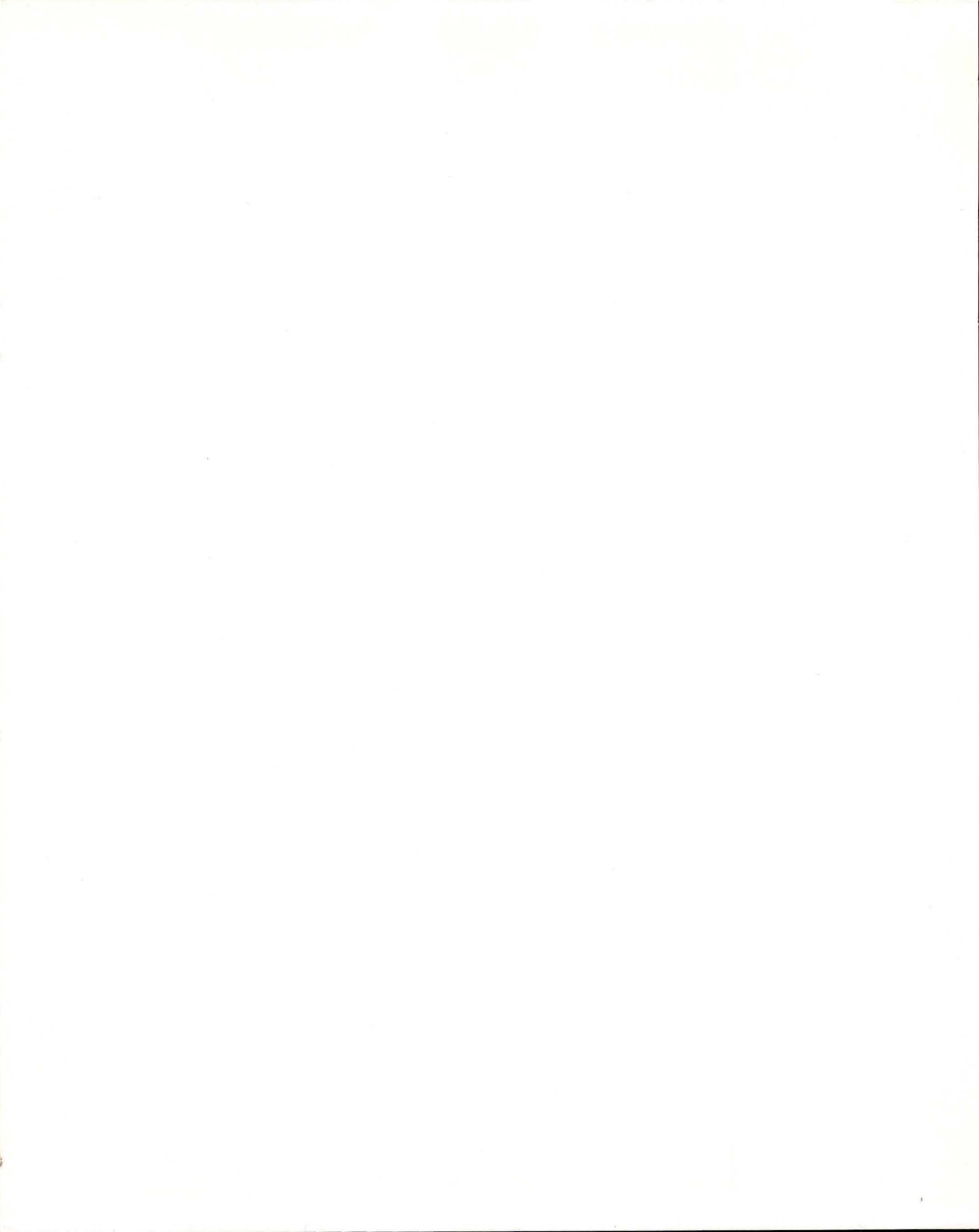


HP 9000

HP MPI User's Guide



HP MPI User's Guide

Second Edition



B6011-96002

Customer Order Number: B6011-90001

October 1997

Printed in: USA

© Copyright 1997 Hewlett-Packard company. All rights reserved.

Revision History

Edition: Second

Document Number: B6011-90001

Remarks: Released with HP MPI V1.3, October, 1997.

Edition: First

Document Number: B6011-90001

Remarks: Released with HP MPI V1.1, January, 1997.

Notice

Reproduction, adaptation, or translation without prior written permission is prohibited, except as allowed under the copyright laws.

The information contained in this document is subject to change without notice.

Hewlett-Packard makes no warranty of any kind with regard to this material, including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose. Hewlett-Packard shall not be liable for errors contained herein or for incidental or consequential damages in connection with the furnishing, performance or use of this material.

Parts of this book came from Cornell Theory Center's web document. That document is copyrighted by the Cornell Theory Center.

Parts of this book came from *MPI: A Message Passing Interface*. That book is copyrighted by the University of Tennessee. These sections were copied by permission of the University of Tennessee.

Parts of this book came from *MPI Primer / Developing with LAM*. That document is copyrighted by the Ohio Supercomputer Center. These sections were copied by permission of the Ohio Supercomputer Center.

Contents

Preface	xi
HP MPI features	xi
System platforms	xii
Notational conventions	xiii
Associated Documents	xiv
1 Introduction	1
Parallel computational models	2
Message passing	3
MPI concepts	4
Point-to-point communication	6
Communicators	6
Sending and receiving messages	7
Blocking communication	8
Nonblocking communication	9
Collective operations	10
Communication	11
Computation	13
Synchronization	14
Noncontiguous data	15
Multilevel parallelism	17
Advanced topics	17
2 Getting started	19
Configuring your environment	20
Building and running your first application	21
Building and running on a single host	22
hello_world output	22
Building and running on multiple hosts	23
hello_world output	23
3 Understanding HP MPI	25
Directory structure	26
Compatibility issues	28
Compiling applications	29
64-bit support	30
Language interoperability	31

Running applications	32
Types of applications.	32
Running SPMD applications.	32
Running MPMD applications	33
Multiprotocol messaging.	34
Run-time environment variables	36
MPI_FLAGS.	37
MPI_GLOBSIZE.	38
MPI_TOPOLOGY	39
MPI_SHMEMCNTL	40
MPI_TMPDIR	40
MPI_XMPI	41
MPI_WORKDIR.	42
MPI_CHECKPOINT	43
MPI_INSTR	43
Run-time utility commands	45
mpirun	45
Creating an appfile.	47
mpijob	48
mpiclean	49
xmpi	50
mpitrget	51
mpitrstat.	51
4 Profiling	53
Using counter instrumentation	54
Creating an instrumentation profile	54
Using XMPI	58
Working with postmortem mode	58
Creating a trace file	59
Viewing a trace file	60
Viewing process information from a trace.	65
Viewing kivi information from a trace file.	69
Working with interactive mode	70
Running an appfile	70
Changing default settings and viewing options.	75
Using CXperf	79
Using the profiling interface	80

5	Tuning	81
	General tuning.....	82
	Message latency and bandwidth	82
	Multiple network interfaces.....	83
	Processor subscription	85
	MPI routine selection	86
	SPP-UX platform tuning.....	87
	Multilevel parallelism	87
	Process placement	87
6	Debugging and troubleshooting	91
	Debugging HP MPI applications	92
	Setting options in MPI_FLAGS.....	92
	Using CXdb.....	93
	Troubleshooting HP MPI applications	94
	Building.....	94
	Starting.....	94
	Running.....	95
	Propagation of environment variables.....	95
	Shared memory	95
	Interoperability	96
	Message buffering	96
	External input and output	97
	Fortran 90 programming features	98
	UNIX open file descriptors	98
	Completing	98
	Frequently asked questions	100
	Appendix A: MPI library routines and extensions	103
	C version of MPI routines	104
	Fortran version of MPI routines	117
	C version of HP MPI library extensions.....	132
	Fortran version of HP MPI library extensions.....	133
	MPI 1.2 extensions	134
	MPI 2.0 extensions	135

Appendix B: Example applications	137
send_receive.f.....	139
send_receive output	140
ping_pong.c.....	141
ping_pong output	143
compute_pi.f.....	144
compute_pi output.....	145
master_worker.f90.....	146
master_worker output	147
cart.C.....	148
cart output.....	150
communicator.c.....	152
communicator output	152
multi_par.f.....	153
Appendix C: XMPI resource file	163
Glossary	165
Index	173

Figures

Figure 1-1	MPI broadcast operation	11
Figure 1-2	MPI scatter operation	12
Figure 3-1	Multiprotocol messaging with an X-Class server	34
Figure 3-2	Multiprotocol messaging with a K-Class server	35
Figure 5-1	Multiple network interfaces	85
Figure 5-2	Default process placement	88
Figure 5-3	Optimal process placement	89
Figure B-1	Array partitioning	154

Tables

Table 1-1	Six commonly used MPI routines	5
Table 1-2	MPI blocking and nonblocking calls	9
Table 3-1	Organization of the /opt/mpi directory	26
Table 3-2	Man page categories	27
Table 3-3	Compilation utilities	29
Table 3-4	Compilation environment variables	29
Table 5-1	Subscription types	85
Table 6-1	Run invocations that support stdin	97
Table A-1	C version of MPI routines	104
Table A-2	Fortran version of MPI routines	117
Table A-3	C version of HP MPI library extensions	132
Table A-4	Fortran version of HP MPI library extensions	133
Table A-5	MPI 1.2 extensions	134
Table A-6	MPI 2.0 extensions	135
Table B-1	Example applications shipped with HP MPI	137

Preface

This guide describes the HP MPI implementation of the Message Passing Interface (MPI) standard. The guide helps you use HP MPI to develop parallel applications.

You should already have experience developing UNIX applications. You should also understand the basic concepts behind parallel processing and be familiar with MPI.

This guide is intended to supplement *MPI: The Complete Reference* (B6011-90003) with information specific to the HP implementation of MPI.

An .html version of this guide is provided with HP MPI. See “Directory structure” on page 26 for more information.

NOTE

Some of the sections in this book contain command-line examples used to demonstrate HP MPI concepts. These examples use the /bin/csh syntax for illustration purposes only.

HP MPI features

HP MPI provides a wide range of features that offer you flexibility in developing parallel applications. These features include:

- Compliance with the MPI version 1.2 standard.
- Support for a small subset of the MPI version 2.0 standard.
- Single program multiple data (SPMD) and multiple program multiple data (MPMD) styles of programming—Allow you to create an application that consists of a single program that is executed by each process or two or more programs where each process can execute a different program. In both cases, processes normally act on different data.
- Profiling—Supports process tracing and monitoring using XMPI and counter instrumentation for collecting cumulative application statistics.
- Multiprotocol support—Supports different communication protocols depending upon where the processes are located and what type of platform is used. The supported protocols include shared memory within a host and TCP/IP between hosts.

Preface

System platforms

- **Data mover**—Accelerates messaging on scalable servers running under SPP-UX.
- **Compliance with the UNIX 95 standard.**
- **Derived data types**—Supports optimized collection and communication of derived data types.

System platforms

HP MPI runs under the HP-UX and SPP-UX operating systems.

The HP-UX operating system is used on:

- **Workstations:** s700 series B-, C-, and J-Class
- **Midrange servers:** s800 series D- and K-Class
- **High-end servers:** V-Class.

The SPP-UX operating system is used on:

- **SPP1600 servers** (single- and multi-hypernode)
- **S-Class servers**
- **X-Class servers**

Notational conventions

This section describes notational conventions used in this book.

bold monospace	In command examples, bold monospace identifies input that must be typed exactly as shown.
monospace	In paragraph text, monospace identifies command names, system calls, and data structures and types. In command examples, monospace identifies command output, including error messages.
<i>italic</i>	In paragraph text, <i>italic</i> identifies titles of documents. In command syntax diagrams, <i>italic</i> identifies variables that you must provide. The following command example uses brackets to indicate that the variable <i>output_file</i> is optional: command <i>input_file</i> [<i>output_file</i>]
Brackets ([])	In command examples, square brackets designate optional entries.

NOTE

A note highlights important supplemental information.

Associated Documents

Associated documents include:

- *MPI: The Complete Reference*, published by MIT Press (B6011-90003)
- *Exemplar Programming Guide for HP-UX Systems* (B6056-90002)
- *MPI-2: Extensions to the Message-Passing Interface*, available from the University of Tennessee
- *Using CXperf* (B6059-90005)
- *Using CXdb* (B6059-90002)
- *Managing Systems and Workgroups* (B2355-90157)

The table below shows World Wide Web sites that contain additional MPI information.

Access ...	To learn more about ...
http://www.mpi-forum.org	MPI standardization forum.
http://www.mcs.anl.gov/Projects/mpi/index.html	Argonne National Laboratory's MPICH implementation of MPI.
http://www.tc.cornell.edu/Edu/Tutor/MPI/	Cornell Theory Center's MPI tutorial and lab exercises.
http://www.osc.edu/Lam.html	Ohio Supercomputer Center's LAM implementation of MPI.
http://www.erc.msstate.edu/mpi/	Mississippi State University's MPI web page.
http://www.mpi-forum.org/	Official site of the MPI forum.
http://www.hp.com/go/mpi	Hewlett-Packard's HP MPI web page.

1

Introduction

This chapter provides introductory information about MPI. The topics covered include:

- Parallel computational models
- Message passing
- MPI concepts

Parallel computational models

Computational models represent a way to look at the types of operations that are available to parallel applications.

These models are independent of the underlying hardware. They can be implemented on any machine designed to run parallel applications. Model performance, however, depends on how optimally a model is implemented on a particular hardware architecture.

The models are categorized by how memory is used (shared versus distributed) and how communication occurs (software versus hardware).

The models include:

- **Shared memory**—Each process can access a shared address space.
- **Message passing**—An application runs as a collection of autonomous processes, each with its own local memory.
- **Remote memory operations**—A local process accesses the memory of a remote process without aid from the remote process. The local process accesses this memory explicitly—not the way it accesses its local memory.
- **Threads**—In a multithreaded process, the values of application variables are shared by all the threads.
- **Hybrid models**—Two or more of the models are used together.

Message passing

In message passing, a parallel application consists of a number of processes that run concurrently. Each process has its own local memory and communicates with other processes by sending and receiving messages. When data is passed in a message, the sending and receiving processes must work to transfer the data from the local memory of one to the local memory of the other.

Message passing is one of the most popular computation models for designing parallel applications. The advantages of using message passing include:

- **Portability**—Message passing is implemented on most parallel platforms.
- **Universality**—Model makes minimal assumptions about underlying parallel hardware. Message-passing libraries exist on computers linked by networks and on shared and distributed multiprocessors.
- **Simplicity**—Model supports explicit control of memory references for easier debugging.

On the other hand, creating message-passing applications may require more effort than letting a parallelizing compiler produce parallel applications.

Applications based on the message-passing model are nondeterministic by default (that is, the arrival order of messages sent from two processes to a third process is not defined). However, when one process sends two or more messages to another process, the transfer is deterministic as the messages are always received in the order sent. You must ensure that process communication and computation are deterministic when required within your application.

MPI concepts

MPI is a standard specification for interfaces to a library of message-passing routines. The goals of MPI are efficient communication, portability, and rich functionality.

Although several message-passing libraries exist on different systems (for example, Intel Paragon's NX, Connection Machine's CMMD, or PVM), MPI has emerged as a favorite for the following reasons:

- Support for full asynchronous communication—Process communication can overlap process computation.
- Group membership—Processes may be grouped based on context.
- Synchronization variables that protect process messaging—When sending and receiving messages, synchronization is enforced by source and destination information, message labeling, and context information.
- Portability—All implementations are based on a published standard that specifies the semantics for usage.

The MPI library contains 128 routines. These routines provide a set of functions that support point-to-point communications, collective operations, process groups and communication contexts, process topologies, and data type manipulation. In addition, HP MPI supports a small number of library extensions, MPI 1.2 extensions, and MPI 2.0 extensions. See "MPI library routines and extensions" on page 103 for more information.

Although the MPI library contains a large number of routines to choose from, you can design a number of applications by only using the six listed in Table 1-1.

Table 1-1 **Six commonly used MPI routines**

MPI routine	Description
MPI_Init	Initializes the MPI environment
MPI_Finalize	Terminates the MPI environment
MPI_Comm_rank	Determines the rank of the calling process within a group
MPI_Comm_size	Determines the size of the group
MPI_Send	Sends messages
MPI_Recv	Receives messages

NOTE

You must include `MPI_Finalize` in your application to prevent deadlock and incorrect message transfers. HP MPI issues a warning when a process exits without calling `MPI_Finalize`.

As your application grows in complexity, you can introduce other routines from the library. For example, `MPI_Bcast` is an often-used routine for sending data from one process to other processes in a single operation. This is much more efficient in terms of performance than using `MPI_Send` and `MPI_Recv` to transfer data from the sending process to each receiving process one by one.

Point-to-point communication

Point-to-point communication involves sending and receiving messages between two processes. This is the simplest form of data transfer in a message-passing model.

The performance of point-to-point communication is measured in terms of total transfer time. The total transfer time is defined as

$$total_transfer_time = latency + (message_size/bandwidth)$$

where

<i>latency</i>	Specifies the time between the initiation of the data transfer in the sending process and the arrival of the first byte in the receiving process.
<i>message_size</i>	Specifies the size of the message in Mbytes.
<i>bandwidth</i>	Denotes the reciprocal of the time needed to transfer a byte. Bandwidth is normally expressed in Mbytes per second.

Obviously, low latencies and high bandwidths lead to better performance.

Communicators

A communicator is an object that represents a group of processes and their communication medium or context. These processes exchange messages with each other to transfer data. In this context, communicators encapsulate their processes such that communication is restricted to processes only within the group.

Communicators are also used to synchronize communications. In this context, they prevent process messages from being mistakenly intercepted by other program layers (for example, a user-defined library that the application calls).

The default communicators provided by MPI are `MPI_COMM_WORLD` and `MPI_COMM_SELF`. `MPI_COMM_WORLD` consists of all processes that are running when an application begins execution. Each process is the single member of its own `MPI_COMM_SELF`.

Communicators that allow processes within a single group to exchange data are termed intracommunicators. Communicators that allow processes between two different groups to exchange data are called intercommunicators.

Many MPI applications depend upon knowing the number of processes and the process rank within a given communicator.

To determine the number of processes in a communicator named `comm`, use

```
MPI_Comm_size (MPI_Comm comm, int *size);
```

To determine the rank of each process in `comm`, use

```
MPI_Comm_rank (MPI_Comm comm, int *rank);
```

where *rank* is an integer between zero and (*size* - 1).

Refer to example “communicator.c” on page 152 for more information about using communicators.

Sending and receiving messages

There are two methods for sending and receiving data: blocking and nonblocking.

Blocking communication means the sending process does not return until the send buffer is available for reuse.

Nonblocking communication means the sending process returns immediately, but the send buffer is not safe for reuse. In this case:

1. The sending routine begins the message transfer and returns immediately.
2. The application does some computation.
3. The application calls a completion routine (for example, `MPI_Test` or `MPI_Wait`) to test or wait for completion of the send operation.

Blocking communication

Blocking communication consists of four send modes and one receive mode.

The four send modes include:

- **Standard mode** (`MPI_Send`)—The sending process returns when the system can buffer the message or when the message is received.
- **Buffered mode** (`MPI_Bsend`)—The sending process returns when the message is buffered in the application-supplied space.

Avoid using the `MPI_Bsend` mode. This mode just forces an additional copy operation.

- **Synchronous mode** (`MPI_Ssend`)—The sending process returns only if a matching receive is posted and the receiving process has started to receive the message.
- **Ready mode** (`MPI_Rsend`)—The sending process assumes that a matching receive is posted. The sending process returns after the message is sent.

The four send modes are all invoked in a similar manner and all pass the same arguments. The only difference is in the routine name used to send the message (that is, `MPI_Send` versus `MPI_Ssend`).

To code a standard blocking send, use

```
MPI_Send (void *buf, int count, MPI_Datatype dtype, int dest, int tag, MPI_Comm comm);
```

where

<i>buf</i>	Specifies the starting address of the buffer.
<i>count</i>	Indicates the number of buffer elements.
<i>dtype</i>	Denotes the data type of the buffer elements.
<i>dest</i>	Specifies the rank of the destination process in the group associated with the communicator <i>comm</i> .
<i>tag</i>	Denotes the message label.
<i>comm</i>	Designates the communication context that identifies a group of processes.

To code a blocking receive, use

```
MPI_Recv (void *buf, int count, MPI_datatype dtype, int source, int tag, MPI_Comm comm,
MPI_Status *status);
```

where

<i>buf</i>	Specifies the starting address of the buffer.
<i>count</i>	Indicates the number of buffer elements.
<i>dtype</i>	Denotes the data type of the buffer elements.
<i>dest</i>	Specifies the rank of the destination process in the group associated with the communicator <i>comm</i> .
<i>tag</i>	Denotes the message label.
<i>comm</i>	Designates the communication context that identifies a group of processes.
<i>source</i>	Specifies the rank of the source process in the group associated with the communicator <i>comm</i> .
<i>status</i>	Returns information about the received message. Status information is useful when wildcards are used or the received message is smaller than expected. Status may also contain error codes.

Refer to examples “send_receive.f” on page 139, “ping_pong.c” on page 141, and “master_worker.f90” on page 146 for more information about using blocking communications.

Nonblocking communication

MPI provides nonblocking versions of the blocking send and receive calls. Table 1-2 lists these calls.

Table 1-2

MPI blocking and nonblocking calls

Blocking mode	Nonblocking mode
MPI_Send	MPI_Isend
MPI_Bsend	MPI_Ibsend
MPI_Ssend	MPI_Issend
MPI_Rsend	MPI_Irsend
MPI_Recv	MPI_Irecv

The nonblocking calls have the same arguments as their blocking counterparts plus an additional argument for a request.

To code a standard nonblocking send, use

```
MPI_Isend (void *buf, int count, MPI_datatype dtype, int dest, int tag, MPI_Comm comm,  
MPI_Request *req);
```

where *req* specifies the request used by a completion routine when called by the application to complete the send operation.

Collective operations

Applications may require coordinated operations among multiple processes. For example, all processes need to cooperate to sum sets of numbers distributed among them.

MPI provides a set of collective operations to coordinate operations among processes. These operations are implemented such that all processes call the same operation with the same arguments. Thus, when sending and receiving messages, one collective operation can replace multiple sends and receives, resulting in lower overhead and higher performance.

Collective operations consist of routines for communication, computation, and synchronization. These routines all specify a communicator argument that defines the group of participating processes and the context of the operation.

NOTE

Collective operations are valid only for intracommunicators. Intercommunicators are not allowed as arguments.

Communication

Collective communication involves the exchange of data among all processes in a group. The communication can be one-to-many, many-to-one, or many-to-many.

The single originating or receiving process in the one-to-many and many-to-one routines is called the root.

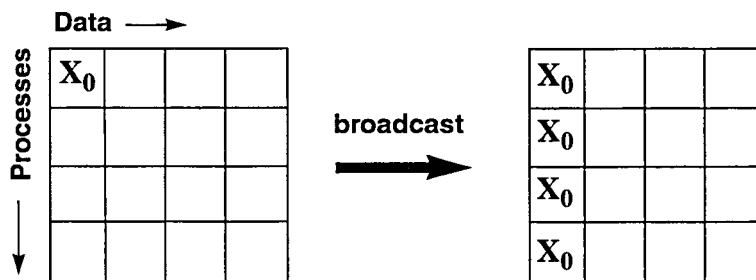
Examples of such communication routines are:

Broadcast

(MPI_Bcast) A one-to-many operation where the root sends its data to all other processes in the communicator, including itself. Figure 1-1 shows the broadcast operation for a four-process application. Each row of boxes represents data locations in one process.

Figure 1-1

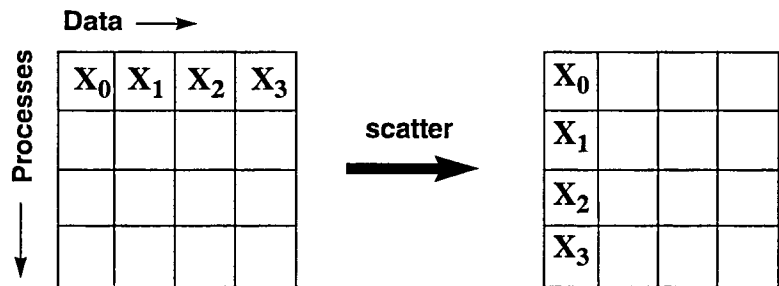
MPI broadcast operation



Scatter

(MPI_Scatter) A one-to-many operation where the root's data is split among all processes in the communicator. Figure 1-2 shows the scatter operation for a four-process application. As before, each row of boxes represents data locations in one process.

Figure 1-2 MPI scatter operation



To code a broadcast, use

`MPI_Bcast (void *buf, int count, MPI_Datatype dtype, int root, MPI_Comm comm);`

where

- buf* Specifies the starting address of the buffer.
- count* Indicates the number of buffer entries.
- dtype* Denotes the data type of the buffer entries.
- root* Specifies the rank of the root.
- comm* Designates the communication context that identifies a group of processes.

To code a scatter, use

```
MPI_Scatter (void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf,  
int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm);
```

where

<i>sendbuf</i>	Specifies the starting address of the send buffer.
<i>sendcount</i>	Specifies the number of elements sent to each process.
<i>sendtype</i>	Denotes the data type of the send buffer.
<i>recvbuf</i>	Specifies the address of the receive buffer.
<i>recvcount</i>	Indicates the number of elements in the receive buffer.
<i>recvtype</i>	Indicates the data type of the receive buffer elements.
<i>root</i>	Denotes the rank of the sending process.
<i>comm</i>	Designates the communication context that identifies a group of processes.

Computation

Computation uses `MPI_Reduce` to apply reduction operations across all processes in a communicator. Reduction operations are binary and are only valid on numeric data. Also, reductions are always associative but may or may not be commutative.

You can select a reduction operation from a predefined list or define your own operation. Examples of predefined operations include `MPI_SUM` and `MPI_PROD`, which apply a summation and a multiplication across all processes respectively.

Introduction
MPI concepts

To implement a reduction, use

```
MPI_Reduce (void *sendbuf, void *recvbuf, int count, MPI_Datatype dtype, MPI_Op op,  
int root, MPI_Comm comm);
```

where

<i>sendbuf</i>	Specifies the address of the send buffer.
<i>recvbuf</i>	Denotes the address of the receive buffer.
<i>count</i>	Indicates the number of elements in the send buffer.
<i>dtype</i>	Specifies the data type of the send and receive buffers.
<i>op</i>	Specifies the reduction operation.
<i>root</i>	Indicates the rank of the root process.
<i>comm</i>	Designates the communication context that identifies a group of processes.

Synchronization

Collective routines return as soon as their participation in a communication is complete. However, the return of the calling process does not guarantee that the receiving processes have completed or even started the operation.

To synchronize the execution of processes, call `MPI_Barrier`. `MPI_Barrier` blocks the calling process until all processes in the communicator have called it. This is a useful approach for separating two stages of a computation so messages from each stage are not overlapped.

To implement a barrier, use

```
MPI_Barrier (MPI_Comm comm);
```

where *comm* identifies a group of processes and a communication context.

Refer to examples “compute_pi.f” on page 144 and “cart.C” on page 148 for more information about using collective operations.

Noncontiguous data

Predefined data types are used for transferring data between two processes using point-to-point communication. This transfer is based on the assumption that the data transferred is stored in contiguous memory (for example, sending an array in a C or Fortran application).

What happens, however, when you want to transfer data that is not stored contiguously? In this case, you can create a derived data type or use `MPI_Pack` and `MPI_Unpack`.

A derived data type specifies a sequence of basic data types and integer displacements describing the data layout in memory. Derived data types are more efficient than using `MPI_Pack` and `MPI_Unpack`, but they cannot handle the case where the data layout varies and is unknown by the receiver beforehand (for example, messages that embed their own layout description).

You create derived data types through the use of type-constructor functions. Once a derived data type is created, you can use it repeatedly in all communicating calls. This allows MPI to pack and unpack the data as necessary and further optimize the data transfer.

The types of constructor functions include:

- **Contiguous**—Allows replication of a data type into contiguous locations.
- **Vector**—Allows replication of a data type into locations that consist of equally spaced blocks.
- **Indexed**—Allows replication of a data type into a sequence of blocks where each block can contain a different number of copies and have a different displacement.
- **Structure**—Allows replication of a data type into a sequence of blocks such that each block consists of replications of different data types, copies, and displacements.

NOTE HP MPI optimizes collection and communication of derived data types.

Introduction
MPI concepts

To create a vector data type, use

```
MPI_Type_Vector (int count, int blocklength, int stride, MPI_Datatype oldtype,  
MPI_Datatype *newtype);
```

where

<i>count</i>	Indicates the number of blocks.
<i>blocklength</i>	Specifies the number of elements in each block.
<i>stride</i>	Denotes the number of elements between the start of two consecutive blocks.
<i>oldtype</i>	Specifies the old data type.
<i>newtype</i>	Specifies the new data type.

You must now commit the derived data type by calling `MPI_Type_commit`.

`MPI_Pack` allows you to store noncontiguous data in contiguous memory locations. `MPI_Unpack` copies data from a contiguous buffer into noncontiguous memory locations. Used together, these routines allow you to transfer heterogeneous data in a single message.

To code a pack, use

```
MPI_Pack (void *inbuf, int incount, MPI_Datatype dtype, void *outbuf, int outsize,  
int *position, MPI_Comm, comm);
```

where

<i>inbuf</i>	Specifies the start of the input buffer.
<i>incount</i>	Indicates the number of input data items.
<i>dtype</i>	Denotes the data type of each input data item.
<i>outbuf</i>	Specifies the start of the output buffer.
<i>outsize</i>	Indicates the output buffer size in bytes.
<i>position</i>	Specifies the current position in the buffer in bytes.
<i>comm</i>	Designates the communication context that identifies a group of processes.

Multilevel parallelism

By default, a process in MPI applications can only do one task at a time. Such processes are known as single-threaded processes. This means that each process has an address space together with a single program counter, a set of registers, and a stack.

A multithreaded process has one address space, but each process thread contains its own counter, registers, and stack.

Multilevel parallelism refers to MPI processes that have multiple threads. Processes become multithreaded through calls to multithreaded libraries, parallel directives and pragmas, and auto-compiler parallelism.

Multilevel parallelism is beneficial for problems you can decompose into logical parts for parallel execution (for example, a looping construct that spawns multiple threads to do a computation and joins after the computation is complete).

See “multi_par.f” on page 153 for an example of multilevel parallelism.

NOTE

HP MPI is not a thread-safe implementation of MPI. This means that more than one thread cannot call MPI functions at the same time.

Advanced topics

This chapter only provides information about basic MPI concepts. Advanced MPI topics include:

- Error handling
- Process topologies
- User-defined data types
- Process grouping
- Attribute caching

To learn more about these advanced topics, see *MPI: The Complete Reference*, the companion to this guide.

Introduction
MPI concepts

This chapter describes how to get started with HP MPI. The topics covered are:

- Configuring your environment
- Building and running your first application

Configuring your environment

Use the following steps to configure your environment before running your first HP MPI application:

Step 1. Verify that HP MPI is installed on your system in the `/opt/mpi` directory.

Step 2. Add `/opt/mpi/bin` to the `PATH` variable by entering:

```
% setenv PATH /opt/mpi/bin:$PATH
```

Step 3. Add `/opt/mpi/share/man` to the `MANPATH` variable by entering:

```
% setenv MANPATH /opt/mpi/share/man:$MANPATH
```

If you ever move the HP MPI installation directory from its default location (`/opt/mpi`), set the `MPI_ROOT` environment variable to point to the new location. For example:

```
% setenv MPI_ROOT /usr/local/mpi
```

Building and running your first application

To quickly gain experience with HP MPI, start with a C version of the familiar hello-world program. This program is called `hello_world.c` and prints out the text string “Hello world! I’m r of s on *host*” where r is a process’s rank, s is the size of the communicator, and *host* is the host on which the program is run.

The source code for `hello_world.c` is stored in `/opt/mpi/help` and is shown below.

```
/* hello_world.c */
#include <stdio.h>
#include <mpi.h>

main(argc, argv)

int          argc;
char        *argv[];

{
    int          rank, size, len;
    char        name[MPI_MAX_PROCESSOR_NAME];

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Get_processor_name(name, &len);
    printf ("Hello world! I'm %d of %d on %s\n", rank, size, name);

    MPI_Finalize();
    exit(0);
}
```

Building and running on a single host

To build and run `hello_world.c` on a local host named `jawbone`:

Step 1. Change to a writable directory.

Step 2. Enter

```
% mpicc -o hello_world /opt/mpi/help/hello_world.c
```

This step builds the `hello_world` executable.

Step 3. Enter

```
% mpirun -np 4 hello_world
```

This step runs the `hello_world` executable using four processes.

hello_world output

The output from running the `hello_world` executable is printed in nondeterministic order and is shown below.

```
Hello world! I'm 1 of 4 on jawbone  
Hello world! I'm 3 of 4 on jawbone  
Hello world! I'm 0 of 4 on jawbone  
Hello world! I'm 2 of 4 on jawbone
```

Building and running on multiple hosts

To build and run `hello_world.c` on a local host named `jawbone` and a remote host named `wizard` (assuming that both machines run under either HP-UX or SPP-UX or `hello_world.c` is built on HP-UX so the same binary can run on both hosts):

Step 1. Edit the `.rhosts` file on `jawbone` and `wizard`. Add an entry for `wizard` in the `.rhosts` file on `jawbone` and an entry for `jawbone` in the `.rhosts` file on `wizard`.

Step 2. Change to a writable directory.

Step 3. Enter

```
% mpicc -o hello_world /opt/mpi/help/hello_world.c
```

This step builds the `hello_world` executable.

Step 4. Copy the `hello_world` executable from `jawbone` to your home directory on `wizard`.

Step 5. Create a text file called `appfile` and add the following two lines:

```
-np 2 hello_world  
-h wizard -np 2 hello_world
```

The `appfile` file contains a separate line for each host, which specifies the name of the executable and the number of processes to run on that host. The `-h` option identifies remote hosts.

Step 6. Enter

```
% mpirun -f appfile
```

This step runs `hello_world` on the hosts specified in the `appfile` file.

hello_world output

The output from running the `hello_world` executable is printed in nondeterministic order and is shown below.

```
Hello world! I'm 1 of 4 on wizard  
Hello world! I'm 3 of 4 on jawbone  
Hello world! I'm 0 of 4 on wizard  
Hello world! I'm 2 of 4 on jawbone
```

Getting started

Building and running your first application

NOTE

Ranks in `MPI_COMM_WORLD` are mapped from top to bottom.

3

Understanding HP MPI

This chapter provides information about the HP MPI implementation of MPI. The topics covered are:

- Directory structure
- Compatibility issues
- Compiling applications
- Running applications

Directory structure

All HP MPI files are stored in the /opt/mpi directory. The directory structure is organized as shown in Table 3-1.

Table 3-1

Organization of the /opt/mpi directory

Subdirectory	Contents
bin	Command files for the HP MPI utilities
doc/html	HTML version of the <i>HP MPI User's Guide</i> . View cover.html to browse the guide.
help	Source files for the example programs
include	Header files
lib/X11/app-defaults	Application default settings for the XMPI trace utility
lib/pa1.1/libfmpi.a	MPI library for 32-bit Fortran applications
lib/pa1.1/libmpi.a	MPI library for 32-bit C and C++ applications
lib/pa1.1/libpmpi.a	MPI 32-bit profiling interface library
lib/pa20_64/libfmpi.a	MPI library for 64-bit Fortran applications
lib/pa20_64/libmpi.a	MPI library for 64-bit C and C++ applications
lib/pa20_64/libpmpi.a	MPI 64-bit profiling interface library
newconfig/	Configuration files and release notes
share/man/man1.Z	Man pages for the HP MPI utilities
share/man/man3.Z	Man pages for HP MPI library and library functions

The man pages located in the /opt/mpi/share/man/man1.Z subdirectory are grouped into three categories: compilation, general, and run time. The compilation and run-time categories correspond to available types of HP MPI utilities. All three categories are described in Table 3-2.

Table 3-2

Man page categories

Man page category	Description
Compilation	Describes the available compilation utilities. Refer to “Compiling applications” on page 29 for more information.
General	Describes the general features of HP MPI. The man page is called MPI.1.
Run time	Describes the available run-time utilities. Refer to “Run-time utility commands” on page 45 for more information.

Compatibility issues

Several compatibility issues exist for HP MPI V1.3:

- HP MPI V1.3 supports 32-bit versions of the MPI library on all platforms running SPP-UX and HP-UX and 64-bit versions on platforms running HP-UX 11.0. However, in the same application, you cannot mix 32-bit and 64-bit executables.
- In order for HP MPI V1.3 to comply with the MPI standard, data structures defined in the `mpi.h` and `mpif.h` header files have changed. Therefore, you must recompile all HP MPI V1.1 files that include `mpi.h` and `mpif.h` before linking your application with HP MPI V1.3. For HP MPI V1.2 files, this recompilation is not necessary.

Compiling applications

The compiler used to build HP MPI applications depends upon which programming language you use. HP MPI provides separate compilation utilities and default compilers for the languages shown in Table 3-3.

Table 3-3 **Compilation utilities**

Language	Utility	Default compiler
C	mpicc	/opt/ansic/bin/cc
C++	mpiCC	/opt/aCC/bin/aCC
Fortran 77	mpif77	/opt/fortran/bin/f77
Fortran 90	mpif90	/opt/fortran90/bin/f90

NOTE If aCC is not available, mpicc uses CC as the default C++ compiler.

NOTE Even though the mpiCC and mpif90 compilation utilities are shipped with HP MPI, all C++ and Fortran 90 applications use C and Fortran 77 bindings respectively.

If you want to use a compiler other than the default one assigned to each utility, you can set the environment variables shown in Table 3-4.

Table 3-4 **Compilation environment variables**

Utility	Environment variable
mpicc	MPI_CC
mpiCC	MPI_CXX
mpif77	MPI_F77
mpif90	MPI_F90

Understanding HP MPI

Compiling applications

To set a compilation environment variable, enter:

```
% setenv compilation_environment_variable path
```

where *compilation_environment_variable* is the name of the variable you want to set and *path* specifies the path to the compiler you want to use.

64-bit support

HP-UX 11.0 is available as a 64-bit operating system for PA2.0 architectures and as a 32-bit operating system for older PA-RISC processor architectures. You must run 64-bit executables on the 64-bit system (though you can build 64-bit executables on the 32-bit system).

HP MPI supports a 64-bit version of the MPI library on platforms running HP-UX 11.0. Both 32- and 64-bit versions of the library are shipped with HP-UX 11.0 (only a 32-bit version is shipped with HP-UX 10.20).

The `mpicc` and `mpiCC` commands link the 64-bit version of the library if you compile with the `+DA2.0W` or `+DD64` options. The `mpif90` command links the 64-bit version of the library if you compile with the `+DA2.0W` option. Otherwise, the 32-bit version is used.

Statically-bound binaries built on HP-UX 10.20 platforms can run on HP-UX 11.0 systems. However, dynamically-bound binaries can only run on the HP-UX platform on which they were built.

Language interoperability

HP MPI complies with the language interoperability requirements of the MPI-2 standard. Language interoperability allows you to write mixed-language applications or applications that call library routines written in another language. For example, you can write applications in Fortran or C that call MPI library routines written in C or Fortran respectively.

MPI provides a special set of conversion functions for converting objects between languages. The types of objects that you can convert include MPI communicators, data types, groups, requests, reduction operations, and status. See “MPI 2.0 extensions” on page 135 for a list of these MPI conversion functions.

Running applications

Most HP MPI applications are run using the `mpirun` command. You should invoke the `mpirun` command with the `-j` option. This option displays the job ID of your job. The job ID is useful during troubleshooting if you want to check for a hung job using the `mpijob` command or want to terminate your job using the `mpiclean` command.

In some cases, you can use the `executable -np #` syntax to start your application. For example, to start an executable named `hello_world` with four processes, enter:

```
% hello_world -j -np 4
```

For multiprotocol applications that span multiple subcomplexes or multiple hosts, you must use `mpirun` together with an `appfile`. For applications that run on a single host and have a single executable, you can use `executable -np #` syntax, although `mpirun` is still recommended.

Types of applications

HP MPI supports two programming styles: SPMD applications and MPMD applications.

Running SPMD applications

A single program multiple data (SPMD) application consists of a single program that is executed by each process in the application. Each process normally acts upon different data. Even though this style simplifies the execution of an application, using SPMD can also make the executable larger and more complicated.

Each process calls `MPI_Comm_rank` to distinguish itself from all other processes in the application. It then determines what processing to do.

To run a SPMD application, use the `mpirun` command like this:

```
% mpirun -np # program
```

where `#` is the number of processors and `program` is the name of your application.

Suppose you want to build a C application called `poisson` and run it using five processes to do the computation. To do this, use the following command sequence:

```
% mpicc -o poisson poisson.c
% mpirun -np 5 poisson
```

Running MPMD applications

A multiple program multiple data (MPMD) application uses two or more separate programs to functionally decompose a problem.

This style can be used to simplify the application source and reduce the size of spawned processes. Each process can execute a different program.

To run an MPMD application, the `mpirun` command must reference an appfile that contains the number of processes to be created from each program and the list of programs to be run.

A simple invocation of an MPMD application looks like this:

```
% mpirun -f appfile
```

where *appfile* is the path name to a file that contains process counts and a list of programs.

Suppose you decompose the `poisson` application into two source files: `poisson_master` (uses a single master process) and `poisson_child` (uses four child processes).

The appfile for the example application contains the two lines shown below:

```
-np 1 poisson_master
-np 4 poisson_child
```

To build and run the example application, use the following command sequence:

```
% mpicc -o poisson_master poisson_master.c
% mpicc -o poisson_child poisson_child.c
% mpirun -f appfile
```

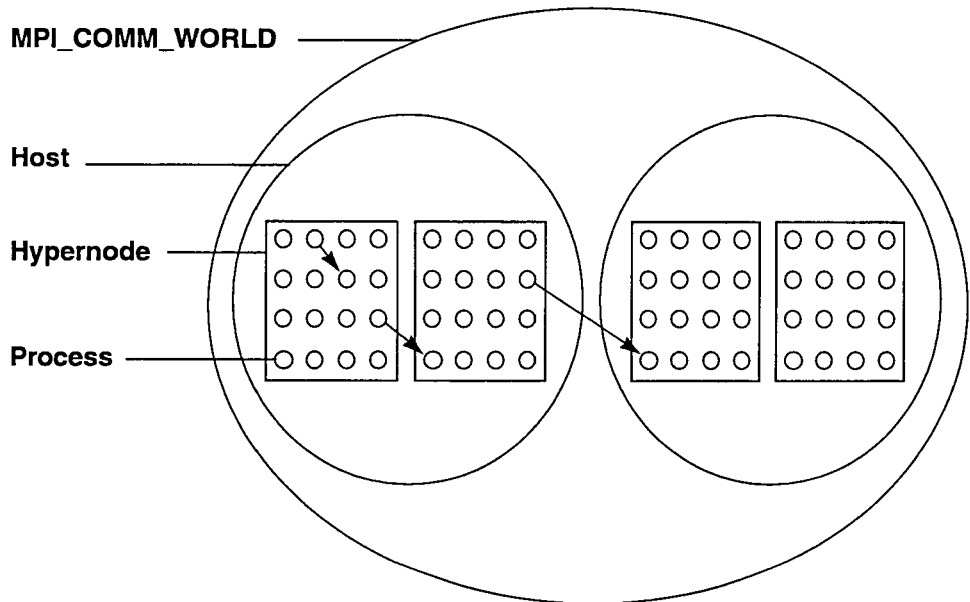
See “Creating an appfile” on page 47 for more information about using appfiles.

Multiprotocol messaging

Multiprotocol messaging refers to process communication that uses different protocols depending upon where the processes are located and what type of Exemplar system is used.

An example configuration for an X-Class server is shown in Figure 3-1.

Figure 3-1 Multiprotocol messaging with an X-Class server



The circles within each hypernode represent processes. The arrows represent message passing. An arrow originates from the sending process and terminates at the receiving process.

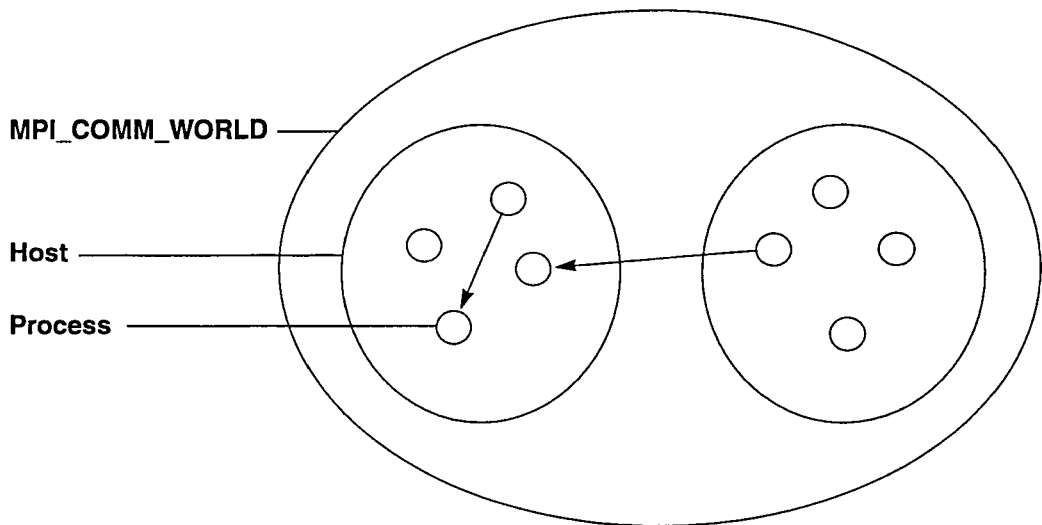
Point-to-point and collective protocols on an X-Class server support messaging between:

- Processes on the same host (on the same or different hypernodes)—Data is transferred using optimized byte-copy and global shared memory.
- Processes on different hosts—Data is transferred using TCP/IP.

The communication speed of protocols for servers running under SPP-UX is fastest for processes on the same hypernode, slower for processes on different hypernodes in the same host, and slowest for processes on different hosts.

An example configuration for a K-Class server is shown in Figure 3-2.

Figure 3-2 Multiprotocol messaging with a K-Class server



The circles within each host represent processes. The arrows represent message passing. An arrow originates from the sending process and terminates at the receiving process.

Point-to-point and collective protocols on servers running under HP-UX support messaging between:

- Processes on the same host (on the same or different hypernodes)—Data is transferred using optimized byte-copy and global shared memory.
- Processes on different hosts—Data is transferred using TCP/IP.

Run-time environment variables

Environment variables are used to alter the way HP MPI executes an application. The variable settings determine how an application behaves and how an application allocates internal resources at run time.

Many applications run without setting any environment variables. However, applications that use a large number of nonblocking messaging requests, require debugging support, or need to control process placement may need a more customized configuration.

Environment variables are always local to the system where `mpirun` is running. To propagate environment variables to remote hosts, you must specify each variable in an appfile using the `-e` option. See “Creating an appfile” on page 47 for more information.

The environment variables listed below affect the behavior of HP MPI at run time:

- `MPI_FLAGS`
- `MPI_GLOBSIZE`
- `MPI_TOPOLOGY`
- `MPI_SHMEMCNTL`
- `MPI_TMPDIR`
- `MPI_XMPI`
- `MPI_WORKDIR`
- `MPI_CHECKPOINT`
- `MPI_INSTR`

MPI_FLAGS

MPI_FLAGS modifies the general behavior of HP MPI. The MPI_FLAGS syntax is shown below:

```
[ecxdb,][edde,][exdb,][egdb,][j,][l,][s[a|p][#],][v,][+E2]
```

where

- | | |
|-------|---|
| ecxdb | Starts a separate CXdb session for each process. The debugger must be in the command search path. This option is only provided for backward compatibility on servers running under SPP-UX. See “Debugging HP MPI applications” on page 92 for more information. |
| edde | Starts the application under the DDE debugger. The debugger must be in the command search path. This option is only supported on servers running under HP-UX. See “Debugging HP MPI applications” on page 92 for more information. |
| exdb | Starts the application under the xdb debugger. The debugger must be in the command search path. This option is only supported on servers running under HP-UX. See “Debugging HP MPI applications” on page 92 for more information. |
| egdb | Starts the application under the gdb debugger. The debugger must be in the command search path. This option is only supported on servers running under HP-UX. See “Debugging HP MPI applications” on page 92 for more information. |
| j | Prints the HP MPI job identifier. |
| l | Reports memory leaks caused by erroneous handling of HP MPI objects. Setting this option may decrease performance. |

Understanding HP MPI
Running applications

`s[a|p][#]` Selects signal and maximum time-delay for guaranteed message progression. The `sa` option selects `SIGALRM`. The `sp` option selects `SIGPROF`. The `#` option is the number of seconds to wait before issuing a signal to trigger message progression. The default value of this option is `sp604800`, which issues a `SIGPROF` once a week.

This mechanism is used to guarantee message progression in applications that use nonblocking messaging requests followed by prolonged periods of time in which HP MPI routines are not called.

NOTE

The `SIGPROF` option is not supported on servers running under `SPP-UX` when your application executable is in Extended Standard Object Module format.

`v` Prints the version number.

`+E2` Sets `-1` as the value of `.TRUE.` and `0` as the value for `FALSE.` when returning logical values from HP MPI routines called within Fortran 77 applications.

MPI_GLOBMEMSIZE

`MPI_GLOBMEMSIZE` specifies the amount of shared memory allocated for all processes in an HP MPI application. `MPI_GLOBMEMSIZE` has the following syntax:

amount

where *amount* specifies the total amount of shared memory in bytes for all processes. The default is 2 Mbytes for up to 64-way applications and 4 Mbytes for larger applications.

NOTE

Be sure that the value specified for `MPI_GLOBMEMSIZE` is less than the amount of global shared memory allocated for the subcomplex when working with X-Class servers. Otherwise, swapping overhead will degrade application performance.

MPI_TOPOLOGY

MPI_TOPOLOGY controls application process placement within a subcomplex on servers running under SPP-UX (the value is ignored on HP-UX systems). MPI_TOPOLOGY has the following syntax:

```
[ [sc] / [hypernode] : ] [topology]
```

where

- | | |
|------------------|---|
| <i>sc</i> | Identifies the name of a subcomplex. |
| <i>hypernode</i> | Specifies the logical hypernode within the subcomplex on which to start the first process. By default, the initial logical hypernode is chosen by the operating system. |
| <i>topology</i> | Is a comma-separated list that specifies the number of processes to start on each logical hypernode in the subcomplex, beginning with logical hypernode 0. |

HP MPI uses logical hypernode numbering. The operating system handles the mapping from physical to logical hypernodes. This mapping follows the lowest-to-highest sorted order of physical hypernode numbers. For example, in a 2-node subcomplex using physical hypernodes 3 and 4, physical hypernode 3 would map to logical hypernode 0, and physical hypernode 4 would map to logical hypernode 1.

An MPI_TOPOLOGY value of `System/3:4,0,4,4` specifies that logical hypernodes zero, two, and three of the subcomplex `System` each run four processes. The first application process is started on logical hypernode three.

When running a multinode application where some processes run different executables, MPI_TOPOLOGY settings in the appfile override any settings you might have specified by setting MPI_TOPOLOGY from the command line. See “Creating an appfile” on page 47 for more information.

The number of processes specified using MPI_TOPOLOGY must match the number of processes specified in `mpirun`. For example, if you set MPI_TOPOLOGY to `2,3` and invoke `mpirun` with `-np 6`, the system generates an error message and terminates your job.

Also, be sure that the number of hypernodes specified in `MPI_TOPOLOGY` matches the number of available hypernodes on the subcomplex you want to use. For example, if you set `MPI_TOPOLOGY` to 6, 2, 3 and System only contains hypernodes 0 and 1, the system will generate an error message and terminate your job. To prevent this, use the `scm` utility to determine the configuration of system subcomplexes before invoking `mpirun`.

NOTE

The default subcomplex on all systems is called System. Use the `mpa` utility to change the default to another subcomplex.

MPI_SHMEMCNTL

`MPI_SHMEMCNTL` controls the subdivision of each process's shared memory for the purposes of point-to-point and collective communications. `MPI_SHMEMCNTL` syntax is shown below:

`[nenv,] [frag,] [generic]`

where

- | | |
|----------------|--|
| <i>nenv</i> | Specifies the number of envelopes per process pair. The default is 8. |
| <i>frag</i> | Denotes the size in bytes of the message-passing fragments region. The default is 87.5 percent of shared memory. |
| <i>generic</i> | Specifies the size in bytes of the generic-shared memory region. The default is 12.5 percent of shared memory. |

MPI_TMPDIR

By default, HP MPI uses the `/tmp` directory to store temporary files needed for its operations. `MPI_TMPDIR` is used to point to a different temporary directory. `MPI_TMPDIR` syntax is shown below:

`directory`

where *directory* specifies an existing directory used to store temporary files.

MPI_XMPI

MPI_XMPI specifies options for run-time raw trace generation. These options represent an alternate way to set tracing rather than using the trace options supplied with `mpirun`.

The argument list for MPI_XMPI contains the prefix name for the file where each process writes its own raw trace data. Each process creates its own filename by concatenating the prefix, a period, and the process's global rank number.

For example, if a process has rank 0 and the prefix is `hello_world`, the process's raw trace file would be `hello_world.0`. If the file prefix name does not begin with a forward slash (/) (for example, `/tmp/test`), the raw trace file is stored in the directory in which the process is executing `MPI_Init`.

MPI_XMPI syntax is shown below:

```
prefix[:bs###] [:nc] [:off] [:s] [:nf] [:k]
```

where

<i>prefix</i>	Specifies the tracing output file prefix. This is a required parameter.
<i>bs###</i>	Denotes the buffering size in kbytes for dumping raw trace data. Actual buffering size may be rounded up by the system. The default buffering size is 4096 kbytes. Specifying a large buffering size reduces the need to flush raw trace data to a file when process buffers reach capacity. Flushing too frequently can cause communication routines to run slower. If this problem occurs, increase the buffering size.
<i>nc</i>	Specifies no clobber, which means that an HP MPI application aborts if a file with the name specified in <i>prefix</i> already exists.
<i>off</i>	Denotes that trace generation is initially turned off and only begins after all processes collectively call <code>MPIHP_Trace_on</code> .
<i>s</i>	Specifies a simpler tracing mode by omitting tracing for <code>MPI_Test</code> , <code>MPI_Testall</code> , <code>MPI_Testany</code> , and <code>MPI_Testsome</code> calls that do not complete a request. This option may reduce the size of trace data so that <code>xmpi</code> runs faster.

Understanding HP MPI

Running applications

- nf** Denotes that a consolidated trace file is not generated. In addition, raw trace files are not deleted. You may want to use this option if your application contains a large number of processes, and you do not want to wait for `MPI_Finalize` to consolidate the raw trace files before your application terminates.
- k** Specifies that raw trace files are kept.

NOTE

Even though you can specify tracing options through the `MPI_XMPI` environment variable, the recommended approach is to use the `mpirun` command with the `-t` option instead. In this case, the specifications you provide with the `-t` option take precedence over any specifications you may have set with `MPI_XMPI`. Using `mpirun` to specify tracing options guarantees that multihost applications do tracing in a consistent manner. See “`mpirun`” on page 45 for more information.

NOTE

Trace-file generation (in conjunction with XMPI) and counter instrumentation are mutually exclusive profiling techniques.

MPI_WORKDIR

By default, HP MPI applications execute in the directory where they are started. `MPI_WORKDIR` changes the execution directory. `MPI_WORKDIR` has the following syntax:

directory

where *directory* specifies an existing directory where you want the application to execute.

MPI_CHECKPOINT

You can checkpoint and restart HP MPI applications running under SPP-UX on a single subcomplex by setting `MPI_CHECKPOINT`. In this case, you cannot start your application using `mpirun`. `MPI_CHECKPOINT` does not require specific arguments. For example, to checkpoint and restart the `hello_world` application:

```
% setenv MPI_CHECKPOINT  
% hello_world -np 4
```

When you use `MPI_CHECKPOINT`, the following limitations apply:

- Your HP MPI job is not assigned a job ID. You cannot monitor or terminate the job using the `mpijob` or `mpiclean` utilities respectively. Also, no job ID is printed when the `j` option is set in `MPI_FLAGS`. If your application crashes or hangs, you must do manual cleanup using the `kill` or `ipcrm` UNIX utilities.
- `MPI_Abort` does not kill peer processes in the communicator. In this case, only the calling process terminates.
- Direct process-to-process byte-copy is disabled. This results in a bandwidth reduction for large message transfers.

MPI_INSTR

`MPI_INSTR` enables counter instrumentation for profiling HP MPI applications. The measurements collected are similar to the reports generated by `mpitrstat`. `MPI_INSTR` has the following syntax:

```
prefix [:b#1, #2] [:nc] [:off] [:n1] [:np] [:nm] [:c]
```

where

- | | |
|------------------------|---|
| <i>prefix</i> | Specifies the instrumentation output file prefix. The rank zero process writes the application's measurement data to <i>prefix</i> .instr. If the prefix does not represent an absolute pathname, the instrumentation output file is opened in the working directory of the rank zero process when <code>MPI_Init</code> is called. |
| <i>b#1</i> , <i>#2</i> | Redefines the instrumentation message bins to include a bin having byte range <i>#1</i> and <i>#2</i> inclusive. The high bound of the range can be infinity, representing the largest possible message size. |

Understanding HP MPI

Running applications

<code>nc</code>	Specifies no clobber. If the instrumentation output file exists, <code>MPI_Init</code> aborts.
<code>off</code>	Denotes that counter instrumentation is initially turned off and only begins after all processes collectively call <code>MPIHP_Trace_on</code> .
<code>nl</code>	Specifies not to dump a long breakdown of the measurement data to the instrumentation output file (in this case, do not dump minimum, maximum, and average time data).
<code>np</code>	Denotes not to dump a per-process breakdown of the measurement data to the instrumentation output file.
<code>nm</code>	Specifies not to dump message-size measurement data to the instrumentation output file.
<code>c</code>	Specifies not to dump time measurement data to the instrumentation output file.

See “Using counter instrumentation” on page 54 for more information.

NOTE

Even though you can specify profiling options through the `MPI_INSTR` environment variable, the recommended approach is to use the `mpirun` command with the `-i` option instead. Using `mpirun` to specify profiling options guarantees that multihost applications do profiling in a consistent manner. See “`mpirun`” on page 45 for more information.

NOTE

Counter instrumentation and trace-file generation (used in conjunction with XMPI) are mutually exclusive profiling techniques.

Run-time utility commands

HP MPI provides a set of utility commands to supplement the MPI library routines. These commands include:

- `mpirun`
- `mpiclean`
- `mpijob`
- `xmpi`
- `mpitrget`
- `mpitrstat`

mpirun

`mpirun` starts an HP MPI application.

`mpirun` syntax has two forms:

- `mpirun [-np #] [-help] [-version] [-djpW] [-t spec] [-i spec] [-h host] [-l user] [-e var[=val]][...] [-sp paths] program [args]`
- `mpirun [-help] [-version] [-jpvW] [-t spec] [-i spec] [-f appfile]`

where

- | | |
|-----------------------|--|
| <code>-np #</code> | Specifies the number of processes to run. |
| <code>-help</code> | Prints usage information for the utility. |
| <code>-version</code> | Prints the version information. |
| <code>-j</code> | Prints the HP MPI job ID. |
| <code>-p</code> | Turns on pretend mode. That is, go through the motions of starting an HP MPI application but do not create any processes. This is useful for debugging and checking whether the <i>appfile</i> (if used) is setup correctly. |
| <code>-v</code> | Turns on verbose mode. |
| <code>-W</code> | Does not wait for the application to terminate before returning. |

Running applications

<code>-t <i>spec</i></code>	Enables run-time raw trace generation for all processes. <i>spec</i> specifies options used when tracing. See “MPI_XMPI” on page 41 for the list of options you can use.
<code>-i <i>spec</i></code>	Enables run-time instrumentation profiling for all processes. <i>spec</i> specifies options used when profiling. See “MPI_INSTR” on page 43 for the list of options you can use.
<code>-h <i>host</i></code>	Starts the processes on <i>host</i> (default is localhost).
<code>-l <i>user</i></code>	Specifies the user name on the target host (default is local username).
<code>-e <i>var</i> [=<i>val</i>]</code>	Sets the environment variable <i>var</i> for the program and gives it the value <i>val</i> if provided. Environment variable substitutions (for example, \$FOO) are supported in the <i>val</i> argument.
<code>-sp <i>paths</i></code>	Sets the target shell PATH environment variable to <i>paths</i> . Search paths are separated by the colon (:) character.
<i>program</i>	Specifies the name of the executable to run.
<i>args</i>	Specifies command-line arguments to the program.
<code>-f <i>appfile</i></code>	Starts the application described in <i>appfile</i> .

The first syntax is used for applications where all processes execute the same program on the same host. For example:

```
% mpirun -j -np 3 send_receive
```

runs the `send_receive` application with three processes and prints out the job ID.

The second syntax must be used for applications that consist of multiple programs or that run on multiple hosts or subcomplexes. In this case, each program called by the application is listed in a file called an `appfile`. For example:

```
% mpirun -t my_trace:k -f my_appfile
```

enables tracing, sets the prefix of the tracing output file to `my_trace`, specifies that the raw trace files are kept, and runs an `appfile` named `my_appfile`.

Creating an appfile

The format of entries in an appfile is line oriented. Lines that end with the backslash (\) character are continued on the next line, forming a single logical line. A logical line starting with the pound (#) character is treated as a comment. Each program, along with its arguments, is listed on a separate logical line.

You can specify the `-h`, `-l`, `-np`, `-e`, and `-sp` options (from the `mpirun` command) in an appfile. Options following a program name are treated as the program's command line arguments and are not processed by `mpirun`.

The ranks of the processes in `MPI_COMM_WORLD` are guaranteed to be ordered according to their sequential order in an appfile.

The general form of an appfile entry is:

```
[-h remote_host] [-e var [=val] [...]] [-l user] [-sp paths] [-np #] program  
[args]
```

where

- `-h remote_host` Specifies the remote host where a remote executable is stored (defaults to local host). *remote_host* is either a host name or an IP address.
- `-e var=val` Sets the environment variable *var* for the program and gives it the value *val* if provided (defaults to not setting environment variables).
- `-l user` Specifies the user name on the target host (default is current user name).
- `-sp paths` Sets the target shell PATH environment variable to *paths*. Search paths are separated by the colon (:) character (default is do not override the path).
- `-np #` Specifies the number of processes to run (defaults to one).
- program* Specifies the name of the executable to run. The executable is searched for in \$PATH.
- args* Specifies command line arguments to the program.

Understanding HP MPI

Running applications

One way to set environment variables on remote hosts is to use the `-e` option in the appfile:

```
-h remote_host -e MPI_TOPOLOGY=val [-np #] program [args]
```

Alternatively, you can set environment variables using the `.cshrc` file on each remote host (only for users that use a `/bin/csh`-based shell).

mpijob

`mpijob` lists the HP MPI jobs running on the system. The `mpijob` syntax is shown below:

```
mpijob [-help] [-a] [-u] [-j id [...]]
```

where

<code>-help</code>	Prints usage information for the utility.
<code>-a</code>	Lists jobs for all users.
<code>-u</code>	Sorts jobs by user name.
<code>-j <i>id</i></code>	Provides process status for job <i>id</i> .

When invoked, `mpijob` reports the following information for each job:

JOB	HP MPI job identifier.
USER	User name of the owner.
NPROCS	Number of processes.
PROGRAMNAME	Program names used in the HP MPI application.

By default, your jobs are listed by job ID in increasing order. However, you can specify the `-a` and `-u` options to change the default behavior.

If you specify the `-j` option, `mpijob` reports the following information for each job:

RANK	Rank for each process in the job.
HOST	Host where the job is running.
PID	Process identifier for each process in the job.
LIVE	Option that indicates whether the process is running (an <code>x</code> is used) or has been terminated.
PROGRAMNAME	Program names used in the HP MPI application.

An `mpijob` output using the `-a` and `-u` options is shown below. The output lists jobs for all users and sorts them by user name.

JOB	USER	NPROCS	PROGNAME
22623	charlie	12	/home/watts
22573	keith	14	/home/richards
22617	mick	100	/home/jagger
22677	ron	4	/home/wood

NOTE

You should invoke `mpijob` on the host on which you initiated `mpirun`.

mpiclean

`mpiclean` kills lingering processes in a running HP MPI application. `mpiclean` syntax has three forms:

- `mpiclean [-help] [-v] -j id [...]`
- `mpiclean [-help] [-v] [-sc name | -scid id] prog [...]`
- `mpiclean [-help] [-v] -m`

where

<code>-help</code>	Prints usage information for the utility.
<code>-v</code>	Turns on verbose mode.
<code>-m</code>	Cleans up your shared-memory segments.
<code>-j <i>id</i></code>	Kills the processes of job number <i>id</i> . You can specify multiple job IDs.
<code>-sc <i>name</i></code>	Restricts the operation to the named subcomplex. This option is mutually exclusive with the <code>-scid</code> option.
<code>-scid <i>id</i></code>	Restricts the operation to subcomplex number <i>id</i> . This option is mutually exclusive with the <code>-sc</code> option.
<i>prog</i>	Specifies the binary filename to kill. You can specify multiple filenames.

The first syntax is used for all servers. The second syntax is provided for backward compatibility on servers running under SPP-UX. The third syntax is used when an application aborts during `MPI_Init`, and the termination of processes does not destroy the allocated shared-memory segments.

The MPI library checks for the abnormal termination of processes while your application is running. In some cases, application bugs can cause processes to deadlock and linger in the system. When this occurs, you can use `mpijob` to identify hung jobs and `mpiclean` to kill all processes in the hung application.

There are two ways to kill an HP MPI application. The preferred way is to provide `mpiclean` with the application's job ID (obtained by using the `-j` option when invoking `mpirun`). However, you can only kill jobs that you own.

The second way is only provided on servers running under SPP-UX for backward compatibility. In this approach, you specify `mpiclean` with a list of binary filenames you own. `mpiclean` locates the matching processes and kills them.

You can restrict the second cleanup method to a single subcomplex by using the `-sc` or `-scid` options. This is helpful in cases where the same code is running independently on several subcomplexes and only one of these applications needs to be killed.

NOTE

You should invoke `mpiclean` on the host on which you initiated `mpirun`.

xmpi

`xmpi` invokes the XMPI utility. The `xmpi` syntax is shown below:

```
xmpi [-h] [-bg arg] [-bd arg] [-bw arg] [-display arg] [-fg arg]  
[-geometry arg] [-iconic] [-title arg]
```

where

- | | |
|-----------------------------------|---|
| <code>-h</code> | Prints usage information for the utility. |
| <code>-bg <i>arg</i></code> | Specifies the background color. |
| <code>-bd <i>arg</i></code> | Denotes the border color. |
| <code>-bw <i>arg</i></code> | Specifies the width of the border in pixels. |
| <code>-display <i>arg</i></code> | Designates the X-window display server to use. |
| <code>-fg <i>arg</i></code> | Specifies the foreground color. |
| <code>-geometry <i>arg</i></code> | Specifies size and position. |
| <code>-iconic</code> | Designates that the application start as an icon. |
| <code>-title <i>arg</i></code> | Specifies the title of the application. |

For more information, see “Using XMPI” on page 58.

mpitrget

`mpitrget` combines raw trace files into a consolidated output file with a `.tr` suffix. This output file is then loaded and reviewed on XMPI.

NOTE

`mpitrget` is obsolete for HP MPI V1.3. Consolidation of raw trace files now occurs automatically when your application calls `MPI_Finalize`.

mpitrstat

`mpitrstat` provides profiling information for HP MPI applications. Counter instrumentation has subsumed the functionality provided by `mpitrstat`. However, if you want more information about this command, see the appropriate man page.

Understanding HP MPI
Running applications

This chapter provides information about utilities used to analyze HP MPI applications. The topics covered are:

- Using counter instrumentation
- Using XMPI
- Using CXperf
- Using the profiling interface

Using counter instrumentation

Counter instrumentation provides cumulative statistics for your applications. Counter instrumentation is the recommended method for collecting profiling data because it is faster and less intrusive than `mpitrstat`.

Creating an instrumentation profile

To create an instrumentation profile, enter:

```
% mpirun -i spec -np # program
```

where

-i spec Enables run-time instrumentation profiling for all processes. *spec* provides options used when profiling. See “MPI_INSTR” on page 43 for information about options you can use.

You must specify the `-i` option before the program name.

-np # Specifies the number of processes to run.

program Specifies the name of the executable to run.

HP MPI provides the nonstandard `MPIHP_Trace_on` and `MPIHP_Trace_off` routines to collect profile information for selected code sections only (by default, the entire application is profiled from `MPI_Init` to `MPI_Finalize`). You insert the `MPIHP_Trace_on` and `MPIHP_Trace_off` pair around code that you want to profile. Then, you build the application and invoke `mpirun` using the appropriate syntax.

A sample instrumentation profile for the compute_pi.f application is shown below. In this case, instrumentation was invoked by entering:

```
% mpirun -i -np 2 compute_pi.exe
```

The overhead time in the profile represents the time a process or routine spends inside MPI. For example, the time a process spends doing message packing.

The blocking time in the profile represents the time a process or routine is blocked, waiting for communication to complete before resuming execution.

Version: HP MPI 01.03.00.00 - HP-UX 10.20

Date: Thu Nov 6 11:12:28 1997

Scale: Wall Clock Seconds

Processes: 2

User: 25.23%

MPI: 74.77% [Overhead:74.77% Blocking:0.00%]

Application Summary by Rank:

Rate	Duration	Overhead	Blocking	User	MPI
0	0.710857	0.504448	0.000000	29.04%	70.96%
1	0.723631	0.537273	0.000000	25.75%	74.25%

Routine Summary:

Routine	Calls	Overhead	Blocking
MPI_Bcast	2	1.898292	0.000000
min		0.157720	0.000000
max		0.249885	0.000000
avg		0.189829	0.000000
MPI_Init	2	1.616565	0.000000

Profiling
Using counter instrumentation

min		0.135646	0.000000
max		0.179480	0.000000
avg		0.161656	0.000000
MPI_Reduce	2	0.944703	0.000000
min		0.000089	0.000000
max		0.128702	0.000000
avg		0.094470	0.000000
MPI_Finalize	2	0.941750	0.000000
min		0.061623	0.000000
max		0.118828	0.000000
avg		0.094175	0.000000

Routine Summary by Rank:

Routine	Rank	Calls	Overhead	Blocking
MPI_Bcast	0	1	0.249885	0.000000
	1	1	0.157720	0.000000
MPI_Init	0	1	0.135646	0.000000
	1	1	0.157633	0.000000
MPI_Reduce	0	1	0.000089	0.000000
	1	1	0.119893	0.000000
MPI_Finalize	0	1	0.118828	0.000000
	1	1	0.102027	0.000000

Message Summary:

Routine	Message Bin	Count
MPI_Bcast	[0.32]	2
MPI_Reduce	[0.32]	2

Message Summary by Rank:

Routine	Message Bin	Count	Rank
MPI_Bcast	[0.32]	1	0
	[0.32]	1	1
MPI_Reduce	[0.32]	1	0
	[0.32]	1	1

Using XMPI

XMPI is an X/Motif graphical user interface for running applications, monitoring processes and messages, and viewing trace files. XMPI provides a graphical display of the state of processes within an HP MPI application.

XMPI is useful when analyzing programs at the application level (for example, examining HP MPI data types and communicators). Unlike other profilers and debuggers, you can run XMPI without having to recompile or relink your application.

XMPI runs in one of two modes: postmortem mode or interactive mode. In postmortem mode, you can view trace information for each process in your application. In interactive mode, you can monitor process communications by taking snapshots while your application is running.

The default X resource settings that determine how XMPI displays on your workstation are stored in `/opt/mpi/lib/X11/app-defaults/XMPI`. See “XMPI resource file” on page 163 for a list of these settings.

Working with postmortem mode

To use XMPI's postmortem mode, you must first create a trace file. Then, you can load this file into XMPI to view state information for each process in your application.

Creating a trace file

To create a trace file, enter:

```
% mpirun -t spec -np # program
```

where

-t spec Enables run-time raw trace generation for all processes. *spec* specifies options used when tracing. See “MPI_XMPI” on page 41 for information about options you can specify.

You must specify the **-t** option before the program name.

-np # Specifies the number of processes to run.

program Specifies the name of the executable to run.

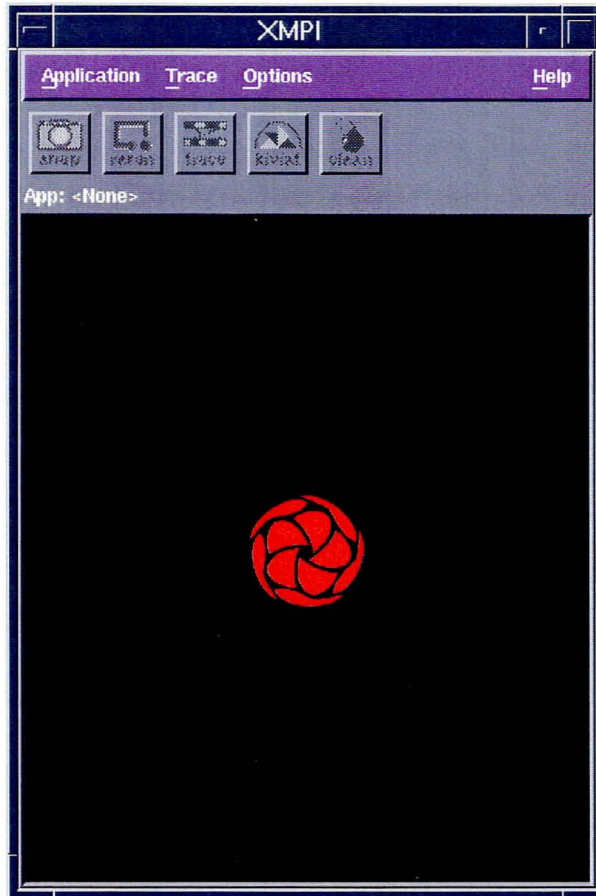
When you use the **-t** option to enable trace generation, you must specify the prefix name used for each raw trace file as part of *spec*. Then, when `mpirun` is invoked, a raw trace dump, *prefix.n*, is created for each application process where *n* ranges from 0 to (# - 1). `MPI_Finalize` consolidates all the raw trace dump files into a single file (*prefix.tr*) that you can load into XMPI.

HP MPI provides the nonstandard `MPIHP_Trace_on` and `MPIHP_Trace_off` routines to help troubleshoot application problems. You insert the `MPIHP_Trace_on` and `MPIHP_Trace_off` pair around suspect code in your application. Then, you build the application and invoke `mpirun` with `-t:off` to enable application tracing. The trace information collected is only for the code between `MPIHP_Trace_on` and `MPIHP_Trace_off`. You can then run the trace file in XMPI to identify problems during application execution.

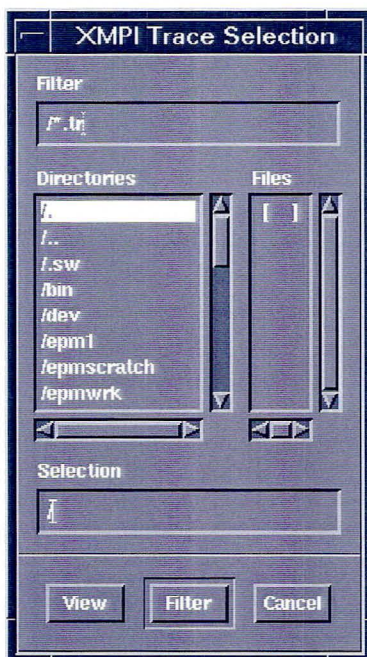
Viewing a trace file

Use these instructions to view a trace file:

- Step 1.** Enter `xmpi` to open the XMPI main window (see “`xmpi`” on page 50 for information about other options you can specify).



Step 2. Select View from the Trace menu to open the XMPI Trace Selection dialog.



Profiling
Using XMPI

Step 3. Type the full path name of the appropriate trace file in the Selection field and choose View to open the XMPI Trace dialog.



NOTE

When viewing trace files containing multiple segments (that is, multiple MPIHP_Trace_on and MPIHP_Trace_off pairs), XMPI prompts you for the number of the segment you want to view. If you want to view a different segment later, simply reload the trace file and specify the new segment number when prompted.

The XPMI Trace dialog consists of an icon bar across the top, the current magnification and dial time just below, and a trace log display area below that.

The icon bar contains icons that (from left to right):

- Increase the magnification of the trace log.
- Decrease the magnification of the trace log.

- Rewind the trace log to the beginning. The dial time is also reset to the beginning.
- Stop playing the trace log.
- Play the trace log.
- Fast forward the trace log.

To set the magnification for viewing a trace file, select the Increase or Decrease icon on the icon bar.

The dial time indicates how long the application has been running in seconds.

The trace log display area shows a separate trace for each process in the application. The dial time is represented as a vertical line. The rank for each process is shown where the dial time line intersects a process trace.

The state of a process at any time is indicated by one of three colors:

Green	Signifies that a process is running outside MPI.
Red	Denotes that a process is blocked, waiting for communication to finish before the process resumes execution.
Yellow	Represents a process's overhead time inside MPI (for example, time spent doing message packing).

Blocking point-to-point communications are represented by a trace for each process showing the time spent in system overhead and time spent blocked waiting for communication. A line is drawn connecting the appropriate send and receive trace segments. The line starts at the beginning of the send segment and ends at the end of the receive segment.

For nonblocking point-to-point communications, a system overhead segment is drawn when a send and receive are initiated. When the communication is completed using a wait or a test, segments are drawn showing system overhead and blocking time. Lines are also drawn between matching sends and receives, except in this case, the line is drawn from the segment where the send was initiated to the segment where the corresponding receive completed.

Profiling Using XMPI

Collective communications are also represented by a trace for each process showing the time spent in system overhead and time spent blocked waiting for communication.

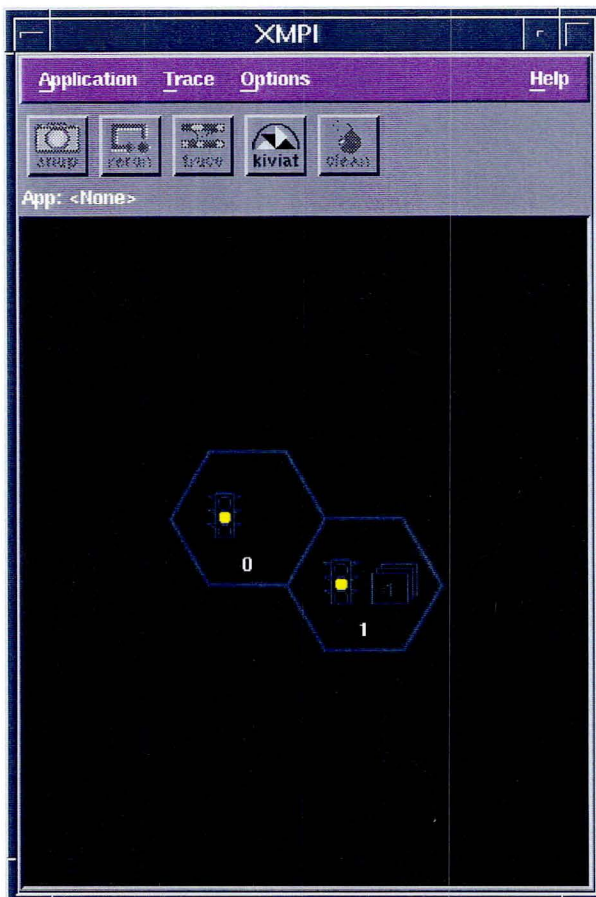
Owing to the use of partial tracing, some send and receive segments may not have a matching segment. In this case, a stub line is drawn out of the send segment or into the receive segment.

To play the trace file, select the Play or Fast forward icons on the icon bar. For any given dial time, the state of the trace file is reflected in the main window, the Focus dialog, the Datatype dialog, and the Kiviat dialog.

Viewing process information from a trace

Use these instructions to view process information from a trace.

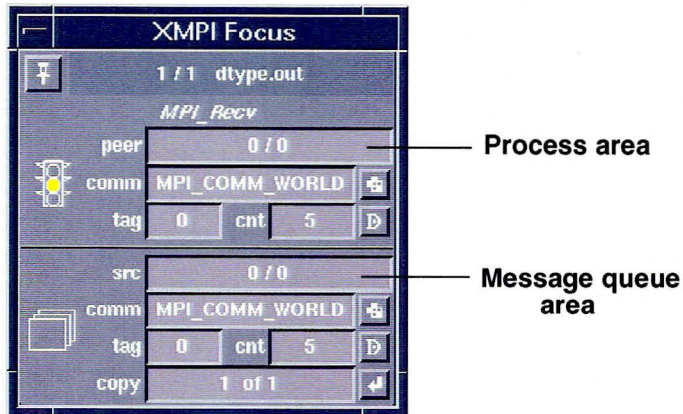
- Step 1.** Start XMPI and open a trace for viewing. The XMPI main window fills with a group of tiled hexagons, each representing the current state of a process and labelled by the process's rank within MPI_COMM_WORLD.



The current state of a process is indicated by the color of the signal light (either green, red, or yellow) in the hexagon. This color corresponds to the elapsed run time (current dial time) of the trace file in the XMPI Trace dialog. As the trace file is played, the color changes as processes communicate with each other.

Profiling
Using XMPI

Step 2. Select the hexagon representing the process you want more information about to open the XMPI Focus dialog.



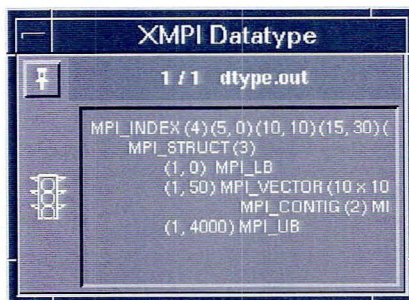
The XMPI Focus dialog consists of a process area and a message queue area.

The values in the process area and message queue area fields correspond to the current dial time of the trace file in the XMPI Trace dialog. As the trace file is played, the values in the fields change as processes communicate with each other.

The process area describes the state of a process together with the name and arguments for the HP MPI function being executed. The fields include:

- | | |
|------|--|
| peer | Displays the rank of the displayed function's peer process. A process is identified by its rank in MPI_COMM_WORLD, a slash (/), and the rank of the process within the current communicator. |
| comm | Shows the communicator being used by the HP MPI function. If you select the icon to the right of the comm field, the hexagons for processes that belong to the communicator are highlighted in the XMPI main window. |
| tag | Displays the value of the tag argument associated with the message. |

cnt Shows the count of the message data elements associated with the message when it was sent. Select the icon to the right of the cnt field to open the XMPI Datatype dialog.



The XMPI Datatype dialog displays the type map of the data type associated with the message when it was sent. This data type can be one of the predefined data types or a user-defined data type.

The data type shown corresponds to the current dial time of the trace file in the XMPI Trace dialog. As the trace file is played, the data type changes as processes communicate with each other.

The message queue area describes the current state of the queue of messages sent to the process but not yet received. The fields include:

- src Displays the rank of the process sending the message. A process is identified by its rank in MPI_COMM_WORLD, a slash (/), and the rank of the process within the current communicator.
- comm Shows the communicator being used by the HP MPI function. If you select the icon to the right of the comm field, the hexagons for processes that belong to the communicator are highlighted in the XMPI main window.
- tag Displays the value of the tag argument associated with the message when it was sent.

Profiling
Using XMPI

- cnt** Shows the count of the message data elements associated with the message when it was sent. If you select the icon to the right of the cnt field, the XMPI Datatype dialog displays. The XMPI Datatype dialog displays the type map of the data type associated with the message when it was sent.
- copy** Displays the number of copies of the message that was sent. For example, if a process sends 10 messages to another process where six of the messages have one type of message envelope and the remaining four have another, the copy field toggles between 6 of 10 and 4 of 10. In this case, a message envelope consists of the sender, the communicator, the tag, the count, and the data type.

This behavior results from treating the six messages that all have the same envelope as one copy and the remaining four messages as a different copy. That way, if a communication involves a hundred messages all having the same envelope, you can work with a single copy rather than a hundred copies.

Step 3. Select Quit from the Application menu to close XMPI.

Viewing kiviatic information from a trace file

Kiviatic graphs are used to display performance data. Use these instructions to view kiviatic information from a trace file.

- Step 1.** Start XMPI and open a trace for viewing.
- Step 2.** Select Kiviatic from the Trace menu to open the XMPI Kiviatic dialog.



The XMPI Kiviatic window shows, in a segmented pie-chart format, the cumulative time up to the current dial time spent by each process in the running, overhead, and blocked states.

The cumulative time for each process corresponds to the current dial time of the trace file in the XMPI Trace dialog. As the trace file is played, the cumulative time changes as processes communicate with each other.

You can use the kiviatic view to determine whether processes are load balanced and applications are synchronized. If an application is load balanced, the amount of time processes spend in each state should be equal. If an application is synchronized, the segments representing each of the three states should be concentric.

- Step 3.** Select Quit from the Application menu to close XMPI.

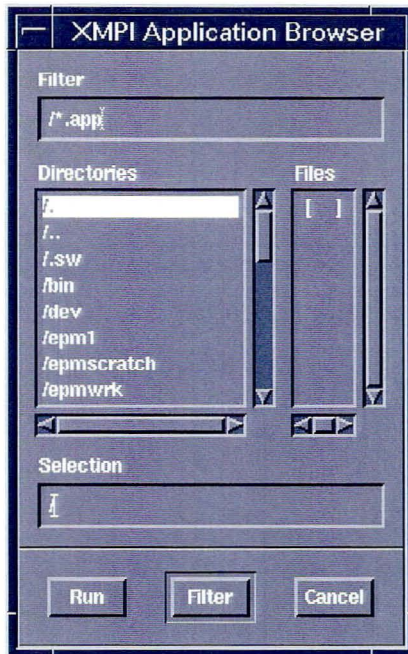
Working with interactive mode

Interactive mode allows you to load and run an existing appfile to view state information for each process in your application.

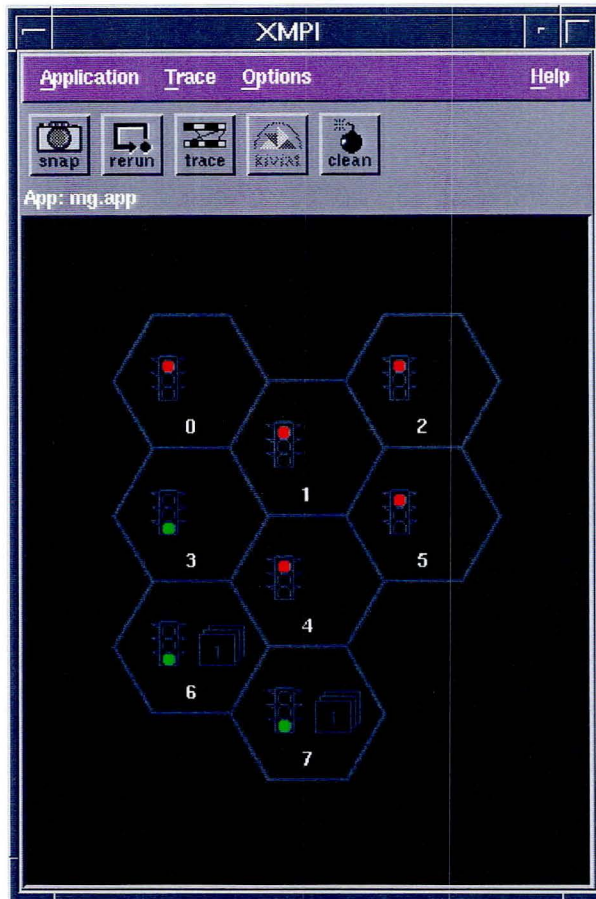
Running an appfile

Use these instructions to run and view an appfile:

- Step 1.** Enter `xmpi` to open the XMPI main window (see “`xmpi`” on page 50 for information about other options you can specify).
- Step 2.** Select Browse&Run from the Application menu to open the XMPI Application Browser dialog.



Step 3. Type the full path name of an existing appfile in the Selection field and choose Run. The XMPI main window fills with a group of tiled hexagons, each representing the current state of a process and labelled by the process's rank within MPI_COMM_WORLD.



The current state of a process is indicated by the color of the signal light (either green, red, or yellow) in the hexagon. These process hexagons disappear when the application has run to completion.

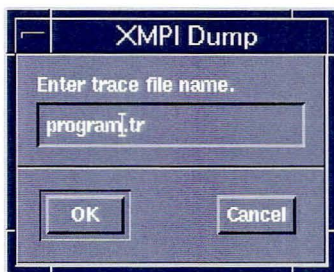
Interactive mode provides the snapshot utility to help debug applications that hang. If automatic snapshot is enabled, XMPI takes periodic snapshots of the application and displays state information for each process on the XMPI main window, the XMPI Focus dialog, and the XMPI Datatype dialog. You can use this information to view the state of each process while the application hangs.

If automatic snapshot is disabled, XMPI displays information for each process when the application begins, but this information is not updated.

Regardless of whether automatic snapshot is enabled, you can take application snapshots manually by selecting Snapshot from the Application menu. In this case, XMPI displays information for each process, but this information is not updated until you take the next snapshot.

You can take snapshots only when an appfile is running. Also, you cannot replay snapshots like trace files.

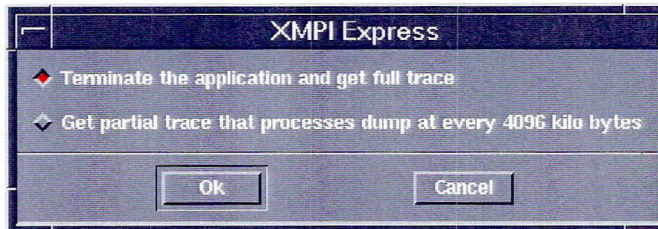
At any time while your application is running, you can select Dump from the Trace menu to open the XMPI Dump dialog.



The Dump option is only available if you have previously selected the Tracing button on the mpirun options trace dialog. Selecting Dump consolidates all raw trace-file data collected up to that point into a single .tr output file.

The single field specifies the name of the consolidated .tr output file. The value you specified for the Prefix field in the mpirun options trace dialog is automatically loaded. You can use this name or choose another. After you have created the .tr output file, you can resume snapshot monitoring.

You can also select Express from the Trace menu while your application is running to open the XMPI Express dialog.



As with the Dump option, the Express option is only available if you have previously selected the Tracing button on the mpirun options trace dialog.

The fields include:

Terminate the application and get full trace

Specifies that the contents of each process buffer (whether partial or full up to that point) are written to a raw trace file. These raw trace files are then consolidated in a .tr output file (previously specified in the Prefix field of the mpirun options trace dialog). Last, the .tr output file is loaded and displayed in the XMPI Trace dialog for viewing.

When you select this field, the XMPI Confirmation dialog displays asking if you are sure you want to terminate the application. You must select Yes before processing will continue.

After the .tr output file is loaded and displayed in the XMPI Trace dialog, you cannot resume snapshot monitoring (the application should have already terminated).

Get partial trace
that processes
dump at every
4096 kilobytes

Specifies that the contents of each process buffer are written to a raw trace file only after the buffer becomes full. These raw trace files are then consolidated to a .tr output file (previously specified in the Prefix field of the mpirun options trace dialog). Last, the .tr output file is loaded and displayed in the XMPI Trace dialog for viewing.

After the .tr output file is loaded and displayed in the XMPI Trace dialog, you cannot resume snapshot monitoring even though the application may still be running.

When using interactive mode, XMPI gathers and displays data from the running appfile or a trace file.

When an application is running, the data source is the appfile, and automatic snapshot is enabled. Even though the application may be creating trace data, the snapshot function does not use it. Instead, the snapshot function acquires data from internal hooks in HP MPI.

At any point in interactive mode, you can load and view a trace file using the View or Express commands under the Trace menu. In this case, the data source switches to the loaded trace file, and the snapshot function is disabled. For HP MPI V1.3, you must rerun your application to switch the data source from a trace file back to an appfile.

- Step 4.** Select Clean from the Application menu at any time to kill the application and close any associated XMPI Focus and XMPI Datatype dialogs. The XMPI Confirmation dialog displays asking if you are sure you want to terminate the application.
- Step 5.** Select Yes to terminate your application and close any associated dialogs. You can then run another application by selecting an appfile from the XMPI Application Browser dialog.

Changing default settings and viewing options

You should initially run your appfile using the XMPI default settings. You can change these default settings and your viewing options later if you like.

Use these instructions to change XMPI's default settings and your viewing options:

- Step 1.** Enter `xmpi` to open the XMPI main window (see “xmpi” on page 50 for information about other options you can specify).
- Step 2.** Select Monitoring from the Options menu to open the XMPI monitor options dialog.



The fields include:

Automatic
snapshot

Enables the automatic snapshot function. If automatic snapshot is enabled, XMPI takes snapshots of the application you are running and displays state information for each process.

If automatic snapshot is disabled, XMPI displays information for each process when the application begins. However, you can only update this information manually. Disabling automatic snapshot may lead to buffer overflow problems because the contents of each process buffer are unloaded every time a snapshot is taken. For communication-intensive applications, process buffers can quickly fill and overflow.

Profiling
Using XMPI

You can enable or disable automatic snapshot while your application is running. This could be useful during troubleshooting when the application has run to a certain point and you want to disable automatic snapshot to study process state information.

Monitor interval
in second

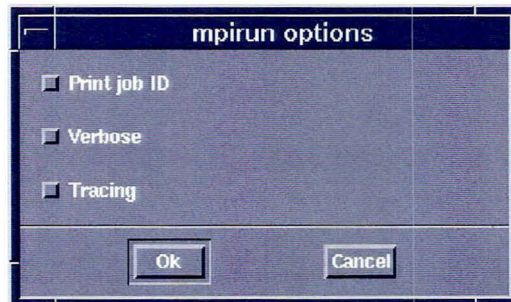
Determines how often XMPI takes a snapshot when automatic snapshot is enabled.

Step 3. Select Buffers from the Options menu to open the XMPI buffer size dialog.



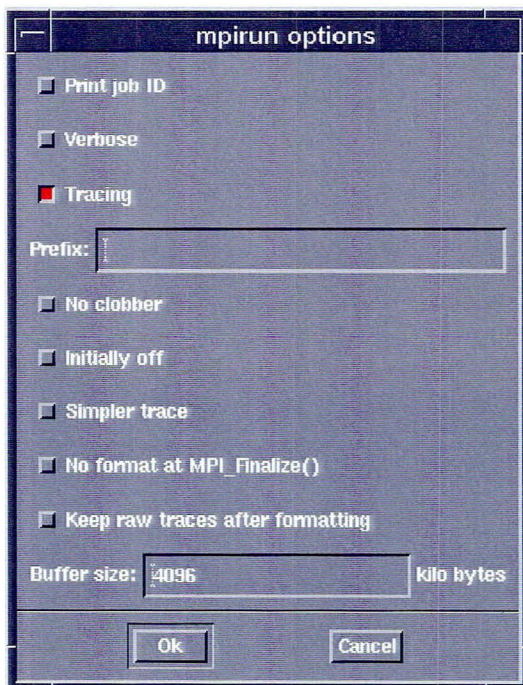
The single field specifies the size of each process buffer. When you run an application, state information for each process is stored in a separate buffer. You may need to increase buffer size if overflow problems occur.

Step 4. Select mpirun from the Options menu to open the mpirun options dialog.



The fields include:

- Print job ID Enables printing of the HP MPI job ID.
- Verbose Enables verbose mode.
- Tracing Enables run-time raw trace generation for all application processes. If you select the Tracing button, the mpirun options trace dialog is opened.



The fields include:

- Prefix Specifies the prefix name for the file where each process writes its own raw trace data. Each process creates its own filename by concatenating the prefix, a period, and the process's global rank number. This is a required field.
- No clobber Specifies no clobber, which means that an HP MPI application aborts if a file with the name specified in the Prefix field already exists.

Profiling
Using XMPI

- Initially off** Specifies that trace generation is initially turned off.
- Simpler trace** Specifies a simpler tracing mode by omitting `MPI_Test`, `MPI_Testall`, `MPI_Testany`, and `MPI_Testsome` calls that do not complete a request.
- No format at `MPI_Finalize()`** Specifies that raw trace files are not consolidated into a single `.tr` output file when `MPI_Finalize` is called. Raw trace-file consolidation can add substantially to the `MPI_Finalize` time when working with large applications.
- Keep raw traces after formatting** Specifies that raw trace files are saved after they are consolidated by `MPI_Finalize`. The default is to delete raw trace files after consolidation.
- Buffer size** Denotes the buffering size in kilobytes for dumping raw trace data. Actual buffering size may be rounded up by the system. The default buffering size is 4096 kilobytes. Specifying a large buffering size reduces the need to flush raw trace data to a file when process buffers reach capacity. Flushing too frequently can increase the overhead for I/O. If this problem occurs, increase the buffering size.

Using CXperf

CXperf allows you to profile each process in an HP MPI application. CXperf replaces the functionality formally provided by CXpa. The profile information is stored in a separate performance data file. During analysis, you merge the data from these separate files into a single performance data file for the application.

With CXperf, you can analyze data using one or more of the following metrics:

- Wall clock time
- CPU time
- Execution counts
- Cache miss counts
- Latency time
- Dynamic call graph

You can display the data as a 3D profile, a 2D profile, a report, or a dynamic call graph. For more information, see *Using CXperf*.

Using the profiling interface

MPI provides a profiling interface for collecting statistics and measuring performance. The profiling interface allows you to intercept calls to the MPI library at link time and perform some action. For example, you may want to measure the time spent in each call to a certain library routine or create a logfile.

All routines in the MPI library begin with the `MPI` prefix. Based on the MPI standard, these routines are also callable using the `PMPI` prefix (for example, `PMPI_Send`).

To use the profiling interface, you write wrapper versions of the MPI library routines you want the linker to intercept. These wrapper routines collect data for some statistic or perform some other action. The wrapper then calls its corresponding routine in the MPI library using its `PMPI` prefix.

For example, suppose you want to measure the elapsed time for each call to `MPI_Send`. In this case, you create a wrapper called `MPI_Send` that uses `MPI_wtime` to measure the elapsed time for each call. After `MPI_wtime` completes, your wrapper then calls `PMPI_Send` from the MPI library to actually send the message.

This chapter provides information about tuning applications to improve performance. The topics covered are:

- General tuning
- SPP-UX platform tuning

General tuning information applies to all applications running on HP-UX and SPP-UX platforms. SPP-UX tuning information applies only to applications running on that platform.

NOTE

The tuning information in this chapter will improve application performance in most but not all cases. You should use the output from counter instrumentation or XMPI to determine which tuning changes are appropriate.

General tuning

When you are developing HP MPI applications, several factors can affect performance. These factors include:

- Message latency and bandwidth
- Multiple network interfaces
- Processor subscription
- MPI routine selection

Message latency and bandwidth

Latency is the time between the initiation of the data transfer in the sending process and the arrival of the first byte in the receiving process.

Latency is often dependent upon the length of messages being sent. An application's messaging behavior can vary greatly based upon whether a large number of small messages or a few large messages are sent.

Message bandwidth is the reciprocal of the time needed to transfer a byte. Bandwidth is normally expressed in megabytes per second. Bandwidth becomes important when message sizes are large.

To improve latency or bandwidth or both:

- Reduce the number of process communications by designing coarse-grained applications.
- Use derived, contiguous data types for dense data structures to eliminate unnecessary byte-copy operations in certain cases. Use derived data types instead of `MPI_Pack` and `MPI_Unpack` if possible. HP MPI optimizes noncontiguous transfers of derived data types.
- Use collective operations whenever possible. This eliminates the overhead of using `MPI_Send` and `MPI_Recv` each time when one process communicates with others. Also, use the HP MPI collectives rather than customizing your own. HP MPI collectives have three-level optimizations when used in a NUMA environment.
- Specify the source process rank whenever possible when calling MPI routines. Using `MPI_ANY_SOURCE` may increase latency.

- Double-word align data buffers if possible. This improves byte-copy performance between sending and receiving processes because of double-word loads and stores.
- Use `MPI_Recv_init` and `MPI_Startall` instead of a loop of `MPI_Irecv` calls in cases where requests may not complete immediately.

For example, suppose you write an application with the following code section:

```
j = 0
for (i=0; i<size; i++) {
    if (i==rank) continue;
    MPI_Irecv(buf[i], count, dtype, i, 0, comm, &requests[j++]);
}
MPI_Waitall(size-1, requests, statuses);
```

Suppose that one of the iterations through `MPI_Irecv` does not complete before the next iteration of the loop. In this case, HP MPI tries to progress both requests. This progression effort could continue to grow if succeeding iterations also do not complete immediately, resulting in a higher latency.

However, you could rewrite the code section as follows:

```
j = 0
for (i=0; i<size; i++) {
    if (i==rank) continue;
    MPI_Recv_init(buf[i], count, dtype, i, 0, comm,
&requests[j++]);
}
MPI_Startall(size-1, requests);
MPI_Waitall(size-1, requests, statuses);
```

In this case, all iterations through `MPI_Recv_init` are progressed just once when `MPI_Startall` is called. This approach avoids the additional progression overhead when using `MPI_Irecv` and can reduce application latency.

Multiple network interfaces

You can use multiple network interfaces for interhost communication while still having intrahost exchanges. In this case, the intrahost exchanges use shared memory between processes mapped to different same-host IP addresses.

Tuning

General tuning

To use multiple network interfaces, you must specify which MPI processes are associated with each IP address in your appfile. To improve performance, you should use the `MPI_TOPOLOGY` environment variable to associate each network interface with the hypernode where it physically resides on SPPUX.

For example, suppose you have two hosts called `host0` and `host1` that each communicate using the two HIPPI cards `hippi0` and `hippi1`. Assume the network interfaces are named:

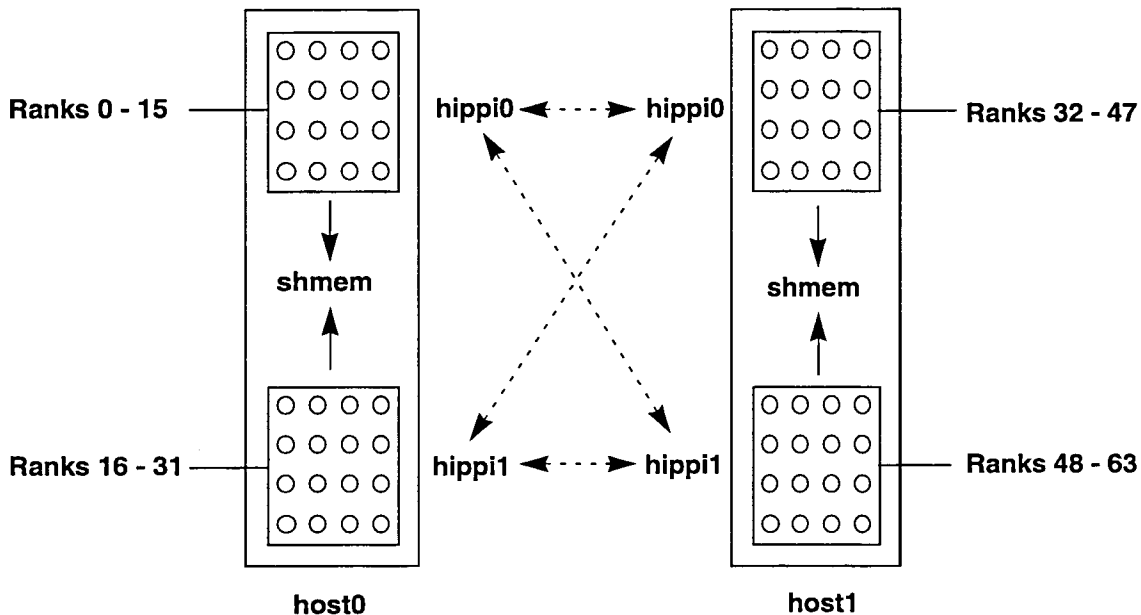
- `host0-hippi0`
- `host0-hippi1`
- `host1-hippi0`
- `host1-hippi1`

If your executable is called `beavis.exe` and uses 64 processes, your appfile should contain the following entries:

```
-h host0-hippi0 -e MPI_TOPOLOGY=/0:16,0 -np 16 beavis.exe  
-h host0-hippi1 -e MPI_TOPOLOGY=/1:0,16 -np 16 beavis.exe  
-h host1-hippi0 -e MPI_TOPOLOGY=/0:16,0 -np 16 beavis.exe  
-h host1-hippi1 -e MPI_TOPOLOGY=/1:0,16 -np 16 beavis.exe
```

Now, when the appfile is run, 32 processes are run on `host0` and 32 processes are run on `host1` as shown in Figure 5-1.

Figure 5-1 Multiple network interfaces



Host0 processes with rank 0 - 15 communicate with processes with rank 16 - 31 through shared memory (shmem). Host0 processes also communicate through the host0-hippi0 network interface with host1 processes.

Processor subscription

Subscription refers to the match of processors and active processes on a host or subcomplex. Table 5-1 lists the possible subscription types.

Table 5-1

Subscription types

Subscription type	Description
Under subscribed	More processors than active processes
Fully subscribed	Equal number of processors and active processes
Over subscribed	More active processes than processors

Tuning

General tuning

When a host or subcomplex is over subscribed, application performance decreases because of increased context switching.

Context switching can degrade application performance by slowing the computation phase, increasing message latency, and lowering message bandwidth. Simulations that use timing-sensitive algorithms can produce unexpected or erroneous results when run on an over-subscribed system.

MPI routine selection

To achieve the lowest message latencies and highest message bandwidths for point-to-point synchronous communications, use the MPI blocking routines `MPI_Send` and `MPI_Recv`. For asynchronous communications, use the MPI nonblocking routines `MPI_Isend` and `MPI_Irecv`.

When using the blocking routines, try to avoid pending requests. MPI must advance nonblocking messages, so calls to blocking receives must advance pending requests occasionally resulting in lower application performance.

For tasks that require collective operations, use the appropriate MPI collective routine. HP MPI takes advantage of shared memory to perform efficient data movement and maximize your application's communication performance.

SPP-UX platform tuning

There are two factors that affect application performance when working with HP MPI applications on the SPP-UX platform. These factors include:

- Multilevel parallelism
- Process placement

Multilevel parallelism

There are several ways to improve the performance of applications that use multilevel parallelism:

- Use the MPI library to provide coarse-grained parallelism and a parallelizing compiler to provide fine-grained (that is, thread-based) parallelism. An appropriate mix of coarse- and fine-grained parallelism provides better overall performance.
- Assign only one multithreaded process per hypernode when placing application processes. This ensures that enough processors are available as different process threads become active.

Process placement

Because messaging bandwidth and latency are better within a hypernode than between hypernodes, you can improve performance by placing HP MPI processes that communicate heavily on the same hypernode. One way to do this is to use the `MPI_TOPOLOGY` environment variable to tell an application the number of processes to run on each available hypernode.

For example, suppose you want to run an application on an X-Class server using a subcomplex called System. This subcomplex spans four hypernodes and contains the 20 processors listed below:

- Hypernode 0: 5 CPUs
- Hypernode 1: 2 CPUs
- Hypernode 2: 5 CPUs
- Hypernode 3: 8 CPUs

Tuning
SPP-UX platform tuning

Suppose the application you want to run contains the 16 processes listed below:

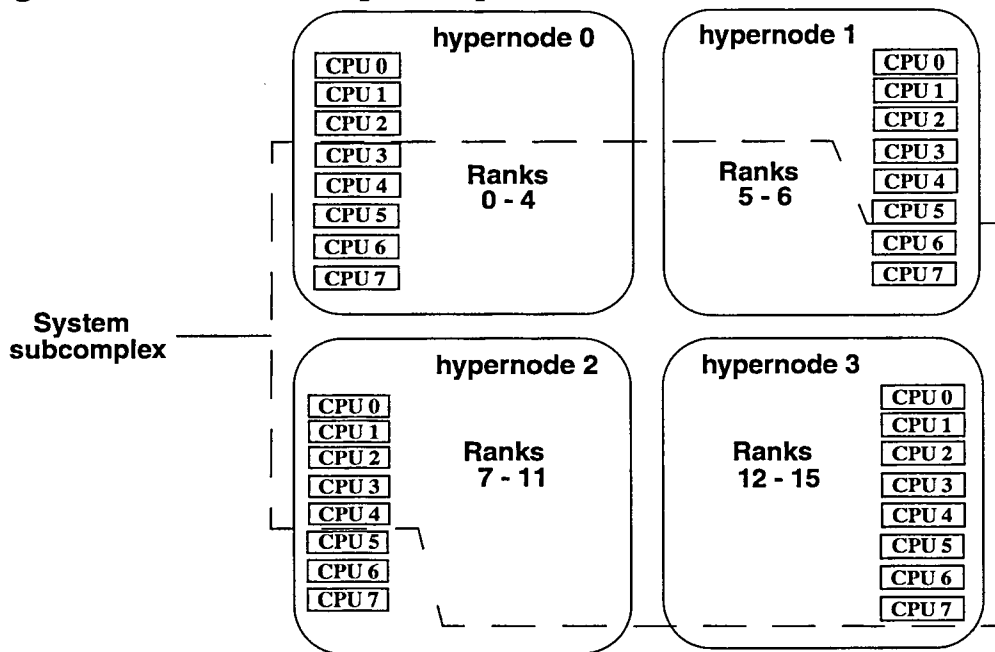
- Set 1: Ranks 0-3
- Set 2: Ranks 4-7
- Set 3: Ranks 8-11
- Set 4: Ranks 12-15

Ideally, you should use a process placement that allows each set of processes to run on a single hypernode to maximize message-passing performance.

By default, HP MPI places processes by fully subscribing each hypernode before moving on to the next. If the processes in your application are placed using this approach, you get the placement shown in Figure 5-2.

Figure 5-2

Default process placement



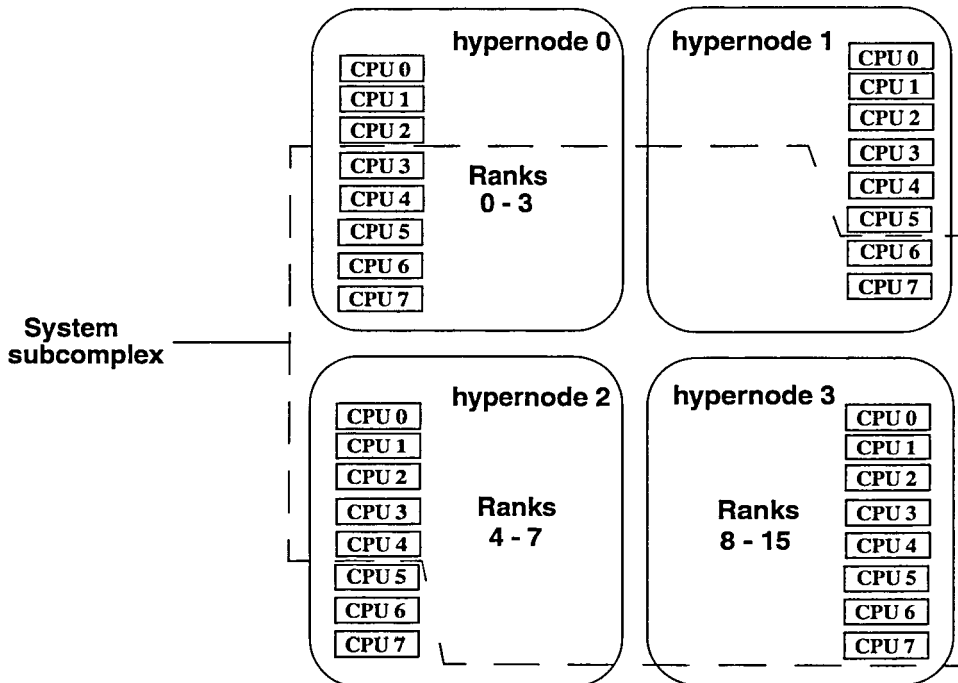
While this distribution prevents processor oversubscription, it does not provide optimum message-passing performance because the processes from sets two and three are split across hypernodes. Communications within these process groups may become a bottleneck when running the application.

You can solve this problem by specifying the number of processes you want to run on each hypernode as shown below:

- Hypernode 0 --> Ranks 0-3
- Hypernode 1 --> unused
- Hypernode 2 --> Ranks 4-7
- Hypernode 3 --> Ranks 8-15

This distribution results in a placement shown in Figure 5-3.

Figure 5-3 Optimal process placement



To specify this process placement, set `MPI_TOPOLOGY` by entering:

```
% setenv MPI_TOPOLOGY 4,0,4,8
```

Tuning
SPP-UX platform tuning

For more information, see “MPI_TOPOLOGY” on page 39.

NOTE

Make sure you use MPI_TOPOLOGY to place processes doing I/O on the hypernodes hosting the appropriate I/O controller. Placing these processes on noncontroller nodes results in lower I/O performance.

6 Debugging and troubleshooting

This chapter describes debugging and troubleshooting HP MPI applications. The topics covered are:

- Debugging HP MPI applications
- Troubleshooting HP MPI applications
- Frequently asked questions

Debugging HP MPI applications

You can debug your applications by setting options in the `MPI_FLAGS` environment variable or, if using `CXdb`, by invoking it directly.

NOTE

Setting options in `MPI_FLAGS` allows you to debug multihost applications and use other debuggers. `CXdb`, however, does not support multihost debugging.

Setting options in `MPI_FLAGS`

HP MPI provides single-process debuggers in `MPI_FLAGS` to debug applications running on SPP-UX and HP-UX platforms.

To use a single-process debugger:

Step 1. Set the `ecxdb`, `edde`, `exdb`, or `egdb` option in the `MPI_FLAGS` environment variable to use the `CXdb`, `DDE`, `XDB`, or `GDB` debugger respectively (refer to “`MPI_FLAGS`” on page 37 for information about these options).

On remote hosts, set `DISPLAY` to point to your console. In addition, use `xhost` to allow remote hosts to redirect their windows to your console.

Step 2. Run your application. When `MPI_Init` is executed, HP MPI starts one debugger session per process.

Step 3. Set a breakpoint anywhere following `MPI_Init` in each session.

Step 4. Set the global variable `MPI_DEBUG_CONT` to one using each session’s command-line interface or graphical user interface. For example, the command for `CXdb` is:

```
(CXdb) fill MPI_DEBUG_CONT \; 1
```

Step 5. Issue the appropriate debugger command in each session to continue program execution. Each process will run and stop at the breakpoint after `MPI_Init` that you set earlier.

Step 6. Debug the execution of each process using the appropriate commands for your debugger.

Using CXdb

CXdb offers a number of features that aid in the debugging of HP MPI applications:

- Consolidated debugging of all processes through one invocation of CXdb
- Automatic attaching of new processes as they are created by the application
- Ability to control execution of individual processes or sets of processes (using CXdb focus groups)

A single session of CXdb cannot debug an HP MPI application that spans multiple machines. However, you can use multiple sessions of CXdb (one for each machine) to debug an application. In this case, you can attach and detach processes to CXdb as needed.

See *Using CXdb* for more information.

Troubleshooting HP MPI applications

This section describes hints and limitations when working with HP MPI applications. You should check this information first when troubleshooting problems. The topics covered are organized by development task and include problems in areas such as:

- Building
- Starting
- Running
- Completing

Building

You can solve most build-time problems by referring to the documentation for the compiler you are using.

If you decide to use your own build script, be sure to specify all necessary input libraries. To determine what libraries are needed, check the contents of the compilation utilities stored in the HP MPI `/opt/mpi/bin` subdirectory.

Avoid using the `+autodbl` option when compiling Fortran 77 applications. This option may lead to unpredictable results.

Starting

When starting multihost applications, make sure that:

- All remote hosts are listed in your `.rhosts` file on each machine
- Application binaries are available on the necessary remote hosts and are executable on those machines
- The `-sp` option is passed to `mpirun` using an appfile if necessary
- The `.cshrc` file does not contain `tty` commands such as `stty` if you are using a `/bin/csh`-based shell.

Running

Run-time problems originate from many sources. These sources include:

- Propagation of environment variables
- Shared memory
- Interoperability
- Message buffering
- External input and output
- Fortran 90 programming features
- UNIX open file descriptors

Propagation of environment variables

When working with applications that run on multiple hosts, you must set values for environment variables on each host that participates in the job.

A recommended way to accomplish this is to set the `-e` option in the appfile:

```
-h remote_host -e MPI_TOPOLOGY=val [-np #] program [args]
```

Alternatively, you can set environment variables using the `.cshrc` file on each remote host if you are using a `/bin/csh`-based shell.

Shared memory

When an MPI application starts, each MPI process attempts to allocate a section of shared memory. This allocation can fail if the system-imposed limit on the maximum number of allowed shared-memory identifiers is exceeded or if the amount of available physical memory is not sufficient to fill the request.

After shared-memory allocation is done, every MPI process attempts to attach to the shared-memory region of every other process residing on the same host. This attachment can fail if the number of shared-memory segments attached to the calling process exceeds the system-imposed limit.

Furthermore, all processes must be able to attach to a shared-memory region at the same virtual address. For example, if the first process to attach to the segment attaches at address `ADR`, then the virtual-memory region starting at `ADR` and of the same size as the shared-memory region must be available to all other processes. Placing `MPI_Init` to execute first can help avoid this problem. A process with a large stack size is also prone to this failure. Choose process stack size carefully.

Interoperability

Depending upon what server resources are available, applications may run on heterogeneous systems such that certain portions run on SPP-UX platforms and other portions run on HP-UX platforms.

For example, suppose you create a MPMD application that calculates the average acceleration of particles in a simulated cyclotron. The application consists of a four-process program called `sum_accelerations` and an eight-process program called `calculate_average`.

Because you have access to a K-Class server called `hpux_server` and an X-Class server called `sppux_server`, you create the following appfile:

```
-h hpux_server -np 4 sum_accelerations  
-h sppux_server -np 8 calculate_average
```

Then, you invoke `mpirun` passing it the name of the appfile you created. In this case, even though the two application programs run on different platforms, all processes can communicate with each other resulting in twelve-way parallelism. The four processes belonging to the `sum_accelerations` application are ranked 0 through 3, and the eight processes belonging to the `calculate_average` application are ranked 4 through 11.

Message buffering

According to the MPI standard, message buffering may or may not occur when processes communicate with each other using `MPI_Send`. Therefore, you should take care when coding communications that depend upon buffering to work correctly.

For example, when two processes use `MPI_Send` to simultaneously send a message to each other and use `MPI_Recv` to receive the messages, the results are unpredictable. If the messages are buffered, communication works correctly. If the messages are not buffered, however, each process hangs in `MPI_Send` waiting for `MPI_Recv` to take the message.

External input and output

Each process in HP MPI applications can read and write data to an external drive. In some applications, however, having one process handle all input and output (and communicate with other processes using collective operations) is more efficient.

You can use `stdin` and `stdout` in your applications to read and write data. `Stdout` is supported regardless of whether your application runs locally or remotely. `Stdin`, however, may or may not be supported depending upon how you run your application, whether the application is run locally or remotely, and whether the `-W` option is used when invoking `mpirun`. The run invocations under which `stdin` is supported are shown in Table 6-1 for the `hello_world` application. All multihost invocations use an appfile called `hello_world`.

Table 6-1 **Run invocations that support `stdin`**

Run invocation	Is <code>stdin</code> supported?
<code>hello_world -np #</code>	Yes
<code>mpirun -np # hello_world</code>	Yes
<code>mpirun -W -np # hello_world</code>	No
<code>mpirun -W -np # -f hello_world</code> (multihost local and remote)	No
<code>mpirun -np # -f hello_world</code> (multihost local)	Yes
<code>mpirun -np # -f hello_world</code> (multihost remote)	No

Fortran 90 programming features

The MPI 1.1 standard defines bindings for Fortran 77 but not Fortran 90.

Although most Fortran 90 MPI applications work using the Fortran 77 MPI bindings, some Fortran 90 features can cause unexpected behavior when used with HP MPI.

In Fortran 90, an array is not always stored in contiguous memory. When noncontiguous array data are passed to an HP MPI subroutine, Fortran 90 copies the data into temporary storage, passes it to the HP MPI subroutine, and copies it back when the subroutine returns. As a result, HP MPI is given the address of the copy but not the original data.

In some cases, this copy-in and copy-out operation can cause a problem. For a nonblocking HP MPI call, the subroutine returns immediately and the temporary storage is deallocated. When HP MPI tries to access the already invalid memory, the behavior is unknown. Moreover, HP MPI operates close to the system level and needs to know the address of the original data (although even if the address is known, HP MPI does not know if the data are contiguous or not).

UNIX open file descriptors

UNIX imposes a limit to the number of file descriptors that application processes can have open at one time. When running a multihost application, each local process opens a socket to each remote process. An HP MPI application with a large amount of off-host processes can quickly reach the file descriptor limit. Ask your system administrator to increase the limit if your applications frequently exceed the maximum.

Completing

In HP MPI, `MPI_Finalize` is a barrier-like collective routine that waits until all application processes have called it before returning. If your application exits without calling `MPI_Finalize`, pending requests may not complete.

When running an application, `mpirun` normally waits until all processes have called `MPI_Finalize`. If an application detects an MPI error that leads to program termination, it calls `MPI_Abort` instead.

You may want to code your error conditions using `MPI_Abort`, which calls `MPI_Finalize`.

Each HP MPI application is identified by a job ID, unique on the server where `mpirun` is invoked. If you use the `-j` option, `mpirun` prints the job ID of the application that it runs. Then, you can invoke `mpijob` with the job ID to display the status of your application.

If your application hangs or terminates abnormally, you can use `mpiclean` to kill any lingering processes and shared-memory segments. In this context, you use the job ID from `mpirun` to specify the application to terminate.

Frequently asked questions

This section describes frequently asked HP MPI questions.

1. How do I find out what hypernode each MPI process is running on in applications running under SPP-UX?

ANSWER: Use the `node_num(3)` interface in the AIL library, the `pot` utility, or the `syspic` utility. The AIL library is searched automatically when you use the HP MPI compilation scripts `mpicc`, `mpif77`, `mpicc`, and `mpif90`.

2. My application times out before invoking `MPI_Init`. I get the message:

```
mpirun: cannot accept connection: Connection timed out
```

I am running on a single system, not a cluster. What might be causing this?

ANSWER: `mpirun` makes some assumptions about how long it will take before a process calls `MPI_Init`. If the process does not call `MPI_Init` in time, `mpirun` assumes there is an error.

In general, MPI programs should call `MPI_Init` before doing anything else. This helps ensure that all processes respond to `mpirun` in time, and that they will be able to attach shared memory at the same addresses.

3. When I build with HP MPI and then turn tracing on, the application takes a long time inside `MPI_Finalize`. This was not happening previously. What is causing this?

ANSWER: `MPI_Finalize` now consolidates the raw trace generated by each process into a single output file (with a `.tr` extension). Previously, you had to invoke `mpitrget` explicitly to consolidate raw traces. You can instruct HP MPI to not merge traces by specifying the `nf` option using the `MPI_XMPI` environment variable. For example:

```
% setenv MPI_XMPI prefix:nf
```

4. How does HP MPI clean up when something goes wrong?

ANSWER: HP MPI uses several mechanisms to clean up job files. Note that all processes in your application must call `MPI_Finalize`.

- a. When a correct HP MPI program (that is, one that calls `MPI_Finalize`) exits successfully, the root host deletes the job file.
- b. If `mpirun` was used, it deletes the job file when the application terminates, whether successfully or not.
- c. When an application calls `MPI_Abort` or when a process gets a signal that is trapped by `SIGINT`, `SIGTERM`, and so forth, any processes present on the root host delete the job file.
- d. If you use `mpijob -j` to get more information on a job, and the processes of that job have all exited, `mpijob` issues a warning that the job has completed, and deletes the job file.

5. The documentation for HP MPI says that the syntax for the `MPI_TOPOLOGY` environment variable is `[[[sc][hypernode]:][topology]`. Based on this, I would expect both of the following two commands to handle a program run on two nodes:

```
% setenv MPI_TOPOLOGY System/0:1,1
```

or

```
% setenv MPI_TOPOLOGY System/0:0,2
```

The first command works fine and gives me:

```
happy04 [89] mpirun -w -np 2 hello_world
Hello world! I'm 1 of 2 on happy04
Hello world! I'm 0 of 2 on happy04
```

But the second command does not work and gives me the following error:

```
mpirun -w -np 2 hello_world
hello_world: Pid 3261: MPI_Init: MPI_TOPOLOGY: Logical startup node 0
must have an assigned process in the topology process list [see MPI.1,
TROUBLESHOOTING]
hello_world: Pid 3261: MPI_Init: MPI_TOPOLOGY=System/0:0,2
hello_world: Pid 3261: MPI_Init: Aborting the application
```

Debugging and troubleshooting
Frequently asked questions

ANSWER: The problem is that the operating system can place processes on any node. In particular, it is biased against placing processes on node 0 because node 0 normally has a higher load than other nodes. So, the operating system spawns the root process of `hello_world` on node 1.

However, in the second command, `MPI_TOPOLOGY` forces both processes to start on virtual node 0, resulting in startup failure. You can avoid this problem if you place the root process on node 0 by entering:

```
% mpa -node 0 hello_world -np 2
```

or

```
% mpa -node 0 mpirun -np 2 hello_world
```

You can also set `MPI_TOPOLOGY` to `System/1:0,2` to use the default subcomplex. In this case, there is no need to use `mpa`, just `mpirun`. `mpirun` will migrate the application to node 1 as follows:

```
% mpirun -np 2 -e MPI_TOPOLOGY=System/1:0,2\  
hello_world
```

Because this command does not depend on `mpa`, it is portable between HP-UX and SPP-UX (`MPI_TOPOLOGY` is ignored on HP-UX, so setting it there has no effect).

6. My MPI application hangs at `MPI_Send`. Why?

ANSWER: Check if your code assumes buffering behavior for standard communication mode. Deadlock situations may occur when standard send operations are used.

7. My MPI processes need to run with `mpa`. How do I do that?

ANSWER: The answer depends on whether you use `mpirun` or the *executable* `-np # syntax` to run your application. If you use the *executable* `-np # syntax`, you can invoke your application with `mpa` as follows:

```
% mpa -DATA n -STACK m executable -np #
```

If you use `mpirun`, you can create a shell script with your `mpa` execution line and include that script in your `mpi` file. For example:

```
% mpirun -np 32 mpa -DATA n -STACK m executable
```

A

MPI library routines and extensions

This appendix describes the MPI library routines, HP MPI library extensions, MPI 1.2 extensions, and MPI 2.0 extensions.

The following tables are included in this appendix:

- Table A-1 on page 104 describes the C version of the library routines.
- Table A-2 on page 117 describes the Fortran version of the library routines.
- Table A-3 on page 132 describes the C version of the library extensions.
- Table A-4 on page 133 describes the Fortran version of the library extensions.
- Table A-5 on page 134 describes the MPI 1.2 extensions.
- Table A-6 on page 135 describes the MPI 2.0 extensions.

C version of MPI routines

Table A-1 **C version of MPI routines**

Routine	Description
int MPI_Abort(MPI_Comm comm, int errorcode)	Aborts all tasks in a communicator.
int MPI_Address(void* location, MPI_Aint *address)	Returns the address of a location.
int MPI_Allgather(void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf, int recvcnts, MPI_Datatype recvtype, MPI_Comm comm)	Sends the contents of each process's send buffer to all other processes.
int MPI_Allgatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf, int *recvcnts, int *displs, MPI_Datatype recvtype, MPI_Comm comm)	Sends a varying count of data from each process's send buffer to all other processes.
int MPI_Allreduce(void* sendbuf, void* recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)	Combines the elements in the input buffer of each process and returns the combined value to the receive buffer of each process.
int MPI_Alltoall(void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf, int recvcount, MPI_Datatype recvtype, MPI_Comm comm)	Sends distinct data from each process to all other processes.
int MPI_Alltoallv(void* sendbuf, int *sendcounts, int *sdispls, MPI_Datatype sendtype, void* recvbuf, int *recvcnts, int *rdispls, MPI_Datatype recvtype, MPI_Comm comm)	Sends distinct data from each process to all other processes. The location of data for the send and the location of the placement of the data on the receive side are specified in the call.
int MPI_Attr_delete(MPI_Comm comm, int keyval)	Deletes attributes from the cache based upon a key.

Routine	Description
int MPI_Attr_get(MPI_Comm comm, int keyval, void *attribute_val, int *flag)	Receives attribute values based upon a key.
int MPI_Attr_put(MPI_Comm comm, int keyval, void* attribute_val)	Stores an attribute value.
int MPI_Barrier(MPI_Comm comm)	Blocks the calling process until all other processes in the group have called the routine.
int MPI_Bcast(void* buffer, int count, MPI_Datatype datatype, int root, MPI_Comm comm)	Broadcasts a message from the root process to all other processes in the group including itself.
int MPI_Bsend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Sends a message in buffered mode.
int MPI_Bsend_init(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_request *request)	Creates a persistent communication request for a buffered mode send.
int MPI_Buffer_attach(void* buffer, int size)	Provides MPI with a buffer in the user's memory that can be used for buffering outgoing messages.
int MPI_Buffer_detach(void* buffer, int* size)	Detaches the buffer currently associated with MPI.
int MPI_Cancel(MPI_Request *request)	Marks a pending, nonblocking send or receive call for cancellation.
int MPI_Cart_coords(MPI_Comm comm, int rank, int maxdims, int *coords)	Translates process ranks to logical process coordinates as the ranks are used in point-to-point communications.

MPI library routines and extensions
C version of MPI routines

Routine	Description
int MPI_Cart_create(MPI_Comm comm_old, int ndims, int *dims, int *periods, int reorder, MPI_Comm *comm_cart)	Returns a handle to a new communicator to which cartesian topology information is attached.
int MPI_Cart_get(MPI_Comm comm, int maxdims, int *dims, int *periods, int *coords)	Returns the cartesian topology information associated with a communicator.
int MPI_Cart_map(MPI_Comm comm, int ndims, int *dims, int *periods, int *newrank)	Computes an optimal placement for the calling process on the physical machine.
int MPI_Cart_rank(MPI_Comm comm, int *coords, int *rank)	Translates logical process coordinates to process ranks as the coordinates are used in point-to-point communications.
int MPI_Cart_shift(MPI_Comm comm, int direction, int disp, int *rank_source, int *rank_dest)	Returns the shifted source and destination ranks given a shift amount and direction.
int MPI_Cart_sub(MPI_Comm comm, int *remain_dims, MPI_Comm *newcomm)	Partitions the communicator group into subgroups that form lower-dimensional cartesian subgrids.
int MPI_Cartdim_get(MPI_Comm comm, int *ndims)	Returns the cartesian topology information associated with a communicator.
int MPI_Comm_compare(MPI_Comm comm1, MPI_Comm comm2, int *result)	Compares two communicators.
int MPI_Comm_create(MPI_Comm comm, MPI_Group group, MPI_Comm *newcomm)	Creates a new communicator and context with a communication group.

Routine	Description
int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)	Duplicates the existing communicator with associated key values.
int MPI_Comm_free(MPI_Comm *comm)	Marks the communication object for deallocation.
int MPI_Comm_group(MPI_Comm comm, MPI_Group *group)	Returns a handle to the group in the communicator.
int MPI_Comm_rank(MPI_Comm comm, int *rank)	Returns the rank of a process in the communicator's group.
int MPI_Comm_remote_group(MPI_Comm comm, MPI_Group *group)	Returns the remote group in the intercommunicator.
int MPI_Comm_remote_size(MPI_Comm comm, int *size)	Returns the size of the remote group in the intercommunicator.
int MPI_Comm_size(MPI_Comm comm, int *size)	Returns the number of processes involved in a communicator.
int MPI_Comm_split(MPI_Comm comm, int color, int key, MPI_Comm *newcomm)	Partitions the group associated with the communicator into disjoint subgroups.
int MPI_Comm_test_inter(MPI_Comm comm, int *flag)	Determines whether a communicator is an intercommunicator or not.
int MPI_Dims_create(int nnodes, int ndims, int *dims)	Helps select a balanced distribution of processes per coordinate direction.
int MPI_Errhandler_create(MPI_Handler_function *function, MPI_Errhandler *errhandler)	Registers a user-specified routine for use as an exception handler.
int MPI_Errhandler_free(MPI_Errhandler *errhandler)	Marks an error handler for deallocation.

MPI library routines and extensions
C version of MPI routines

Routine	Description
int MPI_Errhandler_get(MPI_Comm comm, MPI_Errhandler *errhandler)	Returns a handle to the error handler that is currently associated with the communicator.
int MPI_Errhandler_set(MPI_Comm comm, MPI_Errhandler errhandler)	Associates a new error handler with the communicator at the calling process.
int MPI_Error_class(int errorcode, int *errorclass)	Maps each standard error code onto itself.
int MPI_Error_string(int errorcode, char *string, int *resultlen)	Returns the error string associated with an error code.
int MPI_Finalize(void)	Cleans up all MPI states.
int MPI_Gather(void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)	Sends the contents of each process's send buffer to the root process.
int MPI_Gatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf, int *recvcounts, int *displs, MPI_Datatype recvtype, int root, MPI_Comm comm)	Sends a varying count of data from each process's send buffer to the root process.
int MPI_Get_count(MPI_Status *status, MPI_Datatype datatype, int *count)	Returns the number of entries received in the receive buffer.
int MPI_Get_elements(MPI_Status *status, MPI_Datatype datatype, int *elements)	Returns the number of elements in a data type.
int MPI_Get_processor_name(char *name, int len)	Returns the name of the processor on which it was called at the moment of the call.
int MPI_Graph_create(MPI_Comm comm_old, int nnodes, int *index, int *edges, int reorder, MPI_Comm *comm_graph)	Returns a handle to a new communicator to which the graph topology information is attached.

Routine	Description
int MPI_Graph_get(MPI_Comm comm, int maxindex, int maxedges, int *index, int *edges)	Retrieves graph topology information associated with a communicator.
int MPI_Graph_map(MPI_Comm comm, int nnodes, int *index, int *edges, int *newrank)	Maps process to graph topology information.
int MPI_Graph_neighbors(MPI_Comm comm, int rank, int maxneighbors, int *neighbors)	Returns the neighbors of a node associated with a graph topology.
int MPI_Graph_neighbors_count(MPI_Comm comm, int rank, int *nneighbors)	Returns the number of neighbors of a node associated with a graph topology.
int MPI_Graphdims_get(MPI_Comm comm, int *nnodes, int *nedges)	Retrieves graph topology information associated with a communicator.
int MPI_Group_compare(MPI_Group group1, MPI_Group group2, int *result)	Returns the relationship between two groups.
int MPI_Group_difference(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)	Creates a group by computing the difference between two existing groups.
int MPI_Group_excl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup)	Creates a group by reordering an existing group and only taking unlisted members.
int MPI_Group_free(MPI_Group *group)	Marks a group object for deallocation.
int MPI_Group_incl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup)	Creates a group by reordering an existing group and only taking listed members.
int MPI_Group_intersection(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)	Creates a group by computing the intersection of two existing groups.

MPI library routines and extensions
C version of MPI routines

Routine	Description
int MPI_Group_range_excl(MPI_Group group, int n, int ranges [] [3], MPI_Group *newgroup)	Creates a group by excluding ranges of processes from an existing group.
int MPI_Group_range_incl(MPI_Group group, int n, int ranges [] [3], MPI_Group *newgroup)	Creates a group from ranges of ranks in an existing group.
int MPI_Group_rank(MPI_Group group, int *rank)	Returns the rank of this process in a group.
int MPI_Group_size(MPI_Group group, int *size)	Returns the size or number of processes in the group.
int MPI_Group_translate_ranks(MPI_Group group1, int n, int *ranks1, MPI_Group group2, int *ranks2)	Determines the relative numbering of the same processes in two different groups.
int MPI_Group_union(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)	Creates a group by computing the union of two existing groups.
int MPI_Ibsend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Starts a buffered mode, nonblocking send.
int MPI_Init(int *argc, char ***argv)	Initializes the MPI execution environment.
int MPI_Initialized(int *flag)	Determines whether MPI_Init has been called.
int MPI_Intercomm_create(MPI_Comm local_comm, int local_leader, MPI_Comm peer_comm, int remote_leader, int tag, MPI_Comm *newintercomm)	Creates an intercommunicator.
int MPI_Intercomm_merge(MPI_Comm intercomm, int high, MPI_Comm *newintracomm)	Creates an intracommunicator from the union of two groups.
int MPI_Iprobe(int source, int tag, MPI_Comm comm, int *flag, MPI_Status *status)	Returns flag = true if there is a message that can be received and that matches the pattern specified by the arguments.

Routine	Description
int MPI_Irecv(void* buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Request *request)	Starts a nonblocking receive.
int MPI_Irsend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Starts a ready mode, nonblocking send.
int MPI_Isend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Starts a standard mode, nonblocking send.
int MPI_Issend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Starts a synchronous mode, nonblocking send.
int MPI_Keyval_create(MPI_Copy_function *copy_fn, MPI_Delete_function *delete_fn, int *keyval, void* extra_state)	Generates a new attribute key.
int MPI_Keyval_free(int *keyval)	Frees an existing attribute key.
int MPI_Op_create(MPI_User_function function, int commute, MPI_Op *op)	Binds a user-defined global operation to a handle.
int MPI_Op_free(MPI_Op *op)	Marks a user-defined reduction operation for deallocation.
int MPI_Pack(void* inbuf, int incount, MPI_Datatype datatype, void *outbuf, int outsize, int *position, MPI_Comm comm)	Packs the message in the send buffer into the specified buffer space.
int MPI_Pack_size(int incount, MPI_Datatype datatype, MPI_Comm comm, int *size)	Returns the maximum number of buffer bytes required to hold the data.
int MPI_Pcontrol(const int level, ...)	Calls the specified profiling package.

MPI library routines and extensions
C version of MPI routines

Routine	Description
int MPI_Probe(int source, int tag, MPI_Comm comm, MPI_Status *status)	Returns only after a message is found that matches the pattern specified by the arguments.
int MPI_Recv(void* buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status *status)	Starts a blocking receive.
int MPI_Recv_init(void* buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Request *request)	Creates a persistent communication request for a receive operation.
int MPI_Reduce(void* sendbuf, void* recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)	Combines the elements in the input buffer of each process in the group and returns the combined value in the output buffer of the root process.
int MPI_Reduce_scatter(void* sendbuf, void* recvbuf, int *recvcounts, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)	Combines values and scatters the results.
int MPI_Request_free(MPI_Request *request)	Marks the request object for deallocation.
int MPI_Rsend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Sends a message in ready mode.
int MPI_Rsend_init(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Creates a persistent communication object for a ready mode send operation.
int MPI_Scan(void* sendbuf, void* recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)	Performs a prefix reduction on data distributed across the group.
int MPI_Scatter(void* sendbuf, int sendcount, MPI_Datatype sendtype, void* recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)	Sends the contents of the root process's send buffer to other processes in the group.

Routine	Description
int MPI_Scatterv(void* sendbuf, int *sendcounts, int *displs, MPI_Datatype sendtype, void* recvbuf, int recvcnt, MPI_Datatype recvtpe, int root, MPI_Comm comm)	Sends a varying count of data from the root process's send buffer to other processes in the group.
int MPI_Send(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Starts a standard mode, blocking send.
int MPI_Send_init(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Creates a persistent communication request for a standard mode send operation and binds it to all the arguments of a send operation.
int MPI_Sendrecv(void *sendbuf, int sendcount, MPI_Datatype sendtype, int dest, int sendtag, void *recvbuf, int recvcnt, MPI_Datatype recvtpe, int source, int revctag, MPI_Comm comm, MPI_Status *status)	Executes a blocking send and receive.
int MPI_Sendrecv_replace(void* buf, int count, MPI_Datatype datatype, int dest, int sendtag, int source, int revctag, MPI_Comm comm, MPI_Status *status)	Executes a blocking send and receive where both operations use the same buffer.
int MPI_Ssend(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Sends a message in synchronous mode.
int MPI_Ssend_init(void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Creates a persistent communication object for a synchronous mode send operation.
int MPI_Start(MPI_Request *request)	Initiates a communication with a persistent request handle.
int MPI_Startall(int count, MPI_Request *array_of_requests)	Initiates a collection of requests.

MPI library routines and extensions
C version of MPI routines

Routine	Description
int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status)	Returns flag = true if the request (operation) identified in the call is complete.
int MPI_Test_cancelled(MPI_Status *status, int *flag)	Marks a pending, nonblocking communication operation (send or receive) for cancellation.
int MPI_Testall(int count, MPI_Request *array_of_requests, int *flag, MPI_Status *array_of_statuses)	Returns flag = true if all communications associated with active handles in the array have completed.
int MPI_Testany(int count, MPI_Request *array_of_requests, int *index, int *flag, MPI_Status *status)	Tests for completion of either one or none of the operations associated with active handles.
int MPI_Testsome(int incount, MPI_Request *array_of_requests, int *outcount, int *array_of_indices, MPI_Status *array_of_statuses)	Tests for some given communications to complete.
int MPI_Topo_test(MPI_Comm comm, int *status)	Returns the type of topology that is assigned to a communicator.
int MPI_Type_commit(MPI_Datatype *datatype)	Commits the data type (that is, the formal description of a communication buffer), not the content of that buffer.
int MPI_Type_contiguous(int count, MPI_Datatype oldtype, MPI_Datatype *newtype)	Creates a contiguous data type.
int MPI_Type_extent(MPI_Datatype datatype, MPI_Aint *extent)	Returns the extent of a data type.
int MPI_Type_free(MPI_Datatype *datatype)	Marks the data-type object identified in the call for deallocation.

Routine	Description
<code>int MPI_Type_hindexed(int count, int *array_of_blocklengths, MPI_Aint *array_of_displacements, MPI_Datatype oldtype, MPI_Datatype *newtype)</code>	Creates an indexed data type with offsets in bytes.
<code>int MPI_Type_hvector(int count, int blocklength, MPI_Aint stride, MPI_Datatype oldtype, MPI_datatype *newtype)</code>	Creates a vector (strided) data type with offset in bytes.
<code>int MPI_Type_indexed(int count, int *array_of_blocklengths, int *array_of_displacements, MPI_Datatype oldtype, MPI_Datatype *newtype)</code>	Creates an indexed data type.
<code>int MPI_Type_lb(MPI_Datatype datatype, MPI_Aint *displacement)</code>	Returns the lower bound of a data type.
<code>int MPI_Type_size(MPI_Datatype datatype, int *size)</code>	Returns the total size (in bytes) of the entries in the data type identified in the call.
<code>int MPI_Type_struct(int count, int *array_of_blocklengths, MPI_Aint *array_of_displacements, MPI_Datatype *array_of_types, MPI_Datatype *newtype)</code>	Creates a structured data type.
<code>int MPI_Type_ub(MPI_datatype datatype, MPI_Aint *displacement)</code>	Returns the upper bound of a data type.
<code>int MPI_Type_vector(int count, int blocklength, int stride, MPI_datatype oldtype, MPI_datatype *newtype)</code>	Creates a vector (strided) data type.
<code>int MPI_Unpack(void* inbuf, int insize, int *position, void *outbuf, int outcount, MPI_Datatype datatype, MPI_Comm comm)</code>	Unpacks a message into the receive buffer.
<code>int MPI_Wait(MPI_Request *request, MPI_Status *status)</code>	Waits for an MPI send or receive to complete.

MPI library routines and extensions
C version of MPI routines

Routine	Description
int MPI_Waitall(int count, MPI_Request *array_of_requests, MPI_Status *array_of_statuses)	Blocks until all communication operations associated with active handles in the list complete and return their status.
int MPI_Waitany(int count, MPI_Request *array_of_requests, int *index, MPI_Status *status)	Blocks until one of the operations associated with the active requests in the array has completed.
int MPI_Waitsome(int count, MPI_Request *array_of_requests, int *outcount, int *array_of_indices, MPI_Status *array_of_statuses)	Waits until at least one of the operations associated with active handles in the list have completed.
double MPI_Wtick(void)	Returns the resolution of the MPI_Wtime routine in seconds.
double MPI_Wtime(void)	Returns a floating-point number of seconds representing wall-clock time.

Fortran version of MPI routines

Table A-2 Fortran version of MPI routines

Routine	Description
MPI_Abort (comm, errorcode, ierror) integer comm, errorcode, ierror	Aborts all tasks in a communicator.
MPI_Address(location, address, ierror) <type> location(*) integer address, ierror	Returns the address of a location.
MPI_Allgather(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvttype, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, recvcount, recvttype, comm, ierror	Sends the contents of each process's send buffer to all other processes.
MPI_Allgatherv(sendbuf, sendcount, sendtype, recvbuf, recvcounst, displs, recvttype, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, recvcounst(*), displs(*), recvttype, comm, ierror	Sends a varying count of data from each process's send buffer to all other processes.
MPI_Allreduce(sendbuf, recvbuf, count, datatype, op, comm, ierror) <type> sendbuf(*), recvbuf(*) integer count, datatype, op, comm, ierror	Combines the elements in the input buffer of each process and returns the combined value to the receive buffer of each process.
MPI_Alltoall(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvttype, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, recvcount, recvttype, comm, ierror	Sends distinct data from each process to all other processes.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
<p>MPI_Alltoallv(sendbuf, sendcounts, sdispls, sendtype, recvbuf, recvcounts, rdispls, recvtype, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcounts(*), sdispls(*), sendtype, recvcounts(*), rdispls(*), recvtype, comm, ierror</p>	<p>Sends distinct data from each process to all other processes. The location of data for the send and the location of the placement of the data on the receive side are specified in the call.</p>
<p>MPI_Attr_delete(comm, keyval, ierror) integer comm, keyval, ierror</p>	<p>Deletes attributes from the cache based upon a key.</p>
<p>MPI_Attr_get(comm, keyval, attribute_val, flag, ierror) integer comm, keyval, attribute_val, ierror logical flag</p>	<p>Receives attribute values based upon a key.</p>
<p>MPI_Attr_put(comm, keyval, attribute_val, ierror) integer comm, keyval, attribute_val, ierror</p>	<p>Stores an attribute value.</p>
<p>MPI_Barrier(comm, ierror) integer comm, ierror</p>	<p>Blocks the calling process until all other processes in the group have called the routine.</p>
<p>MPI_Bcast(buffer, count, datatype, root, comm, ierror) <type> buffer(*) integer count, datatype, root, comm, ierror</p>	<p>Broadcasts a message from the root process to all other processes in the group including itself.</p>
<p>MPI_Bsend(buf, count, datatype, dest, tag, comm, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, ierror</p>	<p>Sends a message in buffered mode.</p>
<p>MPI_Bsend_init(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror</p>	<p>Creates a persistent communication request for a buffered mode send.</p>
<p>MPI_Buffer_attach(buffer, size, ierror) <type> buffer(*) integer size, ierror</p>	<p>Provides MPI with a buffer in the user's memory that can be used for buffering outgoing messages.</p>

Routine	Description
MPI_Buffer_detach(buffer, size, ierror) <type> buffer(*) integer size, ierror	Detaches the buffer currently associated with MPI.
MPI_Cancel(request, ierror) integer request, ierror	Marks a pending, nonblocking send or receive call for cancellation.
MPI_Cart_coords(comm, rank, maxdims, coords, ierror) integer comm, rank, maxdims, coords(*), ierror	Translates process ranks to logical process coordinates as the ranks are used in point-to-point communications.
MPI_Cart_create(comm_old, ndims, dims, periods, reorder, comm_cart, ierror) integer comm_old, ndims, dims(*), comm_cart, ierror logical periods(*), reorder	Returns a handle to a new communicator to which cartesian topology information is attached.
MPI_Cart_get(comm, maxdims, dims, periods, coords, ierror) integer comm, maxdims, dims(*), coords(*), ierror logical periods(*)	Returns the cartesian topology information associated with a communicator.
MPI_Cart_map(comm, ndims, dims, periods, newrank, ierror) integer comm, ndims, dims(*), newrank, ierror logical periods(*)	Computes an optimal placement for the calling process on the physical machine.
MPI_Cart_rank(comm, coords, rank, ierror) integer comm, coords(*), rank, ierror	Translates logical process coordinates to process ranks as the coordinates are used in point-to-point communications.
MPI_Cart_shift(comm, direction, disp, rank_source, rank_dest, ierror) integer comm, direction, disp, rank_source, rank_dest, ierror	Returns the shifted source and destination ranks given a shift amount and direction.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Cart_sub(comm, remain_dims, newcomm, ierror) integer comm, newcomm, ierror logical remain_dims(*)	Partitions the communicator group into subgroups that form lower-dimensional cartesian subgrids.
MPI_Cartdim_get(comm, ndims, ierror) integer comm, ndims, ierror	Returns the cartesian topology information associated with a communicator.
MPI_Comm_compare(comm1, comm2, result, ierror) integer comm1, comm2, result, ierror	Compares two communicators.
MPI_Comm_create(comm, group, newcomm, ierror) integer comm, group, newcomm, ierror	Creates a new communicator and context with a communication group.
MPI_Comm_dup(comm, newcomm, ierror) integer comm, newcomm, ierror	Duplicates the existing communicator with associated key values.
MPI_Comm_free(comm, ierror) integer comm, ierror	Marks the communication object for deallocation.
MPI_Comm_group(comm, group, ierror) integer comm, group, ierror	Returns a handle to the group in the communicator.
MPI_Comm_rank(comm, rank, ierror) integer comm, rank, ierror	Returns the rank of a process in the communicator's group.
MPI_Comm_remote_group(comm, group, ierror) integer comm, group, ierror	Returns the remote group in the intercommunicator.
MPI_Comm_remote_size(comm, size, ierror) integer comm, size, ierror	Returns the size of the remote group in the intercommunicator.
MPI_Comm_size(comm, size, ierror) integer comm, size, ierror	Returns the number of processes involved in a communicator.

Routine	Description
MPI_Comm_split(comm, color, key, newcomm, ierror) integer comm, color, key, newcomm, ierror	Partitions the group associated with the communicator into disjoint subgroups.
MPI_Comm_test_inter(comm, flag, ierror) integer comm, ierror logical flag	Determines whether a communicator is an intercommunicator or not.
MPI_Dims_create(nnodes, ndims, dims, ierror) integer nnodes, ndims, dims(*), ierror	Helps select a balanced distribution of processes per coordinate direction.
MPI_Errhandler_create(function, errhandler, ierror) external function integer errhandler, ierror	Registers a user-specified routine for use as an exception handler.
MPI_Errhandler_free(errhandler, ierror) integer errhandler, ierror	Marks an error handler for deallocation.
MPI_Errhandler_get(comm, errhandler, ierror) integer comm, errhandler, ierror	Returns a handle to the error handler that is currently associated with the communicator.
MPI_Errhandler_set(comm, errhandler, ierror) integer comm, errhandler, ierror	Associates a new error handler with the communicator at the calling process.
MPI_Error_class(errorcode, errorclass, ierror) integer errorcode, errorclass, ierror	Maps each standard error code onto itself.
MPI_Error_string(errorcode, string, resultlen, ierror) integer errorcode, resultlen, ierror character*MPI_MAX_ERROR_STRING string	Returns the error string associated with an error code.
MPI_Finalize(ierror) integer ierror	Cleans up all MPI states.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Gather(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, recvcount, recvtype, root, comm, ierror	Sends the contents of each process's send buffer to the root process.
MPI_Gatherv(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtype, root, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, recvcounts(*), displs(*), recvtype, root, comm, ierror	Sends a varying count of data from each process's send buffer to the root process.
MPI_Get_count(status, datatype, count, ierror) integer status(*), datatype, count, ierror	Returns the number of entries received in the receive buffer.
MPI_Get_elements(status, datatype, elements, ierror) integer status(*), datatype, elements, ierror	Returns the number of elements in a data type.
MPI_Get_processor_name(name, resultlen, ierror) character*MPI_MAX_PROCESSOR_NAME name integer resultlen, ierror	Returns the name of the processor on which it was called at the moment of the call.
MPI_Graph_create(comm_old, nnodes, index, edges, reorder, comm_graph, ierror) integer comm_old, nnodes, index(*), edges(*), comm_graph, ierror logical reorder	Returns a handle to a new communicator to which the graph topology information is attached.
MPI_Graph_get(comm, maxindex, maxedges, index, edges, ierror) integer comm, maxindex, maxedges, index(*), edges(*), ierror	Retrieves graph topology information associated with a communicator.
MPI_Graph_map(comm, nnodes, index, edges, newrank, ierror) integer comm, nnodes, index(*), edges(*), newrank, ierror	Maps process to graph topology information.
MPI_Graph_neighbors(comm, rank, maxneighbors, neighbors, ierror) integer comm, rank, maxneighbors, neighbors(*), ierror	Returns the neighbors of a node associated with a graph topology.

Routine	Description
MPI_Graph_neighbors_count(comm, rank, nneighbors, ierror) integer comm, rank, nneighbors, ierror	Returns the number of neighbors of a node associated with a graph topology.
MPI_Graphdims_get(comm, nnodes, nedges, ierror) integer comm, nnodes, nedges, ierror	Retrieves graph topology information associated with a communicator.
MPI_Group_compare(group1, group2, result, ierror) integer group1, group2, result, ierror	Returns the relationship between two groups.
MPI_Group_difference(group1, group2, newgroup, ierror) integer group1, group2, newgroup, ierror	Creates a group by computing the difference between two existing groups.
MPI_Group_excl(group, n, ranks, newgroup, ierror) integer group, n, ranks(*), newgroup, ierror	Creates a group by reordering an existing group and only taking unlisted members.
MPI_Group_free(group, ierror) integer group, ierror	Marks a group object for deallocation.
MPI_Group_incl(group, n, ranks, newgroup, ierror) integer group, n, ranks(*), newgroup, ierror	Creates a group by reordering an existing group and only taking listed members.
MPI_Group_intersection(group1, group2, newgroup, ierror) integer group1, group2, newgroup, ierror	Creates a group by computing the intersection of two existing groups.
MPI_Group_range_excl(group, n, ranges, newgroup, ierror) integer group, n, ranges(3, *), newgroup, ierror	Creates a group by excluding ranges of processes from an existing group.
MPI_Group_range_incl(group, n, ranges, newgroup, ierror) integer group, n, ranges(3, *), newgroup, ierror	Creates a group from ranges of ranks in an existing group.
MPI_Group_rank(group, rank, ierror) integer group, rank, ierror	Returns the rank of this process in a group.
MPI_Group_size(group, size, ierror) integer group, size, ierror	Returns the size or number of processes in the group.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Group_translate_ranks(group1, n, ranks1, group2, ranks2, ierror) integer group1, n, ranks1(*), group2, ranks2(*), ierror	Determines the relative numbering of the same processes in two different groups.
MPI_Group_union(group1, group2, newgroup, ierror) integer group1, group2, newgroup, ierror	Creates a group by computing the union of two existing groups.
MPI_Ibsend(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror	Starts a buffered mode, nonblocking send.
MPI_Init(ierror) integer ierror	Initializes the MPI execution environment.
MPI_Initialized(flag, ierror) logical flag integer ierror	Determines whether MPI_Init has been called.
MPI_Intercomm_create(local_comm, local_leader, peer_comm, remote_leader, tag, newintercomm, ierror) integer local_comm, local_leader, peer_comm, remote_leader, tag, newintercomm, ierror	Creates an intercommunicator.
MPI_Intercomm_merge(intercomm, high, newintracomm, ierror) logical high integer intercomm, newintracomm, ierror	Creates an intracommunicator from the union of two groups.
MPI_Iprobe(source, tag, comm, flag, status, ierror) logical flag integer source, tag, comm, status (*), ierror	Returns flag = true if there is a message that can be received and that matches the pattern specified by the arguments.
MPI_Irecv(buf, count, datatype, source, tag, comm, request, ierror) <type> buf(*) integer count, datatype, source, tag, comm, request, ierror	Starts a nonblocking receive.

Routine	Description
MPI_Irsend(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror	Starts a ready mode, nonblocking send.
MPI_Isend(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf (*) integer count, datatype, dest, tag, comm, request, ierror	Starts a standard mode, nonblocking send.
MPI_Issend(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror	Starts a synchronous mode, nonblocking send.
MPI_Keyval_create(copy_fn, delete_fn, keyval, extra_state, ierror) external copy_fn, delete_fn integer keyval, extra_state, ierror	Generates a new attribute key.
MPI_Keyval_free(keyval, ierror) integer keyval, ierror	Frees an existing attribute key.
MPI_Op_create(function, commute, op, ierror) external function logical commute integer op, ierror	Binds a user-defined global operation to a handle.
MPI_Op_free(op, ierror) integer op, ierror	Marks a user-defined reduction operation for deallocation.
MPI_Pack(inbuf, incount, datatype, outbuf, outsize, position, comm, ierror) <type> inbuf(*), outbuf(*) integer incount, datatype, outsize, position, comm, ierror	Packs the message in the send buffer into the specified buffer space.
MPI_Pack_size(incount, datatype, comm, size, ierror) integer incount, datatype, comm, size, ierror	Returns the maximum number of buffer bytes required to hold the data.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Pcontrol(level) integer level	Calls the specified profiling package.
MPI_Probe(source, tag, comm, status, ierror) integer source, tag, comm, status(*), ierror	Returns only after a message is found that matches the pattern specified by the arguments.
MPI_Recv(buf, count, datatype, source, tag, comm, status, ierror) <type> buf(*) integer count, datatype, source, tag, comm, status(*), ierror	Starts a blocking receive.
MPI_Recv_init(buf, count, datatype, source, tag, comm, request, ierror) <type> buf(*) integer count, datatype, source, tag, comm, request, ierror	Creates a persistent communication request for a receive operation.
MPI_Reduce(sendbuf, recvbuf, count, datatype, op, root, comm, ierror) <type> sendbuf(*), recvbuf(*) integer count, datatype, op, root, comm, ierror	Combines the elements in the input buffer of each process in the group and returns the combined value in the output buffer of the root process.
MPI_Reduce_scatter(sendbuf, recvbuf, recvcnts, datatype, op, comm, ierror) <type> sendbuf(*), recvbuf(*) integer recvcnts(*), datatype, op, comm, ierror	Combines values and scatters the results.
MPI_Request_free(request, ierror) integer request, ierror	Marks the request object for deallocation.
MPI_Rsend(buf, count, datatype, dest, tag, comm, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, ierror	Sends a message in ready mode.
MPI_Rsend_init(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror	Creates a persistent communication object for a ready mode send operation.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Scan (sendbuf, recvbuf, count, datatype, op, comm, ierror) <type> sendbuf(*), recvbuf(*) integer count, datatype, op, comm, ierror	Performs a prefix reduction on data distributed across the group.
MPI_Scatter (sendbuf, sendcount, sendtype, recvbuf, recvcnt, recvtype, root, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, recvcnt, recvtype, root, comm, ierror	Sends the contents of the root process's send buffer to other processes in the group.
MPI_Scatterv (sendbuf, sendcounts, displs, sendtype, recvbuf, recvcnt, recvtype, root, comm, ierror) <type> sendbuf(*), recvbuf(*) integer sendcounts(*), displs(*), sendtype, recvcnt, recvtype, root, comm, ierror	Sends a varying count of data from the root process's send buffer to other processes in the group.
MPI_Send (buf, count, datatype, dest, tag, comm, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, ierror	Starts a standard mode, blocking send.
MPI_Send_init (buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror	Creates a persistent communication request for a standard mode send operation and binds it to all the arguments of a send operation.
MPI_Sendrecv (sendbuf, sendcount, sendtype, dest, sendtag, recvbuf, recvcnt, recvtype, source, recvtag, comm, status, ierror) <type> sendbuf(*), recvbuf(*) integer sendcount, sendtype, dest, sendtag, recvcnt, recvtype, source, recvtag, comm, status(*), ierror	Executes a blocking send and receive.
MPI_Sendrecv_replace (buf, count, datatype, dest, sendtag, source, recvtag, comm, status, ierror) <type> buf(*) integer count, datatype, dest, sendtag, source, recvtag, comm, status(*), ierror	Executes a blocking send and receive where both operations use the same buffer.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Ssend(buf, count, datatype, dest, tag, comm, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, ierror	Sends a message in synchronous mode.
MPI_Ssend_init(buf, count, datatype, dest, tag, comm, request, ierror) <type> buf(*) integer count, datatype, dest, tag, comm, request, ierror	Creates a persistent communication object for a synchronous mode send operation.
MPI_Start(request, ierror) integer request, ierror	Initiates a communication with a persistent request handle.
MPI_Startall(count, array_of_requests, ierror) integer count, array_of_requests(*), ierror	Initiates a collection of requests.
MPI_Test(request, flag, status, ierror) logical flag integer request, status(*), ierror	Returns flag = true if the request (operation) identified in the call is complete.
MPI_Test_cancelled(status, flag, ierror) logical flag integer status(*), ierror	Marks a pending, nonblocking communication operation (send or receive) for cancellation.
MPI_Testall(count, array_of_requests, flag, array_of_statuses, ierror) logical flag integer count, array_of_requests(*), array_of_statuses(MPI_STATUS_SIZE, *), ierror	Returns flag = true if all communications associated with active handles in the array have completed.
MPI_Testany(count, array_of_requests, index, flag, status, ierror) logical flag integer count, array_of_requests(*), index, status(*), ierror	Tests for completion of either one or none of the operations associated with active handles.
MPI_Testsome(incount, array_of_requests, outcount, array_of_indices, array_of_statuses, ierror) integer incount, array_of_requests(*), outcount, array_of_indices(*), array_of_statuses(MPI_STATUS_SIZE, *), ierror	Tests for some given communications to complete.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Topo_test(comm, status, ierror) integer comm, status, ierror	Returns the type of topology that is assigned to a communicator.
MPI_Type_commit(datatype, ierror) integer datatype, ierror	Commits the data type (that is, the formal description of a communication buffer), not the content of that buffer.
MPI_Type_contiguous(count, oldtype, newtype, ierror) integer count, oldtype, newtype, ierror	Creates a contiguous data type.
MPI_Type_extent(datatype, extent, ierror) integer datatype, extent, ierror	Returns the extent of a data type.
MPI_Type_free(datatype, ierror) integer datatype, ierror	Marks the data-type object identified in the call for deallocation.
MPI_Type_hindexed(count, array_of_blocklengths, array_of_displacements, oldtype, newtype, ierror) integer count, array_of_blocklengths(*), array_of_displacements(*), oldtype, newtype, ierror	Creates an indexed data type with offsets in bytes.
MPI_Type_hvector(count, blocklength, stride, oldtype, newtype, ierror) integer count, blocklength, stride, oldtype, newtype, ierror	Creates a vector (strided) data type with offset in bytes.
MPI_Type_indexed(count, array_of_blocklengths, array_of_displacements, oldtype, newtype, ierror) integer count, array_of_blocklengths(*), array_of_displacements(*), oldtype, newtype, ierror	Creates an indexed data type.
MPI_Type_lb(datatype, displacement, ierror) integer datatype, displacement, ierror	Returns the lower bound of a data type.
MPI_Type_size(datatype, size, ierror) integer datatype, size, ierror	Returns the total size (in bytes) of the entries in the data type identified in the call.

MPI library routines and extensions
 Fortran version of MPI routines

Routine	Description
MPI_Type_struct(count, array_of_blocklengths, array_of_displacements, array_of_types, newtype, ierror) integer count, array_of_blocklengths(*), array_of_displacements(*), array_of_types(*), newtype, ierror	Creates a structured data type.
MPI_Type_ub(datatype, displacement, ierror) integer datatype, displacement, ierror	Returns the upper bound of a data type.
MPI_Type_vector(count, blocklength, stride, oldtype, newtype, ierror) integer count, blocklength, stride, oldtype, newtype, ierror	Creates a vector (strided) data type.
MPI_Unpack(inbuf, insize, position, outbuf, outcount, datatype, comm, ierror) <type> inbuf(*), outbuf(*) integer insize, position, outcount, datatype, comm, ierror	Unpacks a message into the receive buffer.
MPI_Wait(request, status, ierror) integer request, status(*), ierror	Waits for an MPI send or receive to complete.
MPI_Waitall(count, array_of_requests, array_of_statuses, ierror) integer count, array_of_requests(*), array_of_statuses(MPI_STATUS_SIZE,*), ierror	Blocks until all communication operations associated with active handles in the list complete and return their status.
MPI_Waitany(count, array_of_requests,index, status, ierror) integer count, array_of_requests(*), index, status(*), ierror	Blocks until one of the operations associated with the active requests in the array has completed.

MPI library routines and extensions
Fortran version of MPI routines

Routine	Description
MPI_Waitsome(incount, array_of_requests, outcount, array_of_indices, array_of_statuses, ierror) integer incount, array_of_requests(*), outcount, array_of_indices(*), array_of_statuses(MPI_STATUS_SIZE,*), ierror	Waits until at least one of the operations associated with active handles in the list have completed.
double precision MPI_Wtick()	Returns the resolution of the MPI_Wtime routine in seconds.
double precision MPI_Wtime()	Returns a floating-point number of seconds representing wall-clock time since some time in the past.

C version of HP MPI library extensions

Table A-3 C version of HP MPI library extensions

Routine	Description
int MPIHP_Comm_id(MPI_Comm comm, int *id)	Returns the communicator context ID.
int MPIHP_Trace_off(void)	Stops HP MPI instrumentation or application tracing.
int MPIHP_Trace_on(void)	Starts HP MPI instrumentation or application tracing.

Fortran version of HP MPI library extensions

Table A-4 Fortran version of HP MPI library extensions

Routine	Description
subroutine MPIHP_Comm_id(comm, id, ierr) integer comm, id, ierr	Returns the communicator context ID.
subroutine MPIHP_Trace_off(ierr) integer ierr	Stops HP MPI instrumentation or application tracing.
subroutine MPIHP_Trace_on(ierr) integer ierr	Starts HP MPI instrumentation or application tracing.

MPI 1.2 extensions

Table A-5 MPI 1.2 extensions

Routine	Description
MPI_Get_version(int *, int *);	Returns the current version of the MPI Standard (for example, 1.2).

MPI 2.0 extensions

Table A-6 **MPI 2.0 extensions**

Routine	Description
MPI_Fint MPI_Comm_c2f(MPI_Comm);	Converts a C communicator handle into a Fortran handle.
MPI_Comm MPI_Comm_f2c(MPI_Fint);	Converts a Fortran communicator handle into a C handle.
MPI_Fint MPI_Type_c2f(MPI_Datatype);	Converts a C data type into a Fortran data type.
MPI_Datatype MPI_Type_f2c(MPI_Fint);	Converts a Fortran data type into a C data type.
MPI_Fint MPI_Group_c2f(MPI_Group);	Converts a C group into a Fortran group.
MPI_Group MPI_Group_f2c(MPI_Fint);	Converts a Fortran group into a C group.
MPI_Fint MPI_Op_c2f(MPI_Op);	Converts a C reduction operation into a Fortran reduction operation.
MPI_Op MPI_Op_f2c(MPI_Fint);	Converts a Fortran reduction operation into a C reduction operation.
MPI_Fint MPI_Request_c2f(MPI_Request);	Converts a C request into a Fortran request.
MPI_Request MPI_Request_f2c(MPI_Fint);	Converts a Fortran request into a C request.
int MPI_Status_c2f(MPI_Status *, MPI_Fint *);	Converts a C status into a Fortran status.
int MPI_Status_f2c(MPI_Fint *, MPI_Status *);	Converts a Fortran status into a C status.
int MPI_Finalized(int *flag);	Checks if MPI_Finalize has completed.

MPI library routines and extensions
MPI 2.0 extensions

Routine	Description
MPI_Finalized(flag, ierror); logical flag integer ierror	Checks if MPI_Finalize has completed.
MPI_Type_get_envelope(MPI_Datatype datatype, int *num_integers, int *num_addresses, int *num_datatypes, int *combiner);	Returns information on the number and type of input arguments used in the call that created a data type.
MPI_Type_get_envelope(datatype, num_integers, num_addresses, num_datatypes, combiner, ierror) integer datatype, num_integers, num_addresses, num_datatypes, combiner, ierror	Returns information on the number and type of input arguments used in the call that created a data type.
MPI_Type_get_contents(MPI_Datatype datatype, int num_integers, int num_addresses, int num_datatypes, int integers [], MPI_Aint addresses [], MPI_Datatype datatypes []);	Returns information about how a data type was constructed.
MPI_Type_get_contents(datatype, num_integers, num_addresses, num_datatypes, integers [], addresses [], datatypes [], ierror) integer datatype, num_integers, num_addresses, num_datatypes, integers [], addresses [], datatypes [], ierror	Returns information about how a data type was constructed.

B

Example applications

This appendix provides example applications that supplement the information in “MPI concepts” on page 4. The examples included are shown in Table B-1.

Table B-1 Example applications shipped with HP MPI

Name	Language	Description	-np argument
send_receive.f	Fortran 77	Illustrates a simple send and receive operation.	-np >= 2
ping_pong.c	C	Measures the time it takes to send and receive data between two processes.	-np = 2
compute_pi.f	Fortran 77	Computes pi by integrating $f(x)=4/(1+x^2)$.	-np >= 1
master_worker.f90	Fortran 90	Distributes sections of an array and performs computation on all sections in parallel.	-np >= 2
cart.C	C++	Generates a virtual topology.	-np = 4
communicator.c	C	Copies the default communicator MPI_COMM_WORLD.	-np = 2
multi_par.f	Fortran 77	Uses the alternating direction iterative (ADI) method on a 2-dimensional compute region. Access to CPSlib is required to run this example. You must also set MPI_TOPOLOGY so the system allocates one process per hypernode. A script file, multi_par.sh, is included to help automate this task.	-np >= 1

These examples and their make file are located in the /opt/mpi/help subdirectory. The examples are presented for illustration purposes only. They may not necessarily represent the most efficient way to solve a given problem.

Example applications

To build and run the examples:

Step 1. Change to a writable directory.

Step 2. Enter

```
% cp /opt/mpi/help/* .
```

Step 3. Enter `make` to build and run all the examples or `make example_name` to build and run a specific application example.

Step 4. Enter

```
% mpirun -j -w -np # program
```

where *program* specifies the path to the executable created in step 3.

send_receive.f

In this Fortran 77 example, process 0 sends an array to other processes in the default communicator MPI_COMM_WORLD.

```

program main

include 'mpif.h'

integer rank, size, to, from, tag, count, i, ierr
integer src, dest
integer st_source, st_tag, st_count
integer status(MPI_STATUS_SIZE)

double precision data(100)

call MPI_Init(ierr)
call MPI_Comm_rank(MPI_COMM_WORLD, rank, ierr)
call MPI_Comm_size(MPI_COMM_WORLD, size, ierr)

if (size .eq. 1) then
  print *, 'must have at least 2 processes'
  call MPI_Finalize(ierr)
  stop
endif

print *, 'Process ', rank, ' of ', size, ' is alive'
dest = size - 1
src = 0

if (rank .eq. src) then
  to = dest
  count = 10
  tag = 2001

  do i=1, 10
    data(i) = 1
  enddo

  call MPI_Send(data, count, MPI_DOUBLE_PRECISION,
+               to, tag, MPI_COMM_WORLD, ierr)
endif

if (rank .eq. dest) then
  tag = MPI_ANY_TAG
  count = 10
  from = MPI_ANY_SOURCE

```

Example applications
send_receive.f

```
call MPI_Recv(data, count, MPI_DOUBLE_PRECISION,  
+           from, tag, MPI_COMM_WORLD, status, ierr)  
call MPI_Get_Count(status, MPI_DOUBLE_PRECISION,  
+           st_count, ierr)  
st_source = status(MPI_SOURCE)  
st_tag = status(MPI_TAG)  
  
print *, 'Status info: source = ', st_source,  
+       ' tag = ', st_tag, ' count = ', st_count  
print *, rank, ' received', (data(i),i=1,10)  
endif  
  
call MPI_Finalize(ierr)  
  
stop  
end
```

send_receive output

The output from running the send_receive executable is shown below.
The application was run with -np = 10.

```
Process 0 of 10 is alive  
Process 1 of 10 is alive  
Process 3 of 10 is alive  
Process 5 of 10 is alive  
Process 9 of 10 is alive  
Process 2 of 10 is alive  
Status info: source = 0 tag = 2001 count = 10  
9 received 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0  
Process 4 of 10 is alive  
Process 7 of 10 is alive  
Process 8 of 10 is alive  
Process 6 of 10 is alive
```

ping_pong.c

This C example is used as a performance benchmark to measure the amount of time it takes to send and receive data between two processes. The buffers are aligned and offset from each other to avoid cache conflicts caused by direct process-to-process byte-copy operations

To run this example:

- Define the CHECK macro to check data integrity.
- Increase the number of bytes to at least twice the cache size to obtain representative bandwidth measurements.

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include <mpi.h>

#define NLOOPS          1000
#define ALIGN          4096

main(argc, argv)

int          argc;
char        *argv[];

{
    int          i, j;
    double      start, stop;
    int         nbytes = 0;
    int         rank, size;
    MPI_Status  status;
    char        *buf;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    if (size != 2) {
        if (!rank) printf("ping_pong: must have two processes\n");
        MPI_Finalize();
        exit(0);
    }

    nbytes = (argc > 1) ? atoi(argv[1]) : 0;
    if (nbytes < 0) nbytes = 0;

    /*
     * Page-align buffers and displace them in the cache to avoid collisions.
     */
    buf = (char *) malloc(nbytes + 524288 + (ALIGN - 1));
```

Example applications
ping_pong.c

```

if (buf == 0) {
    MPI_Abort(MPI_COMM_WORLD, MPI_ERR_BUFFER);
    exit(1);
}

buf = (char *) (((unsigned long) buf) + (ALIGN - 1)) & ~(ALIGN - 1);
if (rank == 1) buf += 524288;
memset(buf, 0, nbytes);
/*
 * Ping-pong.
 */
    if (rank == 0) {
        printf("ping-pong %d bytes ...\n", nbytes);

/*
 * warm-up loop
 */
        for (i = 0; i < 5; i++) {
            MPI_Send(buf, nbytes, MPI_CHAR, 1, 1, MPI_COMM_WORLD);
            MPI_Recv(buf, nbytes, MPI_CHAR,
                    1, 1, MPI_COMM_WORLD, &status);
        }

/*
 * timing loop
 */
        start = MPI_Wtime();
        for (i = 0; i < NLOOPS; i++) {
#ifdef CHECK
            for (j = 0; j < nbytes; j++) buf[j] = (char) (j + i);
#endif
            MPI_Send(buf, nbytes, MPI_CHAR,
                    1, 1000 + i, MPI_COMM_WORLD);
#ifdef CHECK
            memset(buf, 0, nbytes);
            MPI_Recv(buf, nbytes, MPI_CHAR,
                    1, 2000 + i, MPI_COMM_WORLD, &status);
#endif

#ifdef CHECK
            for (j = 0; j < nbytes; j++) {
                if (buf[j] != (char) (j + i)) {
                    printf("error: buf[%d] = %d, not %d\n",
                            j, buf[j], j + i);
                    break;
                }
            }
#endif
        }
        stop = MPI_Wtime();

        printf("%d bytes: %.2f usec/msg\n",
            nbytes, (stop - start) / NLOOPS / 2 * 1000000);
        if (nbytes > 0) {
            printf("%d bytes: %.2f MB/sec\n", nbytes,
                nbytes / 1000000. /
                ((stop - start) / NLOOPS / 2));
        }
    }
}

```

```
        }
    }
    else {
/*
 * warm-up loop
 */
        for (i = 0; i < 5; i++) {
            MPI_Recv(buf, nbytes, MPI_CHAR,
                    0, 1, MPI_COMM_WORLD, &status);
            MPI_Send(buf, nbytes, MPI_CHAR, 0, 1, MPI_COMM_WORLD);
        }

        for (i = 0; i < NLOOPS; i++) {
            MPI_Recv(buf, nbytes, MPI_CHAR,
                    0, 1000 + i, MPI_COMM_WORLD, &status);
            MPI_Send(buf, nbytes, MPI_CHAR,
                    0, 2000 + i, MPI_COMM_WORLD);
        }
    }

    MPI_Finalize();
    exit(0);
}
```

ping_pong output

The output from running the ping_pong executable is shown below. The application was run with `-np = 2`.

```
ping-pong 0 bytes ...
0 bytes: 2.98 3.99 34.99 usec/msg
```

compute_pi.f

This Fortran 77 example computes pi by integrating $f(x) = 4/(1 + x^2)$.
Each process:

- Receives the number of intervals used in the approximation
- Calculates the areas of its rectangles
- Synchronizes for a global summation

Process 0 prints the result and the time it took to complete the calculation.

```
program main

include 'mpif.h'

double precision PI25DT
parameter(PI25DT = 3.141592653589793238462643d0)

double precision mypi, pi, h, sum, x, f, a
integer n, myid, numprocs, i, ierr

C
C Function to integrate
C
  f(a) = 4.d0 / (1.d0 + a*a)

  call MPI_INIT(ierr)
  call MPI_COMM_RANK(MPI_COMM_WORLD, myid, ierr)
  call MPI_COMM_SIZE(MPI_COMM_WORLD, numprocs, ierr)
  print *, "Process ", myid, " of ", numprocs, " is alive"

  sizetype = 1
  sumtype = 2

  if (myid .eq. 0) then
    n = 100
  endif

  call MPI_BCAST(n, 1, MPI_INTEGER, 0, MPI_COMM_WORLD, ierr)

C
C Calculate the interval size.
C
  h = 1.0d0 / n
  sum = 0.0d0
```

```
do 20 i = myid + 1, n, numprocs
    x = h * (dble(i) - 0.5d0)
    sum = sum + f(x)
20 continue

    mypi = h * sum
C
C Collect all the partial sums.
C
    call MPI_REDUCE(mypi, pi, 1, MPI_DOUBLE_PRECISION,
+                 MPI_SUM, 0, MPI_COMM_WORLD, ierr)
C
C Process 0 prints the result.
C
    if (myid .eq. 0) then
        write(6, 97) pi, abs(pi - PI25DT)
97      format(' pi is approximately: ', F18.16,
+            ' Error is: ', F18.16)
    endif

    call MPI_FINALIZE(ierr)

    stop
end
```

compute_pi output

The output from running the compute_pi executable is shown below. The application was run with `-np = 10`.

```
Process 0 of 10 is alive
Process 1 of 10 is alive
Process 3 of 10 is alive
Process 9 of 10 is alive
Process 7 of 10 is alive
Process 5 of 10 is alive
Process 6 of 10 is alive
Process 2 of 10 is alive
Process 4 of 10 is alive
Process 8 of 10 is alive
pi is approximately: 3.1416009869231250
Error is: .0000083333333318
```

master_worker.f90

In this Fortran 90 example, a master task initiates (numtasks - 1) number of worker tasks. The master distributes an equal portion of an array to each worker task. Each worker task receives its portion of the array and sets the value of each element to (the element's index + 1). Each worker task then sends its portion of the modified array back to the master.

```
program array_manipulation
include 'mpif.h'

integer (kind=4) :: status(MPI_STATUS_SIZE)
integer (kind=4), parameter :: ARRAYSIZE = 10000, MASTER = 0
integer (kind=4) :: numtasks, numworkers, taskid, dest, index, i
integer (kind=4) :: arraymsg, indexmsg, source, chunksize, int4, real4
real (kind=4) :: data(ARRAYSIZE), result(ARRAYSIZE)
integer (kind=4) :: numfail

call MPI_Init(ierr)
call MPI_Comm_rank(MPI_COMM_WORLD, taskid, ierr)
call MPI_Comm_size(MPI_COMM_WORLD, numtasks, ierr)
numworkers = numtasks - 1
chunksize = (ARRAYSIZE / numworkers)
arraymsg = 1
indexmsg = 2
int4 = 4
real4 = 4
numfail = 0

! ***** Master task *****
if (taskid .eq. MASTER) then
  data = 0.0
  index = 1
  do dest = 1, numworkers
    call MPI_Send(index, 1, MPI_INTEGER, dest, 0, MPI_COMM_WORLD, ierr)
    call MPI_Send(data(index), chunksize, MPI_REAL, dest, 0, &
      MPI_COMM_WORLD, ierr)
    index = index + chunksize
  end do

  do i = 1, numworkers
    source = i
    call MPI_Recv(index, 1, MPI_INTEGER, source, 1, MPI_COMM_WORLD, &
      status, ierr)
    call MPI_Recv(result(index), chunksize, MPI_REAL, source, 1, &
      MPI_COMM_WORLD, ststus, ierr)
  end do

  do i = 1, numworkers*chunksize
    if (result(i) .ne. (i+1)) then
```

```
        print *, 'element ', i, ' expecting ', (i+1), ' actual is ', result(i)
        numfail = numfail + 1
    endif
enddo

if (numfail .ne. 0) then
    print *, 'out of ', ARRAYSIZE, ' elements, ', numfail, ' wrong answers'
else
    print *, 'correct results!'
endif
end if

! ***** Worker task *****
if (taskid .gt. MASTER) then
    call MPI_Recv(index, 1, MPI_INTEGER, MASTER, 0, MPI_COMM_WORLD, &
        status, ierr)
    call MPI_Recv(result(index), chunksize, MPI_REAL, MASTER, 0, &
        MPI_COMM_WORLD, status, ierr)

    do i = index, index + chunksize
        result(i) = i + 1
    end do

    call MPI_Send(index, 1, MPI_INTEGER, MASTER, 1, MPI_COMM_WORLD, ierr)
    call MPI_Send(result(index), chunksize, MPI_REAL, MASTER, 1, &
        MPI_COMM_WORLD, ierr)
end if
call MPI_Finalize(ierr)

end program array_manipulation
```

master_worker output

The output from running the master_worker executable is shown below.
The application was run with `-np = 2`.

correct results!

cart.C

This C++ program generates a virtual topology. The class Node represents a node in a 2-D torus. Each process is assigned a node or nothing. Each node holds integer data, and the shift operation exchanges the data with its neighbors. Thus, north-east-south-west shifting returns the initial data.

```
#include <stdlib.h>
#include <mpi.h>

#define NDIMS 2

typedef enum { NORTH, SOUTH, EAST, WEST } Direction;

// A node in 2-D torus
class Node {
private:
    MPI_Comm      comm;
    int           dims[NDIMS], coords[NDIMS];
    int           grank, lrank;
    int           data;
public:
    Node(void);
    ~Node(void);
    void profile(void);
    void print(void);
    void shift(Direction);
};

// A constructor
Node::Node(void)
{
    int i, nnodes, periods[NDIMS];

    // Create a balanced distribution
    MPI_Comm_size(MPI_COMM_WORLD, &nnodes);
    for (i = 0; i < NDIMS; i++) { dims[i] = 0; }
    MPI_Dims_create(nnodes, NDIMS, dims);

    // Establish a cartesian topology communicator
    for (i = 0; i < NDIMS; i++) { periods[i] = 1; }
    MPI_Cart_create(MPI_COMM_WORLD, NDIMS, dims, periods, 1, &comm);

    // Initialize the data
    MPI_Comm_rank(MPI_COMM_WORLD, &grank);
    if (comm == MPI_COMM_NULL) {
        lrank = MPI_PROC_NULL;
        data = -1;
    } else {
        MPI_Comm_rank(comm, &lrank);
    }
}
```

```

    data = lrank;
    MPI_Cart_coords(comm, lrank, NDIMS, coords);
}
}

// A destructor
Node::~Node(void)
{
    if (comm != MPI_COMM_NULL) {
        MPI_Comm_free(&comm);
    }
}

// Shift function
void Node::shift(Direction dir)
{
    if (comm == MPI_COMM_NULL) { return; }

    int direction, disp, src, dest;
    if (dir == NORTH) {
        direction = 0; disp = -1;
    } else if (dir == SOUTH) {
        direction = 0; disp = 1;
    } else if (dir == EAST) {
        direction = 1; disp = 1;
    } else {
        direction = 1; disp = -1;
    }

    MPI_Cart_shift(comm, direction, disp, &src, &dest);
    MPI_Status stat;
    MPI_Sendrecv_replace(&data, 1, MPI_INT, dest, 0, src, 0, comm, &stat);
}

// Synchronize and print the data being held
void Node::print(void)
{
    if (comm != MPI_COMM_NULL) {
        MPI_Barrier(comm);
        if (lrank == 0) { puts(""); } // line feed
        MPI_Barrier(comm);
        printf("(%d, %d) holds %d\n", coords[0], coords[1], data);
    }
}

// Print object's profile
void Node::profile(void)
{
    // Non-member does nothing
    if (comm == MPI_COMM_NULL) { return; }

    // Print "Dimensions" at first
    if (lrank == 0) {
        printf("Dimensions: (%d, %d)\n", dims[0], dims[1]);
    }
    MPI_Barrier(comm);
}

```

Example applications

cart.C

```
// Each process prints its profile
printf("global rank %d: cartesian rank %d, coordinate (%d, %d)\n",
      grank, lrank, coords[0], coords[1]);
}

// Program body

//
// Define a torus topology and demonstrate shift operations.
//
void body(void)
{
    Node node;

    node.profile();

    node.print();

    node.shift(NORTH);
    node.print();
    node.shift(EAST);
    node.print();
    node.shift(SOUTH);
    node.print();
    node.shift(WEST);
    node.print();
}

//
// Main program---it is probably a good programming practice to call
// MPI_Init() and MPI_Finalize() here.
//
int main(int argc, char **argv)
{
    MPI_Init(&argc, &argv);
    body();
    MPI_Finalize();
}
```

cart output

The output from running the cart executable is shown below. The application was run with `-np = 4`.

```
Dimensions: (2, 2)
global rank 0: cartesian rank 0, coordinate (0, 0)
global rank 2: cartesian rank 2, coordinate (1, 0)
global rank 3: cartesian rank 3, coordinate (1, 1)
global rank 1: cartesian rank 1, coordinate (0, 1)

(0, 0) holds 0
(0, 1) holds 1
(1, 0) holds 2
```

```
(1, 1) holds 3
(0, 0) holds 2
(0, 1) holds 3
(1, 0) holds 0
(1, 1) holds 1

(0, 0) holds 3
(0, 1) holds 2
(1, 0) holds 1
(1, 1) holds 0

(0, 0) holds 1
(0, 1) holds 0
(1, 0) holds 3
(1, 1) holds 2

(0, 0) holds 0
(1, 1) holds 3
(1, 0) holds 2
(0, 1) holds 1
```

communicator.c

This C example shows how to make a copy of the default communicator MPI_COMM_WORLD.

```
#include <stdio.h>
#include <mpi.h>

main(argc, argv)

int          argc;
char        *argv[];

{
    int          rank, size, data;
    MPI_Status  status;
    MPI_Comm    libcomm;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Comm_dup(MPI_COMM_WORLD, &libcomm);

    if (rank == 0) {
        data = 12345;
        MPI_Send(&data, 1, MPI_INT, 1, 5, MPI_COMM_WORLD);
        data = 6789;
        MPI_Send(&data, 1, MPI_INT, 1, 5, libcomm);
    } else {
        MPI_Recv(&data, 1, MPI_INT, 0, 5, libcomm, &status);
        printf("received libcomm data = %d\n", data);
        MPI_Recv(&data, 1, MPI_INT, 0, 5, MPI_COMM_WORLD, &status);
        printf("received data = %d\n", data);
    }

    MPI_Comm_free(&libcomm);
    MPI_Finalize();
    exit(0);
}
```

communicator output

The output from running the communicator executable is shown below.
The application was run with `-np = 2`.

```
received libcomm data = 6789
received data = 12345
```

multi_par.f

The ADI method is often used to solve differential equations. multi_par.f implements the following logic for a 2-dimensional compute region:

```
DO J=1, JMAX
  DO I=2, IMAX
    A(I, J) = A(I, J) + A(I-1, J)
  ENDDO
ENDDO

DO J=2, JMAX
  DO I=1, IMAX
    A(I, J) = A(I, J) + A(I, J-1)
  ENDDO
ENDDO
```

There are loop-carried dependencies in the first inner DO loop (the array's rows) and the second outer DO loop (the array's columns).

Partitioning the array into column sections supports parallelization of the first outer loop. Partitioning the array into row sections supports parallelization of the second outer loop. However, this approach requires a massive data exchange among processes because of run-time partition changes.

Example applications
multi_par.f

In this case, twisted-data layout partitioning is a better approach because the partitioning used for the parallelization of the first outer loop can accommodate the partitioning of the second outer loop. The partitioning of the array is as follows:

Figure B-1

Array partitioning

		column block			
		0	1	2	3
row block	0	0	1	2	3
	1	3	0	1	2
	2	2	3	0	1
	3	1	2	3	0

In this sample program, the rank n process is assigned to the partition n at distribution initialization. Because these partitions are not contiguous- memory regions, MPI's derived datatype is used to define the partition layout to the MPI system.

Each process starts with computing summations in row-wise fashion. For example, the rank 2 process starts with the block that is on the 0th-row block and 2nd-column block (denoted as [0,2]).

The block computed in the second step is [1,3]. Computing the first row elements in this block requires the last row elements in the [0,3] block (computed in the first step in the rank 3 process). Thus, the rank 2 process receives the data from the rank 3 process at the beginning of the second step. Note that the rank 2 process also sends the last row elements of the [0,2] block to the rank 1 process that computes [1,2] in the second step. By repeating these steps, all processes finish summations in row-wise fashion (the first outer-loop in the illustrated program).

The second outer-loop (the summations in column-wise fashion) is done in the same manner. For example, at the beginning of the second step for the column-wise summations, the rank 2 process receives data from the rank 1 process that computed the [3,0] block. The rank 2 process also sends the last column of the [2,0] block to the rank 3 process. Note that each process keeps the same blocks for both of the outer-loop computations.

This approach is good for distributed memory architectures on which repartitioning requires massive data communications that are expensive. However, on shared memory architectures, the partitioning of the compute region does not imply data distribution. The row- and column-block partitioning method requires just one synchronization at the end of each outer loop.

For distributed shared-memory architectures on X-class servers, the mix of the two methods can be effective. The sample program implements the twisted-data layout method with MPI and the row- and column-block partitioning method with CPSlib. Each MPI process spawns symmetric threads and assigns a block of rows or columns in its computing block to each of threads. Because the computation in the first outer loop does not have dependencies between columns, a thread that is assigned to a block of columns can execute without synchronization. In the second outer-loop, each thread is assigned a block of rows.

Because there is no need for an MPI process to access other processes' memory, it can use node-private memory. On the other hand, threads spawned by a process need to share memory. Given this, the strategy here is to execute MPI processes one per hypernode, assign the computation region on node-private memory, and spawn threads on the node it is running.

```

implicit none
include 'mpif.h'
integer nrow           ! # of rows
integer ncol          ! # of columns
parameter(nrow=1000,ncol=1000)
double precision array(nrow,ncol) ! compute region
c Allocate the compute region in node-private memory
$dir node_private(array)
integer blk           ! block iteration counter
integer rb           ! row block number
integer cb           ! column block number
integer nrb          ! next row block number
integer ncb          ! next column block number
integer rbs(:)       ! row block start subscripts
integer rbe(:)       ! row block end subscripts
integer cbs(:)       ! column block start subscripts

```

Example applications
multi_par.f

```

integer cbe(:)                ! column block end subscripts
integer rdtype(:)            ! row block communication datatypes
integer cdtype(:)           ! column block communication datatypes
integer twdtype(:)          ! twisted distribution datatypes
integer ablen(:)            ! array of block lengths
integer adisp(:)            ! array of displacements
integer adtype(:)           ! array of datatypes
allocatable rbs,rbe,cbs,cbe,rdtype,cdtype,twdtype,ablen,adisp,
*   adtype
integer rank                 ! rank iteration counter
integer comm_size            ! number of MPI processes
integer comm_rank           ! sequential ID of MPI process
integer ierr                 ! MPI error code
integer mstat(mpi_status_size) ! MPI function status
integer src                  ! source rank
integer dest                 ! destination rank
integer dsize                ! size of double precision in bytes
double precision startt,endt,elapsed ! time keepers
integer ncpus                ! # of CPUs on the node
integer acpus                ! total # of CPUs
integer params(4)           ! parameters for thread creation
integer rc                   ! cps funtion return code
external compcolumn,comprow ! subroutines execute in threads

c   CPS functions
integer cps_node_cpus,cps_ppcalln

c$dir sync_routine(cps_node_cpus,cps_ppcalln)
c
c   MPI initialization
c
call mpi_init(ierr)
call mpi_comm_size(mpi_comm_world,comm_size,ierr)
call mpi_comm_rank(mpi_comm_world,comm_rank,ierr)

c
c   CPS setup
c
ncpus=cps_node_cpus()
if (ncpus.le.0) then
  write(7,*) 'More than 1 CPUs must be available on the node ',
*   'on which rank',comm_rank,' is running'
  call mpi_abort(mpi_comm_world,1,ierr)
else
  write(6,*) ncpus,' threads in rank',comm_rank,' process'
  call mpi_reduce(ncpus,acpus,1,mpi_integer,mpi_sum,0,
*   mpi_comm_world,ierr)
endif

c
c   Allocate threads on the same node on which calling processes
c   execute so all of them can access the compute region that is
c   assigned in the node-private memory. The following parameter
c   vector is passed to the symmetric thread-spawn function.
c
params(1)=-1          ! alloc. threads on same node as calling thread

```

```

params(2)=ncpus      ! min. # of threads
params(3)=ncpus      ! max. # of threads
params(4)=1          ! alloc. threads per node

c
c  Data initialization and start up
c
c  if (comm_rank.eq.0) then
      write(6,*) 'Initializing',nrow,' x',ncol,' array...'
      call getdata(nrow,ncol,array)
      write(6,*) 'Start computation with total of',acpus,' threads'
      startt=mpi_wtime()
  endif

c
c  Compose MPI datatypes for row/column send-receive
c
c  Because rows are assigned contiguously in memory, row-wise
c  communication regions are defined as MPI's contiguous datatype
c  while column-wise regions are defined as MPI's vector datatype.
c
  allocate(rbs(0:comm_size-1),rbe(0:comm_size-1),cbs(0:comm_size-1),
*         cbe(0:comm_size-1),rdtype(0:comm_size-1),
*         cdtype(0:comm_size-1),twdtype(0:comm_size-1))
  do blk=0,comm_size-1
      call blockasn(1,nrow,comm_size,blk,rbs(blk),rbe(blk))
      call mpi_type_contiguous(rbe(blk)-rbs(blk)+1,
*                             mpi_double_precision,rdtype(blk),ierr)
      call mpi_type_commit(rdtype(blk),ierr)
      call blockasn(1,ncol,comm_size,blk,cbs(blk),cbe(blk))
      call mpi_type_vector(cbe(blk)-cbs(blk)+1,1,nrow,
*                          mpi_double_precision,cdtype(blk),ierr)
      call mpi_type_commit(cdtype(blk),ierr)
  enddo

c
c  Compose MPI datatypes for gather/scatter
c
c  Each block of the partitioning is defined as a set of fixed-length
c  vectors. Each process's partition is defined as a structure of such
c  blocks.
c
  allocate(adtype(0:comm_size-1),adisp(0:comm_size-1),
*         ablen(0:comm_size-1))
  call mpi_type_extent(mpi_double_precision,dsize,ierr)
  do rank=0,comm_size-1
      do rb=0,comm_size-1
          cb=mod(rb+rank,comm_size)
          call mpi_type_vector(cbe(cb)-cbs(cb)+1,rbe(rb)-rbs(rb)+1,
*                              nrow,mpi_double_precision,adtype(rb),ierr)
          call mpi_type_commit(adtype(rb),ierr)
          adisp(rb)=((rbs(rb)-1)+(cbs(cb)-1)*nrow)*dsize
          ablen(rb)=1
      enddo
      call mpi_type_struct(comm_size,ablen,adisp,adtype,
*                          twdtype(rank),ierr)
      call mpi_type_commit(twdtype(rank),ierr)
  enddo

```

Example applications

multi_par.f

```
do rb=0,comm_size-1
  call mpi_type_free(adtype(rb),ierr)
enddo

enddo
deallocate(adtype,adisp,ablen)
```

```
c
c Scatter initial data using derived datatypes defined above
c for the partitioning. MPI_send() and MPI_recv() determine the
c layout of the data from these datatypes. This saves application
c programs from having to manually pack/unpack the data. More importantly,
c it provides opportunities for optimal communication
c strategies.
c
```

```
if (comm_rank.eq.0) then
  do dest=1,comm_size-1
    call mpi_send(array,1,twdtype(dest),dest,0,mpi_comm_world,
*      ierr)
  enddo
else
  call mpi_recv(array,1,twdtype(comm_rank),0,0,mpi_comm_world,
*      mstat,ierr)
endif
```

```
c
c Computation
```

```
c
c Sum up in each column.
c Each MPI process, or a rank, computes blocks that it is assigned.
c The column block number is assigned in the variable 'cb'. The
c starting and ending subscripts of the column block 'cb' are
c stored in 'cbs(cb)' and 'cbe(cb)' respectively. The row block
c number is assigned in the variable 'rb'. The starting and ending
c subscripts of the row block 'rb' are stored in 'rbs(rb)' and
c 'rbe(rb)' respectively.
c
```

```
src=mod(comm_rank+1,comm_size)
dest=mod(comm_rank-1+comm_size,comm_size)
ncb=comm_rank
do rb=0,comm_size-1
  cb=ncb
```

```
c
c Compute a block with threads.
c cps_ppcalln() spawns symmetric threads, executes a function in
c parallel, and automatically joins the threads. Note that all MPI
c processes execute this code
c
```

```
rc=cps_ppcalln(params,compcolumn,8,ncpus,nrow,ncol,array,
*   rbs(rb),rbe(rb),cbs(cb),cbe(cb))
if (rc.le.0) then
  write(7,*) 'rank',comm_rank,
*   ' cps_ppcalln(compcolumn) returned',rc
  call mpi_abort(mpi_comm_world,1,ierr)
endif
if (rb.lt.comm_size-1) then
```

```

c
c   Send the last row of the block to the rank that computes the
c   block next to the computed block. Receive the last row of the
c   block that the next block being computed depends on.
c
      nrb=rb+1
      ncb=mod(nrb+comm_rank,comm_size)
      call mpi_sendrecv(array(rbe(rb),cbs(cb)),1,cdtype(cb),dest,
*         0,array(rbs(nrb)-1,cbs(ncb)),1,cdtype(ncb),src,0,
*         mpi_comm_world,mstat,ierr)
      endif
    enddo

c
c   Sum up in each row.
c   The same logic as the loop above except rows and columns are
c   switched.
c
      src=mod(comm_rank-1+comm_size,comm_size)
      dest=mod(comm_rank+1,comm_size)
      do cb=0,comm_size-1
        rb=mod(cb-comm_rank+comm_size,comm_size)
        rc=cps_ppcalln(params,comprow,8,ncpus,nrow,ncol,array,
*         rbs(rb),rbe(rb),cbs(cb),cbe(cb))
        if (rc.le.0) then
          write(7,*) 'rank',comm_rank,
*         ' cps_ppcalln(comprow) returned',rc
          call mpi_abort(mpi_comm_world,1,ierr)
        endif
        if (cb.lt.comm_size-1) then
          ncb=cb+1
          nrb=mod(ncb-comm_rank+comm_size,comm_size)
          call mpi_sendrecv(array(rbs(rb),cbe(cb)),1,rdtype(rb),dest,
*         0,array(rbs(nrb),cbs(ncb)-1),1,rdtype(nrb),src,0,
*         mpi_comm_world,mstat,ierr)
        endif
      enddo

c
c   Gather computation results
c
      if (comm_rank.eq.0) then
        do src=1,comm_size-1
          call mpi_recv(array,1,twdtype(src),src,0,mpi_comm_world,
*         mstat,ierr)
        enddo
        endt=mpi_wtime()
        elapsed=endt-startt
        write(6,*) 'Computation took',elapsed,' seconds'
      else
        call mpi_send(array,1,twdtype(comm_rank),0,0,mpi_comm_world,
*         ierr)
      endif

c
c   Dump to a file
c

```

Example applications
multi_par.f

```

c      if (comm_rank.eq.0) then
c          print*, 'Dumping to adi.out...'
c          open(8, file='adi.out', form='unformatted')
c          write(8, *) array
c          close(8, status='keep')
c      endif

c
c      Free the resources
c
c      do rank=0, comm_size-1
c          call mpi_type_free(twdtype(rank), ierr)
c      enddo
c      do blk=0, comm_size-1
c          call mpi_type_free(rdtype(blk), ierr)
c          call mpi_type_free(cdtype(blk), ierr)
c      enddo
c      deallocate(rbs, rbe, cbs, cbe, rdtype, cdtype, twdtype)

c
c      Finalize the MPI system
c
c      call mpi_finalize(ierr)
c      end

c*****
subroutine blockasgn(subs, sube, blockcnt, nth, blocks, blocke)
c
c      This subroutine
c      is given a range of subscripts and the total number of blocks in
c      which the range is to be divided and assigns a subrange to the caller
c      that is n-th member of the blocks.
c
c      implicit none
c      integer subs          ! (in)      subscript start
c      integer sube          ! (in)      subscript end
c      integer blockcnt      ! (in)      block count
c      integer nth           ! (in)      my block (begin from 0)
c      integer blocks        ! (out)     assigned block start subscript
c      integer blocke        ! (out)     assigned block end subscript

c
c      integer dl, m1

c      dl=(sube-subs+1)/blockcnt
c      m1=mod(sube-subs+1, blockcnt)
c      blocks=nth*d1+subs+min(nth, m1)
c      blocke=blocks+d1-1
c      if (m1.gt.nth) blocke=blocke+1
c      end

c*****
subroutine compcolumn(ncpus, nrow, ncol, array, rbs, rbe, cbs, cbe)
c
c      This subroutine
c      does summations of columns in a thread.
c

```

```

implicit none
integer ncpus           ! # of cpus on node
integer nrow            ! length of row
integer ncol            ! length of column
double precision array(nrow,ncol) ! region
integer rbs             ! row block start subscript
integer rbe             ! row block end subscript
integer cbs             ! column block start subscript
integer cbe             ! column block end subscript

c
integer stid,cps_stid   ! ID of the symmetric thread
c$dir sync_routine(cps_stid)
c
c Local variables
c
integer mycbs,mycbe     ! my column start/end index
integer i,j

c
c Assign a range of columns to this thread and compute
c
stid=cps_stid()
call blockasgn(cbs,cbe,ncpus,stid,mycbs,mycbe)
do j=mycbs,mycbe
  do i=max(2,rbs),rbe
    array(i,j)=array(i-1,j)+array(i,j)
  enddo
enddo
end

c*****
subroutine comprow(ncpus,nrow,ncol,array,rbs,rbe,cbs,cbe)

c
c This subroutine
c does summations of rows in a thread.
c
implicit none
integer ncpus           ! # of cpus on node
integer nrow            ! length of row
integer ncol            ! length of column
double precision array(nrow,ncol) ! region
integer rbs             ! row block start subscript
integer rbe             ! row block end subscript
integer cbs             ! column block start subscript
integer cbe             ! column block end subscript

c
integer stid,cps_stid   ! ID of the symmetric thread
c$dir sync_routine(cps_stid)
c
c Local variables
c
integer myrbs,myrbe     ! my row start/end index
integer i,j

```

Example applications
multi_par.f

```
c
c   Assign a range of rows to this thread and compute
c
  stid=cps_stid()
  call blockasgn(rbs,rbe,ncpus,stid,myrbs,myrbe)
  do j=max(2,cbs),cbe
    do i=myrbs,myrbe
      array(i,j)=array(i,j-1)+array(i,j)
    enddo
  enddo
end

c*****
subroutine getdata(nrow,ncol,array)
c
c   Enter dummy data
c
  integer nrow,ncol
  double precision array(nrow,ncol)
c
  do j=1,ncol
    do i=1,nrow
      array(i,j)=(j-1.0)*ncol+i
    enddo
  enddo
end
```

C

XMPI resource file

This appendix displays the contents of the XMPI X resource file stored in `/opt/mpi/lib/X11/app-defaults/XMPI`.

You should make your own copy of the resource file (you can copy it to the `.Xdefaults` file in your home directory) and tailor it accordingly.

To save your changes and rebuild the X resource database from scratch, enter:

```
% xrdb filename
```

To save your changes and merge them into the existing X resource database, enter:

```
% x-rdb -merge filename
```

In both cases, *filename* represents the name of your tailored resource file.

```
XMPI*Title:XMPI
XMPI*IconName:XMPI
XMPI*multiClickTime:500
XMPI*background:lightgray
XMPI*fontList:--*helvetica-bold-r-normal--*-120-*-*-*-*-*-*
XMPI*rankFont:--*helvetica-bold-r-normal--*-120-*-*-*-*-*-*
XMPI*msgFont:--*helvetica-medium-r-normal--*-120-*-*-*-*-*-*
XMPI*fo_func.fontList:--*helvetica-bold-o-normal--*-120-*-*-*-*-*-*
XMPI*dt_dtype.fontList:--*helvetica-medium-r-normal--*-100-*-*-*-*-*-*
XMPI*ctl_bar.topShadowColor:lightslateblue
XMPI*ctl_bar.bottomShadowColor:darkslateblue
XMPI*ctl_bar.background:slateblue
XMPI*ctl_bar.foreground:white
XMPI*banner.background:slateblue
XMPI*banner.foreground:white
XMPI*view_draw.background:black
XMPI*view_draw.foreground:gray
```

XMPI resource file

```
XMPI*trace_draw.background:gray
XMPI*trace_draw.foreground:black
XMPI*kiviat_draw.background:gray
XMPI*kiviat_draw.foreground:black
XMPI*app_list.visibleItemCount:8
XMPI*aschema_list.visibleItemCount:20
XMPI*aschema_text.columns:24
XMPI*prog_mgr*columns:16
XMPI*comCol:cyan
XMPI*rcomCol:plum
XMPI*label_frame.XmLabel.background:#D3B5B5
XMPI*XmToggleButtonGadget.selectColor:red
XMPI*XmToggleButton.selectColor:red
```

Glossary

asynchronous Communication in which sending and receiving processes place no constraints on each other in terms of completion. The communication operation between the two processes may also overlap with computation.

bandwidth Reciprocal of the time needed to transfer a byte. Bandwidth is normally expressed in megabytes per second.

barrier Collective operation used to synchronize the execution of processes. `MPI_Barrier` blocks the calling process until all receiving processes have called it. This is a useful approach for separating two stages of a computation so messages from each stage are not overlapped.

blocking receive Communication in which the receiving process does not return until its data buffer contains the data transferred by the sending process.

blocking send Communication in which the sending process does not return until its associated data buffer is available for reuse. The data transferred can be copied directly into the matching receive buffer or a temporary system buffer.

broadcast One-to-many collective operation where the root process sends a message to all other processes in the communicator including itself.

buffered send mode Form of blocking send where the sending process returns when the message is buffered in application-supplied space or when the message is received.

buffering Amount or act of copying that a system uses to avoid deadlocks. A large amount of buffering can adversely affect performance and make MPI applications less portable and predictable.

cluster Group of computers linked together with an interconnect and software that functions collectively as a parallel machine.

collective communication Communication that involves sending or receiving messages among a group of processes at the same time. The communication can be one-to-many, many-to-one, or many-to-many. The main collective routines are `MPI_Bcast`, `MPI_Gather`, and `MPI_Scatter`.

communicator Global object that groups application processes together. Processes in a communicator can communicate with each other or with processes in another group. Conceptually, communicators define a communication context and a static group of processes within that context.

context Internal abstraction used to define a safe communication space for processes. Within a communicator, context separates point-to-point and collective communications.

data-parallel model Design model where data is partitioned and distributed to each process in an application. Operations are performed on each set of data in parallel and intermediate results are exchanged between processes until a problem is solved.

derived data types User-defined structures that specify a sequence of basic data types and integer displacements for noncontiguous data. You create derived data types through the use of type-constructor functions that describe the layout of sets of primitive types in memory. Derived types may contain arrays as well as combinations of other primitive data types.

domain decomposition Breaking down an MPI application's computational space into regular data structures such that all computation on these structures is identical and performed in parallel.

explicit parallelism Programming style that requires you to specify parallel constructs directly. Using the MPI library is an example of explicit parallelism.

functional decomposition Breaking down an MPI application's computational space into separate tasks such that all computation on these tasks is performed in parallel.

gather Many-to-one collective operation where each process (including the root) sends the contents of its send buffer to the root.

granularity Measure of the work done between synchronization points. Fine-grained applications focus on execution at the instruction level of a program. Such applications are load balanced but suffer from a low computation/communication ratio. Coarse-grained applications focus on execution at the program level where multiple programs may be executed in parallel.

group Set of tasks that can be used to organize MPI applications. Multiple groups are useful for solving problems in linear algebra and domain decomposition.

hypernode Building block of an Exemplar-scalable system. Each hypernode consists of a number of processors, I/O, and memory connected by a crossbar and joined to other hypernodes by a Coherent Toroidal Interconnect link.

implicit parallelism Programming style where parallelism is achieved by software layering (that is, parallel constructs are generated through the software). High performance Fortran is an example of implicit parallelism.

intercommunicators Communicators that allow only processes within the same group or in two different groups to exchange data. These communicators support only point-to-point communication.

intracommunicators Communicators that allow processes within the same group to exchange data. These communicators support both point-to-point and collective communication.

instrumentation Cumulative statistical information collected and stored in ascii format. Instrumentation is the recommended method for collecting profiling data.

latency Time between the initiation of the data transfer in the sending process and the arrival of the first byte in the receiving process.

load balancing Measure of how evenly the work load is distributed among an application's processes. When an application is perfectly balanced, all processes share the total work load and complete at the same time.

locality Degree to which computations performed by a processor depend only upon local data. Locality is measured in several ways including the ratio of local to nonlocal data accesses.

message-passing model Model in which processes communicate with each other by sending and receiving messages. Applications based on message passing are nondeterministic by default. However, when one process sends two or more messages to another, the transfer is deterministic as the messages are always received in the order sent.

MIMD Multiple instruction multiple data. Category of applications in which many instruction streams are applied concurrently to multiple data sets.

MPI Message passing interface. Set of library routines used to design scalable parallel applications. These routines provide a wide range of operations that include computation, communication, and synchronization. MPI 1.2 is the current standard supported by major vendors.

MPMD Multiple data multiple program. Implementations of HP MPI that use two or more separate executables to construct an application. This design style can be used to simplify the application source and reduce the size of spawned processes. Each process may run a different executable.

multilevel parallelism Refers to multithreaded processes that call MPI routines to perform computations. This approach is beneficial for problems that can be decomposed into logical parts for parallel execution (for example, a looping construct that spawns multiple threads to perform a computation and then joins after the computation is complete).

nonblocking receive Communication in which the receiving process returns before a message is stored in the receive buffer. Nonblocking receives are useful when communication and computation can be effectively overlapped in an MPI application. Use of nonblocking receives may also avoid system buffering and memory-to-memory copying.

nonblocking send Communication in which the sending process returns before a message is stored in the send buffer. Nonblocking sends are useful when communication and computation can be effectively overlapped in an MPI application.

NUMA Nonuniform memory access architecture. Amount of time for processes to access memory across hypernodes is nonuniform depending upon where data is stored.

parallel efficiency Speed up in the execution of a parallel application.

point-to-point communication Communication where data transfer involves sending and receiving messages between two processes. This is the simplest form of data transfer in a message-passing model.

polling Mechanism to handle asynchronous events by actively checking to determine if an event has occurred.

process Address space together with a program counter, a set of registers, and a stack. Processes can be single threaded or multithreaded. Single-threaded processes can only perform one task at a time. Multithreaded processes can perform multiple tasks concurrently as when overlapping computation and communication.

race condition Situation in which multiple processes vie for the same resource and receive it in an unpredictable manner. Race conditions can lead to cases where applications do not run correctly from one invocation to the next.

rank Integer between zero and (number of processes - 1) that defines the order of a process in a communicator. Determining the rank of a process is important when solving problems where a master process partitions and distributes work to slave processes. The slaves perform some computation and return the result to the master as the solution.

ready send mode Form of blocking send where the sending process cannot start until a matching receive is posted. The sending process returns immediately.

reduction Binary operations (such as summation, multiplication, and boolean) applied globally to all processes in a communicator. These operations are only valid on numeric data and are always associative but may or may not be commutative.

scalable Ability to deliver an increase in application performance proportional to an increase in hardware resources (normally, adding more processors).

scatter One-to-many operation where the root's send buffer is partitioned into n segments and distributed to all processes such that the i th process receives the i th segment. n represents the total number of processes in the communicator.

send modes Point-to-point communication in which messages are passed using one of four different types of blocking sends. The four send modes include standard mode (MPI_Send), buffered mode (MPI_Bsend), synchronous mode (MPI_Ssend), and ready mode (MPI_Rsend). The modes are all invoked in a similar manner and all pass the same arguments.

shared memory model Model in which each process can access a shared address space. Concurrent accesses to shared memory are controlled by synchronization primitives.

SIMD Single instruction multiple data. Category of applications in which homogeneous processes execute the same instructions on their own data.

SPMD Single program multiple data. Implementations of HP MPI where an application is completely contained in a single executable. SPMD applications begin with the invocation of a single process called the master. The master then spawns some number of identical child processes. The master and the children all run the same executable.

standard send mode Form of blocking send where the sending process returns when the system can buffer the message or when the message is received.

stride Constant amount of memory space between data elements where the elements are stored noncontiguously. Strided data are sent and received using derived data types.

subcomplex Group of processors and their associated memory that may span multiple hypernodes on the same host. Hosts are partitioned into subcomplex configurations to achieve the best mix of hardware and software resources.

synchronization Bringing multiple processes to the same point in their execution before any can continue. For example, `MPI_Barrier` is a collective routine that blocks the calling process until all receiving processes have called it. This is a useful approach for separating two stages of a computation so messages from each stage are not overlapped.

synchronous send mode Form of blocking send where the sending process returns only if a matching receive is posted and the receiving process has started to receive the message.

tag Integer label assigned to a message when it is sent. Message tags are one of the synchronization variables used to ensure that a message is delivered to the correct receiving process.

task Uniquely addressable thread of execution.

thread Smallest notion of execution in a process. All MPI processes have one or more threads. Multithreaded processes have one address space but each process thread contains its own counter, registers, and stack. This allows rapid context switching because threads require little or no memory management.

topologies Process configurations that determine which processes to run on specific hypernodes in a given subcomplex. You can use the `MPI_TOPOLOGY` environment variable in HP MPI or one of the MPI library routines (for example, `MPI_GRAPH_CREATE` or `MPI_CART_CREATE`) to define an application topology.

trace Information collected during program execution that you can use to analyze your application. You can collect trace information and store it in a file for later use or analyze it directly when running your application interactively (for example, when you run an application in the XMPI utility).

Index

Numerics

64-bit support, 30

A

advanced MPI topics, 17
application compiling, 29
array partitioning figure, 154
associated documentation, xiv

B

bandwidth, 6
blocking communication, 8
build problems, 94
building and running example applications, 138
building and running your first application
 building and running on a single host, 22
 building and running on multiple hosts, 23

C

C version of HP MPI library extensions table, 132
C version of MPI routines table, 104
cart output, 150
cart.C application, 148
collective operations
 communication, 11
 computation, 13
 synchronization, 14
communication, 11
communicator output, 152
communicator.c application, 152
communicators, 6
compatibility issues, 28
compilation environment variables table, 29
compilation utilities table, 29

compiling applications, 29
computation, 13
compute_pi output, 145
compute_pi.f application, 144
configuring your environment, 20
counter instrumentation, 54
creating an appfile, 47
creating an instrumentation profile, 54
CXdb, 93
CXperf, 79

D

debugging HP MPI applications
 setting options in MPI_FLAGS, 92
 using CXdb, 93
default process placement figure, 88
dial time, 62
dialogs
 mpirun options, 76, 77
 XMPI Application Browser, 70
 XMPI buffer size, 76
 XMPI Datatype, 67
 XMPI Dump, 72
 XMPI Express, 73
 XMPI Focus, 66
 XMPI Kiviat, 69
 XMPI monitor options, 75
 XMPI Trace, 62
 XMPI Trace Selection, 61
directory structure, 26

E

environment variables
 run time
 MPI_SHMEMCNTL, 40

environment variables

compilation
 MPI_CC, 29
 MPI_CXX, 29
 MPI_F77, 29
 MPI_F90, 29
run time
 MPI_CHECKPOINT, 43
 MPI_FLAGS, 37
 MPI_GLOBBMEMSIZE, 38
 MPI_INSTR, 43
 MPI_TMPDIR, 40
 MPI_TOPOLOGY, 39
 MPI_WORKDIR, 42
 MPI_XMPI, 41
example applications
 cart.C, 148
 communicator.c, 152
 compute_pi.f, 144
 master_worker.f90, 146
 multi_par.f, 153
 ping_pong.c, 141
 send_receive.f, 139
example applications shipped with HP MPI table, 137

F

figures
 array partitioning, 154
 default process placement, 88
 MPI broadcast operation, 11
 MPI scatter operation, 12
 multiple network interfaces, 85
 multiprotocol messaging with a K-Class server, 35
 multiprotocol messaging with an X-Class server, 34
 optimal process placement, 89
Fortran 77 version of HP MPI library extensions table, 133

-
- Fortran 77 version of MPI
 routines table, 117
 frequently asked questions, 100
- G**
 getting started, 19
 glossary, 165
- H**
 HP MPI features
 compliance with the MPI 1.2
 standard, xi
 compliance with the UNIX 95
 standard, xii
 data mover, xii
 derived data types, xii
 multiprotocol support, xi
 profiling, xi
 single program multiple data
 and multiple program
 multiple data, xi
 supports a subset of the MPI
 2.0 standard, xi
- K**
 kiviav view, 69
- L**
 language interoperability, 31
 latency, 6
- M**
 man page categories table, 27
 master_worker output, 147
 master_worker.f90 application,
 146
 message passing, advantages, 3
- MPI**
 advanced topics, 17
 collective operations, 10
 library routines and
 extensions, 103
 message-passing library, 4
 multilevel parallelism, 17
 noncontiguous data, 15
 point-to-point
 communications, 6
 six commonly used routines, 5
 MPI 1.2 extensions, 134
 MPI 1.2 extensions table, 134
 MPI 2.0 extensions, 135
 MPI 2.0 extensions table, 135
 MPI blocking and nonblocking
 calls table, 9
 MPI broadcast operation figure,
 11
 MPI library extensions
 C version of HP MPI library
 routines, 132
 Fortran 77 version of HP MPI
 library extensions, 133
 MPI library routines
 C version of MPI routines, 104
 Fortran 77 version of MPI
 routines, 117
 MPI routine selection, 86
 MPI scatter operation figure, 12
 MPI web sites, xiv
 MPI_CC environment variable,
 29
 MPI_COMM_SELF
 communicator, 6
 MPI_COMM_WORLD
 communicator, 6
 MPI_CXX environment variable,
 29
 MPI_F77 environment variable,
 29
 MPI_F90 environment variable,
 29
- mpiCC utility, 29
 mpicc utility, 29
 mpiclean utility, 49
 mpif77 utility, 29
 mpif90 utility, 29
 mpijob utility, 48
 mpirun options dialog, 76
 mpirun options trace dialog, 77
 mpirun utility, 45
 mpitrget utility, 51
 mpitrstat utility, 51
 multiprotocol messaging with a
 K-Class server figure, 35
 multi_par.f application, 153
 multilevel parallelism, 17, 87
 multiple network interfaces
 figure, 85
 multiprotocol messaging
 K-Class server, 35
 X-Class server, 34
 multiprotocol messaging with an
 X-Class server figure, 34
- N**
 nonblocking communication, 9
 noncontiguous data, 15
 notational conventions, xiii
- O**
 optimal process placement
 figure, 89
 organization of the /opt/mpi
 directory table, 26
- P**
 parallel computational models, 2
 ping_pong output, 143
 ping_pong.c application, 141
 point-to-point communications,
 6
 process placement, 87
 processor subscription, 85
-

profiling
 using counter
 instrumentation, 54
 using CXperf, 79
 using the profiling interface,
 80
 using XMPI, 58
profiling interface, 80

R

resource file, XMPI, 163
run invocations that support
 stdin table, 97
running MPMD applications, 33
run-time problems, 95

S

send modes, blocking
 communication, 8
send_receive output, 140
send_receive.f application, 139
sending and receiving messages
 blocking communication, 8
 nonblocking communication, 9
setting options in MPI_FLAGS,
 92
six commonly used MPI routines
 table, 5
start problems, 94
stdin support, 97
subscription types table, 85
supported platforms
 HP-UX, xii
 SPP-UX, xii
synchronization, 14

T

tables
 C version of HP MPI library
 extensions, 132
 C version of MPI routines, 104
 compilation environment
 variables, 29
 compilation utilities, 29
 example applications shipped
 with HP MPI, 137
 Fortran 77 version of HP MPI
 library extensions, 133
 Fortran 77 version of MPI
 routines, 117
 man page categories, 27
 MPI 1.2 extensions, 134
 MPI 2.0 extensions, 135
 MPI blocking and nonblocking
 calls, 9
 organization of the /opt/mpi
 directory, 26
 run invocations that support
 stdin table, 97
 six commonly used MPI
 routines, 5
 subscription types, 85
total transfer time, 6
troubleshooting applications
 build problems, 94
 completion problems, 98
 run-time problems
 external input and output, 97
 Fortran 90, 98
 interoperability, 96
 message buffering, 96
 propagation of environment
 variables, 95
 shared memory, 95
 UNIX open file descriptors,
 98
 start problems, 94

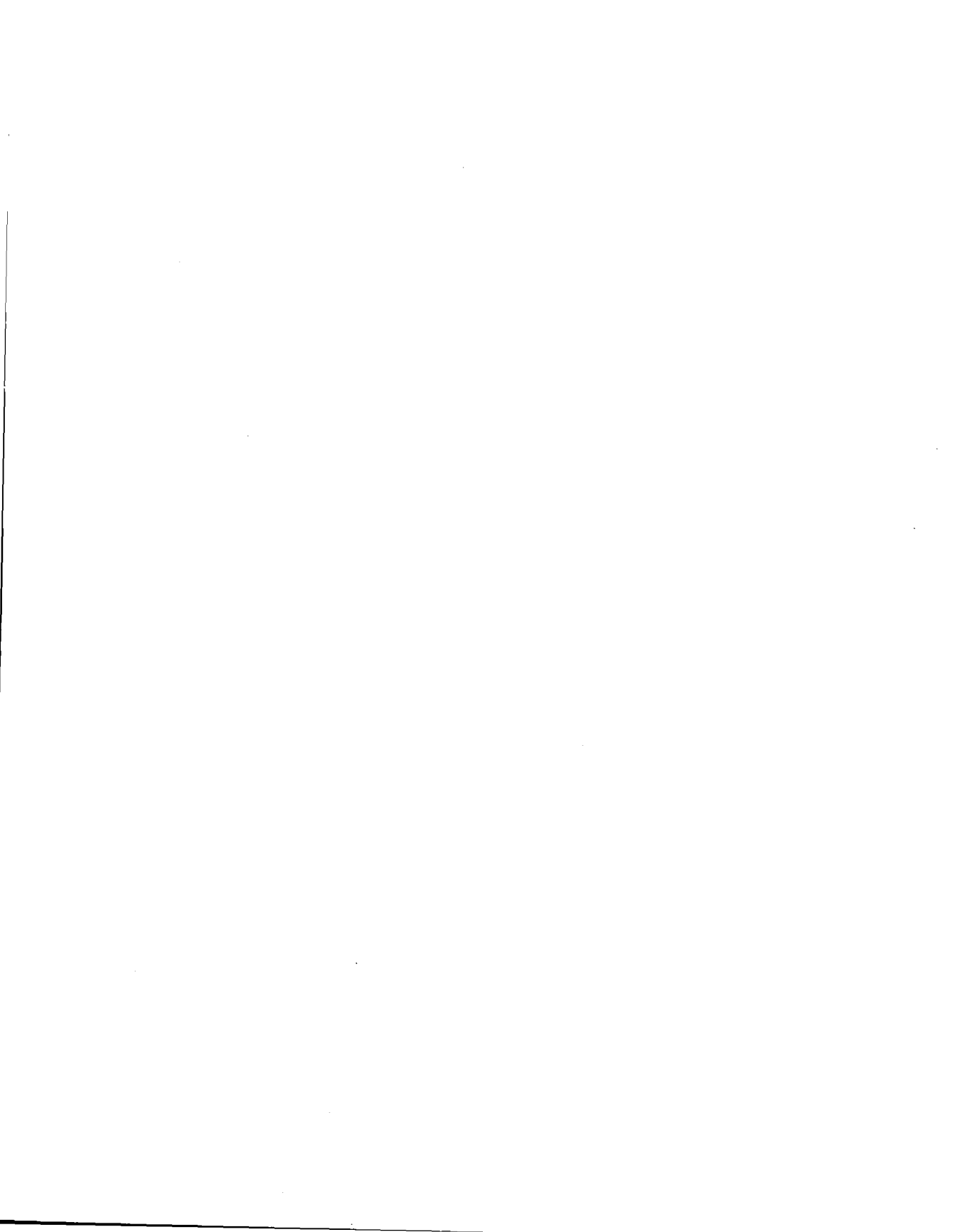
tuning application performance
 general tuning
 message latency and
 bandwidth, 82
 MPI routine selection, 86
 multiple network interfaces,
 83
 processor subscription, 85
 SPP-UX platform tuning
 multilevel parallelism, 87
 process placement, 87
types of applications
 running MPMD applications,
 33
 running SPMD applications,
 32

U

understanding HP MPI, 25
using the profiling interface, 80
using XMPI
 working with interactive mode
 running an appfile, 70
 setting up your viewing
 options, 75
 working with postmortem
 mode
 creating a trace file, 59
 viewing a trace file, 60
utilities
 compilation
 mpiCC, 29
 mpicc, 29
 mpif77, 29
 mpif90, 29
 run time
 mpiclean, 49
 mpijob, 48
 mpirun, 45
 mpitrget, 51
 mpitrstat, 51
 xmpi, 50

X

- XMPI Application Browser dialog, 70
- XMPI buffer size dialog, 76
- XMPI Datatype dialog, 67
- XMPI Dump dialog, 72
- XMPI Express dialog, 73
- XMPI Focus dialog, 66
- XMPI Kiviat dialog, 69
- XMPI main window, 60
- XMPI monitor options dialog, 75
- XMPI resource file, 163
- XMPI Trace dialog, 62
- XMPI Trace Selection dialog, 61
- xmpi utility, 50



Order Number
B6011-90001
E1097



Printed on
Recycled Paper

Printed in USA



B6011-96002
Manufacturing Part Number