



**HAL**  
open science

## Solving dense symmetric indefinite systems using GPUs

Marc Baboulin, Jack Dongarra, Adrien Rémy, Stanimire Tomov, Ichitaro Yamazaki

► **To cite this version:**

Marc Baboulin, Jack Dongarra, Adrien Rémy, Stanimire Tomov, Ichitaro Yamazaki. Solving dense symmetric indefinite systems using GPUs. *Concurrency and Computation: Practice and Experience*, 2017, 29 (9), pp.1 - 17. 10.1002/cpe.4055 . hal-01662358

**HAL Id: hal-01662358**

**<https://inria.hal.science/hal-01662358>**

Submitted on 20 Jan 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Solving Dense Symmetric Indefinite Systems using GPUs

Marc Baboulin\*    Jack Dongarra†    Adrien Rémy\*  
Stanimire Tomov†    Ichitaro Yamazaki†

## Abstract

This paper studies the performance of different algorithms for solving a dense symmetric indefinite linear system of equations on multicore CPUs with a Graphics Processing Unit (GPU). To ensure the numerical stability of the factorization, *pivoting* is required. Obtaining high performance of such algorithms on the GPU is difficult because all the existing pivoting strategies lead to frequent synchronizations and irregular data accesses. Until recently, there has not been any implementation of these algorithms on a hybrid CPU/GPU architecture. To improve their performance on the hybrid architecture, we explore different techniques to reduce the expensive data transfer and synchronization between the CPU and GPU, or on the GPU (e.g., factorizing the matrix entirely on the GPU or in a communication-avoiding fashion). We also study the performance of the solver using iterative refinements along with the factorization without pivoting combined with the preprocessing technique based on Random Butterfly Transformations, or with the mixed-precision algorithm where the matrix is factorized in single precision. This randomization algorithm only has a probabilistic proof on the numerical stability, and for this paper, we only focused on the mixed-precision algorithm without pivoting. However, they demonstrate that we can obtain good performance on the GPU by avoiding the pivoting and using the lower precision arithmetics, respectively. As illustrated with the application in acoustics studied in this paper, in many practical cases, the matrices can be factorized without pivoting. Since the componentwise backward error computed in the iterative refinement signals when the algorithm failed to obtain the desired accuracy, the user can use these potentially-unstable but efficient algorithms in most of the cases, and fall back to a more stable

---

\*University of Paris-Sud, France, [baboulin,aremy]@lri.fr

†University of Tennessee, Knoxville, USA, [dongarra,tomov,iyamazak]@eecs.utk.edu

algorithm with pivoting only in the case of the failure.

**Keywords:** dense symmetric indefinite systems; symmetric pivoting; communication-avoiding; randomization; GPU computation; iterative refinement.

## 1 Introduction

A symmetric matrix  $A$  is called indefinite when the quadratic form  $x^T Ax$  can take both positive and negative values. Dense linear systems of equations with symmetric indefinite matrices appear in many studies of physics, including physics of structures, acoustics, and electromagnetism. For instance, such systems arise in the linear least-squares problem for solving an augmented system [18, p. 77], or in the electromagnetism where the discretization by the Boundary Element Method results in linear systems with dense complex symmetric (non-Hermitian) matrices [28]. The efficient solution of these linear systems demands a high performance implementation of a dense symmetric indefinite solver that can efficiently use the current hardware architecture. In particular, the use of accelerators has become pervasive in scientific computing due to their high-performance capabilities and low energy consumptions. For example, in term of the floating-point operation per second, or flop/s in short, a single K40 NVIDIA GPU has a double precision peak performance of 1,689 Gflop/s for a thermal design power (TDP) of 235 W. According to benchmarks in the MAGMA library [33], optimized large dense matrix computations, e.g., matrix-matrix multiplications, reach 1,200 Gflop/s for a power draw of about 200 W, i.e.,  $\approx 6$  Gflop/W. In contrast, two Sandy Bridge E5-2670 CPUs have about the same TDP ( $2 \times 115 = 230$  W) as the K40 but for a peak of 333 Gflop/s, which translates to only 1.4 Gflop/W for the Sandy Bridge CPU. To achieve the high performance, however, the algorithms must be designed for high parallelism and high “flops to data” ratio, while maintaining a low number of flops and exploiting the hardware features of the hybrid CPU/GPU architecture. A dense symmetric indefinite solver which can efficiently exploit the GPU’s high computing power would be useful for many physical applications.

To solve a symmetric indefinite linear system of equations,  $Ax = b$ , a classical method decomposes the matrix  $A$  into an  $LDL^T$  factorization,

$$PAP^T = LDL^T, \quad (1)$$

where  $L$  is unit lower triangular,  $D$  is block diagonal with either 1-by-1 or 2-by-2 diagonal blocks, and  $P$  is a permutation matrix to ensure the numerical

stability of the factorization. Then the solution  $x$  is computed by successively solving the linear systems with the coefficient matrices  $L$ ,  $D$ , and  $L^T$  along with the permutation. The strategies to compute the permutation matrix  $P$  for the  $LDL^T$  factorization include complete pivoting (Bunch-Parlett algorithm) [13], partial pivoting (Bunch-Kaufman algorithm) [14], rook pivoting (bounded Bunch-Kaufman) [6, p. 523], and fast Bunch-Parlett [6, p. 525]. In particular, the Bunch-Kaufman and rook pivoting strategies are implemented in LAPACK [4], a set of dense linear algebra routines on multicore CPUs, that are extensively used in many scientific and engineering simulations. The routines implemented in LAPACK are based on block algorithms that can exploit the memory hierarchy on modern architectures, using BLAS-3 matrix operations for most of its floating-point operations.

Another promising method for solving a symmetric indefinite linear system is the Aasen’s method [1], which computes the  $LTL^T$  factorization of the matrix  $A$ ,

$$PAP^T = LTL^T, \quad (2)$$

where  $T$  is now a symmetric tridiagonal matrix. The algorithm requires  $\frac{1}{3}n^3 + O(n^2)$  flops [21, p. 166], similarly to the  $LDL^T$  factorization. A block algorithm for computing the  $LTL^T$  factorization was also proposed [30]. Though the block implementation performs slightly more flops (i.e., an additional rank-1 update of the trailing submatrix, see Section 2.3), it can exploit a modern computer’s memory hierarchy and obtain performance similar to the Bunch-Kaufman algorithm implemented in LAPACK [30].

To maintain numerical stability, the pivoting techniques mentioned above involve between  $O(n^2)$  and  $O(n^3)$  comparisons to search for pivots and possible interchanges of selected columns and rows. Hence, factorizing each column of the matrix requires the synchronization for selecting the pivot and the data movement for exchanging the columns and rows, which have become significantly more expensive compared to the arithmetic operations on modern computers. Furthermore, since only either the upper or lower triangular part of the matrix  $A$  is stored, the symmetric pivoting<sup>1</sup> requires irregular data access (i.e., some parts of the pivot column may be stored as the transpose of the corresponding part of the row), which dramatically increases the cost of the data movement. Partially due to these performance challenges, ScaLAPACK [12], which is the extension of LAPACK for distributed-memory machines, does not support the symmetric indefinite factorization, and until recently, there were no implementations of the algo-

---

<sup>1</sup>To maintain the symmetry, both columns and rows must be swapped.

rithm, that could exploit a GPU<sup>2</sup>. This motivated our efforts to review the different factorization algorithms, develop their efficient implementations on multicores with a GPU to address their current limitations, and show the new state-of-the-art outlook for this important problem. For example, recently, a communication-avoiding variant of the Aasen’s algorithm was proposed [11]. However, the pivoting must still be applied symmetrically, leading to expensive irregular data accesses. Another technique studied in this paper is a symmetric version of Random Butterfly Transformations (RBT) [29] on the GPU. RBT can be combined with an  $LDL^T$  factorization to probabilistically improve the stability of the factorization without pivoting. The performance of RBT for symmetric indefinite systems has been studied on multicore systems [10] and distributed-memory systems [8], but its performance has not been investigated on a GPU. Finally, we study the potential of a mixed-precision algorithm to improve the performance of the solver, where the matrix is first factorized in single precision, and the solution is computed through iterative refinement.

This paper is organized as follows. Section 2 describes the three algorithms for solving dense symmetric indefinite systems (i.e., the Bunch-Kaufman and Aasen’s algorithms, and the random butterfly transformations) and their implementations on the hybrid CPU/GPU architecture. It also explains how we can use mixed precision to accelerate the solver. Section 3 shows our experimental results, where Sections 3.1 and 3.2 present the performance and numerical results for random matrices and two acoustic scattering problems, respectively, while Section 3.3 gives performance results of the mixed-precision algorithm applied to random matrices without pivoting. Section 4 contains concluding remarks. In this paper, we use  $a_{i,j}$  and  $a_j$  to denote the  $(i, j)$ -th entry and the  $j$ -th column of the matrix  $A$ , respectively, while  $A_{i_1:i_2, j_1:j_2}$  is the submatrix consisting of the  $i_1$ -th through the  $i_2$ -th rows and the  $j_1$ -th through the  $j_2$ -th columns of  $A$ . We also use  $A_{I,J}$  and  $A_{I_1:I_2, J_1:J_2}$  to denote the  $(I, J)$ -th block and the submatrix consisting of the  $I_1$ -th through the  $I_2$ -th block rows and  $J_1$ -th through the  $J_2$ -th block columns of  $A$ , where the block size is  $n_b$  and the number of block columns/rows in  $A$  is  $n_t$  (i.e.,  $n_t = \lceil \frac{n}{n_b} \rceil$ ).

This paper extends our previous proceedings paper [38] presented at the PPAM 2015 conference. In this extended paper, we describe the current general trends in designing efficient numerical linear algebra libraries on manycore accelerated architectures (Section 2.1) before presenting our spe-

---

<sup>2</sup>A Bunch-Kaufman implementation became recently available in the cuSolver library as part of the CUDA Toolkit v7.5 from NVIDIA.

cific design and optimization of the symmetric indefinite solvers for the GPU architectures (Section 2.2). We also include the time to obtain the solution while the previous paper only showed the factorization time, and give more details about the acoustic scattering problems studied in the paper (Sections 3.1 and 3.2, respectively). Finally, we describe our implementation and its performance of the mixed-precision algorithm which may improve the performance of the solver in practice (Sections 2.5 and 3.3).

## 2 Symmetric Indefinite Factorizations with a GPU

In this section, we describe the existing algorithms for solving a dense symmetric indefinite linear system of equations. First, we describe the general principles for designing an efficient dense linear algebra algorithm on heterogeneous systems, and then, we concentrate on the specifics for the design and optimization of symmetric indefinite solvers for GPU architectures, along with discussion on what design principles can (or can not) be applied for these solvers.

### 2.1 Programming linear algebra solvers on GPUs

The LAPACK’s programming model [4] is based on expressing algorithms in terms of BLAS calls. Subsequently, LAPACK can achieve high efficiency, provided that highly efficient BLAS implementations are provided on the target machine, e.g., by the manufacturer. Since the 1980s this model has turned out to be very successful for cache-based shared-memory vector and parallel processors with multi-layered memory hierarchies.

To account for the deep memory hierarchies today, efficient BLAS implementations feature multilevel blocking where, for example, the Level 3 matrix-matrix computations are split hierarchically into blocks that fit into corresponding levels of the memory hierarchy [40]. In effect, a programming model based on BLAS is still an effective model for exploiting the deep memory hierarchies at the present time [2]. However, the resulting parallelism is *fork-join* – a sequence of BLAS calls is implicitly synchronized after each individual BLAS call (join), though the routines by themselves run in parallel (fork). This brings synchronization overheads and idle time for some processors/cores, especially on the highly parallel current and future heterogeneous system designs [39], motivating the search for improved models where the BLAS routines are broken into small tasks and properly scheduled for execution over the heterogeneous hardware components.

The typical hybrid algorithm splits the overall computation into small tasks to execute on the CPU, and large update tasks to execute on the accelerator [33, 42, 43, 44]. For instance, in LU and QR factorizations, each step is split into a panel factorization of  $n_b$  columns, followed by a trailing matrix update. The panel factorization is assigned to the CPU, and includes such decisions as selecting the maximum pivot in each column or computing a Householder reflector for each column. The trailing matrix update is assigned to the accelerator, and involves some form of matrix-matrix multiply. The block size,  $n_b$ , can be tuned to adjust the amount of work on the CPU vs. on the accelerator. Optimally, during the trailing matrix update, a look-ahead panel is updated first and sent back to the CPU. Asynchronous data transfers are used to copy data between the CPU and accelerator while the accelerator continues computing. The CPU performs the next panel factorization while the accelerator continues with the remainder of the trailing matrix update. In this way, the inputs for the next trailing matrix update are ready when the current update finishes. The goal is to keep the accelerator always busy, which has the highest performance.

Unfortunately, the pivoting required to maintain the numerical stability of the symmetric indefinite factorization leads to the fork-join and prohibits the look-ahead as we describe in the rest of this section.

## 2.2 Bunch-Kaufman Algorithm

One of the most widely used algorithms for solving a symmetric indefinite linear system is based on the block  $LDL^T$  factorization with the Bunch-Kaufman algorithm [14], which is also implemented in LAPACK (i.e., xSYTRF). The pseudo-code of the algorithm is shown in Figure 1a, which is referred to as a *right-looking* algorithm because after the set of  $n_b$  columns, commonly referred to as *panel*, are factorized, the panel is used to update the trailing submatrix, which is on the right of the panel. To select the pivot at each step of the factorization, it scans at most two columns of the trailing submatrix, and depending on the numerical values of the scanned matrix entries, it uses either a 1-by-1 or a 2-by-2 pivot. This algorithm has satisfactory backward stability [27, p. 219]. Then a variant of the Bunch-Kaufman algorithm, also called “rook pivoting”, was proposed in [6] that provides a better accuracy by bounding the triangular factors. However, depending on the matrix, the rook pivoting method could perform  $O(n^3)$  comparisons, as opposed to the  $O(n^2)$  comparisons of the Bunch-Kaufman algorithm. Hence, in this paper, we focus on the Bunch-Kaufman algorithm as a baseline for our performance comparison.

```

 $\alpha = (1 + \sqrt{17})/8$  and  $j = 1$ 
while  $j < n$  do
   $k = j$ 
  {Panel factorization}
  while  $j < k + n_b - 1$  do
    Update column  $a_j$  with previous
    columns
     $r = \arg \max_{i>j} |a_{i,j}|$  and  $\gamma = |a_{r,j}|$ 
    if  $\gamma > 0$  then
      if  $|a_{j,j}| \geq \alpha\gamma$  then
         $s = 1$ 
        Use  $a_{j,j}$  as a  $1 \times 1$  pivot.
      else
        Update  $a_r$  with previous
        columns
         $\omega = \max_{i \geq j, i \neq r} |a_{i,r}|$ 
        if  $|a_{j,j}| \omega \geq \alpha\gamma^2$  then
           $s = 1$ 
          Use  $a_{j,j}$  as a  $1 \times 1$  pivot.
        else
          if  $|a_{r,r}| \geq \alpha\omega$  then
             $s = 1$ 
            Swap rows/columns
             $(j, r)$ 
            Use  $a_{r,r}$  as a  $1 \times 1$  pivot.
          else
             $s = 2$ 
            Swap rows/columns  $(j +$ 
             $1, r)$ 
            Use  $\begin{pmatrix} a_{j,j} & a_{r,j} \\ a_{r,j} & a_{rr} \end{pmatrix}$  as
             $2 \times 2$  pivot.
          end if
        end if
      end if
    end if
     $s = 1$ 
    end if
    Scale the pivot columns to compute
     $L_{j:j+s-1}$ 
     $j = j + s$ 
  end while
  {Right-looking trailing submatrix up-
  date}
   $A_{j:n,j:n} := A_{j:n,j:n} -$ 
   $L_{j:n,k:j} D_{k:j,k:j} L_{j:n,k:j}^T$ 
end while

```

(a) Bunch-Kaufman [14].

```

for  $J = 1, 2, \dots, n_t$  do
  for  $I = 2, 3, \dots, J - 1$  do
     $X = T_{I,I-1} L_{J,I-1}^T$ 
     $Y = T_{I,I} L_{J,I}^T$ 
     $Z = T_{I,I+1} L_{J,I+1}^T$ 
     $W_{I,J} = 0.5Y + Z$ 
     $H_{I,J} = X + Y + Z$ 
  end for

   $A_{J,J} = A_{J,J} - L_{J,2:J-1} W_{2:J-1,J}$ 
   $- W_{2:J-1,J}^T L_{J,2:J-1}^T$ 
   $T_{J,J} = L_{J,J}^{-1} A_{J,J} L_{J,J}^{-T}$ 
  if  $J < n_t$  then
    if  $J > 1$  then
       $H_{J,J} = T_{J,J-1} L_{J,J-1}^T + T_{J,J} L_{J,J}^T$ 
    end if

     $E = A_{J+1:n_t,J} - L_{J+1:n_t,2:J} H_{2:J,J}$ 
     $[L_{J+1:n_t,J+1}, H_{J+1,J}, P^{(J)}] = \text{LU}(E)$ 

     $T_{J+1,J} = H_{J+1,J} L_{J,J}^{-T}$ 

     $L_{J+1:n_t,2:J} = P^{(J)} L_{J+1:n_t,2:J}$ 
     $A_{J+1:n_t,J+1:n_t}$ 
     $P^{(J)} A_{J+1:n_t,J+1:n_t} P^{(J)T}$ 
     $P_{J+1:n_t,1:n_t} = P^{(J)} P_{J+1:n_t,1:n_t}$ 
  end if
end for

```

(b) CA Aasen's [11], where the first block column  $L_{1:n_t,1}$  is the first  $n_b$  columns of the identity matrix and  $[L, U, P] = \text{LU}(A)$  returns the LU factors of  $A$  with partial pivoting such that  $LU = PA$ .

Figure 1: Symmetric indefinite factorization algorithm.



Our implementation of the Bunch-Kaufman algorithm on the hybrid architecture is based on BLAS and LAPACK task representations (as described in subsection 2.1), where the BLAS and LAPACK calls on the CPU are replaced with the corresponding GPU kernels. In addition, our first implementation is based on a hybrid CPU/GPU programming paradigm where the panel is factorized on the CPU (e.g., using the multithreaded MKL library [5]), while the trailing submatrix is updated on the GPU. This is often an effective programming paradigm for many of the LAPACK subroutines because the panel factorization is based on BLAS-1 or BLAS-2, which can be efficiently implemented on the CPU, while BLAS-3 is used for the submatrix updates, which exhibits high data parallelism and can be efficiently implemented on the GPU [32, 33]. Unfortunately, at each step of the panel factorization, the Bunch-Kaufman algorithm may select the pivot from the trailing submatrix. Hence, though copying the panel from the GPU to the CPU can be overlapped with the update of the rest of the trailing submatrix on the GPU, the *look-ahead* – a standard optimization technique to overlap the panel factorization on the CPU with the trailing submatrix update on the GPU – is prohibited. In addition, when the pivot column is on the GPU, this leads to an expensive data transfer between the GPU and the CPU at each step of the factorization. To avoid this expensive data transfer, our second implementation performs the entire factorization on the GPU. Though the CPU may be more efficient at performing the BLAS-1 and BLAS-2 based panel factorization, this implementation often obtains higher performance by avoiding the expensive data transfer (See Figure 3).

When the entire factorization is implemented on the GPU, up to two columns of the trailing submatrix must be scanned to select a pivot – the current column and the column with index corresponding to the row index of the element with the maximum modulus in the first column. This not only leads to the expensive global reduce on the GPU, but also to irregular data accesses since only the lower-triangular part of the submatrix is stored. This makes it difficult to obtain high performance on the GPU. In the next two sections, we describe two other algorithms (i.e., communication-avoiding and randomization algorithms) that aim at reducing this bottleneck.

### 2.3 Aasen’s Algorithm

To solve a symmetric indefinite linear system, Aasen’s algorithm [1] factorizes  $A$  into an  $LTL^T$  decomposition. The algorithm takes advantage of the symmetry of  $A$  and performs  $\frac{1}{3}n^3 + O(n^2)$  flops, which is the same flop count as that of the Bunch-Kaufman algorithm. In addition, like the

Bunch-Kaufman algorithm, it is backward stable subject to a growth factor. To maintain the stability, at each step of the factorization, it uses the largest element of the current column being factorized as the pivot, leading to more regular data access compared to the Bunch-Kaufman algorithm (that may scan an additional column, some part of which may be stored as the transpose of the corresponding part of the row). To exploit the memory hierarchy of modern computers, a blocked version of the algorithm was developed [30], which is based on a *left-looking* panel factorization, followed by a right-looking trailing submatrix update using BLAS-3 routines. Compared to the column-wise algorithm, this blocked algorithm performs slightly more flops, requiring  $\frac{1}{3}(1 + \frac{1}{n_b})n^3 + O(n^2n_b)$  flops with a block size  $n_b$ , but BLAS-3 can be used to perform most of these flops (i.e.,  $\frac{1}{3}(1 + \frac{1}{n_b})n^3$ ). However, the panel factorization is still based on BLAS-1 and BLAS-2, which often obtains only a small fraction of the peak performance. To improve the performance of the panel factorization, another variant of the algorithm was proposed [11]. This other variant computes an  $LTL^T$  factorization of  $A$ , where  $T$  is a banded matrix with its half-bandwidth equal to the block size  $n_b$ , and then uses a banded matrix solver to compute the solution. This algorithm factorizes each panel using an existing LU factorization algorithm, such as recursive LU [22, 26, 34] or communication-avoiding LU (TSLU, for the panel) [23, 24]. In comparison with the panel factorization algorithm used in the block Aasen’s algorithm, these LU factorization algorithms reduce communication, and are likely to speed up the whole factorization process. This is referred to as a communication-avoiding (CA) variant of the Aasen’s algorithm, and its pseudocode is shown in Figure 1b.

In general, a GPU has a greater memory bandwidth than a CPU, but the memory accesses are still expensive compared to the arithmetic operations. Hence, our implementation is based on the CA Aasen’s algorithm. Though this algorithm performs most of the flops using BLAS-3 (e.g., xGEMM), most of the operations are on the submatrices of the dimension  $n_b$ -by- $n_b$ . In order to run these small independent BLAS calls in parallel on the GPU, we use GPU streams. An alternative is to use Batched BLAS, where all independent xGEMMs are grouped together in a single call. Implementations are available in both MAGMA and CUBLAS. However, as shown in [3, Figure 8(f)], the streamed implementation (that we use here) is faster than either the MAGMA Batched or CUBLAS Batched DGEMM for matrices of size above 160 (on K40c GPU for DGEMM on square matrices, i.e.,  $m = n = k$ ), which is the case here. With the GPU streams, the CA Aasen obtained its best performance using  $n_b = 256$  (see Figure 3).

Our CA Aasen’s implementation applies the pivots in two steps: The first step copies all the columns of the trailing submatrix, which need to be swapped, into an  $n$ -by- $2n_b$  workspace. Here, because of the symmetry, the  $k$ -th block column consists of the blocks in the  $k$ -th block row and those in the  $k$ -th block column (each block column consists of the  $n_b$  contiguous columns). Then, in the second step, we copy the columns of the workspace back to a block column of the submatrix after the column pivoting is applied. The two-step implementation is used to exploit the parallelism on multicore CPU [15] and in our non-GPU-resident implementations to factorize the matrices that do not fit in the GPU memory at once [16]. In our experiments, to factorize the panel, we used the LU factorization with partial pivoting, using either the multithreaded MKL library on the CPU or using its native GPU implementation in MAGMA on the GPU. Though the BLAS-1 and BLAS-2 based panel factorization may be more efficient on the CPU, the second approach avoids the expensive data transfer required to copy the panel from the GPU to the CPU (see Section 3 for the performance results).

## 2.4 Random Butterfly Transformations

Random Butterfly Transformation (RBT) is a randomization technique initially described by Parker [29] and recently revisited for dense linear systems, either general [7] or symmetric indefinite [8]. It has also been applied recently to a sparse direct solver in [9]. The procedure to solve  $Ax = b$ , where  $A$  is a symmetric indefinite matrix, using a random transformation and the  $LDL^T$  factorization is summarized in Algorithm 1. The random matrix  $U$  is chosen among a particular class of matrices called *recursive butterfly matrices*. A *butterfly matrix* is an  $n \times n$  matrix of the form

$$B^{<n>} = \frac{1}{\sqrt{2}} \begin{bmatrix} R_0 & R_1 \\ R_0 & -R_1 \end{bmatrix}$$

where  $R_0$  and  $R_1$  are random diagonal  $\frac{n}{2} \times \frac{n}{2}$  matrices. A *recursive butterfly matrix* of size  $n$  and depth  $d$  is defined recursively as

$$W^{<n,d>} = \begin{bmatrix} B_1^{<n/2^{d-1}>} & & & \\ & \ddots & & \\ & & & B_{2^{d-1}}^{<n/2^{d-1}>} \end{bmatrix} \cdot W^{<n,d-1>}, \text{ with } W^{<n,1>} = B^{<n>}$$

where the  $B_i^{<n/2^{d-1}>}$  are butterflies of size  $n/2^{d-1}$ , and  $B^{<n>}$  is a butterfly of size  $n$ . The application of RBT to symmetric indefinite problems was

studied in [17] where it is shown that in practice,  $d = 1$  or  $2$  gives satisfactory results. Note that, as mentioned in [7], the solution can be improved by adding systematically some steps of iterative refinement in the working precision as indicated in [27, p. 232]. It is also shown that random butterfly matrices are cheap to store and apply ( $O(nd)$  and  $O(dn^2)$  respectively). An implementation for the multicore library PLASMA was described in [10].

---

**Algorithm 1** Random Butterfly Transformation Algorithm

---

Generate recursive butterfly matrix  $U$   
 Apply randomization to update the matrix  $A$  and compute the matrix  $A_r = U^T A U$   
 Factorize the randomized matrix using  $LDL^T$  factorization with no pivoting  
 Compute right-hand side  $U^T b$ , solve  $A_r y = U^T b$ , then  $x = U y$

---

For the GPU implementation, we use a recursive butterfly matrix  $U$  of depth  $d = 2$ . Only the diagonal values of the blocks are stored into a vector of size  $2 \times N$  as described in [7]. Applying the depth 2 recursive butterfly matrix  $U$  consists of multiple applications of depth 1 butterfly matrices on different parts of the matrix  $A$ . The application of a depth 1 butterfly matrix is performed using a CUDA kernel where the computed part of the matrix  $A$  is split into blocks. For each of these blocks, the corresponding part of the matrix  $U$  is stored in the shared memory to improve the memory access performance. Matrix  $U$  is small enough to fit into the shared memory due to its packed storage.

To compute the  $LDL^T$  factorization of  $A_r$  without pivoting, we implemented a block factorization algorithm on multicore CPUs with a GPU. In our implementation, the matrix is first copied to the GPU, then the CPU is used to compute the  $LDL^T$  factorization of the diagonal block. Once the resulting  $LDL^T$  factors of the diagonal block are copied back to the GPU, the corresponding off-diagonal blocks of the  $L$ -factor are computed by the triangular solve on the GPU. Finally, we update each block column of the trailing submatrix calling a matrix-matrix multiply on the GPU.

## 2.5 Mixed precision algorithm

On modern computers, single precision 32-bit floating point arithmetic is usually at least twice as fast as double precision 64-bit floating point arithmetic. For example, on a latest NVIDIA GPU (e.g., the GeForce GTX Titan Black), the single precision peak performance is about  $3 \times$  greater

than the double precision peak performance. This gap can be much greater depending on the number of 32-bit and 64-bit CUDA cores (e.g.,  $32\times$  faster on the Titan X). To take advantage of this hardware trend for solving a linear system of equations, the mixed-precision algorithm may compute a solution in single precision and then aims to refine the solution to have double precision accuracy by performing only the critical parts of the algorithm in double precision. Iterative refinement in single/double precision is presented in [31, 35, 36] and has been implemented in so-called mixed precision solvers in [19, 20].

```

compute  $LDL^T := A$       in single
precision
solve  $Ly = b$  for  $y$       in single
precision
solve  $Dz = y$  for  $z$       in single
precision
solve  $L^T x = z$  for  $x$     in single
precision
compute  $r := b - Ax$       in double
precision
while  $\|r\|_2 > \|x\|_2 \|A\|_\infty \epsilon \sqrt{n}$  do
  solve  $Ly = r$  for  $y$       in single
  precision
  solve  $Dz = y$  for  $z$       in single
  precision
  solve  $L^T e = z$  for  $e$     in single
  precision
  compute  $x := x + e$       in double
  precision
  compute  $r := b - Ax$       in double
  precision
end while

```

Figure 2: Fixed precision iterative refinement without pivoting where  $\epsilon$  is the relative machine precision in double precision, given by LAPACK’s DLAMCH. The algorithm can be trivially extended to use pivoting.

Figure 2 shows the pseudocode of such a mixed-precision algorithm, applied to the  $LDL^T$  factorization with no pivoting. Note that this is different

from what is called “mixed precision” in the literature (e.g. [21, p. 127]) since in our case  $x := x + e$  is computed in double precision. The factorization of the coefficient matrix  $A$  is the most computationally-expensive kernel, requiring  $\mathcal{O}(n^3)$  flops, while the other kernels require at most  $\mathcal{O}(n^2)$  flops. To take advantage of the higher performance, the coefficient matrix  $A$  is converted to single precision and factorized in single precision. Then, in order to obtain double-precision accuracy, double-precision arithmetic is used to compute the residual vector and to update the solution vector. To compute the residual vector, the original coefficient matrix  $A$  is needed. Hence, compared to the standard algorithm, which performs all the operations in double precision, the mixed-precision algorithm requires 50% more memory to store  $A$  in single precision. However, the most expensive kernel is handled in single precision, and the mixed-precision algorithm may obtain a higher performance than the standard algorithm does, as long as it requires a small number of iterations. The numerical analysis of the standard or mixed-precision iterative refinements can be found in [21, 25, 31, 35, 36].

### 3 Experimental Results

#### 3.1 Comparison of symmetric indefinite solvers

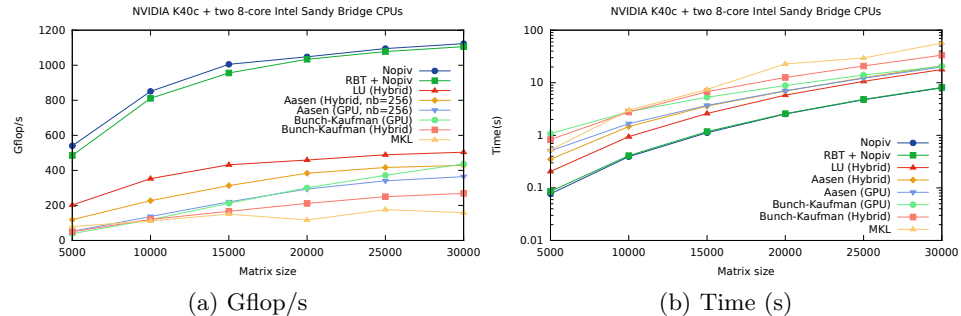


Figure 3: Performance of dense symmetric factorizations (double precision).

Figures 3a and 3b compare respectively the performance in Gflop/s and time for the symmetric indefinite factorizations where the test matrices are random. The “Gflop/s” is computed as the ratio of the number of flops required for the  $LDL^T$  factorization (i.e.,  $n^3/3$ ) over time (in seconds) for the particular dimension of the matrix,  $n$ . Note that, for normalization of

the graph, we also consider the same flop count for  $LU$ , even though it performs twice more flops. The experiments were conducted on two eight-core Intel SandyBridge CPUs with an NVIDIA K40c GPU. The code is compiled using the GNU gcc version 4.4.7 and the nvcc version 7.0 with the optimization flag `-O3` and linked with Intel’s Math Kernel Library (MKL) version `xe_2013_sp1.2.144`. First, when the matrix size is large enough (i.e.,  $n > 10,000$ ), the performance of the Bunch-Kaufman algorithm can be improved using the GPU over the multithreaded MKL implementation (routine `dsytrf`) on the 16 cores of two Sandy Bridge CPUs. In addition, performing the panel factorization on the GPU avoids the expensive data transfer between the CPU and GPU, and may improve the performance of the hybrid CPU/GPU implementation. Next, the communication-avoiding variant of the Aasen’s algorithm further improves the performance of the Bunch-Kaufman by reducing the synchronization and communication costs required for selecting the pivots. The RBT approach outperforms the Bunch-Kaufman and Aasen factorizations but, as mentioned in [10], it may not be numerically stable for some matrices and it requires in general a few steps of iterative refinement in the working precision. However, the performance of all the symmetric factorizations with provable stability was lower than that of the LU factorization. In addition, though our current implementations of the Bunch-Kaufman and Aasen’s algorithms were slower than the LU factorization, they preserve the symmetry which can reduce the runtime or memory requirement for the rest of the software (e.g., sparse symmetric factorization, or any simulation code). The symmetric factorization also preserves the inertia of the coefficient matrix.

After having compared the performance of the factorization, we now compare the performance of solving a linear system using random matrices. Figures 4a and 4b compare respectively the performance in Gflop/s and time for the symmetric indefinite solvers on multicores with a GPU. The “Gflop/s” is computed as the ratio of the number of flops required for the factorization (i.e.,  $n^3/3$ ) plus the number of flops for the solve (i.e.,  $2n^2 + n$ ) over time (in seconds). The time for the transfer of the matrices between CPU and GPU is also taken into account. Here the randomization and the iterative refinement are performed on the GPU, the factorizations are performed with the hybrid CPU/GPU implementations as described previously. The solve is performed on the CPU for Aasen and Bunch-Kaufman and on the GPU for the other implementations. Here the curve for the RBT solver with iterative refinement stops at size 20,000 because the iterative refinement requires a copy of the original matrix and thereby two times more memory on the GPU. Consistently with the previous experiments, the Aasen

solver is slightly faster than the Bunch-Kaufman solver and the no-pivoting solvers outperform those that use pivoting.

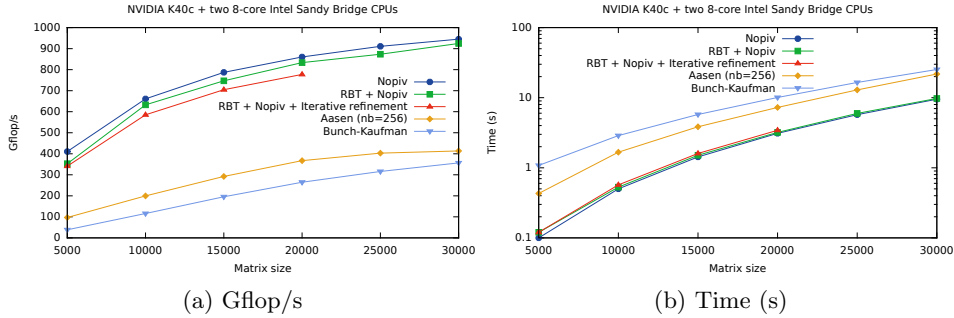


Figure 4: Performance of dense symmetric solvers (double precision).

Let us now study the backward error obtained for the linear system solution computed with the corresponding solvers (on random matrices). We plot in Figure 5 the componentwise backward error given in [4, p. 78] and expressed by

$$\omega = \max_i \frac{|Ax - b|_i}{(|A| \cdot |x| + |b|)_i},$$

where  $x$  is the computed solution. For the RBT solver, we consider the cases without iterative refinement and with one step of iterative refinement in the working precision. We observe that adding one step of iterative refinement is sufficient to obtain a backward error similar to the other solvers (i.e., in the range  $10^{-14} - 10^{-15}$  for the random matrices considered in these experiments).

### 3.2 Experiments and applications for no-pivoting $LDL^T$

In some physical applications involving dense symmetric complex non-Hermitian systems, it is not necessary to pivot in the  $LDL^T$  factorization (see e.g., [27, p. 209] for more information on this class of matrices). These systems are classically solved using an LU factorization since ScaLAPACK does not provide symmetric factorization for this type of matrix. The application considered here is related to the simulation of processes in which acoustic waves are scattered by obstacles. Unless the geometry of the scattering object is very simple, it is generally not possible to find an analytical solution of scattering problems and then numerical schemes are required. A classical



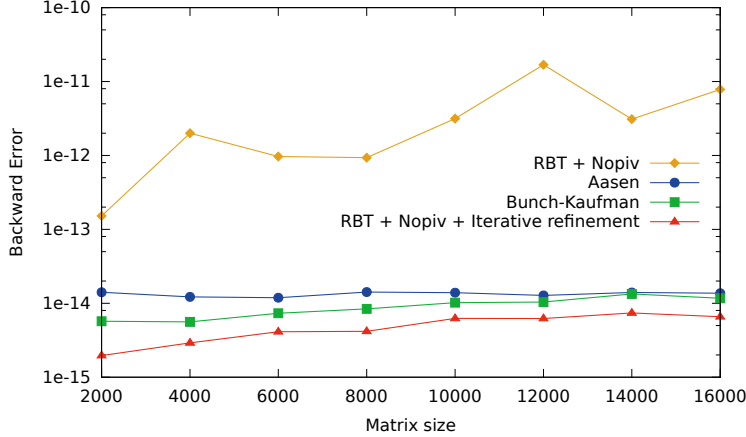


Figure 5: Comparison of componentwise backward error (double precision).

approach is to approximate the solution to time harmonic acoustic problems using the Boundary Element Method (BEM). The BEM discretization leads to linear systems with dense complex symmetric (non-Hermitian) matrices that usually do not require pivoting. Here we consider two test cases where the scattering objects correspond to a human head and a truck engine (see Figure 6).

The matrices (in single complex precision) resulting from the BEM discretization have, respectively, the sizes 10,424 and 15,135. Tables 1 and 2 present numerical results for the solution based on our  $LDL^T$  factorization with no pivoting on the GPU (see end of Section 2.4, here no RBT is used), applied to two sample matrices with comparison to LU factorization. Due to the smaller number of flops, our  $LDL^T$  factorization enables us to accelerate the calculation by about 48%, while keeping a similar accuracy, expressed here by computing the scaled residual  $\|b - Ax\|_\infty / (N\|A\|_\infty \times \|x\|_\infty)$ .

### 3.3 Performance results of mixed-precision iterative refinements

Figure 7 compares the performance of the mixed-precision algorithm (routine ZCHESV) with that of the standard symmetric indefinite solvers using single and double complex precisions (routines CHESV and ZHESV), and on random matrices. We computed the Gflop/s using the flop count needed for the standard algorithm in double complex precision (i.e.,  $\frac{4}{3}n^3 + 8n^2n_{rhs} + o(n^2)$  flops needed to compute the  $LDL^T$  factorization and to perform a pair of

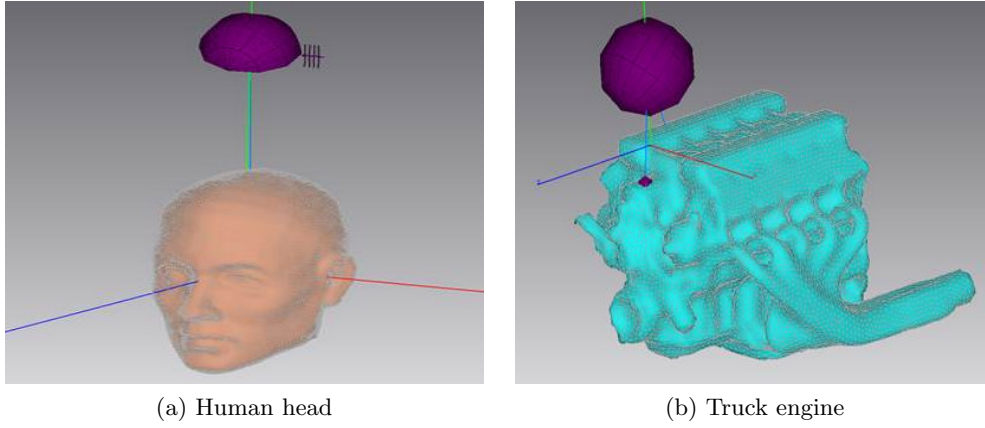


Figure 6: Test cases for acoustic scattering problems

Table 1: Human head (matrix size is 10,424 in single complex precision).

	Time (sec)	Scaled residual
LU	1.34	1.44e-10
$LDL^T$ NoPiv	0.69	1.37e-10

Table 2: Car motor (matrix size is 15,135 in single complex precision).

	Time (sec)	Scaled residual
LU	3.74	7.46e-11
$LDL^T$ NoPiv	1.93	9.28e-11

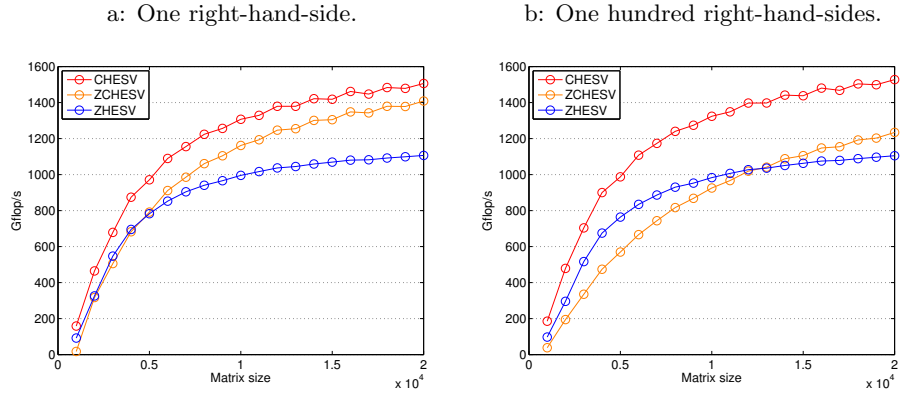


Figure 7: Performance of the standard and mixed-precision solvers: **CHESV** and **ZHESV** is the standard solver in single and double complex precision, while **ZCHESV** is the mixed-precision solver.

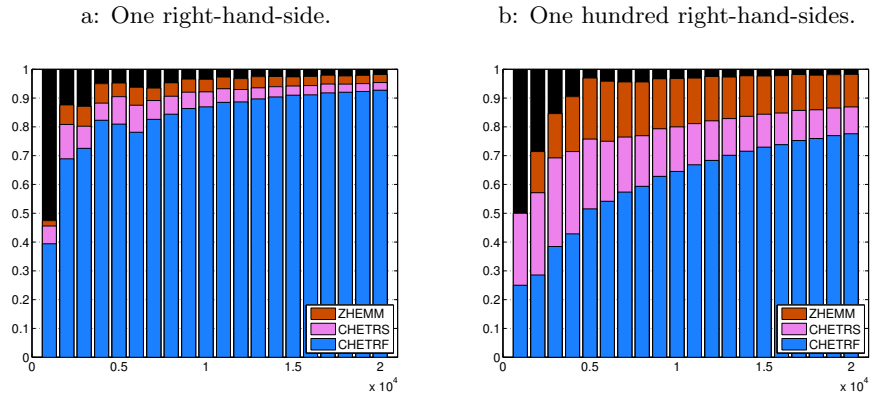


Figure 8: Time breakdown of the mixed-precision solver: **CHETRF** compute the  $LDL^T$  factorization in single complex precision, while **CHETRS** and **ZHEMM** are used for iterative refinement to compute the solution with the  $LDL^T$  factors and to perform the matrix-matrix multiply, respectively. See Figure 2 for the pseudocode.

forward and backward substitutions, where  $n$  is the dimension of  $A$  and  $n_{rhs}$  is the number of right-hand-sides). The iterative refinement converged in two iterations to obtain the accuracy of the double precision. As we expected, for a large enough matrix, the mixed-precision algorithm obtained a performance close to that in single precision (e.g., for  $n = 20,000$ , the single precision and mixed-precision solvers were about  $1.36\times$  and  $1.27\times$  faster than the double precision solver, respectively).

For these experiments, we used the  $LDL^T$  factorization without pivoting. This is motivated by our observation that in many real applications, the pivoting is not needed in most of the cases. In the rare case of the failure, the iterative refinement would not converge, signaling the need for pivoting. When this happens, the user can fall back on a stable algorithm like Bunch and Kaufman's. We can easily integrate the RBT into the mixed-precision solver in order to reduce the probability of encountering small diagonal entries.

In Figure 7, we observe that the performance benefit of using the mixed-precision algorithm decreases as the number of right-hand-sides increases. This is due to the increase in the relative overhead associated with the residual computation in double precision compared to the factorization cost. This can be also seen in Figure 8, where the time spent by the double-precision arithmetic increases (e.g., ZHEMM).

## 4 Conclusion

We presented the performance of dense symmetric indefinite solvers on hybrid GPU+CPU machines for which until recently, there were no implementations of the algorithms that can utilize the GPU. The symmetric pivoting required to maintain the numerical stability of the factorization leads to frequent synchronizations and exhibits irregular memory accesses which are difficult to optimize on a GPU. We investigated several techniques to reduce the expensive communication required for pivoting (e.g., native GPU and communication-avoiding implementations). Unfortunately, the overhead associated with the symmetric pivoting can still be significant. However, these algorithms preserve the symmetry, which is required in several physical applications, and reduces the runtime and memory requirement for the rest of the application software. The randomization using RBT followed by an  $LDL^T$  factorization without pivoting outperforms other algorithms, and is about twice as fast as the LU factorization. We also presented experimental results for acoustic scattering problems where there is no need for pivoting

and how mixed precision can be used to enhance performance. Our current implementations are based on standard BLAS/LAPACK routines, and we are improving the performance of factorization by developing specialized GPU kernels. We point out that low-level optimizations are also provided in vendor libraries (e.g., CuSolver implementation of the Bunch-Kaufman algorithm). Our implementations have been released as a part of MAGMA software package, including the iterative refinements which use the mixed-precision arithmetics.

## Acknowledgments

The authors would like to thank the NSF (grant #ACI-1339822), NVIDIA, and MathWorks for supporting this research effort. The authors are also grateful to Nicolas Zerbib (ESI Group, Compiègne, France) for his help in using test matrices from acoustics.

## References

- [1] J. Aasen, *On the reduction of a symmetric matrix to tridiagonal form*, BIT **11** (1971), 233–242.
- [2] M. Abalenkovs, A. Abdelfattah, J. Dongarra, M. Gates, A. Haidar, J. Kurzak, P. Luszczek, S. Tomov, I. Yamazaki, and A. YarKhan. Parallel programming models for dense linear algebra on heterogeneous systems. *Supercomputing Frontiers and Innovations*, 2(4), 10-2015 2015.
- [3] A. Abdelfattah, A. Haidar, S. Tomov, and J. Dongarra. Performance, Design, and Autotuning of Batched GEMM for GPUs. The International Supercomputing Conference (ISC High Performance 2016), 06-2016, Frankfurt, Germany.
- [4] E. Anderson, Z. Bai, J. J. Dongarra, A. Greenbaum, A. McKenney, J. Du Croz, S. Hammarling, J. W. Demmel, C. Bischof, D. Sorensen, *LAPACK: a portable linear algebra library for high-performance computers*, Proceedings of the 1990 ACM/IEEE conference on Supercomputing.
- [5] Intel, *Math Kernel Library (MKL)*, <http://www.intel.com/software/products/mkl/>.

- [6] C. Ashcraft, R. G. Grimes, and J. G. Lewis. Accurate symmetric indefinite linear equation solvers. *SIAM J. Matrix Anal. and Appl.*, 20(2):513–561, 1998.
- [7] M. Baboulin, J. J. Dongarra, J. Hermann, and S. Tomov, *Accelerating Linear System Solutions using Randomization Techniques*, ACM Transactions on Mathematical Software, 39(2), 2013.
- [8] M. Baboulin, D. Becker, G. Bosilca, A. Danalis, and J. J. Dongarra, *An efficient distributed randomized algorithm for solving large dense symmetric indefinite linear systems*, Parallel Computing, 40(7):213–223, 2014.
- [9] M. Baboulin, X. S. Li, and F-H. Rouet, *Using Random Butterfly Transformations to Avoid Pivoting in Sparse Direct Methods*, Proceedings of International Conference on Vector and Parallel Processing (VecPar 2014), Eugene (OR), USA.
- [10] M. Baboulin, D. Becker, and J. J. Dongarra, *A Parallel Tiled Solver for Dense Symmetric Indefinite Systems on Multicore Architectures*, Parallel & Distributed Processing Symposium (IPDPS), 2012.
- [11] G. Ballard, D. Becker, J. Demmel, J. Dongarra, A. Druinsky, I. Peled, O. Schwartz, S. Toledo, and I. Yamazaki, *Communication-avoiding symmetric-indefinite factorization*, SIAM J. Matrix Anal. Appl. **35** (2014), 1364–1460.
- [12] L. Blackford, J. Choi, A. Cleary, E. D’Azevedo, J. W. Demmel, I. Dhillon, J. J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. Whaley, *ScaLAPACK Users Guide*, SIAM, 1997.
- [13] J. R. Bunch and B. N. Parlett. Direct methods for solving symmetric indefinite systems of linear equations. *SIAM J. Numerical Analysis*, 8:639–655, 1971.
- [14] J. R. Bunch and L. Kaufman. Some stable methods for calculating inertia and solving symmetric linear systems. *Math. Comput.*, 31:163–179, 1977.
- [15] G. Ballard, D. Becker, J. Demmel, J. Dongarra, A. Druinsky, I. Peled, O. Schwartz, S. Toledo, and I. Yamazaki, *Implementing a blocked Aasen’s algorithm with a dynamic scheduler on multicore architectures*, Proceedings of the 27th international symposium on parallel and distributed processing, 2013, pp. 895–907.

- [16] I. Yamazaki, S. Tomov, and J. Dongarra, *Non-GPU-resident Dense Symmetric Indefinite Factorization, Concurrency and Computation: Practice and Experience*, 2016.
- [17] D. Becker, M. Baboulin, and J. J. Dongarra, *Reducing the amount of pivoting in symmetric indefinite systems*, Proceedings of the 9th International Conference on Parallel Processing and Applied Mathematics (PPAM 2011), 133–142, 2012.
- [18] Å. BJÖRCK, *Numerical Methods for Least Squares Problems*, SIAM, 1996.
- [19] A. Buttari, J. Dongarra, J. Langou, J. Langou, P. Luszczek, J. Kurzak, Mixed Precision Iterative Refinement Techniques for the Solution of Dense Linear Systems. *Int. J. High Perform. Comput. Appl.*, 21: 457–466, 2007.
- [20] M. Baboulin and A. Buttari and J. Dongarra and J. Kurzak and J. Langou and J. Langou and P. Luszczek and S. Tomov, *Accelerating scientific computations with mixed precision algorithms*, Computer Physics Communications **180** (2009), no. 12, 2526–2533.
- [21] G. H. Golub and C. F. Van Loan, *Matrix computations*, The Johns Hopkins University Press, Baltimore, 1996, Third edition.
- [22] A. Castaldo and R. Whaley, *Scaling LAPACK panel operations using parallel cache assignment*, Proceedings of the 15th AGM SIGPLAN symposium on principle and practice of parallel programming, 2010, pp. 223–232.
- [23] J. Demmel, L. Grigori, M. Hoemmen, and J. Langou, *Communication-optimal parallel and sequential QR and LU factorizations*, SIAM J. Sci. Comput. **34** (2012), A206–A239, , also available as EECS Department, University of California, Berkeley, Technical report (UCB/EECS-2008-89).
- [24] L. Grigori, J. Demmel, and H. Xiang, *CALU: a communication optimal LU factorization algorithm*, SIAM. J. Matrix Anal. Appl. **32** (2011), no. 4, 1317–1350.
- [25] J. W. Demmel, *Applied Numerical Linear Algebra*, SIAM, 1997

- [26] F. Gustavson, *Recursive leads to automatic variable blocking for dense linear-algebra algorithms*, IBM Journal of Research and Development **41** (1997), 737–755.
- [27] N. J. Higham, *Accuracy and stability of numerical algorithms*, SIAM, 2002.
- [28] J-C. Nédélec, *Acoustic and electromagnetic equations. Integral representations for harmonic problems*, Appl. Math. Sci., vol. 144, Springer-Verlag, New-York, 2001.
- [29] D. S. Parker, *Random Butterfly Transformations with Applications in Computational Linear Algebra*, Technical Report CSD-950023, UCLA Computer Science Department, 1995.
- [30] M. Rozložník, G. Shklarski, and S. Toledo, *Partitioned triangular tridiagonalization*, ACM Trans. Math. Softw. **37** (2011), no. 4, 1–16.
- [31] G. W. Stewart, *Introduction to Matrix Computations*, Academic Press, 1973
- [32] M. Baboulin and J. Dongarra and J. Demmel and S. Tomov and V. Volkov, *Enhancing the performance of dense linear algebra solvers on GPUs in the MAGMA project*, Poster at Supercomputing (SC'08), Austin, 2008.
- [33] S. Tomov and J. Dongarra, and M. Baboulin, *Towards dense linear algebra for hybrid GPU accelerated manycore systems*, Parallel Computing, 36(5&6):232–240, 2010.
- [34] S. Toledo, *Locality of reference in LU decomposition with partial pivoting*, SIAM J. Matrix Anal. Appl. **18** (1997), no. 4, 1065–1081.
- [35] J. H. Wilkinson, *Rounding Errors in Algebraic Processes*, Prentice-Hall, 1963.
- [36] C. B. Moler, *Iterative Refinement in Floating Point*, J. ACM **14** (1967), no. 2, 316–321.
- [37] *PLASMA Users' Guide, Parallel Linear Algebra Software for Multicore Architectures*, Version 2.3, 2010. University of Tennessee.
- [38] M. Baboulin, J. Dongarra, A. Rémy, S. Tomov, I. Yamazaki, *Dense Symmetric Indefinite Factorization on GPU Accelerated Architectures*,



Proceedings of 11th International Conference on Parallel Processing and Applied Mathematics (PPAM 2015), 2015.

- [39] J. Dongarra, J. Kurzak, P. Luszczek, T. Moore, and S. Tomov. Numerical algorithms and libraries at exascale. <http://www.hpcwire.com/2015/10/19/numerical-algorithms-and-libraries-at-exascale/>, October 19 2015. HPCwire.
- [40] R. Nath, S. Tomov, and J. Dongarra. An improved MAGMA GEMM for Fermi graphics processing units. *Int. J. High Perform. Comput. Appl.*, 24(4):511–515, Nov. 2010.
- [41] Y. Yan, B. M. Chapman, and M. Wong. A comparison of heterogeneous and manycore programming models. <http://www.hpcwire.com/2015/03/02/a-comparison-of-heterogeneous-and-manycore-programming-models>, March 2 2015. HPCwire.
- [42] J. Dongarra, M. Gates, A. Haidar, J. Kurzak, P. Luszczek, S. Tomov, and I. Yamazaki. Accelerating numerical dense linear algebra calculations with GPUs. *Numerical Computations with GPUs*, pages 1–26, 2014.
- [43] A. Haidar, C. Cao, I. Yamazaki, J. Dongarra, M. Gates, P. Luszczek, and S. Tomov. Performance and Portability with OpenCL for Throughput-Oriented HPC Workloads Across Accelerators, Coprocessors, and Multicore Processors. In *5th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems (Scala 14)*, New Orleans, LA, 11-2014 2014. IEEE.
- [44] A. Haidar, J. Dongarra, K. Kabir, M. Gates, P. Luszczek, S. Tomov, and Y. Jia. HPC programming on Intel many-integrated-core hardware with MAGMA port to Xeon Phi. *Scientific Programming*, 23, 01-2015 2015.