

Direct solvers for Sparse Matrices

X. Li

Direct solvers for sparse matrices involve much more complicated algorithms than for dense matrices. The main complication is due to the need for efficient handling the *fill-in* in the factors L and U . A typical sparse solver consists of four distinct steps as opposed to two in the dense case:

1. An ordering step that reorders the rows and columns such that the factors suffer little fill, or that the matrix has special structure such as block triangular form.
2. An analysis step or symbolic factorization that determines the nonzero structures of the factors and create suitable data structures for the factors.
3. Numerical factorization that computes the L and U factors.
4. A solve step that performs forward and back substitution using the factors.

There is a vast variety of algorithms associated with each step. The review papers by Duff [9] (see also [8, Chapter 6]) and Heath et al. [18] can serve as excellent reference of various algorithms. Usually steps 1 and 2 involve only the graphs of the matrices, and hence only integer operations. Steps 3 and 4 involve floating-point operations. Step 3 is usually the most time-consuming part, whereas step 4 is about an order of magnitude faster. The algorithm used in step 1 is quite independent of that used in step 3. But the algorithm in step 2 is often closely related to that of step 3. In a solver for the simplest systems, i.e., symmetric and positive definite systems, the four steps can be well separated. For the most general unsymmetric systems, the solver may combine steps 2 and 3 (e.g. SuperLU) or even combine steps 1, 2 and 3 (e.g. UMFPACK) so that the numerical values also play a role in determining the elimination order.

In the past 10 years, many new algorithms and software have emerged which exploit new architectural features, such as memory hierarchy and parallelism. In Table 1, we compose a rather comprehensive list of sparse direct solvers. It is most convenient to organize the software in three categories: the software for serial machines, the software for SMPs, and the software for distributed memory parallel machines.

Code	Technique	Scope	Contact
<i>Serial platforms</i>			
MA27	Multifrontal	Sym	HSL [12]
MA41	Multifrontal	Sym-pat	HSL [1]
MA42	Frontal	Unsym	HSL [13]
MA47	Multifrontal	Sym	HSL [10]
MA48	Right-looking	Unsym	HSL [11]
SPARSE	Right-looking	Unsym	Kundert [21]
SPARSPAK	Left-looking	SPD	George [15]
SPOOLES	Left-looking	Sym and Sym-pat	Ashcraft [4]
SuperLLT	Left-looking	SPD	Ng [24]
SuperLU	Left-looking	Unsym	Li [6]
UMFPACK	Multifrontal	Unsym	Davis [5]
<i>Shared memory parallel machines</i>			
Cholesky	Left-looking	SPD	Rothberg [27]
DMF	Multifrontal	Sym	Lucas [23]
MA41	Multifrontal	Sym-pat	HSL [2]
PanelLLT	Left-looking	SPD	Ng [17]
PARASPAR	Right-looking	Unsym	Zlatev [28]
PARDISO	Left-right looking	Sym-pat	Schenk [26]
SPOOLES	Left-looking	Sym and Sym-pat	Ashcraft [4]
SuperLU_MT	Left-looking	Unsym	Li [7]
<i>Distributed memory parallel machines</i>			
CAPSS	Multifrontal	SPD	Raghavan [19]
DMF	Multifrontal	Sym	Lucas [23]
MUMPS	Multifrontal	Sym and Sym-pat	Amestoy [3]
PaStiX	Left-right looking*	SPD	CEA [20]
PSPASES	Multifrontal	SPD	Gupta [16]
SPOOLES	Left-looking	Sym and Sym-pat	Ashcraft [4]
SuperLU_DIST	Right-looking	Unsym	Li [22]
S+	Right-looking†	Unsym	Yang [14]

Table 1: Software to solve sparse linear systems using direct methods.

* In spite of the title of the paper

† Uses QR storage to statically accommodate any LU fill-in

Abbreviations used in the table:

SPD = symmetric and positive definite

Sym = symmetric and may be indefinite

Sym-pat = symmetric nonzero pattern but unsymmetric values

Unsym = unsymmetric

HSL = Harwell Subroutine Library:

<http://www.cse.clrc.ac.uk/Activity/HSL>

Fair to say, there is no single algorithm or software that is best for all types of linear systems. Some software is targeted for special matrices such as symmetric and positive definite, some is targeted for the most general cases. This is reflected in column 3 of the table, “Scope”. Even for the same scope, the software may decide to use a particular algorithm or implementation technique, which is better for certain applications but not for others. In column 2, “Technique”, we give a high level algorithmic description. For a review of the distinctions between left-looking, right-looking, and multifrontal and their implications on performance, we refer the reader to the papers by Heath et al. [18] and Rothberg [25]. Sometimes the best (or only) software is not in public domain, but available commercially or in research prototypes. This is reflected this in column 4, “Contact”, which could be the name of a company, or the name of the author of the research code.

In the context of shift-and-invert spectral transformation for eigensystem analysis, we need to factorize $A - \sigma I$, where A is fixed. Therefore, the nonzero structure of $A - \sigma I$ is fixed. Furthermore, for the same shift σ , it is common to solve many systems with the same matrix and different right-hand sides. (in which case the solve cost can be comparable to factorization cost.) It is reasonable to spend a little more time in steps 1 and 2 but speed up steps 3 and 4. That is, one can try different ordering schemes and estimate the costs of numerical factorization and solution based on symbolic factorization, and use the best ordering. For instance, in computing the SVD, one has the choice between shift-and-invert on AA^* , A^*A , and $\begin{bmatrix} 0 & A \\ A^* & 0 \end{bmatrix}$, all of which can have rather different factorization costs.

Some solvers have the ordering schemes built in, but others do not. It is also possible that the built-in ordering schemes are not the best for the target applications. It is sometimes better to substitute an external ordering scheme for the built-in one. Many solvers provide well-defined interfaces so that the user can make this substitution easily. One should read the solver documentation to see how to do this, as well as to find out the recommended ordering methods.

References

- [1] P. R. Amestoy and I. S. Duff. Vectorization of a multiprocessor multifrontal code. *Int. J. of Supercomputer Applics.*, 3:41–59, 1989.

- [2] Patrick R. Amestoy and Iain S. Duff. Memory management issues in sparse multifrontal methods on multiprocessors. *Int. J. Supercomputer Applics*, 7:64–82, 1993.
- [3] Patrick R. Amestoy, Iain S. Duff, Jean-Yves L'Excellent, and Jacko Koster. A fully asynchronous multifrontal solver using distributed dynamic scheduling. Technical Report RAL-TR-1999-059, Rutherford Appleton Laboratory, 1999. Also appeared as Report TR/PA/99/28, CERFACS, Toulouse, France. An earlier version appeared as Report RT/APO/99/02, ENSEEIHT-IRIT, Toulouse. Submitted to SIAM J Matrix Analysis and Applications. (<http://www.pallas.de/parasol>).
- [4] C. Ashcraft and R. Grimes. SPOOLES: An object-oriented sparse matrix library. In *Proceedings of the Ninth SIAM Conference on Parallel Processing*, 1999. (<http://www.netlib.org/linalg/spooles>).
- [5] T. A. Davis and I. S. Duff. A combined unifrontal/multifrontal method for unsymmetric sparse matrices. Technical Report TR-95-020, Computer and Information Sciences Department, University of Florida, 1995. (<http://www.netlib.org/linalg/umfpack2.2.tgz>).
- [6] James W. Demmel, Stanley C. Eisenstat, John R. Gilbert, Xiaoye S. Li, and Joseph W. H. Liu. A supernodal approach to sparse partial pivoting. *SIAM J. Matrix Anal. Appl.*, 20(3):720–755, 1999. (<http://www.nersc.gov/~xiaoye/SuperLU>).
- [7] James W. Demmel, John R. Gilbert, and Xiaoye S. Li. An asynchronous parallel supernodal algorithm for sparse Gaussian elimination. *SIAM J. Matrix Anal. Appl.*, 20(4):915–952, 1999. (<http://www.nersc.gov/~xiaoye/SuperLU>).
- [8] J. J. Dongarra, I. S. Duff, D. C. Sorensen, and H. A. van der Vorst. *Numerical Linear Algebra for High-Performance Computers*. SIAM, Philadelphia, PA, 1998.
- [9] Iain S. Duff. Direct methods. Technical Report RAL-98-056, Rutherford Appleton Laboratory, 1998.
- [10] I.S Duff and J. K. Reid. MA47, a Fortran code for direct solution of indefinite sparse symmetric linear systems. Technical Report RAL-95-001, Rutherford Appleton Laboratory, 1995.

- [11] I.S Duff and J. K. Reid. The design of MA48, a code for the direct solution of sparse unsymmetric linear systems of equations. *ACM Trans. Math. Software*, 22:187–226, 1996.
- [12] I.S Duff and J.K Reid. The multifrontal solution of indefinite sparse symmetric linear equations. *ACM Transactions on Mathematical Software*, 9(3):302–325, September 1983.
- [13] I.S Duff and J. A. Scott. The design of a new frontal code for solving sparse unsymmetric systems. *ACM Trans. Math. Software*, 22(1):30–45, 1996.
- [14] Cong Fu, Xiangmin Jiao, and Tao Yang. Efficient sparse LU factorization with partial pivoting on distributed memory architectures. *IEEE Trans. Parallel and Distributed Systems*, 9(2):109–125, 1998. (<http://www.cs.ucsb.edu/research/S+>).
- [15] A. George and J. Liu. *Computer Solution of Large Sparse Positive Definite Systems*. Prentice-Hall Inc., Englewood Cliffs, New Jersey, 1981. (jageorge@sparse1.uwaterloo.ca).
- [16] A. Gupta, G. Karypis, and V. Kumar. Highly scalable parallel algorithms for sparse matrix factorization. *IEEE Trans. Parallel and Distributed Systems*, 8:502–520, 1997. (<http://www.cs.umn.edu/~mjoshi/pspases>).
- [17] A. Gupta, E. Rothberg, E. Ng, and B. W. Peyton. Parallel sparse Cholesky factorization algorithms for shared-memory multiprocessor systems. In R. Vichnevetsky, D. Knight, and G. Richter, editors, *Advances in Computer Methods for Partial Differential Equations–VII*, pages 622–628. IMACS, 1992. (engng@lbl.gov).
- [18] M. Heath, E. Ng, and B. Peyton. Parallel algorithms for sparse linear systems. *SIAM Review*, 33:420–460, 1991.
- [19] Michael T. Heath and Padma Raghavan. Performance of a fully parallel sparse solver. *Int. J. Supercomputer Applications*, 11(1):49–64, 1997. (<http://www.netlib.org/scalapack>).
- [20] P. Henon, P. Ramet, and J. Roman. A mapping and scheduling algorithm for parallel sparse fan-in numerical factorization. In *EuroPar’99 Parallel Processing*, Lecture Notes in Computer Science, No. 1685, pages 1059–1067, Berlin, Heidelberg, New York, 1999. Springer-Verlag.

- [21] Kenneth Kundert. Sparse matrix techniques. In Albert Ruehli, editor, *Circuit Analysis, Simulation and Design*. North-Holland, 1986. (<http://www.netlib.org/sparse>).
- [22] Xiaoye S. Li and James W. Demmel. A scalable sparse direct solver using static pivoting. In *Proceedings of the Ninth SIAM Conference on Parallel Processing for Scientific Computing*, San Antonio, Texas, March 22–24 1999. (<http://www.nersc.gov/~xiaoye/SuperLU>).
- [23] Robert Lucas. Private communication (rflucas@lbl.gov), 2000.
- [24] Esmond G. Ng and Barry W. Peyton. Block sparse Cholesky algorithms on advanced uniprocessor computers. *SIAM J. Sci. Comput.*, 14(5):1034–1056, September 1993. (egng@lbl.gov).
- [25] E. Rothberg. *Exploiting the memory hierarchy in sequential and parallel sparse Cholesky factorization*. PhD thesis, Dept. of Computer Science, Stanford University, December 1992.
- [26] O. Schenk, K. Gärtner, and W. Fichtner. Efficient sparse LU factorization with left–right looking strategy on shared memory multiprocessors. *BIT*, 40(1):158–176, 2000.
- [27] J.P. Singh, W-D. Webber, and A. Gupta. Splash: Stanford parallel applications for shared-memory. *Computer Architecture News*, 20(1):5–44, 1992. (<http://www-flash.stanford.edu/apps/SPLASH>).
- [28] Z. Zlatev, J. Waśniewski, P. C. Hansen, and Tz. Ostromsky. PARASPAR: a package for the solution of large linear algebraic equations on parallel computers with shared memory. Technical Report 95-10, Technical University of Denmark, September 1995.