



ELSEVIER

Parallel Computing 000 (2001) 000–000

PARALLEL
COMPUTINGwww.elsevier.com/locate/parco

HARNESS and fault tolerant MPI

Graham E. Fagg*, Antonin Bukovsky, Jack J. Dongarra

*Department of Computer Science, University of Tennessee, Suite 203,
1122 Volunteer Blvd., Knoxville, TN 37996-3450, USA*

Received 6 February 2001; received in revised form 19 April 2001

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 would have affected their performance. As current HPC systems increase in size with greater potential levels of individual node failure, the need arises for new fault tolerant systems to be developed. Here we present a new implementation of MPI called fault tolerant MPI (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, debugging tools and some performance issues. Also discussed is the experimental HARNESS core (G_HCORE) implementation that FT-MPI is built to operate upon. © 2001 Published by Elsevier Science B.V.

Keywords: Message passing; Parallel computing; Fault tolerant application; Metacomputing middleware

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

* Corresponding author.

E-mail addresses: fagg@cs.utk.edu, snipe@cs.utk.edu (G.E. Fagg).

28 for dynamic process control, although this was limited to the creation of new MPI
29 process groups with separate communicators. These new processes could not be
30 merged with previously existing communicators to form intracommunicators needed
31 for a seamless single application model and were limited to a special set of extended
32 collectives (group) communications.

33 The MPI static process model suffices for small numbers of distributed nodes
34 within the currently emerging masses of clusters and several hundred nodes of
35 dedicated MPPs. Beyond these sizes the mean time between failure (MTBF) of CPU
36 nodes starts becoming a factor. As attempts to build the next generation Peta-flop
37 systems advance, this situation will only become more adverse as individual node
38 reliability becomes out weighted by orders of magnitude increase in node numbers
39 and hence node failures.

40 The aim of FT-MPI is to build a fault tolerant MPI implementation that can
41 survive failures, while offering the application developer a range of recovery options
42 other than just returning to some previous check-pointed state. FT-MPI is built on
43 the HARNESS [1] meta-computing system, and is meant to be used as its default
44 application level message passing interface.

45 2. Check-point and roll back versus replication techniques

46 The first method attempted to make MPI applications fault tolerant was through
47 the use of check-pointing and roll back. Co-Check MPI [2] from the Technical
48 University of Munich being the first MPI implementation built that used the Condor
49 library for check-pointing an entire MPI application. In this implementation, all
50 processes would flush their messages queues to avoid in flight messages getting lost,
51 and then they would all synchronously check-point. At some later stage if either an
52 error occurred or a task was forced to migrate to assist load balancing, the entire
53 MPI application would be rolled back to the last complete check-point and be re-
54 started. This systems main drawback being the need for the entire application having
55 to check-point synchronously, which depending on the application and its size could
56 become expensive in terms of time (with potential scaling problems). A secondary
57 consideration was that they had to implement a new version of MPI known as
58 tuMPI as retro-fitting MPICH was considered too difficult.

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

66 The project closest to FT-MPI known to the author is the implicit fault tolerance
67 MPI project MPI-FT [15] by Paraskevas Evripidou of Cyprus University. This project
68 supports several master–slave models where all communicators are built from grids
69 that contain ‘spare’ processes. These spare processes are utilized when there is a failure.

70 To avoid loss of message data between the master and slaves, all messages are copied to
71 an observer process, which can reproduce lost messages in the event of any failures.
72 This system appears only to support SPMD style computation and has a high overhead
73 for every message and considerable memory needs for the observer process for long
74 running applications. This system is not a full check-point system in that it assumes any
75 data (or state) can be rebuilt using just the knowledge of any passed messages, which
76 might not be the case for non deterministic unstable solvers.

77 FT-MPI has much lower overheads compared to the above check-pointing system,
78 and thus much higher potential performance. These benefits do, however, have con-
79 sequences. An application using FT-MPI has to be designed to take advantage of its
80 fault tolerant features as shown in the next section, although this extra work can be
81 trivial depending on the structure of the application. If an application needs a high level
82 of fault tolerance where node loss would equal data loss, then the application has to be
83 designed to perform some level of user directed check-pointing. FT-MPI does allow for
84 atomic communications much like StarFish, but unlike in StarFish, the level of cor-
85 rectness can be varied on for individual communicators. This provides users the ability
86 to fine tune for coherency or performance as system and application conditions dictate.

87 3. FT-MPI semantics

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

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

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

103 The typical MPI semantics is from OK to Failed which then causes an application
104 abort. By allowing the communicator to be in an intermediate state we allow the
105 application the ability to decide how to alter the communicator and its state as well
106 as how communication within the intermediate state behaves.

107 3.1. Failure modes

108 On detecting a failure within a communicator, that communicator is marked as
109 having a probable error. Immediately as this occurs the underlying system sends a

110 state update to all other processes involved in that communicator. If the error was a
 111 communication error, not all communicators are forced to be updated, if it was a
 112 process exit then all communicators that include this process are changed. Note, this
 113 might not be all current communicators as we support MPI-2 dynamic tasks and
 114 thus multiple MPI_COMM_WORLD.

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

118 3.2. Communicator and communication handling

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

- 124 • SHRINK: The communicator is reduced so that the data structure is contiguous.
 125 The ranks of the processes are **changed**, forcing the application to recall MPI_
 126 COMM_RANK.
- 127 • BLANK: This is the same as SHRINK, except that the communicator can now
 128 contain gaps to be filled in later. Communicating with a gap will cause an invalid
 129 rank error. Note also that calling MPI_COMM_SIZE will return the extent of the
 130 communicator, not the number of valid processes within it.
- 131 • REBUILD: Most complex mode that forces the creation of new processes to fill
 132 any gaps until the size is the same as the extent. The new processes can either
 133 be places in to the empty ranks or the communicator can be shrank and the re-
 134 maining processes filled at the end. This is used for applications that require a cer-
 135 tain size to execute as in power of two FFT solvers.
- 136 • ABORT: Is a mode which affects the application immediately an error is detected
 137 and forces a graceful abort. The user is unable to trap this. If the application need
 138 to avoid this they must set all communicators to one of the above communicator
 139 modes.

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

- 142 • NOP: No operations on error. That is no user level message operations are al-
 143 lowed and all simply return an error code. This is used to allow an application
 144 to return from any point in the code to a state where it can take appropriate action
 145 as soon as possible.
- 146 • CONT: All communication that is NOT to the affected/failed node can continue as
 147 normal. Attempts to communicate with a failed node will return errors until the
 148 communicator state is reset.

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

151 3.3. *Point-to-point versus collective correctness*

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

156 Broadcast, gather and all gather demonstrate this perfectly. In Broadcast even if
 157 there is a failure of a receiving node, the receiving nodes still receive the same data,
 158 i.e., the same end result for the surviving nodes. Gather and all-gather are different in
 159 that the result depends on if the problematic nodes sent data to the gatherer/root or
 160 not. In the case of gather, the root might or might not have gaps in the result. For the
 161 all2all operation, which typically uses a ring algorithm it is possible that some nodes
 162 may have complete information and others incomplete. Thus for operations that
 163 require multiple node input as in gather/reduce type operations any failure causes all
 164 nodes to return an error code, rather than possibly invalid data. Currently, an ad-
 165 dition flag controls how strict the above rule is enforced by utilizing an extra barrier
 166 call at the end of the collective call if required.

167 3.4. *FT-MPI usage*

168 Typical usage of FT-MPI would be in the form of an error check and then some
 169 corrective action such as a communicator rebuild. A typical code fragment is shown
 170 in Example 1, where on an error the communicator is simply rebuilt and reused:

171 **Example 1** (*Simple FT-MPI send usage*).

```
172 rc=MPI_Send(----, com);
173 If (rc == MPI_ERR_OTHER)
174 MPI_Comm_dup (com, newcom);
175 com=newcom; /* continue... */
```

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

183 **Example 2** (*FT-MPI Master-Worker code*).

```
184 rc=MPI_Bcast (initial_work...);
185 if(rc == MPI_ERR_OTHER)reclaim_lost_work(...);
186 while (!all_work_done) {
187   if (work_allocated) {
188     rc=MPI_Recv (buf, ans_size, result_dt,
189                 MPI_ANY_SOURCE, MPI_ANY_TAG, comm, & status);
```

```

190     if (rc == MPI_SUCCESS) {
191         handle_work (buf);
192         free_worker (status.MPI_SOURCE);
193         all_work_done--;
194     }
195     else {
196         reclaim_lost_work (status.MPI_SOURCE);
197         if (no_surviving_workers) { /* !do something ! */ }
198     }
199     /* work allocated */
200     /* Get a new worker as we must have received a result or a
201    death */
202     rank = get_free_worker_and_allocat_work ();
203     if (rank) {
204         rc = MPI_Send (... rank...);
205         if (rc == MPI_OTHER_ERR) reclaim_lost_work (rank);
206         if (no_surviving_workers) { /* !do something ! */ }
207     } /* if free worker */
208     /* while work to do */

```

209 4. FT_MPI implementation details

210 FT-MPI is a partial MPI-2 implementation. It currently contains support for both
 211 C and Fortran interfaces, all the MPI-1 function calls required to run both the
 212 PSTSWM [6] and BLAS [21] applications. BLAS is supported so that SCALAPACK
 213 [20] applications can be tested. Currently only some of the dynamic process control
 214 functions from MPI-2 are supported.

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

219 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.

223 4.1. Derived data type handling

224 MPI-1 introduced extensive facilities for user derived data type (DDT)[11] han-
 225 dling that allows for in effect strongly typed message passing. The handling of these
 226 possibly non-contiguous data types is very important in real applications, and is
 227 often a neglected area of communication library design [17]. Most communications
 228 libraries are designed for low latency and/or high bandwidth with contiguous blocks

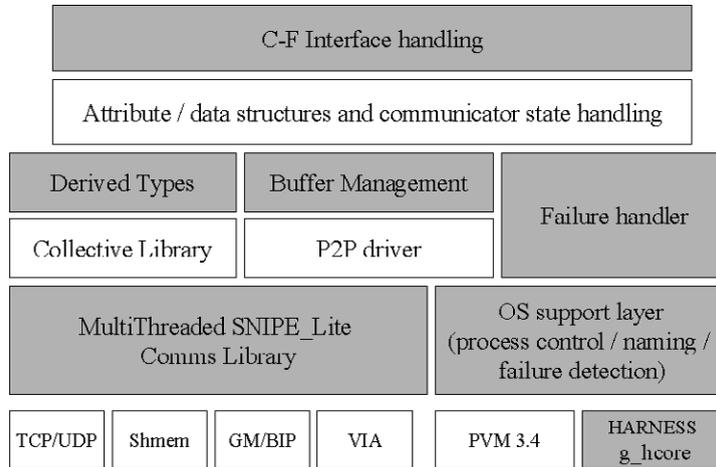


Fig. 1. Overall structure of the FT-MPI implementation.

229 of data [14]. Although this means that they must avoid unnecessary memory copies,
 230 the efficient handling of recursive data structures is often left to simple iterations of a
 231 loop that packs a send/receive buffer.

232 *4.1.1. FT-MPI DDT handling*

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

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

- 246 • *Gather/scatter.* Data is collected into or from recursively structured non-contigu-
 ous memory.
- 248 • *Encoding/decoding.* Data passed between heterogeneous machine architectures
 249 than use different floating point representations need to be converted so that the
 data maintains the original meaning.
- 251 • *Send/receive packetizing.* All of the send or receive cannot be completed in a single
 252 attempt and the data has to be sent in blocks. This is usually due to buffering con-
 straints in the communications library/OS or even hardware flow control.

253 4.1.2. DDT methods and algorithms

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

262 When the DDT is being processed the actual user data itself can also be com-
 263 pacted into/from a contiguous buffer. Several options for this type of buffering are
 264 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.

267
 269 The minimal and no padded methods are used when moving the data type within
 270 a homogeneous set of machines that require no numeric representation encoding or
 271 decoding. The zero padding method benefits slower networks, and alignment padded
 272 can in some cases assist memory copy operations, although its real benefit is when
 273 used with re-ordering.

274
 275 The re-ordered compacting method shown in Fig. 2, for a data type that consists
 276 of characters (C) and integers (I). This type of compacting is designed to be used
 277 when some additional form encoding/decoding takes place. In particular, moving the
 278 re-ordered data type by type through fixed XDR/Swap buffers improves its perfor-
 279 mance considerably. Two types of DDT encoding are supported, the first is the

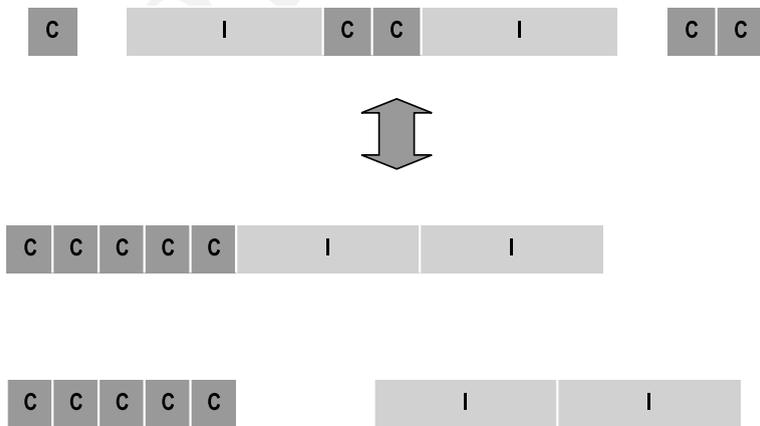


Fig. 2. Compacting storage of re-ordered DDT. Without padding, and with correct alignment.

280 slower generic SUN XDR format and the second is simple byte swapping to convert
281 between little and big endian numbers.

282 4.1.3. FT-MPI DDT performance

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

289 The tests show that the compacted data type handling gives from 10% to 19%
290 improvement for small messages and 78–81% for larger arrays on same numeric
291 representation machines. The benefits of buffer reuse and re-ordered of data elements
292 leads to considerable improvements on heterogeneous networks as shown in the
293 comparison between the default MPICH operations and the DDT ByteSwap and
294 DDT XDR results. It is important to note that all these tests used MPICH to
295 perform the point-to-point communication, and thus the overlapping of the data
296 gather/scatter, encoding/decoding and non-blocking communication is not shown
297 here, and is expected to yield even higher performance.

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

305 4.2. Collective communications

306 The performance of the MPI's collective communications is critical to most MPI-
307 based applications [6]. A general algorithm for a given collective communication
308 operation may not give good performance on all systems due to the differences in

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

Type of operation (arch 2 arch) (method) (encoding)	11,956 bytes B/W MB/s	% compared to MPICH	95,648 bytes B/W MB/s	% compared to MPICH
Sparc 2 Sparc MPICH	5.49		5.47	
Sparc 2 Sparc DDT	6.54	+19	9.74	+78
Linux 2 Linux MPICH	7.11		8.79	
Linux 2 Linux DDT	7.87	+10	9.92	+81
Sparc 2 Linux MPICH	0.855		0.729	
Sparc 2 Linux DDT Byte Swap	5.87	+586	8.20	+1024
Sparc 2 Linux DDT XDR	5.31	+621	6.15	+743

309 architectures, network parameters and the storage capacity of the underlying MPI
 310 implementation [7]. In an attempt to improve over the usual collective library built
 311 on point-to-point communications design as in the logP model [9], we built a col-
 312 lective communications library that is tuned to its target architecture though the use
 313 a limited set of micro-benchmarks. Once the static system is optimized, we then tune
 314 the topology dynamically by re-orders the logical addresses to compensate for
 315 changing run time variations. Other projects that use a similar approach to opti-
 316 mizing include [12,13].

317 4.2.1. Collective communication algorithms and benchmarks

318 The micro-benchmarks are conducted for each of the different classes of MPI
 319 collective operations broadcast, gather, scatter, reduce, etc., individually. That is, the
 320 algorithm that produces the best broadcast might not produce the best scatter even
 321 though they appear similar.

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

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

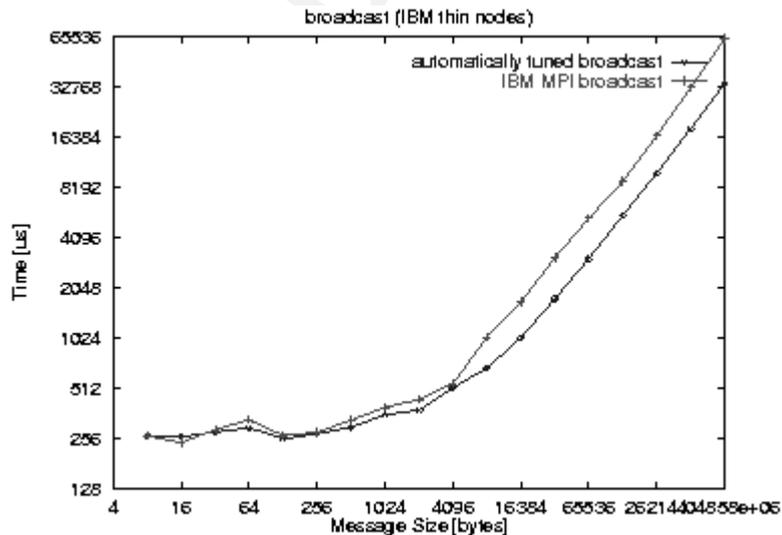


Fig. 3. FT-MPI tuned collective broadcast versus IBM MPI broadcast on IBM SP2 thin node system.

333 4.2.2. Dynamic re-ordering of topologies

334 Most systems rely on all processes in a communicator or process group entering
 335 the collective communication call synchronously for good performance, i.e., all
 336 processes can start the operation without forcing others later in the topology to be
 337 delayed. There are some obvious cases where this is not the case:

338 (1) The application is executed upon heterogeneous computing platforms where
 339 the raw CPU power varies (or load balancing is not optimal).

340 (2) The computational cycle time of the application can be non-deterministic as is
 341 the case in many of the newer iterative solvers that may converge at different rates
 342 continuously.

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

349 FT-MPI can be configured to use a reordering strategy that changes the non-root
 350 ordering of nodes in a tree depending on their availability at the beginning of the
 351 collective operation. Fig. 4 shows the process by example of a binomial tree.

352 In Fig. 5, Case 1 is where all processes within the tree are ready to run immediately
 353 and thus performance is optimal. In Case 2, both processes B and C are delayed
 354 and initially the root A can only send to D. As B and C become available, they are
 355 added to the topology. At this point we have to choose whether to add the nodes
 356 depth first as in Case 2a or breadth first as in Case 2b. Currently, breadth first has
 357 given us the best results. Also note that in CASE 1, if process B is not ready to

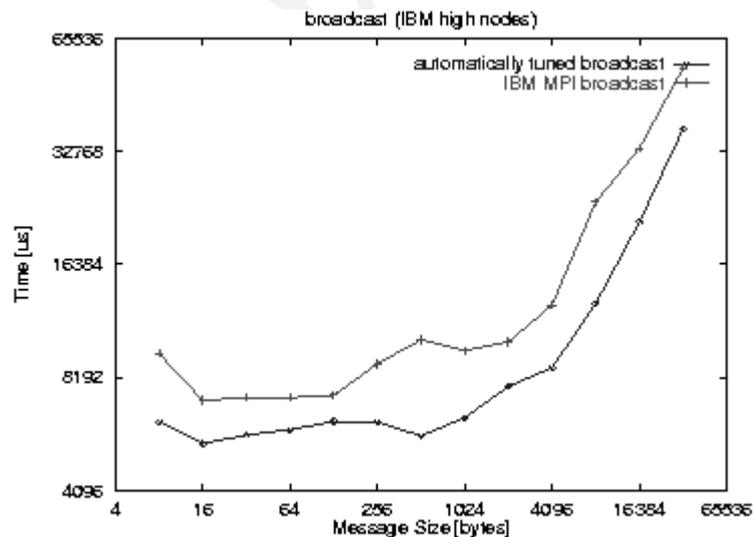


Fig. 4. FT-MPI tuned broadcast versus the IBM MPI broadcast on an SP2 high node system.

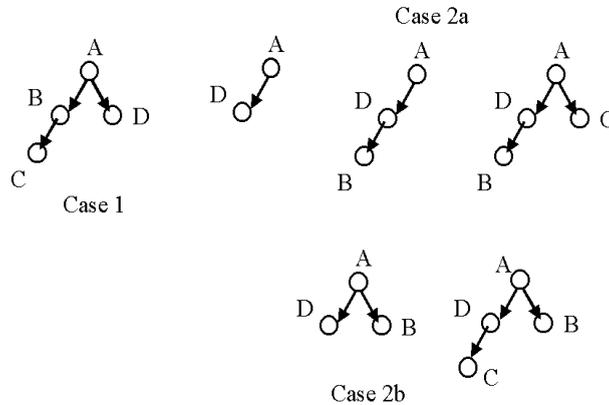


Fig. 5. Re-ordering of a collective topology.

358 receive, it affects not only its own sub-tree, but depending on the message/segment
 359 size, it is possible that it would block any other messages that A might send, such as
 360 to D's sub-tree etc. Faster network protocols might not implement non-blocking
 361 sends in a manner that could overcome this limitation without affecting the syn-
 362 chronous static optimal case, and thus blocking sends are often used instead.

363 *4.3. Point-to-point multi-thread communications*

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

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

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

377 *4.4. Failure detection*

378 It is important to note that the failure handler shown in Fig. 1, gets notification of
 379 failures from both the point-to-point communications libraries as well as the OS
 380 support layer. In the case of communication errors, the notify is usually started by
 381 the communication library detecting a point-to-point message not being delivered to
 382 a failed party rather than the failed parties OS layer detecting the failure. The

383 handler is responsible for notifying all tasks of errors as they occur by injecting
384 notify messages into the send message queues ahead of user level messages.

385 5. OS support and the HARNESS G_HCORE

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

391 As the project has progressed, the primary author developed the G_HCORE, a C
392 based HARNESS core library that uses the same policies as the Java version. This
393 core allows for services to be built that FT-MPI requires.

394 5.1. G_HCORE design and performance

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

- 401 • *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.
- 403 • *Indirect invocation.* The core loads the function and then handles requests to the
404 function on behalf of the calling program, or, it sets the function up as a separate
service and advertises how to access the function.

405 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);
- 410 • H_GCORE forks a new handler that only handles that type of request (multi-in-
vocation service).

411 Remote invocation services only provide very simple marshalling of argument
412 lists. The simplest call format passes the socket of the request caller to the plug-in
413 function which is then responsible for marshalling its own input and output much
414 like skeleton functions under SUN RPC.

415 Currently, the indirect remote invocation services are callable via both the UDP
416 and TCP protocols. Table 2 contains performance details of the G_HCORE com-
417 pared to the Java based Emory DVM system tested on a Linux cluster over 100
418 Mbytes s fast Ethernet.

419 From Table 2, we can see that socket invocation under Java performs poorly,
420 although RMI is comparable to C socket code for remote invocation. The fastest
421 remote invocation method is via UDP on the G_HCORE at just over 300 ms per end
422 to end invocation.

Table 2
Performance of various invocation methods in ms.

	Local (direct invocation)	Local (via TCP/sockets to core)	Local (new thread)	Remote (RMI)	Remote (TCP)	Remote (TCP)
Emory Java DVM	-/-	10.4	0.172	1.406	8.6	-/-
G_HCORE	0.0021	0.58	0.189	-/-	1.17	0.32

423 5.2. *G_HCORE plug-in management*

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

- 428 • The local file system is checked first in a directory constructed from the plug-in
429 name. That is, Package FT_MPI, might have a component TCP_COMS. Thus
430 the G_HCORE would first look in <HARNESS_ROOT>/lib/FT_MPI for a
TCP_COMS shared object.
- 432 • If the plug-in was not in its correct location a search of the temporary cache direc-
433 tory 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.
- 435 • If the plug-in was not local, then the internal system “get by HTTP” routine
436 would be used. This currently functions with either a pure download, as in just
437 a shared object stored in native format on a remote web server or with a complex
438 download that contains a PGP signed MIME encoded plug-in. This later method
439 allows for signed plug-ins that are protected against external tampering once they
are published.

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

445 5.3. *G_HCORE services for FT-MPI*

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

- 448 • *Spawn and Notify service*. This plug-in allows remote processes to be initiated and
449 then monitored. The service notifies other interested processes when a failure or
exit of the invoked process occurs.
- *Naming services*. These allocate unique identifiers in a distributed environment.
- 452 • *Distributed replicated database (DRD)*. This service allows for system state and ad-
453 ditional MetaData to be distributed, with replication specified at the record level.
454 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.

455 **6. FT-MPI tool support**

456 Current MPI debuggers and visualization tools such as totalview, vampir, upshot,
 457 etc., do not have a concept of how to monitor MPI jobs that change their com-
 458 municators on the fly, nor do they know how to monitor a virtual machine. To assist
 459 users in understanding these the author has implemented two monitor tools.
 460 HOSTINFO which displays the state of the Virtual Machine. COMINFO which
 461 displays processes and communicators in colour coded fashion so that users know
 462 the state of an applications processes and communicators. Both tools are currently
 463 built using the X11 libraries but will be rebuilt using the Java SWING system to aid
 464 portability. An example displays during a SHRINK communicator rebuild opera-
 465 tion is shown in Figs. 6–8, where a process (rank 1) exits and the communicator is
 466 reduced in size and extent.

467 **7. Conclusions**

468 FT-MPI is an attempt to provide application programmers with different methods
 469 of dealing with failures within MPI application than just check-point and restart. It

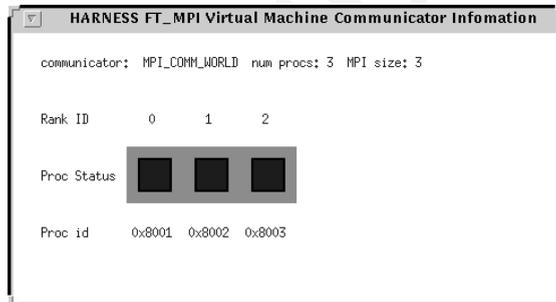


Fig. 6. Cominfo display for a healthy three process MPI application. The colours of the inner boxes indicate the state of the processes and the outer box indicates the communicator state.

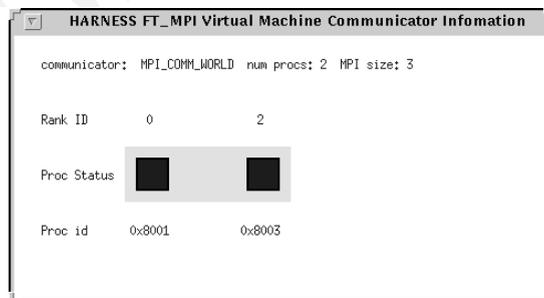


Fig. 7. COMINFO display for an application with an exited process. Note that the number of nodes and size of communicator do not match.

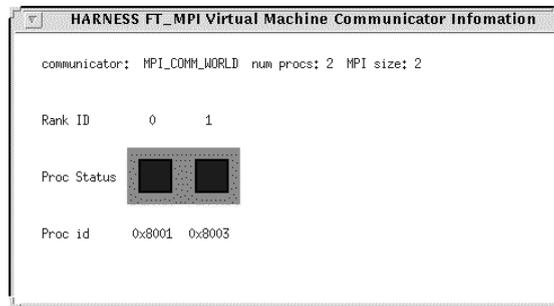


Fig. 8. Cominfo display for the above application after a communicator rebuild using the SHRINK option. Note the communicator status box has changed back to a blue (dark) colour.

470 is hoped that by experimenting with FT-MPI, new applications methodologies and
 471 algorithms will be developed to allow for both high performance and the surviv-
 472 ability required by the next generation of terra-flop and beyond machines.

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

477 Future work in the FT-MPI library system will concentrate on developing both
 478 further implementations that support more restrictive environments such as dedi-
 479 cated MPPs runtimes and embedded clusters. Application development will be
 480 supported by the creation of a number of drop-in library templates to simplify the
 481 construction of fault tolerant applications.

482 References

- 483 [1] Beck, Dongarra, Fagg, Geist, Gray, Kohl, Migliardi, K. Moore, T. Moore, P. Papadopoulos, S.
 484 Scott, V. Sunderam, HARNESS: a next generation distributed virtual machine, *Journal of Future*
 485 *Generation Computer Systems* 15 (1999).
- 486 [2] G. Stellner, CoCheck: checkpointing and process migration for MPI, in: *Proceedings of the*
 487 *International Parallel Processing Symposium*, Honolulu, April 1996, pp. 526–531.
- 488 [3] A. Agbaria, R. Friedman, StarFish: fault-tolerant dynamic MPI programs on clusters of
 489 workstations, in: *The 8th IEEE International Symposium on High Performance Distributed*
 490 *Computing*, 1999.
- 491 [4] G.E. Fagg, K. Moore, J.J. Dongarra, Scalable networked information processing environment
 492 (SNIPE), *Journal of Future Generation Computer Systems* 15 (1999) 571–582.
- 493 [5] M. Migliardi, V. Sunderam, PVM emulation in the HARNESS MetaComputing system: a plug-in
 494 based approach, in: *Lecture Notes in Computer Science*, vol. 1697, September 1999, pp. 117–124.
- 495 [6] P.H. Worley, I.T. Foster, B. Toonen, Algorithm comparison and benchmarking using a parallel
 496 spectral transform shallow water model, in: G.-R. Hoffmann, N. Kreitz (Eds.), *Proceedings of the*
 497 *Sixth Workshop on Parallel Processing in Meteorology*, World Scientific, Singapore, 1995, pp. 277–
 498 289.
- 499 [7] T. Kielmann, H.E. Bal, S. Gorlatch, Bandwidth-efficient collective communication for clustered wide
 500 area systems, in: *IPDPS 2000*, Cancun, Mexico, 1–5 May, 2000.

- 501 [8] L.P. Huse, Collective communication on dedicated clusters of workstations, in: Proceedings of the 6th
502 European PVM/MPI Users' Group Meeting, Lecture Notes in Computer Science, vol. 1697, Springer,
503 Barcelona, September 1999, pp. 469–476.
- 504 [9] D. Culler, R. Karp, D. Patterson, A. Sahay, K.E. Schauser, E. Santos, R. Subramonian, T. von
505 Eicken, LogP: towards a realistic model of parallel computation, in: Proceedings of the Symposium
506 on Principles and Practice of Parallel Programming (PpoPP), San Diego, CA, May 1993, pp. 1–12.
- 507 [10] R. Rabenseifner, A new optimized MPI reduce algorithm. [http://www.hlrs.de/structure/support/
508 parallel_computing/model/myreduce.html](http://www.hlrs.de/structure/support/parallel_computing/model/myreduce.html), 1997.
- 509 [11] M. Snir, S. Otto, S. Huss-Lederman, D. Walker, J. Dongarra, MPI – The Complete Reference. vol. 1,
510 The MPI Core, second ed., 1998.
- 511 [12] M. Frigo, FFTW: an adaptive software architecture for the FFT, in: Proceedings of the ICASSP
512 Conference, vol. 3, 1998, p. 1381.
- 513 [13] R.C. Whaley, J. Dongarra, Automatically Tuned Linear Algebra Software. SC98: High Performance
514 Networking and Computing, <http://www.cs.utk.edu/~rwhaley/ATL/INDEX.HTM>, 1998.
- 515 [14] L. Prylli, B. Tourancheau, BIP: a new protocol designed for high performance networking on
516 myrinet, in: The PC-NOW Workshop, IPPS/SPDP 1998, Orlando, USA, 1998.
- 517 [15] S. Louca, N. Neophytou, A. Lachanas, P. Evripidou, MPI-FT: a portable fault tolerance scheme for
518 MPI, in: Proceedings of PDPTA '98 International Conference, Las Vegas, NV, 1998.
- 519 [16] J.L. Traff, R. Hempel, H. Ritzdort, F. Zimmermann, Flattening on the Fly: efficient handling of MPI
520 derived datatypes, in: Proceedings of the 6th European PVM/MPI Users' Group Meeting, Lecture
521 Notes in Computer Science, vol. 1697, Springer, Barcelona, September 1999, pp. 109–116.
- 522 [17] W.D. Gropp, E. Lusk, D. Swider, Improving the performance of MPI derived datatypes, in:
523 Proceedings of the Third MPI Developer's and User's Conference (MPIDC '99), 1999, pp. 25–30.
- 524 [18] G.E Fagg, K.S. London, J.J. Dongarra, MPI_Connect, managing heterogeneous MPI application
525 interoperation and process control, in: EuroPVM-MPI '98, Lecture Notes in Computer Science, vol.
526 1497, Springer, Berlin, 1998, pp. 93–96.
- 527 [19] E. Gabriel, M. Resch, T. Beisel, R. Keller, Distributed computing in a heterogeneous computing
528 environment, in: EuroPVM-MPI '98, Lecture Notes in Computer Science, vol. 1497, Springer, Berlin,
529 1998, pp.180–187.
- 530 [20] S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling,
531 G. Henry, A. Petitet, K. Stanley, D. Walker, R. Whaley, Scalapack: a linear algebra library for
532 message-passing computers, in: Proceedings of 1997 SIAM Conference on Parallel Processing, May
533 1997.
- 534 [21] J. Choi, J. Dongarra, S. Ostrouchov, A. Petitet, D. Walker, R. Whaley, A proposal for a set of
535 parallel basic linear algebra subprograms, in: LAPACK Working Note #100, CS-95-292, May 1995.
- 536 [22] W. Gropp, E. Lusk, R. Thakur, Using MPI-2: Advanced Features of the Message Passing Interface,
537 first ed., MIT Press, Cambridge, MA, 2000.