

Designing a BSP version of ScaLAPACK*

Guy Horvitz[†] and Rob H. Bisseling[‡]

July 13, 1998

Abstract

The ScaLAPACK library for parallel dense matrix computations is built on top of the BLACS communications layer. In this work, we investigate the use of BSPlib as the basis for a communications layer. We examine the LU decomposition from ScaLAPACK and develop a BSP version which is significantly faster. The savings in communication time are typically 10-15%. The gain in overall execution time is less pronounced, but still significant. We present the main features of a new library, BSP2D, which we propose to develop for porting the whole of ScaLAPACK.

Keywords: Bulk Synchronous Parallel, LU Decomposition, Numerical Linear Algebra

1 Introduction

To obtain the highest performance in parallel computation both computation and communication must be optimised. LAPACK [1] has provided us with highly optimised implementations of state-of-the-art algorithms in the field of numerical linear algebra, in particular for the solution of dense linear systems and eigensystems. Many years of effort have gone into optimising LAPACK, and much of its success is due to the encapsulation of system-dependent optimisations into the Basic Linear Algebra Subprograms (BLAS).

*Preprint 1074 of the Department of Mathematics, Utrecht University, July 1998.

[†]Fritz Haber Research Center for Molecular Dynamics, Hebrew University, Jerusalem 91904, Israel (guyh@fh.huji.ac.il).

[‡]Department of Mathematics, Utrecht University, P.O. Box 80010, 3508 TA Utrecht, The Netherlands (Rob.Bisseling@math.uu.nl).

LAPACK is available for sequential computers, vector supercomputers, and parallel computers with shared memory.

The ScaLAPACK [4] project aims to provide a scalable version of LAPACK for parallel computers with distributed memory. Portability is ensured by building ScaLAPACK on top of the Basic Linear Algebra Communication Subprograms (BLACS). The parallel efficiency depends critically on the communication performance achieved by this library and thus it is natural to ask whether the performance can be further improved.

The bulk synchronous parallel (BSP) model [11] views a parallel algorithm as a sequence of supersteps, each containing computation and/or communication, followed by a global synchronisation of all the processors. (Actually, the original BSP model by Valiant [11] left open the possibility of synchronising a subset of the processors.) This imposes a discipline on the user, thus making parallel programming simpler, but it also provides possibilities for system-level optimisation such as combining and rescheduling of messages. This can be done because the superstep provides a natural context for communication optimisation by the system. The user need not be concerned about such optimisations.

A BSP computer can be characterised by four global parameters:

- p , the number of processors
- s , the computing speed in flop/s (floating point operations per second)
- g , the communication time per data element sent or received, measured in flop time units
- l , the synchronisation time, also measured in flop time units.

Algorithms can be analysed by using the parameters p , g , and l ; the parameter s just scales the time. The time of a superstep with both computation and communication is $w + hg + l$, where w denotes the maximum amount of work (in flops) of a processor, and h is the maximum number of data elements sent or received by a processor. The total execution time of an algorithm (in flops) can be obtained by adding the times of the separate supersteps. This yields an expression of the form $a + bg + cl$. In the following presentation, we consider the architecture as an abstract BSP computer, and therefore we use the term ‘processes’ instead of ‘processors’. In our experiments, only one process executes on each processor, so one may use the terms interchangeably.

BSPLib [8] is a proposed standard which makes it possible to program directly in BSP style. BSPLib is an alternative to PVM [10] and MPI [6], and its communication performance competes with that of MPI. BSPLib provides

both direct remote memory access (i.e., one-sided communications such as put and get), which is suitable for certain regular computations, and bulk synchronous message passing.

BSPPACK [3] is an application package built on top of BSPlib. It is a research and educational library which contains parallel implementations of algorithms for sparse and dense linear system solution, fast Fourier transforms, and other scientific computations.

The aim of this work is to answer the question: can ScaLAPACK be ported to BSPlib and does this improve performance? This may indeed be the case, because we expect ScaLAPACK to benefit from ideas developed within the context of BSPPACK and from the excellent implementation of BSPlib available as the Oxford BSP toolset [7]. We limit ourselves to investigating the ScaLAPACK LU decomposition subroutine `PSGETRF`, but we expect our results to be valid for other ScaLAPACK subroutines as well.

The design philosophy of ScaLAPACK is to use a hierarchy of software layers. The top of the pyramid is ScaLAPACK itself, which calls the Parallel BLAS (PBLAS). The PBLAS use the BLAS for single-process linear algebra computations and the BLACS for communication. The BLACS can be built on top of a basic communications layer such as MPI or PVM. The BLACS perform communication at a higher level: they send complete matrices of all types and they allow us to view the processes as a two-dimensional grid and to perform operations within the scope of a process row or column, or the complete grid.

The data distribution of ScaLAPACK is the two-dimensional block-cyclic distribution with a user determined block size nb . Another parameter is the algorithmic block size nb' . The algorithms in the sequential package LAPACK handle complete blocks of size nb' . ScaLAPACK structures its algorithms in the same way, but it imposes $nb' = nb$. We make the same choice for reasons of convenience, but in our case it is straightforward to relax this constraint so that nb' can be any multiple of nb ; we shall discuss this later.

Since the communication in ScaLAPACK is isolated in the BLACS it would be the most natural choice to construct a BLACS version based on BSPlib. A straightforward BSPlib implementation of the BLACS, however, would be impossible for different reasons; one important reason is the following. The BLACS include pair-wise message passing for communication where the receiver has to wait for the data to arrive in order to continue. The sender can continue as soon as the message is sent off. In BSPlib, a message transfer is completed only after the next global synchronisation. Suppose there is exactly one message to be communicated and hence in the program there is one call to a BLACS send and one to a BLACS receive. The processes

that do not send or receive are not aware of this communication and hence do not synchronise, thus violating the principle of global synchronisation.

Forcing the user to synchronise globally between a send and a receive requires drastic changes in both the syntax and the semantics of the BLACS subroutines. This would turn the BLACS into a different library, which could be called BSP2D. Section 4 outlines how such a library could be constructed in the future. The present work simply removes the BLACS and adapts ScaLAPACK and the PBLAS using direct calls to BSPlib. This alone is not sufficient: it is also necessary to restructure ScaLAPACK and the PBLAS on the basis of supersteps.

The remainder of this paper is organised as follows. Section 2 discusses the changes needed to port the ScaLAPACK LU decomposition to BSPlib. Section 3 presents experimental results, including a comparison of communication based on BSPlib, MPI BLACS, and native BLACS. Section 4 presents BSP2D. Section 5 draws the conclusions.

2 BSP version of ScaLAPACK LU decomposition

Programming in BSPlib requires global synchronisation. For this reason, every process should know when a global synchronisation is needed to perform a certain task. Sometimes, a process also needs to know about resources (such as buffers) provided by remote processes. Such knowledge can be transferred by communication, but this would be inefficient.

Another solution would be to let all the processes call subroutines together and with the same values for the scalar input parameters. This way, each process can deduce the behaviour of the other processes. We adopted this solution for the PBLAS. For example, consider the PBLAS subroutine `PSSWAP` which swaps two rows or columns of distributed matrices. If the swap is local and no communication is needed, the processes do not synchronise. Otherwise, all the processes perform one synchronisation, even if they do not hold any of the related data and do not actively participate in the operation. All processes can distinguish between the two situations, because they all have the necessary information.

2.1 Unblocked LU decomposition and pivot search subroutines

An example of how a ScaLAPACK subroutine and a PBLAS should be altered, is shown in the case of the ScaLAPACK subroutine `PSGETF2`, which performs an unblocked parallel LU decomposition on a block of consecutive columns. The main part of the `PSGETF2` code is given in Fig. 1. This subroutine is called by the main LU decomposition subroutine `PSGETRF`.

In the original subroutine, the main loop (`DO 10 ... 10 CONTINUE`) is executed only by the process column `IACOL` that holds the block to be decomposed. After the decomposition, the pivot indices `IPIV(IIA..IIA+MN-1)` of that block are broadcast to the other process columns by the sending subroutine `IGEB2D` and the receiving subroutine `IGEBR2D`. This structure is inherited from the PBLAS. Since the PBLAS subroutine `PSAMAX`, which finds the pivot of matrix column `J`, returns the result only to the processes of the process column `IACOL` that holds `J`, the other processes cannot evaluate the if-condition `GMAX.NE.ZERO` that tests for singularity.

We mentioned earlier that a BSP based PBLAS should be called by all processes with the same values for the scalar input parameters. The example of `PSGETF2` makes it clear that scalar output parameters must be returned to all processes too. This way, `GMAX` and the pivot index become available to all the processes so they can participate in the main loop, and can call subsequent PBLAS together, as required. Inevitably, sending the output scalars to all processes costs extra communication and synchronisation time.

A clear advantage of the changes in `PSAMAX` and `PSGETF2` is the ability to choose an algorithmic block size that differs from the distribution block size. This is impossible in the current version of ScaLAPACK; e.g. if $nb' = 2nb$, then two process columns should participate in the decomposition of one (algorithmic) block of columns. The subroutine `PSAMAX`, however, returns its results only to one process column, namely to the process column that holds matrix column `J`.

2.2 Collective communication subroutines

Sometimes, we need subroutines to perform collective communications such as broadcasts or reductions. In our case, we need to broadcast data within a process row (or column), and perform this operation for all process rows simultaneously. The method adopted for the PBLAS, global replication of scalar parameters, is not suitable here. The reason is that the size of the broadcast may differ between the process rows. We must allow different sizes, but the number of synchronisations should not depend on them.

```

IF( MYCOL.EQ.IACOL ) THEN
DO 10 J = JA, JA+MN-1
  I = IA + J - JA
*
*   Find pivot and test for singularity.
*
CALL PSAMAX( M-J+JA, GMAX, IPIV( IIA+J-JA ), A, I, J,
$           DESCA, 1 )
IF( GMAX.NE.ZERO ) THEN
*
*   Apply the row interchanges to columns JA:JA+N-1
*
CALL PSSWAP( N, A, I, JA, DESCA, DESCA( M_ ), A,
$           IPIV( IIA+J-JA ), JA, DESCA, DESCA( M_ ) )
*
*   Compute elements I+1:IA+M-1 of J-th column.
*
IF( J-JA+1.LT.M )
$   CALL PSSCAL( M-J+JA-1, ONE / GMAX, A, I+1, J,
$             DESCA, 1 )
ELSE IF( INFO.EQ.0 ) THEN
  INFO = J - JA + 1
END IF
*
*   Update trailing submatrix
*
IF( J-JA+1.LT.MN ) THEN
$   CALL PSGER( M-J+JA-1, N-J+JA-1, -ONE, A, I+1, J, DESCA,
$             1, A, I, J+1, DESCA, DESCA( M_ ), A, I+1,
$             J+1, DESCA )
END IF
10 CONTINUE
*
CALL IGEBS2D( ICTXT, 'Rowwise', ROWBTOP, MN, 1, IPIV( IIA ),
$           MN )
*
ELSE
*
CALL IGEBR2D( ICTXT, 'Rowwise', ROWBTOP, MN, 1, IPIV( IIA ),
$           MN, MYROW, IACOL )
*
END IF

```

DEL

DEL
DEL
DEL
DEL
DEL
DEL
DEL
DEL

Figure 1: Main part of the PSGETF2 source code. Lines marked with DEL are deleted in the BSP version.

The simplest solution is always to use a broadcast with two synchronisations, except when the broadcast is in the scope of one or two processes. For one process no synchronisation is needed, and for two processes a single synchronisation suffices. All processes can take the same decision because the number of participants in the broadcast is known to all of them. The choice of performing two synchronisations in the general case is based on the efficiency of the so-called *two-phase broadcast* [2, 3, 9], which first scatters the elements of a data vector across all the processes, and then lets each process broadcast the data it received. This was shown to be efficient in the LU decomposition program from BSPPACK, see [3].

2.3 Multiple row swap subroutine

The ScaLAPACK subroutine PSLASWP applies a series of row exchanges in a matrix prescribed by a given vector of pivoting indices. This is originally done by pairwise row swaps, each time using the PBLAS subroutine PSSWAP. A direct translation into BSP would imply one superstep for each swap. We change the method so that all the swaps are done in one superstep, in good BSP style. The changes are as follows.

First we translate the representation of the permutation from swaps into cycles. For example, suppose the swap vector is

$$(4, 10)(3, 10)(2, 10)(1, 10),$$

which means: first swap rows 1 and 10, then 2 and 10, etc. In this example, rows 1, 2, 3, 4 are on the same process A and row 10 resides on a different process B. The cycle representation of this permutation is

$$(10, 4, 3, 2, 1),$$

which means: 1 goes to 2, 2 goes to 3, ..., 10 goes to 1. The operations performed by A and B in this case are:

Process A	Process B
Put row 4 in buffer on process B	Put row 10 in buffer on process A
For $i = 3$ to 1 step -1	
copy row i into row $i + 1$	
Sync	Sync
Copy buffer into row 1	Copy buffer into row 10

In this way, only one row is exchanged between A and B. In the original algorithm, which performs the swaps sequentially, four rows are exchanged.

In the general case, the cycles are handled separately, but with one global synchronisation for all of them.

The case of a cycle of three or more rows is not expected to appear often for a randomly constructed matrix, except at the end of the algorithm. Yet, it could occur for certain structured matrices. In that case, our algorithm has an advantage over the original one.

2.4 Registered buffers

Often, we have to communicate noncontiguous data like e.g. a matrix row, which in ScaLAPACK is stored as a strided subarray. (Matrix columns, however, are stored contiguously.) The data elements can of course be sent separately, but even though BSPlib combines small messages, there is still a notable overhead for extremely small messages such as single words. This is because some additional information must be transmitted for each message. If the access pattern is regular, e.g. if it is based on a stride, then the overhead can be avoided by packing messages in buffers.

Put operations are the most efficient means of communication on many architectures, including our test machine. When we use put operations for communications, the locations of the buffers in which we put the packed data must have been registered previously. The purpose of registration is to link the name of a remote variable with its local address. Since registration incurs communication and synchronisation cost, it is more efficient to register the locations only once, at the beginning of the computation. The locations should then be passed to the PBLAS.

For this purpose, we implemented a management system for registered buffers. At the start of the program, we allocate and register buffers of appropriate sizes. When a PBLAS requests a buffer of a certain size, it calls a subroutine which returns a pointer to the smallest buffer of at least the requested size. Similar to the registration procedure of BSPlib, buffers are requested in lock step. All processes participate in all requests, and they ask for a buffer of the same size. Otherwise it can happen that different processes obtain differently named buffers, so that data are put in unexpected locations.

To achieve the ultimate in efficiency, we use the high performance put primitive `bsp_hpput` which is unbuffered on source and unbuffered on destination, instead of `bsp_put` which is doubly buffered. In the case of the high performance primitives, responsibility for buffering rests on the user instead of on the BSPlib system. On our test machine, we found that the improvement in performance was significant. (On other machines, this may not be the case.) The use of `bsp_hpput` means we cannot put data in a remote buffer that holds data which may be used in the same superstep. Such a

conflict may occur when two subsequent PBLAS subroutines use the same buffer, as will be shown below.

Most of the PBLAS we used have this general structure:

1. Pack local data needed for the global operation in a buffer.
2. Communicate the buffer contents to remote buffers.
3. Synchronise.
4. Calculate locally using the data in the buffer.

Now, suppose two subroutines **SR1** and **SR2** are called one after the other, using the same buffer. Let us examine the possible behaviour of two processes A and B. The two processes are synchronised after the communication in **SR1**. In this scenario, process A spends more time in local calculation in **SR1** than B. While A is still using the data in the buffer for local calculation, B has already reached the communication stage of **SR2** and it has put data in the buffer of A. Since the data used by A in the calculation is now overwritten, the calculation will be erroneous. To avoid such a problem, we have to make sure that a buffer is not reused immediately. This is done by associating a *used* flag with each buffer. A used buffer will be ignored by subsequent requests, until it is released by the process and synchronisation occurs.

3 Experimental results

We performed numerical experiments on a CRAY T3E computer with 64 processors, each with a theoretical peak performance of 600 Mflop/s. We measured a sequential speed of $s = 250$ Mflop/s for the matrix multiplication part of the LU decomposition. Normalised for this value of s , we found $g \approx 14$ and $l \approx 16000\text{--}80000$. (We measured these values within the context of the program, not in a separate benchmark. This explains the variation in l .) The aim was to compare the ScaLAPACK performance of three communication layers: BSPLib, a Cray-specific version of the BLACS, and an MPI version. We ran tests for three different process grids (with size 8×8 , 16×4 , 4×16) and four different block sizes ($nb = 16, 32, 48, 64$). The optimal grid size for all three communication layers was 8×8 , and the optimal block size was 32. We used single precision, which is 64 bits on this machine. We ran a test program which generates a square matrix with random elements.

The measured computing rate is given in Table 1. The rate is based on the overall execution time, including all overheads. For small problems, where

communication dominates, the table shows a considerable gain in speed obtained by using BSPlib: about 10% compared to the native BLACS and 50% compared to the MPI BLACS, for $n = 1000$. For large problems, where computation dominates, the differences are less pronounced: about 2% compared to the native BLACS and 4% compared to the MPI BLACS, for $n = 10000$. In all cases, the BSPlib version is faster than the others, with one exception, namely $n = 8000$. Here, the native BLACS are slightly faster than BSPlib. We investigated this further; our findings will be given below.

Table 1: Computing rate in Gflop/s of LU decomposition on a CRAY T3E for three different communication layers. The process grid has size 8×8 ; the block size is 32.

Size	BSPlib	native BLACS	MPI BLACS
500	0.37	0.33	0.22
1000	1.22	1.11	0.81
1500	2.33	2.04	1.60
2000	3.19	3.04	2.42
2500	4.49	4.01	3.36
3000	5.34	4.96	4.15
3500	6.25	5.73	4.93
4000	6.79	6.57	5.61
4500	7.67	7.16	6.34
5000	8.43	8.05	7.20
6000	9.63	9.24	8.56
7000	10.65	10.25	9.71
8000	10.47	10.55	10.04
9000	12.15	11.83	11.47
10000	12.73	12.48	12.20

To understand the savings in execution time, we measured the time used by each part of the program. Using BSPlib we can measure the communication/synchronisation time separately from the local computation time. We then separated the communication time from the synchronisation time by using a theoretical estimate of the synchronisation time provided by the BSP cost model. We also measured the packing time, i.e., the time each process spent in packing and unpacking data, including the time spent in local swaps. Finally, from the total time and other measurements, we could estimate the idle time of each process, which we define as the average time

a process waits for the others to finish their computation. (We do not include the time a process waits while others communicate.) The resulting breakdown of the costs is presented in Fig. 2. As expected, the computation time, which is of the order $O(n^3/p)$, dominates for large n . Note that the synchronisation time, although only linear in n , is still significant compared to the computation time.

The computation and idling time is identical for all three versions, because they differ only in the communication part. By subtracting the computation and idling time from the measured total time, we can obtain the time of the communication part (including packing and synchronisation where needed). The results are presented in Fig. 3. It is clear that the communication time for BSP is significantly less than for the other two versions. The exception is again the case $n = 8000$, for which the native BLACS are slightly faster. For large n , the typical savings compared to the native BLACS are 10–15%.

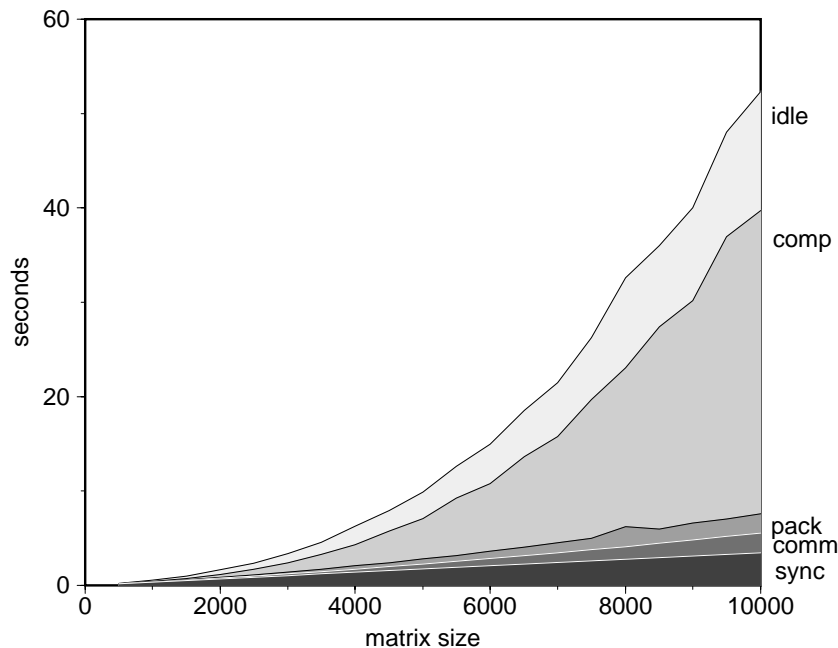


Figure 2: Breakdown of the total execution time for BSP based LU decomposition. The components are: synchronisation, communication, packing, computation, and idling while other processes compute. The process grid has size 8×8 ; the block size is 32.

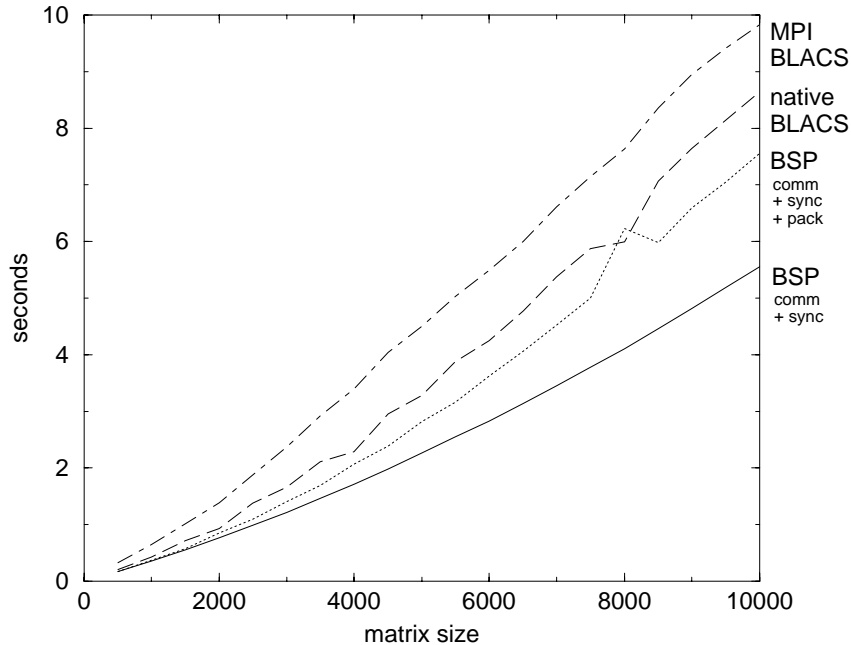


Figure 3: Communication/packing time during LU decomposition for three communication layers: BSPlib, native BLACS and MPI BLACS. For BSPlib, the time without packing is also given. The process grid has size 8×8 ; the block size is 32.

What is the reason for this exception? We found the BLAS `SCOPY` and `SSWAP`, which we used for packing and local swaps in `PSSWAP` and `PSLASWP`, to be highly inefficient on the Cray T3E, and they reduced the performance of our program considerably. For example, about 85% of the time consumed by `PSLASWP` is used for packing of data and local swaps and only 15% for communication. (The native BLACS suffer somewhat less, since 75% of their time is spent packing.) To expose how badly these BLAS implementations perform, we ran some tests using the native BLACS based `PSLASWP`, comparing the case where all swaps are done between different processes and the opposite case where they are all performed locally so that no communication is needed. Surprisingly, the swaps with communication are faster: for $n = 4000$ by 15%; and for $n = 8000$ by an exceptionally high 40%. This is responsible for the better performance of the BLACS in the case $n = 8000$. This abnormality indicates that there is much scope for improvement on this particular machine. Improving the copying would reduce the communica-

tion/packing time for BSPlib to that depicted in the lower line of Fig. 3. The native BLACS version would gain as well, but not as much as BSPlib; for the MPI version we cannot estimate the gain.

4 Proposal for a BSP2D library

When developing a BSP implementation of the whole ScaLAPACK, it would be most convenient to have available a high level BSP based communication layer, called BSP2D. This would save much effort and would also improve modularity. The position of the BSP2D layer in the ScaLAPACK hierarchy is shown in Fig. 4. BSP2D has the functionality of the BLACS, i.e., communicating complete matrices (or vectors) of different types. Like the BLACS it views the processes as a two-dimensional grid. It can be built on top of BSPlib, but it is not limited to BSPlib; in principle it could be implemented using any suitable BSP communication library.

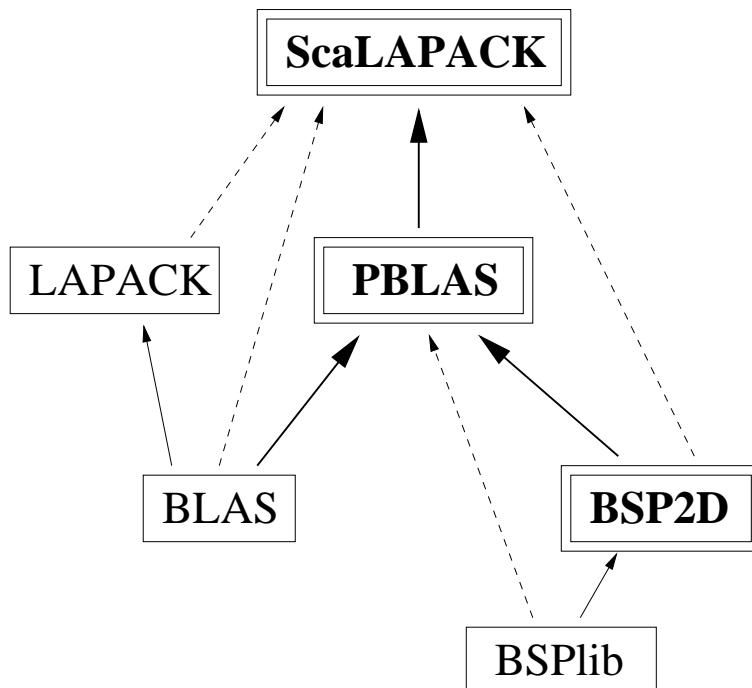


Figure 4: Hierarchical view of a possible BSP based ScaLAPACK, adapted from [5]. The double boxes contain the parts affected by moving to BSP. Solid arrows represent major dependencies; dashed arrows represent minor dependencies. The solid bold arrows represent the main structure of ScaLAPACK.

There are two types of communication operation in BSP2D: pair-wise communications and collective communications. Pair-wise communications should be done by bulk synchronous message passing (using `bsp_send`) and not by direct remote memory access (using `bsp_put`, `bsp_get` or their high performance equivalents). Direct remote memory access cannot be used for the following reason. The communication of noncontiguous data structures involves packing of the data in buffers. Communicating by direct remote memory access requires previous registration of these buffers. Since the size of data each process sends is not always known to the other processes we cannot use the global management system for registered buffers described in Subsection 2.4. (For a general library such as BSP2D we cannot adopt the same solution as for the PBLAS, namely calling each subroutine with the same global parameters. This would render the library hard to use.) An alternative would be to register a buffer for each put operation, but this would be inefficient. A third possibility would be to use static pre-registered buffers, where each process makes $p - 1$ buffers available for use by the other processes; this solution wastes too much memory. Therefore, none of the solutions based on direct remote memory access are satisfactory.

As a consequence, pair-wise communication is done by bulk synchronous message passing. This means that data are sent, and after synchronisation the destination process moves the data from its receive queue. Messages are identified in the following way. Each message consists of a payload and a tag. The payload consists of the matrix to be communicated, packed in a suitable form. The tag consists of the identity of the sending process and the number of messages that were already sent by that process to the receiving process in the current superstep. (This number represents the order in which the send operations occur in the program text, and not the actual order in which BSPlib sends them. BSPlib is still allowed to optimise communication by rescheduling messages within a superstep.) The tag may contain other information, such as type.

In BSP2D, moves of messages originating in the same process must be done in the order those messages were sent; this is similar to the requirement for receives in the BLACS. The messages of the receive queue of BSPlib, however, are in arbitrary order and the queue can only be accessed in this order. Still, this poses no problem since the high performance move operation `bsp_hpmove` can be used to create a list of the message positions in the queue. This operation is done as part of the BSP2D synchronisation subroutine. In an implementation, the list can be sorted in linear time by source process and message number.

The use of `bsp_hpmove` instead of `bsp_move` also enables BSP2D to unpack data straight from the receive buffer, thus saving the time of local copy-

ing. Performance could be improved even more if a high performance send were available, so the data could be sent straight from the source memory. However, the primitive `bsp_hpsend` does not exist in BSPLib.

Collective communications such as broadcasts and reductions involve synchronisation, so they should be called by all processes at the same time. They can be performed in the scope of a process row, a process column, or the whole process grid. To ensure that all the processes perform the same number of synchronisations, these subroutines always have two supersteps, except when the number of processes in the scope is one or two. As already observed in our study of LU decomposition (see Subsection 2.2) the decision on the number of synchronisations cannot rely on the number of data to be communicated, since it may vary between different process rows or columns. Only if the scope of the collective communication is the whole process grid, the decision may depend on the data size, and we can use this to our advantage.

We already described the two-phase broadcast in Subsection 2.2. Two-phase reduction is similar. Suppose the scope of the operation has q processes. In the reduction, each process has a vector of the same size n . Associative and commutative component-wise operations have to be performed on these q vectors, such as taking the minimum or maximum, or adding. This is done as follows. The data on each process are divided into q blocks of size n/q , numbered $0, \dots, q-1$, and each block is sent to a different process, so that process i gets all the blocks numbered i . Then each process performs a local reduction of the blocks, and sends the result to all the other processes. The total cost is about $2ng + 2l$.

In summary, BSP2D will include subroutines for pair-wise and collective communications, for global synchronisation with additional housekeeping, for the creation, initialisation, and the destruction of the process grid, and for retrieving the grid dimensions and process coordinates.

5 Conclusions

In this work, we have demonstrated that it is feasible to produce a bulk synchronous parallel version of an important ScaLAPACK subroutine. The BSP version outperforms two other versions, one based on a vendor-built BLACS communication layer, and the other on MPI BLACS. The savings in execution time were entirely due to a reduction of the communication time; the computation part of the program was left unchanged.

For large problems, e.g. $n = 10000$, communication time was reduced by up to 15% compared to the vendor-built BLACS version and even more compared to the MPI version. Because our test machine has relatively fast

communications the reduction in total execution time is less pronounced: about 2% compared to the native BLACS and 4% compared to the MPI BLACS. For machines with a higher g , i.e., with slower communication relative to computation, the influence of communication on the total time will be larger, and hence the gain we expect to achieve by using BSPlib will be proportionally larger.

For small problems, communication is dominant and the savings in total time are considerable: for $n = 1000$ the gain in overall computing rate is about 10% compared to the native BLACS and 50% compared to the MPI BLACS.

We have outlined how the complete ScaLAPACK package could be ported to BSP. When porting software based on message passing such as ScaLAPACK, the BSP philosophy may sometimes be felt as restrictive, but in developing new software from scratch the BSP constraints help tremendously in simplifying the programming task. Still, we have shown that despite the constraints imposed by BSP we can port ScaLAPACK code with a relatively minor effort and with substantial performance improvements, while maintaining full functionality.

Our practical experience in porting one major ScaLAPACK subroutine led to the formulation of the BSP2D library. Whereas we could build one single routine (and the required PBLAS) directly on top of BSPlib and we could manage the registered buffers within the subroutine, this would not be a feasible solution for the whole of ScaLAPACK. Instead, using an intermediate BSP2D layer would increase modularity and software reuse, at only a slight increase in cost due to copying and other overheads. We emphasise that an efficient implementation of `bsp_send` is crucial for the efficiency of BSP2D. The current work has shown that a public-domain software layer such as BSPlib can outperform a vendor-supplied layer. In principle, we would expect a vendor-supplied version of BSPlib to improve performance even more.

The approach of BSPlib, based on global synchronisation, can be carried over to the PBLAS, and this gives the additional advantage that the algorithmic block size can be decoupled from the distribution block size. This enables a better trade-off between load balance, speed of the BLAS operations in the unblocked part of the algorithm, and speed in the blocked part. This way, we provide further opportunities for improving the performance of ScaLAPACK.

Acknowledgements

The numerical experiments were performed on the Cray T3E of the High Performance Applied Computing centre at the Technical University of Delft, with support from the NCF. The BSP version of the ScaLAPACK LU decomposition is available upon request.

References

- [1] 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, Release 2.0*. SIAM, Philadelphia, PA, second edition, 1995.
- [2] Mike Barnett, Satya Gupta, David G. Payne, Lance Shuler, and Robert van de Geijn. Building a high-performance collective communication library. In *Proceedings Supercomputing 1994*, 1994.
- [3] Rob H. Bisseling. Basic techniques for numerical linear algebra on bulk synchronous parallel computers. In Lubin Vulkov, Jerzy Waśniewski, and Plamen Yalamov, editors, *Workshop Numerical Analysis and its Applications 1996*, volume 1196 of *Lecture Notes in Computer Science*, pages 46–57. Springer-Verlag, Berlin, 1997.
- [4] L. S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitdet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users' Guide*. SIAM, Philadelphia, PA, 1997.
- [5] Jaeyoung Choi, Jack J. Dongarra, L. Susan Ostrouchov, Antoine P. Petitdet, David W. Walker, and R. Clint Whaley. The design and implementation of the ScaLAPACK LU, QR, and Cholesky factorization routines. *Scientific Programming*, 5:173–184, 1996.
- [6] William Gropp, Ewing Lusk, and Anthony Skjellum. *Using MPI: Portable programming with the Message-Passing Interface*. MIT Press, Cambridge, MA, 1994.
- [7] Jonathan M. D. Hill, Stephen R. Donaldson, and Alistair McEwan. Installation and user guide for the Oxford BSP toolset (v1.3) implementation of BSPLib. Technical report, Oxford University Computing Laboratory, Oxford, UK, November 1997.

- [8] Jonathan M. D. Hill, Bill McColl, Dan C. Stefanescu, Mark W. Goudreau, Kevin Lang, Satish B. Rao, Torsten Suel, Thanasis Tsantilas, and Rob Bisseling. BSPlib: The BSP programming library. Technical Report PRG-TR-29-97, Oxford University Computing Laboratory, Oxford, UK, May 1997. To appear in *Parallel Computing*.
- [9] Ben H. H. Juurlink and Harry A. G. Wijshoff. Communication primitives for BSP computers. *Information Processing Letters*, 58:303–310, 1996.
- [10] V. S. Sunderam. PVM: A framework for parallel distributed computing. *Concurrency: Practice and Experience*, 2(4):315–339, 1990.
- [11] Leslie G. Valiant. A bridging model for parallel computation. *Communications of the ACM*, 33(8):103–111, 1990.