

Dissection solver for higher precision arithmetic by inner iterative refinement

Atsushi Suzuki¹

¹Cybermedia Center, Osaka University
atsushi.suzuki@cas.cmc.osaka-u.ac.jp

joint work with
François-Xavier Roux, ONERA, LJLL Sorbonne Université

Outline

- ▶ overview of nested-dissection algorithm for sparse matrix
- ▶ LDU factorization with symmetric pivoting
- ▶ computation of a solution of singular linear system
- ▶ extension with 2x2 pivoting and unsymmetric pivoting
- ▶ kernel detection algorithm
- ▶ higher precision arithmetic with inner iterative refinement
- ▶ numerical results on matrices from a semi-conductor problem

recursive generation of Schur complement by nested-dissection

$$\begin{bmatrix} A_{11} & A_{21} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S_{22} \end{bmatrix} \begin{bmatrix} I_1 & A_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix}$$

$S_{22} = A_{22} - A_{21}A_{11}^{-1}A_{12} = A_{22} - (A_{21}U_{11}^{-1})D_{11}^{-1}L_{11}^{-1}A_{12}$: recursively computed

8 9 a b c d e f 4 5 6 7 2 3 1

88 99 aa bb cc dd ee ff	84 94 a5 b5 c6 d6 e7 f7	82 92 91 a2 b1 b2 b1 c1 d3 d1 e3 e1 f3
48 49 5a 5b 6c 6d 7e 7f	44 55 66 77	42 52 61 73
28 29 2a 2b 3d 3e 3f 19 1b 1c 1d 1e	24 25 37 16	22 21 33 31 12 13 11

Schur complement
by sparse solver

44 55 66 77	42 41 52 51 63 61 73 71
24 25 36 37 14 15 16 17	22 21 33 31 12 13 11

Schur complement
by dense solver

Schur complement
by dense solver

11

dense factorization

22 33	21 31
12 13 11	

sparse part : completely in parallel

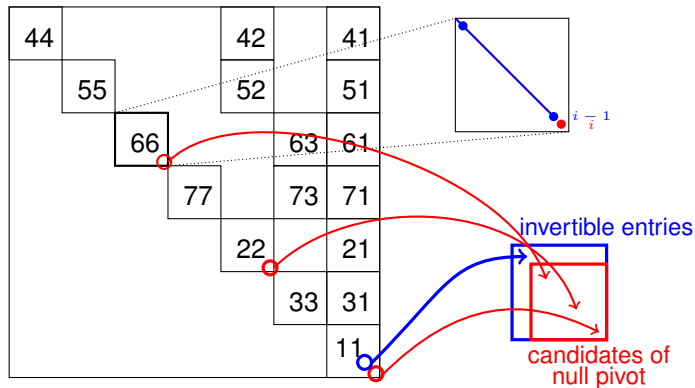
dense part : better use of **BLAS 3**; dgemm, dtrsm

Symmetric pivoting with postponing for block strategy : 1/2

- ▶ nested-dissection decomposition may produce singular sub-matrix for indefinite matrix

τ : given threshold for postponing, 10^{-2}

$|A(i, i)|/|A(i-1, i-1)| < \tau \Rightarrow \{A(k, j)\}_{i \leq k, j}$ are postponed

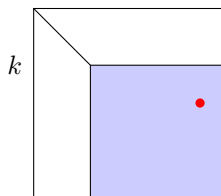


Schur complement matrix from postponed pivots is computed

Pivoting strategy

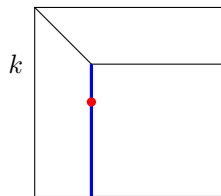
full pivoting : $A = \Pi_L^T L U \Pi_R$

find $\max_{k < i, j \leq n} |A(i, j)|$



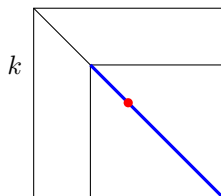
partial pivoting : $A = \Pi L U$

find $\max_{k < i \leq n} |A(i, k)|$



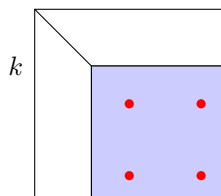
symmetric pivoting : $A = \Pi^T L D U \Pi$

find $\max_{k < i \leq n} |A(k, k)|$



2×2 pivoting : $A = \Pi^T L \tilde{D} U \Pi$

find $\max_{k < i, j \leq n} \det \begin{vmatrix} A(i, i) & A(i, j) \\ A(j, i) & A(j, j) \end{vmatrix}$



sym. pivoting is mathematically not always possible

Computation of a solution of singular linear system

- ▶ $A \in \mathbb{R}^{N \times N}$, $k = \dim \text{Ker} A$.
- ▶ index ordering $\{i_1, i_2, \dots, i_N\}$ $V_{N-k} = \text{span}[\vec{e}_{i_1}, \vec{e}_{i_2}, \dots, \vec{e}_{i_{N-k}}]$

assumption 1 : $V_{N-k} \cap \text{ker} A = \{\vec{0}\} \Rightarrow \exists A_{11}^{-1}$

Solution in the space V_{N-k}

$$\vec{x} = A^\dagger \vec{f} \Leftrightarrow \vec{x} \in V_{N-k} \quad (A\vec{x} - \vec{f}, \vec{v}) = 0 \quad \forall \vec{v} \in V_{N-k}$$

$$N_1 = \begin{bmatrix} A_{11}^{-1} A_{12} \\ -I_2 \end{bmatrix} \quad N_2 = \begin{bmatrix} A_{11}^{-T} A_{21}^T \\ -I_2 \end{bmatrix}, \quad \text{ker} A = \text{span} N_1, \quad \text{ker} A^T = \text{span} N_2$$

assumption 2 :

$\vec{u} = A^\dagger \vec{f}$ is computable by *LDU-factorization with symmetric pivot*.

$$\vec{f} \in \text{Im} A$$

$$\Leftrightarrow \vec{f} \perp \text{Ker} A^T \Leftrightarrow N_2^T \vec{f} = \vec{0} \Leftrightarrow A_{21} A_{11}^{-1} \vec{f}_1 - \vec{f}_2 = \vec{0} \Rightarrow A A^\dagger \vec{f} = \vec{f}$$

$$\forall \vec{\zeta} \in \mathbb{R}^k, \quad A(\vec{u} + N_1 \vec{\zeta}) = \vec{f}$$

- ▶ $\vec{u} + N_1 \vec{\zeta} \in (\text{Ker} A)^\perp \Leftrightarrow N_1^T (\vec{u} + N_1 \vec{\zeta}) = \vec{0}$: algebraic inverse
- ▶ $\vec{u} + N_1 \vec{\zeta} \in V_{N-k} \Leftrightarrow \vec{\zeta} = \vec{0}$: computational
- ▶ $\vec{u} + N_1 \vec{\zeta} \in \text{Im} A \Leftrightarrow N_2^T (\vec{u} + N_1 \vec{\zeta}) = \vec{0} \Leftrightarrow \text{Im} A \cap \text{Ker} A = \{\vec{0}\}$

$\text{Im} A \cap \text{Ker} A = \{\vec{0}\} \Rightarrow N_2^T N_1$ is invertible

$$N_1 \vec{\eta} \in \text{Im} A \cap \text{Ker} A \Leftrightarrow N_2^T N_1 \vec{\eta} = \vec{0}, \quad \text{Im} A \cap \text{Ker} A = \{\vec{0}\} \Leftrightarrow (N_2^T N_1 \vec{\eta} = \vec{0} \Rightarrow \vec{\eta} = \vec{0})$$

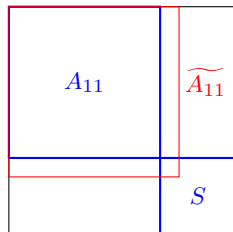
extension with 2x2 pivoting or unsymmetric pivoting

For matrix from PDE, most of part can be factorized with symmetric pivoting.

\widetilde{A}_{11} has LDU -factorization with symmetric pivot entries of $S \Leftarrow$ postponed pivots + some invertible entries of \widetilde{A}_{11}

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S \end{bmatrix} \begin{bmatrix} I_{11} & A_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix}$$

$$S = A_{22} - A_{21}A_{11}^{-1}A_{12}$$



- ▶ applying full pivot when S is (highly) unsymmetric, $\bar{S} = \Pi_L S \Pi_R$
 - ▶ applying 2x2 pivot when S is symmetric + indefinite
- with full pivoting Π_L/Π_R

$$\bar{S} = \Pi_L S \Pi_R = \begin{bmatrix} \bar{S}_{22} & \bar{S}_{23} \\ \bar{S}_{32} & \bar{S}_{33} \end{bmatrix} = \begin{bmatrix} \bar{S}_{22} & 0 \\ \bar{S}_{32} & E_{33} \end{bmatrix} \begin{bmatrix} I_2 & \bar{S}_{22}^{-1} \bar{S}_{23} \\ 0 & I_3 \end{bmatrix}$$

$$E_{33} = 0 \Rightarrow \text{Ker} \bar{S} = \text{span} \begin{bmatrix} \bar{S}_{22}^{-1} \bar{S}_{23} \\ -I_3 \end{bmatrix}, \quad \text{Ker} S = \text{span} \Pi_R \begin{bmatrix} \bar{S}_{22}^{-1} \bar{S}_{23} \\ -I_3 \end{bmatrix}$$

$$\text{Ker} A = \begin{bmatrix} A_{11}^{-1} A_{12}^{(2)} & A_{11}^{-1} A_{12}^{(3)} \\ -I_2^{(2)} & 0 \\ 0 & -I_2^{(3)} \end{bmatrix} \Pi_R \begin{bmatrix} \bar{S}_{22}^{-1} \bar{S}_{23} \\ -I_3 \end{bmatrix}$$

novel kernel detection algorithm based on LDU

$A : N \times N$ unsymmetric, $\dim \text{Ker} A = k \geq 1$, $\dim \text{Im} A \geq m$.

two parameters: l, n , which define size of factorization,

$$\begin{matrix} N-n \\ n \end{matrix} \begin{matrix} \updownarrow \\ \updownarrow \end{matrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S_{22} \end{bmatrix} \begin{bmatrix} I_1 & A_{11}^{-1} A_{12} \\ 0 & I_2 \end{bmatrix} \quad \widetilde{\text{Im}}_n = \text{span} \begin{bmatrix} \bar{A}_{11}^{-1} A_{12} \\ -I_2 \end{bmatrix}^\perp.$$

► projection : $P_n^\perp : \mathbb{R}^N \rightarrow \widetilde{\text{Im}}_n$

► solution in subspace, $\bar{A}_{N-l}^\dagger \vec{b} = \begin{bmatrix} \bar{A}_{11}^{-1} \vec{b}_1 \\ \vec{0} \end{bmatrix}$, $b = \begin{bmatrix} \vec{b}_1 \\ \vec{b}_2 \end{bmatrix}$ $\begin{matrix} \updownarrow \\ \updownarrow \end{matrix} \begin{matrix} N-l \\ l \end{matrix}$

\bar{A}_{11}^{-1} : computed in quadruple-precision with perturbation to simulate double-precision round-off error.

kernel detection algorithm

DOI: 10.1002/nme.4729

n : candidate of dimension of the kernel

compute for $l = n - 1, n, n + 1$

$$\text{err}_l^{(n)} := \max \left\{ \max_{\vec{x} = [\vec{0} \ \vec{x}_l] \neq \vec{0}} \frac{\|P_n^\perp(\bar{A}_{N-l}^\dagger A \vec{x} - \vec{x})\|}{\|\vec{x}\|}, \max_{\vec{x} = [\vec{x}_{N-l} \ \vec{0}] \neq \vec{0}} \frac{\|\bar{A}_{N-l}^\dagger A \vec{x} - \vec{x}\|}{\|\vec{x}\|} \right\}$$

$$n = k + 1 \quad \Leftrightarrow \quad \text{err}_{n-1} \approx 0 \quad \wedge \quad \text{err}_n \approx 0 \quad \wedge \quad \text{err}_{n+1} \sim 1$$

$$n = k \quad \Leftrightarrow \quad \text{err}_{n-1} \gg 0 \quad \wedge \quad \text{err}_n \approx 0 \quad \wedge \quad \text{err}_{n+1} \sim 1$$

$$n = k - 1 \quad \Leftrightarrow \quad \text{err}_{n-1} \gg 0 \quad \wedge \quad \text{err}_n \gg 0 \quad \wedge \quad \text{err}_{n+1} \sim 1$$

Theoretically $\neg A_{N-k+1}^{-1}$, but $\text{err}_{k-1}^{(k)}$ is computable; $\sim \|\bar{A}_N^{-1} A_N - I_N\|$.

higher precision arithmetic with inner iterative refinement

for high condition number matrix with coefficients given by “double”
 A_{11} is moderate part \Leftarrow postponing strategy with given threshold

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S \end{bmatrix} \begin{bmatrix} I_{11} & A_{11}^{-1} A_{12} \\ 0 & I_2 \end{bmatrix}$$

$$S = A_{22} - A_{21} A_{11}^{-1} A_{12}$$

computation of S in higher precision is necessary

- ▶ iterative refinement to compute $X_{12} = A_{11}^{-1} A_{12}$
- ▶ quadruple precision arithmetic to compute $S = A_{22} - A_{21} X_{12}$.
- ▶ quadruple precision LDU -factorization for S with kernel detection

“quadruple” arithmetic is realized by “double-double” in QD library
iterative refinement

$A_{11} X_{12} = A_{12}$ in “double”

$R_{12}^{(1)} = A_{12} - A_{11} X_{12}$ in “quadruple”

loop $n = 1, 2, \dots$

$A_{11} Y_{12}^{(n)} = R_{12}^{(n)}$ in “double”

$X_{12} = X_{12} + Y_{12}^{(n)}$ in “quadruple”

$R_{12}^{(n+1)} = A_{12} - A_{11} X_{12}$ in “quadruple”

SpMV by double-double is 10 times arithmetic intensive than double

Drift-Diffusion system at stationary state

- ▶ φ : electrostatic potential
- ▶ n : electron concentration
- ▶ p : hole concentration

$$\operatorname{div}(\varepsilon E) = q(p - n + C(x))$$

$$E = -\nabla\varphi$$

$$-\operatorname{div}J_n = 0$$

$$J_n = -q(\mu_n n \nabla\varphi - D_n \nabla n)$$

$$\operatorname{div}J_p = 0$$

$$J_p = -q(\mu_p p \nabla\varphi + D_p \nabla p)$$

- ▶ q : positive electron charge
- ▶ ε : dielectric constant of the materials
- ▶ $D_n = \mu_n \theta$, $D_p = \mu_p \theta$: Einstein's relation

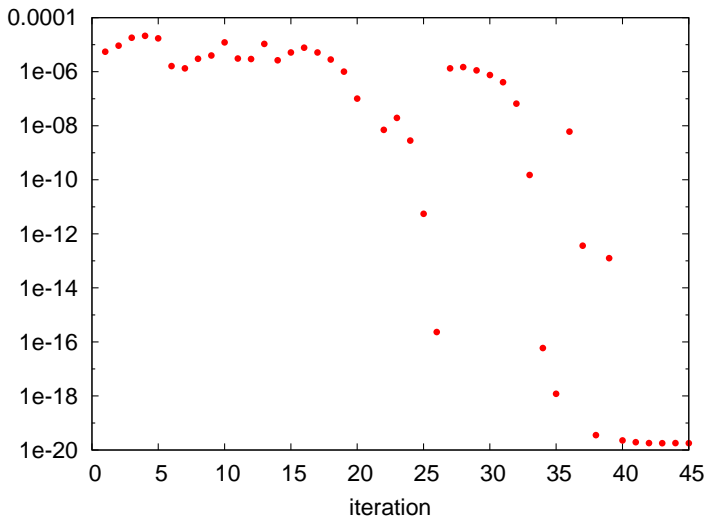
Maxwell-Boltzmann statistics : $p = N_i \exp(\frac{\varphi_p - \varphi}{\theta})$ leads to

$$J_p = -q\mu_p p \nabla\varphi_p$$

- ▶ φ_p : quasi-Fermi level
- ▶ N_i : intrinsic concentration of the semiconductor
- ▶ $\theta = kT/q$: thermal voltage
- ▶ k : Boltzmann constant, T : lattice temperature

Finite Volume with Scharfetter-Gummel scheme + Newton iteration

Newton iteration of a semi-conductor problem



from 1st to 5th iterations and 10th, matrices are invertible, others are recognized as singular

matrices are provided by [Toshiba Memory Corp.](#)

results of kernel detection

$\tilde{A} = W A W$: scaled with $[W]_{ii} = 1\sqrt{[A]_{ii}}$, $n = 568,455$, $nnz = 8,330,429$

matrix id	5th	6th	25th
dimKer \tilde{A}	0	2	2
diagonal entries of S by QR	8.574766e+07	1.996129e+07	1.943695e+08
	8.902046e+04	9.001534e+04	4.834031e+02
	2.063527e-01	3.050581e-01	2.165032e-01
	1.799643e-01	2.632021e-01	2.034221e-01
	5.406636e-11	4.885182e-13	2.373664e-10
	1.286903e-13	2.328076e-15	2.063879e-16
$\ \tilde{A}v_k\ $		4.376530e-16	1.102714e-16
		1.565334e-15	7.982637e-17
$\ \tilde{A}^T w_k\ $		1.982592e-17	3.589934e-18
		1.983189e-17	3.589934e-18
(v_k, w_l)		4.482117e-12	2.802405e-13
		9.361029e-15	9.199130e-14
		2.251826e-13	2.810891e-13
		3.187650e-14	9.222188e-14
error	7.793333e-32	2.147198e-29	6.465044e-23
residual	5.901230e-32	5.658071e-32	1.605194e-30

- ▶ for 5th matrix, by assuming 2 dimensional kernel, $\|\tilde{A}v_1\|=1.286904e-13$
- ▶ error and residual are computed by a manufactured solution

performance comparison

$n = 568,455$, $nnz = 8,330,429$, i7-6770HQ CPU @ 2.60GHz, g++-5, MKL BLAS

matrix id	1st	25th	25th
# postponed	w/o postponing	6	6
$\dim \text{Ker} \tilde{A}$	0	2	0
diagonal entries of S by QR		1.943695e+08 4.834031e+02 2.165032e-01 2.034221e-01 2.373664e-10 2.063879e-16	1.799841e+08 4.161368e+02 2.975576e-01 2.440021e-01 2.158641e-10 1.726983e-16
$\ \tilde{A}v_k\ $		1.102714e-16 7.982637e-17	2.063879e-16 1.188300e-20
$\ \tilde{A}^T w_k\ $		3.589934e-18 3.589934e-18	
error	1.272821e-15	6.465044e-23	8.784892e-23
residual	2.442552e-15	1.605194e-30	1.584077e-30
time for factorization(sec.)	double 22.580	quasi-quadruple 32.010	quadruple 5,620.7

- ▶ in quasi-quadruple computation, $A_{11}^{-1}A_{12}$ is computed by iterative refinement with 3 iterations
- ▶ in whole quadruple computation, S is generated by recursive computation following elimination tree

Summary

- ▶ Iterative refinement process for computing Schur complement matrix consisting postponed pivots can drastically reduce computational time without losing accuracy
- ▶ Full pivoting procedure works as a preconditioner for kernel detection algorithm based on symmetric pivoting
- ▶ Complexity of quadruple precision arithmetic by double-double is 20 times large as double precision arithmetic. However, combination of inner iterative refinement and quadruple arithmetic for the last Schur complement requires only 50% increased elapsed time than double precision, which is numerically verified by matrices in a semi-conductor simulation

ongoing

- ▶ optimizing computation of iterative refinement by replacing multiple SpMV by SpMM

source code of Dissection is accessible within FreeFem++ repository
<https://github.com/FreeFem/FreeFem-sources/tree/master/download/dissection>

under GPL linking-exception / CeCILL-C licenses