

IBM Research Report

Rectangular Full Packed Format for Cholesky's Algorithm: Factorization, Solution, and Inversion

Fred G. Gustavson

IBM Research Division

Thomas J. Watson Research Center

P.O. Box 218

Yorktown Heights, NY 10598

USA

Jerzy Wasniewski

Technical University of Denmark

Richard Petersens Plads, Building 321

DK-2800 Kongens Lyngby

Denmark

Jack J. Dongarra

University of Tennessee

1122 Volunteer Blvd.

Knoxville, TN 37996-3450

USA

Julien Langou

University of Colorado Denver

1250 14th Street - Room 646

Denver, CO 80202

USA



Research Division

Almaden - Austin - Beijing - Cambridge - Haifa - India - T. J. Watson - Tokyo - Zurich

Rectangular Full Packed Format for Cholesky's Algorithm: Factorization, Solution and Inversion

Fred G. Gustavson

IBM T.J. Watson Research Center

and

Jerzy Waśniewski

Technical University of Denmark

and

Jack J. Dongarra

University of Tennessee, Oak Ridge National Laboratory and University of Manchester

and

Julien Langou

University of Colorado Denver

We describe a new data format for storing triangular, symmetric, and Hermitian matrices called RFPF (Rectangular Full Packed Format). The standard two dimensional arrays of Fortran and C (also known as full format) that are used to represent triangular and symmetric matrices waste nearly half of the storage space but provide high performance via the use of Level 3 BLAS. Standard packed format arrays fully utilize storage (array space) but provide low performance as there is no Level 3 packed BLAS. We combine the good features of packed and full storage using RFPF to obtain high performance via using Level 3 BLAS as RFPF is a standard full format representation. Also, RFPF requires exactly the same minimal storage as packed format. Each LAPACK full and/or packed triangular, symmetric, and Hermitian routine becomes a single new RFPF routine based on eight possible data layouts of RFPF. This new RFPF routine usually consists of two calls to the corresponding LAPACK full format routine and two calls to Level 3 BLAS routines. This means *no* new software is required. As examples, we present LAPACK routines for Cholesky factorization, Cholesky solution and Cholesky inverse computation in RFPF to illustrate this new work and to describe its performance on several commonly used computer platforms. Performance of LAPACK full routines using RFPF versus LAPACK full routines using standard format for both serial and SMP parallel processing is about the same while using half the storage. Performance gains are roughly one to a factor of 43 for serial and one to a factor of 97 for SMP parallel times faster using vendor LAPACK full routines with RFPF than with using vendor and/or reference packed routines.

Categories and Subject Descriptors: G.1.3 [**Numerical Analysis**]: Numerical Linear Algebra – Linear Systems (symmetric and Hermitian); G.4 [**Mathematics of Computing**]: Mathematical Software

General Terms: Algorithms, BLAS, Performance, Linear Algebra Libraries

Authors' addresses: F.G. Gustavson, IBM T.J. Watson Research Center, Yorktown Heights, NY-10598, USA, email: fg2@us.ibm.com; J. Waśniewski, Department of Informatics and Mathematical Modelling, Technical University of Denmark, Richard Petersens Plads, Building 321, DK-2800 Kongens Lyngby, Denmark, email: jw@imm.dtu.dk; J.J. Dongarra, Electrical Engineering and Computer Science Department, University of Tennessee, 1122 Volunteer Blvd, Knoxville, TN 37996-3450, USA, email: dongarra@eecs.utk.edu; Julien Langou, Department of Mathematical and Statistical Sciences, University of Colorado Denver, 1250, 14th Street – Room 646, Denver, Colorado 80202, USA, email: julien.langou@ucdenver.edu.

Additional Key Words and Phrases: real symmetric matrices, complex Hermitian matrices, positive definite matrices, Cholesky factorization and solution, recursive algorithms, novel packed matrix data structures, LAPACK, Rectangular Full Packed Format

1. INTRODUCTION

A very important class of linear algebra problems deals with a coefficient matrix A that is symmetric and positive definite [Dongarra et al. 1998; Demmel 1997; Golub and Van Loan 1996; Trefethen and Bau 1997]. Because of symmetry it is only necessary to store either the upper or lower triangular part of the matrix A .

Fig. 1. The **full** format array layout of an order N symmetric matrix required by LAPACK. LAPACK requires $LDA \geq N$. Here we set $LDA=N=7$.

Lower triangular case	Upper triangular case
$\begin{pmatrix} 1 & & & & & & \\ 2 & 9 & & & & & \\ 3 & 10 & 17 & & & & \\ 4 & 11 & 18 & 25 & & & \\ 5 & 12 & 19 & 26 & 33 & & \\ 6 & 13 & 20 & 27 & 34 & 41 & \\ 7 & 14 & 21 & 28 & 35 & 42 & 49 \end{pmatrix}$	$\begin{pmatrix} 1 & 8 & 15 & 22 & 29 & 36 & 43 \\ & 9 & 16 & 23 & 30 & 37 & 44 \\ & & 17 & 24 & 31 & 38 & 45 \\ & & & 25 & 32 & 39 & 46 \\ & & & & 33 & 40 & 47 \\ & & & & & 41 & 48 \\ & & & & & & 49 \end{pmatrix}$

Fig. 2. The **packed** format array layout of an order 7 symmetric matrix required by LAPACK.

Lower triangular case	Upper triangular case
$\begin{pmatrix} 1 & & & & & & \\ 2 & 8 & & & & & \\ 3 & 9 & 14 & & & & \\ 4 & 10 & 15 & 19 & & & \\ 5 & 11 & 16 & 20 & 23 & & \\ 6 & 12 & 17 & 21 & 24 & 26 & \\ 7 & 13 & 18 & 22 & 25 & 27 & 28 \end{pmatrix}$	$\begin{pmatrix} 1 & 2 & 4 & 7 & 11 & 16 & 22 \\ & 3 & 5 & 8 & 12 & 17 & 23 \\ & & 6 & 9 & 13 & 18 & 24 \\ & & & 10 & 14 & 19 & 25 \\ & & & & 15 & 20 & 26 \\ & & & & & 21 & 27 \\ & & & & & & 28 \end{pmatrix}$

1.1 LAPACK full and packed storage formats

The LAPACK library [Anderson et al. 1999] offers two different kinds of subroutines to solve the same problem: POTRF¹ and PPTRF both factorize symmetric,

¹Four names SPOTRF, DPOTRF, CPOTRF and ZPOTRF are used in LAPACK for real symmetric and complex Hermitian matrices [Anderson et al. 1999], where the first character indicates the precision and arithmetic versions: S – single precision, D – double precision, C – complex and Z – double complex. LAPACK95 uses one name LA.POTRF for all versions [Barker et al. 2001]. In this paper, POTRF and/or PPTRF express, any precision, any arithmetic and any language version of the PO and/or PP matrix factorization algorithms.

positive definite matrices by means of the Cholesky algorithm. A major difference in these two routines is the way they access the array holding the triangular matrix (see Figures 1 and 2).

In the POTRF case, the matrix is stored in one of the lower left or upper right triangles of a full square matrix ([Anderson et al. 1999, pages 139 and 140] and [IBM 1997, page 64])², the other triangle is wasted (see Figure 1). Because of the uniform storage scheme, blocked LAPACK and Level 3 BLAS subroutines [Dongarra et al. 1990b; Dongarra et al. 1990a] can be employed, resulting in a fast solution.

In the PPTRF case, the matrix is stored in *packed* storage ([Anderson et al. 1999, pages 140 and 141], [Agarwal et al. 1994] and [IBM 1997, pages 74 and 75]), which means that the columns of the lower or upper triangle are stored consecutively in a one dimensional array (see Figure 2). Now the triangular matrix occupies the strictly necessary storage space but the nonuniform storage scheme means that use of full storage BLAS is impossible and only the Level 2 BLAS packed subroutines [Lawson et al. 1979; Dongarra et al. 1988] can be employed, resulting in a slow solution.

To summarize: LAPACK offers a choice between high performance and wasting half of the memory space (POTRF) versus low performance with optimal memory space (PPTRF).

1.2 Packed Minimal Storage Data Formats related to RFPF

Recently many new data formats for matrices have been introduced for improving the performance of Dense Linear Algebra (DLA) algorithms. The survey article [Elmroth et al. 2004] gives an excellent overview.

Recursive Packed Format (RPF) [Andersen et al. 2001; Andersen et al. 2002]: A new compact way to store a triangular, symmetric or Hermitian matrix called Recursive Packed Format is described in [Andersen et al. 2001] as are novel ways to transform RPF to and from standard packed format. New algorithms, called Recursive Packed Cholesky (RPC) [Andersen et al. 2001; Andersen et al. 2002] that operate on the RPF format are presented. RPF format operates almost entirely by calling Level 3 BLAS GEMM [Dongarra et al. 1990b; Dongarra et al. 1990a] but requires variants of algorithms TRSM and SYRK [Dongarra et al. 1990b; Dongarra et al. 1990a] that are designed to work on RPF. The authors call these algorithms RPTRSM and RPSYRK [Andersen et al. 2001] and find that they do most of their FLOPS by calling GEMM [Dongarra et al. 1990b; Dongarra et al. 1990a]. It follows that almost all of execution time of the RPC algorithm is done in calls to GEMM. There are three advantages of this storage scheme compared to traditional packed and full storage. First, the RPF storage format uses the minimum amount of storage required for symmetric, triangular, or Hermitian matrices. Second, the RPC algorithm is a Level 3 implementation of Cholesky factorization. Finally, RPF requires no block size tuning parameter. A disadvantage of the RPC algorithm was that it had a high recursive calling overhead. The paper [Gustavson and Jonsson 2000] removed this overhead and added other novel features to the RPC algorithm.

Square Block Packed Format (SBPF) [Gustavson 2003]: SBPF is described in Section 4 of [Gustavson 2003]. A strong point of SBPF is that it requires mini-

²In Fortran column major, in C row major.

mum block storage and all its blocks are contiguous and of equal size. If one uses SBPF with kernel routines then data copying is mostly eliminated during Cholesky factorization.

Block Packed Hybrid Format (BPHF) [Andersen et al. 2005; Gustavson et al. 2007]: We consider an efficient implementation of the Cholesky solution of symmetric positive-definite full linear systems of equations using packed storage. We take the same starting point as that of LINPACK [Dongarra et al. 1979] and LAPACK [Anderson et al. 1999], with the upper (or lower) triangular part of the matrix being stored by columns. Following LINPACK [Dongarra et al. 1979] and LAPACK [Anderson et al. 1999], we overwrite the given matrix by its Cholesky factor. The paper [Andersen et al. 2005] uses the BPHF where blocks of the matrix are held contiguously. The paper compares BPHF versus conventional full format storage, packed format and the RPF for the algorithms. BPF is a variant of SBPF in which the diagonal blocks are stored in packed format and so its storage requirement is equal to that of packed storage.

We mention that for packed matrices SBPF and BPHF have become the format of choice for multicore processors when one stores the blocks in register block format [Gustavson et al. 2007]. Recently, there have been many papers published on new algorithms for multicore processors. This literature is extensive. So, we only mention two projects, PLASMA [Buttari et al. 2007] and FLAME [Chan et al. 2007], and refer the interested reader to the literature for additional references.

In regard to other references on new data structures, the survey article [Elmroth et al. 2004] gives an excellent overview. However, since 2005 at least two new data formats for Cholesky type factorizations have emerged, [Herrero 2006] and the subject matter of this paper, RFPF [Gustavson and Waśniewski 2007]. In the next subsection we highlight the main features of RFPF.

1.3 A novel way of representing triangular, symmetric, and Hermitian matrices in LAPACK

LAPACK has two types of subroutines for triangular, symmetric, and Hermitian matrices called packed and full format routines. LAPACK has about 300 these kind of subroutines. So, in either format, a variety of problems can be solved by these LAPACK subroutines. From a user point of view, RFPF can replace both these LAPACK data formats. Furthermore, and this is important, using RFPF does not require any new LAPACK subroutines to be written. Using RFPF in LAPACK only requires the use of already existing LAPACK and BLAS routines. RFPF strongly relies on the existence of the BLAS and LAPACK routines for full storage format.

1.4 Overview of the Paper

First we introduce the RFPF in general, see Section 2. Secondly we show how to use RFPF on symmetric and Hermitian positive definite matrices; e.g., for the factorization (Section 3), solution (Section 4), and inversion (Section 5) of these matrices. Section 6 describes LAPACK subroutines for the Cholesky factorization, Cholesky solution, and Cholesky inversion of symmetric and Hermitian positive definite matrices using RFPF. Section 7 indicates that the stability results of using RFPF is unaffected by this format choice as RFPF uses existing LAPACK

algorithms which are already known to be stable. Section 8 describes a variety of performance results on commonly used platforms both for serial and parallel SMP execution. These results show that performance of LAPACK full routines using RFPF versus LAPACK full routines using standard format for both serial and SMP parallel processing is about the same while using half the storage. Also, performance gains are roughly one to a factor of 43 for serial and one to a factor of 97 for SMP parallel times faster using vendor LAPACK full routines with RFPF than with using vendor and/or reference packed routines. Section 9 explains how some new RFPF routines have been integrated in LAPACK. LAPACK software for Cholesky algorithm (factorization, solution and inversion) using RFPF has been released with LAPACK-3.2 on November 2008. Section 10 gives a short summary and brief conclusions.

2. DESCRIPTION OF RECTANGULAR FULL PACKED FORMAT

We describe Rectangular Full Packed Format (RFPF). It transforms a standard Packed Array AP of size $NT = N(N + 1)/2$ to a full 2D array. This means that performance of LAPACK's [Anderson et al. 1999] packed format routines becomes equal to or better than their full array counterparts. RFPF is a variant of Hybrid Full Packed (HFP) format [Gunnels and Gustavson 2004]. RFPF is a rearrangement of a Standard full format rectangular Array SA of size $LDA * N$ where $LDA \geq N$. Array SA holds a triangular part of a symmetric, triangular, or Hermitian matrix A of order N. The rearrangement of array SA is equal to compact full format Rectangular Array AR of size $LDA1 * N1 = NT$ and hence array AR like array AP uses minimal storage. (The specific values of LDA1 and N1 can vary depending on various cases and they will be specified later during the text.) Array AR will hold a full rectangular matrix A_R obtained from a triangle of matrix A. Note also that the transpose of the rectangular matrix A_R^T resides in the transpose of array AR and hence also represents A. Therefore, Level 3 BLAS [Dongarra et al. 1990b; Dongarra et al. 1990a] can be used on array AR or its transpose. In fact, with the equivalent LAPACK algorithm which uses the array AR or its transpose, the performance is slightly better than standard LAPACK algorithm which uses the array SA or its transpose. Therefore, this offers the possibility to replace all packed or full LAPACK routines with equivalent LAPACK routines that work on array AR or its transpose. For examples of transformations of a matrix A to a matrix A_R see the figures in Section 6.

RFPF is closely related to HFP format, see [Gunnels and Gustavson 2004], which represents A as the concatenation of two standard full arrays whose total size is also NT . A basic simple idea leads to both formats. Let A be an order N symmetric matrix. Break A into a block 2-by-2 form

$$A = \begin{bmatrix} A_{11} & A_{21}^T \\ A_{21} & A_{22} \end{bmatrix} \text{ or } A = \begin{bmatrix} A_{11} & A_{12} \\ A_{12}^T & A_{22} \end{bmatrix} \quad (1)$$

where A_{11} and A_{22} are symmetric. Clearly, we need only store the lower triangles of A_{11} and A_{22} as well as the full matrix $A_{21} = A_{12}^T$ when we are interested in a lower triangular formulation.

When $N = 2k$ is even, the lower triangle of A_{11} and the upper triangle of A_{22}^T

can be concatenated together along their main diagonals into a $(k + 1)$ -by- k dense matrix (see the figures where N is even in Section 6). This last operation is the crux of the basic simple idea. The off-diagonal block A_{21} is k -by- k , and so it can be appended below the $(k + 1)$ -by- k dense matrix. Thus, the lower triangle of A can be stored as a single $(N + 1)$ -by- k dense matrix A_R . In effect, each block matrix A_{11} , A_{21} and A_{22} is now stored in ‘full format’. This means all entries of matrix A_R in array **AR** of size $\text{LDA1} = N + 1$ by $\text{N1} = k$ can be accessed with constant row and column strides. So, the full power of LAPACK’s block Level 3 codes are now available for RFPF which uses the minimum amount of storage. Finally, matrix A_R^T which has size k -by- $(N + 1)$ is represented in the transpose of array **AR** and hence has the same desirable properties. There are eight representations of RFPF. The matrix A can have either odd or even order N , or it can be represented either in standard lower or upper format or it can be represented by either matrix A_R or its transpose A_R^T giving $2^3 = 8$ representations in all.

All eight cases or representations are presented in Section 6. The RFPF matrices are in the upper right part of the figures. We have introduced colors and horizontal lines to try to visually delineate triangles T_1 , T_2 representing lower, upper triangles of symmetric matrices A_{11} , A_{22}^T respectively and square or near square S_1 representing matrices A_{21} . For an upper triangle of A , T_1 , T_2 represents lower, upper triangles of symmetric matrices A_{11}^T , A_{22} respectively and square or near square S_1 representing matrices A_{12} . For both lower and upper triangles of A we have, after each $a_{i,j}$, added its position location in the arrays holding matrices A and A_R .

We now consider performance aspects of using RFPF in the context of using LAPACK routines on triangular matrices stored in RFPF. Let X be a Level 3 LAPACK routine that operates either on full format. X has a full Level 3 LAPACK block 2-by-2 algorithm, call it FX . We write a simple related partition algorithm (SRPA) with partition sizes $n1$ and $n2$ where $n1 + n2 = N$. Apply the new SRPA using the new RFPF. The new SRPA almost always has four major steps consisting entirely of calls to existing full format LAPACK routines in two steps and calls to Level 3 BLAS in the remaining two steps, see Figure 3.

Fig. 3. Simple related partition algorithm (SRPA) of RFPF

call X('L',n1,T1,ldt) ! step 1	call L3BLAS(n1,n2,S,lds,'U',T2,ldt) ! step 3
call L3BLAS(n1,n2,'L',T1,ldt,S,lds) ! step 2	call X('U',n2,T2,ldt) ! step 4

Section 6 shows FX algorithms equal to factorization, solution and inversion algorithms on symmetric positive definite or Hermitian matrices.

3. CHOLESKY FACTORIZATION USING RECTANGULAR FULL PACKED FORMAT

The Cholesky factorization of a symmetric and positive definite matrix A can be expressed as

$$\begin{aligned} A &= LL^T \text{ or } A = U^T U \text{ (in the symmetric case)} \\ A &= LL^H \text{ or } A = U^H U \text{ (in the Hermitian case)} \end{aligned} \quad (2)$$

where L and U are lower triangular and upper triangular matrices.

Break the matrices L and U into 2-by-2 block form in the same way as was done for the matrix A in Equation (1):

$$L = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \text{ and } U = \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix} \quad (3)$$

We now have

$$LL^T = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} L_{11}^T & L_{21}^T \\ 0 & L_{22}^T \end{bmatrix} \text{ and } U^T U = \begin{bmatrix} U_{11}^T & 0 \\ U_{12}^T & U_{22}^T \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix} \quad (4)$$

$$LL^H = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} L_{11}^H & L_{21}^H \\ 0 & L_{22}^H \end{bmatrix} \text{ and } U^H U = \begin{bmatrix} U_{11}^H & 0 \\ U_{12}^H & U_{22}^H \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix}$$

where L_{11} , L_{22} , U_{11} , and U_{22} are lower and upper triangular submatrices, and L_{21} and U_{12} are square or almost square submatrices.

Using Equations (2) and equating the blocks of Equations (1) and Equations (4) gives us the basis of a 2-by-2 block algorithm for Cholesky factorization using RFPPF. We can now express each of these four block equalities by calls to existing LAPACK and Level 3 BLAS routines. An example, see Section 6, of this is the three block equations is $L_{11}L_{11}^T = A_{11}$, $L_{21}L_{11}^T = A_{21}$ and $L_{21}L_{21}^T + L_{22}L_{22} = A_{22}$. The first and second of these block equations are handled by calling LAPACK's POTRF routine $L_{11} \leftarrow A_{11}$ and by calling Level 3 BLAS TRSM via $L_{21} \leftarrow L_{21}L_{11}^{-T}$. In both these block equations the Fortran equality of replacement (\leftarrow) is being used so that the lower triangle of A_{11} is being replaced L_{11} and the nearly square matrix A_{21} is being replaced by L_{21} . The third block equation breaks into two parts: $A_{22} \leftarrow L_{21}L_{21}^T$ and $L_{22} \leftarrow A_{22}$ which are handled by calling Level 3 BLAS SYRK or HERK and by calling LAPACK's POTRF routine. At this point we can use the flexibility of the LAPACK library. In RFPPF A_{22} is in upper format (upper triangle) while in standard format A_{22} is in lower format (lower triangle). Due to symmetry, both formats of A_{22} contain equal values. This flexibility allows LAPACK to accommodate both formats. Hence, in the calls to SYRK or HERK and POTRF we set uplo = 'U' even though the rectangular matrix of SYRK and HERK comes from a lower triangular formulation.

New LAPACK like routine PFTRF performs these four computations. PF was chosen to fit with LAPACK's use of PO and PP. The PFTRF routine covers the Cholesky Factorization algorithm for the eight cases of the RFPPF. Section 6 has Figure 4 with four subfigures. Here we are interested in the first and second subfigure. The first subfigure contains the layouts of matrices A and A_R . The second subfigure has the Cholesky factorization algorithm obtained by simple algebraic manipulations of the three block equalities obtained above.

4. SOLUTION

In Section 3 we obtained the 2-by-2 Cholesky factorization (3) of matrix A . Now, we can solve the equation $AX = B$:

- If A has lower triangular format then

$$\begin{aligned} LY = B \text{ and } L^T X = Y \text{ (in the symmetric case)} \\ LY = B \text{ and } L^H X = Y \text{ (in the Hermitian case)} \end{aligned} \quad (5)$$

- If A has an upper triangular format then

$$\begin{aligned} U^T Y = B \text{ and } UX = Y \text{ (in the symmetric case)} \\ U^H Y = B \text{ and } UX = Y \text{ (in the Hermitian case)} \end{aligned} \quad (6)$$

B , X and Y are either vectors or rectangular matrices. B contains the RHS values. X and Y contain the solution values. B , X and Y are vectors when there is one RHS and matrices when there are many RHS. The values of X and Y are stored over the values of B .

Expanding (5) and (6) using (3) gives the forward substitution equations

$$\begin{aligned} & \text{in the symmetric case:} \\ & \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \text{ and } \begin{bmatrix} U_{11}^T & 0 \\ U_{12}^T & U_{22}^T \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \end{aligned} \quad (7)$$

$$\begin{aligned} & \text{and in the Hermitian case:} \\ & \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \text{ and } \begin{bmatrix} U_{11}^H & 0 \\ U_{12}^H & U_{22}^H \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \end{aligned}$$

and the back substitution equations

$$\begin{aligned} & \text{in the symmetric case:} \\ & \begin{bmatrix} L_{11}^T & L_{21}^T \\ 0 & L_{22}^T \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \text{ and } \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \end{aligned} \quad (8)$$

$$\begin{aligned} & \text{and in the Hermitian case:} \\ & \begin{bmatrix} L_{11}^H & L_{21}^H \\ 0 & L_{22}^H \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \text{ and } \begin{bmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}. \end{aligned}$$

The Equations (7) and (8) gives the basis of a 2×2 block algorithm for Cholesky solution using RFPF format. We can now express these two sets of two block equalities by using existing Level 3 BLAS routines. An example, see Section 6, of the first set of these two block equalities is $L_{11}Y_1 = B_1$ and $L_{21}Y_1 + L_{22}Y_2 = B_2$. The first block equality is handled by Level 3 BLAS TRSM: $Y_1 \leftarrow L_{11}^{-1}B_1$. The second block equality is handled by Level 3 BLAS GEMM and TRSM: $B_2 \leftarrow B_2 - L_{21}Y_1$ and $Y_2 \leftarrow L_{22}^{-1}B_2$. The backsolution routines are similarly derived. One gets $X_2 \leftarrow L_{22}^{-T}Y_2$, $Y_1 \leftarrow Y_1 - L_{21}^T X_2$ and $X_1 \leftarrow L_{11}^{-T}Y_1$.

New LAPACK like routine PFTRS performs these two solution computations for the eight cases of RFPF. PFTRS calls a new Level 3 BLAS TFMSM in the same way that POTRS calls TRSM. The third subfigure in Section 6 gives the Cholesky solution algorithm using RFPF obtained by simple algebraic manipulation of the block Equations (7) and (8).

5. INVERSION

Following LAPACK we consider the following three stage procedure:

- (1) Factorize the matrix A and overwrite A with either L or U by calling PFTRF; see Section 3.
- (2) Compute the inverse of either L or U . Call these matrices W or V and overwrite either L or U with them. This is done by calling new routine new LAPACK like TFTRI.
- (3) Calculate either the product $W^T W$ or $V V^T$ and overwrite either W or V with them.

As in Sections 3 and 4 we examine 2-by-2 block algorithms for the steps two and three above. In Section 3 we obtain either matrices L or U in RFPF. Like LAPACK inversion algorithms for POTRI and PPTRI, this is our starting point for our LAPACK inversion algorithm using RFPF. The LAPACK inversion algorithms for POTRI and PPTRI also follow from steps two and three above by first calling in the full case LAPACK TRTRI and then calling LAPACK LAUUM.

Take the inverse of Equation (2) and obtain

$$\begin{aligned} A^{-1} &= W^T W \text{ or } A^{-1} = V V^T \text{ (in the symmetric case)} \\ A^{-1} &= W^H W \text{ or } A^{-1} = V V^H \text{ (in the Hermitian case)} \end{aligned} \quad (9)$$

where W and V are lower and upper triangular matrices.

Using the 2-by-2 blocking for either L or U in Equation (3) we obtain the following 2-by-2 blocking for W and V :

$$W = \begin{bmatrix} W_{11} & 0 \\ W_{21} & W_{22} \end{bmatrix} \text{ and } V = \begin{bmatrix} V_{11} & V_{12} \\ 0 & V_{22} \end{bmatrix} \quad (10)$$

From the identities $WL = LW = I$ and $VU = UV = I$ and the 2-by-2 block layouts of Equations (3) and (3), we obtain three block equations for W and V which can be solved using LAPACK routines for TRTRI and Level 3 BLAS TRMM. An example, see Figure 4, of these three block equations is $L_{11}W_{11} = I$, $L_{21}W_{11} + L_{22}W_{21} = 0$ and $L_{22}W_{22} = I$. The first and third of these block equations are handled by LAPACK TRTRI routines as $W_{11} \leftarrow L_{11}^{-1}$ and $V_{22} \leftarrow U_{22}^{-1}$. In the second inverse computation we use the fact that L_{22} is equally represented by its transpose L_{22}^T which is U_{22} in RFPF. The second block equation leads to two calls to Level 3 BLAS TRMM via $L_{21} \leftarrow -L_{21}W_{11}$ and $W_{21} = W_{22}L_{21}$. In the last two block equations the Fortran equality of replacement (\leftarrow) is being used so that $W_{21} = -W_{22}L_{21}W_{11}$ is replacing L_{21} .

Now we turn to part three of the three stage LAPACK procedure above. For this we use the 2-by-2 blocks layouts of Equation (10) and the matrix multiplications indicated by following block Equations (11) giving

$$\begin{aligned} &\text{the symmetric case:} \\ W^T W &= \begin{bmatrix} W_{11}^T & W_{21}^T \\ 0 & W_{22}^T \end{bmatrix} \begin{bmatrix} W_{11} & 0 \\ W_{21} & W_{22} \end{bmatrix} \text{ and } V V^T = \begin{bmatrix} V_{11} & V_{12} \\ 0 & V_{22} \end{bmatrix} \begin{bmatrix} V_{11}^T & 0 \\ V_{12}^T & V_{22}^T \end{bmatrix} \end{aligned} \quad (11)$$

$$\begin{aligned} &\text{and the Hermitian case:} \\ W^H W &= \begin{bmatrix} W_{11}^H & W_{21}^H \\ 0 & W_{22}^H \end{bmatrix} \begin{bmatrix} W_{11} & 0 \\ W_{21} & W_{22} \end{bmatrix} \text{ and } V V^H = \begin{bmatrix} V_{11} & V_{12} \\ 0 & V_{22} \end{bmatrix} \begin{bmatrix} V_{11}^H & 0 \\ V_{12}^H & V_{22}^H \end{bmatrix} \end{aligned}$$

where W_{11} , W_{22} , V_{11} , and V_{22} are lower and upper triangular submatrices, and W_{21} and V_{12} are square or almost square submatrices. The values of the indicated block multiplications of W or V in Equation (11) are stored over the block values of W or V .

Performing the indicated 2-by-2 block multiplications of Equation (11) leads to three block matrix computations. An example, see Section 6, of these three block computations is $W_{11}^T W_{11} + W_{21}^T W_{21}$, $W_{22}^T W_{21}$ and $W_{22}^T W_{22}$. Additionally, we want to overwrite the values of these block multiplications on their original block operands. Block operand W_{11} only occurs in the (1,1) block operand computation and hence can be overwritten by a call to LAPACK LAUUM, $W_{11} \leftarrow W_{11}^T W_{11}$, followed by a call to Level 3 BLAS SYRK or HERK, $W_{11} \leftarrow W_{11} + W_{21}^T W_{21}$. Block operand W_{21} now only occurs in the (2,1) block computation and hence can be overwritten by a call to Level 3 BLAS TRMM, $W_{21} \leftarrow W_{22}^T W_{21}$. Finally, block operand W_{22} can be overwritten by a call to LAPACK LAUUM, $W_{22} \leftarrow W_{22}^T W_{22}$.

The fourth subfigure in Section 6 has the Cholesky inversion algorithms using RFPF based on the results of this Section. New LAPACK routine, PFTRI, performs this computation for the eight cases of RFPF.

6. RFP DATA FORMATS AND ALGORITHMS

This section contains three figures.

- (1) The first figure describes the RFPF (Rectangular Full Packed Format) and gives algorithms for Cholesky factorization, solution and inversion of symmetric positive definite matrices, where N is odd, `uplo = 'lower'`, and `trans = 'no transpose'`. This figure has four subfigures.
 - (a) The first subfigure depicts the lower triangle of a symmetric positive definite matrix A in **standard full** and its representation by the matrix A_R in **RFPF**.
 - (b) The second subfigure gives the RFPF Cholesky factorization algorithm and its calling sequences of the LAPACK and BLAS subroutines, see Section 3.
 - (c) The third subfigure gives the RFPF Cholesky solution algorithm and its calling sequences to the LAPACK and BLAS subroutines, see Section 4.
 - (d) The fourth subfigure in each figure gives the RFPF Cholesky inversion algorithm and its calling sequences to the LAPACK and BLAS subroutines, see Section 5.
- (2) The second figure shows the transformation from full to RFPF of all “no transform” cases.
- (3) The third figure depicts all eight cases in RFPF.

The data format for A has `LDA = N`. Matrix A_R has `LDAR = N` if N is odd and `LDAR = N + 1` if N is even and $n1$ columns where $n1 = \lceil N/2 \rceil$. Hence, matrix A_R always has `LDAR` rows and $n1$ columns. Matrix A_R^T always has $n1$ rows and `LDAR` columns and its leading dimension is equal to $n1$. Matrix A_R always has `LDAR × n1 = NT = N(N + 1)/2` elements as does matrix A_R^T .

The order N of matrix A in the first figure is seven and six or seven in the remaining two figures.

Fig. 4. The Cholesky factorization algorithm using the Rectangular Full Packed Format (RFPF) if N is odd, uplo = 'lower', and trans = 'no transpose'.

<p>A of LAPACK full data format LDA=N = 7, memory needed LDA × N = 49</p> $\begin{pmatrix} a_{1,1} & \diamond & \diamond & \diamond & \diamond & \diamond & \diamond \\ a_{2,12} & a_{2,29} & \diamond & \diamond & \diamond & \diamond & \diamond \\ a_{3,13} & a_{3,210} & a_{3,317} & \diamond & \diamond & \diamond & \diamond \\ a_{4,14} & a_{4,211} & a_{4,318} & a_{4,425} & \diamond & \diamond & \diamond \\ a_{5,15} & a_{5,212} & a_{5,319} & a_{5,426} & a_{5,533} & \diamond & \diamond \\ a_{6,16} & a_{6,213} & a_{6,320} & a_{6,427} & a_{6,534} & a_{6,641} & \diamond \\ a_{7,17} & a_{7,214} & a_{7,321} & a_{7,428} & a_{7,535} & a_{7,642} & a_{7,749} \end{pmatrix}$ <p>Matrix A</p>	<p>A_R of Rectangular full packed LDAR=N = 7, memory needed LDAR × n1 = 28</p> $\begin{pmatrix} a_{1,11} & a_{5,58} & a_{6,515} & a_{7,522} \\ a_{2,12} & a_{2,29} & a_{6,616} & a_{7,623} \\ a_{3,13} & a_{3,210} & a_{3,317} & a_{7,724} \\ a_{4,14} & a_{4,211} & a_{4,318} & a_{4,425} \\ a_{5,15} & a_{5,212} & a_{5,319} & a_{5,426} \\ a_{6,16} & a_{6,213} & a_{6,320} & a_{6,427} \\ a_{7,17} & a_{7,214} & a_{7,321} & a_{7,428} \end{pmatrix}$ <p>Matrix A_R</p>
---	---

Cholesky Factorization Algorithm ($n1 = \lceil N/2 \rceil, n2 = N - n1$) :

- | | |
|--|---|
| <p>1) factor $L_{11}L_{11}^T = A_{11}$;
 call POTRF('L', n1, AR, N, & info);</p> <p>2) solve $L_{21}L_{11}^T = A_{21}$;
 call TRSM('R', 'L', 'L', 'T', 'N', n2, & n1, one, AR, N, AR(n1 + 1, 1), N);</p> | <p>3) update $A_{22} := A_{22} - L_{21}L_{21}^T$;
 call SYRK/HERK('U', 'N', n2, n1, & -one, AR(n1 + 1, 1), N, one, AR(1, 2), N);</p> <p>4) factor $U_{22}^TU_{22} = A_{22}$;
 call POTRF('U', n2, AR(1, 2), N, & info);</p> |
|--|---|

Cholesky Solution Algorithm,

where $B(LDB, nr)$ and $LDB \geq N$ (here $LDB = N$) :

- | | |
|---|--|
| <p>$LY = B$</p> <p>1) solve $L_{11}Y_1 = B_1$;
 call TRSM('L', 'L', 'N', 'N', n1, & nr, one, AR, N, B, N);</p> <p>2) Multiply $B_2 = B_2 - L_{21}Y_1$;
 call GEMM('N', 'N', 'N', n2, nr, n1, -one, & AR(n1 + 1, 1), N, B, N, one, & B(n1 + 1, 1), N);</p> <p>3) solve $L_{22}Y_2 = B_2$;
 call TRSM('L', 'U', 'T', 'N', n2, & nr, one, AR(1, 2), N, B(n1 + 1, 1), N);</p> | <p>$L'X = Y$</p> <p>1) solve $L_{22}^TX_2 = Y_2$;
 call TRSM('L', 'U', 'N', 'N', n2, & nr, one, AR(1, 2), N, B(n1 + 1, 1), N);</p> <p>2) Multiply $Y_1 = Y_1 - L_{21}^TX_2$;
 call GEMM('T', 'N', 'N', n1, nr, n2, -one, & AR(n1 + 1, 1), N, B(n1 + 1, 1), & N, one, B, N);</p> <p>3) solve $L_{22}^TX_1 = Y_1$;
 call TRSM('L', 'L', 'T', 'N', n1, & nr, one, AR, N, B, N);</p> |
|---|--|

Cholesky Inversion Algorithm :

- | | |
|--|--|
| <p>Inversion</p> <p>1) invert $W_{11} = L_{11}^{-1}$;
 call TRTRI('L', 'N', n1, AR, N, info);</p> <p>2) Multiply $L_{21} = -L_{21}W_{11}$;
 call TRMM('R', 'L', 'N', 'N', n2, & n1, -one, AR, N, AR(n1 + 1, 1), N);</p> <p>3) invert $V_{22} = U_{22}^{-1}$;
 call TRTRI('U', 'N', n2, AR(1, 2), & N, info);</p> <p>4) invert $V_{22} = U_{22}^{-1}$;
 call TRMM('L', 'U', 'T', 'N', n2, & n1, one, AR(1, 2), N, AR(n1 + 1, 1), N);</p> | <p>Triangular matrix multiplication</p> <p>1) Triang. Mult. $W_{11} = W_{11}^TW_{11}$;
 call LAUUM('L', n1, AR, N, info);</p> <p>2) update $W_{11} = W_{11} + W_{21}^TW_{21}$;
 call SYRK/HERK('L', 'T', n1, n2, & one, AR(n1 + 1, 1), N, one, AR, N);</p> <p>3) Multiply $W_{21} = V_{22}W_{21}$;
 call TRMM('L', 'U', 'N', 'N', n2, & n1, one, AR(1, 2), N, A(n1 + 1, 1), N);</p> <p>4) Triang. Mult. $V_{11} = V_{11}V_{11}^T$;
 call LAUUM('U', n2, AR(1, 2), N, info);</p> |
|--|--|

Fig. 5. Eight two-dimensional arrays for storing the matrices A and A_R that are needed by the LAPACK subroutine POTRF (full format) and PFTRF RFPF respectively. The leading dimension LDA is N for LAPACK, and LDAR for RFPF. $LDAR = N$ for N odd, and $N + 1$ for N even. Here N is 7 or 6. The memory needed is $LDA \times N$ for full format and $LDAR \times n1 = (N + 1)N/2$ for RFPF. Here 49 and 36 for full format and 28 and 21 for RFPF. The column size of RFPF is $n1 = \lceil N/2 \rceil$, here 4 and 3.

5.1 The matrices A of order N and A_R of size $LDAR$ by $n1$, here $N = 7$.

5.1.1 Full Format		5.1.2 RFPF
$\left[\begin{array}{cccc ccc} a_{1,1} & \diamond & \diamond & \diamond & \diamond & \diamond & \diamond \\ a_{2,1} & a_{2,2} & \diamond & \diamond & \diamond & \diamond & \diamond \\ a_{3,1} & a_{3,2} & a_{3,3} & \diamond & \diamond & \diamond & \diamond \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} & \diamond & \diamond & \diamond \\ \hline a_{5,1} & a_{5,2} & a_{5,3} & a_{5,4} & a_{5,5} & \diamond & \diamond \\ a_{6,1} & a_{6,2} & a_{6,3} & a_{6,4} & a_{6,5} & a_{6,6} & \diamond \\ a_{7,1} & a_{7,2} & a_{7,3} & a_{7,4} & a_{7,5} & a_{7,6} & a_{7,7} \end{array} \right],$		$\left[\begin{array}{cccc} a_{1,1} & a_{5,5} & a_{6,5} & a_{7,5} \\ a_{2,1} & a_{2,2} & a_{6,6} & a_{7,6} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{7,7} \\ \hline a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \\ a_{5,1} & a_{5,2} & a_{5,3} & a_{5,4} \\ a_{6,1} & a_{6,2} & a_{6,3} & a_{6,4} \\ a_{7,1} & a_{7,2} & a_{7,3} & a_{7,4} \end{array} \right]$

5.2 The matrices A of order N and A_R of size $LDAR$ by $n1$, here $N = 6$.

5.2.1 Full format		5.2.2 RFPF
$\left[\begin{array}{ccc ccc} a_{1,1} & \diamond & \diamond & \diamond & \diamond & \diamond \\ a_{2,1} & a_{2,2} & \diamond & \diamond & \diamond & \diamond \\ a_{3,1} & a_{3,2} & a_{3,3} & \diamond & \diamond & \diamond \\ \hline a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} & \diamond & \diamond \\ a_{5,1} & a_{5,2} & a_{5,3} & a_{5,4} & a_{5,5} & \diamond \\ a_{6,1} & a_{6,2} & a_{6,3} & a_{6,4} & a_{6,5} & a_{6,6} \end{array} \right],$		$\left[\begin{array}{ccc} a_{4,4} & a_{5,4} & a_{6,4} \\ a_{1,1} & a_{5,5} & a_{6,5} \\ a_{2,1} & a_{2,2} & a_{6,6} \\ \hline a_{3,1} & a_{3,2} & a_{3,3} \\ a_{4,1} & a_{4,2} & a_{4,3} \\ a_{5,1} & a_{5,2} & a_{5,3} \\ a_{6,1} & a_{6,2} & a_{6,3} \end{array} \right]$

5.3 The matrices A of order N and A_R of size $LDAR$ by $n1$, here $N = 7$.

5.3.1 Full format		5.3.2 RFPF
$\left[\begin{array}{ccc cccc} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} & a_{1,5} & a_{1,6} & a_{1,7} \\ \diamond & a_{2,2} & a_{2,3} & a_{2,4} & a_{2,5} & a_{2,6} & a_{2,7} \\ \diamond & \diamond & a_{3,3} & a_{3,4} & a_{3,5} & a_{3,6} & a_{3,7} \\ \hline \diamond & \diamond & \diamond & a_{4,4} & a_{4,5} & a_{4,6} & a_{4,7} \\ \diamond & \diamond & \diamond & \diamond & a_{5,5} & a_{5,6} & a_{5,7} \\ \diamond & \diamond & \diamond & \diamond & \diamond & a_{6,6} & a_{6,7} \\ \diamond & \diamond & \diamond & \diamond & \diamond & \diamond & a_{7,7} \end{array} \right],$		$\left[\begin{array}{cccc} a_{1,4} & a_{1,5} & a_{1,6} & a_{1,7} \\ a_{2,4} & a_{2,5} & a_{2,6} & a_{2,7} \\ a_{3,4} & a_{3,5} & a_{3,6} & a_{3,7} \\ \hline a_{4,4} & a_{4,5} & a_{4,6} & a_{4,7} \\ a_{1,1} & a_{5,5} & a_{5,6} & a_{5,7} \\ a_{1,2} & a_{2,2} & a_{6,6} & a_{6,7} \\ a_{1,3} & a_{2,3} & a_{3,3} & a_{7,7} \end{array} \right]$

5.4 The matrices A of order N and A_R of size $LDAR$ by $n1$, here $N = 6$.

5.4.1 Full format		5.4.2 RFPF
$\left[\begin{array}{ccc ccc} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} & a_{1,5} & a_{1,6} \\ \diamond & a_{2,2} & a_{2,3} & a_{2,4} & a_{2,5} & a_{2,6} \\ \diamond & \diamond & a_{3,3} & a_{3,4} & a_{3,5} & a_{3,6} \\ \hline \diamond & \diamond & \diamond & a_{4,4} & a_{4,5} & a_{4,6} \\ \diamond & \diamond & \diamond & \diamond & a_{5,5} & a_{5,6} \\ \diamond & \diamond & \diamond & \diamond & \diamond & a_{6,6} \end{array} \right],$		$\left[\begin{array}{ccc} a_{1,4} & a_{1,5} & a_{1,6} \\ a_{2,4} & a_{2,5} & a_{2,6} \\ a_{3,4} & a_{3,5} & a_{3,6} \\ \hline a_{4,4} & a_{4,5} & a_{4,6} \\ a_{1,1} & a_{5,5} & a_{5,6} \\ a_{1,2} & a_{2,2} & a_{6,6} \\ a_{1,3} & a_{2,3} & a_{3,3} \end{array} \right]$

Fig. 6. Eight two-dimensional arrays for storing the matrices A_R and A_R^T in RFPF. The leading dimension $LDAR$ of A_R is N when N is odd and $N + 1$ when N is even. For the matrix A_R^T it is $n1 = \lceil N/2 \rceil$. The memory needed for both A_R and A_R^T is $(N + 1)/2 \times N$. This amount is 28 for $N = 7$ and 21 for $N = 6$.

6.1 RFPF for the matrices of rank odd, here $N = 7$ and $n1 = 4$

Lower triangular

$$\begin{array}{c}
 LDAR = N \\
 \left[\begin{array}{cccc|cccc}
 a_{1,1} & a_{5,5} & a_{6,5} & a_{7,5} & & & & \\
 a_{2,1} & a_{2,2} & a_{6,6} & a_{7,6} & & & & \\
 a_{3,1} & a_{3,2} & a_{3,3} & a_{7,7} & & & & \\
 a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} & & & & \\
 \hline
 a_{5,1} & a_{5,2} & a_{5,3} & a_{5,4} & & & & \\
 a_{6,1} & a_{6,2} & a_{6,3} & a_{6,4} & & & & \\
 a_{7,1} & a_{7,2} & a_{7,3} & a_{7,4} & & & &
 \end{array} \right]
 \end{array}
 \quad
 \begin{array}{c}
 \text{transpose, } lda = n1 \\
 \left[\begin{array}{cccc|cccc}
 a_{1,1} & a_{2,1} & a_{3,1} & a_{4,1} & a_{5,1} & a_{6,1} & a_{7,1} & \\
 a_{5,5} & a_{2,2} & a_{3,2} & a_{4,2} & a_{5,2} & a_{6,2} & a_{7,2} & \\
 a_{6,5} & a_{6,6} & a_{3,3} & a_{4,3} & a_{5,3} & a_{6,3} & a_{7,3} & \\
 a_{7,5} & a_{7,6} & a_{7,7} & a_{4,4} & a_{5,4} & a_{6,4} & a_{7,4} &
 \end{array} \right]
 \end{array}$$

Upper triangular

$$\begin{array}{c}
 LDAR = N \\
 \left[\begin{array}{cccc|cccc}
 a_{1,4} & a_{1,5} & a_{1,6} & a_{1,7} & & & & \\
 a_{2,4} & a_{2,5} & a_{2,6} & a_{2,7} & & & & \\
 a_{3,4} & a_{3,5} & a_{3,6} & a_{3,7} & & & & \\
 \hline
 a_{4,4} & a_{4,5} & a_{4,6} & a_{4,7} & & & & \\
 a_{1,1} & a_{5,5} & a_{5,6} & a_{5,7} & & & & \\
 a_{1,2} & a_{2,2} & a_{6,6} & a_{6,7} & & & & \\
 a_{1,3} & a_{2,3} & a_{3,3} & a_{7,7} & & & &
 \end{array} \right]
 \end{array}
 \quad
 \begin{array}{c}
 \text{transpose, } lda = n1 \\
 \left[\begin{array}{ccc|ccc|ccc}
 a_{1,4} & a_{2,4} & a_{3,4} & a_{4,4} & a_{1,1} & a_{1,2} & a_{1,3} & \\
 a_{1,5} & a_{2,5} & a_{3,5} & a_{4,5} & a_{5,5} & a_{2,2} & a_{2,3} & \\
 a_{1,6} & a_{2,6} & a_{3,6} & a_{4,6} & a_{5,6} & a_{6,6} & a_{3,3} & \\
 a_{1,7} & a_{2,7} & a_{3,7} & a_{4,7} & a_{5,7} & a_{6,7} & a_{7,7} &
 \end{array} \right]
 \end{array}$$

6.2 RFPF for the matrices of rank even, here $N = 6$ and $n1 = 3$.

Lower triangular

$$\begin{array}{c}
 LDAR = N + 1 \\
 \left[\begin{array}{ccc|ccc}
 a_{4,4} & a_{5,4} & a_{6,4} & a_{1,5} & & \\
 a_{1,1} & a_{5,5} & a_{6,5} & a_{1,6} & & \\
 a_{2,1} & a_{2,2} & a_{6,6} & & & \\
 a_{3,1} & a_{3,2} & a_{3,3} & & & \\
 \hline
 a_{4,1} & a_{4,2} & a_{4,3} & & & \\
 a_{5,1} & a_{5,2} & a_{5,3} & & & \\
 a_{6,1} & a_{6,2} & a_{6,3} & & &
 \end{array} \right]
 \end{array}
 \quad
 \begin{array}{c}
 \text{transpose, } lda = n1 \\
 \left[\begin{array}{ccc|ccc}
 a_{4,4} & a_{1,1} & a_{2,1} & a_{3,1} & a_{4,1} & a_{5,1} & a_{6,1} \\
 a_{5,4} & a_{5,5} & a_{2,2} & a_{3,2} & a_{4,2} & a_{5,2} & a_{6,2} \\
 a_{6,4} & a_{6,5} & a_{6,6} & a_{3,3} & a_{4,3} & a_{5,3} & a_{6,3}
 \end{array} \right]
 \end{array}$$

Upper triangular

$$\begin{array}{c}
 LDAR = N + 1 \\
 \left[\begin{array}{ccc|ccc}
 a_{1,4} & a_{1,5} & a_{1,6} & & & \\
 a_{2,4} & a_{2,5} & a_{2,6} & & & \\
 a_{3,4} & a_{3,5} & a_{3,6} & & & \\
 \hline
 a_{4,4} & a_{4,5} & a_{4,6} & & & \\
 a_{1,1} & a_{5,5} & a_{5,6} & a_{1,9} & & \\
 a_{1,2} & a_{2,2} & a_{6,6} & & & \\
 a_{1,3} & a_{2,3} & a_{3,3} & & &
 \end{array} \right]
 \end{array}
 \quad
 \begin{array}{c}
 \text{transpose, } lda = n1 \\
 \left[\begin{array}{ccc|ccc}
 a_{1,4} & a_{2,4} & a_{3,4} & a_{4,4} & a_{1,1} & a_{1,2} & a_{1,3} \\
 a_{1,5} & a_{2,5} & a_{3,5} & a_{4,5} & a_{5,5} & a_{2,2} & a_{2,3} \\
 a_{1,6} & a_{2,6} & a_{3,6} & a_{4,6} & a_{5,6} & a_{6,6} & a_{3,3}
 \end{array} \right]
 \end{array}$$

7. STABILITY OF THE RFPF ALGORITHM

The RFPF Cholesky factorization (Section 3), Cholesky solution (Section 4), and Cholesky inversion (Section 5) algorithms are equivalent to the traditional algorithms in the books [Dongarra et al. 1998; Demmel 1997; Golub and Van Loan 1996; Trefethen and Bau 1997]. The whole theory of the traditional Cholesky factorization, solution, inversion and BLAS algorithms carries over to this three Cholesky and BLAS algorithms described in Sections 3, 4, and 5. The error analysis and stability of these algorithms is very well described in the book of [Higham 1996]. The difference between LAPACK algorithms PO, PP and RFPF³ is how inner products are accumulated. In each case a different order is used. They are all mathematically equivalent, and, stability analysis shows that any summation order is stable.

8. A PERFORMANCE STUDY USING RFP FORMAT

The LAPACK library [Anderson et al. 1999] routines POTRF/PPTRF, POTRI/PPTRI, and POTRS/PPTRS are compared with the RFPF routines PFTRF, PFTRI, and PFTRS for Cholesky factorization (PxTRF), Cholesky inverse (PxTRI) and Cholesky solution (PxTRS) respectively. In the previous sentence, the character 'x' can be 'O' (full format), 'P' (packed format), or 'F' (RFPF). In all cases long real precision arithmetic (also called double precision) is used. Sometimes we also show results for long complex precision (also called complex*16). Results were obtained on several different computers using everywhere the vendor Level 3 and Level 2 BLAS. The sequential performance results were done on the following computers:

- **Sun Fire E25K (newton):** 72 UltraSPARC IV+ dual-core CPUs (1800 MHz/ 2 MB shared L2-cache, 32 MB shared L3-cache), 416 GB memory (120 CPUs/368 GB). Further information at "<http://www.gbar.dtu.dk/index.php/Hardware>".
- **SGI Altix 3700 (Freke):** 64 CPUs - Intel Itanium2 1.5 GHz/6 MB L3-cache. 256 GB memory. Peak performance: 384 GFlops. Further information at "<http://www.csc.dk/freke/>".
- **Intel Tigerton computer (zoot):** quad-socket quad-core Intel Tigerton 2.4GHz (16 total cores) with 32 GB of memory. We use Intel MKL 10.0.1.014.
- **DMI Itanium:** CPU Intel Itanium2: 1.3 GHz, cache: 3 MB on-chip L3 cache.
- **DMI NEC SX-6 computer:** 8 CPU's, per CPU peak: 8 Gflops, per node peak: 64 Gflops, vector register length: 256.

The performance results are given in Figures 7 to 15.

The figures from 7 to 10 are paired. Figure 7 (double precision) and Figure 8 (double complex precision) present results for the Sun UltraSPARC IV+ computer. Figure 9 (double precision) and Figure 10 (double complex precision) present results for the SGI Altix 3700 computer. Figure 11 (double precision) presents results for the Intel Itanium2 computer. Figure 12 (double precision) presents results for the NEC SX-6 computer. Figure 13 (double precision) presents results for the quad-socket quad-core Intel Tigerton computer using reference LAPACK-3.2.0

³full, packed and rectangular full packed.

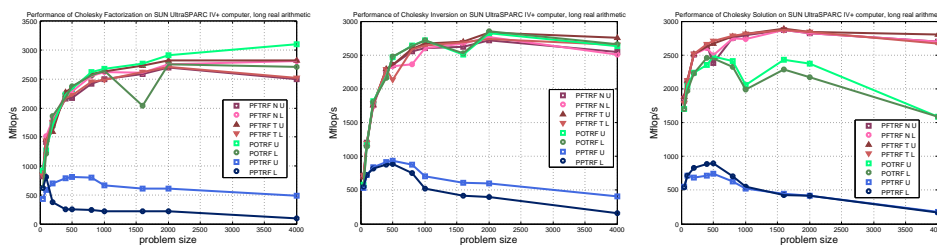


Fig. 7. Performance in Mflop/s of Cholesky Factorization/Inversion/Solution on SUN UltraSPARC IV+ computer, long real arithmetic. For PxTRF, $nrhs = \max(100, n/10)$.

(from netlib.org). Figure 14 (double precision) presents results for the quad-socket quad-core Intel Tigerton computer using vendor LAPACK library (MKL-10.0.1.14).

Figure 15 shows the SMP parallelism of these subroutines on the IBM Power4 (clock rate: 1300 MHz; two CPUs per chip; L1 cache: 128 KB (64 KB per CPU) instruction, 64 KB 2-way (32 KB per CPU) data; L2 cache: 1.5 MB 8-way shared between the two CPUs; L3 cache: 32 MB 8-way shared (off-chip); TLB: 1024 entries) and SUN UltraSPARC-IV (clock rate: 1350 MHz; L1 cache: 64 kB 4-way data, 32 kB 4-way instruction, and 2 kB Write, 2 kB Prefetch; L2 cache: 8 MB; TLB: 1040 entries) computers respectively. They compare SMP times of PFTRF, vendor POTRF and reference PPTRF.

The RFPF packed results greatly outperform the packed and more often than not are better than the full results. Note that our timings do *not* include the cost of sorting any LAPACK data formats to RFPF data formats and vice versa. We think that users will input their matrix data using RFPF. Hence, this is our rationale for not including the data transformation times.

For all our experiments, we use vendor Level 3 and Level 2 BLAS. For all our experiments except Figure 13 and Figure 15, we use the provided vendor library for LAPACK and BLAS.

We include comparisons with reference LAPACK for the quad-socket quad-core Intel Tigerton machine in Figure 13. In this case, the vendor LAPACK library packed storage routines significantly outperform the LAPACK reference implementation from netlib. In Figure 14, you find the same experiments on the same machine but, this time, using the vendor library (MKL-10.0.1.014). We think that MKL is using the reference implementation for Inverse Cholesky (packed and full format). For Cholesky factorization, we see that both packed and full format routines (PPTRF and POTRF) are tuned. But even, in this case, our RFPF storage format results are better.

When we compare RFPF with full storage, results are mixed. However, both codes are rarely far apart. Most of the performance ratios are between 0.95 to 1.05 overall. But, note that the RFPF performance is more uniform over its versions (four presented; the other four are for n odd). For LAPACK full (two versions), the performance variation is greater. Moreover, in the case of the inversion on quad-socket quad-core Tigerton (Figure 13 and Figure 14) RFPF clearly outperforms both variants of the full format.

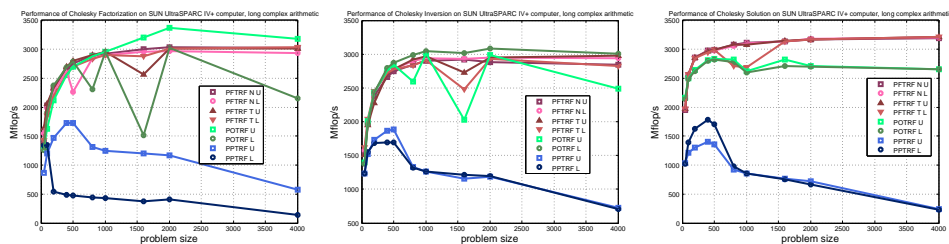


Fig. 8. Performance in Mflop/s of Cholesky Factorization/Inversion/Solution on SUN UltraSPARC IV+ computer, long complex arithmetic. For PxTRF, $nrhs = \max(100, n/10)$.

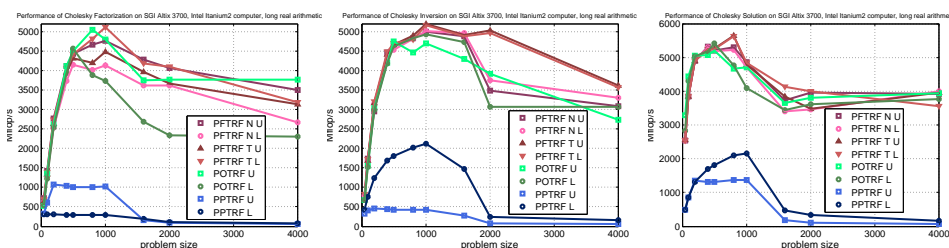


Fig. 9. Performance in Mflop/s of Cholesky Factorization/Inversion/Solution on SGI Altix 3700, Intel Itanium 2 computer, long real arithmetic. For PxTRF, $nrhs = \max(100, n/10)$.

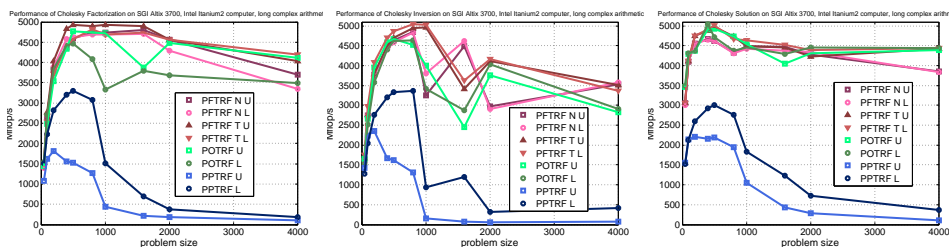


Fig. 10. Performance in Mflop/s of Cholesky Factorization/Inversion/Solution on SGI Altix 3700, Intel Itanium 2 computer, long complex arithmetic. For PxTRF, $nrhs = \max(100, n/10)$.

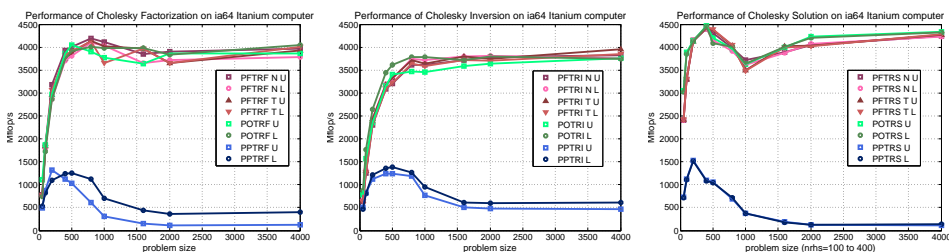


Fig. 11. Performance in Mflop/s of Cholesky Factorization/Inversion/Solution on ia64 Itanium computer, long real arithmetic. For PxTRF, $nrhs = \max(100, n/10)$.

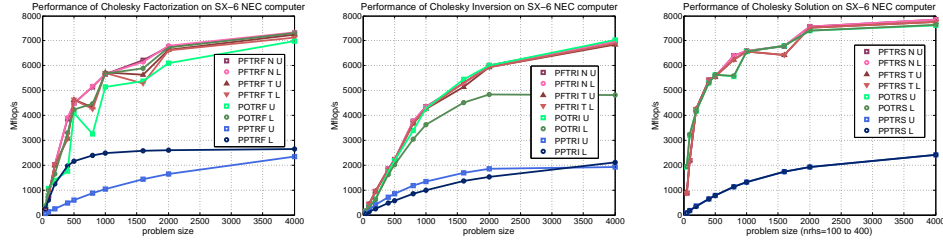


Fig. 12. Performance in Mflop/s of Cholesky Factorization/Inversion/Solution on SX-6 NEC computer, long real arithmetic. For P_xTRF, $n_{rhs} = \max(100, n/10)$.

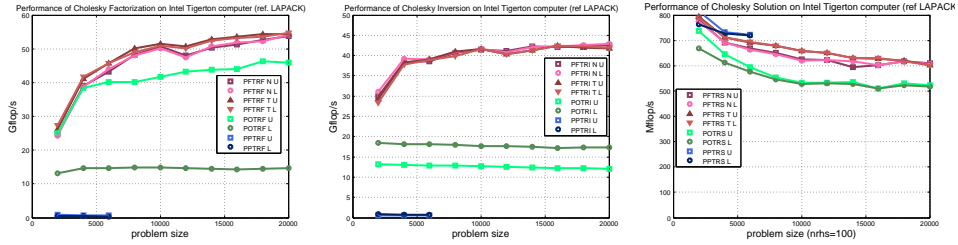


Fig. 13. Performance of Cholesky Factorization/Inversion/Solution on quad-socket quad-core Intel Tigerton computer, long real arithmetic. We use reference LAPACK-3.2.0 (from netlib) and MKL-10.0.1.014 multithreaded BLAS. For the solution phase, n_{rhs} is fixed to 100 for any n . Due to time limitation, the experiment was stopped for the packed storage format inversion at $n = 4000$.

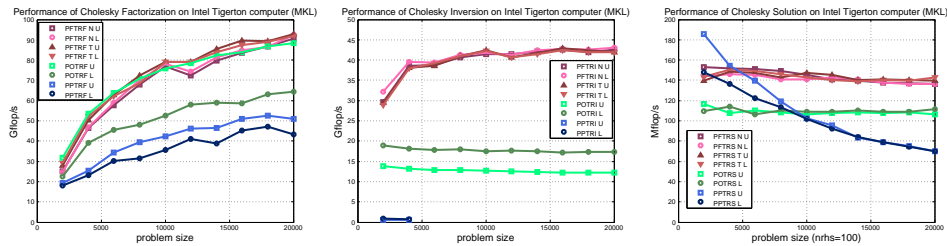


Fig. 14. Performance of Cholesky Factorization/Inversion/Solution on quad-socket quad-core Intel Tigerton computer, long real arithmetic. We use MKL-10.0.1.014 multithreaded LAPACK and BLAS. For the solution phase, n_{rhs} is fixed to 100 for any n . Due to time limitation, the experiment was stopped for the packed storage format inversion at $n = 4000$.

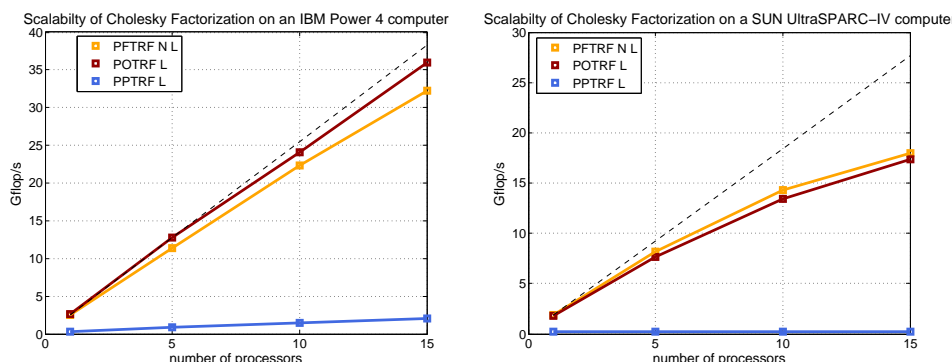


Fig. 15. Performance in Gflop/s of Cholesky Factorization on IBM Power 4 (left) and SUN UltraSPARC-IV (right) computer, long real arithmetic, with a different number of Processors, testing the SMP Parallelism. The implementation of PPTRF of sunperf does not show any SMP parallelism. UPLO = 'L'. $N = 5,000$ (strong scaling experiment).

9. INTEGRATION IN LAPACK

As mentioned in the introduction, as of release 3.2 (November 2008), LAPACK supports a preliminary version of RFPF. Ultimately, the goal would be for RFPF to support as many functionalities as full format or standard packed format does. The 44 routines included in release 3.2 for RFPF are given in Table 1. The names for the RFPF routines follow the naming nomenclature used by LAPACK. We have added the format description letters: PF for Symmetric/Hermitian Positive Definite RFPF (PO for full, PP for packed), SF for Symmetric RFPF (SY for full, SP for packed), HF for Hermitian RFPF (HE for full, HP for packed), and TF for Triangular RFPF (TR for full, TP for packed).

Currently, for the complex case, we assume that the transpose complex-conjugate part is stored whenever the transpose part is stored in the real case. This corresponds to the theory developed in this present manuscript. In the future, we will want to have the flexibility to store the transpose part (as opposed to transpose complex conjugate) whenever the transpose part is stored in the real case. In particular, this feature will be useful for complex symmetric matrices.

10. SUMMARY AND CONCLUSIONS

This paper describes RFPF as a standard minimal full format for representing both symmetric and triangular matrices. Hence, from a user point of view, these matrix layouts are a replacement for both the standard formats of DLA, namely full and packed storage. These new layouts possess three good features: they are efficient, they are supported by Level 3 BLAS and LAPACK full format routines, and they require minimal storage.

11. ACKNOWLEDGMENTS

The results in this paper were obtained on seven computers, an IBM, a SGI, two SUNs, Itanium, NEC, and Intel Tigerton computers. The IBM machine belongs to the Center for Scientific Computing at Aarhus, the SUN machines to the Danish

functionality	routine names and calling sequence			
Cholesky factorization	CPFTRF	DPFTRF	SPFTRF	ZPFTRF (TRANSR,UPLO,N,A,INFO)
Multiple solve after PFTRF	CPFTRS	DPFTRS	SPFTRS	ZPFTRS (TRANSR,UPLO,N,NR,A,B,LDB,INFO)
Inversion after PFTRF	CPFTRI	DPFTRI	SPFTRI	ZPFTRI (TRANSR,UPLO,N,A,INFO)
Triangular inversion	CTRTRI	DTRTRI	STRTRI	ZTRTRI (TRANSR,UPLO,DIAG,N,A,INFO)
Sym/Herm matrix norm	CLANHF	DLANSF	SLANSF	ZLANHF (NORM,TRANSR,UPLO,N,A,WORK)
Triangular solve	CTFSM	DTFSM	STFSM	ZTFSM (TRANSR,SIDE,UPLO,TRANS,DIAG,M,N,ALPHA,A,B,LDB)
Sym/Herm rank- k update	CHFRK	DSFRK	SSFRK	ZHFRK (TRANSR,UPLO,TRANS,N,K,ALPHA,A,LDA,BETA,C)
Conv. from TP to TF	CTPTTF	DTPPTF	STPTTF	ZTPTTF (TRANSR,UPLO,N,AP,ARF,INFO)
Conv. from TR to TF	CTRPTF	DTRPTF	STRPTF	ZTRPTF (TRANSR,UPLO,N,A,LDA,ARF,INFO)
Conv. from TF to TP	CTFTTP	DTFTTP	STFTTP	ZFTTTP (TRANSR,UPLO,N,ARF,AP,INFO)
Conv. from TF to TR	CTFTTR	DTFTTR	STFTTR	ZFTTTR (TRANSR,UPLO,N,ARF,A,LDA,INFO)

Table 1. LAPACK 3.2 RFPF routines.

Technical University, the Itanium and NEC machines to the Danish Meteorological Institute, and the Intel Tigerton machine to the Innovative Computing Laboratory at the University of Tennessee.

We would like to thank Bernd Dammann for consulting on the SUN systems; Niels Carl W. Hansen for consulting on the IBM and SGI systems; and Bjarne Stig Andersen for obtaining the results on the Itanium and NEC computers. We thank IBMers John Gunnels who worked earlier on the HFPP format and JP Fasano who was instrumental in getting the source code released by the IBM Open Source Committee. We thank Allan Backer for discussions about an older version of this manuscript.

REFERENCES

- AGARWAL, R. C., GUSTAVSON, F. G., AND ZUBAIR, M. 1994. Exploiting functional parallelism on power2 to design high-performance numerical algorithms. *IBM Journal of Research and Development* 38, 5 (September), 563–576.
- ANDERSEN, B. S., GUNNELS, J. A., GUSTAVSON, F., AND WAŚNIEWSKI, J. 2002. A Recursive Formulation of the Inversion of symmetric positive definite Matrices in Packed Storage Data Format. In J. FAGERHOLM, J. HAATAJA, J. JÄRVINEN, M. LYLÄ, AND P. R. V. SAVOLAINEN Eds., *Proceedings of the 6th International Conference, PARA 2002, Applied Parallel Computing*, Number 2367 in Lecture Notes in Computer Science (Espoo, Finland, June 2002), pp. 287–296. Springer.
- ANDERSEN, B. S., GUSTAVSON, F. G., REID, J. K., AND WAŚNIEWSKI, J. 2005. A Fully Portable High Performance Minimal Storage Hybrid Format Cholesky Algorithm. *ACM Transactions on Mathematical Software* 31, 201–227.
- ANDERSEN, B. S., GUSTAVSON, F. G., AND WAŚNIEWSKI, J. 2001. A Recursive Formulation of Cholesky Factorization of a Matrix in Packed Storage. *ACM Transactions on Mathematical Software* 27, 2 (Jun), 214–244.

- ANDERSON, E., BAI, Z., BISCHOF, C., BLACKFORD, L. S., DEMMEL, J., DONGARRA, J. J., DU CROZ, J., GREENBAUM, A., HAMMARLING, S., MCKENNEY, A., AND SORENSEN, D. 1999. *LAPACK Users' Guide* (Third ed.). Society for Industrial and Applied Mathematics, Philadelphia, PA.
- BARKER, V. A., BLACKFORD, L. S., DONGARRA, J. J., CROZ, J. D., HAMMARLING, S., MARINOVA, M., WAŚNIEWSKI, J., AND YALAMOV, P. 2001. *LAPACK95 Users' Guide* (first ed.). Society for Industrial and Applied Mathematics, Philadelphia, PA.
- BUTTARI, A., LANGOU, J., KURZAK, J., AND DONGARRA, J. 2007. A class of parallel tiled linear algebra algorithms for multi-core architectures. Tech rep. ut-cs-07-0600, Department of Electrical Engineering and Computer Science of the University of Tennessee.
- CHAN, E., QUINTANA-ORTÍ, E., QUINTANA-ORTÍ, G., AND VAN DE GEIJN, R. 2007. Supermatrix out-of-order scheduling of matrix operations for smp and multi-core architectures. In *SPAA 07, Proceedings of the 19th ACM Symposium on Parallelism in Algorithms and Architecture* (2007), pp. 116–125.
- DEMMEL, J. W. 1997. *Applied Numerical Linear Algebra*. SIAM, Philadelphia.
- DONGARRA, J. J., BUNCH, J. R., MOLER, C. B., AND STEWART, G. W. 1979. *Lapack Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA.
- DONGARRA, J. J., DU CROZ, J., DUFF, I. S., AND HAMMARLING, S. 1990a. Algorithm 679: A set of Level 3 Basic Linear Algebra Subprograms. *ACM Trans. Math. Soft.* 16, 1 (March), 18–28.
- DONGARRA, J. J., DU CROZ, J., DUFF, I. S., AND HAMMARLING, S. 1990b. A Set of Level 3 Basic Linear Algebra Subprograms. *ACM Trans. Math. Soft.* 16, 1 (March), 1–17.
- DONGARRA, J. J., DU CROZ, J., HAMMARLING, S., AND HANSON, R. J. 1988. An Extended Set of Fortran Basic Linear Algebra Subroutines. *ACM Trans. Math. Soft.* 14, 1 (March), 1–17.
- DONGARRA, J. J., DUFF, I. S., SORENSEN, D. C., AND VAN DER VORST, H. A. 1998. *Numerical Linear Algebra for High Performance Computers*. SIAM, Society for Industrial and Applied Mathematics, Philadelphia.
- ELMROTH, E., GUSTAVSON, F. G., KAGSTROM, B., AND JONSSON, I. 2004. Recursive Blocked Algorithms and Hybrid Data Structures for Dense Matrix Library Software. *SIAM Review* 46, 1 (March), 3–45.
- GOLUB, G. AND VAN LOAN, C. F. 1996. *Matrix Computations* (Third ed.). Johns Hopkins University Press, Baltimore, MD.
- GUNNELS, J. A. AND GUSTAVSON, F. G. 2004. A new array format for symmetric and triangular matrices. In J. W. J.J. DONGARRA, K. MADSEN Ed., *Applied Parallel Computing, State of the Art in Scientific Computing, PARA 2004*, Volume LNCS 3732 (Springer-Verlag, Berlin Heidelberg, 2004), pp. 247–255. Springer.
- GUSTAVSON, F. G. 2003. High Performance Linear Algebra Algorithms using New Generalized Data Structures for Matrices. *IBM Journal of Research and Development* 47, 1 (January), 823–849.
- GUSTAVSON, F. G., GUNNELS, J., AND SEXTON, J. 2007. Minimal Data Copy for Dense Linear Algebra Factorization. In *Applied Parallel Computing, State of the Art in Scientific Computing, PARA 2006*, Volume LNCS 4699 (Springer-Verlag, Berlin Heidelberg, 2007), pp. 540–549. Springer.
- GUSTAVSON, F. G. AND JONSSON, I. 2000. Minimal storage high performance cholesky via blocking and recursion. *IBM Journal of Research and Development* 44, 6 (Nov), 823–849.
- GUSTAVSON, F. G., REID, J. K., AND WAŚNIEWSKI, J. 2007. Algorithm 865: Fortran 95 Subroutines for Cholesky Factorization in Blocked Hybrid Format. *ACM Transactions on Mathematical Software* 33, 1 (March), 5.
- GUSTAVSON, F. G. AND WAŚNIEWSKI, J. 2007. Rectangular full packed format for LAPACK algorithms timings on several computers. In *Applied Parallel Computing, State of the Art in Scientific Computing, PARA 2006*, Volume LNCS 4699 (Springer-Verlag, Berlin Heidelberg, 2007), pp. 570–579. Springer.

- HERRERO, J. R. 2006. *A Framework for Efficient Execution of Matrix Computations*. Ph. D. thesis, Universitat Politècnica de Catalunya.
- HIGHAM, N. J. 1996. *Accuracy and Stability of Numerical Algorithms*. SIAM.
- IBM. 1997. *Engineering and Scientific Subroutine Library for AIX* (Version 3, Volume 1 ed.). IBM. Pub. number SA22-7272-0.
- LAWSON, C. L., HANSON, R. J., KINCAID, D., AND KROGH, F. T. 1979. Basic Linear Algebra Subprograms for Fortran Usage. *ACM Trans. Math. Soft.* 5, 308–323.
- TREFETHEN, L. N. AND BAU, D. 1997. *Numerical Linear Algebra*. SIAM, Philadelphia.