

A Set of Level 3 Basic Linear Algebra Subprograms

JACK J. DONGARRA

University of Tennessee and Oak Ridge National Laboratory

JEREMY DU CROZ and SVEN HAMMARLING

Numerical Algorithms Group, Ltd.

and

IAIN DUFF

Harwell Laboratory

This paper describes a set of Level 3 Basic Linear Algebra Subprograms (Level 3 BLAS). The Level 3 BLAS are targeted at matrix-matrix operations, with the aim of providing more efficient, but portable, implementations of algorithms on high-performance computers, especially those with hierarchical memory and parallel processing capability.

Categories and Subject Descriptors: F.2.1 [Analysis of Algorithms and Problem Complexity]: Numerical Algorithms and Problems—*computations on matrices*; G.1.0 [Numerical Analysis]: General—*numerical algorithms*; G.1.3 [Numerical Analysis]: Numerical Linear Algebra—*linear systems (direct and iterative methods)*; G.4 [Mathematics of Computing]: Mathematical Software—*certification and testing; efficiency; portability; reliability and robustness; verification*

General Terms: Algorithms, Measurement, Performance, Reliability, Verification

Additional Key Words and Phrases: Extended BLAS, utilities

1. INTRODUCTION

In 1973, Hanson, Krogh, and Lawson wrote an article in the *SIGNUM Newsletter* (Vol. 8, no. 4, p. 16) describing the advantages of adopting a set of basic routines for problems in linear algebra. The original basic linear algebra subprograms, now commonly referred to as the BLAS and fully described in [27, 28], have been very successful, and have been used in a wide range of software including LINPACK [13] and many of the algorithms published by the *ACM Transactions on Mathematical Software*. In particular, they are an aid to clarity, portability,

This work was supported in part by the Applied Mathematical Sciences subprogram of the Office of Energy Research, U.S. Department of Energy, under contract W-31-109-Eng-38.

Authors' addresses: J. J. Dongarra, Department of Computer Science, University of Tennessee, Knoxville, TN 37996-1301; J. Du Croz and S. Hammarling, Numerical Algorithms Group Ltd., Wilkinson House, Jordan Hill Road, Oxford OX2 8DR, England; I. Duff, Computer Science and Systems Division, Harwell Laboratory, Oxfordshire OX11 0RA, England.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1990 ACM 0098-3500/90/0300-0001 \$01.50

modularity, and maintenance of software; and they have become a de facto standard for the elementary vector operations. An excellent discussion of the *raison d'être* of the BLAS is given in [13].

An extended set of Fortran BLAS, aimed at matrix-vector operations (Level 2 BLAS), was subsequently proposed by Dongarra, Du Croz, Hammarling, and Hanson [14, 15]. The Level 2 BLAS were proposed in order to support the development of software that would be both portable and efficient across a wide range of machine architectures, with emphasis on vector-processing machines.

Many of the frequently used algorithms of numerical linear algebra can be coded so that the bulk of the computation is performed by calls to Level 2 BLAS routines; efficiency can then be obtained by utilizing tailored implementations of the Level 2 BLAS routines. On vector-processing machines, one of the aims of such implementations is to keep the vector lengths as long as possible, and in most algorithms the results are computed one vector (row or column) at a time. In addition, on vector register machines, performance is increased by reusing the results of a vector register, and not storing the vector back into memory.

Unfortunately, this approach to software construction is often not well suited to computers with a hierarchy of memory (such as global memory, cache or local memory, and vector registers) and true parallel-processing computers. (For a description of many advanced computer architectures, see [17].) For those architectures, it is often preferable to partition the matrix or matrices into blocks and to perform the computation by matrix-matrix operations on the blocks. By organizing the computation in this fashion we provide for full reuse of data while the block is held in the cache or local memory. This approach avoids excessive movement of data to and from memory and gives a *surface-to-volume* effect for the ratio of operations to data movement. In addition, on architectures that provide for parallel processing, parallelism can be exploited in two ways: (1) operations on distinct blocks may be performed in parallel; and (2) within the operations on each block, scalar or vector operations may be performed in parallel.

The Level 3 BLAS specified here are targeted at the matrix-matrix operations required for these purposes. If the vectors and matrices involved are of order n , then the original BLAS include operations that are of order $O(n)$, the extended or Level 2 BLAS provide operations of order $O(n^2)$, and the routines specified here provide operations of order $O(n^3)$ —hence our use of the term Level 3 BLAS. Section 9 illustrates how an algorithm for Cholesky factorization can be implemented by calls to the Level 3 BLAS. Such implementations can, we believe, be portable across a wide variety of vector and parallel computers, and are also efficient (assuming that efficient implementations of the Level 3 BLAS are available). There is certainly considerable evidence for the efficiency of such algorithms on particular machines (see, for example, the references cited in Section 9); the question of portability has been much less studied but, by proposing a standard set of building blocks, we hope, to encourage research into this aspect.

The scope of the Level 3 BLAS is intentionally limited. For example, no routines are included for matrix factorization; these are currently provided by LINPACK [13], and will be included in a new linear algebra package currently under development, which will use block algorithms and calls to Level 3 BLAS

wherever possible [10]. Nor are the Level 3 BLAS intended to be a comprehensive set of routines for elementary matrix algebra. They are intended primarily for software developers, and to a lesser extent for experienced applications programmers.

The details of this paper are concerned specifically with defining a set of subroutines for use in Fortran 77 programs. However, the essential features could be adapted to other programming languages.

In a companion paper [19], we present a model implementation of the Level 3 BLAS in Fortran 77 (extended to include a COMPLEX*16 data type), and also a set of rigorous test programs.

2. SCOPE OF THE LEVEL 3 BLAS

The routines described here have been derived in a fairly obvious manner from some of the Level 2 BLAS, by replacing the vectors x and y with matrices B and C . The advantage in keeping the design of the software as consistent as possible with that of the Level 2 BLAS is that it will be easier for users to remember the calling sequences and parameter conventions.

In real arithmetic, the operations proposed for the Level 3 BLAS have the following forms:

(a) Matrix-matrix products

$$\begin{aligned} C &\leftarrow \alpha AB + \beta C \\ C &\leftarrow \alpha A^T B + \beta C \\ C &\leftarrow \alpha AB^T + \beta C \\ C &\leftarrow \alpha A^T B^T + \beta C \end{aligned}$$

Note that these operations are more accurately described as matrix-matrix multiply-and-add operations; they include rank- k updates of a general matrix.

(b) Rank- k and rank- $2k$ updates of a symmetric matrix:

$$\begin{aligned} C &\leftarrow \alpha AA^T + \beta C \\ C &\leftarrow \alpha A^T A + \beta C \\ C &\leftarrow \alpha AB^T + \alpha BA^T + \beta C \\ C &\leftarrow \alpha A^T B + \alpha B^T A + \beta C \end{aligned}$$

(c) Multiplying a matrix by a triangular matrix:

$$\begin{aligned} B &\leftarrow \alpha TB \\ B &\leftarrow \alpha T^T B \\ B &\leftarrow \alpha BT \\ B &\leftarrow \alpha BT^T \end{aligned}$$

(d) Solving triangular systems of equations with multiple right-hand sides:

$$\begin{aligned} B &\leftarrow \alpha T^{-1} B \\ B &\leftarrow \alpha T^{-T} B \\ B &\leftarrow \alpha BT^{-1} \\ B &\leftarrow \alpha BT^{-T} \end{aligned}$$

Here, α and β are scalars, A , B , and C are rectangular matrices (in some cases, square and symmetric), and T is an upper or lower triangular matrix (and nonsingular in (d)).

Analogous operations are proposed in complex arithmetic: conjugate transposition is specified as well as simple transposition, and additional operations in (b) provide for updates of a Hermitian matrix, as follows:

$$\begin{aligned} C &\leftarrow \alpha AA^H + \beta C \\ C &\leftarrow \alpha A^H A + \beta C \end{aligned}$$

with α and β real, and

$$\begin{aligned} C &\leftarrow \alpha AB^H + \bar{\alpha} BA^H + \beta C \\ C &\leftarrow \alpha A^H B + \bar{\alpha} B^H A + \beta C \end{aligned}$$

with β real.

3. NAMING CONVENTIONS

The name of a Level 3 BLAS routine follows the conventions of the Level 2 BLAS. The first character in the name denotes the Fortran data type of the matrix, as follows:

S	REAL
D	DOUBLE PRECISION
C	COMPLEX
Z	COMPLEX*16 or DOUBLE COMPLEX (if available)

Characters two and three in the name refer to the kind of matrix involved, as follows:

GE	All matrices are general rectangular
HE	One of the matrices is Hermitian
SY	One of the matrices is symmetric
TR	One of the matrices is triangular

The fourth and fifth, and in one case sixth, characters in the name denote the type of operation, as follows:

MM	Matrix-matrix product
RK	Rank- k update of a symmetric or Hermitian matrix
R2K	Rank- $2k$ update of a symmetric or Hermitian matrix
SM	Solve a system of linear equations for a matrix of right-hand sides

The permitted combinations are indicated in Table I. In the first column, under *complex*, the initial C may be replaced by Z. In the second column, under *real*, the initial S may be replaced by D; see Appendix B for the full subroutine calling sequences.

The collection of routines can be thought of as being divided into four separate parts: *real*, *double precision*, *complex*, and *complex*16*. The routines can be written in ANSI standard Fortran 77, with the exception of the routines that use COMPLEX*16 variables. These routines are included for completeness and for their usefulness on those systems that support this data type; but because they

Table I

Complex	Real	MM	RK	R2K	SM
CGE	SGE	*			
CSY	SSY	*	*	*	
CHE		*	*	*	
CTR	STR	*			*

do not conform to the Fortran standard, they may not be available on all machines. However, note that rank- k updates of general matrices are provided by the GEMM routines.

4. ARGUMENT CONVENTIONS

We follow a convention for the argument lists similar to that for the Level 2 BLAS, with the necessary adaptations. The order of arguments is as follows:

- (a) Arguments specifying options
- (b) Arguments defining the sizes of the matrices
- (c) Input scalar
- (d) Description of input matrices
- (e) Input scalar (associated with input-output matrix)
- (f) Description of the input-output matrix

Note that not every category is present in each of the routines.

The arguments that specify options are character arguments with the names SIDE, TRANSA, TRANSB, TRANS, UPLO, and DIAG.

SIDE is used by the routines as follows:

Value	Meaning
'L'	Multiply general matrix by symmetric, Hermitian or triangular matrix on the left
'R'	Multiply general matrix by symmetric, Hermitian or triangular matrix on the right

TRANSA, TRANSB, and TRANS are used by the routines as follows:

Value	Meaning
'N'	Operate with the matrix
'T'	Operate with the transpose of the matrix
'C'	Operate with the conjugate transpose of the matrix

In the real case, the values 'T' and 'C' have the same meaning.

UPLO is used by the Hermitian, symmetric, and triangular matrix routines to specify whether the upper or lower triangle is being referenced, as follows:

Value	Meaning
'U'	Upper triangle
'L'	Lower triangle

DIAG is used by the triangular matrix routines to specify whether the matrix is unit triangular, as follows:

Value	Meaning
'U'	Unit triangular
'N'	Nonunit triangular

When DIAG is supplied as 'U', the diagonal elements are not referenced.

Thus, UPLO and DIAG have the same values and meanings as for the Level 2 BLAS; TRANSA, TRANSB, and TRANS have the same values and meanings as TRANS in the Level 2 BLAS, where TRANSA and TRANSB apply to the matrices A or B, respectively.

We recommend that the equivalent lower-case characters be accepted with the same meaning; although, because they are not included in the standard Fortran character set, their use may not be supported on all systems (see Section 7 of [13] for further discussion).

The sizes of the matrices are determined by the arguments M, N, and K. It is permissible to call the routines with M or N = 0, in which case the routines exit immediately without referencing their matrix arguments. If M and N > 0, but K = 0, the operation reduces to $C \leftarrow \beta C$ (this applies to the GEMM, SYRK, HERK, SYR2K, and HER2K routines). The input-output matrix (*B* for the TR routines, *C* otherwise) is always $m \times n$ if rectangular, and $n \times n$ if square.

The description of the matrix consists of the array name (A, B, or C) followed by the leading dimension of the array as declared in the calling (sub)program (LDA, LDB, or LDC).

The scalars always have the dummy argument names ALPHA and BETA.

The following values of arguments are invalid: Any value of the character arguments SIDE, TRANSA, TRANSB, TRANS, UPLO, or DIAG whose meaning is not specified.

- M < 0
- N < 0
- K < 0
- LDA < the number of rows in the matrix A.
- LDB < the number of rows in the matrix B.
- LDC < the number of rows in the matrix C.

If a routine is called with an invalid value for any of its arguments, then it must report the fact and terminate execution of the routine. In the model implementation (see Appendix B), each routine, on detecting an error, calls a common error handling routine XERBLA, passing to it the name of the routine and the number of the first argument that is in error. Specialized implementations may call system-specific exception-handling and diagnostic facilities.

5. STORAGE CONVENTIONS

All matrices are stored conventionally in a two-dimensional array with matrix-element $a_{i,j}$ stored in array-element A(I, J); there is no provision for packed storage for symmetric, Hermitian, or triangular matrices.

For symmetric and Hermitian matrices, only the upper triangle (UPLO = 'U') or the lower triangle (UPLO = 'L') is stored.

For triangular matrices, the argument UPLO serves to define whether the matrix is upper (UPLO = 'U') or lower (UPLO = 'L') triangular.

For a Hermitian matrix, the imaginary parts of the diagonal elements are of course zero, and thus the imaginary parts of the corresponding Fortran array elements need not be set, but are assumed to be zero. In the HERK and HER2K routines, these imaginary parts are set to zero on return, except when $\beta = 1$ and α or $k = 0$, in which case the routines exit immediately.

6. SPECIFICATION OF THE LEVEL 3 BLAS

Type and dimension for variables occurring in the subroutine specifications are as follows:

```
INTEGER M, N, K, LDA, LDB, LDC
CHARACTER*1 SIDE, UPLO, TRANSA, TRANSB, TRANS, DIAG
```

For routines whose first letter is an S:

```
REAL ALPHA, BETA
REAL A(LDA, *), B(LDB, *), C(LDC, *)
```

For routines whose first letter is a D:

```
DOUBLE PRECISION ALPHA, BETA
DOUBLE PRECISION A(LDA, *), B(LDB, *), C(LDC, *)
```

For routines whose first letter is a C:

```
COMPLEX ALPHA, BETA
COMPLEX A(LDA, *), B(LDB, *), C(LDC, *)
```

except that, for CHERK, the scalars α and β are real, so the first declaration above is replaced by

```
REAL ALPHA, BETA
```

and, for CHER2K, α is complex and β is real, so this is replaced by

```
COMPLEX ALPHA
REAL BETA
```

For routines whose first letter is a Z:

```
COMPLEX*16 ALPHA, BETA or DOUBLE COMPLEX ALPHA, BETA
COMPLEX*16 A(LDA, *)      DOUBLE COMPLEX A(LDA, *)
COMPLEX*16 B(LDB, *)      DOUBLE COMPLEX B(LDB, *)
COMPLEX*16 C(LDC, *)      DOUBLE COMPLEX C(LDC, *)
```

except that, for ZHERK, the scalars α and β are real, so the first declaration above is replaced by

```
DOUBLE PRECISION ALPHA, BETA
```

and, for ZHER2K, α is complex and β is real, so this is replaced by

COMPLEX*16 ALPHA
DOUBLE PRECISION BETA

(a) General matrix-matrix products:

_GEMM (TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)

Operation (C is always $m \times n$):

	TRANSA = 'N'	TRANSA = 'T'	TRANSA = 'C'
TRANSB = 'N'	$C \leftarrow \alpha AB + \beta C$ <i>A</i> is $m \times k$, <i>B</i> is $k \times n$	$C \leftarrow \alpha A^T B + \beta C$ <i>A</i> is $k \times m$, <i>B</i> is $k \times n$	$C \leftarrow \alpha A^H B + \beta C$ <i>A</i> is $k \times m$, <i>B</i> is $k \times n$
TRANSB = 'T'	$C \leftarrow \alpha AB^T + \beta C$ <i>A</i> is $m \times k$, <i>B</i> is $n \times k$	$C \leftarrow \alpha A^T B^T + \beta C$ <i>A</i> is $k \times m$, <i>B</i> is $n \times k$	$C \leftarrow \alpha A^H B^T + \beta C$ <i>A</i> is $k \times m$, <i>B</i> is $n \times k$
TRANSB = 'C'	$C \leftarrow \alpha AB^H + \beta C$ <i>A</i> is $m \times k$, <i>B</i> is $n \times k$	$C \leftarrow \alpha A^T B^H + \beta C$ <i>A</i> is $k \times m$, <i>B</i> is $n \times k$	$C \leftarrow \alpha A^H B^H + \beta C$ <i>A</i> is $k \times m$, <i>B</i> is $n \times k$

(In the real case, the values 'T' and 'C' have the same meaning.)

(b) Matrix-matrix products where one matrix is real or complex symmetric or complex Hermitian:

_SYMM (SIDE, UPLO, M, N, ALPHA, A, LDA, B, LDB, BETA, C, LDC)
_HEMM (SIDE, UPLO, M, N, ALPHA, A, LDA, B, LDB, BETA, C, LDC)

Operation (C is always $m \times n$, A is symmetric for the _SYMM routines and Hermitian for the _HEMM routines):

SIDE = 'L'	SIDE = 'R'
$C \leftarrow \alpha AB + \beta C$ <i>A</i> is $m \times m$ <i>B</i> is $m \times n$	$C \leftarrow \alpha BA + \beta C$ <i>B</i> is $m \times n$ <i>A</i> is $n \times n$

(c) Rank- k updates of a real or complex symmetric or complex Hermitian matrix:

_SYRK (UPLO, TRANS, N, K, ALPHA, A, LDA, BETA, C, LDC)
_HERK (UPLO, TRANS, N, K, ALPHA, A, LDA, BETA, C, LDC)

Operation (C is always $n \times n$).

For _SYRK routines (C is symmetric):

TRANS = 'N'	TRANS = 'T'
$C \leftarrow \alpha AA^T + \beta C$ <i>A</i> is $n \times k$	$C \leftarrow \alpha A^T A + \beta C$ <i>A</i> is $k \times n$

For _HERK routines (C is Hermitian):

TRANS = 'N'	TRANS = 'C'
$C \leftarrow \alpha AA^H + \beta C$ <i>A</i> is $n \times k$	$C \leftarrow \alpha A^H A + \beta C$ <i>A</i> is $k \times n$

(In the real case, the values ‘T’ and ‘C’ have the same meaning. In the complex case, TRANS = ‘C’ is not allowed in _SYRK and TRANS = ‘T’ is not allowed in _HERK.)

(d) Rank- $2k$ updates of a real or complex symmetric or complex Hermitian matrix:

_SYRK (UPLO, TRANS, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)
 _HERK (UPLO, TRANS, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)

Operation (C is always $n \times n$).

For _SYRK routines (C is symmetric):

TRANS = ‘N’	TRANS = ‘T’
$C \leftarrow \alpha AB^T + \alpha BA^T + \beta C$	$C \leftarrow \alpha A^T B + \alpha B^T A + \beta C$
A and B are $n \times k$	A and B are $k \times n$

For _HERK routines (C is Hermitian):

TRANS = ‘N’	TRANS = ‘C’
$C \leftarrow \alpha AB^H + \bar{\alpha} BA^H + \beta C$	$C \leftarrow \alpha A^H B + \bar{\alpha} B^H A + \beta C$
A and B are $n \times k$	A and B are $k \times n$

(In the real case, the values ‘T’ and ‘C’ have the same meaning. In the complex case, TRANS = ‘C’ is not allowed in _SYRK and TRANS = ‘T’ is not allowed in _HERK.)

(e) Triangular matrix-matrix products:

_TRMM (SIDE, UPLO, TRANSA, DIAG, M, N, ALPHA, A, LDA, B, LDB)

Operation (B is always $m \times n$, A is triangular):

	TRANSA = ‘N’	TRANSA = ‘T’	TRANSA = ‘C’
SIDE = ‘L’	$B \leftarrow \alpha AB$	$B \leftarrow \alpha A^T B$	$B \leftarrow \alpha A^H B$
	A is triangular $m \times m$	A is triangular $m \times m$	A is triangular $m \times m$
SIDE = ‘R’	$B \leftarrow \alpha BA$	$B \leftarrow \alpha BA^T$	$B \leftarrow \alpha BA^H$
	A is triangular $n \times n$	A is triangular $n \times n$	A is triangular $n \times n$

(In the real case, the values ‘T’ and ‘C’ have the same meaning.)

(f) Solution of triangular systems of equations:

_TRSM (SIDE, UPLO, TRANSA, DIAG, M, N, ALPHA, A, LDA, B, LDB)

Operation (B is always $m \times n$, A is triangular):

	TRANSA = ‘N’	TRANSA = ‘T’	TRANSA = ‘C’
SIDE = ‘L’	$B \leftarrow \alpha A^{-1} B$	$B \leftarrow \alpha A^{-T} B$	$B \leftarrow \alpha A^{-H} B$
	A is triangular $m \times m$	A is triangular $m \times m$	A is triangular $m \times m$
SIDE = ‘R’	$B \leftarrow \alpha BA^{-1}$	$B \leftarrow \alpha BA^{-T}$	$B \leftarrow \alpha BA^{-H}$
	A is triangular $n \times n$	A is triangular $n \times n$	A is triangular $n \times n$

(In the real case, the values ‘T’ and ‘C’ have the same meaning.)

7. NUMERICAL STABILITY

Although it is intended that the Level 3 BLAS be implemented as efficiently as possible, it is essential that efficiency should not be achieved at the cost of sacrificing numerical stability. This is particularly important because the Level 3 BLAS are intended to be used as building blocks in higher level algorithms of linear algebra. If the Level 3 BLAS themselves were unstable, they would destroy the stability of the algorithms that call them.

8. RATIONALE

In the design of all levels of BLAS, one of the main concerns is to keep both the calling sequences simple and the range of options limited, while at the same time maintaining sufficient functionality. This clearly implies a compromise, and a good decision is vital if the BLAS are to be accepted as a useful standard. In this section we discuss the reasoning behind some of the decisions we have made.

The three basic matrix-matrix operations were chosen because they occur in a wide range of linear algebra applications: this is consistent with the criteria used for the Level 1 and Level 2 BLAS. We have again aimed at a reasonable compromise between a much larger number of routines, each performing only one type of operation (e.g., $B \leftarrow L^T B$), and a smaller number of routines with a more complicated set of options. There are in fact, in each precision, 6 real routines performing altogether 48 different combinations of options, and 9 complex routines performing altogether 81 different combinations of options. The number of routines is much smaller than in the Level 2 BLAS.

The routines that we have specified are not intended as high-level matrix algebra routines, but rather as building blocks for the construction of such routines.

In almost every case, where appropriate, we include operations involving a matrix and its transpose (the only exceptions are the `_SYMM` and `_HEMM` routines). We could ask the user to transpose the input matrix, but feel that this would be an imposition, particularly if the BLAS routine is being called from deep within the user's code. It would also increase the amount of data movement, whereas one of the aims of our proposal is to assist the development of software that minimizes data movement.

It could also be argued that algorithms can be rewritten to require only one of the patterns of access for symmetric, Hermitian, or triangular matrices (i.e., upper or lower triangle), but we do not feel that the BLAS should be dictating this to the user.

We do not provide routines for operations involving trapezoidal matrices; all our triangular matrices are square. This is consistent with the Level 2 BLAS. It would be possible to extend the routines for triangular matrices so that they could handle trapezoidal matrices, but at the cost of introducing extra arguments. On the other hand, a trapezoidal matrix can always be partitioned into a triangular matrix and a rectangular matrix.

We have not included specialized routines to take advantage of packed storage schemes for symmetric, Hermitian, or triangular matrices, nor of compact storage schemes for banded matrices, because such storage schemes do not seem to lend

themselves to partitioning into blocks, and hence are not likely to be useful in the type of application we are aiming at. Also, packed storage is required much less with the large memory machines available today, and we wish to keep the set of routines as small as possible.

We also have not specified a set of extended-precision routines analogous to the ES and EC routines in the Level 2 BLAS, since this would require a two-dimensional array in extended precision.

As with the Level 2 BLAS, no check has been included for singularity, or near singularity, in the routines for solving triangular equations. The requirements for such a test depend on the application, and so we felt that this should not be included, but should instead be performed outside the triangular solver.

We have tried to adhere to the convention of, and maintain consistency with, the Level 2 BLAS; however, we have deliberately departed from this approach in a few cases. The input-output matrix C in the matrix-multiply routines is the analogue of the vector y in the matrix-vector product routines. But here, C always has the same dimensions, whereas y was either of length m or n depending on context. In the rank- k update routines, we have included a parameter β which was not present in the Level 2 rank update routines. Here we felt that the parameter β is useful in applications, and since the matrix multiply routines can also be viewed as rank- k update routines, we have consistency between the MM, RK, and R2K routines.

We have also added a parameter α to the routines involving triangular matrices. This was not felt to be needed in the corresponding Level 2 BLAS, since there would be little additional cost in a separate operation to scale the result vector by α . However, in the Level 3 BLAS, where there is a whole matrix to be scaled, it is advantageous to incorporate the scaling within `_TRMM` or `_TRSM`.

Additionally, we have provided for complex symmetric, as well as complex Hermitian, matrices, since they occur sufficiently often in applications.

In our proposed naming scheme, the first character (S, D, C, or Z) indicates the relevant Fortran data type. This conforms to the conventions already established for the Level 1 and Level 2 BLAS, and also other software such as Linpack. However, the fact that single- and double-precision versions of a BLAS routine have different names can be an obstacle to portability, because the actual precision of, say an S routine may differ considerably between machines. For example, SGEMM on a Cray 2 will use arithmetic with similar precision to DGEMM on an IBM 3090. The ideal solution would be to use generic names, not only for single- and double-precision versions, but also for real and complex versions. This option is not available in a standard Fortran 77 environment. However, for implementations in other environments or programming languages that do permit generic names, we propose that the first character of the Fortran 77 names should simply be omitted, giving the generic names GEMM, SYMM, SYRK, SYR2K, HEMM, HERK, HER2K, TRMM, and TRSM.

9. APPLICATIONS

The primary intended application of the Level 3 BLAS is in implementing algorithms of numerical linear algebra in terms of operations on submatrices (or blocks). There is a long history of block algorithms, e.g., [1, 4, 6, 8, 9, 11, 22, 29].

Both the NAG and the IMSL (Edition 9) libraries include such algorithms (F01BTF and F01BXF in NAG; LEQIF and LEQOF in IMSL). The earlier work was usually concerned with submatrices being transferred between the main memory and disk or tape. Similar concerns motivated work designed to exploit common page-swapping algorithms in virtual memory machines. Indeed, the techniques are similar wherever there exists any hierarchy of data storage (in terms of access speed). Additionally, full blocks, and hence the multiplication of full matrices, might appear as a subproblem when handling large sparse systems of equations (for example, [9, 23, 25]).

More recently, several workers have demonstrated the effectiveness of block algorithms on a variety of modern computer architectures, with vector-processing or parallel-processing capabilities, on which potentially high performance can easily be degraded by excessive transfer of data between different levels of memory (vector registers, cache, local memory, main memory, or solid-state disks) [2, 3, 5–7, 19–21, 24, 26, 30–32]. See Demmel et al. [10] for a proposal to develop a new linear algebra library using block algorithms wherever possible and calling Level 3 BLAS.

Here we illustrate how the Level 3 BLAS routines can be used to implement a simple algorithm of numerical linear algebra, namely, Cholesky factorization.

The strategy is to compute at each stage a block of consecutive columns of the result. The size of the block is a parameter, nb , that may be varied to suit the size of the problem and the architecture of the machine. (For transportable software, we shall need some means of determining the block size within the routine, but we set that issue aside here.)

There are other ways to organize the computation: for example, it is equally possible to compute a block of consecutive rows at each stage. The analysis of Dongarra et al. [18] can easily be extended to block algorithms. We have chosen an organization that works by columns rather than rows, one that involves fewest memory references.

Also, we have implemented the algorithms in such a way that submatrices passed to the Level 3 BLAS routines are kept as large as possible (once the block size has been fixed); this gives the greatest scope for achieving efficiency within the Level 3 BLAS. Alternatively, one might explicitly partition the matrix into, say, square blocks of size $nb \times nb$; this would require many more calls to the Level 3 BLAS routines, but might allow a more precise control of memory or of parallelism.

We assume that we are given a positive-definite symmetric matrix A whose lower triangle is stored in the lower triangle of a two-dimensional array. We wish to compute L , overwriting the given elements of A .

We can partition the matrices so that

$$\begin{aligned} \begin{pmatrix} A_{11} & A_{21}^T & A_{31}^T \\ A_{21} & A_{22} & A_{32}^T \\ A_{31} & A_{32} & A_{33} \end{pmatrix} &= \begin{pmatrix} L_{11} & & \\ L_{21} & L_{22} & \\ L_{31} & L_{32} & L_{33} \end{pmatrix} \begin{pmatrix} L_{11}^T & L_{21}^T & L_{31}^T \\ & L_{22}^T & L_{32}^T \\ & & L_{33}^T \end{pmatrix} \\ &= \begin{pmatrix} L_{11}L_{11}^T & & & \\ L_{21}L_{11}^T & L_{21}L_{21}^T + L_{22}L_{22}^T & & \\ L_{31}L_{11}^T & L_{31}L_{21}^T + L_{32}L_{22}^T & L_{31}L_{31}^T + L_{32}L_{32}^T + L_{33}L_{33}^T & \end{pmatrix} \end{aligned}$$

Here, L_{22} and L_{32} constitute the current block of columns of L to be computed, and we assume that L_{11} , L_{21} , and L_{31} constitute the blocks, if any, that have already been computed. Note that the blocks in the above partitioning are not all of equal size—the off-diagonal blocks are, in general, rectangular.

Equating blocks, we have

$$\begin{aligned} A_{22} &= L_{21}L_{21}^T + L_{22}L_{22}^T \\ A_{32} &= L_{31}L_{21}^T + L_{32}L_{22}^T, \end{aligned}$$

so that

$$\begin{aligned} L_{22}L_{22}^T &= A_{22} - L_{21}L_{21}^T \\ L_{32} &= (A_{32} - L_{31}L_{21}^T)(L_{22}^T)^{-1} \end{aligned}$$

Thus the computation of one block-column of the result involves the following operations:

- (1) update the diagonal block:

$$A'_{22} \leftarrow A_{22} - L_{21}L_{21}^T$$

- (2) compute the Cholesky factorization of the diagonal block:

$$A'_{22} \rightarrow L_{22}L_{22}^T$$

- (3) update the subdiagonal block:

$$A'_{32} \leftarrow A_{32} - L_{31}L_{21}^T$$

- (4) compute the subdiagonal block of L :

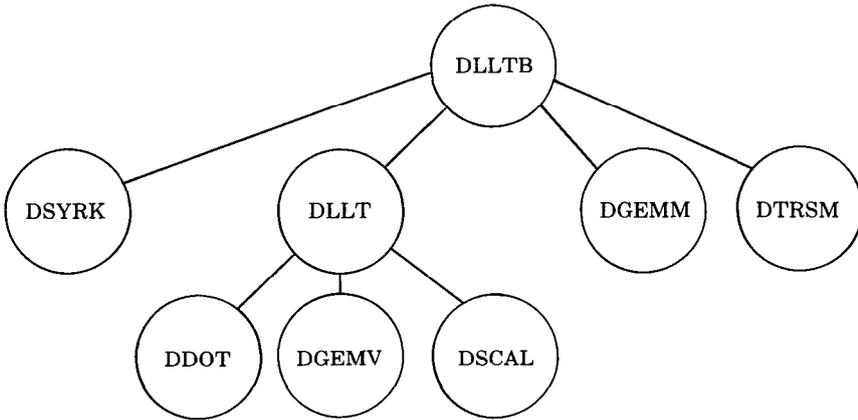
$$L_{32} \leftarrow A'_{32}(L_{22}^T)^{-1}$$

To express the complete algorithm, we adopt a notation in which the matrix is partitioned into blocks A_{ij} of size $nb \times nb$, with $p = \lceil n/nb \rceil$, and use $A_{j,1:j-1}$ to denote, for example, the block row $[A_{j,1}, A_{j,2}, \dots, A_{j,j-1}]$. The algorithm is

$$\begin{array}{ll} \text{for } j = 1 \text{ to } p & \\ \quad A_{jj} \leftarrow A_{jj} - A_{j,1:j-1}A_{j,1:j-1}^T & \quad \text{(-SYRK)} \\ \quad \text{factorize } A_{jj} & \quad \text{(unblocked algorithm)} \\ \quad A_{j+1:p,j} \leftarrow A_{j+1:p,j} - A_{j+1:p,1:j-1}A_{j,1:j-1}^T & \quad \text{(-GEMM)} \\ \quad A_{j+1:p,j} \leftarrow A_{j+1:p,j}(A_{jj}^T)^{-1} & \quad \text{(-TRSM)} \end{array}$$

In Appendix A, we give the Fortran code for a block Cholesky factorization routine DLLTB, calling Level 3 BLAS routines; and also a lower level routine DLLT, which is called by DLLTB and calls Level 1 and 2 BLAS routines. A separate lower level routine is needed, since current standard Fortran forbids recursion. The structure of DLLTB has been kept as similar as possible to that

of DLLT. The call-tree is



APPENDIX A BLOCKED CHOLESKY FACTORIZATION

```

SUBROUTINE DLLTB(N,A,LDA,INFO)
*
* Computes an L*L**T factorization of a symmetric positive-definite
* matrix A.
* Blocked version, calling Level 3 BLAS.
*
  INTEGER          INFO, LDA, N
  DOUBLE PRECISION A(LDA,*)
  INTEGER          J, JB
  INTEGER          NB
  PARAMETER       (NB=64)
  EXTERNAL        DGEMM, DLLT, DSYRK, DTRSM
*
  INFO = 0
  DO 10 J = 1, N, NB
    JB = MIN(NB,N-J+1)
*
*   Update diagonal block.
*
    CALL DSYRK('Lower','No transpose',JB,J-1,-1.0D0,A(J,1),LDA,
$      1.0D0,A(J,J),LDA)
*
*   Factorize diagonal block and test for
*   non-positive-definiteness.
*
    CALL DLLT(JB,A(J,J),LDA,INFO)
    IF (INFO.NE.0) GO TO 20
*
    IF (J+JB.LE.N) THEN
*
*   Update subdiagonal block.
*
      CALL DGEMM('No transpose','Transpose',N-J-JB+1,JB,J-1,
$        -1.0D0,A(J+JB,1),LDA,A(J,1),LDA,1.0D0,A(J+JB,J),
$        LDA)
*
*   Compute subdiagonal block of L.
*
      CALL DTRSM('Right','Lower','Transpose','Non-unit',N-J-JB+1,
$        JB,1.0D0,A(J,J),LDA,A(J+JB,J),LDA)
    END IF

```

```

10 CONTINUE
   RETURN
*
20 INFO = INFO + J - 1
   RETURN
   END
   SUBROUTINE DLLT(N,A,LDA,INFO)
*
*   Computes an L*L**T factorization of a symmetric positive-definite
*   matrix A.
*   Unblocked version, calling Level 2 and Level 1 BLAS.
*
   INTEGER          INFO, LDA, N
   DOUBLE PRECISION A(LDA,*)
   INTEGER          J
   DOUBLE PRECISION DDOT
   EXTERNAL         DDOT
   EXTERNAL         DGEMV, DSCAL

   INFO = 0
   DO 10 J = 1, N
*
*       Update a(j,j).
*
*       A(J,J) = A(J,J) - DDOT(J-1,A(J,1),LDA,A(J,1),LDA)
*
*       Compute l(j,j) and test for non-positive-definiteness.
*
*       IF (A(J,J).LE.0.0D0) GO TO 20
*       A(J,J) = SQRT(A(J,J))
*
*       IF (J.LT.N) THEN
*
*           Update elements j+1:n of j-th column.
*
*           CALL DGEMV('No transpose',N-J,J-1,-1.0D0,A(J+1,1),LDA,
$           A(J,1),LDA,1.0D0,A(J+1,J),1)
*
*           Compute elements j+1:n of j-th column of L.
*
*           CALL DSCAL(N-J,1.0D0/A(J,J),A(J+1,J),1)
*       END IF
10 CONTINUE
   RETURN
*
20 INFO = J
   RETURN
   END

```

APPENDIX B CALLING SEQUENCES FOR ALL THE LEVEL 3 BLAS

name	options	dim	scalar matrix	matrix	scalar matrix
_GEMM (TRANSA, TRANSB,	M, N, K,	ALPHA, A,	LDA, B, LDB,	BETA, C, LDC)
_SYMM (SIDE, UPLO,		M, N,	ALPHA, A,	LDA, B, LDB,	BETA, C, LDC)
_HEMM (SIDE, UPLO,		M, N,	ALPHA, A,	LDA, B, LDB,	BETA, C, LDC)
_SYRK (UPLO, TRANS,	N, K,	ALPHA, A,	LDA,	BETA, C, LDC)
_HERK (UPLO, TRANS,	N, K,	ALPHA, A,	LDA,	BETA, C, LDC)
_SYR2K(UPLO, TRANS,	N, K,	ALPHA, A,	LDA, B, LDB,	BETA, C, LDC)
_HER2K(UPLO, TRANS,	N, K,	ALPHA, A,	LDA, B, LDB,	BETA, C, LDC)
_TRMM (SIDE, UPLO, TRANSA,	DIAG,	M, N,	ALPHA, A,	LDA, B, LDB)	
_TRSM (SIDE, UPLO, TRANSA,	DIAG,	M, N,	ALPHA, A,	LDA, B, LDB)	

ACKNOWLEDGMENTS

Draft proposals for the Level 3 BLAS were discussed initially at a workshop at Argonne National Laboratory, on January 26–27, 1987, and subsequently at various meetings. We thank all the participants in those discussions for their comments and encouragement. We thank John Lewis in particular, for urging us to provide for complex symmetric matrices.

REFERENCES

1. BARRON, D. W., AND SWINNERTON-DYER, H. P. F. Solution of simultaneous linear equations using a magnetic-tape store. *Comput. J.* 3 (1960), 28–33.
2. BERRY, M., GALLIVAN, K., HARROD, W., JALBY, W., LO, S., MEIER, U., PHILIPPE, B., AND SAMEH, A. Parallel algorithms on the CEDAR system. CSRD Report 581, 1986.
3. BISCHOF, C., AND VAN LOAN, C. The WY representation for products of Householder matrices. *SIAM J. Sci. Stat. Comput.* 8, 1 (Jan. 1987), s2–s13.
4. BRONLUND, O. E., AND JOHNSEN, T. QR-factorization of partitioned matrices. *Comput. Meth. Appl. Mech. Eng.*, vol. 3, pp. 153–172, 1974.
5. BUCHER, I., AND JORDAN, T. Linear algebra programs for use on a vector computer with a secondary solid state storage device. In *Advances in Computer Methods for Partial Differential Equations*, R. Vichnevetsky and R. Stepleman, Eds. IMACS, 1984, 546–550.
6. CALAHAN, D. A. Block-oriented local-memory-based linear equation solution on the CRAY-2: Uniprocessor algorithms. In *Proceedings International Conference on Parallel Processing* (Aug. 1986). IEEE Computer Society Press, New York, 1986.
7. CARNEVALI, P., RADICATI DI BROZOLO, G., ROBERT, Y., AND SGUAZZERO, P. Efficient Fortran implementation of the Gaussian elimination and Householder reduction algorithms on the IBM 3090 vector multiprocessor. IBM ECSEC Rep. ICE-0012, 1987.
8. CHARTRES, B. Adaption of the Jacobi and Givens methods for a computer with magnetic tape backup store. Univ. of Sydney Tech. Rep. 8, 1960.
9. DAVE, A. K., AND DUFF, I. S. Sparse matrix calculations on the CRAY-2. *Parallel Comput.* 5 (July 1987), 55–64.
10. DEMMEL, J., DONGARRA, J. J., DUCROZ, J., GREENBAUM, A., HAMMARLING, S., AND SORENSEN, D. Prospectus for the development of a linear algebra library for high-performance computers. Argonne National Lab. Rep. ANL-MCS-TM-97, Sept. 1987.
11. DIETRICH, G. A new formulation of the hypermatrix Householder QR-decomposition. *Comput. Meth. Appl. Mech. Eng.* 9 (1976), 273–280.
12. DODSON, D., AND LEWIS, J. Issues relating to extension of the basic linear algebra subprograms. *ACM SIGNUM Newsl.* 20, 1 (1985), 2–18.
13. DONGARRA, J. J., BUNCH, J., MOLER, C., AND STEWART, G. *LINPACK Users' Guide*. SIAM, Philadelphia, Pa., 1979.
14. DONGARRA, J. J., DUCROZ, J., HAMMARLING, S., AND HANSON, R. An extended set of Fortran basic linear algebra subprograms. *ACM Trans. Math. Softw.* 14, 1 (Mar. 1988), 1–17.
15. DONGARRA, J. J., DUCROZ, J., HAMMARLING, S., AND HANSON, R. An extended set of Fortran basic linear algebra subprograms: Model implementation and test programs. *ACM Trans. Math. Softw.* 14, 1 (Mar. 1988), 18–32.
16. DONGARRA, J. J., DUCROZ, J., DUFF, I. S., AND HAMMARLING, S. A set of level 3 basic linear algebra subprograms: Model implementation and test programs. This issue, pp. 18–37.
17. DONGARRA, J. J., AND DUFF, I. S. Advanced architecture computers. Univ. of Tennessee, Rep. CS-89-90, Nov. 1989.
18. DONGARRA, J. J., GUSTAVSON, F., AND KARP, A. Implementing linear algebra algorithms for dense matrices on a vector pipeline machine. *SIAM Rev.* 26, 1 (1984), 91–112.
19. DONGARRA, J. J., HAMMARLING, S., AND SORENSEN, D. C. Block reduction of matrices to condensed forms for eigenvalue computations. Argonne National Lab. Rep. ANL-MCS-TM-99, Sept. 1987.
20. DONGARRA, J. J., AND HEWITT, T. Implementing dense linear algebra using multitasking on the CRAY X-MP-4. *J. Comput. Appl. Math.* 27 (1989), 215–227.

21. DONGARRA, J. J., AND SORENSEN, D. C. Linear algebra on high-performance computers. In *Proceedings Parallel Computing 85*, U. Schendel, Ed. North Holland, Amsterdam, 1986, 113–136.
22. DUCROZ, J., NUGENT, S., REID, J., AND TAYLOR, D. Solving large full sets of linear equations in a paged virtual store. *ACM Trans. Math. Softw.* 7, 4 (1981), 527–536.
23. DUFF, I. S. Full matrix techniques in sparse Gaussian elimination. In *Numerical Analysis Proceedings, Dundee 1981, Lecture Notes in Mathematics 912*. Springer-Verlag, New York, 1981, 71–84.
24. GALLIVAN, K., JALBY, W., AND MEIER, U. The use of BLAS3 in linear algebra on a parallel processor with a hierarchical memory. *SIAM J. Sci. Stat. Comput.* 8, 6 (Nov. 1987), 1079–1084.
25. GEORGE, A., AND RASHWAN, H. Auxiliary storage methods for solving finite element systems. *SIAM J. Sci. Stat. Comput.* 6, 4 (Oct. 1985), 882–910.
26. IBM. Engineering and scientific subroutine library. Program 5668-863, 1986.
27. LAWSON, C., HANSON, R., KINCAID, D., AND KROGH, F. Basic linear algebra subprograms for Fortran usage. *ACM Trans. Math. Softw.* 5 (1979), 308–323.
28. LAWSON, C., HANSON, R., KINCAID, D., AND KROGH, F. Algorithm 539: Basic linear algebra subprograms for Fortran usage. *ACM Trans. Math. Softw.* 5 (1979), 324–325.
29. MCKELLAR, A. C., AND COFFMAN, E. G., JR. Organizing matrices and matrix operations for paged memory systems. *Commun. ACM* 12, 3 (1969), 153–165.
30. ROBERT, Y., AND SGUAZZERO, P. The LU decomposition algorithm and its efficient Fortran implementation on the IBM 3090 vector multiprocessor. IBM ECSEC Rep. ICE-0006, 1987.
31. SCHREIBER, R. Module design specification (Version 1.0). SAXPY Computer Corp., 255 San Geronimo Way, Sunnyvale, CA 94086, 1986.
32. SCHREIBER, R., AND PARLETT, B. Block reflectors: Theory and computation. *SIAM J. Numer. Anal.* 25, 1 (Feb. 1988), 189–205.

Received October 1988; revised February 1989; accepted March 1989