

# Direct Solvers for Sparse Matrices

X. Li

April 2022

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 [16] (see also [15, Chapter 6]) and Heath et al. [27] 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.

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. [27] and Rothberg [38]. 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  $\sigma$ , 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

Code	Technique	Scope	Contact
<i>Serial platforms (possibly on GPU)</i>			
CHOLMOD	Left-looking	SPD	Davis [8]
GLU3.0	Left-looking	Unsym (GPU)	Peng [36]
KLU	Left-looking	Unsym	Davis [11]
MA57	Multifrontal	Sym	HSL [19]
MA41	Multifrontal	Sym-pat	HSL [1]
MA42	Frontal	Unsym	HSL [20]
MA67	Multifrontal	Sym	HSL [17]
MA48	Right-looking	Unsym	HSL [18]
Oblio	Left/right/Multifr.	sym, Out-core	Dobrian [14]
SPARSE	Right-looking	Unsym	Kundert [32]
SPARSPAK	Left-looking	SPD, Unsym, QR	George et al. [22]
SPOOLES	Left-looking	Sym, Sym-pat, QR	Ashcraft [5]
SSIDS	Multifrontal	Sym (GPU)	Hogg [28]
SuperLLT	Left-looking	SPD	Ng [35]
SuperLU	Left-looking	Unsym	Li [12]
UMFPACK	Multifrontal	Unsym	Davis [9]
<i>Shared memory parallel machines (possibly on GPU)</i>			
BCSLIB-EXT	Multifrontal	Sym, Unsym, QR	Ashcraft et al. [6]
Cholesky	Left-looking	SPD	Rothberg [31]
MF2	Multifrontal	Sym, Sym-pat, Out-core (GPU)	Lucas [34]
MA41	Multifrontal	Sym-pat	HSL [2]
MA49	Multifrontal	QR	HSL [4]
PanelLLT	Left-looking	SPD	Ng [24]
PARASPAR	Right-looking	Unsym	Zlatev [41]
PARDISO	Left-Right looking	Sym-pat	Schenk [39]
SPOOLES	Left-looking	Sym, Sym-pat	Ashcraft [5]
SuiteSparseQR	Multifrontal	Rank-revealing QR	Davis [10]
SuperLU_MT	Left-looking	Unsym	Li [13]
TAUCS	Left/Multifr.	Sym, Unsym, Out-core	Toledo [7]
WSMP	Multifrontal	SPD, Unsym	Gupta [25]
<i>Distributed memory parallel machines</i>			
Clique	Multifrontal	Sym (no pivoting)	Poulson [37]
MF2	Multifrontal	Sym, Sym-pat, Out-core, GPU	Lucas [34]
DSCPACK	Multifrontal	SPD	Raghavan [26]
MUMPS	Multifrontal	Sym, Sym-pat	Amestoy [3]
PARDISO	Left-Right looking	Sym-pat, Unsym	Schenk [39]
PaStiX	Left-Right looking	SPD, Sym, Sym-pat	Ramet [29]
PSPASES	Multifrontal	SPD	Gupta [23]
SPOOLES	Left-looking	Sym, Sym-pat, QR	Ashcraft [5]
STRUMPACK	Multifrontal	Unsym, Sym-pat (GPU)	Ghysels [40]
SuperLU_DIST	Right-looking	Unsym (GPU)	Li [33]
symPACK	Left-Right looking	SPD	Jacquelin [30]
S+	Right-looking†	Unsym	Yang [21]
WSMP	Multifrontal	SPD, Unsym	Gupta [25]

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

† 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>

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. *The International Journal of Supercomputer Applications*, 3:41–59, 1989.
- [2] P. R. Amestoy and I. S. Duff. Memory management issues in sparse multifrontal methods on multiprocessors. *Int. J. Supercomputer Appl.*, 7(1):64–82, Spring 1993.
- [3] P. R. Amestoy, I. S. Duff, J.-Y. L’Excellent, and J. Koster. A fully asynchronous multifrontal solver using distributed dynamic scheduling. *SIAM Journal on Matrix Analysis and Applications*, 23(1):15–41, 2001.
- [4] P. R. Amestoy, I. S. Duff, and C. Puglisi. Multifrontal QR factorization in a multiprocessor environment. *Numer. Linear Algebra Appl.*, 3(4):275–300, 1996.
- [5] C. Ashcraft and R. G. Grimes. SPOOLES: An object oriented sparse matrix library. In *Proceedings of the Ninth SIAM Conference on Parallel Processing for Scientific Computing*, San Antonio, Texas, March 22–24, 1999. <http://www.netlib.org/linalg/spooles>.
- [6] BASLIB-EXT: Sparse matrix software. <http://www.aanalytics.com/products.htm>.
- [7] D. Chen, V. Rotkin, and S. Toledo. TAUCS: A Library of Sparse Linear Solvers. Tel-Aviv University. <http://www.tau.ac.il/~stoledo/taucs/>.
- [8] Y. Chen, T. A. Davis, W. W. Hager, and S. Rajamanickam. Algorithm 887: CHOLMOD, supernodal sparse Cholesky factorization and update/downdate. *ACM Trans. Mathematical Software*, 35(3), January 2009. <http://www.cise.ufl.edu/research/sparse/cholmod/>.
- [9] T. A. Davis. Algorithm 832: UMFPACK V4.3, an unsymmetric-pattern multifrontal method with a column pre-ordering strategy. *ACM Trans. Mathematical Software*, 30(2):196–199, 2004. <https://people.engr.tamu.edu/davis/suitesparse.html>.
- [10] T. A. Davis. Algorithm 915, SuiteSparseQR: Multifrontal multithreaded rank-revealing sparse QR factorization. *ACM Trans. Mathematical Software*, 38(1), 2011. <https://people.engr.tamu.edu/davis/suitesparse.html>.
- [11] T. A. Davis and E. Palamadai Natarajan. Algorithm 907: KLU, A Direct Sparse Solver for Circuit Simulation Problems. *ACM Trans. Mathematical Software*, 37(3), 2011. <https://people.engr.tamu.edu/davis/suitesparse.html>.
- [12] J. W. Demmel, S. C. Eisenstat, J. R. Gilbert, X. S. Li, and J. W. H. Liu. A supernodal approach to sparse partial pivoting. *SIAM J. Matrix Analysis and Applications*, 20(3):720–755, 1999. <https://portal.nersc.gov/project/sparse/superlu/>.

- [13] J. W. Demmel, J. R. Gilbert, and X. S. Li. An asynchronous parallel supernodal algorithm for sparse gaussian elimination. *SIAM J. Matrix Analysis and Applications*, 20(4):915–952, 1999. <https://portal.nersc.gov/project/sparse/superlu/>.
- [14] F. Dobrian and A. Pothen. Oblio: a sparse direct solver library for serial and parallel computations. Technical report, Old Dominion University, 2000.
- [15] J. J. Dongarra, I. S. Duff, D. C. Sorensen, and H. A. van der Vorst. *Numerical Linear algebra for high-performance computers*. SIAM, Philadelphia, 1998.
- [16] I. S. Duff. Direct methods. Technical Report RAL-98-054, Rutherford Appleton Laboratory, 1998.
- [17] 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.
- [18] 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. Mathematical Software*, 22:187–226, 1996.
- [19] I.S Duff and J.K Reid. The multifrontal solution of indefinite sparse symmetric linear equations. *ACM Trans. Mathematical Software*, 9(3):302–325, September 1983.
- [20] I.S Duff and J. A. Scott. The design of a new frontal code for solving sparse unsymmetric systems. *ACM Trans. Mathematical Software*, 22(1):30–45, 1996.
- [21] C. Fu, X. Jiao, and T. Yang. Efficient sparse LU factorization with partial pivoting on distributed memory architectures. *IEEE Trans. Parallel and Distributed Systems*, 9(2):109–125, 1998. <https://sites.cs.ucsb.edu/projects/s+/>.
- [22] A. George and J. W. H. Liu. *Computer Solution of Large Sparse Positive Definite Systems*. Prentice Hall, Englewood Cliffs, NJ, 1981.
- [23] A. Gupta, G. Karypis, and V. Kumar. Scalable parallel algorithms for sparse matrix factorization. *IEEE Trans. Parallel Distrib. Syst.*, 8(5):502–520, 1997. <http://glaros.dtc.umn.edu/gkhome/pspases/overview>.
- [24] 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, New Brunswick, NJ., 1992.
- [25] Anshul Gupta. WSMP: Watson Sparse Matrix Package. IBM T.J. Watson Research Center, Yorktown Heights. <http://www-users.cs.umn.edu/~agupta/wsmp.html>.
- [26] M. T. Heath and P. Raghavan. Performance of a fully parallel sparse solver. *Int. J. Supercomputer Applications*, 11(1):49–64, 1997. <http://www.cse.psu.edu/~raghavan>.
- [27] M.T. Heath, E. Ng., and B.W. Peyton. Parallel algorithms for sparse linear systems. *SIAM Review*, 33(3):420–460, September 1991.
- [28] J.D. Hogg, E. Ovtchinnikov, and J.A. scott. A Sparse Symmetric Indefinite Direct Solver for GPU Architectures. *ACM Trans. Mathematical Software*, 42(1), January 2016.

- [29] P. Hénon, P. Ramet, and J. Roman. PaStiX: A High-Performance Parallel Direct Solver for Sparse Symmetric Definite Systems. *Parallel Computing*, 28(2):301–321, 2002. <https://solverstack.gitlabpages.inria.fr/pastix/index.html>.
- [30] M. Jacquelin and E. Ng. Solver for sparse symmetric matrices. Lawrence Berkeley National Laboratory. <http://www.sympack.org>.
- [31] W-D. Webber J.P. Singh and A. Gupta. Splash: Stanford parallel applications for shared-memory. *Computer Architecture News*, 20(1):5–44, 1992.
- [32] Kenneth Kundert. Sparse matrix techniques. In Albert Ruehli, editor, *Circuit Analysis, Simulation and Design*. North-Holland, 1986. <https://arxiv.org/abs/1908.00204v3>.
- [33] X. S. Li and J. W. Demmel. SuperLU\_DIST: A scalable distributed-memory sparse direct solver for unsymmetric linear systems. *ACM Trans. Mathematical Software*, 29(2):110–140, June 2003. <https://portal.nersc.gov/project/sparse/superlu/>.
- [34] Robert Lucas, 2015. Private communication (rflucas@isi.edu).
- [35] 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.
- [36] S. Peng and S. X.-D. Tan. GLU3.0: Fast GPU-based Parallel Sparse LU Factorization for Circuit Simulation. IBM T.J. Watson Research Center, Yorktown Heights, 2020. <https://arxiv.org/abs/1908.00204v3>.
- [37] Jack Poulson. Clique: Sparse direct solver. <http://www-users.cs.umn.edu/~agupta/wsmp.html>.
- [38] Edward E. Rothberg. *Exploiting the memory hierarchy in sequential and parallel sparse Cholesky factorization*. PhD thesis, Dept. of Computer Science, Stanford University, December 1992.
- [39] 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. <https://pardiso-project.org/>.
- [40] STRUMPACK: STRUctured Matrices PACKages. <http://portal.nersc.gov/project/sparse/strumpack/>.
- [41] 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.