

# Questions and Answers about BSP

---

D.B. SKILLICORN,<sup>1</sup> JONATHAN M.D. HILL,<sup>2</sup> AND W.F. McCOLL<sup>2</sup>

<sup>1</sup> Department of Computing and Information Science, Queen's University, Kingston, Canada; e-mail: skill@qcis.queensu.ca

<sup>2</sup> Computing Laboratory, University of Oxford, Oxford, UK; e-mail: {Jonathan.Hill, Bill.McColl}@comlab.ox.ac.uk

## ABSTRACT

Bulk Synchronous Parallelism (BSP) is a parallel programming model that abstracts from low-level program structures in favour of *supersteps*. A superstep consists of a set of independent local computations, followed by a global communication phase and a barrier synchronisation. Structuring programs in this way enables their costs to be accurately determined from a few simple architectural parameters, namely the permeability of the communication network to uniformly-random traffic and the time to synchronise. Although permutation routing and barrier synchronisations are widely regarded as inherently expensive, this is not the case. As a result, the structure imposed by BSP does not reduce performance, while bringing considerable benefits for application building. This paper answers the most common questions we are asked about BSP and justifies its claim to be a major step forward in parallel programming.

## 1 Why Is Another Model Needed?

In the 1980s, a large number of different types of parallel architectures were developed. This variety may have been necessary to thoroughly explore the design space but, in retrospect, it had a negative effect on the commercial development of parallel applications software. To achieve acceptable performance, software had to be carefully tailored to the specific architectural properties of each computer, making portability almost impossible. Each new generation of processors appeared in strikingly-different parallel architectural frameworks, forcing performance-driven software developers to redesign their applications from the ground up. Understandably, few were keen to join this process.

Today, the number of parallel computation models and languages probably exceeds the number of different architectures with which parallel programmers had to contend ten years ago. Most make it hard to achieve portability, hard to achieve performance, or both.

The two largest classes of models are based on message passing, and on shared memory. Those based on message passing are inadequate for three reasons. First, messages require paired actions at the sender and receiver, which it is difficult to ensure are correctly matched. Second, messages blend communication and synchronisation so that sender and receiver must be in appropriately-consistent states when the communication takes place. This is appallingly difficult to ensure in most models, and programs are prone to deadlock as a result. Third, the performance of such programs is impossible to predict because the interaction of large numbers of individual messages in the interconnection mechanism makes the variance in their delivery times large.

The argument for shared-memory models is that they are easier to program because they provide the abstraction of a single, shared address space. A whole class of placement decisions are avoided. This is true, but is only half of the issue. When memory is shared, simultaneous access to the same location must be prevented. This requires either PRAM-style discipline by the programmer, or expensive lock management (and locks *are* expensive on today's parallel computers [16]). In both cases, the benefits are counterbalanced by quite serious drawbacks.

From an architectural point of view, shared-memory abstractions limit the size of computer that can be built because a larger and larger fraction of the computer's resources must be devoted to communication and the maintenance of coherence. Even worse, this part of the computer is most likely to be highly customized, and hence to be proportionally more expensive. Thus even the proponents of shared memory agree that, with our current understanding, such architectures can contain no more than, say, fifty processors. Whether this is sufficient for the application demands of the next decade is debatable.

The Bulk Synchronous Parallel (BSP) model [36] is a distributed-memory abstraction that treats communication as a bulk action of a program, rather than as the aggregate of a set of individual, point-to-point messages. It provides software developers with an attractive escape route from the world of architecture-dependent parallel software. The emergence of the model has coincided with the convergence of commercial parallel machine designs to a standard architectural form with which it is compatible. These developments have been enthusiastically welcomed by a rapidly-growing community of software engineers who produce scalable and portable parallel applications. However, while the parallel-applications community has welcomed the approach, there is a degree of skepticism amongst parts of the computer science research community. Some people seem to regard some of the claims made in support of the BSP approach as "too good to be true". We will make these claims, and back them up, in what follows.

The only sensible way to evaluate an architecture-independent model of parallel computation such as BSP is to consider it in terms of *all* of its properties, that is

- (a) its usefulness as a basis for the design and analysis of algorithms,
- (b) its applicability across the whole range of general-purpose architectures and its ability to provide efficient, scalable performance on them, and
- (c) its support for the design of fully-portable programs with analytically-predictable performance.

To focus on only one of these at a time, is simply to replace the zoo of parallel architectures in the 1980s by a new zoo of parallel models in the 1990s. A fully-rounded viewpoint on the nature and role of models seems more appropriate as we move from the straightforward world of parallel algorithms to the much more complex world of parallel software systems.

## 2 What is Bulk Synchronous Parallelism?

Bulk Synchronous Parallelism is a style of parallel programming intended for parallelism across all application areas and a wide range of architectures [25]. Its goals are more ambitious than most parallel-programming systems which are aimed at particular kinds of applications, or work well only on particular classes of parallel architectures [26].

BSP's most fundamental properties are that:

1. *It is simple to write.* BSP imposes a high-level *series-parallel* structure on programs which makes them easy to write, and to read. Existing BSP languages are SPMD, making programs even simpler, since the parallelism is largely implicit.
2. *It is independent of target architectures.* Unlike many parallel programming systems, BSP is designed to be architecture-independent, so that programs run unchanged when they are moved from one architecture to another. Thus BSP programs are portable in a strong sense.
3. *The performance of a program on a given architecture is predictable.* The execution time of a BSP program can be computed from the text of the program and a few simple parameters of the target architecture. This makes engineering design possible, since the effect of a decision on performance can be determined at the time it is made.

BSP achieves these properties by raising the level of abstraction at which programs are written and implementation decisions made. Rather than considering individual processes and individual communication actions, BSP considers computation and communication at the level of the entire program, and the entire executing computer and its interconnection mechanism. Determining the *bulk* properties of a program, and the *bulk* ability of a particular computer to satisfy them makes it possible to design with new clarity.

One way in which BSP is able to achieve this abstraction is by renouncing locality as a performance optimisation. This simplifies many aspects of both program and implementation design, and in the end does not adversely affect performance for most application domains. There will always be some application domains for which locality is critical, for example low-level image processing, and for these BSP may not be the best choice.

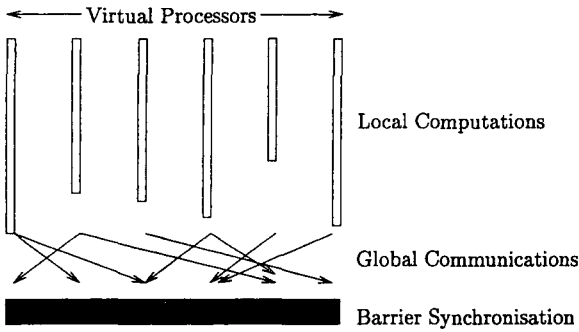


FIGURE 1 A superstep.

### 3 What Does the BSP Programming Style Look Like?

BSP programs have both a vertical structure and a horizontal structure. The vertical structure arises from the progress of a computation through time. For BSP, this is a sequential composition of global *supersteps*, which conceptually occupy the full width of the executing architecture. A superstep is shown in Figure 1.

Each superstep is further subdivided into three ordered phases consisting of:

1. simultaneous local computation in each process, using only values stored in the memory of its processor;
2. communication actions amongst the processes, causing transfers of data between processors;
3. a barrier synchronisation, which waits for all of the communication actions to complete, and which then makes any data transferred visible in the local memories of the destination processes.

The horizontal structure arises from concurrency, and consists of a fixed number of virtual processes. These processes are **not** regarded as having a particular linear order, and may be mapped to processors in any way. Thus locality plays no role in the placement of processes on processors.

We will use  $p$  to denote the virtual parallelism of a program, that is the number of processes it uses. If the target parallel computer has fewer processors than the virtual parallelism, an extension of Brent's theorem [5] can be used to transform any BSP program into a slimmer version.

### 4 How Does BSP Communication Work?

Most parallel programming systems treat communication, both conceptually and in implementations, at the

level of individual actions: memory-to-memory transfers, sends and receives, or active messages. This level is difficult to work with because parallel programs contain many simultaneous communication actions, and their interactions are complex. For example, congestion in the interconnection mechanism is typically very sensitive to the applied load. This makes it hard to discover much about the time any single communication action will take to complete, because it depends so much on what else is happening in the computer at the same time.

Considering communication actions *en masse* both simplifies their treatment, and makes it possible to bound the time it takes to deliver a whole set of data. BSP does this by considering all of the communication actions of a superstep as a unit. For the time being, imagine that all messages have a fixed size. During a superstep, each process has designated some set of outgoing messages and is expecting to receive some set of incoming messages. If the maximum number of incoming or outgoing messages per processor is  $h$ , then such a communication pattern is called an  $h$ -relation. The communication pattern in Figure 1 is a 2-relation.

Many communication topologies deliver almost all message patterns well, but perform badly for a particular, small set of patterns. The patterns in this set are typically regular ones. In other words, a random message pattern is unlikely to be in this set of 'bad' patterns unless it has some regular structure. One of the attractions of adaptive routing techniques is that they reduce the likelihood of such 'bad' patterns. BSP randomises the placement of processes on processors so that regularities from the problem domain, which are often reflected in programs, are destroyed in the implementation. This tends to make the destination processor addresses of an  $h$ -relation approximate a random permutation. This, in turn, makes it unlikely that each  $h$ -relation will be a 'bad' pattern. The performance advantage of avoiding patterns that take the network a long time to deliver outweighs any advantage gained by exploiting locality in placement.

The ability of a communication network to deliver data is captured by a BSP parameter,  $g$ , that measures the permeability of the network to continuous traffic addressed to uniformly-random destinations. As we have seen, BSP programs randomise to approximate such traffic. The parameter  $g$  is defined such that an  $h$ -relation will be delivered in time  $hg$ . Subject to some small provisos, discussed later,  $hg$  is an accurate measure of communication performance over a large range of architectures. The value of  $g$  is normalised with respect to the clock rate of each architecture so that it is in the same units as the time for executing sequences of instructions.

Sending a message of length  $m$  clearly takes longer than sending a message of size 1. For reasons that will become clear later, BSP does not distinguish between a

message of length  $m$  and  $m$  messages of length 1 – the cost in either case is  $mhg$ . So messages of varying lengths may either be costed using the form  $mhg$  where  $h$  is the number of messages, or the message lengths can be folded into  $h$ , so that it becomes the number of units of data to be transferred.

The parameter  $g$  is related to the bisection bandwidth of the communication network but they are not equivalent –  $g$  also depends on factors such as:

1. the protocols used to interface with, and within, the communication network;
2. buffer management by both the processors and the communication network;
3. the routing strategy used in the communication network; and
4. the BSP runtime system.

So  $g$  is bounded below by the ratio of  $p$  to the bisection bandwidth, suitably normalised, but may be much larger because of these other factors. Only a very unusual network would have a bisection bandwidth that grew faster than  $p$ , so  $g$  is a monotonically increasing function of  $p$ . The precise values of  $g$  is, in practice, determined empirically for each parallel computer, by running suitable benchmarks. A BSP benchmarking protocol is given in Appendix B.

Note that  $g$  is not the single-word delivery time, but the single-word delivery time under continuous traffic conditions. This difference is subtle but crucial.

## 5 Surely This Isn't a Very Precise Measure of How Long Communication Takes? Don't Hotspots and Congestion Make It Very Inaccurate?

One of the most difficult problems of determining the performance of conventional messaging systems is precisely that congestion makes upper bounds hard to determine and quite pessimistic. BSP largely avoids this difficulty.

An apparently-balanced communication pattern may always generate hotspots in some region of the interconnection network. BSP prevents this in several ways. First, the random allocation of processes to processors breaks up patterns arising from the problem domain. Second, the BSP runtime system uses routing techniques that avoid localized congestion. These include randomized routing [37], in which particular kinds of randomness are introduced into the choice of route for each communication action, and adaptive routing [4], in which data are diverted from their normal route in a controlled way to avoid congestion. If congestion occurs, as when an architecture has only a limited range of deterministic routing

techniques for the BSP runtime system to choose from, this limitation on continuous message traffic is reflected in the measured value of  $g$ .

Notice also that the definition of an  $h$ -relation distinguishes the cost of a balanced communication pattern from one that is skewed. A communication pattern in which each processor sends a single message to some other (distinct) processor counts as a 1-relation. However, a communication pattern that transfers the same number of messages, but in the form of a broadcast from one processor to all of the others, counts as a  $p$ -relation. Hence, unbalanced communication, which is the most likely to cause congestion, is charged a higher cost. Thus the cost model does take into account congestion phenomena arising from the limits on each processor's capacity to send and receive data, and from extra traffic that might occur on the communication links near a busy processor.

Experiments have shown that  $g$  is an accurate measure of the cost of moving large amounts of data on a wide range of existing parallel computers. The reason that  $g$  works so well is that, while today's interconnection networks do have non-uniform latencies, these are quite flat. Once a message has entered the network, the latency to an immediate neighbour is not very much smaller than the latency to the other side of the network. Almost all of the end-to-end latency arises on the path from the processor to the network itself, and is caused by operating system overheads, protocol overheads, and limited bandwidth into the network.

## 6 Isn't It Expensive to Give up Locality?

There will always be application domains where exploiting locality is the key to achieving good performance. However, there are not as many of them as a naive analysis might suggest.

There are two reasons why locality is of limited importance. The first is that the communication networks of today's parallel computers seldom have the regular topologies that are often assumed. They are far more likely to have a hierarchical, cluster-based topology (the important exceptions being the Cray T3D and T3E which have a torus topology). Hence each processor has a few neighbours in its cluster, a lot more neighbours slightly further away, and then all of the other nodes at the same effective distance. Furthermore, these distances vary only slightly. So there is just not much advantage to locality in the architecture, since it makes very little difference to latencies once in the network.

The second reason why locality is of limited importance is that most performance-limited problems work with large amounts of data, and can therefore exploit large amounts of virtual parallelism. However, most existing

parallel computers have only modest numbers of processors. When highly-parallel programs are mapped to much less parallel architectures, many virtual processes must be multiplexed onto each physical processor by the programmer. Almost all of the locality is lost when this is done, unless the application domain is highly-regular and matches the structure of the communication topology very closely. Most interesting applications have locality arising from the three-dimensional nature of the world, while most communication networks have two-dimensional locality. For example, finite element applications typically triangulate a three-dimensional surface, and there is no obvious way to map such triangulations onto, say, a 2D torus, while preserving all of the locality. So, while there are applications where locality can be exploited, they are, in practice, less frequent than is commonly supposed.

**7 Most Parallel Computers Have a Considerable Cost Associated with Starting up Communicaton. Doesn't This Mean that the Cost Model Is Inaccurate for Small Messages, Since  $g$  Doesn't Account for Start-up Costs?**

The cost model can be inaccurate, but only in rather special circumstances. Recall that all of the communications in a superstep are regarded as taking place at the end of the superstep. This semantics makes it possible for implementations to wait until the end of the computation part of each superstep to begin the communication actions that have been requested. The implementation can then package the data to be transferred into larger message units. The cost of starting up a data transfer is thus only paid once per destination per superstep.

However, if the **total** amount of communication in a superstep is small, then start-up effects may make a noticeable difference to the performance. We address this quantitatively later.

**8 Aren't Barrier Synchronisations Expensive? How Are Their Costs Accounted for?**

Barriers are often expensive on today's architectures. The reasons can usually be traced back to naive implementations based on, say, trees of pairwise synchronisations, which are themselves expensive on most machines because of poor implementations of semaphores and locks [16]. There is nothing inherently expensive about barriers, and there are signs that future architecture developments will make them much cheaper.

The cost of a barrier synchronisation comes in two parts:

1. The cost caused by the variation in the completion times of the computation steps that participate. There is not much that an implementation can do about this, but it does suggest that balance in the computation parts of a superstep is a good thing.
2. The cost of reaching a globally-consistent state in all of the processors. This depends, of course, on the communication network, but also on whether or not special-purpose hardware is available for synchronizing, and on the way in which interrupts are handled by processors.

For each architecture, the cost of a barrier synchronisation is captured by a parameter,  $l$ . The diameter of the communication network, or at least the length of the longest path that allows state to be moved from one processor to another clearly imposes a lower bound on  $l$ . However, it is also affected by many other factors, so that, in practice, an accurate value of  $l$  for each parallel architecture is obtained empirically.

Notice that barriers, although potentially costly, have a number of attractive features. They make it possible for communication and synchronisation to be logically separated. Communication patterns can no longer accidentally introduce circular state dependencies, so there is no possibility of deadlock or livelock in a BSP program. This makes software easier to build and to understand, and completely avoids the complex debugging needed to find state errors in traditional parallel programs. Barriers also permit novel forms of fault tolerance.

**9 How Do These Parameters Allow the Cost of Programs to Be Determined?**

The cost of a single superstep is the sum of three terms: the (maximum) cost of the local computations on each processor, the cost of the global communication of an  $h$ -relation, and the cost of the barrier synchronisation at the end of the superstep. Thus the cost is given by

$$\text{cost of a superstep} = \text{MAX}_{\text{processes}} w_i + \text{MAX}_{\text{processes}} h_i g + l,$$

where  $i$  ranges over processes, and  $w_i$  is the time for the local computation in process  $i$ . Often the maxima are assumed and BSP costs are expressed in the form  $w + hg + l$ . The cost of an entire BSP program is just the sum of the cost of each superstep. We call this the standard cost model. At this point we emphasize that the standard cost model is not simply a theoretical construct. It provides an accurate model for the cost of real programs of all sizes, across a wide range of real parallel computers. Hill *et al.* [18] illustrates the use of the cost model to predict the cost of a computational fluid dynamics code running on

one architecture when it is moved to another. In contrast, [33] uses the cost model to compare the predicted and actual speedup of an electromagnetics application.

To make this summation of costs meaningful, and to allow comparisons between different parallel computers, the parameters  $w$ ,  $g$ , and  $l$  are expressed in terms of the basic instruction execution rate,  $s$ , of the target architecture. Since this will only vary by a constant factor across architectures, asymptotic complexities for programs are often given unless the constant factors are critically important. Note that we are assuming that the processors are homogeneous, although it is not hard to avoid that assumption by expressing performance factors in any common unit.

The existence of a cost model that is both tractable and accurate makes it possible to truly design BSP programs, that is to consciously and justifiably make choices between different implementations of a specification. For example, the cost model makes it clear that the following strategies should be used to write efficient BSP programs:

1. balance the computation in each superstep between processes, since  $w$  is a *maximum* over computation times, and the barrier synchronisation must wait for the slowest process;
2. balance the communication between processes, since  $h$  is a *maximum* over fan-in and fan-out of data; and
3. minimise the number of supersteps, since this determines the number of times  $l$  appears in the final cost.

The cost model also shows how to predict performance across target architectures. The values of  $p$ ,  $w$ , and  $h$  for each superstep, and the number of supersteps can be determined by inspection of the program code, subject to the usual limits on determining the cost of sequential programs. Values of  $g$  and  $l$  can then be inserted into the cost formula to estimate execution time before the program is executed. The cost model can be used

1. as part of the design process for BSP programs;
2. to predict the performance of programs ported to new parallel computers; and
3. to guide buying decisions for parallel computers if the BSP program characteristics of typical workloads are known.

Other cost models for BSP have been proposed, incorporating finer detail. For example, communication and computation could conceivably be overlapped, giving a superstep cost of the form

$$\max(w, hg) + l,$$

although this optimisation is not usually a good idea on today's architectures [17, 32]. It is also sometimes argued that the cost of an  $h$ -relation is limited by the time taken to send  $h$  messages and then receive  $h$  messages, so that the communication term should be of the form

$$(h_{\text{in}} + h_{\text{out}})g.$$

All of these variations alter costs by no more than small constant factors, so we will continue to use the standard cost model in the interests of simplicity and clarity.

A more important omission from the standard cost model is any restriction on the amount of memory required at each processor. While the existing cost model encourages balance in communication and limited barrier synchronisation, it encourages profligate use of memory. An extension to the cost model to bound the memory associated with each processor is being investigated.

The cost model also makes it possible to use BSP to design algorithms, not just programs. Here the goal is to build solutions that are optimal with respect to total computation, total communication, and total number of supersteps *over the widest possible range of values of  $p$* . Designing a particular program then becomes a matter of choosing among known algorithms for those that are optimal for the range of machine sizes envisaged for the application.

For example two BSP algorithms for matrix multiplication have been developed. The first, a block parallelization of the standard  $n^3$  algorithm [26], has (asymptotic) BSP complexity

$$\text{Block MM cost} = n^3/p + (n^2/p^{1/2})g + p^{1/2}l,$$

requiring memory at each processor of size  $n^2/p$ . This is optimal in computation time and memory requirement.

A more sophisticated algorithm (McCull and Valiant [23]) has BSP complexity

$$\text{Block and Broadcast MM cost} = n^3/p + (n^2/p^{2/3})g + l,$$

requiring memory at each processor of size  $n^2/p^{2/3}$ . This is optimal in time, communication, and supersteps, but requires more memory at each processor. Therefore the choice between these two algorithms in an implementation may well depend on the relationship between the size of problem instances and the memory available on processors of the target architecture.

## 10 Is BSP a Programming Discipline, or a Programming Language, or Something else?

BSP is a model of parallel computation. It is concerned with high-level structure of computations. Therefore it

**Table 1. Core BSP Operations**

Class	Operation	Meaning
Initialisation	<code>bsp_init</code>	Simulate dynamic processes
	<code>bsp_begin</code>	Start of SPMD code
	<code>bsp_end</code>	End of SPMD code
Enquiry	<code>bsp_pid</code>	Find my process id
	<code>bsp_nprocs</code>	Number of processes
	<code>bsp_time</code>	Local time
Synchronisation	<code>bsp_sync</code>	Barrier synchronisation
DRMA	<code>bsp_pushregister</code>	Make region globally visible
	<code>bsp_popregister</code>	Remove global visibility
	<code>bsp_put</code>	Push to remote memory
	<code>bsp_get</code>	Pull from remote memory
BSMP	<code>bsp_set_tag_size</code>	Choose tag size
	<code>bsp_bsmpt_info</code>	Number of packets in queue
	<code>bsp_send</code>	Send to remote queue
	<code>bsp_get_tag</code>	Get tag of 1st message
	<code>bsp_move</code>	Fetch from queue
Halt	<code>bsp_abort</code>	One process halts all
High Performance	<code>bsp_hpput</code>	Unbuffered versions of communication primitives
	<code>bsp_hpget</code>	
	<code>bsp_hpmove</code>	

does not prescribe the way in which local computations are carried out, nor how communication actions are expressed. All existing BSP languages are imperative, but there is no intrinsic reason why this need be so.

BSP can be expressed in a wide variety of programming languages and systems. For example, BSP programs could be written using existing communication libraries such as PVM [9], MPI [27], or Cray’s SHMEM. All that is required is that they provide non-blocking communication mechanisms and a way to implement barrier synchronisation. Indeed, experienced programmers may already find themselves writing in a style reminiscent of BSP precisely to avoid the deadlock potential of the unrestricted message passing style.

There are two advantages to explicitly adopting the BSP framework. First, the values of  $g$  and  $l$  depend not only on the hardware performance of the target architecture but also on the amount of software overhead required to achieve the necessary behaviour. Systems not designed with BSP in mind may not deliver good values of  $g$  and  $l$ . Second, use of the cost model as a design tool can guide software development and increase confidence that good choices have been made.

The most common approach to BSP programming is SPMD imperative programming using Fortran or C, with BSP functionality provided by library calls. Two BSP libraries have been in use for some years: the Oxford BSP

Library [26] and the Green BSP Library [11, 12]. A standard has recently been agreed for a library called BSPLib [13]. BSPLib contains operations for delimiting supersteps, and two variants of communication, one based on direct memory transfer, and the other on buffered message passing.

Other BSP languages have been developed. These include GPL [24] and Opal [21].

### 11 How Easy Is It to Program Using the BSPLib Library?

The BSPLib library provides the operations shown in Table 1. There are operations to:

1. set up a BSP program;
2. discover properties of the environment in which each process is executing;
3. communicate, either directly into or out of a remote memory, or using a message queue;
4. participate in a barrier synchronisation;
5. abort a computation from anywhere inside it; and
6. communicate in a high-performance unbuffered mode.

The BSPLib library is freely available in both Fortran and C from <http://www.bsp-worldwide.org/implmnts/oxtool.htm>. A more complete description of the library can be found in Appendix A.

Another higher-level library provides specialised collective-communication operations. These are not considered as part of the core library, but they can be easily realised in terms of the core. These include operations for broadcast, scatter, gather, and total exchange.

## 12 In what Application Domains Has BSP Been Used?

BSP has been used in a number of application areas, primarily in scientific computing. Much of this work has been done as part of contracts involving Oxford Parallel (<http://www.comlab.ox.ac.uk/oxpara/>).

Computational fluid dynamics applications of BSP include:

- an implementation of a BSP version of the OPlus library for solving 3D multigrid viscous flows, used for computation of flows around aircraft or complex parts of aircraft in a project with Rolls Royce [6];
- a BSP version of FLOW3D, a computational fluid dynamics code;
- oil reservoir modelling in the presence of discontinuities and anisotropies in a project with Schlumberger Geoquest Ltd.

Computational electromagnetics applications of BSP [30] include:

- 3D modelling of electromagnetic interactions with complex bodies using unstructured 3D meshes, in a project with British Aerospace;
- parallelisation of the TOSCA, SCALA, and ELEKTRA codes, and demonstrations on problems such as design of electric motors and permanent magnets for MRI imaging;
- a parallel implementation of a time domain electromagnetic code ParEMC3d with absorbing boundary conditions;
- parallelisation of the EMMA-T2 code for calculating electromagnetic properties of microstrips, wires and cables, and antennae [33].

BSP has been used to parallelise the MERLIN code in a project with Lloyds Register of Shipping and Ford Motor Company. It has been applied to plasma simulation at Rensselaer Polytechnic Institute in New York [31]. It is being used to build neural network systems for data mining at Queen's University in Kingston, Canada.

## 13 What Do BSP Programs Look Like?

Most BSP programs for real problems are large and it is impractical to include their source here. Instead we include some small example programs to show how the BSPLib interface can be used. We illustrate some different possibilities using the standard parallel prefix or scan operation: given  $x_0, \dots, x_{p-1}$  (with  $x_i$  stored on process  $i$ ), compute  $x_0 + \dots + x_i$  on each process  $i$ .

### All Sums: Version 1

The function `bsp_allsums1` calculates the partial sums of  $p$  integers stored on  $p$  processors. The algorithm uses the logarithmic technique that performs  $\lceil \log p \rceil$  supersteps, such that during the  $k$ th superstep, the processes in the range  $2^{k-1} \leq i < p$  each combine their local partial sums with process  $i - 2^{k-1}$ . Figure 2 shows the steps involved in summing the values `bsp_pid() + 1` using 4 processors.

```
int bsp_allsums1(int x) {
    int i, left, right;
    bsp_pushregister(&left, sizeof(int));
    bsp_sync();

    right = x;
    for(i=1; i<bsp_nprocs(); i*=2) {
        if (bsp_pid()+i < bsp_nprocs())
            bsp_put(bsp_pid()+i, &right, &left,
                    0, sizeof(int));
        bsp_sync();
        if (bsp_pid()>=i) right=left+right;
    }
    bsp_popregister(&left);
    return right;
}
```

A process called *registration* is used to enable references to a data structure on one processor to be correctly mapped to locations on other processors. BSPLib does not assume that processors are homogeneous. In any case, heap-allocated data structures need not have the same addresses on different processors, so some mechanism for associating names to addresses is required. The procedure

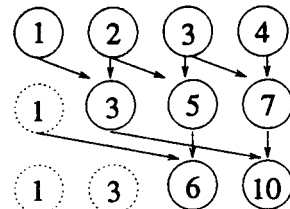


FIGURE 2 All sums using the logarithmic technique.



`bsp_pushregister` allows all processors to declare that the variable `left` is willing to have data put into it by a DRMA operation.

When

```
bsp_put(bsp_pid()+i, &right, &left,
        0, sizeof(int))
```

is executed on process `bsp_pid()`, then a single integer `right` is copied into the memory of processor `bsp_pid()+i` at the address `&left+0`.

The cost of the algorithm is  $\lceil \log p \rceil(1 + g + l) + l$  as there are  $\lceil \log p \rceil + 1$  supersteps (including one for registration); during each superstep a local addition is performed (which costs 1 flop), and at most one message of size 1 word enters and exits each process.

### All Sums: Version 2

An alternative implementation of the prefix sums function can be achieved in a single superstep by using a temporary data structure containing up to  $p$  integers. Each process  $i$  puts the data to be summed into the  $i$ th element of the temporary array on processes  $j$  (where  $0 \leq j \leq i$ ). After all communications have been completed, a local sum is then performed on the accumulated data. The cost of the algorithm is  $p + pg + 2l$ .

```
int bsp_allsums2(int x) {
    int i, result, *array =
        calloc(bsp_nprocs(), sizeof(int));
    if (array==NULL)
        bsp_abort("Unable to allocate %d
                  element array", bsp_nprocs());
    bsp_pushregister(array, bsp_nprocs()
                    *sizeof(int));
    bsp_sync();

    for(i=bsp_pid(); i<bsp_nprocs(); i++)
        bsp_put(i, &x, array, bsp_pid()
                *sizeof(int), sizeof(int));
    bsp_sync();

    result = array[0];
    for(i=1; i<=bsp_pid(); i++)
        result += array[i];
    free(array);
    bsp_popregister(array);
    return result;
}
```

The first algorithm performs a logarithmic number of additions and supersteps, while the second algorithm performs a linear number of additions but a constant number of supersteps. If the operation being performed at each iteration of the algorithm were changed from addition to another, more costly, associative operator, then BSP cost analysis provides a simple mechanism for determining which is the better implementation.

### All Sums on an Array

Either of the routines defined above can be used to sum  $n$  values held in  $n/p$  blocks distributed among  $p$  processors. The algorithm proceeds in four phases:

1. The running sum of each  $n/p$  block of integers is computed locally on each processor.
2. As the last element of each  $n/p$  block contains the sum of each  $(n/p)$ -element segment, then either of the two simple algorithms can be used to calculate the running sums of the last element in each block (call this `last`).
3. Each processor gets the value of `last` from its left neighbouring processor (we call this `lefts_last`).
4. Adding `lefts_last` to each of the locally-summed  $n/p$  elements produces the desired effect of the running sums of all  $n$  elements.

```
void bsp_allsums(int*array,
                 int n_over_p) {
    int i, last, lefts_last;
    bsp_pushregister(&last, sizeof(int));

    for (i=1; i<n_over_p; i++)
        array[i] += array[i-1];

    last = bsp_allsums2
           (array[n_over_p-1]);

    if (bsp_pid()==0) lefts_last=0;
    else
        bsp_get(bsp_pid()-1, &last, 0,
                &lefts_last, sizeof(int));
    bsp_sync();
    for(i=0; i<n_over_p; i++)
        array[i] += lefts_last;

    bsp_popregister(&last);
}

void main() {
    int i, j, n_over_p, *xs;
    bsp_begin(bsp_nprocs());

    n_over_p = 100;
    xs = calloc(n_over_p, sizeof(int));
    for (i=0; i<n_over_p; i++) xs[i]=1;
    bsp_allsums(xs, n_over_p);

    for(i=0; i<bsp_nprocs(); i++) {
        if (bsp_pid()==i) {
            printf("On process %d: ",
```

```

        bsp_pid() );
    for (j=0; j<n_over_p; j++)
        printf("%d ", xs[j]);
    printf("\n");
    fflush(stdout);
}
bsp_sync();
}
bsp_end();
}

```

## 14 What Are Typical Values of $g$ and $l$ for Common Parallel Computers?

Values of the BSP cost model parameters are shown in Table 2. The values of the  $g$  and  $l$  parameters are normalised by the instruction rate,  $s$ , of each processor (to aid comparisons between machines, raw rates are also given in microseconds). Because this instruction rate depends heavily upon the kind of computations being done, the average of two different measured values are used:

[ $s$ ] measures the cost of an inner product, where  $\mathcal{O}(n)$  operations are performed on a data structure of size  $n$ . The value of  $n$  is chosen to be far greater than the cache size on each processor. This benchmark therefore gives a lower-bound megaflop rate for the processor as each arithmetic operation induces a cache miss.

[ $s$ ] measures the cost of a dense matrix multiplication, where  $\mathcal{O}(n^3)$  operations are performed on a data structures of size  $n^2$ . Because a large percentage of the computation can be kept in cache, this benchmark gives an upper-bound megaflop rate for the processor.

As we have already mentioned, good BSP algorithm design is often based around balanced patterns of communication. We illustrate the communication capacity,  $g$ , using two balanced communications. The first is a particularly easy 1-relation, a local communication that performs a cyclic shift of data between neighbouring processors. This benchmark provides an upper-bound rate for communication as there are only  $p$  messages injected into the communication network during a superstep.

Parallel computers have far greater difficulty in achieving scalable communication for patterns of communication that move lots of data to many destinations. As an extreme example, we consider the total exchange global communication that injects  $p^2$  messages into the network and realises a  $p$ -relation. As no scalable architecture can provide  $p^2$  dedicated wires because it is too expensive, sparser interconnection networks are used in practice. For

example, the Cray T3D uses a 3D Torus, while the IBM SP2 uses a hierarchy of 8-node fully-connected crossbar switches. The value of  $g$  for a total exchange therefore provides a good measure of the lower-bound rate of communication of an architecture.

Not very surprisingly, the two values of  $g$ , derived directly from a 1-relation, and from the  $pg$  cost of a  $p$ -relation total exchange can be quite different. This might mean that the 1-relation performance of the network is not very good (for example, a ring takes time proportional to  $p$  to deliver both a 1-relation and a  $p$ -relation), but usually means that the network's effective capacity is not as large as the per-link bandwidth would suggest. When cost modelling algorithms, it is advisable to use the value of  $g$  produced by the global communication (total exchange) benchmark.

Appendix B shows how these figures were obtained. The meaning of  $n_{1/2}$  is explained in Section 16.

## 15 How Can the BSPLib Be Implemented Efficiently on Today's Architectures?

The semantics of supersteps separates local computation from communication, and the Oxford implementation of BSPLib keeps these two phases separate in the implementation also. Thus while the semantics of calls to `put` and `get` permits them to begin executing concurrently with the local process's computation, calls to these functions in fact buffer the data for later transfer. Not overlapping computation and communication contradicts conventional wisdom, but it turns out that the performance advantages of postponing communication are larger than of exploiting the potential overlap [17].

We begin by noting that overlapping computation and communication can give at best a factor of two performance improvement, and then only when the computation and communication times are precisely equal. This equality is neither a scalable nor portable property, so we must expect an appropriate balance to be quite rare. Thus the performance improvement factor due to overlapping is likely to be much less than two in practice.

On the other hand, postponing communication is a big performance win because it permits two major optimisations:

1. Combining all of the transfers between a pair of processors into a single messages, so that the overhead of message startup is paid only once. The benefits of doing this are discussed in the next section.
2. Reordering communications so that the load they generate is applied to the communication network effectively, rather than in the order in which the

Table 2. BSP Machine Parameters

Machine	computation			<i>p</i>	barrier		local comm.		global comm.		<i>n</i> <sub>1/2</sub> words
	[ <i>s</i> ]	[ <i>s</i> ]	<i>s</i>		<i>l</i>	<i>g</i>	<i>g/s</i>	<i>g</i>	<i>g/s</i>		
	Mflops				flops	$\mu$ s	flop/word	$\mu$ s/word	flop/word	$\mu$ s/word	
SGI PowerChallenge	53	94	74	1	226	3.1	0.5	0.007	0.5	0.007	80
				2	1132	15.3	9.8	0.13	10.2	0.14	12
				3	1496	20.2	8.9	0.12	9.5	0.13	12
				4	1902	25.7	9.8	0.13	9.3	0.13	12
Cray T3E	4.3	89.2	46.7	1	86	1.8	2.12	0.05	2.14	0.05	9
				2	269	5.7	0.87	0.02	2.61	0.07	33
				3	296	6.3	0.86	0.02	2.11	0.04	35
				4	357	7.6	0.87	0.02	1.77	0.04	40
				8	506	10.8	0.81	0.02	1.64	0.03	40
				9	552	11.7	0.82	0.02	1.57	0.03	42
				16	751	16.0	1.04	0.02	1.66	0.04	38
				20	880	18.7	0.96	0.02	1.63	0.03	38
Cray T3D	5	19	12	24	1013	21.6	1.39	0.03	1.70	0.04	36
				1	68	5.6	0.3	0.02	0.3	0.02	94
				2	164	13.5	0.7	0.06	1.0	0.08	71
				4	168	13.9	0.7	0.06	0.8	0.65	66
				8	175	14.4	0.8	0.07	0.8	0.65	59
				9	383	31.7	0.9	0.07	1.2	0.10	39
				16	181	14.9	0.9	0.07	1.0	0.08	61
				25	486	40.2	1.1	0.09	1.5	0.13	26
IBM SP2 (switch)	25	27	26	32	201	16.6	1.1	0.09	1.4	0.12	28
				64	148	12.3	1.0	0.09	1.7	0.14	27
				128	301	24.9	1.1	0.09	1.8	0.15	20
				256	387	32.1	1.2	0.11	2.4	0.19	15
				1	244	9.4	1.3	0.05	1.3	0.05	7
				2	1903	73.2	6.3	0.24	7.8	0.30	6
				4	3583	137.8	6.4	0.25	8.0	0.31	7
				8	5412	208.2	6.9	0.27	11.4	0.43	6
Multiprocessor Sun	3.8	16.4	10.1	1	24	2.4	0.4	0.04	0.4	0.04	7
				2	54	5.3	3.0	0.29	3.4	0.34	7
				3	74	7.4	2.9	0.29	4.1	0.41	8
				4	118	11.7	3.3	0.32	4.1	0.41	11
Parsytec GC			19.3	1	98	5.1	1.0	0.05	1.0	0.05	16
				2	6309	325	109	5.6	113	5.9	3
				4	23538	1219	190	9.9	143	7.4	3
				8	29080	1506	252	13.1	254	13.2	3
				16	224977	11600	253	13.1	342	17.7	3
IBM SP2 (ethernet)	25	27	26	32	130527	6700	272	14.1	658	34.1	3
				1	241	9.3	1.3	0.05	1.3	0.05	8
				2	18759	721.5	182.1	7.0	183.6	7.1	3
				4	39025	1500.9	388.2	14.9	628.2	24.2	5
				8	88795	3415.2	1246.6	47.3	1224.1	47.1	2

(1) All values for *g* are for communications of 32-bit words; (2) benchmarks were performed at the -O3 optimisation level; (3) the Cray T3D, SGI PowerChallenge, IBM SP2, and Parsytec GC used native implementations of the toolset; (4) the toolset used on the multiprocessor Sun was built using generic System V shared-memory facilities

particular puts and gets appears in the program. Patterns guaranteed to avoid congestion can be set up in software, rather than requiring expensive hardware solutions operating during the data transfers.

These results are counter-intuitive, since they appear to increase congestion in the network that could be avoided by allowing some messages to begin transmission early. This effect is undoubtedly present, but it is dwarfed by the size of the improvements which postponement makes possible. Reordering communication, for example, gives performance improvements of a factor between about 2 and  $p$ , while combining multiple transfers into single messages can give improvements of several orders of magnitude. One reason why this tradeoff has not been noted previously is that message-passing interfaces that operate at the level of single messages cannot naturally conceive of postponing transmission since there is no clear moment to postpone transmission to.

The performance gains of delaying communication are so large that even the high-performance versions of the put and get operations, which are designed so that computation and communication can be overlapped without buffering, postpone transmissions until the end of the computation phase of each superstep. Congestion within the network is much less important, in practice, than congestion at the network boundaries. A processor that simultaneously receives messages from several other processors has no choice but to sequentialise their removal from the network.

Regardless of the type of parallel architecture, the ability to reorder messages before transmission is crucial to creating a consistent bulk-communication behaviour without increasing the value of  $g$ . Two mechanisms used are:

1. randomly ordering the messages to reduce the likelihood of troublesome patterns, and
2. using a latin square to schedule transmissions in a guaranteed contention-free way.

Which of these mechanisms is to be preferred is architecture-dependent.

Recall that a latin square is a  $p \times p$  square in which each of the values from 1 to  $p$  appears  $p$  times, with no repetition in any row or column. Such a square can be used as a schedule for the routing of the  $h$ -relation, using row  $i$  as the schedule for processor  $i$ , with the contents of the row regarded as the destinations for each communication time step.

The use of such mechanisms has a major effect on performance. For example, consider a total exchange algorithm shown in Figure 3 where each processor  $i$  has

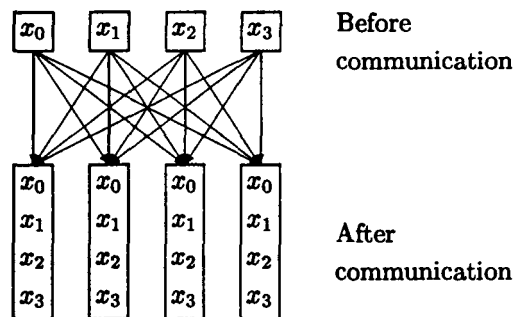


FIGURE 3 Total exchange between four processors.

data  $x_i$  of size  $n$  that is to be exchanged with every other processor. After the communication, each processor will contain a data structure of size  $np$  containing all of the  $x_j$ , where  $1 \leq j < p$ . The BSP cost of the algorithm is  $png + l$  because  $p$  messages enter and exit each processor. However, a naive implementation may have each processor send a message to processor 0 on the first time step, to processor 1 on the second, and so on. This causes  $p$  messages to contend at process 0, then  $p$  to contend at process 1, and so on. The cost of this communication will be  $\mathcal{O}(p^2)$  rather than the linear cost predicted by the BSP cost formula  $png + l$ . An alternative ordering that does not cause contention is for processors to send their data in the order  $\text{mod}(i + j, p)$ ; where  $1 \leq j < p$ , and  $i$  is the processor identifier, using a simple latin square. The expected linear (in  $p$ ) cost can then be achieved.

Table 3 shows the results of an implementation that routes total exchanges. Column 1 shows the performance of a system in which messages are despatched as soon as the puts are encountered, and in which the order of the puts causes contention. The second column shows the performance when messages are immediately despatched, but the programmer has carefully hand-crafted the order of puts to minimise contention. The third and fourth columns show the performance when both of these programs are run with puts postponed until the end of the superstep and reordered by the runtime system using a latin square. The performance is very slightly worse than the best hand-coded program, because of the overhead of the runtime system managing the reordering. Far more importantly, the effect of the programmer's ordering of the puts has been completely removed. In other words, reordering provides consistent performance over varying orderings of the data transfer instructions, at the expense of a very small decrease in best case performance. Note also that reordering provides almost a factor of two performance improvement, enough by itself to make up for any performance loss caused by not overlapping computation and communication.

The precise details of handling communication and building barriers differs depending on the specifics of target architectures:

**Table 3. The Effects of Node Contention on the Cray T3D.** Entries in the table are in seconds for routing a 4,000,000-relation, e.g., for 128 processors, 15625 integers per process

Procs	immediate transmission		BSPLib reordering	
	contention	latin square	contention	latin square
2	.168	.157	.157	.157
4	.392	.194	.191	.191
8	.461	.239	.228	.229
16	.598	.289	.344	.345
32	.784	.413	.465	.456
64	.903	.529	.548	.546
128	.961	.575	.599	.599

*Distributed-memory machines with remote-memory access (Cray T3D and Cray T3E).* A barrier synchronisation is performed to ensure that each process has finished its local computation. Once all the processors have passed the barrier, one-sided memory accesses are used to route messages into the memories of the remote processors. Combining is not used, because there is little to be gained when the actual data transfer mechanism is DRMA. The communication phase of a superstep is completed by performing a further barrier synchronisation.

*Distributed-memory machines with message-passing (IBM SP2, Parsytec GC, Generic TCP/IP).* On architectures that provide native non-blocking send and blocking receive message-passing primitives, the  $h$ -relation is routed through the communication network in three phases:

1. a total exchange is performed, exchanging information about the number, sizes, and destination addresses of messages. This total exchange is considered to be the barrier synchronisation for the superstep.
2. gets are translated into puts and the data they refer to is buffered at the source processor.
3. after the total exchange, each processor knows how many messages, from every other process, it is expecting. Each process therefore knows when the communication phase of the superstep is complete by counting the incoming messages. Communication is performed by interleaving the outgoing and incoming messages, so that minimum buffering requirements are placed on the underlying message-passing system.

*Shared-memory architectures (SGI Power Challenge, Sun).* The implementation on shared-memory architectures combines features from both of the implementations above. The information about the number and size of messages to be sent between each processor pair is constructed in a region of shared memory by each call to put and get. After the computation phase, a barrier synchronisation takes place to ensure that this information is frozen. Because the message information is in shared memory, an implicit total exchange can be considered to have occurred at this point. The actual exchange of data is performed in a message-passing style. First messages are copied into buffers associated with each process in shared memory. These buffers are then inspected by the remote process, and their contents copied into the remote processor's memories. Using a contention-limiting order for messages, the number of message passing buffers associated with each process can be minimised. Finally, the message information region is cleared and a further barrier synchronisation takes place to allow renewed access to it.

### 16 How Much Effect Does Message Size Have on the Value of $g$ ?

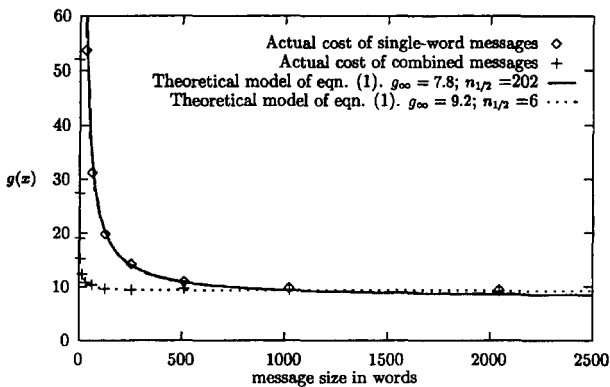
As we have already seen, the way in which BSPLib delays communication until the end of each superstep and then combines messages into the largest possible units reduces the importance of message size. The cost model makes no distinction between the cost of a process sending  $h$  messages of size one or a single message of size  $h$ ; both communications have an  $h$ -relation cost of  $hg$ . However, a superstep in which very little *total* communication occurs may still deviate from the cost model because of the effects of startup costs for message transmission.

Miller refined the standard cost model [29] using a technique of Hockney [20] to model the effect of message granularity on communication cost. In the refined model,  $g$  is defined as a function of the message size  $x$ :

$$g(x) = \left(\frac{n_{1/2}}{x} + 1\right)g_{\infty}, \tag{1}$$

where  $g_{\infty}$  is the asymptotic communication cost for very large messages (that is, the  $g$  reported in Table 2) and  $n_{1/2}$  is the size of message that produces half the optimal bandwidth of the machine so  $g(n_{1/2}) = 2g_{\infty}$ .

The value of  $n_{1/2}$  in Equation (1) is determined experimentally for each machine configuration by fitting a curve to actual values of  $g(x)$ . Figure 4 shows the actual values of  $g(x)$  on an 8-processor IBM SP2. Because messages are combined in each superstep, the value of  $n_{1/2}$  is effectively reduced to 6 words. For comparison purposes, the effect of naively communicating messages separately



**FIGURE 4** Fitting experimental values of  $g(x)$  flops/word to Equation (1) using an 8-processor IBM SP2 with switch communication. The messages are communicated using one-sided *put* communication where a process puts data into another processor's memory. The top curve represents single-word messages and the bottom curve uses a message-combining scheme.

is shown by the data points labeled "actual cost of single-word messages" in the figure. Fitting a curve to this data gives  $n_{1/2} = 202$  words.

The  $n_{1/2}$  parameter can be used to discover the minimum message size for which the standard cost model is within a given percentage of the more-detailed cost model. For the standard model to be within  $y\%$  accuracy of the cost attributed by the model that includes message granularity, then:

$$\left(\frac{100+y}{100}\right)h_0g_\infty = h_0g(h_0) = \left(\frac{n_{1/2}}{h_0} + 1\right)h_0g_\infty, \quad (2)$$

where  $h_0$  words is Valiant's parameter [36] that measures the minimum size of  $h$ -relation to achieve  $n_{1/2}$  throughput. Thus the percentage error in the communication cost  $h_0g_\infty$  is

$$y = \left(\frac{100n_{1/2}}{h_0}\right)\%. \quad (3)$$

So on the IBM SP2 with switch communication the error in the standard BSP model for communicating  $h_0 = 60$  32-bit words is 10%. Moreover, as would be expected, as the size of  $h$ -relation increases, the error in the standard BSP model decreases.

These data show that combining the messages sent between each pair of processors has a significant effect on the achieved value of  $g$ , and so provides further justification for not overlapping computation and communication.

## 17 What Tools Are Available to Help with Building and Tuning BSP Programs?

The *intensional properties* of a parallel program (i.e., how it computes a result) can often be hard to understand. The BSP model goes some way towards alleviating this problem if cost analysis is used to guide program development. Unfortunately, in large-scale problems, cost analysis is rarely used at the time of program development. The role of current BSP tools [18] is to aid programmers in understanding the intensional properties of their programs by graphically providing profiling and cost information. The tools may be used both to analyse the communication properties of a program, and to analyse the predicted performance of the code on a real machine.

A central problem with any parallel-profiling systems is effective visualisation of large amounts of profiling data. In contrast to conventional parallel-profiling tools, which highlight the patterns of communication between individual sender-receiver pairs in a message passing system, the BSP approach significantly simplifies visualisation because all of the communications from a superstep can be visualised as a single monolithic unit.

Figure 5 is an example of the results from a BSP profiling tool running on the IBM SP2. It shows a *communication* profile for the parallel prefix algorithm (with  $n > p$ ) developed on page 260.

The top and bottom graphs in Figure 5 show, on the  $y$ -axis, the volume of data moved, and on the  $x$ -axis, the elapsed time. Each pair of vertically-aligned bars in the two graphs represents the total communication during a superstep. The upper bars represent the output from processors, and the lower bars the input. Within each communication bar is a series of bands. The height of each band represents the amount of data communicated by a particular process, identified by the band's shade. The sum of all the bands (the height of the bar) represents the total amount of communication during a superstep. The width represents the elapsed time spent in *both* communication and barrier synchronisation. The label found at the top left-hand corner of each bar can be used in conjunction with the legend in the right of the graph to identify the end of each superstep (i.e., the call to `bsp_sync`) in the user's code. The white space in the figure represents the computation time of each superstep.

In Figure 5, the start and end of the running sums is identified by the points labelled 0 and 4. The white space in the graphs between supersteps 0 and 1 shows the computation of the running sums executed locally in each process on a block of size  $n/p$ . The first superstep, which is hidden by the label 1 at this scale, shows the synchronisation that arises due to registration in the function `bsp_allsums1`. The three successively-smaller bars represent the logarithmic number of communication phases of the parallel prefix technique. Contrast-

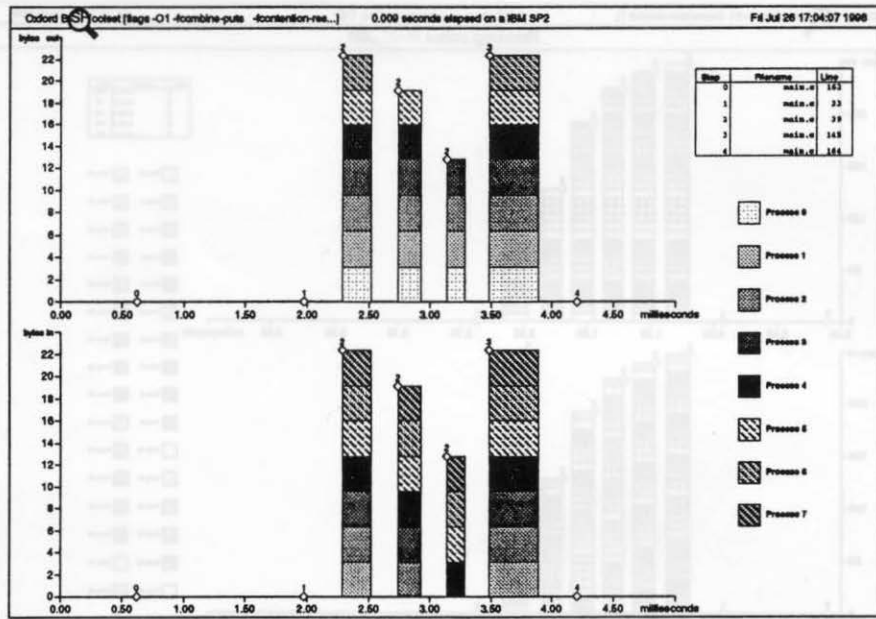


FIGURE 5 All sums of 32,000 elements using the logarithmic technique on an 8-processor IBM SP2.

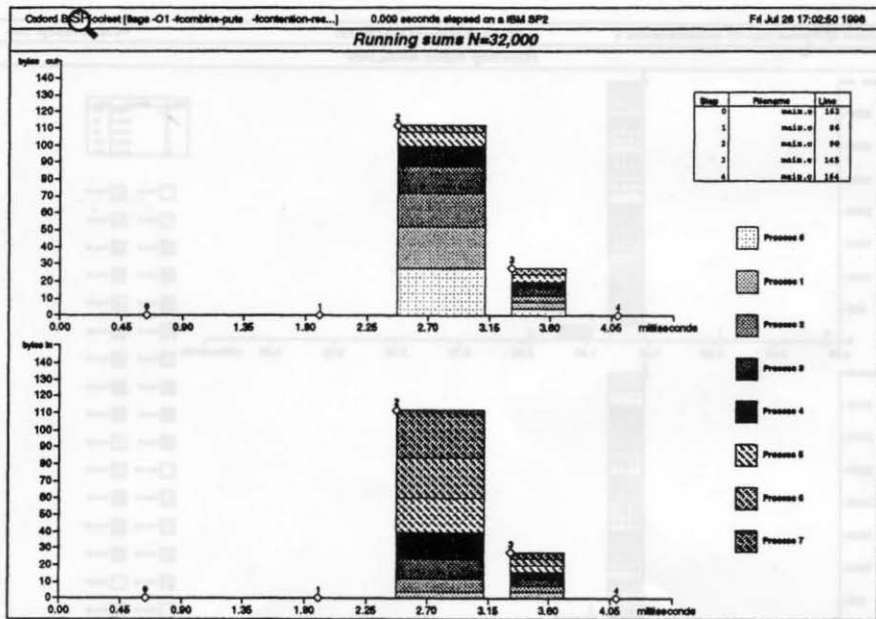


FIGURE 6 All sums of 32,000 elements using total exchange on an 8-processor IBM SP2.

ing the sizes of the communication bars in Figure 5 with the schematic diagram of Figure 2 graphically shows the diminishing numbers of processors involved in communication as the parallel prefix algorithm proceeds. Contrasting this method of running sums with the total-exchange-based algorithm in Figure 6 shows that although the number of synchronisations within the algorithm is reduced from  $\lceil \log p \rceil$  to 1, the time spent in the total exchange of `bsp_allsums2` is approximately the same as the algo-

rithm based upon the logarithmic technique. This is due to the larger amount of data transferred, i.e., 1.51 milliseconds spent in summing  $p$  values in  $p$  processes using the parallel prefix technique, compared to 1.42 milliseconds when the total exchange is used.

Figures 7 and 8 show profiles of the same two algorithms running on a 32-processor Cray T3D, with the same data-set size as the IBM SP2. Although the T3D has a lower value for the barrier synchronisation latency than

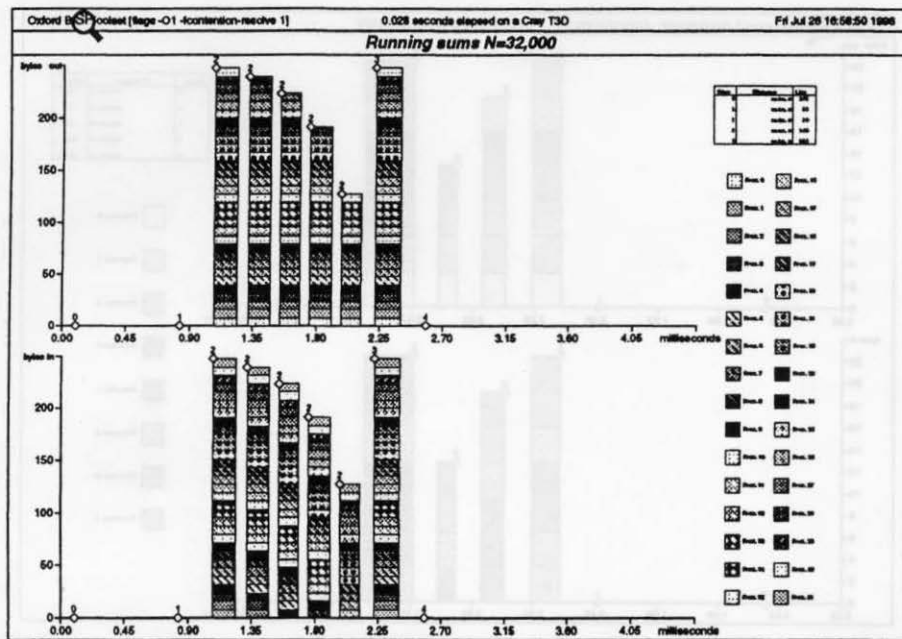


FIGURE 7 All sums of 32,000 elements using the logarithmic technique on a 32-processor Cray T3D.

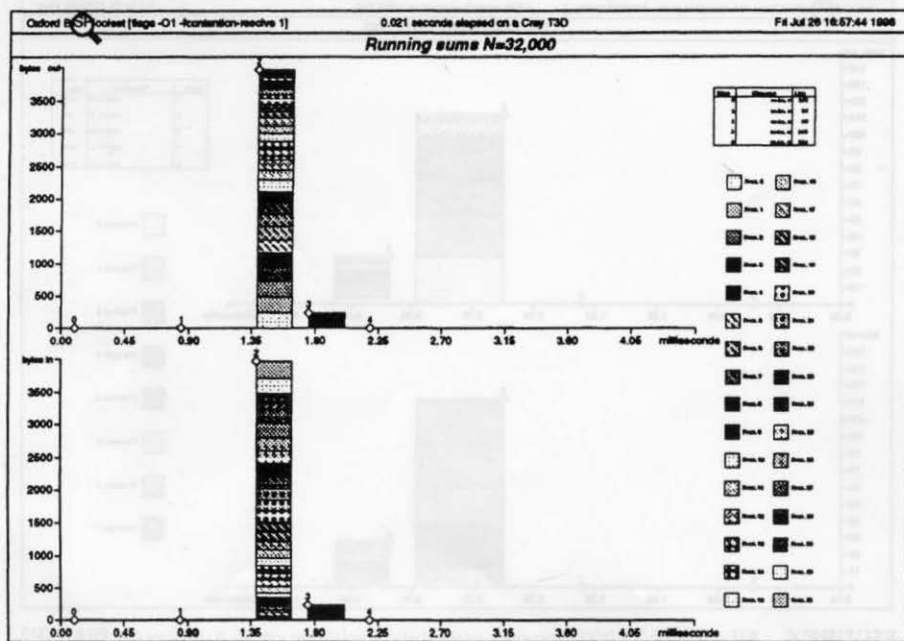


FIGURE 8 All sums of 32,000 elements using a total exchange on a 32-processor Cray T3D.

the IBM SP2 (see Table 2), reducing the number of supersteps from  $\lceil \log 32 \rceil = 6$  supersteps to 1 has a marked effect on the efficiency. The version `bsp_allsums1` (i.e., logarithmic) takes 1.39 milliseconds compared to 0.91 milliseconds for `bsp_allsums2` (i.e., total exchange).

These data show that, for today's parallel computers, it is often better to reduce the number of supersteps, even at the expense of requiring more communication.

## 18 How Does BSPLib Compare with Other Communication Systems such as PVM or MPI?

In recent years, the PVM message-passing library [1, 2, 10] has been widely implemented and widely used. In that respect, the goal of source code portability in parallel computing has already been achieved by PVM. What then, are the advantages of BSP programming, if



any, over a message-passing framework such as PVM? First, PVM and all other message-passing systems based on pairwise, rather than barrier, synchronisation have no simple analytic cost model for performance prediction, and no simple means of examining the global state of a computation for debugging. Second, taking a global view of communication introduces opportunities for optimisation that can improve performance substantially [17] and these are inaccessible to systems such as PVM.

MPI [14] has been proposed as a new standard for those who want to write portable message-passing programs in Fortran and C. At the level of point-to-point communications (send, receive etc.), MPI is similar to PVM, and the same comparisons apply. The MPI standard is very general and is very complex relative to the BSP model. However, one could use some carefully-chosen combination of the various non-blocking communication primitives available in MPI, together with its barrier synchronisation primitive, to produce an MPI-based BSP programming model. At the higher level of collective communications, MPI provides support for various specialised communication patterns which arise frequently in message-passing programs. These include broadcast, scatter, gather, total exchange, reduction, and scan. These standard communication patterns are also provided for BSP in a higher-level library. There have been two comparisons of the performance of BSP and MPI. One by Szymanski on a network of workstations [31] showed performance differences of the order of a few percent. Another by Hyaric (<http://merry.comlab.ox.ac.uk/users/hyaric/doc/BSP/NASfromMPItoBSP>) used the NAS benchmarks. BSP outperformed MPI on four out of five of these, performing ten percent better in some cases. Only on LU did BSP perform about five percent worse.

Compared to PVM and MPI, the BSP approach offers

- (a) a simple programming discipline (based on supersteps) that makes it easier to determine the correctness of programs;
- (b) a cost model for performance analysis and prediction which is simpler and compositional; and
- (c) more efficient implementations on many machines.

### 19 How Is BSP Related to the LogP Model?

LogP [7] differs from BSP in three ways:

- 1. It uses a form of message passing based on pairwise synchronisation.

- 2. It adds an extra parameter representing the *overhead* involved in sending a message. This has the same general purpose as the  $n_{1/2}$  parameter in BSP, except that it applies to *every* communication, whereas the BSP parameter can be ignored except for a few unusual programs.
- 3. It defines  $g$  in local terms. The  $g$  parameter in BSP is regarded as capturing the throughput of an architecture when every processor inserts a message (to a uniformly-distributed address) on every step. It takes no account of the actual capacity of the network, and does not distinguish between delays in the network itself and those caused by inability to actually enter the network (blocking back at the sending processor). In contrast, LogP regards the network as having finite capacity, and therefore treats  $g$  as the minimal permissible *gap* between message sends from a single process. This amounts to the same thing in the end, that is  $g$  in both cases is the reciprocal of the available per-processor network bandwidth, but BSP takes a global view of the meaning of  $g$ , while LogP takes a more local view.

Experience in developing software using the LogP model has shown that, to analyse the correctness and efficiency of LogP programs, it is often necessary, or at least convenient, to use barriers. Also, major improvements in network hardware and in communications software have greatly reduced the overhead associated with sending messages. In early multiprocessors, this overhead could be substantial, since a single processor handled both the application and its communication. Manufacturers have learned that this is a bad idea, and most newer multiprocessors provide either a dedicated processor to handle message traffic at each node or direct remote-memory access. In this new scenario, the only overhead for the application processor in sending or receiving a message is the time to move it from user address space to a system buffer. This is likely to be small and relatively machine-independent, and may even disappear as communication processors gain access to user address space directly. The importance of the overhead parameter in the long term seems negligible.

Given that

$$\text{LogP} + \text{barriers} - \text{overhead} = \text{BSP},$$

the above points would suggest that the LogP model does not improve upon BSP in any significant way. However, it is natural to ask whether or not the more “flexible” LogP model enables a designer to produce a more efficient algorithm or program for some particular problem, at the expense of a more complex style of programming. Recent

results show that this is not the case. In [3] it is shown that the BSP and LogP models can efficiently emulate one another, and that there is therefore no loss of performance in using the more-structured BSP programming style.

## 20 How Is BSP Related to the PRAM Model?

The BSP model can be regarded as a generalisation of the PRAM model which permits the frequency of barrier synchronisation, and hence the demands on the routing network, to be controlled. If a BSP architecture has a very small value of  $g$ , e.g.,  $g = 1$ , then it can be regarded as a PRAM and we can use hashing to automatically achieve efficient memory management. The value of  $l$  determines the degree of parallel slackness required to achieve optimal efficiency. The case  $l = g = 1$  corresponds to the idealised PRAM, where no parallel slackness is required.

## 21 How Is BSP Related to Data Parallelism?

Data parallelism is an important niche within the field of scalable parallel computing. A number of interesting programming languages and elegant theories have been developed in support of the data-parallel style of programming, see, e.g., [34]. High Performance Fortran [22] is a good example of a practical data-parallel language. Data parallelism is particularly appropriate for problems in which locality is crucial.

The BSP approach, in principle, offers a more flexible and general style of programming than is provided by data parallelism. However, the current SPMD language implemented by BSPLib is very much like a large-grain data parallel language, in which locality is not considered and programmers have a great deal of control over partitioning of functionality. In any case, the two approaches are not incompatible in any fundamental way. For some applications, the flexibility provided by the BSP approach may not be required and the more limited data-parallel style may offer a more attractive and productive setting for parallel software development, since it frees the programmer from having to provide an explicit specification of the various processor scheduling, communication and memory management aspects of the parallel computation. In such a situation, the BSP cost model can still play an important role in terms of providing an analytic framework for performance prediction of the data-parallel program.

## 22 Can BSP Handle Synchronisation among a Subset of the Processes?

Synchronising a subset of executing processes is a complex issue because the ability of an architecture to synchronise is not a bulk property in the same sense that its processing power and communication resources are. Certain architecture provide a special hardware mechanism for barrier synchronisation across all of the processors. For example the Cray T3D provides an add-and-broadcast tree, and work at Purdue [8] has created generic, fast, and cheap barrier synchronisation hardware for a wide range of architectures. Sharing this single synchronisation resource among several concurrent subsets that may wish to use it at any time seems difficult. We are currently exploring this issue, but the current version of the library synchronises only across the entire machine.

Architectures in which barrier synchronisation is implemented in software do not have any difficulty in implementing barriers for subsets of the processors. The remaining difficulty here is a language design one – it is not yet clear what an MIMD, subset-synchronising language should be like if it is to retain the characteristics of BSP, such as accurate predictability.

## 23 Can BSP be Used on Vector, Pipelined, or VLIW Architectures?

Nothing about BSP presupposes how the sequential parts of the computation, that is the processes within each processor, are computed. Thus architectures in which the processor uses a specialised technique to improve performance might make it harder to determine the value of  $w$  for a particular program, but they do not otherwise affect the BSP operation or cost modelling. The purpose of normalising  $g$  with respect to processor speed is to enable terms of the form  $hg$  to be compared to computation times so that the balance between computation and communication in a program is obvious. Architectures that issue multiple instructions per cycle might require a more sophisticated normalisation to keep these quantities comparable in useful ways.

## 24 BSP Doesn't Seem to Model Either Input/Output or Memory Hierarchy?

Both properties can be modelled as part of the cost of executing the computation part of a superstep. Modelling the latency of deep storage hierarchies fits naturally into BSP's approach to the latency of communication, and investigations of extensions to the BSP cost model applicable to databases are underway [35].

## 25 Does BSP Have a Formal Semantics?

Several formal semantics for BSP have been developed. He *et al.* [15] show how these may be used to give algebraic laws for developing BSP programs. BSP is used as a semantics case study in a forthcoming book [19].

## 26 Will BSP Influence the Design of Architectures for the Next Generation of Parallel Computers?

The contribution of BSP to architecture design is that it clarifies those factors that are most important for performance on problems without locality. It suggests that the critical properties of an architecture are:

1. high permeability of the communication system, that is the ability to move arbitrary patterns of data quickly; and
2. the ability to reach a consistent global state quickly by barrier synchronisation.

More subtly, it also suggests that predictability of delivery for a wide range of communication patterns is more important than high performance for some special communication patterns, and low performance for others. In other words, low variance is more significant than low mean.

The two parameters  $l$  and  $g$  capture, in a direct way, how well an architecture achieves these two major performance properties. Details of exactly which topology to use, what routing technology, and what congestion control scheme are all subsumed in the single consideration of total throughput.

When the BSP model was first considered, it was often felt to be necessarily inefficient because of its use of permutation routing. After a while, it came to be appreciated that permutation routing is not necessarily expensive, and architectures that do it well were developed. Next the BSP model was considered inefficient because of its requirement for barrier synchronisation. It is now understood that barriers need not be expensive, and architectures that handle them well are being developed. It may be that total exchange is the next primitive to be made central to BSP and the same arguments about its necessary inefficiency may well be made. New communication technologies, such as ATM, repay foreknowledge of communication patterns, and total exchange may turn out to be a reasonable standard building block for parallel architectures as well.

BSP's structured use of machine resources also suggests functions that could be usefully migrated to hardware. We have already seen this possibility for barrier synchronisation. Hardware support for message combining and scheduling would appear to be cost-effective also.

## 27 How Can I Find out More about BSP?

Development of BSP is coordinated by *BSP Worldwide*, and organisation of researchers and users. Information about it can be found at the web site <http://www.bsp-worldwide.org/>. The BSPLib library described here is a BSP Worldwide standard. Other general papers about BSP are [23, 36].

There are groups of BSP researchers at:

1. Oxford – <http://www.comlab.ox.ac.uk/oucl/groups/bsp/>;
2. Harvard – <http://das-www.harvard.edu/cs/research/bsp.html>;
3. Utrecht – <http://www.math.ruu.nl/people/bisseling.html>;
4. Carleton – <http://www.scs.carleton.ca/~palepu/BSP.html>;
5. Central Florida – <http://longwood.cs.ucf.edu/csdept/faculty/goudreau.html>;

as well as individuals working on BSP at a number of other universities.

## ACKNOWLEDGEMENTS

We appreciate the helpful comments made on earlier drafts of this paper by David Burgess, Dave Dove, Gaétan Hains, Jifeng He, Owen Rogers, Heiko Schröder, Bolek Szymanski, and Alexandre Tiskin.

D. B. Skillicorn was supported in part by EPSRC Research Grant GR/K63740 “A Unified Framework for Parallel Programming”.

J. M. D. Hill and W. F. McColl were supported in part by EPSRC Research Grant GR/K40765 “A BSP Programming Environment”.

## REFERENCES

- [1] A. Beguelin, J. Dongarra, A. Geist, R. Manchek, and V. Sunderam, “Recent enhancements to PVM,” *International Journal of Supercomputing Applications and High Performance Computing*, 1995.
- [2] A. Beguelin, J. Dongarra, A. Geist, R. Manchek, and V. Sunderam, “A users’ guide to PVM parallel virtual machine,” University of Tennessee, Tech. Rep. CS-91-136, July 1991.
- [3] G. Bilardi, K. T. Herley, A. Pietracaprina, G. Pucci, and P. Spirakis, “BSP vs LogP,” in *Proc. 8th Ann. Symp. Parallel Algorithms and Architectures*, June 1996, pp. 25–32.
- [4] R. V. Boppana and S. Chalasani, “A comparison of adaptive wormhole routing algorithms,” in *Proc. 20th Ann. Symp. Computer Architecture*, May 1993.

- [5] R. P. Brent, "The parallel evaluation of general arithmetic expressions," *Journal of the ACM*, vol. 21, no. 2, pp. 201–206, April 1974.
- [6] P. I. Crumpton and M. B. Giles, "Multigrid aircraft computations using the OPlus parallel library," in *Parallel Computational Fluid Dynamics: Implementation and Results using Parallel Computers, Proc. Parallel CFD'95*, Pasadena, CA, USA. Elsevier/North-Holland, June 1995, pp. 339–346.
- [7] D. E. Culler, R. M. Karp, D. A. Patterson, A. Sahay, K. E. Schauer, E. Santos, R. Subramonian, and T. von Eicken, "LogP: Towards a realistic model of parallel computation," in *Proc. 4th ACM SIGPLAN Symp. Principles and Practice of Parallel Programming*, San Diego, CA, May 1993.
- [8] H. G. Dietz, T. Muhammad, J. B. Sponaugle, and T. Mattox, "PAPERS: Purdue's adapter for parallel execution and rapid synchronisation," Purdue School of Electrical Engineering, Tech. Rep. TR-EE-94-11, March 1994.
- [9] A. Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manckek, and V. Sunderam, *PVM 3 Users Guide and Reference Manual*. Oak Ridge National Laboratory, Oak Ridge, Tennessee 37831, May 1994.
- [10] G. A. Geist, "PVM3: Beyond network computing," in J. Volkert, ed., *Parallel Computation*, Lecture Notes in Computer Science 734, Springer, 1993, pp. 194–203.
- [11] M. Goudreau, K. Lang, S. Rao, T. Suel, and T. Tsantilas, "Towards efficiency and portability: Programming the BSP model," in *Proc. 8th Ann. Symp. Parallel Algorithms and Architectures*, June 1996, pp. 1–12.
- [12] M. W. Goudreau, K. Lang, S. B. Rao, and T. Tsantilas, "The Green BSP Library," University of Central Florida, Tech. Rep. 95-11, August, 1995.
- [13] M. W. Goudreau, J. M. D. Hill, K. Lang, W. F. McColl, S. D. Rao, D. C. Stefanescu, T. Suel, and T. Tsantilas, "A proposal for a BSP Worldwide standard," BSP Worldwide, <http://www.bsp-worldwide.org/>, April 1996.
- [14] W. Gropp, E. Lusk, and A. Skjellum, *Using MPI: Portable Parallel Programming*. Cambridge, MA: MIT Press, 1994.
- [15] J. He, Q. Miller, and L. Chen, "Algebraic laws for BSP programming," in *Proc. Europar'96*, Lecture Notes in Computer Science, vol. 1124, Springer-Verlag, August 1996.
- [16] J. M. D. Hill and D. B. Skillicorn, "Practical barrier synchronisation," Oxford University Computing Laboratory, Tech. Rep. TR-96-16, October 1996.
- [17] J. M. D. Hill and D. B. Skillicorn, "Lessons learned from implementing BSP," in *High-Performance Computing and Networks*, Lecture Notes in Computer Science, vol. 1225, Springer, April 1997, pp. 762–771. Also appears as Oxford University Computing Laboratory Tech. Rep. TR-96-21.
- [18] J. M. D. Hill, P. I. Crumpton, and D. A. Burgess, "The theory, practice, and a tool for BSP performance prediction," in *Proc. Europar'96, LNCS*, vol. 1124, Springer-Verlag, August 1996, pp. 697–705.
- [19] C. A. R. Hoare and J. He, *Unified Theories of Programming*. Prentice-Hall International, 1997 (to appear).
- [20] R. W. Hockney, "Performance parameters and benchmarking of supercomputers," *Parallel Computing*, vol. 17, pp. 1111–1130, 1991.
- [21] S. Knee, "Program development and performance prediction on BSP machines using Opal," Oxford University Computing Laboratory, Tech. Rep. PRG-TR-18-94, August 1994.
- [22] C. H. Koelbel, D. B. Loveman, R. S. Schreiber, G. L. Steele Jr., and M. E. Zosel, *The High Performance Fortran Handbook*. Cambridge, MA: MIT Press, 1994.
- [23] W. F. McColl, "Scalable computing," in J. van Leeuwen, ed., *Computer Science Today: Recent Trends and Developments*. Lecture Notes in Computer Science, vol. 1000, Springer-Verlag, 1995, pp. 46–61.
- [24] W. F. McColl and Q. Miller, "The GPL language: Reference manual," ESPRIT GEPPCOM Project, Oxford University Computing Laboratory, Tech. Rep., October 1995.
- [25] W. F. McColl, "General purpose parallel computing," in A. M. Gibbons and P. Spirakis, eds, *Lectures on Parallel Computation*. Cambridge International Series on Parallel Computation, Cambridge: Cambridge University Press, 1993, pp. 337–391.
- [26] W. F. McColl, "Special purpose parallel computing," in A. M. Gibbons and P. Spirakis, eds, *Lectures on Parallel Computation*. Cambridge International Series on Parallel Computation, Cambridge: Cambridge University Press, 1993, pp. 261–336.
- [27] Message Passing Interface Forum, "MPI: A message passing interface," in *Proc. Supercomputing'93*, IEEE Computer Society, 1993, pp. 878–883.
- [28] R. Miller, "A library for Bulk Synchronous Parallel programming," in *Proc. BCS Parallel Processing Specialist Group workshop on General Purpose Parallel Computing*, December 1993, pp. 100–108.
- [29] R. Miller, *Two approaches to architecture-independent parallel computation*, Oxford University Computing Laboratory, Wolfson Building, Parks Road, Oxford OX1 3QD, Ph.D. thesis, 1994.
- [30] P. B. Monk, A. K. Parrott, and P. J. Wesson, "A parallel finite element method for electromagnetic scattering," *COMPEL*, vol. 13, Supp. A, pp. 237–242, 1994.
- [31] M. Nibhanupudi, C. Norton, and B. Szymanski, "Plasma simulation on networks of workstations using the bulk synchronous parallel model," in *Proc. Int. Conf. Parallel and Distributed Processing Techniques and Applications*, Athens, GA, November 1995.
- [32] M. J. Quinn and P. J. Hatcher, "On the utility of communication-computation overlap in data-parallel programs," *J. Parallel and Distributed Computing*, vol. 33, pp. 197–204, 1996.
- [33] J. Reed, K. Parrott, and T. Lanfear, "Portability, predictability and performance for parallel computing: BSP in practice," *Concurrency: Practice and Experience* (to appear).
- [34] D. B. Skillicorn, *Foundations of Parallel Programming*. Cambridge Series in Parallel Computation, vol. 6, Cambridge University Press, 1994.
- [35] K. R. Sujithan and J. M. D. Hill, "Collection types for database programming in the BSP model," in *Proc. Eu-*

*romicro Workshop on Parallel and Distributed Processing*, IEEE CS Press, January 1997.

- [36] L. G. Valiant, "A bridging model for parallel computation," *Comm. ACM*, vol. 33, pp. 103–111, August 1990.
- [37] L. G. Valiant, "General purpose parallel architectures," in J. van Leeuwen, ed., *Handbook of Theoretical Computer Science*, vol. A, Elsevier and MIT Press, 1990.

## APPENDIX A THE BSPLib LIBRARY

This Appendix provides slightly more detail about the current major BSP system, BSPLib. We describe C interfaces to the library, but a Fortran version is also available.

### A.1 Initialisation

Processes are created in a BSPLib program by the operations `bsp_begin` and `bsp_end`. There can only be one instance of a `bsp_begin/bsp_end` pair within a program. There are two different ways to start a BSPLib program. If `bsp_begin` and `bsp_end` are the first and last statements in a program, then the entire BSPLib computation is SPMD. In an alternative mode, a single process starts execution and determines the number of parallel processes required for the calculation. It then spawns the required number of processes using `bsp_begin`. Execution of the spawned processes then continue in an SPMD manner, until `bsp_end` is encountered by all the processes. At that point, all processes except process zero are terminated, and process zero is left to continue the execution of the rest of the program sequentially.

One problem with providing this mode is that some parallel machines available today, for example almost all distributed-memory machines, e.g., IBM SP2, Cray T3D, Meiko CS-2, Parsytec GC, Hitachi SR2001, do not provide dynamic process creation. Therefore we *simulate* dynamic spawning using an operation `bsp_init` which takes as its argument a procedure name. The procedure passed as an argument to `bsp_init` must contain `bsp_begin` and `bsp_end` as its first and last statements.

The interface for these library operations is:

```
void bsp_init(void (*startproc)(void),
              int argc, char**argv);
void bsp_begin(int maxprocs);
void bsp_end();
```

`maxprocs` is the number of processes requested by the user.

`startproc` is the name of a procedure that contains `bsp_begin` and `bsp_end` as its first and last statements.

`argc` and `argv` are command line size and arguments.

### A.2 Enquiry

There are also operations to determine the total number of processes, and for each process to find out its process identifier. The interface for these operations is:

```
int bsp_nprocs();
int bsp_pid();
```

If the function `bsp_nprocs` is called before `bsp_begin`, then it returns the number of processors which are available. If it is called after `bsp_begin` it returns `n`, the actual number of processes allocated to the program, where  $1 \leq n \leq \text{maxprocs}$ , and `maxprocs` is the number of processes requested in `bsp_begin`. Each of the `n` processes created by `bsp_begin` has a unique associated value `m` in the range  $0 \leq m \leq n - 1$ . The function `bsp_pid` returns the associated value of the process executing the function call.

### A.3 Synchronisation

A BSPLib calculation consists of a sequence of supersteps. The end of one superstep and the start of the next is identified by a call to the library procedure `bsp_sync` with interface:

```
void bsp_sync();
```

### A.4 DRMA

There are two ways of communicating among processes: one using direct remote-memory access (DRMA), and the other using a BSP version of message passing.

The DRMA communication operations are defined for stack- and heap-allocated data structures as well as for static data. This is achieved by allowing a process to reference only certain *registered* areas of a remote memory. In a registration procedure, processes use the operation `bsp_pushregister` to announce the address of the start of a local area which is available for global remote use. This makes it possible to execute BSP programs using heterogeneous processor architectures. Registration takes effect at the next barrier synchronisation.

```
void bsp_pushregister (void*region,
                      int nbytes);
void bsp_popregister (void*region);
```

`region` is the starting address of the region to be registered or unregistered. The name `region` must be the same for all logically-related calls to `bsp_pushregister` or `bsp_popregister`, and implementations may check that this is true.

`nbytes` is the size of the region (used for range checking).

Each processor maintains a stack of registration slots. Logically-related calls to `bsp_pushregister` in different processes (the  $i$ th call in each process is related to the  $i$ th call in all of the others) associate a variable name and the addresses to which it is mapped in each process with the next available slot. Registration slots can be deallocated using `bsp_popregister`, which invalidates the last slot associated with the variable name passed as an argument – deregistration does not impose the strict nesting of push–pop pairs that is normally associated with a stack; the scheme allows the popping of registrations to occur in an arbitrary order. This provides the benefits of encapsulation provided by a stack, whilst providing the flexibility associated with a heap-based discipline. However, the registration slot of the argument to `popregister` *must be the same across all the processing elements*.

The intent of registration is to make it simple to refer to variables in other processes without requiring their locations to be explicitly known. A reference to a registered name in a put or get is translated to the address corresponding to the remote variable with the same name. Here is an example:

Process 0:

```
int x;
bsp_pushregister(&x, sizeof(int));
bsp_sync();
x = 3;
bsp_put(1, &x, &x, 0, sizeof(int));
bsp_sync();
```

Process 1

```
int x;
bsp_pushregister(&x, sizeof(int));
bsp_sync();
bsp_sync();
```

Process 0 and Process 1 register `x` in the first slot. When Process 0 executes a put, using `x` as the destination region name, this is mapped to the region whose address is associated with the first slot in Process 1. Therefore, the variable `x` in Process 1 has the value 3 placed in it after the second `sync` as the result of the put.

The operation `bsp_put` pushes locally-held data into a registered remote-memory area on a target process, without the active participation of the target process. The operation `bsp_get` reaches into the registered local memory of another process to copy data values held there into a data structure in its own local memory. All gets are executed before all puts at the end of a superstep, consistent with the semantics that communications do not take effect locally until the end of a superstep. Their interfaces are:

```
void bsp_[hp]put(int pid,
                const void *src,
                void *dst,
                int offset,
                int nbytes);
```

`pid` is the identifier of the process where data is to be stored.

`src` is the location of the first byte to be transferred by the put operation. The calculation of `src` is performed on the process that initiates the put.

`dst` is the base address of the area where data is to be stored. It must be a previously-registered data area.

`offset` is the displacement in bytes from `dst` to which `src` will copy. The calculation of `offset` is performed by the process that initiates the put.

`nbytes` is the number of bytes to be transferred from `src` into `dst`. It is assumed that `src` and `dst` are addresses of data structures that are at least `nbytes` in size.

```
void bsp_[hp]get(int pid,
                const void *src,
                int offset,
                void *dst,
                int nbytes);
```

`pid` is the identifier of the process from which data is to be obtained.

`src` is the base address of the area from which data will be obtained. `src` must be a previously-registered data structure.

`offset` is an offset from `src`. The calculation of `offset` is performed by the process that initiates the get.

`dst` is the location of the first byte where the data obtained is to be placed. The calculation of `dst` is performed by the process that initiates the get.

`nbytes` is the number of bytes to be transferred from `src` into `dst`. It is assumed that `src` and `dst` are addresses of data structures that are at least `nbytes` in size.

The semantics adopted for BSPLib `bsp_put` communication is *buffered-locally/buffered-remotely*. When a `bsp_put` is executed, the data to be transferred is copied out of user address space immediately. The executing process is free to alter the contents of those locations after return from the call to `bsp_put`. While the semantics is clean and safety is maximized, puts may unduly tax the memory resources of an implementation, thus preventing large data transfers. Consequently, BSPLib also provides a *high-performance put*

operation `bsp_hput` whose semantics is *unbuffered-locally/unbuffered-remotely*. The use of this operation requires care, as correct data delivery is only guaranteed if neither communication nor local/remote computations modify either the source or the destination areas during a superstep. The main advantage of this operation is its economical use of memory. It is therefore particularly useful for applications which repeatedly transfer large data sets.

The `bsp_get` and `bsp_hpget` operations reach into the local memory of another process and copy previously-registered remote data held there into a data structure in the local memory of the process that initiated them.

### A.5 BSMP

Bulk synchronous remote-memory access is a convenient style of programming for BSP computations that can be statically analysed in a straightforward way. It is less convenient for computations in which the volumes of data being communicated are irregular and data-dependent, or where the computation to be performed in a superstep depends on the quantity and form of data received at its start. A more appropriate style of programming in such cases is bulk-synchronous message passing (BSMP).

In BSMP, a non-blocking send operation delivers messages to a system buffer associated with the destination process. The message is guaranteed to be in the destination buffer at the beginning of the subsequent superstep, and can be accessed by the destination process only during that superstep. A collection of messages sent to the same process has no implied ordering at the receiving end. However, since messages may be tagged, the programmer can identify them by their tag.

In BSPLib, bulk-synchronous message passing is based on the idea of two-part messages, a fixed-length part carrying tagging information that will help the receiver to interpret the message, and a variable-length part containing the main data payload. We will call the fixed-length portion the *tag* and the variable-length portion the *payload*. In C programs, either part could be a complicated structure. The length of the tag is required to be fixed during any particular superstep, but may vary between supersteps. The buffering mode of the BSMP operations is *buffered-locally/buffered-remotely*.

The procedure to set tag size must be called collectively by all processes. Moreover, in any superstep where `bsp_set_tag_size` is called, it must be called before sending any messages.

```
void bsp_set_tag_size(int *tag_bytes);
```

`tag_bytes`, on entry to the procedure, specifies the size of the fixed-length portion of every message from the current superstep until it is updated; the

default tag size is zero. On return from the procedure, `tag_bytes` is changed to reflect the *previous* value of the tag size to allow for its use inside procedures.

The tag size of incoming messages is prescribed by the outgoing tag size of the previous step.

The procedure `bsp_bsmpt_info` is an enquiry operation that returns information concerning how many BSMP packets were sent to the process calling the operation in the prior superstep. This information is intended to help the user to allocate an appropriate sized data structure to hold any incoming BSMP messages.

```
void bsp_bsmpt_info(int*packets,
                   int*accum_nbytes);
```

`packets` becomes the number of packets sent using `bsp_send` in the previous superstep.

`accum_nbytes` is the accumulated size of all the packets.

The `bsp_send` operation is used to send a message that consists of a tag and a payload to a specified destination process. The destination process will be able to access the message during the subsequent superstep. Its interface is:

```
void bsp_send(int pid,
              const void*tag,
              const void*payload,
              int payload_bytes);
```

`pid` is the identifier of the process where data is to be sent.

`tag` is a token that can be used to identify the message. Its size is determined by the value specified in `bsp_set_size_tag`.

`payload` is the location of the first byte of the payload to be communicated.

`payload_bytes` is the size of the payload.

The `bsp_send` operation copies both the tag and the payload of the message out of user space into the system before returning. The `tag` and `payload` inputs may be changed by the user immediately after the `bsp_send` returns.

To receive a message, the operations `bsp_get_tag` and `bsp_move` are used. The operation `bsp_get_tag` returns the tag of the first message in the buffer. The operation `bsp_move` copies the payload of the first message in the buffer into `payload`, and removes that message from the buffer. Its interface is:

```
void bsp_get_tag(int *status,
                 void *tag);
```

`status` returns `-1` if the system buffer is empty. Otherwise it returns the length of the payload of the first message in the buffer. This length can be used to allocate an appropriately-sized data structure for copying the payload using `bsp_move`.

`tag` is unchanged if the system buffer is empty. Otherwise it is assigned the tag of the first message in the buffer.

```
void bsp_move(void *payload,
              int reception_nbytes);
```

`payload` is an address to which the message payload will be copied. The buffer is then advanced to the next message.

`reception_nbytes` specifies the size of the reception area where the payload will be copied into. At most `reception_nbytes` will be copied into payload.

```
int bsp_hpmove(void**tag_ptr_buf,
               void**payload_ptr_buf);
```

`bsp_hpmove` is a function which returns `-1`, if the system buffer is empty. Otherwise it returns the length of the payload of the first message in the buffer and

- places a pointer to the tag in `tag_ptr_buf`;
- places a pointer to the payload in `payload_ptr_buf`; and
- conceptually removes the message (by advancing a pointer representing the head of the buffer).

Note that `bsp_move` flushes the corresponding message from the buffer, while `bsp_get_tag` does not. This allows a program to get the tag of a message (as well as the payload size in bytes) before obtaining the payload of the message. It does, however, require that even if a program only uses the fixed-length tag of incoming messages the program must call `bsp_move` to get successive message tags.

`bsp_get_tag` can be called repeatedly and will always return the same tag until a call to `bsp_move`.

## A.6 Halt

The function `bsp_abort` can be used to print an error message followed by a halt of the entire BSPLib program. The routine is designed *not* to require a barrier synchronisation of all processes. A single process can therefore halt the entire BSPLib program.

```
void bsp_abort(char*format, ...);
```

`format` is a C-style format string as used by `printf`.

Any other arguments are interpreted in the same way as the variable number of arguments to `printf`.

The function `bsp_time` provides access to a high-precision timer – the accuracy of the timer is implementation-specific. The function is a local operation of each process, and can be issued at any point after `bsp_begin`. The result of the timer is the time in seconds since `bsp_begin`. The semantics of `bsp_time` is as though there were `bsp_nprocs` timers, one per process. BSPLib does *not* impose any synchronisation requirements between the timers in each process.

```
double bsp_time();
```

## APPENDIX B BENCHMARKING

The BSP parameter  $l$  measures the minimum time for all processors to barrier synchronise. It is benchmarked by repeatedly over-sampling barrier synchronisation, and measuring the wall-clock time. Repeated barrier synchronisation produces a pessimistic value for  $l$  as it models the case where the computation part of each superstep completes in each processor at the same moment. This produces most contention in whatever resources are used for synchronising.

Two values for the BSP parameter  $g$  are calculated. The first is the value of  $g$  experienced when routing a local communication (a cyclic shift), and the second a global communication using a total exchange. As well as calculating the value of  $g$ , the benchmark also calculates the value for  $n_{1/2}$  used in Equation (1). This is done by routing a fixed-sized  $h$ -relation (an over-sampling of 10 iterations is performed for each  $h$ -relation) for large  $h$  and measuring the elapsed time of a superstep containing no computation.

Sophisticated profiling tools are available to examine how much this measured value of  $g$  is affected by particular properties of the target computer. We have already mentioned some such factors, for example the overhead of message startup and the extra data that must be transferred as control information. This can be clearly seen in Figures 9–11.

These figures shows the amount of data transferred and the effective value of  $g$  in two phases. The first half of each figure shows a cyclic shift; the second half a total exchange. All supersteps in each half send an  $h$ -relation, in sets of size 10 for oversampling, varying the granularity for each set – first using single messages of size  $h$ , then using two messages of size  $h/2$ , and so on. According to the theory, the measured value of  $g$  should be the same for all of these granularities, since the same total volume of data is moved into and out of each processor. The top half of each figure shows the volume of data being moved. The second half shows the measured value of  $g$  for each superstep.



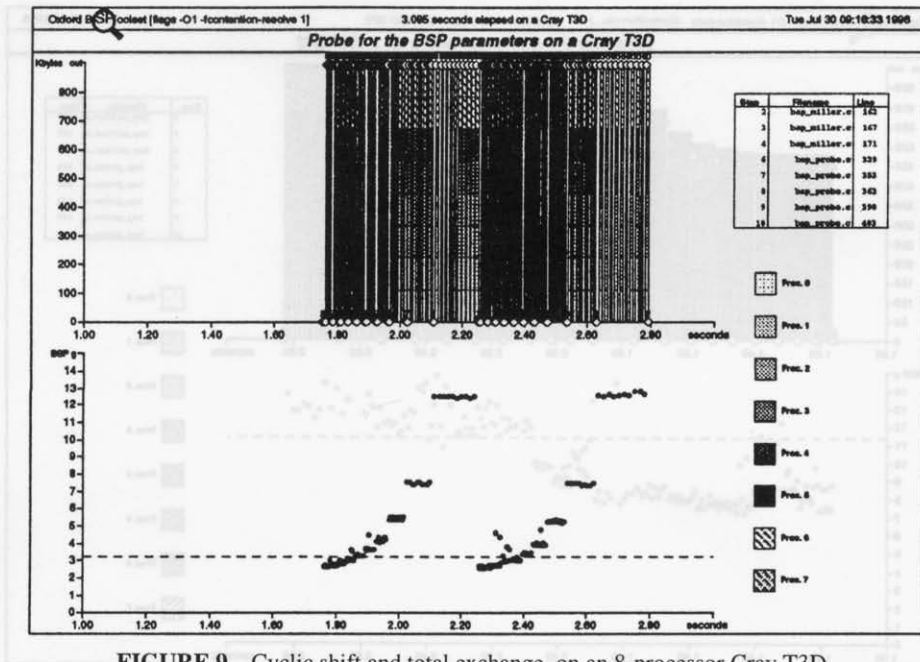


FIGURE 9 Cyclic shift and total exchange, on an 8-processor Cray T3D.

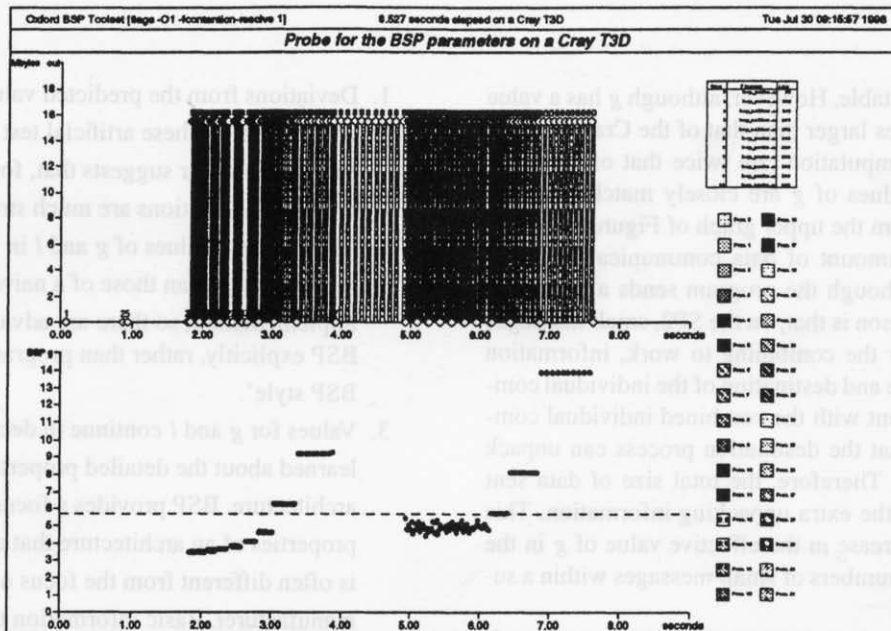


FIGURE 10 Cyclic shift and total exchange, on a 32-processor Cray T3D.

The curves in the lower half of Figure 9 show the value of  $g$  during the cyclic shift phase and the total exchange phase. The curves are not the expected horizontal lines because of the overhead of message startup. The implementation on the CRAY delays messages and re-orders them, but does not combine them because the communication mechanism is DRMA. The curves are good matches for Equation (1) which uses the  $n_{1/2}$  param-

eter to model the extra cost of communicating small messages. The dotted line in the graph shows the value of  $g$  obtained from the benchmark, and given in Table 2. It is very close to the asymptote of the curves. The same structure can be seen for larger numbers of processors (Figure 10).

Figure 11 shows the same benchmark running on an eight-processor IBM SP2. Unlike the Cray, the value of

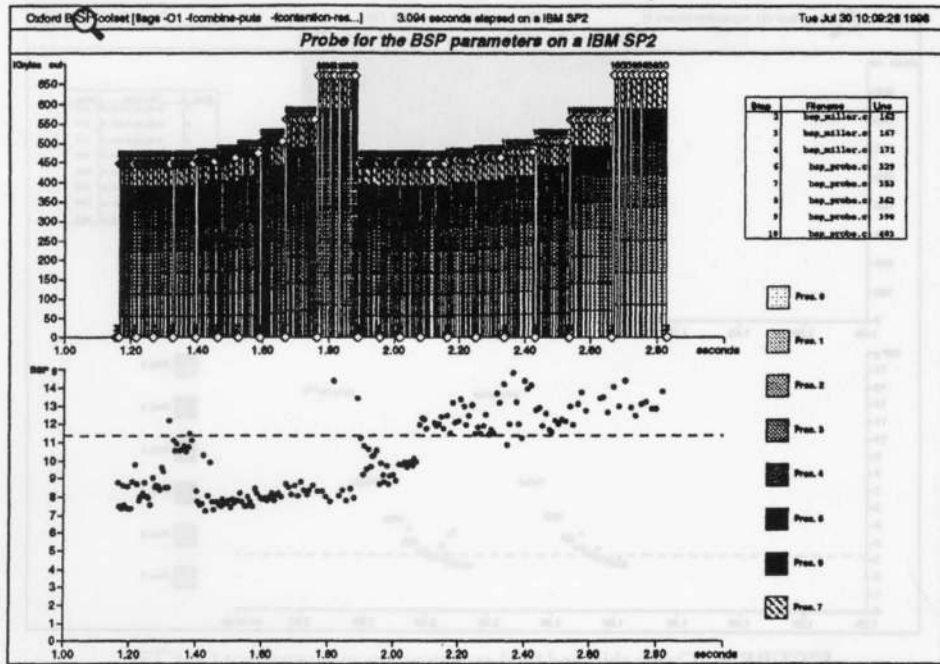
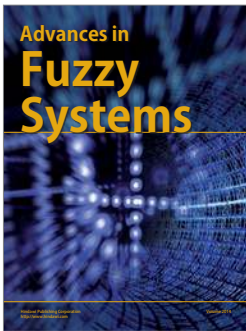
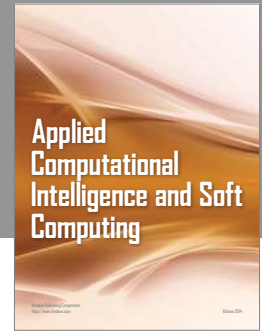
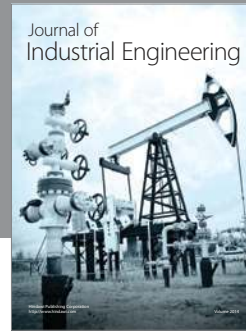
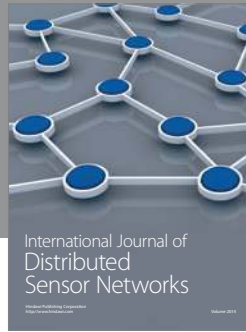
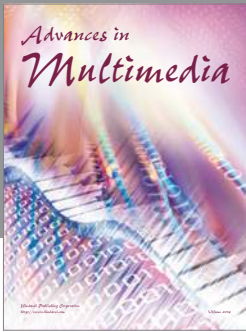


FIGURE 11 Cyclic shift and total exchange on an 8-processor IBM SP2.

$g$  is more unpredictable. However, although  $g$  has a value which is three times larger than that of the Cray, the SP2 has a per-node computation rate twice that of the T3D, so the absolute values of  $g$  are closely matched on the two machines. From the upper graph of Figure 11, it can be seen that the amount of data communicated gradually grows, even though the program sends a fixed size  $h$ -relation. The reason is that, on the SP2, small messages are combined. For the combining to work, information concerning the size and destination of the individual communications are sent with the combined individual communications, so that the destination process can unpack the data correctly. Therefore, the total size of data sent may triple due to the extra unpacking information. This causes a slight increase in the effective value of  $g$  in the presence of large numbers of small messages within a superstep.

These examples show that  $g$  can deviate from the values predicted by the cost model because of properties of the target computers, and unavoidable overheads in the implementation of the library operations. Three observations seem relevant:

1. Deviations from the predicted values are relatively small, even for these artificial test programs, and experience so far suggests that, for practical programs, deviations are much smaller.
2. The achieved values of  $g$  and  $l$  in BSPLib are much smaller than those of a naive implementation, so there are advantages to using BSP explicitly, rather than programming 'in the BSP style'.
3. Values for  $g$  and  $l$  continue to decrease as more is learned about the detailed properties of each architecture. BSP provides a focus for the properties of an architecture that are critical, but it is often different from the focus of the manufacturer. Basic information that would make better implementations possible is hard to obtain, sometimes because even the manufacturers do not know it.



**Hindawi**

Submit your manuscripts at  
<http://www.hindawi.com>

