

- [10] M. A. ELIIS, B. STROUSTRUP, *The Annotated C++ Reference Manual*, Addison-Wesley, 1990.
- [11] K. E GOREN, S MOROW, P. S PEIXO, *Data Abstraction and Object-Oriented Programming in C++*, John Wiley & Sons, Chichester, England, 1990.
- [12] ROGUE WAVE SOFTWARE, *Mith.h++*, Corvallis, Oregon, 1992.
- [13] B T. SMITH, J. M BYRNE, J. J. DONGARRA, B S GARBOV, Y. IKEBE, V. C. KEMM, AND C B MIR, *Matrix Eigen system Routines – EISPACK Guide*, vols. 6 of Lecture Notes in Computer Science, Springer-Verlag Berlin, 2 ed., 1976.

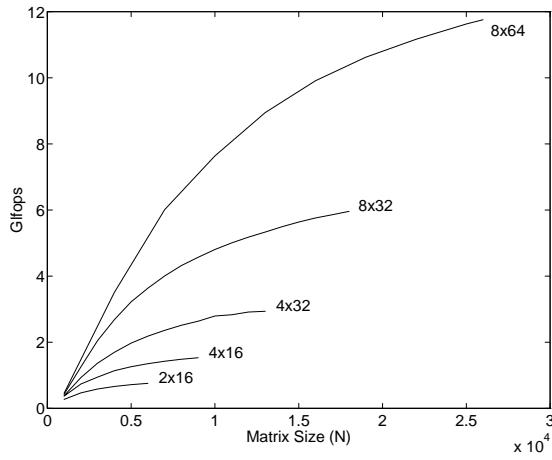


Figure 8: Performance of distributed LU factorization using an Square Block Scattered (SBS) matrix decomposition on the Intel Tualstar Delta system. Results are provided for various processor-grid configurations.

both conventional and distributed matrices as fundamental objects.

In both the parallel and sequential aspects of IA **PACK++** decoupling the matrix algorithm from a specific data decomposition provides three important attributes: (1) it results in simpler code which more closely matches the underlying mathematical formulation, (2) it allows for a “universal” algorithm rather than supporting one version for each data decomposition needed, and (3) it allows one to postpone the data decomposition decision until runtime.

We have used the **inheritance** mechanism of the object oriented design to provide a common source code for both the parallel and sequential versions. Because the parallelism is embedded in the parallel BLAS library, the sequential and parallel high level matrix algorithms in **ScalAPACK++** look the same. This tremendously simplifies the complexity of the parallel libraries.

We have used **polymorphism** or dynamic binding mechanism to achieve a truly portable matrix library in which the data decomposition may be dynamically changed at runtime.

We utilized operator and function **overloading** capabilities of C++ to simplify the syntax and user interface into the LAPACK and BLAS functions.

We utilized the **function inlining** capabilities of C++ to reduce the function call overhead usually associated with interfacing Fortran or assembly kernels.

In short, we have used various important aspects of object oriented mechanisms and C++ in the design of **LAPACK++**. These attributes were utilized not because of novelty but out of necessity to incorporate a design which provides scalability, portability, flexibility and ease-of-use.

## References

- [1] E. ADERSON, Z. BI, C. BSCHE, J. W. DUMMEL, J. J. DONGARRA, J. D. GÓZ, A. GREENBAUM, S. HAMMING, A. MÉNÉZY, S. OSTRouchov, AND D. SENSEN, *LAPACK User's Guide*, SIAM, Philadelphia, 1992.
- [2] J. GÓZ, J. J. DONGARRA, R. RIZO, D. W. WILCOX, *ScalAPACK: A Scalable Linear Algebra Library for Distributed Memory Concurrent Computers*, Fatigues of Massively Parallel Computing Mean, Virginia, October 1992.
- [3] J. GÓZ AND J. J. DONGARRA AND D. W. WILCOX, *PB-BLAS: Parallel Block Basic Linear Algebra Subroutines on Distributed Memory Concurrent Computers*, Oak Ridge National Laboratory Mathematical Sciences Section, in preparation, 1993.
- [4] R. DAMES, *newmat 07, an experimental matrix package in C++*, [robertd@kauri.vuw.ac.nz](mailto:robertd@kauri.vuw.ac.nz), 1993.
- [5] J. J. DONGARRA, J. R. BENCH, C. B. MÉNÉZY, AND G. W. STEWART, *LINPACK Users' Guide*, SIAM, Philadelphia, PA, 1979.
- [6] J. J. DONGARRA, J. D. GÓZ, I. S. DUFF, AND S. HAMMING, *A set of Level 3 Basic Linear Algebra Subprograms*, ACM Trans. Math. Soft., 16 (1990), pp. 1–17.
- [7] J. J. DONGARRA, R. RIZO, D. W. WILCOX, *LAPACK++: Object Oriented Extensions for High Performance Linear Algebra*, in preparation.
- [8] J. J. DONGARRA, R. RIZO, D. W. WILCOX, *An Object Oriented Design for High Performance Linear Algebra on Distributed Memory Architectures*, Object Oriented Numerics Conference (OON), Striver, Oregon, May 26-27, 1993.
- [9] IBM SCIENTIFIC CORPORATION, *M++ Class Library*, Elleville, Winter, 1991.

0	1	2	3	0	1	2	3	0	1	2	3
4	5	6	7	4	5	6	7	4	5	6	7
0	1	2	3	0	1	2	3	0	1	2	3
4	5	6	7	4	5	6	7	4	5	6	7
0	1	2	3	0	1	2	3	0	1	2	3
4	5	6	7	4	5	6	7	4	5	6	7
0	1	2	3	0	1	2	3	0	1	2	3
4	5	6	7	4	5	6	7	4	5	6	7
0	1	2	3	0	1	2	3	0	1	2	3
4	5	6	7	4	5	6	7	4	5	6	7
0	1	2	3	0	1	2	3	0	1	2	3
4	5	6	7	4	5	6	7	4	5	6	7

Figure 6: An example of block scattered decomposition over an 2x4 processor grid

ization algorithms. Each node of the multicopter runs a sequential LAPACK++ library that provides the object oriented framework to describe block algorithms on conventional matrices in each individual processor.

We can view the block scattered decomposition as starting a  $P \times Q$  processor grid, or template, over the matrix, where each cell of the grid covers  $r \times s$  data items and is labeled by its position in the template (Figure 6). The block and scattered decompositions may be regarded as special cases of the block scattered (SBS) decomposition. This scheme is practical and sufficiently general-purpose for most, if not all, dense linear algebra computations. Furthermore, in problems such as LU factorization, in which rows and/or columns are eliminated in successive steps, the SBS decomposition enhances scalability by ensuring processor load balance. Preliminary experiments of an object-based LU factorization algorithm using an SBS decomposition [2] suggest these algorithms scale well on multicopters. For example, on the Intel Tetherless Delta System 50 node i800-based multicopter, such algorithms can achieve nearly twelve Gflops (Figure 8).

Parallelism is exploited through the use of distributed memory versions of the Basic Linear Algebra Subprogram (BLAS) [6] [3] that perform the basic computational units of the block algorithms. This,

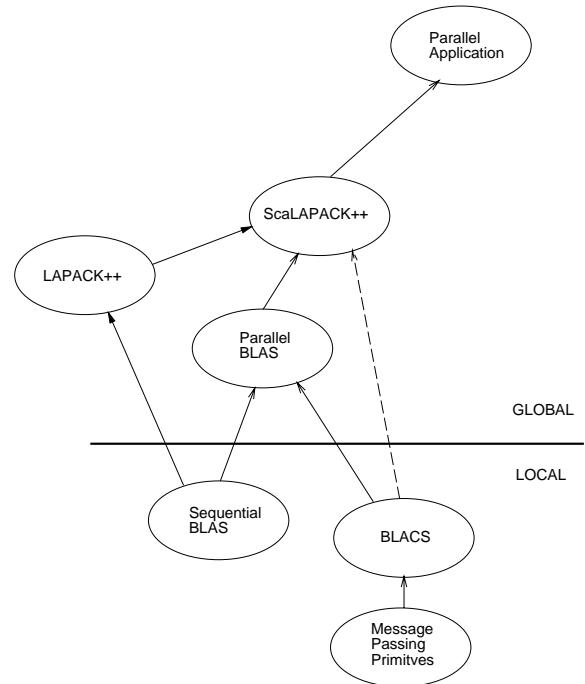


Figure 7: Design Hierarchy of ScalAPACK++. In an SIMD environment, components above the horizontal reference line represent a global viewpoint (a single distributed structure), while elements below represent a per-node local viewpoint of data

at a higher level, the block algorithms look the same for the parallel and sequential versions, and only one version of each needs to be maintained.

The benefits of an object-oriented design (Figure 7) include the ability to hide the implementation details of distributed matrices from the application programmer, and the ability to support a generic interface to basic computational kernels (BLAS), such as matrix multiply without specifying the details of the matrix storage class.

## 7 Conclusion

We have presented a design overview for object-oriented linear algebra on high performance architectures. We have also described extensions for distributed memory architectures. These designs treat

```

void poly_fit(LaVector<double> &x,
              LaVector<double> &y, LaVector<double> &p)
{
    int N = min(x.size(), y.size());
    int d = p.size();

    LaGenMatDouble P(N,d);
    LaVectorDouble a(d);
    double x_to_the_j;

    // construct Vandermonde matrix
    for (i=0; i<N; i++)
    {
        x_to_the_j = 1;
        for (j=0; j<d; j++)
        {
            P(i,j) = x_to_the_j;
            x_to_the_j *= x(i);
        }
    }
    // solve Pa = y using linear least squares
    LaLinSolveIP(P, p, y);
}

```

Figure 5: LAPACK++ code example: polynomial data fitting

degree polynomial equation

$$p(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_d x^d$$

using QR factorization. It is intended for illustrative purposes only, there are more effective methods to solve this problem.

Given the two vectors  $x$  and  $y$ , it returns the vector of coefficients  $a = \{a_0, a_1, a_2, \dots, a_{d-1}\}$ . It is assumed that  $N \gg d$ . The solution arises from solving the overdetermined Vandermonde system  $Xa = y$ :

$$\begin{bmatrix} 1 & x_0^1 & x_0^2 & \dots & x_0^d \\ 1 & x_1^1 & x_1^2 & \dots & x_1^d \\ \vdots & & & \ddots & \\ 1 & x_{N-1}^1 & x_{N-1}^2 & \dots & x_{N-1}^d \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{N-1} \end{bmatrix}$$

in the least squares sense, i.e., minimizing  $\|Xa - y\|_2$ . The resulting code is shown in figure 5.

## 6 ScALAPACK++: an extension for distributed architectures

There are various ways to extend LAPACK++. Here we discuss one such extension, ScALAPACK++ [2], for linear algebra on distributed memory architectures. The intent is that for large scale problems ScALAPACK++ should effectively exploit the computational hardware of modern ungrainsized multicore processors with up to a few thousand processors, such as the Intel Paragon and Tukwila Milles Corporation's CM.

Achieving these goals while maintaining portability, flexibility and ease-of-use can present serious challenges, since the layout of an application's data within the hierarchical memory of a concurrent computer is critical in determining the performance and scalability of the parallel code. To ensure the programmability of the library we would like details of the parallel implementation to be hidden as much as possible, but still provide the user with the capability to control the data distribution.

The design of ScALAPACK++ includes a general two-dimensional matrix decomposition that supports the most common block scattered encodings in the current literature. ScALAPACK++ can be extended to support arbitrary matrix decompositions by providing the specific parallel BLAS library to operate on such matrices.

Decoupling the matrix operations from the details of the decomposition not only simplifies the encoding of an algorithm, it also allows the possibility of postponing the decomposition until runtime. This is only possible with object-oriented programming languages, like C++, that support *dynamic binding*, or *polymorphism*—the ability to examine an object's type and dynamically select the appropriate action [10]. In many applications the optimal matrix decomposition is strongly dependent on how the matrix is utilized in other parts of the program. Furthermore, it is often necessary to dynamically alter the matrix decomposition at runtime to accommodate special routines. The ability to support dynamic run-time decomposition strategies is one of the key features of ScALAPACK++ that makes it integrable with scalable applications.

The currently supported decompositions include fine-grained matrix objects which are distributed across a  $P \times Q$  logical grid of processors. Matrices are mapped to processes using a block-scattered class of decompositions (Figure 6) that allows a wide variety of matrix mappings while enhancing scalability and maintaining good load balance for various factor-

and if the contents of A are (6) assigns the 2x2 submatrix of B to C. Note that C can also be referred as A( LaIndex(5,9,2), LaIndex(3,7,2)).

Although LAPACK+ submatrix expressions allow one to access non-contiguous rows or columns, many of the LAPACK routines only allow submatrices with unit stride in the column direction. Calling an LAPACK routine with a non-contiguous submatrix column *may* cause data to be copied into contiguous submatrix and can optionally generate a runtime warning to advise the programmer that data copying has taken place. (In Fortran, the user would need to code this by hand.)

## 4 Driver Routines

This section discusses LAPACK+ routines for solving linear systems of linear equations:

$$Ax = b,$$

where *A* is the **coefficient matrix**, *b* is the **right hand side**, and *x* is the **solution**. *A* is assumed to be a square matrix of order *n*, although underlying computational routines allow for *A* to be rectangular. For several right hand sides, we write

$$AX = B,$$

where the columns of *B* are individual right hand sides, and the columns of *X* are the corresponding solutions. The task is to find *X*, given *A* and *B*. The coefficient matrix *A* can be of the types shown in Figure 4.

The basic syntax for a linear equation driver in LAPACK+ is given by

```
LaLinSolve(op(A), X, B);
```

The matrices *A* and *B* are input, and *X* is the output. *A* is an  $M \times N$  matrix of one of the above types. Letting *rhs* denote the number of right hand sides in eq. 4, *X* and *B* are both rectangular matrices of size  $N \times \text{rhs}$ . The syntax *op(A)* can denote either *A* or the transpose of *A*, expressed as *transp(A)*.

This version requires intermediate storage of approximately  $M * (N + \text{rhs})$  elements.

In cases where no additional information is supplied, the LAPACK+ routines will attempt to follow an intelligent course of action. For example, if *LaLinSolve(A, X, B)* is called with a non-square  $M \times N$  matrix, the solution returned will be the linear least square that minimizes  $\|Ax - b\|_2$  using a QR

factorization. (If *A* is declared as *SH* then a Cholesky factorization will be used. Alternatively one can directly specify the exact factorization method such as *LaLUFactor(F, A)*). In this case, if *A* is non-square, the factors return only a partial factorization of the upper square portion of *A*.

Error conditions in performing the *LaLinSolve()* operations can be retrieved via the *LaLinSolveInfo()* function, which returns information about the last call to *LaLinSolve()*. A zero value denotes a successful completion. A value of  $-i$  denotes that the *i*th argument was somehow invalid or inappropriate. A positive value of *i* denotes that in the LU decomposition  $U(i, i) = 0$ ; the factorization has been completed but the factor *U* is exactly singular, so the solution could not be computed. In this case, the value returned by *LaLinSolve()* is a null (0x0) matrix.

### 4.1 Memory Optimizations: Factorizing in place

When using large matrices that consume a significant portion of available memory, it may be beneficial to remove the requirement of storing intermediate factorization representations at the expense of destroying the contents of the input matrix *A*. For most matrix factorizations we require temporary data structures roughly equal to the size of the original input matrix. (For general banded matrices, one needs slightly more storage due to pivoting which causes fill in additional bands.) For example, the temporary memory requirement of a square  $N \times N$  dense non-symmetric factorization can be reduced from  $N \times (N + \text{rhs})$  elements to  $N \times 1$ . Such memory-efficient factorizations are performed with the *LaLinSolveIP()* routine:

```
LaLinSolveIP(A, X, B);
```

Here the contents of *A* are overwritten (with the respective factorization). These “in-place” functions are intended for advanced programmers and are not recommended for general use. They assume the programmer’s responsibility to recognize that the contents of *A* have been destroyed; however, they can allow a large numerical problem to be solved on a machine with limited memory.

## 5 Programming Examples

This code example solves the linear least squares problem fitting *N* data points  $(x_i, y_i)$  with a *d*th

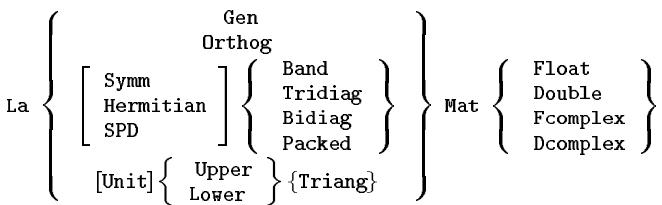


Figure 4: LAPACK++ matrix nomenclature. Brackets in square brackets are optional.

rectangular subregions of a general matrix. These regions are accessed by *reference*, that is, without copying data, and can be used in any matrix expression.

Ideally, one would like to use familiar colon notation of Fortran 90 or Matlab for expressing submatrices. However, this modification of the C++ syntax is not possible without redefining the language specifications. As a reasonable compromise, LAPACK++ denotes submatrices by specifying a subscript range through the `LaIndex()` function. For example, the  $3 \times 3$  matrix in the upper left corner of `A` is denoted

for Lapack General Matrix. The *type* suffix can be `float`, `double`, `fcomplex`, or `dcomplex`. Matrices in this category have the added property that `submat[i:j]`—the `A(0:2,0:2)` colon notation used elsewhere—can be efficiently accessed and referenced. `LaIndex` submatrix expressions may be also be used in structured algorithms. This is a necessity for describing parallelization for assignment, as in

```
A( LaIndex(0,2), LaIndex(0,2) ) = 0.0;
```

which sets the  $3 \times 3$  submatrix of `A` to zero. Following the Fortran 90 conventions, the index notation has an optional third argument denoting the stride value

`LaIndex(start, end, increment)`

If the *increment* value is not specified it is assumed to be one. The expression `LaIndex(s, e, i)` is equivalent to the index sequence

$$s, s+i, s+2i, \dots, s+\lfloor \frac{e-s}{i} \rfloor i$$

The internal representation of an index is not

expanded to a full vector, but kept in its compact triangular form. Line (1) declares `A` to be a rectangular  $200 \times 100$  matrix, with all of its elements initialized to `0.0`. Line (2) declares `B` to be an empty (uninitialized) matrix. Until `B` becomes initialized, any attempt to difference such as in `(10,7,-1)` to denote the sequence of its elements will result in a runtime error. Line (3) illustrates an equivalent way of specifying this at the time of a new object construction. Finally, line (4) demonstrates how one can initialize a  $2 \times 2$  matrix with the data from a standard C++ vector. The values are initialized in column-major form so that the first column of `E` contains  $\{1.0, 2.0\}$  and the second column contains  $\{3.0, 4.0\}$ .

```
LaGenMat<double> A(10,10), B, C; // 1
LaIndex I(1,9,2), J(1,3,2); // 2
B.ref(A(I,J)); // 3
B(2,3) = 3.1; // 4
C = B(LaIndex(2,4,2), J); // 5
```

In lines (2) and (3) we declare indices `I = {1, 3, 5, 7, 9}` and `J = {1, 3}`. Line (4) sets `B` to the specified  $5 \times 5$  submatrix of `A`. The matrix `B` can be used in any matrix expression as their basic unit of computation. It is recommended that submatrix operations be highly optimized. This, LAPACK++ provides mechanisms for accessing a submatrix as `A(5,7)`, so that a change to `B` will also affect

### 3.1.2 Submatrices

Blocked linear algebra algorithms utilize submatrices as their basic unit of computation. It is recommended that submatrix operations be highly optimized. This, LAPACK++ provides mechanisms for accessing a submatrix as `A(5,7)`, so that a change to `B` will also affect

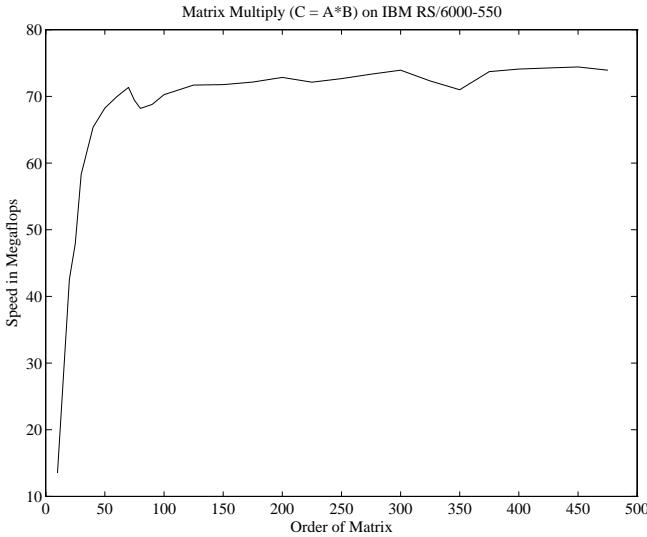


Figure 1: Performance of matrix multiply in LAPACK++ on the IBM RS/6000 Model 550 workstation. GNU g++ v. 2.3.1 was used together with the ESSL Level 3 routine `dgemm`.

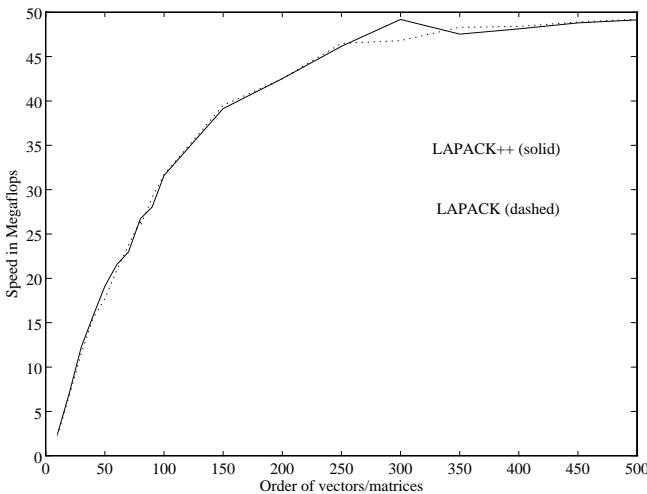


Figure 2: Performance of LAPACK++ LU factorization on the IBMRS/6000 Model 550 workstation, using GNU g++ v. 2.3.1 and BLAS routines from the IBM ESSL library. The results are indistinguishable between the Fortran and C++ interfaces.

are typically stored in column-order for compatibility with Fortran subroutines and libraries.

Various types of matrices are supported: banded, symmetric, Hermitian, packed, triangular, tridiagonal, bidiagonal, and non-symmetric. Rather than have an unstructured collection of matrix classes, LAPACK++ maintains a class hierarchy (Figure 3) to exploit commonality in the derivation of the fundamental matrix types. This limits much of the code redundancy and allows for an open-ended design which can be extended as new matrix storage structures are introduced.

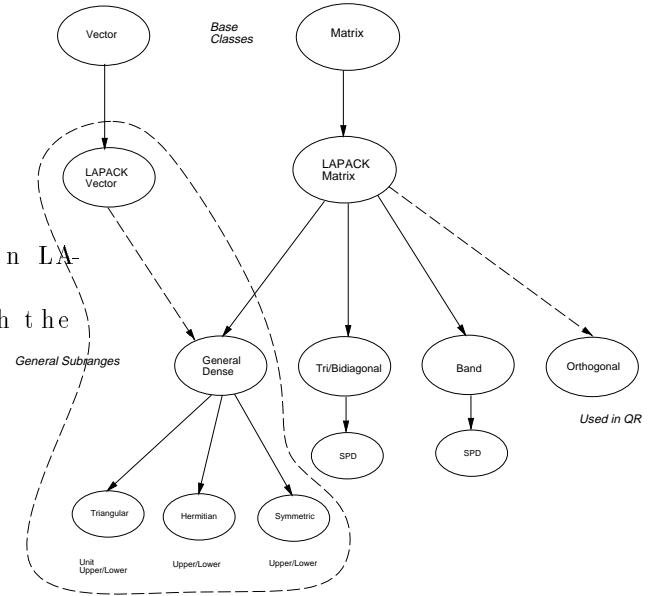


Figure 3: The Matrix Hierarchy of LAPACK++. Rectangular and Vector classes allow indexing of a rectangular submatrix by reference.

Matrix classes and other related data types specific to LAPACK++ begin with the prefix “La” to avoid naming conflicts with user-defined or third-party matrix packages. The list of valid names is a subset of the nomenclature shown in Figure 4.

### 3.1 General Matrices

One of the fundamental matrix types in LAPACK++ is a general (nonsymmetric) rectangular matrix. The possible data element types of this matrix include single and double precision of real numbers, complex numbers. The corresponding LAPACK++ names are given as

`LaGenMatType`

## 2.1 A simple code example

To illustrate how LAPACK++ simplifies the three Megaflop interface, we present a small code fragment to solve linear systems. The examples are incomplete and are

meant to merely illustrate the interface style. The next few sections discuss the details of the implementation used GNU *g++* v. 2.3.1 and utilized

Consider solving the linear LU factorization in LAPACK++:

```
#include <lapack++.h>

LaGenMatDouble A(N,N);
LaVectorDouble x(N), b(N);

// ...

LaLinSolve(A,x,b);
```

The first line includes the LAPACK++ object and numbers are very near the machine peak and function declarations. The second line demonstrates that using C++ with optimized computation be a square  $N \times N$  coefficient matrix, while the third line provides an elegant high-level interface. The fourth line declares the right-hand-side and solution vectors, sacrificing performance. Finally, the `LaLinSolve()` function in the last figure illustrates performance characteristics and calls the underlying LAPACK driver routine for LU factorization of various matrices on the same general linear equations.

Consider now solving a system with a diagonal coefficient matrix:

```
#include <lapack++.h>

LaTridiagMatDouble A(N,N);
LaVectorDouble x(N), b(N);

// ...

LaLinSolve(A, x, b);
```

The only code modification is in the declarations of fundamental objects in LAPACK++ are numerous. In this case `LaLinSolve()` calls the driver routine vectors and matrices; however, LAPACK++ uses `DGTSV()` for tridiagonal linear systems. **not** The general-purpose array package. Rather, `LaLinSolve()` function has been overloaded to provide a self-contained interface consisting of many different tasks depending on the type of object you want to solve. There are many classes to support the matrix  $A$ . There is no runtime overhead as for compatibility of the LAPACK algorithms and data with this; it is resolved by C++ at compile time.

## 2.2 Performance

The elegance of the LAPACK++ matrix class is evident, for example, in referencing a matrix element  $A(i,j)$ . By default, matrix subscripts begin at zero, keeping time performance overhead compared to similar indexing conventions of C++; however, they can be changed to start at one if desired. (Fortran programmers may prefer this convention.) Internally, LAPACK++ matrices are stored in column-major order.

performance and can utilize the BLAS kernels as efficiently as Fortran. Figure 1, for example, illus-

the **M**egaflop rating of the simple code  
is to solve  
 $C = A^*B$ ;

The square matrices of various sizes on the IBM RS/6000 Model 550 workstation. This particular implementation used GNU g++ v. 2.3.1 and utilized

presentation used C++ g++ v. 2.3.1 and utilized the Level 3 BLAS routines from the native ESSL library. The performance results are nearly identical with those of optimized Fortran calling the library. This is accomplished by *inlining* the LAPACK++BLAS kernels. That is, these functions are expanded at the point of their call by C++ compiler saving the runtime overhead of an explicit function call.

In this case the above expression calls the underlying DGEMMBLAS 3 routine. This occurs at *compile time*, without any runtime overhead. The perfor-

large numbers are very near the machine peak and illustrate that using C++ with optimized computation kernels provides an elegant high-level interface without sacrificing performance.

$\mathbb{E}[X] = \text{EP}(A, B)$

`routine`, which overwrites `A` with its LU factors. This routine essentially inlines to the underlying LAPACK routine `DGETRF()` and incurs no runtime overhead. (The `IP` suffix stands for “In Place” factorization. See Section 4.1 for details.)

### 3 LAPACK++ Matrix Objects

LAPACK++ matrices can be referenced, assigned, and used in mathematical expressions as naturally if they were an integral part of C++; the matrix expressions are evaluated at run time.

## to the FORTRAN library; some of the key routines in **Over view**

such as the matrix factorizations, are actually implemented in C++ so that the general algorithmic underlying philosophy of the LAPACK++ de-

applied to derived matrix classes, such as **singular value decomposition**, provide an interface which is simple, memory matrix objects.

powerful enough to express the sophisticated numer-

LAPACK++ provides speed and efficiency computational algorithms within LAPACK, including those which are competitive with native Fortran codes (see Section 2) for performance and/or storage. Programmers while allowing programmers to capitalize on the software engineering benefits of object oriented programming need not be concerned with the intricacies. Replacing the Fortran 77 interface of LAPACK with an object-oriented framework implies following the framework of LAPACK, the C++ exception style and allows for a more flexible and extensible interface to utilize LAPACK++ as a black box to software platform. In Section 6, for example, different types of problems, **computational routines** discuss extensions to support distributed matrices for distinct computational tasks, and **auxiliary routines** on scalable architectures [2] to perform certain subtasks or common

The motivation and design goals for LAPACK++ level computations. Each driver routine typically includes

calls a sequence of computational routines. Taken

a whole, the computational routines can perform a wide range of tasks than are covered by the driver routines.

- Maintaining competitive performance with Fortran 77. Utilizing function overloading and object inheritance details of various matrix storage schemes in C++, the procedural interface to LAPACK and their corresponding factorization subroutines simplified: two fundamental drivers and their variants, `LaLinSolve()` and `LaEigenSolve()`, replace
- Providing a simple interface that hides implementing function overloading and object inheritance details of various matrix storage schemes in C++, the procedural interface to LAPACK and their corresponding factorization subroutines simplified: two fundamental drivers and their variants, `LaLinSolve()` and `LaEigenSolve()`, replace
- Providing a universal interface and opens up hundred subroutines in the original Fortran design for integration into user-defined data structures and third-party matrix packages. LAPACK++ supports various algorithms for solving linear equations and eigenvalue problems:
- Replacing static work array limitations of Fortran with more flexible and type-safe dynamic memory allocation schemes.
  - LU Factorization
- Providing an efficient indexing scheme for matrix elements that has minimal overhead and can be optimized in most application code loops.
  - Cholesky (T) L Factorization
  - QR Factorization (linear least squares)
- Utilizing function and operator overloading in C++ to simplify and reduce the number of interface entry points to LAPACK.
  - Singular Value Decomposition (SVD)
  - Eigenvalue problems (as included in LAPACK)
- Providing the capability to access submatrices efficiently.
  - Storage Classes
  - reference, rather than by value, and perform rectangular matrices factorizations “in place” – vital for implementing blocked algorithms efficiently.
    - symmetric and symmetric positive definite (SPD)
- Providing more meaningful naming conventions for variables and functions (e.g. names no longer limited to six alphanumeric characters, and `std::` / `bi` diagonal matrices on).
- Element Data Types

LAPACK++ also provides an object-oriented interface to the Basic Linear Algebra Subprograms (BLAS) [6] allowing programmers to utilize these optimized computational kernels in their own applications. In this paper we focus on matrix factorizations, linear equations and linear least squares.

# LAPACK++: A Design Overview of Object-Oriented Extensions for High Performance Linear Algebra

Jack J. Dongarra <sup>§‡</sup>, Roldan Pozo<sup>‡</sup>, and David W. Walker<sup>§</sup>

<sup>§</sup>Oak Ridge National Laboratory  
Mathematical Sciences Section  
P. O. Box 2008, Bldg. 6012  
Oak Ridge, TN 37831-6367

<sup>‡</sup>University of Tennessee  
Department of Computer Science  
107 Ayres Hall  
Knoxville, TN 37996-1301

## Abstract

*LAPACK++* is an object-oriented C++ extension of the LAPACK (Linear Algebra PACKage) library for solving the common problem of numerical linear algebra: linear systems, linear least squares, and eigenvalue problems on high-performance computer architectures. The advantages of an object-oriented approach include the ability to encapsulate various matrix representations, hide their implementation details, reduce the number of subroutines, simplify their calling sequences, and provide an extendible software framework that can incorporate future extensions of LAPACK such as ScalAPACK++ for distributed memory architectures. We present an overview of the object-oriented design of the matrix and decomposition classes in C++ and discuss its impact on elegance, generality, and performance.

## 1 Introduction

LAPACK++ is an object-oriented C++ extension to the Fortran LAPACK library for numerical linear algebra. This package includes state-of-the-art implementations for the more common linear algebra problems encountered in scientific and engineering applications. It is based on the LINPACK [5] and EISPACK [13] libraries for solving with multi-dimensional arrays. There are several good public domain and commercial C++ packages for these problems [4][11][12]. The classes in LAPACK++, however, can easily integrate with these or with any other C++ matrix interface. These objects have been explicitly designed with block matrix algorithms and make extensive use of the level 3 BLAS. Furthermore, LAPACK++ is more than just a shell

linear equations, linear least squares, and eigenvalue problems for dense and banded systems. The current LAPACK software consists of over 1,000 routines and 600,000 lines of Fortran 77 source code.

The numerical algorithms in LAPACK utilize block-matrix operations, such as matrix multiplication, the innermost loops to achieve high performance on cached and hierarchical memory architectures. These operations, standardized as a set of subroutines called the Level 3 BLAS (Basic Linear Algebra Subprograms) [9], improve performance by increasing the granularity of the computations and keeping the most frequently accessed subregions of a matrix in the fastest level of memory. The result is that these blocked matrix versions of the fundamental algorithms typically show performance improvements of a factor of three over non-blocked versions [1].

LAPACK++ provides a framework for describing general block matrix computations in C++. Without a proper design of fundamental matrix and factorization classes, the performance benefits of blocked computation can be easily lost due to unnecessary data copying, inefficient access of submatrices, and excessive run-time overhead in the dynamic binding mechanisms of C++. LAPACK++, however, is not a general purpose arithmetic package. There are no functions, for example, for performing trigonometric operations on matrices, or for solving with multi-dimensional arrays. There are several good public domain and commercial C++ packages for these problems [4][11][12]. The classes in LAPACK++, however, can easily integrate with these or with any other C++ matrix interface. These objects have been explicitly designed with block matrix algorithms and make extensive use of the level 3 BLAS. Furthermore, LAPACK++ is more than just a shell

\*This project was supported in part by the Defense Advanced Research Projects Agency under contract DAAL03-91-C-0047, administered by the Army Research Office, the Applied Mathematical Sciences subprogram of the Office of Energy Research, U.S. Department of Energy, under Contract DE-AC05-84OR21400, and by the National Science Foundation Science and Technology Center Cooperative Agreement No. CCR-8809615.