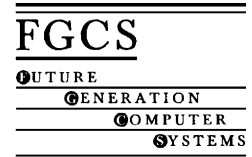




ELSEVIER

Future Generation Computer Systems 18 (2002) 1127–1142



www.elsevier.com/locate/future

HARNESS fault tolerant MPI design, usage and performance issues

Graham E. Fagg^{a,*}, Jack J. Dongarra^b

^a High Performance Computing Center Stuttgart (HLRS), Parallel and Distributed Systems, Allmandring 30, D-70550 Stuttgart, Germany

^b Department of Computer Science, Suite 413, 1122 Volunteer Blvd., University of Tennessee, Knoxville, TN 37996-3450, USA

Abstract

Initial versions of MPI were designed to work efficiently on multi-processors which had very little job control and thus static process models. Subsequently forcing them to support a dynamic process model suitable for use on clusters or distributed systems would have reduced their performance. As current HPC collaborative applications increase in size and distribution the potential levels of node and network failures increase. This is especially true when MPI implementations are used as the communication media for GRID applications where the GRID architectures themselves are inherently unreliable thus requiring new fault tolerant MPI systems to be developed. Here we present a new implementation of MPI called FT-MPI that allows the semantics and associated modes of failures to be explicitly controlled by an application via a modified MPI API. Given is an overview of the FT-MPI semantics, design, example applications and some performance issues such as efficient group communications and complex data handling. Also briefly described is the HARNESS g_hcore system that handles low-level system operations on behalf of the MPI implementation. This includes details of plug-in services developed and their interaction with the FT-MPI runtime library.

© 2002 Elsevier Science B.V. All rights reserved.

Keywords: MPI implementation; HARNESS; FT-MPI; Meta computing; Fault tolerant message passing

1. Introduction

Although MPI [11] is currently the de-facto standard system used to build high performance applications for both clusters and dedicated MPP systems, it is not without its problems. Initially, MPI was designed to allow for very high efficiency and thus performance on a number of early 1990s MPPs, that at the time had limited OS runtime support. This led to the current MPI design of a static process model. While this

model was possible to implement for MPP vendors, easy to program for, and more importantly something that could be agreed upon by a standards committee. The second version of MPI standard known as MPI-2 [22] did include some support for dynamic process control, although this was limited to the creation of new MPI process groups with separate communicators. These new processes could not be merged with previously existing communicators to form intracommunicators needed for a seamless single application model and were limited to a special set of extended collectives (group) communications.

The MPI static process model suffices for small numbers of distributed nodes within the currently emerging masses of clusters and several hundred nodes of dedicated MPPs. Beyond these sizes the

* Corresponding author. Tel.: +49-711-6855990;

fax: +49-711-6787626.

E-mail addresses: fagg@hlrs.de (G.E. Fagg), dongarra@cs.utk.edu (J.J. Dongarra).

URL: <http://www.hlrs.de/people/fagg>

mean time between failure (MTBF) of CPU nodes starts becoming a factor. As attempts to build the next generation Peta-flop systems advance, this situation will only become more adverse as individual node reliability becomes out weighted by orders of magnitude increase in node numbers and hence node failures. Current GRID [30] technologies such as GLOBUS [27] also provide for middleware services such as naming, resource discovery, and book keeping that is robust and handle expected failures gracefully. Unfortunately, the defacto MPI message passing library for Globus, MPICH-G [28] is not expected to handle loss of MPI processes or partitioning of networks gracefully and failures still lead to pathological failure of applications unless special precautions are taken such application check-pointed discussed further in the next section.

The aim of FT-MPI is to build a fault tolerant MPI implementation that can survive failures, while offering the application developer a range of recovery options other than just returning to some previous check-pointed state. FT-MPI is built on the HARNESS [1] meta-computing system, and is meant to be used as the HARNESS default application level message passing interface. Its design allows it to easily ported to other GRID environments by porting of its modular services that are implemented in the form of short-lived daemons.

2. Check-point and roll back versus replication techniques

The first method attempted to make MPI applications fault tolerant was through the use of check-pointing and roll back. Co-Check MPI [2] from the Technical University of Munich being the first MPI implementation built that used the Condor library for check-pointing an entire MPI application. In this implementation, all processes would flush their messages queues to avoid in flight messages getting lost, and then they would all synchronously check-point. At some later stage if either an error occurred or a task was forced to migrate to assist load balancing, the entire MPI application would be rolled back to the last complete check-point and be restarted. This systems main drawback being the need for the entire application having to check-point synchronously,

which depending on the application and its size could become expensive in terms of time (with potential scaling problems). A secondary consideration was that they had to implement a new version of MPI known as tuMPI as updating MPICH was considered too difficult.

Another system that also uses check-pointing but at a much lower level is StarFish MPI [3]. Unlike Co-Check MPI which relies on Condor, Starfish MPI uses its own distributed system to provide built in check-pointing. The main difference with Co-Check MPI is how it handles communication and state changes which are managed by StarFish using strict atomic group communication protocols built upon the Ensemble system [4], and thus avoids the message flush protocol of Co-Check. Being a more recent project StarFish supports faster networking interfaces than tuMPI.

The project closest to FT-MPI known to the author is the Implicit Fault Tolerance MPI project MPI-FT [15] by Paraskevas Evripidou of Cyprus University. This project supports several master–slave models where all communicators are built from grids that contain ‘spare’ processes. These spare processes are utilized when there is a failure. To avoid loss of message data between the master and slaves, all messages are copied to an observer process, which can reproduce lost messages in the event of any failures. This system appears only to support SPMD style computation and has a high overhead for every message and considerable memory needs for the observer process for long running applications. This system is not a full check-point system in that it assumes any data (or state) can be rebuilt using just the knowledge of any passed messages, which might not be the case for non-deterministic unstable solvers.

FT-MPI has much lower overheads compared to the above check-pointing systems, and thus much higher potential performance. These benefits do, however, have consequences. An application using FT-MPI has to be designed to take advantage of its fault tolerant features as shown in the next section, although this extra work can be trivial depending on the structure of the application. If an application needs a high level of fault tolerance where node loss would equal data loss then the application has to be designed to perform some level of user directed check-pointing. FT-MPI does allow for atomic communications much like Starfish, but

unlike Starfish, the level of correctness can be varied on for individual communicators. This provides users the ability to fine tune for coherency or performance as system and application conditions dictate.

Currently, GRID application efforts such as GrADS [25] primarily focus on gaining high performance from GRIDs rather than handling failures, although current efforts at the University of Tennessee [26] involve check-pointing distributed applications to improve fault tolerance. Unlike the above check-pointing systems that rely on local disks for check-pointed data storage, the current GRADS effort is experimenting with replicated distributed storage built on top of the IBP [29] system to improve both availability and performance. This avoids the potential for an unreliable GRID to also make the check-point unavailable and thus thwart any application restarts after a likely failure has occurred.

3. FT-MPI semantics

Current semantics of MPI indicate that a failure of an MPI process or communication causes all communicators associated with them to become *invalid*. As the standard provides no method to reinstate them (and it is unclear if we can even *free* them), we are left with the problem that this causes MPI_COMM_WORLD itself to become invalid and thus the entire MPI application will grid to a halt.

FT-MPI extends the MPI communicator states from {valid, invalid} to a range {FT_OK, FT_DETECTED, FT_RECOVER, FT_RECOVERED, FT_FAILED}. In essence this becomes {OK, PROBLEM, FAILED}, with the other states mainly of interest to the internal fault recovery algorithm of FT-MPI. Processes also have typical states of {OK, FAILED} which FT-MPI replaces with {OK, Unavailable, Joining, Failed}. The *Unavailable* state includes unknown, unreachable or “we have not voted to remove it yet” states.

A communicator changes its state when either an MPI process changes its state, or a communication within that communicator fails for some reason. Some details of failure detection is given in Section 4.4.

The typical MPI semantics is from OK to Failed which then causes an application abort. By allowing the communicator to be in an intermediate state we allow the application the ability to decide how to alter

the communicator and its state as well as how communication within the intermediate state behaves.

3.1. Failure modes

On detecting a failure within a communicator, that communicator is marked as having a probable error. Immediately as this occurs the underlying system sends a state update to all other processes involved in that communicator. If the error was a communication error, not all communicators are forced to be updated, if it was a process exit then all communicators that include this process are changed. Note, this might not be all current communicators as we support MPI-2 dynamic tasks and thus multiple MPI_COMM_WORLDs.

How the system behaves depends on the communicator failure mode chosen by the application. The mode has two parts, one for the communication behavior and one for the how the communicator reforms if at all.

3.2. Communicator and communication handling

Once a communicator has an error state it can only recover by rebuilding it, using a modified version of one of the MPI communicator build functions such as MPI_Comm_{create, split or dup}. Under these functions the new communicator will follow the following semantics depending on its failure mode:

- *SHRINK*. The communicator is reduced so that the data structure is contiguous. The ranks of the processes are *changed*, forcing the application to recall MPI_COMM_RANK.
- *BLANK*. This is the same as SHRINK, except that the communicator can now contain gaps to be filled in later. Communicating with a gap will cause an invalid rank error. Note also that calling MPI_COMM_SIZE will return the extent of the communicator, not the number of valid processes within it.
- *REBUILD*. Most complex mode that forces the creation of new processes to fill any gaps until the size is the same as the extent. The new processes can either be placed in to the empty ranks, or the communicator can be shrank and the remaining processes filled at the end. This is used for applications that

```

rc= MPI_Send (----, com);
If (rc==MPI_ERR_OTHER) {
    MPI_Comm_dup (com, newcom);      /* collective recovery occurs here! */
    MPI_Comm_free (com);
    com = newcom;
}
/* continue.. */

```

Scheme 1. Simple FT-MPI send usage.

require a certain size to execute as in power of two FFT solvers.

- *ABORT*. Is a mode which affects the application immediately an error is detected and forces a graceful abort. The user is unable to trap this. If the application needs to avoid this they must set all communicators to one of the above communicator modes.

Communications within the communicator are controlled by a message mode for the communicator which can be either of:

- *NOP*. No operations on error, i.e. no user level message operations are allowed and all simply return an error code. This is used to allow an application to return from any point in the code to a state where it can take appropriate action as soon as possible.
- *CONT*. All communication that is NOT to the affected/failed node can continue as normal. Attempts to communicate with a failed node will return errors until the communicator state is reset.

The user discovers any errors from the return code of any MPI call, with a new fault indicated by `MPI_ERR_OTHER`. Details as to the nature and specifics of an error is available through the cached attributes interface in MPI.

3.3. Point to point versus collective correctness

Although collective operations pertain to point to point operations in most cases, extra care has been taken in implementing the collective operations so that if an error occurs during an operation, the result of the operation will still be the same as if there had been no error, or else the operation is aborted.

Broadcast, gather and all gather demonstrate this perfectly. In broadcast even if there is a failure of a receiving node, the receiving nodes still receive the same data, i.e. the same end result for the surviving

nodes. Gather and all-gather are different in that the result depends on if the problematic nodes sent data to the gatherer/root or not. In the case of gather, the root might or might not have gaps in the result. For the all2all operation, which typically uses a ring algorithm it is possible that some nodes may have complete information and others incomplete. Thus for operations that require multiple node input as in gather/reduce type operations any failure causes all nodes to return an error code, rather than possibly invalid data. Currently, an addition flag controls how strict the above rule is enforced by utilizing an extra barrier call at the end of the collective call if required.

3.4. FT-MPI usage

Typical usage of FT-MPI would be in the form of an error check and then some corrective action such as a communicator rebuild. A typical code fragment is shown below in [Scheme 1](#), where on an error the communicator is simply rebuilt and reused.

Some types of computation such as SPMD master-worker codes only need the error checking in the master code if the user is willing to accept the master as the only point of failure. [Scheme 2](#) below shows how complex a master code can become. In this example the communicator mode is `BLANK` and communications mode is `CONT`. The master keeps track of work allocated, and on an error just reallocates the work to any ‘free’ surviving processes. Note, the code has to check to see if there are any surviving worker processes left after each death is detected.

4. FT-MPI Implementation details

FT-MPI is a partial MPI-2 implementation. It currently contains support for both C and Fortran

```

rc = MPI_Bcast ( initial_work...);
if(rc==MPI_ERR_OTHER)reclaim_lost_work(...);
while ( ! all_work_done) {
  if (work allocated) {
    rc = MPI_Recv ( buf, ans_size, result_dt,
                  MPI_ANY_SOURCE, MPI_ANY_TAG, comm, &status);
    if (rc==MPI_SUCCESS) {
      handle_work (buf);
      free_worker (status.MPI_SOURCE);
      all_work_done--;
    }
  }
  else {
    reclaim_lost_work(status.MPI_SOURCE);
    if (no_surviving_workers) { /* ! do something ! */ }
  }
} /* work allocated */
/* Get a new worker as we must have received a result or a death */
rank=get_free_worker_and_allocate_work();
if (rank) {
  rc = MPI_Send (... rank... );
  if (rc==MPI_OTHER_ERR) reclaim_lost_work (rank);
  if (no_surviving_workers) { /* ! do something ! */ }
} /* if free worker */
} /* while work to do */

```

Scheme 2. FT-MPI master-worker code.

interfaces, all the MPI-1 function calls required to run both the PSTSWM [6] and BLAS [21] applications. BLAS is supported so that SCALAPACK [20] applications can be tested. Currently, only some of the dynamic process control functions from MPI-2 are supported.

The current implementation is built as a number of layers as shown in Fig. 1. Operating system support is provided by either PVM or the C HARNESS *G_HCORE*. Although point to point communication is provided by a modified SNIPE.Lite communication library taken from the SNIPE project [4].

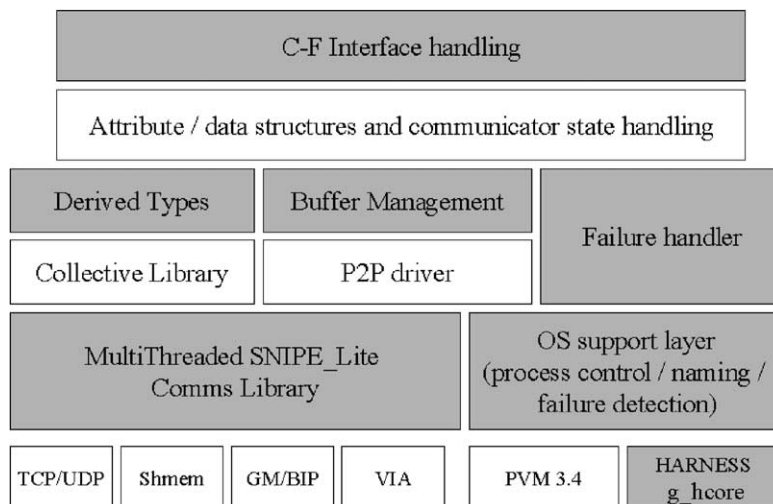


Fig. 1. Overall structure of the FT-MPI implementation.

A number of components have been extensively optimized, these include:

- Derived data types and message buffers.
- Collective communications.
- Point to point communication using multi-threading.

4.1. Derived Data Type handling

MPI-1 introduced extensive facilities for user Derived Data Type (DDT) [11] handling that allows for in effect strongly typed message passing. The handling of these possibly non-contiguous data types is very important in real applications, and is often a neglected area of communication library design [17]. Most communications libraries are designed for low latency and/or high bandwidth with contiguous blocks of data [14]. Although this means that they must avoid unnecessary memory copies, the efficient handling of recursive data structures is often left to simple iterations of a loop that packs a send/receive buffer.

4.1.1. FT-MPI DDT handling

Having gained experience with handling DDTs within a heterogeneous system from the PVMPI/MPI_Connect library [18] the authors of FT-MPI redesigned the handling of DDTs so that they would not just handle the recursive data types flexibly but also take advantage of internal buffer management structure to gain better performance. In a typical system the DDT would be collected/gathered into a single buffer and then passed to the communications library, which may have to encode the data using XDR, e.g., and then segment the message into packets for transmission. These steps involving multiple memory copies across program modules (reducing cache effectiveness) and possibly precluding overlapping (concurrency) of operations.

The DDT system used by FT-MPI was designed to reduce memory copies while allowing for overlapping in the three stages of data handling:

- *Gather/scatter*. Data is collected into or from recursively structured non-contiguous memory.
- *Encoding/decoding*. Data passed between heterogeneous machine architectures than use different floating point representations need to be converted so that the data maintains the original meaning.

- *Send/receive packetizing*. All of the send or receive cannot be completed in a single attempt and the data has to be sent in blocks. This is usually due to buffering constraints in the communications library/OS or even hardware flow control.

4.1.2. DDT methods and algorithms

Under FT-MPI data can be gathered/scattered by compressing the data type representation into a compacted format that can be efficiently transversed (not to be confused with compressing data discussed below). The algorithm used to compact data type representation would break down any recursive data type into an optimized maximum length new representation. FT-MPI checks for this optimization when the users application commits the data type using the MPI_Type_commit API call. This allows FT-MPI to optimize the data type representation before any communication is attempted that uses them.

When the DDT is being processed the actual user data itself can also be compacted into/from a contiguous buffer. Several options for this type of buffering are allowed that include:

- *Zero padding*. Compacting into the smallest buffer space.
- *Minimal padding*. Compacting into smallest space but maintaining correct word alignment.
- *Re-ordering pack*. Re-arranging the data so that all the integers are packed first, followed by floats, etc., i.e. type by type.

The minimal and no padded methods are used when moving the data type within a homogeneous set of machines that require no numeric representation encoding or decoding. The zero padding method benefits slower networks, and alignment padded can in some cases assist memory copy operations, although its real benefit is when used with re-ordering.

The re-ordered compacting method shown in Fig. 2 for a data type that consists of characters (*C*) and integers (*I*). This type of compacting is designed to be used when some additional form encoding/decoding takes place. In particular, moving the re-ordered data type by type through fixed XDR/Swap buffers improves its performance considerably. Two types of DDT encoding are supported, the first is the slower generic SUN XDR format and the second is simple byte swapping to convert between little and big endian numbers.

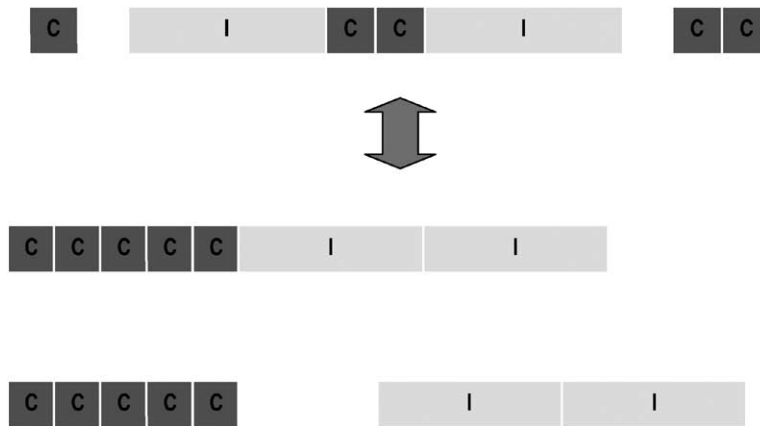


Fig. 2. Compacting storage of re-ordered DDT. Without padding and with correct alignment.

4.1.3. FT-MPI DDT performance

Tests comparing the DDT code to MPICH (1.3.1) on a 93 element DDT taken from a fluid dynamic code were performed between Sun SPARC Solaris and Red Hat (6.1) Linux machines as shown in Table 1. The tests were on small and medium arrays of this data type. All the tests were performed using MPICH MPI_Send and MPI_Recv operations, so that the point to point communications speeds were not a factor, and only the handling of the data types was compared.

The tests show that the compacted data type handling gives from 10 to 19% improvement for small messages and 78–81% for larger arrays on same numeric representation machines. The benefits of buffer reuse and re-ordered of data elements leads to considerable improvements on heterogeneous networks as shown in the comparison between the default MPICH

operations and the DDT ByteSwap and DDT XDR results. It is important to note that all these tests used MPICH to perform the point to point communication, and thus the overlapping of the data gather/scatter, encoding/decoding and non-blocking communication is not shown here, and is expected to yield even higher performance.

The above tests were performed using the DDT software as a standalone library that can be used to improve any MPI implementation. This software is currently being made into an MPI profiling library so that its use will be completely transparent. Two other efforts closely parallel this section of work on DDTs. PACX [19] from HLRS, RUS Stuttgart, requires the heterogeneous data conversion facilities and a project from NEC Europe [16] concentrates on efficient data type representation and transmission in homogeneous systems.

Table 1
Performance of the FT-MPI DDT software compared to MPICH

Type of operation (arch 2 arch) (method) (encoding)	11956 bytes B/W MB/s	Percentage compared to MPICH	95648 bytes B/W MB/s	Percentage compared to MPICH
Sparc 2 Sparc MPICH	5.49		5.47	
Sparc 2 Sparc FTMPI	6.54	+19	9.74	+78
Linux 2 Linux MPICH	7.11		8.79	
Linux 2 Linux FTMPI	7.87	+10	9.92	+81
Sparc 2 Linux MPICH	0.855		0.729	
Sparc 2 Linux FTMPI Byte Swap	5.87	+586	8.20	+1024
Sparc 2 Linux FTMPI XDR	5.31	+621	6.15	+743

4.2. Collective communications

The performance of the MPIs collective communications is critical to most MPI-based applications [6]. A general algorithm for a given collective communication operation may not give good performance on all systems due to the differences in architectures, network parameters and the storage capacity of the underlying MPI implementation [7]. In an attempt to improve over the usual collective library built on point to point communications design as in the $\log P$ model [9], we built a collective communications library that is tuned to its target architecture though the use of a limited set of micro-benchmarks. Once the static system is optimized we then tune the topology dynamically by re-orders the logical addresses to compensate for changing run time variations. Other projects that use a similar approach to optimizing include [12,13].

4.2.1. Collective communication algorithms and benchmarks

The micro-benchmarks are conducted for each of the different classes of MPI collective operations broadcast, gather, scatter, reduce, etc. individually, i.e.

the algorithm that produces the best broadcast might not produce the best scatter even though they appear similar.

The algorithms tested are different variations of standard topologies and methods such as sequential, Rabenseifner [10], binary and binomial trees, using different combinations of blocking/non-blocking send and receives. Each test is varied over a number of processors, message sizes and segmentation sizes. The segmenting of messages was found to improve bi-section bandwidth obtained depending on the target network.

These tests produce an optimal topology and segment size for each MPI collective of interest. Tests against vendor MPI implementations have shown that our collective algorithms are comparable or even faster as shown in Figs. 3 and 4. These tests were comparing the FT-MPI tuned broadcast versus the IBM broadcast on different configurations of an IBM SP2 for eight nodes.

4.2.2. Dynamic re-ordering of topologies

Most systems rely on all processes in a communicator or process group entering the collective

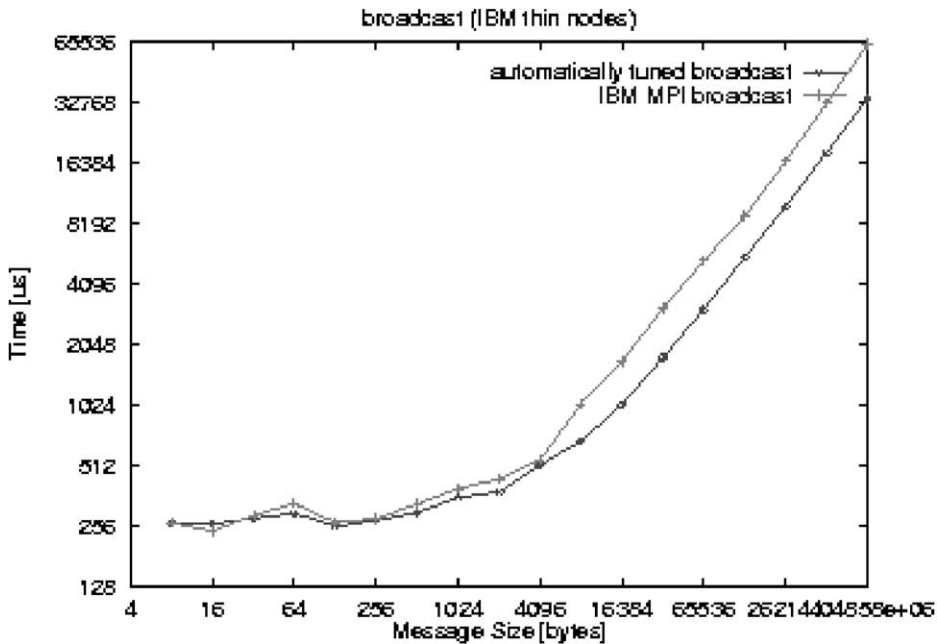


Fig. 3. FT-MPI tuned collective broadcast versus IBM MPI broadcast on IBM SP2 thin node system.

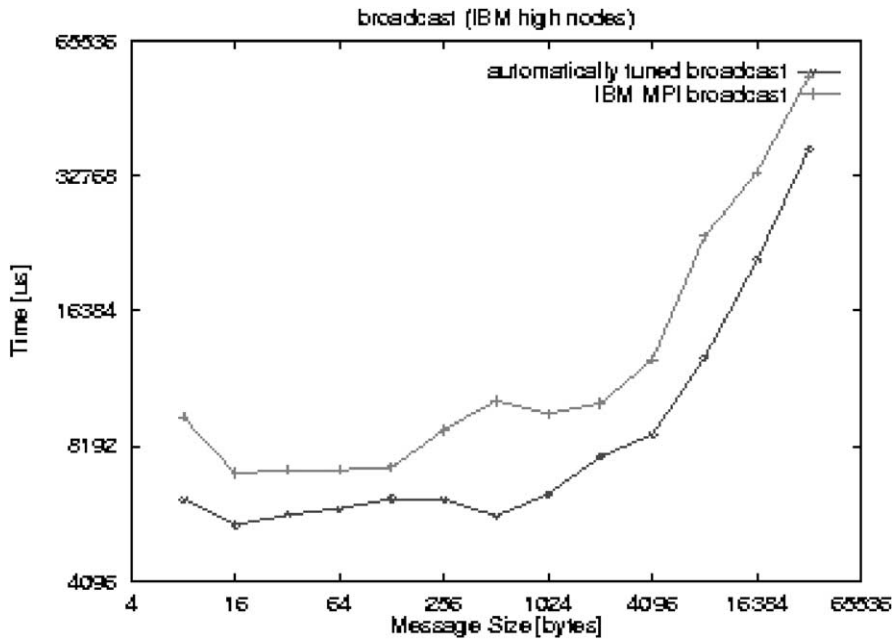


Fig. 4. FT-MPI tuned broadcast versus the IBM MPI broadcast on an SP2 high node system.

communication call synchronously for good performance, i.e. all processes can start the operation without forcing others later in the topology to be delayed. There are some obvious cases where this is not the case:

- (1) The application is executed upon heterogeneous computing platforms where the raw CPU power varies (or load balancing is not optimal).
- (2) The computational cycle time of the application can be non-deterministic as is the case in many of the newer iterative solvers that may converge at different rates continuously.

Even when the application executes in a regular pattern, the physical network characteristics can cause problems with the simple $\log P$ model, such as when running between dispersed clusters. This problem becomes even more acute when the target systems latency is so low that any buffering, while waiting for slower nodes, drastically changes performance characteristics as is the case with BIP-MPI [14] and SCI MPI [8].

FT-MPI can be configured to use a re-ordering strategy that changes the non-root ordering of nodes in a

tree depending on their availability at the beginning of the collective operation. Fig. 5 shows the process by example of a binomial tree.

In Fig. 5, Case 1 is where all processes within the tree are ready to run immediately and thus performance is optimal. In Case 2, both processes B and C are delayed and initially the root A can only send to D. As B and C become available, they are added to the topology. At this point we have to choose whether

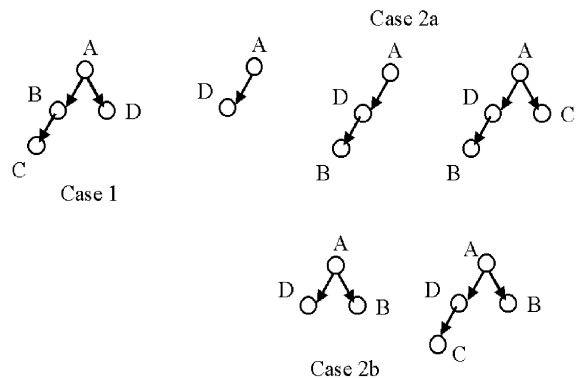


Fig. 5. Re-ordering of a collective topology.

to add the nodes depth first as in Case 2a or breadth first as in Case 2b. Currently, breadth first has given us the best results. Also note that in Case 1, if process B is not ready to receive, it affects not only its own sub-tree, but depending on the message/segment size, it is possible that it would block any other messages that A might send, such as to Ds sub-tree etc. Faster network protocols might not implement non-blocking sends in a manner that could overcome this limitation without affecting the synchronous static optimal case, and thus blocking sends are often used instead.

4.3. Point to point multi-thread communications

FT-MPIs requirements for communications have forced us to use a multi-threaded communications library. The three most important criteria were:

- High performance networking is not affected by concurrent use of slower networking (Myrinet versus Ethernet).
- Non-blocking calls make progress outside of API calls.
- Busy wait (CPU spinning) is avoided within the runtime library.

To meet these requirements, in general, communication requests are passed to a thread via a shared queue to be completed unless the calling thread can complete the operation immediately. Receives are placed into a pending queue by a separate thread. There is one sending and receiving thread per type of communication media, i.e. a thread for TCP communications, a thread for VIA and a thread for handling GM message events. The collective communications are built upon this point to point library.

4.4. Failure detection

It is important to note that the failure handler shown in Fig. 1 gets notification of failures from both the point to point communications libraries as well as the OS support layer. In the case of communication errors, the notify is usually started by the communication library detecting a point to point message not being delivered to a failed party rather than the failed parties OS layer detecting the failure. The handler is responsible for notifying all tasks of errors as they occur by injecting notify messages into the send message queues

ahead of user level messages. An additional daemon know as the FTMPI_NOTIFIER can be used to guarantee ordered delivery of failure notification messages and thus aid in complex debugging.

5. OS support and the HARNESS G_HCORE

When FT-MPI was first designed the only HARNESS Kernel available was an experiment Java implementation from Emory University [5]. Tests were conducted to implement required services on this from C in the form of C-Java wrappers that made RMI calls. Although they worked, they were not very efficient and so FT-MPI was instead initially developed using the readily available PVM system [23].

As the project has progressed, the primary author developed the General HARNESS Core ‘G_HCORE’, a C-based HARNESS core library that uses the same policies as the Java version. This core allowed for services to be built that FT-MPI required as discussed below. Since the G_HCORE, a team at ORNL has developed another C-based H_CORE based on the concept of message handlers [24]. It is expected that the FT-MPI service plug-ins described below will ported to this new H_CORE during early 2002.

5.1. G_HCORE design and performance

The core is built as a daemon wrote in C code that provides a number of very simple services that can be dynamically added to [1]. The simplest service is the ability to load additional code in the form of a dynamic library (shared object) and make this available to either a remote process or directly to the core itself. Once the code is loaded it can be invoked using a number of different techniques such as:

- *Direct invocation.* The core calls the code as a function, or a program uses the core as a runtime library to load the function, which it then calls directly itself.
- *Indirect invocation.* The core loads the function and then handles requests to the function on behalf of the calling program, or, it sets the function up as a separate service and advertises how to access the function.

Table 2
Performance of various invocation methods in milliseconds

	Local (direct invocation)	Local (via TCP/ sockets to core)	Local (new thread)	Remote (RMI)	Remote (TCP)	Remote (UDP)
Emory Java DVM	–	10.4	0.172	1.406	8.6	–
G.HCORE	0.0021	0.58	0.189	–	1.17	0.32

The indirect invocation method allows a range of options such as:

- The H.GCOREs main thread calls the function directly.
- The H.GCORE hands the function call over to a separate thread per invocation.
- H.GCORE forks a new process to handle the request (once per invocation).
- H.GCORE forks a new handler that only handles that type of request (multi-invocation service).

Remote invocation services only provide very simple marshalling of argument lists. The simplest call format passes the socket of the request caller to the plug-in function which is then responsible for marshalling its own input and output much like skeleton functions under SUN RPC. Currently, the indirect remote invocation services are callable via both the UDP and TCP protocols. Table 2 contains performance details of the G.HCORE compared to the Java-based Emory DVM system tested on a Linux cluster over 100 MB Second fast Ethernet.

From Table 2 we can see that socket invocation under Java performs poorly, although RMI is comparable to C socket code for remote invocation. The fastest remote invocation method is via UDP on the G.HCORE at just over 300 μ s per end to end invocation.

5.2. G.HCORE plug-in management and invocation

The plug-ins used by the G.HCORE can either be located locally via a mounted file system or downloaded via HTTP from a web repository. This loading scheme is very similar to that used by JAVA. The search actions are as follows:

- The local file system is checked first in a directory constructed from the plug-in name, i.e. package FT_MPI, might have a component TCP_COMS. Thus the G.HCORE would first

look in <HARNESS_ROOT>/lib/FT_MPI for a TCP_COMS shared object.

- If the plug-in was not in its correct location a search of the temporary cache directory would occur, i.e. <HARNESS_ROOT>/cache/lib/FT_MPI. If the plug-in was found its time to live index would be checked to see if it was still current.
- If the plug-in was not local, then the internal system “get by HTTP” routine would be used. This currently functions with either a pure download, as in just a shared object stored in native format on a remote web server, or with a complex download that contains a PGP signed MIME encoded plug-in. This later method allows for signed plug-ins that are protected against external tampering once they are published.

The locations and calling parameters of the plug-ins available are stored within a distributed replicated database (DRD). This information is pushed to the database when plug-ins are ‘published’ at the individual web servers. The use of standard web servers for plug-in distribution was chosen to aid in individual site deployment of HARNESS, as existing servers can be used without modification.

The API to manage and invoke plug-ins know as the Hlib is based on five classes of operation:

- Locating and loading of a package (set of functions/components) into either local memory or the memories of a group of HCORES.
- Accessing a calling method, which involves obtaining a handle to the required functions.
- Status and enquiry calls. An example enquiry function is to ask if a plug-in component of name X has been loaded into a HARNESS virtual machine Y.
- Invoking a method. This could be done either directly as a pointer reference to a function call, or via the Hcore library. Calling via the library allows for synchronous, asynchronous and parallel execution on multiple Hcores.

- Data marshalling routines to handle the transfer of data to and from remote methods. Current implementation supports three types of data encoding such none (when a call returns a socket for passing of arguments), user controlled pack/unpack (like PVM messages), automatic (like VARARGS/RPC).

The following code example demonstrates how a plug-in is loaded and invoked. In this example it is important to note that the `load_package()` call needs a target virtual machine, and after it has located and loaded the package it runs any init routines provided by the package. These routines can include Hlib query routines that check if other required routines are loaded and if not load them as required forming a simple load dependency tree. The load routine also needs to know that it is loading only a single instance on a single host. The get function reference routine also states that it wants just a single reference, and in this example the first routine was found. The call function invokes the method and expects to receive its output via the socket descriptor returned (Scheme 3).

5.3. G-HCORE services for FT-MPI

Current services required by FT-MPI break down into six categories:

- *Spawn and notify service.* This plug-in allows remote processes to be initiated and then monitored. The service notifies other interested processes when a failure or exit of the invoked process occurs.
- *Event distribution.* This is a centralized service that can distribute serialized events between Hcores and other interested individual processes. This is used by FT-MPI to simplify failure notification so that changes to notify lists (which match MPI communicator membership) only has to be stored in a few locations rather than at every notification service daemon.

- *Naming services.* These allocate unique identifiers in a distributed environment.
- *DRD.* This service allows for system state and additional MetaData to be distributed, with replication specified at the record level. This plug-in has a secondary benefit as it can be used by the Emory DVMs PVM plug-in to implement the PVM 3.4 Mailbox features directly.
- *Specialized communication library access and initialization.* For example, TCP is built into the FT-MPI library, but shared memory and myrinet access is not. These are built as separate objects that are loaded only if needed.

FT-MPI uses the Hcore loading of plug-ins at two separate levels. The first is at the Hcore daemon level and are accessed via remote invocation. The spawn and notify services are such an example. Each Hcore that is expected to run an FT-MPI job would load this service via a request from the Hlib. For example, the FT-MPIRUN utility queries the Hcore for all instances of the spawn and notify services when deciding where to run an FT-MPI application, before it sends requests to the spawn services to start execution. The naming services, e.g. would also execute under a Hcore daemon, but in this case the service would enforce that only one instance is loaded and executes.

The second level of plug-in are the specialized communications plug-ins. These plug-ins are loaded directly into the FT-MPI applications memory space so that they can be accessed directly by the FT-MPI runtime library. Without this direct access, all MPI message data using these devices would have to pass through the `g_hcore` daemons incurring serious performance penalties. Run time swapping of these plug-ins requires the use of message queue flush and hold protocols to prevent message loss. The same protocols are used for some of the FT-MPI communicator rebuild operations. Fig. 6 gives the overall scheme for data and control message flow within an FT-MPI application.

```
rc = load_package (myVM, HLIB_SINGLE, "package", "component",
&package, &ncalls);
rc = get_function ("function long_name", package, HLIB_FIRST, &f_info);
rc = call_function (f_info.call_id, &s, NULL);
```

Scheme 3. Loading a package dynamically onto a HARNESS VM via the Hlib API.

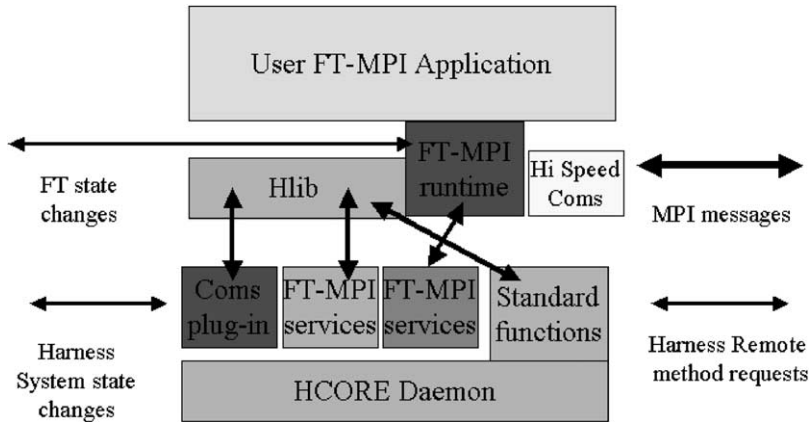


Fig. 6. Overview of FT-MPI application use of HARNCESS services and runtime library.

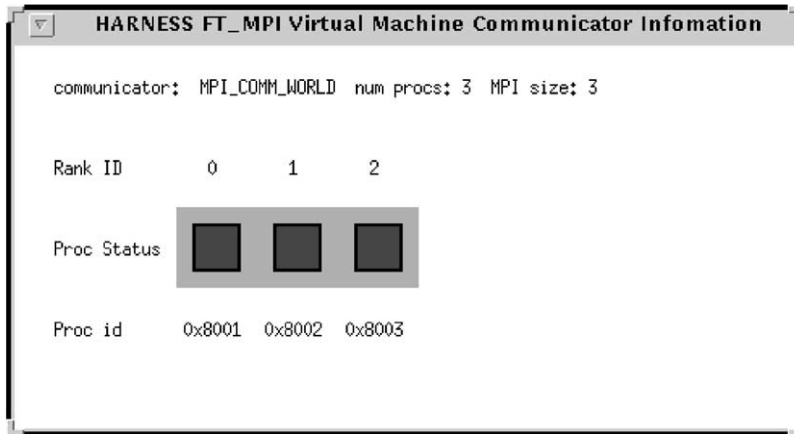


Fig. 7. COMINFO display for a healthy three process MPI application. The colors of the inner boxes indicate the state of the processes and the outer box indicates the communicator state.

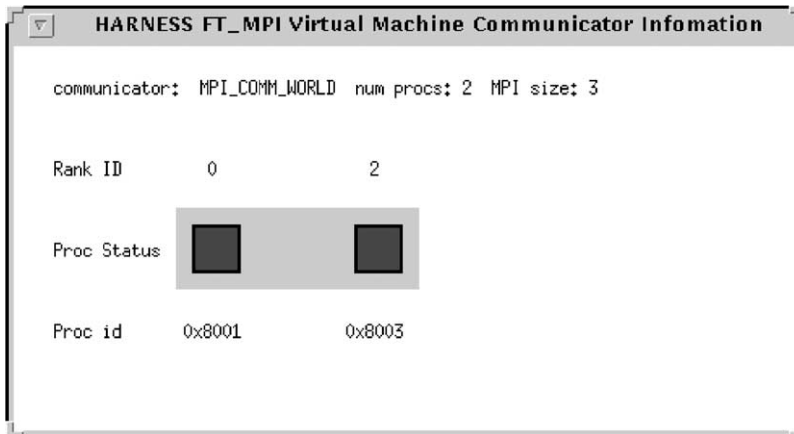


Fig. 8. COMINFO display for an application with an exited process. Note that the number of nodes and size of communicator do not match.

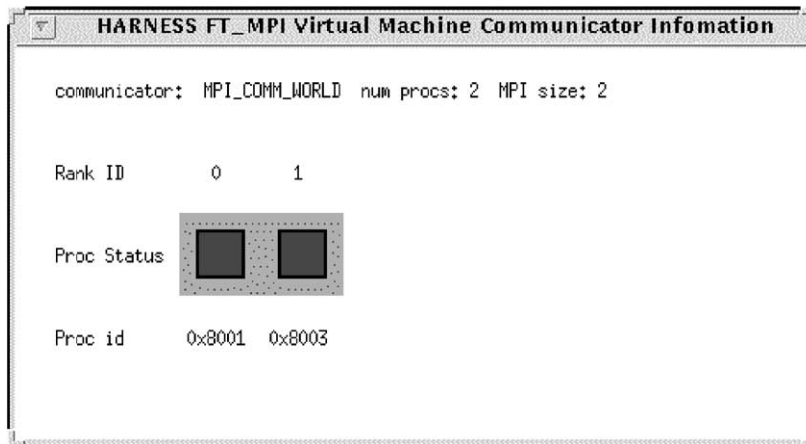


Fig. 9. COMINFO display for the above application after a communicator rebuild using the SHRINK option. Note the communicator status box has changed back to a blue (dark) color.

6. FT-MPI tool support

Current MPI debuggers and visualization tools such as totalview, vampir, upshot, etc. do not have a concept of how to monitor MPI jobs that change their communicators on the fly, nor do they know how to monitor a virtual machine. To assist users in understanding these the author has implemented two monitor tools. HOSTINFO which displays the state of the virtual machine. COMINFO which displays processes and communicators in color-coded fashion so that users know the state of an applications processes and communicators. Both tools are currently built using the X11 libraries but will be rebuilt using the Java SWING system to aid portability. An example displays during a SHRINK communicator rebuild operation is shown in Figs. 7–9, where a process (rank 1) exits and the communicator is reduced in size and extent.

7. Conclusions

FT-MPI is an attempt to provide application programmers with different methods of dealing with failures within MPI application than just check-point and restart. It is hoped that by experimenting with FT-MPI, new applications methodologies and algorithms will be developed to allow for both high performance and the survivability required by both unreliable GRIDs and the next generation of terra-flop

and beyond machines. In the case of GRIDs, newly developed applications could be designed to require less storage of state thus removing the requirement that any check-pointed application state has to be maintained on both highly reliable and highly available storage further burdening a potentially unreliable system.

FT-MPI in itself is already proving to be a useful vehicle for experimenting with self-tuning collective communications, distributed control algorithms, various dynamic library download methods and improved sparse data handling subsystems, as well as being the default MPI implementation for the HARNESS project.

Future work in the FT-MPI library system will concentrate on developing both further implementations that support both common GRID services as well as more restrictive environments such as dedicated MPPs runtimes and embedded clusters. Application development will be supported by the creation of a number of drop-in library templates or skeletons to simplify the construction of fault tolerant applications.

References

- [1] M. Beck, J.J. Dongarra, G.E. Fagg, G.A. Geist, P. Gray, J. Kohl, M. Migliardi, K. Moore, T. Moore, P. Papadopoulos, S. Scott, V. Sunderam, HARNESS: a next generation distributed virtual machine, *J. Future Generation Comput. Syst.* 15 (1999) 571–582.

- [2] G. Stellner, CoCheck: Checkpointing and process migration for MPI, in: Proceedings of the International Parallel Processing Symposium, Honolulu, April 1996, pp. 526–531.
- [3] A. Agbaria, R. Friedman, Starfish: Fault-tolerant dynamic MPI programs on clusters of workstations, in: Proceedings of the Eighth IEEE International Symposium on High Performance Distributed Computing, 1999.
- [4] G.E. Fagg, K. Moore, J.J. Dongarra, Scalable networked information processing environment (SNIPE), *J. Future Generation Comput. Syst.* 15 (1999) 595–605.
- [5] M. Migliardi, V. Sunderam, PVM emulation in the HARNESS MetaComputing system: a plug-in based approach, *Lecture Notes in Computer Science*, Vol. 1697, September 1999, pp. 117–124.
- [6] P.H. Worley, I.T. Foster, B. Toonen, Algorithm comparison and benchmarking using a parallel spectral transform shallow water model, in: G.-R. Hoffmann, N. Kreitz (Eds.), Proceedings of the Sixth Workshop on Parallel Processing in Meteorology, World Scientific, Singapore, 1995, pp. 277–289.
- [7] T. Kielmann, H.E. Bal, S. Gorlatch, Bandwidth-efficient collective communication for clustered wide area systems, in: Proceedings of the IPDPS 2000, Cancun, Mexico, May 1–5, 2000.
- [8] L.P. Huse, Collective communication on dedicated clusters of workstations, in: Proceedings of the Sixth European PVM/MPI Users' Group Meeting, *Lecture Notes in Computer Science*, Vol. 1697, Springer, Barcelona, September 1999, pp. 469–476.
- [9] D. Culler, R. Karp, D. Patterson, A. Sahay, K.E. Schauer, E. Santos, R. Subramonian, T. von Eicken, *LogP*: towards a realistic model of parallel computation, in: Proceedings of the Symposium on Principles and Practice of Parallel Programming (PpoPP), San Diego, CA, May 1993, pp. 1–12.
- [10] R. Rabenseifner, A new optimized MPI reduce algorithm, 1997. http://www.hlr.de/structure/support/parallel_computing/models/mpl/myreduce.html.
- [11] M. Snir, S. Otto, S. Huss-Lederman, D. Walker, J. Dongarra, MPI: The Complete Reference, Vol. 1, The MPI Core, second ed., 1998.
- [12] M. Frigo, FFTW: an adaptive software architecture for the FFT, in: Proceedings of the ICASSP Conference, Vol. 3, 1998, p. 1381.
- [13] R.C. Whaley, J. Dongarra, Automatically tuned linear algebra software, SC98: high performance networking and computing, 1998. <http://www.cs.utk.edu/~rwhaley/ATL/INDEX.HTM>.
- [14] L. Prylli, B. Tourancheau, BIP: a new protocol designed for high performance networking on myrinet, in: Proceedings of the PC-NOW workshop, IPPS/SPDP 1998, Orlando, USA, 1998.
- [15] S. Louca, N. Neophytou, A. Lachanas, P. Evripidou, MPI-FT: a portable fault tolerance scheme for MPI, in: Proceedings of PDPTA'98 International Conference, Las Vegas, Nevada, 1998.
- [16] J.L. Traff, R. Hempel, H. Ritzdort, F. Zimmermann, Flattening on the fly: efficient handling of MPI derived datatypes, in: Proceedings of the Sixth European PVM/MPI Users' Group Meeting, *Lecture Notes in Computer Science*, Vol. 1697, Springer, Barcelona, September 1999, pp. 109–116.
- [17] W.D. Gropp, E. Lusk, D. Swider, Improving the performance of MPI derived datatypes, in: Proceedings of the Third MPI Developer's and User's Conference (MPIDC'99), 1999, pp. 25–30.
- [18] G.E. Fagg, K.S. London, J.J. Dongarra, MPIConnect, managing heterogeneous MPI application interoperation and process control, EuroPVM-MPI 98, *Lecture Notes in Computer Science*, Vol. 1497, Springer, 1998, pp. 93–96.
- [19] E. Gabriel, M. Resch, T. Beisel, R. Keller, Distributed computing in a heterogeneous computing environment, EuroPVM-MPI 98, *Lecture Notes in Computer Science*, Vol. 1497, Springer, Berlin, 1998, pp. 180–187.
- [20] S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, R. Whaley, Scalapack: a linear algebra library for message-passing computers, in: Proceedings of the 1997 SIAM Conference on Parallel Processing, May 1997.
- [21] J. Choi, J. Dongarra, S. Ostrouchov, A. Petitet, D. Walker, R. Whaley, A proposal for a set of parallel basic linear algebra subprograms, LAPACK Working Note #100, CS-95-292, May 1995.
- [22] W. Gropp, E. Lusk, R. Thakur, Using MPI-2: Advanced Features of the Message Passing Interface, first ed., MIT Press, Cambridge, MA, February 2000.
- [23] Al Geist, A. Beguelin, J. Dongarra, W. Jiang, R. Manchek, V. Sunderam, PVM—parallel virtual machine: a users' guide and tutorial for networked parallel computing, Scientific and Engineering Computation, MIT Press, Cambridge, MA, 1994.
- [24] W.R. Elwasif, D.E. Bernholdt, J.A. Kohl, G.A. Geist, An architecture for a multi-threaded harness kernel, in: Proceedings of the Eighth EuroPVM/MPI Users Group Meeting, LNCS, Vol. 2131, Springer, Berlin, 2001, pp. 126–134.
- [25] F. Berman, A. Chen, K. Cooper, J. Dongarra, I. Foster, D. Gannon, L. Johnsson, K. Kennedy, C. Kesselman, J. Mellow-Crummey, D. Reed, L. Torczon, R. Wolski, The GrADS Project, *Int. J. High Perform. Comput. Appl.* 15 (4) (2001) 327–344.
- [26] A. Petitet, S. Blackford, J. Dongarra, B. Ellis, G. Fagg, K. Roche, S. Vadhiyar, Numerical libraries and the grid, *Int. J. High Perform. Comput. Appl.* 15 (4) (2001) 359–374.
- [27] I. Foster, C. Kesselmann, The globus toolkit, in: I. Foster, C. Kesselmann (Eds.), *The GRID: Blueprint for a New Computing Infrastructure*, Morgan Kaufmann, San Francisco, 1999, pp. 259–278.
- [28] I. Foster, N. Karonis, A grid enabled MPI: message passing in heterogeneous distributed computing systems, in: Proceedings of the SuperComputing'98 (SC98), Orlando, FL.
- [29] J.S. Plank, M. Beck, W.R. Elwasif, T. Moore, M. Swany, R. Wolski, The internet backplane protocol: storage in the network, in: Proceedings of the NetStore99: The Network Storage Symposium, Seattle, WA, 1999.
- [30] I. Foster, C. Kesselmann, *The GRID: Blueprint for a New Computing Infrastructure*, Morgan Kaufmann, San Francisco, 1999.



Graham Fagg received his BSc in Computer Science and Cybernetics from the University of Reading (UK) in 1991 and a PhD in Computer Science in 1998. From 1991 to 1993, he worked on CASE tools for interconnecting array processors and Transputer MIMD systems. From 1994 to 1995 he was a research assistant at the Cluster Computing Laboratory at the University of Reading working on code generation tools for collective/group communication.

From 1996 to 1999 he worked as a senior research associate at the University of Tennessee. Since 1999 he has been a Research Assistant Professor. Currently he is a visiting guest scientist at the High Performance Computing Center Stuttgart (HLRS). His current research interests include distributed scheduling, resource management, performance prediction, benchmarking, cluster management tools, and high speed networking. He is currently involved in the development of a number of metacomputing and GRID middle-ware systems including SNIPE, MPI_Connect, HARNES, and a fault tolerant MPI implementation (FT-MPI).



Jack Dongarra holds an appointment as University Distinguished Professor of Computer Science in the Computer Science Department at the University of Tennessee. He specializes in numerical algorithms in linear algebra, parallel computing, use of advanced computer architectures, programming methodology, and tools for parallel computers. Other current research involves the development, testing

and documentation of high quality mathematical software. He was involved in the design and implementation of the open source software packages EISPACK, LINPACK, the BLAS, LAPACK, ScaLAPACK, Netlib, PVM, MPI, NetSolve, ATLAS, PAPI, and HARNES; and is currently involved in the design of algorithms and techniques for high performance computer architectures. He is a fellow of the AAAS, ACM, and IEEE and a member of the National Academy of Engineering.