A Practical Randomized GMRES Algorithm for Solving Linear Equation System In Circuit Simulation

Baiyu Chen, Jiawen Cheng, Wenjian Yu* Department of Computer Science and Technology Tsinghua University Presenter: BAIYU CHEN

Contents

Background

- Proposed Method
- Experimental Results
- Conclusion

Contents

Background

- Proposed Method
- Experimental Results
- Conclusion

Circuit Simulation

• With the advance of chip technology, circuits with billions or even more

nodes need to be simulated efficiently and effectively.

Simulation for large-scale integrated circuits is of significance.



Circuit Simulation

• A typical framework:



(3) Solve the constructed **linear equations**. (the key

process)

Circuit Simulation

- The two key properties of linear equations from circuit simulation problem:
- Sparse: while the nodes of circuits are enormous, their connections are very sparse.
 Therefore, the coefficient matrices are sparse.
- (2) Unsymmetric: the dissymmetry can be caused by many factors. The increasing

dissymmetry can be observed if the integrated chips grow larger.

• We should leverage the above two properties for more efficient solution.

Solution for Linear Equations

• The two typical types of methods for solving linear equations:

(1) Direct methods: decompose the matrix and solve by substitution.

Two main weaknesses: because of the fill-ins in the decomposition process and we cannot control the accuracy of the solution.

(2) Iterative methods: transform the solution process into many-times matrix-vector multiplication (keep the sparsity), and return the solution when the demanded accuracy is met.



• The two most commonly used iterative methods for solving linear equations:

(1) conjugate gradient (**CG**) method: O(n) time complexity per iteration, but only applicable for symmetric positive definite (SPD) matrix.

Only 235 cases in the 2893 cases in SuiteSparse Matrix Collection are SPD.

0 case after 2015 in <u>SuiteSparse Matrix Collection</u> is SPD.

(1) generalized minimal residual (GMRES) method: applicable for general linear equations.

GMRES

- The incredible popularity: Since the birth of GMRES, it has been cited more than 15000 times (the most cited numerical algorithm in our best knowledge).
- The main shortcoming: increasing computational cost with iterations,

caused by its core: the Arnoldi process.

Algorithm 1 The Arnoldi Process with Modified Gram-Schmidt

Input: $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$, $x_0 \in \mathbb{R}^n$, m. **Output:** $V_{m+1} \in \mathbb{R}^{n \times (m+1)}, \overline{H}_m \in \mathbb{R}^{(m+1) \times m}$. 1: Calculate residual $r \leftarrow b - Ax_0$ 2: $v_1 \leftarrow \frac{r}{\|r\|}$ 3: **for** $j = 1, 2, \dots, m$ **do** $w_i \leftarrow Av_i$ 4: for $i = 1, \cdots, j$ do ▶ modified Gram-Schmidt (MGS) 5: $h_{i,j} \leftarrow w_j^T v_i$ 6: $w_i \leftarrow w_i - h_{i,j} v_i$ 7: end for 8: $h_{i+1,i} \leftarrow \|w_i\|$ 9: $v_{j+1} \leftarrow \frac{w_j}{h_{j+1}}$ 10: 11: end for 12: **Return** $V_{m+1} = [v_1, v_2, \cdots, v_{m+1}], \overline{H}_m = (h_{i,j})_{(m+1) \times m}$.

GMRES

 The efficient and effective of orthogonalization is of significance: the orthogonalization is the rooted reason for the high complexity of GMRES and its effectiveness directly influence its convergency.

The efficient and effective relative error
 estimation is of significance: we should know

the error to know when to return the solution.

Algorithm 2 The Restarted GMRES Algorithm

Input: $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$, $x_0 \in \mathbb{R}^n$, m_{max} , tol. **Output:** $x_{appro} \in \mathbb{R}^n$. 1: $v \leftarrow +\infty$ 2: while |y| > tol do $r \leftarrow b - Ax_0, \quad g \leftarrow ||r||, \quad v_1 \leftarrow \frac{r}{a}$ 3: for $i = 1, 2, \cdots, m_{\max}$ do 4: $w_i \leftarrow Av_i$, make q a vector $\in \mathbb{R}^{j+1}$ by appending a zero 5: Run the MGS process to obtain $h_{1:i,i}$ and update w_i 6: $h_{i+1,i} \leftarrow ||w_i|| > h_{i,i}$ denotes element of matrix \overline{H}_m 7: $v_{j+1} \leftarrow \frac{w_j}{h_{j+1,j}}$ 8: Apply the Givens transformations G_1, \dots, G_{i-1} to h_{i} 9: Compute the Givens matrix G_i for $h_{i,j}$ to make $h_{i+1,j} = 0$ 10: Apply the Givens transformation G_i to q11: \triangleright last element of vector *q* $\gamma \leftarrow g_{i+1}$ 12: if $|y| \leq tol$ then 13: Break 14: end if 15: end for 16: $m \leftarrow i$ 17: Solve y_m with the transformed matrix \overline{H}_m and q18: $x_0 \leftarrow x_0 + V_m y_m$, with $V_m = [v_1, \cdots, v_m]$ 19: 20: end while 21: Return $x_{appro} = x_0$.

Randomized GMRES

 The theoretical analysis shows that randomized Gram-Schmidt process can promote effectiveness of Arnoldi process.

- However, it is not practical, for the following three oblivious drawbacks:
- (1) High computational cost of sketching.
- (2) Low efficiency of residual estimation.
- (3) Low efficiency of solving sketched LS.

Algorithm 3 The Randomized Arnoldi Process [3] **Input:** $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$, $x_0 \in \mathbb{R}^n$, *m*, random sketching $\Theta \in \mathbb{R}^{k \times n}$. **Output:** $V_{m+1} \in \mathbb{R}^{n \times (m+1)}, \overline{H}_m \in \mathbb{R}^{(m+1) \times m}$. 1: Calculate residual $r \leftarrow b - Ax_0$ 2: Initialize \overline{H} to be an $(m + 1) \times m$ zero matrix 3: $s_1 \leftarrow \frac{\Theta r}{\|\Theta r\|}, \quad v_1 \leftarrow \frac{r}{\|\Theta r\|}$ \triangleright normalize s_1 4: **for** $i = 1, 2, \dots, m$ **do** $w_i \leftarrow Av_i$ 5: ▶ random sketching 6: $p_i \leftarrow \Theta w_i$ Solve $z_i = \arg \min_z ||S_i z - p_i||$, with $S_i = [s_1, \dots, s_i]$ 7: $v_{j+1} \leftarrow w_j - V_j z_j$, with $V_j = [v_1, \cdots, v_j]$ 8: $s_{j+1} \leftarrow \Theta v_{j+1}$ ▶ random sketching 9: $\overline{H}_{1:j+1,j} \leftarrow [z_j^T, \|s_{j+1}\|]^T$ 10: $v_{j+1} \leftarrow \frac{v_{j+1}}{\|s_{j+1}\|}, \quad s_{j+1} \leftarrow \frac{s_{j+1}}{\|s_{j+1}\|}$ \triangleright normalize s_{i+1} 11: 12: end for 13: **Return** $V_{m+1} = [v_1, v_2, \cdots, v_{m+1}], \overline{H}_m = \overline{H}.$

The Aim of This Work

- We aim to develop a practical and more efficient randomized GMRES algorithm,
 and can replace GMRES algorithm at least in circuit simulation.
- The contribution of this work can be summarized as follows:
- (1) Propose a practical randomized Arnoldi process, which can efficiently and effectively estimate residual in each iteration.
- (2) **Propose a linear-time complexity sketching algorithm.**
- (3) Develop a practical randomized GMRES (PRGMRES) solvers.
- (4) Theoretical analysis of PRGMRES algorithm.

Contents

Background

Proposed Method

Experimental Results

Conclusion

The Practical Randomized Arnoldi Process

$$\|r_{m}\| = \|b - Ax_{m}\| = \|r_{0} - AV_{m}y_{m}\| = \|r_{0} - V_{m+1}\overline{H}_{m}y_{m}\|$$
$$= \|V_{m+1}(\beta e_{1} - \overline{H}_{m}y_{m})\| = \|\beta e_{1} - \overline{H}_{m}y_{m}\| \quad (2)$$
$$= |\gamma|,$$

- The norm of S and V is not equal with high probability due to the low rank of the sketching matrices. Therefore, the (2) does not hold and original residual estimation does not work for randomized GMRES in [3].
- To make the efficient residual estimation, we should make (2) work for randomized GMRES.

The key: Enforce orthonormality of V.

The Practical Randomized Arnoldi Process



(a) all s_j is normalized (b) all v_j is normalized Figure 1: Two different sketching schemes in (a) the RGMRES algorithm in [3] and (b) the proposed algorithm. Notice that the two sets of orthogonal bases satisfy $s_i = \Theta v_i$, $i = 1, 2, \cdots$.

The Practical Randomized Arnoldi Process

$$||r_{m}|| = ||b - Ax_{m}|| = ||r_{0} - AV_{m}y_{m}|| = ||r_{0} - V_{m+1}\overline{H}_{m}y_{m}||$$

= $||V_{m+1}(\beta e_{1} - \overline{H}_{m}y_{m})|| = ||\beta e_{1} - \overline{H}_{m}y_{m}||$ (2)
= $|\gamma|$,

• The sketching process cannot guarantee the fully orthogonality of V.

Therefore, the residual estimation is still not accurate in some extreme cases.

• The key: implement an inner tolerance (double tolerance), which is smaller than demanded tolerance to avoid redundant restarting.

The Linear-Time Sketching

- The Sketching: compute $y = \tilde{\Theta}x$
- The theoretical bound of Rademacher sketching is better than P-SRHT, while the time complexity of Rademacher sketching is O(nk) and the time complexity of P-SRHT is O(nlogn) or O(nlogk).
- We aim to implement **sketching in linear complexity** by leveraging special structure.

The Linear-Time Sketching: Implicit Generation

Algorithm 4 Generate the Matrices for Linear-time Sketching

Input: n, k, ξ, ζ . **Output:** $D \in \mathbb{N}^{k \times u}$, $E \in \mathbb{N}^{k \times u}$, $F \in \mathbb{N}^{k \times u}$, $P \in \mathbb{R}^{\xi \times n}$. 1: $u \leftarrow \lfloor \zeta \frac{n}{k} \rfloor$ 2: Generate random matrix $P \in \mathbb{R}^{\xi \times n}$ whose element is $\frac{1}{\sqrt{L}}$ with probability 0.5 and $-\frac{1}{\sqrt{k}}$ with probability 0.5 3: **for** $i = 1, 2, \cdots, k$ **do** Partition integer array 1 : *n* into *u* segments 4: for $j = 1, 2, \dots u$ do 5: $d_{i,j} \leftarrow$ the left endpoint of *j*-th segment 6: $e_{i,j} \leftarrow$ the right endpoint of *j*-th segment 7: $f_{i,i} \leftarrow$ a uniformly random integer between 1 and ξ . 8: end for 9: 10: **end for** 11: **Return** $D = (d_{i,j})_{k \times u}, E = (e_{i,j})_{k \times u}, F = (f_{i,j})_{k \times u}, P.$

18

The Linear-Time Sketching: Sketching in Linear-Time

Algorithm 5 Perform Random Sketching in Linear Time

Input: $x \in \mathbb{R}^n$, $D \in \mathbb{Z}^{k \times u}$, $E \in \mathbb{Z}^{k \times u}$, $F \in \mathbb{Z}^{k \times u}$, $P \in \mathbb{R}^{\xi \times n}$. ζ . **Output:** $y \in \mathbb{R}^k$. \triangleright compute $y = \tilde{\Theta}x$ 1: **for** $i = 1, 2, \cdots, \xi$ **do** 2: **for** $j = 1, 2, \dots, n$ **do** Compute t(i, j) with (6) 3: end for 4: 5: end for 6: **for** $i = 1, 2, \dots, k$ **do** Compute y_i with (7), where $u = \lfloor \zeta \frac{n}{k} \rfloor$ 7: 8: end for 9: **Return** $y = [y_1, y_2, \cdots, y_k]^T$.

9:

Incremental Solution of the sketched least squares problem in randomized Arnoldi process with Householder transforms.

• The time complexity of solving sketched least-squares problems: $O(kj^2)$ to O(kj).

$$(\Theta V_j)z \approx \Theta w_j$$
, i.e. $\min_z \|(\Theta V_j)z - \Theta w_j\|$,

Solve
$$z_j = \arg \min_z ||S_j z - p_j||$$
, where $S_j = [s_1, \dots, s_j]$

The Overall Algorithm

• With the above

techniques, we can obtain

the overall algorithm for

PRGMRES algorithm.

Algorithm 6 The Practical Randomized GMRES Algorithm **Input:** $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$, $x_0 \in \mathbb{R}^n$, m_{max} , tol, α , k, ξ , ζ . **Output:** $x_{appro} \in \mathbb{R}^n$. 1: Generate matrices *D*. *E*. *F* and *P* with Alg. 4. 2: $r \leftarrow b - Ax_0$, $a \leftarrow ||r||$ 3: $v_1 \leftarrow \frac{r}{a}$ \triangleright normalize v_1 ▶ linear-time sketching 4: $s_1 \leftarrow$ sketching of v_1 with Alg. 5 ▶ check the accurate residual 5: while ||r|| > tol do for $i = 1, 2, \dots, m_{\max}$ do 6: $w_i \leftarrow Av_i$, make *q* a vector $\in \mathbb{R}^{j+1}$ by appending a zero 7: $p_i \leftarrow$ sketching of w_i with Alg. 5 8: Solve $z_i = \arg \min_z ||S_i z - p_i||$, where $S_i = [s_1, \dots, s_i]$ 9: $v_{i+1} \leftarrow w_i - V_i z_i$, with $V_i = [v_1, \cdots, v_i]$ 10: $s_{i+1} \leftarrow$ sketching of v_{i+1} with Alg. 5 11: $h_{1:j+1,j} \leftarrow [z_{j}^{T}, ||v_{j+1}||]^{T}$ 12: $v_{j+1} \leftarrow \frac{v_{j+1}}{\|v_{j+1}\|}, \quad s_{j+1} \leftarrow \frac{s_{j+1}}{\|v_{j+1}\|} \qquad \triangleright \text{ normalize } v_{j+1}$ 13: Apply the Givens transformations G_1, \dots, G_{j-1} to $h_{:,j}$ 14: Compute the Givens matrix G_i for $h_{i,i}$ to make $h_{i+1,i} = 0$ 15: Apply the Givens transformation G_i to q16: \triangleright last element of vector *q* $\gamma \leftarrow g_{i+1}$ 17: if $|y| \leq \alpha \cdot tol$ then ▶ a slightly tighter criterion 18: Break 19: end if 20: end for 21: $m \leftarrow i$ 22: Obtain y_m with the transformed matrix \overline{H}_m and q23: $x_0 \leftarrow x_0 + V_m y_m$, with $V_m = [v_1, \cdots, v_m]$ 24: $r \leftarrow b - Ax_0$ ▶ an extra computation of residual 25: $g \leftarrow ||r||, \quad v_1 \leftarrow \frac{r}{a}$ 26: 27: end while 28: **Return** $x_{appro} = x_0$.

21

Theoretical Analysis

Results on proposed sketching algorithm.

THEOREM 1. Suppose $\tilde{\Theta} \in \mathbb{R}^{k \times n}$ is a random matrix generated with the above method, then the expectation

$$\mathbb{E}[\tilde{\Theta}^T \tilde{\Theta}] = I, \qquad (8)$$

where I denotes the identity matrix. Moreover, every diagonal element of $\tilde{\Theta}^T \tilde{\Theta}$ is always 1.

Theoretical Analysis

• Results supporting smaller sketching size can be summarized in brief.

(1) Oblivious subspace embedding for m-dimensional Krylov subspace is not worse than any m-dimensional subspace of \mathbb{R}^n .

(2) Oblivious subspace embedding for the first k (k is a small number) basis vectors of Krylov subspace may have good results as well.

• The smaller sketching dimensions can be just fine.

Contents

Background

- Proposed Method
- Experimental Results
- Conclusion

Experiments Results

Experimental setting:

(1) Machine: Intel Xenon Gold 6230R CPU@2.10GHz and 128 RAM.

(2) Preconditioning: ILU(k) (k=3)

(3) Tolerance: 1E-8 (to show the stability of our PRGMRES algorithm)

Experiments Results

	Case	n	nnz	GMRES		RGMRES [3]		PRGMRES (ours)				
				iter	$T_{tot}(\mathbf{s})$	iter	$T_{tot}(\mathbf{s})$	iter	$T_{tot}(\mathbf{s})$	Sp ₁	Sp_2	
	ibmpg3t	1.0E6	4.2E6	202	41.9	200	33.6	198	25.7	1.6	1.3	
	ibmpg4t	1.2E6	4.8E6	146	29.0	150	22.6	149	18.5	1.6	1.2	
	ibmpg5t	1.6E6	5.7E6	645	206	670	128	701	122	1.7	1.0	
	thupg1	5.3E6	2.2E7	137	132	140	103	141	96.9	1.4	1.1	
	thupg2	9.5E6	4.1E7	140	292	150	213	143	180	1.6	1.2	
	thupg3	1.2E7	5.3E7	140	402	150	290	143	245	1.6	1.2	
	thupg4	1.6E7	6.9E7	143	583	150	386	146	317	1.8	1.2	
	thupg5	2.0E7	8.8E7	150	972	150	477	154	452	2.2	1.1	
	scircuit	1.7E5	9.6E5	642	19.2	650	15.2	664	12.1	1.6	1.3	
	rajat31	4.7E6	2.0E7	113	84.4	120	68.0	115	61.4	1.4	1.1	
	FullChip	3.0E6	2.7E7	17	7.44	210	98.8	20	9.60	0.8	10.3	
	Freescale1	3.4E6	1.9E7	275	188	280	119	284	118	1.6	1.0	
	nxp1	4.1E5	2.7E6	151	11.7	160	9.86	157	8.49	1.4	1.2	
	nlcircuit	1.6E5	7.3E5	21	0.535	30	0.636	22	0.449	1.2	1.4	
	Average									1.5	1.8	

Table 1: Computational results of solving sparse linear equations. The convergence tolerance of GMRES method is all set to 1E-8, the restarting parameter is $m_{\text{max}} = 200$, and the sketching size k of RGMRES and our PRGMRES are 1500 and 500 respectively.

n and *nnz* are the dimension and the number of nonzeros in coefficient matrix, respectively. *iter* and T_{tot} denote the number of iterations and runtime of GMRES iterations, respectively. Sp_1 denotes the speedup ratio of PRGMRES to GMRES, while Sp_2 denotes the speedup ratio of PRGMRES to RGMRES.

Experiments Results



within the proposed PRGMRES algorithm for case "ibmpg3t".

Contents

Background

- Proposed Method
- Experimental Results
- Conclusion



• A practical randomized Arnoldi process which can efficiently estimate residual.

- A linear-time sketching algorithm.
- The overall PRGMRES algorithm
- The theoretical analysis on PRGMRES
- The practical ILU-PRGMRES solver.

Other Possible Applications



- Below are some of other possible applications in EDA:
- (I) Capacitance Extraction
- (2) Multi-physics simulation
- (3) Other EDA applications as long as linear equation solvers are needed
- Below are some of other possible applications besides EDA:
- (I) Solving eigenvalue: proposed practical randomized Arnoldi process.
- (2) Randomized singular value decomposition: proposed sketching algorithm.



Presenter: BAIYU CHEN