



New robust limited-memory incomplete Cholesky factorizations

Jennifer Scott

STFC Rutherford Appleton Laboratory

Miroslav Tůma

Institute of Computer Science
Academy of Sciences of the Czech Republic

MNR13 Workshop, Manchester, 23rd October 2013

Numerical Analysis Group at RAL

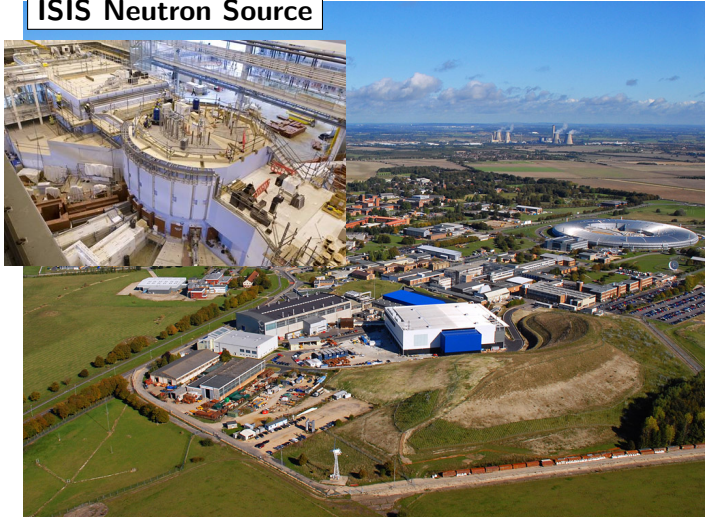
- ▶ Part of **Scientific Computing Department** (SCD) of the **Science and Technology Facilities Council** (STFC).
- ▶ STFC is a research council with the mission to
 - ▶ “maximise the impact of our knowledge, skills, facilities and resources for the benefit of the United Kingdom and its people.”
 - ▶ World-class research
 - ▶ World-class innovation
 - ▶ World-class skills

STFC Rutherford Appleton Laboratory



STFC Rutherford Appleton Laboratory

ISIS Neutron Source



STFC Rutherford Appleton Laboratory

ISIS Neutron Source



Central Laser Facility

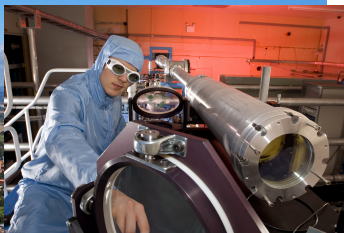


STFC Rutherford Appleton Laboratory

ISIS Neutron Source



Central Laser Facility



Diamond Light Source



STFC Rutherford Appleton Laboratory

ISIS Neutron Source



Central Laser Facility



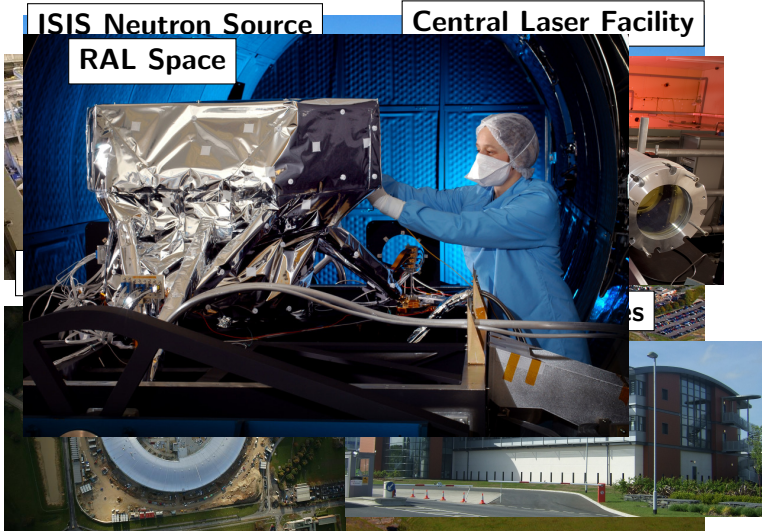
Diamond Light Source



HPC facilities



STFC Rutherford Appleton Laboratory



STFC Rutherford Appleton Laboratory

ISIS Neutron Source

Central Laser Facility

RAL Space

Research Complex at Harwell



Numerical Analysis Group at RAL

- ▶ **SCD** provides world-class expertise and support for UK theoretical and computational science communities, in both academia and industry.
- ▶ Approx. 160 staff based at Daresbury (Warrington) and Rutherford.
- ▶ NA Group is part of the **Technology Division**.
- ▶ Group has existed since late 1950s (originally at Harwell, moved to RAL in 1990).
- ▶ Group currently mainly funded by 4-year **ESPRC grant**.

Numerical Analysis Group at RAL

Jennifer Scott (Group leader) Sparse linear algebra, high-performance computing.

Iain Duff Sparse linear systems, high-performance computing.

Nick Gould Nonlinear systems, iterative methods, optimization.

Jonathan Hogg Parallel linear algebra, optimization.

Evgueni Ovtchinnikov Sparse linear algebra, GPUs.

Tyrone Rees Numerical linear algebra, PDE constrained optimization

Sue Thorne Numerical linear algebra, iterative methods, preconditioning.

Numerical Analysis Group at RAL

Plus

Mario Arioli Numerical linear algebra, numerical solution of PDEs, error analysis. **Retiring but still working!**

John Reid Sparse matrices, automatic differentiation, Fortran. **Honorary Scientist.**

Numerical Analysis Group at RAL

Main areas of expertise:

- ▶ Sparse linear algebra $Ax = b$, $Ax = \lambda Bx$
- ▶ Direct and iterative methods
- ▶ Large-scale optimization
- ▶ High-quality software
 - ▶ **HSL** - emphasis on sparse linear algebra (especially sparse direct solvers). **Celebrating 50 years!**
 - ▶ **GALAHAD** - nonlinear optimization
 - ▶ **CUTEr** - testing environment for optimization and linear algebra solvers

Sparse Solvers

We want to solve

$$Ax = b$$

where $A \in R^{n \times n}$ is large and sparse.

Direct Methods: Factorize $A = LU$, solve $Ly = b$, $Ux = y$.

Black-box, robust.

Memory-hungry \Rightarrow unsuitable for very large problems.

Iterative Methods: CG, GMRES, BiCGStab, etc.

Matrix-free. Fast? Efficient? .

Non-robust, performance depends on preconditioner.

Hybrid Methods: Lots of combinations ...

Background

An ideal preconditioner should be:

- ▶ cheap to compute (time and memory)
- ▶ sparse and fast to apply
- ▶ provide sufficient approximation of the algebraic problem
- ▶ result in rapidly converging preconditioned iterative method

Key target for library software is **robustness**

Background

Incomplete Cholesky factorization

$$A \simeq LL^T$$

Some entries that occur in complete factorization are ignored.

But which to keep?

Long history (> 50 years) and many possible variants:

- ▶ **Structure-based** $IC(\ell)$: potential fill entries allowed only if their level of fill is less than ℓ .
- ▶ **Threshold-based** $IC(\tau)$: entries smaller than τ dropped.
- ▶ **Memory-based** $IC(p)$: dropping of entries based on memory available.

$IC(\ell)$

- ▶ Location of permissible fill entries using sparsity pattern of A prescribed in advance.
- ▶ Aim to mimic how pattern of A is developed during complete factorization.
- ▶ **But** although entries of $E = A - LL^T$ are zero inside prescribed sparsity pattern, outside can be **large**.
- ▶ **Increasing ℓ can be prohibitive** (storage requirements and time to compute and apply the preconditioner).

$$IC(\tau)$$

- ▶ Entries of computed factors or intermediate quantities that are less than **drop tolerance** τ discarded.
- ▶ Success depends on suitable τ : highly problem dependent.
- ▶ Trade-off between sparsity and quality.
- ▶ Memory **not** predictable.

$IC(p)$

- ▶ Prescribe maximum number of entries allowed in each column of L and retain only largest entries.
- ▶ **Memory predictable.**
- ▶ Example is widely-used dual threshold $ILUT(p, \tau)$ (Saad '94).
 - ▶ Designed for non symmetric problems.
 - ▶ Combines use of drop tolerance τ with prescribed maximum column and row counts.
 - ▶ Ignores symmetry in A (if A symmetric, patterns of L and U^T normally different).

ICFS

ICFS code of Lin and Moré '99:

- ▶ Widely used for large-scale trust region subproblems.
- ▶ Given p , retains $n_j + p$ largest entries in the lower triangular part of L_j , where n_j is number of entries in lower triangular part of A_j .
- ▶ **But**, as we will see, efficiency of resulting preconditioner not very sensitive to choice of p .

Breakdown?

- ▶ Problem with *IC* factorization: it may **breakdown** (that is, reach a point where the diagonal entry of the pivot column is zero).
- ▶ Kershaw '78 **locally perturbed** zero or negative diagonal entries to prevent breakdown so method more widely applicable. Straightforward but can give large growth and unstable preconditioner.
- ▶ Alternative remedy: use global diagonal shift so that $A + \alpha I$ factorized for some $\alpha > 0$ (Manteuffel '80).
- ▶ This is used in ICFS.

Are there other things we can do?

Positive semi-definite modifications I

- ▶ Diagonal modification scheme first introduced by Jennings and Malik '77,'78 (also Ajiz and Jennings '84).
- ▶ Every time off-diagonal entry discarded, corresponding diagonal entries modified by adding **SPSD** matrix

$$\begin{matrix} & & i & & j & & \\ & & & & & & \\ i & & & & & & \\ & & & & & & \\ j & & & & & & \\ & & & & & & \end{matrix} \begin{pmatrix} \ddots & & & & & & \\ & |a_{ij}| & & & -|a_{ij}| & & \\ & & \ddots & & & & \\ & -|a_{ij}| & & & |a_{ij}| & & \\ & & & & & \ddots & \\ & & & & & & \ddots \end{pmatrix}$$

Jennings-Malik approach

- ▶ **Breakdown-free** factorization that can be expressed as

$$A = LL^T - E$$

where error matrix E is sum of SPSD matrices.

- ▶ **But** modifications to A can be significant.
- ▶ Popular in some engineering applications.

Positive semi-definite modifications II

- ▶ More sophisticated modification scheme due to Tismenetsky '91 (and Kaporin '98).
- ▶ Introduces use of **intermediate memory** that is employed during construction of L but then discarded.
- ▶ Shown to be very **robust** but it “has unfortunately attracted surprisingly little attention” (Benzi '02).
- ▶ Suffers from a serious drawback: **memory requirements can be prohibitively high**.

Our aims

- ▶ Develop **generalisation of ICFS** such that efficiency of preconditioner improves with prescribed memory.
- ▶ Develop **memory-efficient** variant of Tismenetsky-Kaporin approach using **global shifts** to avoid breakdown.
- ▶ Combine in “black-box” *IC* factorization code that is demonstratively robust, efficient and flexible.

New package is **HSL_MI28**.

Tismenetsky approach

Based on matrix decomposition of form

$$A = (L + R)(L + R)^T - E,$$

- ▶ L is lower triangular with positive diagonal entries used for preconditioning,
- ▶ R is strictly lower triangular with small entries that is used to stabilise the factorization process, and
- ▶ E has the structure

$$E = RR^T.$$

Tismenetsky approach

- ▶ On j -th step, decompose col. 1 of Schur complement S into

$$l_j + r_j \quad \text{with} \quad |l_j|^T |r_j| = 0,$$

where entries of l_j are retained in incomplete factorization and those in r_j are discarded.

- ▶ On next step, S updated by subtracting

$$(l_j + r_j)(l_j + r_j)^T.$$

- ▶ Tismenetsky omits the term

$$E_j = r_j r_j^T. \tag{1}$$

- ▶ Thus, SPSD matrix implicitly added to A .

Kaporin's use of drop tolerances

- ▶ Obvious choice for r_j are smallest off-diagonal entries in col j .
- ▶ Controls size of L but **not** memory required to compute it.
- ▶ Kaporin '98: entries of magnitude at least τ_1 kept in L and those smaller than τ_2 are dropped from R .
- ▶ Now E has structure

$$E = RR^T + F + F^T,$$

F strictly lower triangular matrix that is **not computed**;
 R used in computation of L but **discarded**.

Problems of Tismenetsky-Kaporin approach

- ▶ How to choose tolerances τ_1 and τ_2 ? Problem dependent.
- ▶ Method not guaranteed breakdown free ... combine with diagonal compensation or global shift.
- ▶ With no restriction on size of L and R , can achieve **high quality** preconditioner but memory demands **high**.
- ▶ Also too **expensive**. Impractical for the very large problems iterative methods designed for.

Remedy: impose memory limit on L and R .

Limited memory Tismenetsky-Kaporin approach

- ▶ **lsize**: max. number of fill entries in each col. of L

$$nz(L) \leq nz(A) + \text{lsize} * (n - 1)$$

- ▶ **rsize**: max. number of entries in each col. of R .
Amount of intermediate memory and work involved in computing preconditioner depends on **rsize**.
Note: if **rsize** = 0, R not used.
- ▶ Retain **largest** entries in l_j , provided at least τ_1 in magnitude.
- ▶ Retain next largest entries in r_j , provided at least τ_2 in magnitude.

Left-looking algorithm outline

Input: A , $lsize$, $rsize$, τ_1 , τ_2

Set $w(1:n) = 0$

for $j = 1 : n$ **do**

Scatter col. A_j into w

Apply $LL^T + RL^T + LR^T$ updates from columns $1 : j - 1$ to w

(Partially) sort entries in w by magnitude

Keep $n_j + lsize$ entries of largest magnitude in l_j provided
they are at least τ_1

Keep $rsize$ additional entries that are next largest in magnitude
in r_j provided they are at least τ_2

Reset entries of w to zero

end do

end do

Output: L

Coping with breakdown

- ▶ When using limited memory (and/or dropping), factorization may **breakdown**.
- ▶ We hold a copy of diagonal entries and, at each step j , keep them updated. If any becomes zero or negative, restart factorization with

$$A \leftarrow A + \alpha I$$

for some positive α .

- ▶ More than one restart may be required.

Test environment

- ▶ Problems from University of Florida Collection.
- ▶ Selected all non-diagonal SPD matrices with $n > 1000$.
- ▶ Removed those with duplicate sparsity patterns.
- ▶ Following initial experiments, 8 problems discarded as unable to achieve convergence without large amount of fill.
- ▶ **Test set of 139 problems.**
- ▶ CG used with $x_0 = 0$, b computed so that $x = 1$, and stopping criteria

$$\|Ax_k - b\| \leq 10^{-10} \|b\|$$

with limit of **2000** iterations.

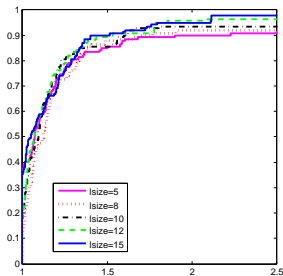
Test environment (continued)

- ▶ What to measure? iteration counts? timings? sparsity of L ?
- ▶ We define the **efficiency** of preconditioner to be

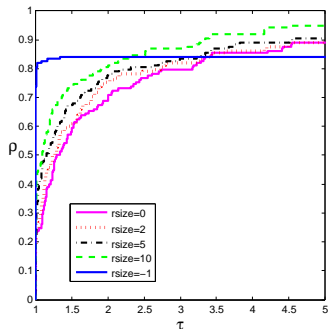
$$iter \times nz(L)$$

- ▶ **Performance profiles** (Moré, Dolan '02) used to assess performance.
- ▶ All software written in Fortran.

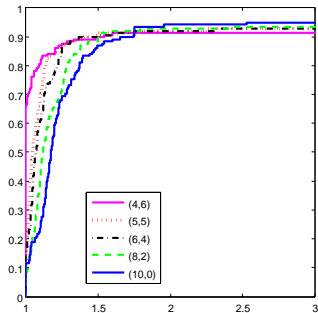
Efficiency performance profile, $r_{\text{size}}=0$



Note: rather insensitive to choice of l_{size} (ICFS).

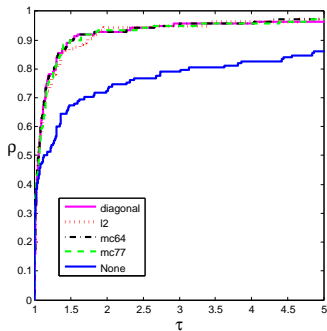
Efficiency (= iteration) performance profile, $lsize=5$ 

$rsize=-1$ is unlimited memory for R (not practical).

Efficiency performance profile $lsize+rsize$ constant

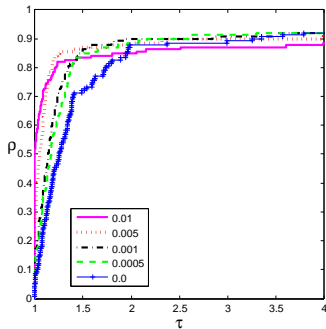
Pairs $(lsize+rsize)=10$

Intermediate memory ($rsize > 0$) can compensate for $lsize$.

Effect of scaling on efficiency ($lsize = rsize = 10$)

With scaling: 1 failure. No scaling: 10 failures.

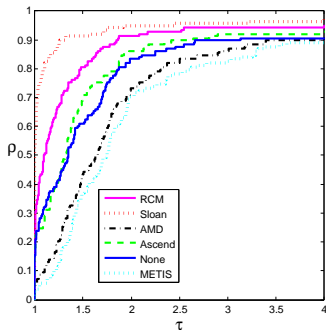
HSL_MI28 default is l_2 scaling.

Effect of dropping on efficiency ($lsize = rsize = 5$)

Often advantageous to use small drop tolerance.

Default $\tau_1 = 0.001$.

Effect of ordering on efficiency



Sloan profile-reduction ordering is the winner.

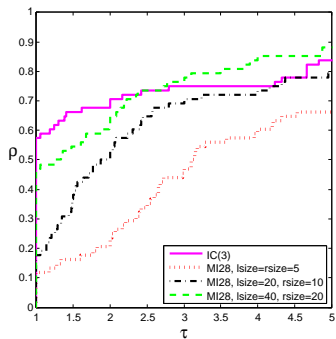
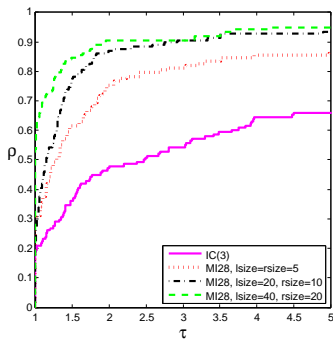
How many shifts? How to pick?

Problem	$\alpha_0 = 0.01$	$\alpha_0 = 0.001$	$\alpha_0 = 0.0001$
mcs01440	40 (4, $1.56 * 10^{-4}$)	40 (3, $2.50 * 10^{-4}$)	39 (2, $2.00 * 10^{-4}$)
ex33	427 (2, $1.00 * 10^{-2}$)	415 (3, $8.00 * 10^{-3}$)	460 (5, $1.28 * 10^{-2}$)
bcsstk17	112 (3, $2.50 * 10^{-3}$)	151 (3, $8.00 * 10^{-3}$)	136 (5, $6.40 * 10^{-3}$)
oilpan	817 (3, $2.50 * 10^{-3}$)	698 (2, $1.00 * 10^{-3}$)	708 (3, $1.60 * 10^{-3}$)
offshore	69 (2, $1.00 * 10^{-2}$)	71 (3, $1.60 * 10^{-2}$)	72 (4, $6.40 * 10^{-3}$)

- ▶ α_0 is first non-zero shift.
- ▶ How to pick? Too large: try and reduce.
Too small: must increase.
- ▶ Generally, want shift as **small** as possible

Comparison with level-based approach ($IC(3)$)

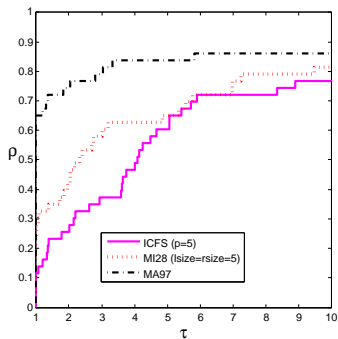
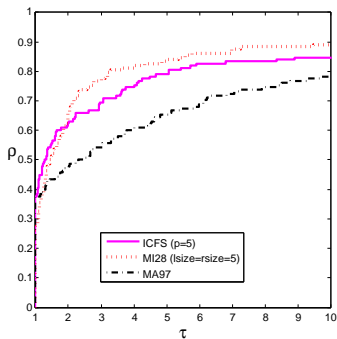
Efficiency (left) and iterations (right).



HSL_MI28 solved all problems; $IC(3)$ failed to give convergence for 19 problems

Comparison with direct solver HSL_MA97

Total time: all problems (left) and large problems (right).



HSL_MI28 can sometimes compete with direct solver
(and succeeds when HSL_MA97 runs out of memory).

Concluding remarks

- ▶ We have developed a new *IC* code **HSL_MI28** that may be used as a “black box” or tuned for a particular problem.
- ▶ Memory usage is under the user’s control.
- ▶ Using restricted **intermediate memory** improves efficiency.
- ▶ The intermediate memory can **compensate** for the preconditioner size.
- ▶ Based on extensive experimentation, HSL_MI28 appears **robust and efficient**.
- ▶ **Next step:** the indefinite case. **Pivoting challenge**



Thank you!

HSL_MI28 is available as part of HSL 2013.

Technical Reports RAL-P-2013-004 and RAL-P-2013-005
(to appear in SISC and TOMS)

Supported by EPSRC grant EP/I013067/1

Grant Agency of the Czech Republic Project No. P201/13-06684