

# Algorithmic Redistribution Methods for Block-Cyclic Decompositions

Antoine P. Petitet and Jack J. Dongarra, *Fellow, IEEE*

**Abstract**—This article presents various data redistribution methods for block-partitioned linear algebra algorithms operating on dense matrices that are distributed in a block-cyclic fashion. Because the algorithmic partitioning unit and the distribution blocking factor are most often chosen to be equal, severe alignment restrictions are induced on the operands, and optimal values with respect to performance are architecture dependent. The techniques presented in this paper redistribute data “on the fly,” so that the user’s data distribution blocking factor becomes independent from the architecture dependent algorithmic partitioning. These techniques are applied to the matrix-matrix multiplication operation. A performance analysis along with experimental results shows that alignment restrictions can then be removed and that high performance can be maintained across platforms independently from the user’s data distribution blocking factor.

**Index Terms**—Algorithmic blocking, redistribution, block-cyclic decomposition.

## 1 INTRODUCTION

IN a serial computational environment, *transportable efficiency* is the essential motivation for developing blocking strategies and block-partitioned algorithms [5], [19], [32]. An algorithmic blocking factor adjusts the granularity of the subtasks to maximize the utilization of hardware resources. In a distributed-memory environment, *load balance* is the essential motivation for distributing the data over a collection of processes according to the block-cyclic decomposition scheme [1], [12], [29], [21], [11], [17], [22]. A distribution blocking factor is used to partition an array into blocks that are then mapped onto the processes. Optimal values of these algorithmic and distribution blocking factors often differ for a given algorithm and target architecture. Despite this fact, most of the parallel algorithms proposed in the literature assume the values of these blocking factors to be identical [16], [17], [40], [47]. This assumption simplifies the expression of such algorithms, but limits their flexibility and ease of use. The application’s scope of these algorithms is thus limited in a way that does not satisfy general purpose library requirements. High performance is nevertheless achievable on a wide range of distributed-memory concurrent computers, but depends on the chosen value of the distribution and algorithmic blocking factors. General redistribution packages [33], [42] can alleviate this constraint at a possibly high memory cost.

This paper presents and discusses various methods for redistributing the data into an appropriate form as the parallel algorithm progresses, so that the algorithmic and distribution blocking factors can be selected independently. The data distribution blocking factor is chosen by the user, and the selection of an efficient and machine specific

algorithmic blocking factor is left to the software library designer. In this paper, we focus on elementary blocked linear algebra computations and assume that the user’s data is distributed onto the processes according to the general block-cyclic decomposition scheme. We show that such redistribution methods alleviate the alignment restrictions imposed by the block-cyclic data layout. We analyze the efficiency and scalability of these methods and illustrate our study by presenting experimental results on two distributed-memory concurrent computers.

Algorithmic redistribution methods attempt to reorganize logically the computations and communications within an algorithmic context. In order to derive such methods, some properties of the block-cyclic data distribution are first exhibited in Section 2. Various algorithmic redistribution methods are then presented and applied to the representative outer product matrix-matrix multiply algorithm in Section 3. These techniques are presented within a single framework, making them suitable for their integration into a software library. For some of these strategies, little is known in terms of their impact on efficiency and/or ease of modular implementation. To our knowledge, few practical experiments have thus far been reported in the literature. A scalability analysis of these algorithmic redistribution methods as well as a number of experimental performance results are finally presented in Section 4.

## 2 PROPERTIES OF THE BLOCK-CYCLIC DATA DISTRIBUTION

Since the data decomposition largely determines the performance and scalability of a concurrent algorithm, a great deal of research [21], [24], [26], [30], [11] has focused on studying various data decompositions [6], [12], [31]. As a result, the two-dimensional block-cyclic distribution [36] has been suggested as a possible general purpose basic decomposition for parallel dense linear algebra software libraries [18], [29], [38]. This section presents important

• The authors are with the Department of Computer Science, University of Tennessee, Knoxville, TN 37996. E-mail: {petitet, dongarra}@cs.utk.edu.

Manuscript received 12 Mar. 1998; revised 01 June 1999.

For information on obtaining reprints of this article, please send e-mail to: tpd@computer.org, and reference IEEECS Log Number 106471.

properties of the two-dimensional, block-cyclic data distribution. These properties are the basis of efficient algorithms for address generation, fast indexing techniques, and communication scheduling.

## 2.1 Analytical Definition of the Block-Cyclic Data Distribution

In general, there may be several processes executed by one processor; therefore, without loss of generality, the underlying concurrent computer is regarded as a collection of *processes*, rather than physical processors. Consider a  $P \times Q$  grid of processes, where  $\Gamma$  denotes the set of all process coordinates  $(p, q)$  in this grid:

$$\Gamma = \{(p, q) \in \{0 \dots P-1\} \times \{0 \dots Q-1\}\}.$$

Consider an  $M \times N$  matrix partitioned into blocks of size  $r \times s$ . Each matrix entry  $a_{ij}$  is uniquely identified by the integer pair  $(i, j)$  of its row and column indexes. Let  $\Delta$  be the set constructed from all these pairs:

$$\begin{aligned} \Delta &= \{(i, j) \in \{0 \dots M-1\} \times \{0 \dots N-1\}\} \\ &= \{((lP+p)r+x, (mQ+q)s+y), \\ &\quad ((p, q), (l, m), (x, y)) \in \Gamma \times \Lambda \times \Theta\}. \end{aligned}$$

with

$$\Lambda = \{(l, m) \in \{0 \dots \lfloor \frac{M-1}{P} \rfloor\} \times \{0 \dots \lfloor \frac{N-1}{Q} \rfloor\}\}$$

and

$$\Theta = \{(x, y) \in \{0 \dots r-1\} \times \{0 \dots s-1\}\}.$$

**Definition 2.1.** *The block-cyclic distribution is defined by the three following mappings associating to a matrix entry index pair  $(i, j)$ :*

- the coordinates  $(p, q)$  of the process into which the matrix entry resides

$$\begin{cases} \Delta \longrightarrow \Gamma \\ (i, j) \longmapsto (p, q). \end{cases} \quad (1)$$

- the coordinates  $(l, m)$  of the local block in which the matrix entry resides

$$\begin{cases} \Delta \longrightarrow \Lambda \\ (i, j) \longmapsto (l, m). \end{cases} \quad (2)$$

- the local row and column offsets  $(x, y)$  within this local block  $(l, m)$

$$\begin{cases} \Delta \longrightarrow \Theta \\ (i, j) \longmapsto (x, y). \end{cases} \quad (3)$$

A few particular occurrences of the above definition are worth mentioning. First, **the blocked distribution** is determined by Definition 2.1 with  $r = \lceil \frac{M}{P} \rceil$  and  $s = \lceil \frac{N}{Q} \rceil$ , i.e.,  $\Lambda = \{(0, 0)\}$ . Second, **the square block-cyclic distribution** is a special case of Definition 2.1 with  $r = s$ . Finally, **the purely scattered or cyclic decomposition** is a particular

instance of the square block-cyclic distribution with  $r = s = 1$ , i.e.,  $\Theta = \{(0, 0)\}$ .

## 2.2 Properties of the Block-Cyclic Data Distribution and LCM Tables

The properties of the block-cyclic data distribution presented below are centered around the concept of the *k*-diagonal of a matrix, which is defined as the set of entries  $a_{ij}$  such that  $i - j = k$ . Indeed, if we consider the matrix induced by a column and a row vector, the one-dimensional redistribution operation can be formulated as first; projecting the source column vector entries onto the 0-diagonal of this matrix and second; projecting this 0-diagonal onto the target row vector. The complexity of such a redistribution operation is thus closely related to the 0-diagonal of the induced matrix. This suggests the following definition:

**Definition 2.** *Given a k-diagonal, the k-LCM table (LCMT) is a two-dimensional infinite array of integers local to each process  $(p, q)$  defined by*

$$LCMT_{l,m}^{p,q} = (mQ+q)s - (lP+p)r + k,$$

for  $(l, m)$  in  $\mathbb{N}^2$ .

The equation for the *k*-diagonal is thus given by

$$LCMT_{l,m}^{p,q} = x - y, \quad (4)$$

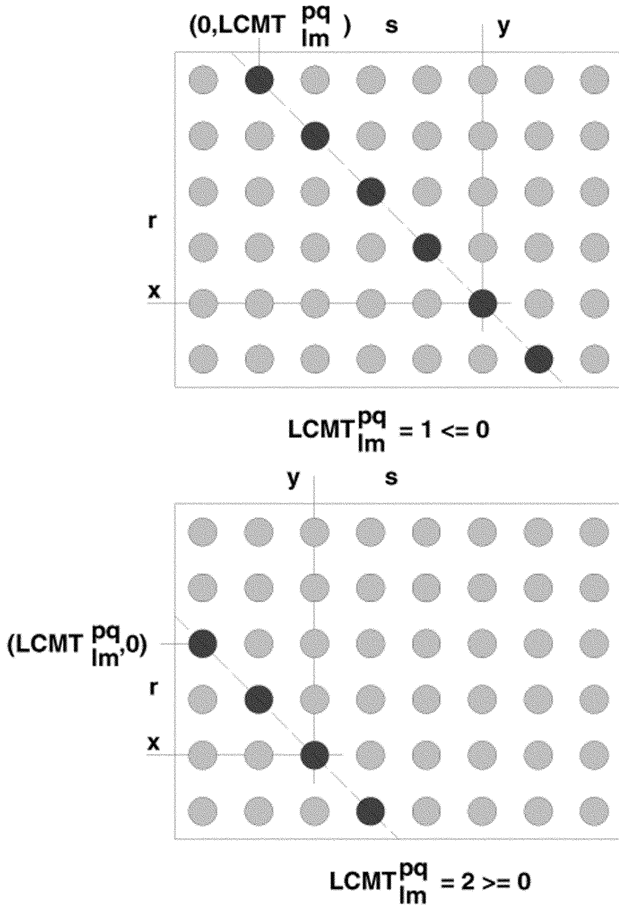
with  $(x, y)$  in  $\Theta$ . Blocks owning the *k*-diagonal entries are therefore such that

$$1 - s \leq LCMT_{l,m}^{p,q} \leq r - 1. \quad (5)$$

In addition the value of  $LCMT_{l,m}^{p,q}$  specifies where the diagonal starts within a block owning diagonals as illustrated in Fig. 1.

It follows from Definition 2 that the local blocks in process  $(p, q)$  such that  $LCMT_{l,m}^{p,q} \leq 0$  own matrix entries  $a_{ij}$  that are globally *below* the *k*-diagonal. The local blocks such that  $LCMT_{l,m}^{p,q} \leq -s$  correspond globally to *strictly lower* blocks of the matrix. Moreover, within each process, if the  $r \times s$  block of local coordinates  $(l, m)$  owns *k*-diagonal entries, the block of local coordinates  $(l+1, m)$  owns either *k*-diagonals or matrix entries that are *strictly below* the *k*-diagonal. Furthermore, within each process, if the  $r \times s$  blocks of local coordinates  $(l, m)$  and  $(l+1, m)$  own *k*-diagonals, then the block of local coordinates  $(l, m+1)$  owns matrix entries that are *strictly above* the *k*-diagonal.

Similarly, the local blocks in process  $(p, q)$  such that  $LCMT_{l,m}^{p,q} \geq 0$  own matrix entries  $a_{ij}$  that are globally *above* the *k*-diagonal. The local blocks in process  $(p, q)$  such that  $LCMT_{l,m}^{p,q} \geq r$  correspond globally to *strictly upper* blocks of the matrix. Moreover, within each process, if the  $r \times s$  block of local coordinates  $(l, m)$  owns *k*-diagonal entries, the block of local coordinates  $(l, m+1)$  owns either *k*-diagonals or matrix entries that are *strictly above* the *k*-diagonal. Furthermore, within each process, if the  $r \times s$  blocks of local coordinates  $(l, m)$  and  $(l, m+1)$  own *k*-diagonals, then the block of local coordinates  $(l+1, m)$  owns matrix entries that are *strictly below* the *k*-diagonal.


 Fig. 1. Meaning of different values of  $LCMT_{lm}^{pq}$  with  $r = 6$ ,  $s = 8$ .

Let  $L = \text{lcm}(Pr, Qs)$  and  $g = \text{gcd}(Pr, Qs)$  be respectively the least common multiple and greatest common divisor of the quantities  $Pr$  and  $Qs$ . The index pairs  $(i, j)$ ,  $(i + L, j)$ ,  $(i, j + L)$ , and  $(i + L, j + L)$  refer to array entries that are residing in the same process  $(p, q)$ . The distribution pattern has therefore a periodicity of length  $L$  along the  $k$ -diagonal. A block of order  $L$  is referred to as an LCM-block in what follows. Each process owns exactly  $L/P \times L/Q$  entries of such an LCM-block. This larger partitioning unit has been originally introduced in the restricted context of square block-cyclic mappings in [15], [16]. The meaningful part of the LCM tables to be considered in each process is thus of size  $L/(Pr) \times L/(Qs)$ .

Fig. 2 shows an LCM-block for a given set of distribution parameters  $P$ ,  $r$ ,  $Q$ , and  $s$ . The corresponding 1-LCM tables are shown in Fig. 3. Each of these tables is associated to a distinct process of coordinates  $(p, q)$  as indicated in the upper left corner of each table. Examine, for example, the table corresponding to process  $(0, 0)$ . The value of the LCM table entry  $(0, 0)$  is 1. Since this value is greater than  $-s = -3$ , and less than  $r = 2$ , it follows that this block  $(0, 0)$  owns 1-diagonals. Moreover, within this block this diagonal starts in position  $(LCMT_{00}^{00}, 0) = (1, 0)$ . The distribution periodicity within this process is illustrated by the block of local coordinates  $(3, 2)$  which is such that

	0	1	2	3	4	5	6	7	8	9	10	11
0	○	○	○	○	○	○	○	○	○	○	○	○
1	●	○	○	○	○	○	○	○	○	○	○	○
2	○	●	○	○	○	○	○	○	○	○	○	○
3	○	○	○	●	○	○	○	○	○	○	○	○
4	○	○	○	○	○	○	○	○	○	○	○	○
5	○	○	○	○	○	○	○	○	○	○	○	○
6	○	○	○	○	○	○	○	○	○	○	○	○
7	○	○	○	○	○	○	○	○	○	○	○	○
8	○	○	○	○	○	○	○	○	○	○	○	○
9	○	○	○	○	○	○	○	○	○	○	○	○
10	○	○	○	○	○	○	○	○	○	○	○	○
11	○	○	○	○	○	○	○	○	○	○	○	○
12	○	○	○	○	○	○	○	○	○	○	○	○
13	○	○	○	○	○	○	○	○	○	○	○	○

 Fig. 2. LCM-block for  $P = 2$ ,  $r = 2$ ,  $Q = 2$ , and  $s = 3$ .

$$LCMT_{00}^{00} = LCMT_{32}^{00} = 1.$$

One can also verify that a block of local coordinates  $(l, m)$  in this table corresponds to a strictly lower (respectively upper) block in the original LCM-block if and only if  $LCMT_{lm}^{00} \leq -s$  (respectively  $LCMT_{lm}^{00} \geq r$ ). These same remarks apply to all of the other LCM tables shown in Fig. 3.

**Property 1.** The number of  $r \times s$  blocks owning  $L$  consecutive  $k$ -diagonals is given by

$$\begin{cases} \frac{L(r+s-\text{gcd}(r,s))}{rs} & \text{if } \text{gcd}(r,s) \text{ divides } k, \\ \frac{L(r+s)}{rs} & \text{otherwise.} \end{cases}$$

**Proof.** Because of Inequality 5, one can assume that  $0 \leq k < \text{gcd}(r, s)$  by simply renumbering the processes with appropriate relative coordinates. Consider, then, a square array of size  $\text{lcm}(r, s)$ . If  $\text{gcd}(r, s)$  divides  $k$  (i.e.,  $k$  is zero), there is exactly one  $r \times s$  subblock such that its  $(r-1, s-1)$  entry belongs to the  $k$ -diagonal: there is obviously one, namely the right bottom subblock; there is exactly one by definition of the least common multiple of  $r$  and  $s$ . If  $\text{gcd}(r, s)$  does not divide  $k$ , such a subblock does not exist: if it did, one could have chosen  $k$  to be zero, by considering the next block owning  $k$ -diagonals, which is a contradiction. In other words, if  $\text{gcd}(r, s)$  does not divide  $k$ , the  $k$ -diagonal never cut the right lower corner of a  $r \times s$  subblock; otherwise, the  $k$ -diagonal cut the right lower corner of exactly one  $r \times s$  subblock.

In this  $\text{lcm}(r, s)$  array, the column edges of the  $r \times s$  subblocks will be cut exactly  $\text{lcm}(r, s)/s$ , that is  $r/\text{gcd}(r, s)$  times by the  $k$ -diagonal. Similarly, the row edges of the  $r \times s$  subblocks will be cut exactly

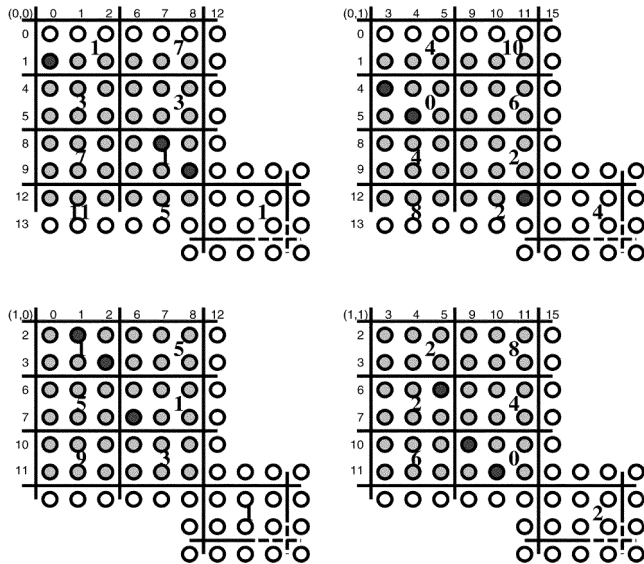


Fig. 3. 1-LCM tables for  $P = 2$ ,  $r = 2$ ,  $Q = 2$ , and  $s = 3$ .

$\text{lcm}(r, s)/r$ , that is  $s/\text{gcd}(r, s)$  times by the  $k$ -diagonal. Therefore, when  $\text{gcd}(r, s)$  divides  $k$ , there are exactly  $(r + s)/\text{gcd}(r, s) - 1$  blocks of size  $r \times s$  owning  $\text{lcm}(r, s)$  consecutive  $k$ -diagonals; otherwise, there are exactly  $(r + s)/\text{gcd}(r, s)$  such blocks. Finally, there are exactly  $L/\text{lcm}(r, s)$  such blocks in an LCM block. To see that  $L/\text{lcm}(r, s)$  is indeed an integer, one may observe that this quantity can be rewritten as  $((uQ)Pr + (tP)Qs)/g$  with  $u$  and  $t$  in  $\mathbb{Z}$ .  $\square$

Property 1 quantifies the complexity of general redistribution operations as a function of the distribution parameters, namely the perimeter  $r + s$  of the  $r \times s$  partitioning unit and the quantities  $\text{gcd}(r, s)$  and  $g = \text{gcd}(Pr, Qs)$ .

**Property 2.** A necessary and sufficient condition for every process to own  $k$ -diagonal entries is given by  $r + s - \text{gcd}(r, s) \geq g$  if  $\text{gcd}(r, s)$  divides  $k$  and  $r + s \geq g$  otherwise.

**Proof.** The condition is sufficient: remark that  $\text{gcd}(r, s)$  divides  $g$ . If  $\text{gcd}(r, s)$  divides  $k$  (note that this will always be the case if  $\text{gcd}(r, s) = 1$ ), then  $\frac{r + s}{\text{gcd}(r, s)} - 1$  is the number of multiples of  $\text{gcd}(r, s)$  in the interval

$$I_{p,q} = (pr - (q - 1)s \dots (p + 1)r - qs).$$

The number of multiples of  $g$  in the interval  $I_{p,q}$  is  $\frac{g}{\text{gcd}(r, s)}$ . Thus, the inequality

$$\frac{r + s}{\text{gcd}(r, s)} - 1 \geq \frac{g}{\text{gcd}(r, s)}$$

is a sufficient condition for a multiple of  $g$  to be in this interval  $I_{p,q}$ . Otherwise, when  $\text{gcd}(r, s)$  does not divide  $k$ , (5) can be rewritten as

$$pr - qs - s < mQs - lPr + k < pr - qs + r. \quad (6)$$

For any given process of coordinates  $(p, q)$ , there must exist a  $t \in \mathbb{Z}$  such that  $mQs - lPr = tg$  verifying Inequality 6. Moreover, the interval of interest  $I_{p,q}$  is of length  $r + s - 1$ . A sufficient condition for all processes to have  $k$ -diagonals is given by  $r + s - 1 \geq g$ . Since  $\text{gcd}(r, s) \neq 1$  and  $\text{gcd}(r, s)$  divides  $g$ , this sufficient condition can be equivalently written as  $r + s \geq g$ .

The condition is necessary: suppose there exists a process  $(p, q)$  having two distinct blocks owning  $k$ -diagonals. Then,  $r + s - \text{gcd}(r, s) \geq g$  if  $\text{gcd}(r, s)$  divides  $k$ , and  $r + s \geq g$  otherwise. There are two multiples of  $g$  in some interval  $I_{p,q}$ . Otherwise, each process owns at most one  $r \times s$  block in which  $k$ -diagonals reside. Therefore, the number of blocks owning  $k$ -diagonals is equal to the number of processes owning these diagonals. The result then follows from Property 1.  $\square$

**Property 3.** The number of processes owning  $k$ -diagonal entries is given by

$$PQ \max\left(\frac{r + s - \text{gcd}(r, s)}{g}, 1\right)$$

if  $\text{gcd}(r, s)$  divides  $k$ , and by

$$PQ \max\left(\frac{r + s}{g}, 1\right)$$

otherwise.

**Proof.** The result follows from the fact that  $Lg = (Pr)(Qs)$  and Properties 1 and 2.  $\square$

Properties 2 and 3 say that the distribution of the  $k$ -diagonals essentially depends on the perimeter of the  $r \times s$  partitioning unit as opposed to its shape. In other words, restricting the data decomposition to a square block-cyclic mapping does not affect in any way the problem of locating the  $k$ -diagonals, and consequently the complexity of redistribution operations. To reduce this complexity, it is necessary to choose small values of the distribution parameters  $r$  and  $s$ .

The end of this section aims at determining the probability that the quantities  $r + s - \text{gcd}(r, s)$  or  $r + s$  are greater or equal to  $g$ , that is, the probability that every process owns  $k$ -diagonal entries. It is difficult to analytically compute the probability that all processes will own  $k$ -diagonal entries. However, it is likely that this probability, if it exists, converges rapidly [41]. It is possible to rely on a computer to enumerate all 4-tuples in a finite and practical range such that the quantities  $r + s - \text{gcd}(r, s)$  or  $r + s$  are greater or equal to  $g$ . The results are presented in Fig. 4. In practice, i.e., for a finite range of values ( $1 \leq P, r, Q, s \leq n$ ), there is almost no difference between the finite ratios of all 4-tuples verifying these inequalities over  $[1..n]^4$ .

Fig. 4 does not prove the existence of the limit, and therefore, of the probability. However, if it exists, its value is very close to one. In other words, if one picks random distribution parameters, it is very likely that all processes in the grid will own  $k$ -diagonals. Fig. 4 also shows that the

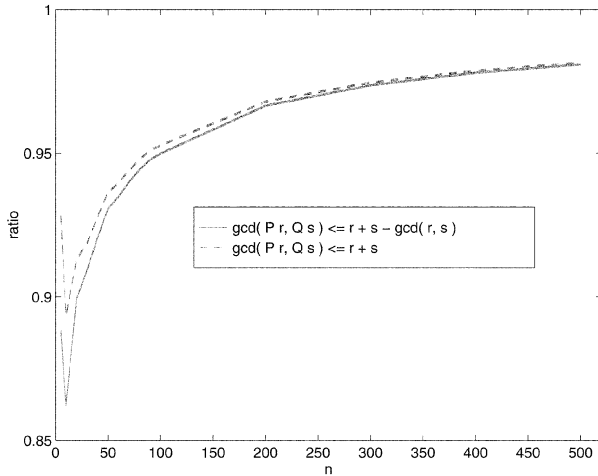


Fig. 4. Ratios of tuples  $(P, r, Q, s)$  in  $[1..n]^4$  such that  $r + s - \gcd(r, s) \geq g$  or  $r + s \geq g$ .

ratios of distribution parameters such that  $k$ -diagonals are evenly distributed tends towards one. More interesting is the fact that this function  $R$  increases very rapidly ( $R(10) \approx .88$ ,  $R(20) \approx .90$ ,  $R(50) \approx .93$ ). Therefore, it is very likely that all processes in the grid will own  $k$ -diagonals.

### 3 ALGORITHMIC REDISTRIBUTION METHODS

This section presents different kinds of blocking strategies for distributed-memory hierarchies. Most of these algorithmic redistributed operations can be expressed in terms of locating diagonals of a distributed matrix. All of these techniques are therefore presented within a single framework. For some of these strategies, little is known in terms of their impact on efficiency and/or ease of modular implementation. To our knowledge, few practical experiments have thus far been reported in the literature.

The same example operation called a rank- $K$  update is used to illustrate the differences between all blocking strategies presented below. This operation produces an  $M \times N$  matrix  $C$  by adding to itself the product of an  $M \times K$  matrix  $A$  and a  $K \times N$  matrix  $B$ :

$$C \leftarrow C + AB.$$

The *physical blocking* strategy uses the distribution blocking factors as a unit for the computational blocks. In other words, the computations are partitioned accordingly to the data distribution parameters. No attempts are made to either gather rows or columns residing in distinct processes, or scatter rows or columns residing in a single process row or column. It is assumed that the distribution parameters have been determined a priori, presumably by the user. Ideally, this choice has been influenced by its performance implications. Most of the parallel algorithms presented in the literature [2], [6], [9], [10], [17], [18], [22], [25], [26], [27], [31], [37], [48] rely on this strategy. The algorithm performing the rank- $K$  update operation using a physical blocking strategy is relatively easy to express.

Strong alignment and distribution assumptions are made on the matrix operands. In particular, the distribution blocking factors used to decompose the columns of  $A$  and the rows of  $B$  must be equal. Moreover, the rows of  $A$  (respectively the columns of  $B$ ) must be aligned to the rows (respectively columns) of  $C$ . The pseudocode for this algorithm is given below.

```

for  $kk = 1, K, NB_{dis}$ 
   $kb = \min(K - kk + 1, NB_{dis})$ 
  Broadcast  $\hat{A} = A(:, kk : kk + kb - 1)$ 
  within process rows;
  Broadcast  $\hat{B} = B(kk : kk + kb - 1, :)$ 
  within process columns;
   $C \leftarrow C + \hat{A} * \hat{B}$ ;
end for
    
```

It is possible to take advantage of communication pipelines in both directions of the process grid. However, the data allocation imposes that the source process of the broadcasts changes at each iteration in a cyclic fashion. That is, a given process broadcasts all of its columns of  $A$  or rows of  $B$  in multiple pieces of size proportional to the value of the distribution blocking factor  $NB_{dis}$ . The smaller this value is, the larger the number of messages and the lower the possible data reuse during each computational phase. In other words, the performance degrades as the value of the distribution blocking factor is decreased. If the value of this factor is very large, the communication computation overlap decreases, causing a performance degradation. However, high performance and efficiency can still be achieved for a wide range of blocking factors. This has been reported in [2], [21], [40], [47], [14].

Three alternatives to the physical blocking strategy are first presented in this section. Then, a few other related applications of those methods are outlined. The originality of the algorithms presented here is their systematic derivation from the properties of the underlying mapping. These blocking strategies are expressed within a single framework using LCM tables. The resulting blocked operations are appropriate for library software. They indeed feature potential for high performance without any specific alignment restrictions on their operands. This says that the antagonism between efficiency and flexibility is not a property of the block-cyclic mapping, but merely a characteristic of the algorithms that have been thus far proposed to deal with a distributed-memory hierarchy.

#### 3.1 Aggregation and Disaggregation

The *aggregation or algorithmic blocking* strategy operates on a panel of rows or columns that are globally contiguous. The local components of this panel before aggregation are also contiguous. The size of this panel is an algorithmic blocking factor. Its optimal value depends on the target machine characteristics. If this logical value is equal to the distribution blocking factors, then aggregation and physical blocking are the same. Otherwise, a few rows or columns that are globally contiguous and residing in distinct processes, are gathered into a single process row or column, and this

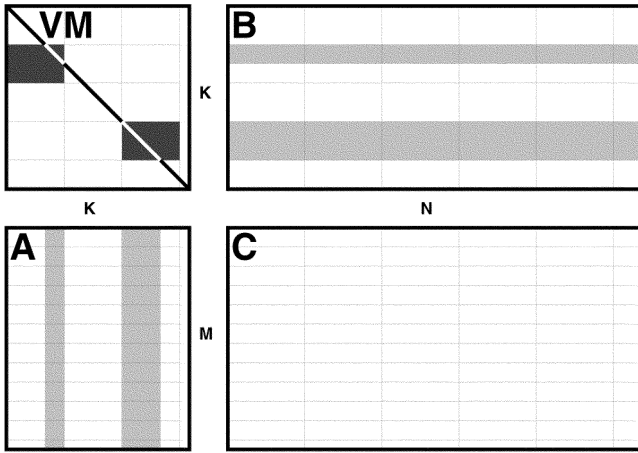


Fig. 5. Global view of the LCM blocked rank- $K$  update.

panel becomes the matrix operand. This strategy is particularly efficient when the distribution blocking factor is so small that Level 3 BLAS performance cannot be achieved locally on each process. Obviously, the aggregation phase induces some communication overhead. However, this must be weighted against the local computational gain. The problem is then to determine an algorithmic blocking factor  $NB_{alg}$  that keeps this overhead as low as possible and simultaneously optimizes the time spent in local computation. The feasibility and performance characteristics of this approach have been previously illustrated for the numerical resolution of a general linear system of equations and the symmetric eigenvalue problem [12], [7], [28], [43], [44] for the purely scattered distribution. Similarly, it is sometimes beneficial to disaggregate a panel into multiple panels in order to overlap communication and computation phases. When applicable, this last strategy also presents the advantage of requiring a smaller amount of workspace. The pseudocode of the rank- $K$  update operation using aggregation follows.

```

for  $kk = 1, K, NB_{alg}$ 
   $kb = \min(K - kk + 1, NB_{alg})$ ;
  Aggregate  $\hat{A} = A(:, kk : kk + kb - 1)$ ;
  Broadcast  $\hat{A}$  within process rows;
  Aggregate  $\hat{B} = B(kk : kk + kb - 1, :)$ ;
  Broadcast  $\hat{B}$  within process columns;
   $C \leftarrow C + \hat{A} * \hat{B}$ ;
end for

```

The aggregation and disaggregation techniques attempt to address the cases where the physical blocking strategy is not very efficient (i.e., for very small or large distribution blocking factors). In both techniques, the consecutive order of matrix columns or rows is preserved. It is therefore possible to use these techniques for algorithms that feature dependent steps, such as a triangular solve or the LU factorization with partial pivoting. The disaggregation technique, however, can only be applied efficiently for operations that do not feature any dependence between steps, such as a matrix-multiply. The disaggregated data remains consecutively ordered. Therefore, it cannot sig-

nificantly improve the load imbalance caused by consecutive allocation and consecutive elimination [38].

### 3.2 LCM Blocking

The *LCM blocking* strategy operates on a panel of rows or columns that are locally contiguous. The size of this panel is also an algorithmic blocking factor. Its optimal value depends on the target machine characteristics. However, rows or columns that may not be locally contiguous are packed according to an external criterion, typically the distribution parameters of another operand.

Consider the rank- $K$  update operation illustrated in Fig. 5. The LCM blocking strategy proceeds as follows: One is interested in finding the columns of  $A$  residing in a particular process column  $q$  and the rows of  $B$  residing in a particular process row  $p$  that could be multiplied together in order to update the matrix  $C$ . In Fig. 5, these columns of  $A$  and rows of  $B$  are indicated in gray. To accomplish this, one can consider the virtual matrix, denoted  $VM$  in the figure, defined by the column distribution parameters of  $A$  and the row distribution parameters of  $B$ . Locating the 0-diagonals of this  $VM$  in the process of coordinates  $(p, q)$  exactly solves the problem as illustrated in the figure. This can be realized by using LCM tables. As opposed to the physical blocking strategy, this technique does not assume the distribution equivalence of the columns of  $A$  and rows of  $B$  as suggested in Fig. 5. Moreover, the packing of these columns of  $A$  and rows of  $B$  is a local data copy operation (i.e., without communication overhead). For a given  $q$ , one just needs to go over all process rows, and thus treat all of the columns of  $A$  residing in this process column  $q$ . This algorithm can be regarded as a generalization of the physically blocked version. It presents, however, some advantages. First, as mentioned above, it does not assume an equivalent distribution of the columns of  $A$  and rows of  $B$ . Second, the communication overhead of the physically blocked variants has been partially replaced by a local data copy into a buffer that was needed, anyway. The communication pipeline stages in the row direction have been shortened. The cost of this pipeline startup has also been reduced by having the process column emitting the broadcasts remaining fixed as long as possible. Furthermore, this operation can also be logically blocked by limiting the number of columns of  $A$  in process column  $q$  and corresponding rows of  $B$  in the process row  $p$  that will be locally packed and broadcast at each step. The pseudocode for the LCM blocking strategy is given below.

```

for  $q = 0, Q - 1$ 
  for  $p = 0, P - 1$ 
     $n_{pq} =$  number of 0-diagonals process
       $(p, q)$  owns;
    Process column  $q$  packs and broadcasts
      those  $n_{pq}$  columns of  $A$  within process rows;
    Process row  $p$  packs and broadcasts those
       $n_{pq}$  rows of  $B$  within process columns;
    Perform local matrix-matrix multiply;
  end for
end for

```

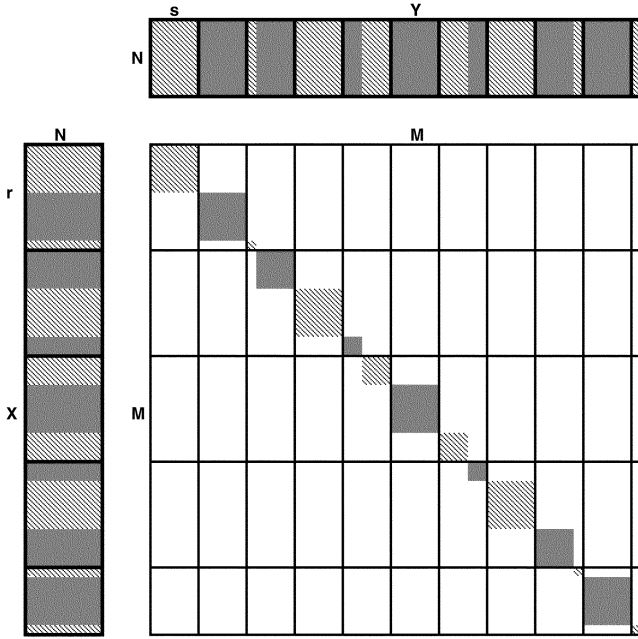


Fig. 6. Global view of one-dimensional redistribution.

This approach presents the advantage that the cost of the gathering phase is put on the processor as opposed to the interconnection network. However, it cannot be used for algorithms where each step depends on the previous one. Typically, the LCM blocking strategy is well-suited for multiplying two matrices, where each contribution to the resulting matrix entries can be added in any order. The LCM blocking strategy is a typical algorithmic redistribution operation, since it rearranges both logically and physically the communication and computation phases for increased efficiency and flexibility.

### 3.3 Aggregated LCM Blocking

The *aggregated LCM blocking* strategy is a hybrid scheme that combines the aggregation and LCM blocking strategies. In the aggregation scheme described earlier, the blocks to be aggregated were globally contiguous. It is, however, possible to use the same strategy for the local blocks obtained via LCM blocking. Furthermore, disaggregated LCM blocking is also possible, as previously noted. This more elaborate algorithmic blocking method maintains the local computational granularity even if the number of diagonals residing in a process is or becomes too small. The algorithm goes as follows: A process of coordinates  $(p, q)$  is considered, and the number of diagonals  $n_{pq}$  that process owns are handled by chunks of size  $NB_{alg}$ . If  $n_{pq}$  is a multiple of  $NB_{alg}$ , then the algorithm proceeds to the next process in the grid, either  $(p + 1, q)$  or  $(0, q + 1)$ . Otherwise, if  $(p + 1, q)$  is the next process, the remaining rows of  $B$  in process  $(p, q)$  are sent to the process  $(p + 1, q)$ , and the LCM blocking method proceeds to this process, taking into account the remainder of the previous step. If  $(0, q + 1)$  is the next process, this last procedure is applied to the remaining rows of  $B$  and columns of  $A$ . This algorithm therefore maintains the local computational granularity at a low communication overhead.

### 3.4 Redistribution and Static Blocking

The above framework can be used to tackle the run-time array redistribution problem when those arrays are distributed in a block-cyclic fashion over a multidimensional process grid. Solving this redistribution problem requires that 1) we generate the messages to be exchanged, and 2) we schedule these messages so that communication overhead is minimized. A comparative survey of the available literature can be found in [50].

Thus far, most of the focus regarding the run-time array distribution problem has been on the message generation phase [13], [34], [4], [46], and only a few papers deal with the communication scheduling phase [33], [42], [49], [39]. It turns out that the properties of the block-cyclic distribution presented in this paper can be used to further study this scheduling problem as it is shown in [20]. Moreover, the message generation phase can also be addressed with the help of the LCM tables. Figs. 6 and 7 illustrate this fact in the one-dimensional case.  $X$  (respectively  $Y$ ) is a  $M \times N$  one-dimensional array distributed over  $P$  (respectively  $Q$ ) processes with a distribution blocking factor of  $r$  (respectively  $s$ ). These figures show the global and local view of the redistribution mapping, as well as the  $M \times M$  distributed matrix induced by  $X$  and  $Y$ . Locating the diagonals of this matrix using LCM tables naturally provides a possible message generation algorithm. These figures also show that the general complexity of the redistribution problem is related to the number of processes in the  $P \times Q$  grid owning diagonals. Furthermore, the transpose and shift operations can be handled similarly within this framework. Finally, this approach can be generalized to handle the multidimensional case and the problem of accessing array entries with a nonunit stride [41]. It follows that efficient algorithms for the re-alignment of block cyclically distributed operands can also be expressed using the LCM tables and the above properties. In other words, flexible and efficient basic linear algebra kernels for distributed-memory concurrent computers can be expressed within the same framework.

LCM tables can be also used to derive another algorithmic blocking method, called the *static blocking* strategy hereafter, which deals only with purely local computational phases. It is assumed that the operation has reached a stage where the operands have already been redistributed if necessary by other techniques. Only local remaining computations need to be performed. It may, however, be the case that a local output operand has to be redistributed subsequently. Within this context, the symmetric rank- $K$  update operation  $C \leftarrow C + AA^T$  is easy to describe.  $C$  is an  $N \times N$  symmetric matrix for which only the upper or lower triangle should be referenced, and  $A$  is a matrix of dimension  $N \times K$ . The matrix  $A$  has been replicated in every process column and the matrix  $A^T$  replicated in every process row. The distributed matrix  $C$  is partitioned into diagonal and strictly upper or lower LCM blocks, as shown in Fig. 8. This figure shows the LCM block-partitioned matrices  $A$  and  $C$ , and the  $r \times s$ ,  $r \times K$  and  $K \times s$  blocks of these matrices that reside in the process of coordinates  $(p, q)$ . The arrangement of these blocks in process  $(p, q)$  is also represented and denoted by the local

arrays in process  $(p, q)$ . Depending on their relative position to the diagonal, the  $r \times s$  blocks of  $C$  are identified by a different shade of color. It is usually easy to deal with the strict upper or lower part using the BLAS matrix-matrix multiply. The diagonal LCM-block, however, requires particular attention.

Fig. 8 shows that the local update can be expressed in terms of symmetric rank- $K$  updates and matrix-matrix multiplies. The LCM tables provide the necessary information to organize the local computation in such a way that one can take advantage of the high efficiency of the matrix-matrix multiply kernel [8], [51]. A similar approach has been proven highly efficient for *GEMM-based* BLAS implementations [19], [32]. The static blocking strategy, even in its simplest form, imposes strong restrictions on the alignment and distribution of the operands. This is, nevertheless, the last opportunity for a large operation to logically rearrange the computations. This suggests that an efficient implementation of the symmetric rank- $K$  update when  $N \gg K$  would use the LCM blocking strategy to replicate  $A$  over  $C$ , and the static blocking technique to perform the local update. The algorithmic blocking factors for both phases can be chosen independently.

The use of physical blocking in conjunction with static blocking can lead to a comprehensive and scalable dense linear algebra software library. Existing serial software such as LAPACK [5] can be reused. The ScaLAPACK [11] software library is the result of this reasoning. As suggested above, if one limits oneself to static and physical blocking, strong alignment restrictions must be met by the matrix operands. It is argued that these restrictions are reasonable because 1) general redistribution software is available, 2) the user is ultimately responsible for choosing the initial data layout, and 3) the majority of practical cases are covered by this approach.

## 4 PERFORMANCE ANALYSIS AND EXPERIMENTAL RESULTS

A theoretical model of a distributed-memory computer is presented early in this section. It is an abstraction of physical models that provides a convenient framework for developing and analyzing parallel distributed, dense linear algebra algorithms without worrying about the implementation details or physical constraints. The model is then applied to the algorithmic blocking strategies presented in Section 3 in order to analyze their scalability. Finally, a number of experimental results are presented and summarized.

### 4.1 The Machine Model

Distributed-memory computers consist of processors that are connected using a message passing interconnection network. Each processor has its own memory, called the *local memory*, which is accessible only to that processor. As the time to access a remote memory is longer than the time to access a local one, such computers are often referred to as Non-Uniform Memory Access (NUMA) machines [36]. The interconnection network of our machine model is static, meaning that it consists of point-to-point communication links among processors. This type of network is also

referred to as a *direct network* as opposed to dynamic networks. The latter are constructed from switches and communication links. These links are dynamically connected to one another by the switching elements to establish, at run time, the paths between processors' memories. The interconnection network of the machine model considered here is a static, two-dimensional  $P \times Q$  rectangular mesh with wraparound connections. It is assumed that all processors can be treated equally in terms of local performance, and the communication rate between two processors is independent from the processors considered. In the two-dimensional mesh, each processor has four communication ports; however, the model assumes that a processor can send or receive data on only one of its ports at a time. This assumption is also referred to as the *one-port communication model* [36].

The time spent to communicate a message between two processors is called the *communication time*  $T_c$ . In our machine model,  $T_c$  is approximated by a linear function of the number  $L$  of items communicated.  $T_c$  is the sum of the time to prepare the message for transmission  $\alpha$  and the time  $\beta L$  taken by the message of length  $L$  to traverse the network to its destination, i.e.,

$$T_c = \alpha + \beta L.$$

This approximation of the communication time assumes that any two processors are equidistant from a communication point of view (cut-through or worm-hole routing). This approximation is reasonable for most current distributed-memory, concurrent computers. Finally, the model assumes that the communication links are bidirectional, that is, the time for two processors to send each other a message of length  $L$  is also  $T_c$ . A processor can send and/or receive a message on only one of its communication links at a time. In particular, a processor can send a message while receiving another message on the same or different link at the same time.

Since this paper is only concerned with a single, regular local operation, namely the matrix-matrix multiplication, the time taken to perform one floating point operation is assumed to be a constant  $\gamma$  in our model. This very crude approximation summarizes in a single number all the steps performed by a processor to achieve such a computation. Obviously, such a model neglects all the phenomena occurring in the processor components, such as cache misses, pipeline startups, memory load or store, floating point arithmetic, etc. that may influence the value of  $\gamma$  as a function of the problem size, for example. Similarly, the model does not make any assumption on the amount of physical memory per node. This machine model is a very crude approximation that is designed specifically to illustrate the cost of the dominant factors to our particular case. More realistic models are described, for example, in [36] and the references therein.

### 4.2 Scalability Analysis

The rank- $K$  update operation produces an  $M \times N$  matrix  $C$  by adding to itself the product of an  $M \times K$  matrix  $A$  and a  $K \times N$  matrix  $B$ ,



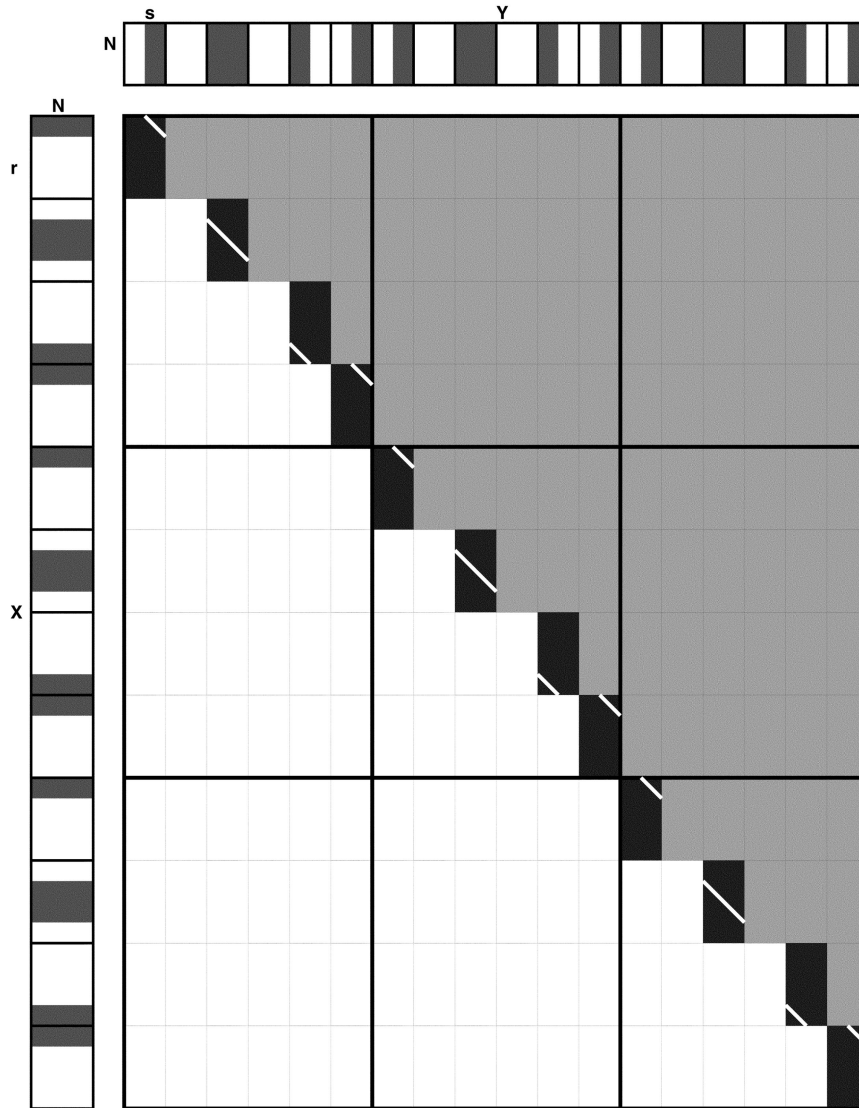


Fig. 7. Local view of one-dimensional redistribution.

$$C \leftarrow C + AB.$$

In the following, we assume for simplicity that  $M = N = K$ . The number of floating point operations is assumed to be equal to  $2N^3$ . All three matrices are distributed onto the same square process grid. Moreover, we also assume that the matrix operands are distributed according to the square block-cyclic data distribution (see Definition 2), and that the distribution blocking factors are the same for all operands. Therefore, no realignment phase is necessary to be performed. This distribution blocking factor is denoted by  $NB_{dis}$  in the following.

Four algorithms are considered, denoted by PHY, AGG, LCM, and RED. PHY denotes the physically blocked variant, AGG identifies the aggregation algorithm, and the LCM blocking algorithm is denoted by LCM. Finally, a fourth variant RED is considered where the matrices  $A$  and  $B$  are completely redistributed beforehand. For the algorithmic blocking variants AGG and LCM, the algorithmic blocking factor is denoted by  $NB_{alg}$ . Parallel efficiency,  $E(n, p)$ , for a problem of size  $n$  on  $p$  processors is defined in the usual way [24] by

$$E(n, p) = \frac{1}{p} \frac{T_{seq}(n)}{T(n, p)}, \quad (7)$$

where  $T(n, p)$  is the runtime of the parallel algorithm, and  $T_{seq}(n)$  is the runtime of the best sequential algorithm. An implementation is said to be *scalable* if the efficiency is an increasing function of  $n/p$ , the problem size per processor (in the case of dense matrix computations,  $n = N^2$ , the number of words in the input). The parallel runtime and efficiency on our machine model for the four algorithms PHY, AGG, LCM, and RED are computed below as a function of the local computational speed  $\gamma$ , the communication parameters  $\alpha$  and  $\beta_d$  (the time for a floating point number to traverse the network), and the total number of processors  $p$ .

The key-factor of this performance analysis is to model the cost of a sequence of broadcasts on a ring [2], [47] where the source either remains the same or is incremented by one after each broadcast. In the physical blocking strategy, the source process of the broadcast sequence is incremented at each step. The parallel runtime of the physically blocked variant algorithm is given by

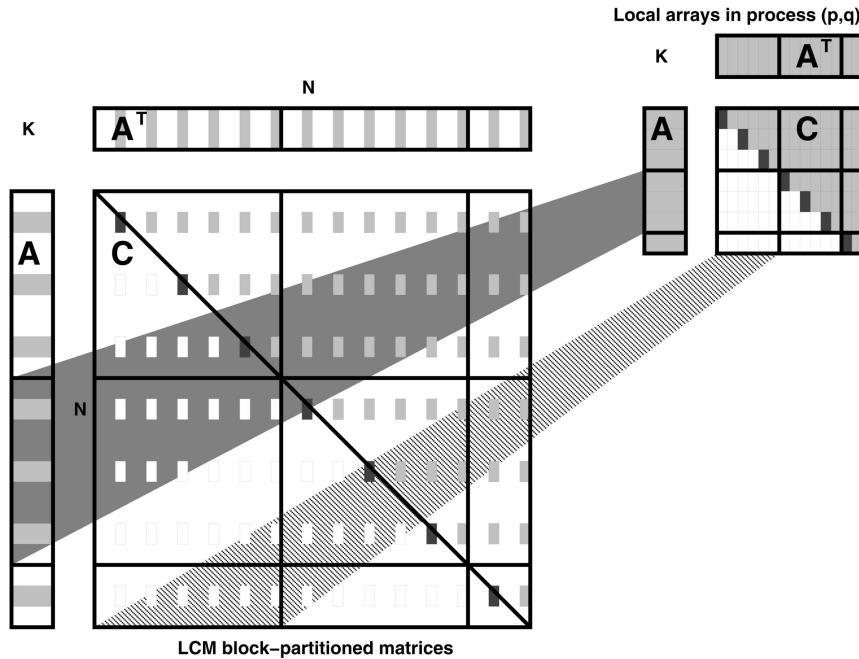


Fig. 8. Static symmetric rank- $K$  update.

$$T_{PHY}(N, p) \approx \frac{2N^3\gamma}{p} \left( 1 + \frac{2}{\gamma} \left( \frac{p\alpha}{NB_{dis}N^2} + \frac{\sqrt{p}\beta_d}{N} \right) \right)$$

when  $\frac{N}{NB_{dis}} \gg \sqrt{p}$ . A similar analysis for the physical blocking variant can also be found in [2], [47]. The aggregation blocking strategy essentially performs a sequence of accumulations followed by a ring broadcast. For the sake of simplicity, it is assumed that  $k$  blocks of the same size are aggregated ( $k \geq 2$ ). In practice, the blocks are only approximately of the same size.  $k$  is clearly bounded above by  $\sqrt{p}$ . In addition, the algorithmic blocking factor  $NB_{alg}$  is used to partition the communication and computation. It follows that the estimated execution time on our machine model for the aggregation strategy is given by

$$T_{AGG}(N, p) \approx \frac{2N^3\gamma}{p} \left( 1 + \frac{k}{\gamma} \left( \frac{p\alpha}{NB_{alg}N^2} + \frac{\sqrt{p}\beta_d}{N} \right) \right)$$

when  $\frac{N}{NB_{alg}} \gg \sqrt{p}$ . In the LCM blocking strategy, one looks at the diagonals of the virtual distributed matrix induced by the columns of  $A$  and rows of  $B$  residing in all process column and row pairs. It is assumed in this section that each process in the grid owns a number of diagonals that is proportional to  $NB_{alg}$ . With these assumptions, the estimated execution time of the LCM blocking strategy is given by

$$T_{LCM}(N, p) \approx \frac{2N^3\gamma}{p} \left( 1 + \frac{3}{2\gamma} \left( \frac{p\alpha}{NB_{alg}N^2} + \frac{\sqrt{p}\beta_d}{N} \right) \right)$$

when  $\frac{N}{NB_{alg}} \gg \sqrt{p}$ . Finally, the parallel run time of the RED variant is obtained by adding to the quantity  $T_{PHY}(N, p)$

computed above the approximated time to redistribute two square matrices of order  $N$ , that is  $2(p\alpha + \frac{N^2\beta_d}{p})$ .

Table 1 summarizes the estimated parallel efficiency for each variant studied in this Section. The LCM blocking variant features a slightly higher efficiency than the physical blocking strategy. This theoretical analysis also explains why one expects to observe better performance for the physical strategy than the aggregation variant when a “good” value of the distribution blocking factor  $NB_{dis}$  is selected. The physical blocking algorithm is thus *scalable* in the sense that if the memory use per process ( $\frac{N^2}{p}$ ) is maintained constantly, this algorithm maintains efficiency. The physical block size  $NB_{dis}$  can be used to lower the importance of the latency  $\alpha$ . The aggregation algorithm is also *scalable*. The value of  $k$  is a constant that only depends on the ratio between the algorithmic  $NB_{alg}$  and distribution  $NB_{dis}$  blocking factors. This formula shows the communication overhead induced by the aggregation strategy in terms of the number of messages, as well as the communication volume. When the distribution blocking factor is larger than the algorithmic blocking factor, the physical blocks are split into smaller logical blocks. Therefore, the estimated execution time of the disaggregation variant is always less than the estimated execution time of the aggregation method. The LCM blocking variant is also *scalable* for aligned matrix operands. This variant is slightly more efficient than the physical and aggregation strategies.

It should be noted, however, that our machine model assumes that the local data copy operation is free. In reality, such an assumption depends on the target machine, and may affect the results presented in Table 1. The RED algorithm, however, is not scalable because of the latency term.

### 4.3 Experimental Performance Results

The purpose of this section is to illustrate the general behavior of algorithmically redistributed operations as opposed to presenting a collection of particular performance numbers. The presentation style aims at facilitating the comparison of the different blocking strategies for a set of illustrative and particular cases. Experimental performance results are presented in Table 2 below for two distinct distributed-memory concurrent computers, namely the Intel XP/S Paragon and the IBM Scalable POWER-parallel System [3], [45].

Table 2 contains the specifications of the experiments that have been performed. In all of the experiments, the matrix operands were square of order  $N$ . The values of  $N$  used for all experiments are 100, 250, 500, 1,000, 1,500, 2,000 and 3,000. Due to memory size constraints, it was not always possible to perform the experiments for all of these values. Results are reported on a  $4 \times 4$  Intel XP/S Paragon and a  $4 \times 8$  IBM SP. Each experiment has been given an encoded name of the form  $XX\_A\#$ .  $XX$  identifies which target machine the experiment was run on, either XP for the Intel XP/S Paragon or SP for the IBM SP.  $\#$  is a number or a string distinguishing each experiment. For each experiment, the distribution parameters of the matrix operands  $A$ ,  $B$ , and  $C$  are the same. All of our experiments were performed in double precision arithmetic. The local rank update operation was performed by calling the appropriate subprogram of the vendor-supplied BLAS. The communication operations were implemented by explicit calls to the Basic Linear Algebra Communications Subprograms (BLACS). The BLACS [23] are a message passing library specifically designed for distributed linear algebra communication operations. The computational model consists of a one- or two-dimensional grid of processes, where each process stores matrices and vectors. The

BLACS include synchronous send/receive routines to send a matrix or submatrix from one process to another, to broadcast submatrices, or to compute global reductions (sums, maxima, and minima). There are also routines to establish, change, or query the process grid. The BLACS provide an adequate interface level for linear algebra communication operations. The performance of our algorithms is measured in Mflops/s. This is appropriate for large, dense linear algebra computations, since floating point dominates communication.

The matrix operands used for our experiments were distributed in such a way that no realignment phase was necessary, as explained in Section 4.2. Experimental performance results for nonaligned operands have been reported in [41]. Different values of the distribution blocking factor have been used. A machine dependent value of the algorithmic blocking factor  $NB_{alg}$  used by the AGG and LCM variants has first been determined for each platform and used for all of the experiments. On our Intel XP/S Paragon, we found that a reasonable value for this algorithmic blocking factor was 14. On the IBM SP, the value of 70 has been selected.

Fig. 9 shows the performance of the physical blocking (PHY), aggregation (AGG), and the LCM blocking (LCM) strategies using the value of  $NB_{alg}$  as the distribution blocking factors for the three matrix operands. According to the conclusions of the previous section, the performance of the three variants is almost identical on each platform with a slight advantage to the LCM blocking variant. In the rest of this section, the performance curves shown in Fig. 9 are considered as a reference. The combined maximum of these curves has been replicated on all of the other plots presented. This maximal curve is hereafter always represented as a bold solid line. Ideally, one would like to observe no difference between the performance obtained for this “good” physical layout and the performance achieved by distributions induced by other distribution blocking factors.

Fig. 10 shows the performance results obtained by the physical blocking strategy. The physical blocking variant uses the distribution blocking factors as the computational unit. When the distribution parameters are very small, the performance is dramatically degraded because of the local performance of the rank- $k$  update for small values of  $k$ . This is the difference that one should expect

TABLE 1

Estimated Parallel Efficiencies for Various Blocking Variants

$E_{PHY}(N, p)$	$(1 + \frac{2}{\gamma} (\frac{p\alpha}{NB_{dis}N^2} + \frac{\sqrt{p}\beta_d}{N}))^{-1}$
$E_{AGG}(N, p)$	$(1 + \frac{k}{\gamma} (\frac{p\alpha}{NB_{alg}N^2} + \frac{\sqrt{p}\beta_d}{N}))^{-1}$ with $k \approx \lceil \frac{NB_{alg}}{NB_{dis}} \rceil$
$E_{LCM}(N, p)$	$(1 + \frac{3}{2\gamma} (\frac{p\alpha}{NB_{alg}N^2} + \frac{\sqrt{p}\beta_d}{N}))^{-1}$
$E_{RED}(N, p)$	$(1 + \frac{1}{\gamma} ((2 + \frac{pNB_{dis}}{N}) \frac{p\alpha}{NB_{dis}N^2} + \frac{(2\sqrt{p} + 1)\beta_d}{N}))^{-1}$

 TABLE 2  
 Specification of the Experiments

Experiment #	Distribution parameters
XP_A0, SP_A0	$NB_{dis} = NB_{alg}$ for all operands.
XP_A1, SP_A1	$NB_{dis} = 1$ for all operands.
XP_A10	$NB_{dis} = 10$ for all operands.
SP_A20	$NB_{dis} = 20$ for all operands.
XP_A40	$NB_{dis} = 40$ for all operands.
XP_A100	$NB_{dis} = 100$ for all operands.
SP_A200	$NB_{dis} = 200$ for all operands.

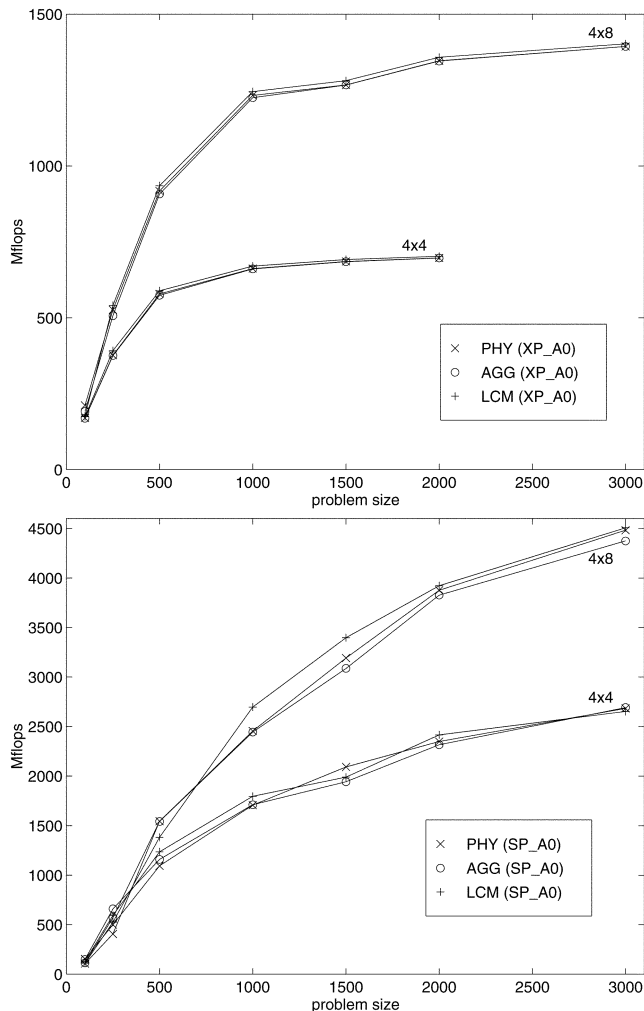


Fig. 9. Performance in Mflops/s of algorithmic blocking variants for a "good" physical data layout case, and various process grids on the Intel XP/S Paragon and the IBM SP.

when using Level 1 or 2 BLAS based algorithms on such computers, as opposed to Level 3 BLAS based algorithms. Very large distribution parameters increase the load imbalance, which is characterized by the ragged curves. For the experiment SP\_A200 and  $N = 1,500$ , each processor has almost the same amount of data. However, for  $N = 2,000$ , some processes have three times as much work to perform than others.

Fig. 11 shows the performance results obtained for the aggregation strategy. The dependence of the performance from the physical distribution parameters is largely decreased. The performance results are pushed towards the result of reference. For very small values of the distribution parameters, one expects a large performance improvement compared to the physical blocking strategy. This aspect is particularly evident for both target platforms. The aggregation phase induces some communication overhead that somewhat limits the potential of this strategy. This phenomenon is not particularly well illustrated on the Intel XP/S Paragon, due to the high speed of the interconnection network compared to the local computational performance. However, on the IBM SP, the performance of

Experiment SP\_A1 has been improved, but remains lower than the reference performance because of its less favorable communication-computation performance ratio.

Fig. 12 shows the performance results obtained for the LCM blocking strategy. The LCM blocking variant decouples the performance results from a poor choice of the distribution blocking factor. The LCM results are slightly better than the ones shown above for the aggregation variant. In particular, the performance results observed for experiments XP\_A1 and SP\_A1 have been considerably improved. On the Intel XP/S Paragon, the performance obtained for very small distribution blocking factors is now superior to the performance observed for distribution blocking factors slightly larger than  $N_{B_{alg}}$  (XP\_A40). On the IBM SP, there is virtually no performance difference between experiments SP\_A1 and SP\_A20. The impact of the less favorable communication-computation performance ratio of this particular machine is somewhat hidden by the algorithmic blocking strategy. This relatively low ratio is, however, the reason for the performance difference between the reference case and the experiments SP\_A1 and SP\_A20. The LCM blocking strategy builds panels of  $N_{B_{alg}}$  rows and columns with less communication overhead, because it essentially determines and regroups the columns of  $A$  and rows of  $B$  that belong to a given process column and process row pair. This phase is communication free. These results show that for aligned data and uniform data distributions, the performance difference due to various distribution blocking factors is no more than a few percentage points from the reference, as previously defined.

Fig. 13 shows the performance results when the matrix operands  $A$  and  $B$  are aligned but redistributed (RED) for efficiency reasons. To perform the complete redistribution of a two-dimensional block-cyclically distributed matrix, the appropriate component of the ScaLAPACK [11] software library [42] has been used. Even if these plots show the performance obtained for the same experiments as above, one could argue that complete redistribution (RED) should only be used for the extreme cases. A major feature of redistributing the entire matrix operands  $A$  and  $B$  at once is the large memory cost required by this operation. This increases the chances of the possible use of virtual memory by a large factor. Fig. 13 illustrates the dramatic performance consequences of using virtual memory on the Intel XP/S Paragon. On this particular machine, the complete redistribution beforehand leads to lower performance than the one obtained by the LCM blocking variant. In other words, the cost of redistributing when needed beforehand is larger than the cost induced by the algorithmically redistributed LCM strategy. In both variants, the amount of computation is performed at the same speed. On the IBM SP, the complete redistribution beforehand leads to slightly higher performance than the LCM blocking strategy. The lower total number of redistribution messages of the complete redistribution strategy takes better advantage of the low communication-computation performance ratio of this machine. It is clear that the IBM SP may need to use virtual memory for sufficiently large problem sizes. However, the nodes of the machine we used for our experiments each had at least 128 Megabytes of physical

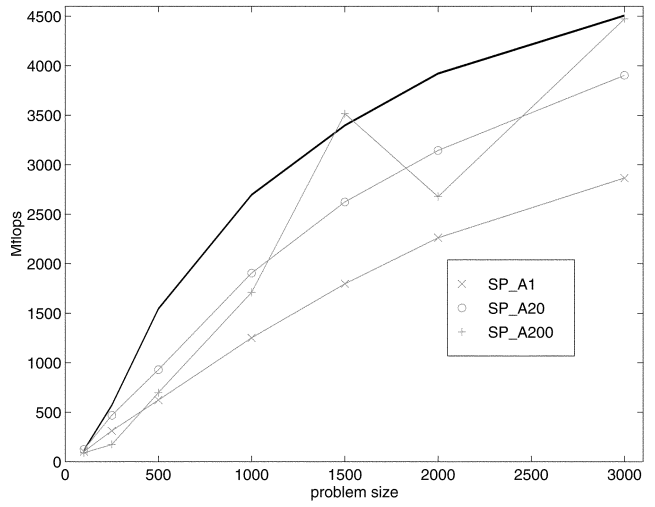
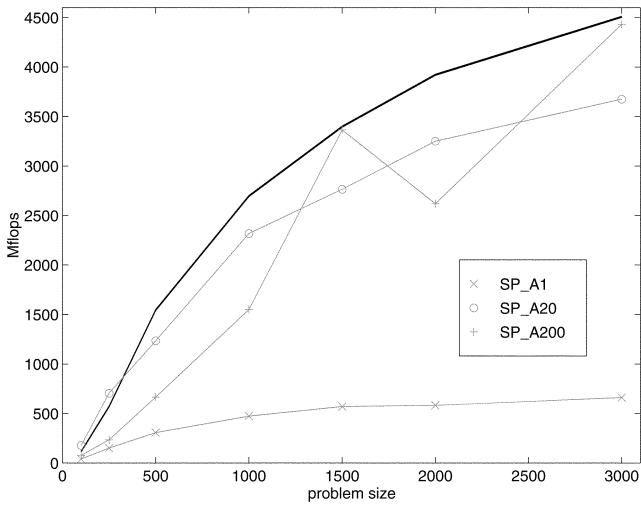
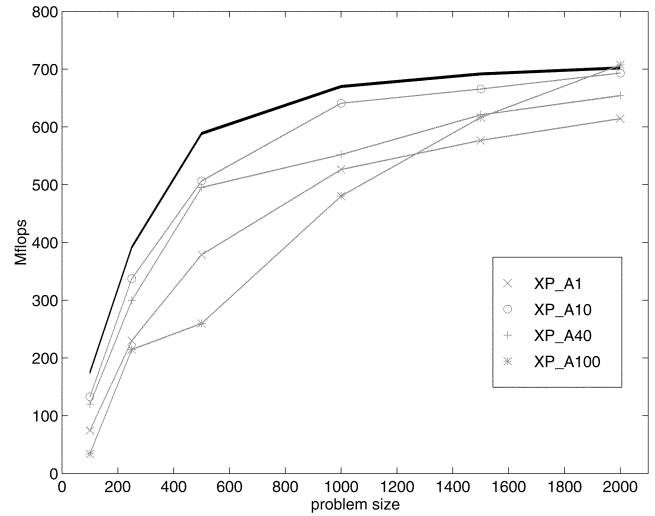
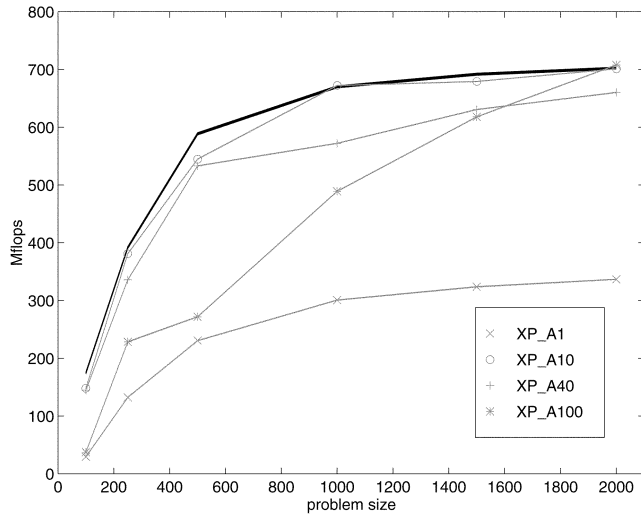


Fig. 10. Performance of the PHY variant on a  $4 \times 4$  Intel XP/S Paragon and on a  $4 \times 8$  IBM SP.

Fig. 11. Performance of the AGG variant on a  $4 \times 4$  Intel XP/S Paragon and on a  $4 \times 8$  IBM SP.

memory. It was not feasible to estimate the impact of the use of virtual memory in a reasonable amount of time.

For the aligned experiments on both platforms, it is legitimate to use algorithmic redistribution variants. By doing so, one can obtain high performance and efficiency independently from the distribution parameters. Moreover, the performance numbers obtained by the aggregation and LCM blocking techniques show a slight superiority for the latter. However, both techniques are complementary in the sense that it is not always possible to use the LCM blocking strategy as mentioned in Section 3.1. In order to address the problems induced by badly balanced computations, it is always possible to redistribute the matrix operand  $C$ , even if this somewhat contradicts the “owner’s compute” rule. The larger the operands, the more benefits one should obtain from a complete redistribution. However, the amount of memory necessary to perform such an operation grows with the number of items to be redistributed. This prevents the redistribution of the largest operands. Such an argumentation was at the forefront of our motivation for developing algorithmically redistributed operations that require a much smaller amount of memory.

## 5 CONCLUSION

Most of the parallel algorithms for basic linear algebra operations proposed in the literature thus far focus on the naturally aligned cases and rely on the physical blocking strategy to efficiently use a distributed memory hierarchy. This restricted interest prevents one from providing the necessary flexibility that a parallel software library requires to be truly usable. These restrictions also considerably handicap the ease-of-use of such a library, since one often needs to reformulate general operations to match obscure alignment restrictions that are difficult to document and to explain. In this paper, we presented various algorithmic redistribution methods that can be efficiently used to address this problem. These techniques were systematically derived from properties of the underlying mapping. Such a feature is particularly attractive from the software library design point of view. Furthermore, this approach can be extended to the more general family of Cartesian mappings [9] by generalizing the definition of an LCM table, as shown in [41]. Such a generalization is convenient to allow for the specification of submatrix operands where the upper left corner is not aligned on block boundaries [35].

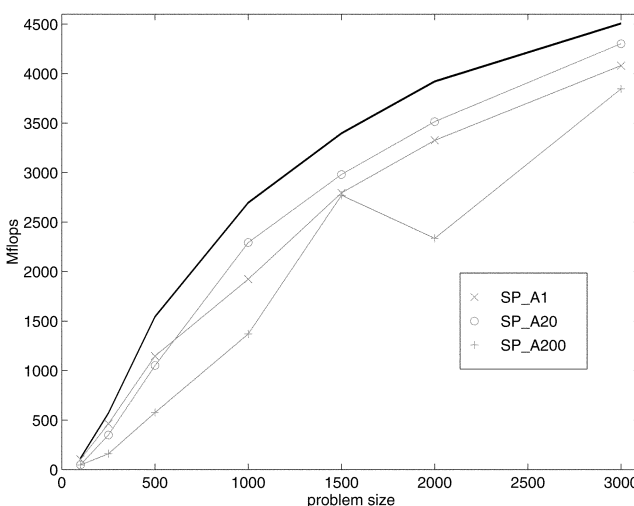
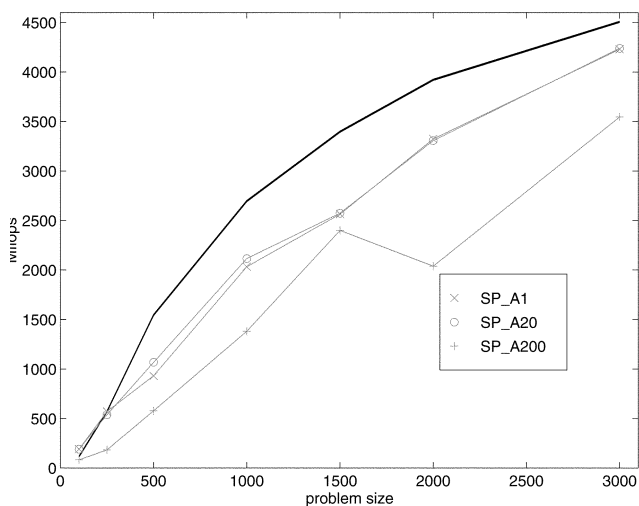
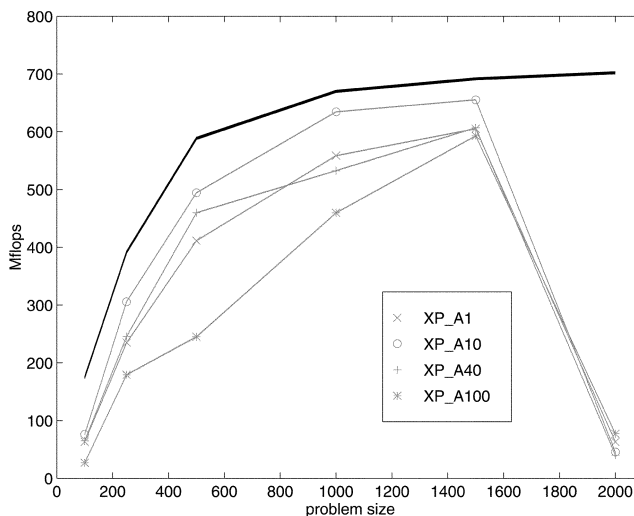
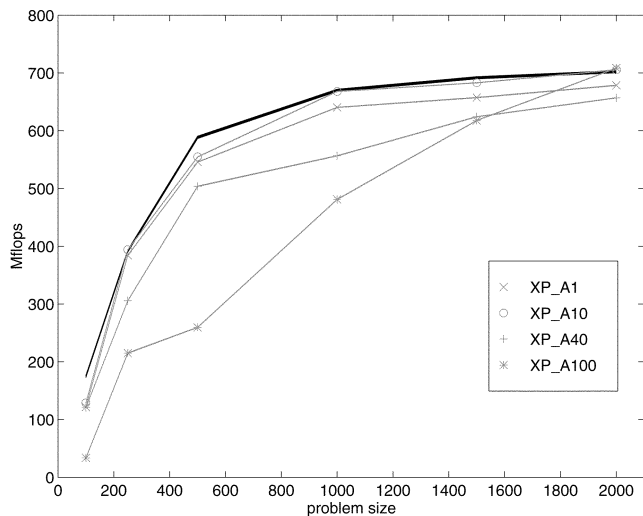


Fig. 12. Performance of the LCM variant on a  $4 \times 4$  Intel XP/S Paragon and on a  $4 \times 8$  IBM SP.

Fig. 13. Performance of the RED variant on a  $4 \times 4$  Intel XP/S Paragon and on a  $4 \times 8$  IBM SP.

The performance results presented in this paper show that when the matrix operands are aligned, the algorithmically redistributed operations based on the aggregation and the LCM blocking strategies are competitive in terms of performance with the beforehand complete redistribution variant (RED). For a variety of distribution and machine parameters, one can thus afford to redistribute the matrix operands "on the fly" without a significant performance degradation. This conclusion must be refined when the matrix operands have to be redistributed before the aligned operation can take place [41]. Nevertheless, for concurrent computers featuring slow communication performance compared to their computational power, it is necessary to preserve the possibility of redistributing the data beforehand despite the high memory cost. This problem can be resolved in two ways: 1) It is conceivable to redistribute the operands in two steps. At each step the same workspace can be reused and only part of the computation performed. This approach is viable, even if it is problematic from a software point of view to estimate at run-time the amount of usable memory on each process. 2) Redistribution in place is also possible, assuming enough memory has been initially allocated. The redistribution methods presented in this

paper also apply to networks or clusters of workstations and PCs. The efficiency of such methods clearly depends on the connectivity and performance of the network. When the nodes of such a computer are heterogeneous in terms of computational power, redistribution methods can still be efficiently utilized if the performance imbalance is compensated by allocating distinct number of processes on each node.

Algorithmic redistribution methods can alleviate natural alignment restrictions at a low, sometimes negligible, performance cost for basic operations and various block-cyclic distributions. In addition, these techniques considerably reduce and often completely remove the complicated dependence between the performance of parallel basic linear algebra operations and the physical distribution parameters. Indeed, it says that these algorithms facilitate the implementation of a general purpose and flexible parallel software library of basic linear algebra subprograms. These algorithms have been shown to achieve high performance independently from the actual block-cyclic distribution parameters. Efficiency and flexibility are not antagonistic objectives for basic dense linear algebra operations, but merely a characteristic of the algorithms

that have been so far proposed to deal with a distributed memory hierarchy.

## ACKNOWLEDGMENTS

The authors acknowledge using the Intel Paragon XP/S 5 computer, located in the Oak Ridge National Laboratory Center for Computational Sciences (CCS), funded by the Department of Energy's Mathematical, Information, and Computational Sciences (MICS) Division of the Office of Computational and Technology Research. This work also was supported by the U.S. Defense Advanced Research Projects Agency under contract DAAH04-95-1-0077, administered by the Army Research Office. This research was conducted, as well, by using the resources of the Cornell Theory Center, which receives major funding from the U.S. National Science Foundation (NSF) and New York State, with additional support from the Advanced Research Projects Agency (ARPA), the National Center for Research Resources at the National Institutes of Health (NIH), IBM Corporation, and other members of the center's Corporate Partnership Program.

A complete set of parallel basic linear algebra subprograms (PBLAS), along with the testing and timing programs for distributed-memory computers relying heavily on the algorithmic redistribution methods presented in this paper are available at <http://www.netlib.org/scalapack/prototype>.

## REFERENCES

- [1] M. Aboelaze, N. Chrisochoides, and E. Houstis, "The Parallelization of Level 2 and 3 BLAS Operations on Distributed-Memory Machines," Technical Report CSD-TR-91-007, Purdue Univ., West Lafayette, Ind., 1991.
- [2] R. Agarwal, F. Gustavson, and M. Zubair, "A High Performance Matrix Multiplication Algorithm on a Distributed-Memory Parallel Computer, Using Overlapped Communication," *IBM J. Research and Development*, vol. 38, no. 6, pp.673-681, 1994.
- [3] T. Agerwala, J. Martin, J. Mirza, D. Sadler, D. Dias, and M. Snir, "SP2 System Architecture," *IBM Systems J.*, vol. 34, no. 2, pp. 153-184, 1995.
- [4] C. Ancourt, F. Coelho, F. Irigoien, R. Keryell, "A linear Algebra Framework for Static HPF Code Distribution," Technical Report A-278-CRI, CRI-Ecole des Mines, Fontainebleau, France, 1995. (Available at <http://www.cri.ensmp.fr>.)
- [5] E. Anderson, Z. Bai, C. Bischof, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, S. Ostrouchov, and D. Sorensen, *LAPACK Users' Guide*. Philadelphia, Penn.: SIAM, 1995.
- [6] C. Ashcraft, "The Distributed Solution of Linear Systems Using the Torus-Wrap Data Mapping," Technical Report ECA-TR-147, Boeing Computer Services, Seattle, Wash., 1990.
- [7] P. Bangalore, "The Data-Distribution-Independent Approach to Scalable Parallel Libraries," master's thesis, Mississippi State Univ., 1995.
- [8] J. Bilmes, K. Asanovic, J. Demmel, D. Lam, and C. Chin, "Optimizing Matrix Multiply using PhiPAC: A Portable, High-Performance, ANSI C Coding Methodology," Technical Report UT CS-96-326, LAPACK Working Note 111, Univ. Tennessee, 1996.
- [9] R. Bisseling and J. van der Vorst, "Parallel LU Decomposition on a Transputer Network," *Lecture Notes in Computer Sciences*, G. van Zee and J. van der Vorst, eds., vol. 384, pp. 61-77, 1989.
- [10] R. Bisseling and J. van der Vorst, "Parallel Triangular System Solving on a Mesh Network of Transputers," *SIAM J. Scientific and Statistical Computing*, vol. 12, pp. 787-799, 1991.
- [11] L. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R.C. Whaley, *ScaLAPACK Users' Guide*. Philadelphia, Penn.: SIAM, 1997.
- [12] R. Brent and P. Strazdins, "Implementation of BLAS Level 3 and LINPACK Benchmark on the AP1000," *Fujitsu Scientific and Technical J.*, vol. 5, no. 1, pp. 61-70, 1993.
- [13] S. Chatterjee, J. Gilbert, F. Long, R. Schreiber, and S. Tseng, "Generating Local Addresses and Communication Sets for Data Parallel Programs," *J. Parallel and Distributed Computing*, vol. 26, pp. 72-84, 1995.
- [14] J. Choi, "A New Parallel Matrix Multiplication Algorithm on Distributed-Memory Concurrent Computers," Technical Report UT CS-97-369, LAPACK Working Note 129, Univ. Tennessee, 1997.
- [15] J. Choi, J. Dongarra, and D. Walker, "PUMMA: Parallel Universal Matrix Multiplication Algorithms on Distributed-Memory Concurrent Computers," *Concurrency: Practice and Experience*, vol. 6, no. 7, pp. 543-570, 1994.
- [16] J. Choi, J. Dongarra, and D. Walker, "PB-BLAS: A Set of Parallel Block Basic Linear Algebra Subroutines" *Concurrency: Practice and Experience*, vol. 8, no. 7, pp. 517-535, 1996.
- [17] A. Chtchelkanova, J. Gunnels, G. Morrow, J. Overfelt, and R. van de Geijn, "Parallel Implementation of BLAS: General Techniques for Level 3 BLAS," *Concurrency: Practice and Experience*, vol. 9, no. 9, pp. 837-857, 1997.
- [18] E. Chu and A. George, "QR Factorization of a Dense Matrix on a Hypercube Multiprocessor," *SIAM J. Scientific and Statistical Computing*, vol. 11, pp. 990-1,028, 1990.
- [19] M. Daye, I. Duff, and A. Petitet, "A Parallel Block Implementation of Level 3 BLAS for MIMD Vector Processors," *ACM Trans. Mathematical Software*, vol. 20, no. 2, pp. 178-193, 1994.
- [20] F. Desprez, J. Dongarra, and A. Petitet, C. Randriamaro, Y. Robert, "Scheduling Block-Cyclic Array Redistribution," *IEEE Trans. Parallel and Distributed Systems*, vol. 9, no. 2, pp. 192-205 1998.
- [21] J. Dongarra, R. van de Geijn, and D. Walker, "Scalability Issues in the Design of a Library for Dense Linear Algebra," *J. Parallel and Distributed Computing*, vol. 22, no. 3, pp. 523-537, 1994.
- [22] J. Dongarra and D. Walker, "Software Libraries for Linear Algebra Computations on High Performance Computers," *SIAM Review*, vol. 37, no. 2, pp. 151-180, 1995.
- [23] J. Dongarra and R.C. Whaley, "A User's Guide to the BLACS v1.0," Technical Report UT CS-95-281, LAPACK Working Note 94, Univ. Tennessee, 1995. (<http://www.netlib.org/blacs/>)
- [24] G. Fox, M. Johnson, G. Lyzenga, S. Otto, J. Salmon, and D. Walker, *Solving Problems on Concurrent Processors*. Englewood Cliffs, N.J.: Prentice Hall, 1988.
- [25] G. Fox, S. Otto, and A. Hey, "Matrix Algorithms on a Hypercube I: Matrix Multiplication," *Parallel Computing*, vol. 3, pp. 17-31, 1987.
- [26] G. Geist and C. Romine, "LU Factorization Algorithms on Distributed-Memory Multiprocessor Architectures," *SIAM J. Scientific and Statistical Computing*, vol. 9, pp. 639-649, 1988.
- [27] M. Heath and C. Romine, "Parallel Solution Triangular Systems on Distributed-Memory Multiprocessors," *SIAM J. Scientific and Statistical Computing*, vol. 9, pp. 558-588, 1988.
- [28] B. Hendrickson, E. Jessup, and C. Smith, "A Parallel Eigensolver for Dense Symmetric Matrices," personal communication, 1996.
- [29] B. Hendrickson and D. Womble, "The Torus-Wrap Mapping for Dense Matrix Calculations on Massively Parallel Computers," *J. Scientific and Statistical Computing*, vol. 15, no. 5, pp. 1,201-1,226, Sept. 1994.
- [30] G. Henry and R. van de Geijn, "Parallelizing the QR Algorithm for the Unsymmetric Algebraic Eigenvalue Problem: Myths and Reality," Technical Report UT CS-94-244, LAPACK Working Note 79, Univ. Tennessee, 1994.
- [31] S. Huss-Lederman, E. Jacobson, A. Tsao, and G. Zhang, "Matrix Multiplication on the Intel Touchstone DELTA," *Concurrency: Practice and Experience*, vol. 6, no. 7, pp. 571-594, 1994.
- [32] B. Kågström, P. Ling, and C. van Loan, "GEMM-Based Level 3 BLAS: High-Performance Model Implementations and Performance Evaluation Benchmark," Technical Report UMINF 95-18, Dept. Computing Science, Umeå Univ., 1995.
- [33] E. Kalns and L. Ni, "Processor Mapping Techniques towards Efficient Data Redistribution," *IEEE Trans. Parallel and Distributed Systems*, vol. 12, no. 6, pp. 1,234-1,247, 1995.

- [34] K. Kennedy, N. Nedeljković, and A. Sethi, "A Linear-Time Algorithm for Computing the Memory Access Sequence in Data Parallel Programs," *Proc. Fifth ACM SIGPLAN, Symp. Principles and Practice of Parallel Programming*, 1995.
- [35] C. Koebel, D. Loveman, R. Schreiber, G. Steele, and M. Zosel, *The High Performance Fortran Handbook*. Cambridge, Mass.: MIT Press, 1994.
- [36] V. Kumar, A. Grama, A. Gupta, and G. Karypis, *Introduction to Parallel Computing*. Redwood City, Calif.: Benjamin/Cummings Publishing Company, Inc., 1994.
- [37] G. Li and T. Coleman, "A New Method for Solving Triangular Systems on Distributed-Memory Message-Passing Multiprocessor," *SIAM J. Scientific and Statistical Computing*, vol. 10, no. 2, pp. 382–396, 1989.
- [38] W. Lichtenstein and S.L. Johnsson, "Block-Cyclic Dense Linear Algebra," *SIAM J. Scientific and Statistical Computing*, vol. 14, no. 6, pp. 1,259–1,288 1993.
- [39] Y. Lim, P. Bhat, and V. Prasanna, "Efficient Algorithms for Block-Cyclic Redistribution of Arrays," Technical Report CENG 97-10, Dept. Electrical Engineering-Systems, Univ. Southern California, Los Angeles, Calif., 1997.
- [40] K. Mathur, S.L. Johnsson, "Multiplication of Matrices of Arbitrary Shapes on a Data Parallel Computer," *Parallel Computing*, vol. 20, pp. 919–951, 1994.
- [41] A. Pettit, *Algorithmic Redistribution Methods for Block Cyclic Decompositions*, doctoral thesis, Univ. Tennessee, Knoxville, 1996.
- [42] L. Prylli and B. Tourancheau, "Fast Runtime Block Cyclic Data Redistribution on Multiprocessors," *J. Parallel and Distributed Computing*, vol. 45, 1997.
- [43] P. Strazdins, "Matrix Factorization using Distributed Panels on the Fujitsu AP1000," *Proc. IEEE First Int'l Conf. Algorithms and Architectures for Parallel Processing (ICA3PP-95)*, 1995.
- [44] P. Strazdins and H. Koesmarno, "A High Performance Version of Parallel LAPACK: Preliminary Report," *Proc. Sixth Parallel Computing Workshop*, Fujitsu Parallel Computing Center, 1996.
- [45] C. Stunkel, D. Shea, B. Abali, M. Atkins, C. Bender, D. Grice, P. Hochschild, D. Joseph, B. Nathanson, R. Swetz, R. Stucke, M. Tsao, and P. Varker, "The SP2 High-Performance Switch," *IBM Systems J.*, vol. 34, no. 2, pp. 185–204, 1995.
- [46] A. Thirumalai and J. Ramanujam, "Fast Address Sequence Generation for Data Parallel Programs Using Integer Lattices," *Languages and Compilers for Parallel Computing: Lecture Notes in Computer Science*. P. Sadayappan et al., eds., Springer-Verlag, 1996.
- [47] R. van de Geijn and J. Watts, "SUMMA: Scalable Universal Matrix Multiplication Algorithm," *Concurrency: Practice and Experience*, vol. 9, no. 4, pp. 255–274, 1997.
- [48] E. van de Velde, "Experiments with Multicomputer LU Decomposition," *Concurrency: Practice and Experience*, vol. 2, pp. 1–26, 1990.
- [49] D. Walker and S. Otto, "Redistribution of Block-Cyclic Data Distributions Using MPI," *Concurrency: Practice and Experience*, vol. 8, no. 9, pp. 707–728, 1996.
- [50] L. Wang, J. Stichnoth, S. Chatterjee, "Runtime Performance of Parallel Array Assignment: An Empirical Study," *Proc. Supercomputing*, 1996. (<http://www.supercomp.org/sc96/proceedings/>).
- [51] R. Whaley and J. Dongarra, "Automatically Tuned Linear Algebra Software," Technical Report UT CS-97-366, LAPACK Working Note 131, Univ. Tennessee, 1997.



**Antoine P. Pettit** is a research scientist in the Computer Science Department at the University of Tennessee, Knoxville. His research interests primarily focus on parallel computing, numerical linear algebra, and the design of scientific parallel numerical software libraries for distributed-memory concurrent computers. He was involved in the design and implementation of the software packages ScaLAPACK and ATLAS, and is currently working on the design of algorithms and techniques for high performance computer architectures.



**Jack J. Dongarra** holds a joint appointment as distinguished professor of computer science in the Computer Science Department at the University of Tennessee (UT) and as a distinguished scientist in the Mathematical Sciences Section at Oak Ridge National Laboratory (ORNL) under the UT/ORNL Science Alliance Program. He specializes in numerical algorithms in linear algebra, parallel computing, use of advanced-computer architectures, programming methodology, and tools for parallel computers. Other current research interests include the development, testing, and documentation of high quality mathematical software. He was involved in the design and implementation of the software packages EISPACK, LINPACK, the BLAS, LAPACK, ScaLAPACK, Netlib, PVM/Hence, MPI, the National High-Performance Software Exchange, NetSolve, and ATLAS. He is currently involved in the design of algorithms and techniques for high performance computer architectures. Dr. Dongarra's professional activities include membership in SIAM, the IEEE, the ACM, and as a fellow of the AAAS. He has published numerous articles, papers, reports, and technical memoranda, and has given many presentations on his research interests.