

***Performance Evaluation of Mixed Precision
Algorithms for Solving Sparse Linear Systems***

Zounon, Mawussi and Higham, Nicholas J. and Lucas,
Craig and Tisseur, Françoise

2020

MIMS EPrint: **2020.21**

Manchester Institute for Mathematical Sciences
School of Mathematics

The University of Manchester

Reports available from: <http://eprints.maths.manchester.ac.uk/>

And by contacting: The MIMS Secretary
School of Mathematics
The University of Manchester
Manchester, M13 9PL, UK

ISSN 1749-9097

Performance Evaluation of Mixed Precision Algorithms for Solving Sparse Linear Systems

MAWUSSI ZOUNON, Numerical Algorithms Group, UK
NICHOLAS J. HIGHAM, The University of Manchester, UK
CRAIG LUCAS, Numerical Algorithms Group, UK
FRANÇOISE TISSEUR, The University of Manchester, UK

It is well established that mixed precision algorithms that factorize a matrix at a precision lower than the working precision can reduce the execution time and the energy consumption of parallel solvers for dense linear systems. Much less is known about the efficiency of mixed precision parallel algorithms for sparse linear systems, and existing work focuses on single core experiments. We evaluate the benefits of using single precision arithmetic in solving a double precision sparse linear systems using multiple cores, focusing on the key components of LU factorization and matrix–vector products. We find that single precision sparse LU factorization is prone to a severe loss of performance due to the intrusion of subnormal numbers. We identify a mechanism that allows cascading fill-ins to generate subnormal numbers and show that automatically flushing subnormals to zero avoids the performance penalties. Our results show that the anticipated speedup of 2 over a double precision LU factorization is obtained only for the very largest of our test problems. For iterative solvers, we find that for the majority of the matrices computing or applying incomplete factorization preconditioners in single precision does not present sufficient performance benefits to justify the loss of accuracy compared with the use of double precision. We also find that using single precision for the matrix–vector product kernels provides an average speedup of 1.5 over double precision kernels, but new mixed precision algorithms are needed to exploit this benefit without losing the performance gain in the process of refining the solution to double precision accuracy.

CCS Concepts: • **Mathematics of computing** → **Mathematical software performance**; **Solvers**.

Additional Key Words and Phrases: Mixed precision algorithms, iterative refinement, sparse linear algebra, subnormal numbers

ACM Reference Format:

Mawussi Zounon, Nicholas J. Higham, Craig Lucas, and Françoise Tisseur. 2020. Performance Evaluation of Mixed Precision Algorithms for Solving Sparse Linear Systems. 1, 1 (September 2020), 18 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Ever since early versions of Fortran offered real and double precision data types, we have been able to choose between single and double precision floating-point arithmetics. Although single precision was no faster than double precision on most processors up to the early 2000s, on modern processors

Authors' addresses: Mawussi Zounon, Numerical Algorithms Group, Suite 21A, Manchester One, Portland Street, Manchester, M1 3LD, UK, mawussi.zounon@nag.co.uk; Nicholas J. Higham, The University of Manchester, Department of Mathematics, Manchester, M13 9PL, UK, nick.higham@manchester.ac.uk; Craig Lucas, Numerical Algorithms Group, Suite 21A, Manchester One, Portland Street, Manchester, M1 3LD, UK, craig.lucas@nag.co.uk; Françoise Tisseur, The University of Manchester, Department of Mathematics, Manchester, M13 9PL, UK, francoise.tisseur@manchester.ac.uk.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

XXXX-XXXX/2020/9-ART \$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

it executes twice as fast as double precision and has the additional benefit of halving the data movement. As a result, single precision (as well as half precision) is starting to be used in applications such as weather and climate modelling [Dawson et al. 2018], [Vána et al. 2017] and seismic modeling [Fabien-Ouellet 2020], where traditionally double precision was used. Mixed precision algorithms, which use some combination of half, single, double, and perhaps even quadruple precisions, are increasingly being developed and used in high performance computing [Abdelfattah et al. 2020].

In 2006, Buttari et al. [Buttari et al. 2007], [Langou et al. 2006] drew the attention of the HPC community to the potential of mixed precision iterative refinement algorithms for solving dense linear systems with unprecedented efficiency. The underlying principle is to carry out the most expensive part of the computation, the LU factorization or Cholesky factorization, in single precision instead of double precision (the working precision) and then refine the initial computed solution using residuals computed in double precision. This contrasts with traditional iterative refinement, in which only a precision higher than the working precision is used. The resulting algorithms are now implemented in LAPACK [Anderson et al. 1999] (as DSGETRS, and DSPOTRS for general and symmetric positive definite problems, respectively), and are generally twice as fast as a full double precision solve for sufficiently well conditioned matrices.

A decade after the two-precision iterative refinement work by Buttari et al., Carson and Higham introduced a GMRES-based iterative refinement algorithm that uses up to three precisions for the solution of linear systems [Carson and Higham 2017], [Carson and Higham 2018]. This algorithm enabled Haidar et al. [Haidar et al. 2018a], [Haidar et al. 2020], [Haidar et al. 2018b] to successfully exploit the half-precision floating-point arithmetic units of NVIDIA tensor cores in the solution of linear systems. Compared with linear solvers using exclusively double precision, their implementation shows up to a 4x–5x speedup while still delivering double precision accuracy [Haidar et al. 2020], [Haidar et al. 2018b]. This algorithm is now implemented in the MAGMA library [Agullo et al. 2009], [Magma [n.d.]] (routine `magma_dhgesv_iterref_gpu`) and in `cuSOLVER`, the NVIDIA library that provides LAPACK-like routines (routine `cusolverDnDhgesv`).

Mixed precision iterative refinement algorithms can be straightforwardly applied to parallel sparse direct solvers. But the variability of sparse matrix patterns and the complexity of sparse direct solvers make the estimation of the performance speedup difficult to predict. The primary aim of this work is to design a performance benchmark that will provide insight into the speedup of mixed precision parallel solvers of sparse linear systems.

The rest of this paper is structured as follows. We discuss existing work and the need for new studies in Section 2. Section 3 presents experimental settings, including details of the sparse matrices and hardware selected. Section 4 introduces the issue of appearance of subnormal numbers in single precision sparse LU factorization, explains how the subnormal numbers can be generated and proposes different mitigation strategies. In Section 5, we discuss our approach for assessing the performance gain of using reduced precision in parallel sparse direct solvers and give experimental results. Section 6 is dedicated to a similar study for iterative methods. In Section 7, we provide an advanced performance profiling to explain why mixed precision iterative refinement algorithms show lesser performance gains for sparse matrices than for dense matrices. Concluding remarks are given in Section 8.

2 DISCUSSION OF EXISTING STUDIES

The performance benefits of mixed precision iterative refinement have been widely demonstrated for dense linear systems. The few such performance studies for sparse linear systems are summarized below, with an emphasis on the performance metrics reported.

2.1 Mixed Precision Iterative Refinement for Sparse Direct Solvers

In 2008, [Buttari et al. \[2008\]](#) studied the performance of mixed precision iterative refinement algorithms for sparse linear systems. They used Algorithm 1, in which the precision that each line should be executed in is shown at the end of the line, with FP32 denoting single precision and FP64 double precision. To implement Algorithm 1 they selected two existing sparse direct solvers: a multifrontal sparse direct solver MUMPS, by Amestoy et al. [[Amestoy et al. 2000](#)] and a supernodal sparse direct solver SuperLU, by [Li and Demmel \[2003\]](#). Multifrontal and supernodal methods are the two main variants of sparse direct methods; for a full description and a performance comparison see [[Amestoy et al. 2001](#)].

Algorithm 1 Mixed-precision iterative refinement. Given a sparse matrix $A \in \mathbb{R}^{n \times n}$, and a vector $b \in \mathbb{R}^n$, this algorithm solves $Ax = b$ using a single precision sparse LU factorization of A then refines x to double precision accuracy.

- 1: Carry out the reordering and analysis for A .
 - 2: $LU \leftarrow \text{sparse_lu}(A)$ ▷ (FP32)
 - 3: Solve $Ax = b$ using the LU factors. ▷ (FP32)
 - 4: **while** not converged **do**
 - 5: $r \leftarrow b - Ax$ ▷ (FP64)
 - 6: Solve $Ad = r$ using the LU factors. ▷ (FP32)
 - 7: $x \leftarrow x + d$ ▷ (FP64)
 - 8: **end while**
-

Buttari et al. showed that the version of SuperLU used in their study does not benefit from using low-precision arithmetic. Put differently, the time spent in matrix factorization, which is the most time-consuming part of the algorithm, is hardly reduced when single precision arithmetic is used in place of double precision. They concluded that a mixed precision iterative refinement based on SuperLU would be no faster than the standard double precision algorithm.

For MUMPS, their experimental results showed that the mixed precision version can be up to two times faster than the standard double precision MUMPS. While this result is consistent with the performance observed for dense linear systems, there is an important difference to point out here: all the experimental results in [[Buttari et al. 2008](#)] were obtained using a single core.

In 2010, [Hogg and Scott \[2010\]](#) designed a mixed precision iterative solver for the solution of sparse symmetric linear systems. The algorithm is similar to Algorithm 1, except they perform LDL^T factorization instead of LU factorization and they also considered flexible GMRES [[Saad 1993](#)] for the refinement process. Their experimental results show that the advantage of mixed precision is limited to very large problems, where the computation time can be reduced up to a factor of two. But the results of this study are again based on single core benchmarks and also involve out-of-core techniques.

As these existing works are limited to a single core, further study is required to evaluate how the performance will be affected in fully-featured parallel sparse direct solvers using many cores. The main objective of using single precision arithmetic in sparse direct solvers is to reduce the time to solution. A safe way to improve performance without risking accuracy loss or inducing numerical stability is by exploiting the thread-level parallelism available in modern multicore processors. It is then sensible to first take advantage of core parallelism before using mixed precision algorithms for further performance enhancement. We aim to provide new insights into how far the mixed precision algorithms can advance the performance of parallel sparse solvers.

2.2 Mixed Precision Methods for Iterative Solvers

Here we summarize studies that use mixed precision arithmetic to improve the performance of iterative solvers. The existing works can be classified in three categories.

The first approach consists of using a single precision preconditioner or a few steps of a single precision iterative scheme as a preconditioner in a double precision iterative method. [Buttari et al. \[2008\]](#) have demonstrated the performance potential of this method using a collection of five sparse matrices, with a speedup ranging from 1.5x to 2.x. But the experiment has been performed on a single core using a diagonal preconditioner with an unvectorized sparse matrix–vector multiplication (SpMV) kernel.

The second approach, proposed by [Anzt et al. \[2019\]](#), uses low precision data storage whenever possible to accelerate data movement while performing all the computation in high precision. This concept is appealing, but hard to implement in practice as it requires an optimized data conversion routine and knowledge of key numerical properties of the matrices, such as the condition number. To illustrate this idea the authors of [\[Anzt et al. 2019\]](#) designed a mixed precision block-Jacobi preconditioning method where the explicit inversion of the block diagonals is required.

The third category consists of studies that focus on designing a mixed precision SpMV kernel for iterative solvers. This approach has been implemented by [Ahmad et al. \[2019\]](#) by proposing a new sparse matrix format that stores selected entries of the input matrix in single precision and the remainder in double precision. Their algorithm accelerates data movement and computation with a small accuracy loss compared with double precision SpMV. Their implementation demonstrates up to 2x speedup in the best case, but hardly achieves any speedup on most of the matrices due to data format conversion overhead. A similar approach has been implemented by [Grigoraş et al. \[2016\]](#) with a better speedup for FPGA architectures.

Our contribution is to assess the benefit of using mixed precision in iterative solvers from a practical point of view by evaluating optimized vendor kernels used in applications.

3 EXPERIMENTAL SETUP

We consider three different processors for the benchmarks: the AMD dual-socket EPYC Naples system with 64 cores, the Intel dual-socket Haswell with 20 cores, and the Intel dual-socket Skylake with 40 cores. The hardware also include two NVIDIA GPUs, the V100 and the P100. However in this paper, we report only the results from the Intel Skylake and the NVIDIA V100 GPUs because the results from others architectures are similar.

The selected sparse matrices for the benchmark are from various scientific and engineering applications and are summarized in [Table 1](#). The Intel Skylake node has 50 gigabytes of main memory, and consequently sparse matrices whose factors require more than 50 gigabytes storage are not included. The matrices are divided in two groups. The first 21 matrices are from the medium size group with 700, 000 to 5, 000, 000 nonzero elements. It takes few seconds on average to factorize these matrices. The second group contains larger matrices with 7,000,000 to 64,000,000 nonzeros and it takes on average a few minutes to factorize most of the matrices in this group. For each matrix, the largest absolute value $\max_{i,j} |a_{ij}|$ and the smallest nonzero absolute value $\min_{i,j} \{ |a_{ij}| : a_{ij} \neq 0 \}$ of the elements are reported. For medium size matrices, an estimate for the 1-norm condition number, $\kappa_1(A) = \|A^{-1}\|_1 \|A\|_1$, computed using the MATLAB `condst` routine, is also provided.

4 APPEARANCE OF SUBNORMAL NUMBERS IN SINGLE PRECISION SPARSE LU AND MITIGATION TECHNIQUES

From [Table 1](#), one can observe that the entries of the matrices fit in the range of single precision arithmetic, which from [Table 2](#) we see comprises numbers of modulus roughly between 10^{-45} and

Table 1. Selected matrices from the SuiteSparse Matrix Collection [Davis [n.d.]], [Davis and Hu 2011]. The first 21 matrices are of medium size and each can be factorized in a few seconds. Matrices 22 to 36 are larger and require more time and memory to solve.

	Matrix	Size	nnz	$\kappa_1(A)$	$\max_{i,j} a_{ij} $	$\min_{i,j} \{ a_{ij} : a_{ij} \neq 0\}$
1	2cubes_sphere	101,492	1,647,264	2.93e+09	2.52e+10	6.68e-15
2	ASIC_320ks	321,67	1,316,085	5.06e+22	1.00e+06	1.26e-39
3	Baumann	112,211	748,331	1.368+09	1.29e+04	5.00e-02
4	cf2	123,440	3,085,406	3.66e+06	1.00e+00	6.66e-09
5	crashbasis	160,000	1,750,416	1.78e+03	4.08e+02	6.42e-11
6	ct20stif	52,329	2,600,295	2.22e+14	8.86e+11	3.02e-34
7	dc1	116,835	861,071	1.01e+10	5.67e+4	3.00e-12
8	Dubcova3	146,689	3,636,643	1.14e+04	2.66e+00	8.47e-22
9	ecology2	999,999	4,995,991	6.66e+07	4.00e+01	1.00e+01
10	FEM_3D_thermal2	147,900	3,489,300	1.66e+03	2.92e-01	1.16e-05
11	G2_circuit	150,102	726,674	1.97e+07	2.22e+04	3.27e-01
12	Goodwin_095	100,037	3,226,066	3.43e+07	1.00e+00	1.41e-21
13	matrix-new_3	125,329	893,984	3.47e+22	1.00e+00	1.27e-21
14	offshore	259,789	4,242,673	2.32e+13	7.47e+14	7.19e-21
15	para-10	155,924	2,094,873	8.13e+18	6.44e+11	2.26e-20
16	parabolic_fem	525,825	3,674,625	2.11e+05	4.00e-01	3.18e-07
17	ss1	205,282	845,089	1.29e+01	1.00e+00	1.06e-11
18	stomach	213,360	3,021,648	8.01e+1	1.38e+00	1.47e-09
19	thermomech_TK	102,158	711,558	1.62e+20	1.96e+02	4.83e-03
20	tmt_unsym	917,825	4,584,801	2.26e+09	4.00e+00	1.00e+00
21	xenon2	157,464	3,866,688	1.76e+05	3.17e+28	5.43e+23
22	af_shell10	1,508,065	52,259,885		5.72e+05	1.00e-06
23	af_shell12	504,855	17,588,875		1.51e+06	4.55e-13
24	atmosmodd	1,270,432	8,814,880		2.22e+04	3.19e+03
25	atmosmodl	1,489,752	10,319,760		7.80e+04	3.96e+04
26	cage13	445,315	7,479,343		9.31e-01	1.15e-02
27	CurlCurl_2	806,529	8,921,789		4.42e+10	8.84e+06
28	dielFilterV2real	1,157,456	48,538,952		6.14e+01	3.25e-13
29	Geo_1438	1,437,960	60,236,322		6.69e+12	4.75e-07
30	Hook_1498	1,498,023	59,374,451		1.58e+05	5.17e-26
31	ML_Laplace	377,002	27,689,972		1.22e+07	1.24e-09
32	nlpkkt80	1,062,400	28,192,672		2.00e+02	4.08e-01
33	Serena	1,391,349	64,131,971		5.51e+13	2.19e-01
34	Si87H76	240,369	10,661,631		1.83e+01	2.57e-13
35	StocF-1465	1,465,137	21,005,389		3.10e+11	9.57e-09
36	Transport	1,602,111	23,487,281		1.00e+00	1.62e-12

10^{38} . There is no risk of underflow or overflow in converting these matrices to single precision format. However, the smallest absolute value of matrix ASIC_320ks, 1.26×10^{-39} , is a subnormal number in single precision. A subnormal floating-point number is a nonzero number with magnitude less than the absolute value of the smallest normalized number [Higham 2002, Chap. 2], [Muller

Table 2. Parameters for IEEE single and double precision point arithmetic. $x_{\min,s}$ is the smallest nonzero subnormal number and x_{\min} and x_{\max} are the smallest and largest normalized floating-point numbers.

	$x_{\min,s}$	x_{\min}	x_{\max}	Unit roundoff
FP32	1.4×10^{-45}	1.2×10^{-38}	3.4×10^{38}	6.0×10^{-8}
FP64	4.9×10^{-324}	2.2×10^{-308}	1.8×10^{308}	1.1×10^{-16}

et al. 2018, Chap. 2]. Floating-point operations on subnormals can be very slow because they are usually processed at the software level, which induces a high overhead.

The risk of underflow, overflow or generating subnormal numbers during the conversion from higher precision to lower precision can be reduced using scaling techniques proposed by Higham et al. [2019]. However, even if matrices have been safely converted from double to normalized single precision numbers, subnormal numbers may still be generated during the computation. We first suspected this behavior in our benchmark when some single precision computations took significantly more time than the corresponding double precision computations. For example, the sparse direct solver MUMPS computed the double precision LU decomposition of the matrix Baumann (#3 in Table 1) in 1.6251 seconds, while the single precision factorization took 3.586 seconds. Instead of being two times faster than the double precision computation, the single precision computation is two times slower. A further analysis reveals that the smallest magnitude entries of the single precision factors L and U are of the order of 10^{-88} , which is a subnormal number in single precision but a normalized number in double precision. The appearance of subnormal numbers in the single precision factors may be surprising since the absolute values of the entries of this matrix range from 5×10^{-2} to 1.29×10^4 , which appears to be innocuous for single precision.

This phenomenon of LU factorization generating subnormal numbers does not appear to have been observed before. How can it happen? The elements at the $(k+1)$ st stage of Gaussian elimination are generated from the formula

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik}a_{kj}^{(k)}, \quad m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}$$

where m_{ik} is a multiplier. If A is a dense matrix of normalized floating-point numbers with norm of order 1, it is extremely unlikely that any of the $a_{ij}^{(k)}$ will become subnormal. However, for sparse matrices we can identify a mechanism whereby fill-in cascades down a column and small multipliers combine multiplicatively. Consider the upper Hessenberg matrix

$$A = \begin{bmatrix} d_1 & 0 & \dots & \dots & 0 & 1 \\ -a_1 & d_2 & 0 & \dots & 0 & 0 \\ & -a_2 & d_3 & 0 & \dots & \vdots \\ & & -a_3 & d_4 & \ddots & \vdots \\ & & & \ddots & \ddots & 0 \\ & & & & -a_{n-1} & d_n \end{bmatrix}.$$

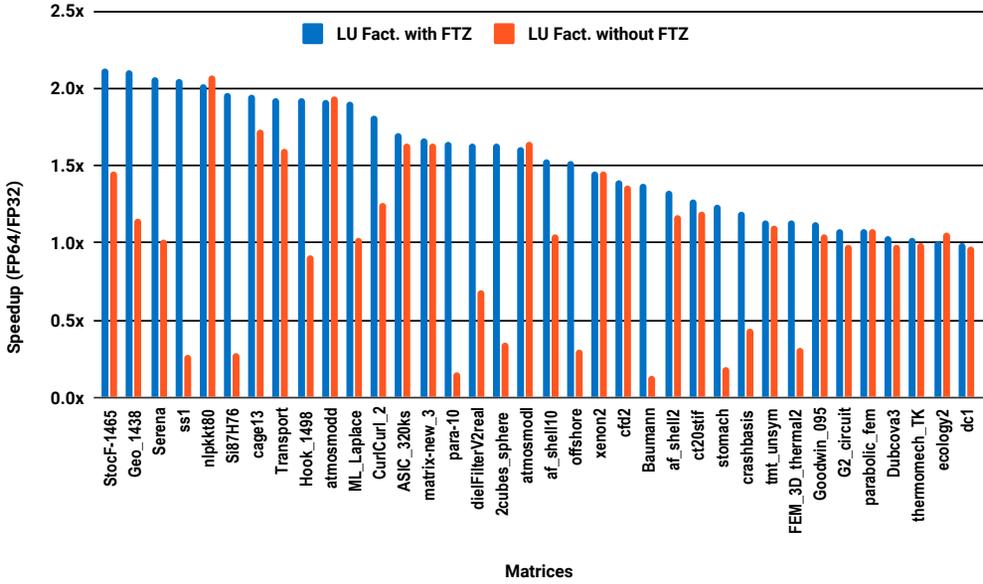


Fig. 1. Single precision speedup over double precision for sparse LU factorization using PARDISO on a single Intel Skylake core.

to decide whether low precision is beneficial. Note that in the case of dense linear systems, the factorization step speedup is usually close to 2x.

In addition to SuperLU and MUMPS, we have added PARDISO [Schenk et al. 2001], which is available in the Intel Math Kernel Library (MKL), to the set of sparse direct solvers for the benchmarks. PARDISO combines left- and right-looking level 3 BLAS supernodal algorithms for better parallelism. The solvers also include the multithreaded version of SuperLU, called SuperLU_MT [Li 2005]. We will refer to both packages as SuperLU unless there is ambiguity. We also considered adding UMFPACK [Davis 2004], but this package does not have support for single precision.

For each sparse direct solver, we report the factorization speedup for both sequential and parallel runs. Even though the Intel Skylake has 40 cores, we report parallel results with 10 cores as for most of the experiments the performance stagnates and sometimes declines beyond 10 cores. To stress the performance penalty induced by subnormals in the single precision computations, the results with and without FTZ are reported.

The experimental results with serial PARDISO are summarized in Figure 1. For each matrix two bars are shown, which give the speedup for LU factorization with and without FTZ. Without FTZ, up to 12 matrices out of 36 show a speedup below 1. In other words, single precision decreases the performance for 30% of the problems compared with double precision. This anomaly is corrected by flushing subnormals to zero. By comparing the results with FTZ with results without FTZ, we see that more than half of the problems generated subnormals during the single precision computation. As for the performance benefit of using single precision for the matrix factorization, half of the matrices show a speedup above the 1.5x threshold. The matrices that did not exceed 1.5x speedup are predominately of medium size. The parallel results in Figure 2 show that with 10 cores the proportion of problems that reach 1.5x speedup drops from 50% to 30%. The problems that still reach 1.5x speedup with 10 cores are exclusively from the large matrices and represent 65% of them.

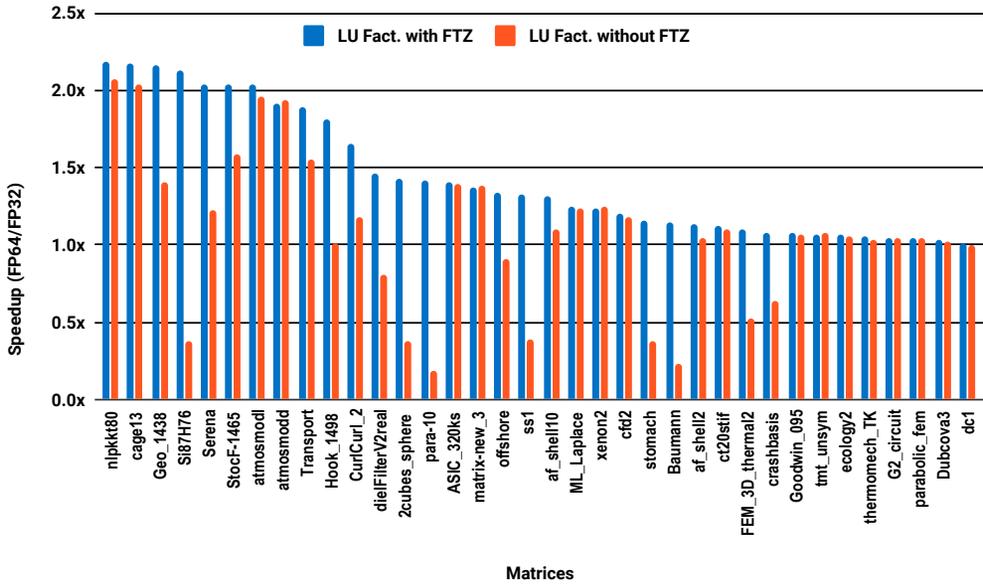


Fig. 2. Single precision speedup over double precision for sparse LU factorization using PARDISO on 10 Intel Skylake cores.

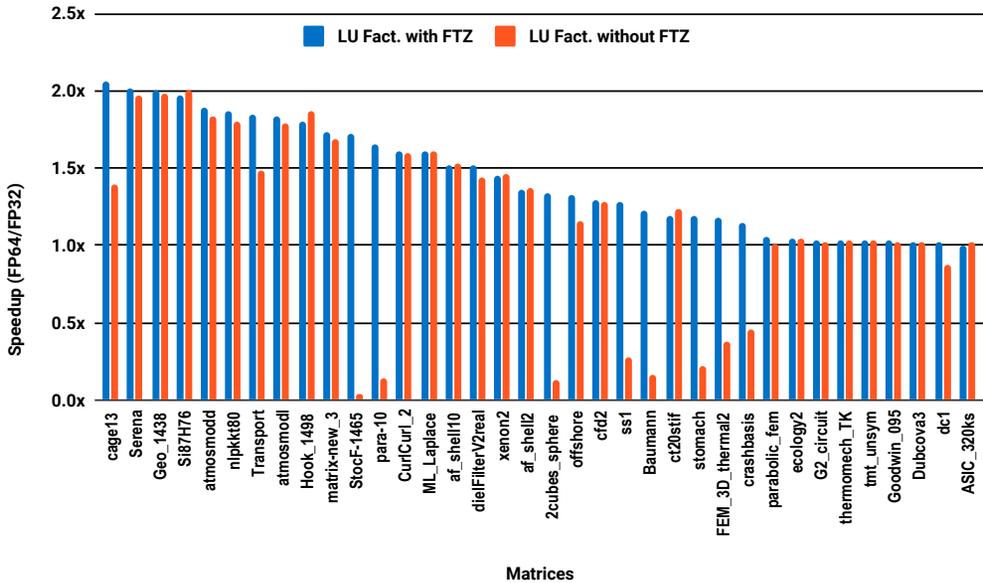


Fig. 3. Single precision speedup over double precision for sparse LU factorization using MUMPS on a single Intel Skylake core.

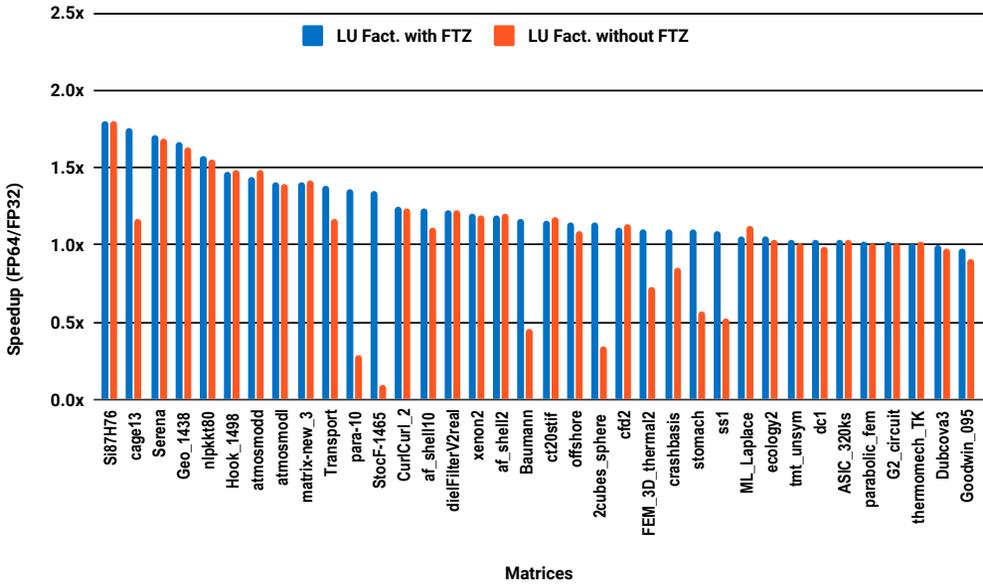


Fig. 4. Single precision speedup over double precision for sparse LU factorization using MUMPS on 10 Intel Skylake cores.

The results for serial MUMPS are summarized in Figure 3. The matrices that suffered performance degradation due to subnormals in the PARDISO experiments exhibit similar behavior with MUMPS. Similarly, half of the matrices did not reach the threshold of 1.5x, and the matrices beyond 1.5x are mainly the large ones. The parallel results in Figure 4 are less attractive as only five matrices deliver a speedup beyond 1.5x. These matrices are from the large size group.

Unlike PARDISO and MUMPS, the multithreaded SuperLU ran out of memory for 15 problems out of the 36, predominantly the large size ones. Results are reported for only the 21 remaining matrices. The serial results in Figure 5 show that only 33% of the 21 problems, successfully solved exceed 1.5x speedup, against 24% for the parallel results in Figure 6.

These results show that mixed precision iterative refinement may only be beneficial for large sparse matrices. However, a large matrix size and higher density are not enough to predict the speedup, as matrix `diefFilterV2real` is much larger and denser than `cage13` but its speedup is lower than `cage13`'s speedup in all the experiments. We note the contrast with dense linear systems, where a 2x speedup is often achieved even for matrices of size as small as 200×200 .

6 PERFORMANCE OF MIXED PRECISION SPARSE ITERATIVE SOLVERS

The performance of an iterative solver depends not only on the algorithm implemented but also on the eigenvalue distribution and condition number of the matrix, the choice of preconditioner, and the accuracy targeted. It is therefore hard to make general statements about how mixed precision techniques will affect the performance of an iterative solver. Therefore in this section we focus instead on analyzing the impact of low precision in sparse matrix–vector multiplication kernels and preconditioners, as they are the main building blocks of iterative solvers.

The two classes of preconditioning techniques commonly used to accelerate iterative solvers are incomplete factorizations and iterative schemes such as algebraic multigrid preconditioners.

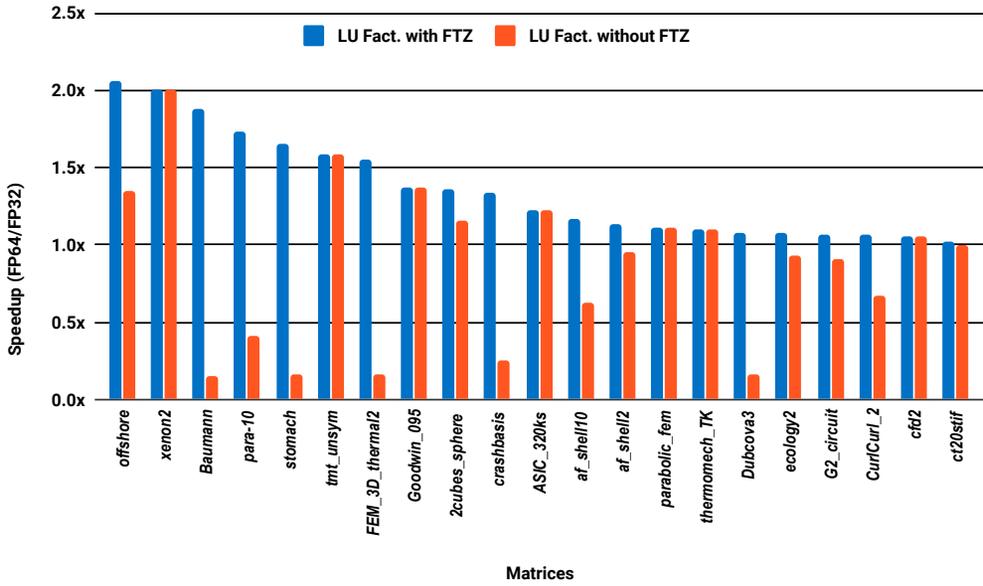


Fig. 5. Single precision speedup over double precision of sparse LU factorization using SuperLU on a single Intel Skylake core. SuperLU ran out of memory for 15 problems.

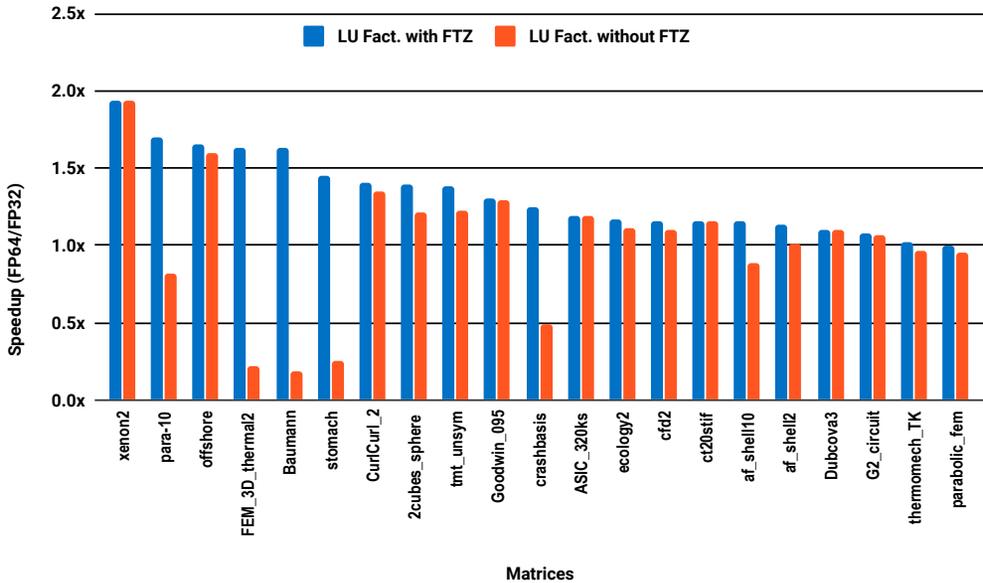


Fig. 6. Single precision speedup over double precision for sparse LU factorization using SuperLU on a 10 Intel Skylake cores. SuperLU ran out of memory for 15 problems.

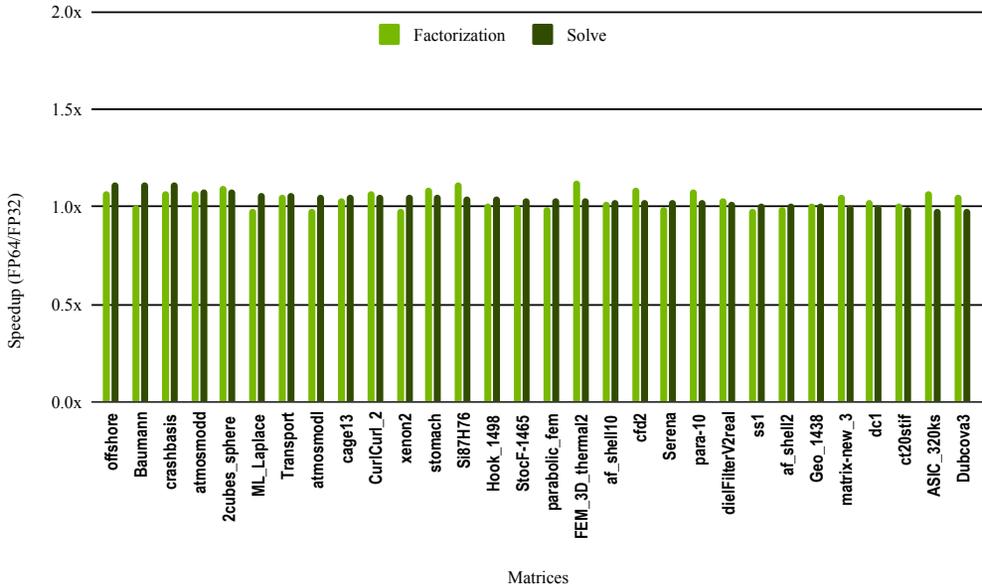


Fig. 7. Speedup of single precision versus double precision for sparse incomplete LU factorization (ILU0) using cuSPARSE on NVIDIA V100 GPU.

Iterative preconditioners are dominated by SpMV kernels, so in this section we focus on the impact of low precision arithmetic on the performance of incomplete factorization and SpMV kernels.

The results in Figure 7 illustrate the speedup from using single precision incomplete LU factorization (ILU0) from the cuSPARSE³ library on an NVIDIA V100 GPU. The cuSPARSE library provides an optimized implementation of a set of sparse linear algebra routines for NVIDIA GPUs. For the sake of readability, the matrices are sorted in a decreasing order of the solve step speedup.

The most critical part of the preconditioner application is the forward and backward solve, because it is executed at each iteration and can easily become the most time consuming part of iterative solvers. The dark green bars in Figure 7 represent the speedup of the single precision ILU0 preconditioner application. The performance shows that lowering the precision in the preconditioner application did not enhance the performance. The same is true for the incomplete factorization itself, so there is no benefit to using single precision in place of double precision. The results from SuperLU ILU in Figure 8 show a better speedup for the solve step compared with the results from cuSPARSE ILU0. However, the speedup is still under the threshold of 1.5x speedup, except for one matrix (Transport). For the incomplete LU factorization step itself, the performance gain from using single precision is insignificant. As the factorization step is more time consuming than the solve steps, the overall speedup of the preconditioner computation and application remains very small and does not seem to present enough potential to accelerate parallel iterative solvers. Note that from the three libraries evaluated in this work only SuperLU provides incomplete LU factorization for preconditioning.

To evaluate how low precision can accelerate SpMV kernels, we have considered the compressed row storage (CSR) format, as it is widely used in applications. In the CSR format, a double precision sparse matrix with nnz nonzero elements requires approximately $12nnz$ bytes for the storage

³<https://docs.nvidia.com/cuda/cusparse>

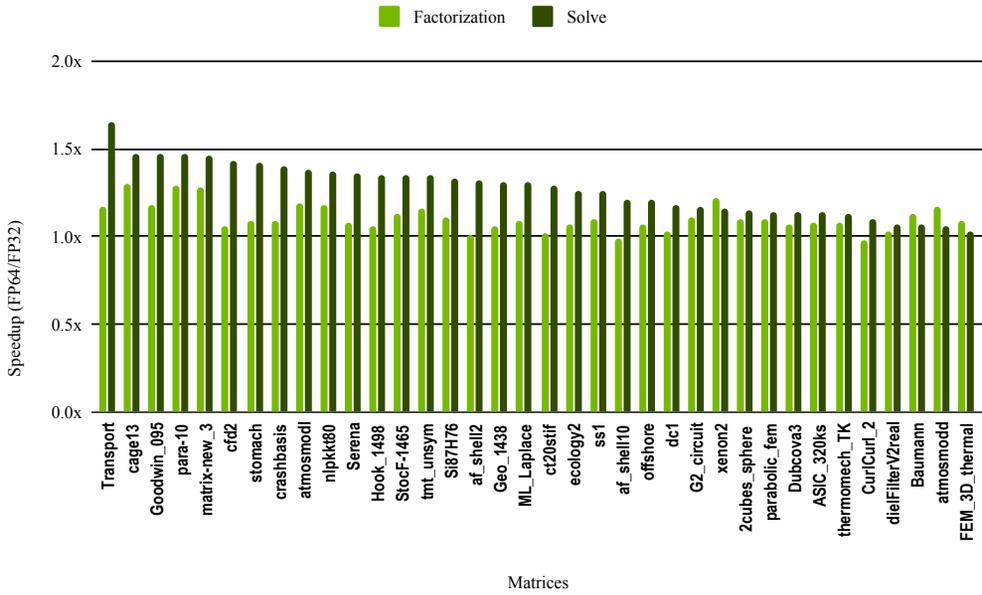


Fig. 8. Speedup of single precision versus double precision for sparse Incomplete LU factorization (ILU) using SuperLU Intel Skylake. The SuperLU ILU implementation is serial but it has been compiled against a multithreaded MKL BLAS and run with 10 cores.

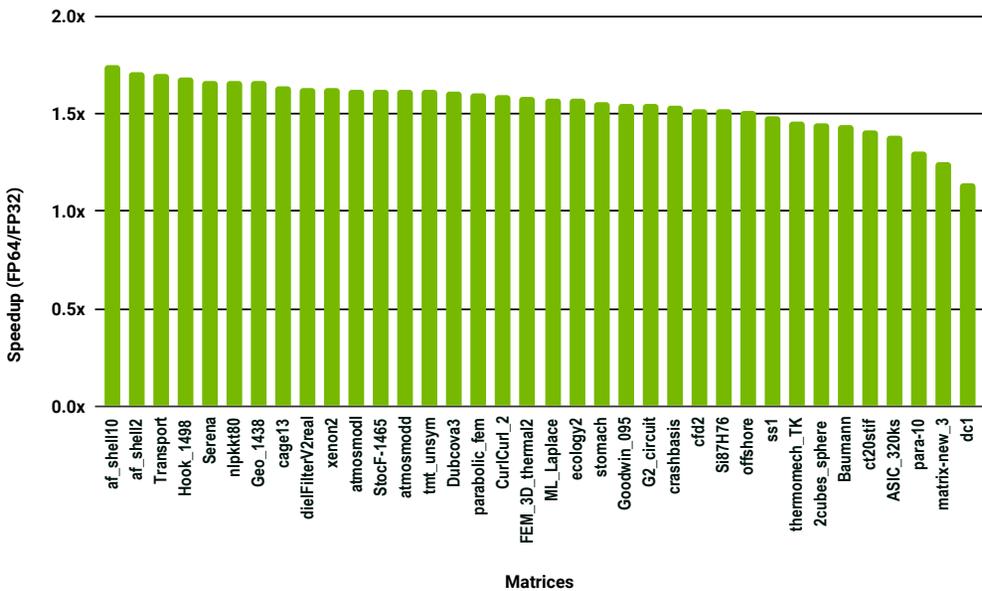


Fig. 9. Speedup of single precision versus double precision for SpMV using cuSPARSE on NVIDIA V100 GPU.

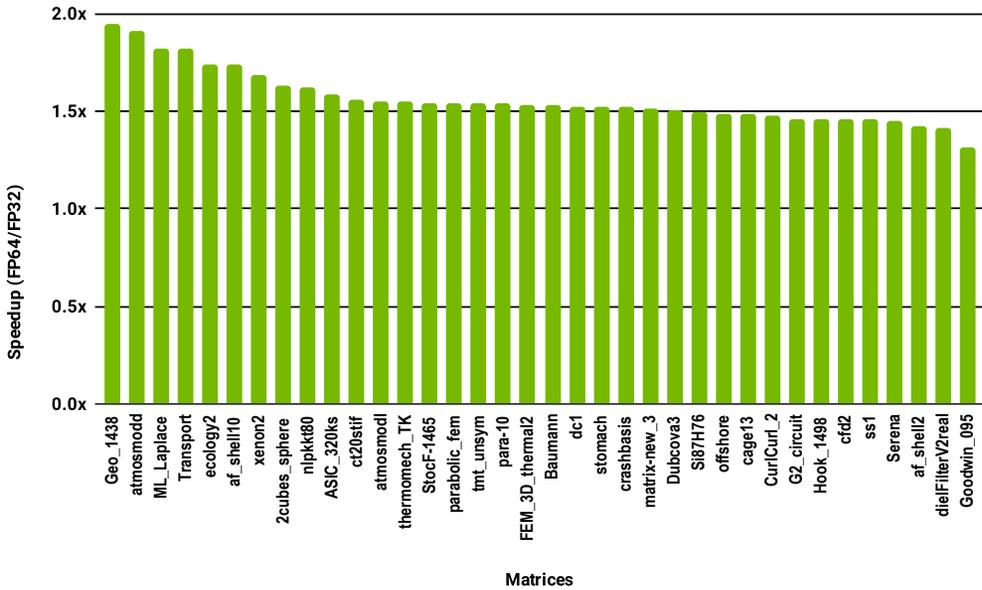


Fig. 10. Speedup of single precision versus double precision for SpMV using MKL on 10 Intel Skylake cores.

(each nonzero element requires 8 bytes for its value and 4 bytes for its column index). In single precision the matrix will occupy approximately $8nnz$ bytes of memory. As SpMV kernels are memory bandwidth bound, the use of single precision will only provide a $1.5x$ ($12nnz$ divided by $8nnz$) speedup in theory. Note that, for simplicity we have ignored the $4n$ bytes for row indices, where n is the number of rows, and the extra memory for left- and right-hand side vectors. The results in Figure 9 for the optimized cuSPARSE SpMV on the NVIDIA V100 GPU show that the speedup is oscillating around $1.5x$. Similarly, the benchmark of the MKL SpMV in Figure 10 shows that the single precision kernel has approximately $1.5x$ speedup over the double precision kernel.

While more experiments and analysis may be necessary to fully understand the benefit of low precision for incomplete factorization, this study shows that applying the preconditioner in low precision does not offer enough performance advantage to take the risk of lowering the computation accuracy. A mixed precision iterative preconditioner may be accelerated by taking advantage of efficient single precision SpMV kernels, but the overall speedup might be far less than $1.5x$ if a double precision solution accuracy is expected, because of the need to refine the solution.

7 FURTHER PERFORMANCE ANALYSIS

Apart from the unforeseen high occurrence of subnormal numbers in single precision sparse LU factorization, two other unexpected observations require further explanation. These are the poor speedup of the matrices from the medium size group, and the fact that many matrices show better speedup in single core experiments than with parallel execution. This section aims to address these questions.

Sparse direct solvers employ more elaborate algorithms than dense solvers. Given a sparse linear system to solve, the rows and the columns of the sparse matrix are first reordered to reduce the number of nonzero elements in the factors, or such that the matrix has dense clusters to take advantage of BLAS 3 kernels. This pre-processing step is called reordering, and it is critical for the overall performance and the memory consumption. After the ordering, the resulting matrix

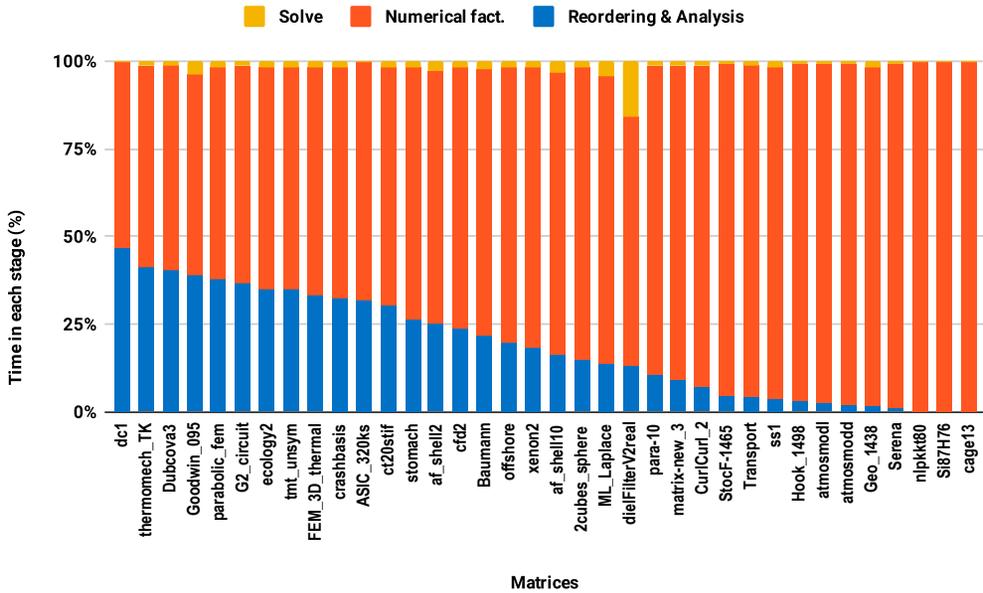


Fig. 11. Time spent by double precision sequential PARDISO LU in each step on a single Intel Skylake core.

is analyzed to determine the nonzero structures of the factors and allocate the required memory accordingly. This step is called symbolic factorization. It is followed by the numerical factorization step that computes the LU factors, and finally the solve step.

The reordering and the analysis steps do not involve floating-point arithmetic. Therefore, they do not benefit from lowering the arithmetic precision. If the reordering and the analysis represent 50% of the overall factorization time, for example, then using single precision instead of double will only reduce the overall time by a quarter in the best case. This explains the poor speedup on average size matrices compared with the large size group. This is illustrated in Figure 11 where one can observe that the majority of average size matrices spend more than 25% of the overall time in the reordering and analysis steps. Here, the bars are sorted by decreasing time associated with the reordering and analysis step. The matrices for which the reordering and analysis time is negligible are the ones that reach up to 2x speedup with single precision.

The second issue, the decrease of speedup in parallel experiments compared with single core executions, is due to the lack of parallelism in the reordering and analysis steps. For example in this work, all the sparse solvers except PARDISO use sequential reordering and analysis algorithms on shared memory multicore architectures. PARDISO provides the parallel version of the nested dissection algorithm for reordering, but compared with the sequential version, it reduces the reordering time only by a factor of two while the numerical factorization time decreases significantly, up to a factor of 8 using 10 cores. Consequently, by increasing the number of cores, the proportion of time spent in reordering and analysis steps increases as illustrated in Figure 12, with the time spent in reordering and analysis step in decreasing order. One can observe that in the parallel experiment, half of the matrices spent more than 50% of the overall factorization time in reordering and analysis, which explains the limited acceleration from lowering the precision.

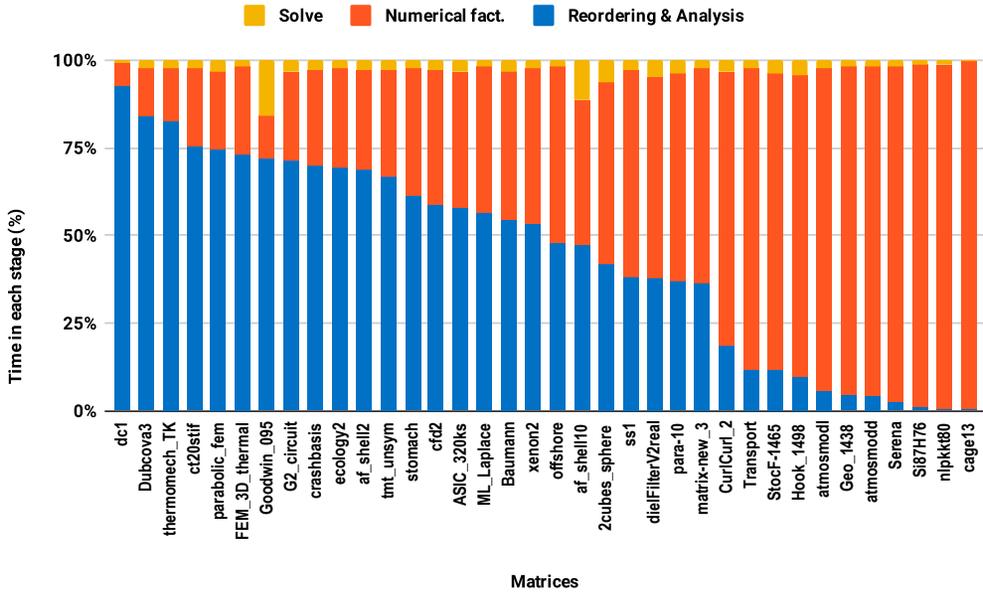


Fig. 12. Time spent by double precision parallel PARDISO LU in each step on 10 Intel Skylake cores.

8 CONCLUSION

The benefits of using mixed precision algorithms for solving dense linear systems are well documented in the HPC community. Much less is known about the efficiency of mixed precision parallel algorithms for sparse linear systems, and existing work focuses on single core experiments. In this work, we have assessed the benefit of using single precision arithmetic in solving double precision sparse linear systems on multicore architectures. We have evaluated two classes of algorithms: iterative refinement based on single precision LU factorization and iterative methods using single precision for the matrix–vector product kernels or preconditioning.

Our first finding is that a limiting factor in the performance of single precision sparse LU factorization is the generation of subnormal numbers, which occurs for the majority of our test matrices. We have identified a mechanism whereby fill-in can cascade down a column, creating and then propagating subnormal numbers with it. We have demonstrated the severe performance drop that can result and have shown how flushing subnormals to zero can mitigate it.

Our second finding is that the anticipated speedup of 2 from mixed precision iterative refinement is obtained only for the very largest of our test problems, where the analysis and reordering time is negligible compared with numerical factorization time.

Our last finding concerns iterative solvers. Our results show that the performance gain in computing or applying incomplete factorization preconditioners in single precision is not appealing enough to justify the accuracy sacrifice, but we have observed a speedup of 1.5 from matrix–vector product kernels by using single precision. In future work, we will explore new approaches to integrate efficiently single precision matrix–vector product kernels and single precision preconditioners in double precision iterative solvers without accuracy loss.

ACKNOWLEDGMENTS

This work was supported by the Innovate UK under grant number KTP011064, by the Engineering and Physical Sciences Research Council under grant number EP/P020720/1, and by the Royal Society.

REFERENCES

- Ahmad Abdelfattah, Hartwig Anzt, Erik G. Boman, Erin Carson, Terry Cojean, Jack Dongarra, Mark Gates, Thomas Grützmacher, Nicholas J. Higham, Sherry Li, Neil Lindquist, Yang Liu, Jennifer Loe, Piotr Luszczek, Pratik Nayak, Sri Pranesh, Siva Rajamanickam, Tobias Ribizel, Barry Smith, Kasia Swirydowicz, Stephen Thomas, Stanimire Tomov, Yaohung M. Tsai, Ichitaro Yamazaki, and Urike Meier Yang. 2020. *A Survey of Numerical Methods Utilizing Mixed Precision Arithmetics*. ArXiv:2007.06674. 44 pages. <https://arxiv.org/abs/2007.06674>
- Emmanuel Agullo, Jim Demmel, Jack Dongarra, Bilel Hadri, Jakub Kurzak, Julien Langou, Hatem Ltaief, Piotr Luszczek, and Stanimire Tomov. 2009. Numerical Linear Algebra on Emerging Architectures: the PLASMA and MAGMA Projects. *Journal of Physics: Conference Series* 180, 1 (2009), 012037. <https://doi.org/10.1088/1742-6596/180/1/012037>
- Khalid Ahmad, Hari Sundar, and Mary Hall. 2019. Data-driven Mixed Precision Sparse Matrix Vector Multiplication for GPUs. *ACM Trans. Archit. Code Optim.* 16, 4 (Dec. 2019), 51:1–51:24. <https://doi.org/10.1145/3371275>
- Patrick R. Amestoy, Iain S. Duff, and Jean-Yves L'Excellent. 2000. Multifrontal Parallel Distributed Symmetric and Unsymmetric Solvers. *Comput. Methods Appl. Mech. Engrg.* 184, 2-4 (2000), 501–520. [https://doi.org/10.1016/S0045-7825\(99\)00242-X](https://doi.org/10.1016/S0045-7825(99)00242-X)
- Patrick R. Amestoy, Iain S. Duff, Jean-Yves L'Excellent, and Xiaoye S. Li. 2001. Analysis and Comparison of Two General Sparse Solvers for Distributed Memory Computers. *ACM Trans. Math. Software* 27, 4 (Dec. 2001), 388–421. <https://doi.org/10.1145/504210.504212>
- E. Anderson, Z. Bai, C. H. Bischof, S. Blackford, J. W. Demmel, J. J. Dongarra, J. J. Du Croz, A. Greenbaum, S. J. Hammarling, A. McKenney, and D. C. Sorensen. 1999. *LAPACK Users' Guide* (third ed.). Society for Industrial and Applied Mathematics, Philadelphia, PA, USA. xxvi+407 pages. <http://www.netlib.org/lapack/lug/>
- Hartwig Anzt, Jack Dongarra, Goran Flegar, Nicholas J. Higham, and Enrique S. Quintana-Ortí. 2019. Adaptive Precision in Block-Jacobi Preconditioning for Iterative Sparse Linear System Solvers. *Concurrency Computat.: Pract. Exper.* 31, 6 (2019), e4460. <https://doi.org/10.1002/cpe.4460>
- Alfredo Buttari, Jack Dongarra, Jakub Kurzak, Piotr Luszczek, and Stanimir Tomov. 2008. Using Mixed Precision for Sparse Matrix Computations to Enhance the Performance While Achieving 64-Bit Accuracy. *ACM Trans. Math. Software* 34, 4 (July 2008), 17:1–17:22. <https://doi.org/10.1145/1377596.1377597>
- Alfredo Buttari, Jack Dongarra, Julie Langou, Julien Langou, Piotr Luszczek, and Jakub Kurzak. 2007. Mixed Precision Iterative Refinement Techniques for the Solution of Dense Linear Systems. *Int. J. High Performance Computing Applications* 21, 4 (2007), 457–466. <https://doi.org/10.1177/1094342007084026>
- Erin Carson and Nicholas J. Higham. 2017. A New Analysis of Iterative Refinement and its Application to Accurate Solution of Ill-Conditioned Sparse Linear Systems. *SIAM J. Sci. Comput.* 39, 6 (2017), A2834–A2856. <https://doi.org/10.1137/17M1122918>
- Erin Carson and Nicholas J. Higham. 2018. Accelerating the Solution of Linear Systems by Iterative Refinement in Three Precisions. *SIAM J. Sci. Comput.* 40, 2 (2018), A817–A847. <https://doi.org/10.1137/17M1140819>
- Timothy A. Davis. [n.d.]. SuiteSparse: A Suite of Sparse Matrix Software. <http://faculty.cse.tamu.edu/davis/suitesparse.html>.
- Timothy A. Davis. 2004. Algorithm 832: UMFPACK V4.3—An Unsymmetric-Pattern Multifrontal Method. *ACM Trans. Math. Software* 30, 2 (June 2004), 196–199. <https://doi.org/10.1145/992200.992206>
- Timothy A. Davis and Yifan Hu. 2011. The University of Florida Sparse Matrix Collection. *ACM Trans. Math. Software* 38, 1 (2011), 1:1–1:25. <https://doi.org/10.1145/2049662.2049663>
- Andrew Dawson, Peter D. Düben, David A. MacLeod, and Tim N. Palmer. 2018. Reliable Low Precision Simulations in Land Surface Models. *Climate Dynamics* 51, 7 (2018), 2657–2666. <https://doi.org/10.1007/s00382-017-4034-x>
- Gabriel Fabien-Ouellet. 2020. Seismic Modeling and Inversion Using Half-Precision Floating-Point Numbers. *GEOPHYSICS* 85, 3 (2020), F65–F76. <https://doi.org/10.1190/geo2018-0760.1>
- Paul Grigoraş, Pavel Burovskiy, Wayne Luk, and Spencer Sherwin. 2016. Optimising Sparse Matrix Vector multiplication for large scale FEM problems on FPGA. In *2016 26th International Conference on Field Programmable Logic and Applications (FPL)*, 1–9. <https://doi.org/10.1109/FPL.2016.7577352>
- Azzam Haidar, Ahmad Abdelfattah, Mawussi Zounon, Panruo Wu, Srikara Pranesh, Stanimire Tomov, and Jack Dongarra. 2018a. The Design of Fast and Energy-Efficient Linear Solvers: On the Potential of Half-Precision Arithmetic and Iterative Refinement Techniques. In *Computational Science—ICCS 2018*, Yong Shi, Haohuan Fu, Yingjie Tian, Valeria V. Krzhizhanovskaya, Michael Harold Lees, Jack Dongarra, and Peter M. A. Sloot (Eds.). Springer, Cham, Switzerland, 586–600. https://doi.org/10.1007/978-3-319-93698-7_45

- Azzam Haidar, Harun Bayraktar, Stanimire Tomov, Jack Dongarra, and Nicholas J. Higham. 2020. *Mixed-Precision Solution of Linear Systems Using Accelerator-Based Computing*. Technical Report ICL-UT-20-05. Innovative Computing Laboratory, University of Tennessee, Knoxville, TN, USA. 30 pages. <https://www.icl.utk.edu/publications/mixed-precision-solution-linear-systems-using-accelerator-based-computing>
- Azzam Haidar, Stanimire Tomov, Jack Dongarra, and Nicholas J. Higham. 2018b. Harnessing GPU Tensor Cores for Fast FP16 Arithmetic to Speed up Mixed-Precision Iterative Refinement Solvers. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis (SC '18 (Dallas, TX))*. IEEE Press, Piscataway, NJ, USA, 47:1–47:11. <https://doi.org/10.1109/SC.2018.00050>
- Nicholas J. Higham. 2002. *Accuracy and Stability of Numerical Algorithms* (second ed.). Society for Industrial and Applied Mathematics, Philadelphia, PA, USA. xxx+680 pages. <https://doi.org/10.1137/1.9780898718027>
- Nicholas J. Higham, Srikara Pranesh, and Mawussi Zounon. 2019. Squeezing a Matrix Into Half Precision, with an Application to Solving Linear Systems. *SIAM J. Sci. Comput.* 41, 4 (2019), A2536–A2551. <https://doi.org/10.1137/18M1229511>
- J. D. Hogg and J. A. Scott. 2010. A Fast and Robust Mixed-Precision Solver for the Solution of Sparse Symmetric Linear Systems. *ACM Trans. Math. Software* 37, 2, Article 17 (April 2010), 17:1–17:24 pages. <https://doi.org/10.1145/1731022.1731027>
- Julie Langou, Julien Langou, Piotr Luszczek, Jakub Kurzak, Alfredo Buttari, and Jack Dongarra. 2006. Exploiting the Performance of 32 bit Floating Point Arithmetic in Obtaining 64 bit Accuracy (Revisiting Iterative Refinement for Linear Systems). In *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing*. <https://doi.org/10.1109/SC.2006.30>
- Xiaoye S. Li. 2005. An Overview of SuperLU: Algorithms, Implementation, and User Interface. *ACM Trans. Math. Software* 31, 3 (2005), 302–325. <https://doi.org/10.1145/1089014.1089017>
- Xiaoye S. Li and James W. Demmel. 2003. SuperLU_DIST: A Scalable Distributed-Memory Sparse Direct Solver for Unsymmetric Linear Systems. *ACM Trans. Math. Software* 29, 2 (2003), 110–140. <https://doi.org/10.1145/779359.779361>
- Magma [n.d.]. Matrix Algebra on GPU and Multicore Architectures (MAGMA). <http://icl.cs.utk.edu/magma/>.
- Jean-Michel Muller, Nicolas Brunie, Florent de Dinechin, Claude-Pierre Jeannerod, Mioara Joldes, Vincent Lefèvre, Guillaume Melquiond, Nathalie Revol, and Serge Torres. 2018. *Handbook of Floating-Point Arithmetic* (second ed.). Birkhäuser, Boston, MA, USA. xxv+627 pages. <https://doi.org/10.1007/978-3-319-76526-6>
- Yousef Saad. 1993. A Flexible Inner-Outer Preconditioned GMRES Algorithm. *SIAM J. Sci. Comput.* 14, 2 (1993), 461–469. <https://doi.org/10.1137/0914028>
- Olaf Schenk, Klaus Gärtner, Wolfgang Fichtner, and Andreas Stricker. 2001. PARDISO: A High-Performance Serial and Parallel Sparse Linear Solver in Semiconductor Device Simulation. *Future Generation Computer Systems* 18, 1 (2001), 69–78. [https://doi.org/10.1016/S0167-739X\(00\)00076-5](https://doi.org/10.1016/S0167-739X(00)00076-5)
- Filip Váňa, Peter Düben, Simon Lang, Tim Palmer, Martin Leutbecher, Deborah Salmond, and Glenn Carver. 2017. Single Precision in Weather Forecasting Models: An Evaluation with the IFS. *Monthly Weather Review* 145, 2 (2017), 495–502. <https://doi.org/10.1175/MWR-D-16-0228.1>