

Programming Matrix Algorithms-by-Blocks for Thread-Level Parallelism

GREGORIO QUINTANA-ORTÍ

and

ENRIQUE S. QUINTANA-ORTÍ

Universidad Jaume I

and

ROBERT A. VAN DE GEIJN

and

FIELD G. VAN ZEE

and

ERNIE CHAN

The University of Texas at Austin

With the emergence of thread-level parallelism as the primary means for continued performance improvement, the programmability issue has reemerged as an obstacle to the use of architectural advances. We argue that evolving legacy libraries for dense and banded linear algebra is not a viable solution due to constraints imposed by early design decisions. We propose a philosophy of abstraction and separation of concerns that provides a promising solution in this problem domain. The first abstraction, FLASH, allows algorithms to express computation with matrices consisting of contiguous blocks, facilitating algorithms-by-blocks. Operand descriptions are registered for a particular operation *a priori* by the library implementor. A runtime system, SuperMatrix, uses this information to identify data dependencies between suboperations, allowing them to be scheduled to threads out-of-order and executed in parallel. But not all classical algorithms in linear algebra lend themselves to conversion to algorithms-by-blocks. We show how our recently proposed LU factorization with incremental pivoting and a closely related algorithm-by-blocks for the QR factorization, both originally designed for out-of-core computation, overcome this difficulty. Anecdotal evidence regarding the development of routines with a core functionality demonstrates how the methodology supports high productivity while experimental results suggest that high performance is abundantly achievable.

Categories and Subject Descriptors: G.4 [Mathematical Software]: —*Efficiency*

General Terms: Algorithms; Performance

Additional Key Words and Phrases: linear algebra, libraries, high-performance, multithreaded architectures

Authors' addresses: Gregorio Quintana-Ortí, Enrique S. Quintana-Ortí, Departamento de Ingeniería y Ciencia de Computadores, Universidad Jaume I, Campus Riu Sec, 12.071, Castellón, Spain, {gquintan,quintana}@icc.uji.es. Robert A. van de Geijn, Field G. Van Zee, Ernie Chan, Department of Computer Sciences, The University of Texas at Austin, Austin, TX 78712, {rvdg,field,echan}@cs.utexas.edu.

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 20YY ACM 0098-3500/20YY/1200-0001 \$5.00

1. INTRODUCTION

For the first time in history, computer architectures are approaching physical and technological barriers which make increasing the speed of a single core exceedingly difficult and economically infeasible. As a result, hardware architects have begun to design microprocessors with multiple processing cores that operate independently and share the same address space. It appears that the advent of these multi-core architectures will finally force a radical change in how applications are programmed. Specifically, developers must consider how to direct the collaboration of many concurrent threads of execution to solve a single problem. In this paper, we present what we believe is a promising solution to what is widely perceived to be a very difficult problem, targeting the domain of dense and banded linear algebra libraries. The approach views submatrices (blocks) as units of data, algorithms as operating on these blocks (algorithms-by-blocks), and schedules the operations on blocks using out-of-order techniques similar to how superscalar processors schedule instructions and resolve dependencies on individual data.

Experience gained from parallel computing architectures with complex hierarchical memories has shown that an intricate and interdependent set of requirements must be met in order to achieve the level of performance that scientific applications demand of linear algebra libraries. These requirements include data locality, load balance, and careful attention to the critical path of computation, to name a few. While it may be possible to satisfy this set of constraints when implementing a single operation on a single architecture, addressing them for an entire library of commonly used operations requires one to face the additional challenge of programmability.

We propose abandoning essentially all programming conventions that were adopted when widely-used libraries like LINPACK [Dongarra et al. 1979], LAPACK [Anderson et al. 1999], and ScaLAPACK [Choi et al. 1992], denoted by LIN(Sca)LAPACK hereafter, were designed in the last quarter of the 20th century.¹ To us, nothing associated with programming these libraries is sacred: not the algorithms, not the notation used to express algorithms, not the data structures used to store matrices, not the APIs used to code the algorithms, and not the runtime system that executes the implementations (if one ever existed). Instead we build on the notation, APIs, and tools developed as part of the FLAME project [van de Geijn and Quintana-Ortí 2008], which provide modern object-oriented abstractions to increase developer productivity and user-friendliness alike. Perhaps the most influential of these abstractions is one that provides advanced shared-memory parallelism by borrowing out-of-order scheduling techniques from sequential superscalar architectures.

From the outset, we have been focused on developing a prototype library that encompasses much of the functionality of the core of LAPACK, including LU, QR and Cholesky factorizations, and related solvers. The LU factorization with pivoting and Householder QR factorizations are key since algorithms-by-blocks for these

¹We do not mean to deminish the contributions of the LINPACK and LAPACK projects. We merely suggest that the programming styles used by those packages, while cutting-edge at the time they were developed, need to be reconsidered. We fully recognize that these projects made tremendous contributions to the field of numerical linear algebra beyond the packages that they delivered.

operations require new algorithms. This is due to the fact that when traditional partial pivoting or traditional Householder transformations are employed, entire columns have to be examined and/or updated, which becomes a synchronization point in the computation [Quintana-Ortí et al. 008a; Quintana-Ortí et al. 008b]. Fortunately, we had already developed high-performance algorithms for out-of-core computation that were algorithm-by-blocks, except that the block size targeted the movement of data from disk to main memory rather than thread-level parallelism [Gunter and van de Geijn 2005; Joffrain et al. 2004]. Thus, we knew from the start that algorithms-by-blocks were achievable for all of these important operations.

We focus on the programmability issue that developers face given that parallelism will be required in order to exploit the performance potential of multi-core and many-core systems. We point to the fact that the parallelization effort described in this paper was conceived in early September 2006. Remarkably, all results presented in this paper were already achieved by September 2007. The key has been a clear separation of concerns between the code that implements the algorithms and the runtime system that schedules tasks.

We are not alone in recognizing the utility of a runtime system that dynamically schedules subproblems for parallel execution. The PLASMA project [Buttari et al. 2007; Buttari et al. 2008], independent of our efforts but with LU and QR factorization algorithms that are similarly based on our earlier work on out-of-core algorithms, has developed a similar mechanism in the context of the QR factorization. Like the SuperMatrix system described here and in [Chan et al. 007a; Chan et al. 2008; Chan et al. 007b], the PLASMA system enqueues operation subproblems as “tasks” on a queue, builds a directed acyclic graph (DAG) to encode dependencies, and then executes the subproblems as task dependencies become satisfied [Buttari et al. 2008]. However, as part of [Buttari et al. 2008] they do not provide any source code—neither example code implementing the runtime-aware algorithm nor code implementing the runtime system itself. The same authors expand on their work in [Buttari et al. 2007] to include remarks and results for the Cholesky factorization and LU factorization with incremental pivoting based on our out-of-core algorithm [Joffrain et al. 2004], which we encourage the reader to study. In contrast to the PLASMA project, we directly address the programmability issue.

The primary contribution of the present paper lies with the more comprehensive description of how abstraction can be exploited to solve the programmability problem that faces linear algebra library development with the advent of multi-core architectures. In doing so, this paper references a number of conference publications [Chan et al. 007a; Chan et al. 2008; Chan et al. 007b; Quintana-Ortí et al. 008a; Quintana-Ortí et al. 008b; Quintana-Ortí et al. 008c] that provide evidence in support of claims that we make. As such, this is also survey paper.

The paper can be viewed as consisting of two parts. The first part describes the methodology. Section 2 provides a motivating example in the form of the LU factorization (without pivoting) and makes the case that the LIN(Sca)LAPACK design philosophies are conducive neither to ease-of-programming nor efficient parallelism. Section 3 discusses algorithms-by-blocks and the challenges they present, and provides an overview of the FLASH API, including an interface for filling

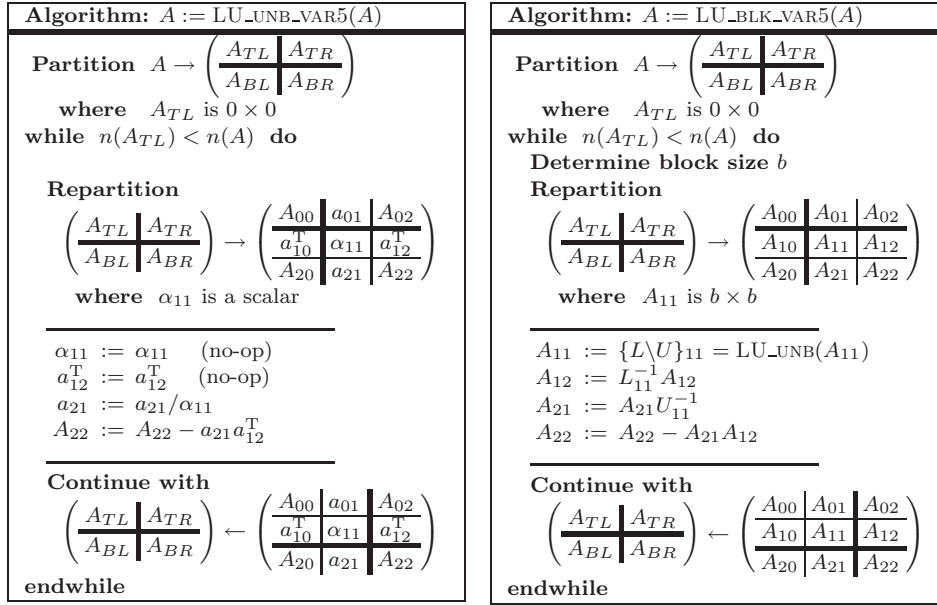


Fig. 1. Unblocked and blocked algorithms (left and right, respectively) for computing the LU factorization (without pivoting). Here, $n(B)$ stands for the number of columns of B .

matrices that uses storage-by-blocks. Section 4 gives an overview of the SuperMatrix runtime parallel execution mechanism in the context of the LU factorization. Upon finishing this first part, the reader will understand the general approach and, hopefully, the opportunities that it enables. The second part of the paper argues that the methodology has the potential of solving the programmability problem in the domain of dense and banded linear algebra libraries. It does so primarily by expanding upon results from our other recent conference publications related to the FLAME project, thereby providing better perspective on the new techniques. Section 5 recounts the authors' work in parallelizing the level-3 Basic Linear Algebra Subprograms (BLAS) operations using SuperMatrix. Section 6 expands the discussion of operations to the LU factorization with partial and incremental pivoting. This section goes into quite a bit more detail because it is important to show how new algorithms can be formulated so that the methodology can be exploited. Section 7 briefly summarizes how these results extend to advanced linear algebra operations. Performance results are reported in Section 8. Finally, Section 9 contains concluding remarks including a discussion of possible future extensions.

2. A MOTIVATING EXAMPLE: THE LU FACTORIZATION WITHOUT PIVOTING

The LU factorization (without pivoting) of an $n \times n$ matrix A is given by $A = LU$ where L is $n \times n$ unit lower triangular and U is $n \times n$ upper triangular. In traditional algorithms for this factorization, the triangular factors overwrite A , with the strictly lower triangular part of L stored on the subdiagonal elements of A and the upper triangular part of U stored on those elements of A on and above the diagonal. We

<pre> FLA_Error FLA_LU_blk_var5(FLA_Obj A, int nb_alg) { FLA_Obj ATL, ATR, A00, A01, A02, ABL, ABR, A10, A11, A12, A20, A21, A22; int b; FLA_Part_2x2(A, &ATL, &ATR, &ABL, &ABR, 0, 0, FLA_TL); while (FLA_Obj_width(ATL) < FLA_Obj_width(A)) { b = min(FLA_Obj_length(ABR), nb_alg); FLA_Repart_2x2_to_3x3(ATL, /**/ ATR, &A00, /**/ &A01, &A02, /* ***** */ /* ***** */ ABL, /**/ ABR, &A10, /**/ &A11, &A12, b, b, FLA_BR); /*-----*/ FLA_LU_unb_var5(A11); FLA_Trsm(FLA_LEFT, FLA_LOWER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_UNIT_DIAG, FLA_ONE, A11, A12); FLA_Trsm(FLA_RIGHT, FLA_UPPER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_NONUNIT_DIAG, FLA_ONE, A11, A21); FLA_Gemm(FLA_NO_TRANSPOSE, FLA_NO_TRANSPOSE, FLA_MINUS_ONE, A21, A12, FLA_ONE, A22); /*-----*/ FLA_Cont_with_3x3_to_2x2(&ATL, /**/ &ATR, A00, A01, /**/ A02, A10, A11, /**/ A12, /* ***** */ /* ***** */ &ABL, /**/ &ABR, A20, A21, /**/ A22, FLA_TL); } return FLA_SUCCESS; } </pre>	<pre> FLA_Error FLASH_LU_by_blocks_var5(FLA_Obj A) { FLA_Obj ATL, ATR, A00, A01, A02, ABL, ABR, A10, A11, A12, A20, A21, A22; FLA_Part_2x2(A, &ATL, &ATR, &ABL, &ABR, 0, 0, FLA_TL); while (FLA_Obj_width(ATL) < FLA_Obj_width(A)) { FLA_Repart_2x2_to_3x3(ATL, /**/ ATR, &A00, /**/ &A01, &A02, /* ***** */ /* ***** */ ABL, /**/ ABR, &A10, /**/ &A11, &A12, 1, 1, FLA_BR); /*-----*/ FLA_LU_unb_var5(FLASH_MATRIX_AT(A11)); FLASH_Trsm(FLA_LEFT, FLA_LOWER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_UNIT_DIAG, FLA_ONE, A11, A12); FLASH_Trsm(FLA_RIGHT, FLA_UPPER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_NONUNIT_DIAG, FLA_ONE, A11, A21); FLASH_Gemm(FLA_NO_TRANSPOSE, FLA_NO_TRANSPOSE, FLA_MINUS_ONE, A21, A12, FLA_ONE, A22); /*-----*/ FLA_Cont_with_3x3_to_2x2(&ATL, /**/ &ATR, A00, A01, /**/ A02, A10, A11, /**/ A12, /* ***** */ /* ***** */ &ABL, /**/ &ABR, A20, A21, /**/ A22, FLA_TL); } return FLA_SUCCESS; } </pre>
---	---

Fig. 2. Left: FLAME/C implementation of the blocked algorithm in Figure 1 (right). Right: FLASH implementation of the algorithm-by-blocks, which is described in Section 3.3.

denote this as $A := \{L \setminus U\}$.

2.1 A typical algorithm

In Figure 1 we give unblocked and blocked algorithms, in FLAME notation [Gunnels et al. 2001], for overwriting a matrix A with the triangular factors L and U . The unblocked algorithm on the left involves vector-vector and matrix-vector operations, which perform $O(1)$ floating-point arithmetic operations (flops) for every memory operation (memop). This ratio renders low performance on current cache-based processors as memops are considerably slower than flops on these architectures. The blocked algorithm on the right of that figure is likely to attain high performance since most computation is cast in terms of the *matrix-matrix product* (GEMM) $A_{22} := A_{22} - A_{21}A_{12}$, which performs $O(b)$ flops for every memop. In [Gunnels et al. 2001] five algorithmic variants for computing the LU factorization are identified. The VAR5 that is part of the algorithm name indicates the given algorithm corresponds to Variant 5 in that paper.

Using the FLAME/C API [Bientinesi et al. 2005], an equivalent blocked algorithm can be represented in code as presented in Figure 2 (left). Comparing and contrasting Figures 1 and 2 (left) shows how the FLAME notation, which departs from the commonly encountered loop-based algorithms, translates more naturally into code when an appropriate API is defined for the target programming language. And thus we abandon conventional algorithm notation and the LIN(Sca)LAPACK style of programming.

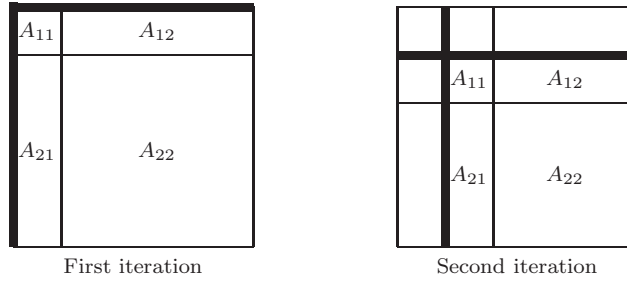


Fig. 3. First two iterations of the blocked algorithm in Figure 1 (right).

2.2 The trouble with evolving legacy code to multithreaded architectures

A commonly employed approach to parallelizing dense linear algebra operations on multithreaded architectures has been to push parallelism into multithreaded versions of the BLAS [Dongarra et al. 1990; Dongarra et al. 1988; Lawson et al. 1979]. The rationale behind this is to make minimal changes to existing codes.

In the case of the LU factorization, this means parallelism is attained only within the two TRSM and the GEMM operations:

$$\begin{aligned} A_{12} &:= L_{11}^{-1}A_{12}, \\ A_{21} &:= A_{21}U_{11}^{-1}, \\ A_{22} &:= A_{22} - A_{21}A_{12}. \end{aligned}$$

While we will see that this works well when the matrix is large and there are relatively few processors, a bottleneck forms when the ratio of the matrix dimension to the number of processors is low. In particular, the block size (variable b in Figures 1 (right) and 2 (left), respectively) must be relatively large (in practice, in the 128–256 range) so that the GEMM subproblems, which form the bulk of the LU computation, deliver high performance [Goto and van de Geijn 2008]. As a result, the LU factorization of A_{11} , typically computed by only a single processor, leaves other threads idle and therefore hinders parallel efficiency. Thus, this approach to extracting parallelism is inherently limited.

One technique that attempts to overcome such a bottleneck is to “compute ahead.” Consider the illustration in Figure 3 of the partitionings of A at the beginning of the first two iterations of the blocked algorithm for the LU factorization. In this technique, the update of A_{22} during the first iteration is broken down into the update of the part of A_{22} that will become A_{11} in the next iteration (see Figure 3), followed by the update of the rest of A_{22} . This then allows the factorization of the next A_{11} to be scheduled before the update of the remaining parts of the current A_{22} , thus overcoming the bottleneck. Extensions of this idea compute ahead several iterations in a similar manner.

The problem with this idea is that it greatly complicates the code that implements the algorithm if coded in a traditional style [Addison et al. 2003; Kurzak and Dongarra 2006; Strazdins 2001]. While feasible for a single, relatively simple algorithm like the LU factorization without pivoting or the Cholesky factorization, re-implementing a linear algebra library like LAPACK would become a daunting task if this strategy were employed.

3. ALGORITHMS-BY-BLOCKS

Fred Gustavson (IBM) has long advocated an alternative to the blocked algorithms in LAPACK [Agarwal and Gustavson 1989; Elmroth et al. 2004; Gustavson et al. 2007]. The solution, *algorithms-by-blocks*, proposes algorithms that view matrices as collections of submatrices and express their computation in terms of these submatrix blocks.

3.1 Basic idea

The idea is simple. When moving from algorithms that cast most computation in terms of matrix-vector operations to algorithms that mainly operate in terms of matrix-matrix computations, rather than improving performance by aggregating the computation into matrix-matrix computations, the developer should raise the granularity of the data by replacing each element in the matrix by a submatrix (block). Algorithms are then written as before, except with scalar operations replaced by operations on the blocks.

For example, consider the LU factorization of the partitioned matrix:

$$A \rightarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^\top & \alpha_{11} & a_{12}^\top \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right) = \left(\begin{array}{cc|c|c} \bar{\alpha}_{00} & \dots & \bar{\alpha}_{0k} & \dots & \bar{\alpha}_{0,n-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \hline \bar{\alpha}_{k0} & \dots & \bar{\alpha}_{kk} & \dots & \bar{\alpha}_{k,n-1} \\ \hline \vdots & \ddots & \vdots & \ddots & \vdots \\ \bar{\alpha}_{n-1,0} & \dots & \bar{\alpha}_{n-1,k} & \dots & \bar{\alpha}_{n-1,n-1} \end{array} \right)$$

where α_{11} and $\bar{\alpha}_{ij}$, $0 \leq i, j < n$, are all scalars. The unblocked algorithm in Figure 1 (left) can be turned into an algorithm-by-blocks by recognizing that if each element in the matrix is itself a matrix, as in

$$A \rightarrow \left(\begin{array}{c|c|c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right) = \left(\begin{array}{cc|c|c} \bar{A}_{00} & \dots & \bar{A}_{0K} & \dots & \bar{A}_{0,N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \hline \bar{A}_{K0} & \dots & \bar{A}_{KK} & \dots & \bar{A}_{K,N-1} \\ \hline \vdots & \ddots & \vdots & \ddots & \vdots \\ \bar{A}_{N-1,0} & \dots & \bar{A}_{N-1,K} & \dots & \bar{A}_{N-1,N-1} \end{array} \right)$$

with A_{11} and \bar{A}_{ij} , $0 \leq i, j < N$, are all $b \times b$ blocks, then:

- (1) $\alpha_{11} := \alpha_{11}(\text{no-op})$ becomes the LU factorization of the matrix element A_{11} :

$$A_{11} := \{L \setminus U\}_{11} = \{\bar{L} \setminus \bar{U}\}_{KK}.$$

- (2) a_{21} becomes the column vector of blocks A_{21} so that $a_{21} := a_{21}/\alpha_{11}$ is replaced by a triangular solve with multiple right-hand sides with the updated upper triangular matrix in A_{11} and each of the blocks in A_{21} :

$$A_{21} := A_{21}U_{11}^{-1} = \begin{pmatrix} \bar{A}_{K+1,K} \\ \vdots \\ \bar{A}_{N-1,K} \end{pmatrix} \bar{U}_{KK}^{-1} = \begin{pmatrix} \bar{A}_{K+1,K} \bar{U}_{KK}^{-1} \\ \vdots \\ \bar{A}_{N-1,K} \bar{U}_{KK}^{-1} \end{pmatrix}.$$

- (3) a_{12}^T becomes a row vector of blocks A_{12} so that $a_{12}^T := a_{12}^T(\text{no-op})$ becomes a triangular solve with multiple right-hand sides with the updated lower triangular matrix in A_{11} and each of the blocks in A_{12} :

$$\begin{aligned} A_{12} &:= L_{11}^{-1} A_{12} = \bar{L}_{KK}^{-1} (\bar{A}_{K,K+1} \dots \bar{A}_{K,N-1}) \\ &= (\bar{L}_{KK}^{-1} \bar{A}_{K,K+1} \dots \bar{L}_{KK}^{-1} \bar{A}_{K,N-1}). \end{aligned}$$

- (4) Each element in A_{22} describes a block that needs to be updated via a matrix-matrix product using blocks from the updated vectors of blocks A_{21} and A_{12} :

$$\begin{aligned} A_{22} &:= A_{22} - A_{21} A_{12} \\ &= \begin{pmatrix} \bar{A}_{K+1,K+1} & \dots & \bar{A}_{K+1,N-1} \\ \vdots & \ddots & \vdots \\ \bar{A}_{N-1,K+1} & \dots & \bar{A}_{N-1,N-1} \end{pmatrix} - \begin{pmatrix} \bar{A}_{K+1,K} \\ \vdots \\ \bar{A}_{N-1,K} \end{pmatrix} (\bar{A}_{K,K+1} \dots \bar{A}_{K,N-1}) \\ &= \begin{pmatrix} \bar{A}_{K+1,K+1} - \bar{A}_{K+1,K} \bar{A}_{K,K+1} & \dots & \bar{A}_{K+1,N-1} - \bar{A}_{K+1,K} \bar{A}_{K,N-1} \\ \vdots & \ddots & \vdots \\ \bar{A}_{N-1,K+1} - \bar{A}_{N-1,K} \bar{A}_{K,K+1} & \dots & \bar{A}_{N-1,N-1} - \bar{A}_{N-1,K} \bar{A}_{K,N-1} \end{pmatrix}. \end{aligned}$$

Below we will show that the algorithm-by-blocks approach also facilitates the high-performance implementation and parallel execution of matrix operations on SMP and multi-core architectures.

3.2 Obstacles

A major obstacle to algorithms-by-blocks lies with the complexity that is introduced into the code when matrices are manipulated and stored by blocks. A number of solutions have been proposed to solve this problem, ranging from storing the matrix in arrays with four or more dimensions and explicitly exposing intricate indexing into the individual elements [Guo et al. 2008], to template programming using C++ [Valsalam and Skjellum 2002], and to compiler-based solutions [Wise et al. 2001]. None of these have yielded a consistent methodology that allows the development of high-performance libraries with functionality that rivals LAPACK or FLAME. The problem is programmability.

3.3 The FLASH API for algorithms-by-blocks

Several recent efforts [Elmroth et al. 2004; Herrero 2006; Low and van de Geijn 2004] follow an approach different from those mentioned above. They view the matrix as a matrix of smaller matrices, just as it is conceptually described. Among these, the FLASH API [Low and van de Geijn 2004], which is an extension of the FLAME API, exploits the fact that FLAME encapsulates matrix information in objects by allowing elements of a matrix to themselves be descriptions of matrices. This approach supports (multi-level) hierarchical storage of matrices by submatrices (blocks). We note that, conceptually, these ideas are by no means new. The notion of storing matrices hierarchically goes back to the early 1970s [Skagestein 1972] and was rediscovered in the 1990s [Collins and Browne 1995].

Using the FLASH API, code for an algorithm-by-blocks for the LU factorization is given in Figure 2 (right). Note the similarity between that implementation and the blocked implementation for the LU factorization on the left of the same figure. That code maintains the traditional layering of subroutine calls that im-

plement linear algebra operations, which illustrates that we are willing to preserve conventions from the BLAS/LIN(Sca)LAPACK efforts that continue to benefit programmability. However, it is worth pointing out that we actually advocate a deeper layering, one that allows the programmer to invoke routines that assume certain matrix shapes. This design yields potential performance benefits when the underlying BLAS implementation provides interfaces to low-level kernels [Goto and van de Geijn 2008; Marker et al. 2007]. Such an extended matrix multiplication interface can be seen in our use and implementation of `FLASH_Gebp_nn` in Figure 4, which assumes a matrix-matrix product where A is a block and B is a row-panel.

It may seem that complexity is merely hidden in the routines `FLASH_Trsm`, and `FLASH_Gemm`. The abbreviated implementations of these operations shown in Figure 4 demonstrate how the FLASH API is used in the implementation of those routines as well. The reader can see here that many of the details of the FLASH implementation have been buried within the FLASH-aware FLAME object definition. The fact that these algorithms operate on hierarchical matrices (which use storage-by-blocks) manifests itself only through the unit block size, the use of alternative `FLASH_` routines to further break subproblems into tasks with block operands, and an additional `FLASH_MATRIX_AT` macro to extract the appropriate submatrix when wrappers to external level-3 BLAS are invoked.

As a result, transforming blocked algorithms into algorithms-by-blocks and/or developing algorithms-by-blocks from scratch using the FLASH API is straightforward.

3.4 Filling the matrix

A nontrivial matter that has prevented acceptance of alternative data structures for storing matrices has been the interface to the application. Historically, the LIN(Sca)LAPACK approach has granted application direct access to the data. This requires the application programmer to understand how data is stored, which greatly increases the programming burden on the user [Edwards and van de Geijn 2006] particularly when matrices are distributed across processors, as is true for ScaLAPACK, or stored by blocks as discussed in the current paper.

Our approach currently supports three alternative solutions.

Referencing conventional arrays. Recall that the FLASH API allows matrix elements to contain submatrices. Leaf objects in this hierarchy encapsulate the actual numerical matrix data. Given a matrix stored in conventional column-major order, a FLASH matrix object can be constructed such that the leaf matrices simply refer to submatrices of the user-supplied matrix. Notice that this means the user can access the elements of the matrix as one would when interfacing with a conventional library like LAPACK. The main disadvantage is that the leaf matrices are not stored contiguously.

Contributions to a matrix object. We will see that there is a distinct performance benefit to storing leaf matrices contiguously. Also, applications often naturally generate matrices by computing submatrices which are contributed to a larger overall matrix, possibly by adding to a partial result [Edwards and van de Geijn 2006].

For this scenario we provide routines for contribution to a FLASH matrix object. For example, FLASH provides a function whose signature is given by:

<pre> void FLASH_Trsm_llnu(FLA_Obj alpha, FLA_Obj L, FLA_Obj B) /* Special case with mode parameters FLASH_Trsm(FLA_LEFT, FLA_LOWER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_UNIT_DIAG, ...) Assumption: L consists of one block and B consists of a row of blocks */ { FLA_Obj BL, BR, BO, B1, B2; FLA_Part_1x2(B, &BL, &BR, 0, FLA_LEFT); while (FLA_Obj_width(BL) < FLA_Obj_width(B)) { FLA_Repart_1x2_to_1x3(BL, /**/ BR, &BO, /**/ &B1, &B2, 1, FLA_RIGHT); /*-----*/ FLA_Trsm(FLA_LEFT, FLA_LOWER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_UNIT_DIAG, alpha, FLASH_MATRIX_AT(L), FLASH_MATRIX_AT(B1)); /*-----*/ FLA_Cont_with_1x2_to_1x2(&BL, /**/ &BR, BO, B1, /**/ B2, FLA_LEFT); } } </pre>	<pre> void FLASH_Trsm_runn(FLA_Obj alpha, FLA_Obj U, FLA_Obj B) /* Special case with mode parameters FLASH_Trsm(FLA_RIGHT, FLA_UPPER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_NONUNIT_DIAG, ...) Assumption: U consists of one block and B consists of a column of blocks */ { FLA_Obj BT, BO, BB, B1, B2; FLA_Part_2x1(B, &BT, &BB, 0, FLA_TOP); while (FLA_Obj_length(BT) < FLA_Obj_length(B)) { FLA_Repart_2x1_to_3x1(BT, &BO, /** **/ /** **/ &B1, BB, &B2, 1, FLA_BOTTOM); /*-----*/ FLA_Trsm(FLA_RIGHT, FLA_UPPER_TRIANGULAR, FLA_NO_TRANSPOSE, FLA_NONUNIT_DIAG, alpha, FLASH_MATRIX_AT(U), FLASH_MATRIX_AT(B1)); /*-----*/ FLA_Cont_with_3x1_to_2x1(&BT, BO, B1, /** **/ /** **/ &BB, B2, FLA_TOP); } } </pre>
<pre> void FLASH_Gemm_nn(FLA_Obj alpha, FLA_Obj A, FLA_Obj B, FLA_Obj beta, FLA_Obj C) /* Special case with mode parameters FLASH_Gemm(FLA_NO_TRANSPOSE, FLA_NO_TRANSPOSE, ...) Assumption: A is a column of blocks (column panel) and B is a row of blocks (row panel) */ { FLA_Obj AT, AO, CT, CO, AB, A1, CB, C1, A2, C2; FLA_Part_2x1(A, &AT, &AB, 0, FLA_LEFT); FLA_Part_2x1(C, &CT, &CB, 0, FLA_TOP); while (FLA_Obj_length(AL) < FLA_Obj_length(A)) { FLA_Repart_2x1_to_3x1(AT, &AO, /** **/ /** **/ &A1, AB, &A2, 1, FLA_BOTTOM); FLA_Repart_2x1_to_3x1(CT, &CO, /** **/ /** **/ &C1, CB, &C2, 1, FLA_BOTTOM); /*-----*/ FLASH_Gebp_nn(FLA_NO_TRANSPOSE, FLA_NO_TRANSPOSE, alpha, A1, B, beta, C1); /*-----*/ FLA_Cont_with_3x1_to_2x1(&AT, AO, A1, /** **/ /** **/ &AB, A2, FLA_TOP); FLA_Cont_with_3x1_to_2x1(&CT, CO, C1, /** **/ /** **/ &CB, C2, FLA_TOP); } } </pre>	<pre> void FLASH_Gebp_nn(FLA_Obj alpha, FLA_Obj A, FLA_Obj B, FLA_Obj beta, FLA_Obj C) /* Special case with mode parameters FLASH_Gebp(FLA_NO_TRANSPOSE, FLA_NO_TRANSPOSE, ...) Assumption: A is a block and B, C are rows of blocks (row panels) */ { FLA_Obj BL, BR, BO, B1, B2, CL, CR, CO, C1, C2; FLA_Part_1x2(B, &BL, &BR, 0, FLA_LEFT); FLA_Part_1x2(C, &CL, &CR, 0, FLA_LEFT); while (FLA_Obj_width(BL) < FLA_Obj_width(B)) { FLA_Repart_1x2_to_1x3(BL, /**/ BR, &BO, /**/ &B1, &B2, 1, FLA_RIGHT); FLA_Repart_1x2_to_1x3(CL, /**/ CR, &CO, /**/ &C1, &C2, 1, FLA_RIGHT); /*-----*/ FLA_Gemm(FLA_NO_TRANSPOSE, FLA_NO_TRANSPOSE, alpha, FLASH_MATRIX_AT(A), FLASH_MATRIX_AT(B1), beta, FLASH_MATRIX_AT(C1)); /*-----*/ FLA_Cont_with_1x3_to_1x2(&BL, /**/ &BR, BO, B1, /**/ B2, FLA_LEFT); FLA_Cont_with_1x3_to_1x2(&CL, /**/ &CR, CO, C1, /**/ C2, FLA_LEFT); } } </pre>

Fig. 4. Code for the routines `FLASH_Trsm` and `FLASH_Gemm` that are needed to complete the LU factorization algorithm-by-blocks in Figure 2 (right). Top: `FLASH` implementations of triangular system solve with multiple right-hand sides. Bottom: `FLASH` implementations of matrix-matrix product.

```
void FLASH_Axpy_submatrix_to_global( FLA_Obj alpha,
                                     int m, int n,
                                     void* B, int ldb,
                                     int i, int j, FLA_Obj H );
```

This call accepts an $m \times n$ matrix B , stored at address B with leading dimension ldb , scales it by scalar α , and adds the result to the submatrix of H that has as its top-left element the (i, j) element of H . Note that matrix descriptor H refers to a hierarchically stored matrix object while B is a matrix created and stored in conventional column-major order. In (MATLAB) Mscript notation, this operation is given by

```
H( i:i+m-1, j:j+n-1 ) = alpha * B + H( i:i+m-1, j:j+n-1 );
```

A complementary routine allows submatrices to be extracted. Notice that given such an interface the user does not need to know how matrix H is actually stored.

This approach has been highly successful for interface applications to our Parallel Linear Algebra Package (PLAPACK) library for distributed-memory architectures [Edwards and van de Geijn 2006; van de Geijn 1997] where filling distributed matrices poses a similar challenge. Analogous interfaces are also used by the Global Array abstraction [Nieplocha et al. 1996] and PETSc [Balay et al. 2004]. We strongly recommend this departure from the LIN/(Sca)LAPACK interface to applications.

Converting whole matrices. It is also possible to allow the user to construct whole matrices in column-major order, which then may be used to build hierarchical matrices that contain the equivalent data. The submission process described above can be used for this conversion.

4. SUPERMATRIX OUT-OF-ORDER SCHEDULING

In this section we discuss how techniques used in superscalar processors can be adopted to systematically expose parallelism in algorithms-by-blocks without obfuscating the coded algorithms with further complexity.

4.1 SuperMatrix dynamic scheduling and out-of-order execution

In order to illustrate the scheduling mechanism during this subsection, we consider the matrix of 3×3 blocks

$$A \rightarrow \begin{pmatrix} \bar{A}_{00} & \bar{A}_{01} & \bar{A}_{02} \\ \bar{A}_{10} & \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{20} & \bar{A}_{21} & \bar{A}_{22} \end{pmatrix}$$

where all blocks are $b \times b$. First, the code in Figure 2 (right) is linked to the SuperMatrix runtime library and executed sequentially. As suboperations are encountered, the information associated with each suboperation is encapsulated and placed onto an internal queue. Once all operations are enqueued, the initial *analyzer stage* of execution is complete. Figure 5 provides a human-readable list corresponding to the full queue generated for a 3×3 matrix of blocks.

For example, during the first iteration of the code, the call

```
FLA_LU_unb_var5( FLASH_MATRIX_AT( A11 ) );
```

Operation	Result	Operands	
		In	In/out
1. FLA_LU_unb_var5(\bar{A}_{00})	$\bar{A}_{00} := \{\bar{L}_{00} \setminus \bar{U}_{00}\} = LU(\bar{A}_{00})$		$\bar{A}_{00} \checkmark$
2. FLA_Trsm($\dots, \bar{A}_{00}, \bar{A}_{01}$)	$\bar{A}_{01} := \bar{L}_{00}^{-1} \bar{A}_{01}$	\bar{A}_{00}	$\bar{A}_{01} \checkmark$
3. FLA_Trsm($\dots, \bar{A}_{00}, \bar{A}_{02}$)	$\bar{A}_{02} := \bar{L}_{00}^{-1} \bar{A}_{02}$	\bar{A}_{00}	$\bar{A}_{02} \checkmark$
4. FLA_Trsm($\dots, \bar{A}_{00}, \bar{A}_{10}$)	$\bar{A}_{10} := \bar{A}_{10} \bar{U}_{00}^{-1}$	\bar{A}_{00}	$\bar{A}_{10} \checkmark$
5. FLA_Trsm($\dots, \bar{A}_{00}, \bar{A}_{20}$)	$\bar{A}_{20} := \bar{A}_{20} \bar{U}_{00}^{-1}$	\bar{A}_{00}	$\bar{A}_{20} \checkmark$
6. FLA_Gemm($\dots, \bar{A}_{10}, \bar{A}_{01}, \dots, \bar{A}_{11}$)	$\bar{A}_{11} := \bar{A}_{11} - \bar{A}_{10} \bar{A}_{01}$	$\bar{A}_{10} \bar{A}_{01}$	$\bar{A}_{11} \checkmark$
7. FLA_Gemm($\dots, \bar{A}_{10}, \bar{A}_{02}, \dots, \bar{A}_{12}$)	$\bar{A}_{12} := \bar{A}_{12} - \bar{A}_{10} \bar{A}_{02}$	$\bar{A}_{10} \bar{A}_{02}$	$\bar{A}_{12} \checkmark$
8. FLA_Gemm($\dots, \bar{A}_{20}, \bar{A}_{01}, \dots, \bar{A}_{21}$)	$\bar{A}_{21} := \bar{A}_{21} - \bar{A}_{20} \bar{A}_{01}$	$\bar{A}_{20} \bar{A}_{01}$	$\bar{A}_{21} \checkmark$
9. FLA_Gemm($\dots, \bar{A}_{20}, \bar{A}_{02}, \dots, \bar{A}_{22}$)	$\bar{A}_{22} := \bar{A}_{22} - \bar{A}_{20} \bar{A}_{02}$	$\bar{A}_{20} \bar{A}_{02}$	$\bar{A}_{22} \checkmark$
10. FLA_LU_unb_var5(\bar{A}_{11})	$\bar{A}_{11} := \{\bar{L}_{11} \setminus \bar{U}_{11}\} = LU(\bar{A}_{11})$		\bar{A}_{11}
11. FLA_Trsm($\dots, \bar{A}_{11}, \bar{A}_{12}$)	$\bar{A}_{12} := \bar{L}_{11}^{-1} \bar{A}_{12}$	\bar{A}_{11}	\bar{A}_{12}
12. FLA_Trsm($\dots, \bar{A}_{11}, \bar{A}_{21}$)	$\bar{A}_{21} := \bar{A}_{21} \bar{U}_{11}^{-1}$	\bar{A}_{11}	\bar{A}_{21}
13. FLA_Gemm($\dots, \bar{A}_{21}, \bar{A}_{12}, \dots, \bar{A}_{22}$)	$\bar{A}_{22} := \bar{A}_{22} - \bar{A}_{21} \bar{A}_{12}$	$\bar{A}_{21} \bar{A}_{12}$	\bar{A}_{22}
14. FLA_LU_unb_var5(\bar{A}_{22})	$\bar{A}_{22} := \{\bar{L}_{22} \setminus \bar{U}_{22}\} = LU(\bar{A}_{22})$		\bar{A}_{22}

Fig. 5. Complete list of operations to be performed on blocks for the LU factorization (without pivoting) of a 3×3 matrix of blocks using algorithm-by-blocks. The “ \checkmark ”-marks denote those operands that are available immediately at the beginning of the algorithm (ie: those operands that are not dependent upon other operations).

inserts the LU factorization of block \bar{A}_{00} onto the list. During the same iteration, suboperations encountered inside `FLASH_Trsm` or `FLASH_Syrk` also enqueue their corresponding task entries. The order of the operations in the list, together with the operands that are read (input operands) and written (output operands) in the operation, determine the dependencies among matrix operations. Thus, the second operation in the list, which has \bar{A}_{00} as an input operand and \bar{A}_{01} as both an input and an output operand, requires the first operation to be completed before it may begin. The list denotes available operands with a “ \checkmark ”-mark; these operands are not dependent upon the completion of any other operations. They also happen to represent the operands that are available at the beginning of the algorithm-by-blocks since the list captures the state of the queue before execution.

During the *scheduler/dispatcher stage*, operations that have all operands available are scheduled for execution. As computation progresses, dependencies are satisfied and new operands become available, allowing more operations to be dequeued and executed (see Figure 6). The overhead of this runtime mechanism is amortized over a large amount of computation, and therefore its overall cost is minor.

Thus, we combine two techniques from superscalar processors, dynamic scheduling and out-of-order execution, while hiding the management of data dependencies from both library developers and users. This approach is similar in philosophy to the inspector-executor paradigm for parallelization [Lu et al. 1997; von Hanxleden et al. 1992], but that work solves a very different problem. This approach also re-

Operation/Result	Original table		After 1st oper.		After 5th oper.		After 9th oper.	
	In	In/out	In	In/out	In	In/out	In	In/out
1. $LU(\bar{A}_{00})$		$\bar{A}_{00}\checkmark$						
2. $\text{TRISL}(\bar{A}_{00})^{-1}\bar{A}_{01}$	\bar{A}_{00}	$\bar{A}_{01}\checkmark$	$\bar{A}_{00}\checkmark$	$\bar{A}_{01}\checkmark$				
3. $\text{TRISL}(\bar{A}_{00})^{-1}\bar{A}_{02}$	\bar{A}_{00}	$\bar{A}_{02}\checkmark$	$\bar{A}_{00}\checkmark$	$\bar{A}_{02}\checkmark$				
4. $\bar{A}_{10}\text{TRIU}(\bar{A}_{00})^{-1}$	\bar{A}_{00}	$\bar{A}_{10}\checkmark$	$\bar{A}_{00}\checkmark$	$\bar{A}_{10}\checkmark$				
5. $\bar{A}_{20}\text{TRIU}(\bar{A}_{00})^{-1}$	\bar{A}_{00}	$\bar{A}_{20}\checkmark$	$\bar{A}_{00}\checkmark$	$\bar{A}_{20}\checkmark$				
6. $\bar{A}_{11}-\bar{A}_{10}\bar{A}_{01}$	\bar{A}_{10} \bar{A}_{01}	$\bar{A}_{11}\checkmark$	\bar{A}_{10} \bar{A}_{01}	$\bar{A}_{11}\checkmark$	$\bar{A}_{10}\checkmark$ $\bar{A}_{01}\checkmark$	$\bar{A}_{11}\checkmark$		
7. $\bar{A}_{12}-\bar{A}_{10}\bar{A}_{02}$	\bar{A}_{10} \bar{A}_{02}	$\bar{A}_{12}\checkmark$	\bar{A}_{10} \bar{A}_{02}	$\bar{A}_{12}\checkmark$	$\bar{A}_{10}\checkmark$ $\bar{A}_{02}\checkmark$	$\bar{A}_{12}\checkmark$		
8. $\bar{A}_{21}-\bar{A}_{20}\bar{A}_{01}$	\bar{A}_{20} \bar{A}_{01}	$\bar{A}_{21}\checkmark$	\bar{A}_{20} \bar{A}_{01}	$\bar{A}_{21}\checkmark$	$\bar{A}_{20}\checkmark$ $\bar{A}_{01}\checkmark$	$\bar{A}_{21}\checkmark$		
9. $\bar{A}_{22}-\bar{A}_{20}\bar{A}_{02}$	\bar{A}_{20} \bar{A}_{02}	$\bar{A}_{22}\checkmark$	\bar{A}_{20} \bar{A}_{02}	$\bar{A}_{22}\checkmark$	$\bar{A}_{20}\checkmark$ $\bar{A}_{02}\checkmark$	$\bar{A}_{22}\checkmark$		
10. $LU(\bar{A}_{11})$		\bar{A}_{11}		\bar{A}_{11}		\bar{A}_{11}		$\bar{A}_{11}\checkmark$
11. $\text{TRISL}(\bar{A}_{11})^{-1}\bar{A}_{12}$	\bar{A}_{11}	\bar{A}_{12}	\bar{A}_{11}	\bar{A}_{12}	\bar{A}_{11}	\bar{A}_{12}	\bar{A}_{11}	$\bar{A}_{12}\checkmark$
12. $\bar{A}_{21}\text{TRIU}(\bar{A}_{11})^{-1}$	\bar{A}_{11}	\bar{A}_{21}	\bar{A}_{11}	\bar{A}_{21}	\bar{A}_{11}	\bar{A}_{21}	\bar{A}_{11}	$\bar{A}_{21}\checkmark$
13. $\bar{A}_{22}-\bar{A}_{21}\bar{A}_{12}$	\bar{A}_{21} \bar{A}_{12}	\bar{A}_{22}	\bar{A}_{21} \bar{A}_{12}	\bar{A}_{22}	\bar{A}_{21} \bar{A}_{12}	\bar{A}_{22}	\bar{A}_{21} \bar{A}_{12}	$\bar{A}_{22}\checkmark$
14. $LU(\bar{A}_{22})$		\bar{A}_{22}		\bar{A}_{22}		\bar{A}_{22}		$\bar{A}_{22}\checkmark$

Fig. 6. An illustration of the scheduling of operations for the LU factorization (without pivoting) of a 3×3 matrix of blocks using algorithm-by-blocks. Here, $\text{TRIU}(B)$ stands for the upper triangular part of B while $\text{TRISL}(B)$ denotes the matrix consisting of the lower triangular part of B with the diagonal entries replaced by ones.

reflects a shift from control-level parallelism, specified strictly by the order in which operations appear in the code, to data-flow parallelism, which is restricted only by true data dependencies and availability of compute resources.

5. AN EXPERIMENT IN PROGRAMMABILITY: THE LEVEL-3 BLAS

In [Chan et al. 007a] we report on the implementation of the level-3 BLAS using FLASH and SuperMatrix. In this section we briefly summarize the insights from that paper with a primary focus on what it tells us about how the approach addresses the programmability issue.

When we commenced parallelizing the level-3 BLAS, we had a full set of sequential level-3 BLAS implemented using the FLAME API. It is important to realize that a “full set” entails all datatypes² and all unblocked and blocked algorithm variants³ for all operations that constitute the level-3 BLAS. We also had an implementation of the FLASH extension to FLAME and the SuperMatrix runtime system for scheduling the tasks used for the execution of a SuperMatrix-enabled algorithm. This implementation had previously been used only for the Cholesky factorization.

Two of the authors, Ernie Chan and Field Van Zee, spent a weekend implementing and testing the parallelization, yielding a full set of multithreaded level-3 BLAS

²This includes single-precision and double-precision for real and complex operations.

³The FLAME methodology often yields half a dozen or more algorithms for each operation.

using FLASH and SuperMatrix. This productivity attests as to how effectively the methodology addresses the programmability issue. Impressive performance is reported in [Chan et al. 007b] despite the absence of dependencies in many of the reported operations, the lack of which reduces much of the SuperMatrix system to overhead.

6. LAPACK-LEVEL OPERATIONS: DENSE FACTORIZATIONS

The LAPACK library provides functionality one level above the level-3 BLAS. The subset with which we will primarily concern ourselves in this section includes the LU with pivoting, QR, and Cholesky factorizations.

6.1 An algorithm-by-blocks for the LU factorization with pivoting

It becomes immediately obvious that algorithms-by-blocks for the LU factorization with partial pivoting and the QR factorization based on Householder transformations require us to abandon the tried-and-trusted algorithms incorporated in LAPACK since pivoting information for the former and the computation of Householder transformations for the latter require access to columns that span multiple blocks. Unblocked and blocked algorithms for the LU factorization with partial pivoting, in FLAME notation, are given in Figure 7.

Thus, a second major obstacle to algorithms-by-blocks is that not all operations lend themselves nicely to this class of algorithms, with a clear example being the LU factorization when pivoting for stability enters the picture. We next describe our solution to this problem, inspired by out-of-core tiled algorithms for the QR and LU factorizations [Gunter and van de Geijn 2005; Joffrain et al. 2004; Quintana-Ortí and van de Geijn 2009].

Traditional algorithms for the LU factorization with partial pivoting exhibit the property that an updated column is required for a critical computation; in order to compute which row to pivot during the k -th iteration, the k -th column must have been updated with respect to all previous computation. This greatly restricts the order in which the computations can be performed. The problem is compounded by the fact that the column needed for computing which row to pivot, as well as the row to be pivoted, likely spans multiple blocks. This need for viewing and/or storing matrices by blocks was also observed for out-of-core dense linear algebra computations [Toledo 1999] and the implementation of dense linear algebra operations on distributed-memory architectures [Choi et al. 1992; van de Geijn 1997].

We will describe how partial pivoting can be modified to facilitate an algorithm-by-blocks. We do so by first reviewing results from [Joffrain et al. 2004; Quintana-Ortí and van de Geijn 2009] that show how an LU factorization can be updated while incorporating pivoting. Afterward, we generalize the insights to the desired algorithm-by-blocks.

6.1.1 Updating an LU factorization. We briefly review how to compute the LU factorization of a matrix A of the form

$$A = \left(\begin{array}{c|c} B & C \\ \hline D & E \end{array} \right) \quad (1)$$

Algorithm: $[A, p] := \text{LUPP_UNB_VAR5}(A)$	Algorithm: $[A, p] := \text{LUPP_BLK_VAR5}(A)$
<p>Partition</p> $A \rightarrow \left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right), p \rightarrow \left(\begin{array}{c} p_T \\ \hline p_B \end{array} \right)$ <p>where A_{TL} is 0×0 and p_T has 0 elements</p> <p>while $n(A_{TL}) < n(A)$ do</p> <p>Repartition</p> $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c c c} A_{00} & a_{01} & A_{02} \\ \hline a_{10} & \alpha_{11} & a_{12} \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right),$ $\left(\begin{array}{c} p_T \\ \hline p_B \end{array} \right) \rightarrow \left(\begin{array}{c} p_0 \\ \hline \pi_1 \\ \hline p_2 \end{array} \right)$ <p>where α_{11} and π_1 are scalars</p> <hr style="width: 30%; margin: 10px auto;"/> $\left[\left(\begin{array}{c} \alpha_{11} \\ \hline a_{21} \end{array} \right), \pi_1 \right]$ <p style="text-align: center;">:= Pivot $\left(\begin{array}{c} \alpha_{11} \\ \hline a_{21} \end{array} \right)$</p> $\left(\begin{array}{c c} a_{10} & a_{12} \\ \hline A_{20} & A_{22} \end{array} \right)$ <p style="text-align: center;">:= $P(\pi_1) \left(\begin{array}{c c} a_{10} & a_{12} \\ \hline A_{20} & A_{22} \end{array} \right)$</p> $a_{21} := l_{21} = a_{21} / \alpha_{11}$ $A_{22} := A_{22} - l_{21} u_{12}^T$ $= A_{22} - a_{21} a_{12}^T$ <hr style="width: 30%; margin: 10px auto;"/> <p>Continue with</p> $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c c c} A_{00} & a_{01} & A_{02} \\ \hline a_{10} & \alpha_{11} & a_{12} \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right),$ $\left(\begin{array}{c} p_T \\ \hline p_B \end{array} \right) \leftarrow \left(\begin{array}{c} p_0 \\ \hline \pi_1 \\ \hline p_2 \end{array} \right)$ <p>endwhile</p>	<p>Partition</p> $A \rightarrow \left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right), p \rightarrow \left(\begin{array}{c} p_T \\ \hline p_B \end{array} \right)$ <p>where A_{TL} is 0×0 and p_T has 0 elements</p> <p>while $n(A_{TL}) < n(A)$ do</p> <p>Determine block size b</p> <p>Repartition</p> $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c c c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right),$ $\left(\begin{array}{c} p_T \\ \hline p_B \end{array} \right) \rightarrow \left(\begin{array}{c} p_0 \\ \hline p_1 \\ \hline p_2 \end{array} \right)$ <p>where A_{11} is $b \times b$ and p_1 has b elements</p> <hr style="width: 30%; margin: 10px auto;"/> $\left[\left(\begin{array}{c} A_{11} \\ \hline A_{21} \end{array} \right), p_1 \right] := \left[\left(\begin{array}{c} \{LU\}_{11} \\ \hline L_{21} \end{array} \right), p_1 \right]$ <p style="text-align: center;">= $\text{LUPP_UNB} \left(\begin{array}{c} A_{11} \\ \hline A_{21} \end{array} \right)$</p> $\left(\begin{array}{c c} A_{10} & A_{12} \\ \hline A_{20} & A_{22} \end{array} \right)$ <p style="text-align: center;">:= $P(p_1) \left(\begin{array}{c c} A_{10} & A_{12} \\ \hline A_{20} & A_{22} \end{array} \right)$</p> $A_{12} := U_{12} = L_{11}^{-1} A_{12}$ $A_{22} := A_{22} - L_{21} U_{12}$ $= A_{22} - A_{21} A_{12}$ <hr style="width: 30%; margin: 10px auto;"/> <p>Continue with</p> $\left(\begin{array}{c c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c c c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right),$ $\left(\begin{array}{c} p_T \\ \hline p_B \end{array} \right) \leftarrow \left(\begin{array}{c} p_0 \\ \hline p_1 \\ \hline p_2 \end{array} \right)$ <p>endwhile</p>

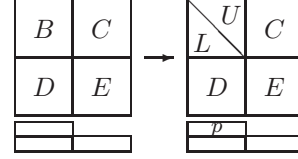
Fig. 7. Unblocked and blocked algorithms (left and right, respectively) for computing the LU factorization with partial pivoting. $\text{Pivot}(v)$ refers to a function that returns the index of the entry of largest magnitude of a vector v and interchanges that element with the first entry of v . $P(\pi_1)$ and $P(p_1)$ denote permutation matrices formed from the rows interchanges registered in π_1 and p_1 , respectively.

in such a way that the LU factorization with partial pivoting of B can be reused if D , C , and E change. In our description, we assume that both B and E are square matrices.

The following procedure [Joffrain et al. 2004; Quintana-Ortí and van de Geijn 2009], consisting of 5 steps, computes an *LU factorization with incremental pivoting* of the matrix in (1):

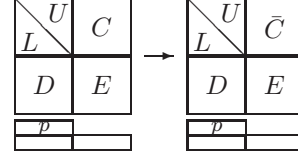
Step 1: Factor B . Compute the LU factorization with partial pivoting of B :

$$[B, p] := [\{L \setminus U\}, p] = \text{LUPP_BLK}(B).$$



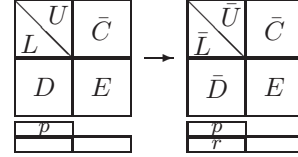
Step 2: Update C consistent with the factorization of B (using forward substitution):

$$\bar{C} := L^{-1}P(p)C = \text{TRSM_LLNU}(L, P(p)C).$$



Step 3: Factor $\begin{pmatrix} U \\ D \end{pmatrix}$. Compute the LU factorization with partial pivoting:

$$\begin{aligned} \left[\begin{pmatrix} U \\ D \end{pmatrix}, r \right] &:= \left[\begin{pmatrix} \bar{L} \setminus \bar{U} \\ D \end{pmatrix}, r \right] \\ &= \text{LUPP_SA_BLK} \left(\begin{pmatrix} U \\ D \end{pmatrix} \right). \end{aligned}$$

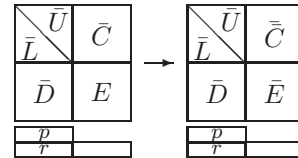


Here, \bar{U} overwrites the upper triangular part of B (where U was stored before this operation). The lower triangular matrix \bar{L} that results needs to be stored separately since both L , computed in Step 1 and used at Step 2, and \bar{L} are needed during the forward substitution stage when solving a linear system.

Care must be taken in this step not to completely fill the zeroes below U , which would greatly increase the computational cost of the next step and the storage costs of both this and the next step. The procedure computes a “structure-aware” (SA) LU factorization with partial pivoting, employing a blocked algorithm that combines the LINPACK and LAPACK styles of pivoting. For details, see algorithm $\text{LU}_{\text{BLK}}^{\text{SA-LIN}}$ in [Quintana-Ortí and van de Geijn 2009].

Step 4: Update $\begin{pmatrix} \bar{C} \\ E \end{pmatrix}$ consistent with the factorization of $\begin{pmatrix} U \\ D \end{pmatrix}$:

$$\begin{aligned} \begin{pmatrix} \bar{\bar{C}} \\ \bar{E} \end{pmatrix} &:= \begin{pmatrix} \bar{L} & | & 0 \\ \bar{L} & | & I \end{pmatrix}^{-1} P(r) \begin{pmatrix} \bar{C} \\ E \end{pmatrix} \\ &= \text{TRSM_SA_LLNU} \left(\begin{pmatrix} \bar{L} & | & 0 \\ \bar{L} & | & I \end{pmatrix}, P(r) \begin{pmatrix} \bar{C} \\ E \end{pmatrix} \right). \end{aligned}$$




```

for  $k = 0 : N - 1$ 
     $[A_{kk}, p_{kk}] := \text{LUPP\_BLK}(A_{kk})$ 
    for  $j = k + 1 : N - 1$ 
         $A_{kj} := \text{TRSM\_LLNU}(A_{kk}, P(p_{kk})A_{kj})$ 
    endfor
    for  $i = k + 1 : N - 1$ 
         $\left[ \begin{pmatrix} A_{kk} \\ A_{ik} \end{pmatrix}, L_{ik}, p_{ik} \right] := \text{LUPP\_SA\_BLK} \left( \frac{A_{kk}}{A_{ik}} \right)$ 
        for  $j = k + 1 : N - 1$ 
             $\left[ \begin{pmatrix} A_{kj} \\ A_{ij} \end{pmatrix} \right] := \text{TRSM\_SA\_LLNU} \left( \text{TRIL} \left( \left( \frac{L_{ik}}{A_{ik}} \middle| 0 \right) \right), P(p_{ik}) \left( \frac{A_{kj}}{A_{ij}} \right) \right)$ 
        endfor
    endfor
endfor
    
```

Fig. 8. Algorithm-by-blocks for the LU factorization with incremental pivoting. Here, $\text{TRIL}(B)$ stands for the lower triangular part of B . The actual implementation is similar to those in Figures 2 (right) and 4, but for conciseness we present it as loops.

Again, care must be taken in this step to exploit the zeroes below the diagonal of the upper triangular matrix produced in the previous step. This structure-aware procedure, though not equivalent to a clean triangular system solve (plus the application of the corresponding permutations), can be performed in terms of level-3 BLAS and presents essentially the same computational cost, modulo a lower order term. For details, see algorithm $\text{FS}_{\text{BLK}}^{\text{SA-LIN}}$ in [Quintana-Ortí and van de Geijn 2009].

Step 5: Factor E . Finally, compute the LU factorization with partial pivoting

$$\left[\begin{array}{c} \bar{\bar{E}}, s \end{array} \right] := \left[\{ \tilde{L} \setminus \tilde{U} \}, s \right] = \text{LUPP_BLK}(\bar{E}).$$

Overall, the five steps of the procedure apply Gauss transforms and permutations to reduce A to the upper triangular matrix

$$\left(\begin{array}{c|c} \bar{U} & \bar{C} \\ \hline 0 & \bar{U} \end{array} \right).$$

6.1.2 *An Algorithm-By-Blocks.* The insights from the previous section naturally extend to an algorithm-by-blocks for the LU factorization with incremental pivoting [Joffrain et al. 2004; Quintana-Ortí et al. 008a]. Consider the partitioning by *blocks*

$$A = \begin{pmatrix} A_{00} & A_{01} & \dots & A_{0,N-1} \\ A_{10} & A_{11} & \dots & A_{1,N-1} \\ \vdots & \vdots & \ddots & \vdots \\ A_{N-1,0} & A_{N-1,1} & \dots & A_{N-1,N-1} \end{pmatrix} \quad (2)$$

where, for simplicity, A_{ij} , $0 \leq i, j < N$, are considered to be of size $b \times b$. Then the algorithm in Figure 8 is a generalization of the algorithm described in Section 6.1.1.

While there is some flexibility in the order in which the loops are arranged, the SuperMatrix runtime system, described in Section 4, rearranges the operations, and therefore the exact order of the loops is not important.

6.1.3 Stability. Strictly speaking, the LU factorization with partial pivoting is not numerically stable; theory predicts that so-called element growth proportional to 2^n may occur. It is *practice* that taught us to rely on this method. In [Quintana-Ortí and van de Geijn 2009] we discuss how the stability of incremental pivoting relates to that of partial pivoting and pairwise pivoting [Sorensen 1985]. In summary, incremental pivoting is a blocked variant of pairwise pivoting, being equivalent to partial pivoting when the block size equals the matrix dimension and to pairwise pivoting when the block size is 1. Element growth for partial pivoting is known to be bounded by 4^n . Therefore, element growth for incremental pivoting can be expected to be bounded by a factor proportional to 2^n and 4^n , depending on the block size. The results in [Quintana-Ortí and van de Geijn 2009] provide evidence in support of this observation. However, as was the case for partial pivoting, further practical experience will be needed to establish incremental pivoting as being a numerically stable method.

6.2 An algorithm-by-blocks for the QR factorization

The QR factorization of an $m \times n$ matrix is given by $A = QR$ where Q is an $m \times m$ orthogonal matrix and R is an $m \times n$ upper triangular matrix. Although there exist several approaches to compute this factorization, the algorithm based on Householder reflectors [Golub and Loan 1996] is usually chosen when seeking high performance.

The QR factorization based on Householder reflectors and the LU factorization with partial pivoting share the property that an updated column is required for a critical computation; in the case of the QR factorization, the k -th column must have been updated with respect to all previous computation before the Householder reflectors that annihilate subdiagonal entries in this column can be computed. This greatly restricts the order in which the computations can be performed.

An algorithm-by-blocks for the QR factorization can be obtained following the out-of-core algorithm in [Gunter and van de Geijn 2005].

The algorithm-by-blocks for the QR factorization incurs a certain extra cost when compared with the traditional implementation of the QR factorization via Householder reflectors. This overhead is negligible for matrices of medium and large size. The use of orthogonal transformations ensures that the algorithm-by-blocks and the traditional QR factorization are numerically stable.

For details, see [Quintana-Ortí et al. 008b].

6.3 An algorithm-by-blocks for the Cholesky factorization

Given a symmetric positive-definite (SPD) matrix A , its Cholesky factorization is given by $A = LL^T$ (or $A = U^T U$) where L is lower triangular (or U is upper triangular). The construction of an algorithm-by-blocks to obtain the Cholesky factorization is straightforward. The algorithm is illustrated in Figure 9 for an

```

for  $k = 0 : N - 1$ 
     $[A_{kk}] := \{L_{kk} \setminus A_{kk}\} = \text{CHOL\_BLK}(A_{kk})$ 
    for  $i = k + 1 : N - 1$ 
         $A_{ik} := L_{ik} = A_{ik} L_{kk}^{-T}$ 
        for  $j = k + 1 : i$ 
             $A_{ij} := A_{ij} - A_{ik} A_{jk}^T$ 
        endfor
    endfor
endfor
    
```

Fig. 9. Algorithm-by-blocks for the Cholesky factorization. $\text{CHOL_BLK}(B)$ refers to a blocked algorithm to compute the Cholesky factorization of B . On completion, this algorithm overwrites the lower triangular part of B with its Cholesky factor. The actual implementation is similar to those in Figures 2 (right) and 4, but for conciseness we present it as loops.

SPD matrix partitioned as in (2). On completion, the lower triangular part of A is overwritten by the Cholesky factor L while the strictly upper triangular part of the matrix is not modified.

This algorithm incurs in the same flop count as the traditional implementation for the Cholesky factorization and the two exhibit the same numerical stability properties.

For details, see [Chan et al. 007a].

7. ADVANCED LAPACK-LEVEL OPERATIONS: BANDED FACTORIZATION AND INVERSION OF MATRICES

In this section, we offer a few comments on the development of algorithms-by-blocks for slightly more complex operations: factorization of a banded matrix and matrix inversion.

7.1 Cholesky factorization of banded matrices

Consider a banded SPD matrix with bandwidth k_d , partitioned into $b \times b$ blocks as in (2) so that nonzero entries only appear in the diagonal blocks A_{kk} , the subdiagonal blocks $A_{k+1,k}, \dots, A_{\min(N-1, k+D), k}$ and, given the symmetry of the matrix, $A_{k, k+1}, \dots, A_{k, \min(N-1, k+D)}$. (If we assume for simplicity that $k_d + 1$ is an exact multiple of b then $D = (k_d + 1)/b - 1$.) An algorithm-by-blocks for the Cholesky factorization of this matrix is easily obtained from the algorithm in Figure 9 by changing the upper limit of the second loop to $\min(N - 1, k + D)$.

The ideas extend to provide algorithm-by-blocks for the LU factorization with incremental pivoting and the QR factorization of a banded matrix.

One of the advantages of using FLASH in the implementation of banded algorithm-by-blocks is that storage of a band matrix does not differ from that of a dense matrix. We can still view the matrix as a matrix of matrices but store only those blocks that contain nonzero entries into the structure. Thus, FLASH easily provides compact storage schemes for banded matrices.

For further details, see [Quintana-Ortí et al. 008c].

7.2 Inversion of SPD matrices

Traditionally, the inversion of an SPD matrix A is performed as a sequence of three stages: compute the Cholesky factorization of the matrix $A = LL^T$; invert the

Cholesky factor $L \rightarrow L^{-1}$; and form $A^{-1} := L^{-T}L^{-1}$. An algorithm-by-blocks has been given above for the first stage and two more algorithms-by-blocks can be easily formulated for the second and third stages. The result is an alternative algorithm-by-blocks that yields much higher performance than one which synchronizes all computation after each stage. For further details, see [Chan et al. 2008].

The same approach provides an algorithm-by-blocks for the inversion of a general matrix via the LU factorization with incremental pivoting.

The authors in [Bientinesi et al. 2008] show it is possible to compute these three stages concurrently, and that doing so enhances load-balance on distributed-memory architectures. Since the run-time system performs the operations on blocks out-of-order, no benefit results from a one-sweep algorithm.

8. EXPERIMENTAL RESULTS

In this section, we examine various multithreaded codes in order to assess the potential performance benefits offered by algorithms-by-blocks. All experiments were performed using double-precision floating-point arithmetic on two architectures:

- A ccNUMA SGI Altix 350 server consisting of eight nodes, each with two 1.5 GHz Intel Itanium2 processors, providing a total of 16 CPUs and a peak performance of 96 GFLOPs (96×10^9 floating point operations per second). The nodes are connected via an SGI NUMALink connection ring and collectively access 32 GBytes of general-purpose physical RAM, with 2 GBytes local to each node. Performance was measured by linking to the BLAS in Intel’s Math Kernel Library (MKL) version 8.1.
- An SMP server with eight AMD Opteron processors, each one with two cores clocked at 2.2 GHz, providing a total of 16 cores and a peak performance of 70.4 GFLOPs. The cores in this platform share 64 GBytes of general-purpose physical RAM. Performance was measured by linking to the BLAS in Intel Math Kernel Library (MKL) version 9.1.

We report the performance of the following three parallel implementations in Figure 10:

- LAPACK: Routines `dpotrf` (Cholesky factorization), `dgeqrf` (QR factorization), `dgetrf` (LU factorization with partial pivoting), and `dpbtrf` (Cholesky factorization of band matrices) in LAPACK 3.0 linked to multithreaded BLAS in MKL.
- MKL: Multithreaded implementation of routines `dpotrf`, `dgeqrf`, and `dgetrf` in MKL.
- AB: Our implementation of algorithm-by-blocks, with matrices stored hierarchically using the FLASH API, scheduled with the SuperMatrix runtime system and linked to serial BLAS in MKL. The OpenMP implementation provided by the Intel C Compiler served as the underlying threading mechanism used by SuperMatrix on both platforms.

We consider the usual flop counts for the factorizations: $n^3/3$, $4n^3/3$ and $2n^3/3$, respectively, for the Cholesky, QR and LU factorizations of a square matrix of order n . The cost of the Cholesky factorization of a matrix with bandwidth k_d is

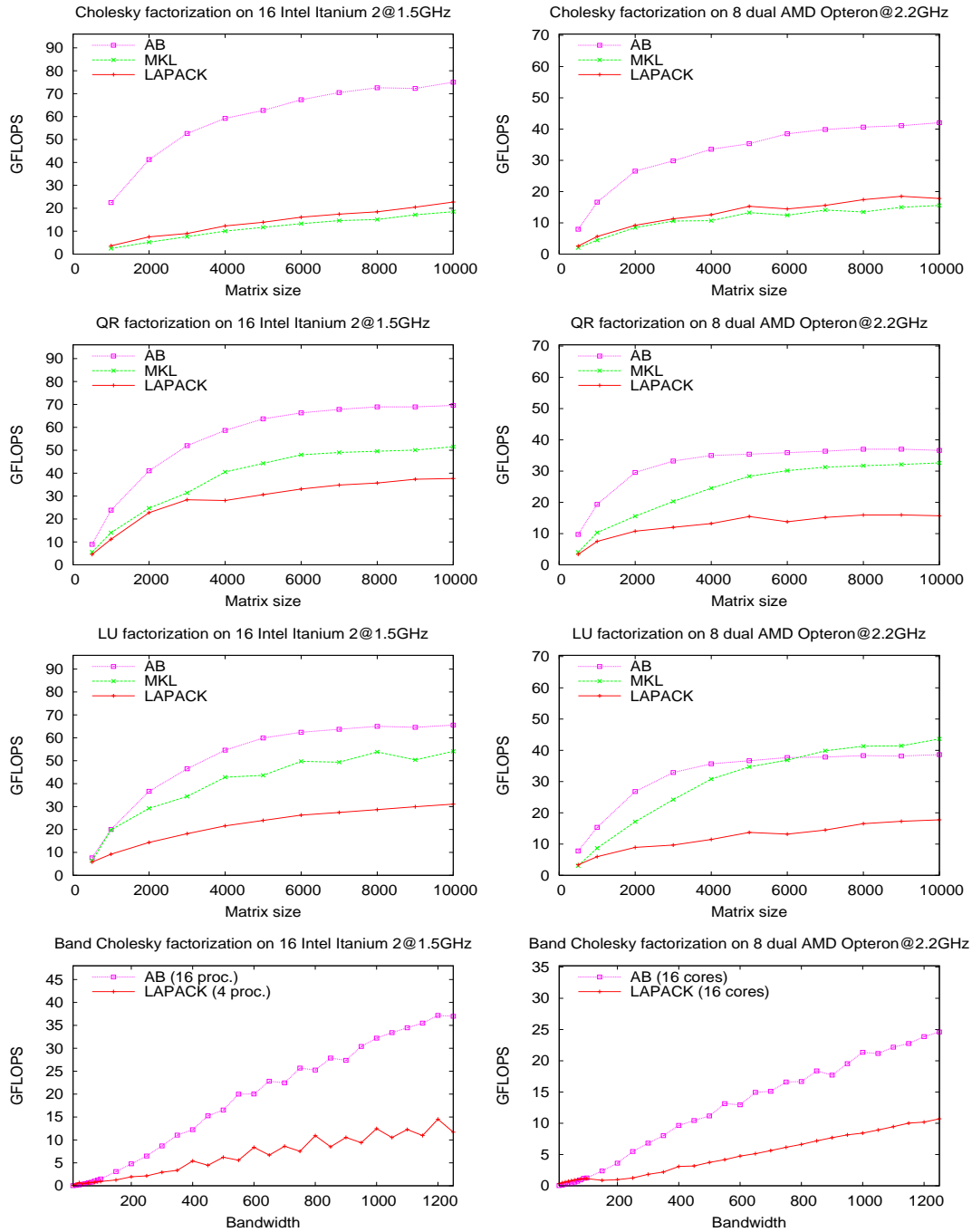


Fig. 10. Performance of the multithreaded factorization algorithms.

computed as $n(k_d^2 + 3k_d)$ flops. Note that the algorithms-by-blocks for the QR and LU factorizations actually perform a slightly higher number of flops that represent a lower order term in the overall cost.

When hand-tuning block sizes, an effort was made to determine the best values for all combinations of parallel implementations and BLAS. In the evaluation of the band factorization case, the dimension of the matrix was set to 5000. In this case, we report those results corresponding to the most favorable number of processors/cores for each implementation since using a lower number of resources in some cases resulted in a lower execution time.

The results show that algorithms-by-blocks clearly outperform the codes in LAPACK and are competitive with highly tuned implementations provided by libraries such as MKL.

An interesting question is whether on multithreaded architectures it would be appropriate to instead use libraries such as ScaLAPACK or PLAPACK, which were designed for distributed-memory parallel architectures. In [Bientinesi et al. 2008] we present evidence that, on multithreaded architectures, implementations that extract parallelism only within the BLAS outperform PLAPACK and ScaLAPACK, which use MPI [Gropp et al. 1994]. In this paper we show that extracting parallelism only within the BLAS is inferior to the proposed approach. Thus, by transitivity, the proposed approach can be expected to outperform MPI based libraries like PLAPACK and ScaLAPACK on multithreaded architectures.

9. CONCLUSION

While architectural advances promise to deliver a high level of parallelism in the form of many-core platforms, we argue that it is programmability that will determine the success of these architectures. In this paper we have illustrated how the notation, APIs, and tools that are part of the FLAME project provide modern object-oriented abstractions to increase developer productivity and user-friendliness alike in the context of dense and banded linear algebra libraries. One of these abstractions targets multi-core architectures by borrowing dynamic out-of-order scheduling techniques from sequential superscalar architectures. Results for the most significant (dense) matrix factorizations on two shared-memory parallel platforms consisting of a relatively large number of processors/cores illustrate the benefits of our approach.

The FLAME project strives to remain forward-looking. By maintaining a clean API design and clear separation of concerns, we streamline the process of taking a new algorithm from whiteboard concept to high-performance parallel implementation. The base FLAME/C API, the FLASH hierarchical matrix extension, and the SuperMatrix runtime scheduling and execution mechanism compliment each other through friendly abstractions that facilitate a striking increase in developer-level productivity as well as uncompromising end-user performance.

From the beginning, we have separated the SuperMatrix heuristic used for scheduling tasks from the library that implements the linear algebra operations. In [Chan et al. 007a] we demonstrated the benefits of using different heuristics to schedule suboperations to threads. As part of ongoing efforts, we continue to investigate the effects of different scheduling strategies on overall performance. We do not

discuss this topic in the present paper because our desire to focus the present paper squarely on the issue of programmability. Another topic for future research is how to take advantage of the FLASH API's ability to capture multiple levels of hierarchy.

Additional information

For additional information on FLAME visit

<http://www.cs.utexas.edu/users/flame/>.

Acknowledgements

This research was partially sponsored by NSF grants CCF-0540926 and CCF-0702714. Additional support came from the *J. Tinsley Oden Faculty Fellowship Research Program* of the Institute for Computational Engineering and Sciences (ICES) at UT-Austin. Gregorio Quintana-Ortí and Enrique S. Quintana-Ortí were supported by the CICYT project TIN2005-09037-C02-02 and FEDER, and projects P1B-2007-19 and P1B-2007-32 of the *Fundación Caixa-Castellón/Bancaixa* and UJI.

We thank the other members of the FLAME team for their support. We thank John Gilbert and Vikram Aggarwal from the University of California at Santa Barbara for granting access to the NEUMANN platform.

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation (NSF).

REFERENCES

- ADDISON, C., REN, Y., AND VAN WAVEREN, M. 2003. OpenMP issues arising in the development of parallel BLAS and LAPACK libraries. *Scientific Programming* 11, 2 (April), 95–104.
- AGARWAL, R. C. AND GUSTAVSON, F. G. 1989. Vector and parallel algorithms for Cholesky factorization on IBM 3090. In *SC '89: Proceedings of the 1989 ACM/IEEE Conference on Supercomputing*. New York, NY, USA, 225–233.
- ANDERSON, E., BAI, Z., BISCHOF, C., BLACKFORD, L. S., DEMMEL, J., DONGARRA, J. J., CROZ, J. D., HAMMARLING, S., GREENBAUM, A., MCKENNEY, A., AND SORENSEN, D. 1999. *LAPACK Users' Guide*, 3rd ed. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA.
- BALAY, S., BUSCHELMAN, K., EIJKHOUT, V., GROPP, W. D., KAUSHIK, D., KNEPLEY, M. G., MCINNES, L. C., SMITH, B. F., AND ZHANG, H. 2004. PETSc users manual. Tech. Rep. ANL-95/11 - Revision 2.1.5, Argonne National Laboratory.
- BIENTINESI, P., GUNTER, B., AND DE GEIJN, R. A. V. 2008. Families of algorithms related to the inversion of a symmetric positive definite matrix. *ACM Transactions on Mathematical Software* 35, 1 (July), 1–22.
- BIENTINESI, P., QUINTANA-ORTÍ, E. S., AND VAN DE GEIJN, R. A. 2005. Representing linear algebra algorithms in code: The FLAME application programming interfaces. *ACM Transactions on Mathematical Software* 31, 1 (March), 27–59.
- BUTTARI, A., LANGOU, J., KURZAK, J., , AND DONGARRA, J. 2007. A class of parallel tiled linear algebra algorithms for multicore architectures. LAPACK Working Note 191 UT-CS-07-600, University of Tennessee. September.
- BUTTARI, A., LANGOU, J., KURZAK, J., AND DONGARRA, J. 2008. Parallel tiled QR factorization for multicore architectures. *Concurrency and Computation: Practice and Experience* 20, 13 (September), 1573–1590.
- CHAN, E., QUINTANA-ORTÍ, E. S., QUINTANA-ORTÍ, G., AND VAN DE GEIJN, R. 2007a. Super-Matrix out-of-order scheduling of matrix operations for SMP and multi-core architectures. In *ACM Transactions on Mathematical Software*, Vol. V, No. N, Month 20YY.

- SPAA '07: Proceedings of the Nineteenth ACM Symposium on Parallelism in Algorithms and Architectures*. San Diego, CA, USA, 116–125.
- CHAN, E., VAN ZEE, F. G., BIENTINESI, P., QUINTANA-ORTÍ, E. S., QUINTANA-ORTÍ, G., AND VAN DE GEIJN, R. 2008. SuperMatrix: A multithreaded runtime scheduling system for algorithms-by-blocks. In *PPoPP '08: Proceedings of the Thirteenth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. Salt Lake City, UT, USA, 123–132.
- CHAN, E., VAN ZEE, F. G., QUINTANA-ORTÍ, E. S., QUINTANA-ORTÍ, G., AND VAN DE GEIJN, R. 2007b. Satisfying your dependencies with SuperMatrix. In *Cluster '07: Proceedings of the 2007 IEEE International Conference on Cluster Computing*. Austin, TX, USA, 91–99.
- CHOI, J., DONGARRA, J. J., POZO, R., AND WALKER, D. W. 1992. ScaLAPACK: A scalable linear algebra library for distributed memory concurrent computers. In *Proceedings of the Fourth Symposium on the Frontiers of Massively Parallel Computation*. McLean, VA, USA, 120–127.
- COLLINS, T. AND BROWNE, J. C. 1995. Matrix++: An object-oriented environment for parallel high-performance matrix computations. In *HICSS '95: Proceedings of the 28th Annual Hawaii International Conference on System Sciences*. Maui, HI, USA, 202–211.
- DONGARRA, J. J., BUNCH, J. R., MOLER, C. B., AND STEWART, G. W. 1979. *LINPACK Users' Guide*. SIAM, Philadelphia, PA, USA.
- DONGARRA, J. J., DU CROZ, J., HAMMARLING, S., AND DUFF, I. 1990. A set of level 3 basic linear algebra subprograms. *ACM Transactions on Mathematical Software* 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 subprograms. *ACM Transactions on Mathematical Software* 14, 1 (March), 1–17.
- EDWARDS, H. C. AND VAN DE GEIJN, R. A. 2006. On application interfaces to parallel dense matrix libraries: Just let me solve my problem! FLAME Working Note #18 TR-2006-15, The University of Texas at Austin, Department of Computer Sciences. February.
- ELMROTH, E., GUSTAVSON, F., JONSSON, I., AND KAGSTROM, B. 2004. Recursive blocked algorithms and hybrid data structures for dense matrix library software. *SIAM Review* 46, 1, 3–45.
- GOLUB, G. H. AND LOAN, C. F. V. 1996. *Matrix Computations*, 3rd ed. The Johns Hopkins University Press, Baltimore, MD, USA.
- GOTO, K. AND VAN DE GEIJN, R. A. 2008. Anatomy of a high-performance matrix multiplication. *ACM Transactions on Mathematical Software* 34, 3 (May), 1–25.
- GROPP, W., LUSK, E., AND SKJELLUM, A. 1994. *Using MPI*. The MIT Press.
- GUNNELS, J. A., GUSTAVSON, F. G., HENRY, G. M., AND VAN DE GEIJN, R. A. 2001. FLAME: Formal linear algebra methods environment. *ACM Transactions on Mathematical Software* 27, 4 (December), 422–455.
- GUNTER, B. C. AND VAN DE GEIJN, R. A. 2005. Parallel out-of-core computation and updating the QR factorization. *ACM Transactions on Mathematical Software* 31, 1 (March), 60–78.
- GUO, J., BIKSHANDI, G., FRAGUELA, B., GARZARAN, M., AND PADUA, D. 2008. Programming with tiles. In *PPoPP '08: Proceedings of the Thirteenth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. Salt Lake City, UT, USA, 111–122.
- GUSTAVSON, F. G., KARLSSON, L., AND KAGSTROM, B. 2007. Three algorithms for Cholesky factorization on distributed memory using packed storage. In *PARA '07: Proceedings of the Workshop on State-of-the-Art in Scientific Computing*. Volume 4699 in Lecture Notes in Computer Science. Springer Berlin / Heidelberg, 550–559.
- HERRERO, J. R. 2006. A framework for efficient execution of matrix computations. Ph.D. thesis, Polytechnic University of Catalonia, Spain.
- JOFFRAIN, T., QUINTANA-ORTÍ, E. S., AND VAN DE GEIJN, R. A. 2004. Rapid development of high-performance out-of-core solvers. In *PARA '04: Proceedings of the Workshop on State-of-the-Art in Scientific Computing*. Volume 3732 in Lecture Notes in Computer Science. Springer Berlin / Heidelberg, 413–422.
- KURZAK, J. AND DONGARRA, J. 2006. Implementing linear algebra routines on multi-core processors with pipelining and a look ahead. LAPACK Working Note 178 UT-CS-06-581, University of Tennessee. September.

- LAWSON, C. L., HANSON, R. J., KINCAID, D. R., AND KROGH, F. T. 1979. Basic linear algebra subprograms for Fortran usage. *ACM Transactions on Mathematical Software* 5, 3 (September), 308–323.
- LOW, T. M. AND VAN DE GEIJN, R. 2004. An API for manipulating matrices stored by blocks. FLAME Working Note #12 TR-2004-15, The University of Texas at Austin, Department of Computer Sciences. May.
- LU, H., COX, A. L., DWARKADAS, S., RAJAMONY, R., AND ZWAENEPOEL, W. 1997. Compiler and software distributed shared memory support for irregular applications. In *PPOPP '97: Proceedings of the Sixth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. Las Vegas, NV, USA, 48–56.
- MARKER, B. A., VAN ZEE, F. G., GOTO, K., QUINTANA-ORTÍ, G., AND VAN DE GEIJN, R. A. 2007. Toward scalable matrix multiply on multithreaded architectures. In *Euro-Par '07: Proceedings of the Thirteenth International European Conference on Parallel and Distributed Computing*. Rennes, France, 748–757.
- NIEPLOCHA, J., HARRISON, R., AND LITTLEFIELD, R. 1996. Global arrays: A nonuniform memory access programming model for high-performance computers. *The Journal of Supercomputing* 10, 2 (June), 197–220.
- QUINTANA-ORTÍ, E. S. AND VAN DE GEIJN, R. 2009. Updating an LU factorization with pivoting. *ACM Transactions on Mathematical Software*. To appear.
- QUINTANA-ORTÍ, G., QUINTANA-ORTÍ, E. S., CHAN, E., VAN DE GEIJN, R., AND VAN ZEE, F. G. 2008a. Design of scalable dense linear algebra libraries for multithreaded architectures: the LU factorization. In *MTAAP '08: Proceedings of the 2008 Workshop on Multithreaded Architectures and Applications*. Miami, FL, USA, 1–8.
- QUINTANA-ORTÍ, G., QUINTANA-ORTÍ, E. S., CHAN, E., VAN ZEE, F. G., AND VAN DE GEIJN, R. A. 2008b. Scheduling of QR factorization algorithms on SMP and multi-core architectures. In *PDP '08: Proceedings of the Sixteenth Euromicro International Conference on Parallel, Distributed and network-based Processing*. Toulouse, France, 301–307.
- QUINTANA-ORTÍ, G., QUINTANA-ORTÍ, E. S., REMÓN, A., AND VAN DE GEIJN, R. 2008c. An algorithm-by-blocks for SuperMatrix band Cholesky factorization. In *VECPAR '08: Proceedings of the Eighth International Meeting on High Performance Computing for Computational Science*. Toulouse, France, 1–13.
- SKAGESTEIN, G. 1972. Rekursiv unterteilte matrisen sowie methoden zur erstellung von rechnerprogrammen fur ihre verarbeitung. Ph.D. thesis, Universitat Stuttgart, Germany.
- SORENSEN, D. C. 1985. Analysis of pairwise pivoting in Gaussian elimination. *IEEE Transactions on Computers* 34, 3 (March), 274–278.
- STRAZDINS, P. 2001. A comparison of lookahead and algorithmic blocking techniques for parallel matrix factorization. *International Journal of Parallel and Distributed Systems and Networks* 4, 1 (June), 26–35.
- TOLEDO, S. 1999. A survey of out-of-core algorithms in numerical linear algebra. In *External Memory Algorithms*, J. Abello and J. S. Vitter, Eds. American Mathematical Society, Boston, MA, USA, 161–179.
- VALSALAM, V. AND SKJELLUM, A. 2002. A framework for high-performance matrix multiplication based on hierarchical abstractions, algorithms and optimized low-level kernels. *Concurrency and Computation: Practice and Experience* 14, 10 (August), 805–840.
- VAN DE GEIJN, R. A. 1997. *Using PLAPACK: Parallel Linear Algebra Package*. The MIT Press.
- VAN DE GEIJN, R. A. AND QUINTANA-ORTÍ, E. S. 2008. *The Science of Programming Matrix Computations*. www.lulu.com.
- VON HANXLEDEN, R., KENNEDY, K., KOELBEL, C. H., DAS, R., AND SALTZ, J. H. 1992. Compiler analysis for irregular problems in Fortran D. In *LCPC '92: Proceedings of the Fifth Workshop on Languages and Compilers for Parallel Computing*. Volume 757 in Lecture Notes in Computer Science. Springer Berlin / Heidelberg, 97–111.
- WISE, D. S., FRENS, J. D., GU, Y., AND ALEXANDER, G. A. 2001. Language support for Morton-order matrices. In *PPoPP '01: Proceedings of the Eighth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. Snowbird, UT, USA, 24–33.

Received Month Year; revised Month Year; accepted Month Year