbb{C}^j}{\operatorname{argmin}} \left\| \beta e_1 - \underline{H}_j y \right\|_2$  with  $\beta = \|r_0\|_2$  (2)



• Generate **block-Arnoldi relation** (Arnoldi algorithm):

 $A\mathscr{V}_{j} = \mathscr{V}_{j+1} \underline{\mathscr{H}}_{j}, \quad \operatorname{span}(\mathscr{V}_{j+1}) = \operatorname{span}(V_{1}, V_{2}, \cdots, V_{j+1}) \quad (3)$ • Generate nested and **orthonormal basis**:

 $\begin{aligned} \mathscr{V}_{j+1}^{H} \mathscr{V}_{j+1} &= I_{n_{j+1}} \text{ with } n_{j+1} = (j+1) \times p, \\ \mathcal{K}_{j+1}(A, V_1) &= \operatorname{span}(R_0, AR_0, \cdots, A^j R_0) = \operatorname{span}(\mathscr{V}_{j+1}), \end{aligned}$ 

where *p* is the number of right-hand sides, and  $R_0 = V_1 \Lambda_1$ ;  $\mathcal{H}_i$  is upper Hessenberg matrix with  $p \times p$  block elements



• Generate **block-Arnoldi relation** (Arnoldi algorithm):

 $A\mathscr{V}_{j} = \mathscr{V}_{j+1}\underline{\mathscr{H}}_{j}, \quad \operatorname{span}(\mathscr{V}_{j+1}) = \operatorname{span}(V_{1}, V_{2}, \cdots, V_{j+1})$ (3)

• Generate nested and orthonormal basis:

$$\mathscr{V}_{j+1}^{H}\mathscr{V}_{j+1}=\mathit{I}_{n_{j+1}}$$
 with  $n_{j+1}=(j+1) imes p_{j}$ 

• Solve a small dimensional **least-squares problem** (Minimal Residual approach):

 $X_j = \operatorname*{argmin}_{X \in X_0 + \mathcal{K}_j(A, V_1)} \|B - AX\|_F \text{ with } X_j = X_0 + \mathscr{V}_j Y_j, \text{ thus}$ 

$$\underbrace{\|B - AX_j\|_F}_{X_j \in \mathbb{C}^{n \times p}} = \left\| \mathscr{V}_{j+1}(\mathscr{V}_{j+1}^H R_0 - \underline{\mathscr{H}}_j Y_j) \right\|_F = \underbrace{\|\Lambda_j - \underline{\mathscr{H}}_j Y_j\|_F}_{Y_j \in \mathbb{C}^{n_j \times p}},$$

where  $n_j \ll n$  and

$$Y_{j} = \operatorname{argmin}_{Y \in \mathbb{C}^{n_{j} \times p}} \left\| \Lambda_{j} - \underline{\mathscr{H}}_{j} Y \right\|_{F} \text{ with } \Lambda_{j} = [\Lambda_{1}^{T}, \mathbf{0}_{p \times n_{j}}]^{T} (4)$$



Two issues in **BGMRES** 

• Inexact Breakdown (IB) caused by partial convergence or linear combination

$$A\mathscr{V}_{j} = \left[\mathscr{V}_{j}, [P_{j-1}, \widetilde{W}_{j}]\right]\mathscr{F}_{j},$$
(5)

then  $\mathscr{V}_{j+1} = [\mathscr{V}_j, \mathbb{V}_{j+1}]$  where  $\mathbb{V}_{j+1} \in \mathbb{C}^{n \times p_{j+1}} (p_{j+1} \le p)$ 

- > Eq. (5) comes from block-Arnoldi relation in BGMRES:  $A\mathscr{V}_{j} = [\mathscr{V}_{j}, V_{j+1}] \underline{\mathscr{H}}_{j}$  with  $V_{j+1} \in \mathbb{C}^{n \times p}$
- >  $\mathbb{V}_{j+1}$  corresponds to the components within  $[P_{j-1}, W_j] \in \mathbb{C}^{n \times p}$  that contribute the most to the residual norms
- > *P<sub>j</sub>* are the abandoned directions at iteration *j* that contribute the less to the residuals.
- $> \mathscr{F}_j$  is no longer upper Hessenberg as soon as one IB occurs

IB-BGMRES [Robbé and Sadkane, 2006]



Two issues in **BGMRES** 

• Inexact Breakdown (IB) caused by partial convergence or linear combination

$$A\mathscr{V}_{j} = \left[\mathscr{V}_{j}, [P_{j-1}, \widetilde{W}_{j}]\right]\mathscr{F}_{j},$$
(5)

then  $\mathscr{V}_{j+1} = [\mathscr{V}_j, \mathbb{V}_{j+1}]$  where  $\mathbb{V}_{j+1} \in \mathbb{C}^{n \times p_{j+1}} (p_{j+1} \le p)$ IB-BGMRES [Robbé and Sadkane, 2006]

 The growth of memory and computational requirements deflated restarting/subspace recycling strategy for accelerating convergence – BGMRES-DR [Morgan, 2005], IB-BGMRES-DR [Agullo et al., 2014], BGCRO-DR [Parks et al., 2016]



- BGMRES-DR: Restarted block-GMRES with deflation of eigenvalues. [R. B. Morgan. Applied Numerical Mathematics, 2005, 54:(2), pp. 222-236]
- IB-BGMRES: Exact and inexact breakdowns in the block GMRES method. [M. Robbé and M. Sadkane. *Linear Algebra* and its Applications, 2006, 419:(1), pp. 265-285]
- IB-BGMRES-DR: Block GMRES Method with Inexact Breakdowns and Deflated Restarting. [E. Agullo, L. Giraud and Y.-F. Jing. *SIAM Journal on Matrix Analysis and Applications*, 2014, 35:(4), pp. 1625-1651]
- BGCRO-DR: A block Recycled GMRES method with investigations into aspects of solver performance. [M. L. Parks, K. M. Soodhalter and D. B. Szyld. *arXiv*, 2016]





### **IB-BGCRO-DR**



10 - IB-BGCRO-DR- Yanfei Xiang

Motivation: devise **block solver** for **sequences of linear systems** with multiple left and right hand sides:  $A^{(\ell)}X^{(\ell)} = B^{(\ell)}, \ell = 1, 2, ...$  families

- 1. For block method: Block GMRES method with inexact breakdown and deflated restarting IB-BGMRES-DR
  - Combine the <u>inexact breakdown</u> mechanism aiming for handling rank deficiency in block iterations and the <u>deflated restarting</u> augmenting with some approximated eigenvectors for accelerating convergence of the next cycles.



Motivation: devise **block solver** for **sequences of linear systems** with multiple left and right hand sides:  $A^{(\ell)}X^{(\ell)} = B^{(\ell)}, \ell = 1, 2, ...$  families

1. For block method: Block GMRES method with inexact breakdown and deflated restarting – IB-BGMRES-DR

2. For sequences of linear systems: A block recycled GMRES method- BGCRO-DR

> Subspace recycling: accelerate convergence rate of the next cycles as well as the <u>next families</u> to be solved.



Motivation: devise **block solver** for **sequences of linear systems** with multiple left and right hand sides:  $A^{(\ell)}X^{(\ell)} = B^{(\ell)}, \ell = 1, 2, ...$  families

1. For block method: Block GMRES method with inexact breakdown and deflated restarting – IB-BGMRES-DR

How to combine <u>inexact breakdown</u> mechanism and <u>subspace recycling</u> strategy for solving sequences of linear systems with multiple left and right-hand sides ?

 For sequences of linear systems: A block recycled GMRES method– BGCRO-DR



• Generate block-Arnoldi like relation  $A[U_k, \mathscr{V}_j] = \left[C_k, \mathscr{V}_j, [P_{j-1}, \widetilde{W}_j]\right] \mathscr{F}_j = \widehat{\mathscr{V}_{j+1}} \mathscr{F}_j, \quad (6)$ > IB:  $A\mathscr{V}_j = * + \left[\mathscr{V}_j, [P_{j-1}, \widetilde{W}_j]\right] \mathscr{F}_j$  from IB-BGMRES > BGCRO-DR:  $AU_k = C_k, \ C_k^H C_k = I_k \text{ (computed by Harmonic-Ritz, Rayleigh-Ritz or other techniques); and <math>\mathscr{F}_j = \left[\begin{matrix}I_k & \mathscr{B}_j\\0_{(n_j+p)\times k} & \mathscr{F}_j\end{matrix}\right] \in \mathbb{C}^{(k+n_j+p)\times(k+n_j)}, \ \mathcal{B}_j = C_k^H A\mathscr{V}_j \in \mathbb{C}^{k\times n_j}$ 



i

- Generate **block-Arnoldi like relation**  $A[U_k, \mathscr{V}_j] = \left[C_k, \mathscr{V}_j, [P_{j-1}, \widetilde{W}_j]\right] \underline{\mathscr{F}}_j = \widehat{\mathscr{V}}_{j+1} \underline{\mathscr{F}}_j, \quad (6)$
- Generate orthonormal basis:

$$\widehat{\mathscr{V}}_{j+1}^{H}\widehat{\mathscr{V}}_{j+1} = I_{(k+n_j+p)}$$
 with  $n_j = \sum_{i=1}^{j} p_i$ 



- Generate block-Arnoldi like relation  $A[U_k, \mathscr{V}_j] = \left[ C_k, \mathscr{V}_j, [P_{j-1}, \widetilde{W}_j] \right] \underline{\mathscr{F}}_j = \widehat{\mathscr{V}}_{j+1} \underline{\mathscr{F}}_j, \quad (6)$
- Generate orthonormal basis:  $\widehat{\mathscr{V}}_{j+1}^{H} \widehat{\mathscr{V}}_{j+1} = I_{(k+n_j+p)}$  with  $n_j = \sum_{i=1}^{j} p_i$
- Solve a small dimensional least-squares problem:

 $X_j = \operatorname*{argmin}_{X \in X_0 + \mathsf{Range}([U_k, \mathscr{V}_j])} \|B - AX\|_F \text{ with } X_j = X_0 + [U_k, \mathscr{V}_j]Y_j,$ 

where 
$$X_j \in \mathbb{C}^{n \times p}$$
,  $Y_j \in \mathbb{C}^{(k+n_j) \times p}$  with  $k + n_j \ll n$  and  
 $Y_j = \underset{Y \in \mathbb{C}^{(k+n_j) \times p}}{\operatorname{argmin}} \left\| \Lambda_j - \underline{\mathscr{F}}_j Y \right\|_F$ ,  $\Lambda_j \in \mathbb{C}^{(k+n_j+p) \times p}$  (7)



### Flowchart of reusing spectral information





Different search space expansion policies can be defined

- Residual based (IB-BGMRES, Robbé-Sadkane, 2006): all the residual norms will be smaller than prescribed threshold
- Backward error on *b* based (IB-BGMRES-DR, Agullo et al., 2014): all the backward error  $\eta_b$  will be smaller than prescribed threshold

$$\begin{split} \eta_b(x_j) &= \min_{\Delta b} \left\{ \tau > 0 : \|\Delta b\|_2 \le \tau \|b\|_2 \text{ and } Ax_j = b + \Delta b \right\} \\ &= \frac{\|Ax_j - b\|_2}{\|b\|_2} \end{split}$$

• Backward error on A and b based (IB-BGCRO-DR): all the backward error  $\eta_{A,b}$  will be smaller than prescribed threshold

$$\begin{split} \eta_{A,b}(x_j) &= \min_{\Delta A, \Delta b} \left\{ \tau > 0 : \|\Delta A\|_2 \le \tau \|A\|_2 , \, \|\Delta b\|_2 \le \tau \|b\|_2 \\ &\text{and} \, (A + \Delta A) x_j = b + \Delta b \right\} \\ &= \frac{\|Ax_j - b\|_2}{\|A\|_2 \|x_j\|_2 + \|b\|_2} \end{split}$$





### Experimental results



15 - IB-BGCRO-DR- Yanfei Xiang

### Benefits of IB and DR between the families



• The solution of the second family benefits from recycled information from the previous family solution



### Benefits of IB and DR between the families



- IB-BGMRES-DR history overlaps BGCRO-DR until the first IB
- IB-BGMRES-DR history is the same for the two families



16 -IB-BGCRO-DR- Yanfei Xiang

### Benefits of IB and DR between the families



- IB-BGCRO-DR history overlaps IB-BGMRES-DR for the first family
- IB-BGCRO-DR inherits all the good genes of IB, DR and GCRO



16 -IB-BGCRO-DR- Yanfei Xiang

### The backward error stopping criterion: $\eta_{A,b}$



• IB-BGCRO-DR is able to decrease  $\eta_{A,b}$  to a very low value close to the machine epsilon:  $\mathcal{O}(10^{-16})$ 



\*-VA : individual convergence threshold



Once converged:

- the ones with largest thresholds do not progress anymore
- computation effort focuses on the ones with smallest thresholds



### \*-VA : individual convergence threshold



• Significant gains compared to converging all solutions with the most stringent threshold



### Test system with slowly varying left-hand sides



• 
$$A^{(\ell)}X^{(\ell)} = B^{(\ell)}$$
 with  $\ell = 1, 2, 3$ 

• Obvious gains compared to another related block iterative solver



19 - IB-BGCRO-DR- Yanfei Xiang





20 -IB-BGCRO-DR- Yanfei Xiang

**IB-BGCRO-DR**: **subspace recycling** block GCRO-DR method with **partial convergence** management

• **Subspace recycling** strategies : Harmonic-Ritz, Rayleigh-Ritz or other techniques can be used to define the space to be recycled



**IB-BGCRO-DR**: **subspace recycling** block GCRO-DR method with **partial convergence** management

- Subspace recycling strategies
- Search space expansion policies governed by
  - > the partial convergence management (IB-, \*-VA and \*-CB)
  - > the backward error stopping criterion ( $\eta_b$  and  $\eta_{A,b}$ )



**IB-BGCRO-DR**: **subspace recycling** block GCRO-DR method with **partial convergence** management

- Subspace recycling strategies
- Search space expansion policies governed by
  - > the partial convergence management (IB-, \*-VA and \*-CB)
  - > the backward error stopping criterion ( $\eta_b$  and  $\eta_{A,b}$ )
- Flexible preconditioning (mixed-precision): IB-BFGCRO-DR



**IB-BGCRO-DR**: **subspace recycling** block GCRO-DR method with **partial convergence** management

- Subspace recycling strategies
- Search space expansion policies governed by
  - > the partial convergence management (IB-, \*-VA and \*-CB)
  - > the backward error stopping criterion ( $\eta_b$  and  $\eta_{A,b}$ )
- Flexible preconditioning (mixed-precision): IB-BFGCRO-DR

• Remarks on some computational and algorithmic aspects

- > Partial convergence detecting in the initial residuals  $\rightarrow$  address possible initial rank deficiency & reduce computation
- > The MGS re-orthogonalization among k + p columns (between the columns of recycling space and initial block Arnoldi basis) at restart → generate a good enough orthonormal basis



**IB-BGCRO-DR**: **subspace recycling** block GCRO-DR method with **partial convergence** management

- Subspace recycling strategies
- Search space expansion policies governed by
  - > the partial convergence management (IB-, \*-VA and \*-CB)
  - > the backward error stopping criterion ( $\eta_b$  and  $\eta_{A,b}$ )
- Flexible preconditioning (mixed-precision): IB-BFGCRO-DR

• Remarks on some computational and algorithmic aspects

- > Partial convergence detecting in the initial residuals  $\rightarrow$  address possible initial rank deficiency & reduce computation
- > The MGS re-orthogonalization among k + p columns (between the columns of recycling space and initial block Arnoldi basis) at restart → generate a good enough orthonormal basis

Thanks for your attention. Questions ?



### IB detecting in the initial basis or after ?



 Execute IB detecting in the initial basis helps reduce (even though slightly) the consuming mvps



### IB detecting in the initial basis or after ?



• Such slight gains comes from the benefits of addressing the initial rank deficiency accumulated at restart



#### Diagonal scaled IB vs min-IB

 $\|(B - AX_j)D_{\varepsilon}\|_2 \leq 1 \text{ vs } \|b^{(i)} - Ax_j^{(i)}\|_2 < \|B - AX_j\|_2 \leq \varepsilon \min_{i=1,...,\rho} \|b^{(i)}\|_2$ 



•  $D_{\varepsilon} = \varepsilon^{-1} \operatorname{diag}(\|(b^{(1)})^{-1}\|_2, \dots, \|(b^{(p)})^{-1}\|_2); B = \operatorname{rand}(n, p)$ 



23 - IB-BGCRO-DR- Yanfei Xiang

#### Diagonal scaled IB vs min-IB

 $\|(B - AX_j)D_{\varepsilon}\|_2 \leq 1 \text{ vs } \|b^{(i)} - Ax_j^{(i)}\|_2 < \|B - AX_j\|_2 \leq \varepsilon \min_{i=1,...,p} \|b^{(i)}\|_2$ 



• B = rand(n, p) and then multiply 20 to the first  $\frac{p}{2}$  columns of B

### Loss of orthogonality of the Krylov basis





• Without re-orthogonalization of  $C_k$  to  $\mathbb{V}_1$  (or  $V_1$ ) at restart



#### Loss of orthogonality of the Krylov basis

 $\text{Loss-Orth} = \left\| \widehat{\mathscr{V}}_{j+1}^{H} \widehat{\mathscr{V}}_{j+1} - I_{j+1} \right\|_{2}$ 



With re-orthogonalization of C<sub>k</sub> to V<sub>1</sub> (or V<sub>1</sub>) at restart (by V<sub>1</sub> = (I − C<sub>k</sub> \* C<sup>H</sup><sub>k</sub>)V<sub>1</sub> implemented in column-column way)



#### Loss of orthogonality of the Krylov basis

 $\text{Loss-Orth} = \left\| \widehat{\mathscr{V}}_{j+1}^{H} \widehat{\mathscr{V}}_{j+1} - I_{j+1} \right\|_{2}$ 



 Apply re-orthogonalization to all the columns of [C<sub>k</sub>, [𝒱<sub>1</sub>, P<sub>0</sub>]] by the modified Gram-Schmidt (MGS) process at restart

