



# Sun™ MPI 6.0 Software Programming and Reference Manual

---

Sun Microsystems, Inc.  
4150 Network Circle  
Santa Clara, CA 95054 U.S.A.  
650-960-1300

Part No. 817-0085-10  
February 2003, Revision 01

Send comments about this document to: [docfeedback@sun.com](mailto:docfeedback@sun.com)

Copyright 2003 Sun Microsystems, Inc., 4150 Network Circle, Santa Clara, California 95054, U.S.A. All rights reserved.

Sun Microsystems, Inc. has intellectual property rights relating to technology embodied in the product that is described in this document. In particular, and without limitation, these intellectual property rights may include one or more of the U.S. patents listed at <http://www.sun.com/patents> and one or more additional patents or pending patent applications in the U.S. and in other countries.

This document and the product to which it pertains are distributed under licenses restricting their use, copying, distribution, and decompilation. No part of the product or of this document may be reproduced in any form by any means without prior written authorization of Sun and its licensors, if any.

Third-party software, including font technology, is copyrighted and licensed from Sun suppliers.

Parts of the product may be derived from Berkeley BSD systems, licensed from the University of California. UNIX is a registered trademark in the U.S. and in other countries, exclusively licensed through X/Open Company, Ltd.

Sun, Sun Microsystems, the Sun logo, AnswerBook2, docs.sun.com, Solaris, Sun HPC ClusterTools, Prism, Sun Performance Library, Sun Performance Workshop, Sun ONE, and RSM are trademarks or registered trademarks of Sun Microsystems, Inc. in the U.S. and in other countries.

All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. in the U.S. and in other countries. Products bearing SPARC trademarks are based upon an architecture developed by Sun Microsystems, Inc.

The OPEN LOOK and Sun™ Graphical User Interface was developed by Sun Microsystems, Inc. for its users and licensees. Sun acknowledges the pioneering efforts of Xerox in researching and developing the concept of visual or graphical user interfaces for the computer industry. Sun holds a non-exclusive license from Xerox to the Xerox Graphical User Interface, which license also covers Sun's licensees who implement OPEN LOOK GUIs and otherwise comply with Sun's written license agreements.

U.S. Government Rights—Commercial use. Government users are subject to the Sun Microsystems, Inc. standard license agreement and applicable provisions of the FAR and its supplements.

DOCUMENTATION IS PROVIDED "AS IS" AND ALL EXPRESS OR IMPLIED CONDITIONS, REPRESENTATIONS AND WARRANTIES, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE OR NON-INFRINGEMENT, ARE DISCLAIMED, EXCEPT TO THE EXTENT THAT SUCH DISCLAIMERS ARE HELD TO BE LEGALLY INVALID.

---

Copyright 2003 Sun Microsystems, Inc., 4150 Network Circle, Santa Clara, California 95054, Etats-Unis. Tous droits réservés.

Sun Microsystems, Inc. a les droits de propriété intellectuels relatants à la technologie incorporée dans le produit qui est décrit dans ce document. En particulier, et sans la limitation, ces droits de propriété intellectuels peuvent inclure un ou plus des brevets américains énumérés à <http://www.sun.com/patents> et un ou les brevets plus supplémentaires ou les applications de brevet en attente dans les Etats-Unis et dans les autres pays.

Ce produit ou document est protégé par un copyright et distribué avec des licences qui en restreignent l'utilisation, la copie, la distribution, et la décompilation. Aucune partie de ce produit ou document ne peut être reproduite sous aucune forme, par quelque moyen que ce soit, sans l'autorisation préalable et écrite de Sun et de ses bailleurs de licence, s'il y en a.

Le logiciel détenu par des tiers, et qui comprend la technologie relative aux polices de caractères, est protégé par un copyright et licencié par des fournisseurs de Sun.

Des parties de ce produit pourront être dérivées des systèmes Berkeley BSD licenciés par l'Université de Californie. UNIX est une marque déposée aux Etats-Unis et dans d'autres pays et licenciée exclusivement par X/Open Company, Ltd.

Sun, Sun Microsystems, le logo Sun, AnswerBook2, docs.sun.com, Solaris, Sun HPC ClusterTools, Prism, Sun Performance Library, Sun Performance Workshop, Sun ONE, et RSM sont des marques de fabrique ou des marques déposées de Sun Microsystems, Inc. aux Etats-Unis et dans d'autres pays.

Toutes les marques SPARC sont utilisées sous licence et sont des marques de fabrique ou des marques déposées de SPARC International, Inc. aux Etats-Unis et dans d'autres pays. Les produits portant les marques SPARC sont basés sur une architecture développée par Sun Microsystems, Inc.

L'interface d'utilisation graphique OPEN LOOK et Sun™ a été développée par Sun Microsystems, Inc. pour ses utilisateurs et licenciés. Sun reconnaît les efforts de pionniers de Xerox pour la recherche et le développement du concept des interfaces d'utilisation visuelle ou graphique pour l'industrie de l'informatique. Sun détient une licence non exclusive de Xerox sur l'interface d'utilisation graphique Xerox, cette licence couvrant également les licenciées de Sun qui mettent en place l'interface d'utilisation graphique OPEN LOOK et qui en outre se conforment aux licences écrites de Sun.

LA DOCUMENTATION EST FOURNIE "EN L'ÉTAT" ET TOUTES AUTRES CONDITIONS, DECLARATIONS ET GARANTIES EXPRESSES OU TACITES SONT FORMELLEMENT EXCLUES, DANS LA MESURE AUTORISEE PAR LA LOI APPLICABLE, Y COMPRIS NOTAMMENT TOUTE GARANTIE IMPLICITE RELATIVE A LA QUALITE MARCHANDE, A L'APTITUDE A UNE UTILISATION PARTICULIERE OU A L'ABSENCE DE CONTREFAÇON.



# Contents

---

**Preface** vii

**1. Introduction to Sun MPI** 1

Sun MPI Features 1

Sun MPI I/O 2

**2. Sun MPI Library** 3

Types of Libraries 3

Sun MPI Routines 4

Point-to-Point Communication Routines 4

One-Sided Communication Routines 5

Collective Communication Routines 12

Using the In-Place Option 13

Using Persistent Communication Requests 14

Managing Process Topologies 14

Name-Publishing Routines 14

Environmental Inquiry Routines 15

Packing and Unpacking Functions 15

Managing Communicators, Groups, and Contexts 15

Data Types 17

Resource Reservation for Batch Processing	20
Programming With Sun MPI	20
Fortran Support	21
Recommendations for All-to-All and All-to-One Communication	22
Signals and MPI	22
Multithreaded Programming	23
Guidelines for Thread-Safe Programming	23
MPI_Wait(), MPI_Waitall(), MPI_Waitany(), MPI_Waitsome()	24
MPI_Cancel()	24
MPI_Probe(), MPI_Iprobe()	24
Collective Calls	25
Communicator Operations	25
Error Handlers	25
Profiling Interface	26
MPE: Extensions to the Library	29
▼ To Obtain and Build the MPE	29
<b>3. Getting Started</b>	<b>31</b>
Header Files	31
Sample Code	32
Compiling and Linking	34
Choosing a Library Path	37
Stubbing Thread Calls	37
Profiling With <code>mpprof</code>	38
Basic Job Execution	38
Executing With CRE	39
Executing With LSF Suite	39
Debugging	39
Debugging With the Prism Environment	40

Starting Up Prism	40
Debugging With TotalView	41
Limitations	41
Related Documentation	42
Starting a New Job Using TotalView	42
Attaching to an mprun Job	44
Launching Sun MPI Batch Jobs Using TotalView	45
Debugging With MPE	47
<b>4. Programming With Sun MPI I/O</b>	<b>49</b>
Data Partitioning and Data Types	49
Definitions	50
Note for Fortran Users	51
Routines	51
File Manipulation	52
File Hints	52
File Views	53
Data Access	53
Data Access With Explicit Offsets	54
Data Access With Individual File Pointers	55
Pointer Manipulation	55
Data Access With Shared File Pointers	56
File Interoperability	57
File Consistency and Semantics	58
Sample Code	59
Partitioned Writing and Reading in a Parallel Job	59
Data Access Styles	63
Overlapping I/O With Computation and Communication	64

<b>A. Sun MPI and Sun MPI I/O Routines</b>	<b>67</b>
Sun MPI Routines	67
Sun MPI I/O Routines	95
<b>B. Environment Variables</b>	<b>103</b>
Informational Variables	103
General Performance Tuning	104
Tuning Memory for Point-to-Point Performance	106
Numerics	109
Synchronization of One-Sided Communications	109
Tuning Rendezvous	111
MPPProf	112
Miscellaneous	113
<b>C. Troubleshooting</b>	<b>117</b>
MPI Messages	117
MPI I/O Error Handling	120
Adjusting System V Shared-Memory Limits	122

# Preface

---

The *Sun MPI Software Programming and Reference Manual* describes the Sun™ MPI library of message-passing routines and explains how to develop an MPI program on a Sun HPC system.

For the most part, this guide does not repeat information that is available in detail in the *MPI Standard*; it focuses instead on what is specific to the Sun MPI implementation.

You should be familiar with programming in C or Fortran, with parallel programming, and with the message-passing model.

---

## Before You Read This Book

For general information about writing MPI programs, see “Related Documentation” on page x. Sun MPI is part of the Sun HPC ClusterTools™ suite of software. For more information about running Sun MPI jobs, see the *Sun HPC ClusterTools Software User's Guide*. Product notes for Sun MPI are included in *Sun HPC ClusterTools Software Release Notes*.

---

## Using UNIX Commands

This document does not describe how to use basic UNIX® commands. For that type of information, see:

- *Solaris Handbook for Sun Peripherals*
- Online documentation for the Solaris™ operating environment

- Other software documentation that you received with your system



---

# Typographic Conventions

Typeface*	Meaning	Examples
AaBbCc123	The names of commands, files, and directories; on-screen computer output	Edit your <code>.login</code> file. Use <code>ls -a</code> to list all files. % You have mail.
<b>AaBbCc123</b>	What you type, when contrasted with on-screen computer output	% <b>su</b> Password:
<i>AaBbCc123</i>	Book titles, new words or terms, words to be emphasized. Replace command-line variables with real names or values.	Read Chapter 6 in the <i>User's Guide</i> . These are called <i>class</i> options. You <i>must</i> be superuser to do this. To delete a file, type <code>rm filename</code> .

\* The settings on your browser might differ from these settings.

---

# Shell Prompts

Shell	Prompt
C shell	<i>machine-name%</i>
C shell superuser	<i>machine-name#</i>
Bourne shell and Korn shell	\$
Bourne shell and Korn shell superuser	#

---

## Related Documentation

Application	Title	Part Number
Sun HPC ClusterTools Documentation	<i>Read Me First: Guide to Sun HPC ClusterTools 5 Software Documentation</i>	817-0096-10
Sun HPC ClusterTools Software	<i>Sun HPC ClusterTools 5 Software Release Notes</i>	817-0081-10
	<i>Sun HPC ClusterTools 5 Software Installation Guide</i>	817-0082-10
	<i>Sun HPC ClusterTools 5 Software Performance Guide</i>	817-0090-10
	<i>Sun HPC ClusterTools 5 Software User's Guide</i>	817-0084-10
	<i>Sun HPC ClusterTools 5 Software Administrator's Guide</i>	817-0083-10
Sun S3L	<i>Sun S3L 4.0 Software Programming Guide</i>	817-0086-10
	<i>Sun S3L 4.0 Software Reference Manual</i>	817-0087-10
Prism™ graphical programming environment	<i>Prism 7.0 Software User's Guide</i>	817-0088-10
	<i>Prism 7.0 Software Reference Manual</i>	817-0089-10

A wealth of documentation on MPI is available on the World Wide Web. It includes the following:

- The MPI home page, with links to specifications for the *MPI-2 Standard*:  
<http://www.mpi-forum.org>
- Additional Web sites that provide links to papers, talks, the *MPI Standard*, implementations, information about MPI-2, tutorials, and pointers to many other sources:  
<http://www.erc.msstate.edu/mpi/>  
<http://www.arc.unm.edu/>

---

## Man Pages

Man pages are also available online for all the Sun MPI and MPI I/O routines and are accessible by means of the Solaris `man` command. These man pages are usually installed in `/opt/SUNWhpc/man`. Ask your system administrator for their location at your site.

---

## Accessing Sun Documentation

You can view, print, or purchase a broad selection of Sun documentation, including localized versions, at:

<http://www.sun.com/documentation>

---

## Sun Welcomes Your Comments

Sun is interested in improving its documentation and welcomes your comments and suggestions. You can email your comments to Sun at:

[docfeedback@sun.com](mailto:docfeedback@sun.com)

Please include the part number (817-0085-10) of your document in the subject line of your email.



# Introduction to Sun MPI

---

The MPI specification was developed by the MPI Forum, a group of software developers, computer vendors, academics, and computer-science researchers whose goal was to develop a standard for writing message-passing programs that would be efficient, flexible, and portable.

The outcome, known as the *MPI Standard*, was first published in 1993. The most recent version (MPI-2) was published in 1997. It was well received, and several implementations are available publicly.

Sun™ MPI is a complete library of message-passing routines, including all MPI 1.2-compliant and MPI 2-compliant routines. Chapter 2 provides an overview of the routines, Appendix A summarizes them, and the Sun `man` pages provide detailed descriptions.

---

## Sun MPI Features

- Optimization for running with Sun HPC ClusterTools 5 software using C, C++, Fortran 77, or Fortran 90.
- Integration with the Sun HPC ClusterTools™ Runtime Environment (CRE) and the Sun GridEngine, LSF, and PBS distributed resource-management systems.
- Support for multithreaded programming.
- Seamless use of various network protocols.
- Multiprotocol support so that MPI selects the fastest available medium for each type of connection (such as shared memory, RSM, and ATM).
- Communication by shared memory for fast performance on clusters of SMPs
- Finely tunable shared-memory communication.
- Optimized collectives for symmetric multiprocessors (SMPs) and clusters of SMPs.

- MPI I/O support for parallel file I/O.
- Prism support, so that users can develop, run, and debug programs in the Prism programming environment.
- Implicit co-scheduling, whereby the Sun HPC `spind` daemon enables certain processes of a given MPI job on a shared-memory system to be scheduled at approximately the same time as other related processes. This co-scheduling reduces the load on the processors, thus reducing the effect that MPI jobs have on each other.
- Complete support of one-sided communication routines and name-publishing routines.
- Dynamic library support.
- MPI-2 dynamic support.

Sun MPI and MPI I/O provide full F77, C, and C++ support, as well as Basic F90 support.

---

## Sun MPI I/O

File I/O in Sun MPI uses MPI 2-compliant routines for parallel file I/O. Chapter 4 describes these routines. Their `man` pages are provided online, and the routines are summarized in Appendix A.

# Sun MPI Library

---

This chapter describes the Sun MPI library:

- “Types of Libraries” on page 3
- “Sun MPI Routines” on page 4
- “Managing Communicators, Groups, and Contexts” on page 15
- “Data Types” on page 17
- “Resource Reservation for Batch Processing” on page 20
- “Programming With Sun MPI” on page 20
- “Multithreaded Programming” on page 23
- “Profiling Interface” on page 26
- “MPE: Extensions to the Library” on page 29

---

**Note** – Sun MPI I/O is described separately, in Chapter 4.

---

---

## Types of Libraries

Sun MPI contains four types of libraries, which represent two categories.

- *32- and 64-bit libraries* —If you want to take advantage of the 64-bit capabilities of Sun MPI, you must explicitly link to the 64-bit libraries. The 32-bit libraries are the default in each category.
- *Thread-safe and non-thread-safe libraries* —For multithreaded programs, link with the thread-safe library in the appropriate category unless the program has only one thread calling MPI. For programs that are not multithreaded, you can link against either the thread-safe or the default (non-thread-safe) library. However, nonmultithreaded programs have better performance using the default library, as it does not incur the extra overhead of providing thread safety. Therefore, use the default libraries whenever possible for maximum performance.

For full information about linking to libraries, see “Compiling and Linking” on page 34.

---

## Sun MPI Routines

This section gives a brief description of the routines in the Sun MPI library. All the Sun MPI routines are listed in Appendix A with brief descriptions and their C syntax. For detailed descriptions of individual routines, see the man pages or the *MPI Standard*. The routines are divided into these categories:

- “Point-to-Point Communication Routines” on page 4
- “One-Sided Communication Routines” on page 5
- “Collective Communication Routines” on page 12
- “Name-Publishing Routines” on page 14
- “Environmental Inquiry Routines” on page 15
- “Packing and Unpacking Functions” on page 15

### Point-to-Point Communication Routines

Point-to-point communication routines include the basic send and receive routines in both blocking and nonblocking forms and in four modes.

A *blocking send* blocks until its message buffer can be written with a new message.

A *blocking receive* blocks until the received message is in the receive buffer.

*Nonblocking sends and receives* differ from blocking sends and receives in that they return immediately and their completion must be waited or tested for. It is expected that eventually nonblocking send and receive calls will allow the overlap of communication and computation.

The four modes for MPI point-to-point communication are as follows:

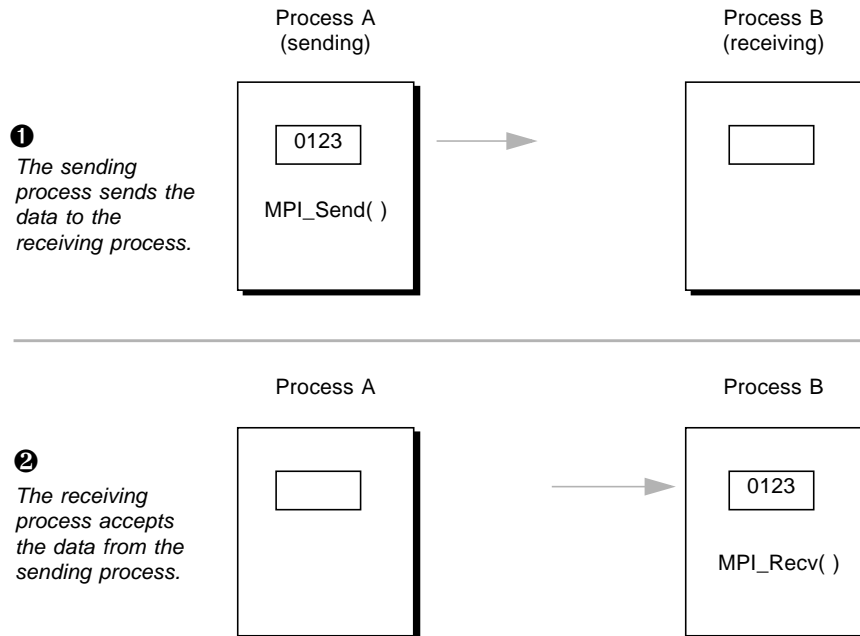
- *Standard* – The completion of a send implies that the message either is buffered internally or has been received. Users are free to overwrite the buffer that they passed in with any of the blocking send or receive routines, after the routine returns.
- *Buffered* – The user guarantees a certain amount of buffering space.
- *Synchronous* – Rendezvous semantics occur between sender and receiver; that is, a send blocks until the corresponding receive has occurred.



- *Ready* – A send can be started only if the matching receive is already posted. The ready mode for sends is a way for the programmer to notify the system that the receive has been posted, so that the underlying system can use a faster protocol if it is available.

## One-Sided Communication Routines

Standard MPI communication is two-sided. To complete a transfer of information, both the sending and receiving processes must call appropriate functions. The operation proceeds in two stages, as shown in the following figure.



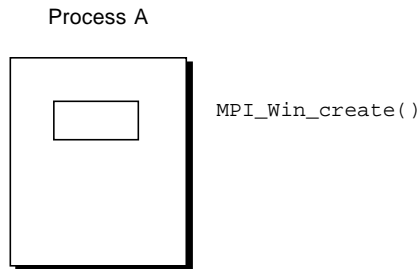
This form of communication requires regular synchronization between the sending and receiving processes. That synchronization can become complicated if the receiving process does not know which process is sending it the data it needs. One-sided communication was developed to solve this problem and to reduce the amount of synchronization required even when both sending and receiving processes know each other's identities.

In one-sided communication, a process opens a window in memory and exposes it to all processes that belong to a particular communicator, provided they reside on the

same node. As long as that window is open, any process in the communicator and node can *put* data into it and *get* data out of it.

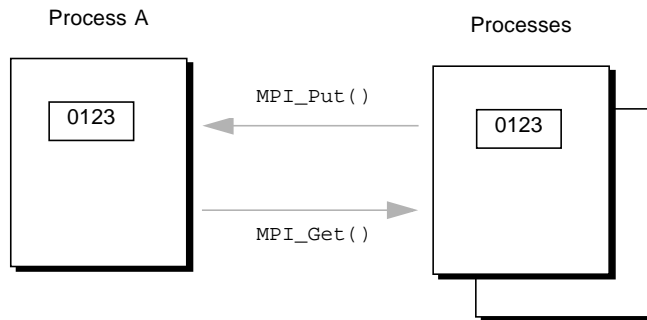
❶

*A process opens a communications window and exposes it to the other processes in the same node and communicator.*



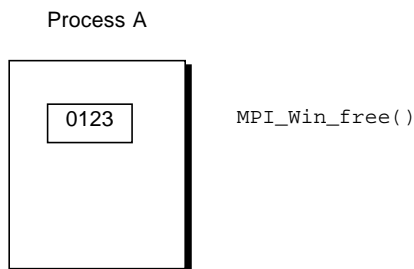
❷

*Any process within the same communicator and node can transfer data directly into or out of that window as long as it is open.*



❸

*Then the original process closes the window.*



The *put* requires no complementary operation from the process that opened the window, and is equivalent to the combination of a *send* and *receive* operation in two-

sided MPI communication.

The functions used to implement one-sided MPI communication fall into three categories and are summarized in TABLE 2-1. You can find their definitions in the *MPI Standard*. Also, Appendix A of this document provides syntax summaries.

**TABLE 2-1** One-Sided Communication Routines

<b>Window Creation</b>	
<code>MPI_Win_create()</code>	Creates a window in memory and exposes it to all processes in the communicator and node.
<code>MPI_Win_free()</code>	Closes the window created with <code>MPI_Win_create()</code> . Requires barrier synchronization.
<code>MPI_Win_get_group()</code>	Returns a duplicate of the group of the communicator used to create the window.
<b>Data Transfer</b>	
<code>MPI_Accumulate()</code>	Combines data from a process with data already in the window. Different from <code>MPI_Put()</code> in that new data is appended to existing data instead of replacing it.
<code>MPI_Get()</code>	The calling process takes data directly from the window. The opposite of <code>MPI_Put()</code> . Equivalent to the combination of a send and receive operation originated by the receiving process.
<code>MPI_Put()</code>	The calling process loads its data directly into the buffer of the target process. Equivalent to the combination of a send and receive operation originated by the sending process. The opposite of <code>MPI_Get()</code> .
<b>Synchronization</b>	
<code>MPI_Win_fence()</code>	Blocks any process from operating on a particular window until all operations relating to that window have completed. Similar to <code>MPI_Barrier()</code> , but applies to a window instead of a communicator.
<code>MPI_Win_lock()</code>	Starts an RMA access epoch. While the lock is in place, only the process whose rank is specified in the function call can be accessed by RMA operations on the window.
<code>MPI_Win_unlock()</code>	Closes an RMA access epoch begun with a call to <code>MPI_Win_lock()</code> .
<code>MPI_Win_start()</code>	Starts an RMA access epoch for <i>window</i> .
<code>MPI_Win_complete()</code>	Completes an RMA access epoch on <i>window</i> started by a call to <code>MPI_Win_start()</code> .

**TABLE 2-1** One-Sided Communication Routines (*Continued*)

---

<code>MPI_Win_post()</code>	Starts an RMA exposure epoch for the local window associated with <i>window</i> .
<code>MPI_Win_wait()</code>	Completes an RMA exposure epoch started by a call to <code>MPI_Win_post()</code> on <i>window</i> .
<code>MPI_Win_test()</code>	Attempts to complete an RMA exposure epoch; a nonblocking version of <code>MPI_Win_wait()</code> .

---

Some special considerations apply to allocating memory for one-sided communications. For example:

- `MPI_Alloc_mem` allocates memory by means of SysV shared memory. If you find that you cannot allocate any more memory because of system-imposed limits, you can either increase `shmsys:shminfo_shmseg` in `/etc/system` or try to preallocate a large segment and then use parts of it for your needs.
- If the memory used for a communication window has been allocated by means of the `MPI_Alloc_mem` function, the memory region specified by the base and size parameters passed to `MPI_Win_create` should not exceed the limits of the memory region allocated using a call either to `malloc` or to the `MPI_Alloc_mem` function. If the memory has been allocated using a call to `MPI_Alloc_mem`, these limits are checked internally and the call to `MPI_Win_create()` returns `MPI_ERR_OTHER`.

Several one-sided communications routines support info keys. These keys, and their descriptions, are listed in TABLE 2-2.

**TABLE 2-2** Info Keys Supported By One-Sided Communications Routines

<b>Routine</b>	<b>Info Key</b>	<b>Description</b>
MPI_Win_create()	no_locks	If MPI_Win_lock is called on a window created with this info key, the call fails. If this info key is present, it is assumed that the local window is never locked, allowing several internal checks to be skipped, permitting a more efficient implementation.
	sun_noexpose	An info key (unique to Sun MPI) that is interpreted only by the remote shared memory (RSM™) protocol module when passed using MPI_Win_create (and while the RSM protocol is in use). When explicitly passed with a base address that was not allocated using MPI_Alloc_mem, sun_noexpose notifies the library not to expose address spaces not specified by the application.
	sun_rsmrexit_hw	An info key (unique to Sun MPI) that is interpreted only by the RSM protocol module when passed using MPI_Win_create (and while the RSM protocol is in use). The sun_rsmrexit_hw info key defines a high-water mark, in bytes, of data that will be sent before any error checking and retransmission of data is done. By default the RSM protocol module assumes the high-water mark to be 16KB.
	sun_shmeuid	The effective user ID (UID) to which the shared-memory segment is set for memory-based protocols, including shared memory (SHM) and RSM. Set this key only when you anticipate connections from programs run by other users. Valid only for server programs run as root.

**TABLE 2-2** Info Keys Supported By One-Sided Communications Routines (*Continued*)

<b>Routine</b>	<b>Info Key</b>	<b>Description</b>
	sun_shmegid	The effective GID to which the shared-memory segment is set for memory-based protocols, including SHM and RSM. Set this key only when you anticipate connections from programs run by other users.
	sun_shmperm	The permissions to which the shared-memory segment is set, in octal, for memory-based protocols, including SHM and RSM. Set this key only when you anticipate connections from programs run by other users.
MPI_Alloc_mem	sun_shmeuid	The effective UID to which the shared-memory segment is set for memory-based protocols, including SHM and RSM. Set this key only when you anticipate connections from programs run by other users. Valid only for server programs run as <code>root</code> .
	sun_shmegid	The effective GID to which the shared-memory segment is set for memory-based protocols, including SHM and RSM. Set this key only when you anticipate connections from programs run by other users.
	sun_shmperm	The permissions to which the shared-memory segment is set, in octal, for memory-based protocols, including SHM and RSM. Set this key only when you anticipate connections from programs run by other users.

Several one-sided communications routines support assertions. These assertions, and their descriptions, are listed in TABLE 2-3.

**TABLE 2-3** Assertions Supported by One-Sided Communications Routines

<b>Routine</b>	<b>Assertion Values</b>	<b>Description</b>
<code>MPI_Win_fence()</code>	<code>MPI_MODE_NOPRECEDE</code>	The RSM protocol module recognizes this value and does not issue any RSM close barriers on the window passed into the <code>MPI_Win_fence()</code> call. The effect of this assertion value is that fence calls execute faster.
	<code>MPI_MODE_NOSTORE</code> <code>MPI_MODE_NOPUT</code> <code>MPI_MODE_NOSUCCEED</code>	All native Sun HPC ClusterTools protocol modules and generic one-sided functions ignore these values. However, they are available for use by third-party protocol modules
<code>MPI_Win_start()</code>	<code>MPI_MODE_NOCHECK</code>	When this value is passed in to this call, the library assumes that the post call on the target has been called and it is not necessary for the library to check to see if such a call has been made.
<code>MPI_Win_post()</code>	<code>MPI_MODE_NOCHECK</code>	When this value is passed in to this call, the library assumes that the post call on the target has been called and it is not necessary for the library to check to see if such a call has been made.

**TABLE 2-3** Assertions Supported by One-Sided Communications Routines (*Continued*)

<b>Routine</b>	<b>Assertion Values</b>	<b>Description</b>
<code>MPI_Win_lock()</code>	<code>MPI_MODE_NOCHECK</code>	The Sun MPI library supports <code>MPI_MODE_NOCHECK</code> for the RSM protocol module only. When this lock is set, the lock and unlock calls in the RSM protocol module do not acquire the lock. However, data synchronization does occur.

## Collective Communication Routines

Collective communication routines are blocking routines that involve all processes in a communicator and, in most cases, an intercommunicator. Collective communication includes broadcasts and scatters, reductions and gathers, all-gathers and all-to-alls, scans, and a synchronizing barrier call.

**TABLE 2-4** Collective Communication Routines

<b>Routine</b>	<b>Description</b>
<code>MPI_Bcast()</code>	Broadcasts from one process to all others in a communicator or intercommunicator.
<code>MPI_Scatter()</code>	Scatters from one process to all others in a communicator or intercommunicator.
<code>MPI_Scatterv()</code>	Scatters from all processes to all others in a communicator or intercommunicator.
<code>MPI_Reduce()</code>	Reduces from all to one in a communicator or intercommunicator.
<code>MPI_Allreduce()</code>	Reduces and then broadcasts result to all nodes in a communicator or intercommunicator.
<code>MPI_Allreducev()</code>	Reduces from all processes and then broadcasts result to all nodes in a communicator or intercommunicator.
<code>MPI_Reduce_scatter()</code>	Scatters a vector that contains results across the nodes in a communicator.
<code>MPI_Gather()</code>	Gathers from all to one in a communicator or intercommunicator.
<code>MPI_Gatherv()</code>	Gathers information from all processes in a communicator or intercommunicator.



**TABLE 2-4** Collective Communication Routines (*Continued*)

Routine	Description
<code>MPI_Allgather()</code>	Gathers and then broadcasts the results of the gather in a communicator or intercommunicator.
<code>MPI_Allgatherv()</code>	Gathers from all processes and then broadcasts the results of the gather in a communicator or intercommunicator.
<code>MPI_Alltoall()</code>	All processes send data to, and receive data from, all other processes in a communicator or intercommunicator.
<code>MPI_Alltoallv()</code>	Like <code>MPI_Alltoall()</code> , but user can use vector style (displacement and element count) to specify what data to send and receive.
<code>MPI_Alltoallw()</code>	Like <code>MPI_Alltoallv()</code> , but user can specify database of individual datablocks, in addition to displacement and element count.
<code>MPI_Scan()</code>	Scans (performs a parallel prefix) across processes in a communicator or intercommunicator.
<code>MPI_Exscan()</code>	Performs an exclusive prefix reduction on data distributed across the calling processes.
<code>MPI_Barrier()</code>	Synchronizes processes in a communicator or intercommunicator (no data is transmitted).

Many of the collective communication calls have alternative vector forms, with which various amounts of data can be sent to or received from various processes. In addition, `MPI_Alltoallw()` accepts a database of individual datablocks.

The syntax and semantics of these routines are basically consistent with the point-to-point routines (upon which they are built), but there are restrictions to keep them from becoming too complex:

- The amount of data sent must exactly match the amount of data specified by the receiver.
- There is only one mode, a mode analogous to the standard mode of point-to-point routines.

## Using the In-Place Option

Several collectives can pass `MPI_IN_PLACE` as the value of *send-buffer* at the root. When they do, *sendcount* and *sendtype* are ignored, and the contribution of the root to the gathered vector is assumed to be already in the correct location in the receive buffer. The collectives are as follows:

- `MPI_Gather()`
- `MPI_Gatherv()`
- `MPI_Scatter()`
- `MPI_Scatterv()`

- `MPI_Allgather()`
- `MPI_Allgatherv()`
- `MPI_Reduce()`
- `MPI_AllReduce()`
- `MPI_Reduce_scatter()`
- `MPI_Scan()`

## Using Persistent Communication Requests

Sometimes within an inner loop of a parallel computation, a communication with the same argument list is executed repeatedly. The communication can be slightly improved by using a *persistent* communication request, which reduces the overhead for communication between the process and the communication controller. A persistent request can be thought of as a communication port or *half-channel*.

## Managing Process Topologies

Process topologies are associated with communicators; they are optional attributes that can be given to an intracommunicator (not to an intercommunicator).

Recall that processes in a group are ranked from 0 to  $n-1$ . This linear ranking often reflects nothing of the logical communication pattern of the processes, which may be, for instance, a two- or three-dimensional grid. The logical communication pattern is referred to as a *virtual topology* (separate and distinct from any hardware topology). In MPI, two types of virtual topologies can be created: Cartesian (grid) topology and graph topology.

You can use virtual topologies in your programs by taking physical processor organization into account to provide a ranking of processors that optimizes communication.

## Name-Publishing Routines

Name-publishing routines enable client applications to retrieve system-supplied port names. A server calls the `MPI_Publish_name()` function to publish the name of the service associated with a particular port name. A client application calls the `MPI_Lookup_name()`, passing it the published service name, and in return gets its associated port name. The server can also call the `MPI_Unpublish_name()` function to stop publishing names.

Sun's implementation of the *MPI Standard* does not provide a *scope* for the published names and does not allow a server to publish the same service name twice. The implementation consists of three routines:

- `MPI_Publish_name()`
- `MPI_Unpublish_name()`
- `MPI_Lookup_name()`

## Environmental Inquiry Routines

Environmental inquiry routines are used for starting up and shutting down error-handling routines and timers.

Few MPI routines can be called before `MPI_Init()` or after `MPI_Finalize()`. Examples include `MPI_Initialized()` and `MPI_Version()`. `MPI_Finalize()` can be called only if there are no outstanding communications involving that process.

The set of errors handled by MPI depends upon the implementation. See Appendix C for tables listing the Sun MPI error classes.

## Packing and Unpacking Functions

Sun's implementation of the *MPI Standard* provides functions for packing and unpacking messages to be exchanged within an MPI implementation, and in the external32 format used to exchange messages between MPI implementations.

- `MPI_Pack()`
- `MPI_Unpack()`
- `MPI_Pack_size()`
- `MPI_Pack_external()`
- `MPI_Unpack_external()`
- `MPI_Pack_external_size()`

---

## Managing Communicators, Groups, and Contexts

A distinguishing feature of the *MPI Standard* is that it includes a mechanism for creating separate worlds of communication, accomplished through *communicators*, *groups*, and *contexts*.

A *communicator* specifies a group of processes that will conduct communication operations within a specified context without affecting or being affected by operations occurring in other groups or contexts elsewhere in the program. A communicator also guarantees that, within any group and context, point-to-point and collective communication are isolated from each other.

A *group* is an ordered collection of processes. Each process has a rank in the group; the rank runs from 0 to  $n-1$ . A process can belong to more than one group; its rank in one group has nothing to do with its rank in any other group.

A *context* is the internal mechanism by which a communicator guarantees safe communication space to the group.

At program startup, two default communicators are defined:

- `MPI_COMM_WORLD`, which has as a process group all the processes of the job
- `MPI_COMM_SELF`, which is equivalent to an identity communicator

The process group that corresponds to `MPI_COMM_WORLD` is not predefined, but can be accessed using `MPI_COMM_GROUP`. One `MPI_COMM_SELF` communicator is defined for each process, each of which has rank zero in its own communicator. For many programs, these are the only communicators needed.

Communicators are of two kinds: *intra*communicators, which conduct operations within a given group of processes; and *inter*communicators, which conduct operations between two groups of processes.

Communicators provide a *caching* mechanism, which allows an application to attach attributes to communicators. Attributes can be user data or any other kind of information.

New groups and new communicators are constructed from existing ones. Group constructor routines are local, and their execution does not require interprocessor communication. Communicator constructor routines are collective, and their execution can require interprocess communication.

You can also create an intercommunicator from two MPI processes that are connected by a socket. Use the `MPI_Comm_join()` function.

---

**Note** – Users who do not need any communicator other than the default `MPI_COMM_WORLD` communicator—that is, who do not need any sub- or supersets of processes—can plug in `MPI_COMM_WORLD` wherever a communicator argument is requested. In these circumstances, users can ignore this section and the associated routines. (These routines can be identified from the listing in Appendix A.)

---

---

# Data Types

All Sun MPI communication routines have a data type argument. They can be primitive data types, such as integers or floating-point numbers, or they can be user-defined, derived data types that are specified in terms of primitive types.

Derived data types enable users to specify more general, mixed, and noncontiguous communication buffers, such as array sections and structures that contain combinations of primitive data types.

Fortran data types are listed in TABLE 2-5. Data types of Fortran used with the `-r8` flag are listed in TABLE 2-6. C data types are listed in TABLE 2-7.

**TABLE 2-5** Fortran Data Types

<b>MPI Data Type</b>	<b>Fortran Data Type</b>
MPI_INTEGER	INTEGER INTEGER*4
MPI_INTEGER1	INTEGER*1 (Fortran 90 only)
MPI_INTEGER2	INTEGER*2
MPI_INTEGER4	INTEGER*4
MPI_INTEGER8	INTEGER*8
MPI_REAL	REAL REAL*4
MPI_REAL4	REAL*4
MPI_REAL8	REAL*8
MPI_REAL16	REAL*16
MPI_DOUBLE_PRECISION	REAL*8 DOUBLE PRECISION
MPI_2DOUBLE_PRECISION	Pair of DOUBLE PRECISION variables*
MPI_COMPLEX	COMPLEX
MPI_LOGICAL	LOGICAL
MPI_CHARACTER	CHARACTER(1)
MPI_SIGNED_CHAR	INTEGER*1

**TABLE 2-5** Fortran Data Types (Continued)

MPI Data Type	Fortran Data Type
MPI_DOUBLE_COMPLEX	DOUBLE COMPLEX
MPI_2REAL	Pair of REALS*
MPI_INTEGER2	INTEGER*2
MPI_INTEGER4	INTEGER*4
MPI_2INTEGER	Pair of INTEGERS*
MPI_BYTE	no corresponding Fortran data type
MPI_PACKED	no corresponding Fortran data type

\* For use with MINLOC and MAXLOC

**TABLE 2-6** Fortran -r8 Data Types

MPI Data Type	Fortran -r8 Data Type
MPI_INTEGER	INTEGER*4
MPI_INTEGER1	INTEGER*1 (Fortran 90 only)
MPI_INTEGER2	INTEGER*2
MPI_INTEGER4	INTEGER*4
MPI_INTEGER8	INTEGER INTEGER*8
MPI_REAL	REAL*4
MPI_REAL4	REAL*4
MPI_REAL8	REAL REAL*8
MPI_REAL16	REAL*16 DOUBLE PRECISION
MPI_DOUBLE_PRECISION	REAL REAL*8
MPI_2DOUBLE_PRECISION	Pair of REAL*8*
MPI_COMPLEX	COMPLEX*4
MPI_LOGICAL	LOGICAL
MPI_CHARACTER	CHARACTER(1)
MPI_SIGNED_CHAR	INTEGER*1

**TABLE 2-6 Fortran -r8 Data Types (Continued)**

<b>MPI Data Type</b>	<b>Fortran -r8 Data Type</b>
MPI_DOUBLE_COMPLEX	COMPLEX
MPI_2REAL	Pair of REAL*4*
MPI_INTEGER2	INTEGER*2
MPI_INTEGER4	INTEGER*4
MPI_2INTEGER	Pair of INTEGER*4
MPI_BYTE	No corresponding Fortran data type
MPI_PACKED	No corresponding Fortran data type

\* For use with MINLOC and MAXLOC

**TABLE 2-7 C Data Types**

<b>MPI Data Type</b>	<b>C Data Type</b>
MPI_BYTE	No corresponding C data type
MPI_PACKED	No corresponding C data type
MPI_CHAR	signed char
MPI_SIGNED_CHAR	signed char
MPI_UNSIGNED_CHAR	unsigned char
MPI_SHORT	signed short int
MPI_UNSIGNED_SHORT	unsigned short int
MPI_INT	signed int
MPI_UNSIGNED	unsigned int
MPI_LONG	signed long int
MPI_UNSIGNED_LONG	unsigned long int
MPI_LONG_LONG_INT	long long int
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double
MPI_WCHAR	wchar_t
MPI_UNSIGNED_LONG_LONG	unsigned long long int
MPI_2INT	Pair of int*

**TABLE 2-7** C Data Types (Continued)

MPI Data Type	C Data Type
MPI_FLOAT_INT	float and int*
MPI_DOUBLE_INT	double and int*
MPI_LONG_DOUBLE_INT	long double and int*
MPI_LONG_INT	long and int*
MPI_SHORT_INT	short and int*

\* For use with MINLOC and MAXLOC

---

## Resource Reservation for Batch Processing

If you plan to launch a job that uses the `MPI_Comm_spawn` or `MPI_Comm_spawn_multiple` functions, you must first reserve the resources with the resource manager that will run the job. As explained in the `mprun.1` man page and the *Sun HPC ClusterTools Software User's Guide*, you can reserve those resources by adding the `-nr` flag to the `mprun` command.

When you launch a job with the `mprun` command from within a resource manager, the number of processes allocated for that job are stored in the environment variable `MPI_UNIVERSE_SIZE`. It is the sum of the processes allocated with the `mprun` command's `-np` flag, and reserved with its `-nr` flag.

---

## Programming With Sun MPI

Although there are about 190 non-I/O routines in the Sun MPI library, you can write programs for a wide range of problems using only six routines, as described in TABLE 2-8.



**TABLE 2-8** Six Basic MPI Routines

Routine	Description
<code>MPI_Init()</code>	Initializes the MPI library.
<code>MPI_Finalize()</code>	Finalizes the MPI library. This includes releasing resources used by the library.
<code>MPI_Comm_size()</code>	Determines the number of processes in a specified communicator.
<code>MPI_Comm_rank()</code>	Determines the rank of calling process within a communicator.
<code>MPI_Send()</code>	Sends a message.
<code>MPI_Recv()</code>	Receives a message.

This set of six routines includes the basic send and receive routines. Programs that depend heavily on collective communication might also include `MPI_Bcast()` and `MPI_Reduce()`.

The functionality of these routines means you can have the benefit of parallel operations without having to learn the whole library at once. As you become more familiar with programming for message passing, you can start learning the more complex and esoteric routines and add them to your programs as needed.

See Appendix A for a complete list of Sun MPI routines.

## Fortran Support

Sun MPI provides extended Fortran support, as described in Section 10.2 of the *MPI-2 Standard*. In other words, it provides basic Fortran support, plus additional functions that specifically support Fortran 90:

- `MPI_Type_create_f90_complex()`
- `MPI_Type_create_f90_integer()`
- `MPI_Type_create_f90_real()`
- `MPI_Type_match_size()`
- `MPI_Sizeof()`

Basic Fortran support provides the original Fortran bindings and an `mpif.h` file specified in the *MPI-1 Standard*. The `mpif.h` file is valid for both fixed- and free-source forms, as specified in the *MPI-2 Standard*.

The MPI interface is known to violate the Fortran standard in several ways, but it causes few problems for FORTRAN 77 programs. Violations of the standard can cause more significant problems for Fortran 90 programs, however, if you do not follow the guidelines recommended in the standard. If you are programming in

Fortran, and particularly if you are using Fortran 90, you should consult Section 10.2 of the *MPI-2 Standard* for detailed information about basic Fortran support in an MPI implementation.

## Recommendations for All-to-All and All-to-One Communication

The Sun MPI library uses the TCP protocol to communicate over a variety of networks. MPI depends on TCP to ensure reliable, correct data flow. TCP's reliability compensates for unreliability in the underlying network, as the TCP retransmission algorithms handle any segments that are lost or corrupted. In most cases, this works well with good performance characteristics. However, when doing all-to-all and all-to-one communication over certain networks, a large number of TCP segments can be lost, resulting in poor performance.

You can compensate for this diminished performance over TCP in these ways:

- When writing your own algorithms, avoid flooding one node with a lot of data.
- If you need to do all-to-all or all-to-one communication, use one of the Sun MPI routines to do so. They are implemented in a way that avoids congesting a single node with lots of data. The following routines fall into this category:
  - `MPI_Alltoall()`, `MPI_Alltoallv()`, and `MPI_Alltoallw()` – These have been implemented using a pairwise communication pattern, so that every rank is communicating with only one other rank at a given time.
  - `MPI_Gather()` and `MPI_Gatherv()` – The root process sends ready-to-send packets to each nonroot-rank process to tell the processes to send their data. In this way, the root process can regulate how much data it is receiving at any one time. Using this ready-to-send method is associated with a minor performance cost, however. For this reason, you can override this method by setting the `MPI_TCPSAFEGATHER` environment variable to 0. (See Appendix B for information about environment variables.)

## Signals and MPI

When running the MPI library over TCP, nonfatal `SIGPIPE` signals can be generated. To handle them, the library sets the signal handler for `SIGPIPE` to `ignore`, overriding the default setting (terminate the process). In this way the MPI library can recover in certain situations. You should therefore avoid changing the `SIGPIPE` signal handler.

The Sun MPI Fortran and C++ bindings are implemented as wrappers on top of the C bindings. The profiling interface is implemented using weak symbols. This means a profiling library need contain only a profiled version of C bindings.

The SIGPIPERs can occur when a process first starts communicating over TCP. This happens because the MPI library creates connections over TCP only when processes actually communicate with one another. There are some unavoidable conditions where SIGPIPERs can be generated when two processes establish a connection. If you want to avoid any SIGPIPERs, set the environment variable `MPI_FULLCONNINIT`, which creates all connections during `MPI_Init()` and avoids any situations that might generate a SIGPIPER. For more information about environment variables, see Appendix B.

---

## Multithreaded Programming

When you are linked to one of the thread-safe libraries, Sun MPI calls are thread safe, in accordance with basic tenets of thread safety for MPI mentioned in the MPI-2 specification. As a result:

- When two concurrently running threads make MPI calls, the outcome is as if the calls executed in some order.
- Blocking MPI calls block the calling thread only. A blocked calling thread does not prevent progress of other runnable threads on the same process, nor does it prevent them from executing MPI calls. Thus, multiple sends and receives are concurrent.

Use `MPI_Init_thread()` in place of `MPI_Init()` to initialize the MPI execution environment with a predetermined level of thread support. Use the `MPI_Is_thread_main()` function to find out whether a thread is the one that called `MPI_Init_thread()`.

## Guidelines for Thread-Safe Programming

Each thread within an MPI process can issue MPI calls; however, threads are not separately addressable. That is, the rank of a send or receive call identifies a process, not a thread, which means that no order is defined for the case in which two threads call `MPI_Recv()` with the same tag and communicator. Such threads are said to be *in conflict*.

If threads within the same application post conflicting communication calls, data races will result. You can prevent such data races by using distinct communicators or tags for each thread.

In general, adhere to these guidelines:

- You must not have a request serviced by more than one thread. Although you can have an operation posted in one thread and then completed in another, you cannot have the operation completed in more than one thread.
- A data type or communicator must not be freed by one thread while it is in use by another thread.
- Once `MPI_Finalize()` is called, subsequent calls in any thread will fail.
- You must ensure that a sufficient number of lightweight processes (LWPs) are available for your multithreaded program. Failure to do so can degrade performance or even result in deadlock.
- You cannot stub the thread calls in your multithreaded program by omitting the threads libraries in the link line. The `libmpi.so` library automatically calls in the threads libraries, which effectively override any stubs.

The following sections describe more specific guidelines that apply for some routines. They also include some general considerations for collective calls and communicator operations that you should be aware of.

`MPI_Wait()`, `MPI_Waitall()`, `MPI_Waitany()`,  
`MPI_Waitsome()`

In a program in which two or more threads call one of these routines, you must ensure that they are not waiting for the same request. Similarly, the same request cannot appear in the array of requests of multiple concurrent wait calls.

`MPI_Cancel()`

One thread must not cancel a request while that request is being serviced by another thread.

`MPI_Probe()`, `MPI_Iprobe()`

A call to `MPI_Probe()` or `MPI_Iprobe()` from one thread on a given communicator should not have a source rank and tags that match those of any other probes or receives on the same communicator. Otherwise, correct matching of message to probe call might not occur.

## Collective Calls

Collective calls are matched on a communicator according to the order in which the calls are issued at each processor. All the processes on a given communicator must make the same collective call. You can avoid the effects of this restriction on the threads on a given processor by using a different communicator for each thread.

No process that belongs to the communicator may omit making a particular collective call; that is, none should be left “dangling.”

## Communicator Operations

Each of the communicator (or intercommunicator) functions operates simultaneously with each of the noncommunicator functions, regardless of what the parameters are and whether the functions are on the same or different communicators. However, if you are using multiple instances of the same communicator function on the same communicator where all parameters are the same, it cannot be determined which threads belong to which resultant communicator. Therefore, when concurrent threads issue such calls, you must ensure that the calls are synchronized in such a way that threads in separate processes participating in the same communicator operation are grouped together. Do this either by using a different base communicator for each call or by making the calls in single-thread mode before actually using them within the separate threads.

Note also these special situations:

- If you are using multiple instances of the same function with differing parameters and multiple threads, you must use separate communicators.
- When using splits with multiple instances of the same function with the same parameters, but with different threads at the split, you must use separate communicators.

For example, you might want to produce several communicators in separate sets of threads by performing `MPI_Comm_split()` on a base communicator. To ensure proper, thread-safe operation, you should replicate the base communicator with `MPI_Comm_dup()` (in the root thread or in one thread) and then perform `MPI_Comm_split()` on the resulting duplicate communicators.

- Do not free a communicator in one thread if it is still being used by another thread.

## Error Handlers

When an error occurs as a result of an MPI call, the handler might not run on the same thread as the thread that made the error-raising call. In other words, you cannot assume that the error handler will execute in the local context of the thread

that made the error-raising call. The error handler can be executed by another thread on the same process, distinct from the one that returns the error code. Therefore, you cannot rely on local variables for error handling in threads; instead, use global variables from the process.

---

## Profiling Interface

The Sun HPC ClusterTools software suite includes MPPProf, a profiling tool to be used with applications that call Sun MPI library routines. When enabled, MPPProf collects information about a program's message-passing activities in a set of intermediate files, one file per MPI process. Once the information is collected, you can invoke the MPPProf command-line utility `mppprof`, which generates a report based on the profiling data stored in the intermediate files. You must enable MPPProf before starting an MPI program. You do this by setting the environment variable `MPI_PROFILE` to 1.

If MPPProf is enabled, it creates and initializes the intermediate files with header information when the program's `MPI_Init` call ends. It also creates an index file that contains a map of the intermediate files. `mppprof` uses this index file to find the intermediate files.

`mppprof` includes an interface for interacting with loadable protocol modules (loadable PMs). If an MPI program uses a loadable PM, this interface allows MPPProf to collect profiling data that is specific to loadable PM activities.

An `mppprof` report contains the following classes of performance information:

- The percentage of total execution time spent in MPI calls across all processes
- The percentage of time each process spent in MPI calls
- The number of calls, amount of time spent, and number of bytes sent or received per MPI routine, averaged over all processes, with percent variation among processes
- Connectivity statistics (message count and volume) between processor pairs
- The settings of environment variables that have performance implications

You can control aspects of `mppprof` behavior with the following environment variables:

- `MPI_PROFINTERVAL` – Use this environment variable to specify a data sampling period. When this value is greater than 0, a sequence of snapshots is recorded at the prescribed intervals. Each snapshot represents the MPI activity that occurred since the previous snapshot. The default behavior is to record a single snapshot at the time of the `MPI_Finalize` call.
- `MPI_PROFDATADIR` – Use this environment variable to specify the location where the intermediate files created for each process rank will be stored.

- `MPI_PROFINDEXDIR` – Use this environment variable to specify the location where the index file created for each profiled job will be stored.
- `MPI_PROFMAXFILESIZE` – Use this environment variable to specify the maximum size of intermediate files.

The Sun HPC ClusterTools software suite also provides a conversion utility, `mpdump`, which converts the data from each intermediate file into a raw (unevaluated) user-readable format. You can use the ASCII files generated by `mpdump` as input to a report generator of your choice.

Once you’ve enabled MPPProf profiling by setting `MPI_PROFILE` to 1 (and run a job using `mprun`) you will find a file in your working directory of the form

```
mpprof.index.rm.jid
```

Type

```
% mpprof mpprof.index.rm.jid
```

to view the profiling report.

Further instructions for using `mpprof` and `mpdump` are provided in the *Sun HPC ClusterTools Software User’s Guide*.

Sun MPI meets the profiling interface requirements described in Chapter 8 of the *MPI-1 Standard*. This means you can write your own profiling library or choose from a number of available profiling libraries, such as those in the multiprocessing environment (MPE) from Argonne National Laboratory. (See “MPE: Extensions to the Library” on page 29 for more information.) The *User’s Guide for mpich, a Portable Implementation of MPI* includes more detailed information about using profiling libraries.

FIGURE 2-1 provides a generic illustration of how the software fits together. In this example, the user is linking against a profiling library that collects information on `MPI_Send()`. No profiling information is being collected for `MPI_Recv()`.

C profiling interfaces are needed even for Fortran programs. If there is profiling for both the Fortran and C version of an MPI function, then a Fortran call will encounter both profilings.

Be sure you make the library dynamic. A static library can experience the linker problems described in Section 8.4.3 of the *MPI 1.1 Standard*.

For compiling the program, the user’s link line would look like this:

```
# cc ..... -llibrary-name -lmpi
```

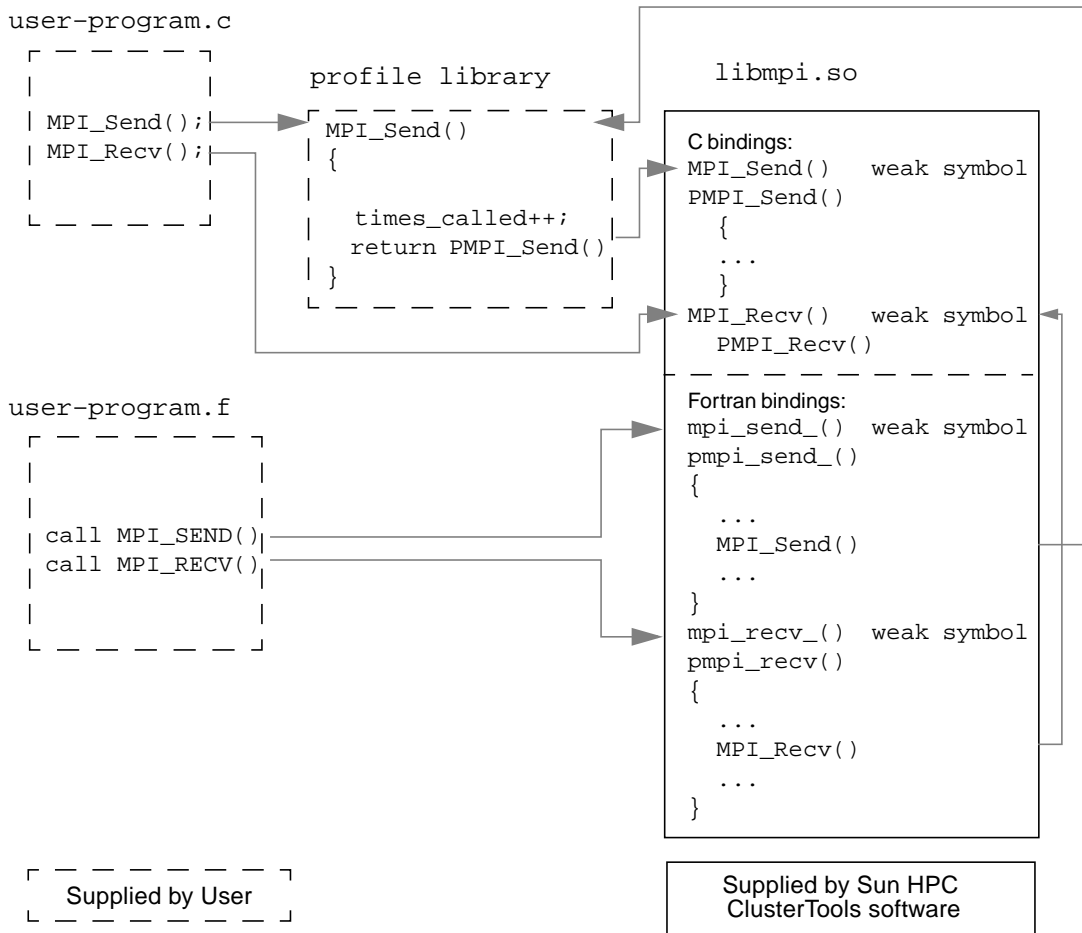


FIGURE 2-1 Sun MPI Profiling Interface

To clarify the layering of PMPI profiling, users need to understand the role of *weak symbols*. A weak symbol is such that, if a user defines the symbol, the user's definition is used. Otherwise, the associated function is used. The relation of weak symbols to associated functions is illustrated in FIGURE 2-2.



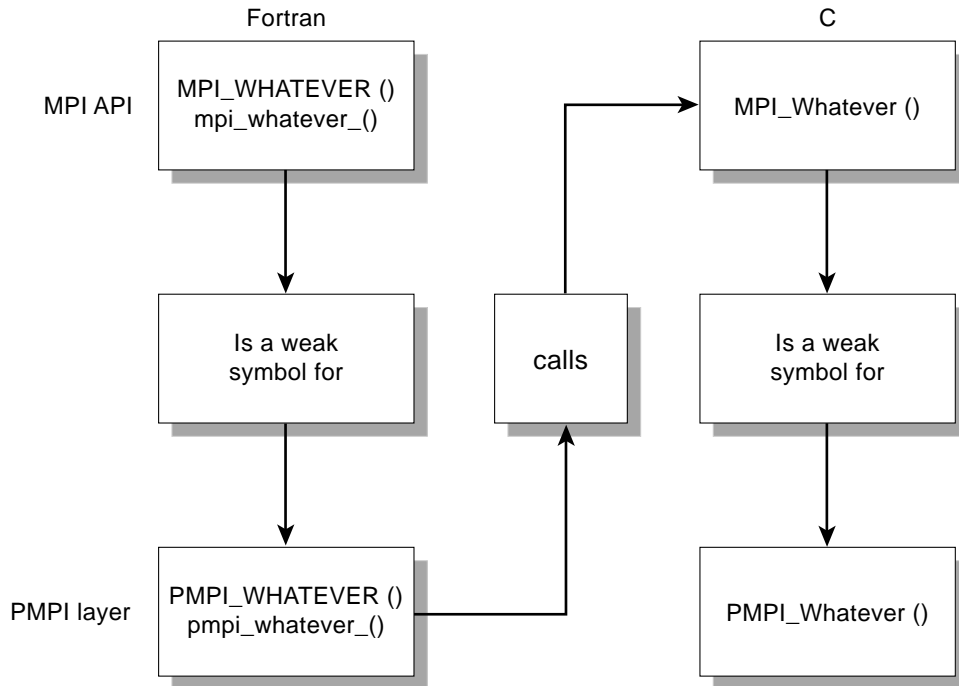


FIGURE 2-2 Layering in PMPI Profiling

## MPE: Extensions to the Library

Although the Sun MPI library does not include or support the multiprocessing environment (MPE) available from Argonne National Laboratory (ANL), it is compatible with MPE. If you would like to use these extensions to the MPI library, see the following instructions for downloading them from ANL and building MPE yourself. Note that this procedure may change if ANL makes changes to MPE.

### ▼ To Obtain and Build the MPE

The MPE software is available from Argonne National Laboratory. The `mpe.tar.gz` file is about 240 Kbytes.

**1. Use ftp to obtain the file.**

```
ftp://ftp.mcs.anl.gov/pub/mpi/misc/mpe.tar.gz
```

**2. Use gunzip and tar to decompress the software.**

```
# gunzip mpe.tar.gz  
# tar xvf mpe.tar
```

**3. Change your current working directory to the mpe directory, and execute configure with the arguments shown.**

```
# cd mpe  
# configure -cc=cc -fc=f77 -opt=-I/opt/SUNWhpc/include
```

**4. Execute a make.**

```
# make
```

---

**Note** – Sun MPI does not include the MPE error handlers. You must call the debug routines `MPE_Errors_call_dbx_in_xterm()` and `MPE_Signals_call_debugger()` yourself.

---

Refer to the *User's Guide for mpich, a Portable Implementation of MPI* for information on how to use MPE. It is available at the Argonne National Laboratory web site:

<http://www.mcs.anl.gov/mpi/mpich/>

## Getting Started

---

This chapter explains how to develop, compile and link, execute, and debug a Sun MPI program. The chapter focuses on what is specific to the Sun MPI implementation and, for the most part, does not repeat information that can be found in related documents. Information about programming with the Sun MPI I/O routines is in Chapter 4.

---

### Header Files

Include syntax must be placed at the top of any program that calls Sun MPI routines.

- For C and C++, use

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

- For Fortran, use

```
INCLUDE 'mpif.h'
```

These lines enable the program to access the Sun MPI version of the `mpi` header file, which contains the definitions, macros, and function prototypes required when compiling the program. Ensure that you are referencing the *Sun MPI* include file.

The include files are usually found in `/opt/SUNWhpc/include/` or `/opt/SUNWhpc/include/v9/`. If the compiler cannot find them, verify that they exist and are accessible from the machine on which you are compiling your code. The location of the include file is specified by a compiler option (see “Compiling and Linking” on page 34).

# Sample Code

Two simple Sun MPI programs are available in `/opt/SUNWhpc/examples/mpi` and are included here in their entirety. In the same directory you will find the `Readme` file, which provides instructions for using the examples, and the `make file` `Makefile`.

## CODE EXAMPLE 3-1 Simple Sun MPI Program in C: `connectivity.c`

```
/*
 * Test the connectivity between all processes.
 */

#pragma ident "@(#)connectivity.c 1.1 99/02/02"

#include <errno.h>
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <netdb.h>
#include <unistd.h>

#include <mpi.h>

int
main(int argc, char **argv)
{
    MPI_Status  status;
    int         verbose = 0;
    int         rank;
    int         np;           /* number of processes in job */
    int         peer;
    int         i;
    int         j;

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

    if (argc>1 && strcmp(argv[1], "-v")==0)
        verbose = 1;

    for (i=0; i<np; i++) {
        if (rank==i) {
            /* rank i sends to and receives from each higher rank */
            for(j=i+1; j<np; j++) {
                if (verbose)
                    printf("checking connection %4d <-> %-4d\n", i, j);
                MPI_Send(&rank, 1, MPI_INT, j, rank, MPI_COMM_WORLD);
            }
        }
    }
}
```

**CODE EXAMPLE 3-1** Simple Sun MPI Program in C: connectivity.c (Continued)

```

        MPI_Recv(&peer, 1, MPI_INT, j, j, MPI_COMM_WORLD, &status);
    }
} else if (rank>i) {
    /* receive from and reply to rank i */
    MPI_Recv(&peer, 1, MPI_INT, i, i, MPI_COMM_WORLD, &status);
    MPI_Send(&rank, 1, MPI_INT, i, rank, MPI_COMM_WORLD);
}
}

MPI_Barrier(MPI_COMM_WORLD);
if (rank==0)
    printf("Connectivity test on %d processes PASSED.\n", np);

MPI_Finalize();
return 0;
}

```

**CODE EXAMPLE 3-2** Simple Sun MPI Program in Fortran: monte.f

```

!
! Estimate pi via Monte-Carlo method.
!
! Each process sums how many of samplesize random points generated
! in the square (-1,-1),(-1,1),(1,1),(1,-1) fall in the circle of
! radius 1 and center (0,0), and then estimates pi from the formula
! pi = (4 * sum) / samplesize.
! The final estimate of pi is calculated at rank 0 as the average of
! all the estimates.
!
    program monte

    include 'mpif.h'

    double precision drand
    external drand

    double precision x, y, pi, pisum
    integer*4 ierr, rank, np
    integer*4 incircle, samplesize

    parameter(samplesize=2000000)

    call MPI_INIT(ierr)
    call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
    call MPI_COMM_SIZE(MPI_COMM_WORLD, np, ierr)

!    seed random number generator

```

**CODE EXAMPLE 3-2** Simple Sun MPI Program in Fortran: `monte.f` (Continued)

```
x = drand(2 + 11*rank)

incircle = 0
do i = 1, samplesize
  x = drand(0)*2.0d0 - 1.0d0      ! generate a random point
  y = drand(0)*2.0d0 - 1.0d0

  if ((x*x + y*y) .lt. 1.0d0) then
    incircle = incircle+1      ! point is in the circle
  endif
end do

pi = 4.0d0 * DBLE(incircle) / DBLE(samplesize)

! sum estimates at rank 0
  call MPI_REDUCE(pi, pisum, 1, MPI_DOUBLE_PRECISION, MPI_SUM,
&    0, MPI_COMM_WORLD, ierr)

if (rank .eq. 0) then
!   final estimate is the average
  pi = pisum / DBLE(np)
  print '(A,I4,A,F8.6,A)', 'Monte-Carlo estimate of pi by ', np,
&    ' processes is ', pi, '.'
endif

call MPI_FINALIZE(ierr)
end
```

---

## Compiling and Linking

Sun MPI programs are compiled with ordinary C, C++, or Fortran compilers, just like any other C, C++, or Fortran program, and linked with the Sun MPI library.

The `mpf77`, `mpf90`, `mpcc`, and `mpCC` utilities can be used to compile Fortran 77, Fortran 90, C, and C++ programs, respectively. For example, you might use the following entry to compile a Fortran 77 program that uses Sun MPI:

```
% mpf77 -fast -xarch=v8plusa -o a.out a.f -lmpi
```

See the man pages for more information on these utilities.

For performance, the single most important compilation switch is `-fast`. This is a macro that expands to settings appropriate for high performance for a general set of circumstances. Because its expansion varies from one compiler release to another, you might prefer to specify the underlying switches. To see what `-fast` expands to, use `-v` for “verbose” compilation output in Fortran, and `-#` for C. Also, `-fast` assumes native compilation, so you should compile on UltraSPARC™ processors.

The next important compilation switch is `-xarch`. Sun ONE Studio 7 Compiler Collection compilers set `-xarch` by default when you select `-fast` for native compilations. If you plan to compile on one type of processor and run the program on another type (nonnative compilation), be sure to use the `-xarch` flag. Also use it to compile in 64-bit mode. For UltraSPARC II processors, specify:

```
-xarch=v8plusa
```

or

```
-xarch=v9a
```

after `-fast` for 32-bit or 64-bit binaries, respectively. This version is only supported on Solaris 8 software. For UltraSPARC III processors, specify:

```
-xarch=v8plusb
```

or

```
-xarch=v9b
```

The `v8plusb` and `v9b` flags apply only to running on UltraSPARC III processors, and do not work when running on UltraSparc II processors. For more information, see the *Sun HPC ClusterTools Software Performance Guide* and the documents that came with your compiler.

Sun MPI programs compiled using the Sun ONE Studio 7 Compiler Collection Fortran compiler should be compiled with `-xalias=actual`. The `-xalias=actual` workaround requires patch 111718-01 (which requires 111714-01).

This recommendation arises because the MPI Fortran binding is inconsistent with the Fortran 90 standard in several respects. Specifically, this is documented in the MPI 2 standard, which you can find on the World Wide Web:

```
http://www-unix.mcs.anl.gov/mpi/mpi-standard/  
mpi-report-2.0/node19.htm#Node19
```

This recommendation applies to the use of high levels of compiler optimization. A highly optimizing Fortran compiler could break MPI codes that use nonblocking operations.

The failure modes are varied and insidious and include the following:

- Silently incorrect answers
- Intermittent and mysterious floating-point exceptions
- Intermittent and mysterious hangs

If you will be using the Prism debugger, you must compile your program with compilers from the Forte™ 6 update 2, or Sun ONE Studio 7 (formerly Forte Development 7 software) Compiler Collections (see “Debugging” on page 39).

**TABLE 3-1** Compile and Link Line Options for Sun MPI and Sun MPI I/O

Program	Options
C (nonthreaded example)	Use <code>mpcc</code> (below), or if you prefer: <pre>% cc filename.c -o filename \ -I/opt/SUNWhpc/include -L/opt/SUNWhpc/lib \ -R/opt/SUNWhpc/lib -lmpi</pre>
C++ Note that <code>x.0</code> represents the version of your C++ compiler (6.0 for Forte 6 update 2, and 7.0 for Sun ONE Studio 7).	Use <code>mpCC</code> (below), or if you prefer: <pre>% CC filename.cc -o filename \ -I/opt/SUNWhpc/include -L/opt/SUNWhpc/lib \ -R/opt/SUNWhpc/lib -L/opt/SUNWhpc/lib/SCx.0 \ -R/opt/SUNWhpc/lib/SCx.y -mt -lmpi++ -lmpi</pre>
<code>mpcc</code> , <code>mpCC</code>	<pre>% mpcc -o filename filename.c -lmpi % mpCC -o filename filename.cc -mt -lmpi</pre>
Fortran 77 (nonthreaded example)	Use <code>mpf77</code> (below), or if you prefer: <pre>% f77 -dalign filename.f -o filename \ -I/opt/SUNWhpc/include -L/opt/SUNWhpc/lib \ -R/opt/SUNWhpc/lib -lmpi</pre>
Fortran on a 64-bit system	<pre>% f77 -dalign filename.f -o filename \ -I/opt/SUNWhpc/include/v9 \ -L/opt/SUNWhpc/lib/sparcv9 \ -R/opt/SUNWhpc/lib/sparcv9 -lmpi</pre>
Fortran 90 <code>mpf90</code> , <code>mpf95</code>	Replace <code>mpf77</code> with <code>mpf90</code> , or <code>f77</code> with <code>f90</code> : <pre>% mpf90 -o filename -dalign filename.f -lmpi % mpf95 -o filename -dalign filename.f -lmpi</pre>
Multithreaded programs and programs containing nonblocking MPI I/O routines	To support multithreaded code, replace <code>-lmpi</code> with <code>-lmpi_mt</code> . This change also supports programs with nonblocking MPI I/O routines.  Note that <code>-lmpi</code> can be used for programs containing nonblocking MPI I/O routines, but <code>-lmpi_mt</code> <i>must</i> be used for multithreaded programs.

**Note** – For the Fortran interface, the `-dalign` option is necessary to avoid the possibility of bus errors. (The underlying C or C++ routines in Sun MPI internals assume that parameters and buffer types passed as `REALs` are double-aligned.)



---

**Note** – If your program has previously been linked to any static libraries, you must relink it to `libmpi.so` before executing it.

---

## Choosing a Library Path

The paths for the MPI libraries, which you must specify when you are compiling and linking your program, are listed in TABLE 3-2.

**TABLE 3-2** Sun MPI Libraries

Category	Description	Path: /opt/SUNWhpc/lib/...
32-bit libraries	Default, not thread-safe	<code>libmpi.so</code>
	C++ (in addition to <code>libmpi.so</code> )	<code>SC6.0/libmpi++.so</code>
	Thread-safe	<code>libmpi_mt.so</code>
64-bit libraries	Default, not thread-safe	<code>sparcv9/libmpi.so</code>
	C++ (in addition to <code>sparcv9/libmpi.so</code> )	<code>sparcv9/SC6.0/libmpi++.so</code>
	Thread-safe	<code>sparcv9/libmpi_mt.so</code>

## Stubbing Thread Calls

The `libthread.so` libraries are automatically linked into the respective `libmpi.so` libraries. This means that any thread-function calls in your program can be resolved by the `libthread.so` library. Simply omitting `libthread.so` from the link line does not cause thread calls to be stubbed out; you must remove the thread calls yourself. For more information about the `libthread.so` library, see its man page. (For the location of Solaris man pages at your site, see your system administrator.)

---

# Profiling With `mpprof`

If you plan to extract MPI profiling information from the execution of a job, you need to set the `MPI_PROFILE` environment variable to 1 before you start the job execution.

```
% setenv MPI_PROFILE 1
```

If you want to set any other `mpprof` environment variables, you must set them also before starting the job. See Appendix B for detailed descriptions of the `mpprof` environment variables.

---

# Basic Job Execution

The CRE environment provides close integration with batch-processing systems, also known as resource managers. You can launch parallel jobs from a batch system to control resource allocation, and continue to use the CRE environment to monitor job status. For a list of currently supported resource managers, see TABLE 3-3.

**TABLE 3-3** Currently Supported Resource Managers

Resource manager	Name used with -x option to <code>mprun</code>	Version	Man page
Sun Grid Engine	<code>sge</code>	SGE 5.3	<code>sge_cre.1</code>
PBS	<code>pbs</code>	PBS 2.3.15	<code>pbs_cre.1</code>
	<code>pbs</code>	PBS Pro 5.x.x	<code>pbs_cre.1</code>
LSF	<code>lsf</code>	LSF 4.x	<code>lsf_cre.1</code>

To enable the integration between the CRE environment and the supported resource managers, you must call `mprun` from a script in the resource manager. Use the `-x` flag to specify the resource manager, and the `-np` and `-nr` flags to specify the resources you need. Instructions and examples for each resource manager are provided in the *Sun HPC ClusterTools Software User's Guide*.

Before starting your job, you might want to set one or more environment variables, which are also described in Appendix B and in the *Sun HPC ClusterTools Software Performance Guide*.

## Executing With CRE

When using CRE software, parallel jobs are launched using the `mprun` command. For example, to start a job with six processes named `mpijob`, use this command:

```
% mprun -np 6 mpijob
```

## Executing With LSF Suite

Parallel jobs can be either launched by the LSF Parallel Application Manager (PAM) or submitted in queues configured to run PAM as the parallel job starter. LSF's `bsub` command launches both parallel interactive and batch jobs. For example, to start a batch job named `mpijob` on four CPUs, use this command:

```
% bsub -n 4 pam mpijob
```

To launch an interactive job, add the `-I` argument to the command line. For example, to launch an interactive job named `earth` on a single CPU in the queue named `sun`, which is configured to launch jobs with PAM, use this command:

```
% bsub -q sun -Ip -n 1 earth
```

---

## Debugging

Debugging parallel programs is notoriously difficult, because you are in effect debugging a program potentially made up of many distinct programs executing simultaneously. Even if the application is an SPMD (single-program, multiple-data) application, each instance can be executing a different line of code at any instant. The Prism development environment eases the debugging process considerably and is recommended for debugging with Sun HPC ClusterTools software.

# Debugging With the Prism Environment

---

**Note** – To run the graphical version of the Prism environment, you must be running the Solaris 8 operating environment with either OpenWindows™ software or the Common Desktop Environment (CDE), and with your `DISPLAY` environment variable set correctly. See the *Prism Software User's Guide* for information.

---

This section provides a brief introduction to the Prism development environment.

You can use a Prism session to debug more than one Sun MPI job at a time. To debug a child or client program it is necessary to launch an additional Prism session. If the child program is spawned using calls to `MPI_Comm_spawn()` or `MPI_Comm_spawn_multiple()`, Prism can (if enabled) debug the child program as well.

However, if an MPI job connects to another job, the current Prism session has control only of the parent or server MPI job. It cannot debug the children or clients of that job. This might occur, for example, when an MPI job sets up a client/server connection to another MPI job with `MPI_Comm_accept` or `MPI_Comm_connect`.

With the exception of programs using calls to `MPI_Comm_spawn()` or `MPI_Comm_spawn_multiple()`, to use the Prism environment to debug a Sun MPI program the program must be written in the SPMD style. In other words, all processes that make up a Sun MPI program must be running the same executable.

`MPI_Comm_spawn_multiple` can create multiple executables with only one job ID. Therefore, you can use the Prism environment to debug jobs with various executables that have been spawned with this command.

## Starting Up Prism

To start Prism with a Sun MPI program, launch it from within the `mprun` command.

For example,

```
% mprun -np 4 -x lsf prism -np 4 foo
```

launches Prism on executable `foo` with four processes.

This starts up a graphical version of Prism with your program loaded. You can then debug and visualize data in your Sun MPI program.

You can also attach Prism to running processes. First determine the job ID (not the individual process ID), jobname (or *jid*), using `mpps`. (See the *Sun HPC ClusterTools Software User's Guide* for further information about `mpps`.) Then specify the `jid` at the command line:

```
% prism foo 12345
```

This launches Prism and attaches it to the processes running in job 12345.

One important feature of the Prism environment is that it enables you to debug the Sun MPI program at any level of detail. You can look at the program as a whole, or you can look at subsets of processes within the program (for example, those that have an error condition) or at individual processes, all within the same debugging session. For complete information, see the *Prism Software User's Guide*.

## Debugging With TotalView

TotalView is a third-party multiprocess debugger from Etnus that runs on many platforms. Support for using the TotalView debugger on Sun MPI applications includes:

- Making Sun HPC ClusterTools software compatible with the Etnus debugger TotalView
- Allowing Sun MPI jobs to be debugged by TotalView using the Sun Grid Engine (SGE), the Portable Batch System (PBS), and Platform Computing's Load Sharing Facility (LSF)
- Displaying Sun MPI message queues
- Allowing multiple instantiations of TotalView on a single cluster
- Supporting TotalView in Sun HPC ClusterTools software

The following sections provide a brief description of how to use the TotalView debugger with Sun MPI applications, including:

- "Limitations" on page 41
- "Related Documentation" on page 42
- "Starting a New Job Using TotalView" on page 42
- "Attaching to an `mprun` Job" on page 44
- "Launching Sun MPI Batch Jobs Using TotalView" on page 45

Refer to your TotalView documentation for more information about using TotalView.

### Limitations

- Debuggable job restricted according to the license with Etnus. Contact the system administrator who installed TotalView for more details.

- Supports the SPARC platform *only*.
- TotalView 5 supports 32-bit application debugging only. However, TotalView 6 supports both 32-bit and 64-bit application debugging.
- Does *not* support `MPI_Comm_spawn` and `MPI_Comm_spawn_multiple` function calls. Use Prism to debug these function calls.
- Displays `MPI_COMM_WORLD` in the message queue graph only after `MPI_Init()` has occurred. Displays *neither* collective communicators *nor* `MPI_COMM_SELF`. Refer to the Etnus TotalView user's guide for more information.
- Does not display any buffer contents for unexpected messages in the message queue window.

## Related Documentation

For more information, refer to the following related documentation:

- *Sun HPC ClusterTools Software User's Guide*
- *Sun HPC ClusterTools Software Administrator's Guide*
- Sun HPC ClusterTools software man pages
  - `lsf_cre(1)`
  - `lsf_cre_admin(1M)`
  - `mpps(1M)`
  - `mprun(1M)`
  - `pbs_cre(1)`
  - `pbs_cre_admin(1M)`
  - `sge_cre(1)`
  - `sge_cre_admin(1M)`
  - `totalview_mprun(1M)`
- Etnus TotalView documentation

## Starting a New Job Using TotalView

You can start a new job from the Total View Graphical User Interface (GUI) using:

- GUI method 1
- GUI method 2
- Command-line interface (CLI)

## ▼ To Start a New Job Using GUI Method 1

### 1. Type:

```
% totalview mprun [totalview args] -a [mprun args]
```

For example:

```
% totalview mprun -bg blue -a -np 4 /opt/SUNWhpc/mpi/conn.x
```

### 2. When the GUI appears, type `g` for go, or click **Go** in the TotalView window.

TotalView may display a dialog box:

```
Process mprun is a parallel job. Do you want to stop the job now?
```

### 3. Click **Yes** to open the TotalView debugger window with the Sun MPI source window, if compiled with option `-g`, and to leave all processes in a traced state.

## ▼ To Start a New Job Using GUI Method 2

### 1. Type:

```
% totalview
```

### 2. Select the menu option **File** and then **New Program**.

### 3. Type `mprun` as the executable name in the dialog box.

### 4. Click **OK**.

TotalView displays the main debug window.

### 5. Select the menu option **Process** and then **Startup Parameters**, which are the *`mprun args`*.

## ▼ To Start a New Job Using the CLI

### 1. Type:

```
% totalviewcli mprun [totalview args] -a [mprun args]
```

For example:

```
% totalviewcli mprun -a -np 4 /opt/SUNWhpc/mpi/conn.x
```

### 2. When the job starts, type `dgo`.

TotalView displays this message:

```
Process mprun is a parallel job. Do you want to stop the job now?
```

### 3. Type `y` to start the MPI job, attach TotalView, and leave all processes in a traced state.

## Attaching to an `mprun` Job

This section describes how to attach to an already running `mprun` job from both the TotalView GUI and CLI.

## ▼ To Attach to a Running Job from the GUI

### 1. Find the host name and process identifier (PID) of the `mprun` job by typing:

```
% mpps -b
```

`mprun` displays the PID and host name in a similar manner to this example:

JOBNAME	MPRUN_PID	MPRUN_HOST
cre.99	12345	hpc-u2-9
cre.100	12601	hpc-u2-8

For more information, refer to the `mpps(1M)` man page, option `-b`.

### 2. In the TotalView GUI, select **File** and then **New Program**.



3. **Type the PID in Process ID.**
4. **Type `mprun` in the field Executable Name.**
5. **Do one of the following:**
  - Leave Remote Host blank if TotalView is running on the same node as the `mprun` job.
  - Enter the host name in Remote Host.
6. **Click OK.**

## ▼ To Attach to a Running Job From the CLI

1. **Find the process identifier (PID) of the launched job.**

See the example under the preceding GUI procedure. For more information, refer to the `mpps(1M)` man page, option `-b`.
2. **Start `totalviewcli` by typing:**

```
% totalviewcli
```

3. **Attach the executable program to the `mprun` PID:**

```
% dattach mprun mprun_pid
```

For example:

```
% dattach mprun 12601
```

## Launching Sun MPI Batch Jobs Using TotalView

This section describes how to launch Sun MPI batch jobs, including:

- TotalView's control using the GUI only
- Interactive sessions using the GUI
- Interactive sessions using the CLI

This section provides examples of launching batch jobs in Sun Grid Engine (SGE). Refer to Chapter 5 of the *Sun HPC ClusterTools Software User's Guide* for descriptions of launching batch jobs in the Load Sharing Facility (LSF) and the Portable Batch System (PBS).

## ▼ To Execute Startup in Batch Mode for the TotalView GUI

Executing startup in batch mode for the TotalView CLI is not practical, because there is no controlling terminal for input and output. This procedure describes executing startup in batch mode for the TotalView GUI:

1. Write a batch script, which contains a line similar to the following:

```
% totalview mprun -a -x SGE /opt/SUNWhpc/mpi/conn.x
```

2. Then submit the script to SGE for execution with a command similar to the following:

```
% qsub -l crfe 4 batch_script
```

The TotalView GUI appears upon successful allocation of resources and execution of the batch script in SGE.

## ▼ To Use the Interactive Mode

The interactive mode creates an `xterm` window for your use, so you can use either the TotalView GUI or the CLI.

1. Submit an interactive mode job to SGE with a command similar to the following:

```
% qsh -l cre 4
```

The system displays an `xterm` window.

2. Run the following, or an equivalent path, to source the SGE environment:

```
% source /opt/sge/default/common/settings.csh
```

### 3. Execute a typical `totalview` or `totalviewcli` command.

TotalView GUI example:

```
% totalview mprun -a -x SGE /opt/SUNWhpc/mpi/conn.x
```

TotalView CLI example:

```
% totalviewcli mprun -a -x SGE /opt/SUNWhpc/mpi/conn.x
```

## Debugging With MPE

The multiprocessing environment (MPE) available from Argonne National Laboratory includes a debugger that can also be used for debugging at the thread level. For information about obtaining and building MPE, see “MPE: Extensions to the Library” on page 29.



## Programming With Sun MPI I/O

---

File I/O in Sun MPI is fully MPI-2 compliant. MPI I/O is specified as part of that standard, which was published in 1997. Its goal is to provide a library of routines featuring a portable parallel file system interface that is an extension of the MPI framework. See “Related Documentation” on page x for more information about the *MPI-2 Standard*.

MPI I/O models file I/O on message passing; that is, writing to a file is analogous to sending a message, and reading from a file is analogous to receiving a message. The MPI library provides a high-level way of partitioning data among processes, which saves you from having to specify the details involved in making sure that the right pieces of data go to the right processes. This section describes basic MPI I/O concepts and the Sun MPI I/O routines.

---

## Data Partitioning and Data Types

MPI I/O uses the MPI model of communicators and derived data types to describe communication between processes and I/O devices. MPI I/O determines which processes are communicating with a particular I/O device. Derived data types can be used to define the layout of data in memory and of data in a file on the I/O device. (For more information about derived data types, see “Data Types” on page 17.) Because MPI I/O builds on MPI concepts, it’s easy for a knowledgeable MPI programmer to add MPI I/O code to a program.

Data is stored in memory and in the file according to MPI data types. Herein lies one of MPI and MPI I/O’s advantages: Because they provide a mechanism whereby you can create your own data types, you have more freedom and flexibility in specifying data layout in memory and in the file.

The library also simplifies the task of describing how your data moves from processor memory to the file and back again. You create derived data types that describe how the data is arranged in the memory of each process and how it should be arranged in that part of the disk file associated with the process.

Three functions are provided to handle the external32 format. This format, defined by the MPI Forum, represents data in a universal format that is useful for exchanging data between implementations or for writing it to a file. The functions are:

- `MPI_Pack_external()`
- `MPI_Unpack_external()`
- `MPI_Pack_external_size()`

The Sun MPI I/O routines are described in “Routines” on page 51. But first, to be able to define a data layout, you will need to understand some basic MPI I/O data-layout concepts. The next section explains some of the fundamental terms and concepts.

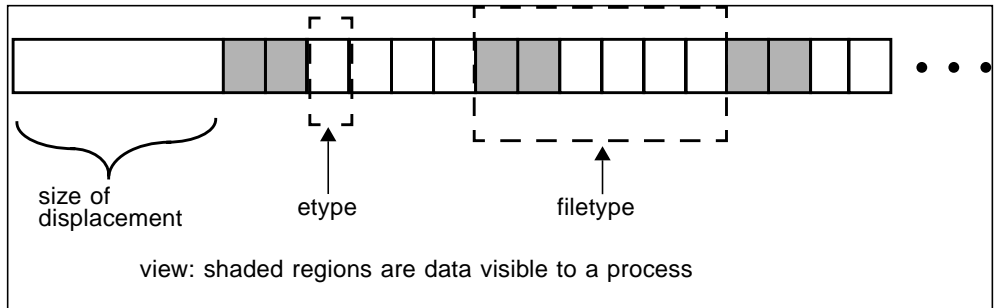
---

## Definitions

The following terms are used to describe partitioning data among processes.

FIGURE 4-1 illustrates some of these concepts.

- An *elementary data type* (or `etype`) is the unit of data access and positioning. It can be any MPI basic or derived data type. Data access is performed in elementary-data-type units, and offsets (defined later in this list) are expressed as a count of elementary data types.
- The *file type* (or `filetype`) is used to partition a file among processes; that is, a file type defines a template for accessing the file. It is either a single elementary data type or a derived MPI data type constructed from elementary data types. A file type can contain “holes,” or extents of bytes that are not accessed by this process.
- A file *displacement* (or `disp`) is an absolute byte position counted from the beginning of a file. The displacement defines the location where a view begins (see FIGURE 4-1).
- A *view* defines the current set of data visible and accessible by a process from an open file in terms of a displacement, an elementary data type, and a file type. The pattern described by a file type is repeated, beginning at the displacement, to define the view.
- An *offset* is a position relative to the current view, expressed as a count of elementary data types. Holes in the view’s file type are ignored when calculating this position.



**FIGURE 4-1** Displacement, the Elementary Data Type, the File Type, and the View

For a more detailed description of MPI I/O, see Chapter 9, “I/O,” of the *MPI-2 Standard*.

## Note for Fortran Users

When writing a Fortran program, you must declare the variable `ADDRESS` as follows:

```
INTEGER*MPI_ADDRESS_KIND ADDRESS
```

`MPI_ADDRESS_KIND` is a constant defined in `mpi.h`. This constant defines the length of the declared integer.

---

## Routines

This release of Sun MPI includes all the MPI I/O routines, which are defined in Chapter 9, “I/O,” of the *MPI-2 Standard*.

Code samples that use many of these routines are provided in “Sample Code” on page 59.

# File Manipulation

Collective coordination	Noncollective coordination
<code>MPI_File_open()</code>	<code>MPI_File_delete()</code>
<code>MPI_File_close()</code>	<code>MPI_File_get_size()</code>
<code>MPI_File_set_size()</code>	<code>MPI_File_get_group()</code>
<code>MPI_File_preallocate()</code>	<code>MPI_File_get_amode()</code>

`MPI_File_open()` and `MPI_File_close()` are collective operations that open and close a file, respectively; that is, all processes in a communicator group must together open or close a file. To achieve a single-user, UNIX-like open, set the communicator to `MPI_COMM_SELF`.

`MPI_File_delete()` deletes a specified file.

The routines `MPI_File_set_size()`, `MPI_File_get_size()`, `MPI_File_get_group()`, and `MPI_File_get_amode()` get and set information about a file. When using the collective routine `MPI_File_set_size()` on a UNIX file, if the size that is set is smaller than the current file size, the file is truncated at the position defined by *size*. If *size* is set to be larger than the current file size, the file size becomes *size*.

When the file size is increased this way with `MPI_File_set_size()`, new regions are created in the file with displacements between the old file size and the larger, newly set file size. Sun MPI I/O does not necessarily allocate file space for such new regions. You can reserve file space either by using `MPI_File_preallocate()` or by performing a read or write to unallocated bytes. `MPI_File_preallocate()` ensures that storage space is allocated for a set quantity of bytes for the specified file; however, its use is very “expensive” in terms of performance and disk space.

The routine `MPI_File_get_group()` returns a communicator group, but it does not free the group.

## File Hints

The opaque `info` object enables you to provide hints for optimization of your code, making it run faster or more efficiently, for example. These hints are set for each file, using the `MPI_File_open()`, `MPI_File_set_view()`, `MPI_File_set_info()`, and `MPI_File_delete()` routines. `MPI_File_set_info()` sets new values for the specified file’s hints. `MPI_File_get_info()` returns all the hints that the system currently associates with the specified file.



When using UNIX files, Sun MPI I/O provides four hints for controlling how much buffer space it uses to satisfy I/O requests: `noncoll_read_bufsize`, `noncoll_write_bufsize`, `coll_read_bufsize`, and `coll_write_bufsize`. These hints can be tuned for your particular hardware configuration and application to improve performance for both noncollective and collective data accesses. For example, if your application uses a single MPI I/O call to request multiple noncontiguous chunks that form a regular strided pattern in the file, you can adjust the `noncoll_write_bufsize` to match the size of the stride. Note that these hints limit the size of MPI I/O's underlying buffers but do not limit the amount of data a user can read or write in a single request.

## File Views

The `MPI_File_set_view()` routine changes the view the process has of the data in the file, specifying its displacement, elementary data type, and file type, as well as setting the individual file pointers and shared file pointer to 0.

`MPI_File_set_view()` is a collective routine; all processes in the group must pass identical values for the file handle and the elementary data type, although the values for the displacement, the file type, and the info object can vary. However, if you use the data-access routines that use file positioning with a shared file pointer, you must also give the displacement and the file type identical values. The data types passed in as the elementary data type and the file type must be committed.

You can also specify the type of data representation for the file. See “File Interoperability” on page 57 for information about registering data representation identifiers.

---

**Note** – Displacements within the file type and the elementary data type must be monotonically nondecreasing.

---

## Data Access

The 35 data-access routines are categorized according to file positioning. Data access can be achieved by any of these methods of file positioning:

- Explicit offset
- Individual file pointer
- Shared file pointer

This section discusses each of these methods in more detail.

While *blocking* I/O calls do not return until the request is completed, *nonblocking* calls do not wait for the I/O request to complete. A separate “request complete” call, such as `MPI_Test()` or `MPI_Wait()`, is necessary to confirm that the buffer is ready to be used again. Nonblocking routines have the prefix `MPI_File_i`, where the *i* stands for *immediate*.

All the nonblocking collective routines for data access are “split” into two routines, each with `_begin` or `_end` as a suffix. These *split collective* routines are subject to the semantic rules described in Section 9.4.5 of the *MPI-2 Standard*.

## Data Access With Explicit Offsets

Synchronism	Noncollective coordination	Collective coordination
Blocking	<code>MPI_File_read_at()</code> <code>MPI_File_write_at()</code>	<code>MPI_File_read_at_all()</code> <code>MPI_File_write_at_all()</code>
Nonblocking or split collective	<code>MPI_File_iread_at()</code> <code>MPI_File_iwrite_at()</code>	<code>MPI_File_read_at_all_begin()</code> <code>MPI_File_read_at_all_end()</code> <code>MPI_File_write_at_all_begin()</code> <code>MPI_File_write_at_all_end()</code>

To access data at an explicit offset, specify the position in the file where the next data access for each process should begin. For each call to a data-access routine, a process attempts to access a specified number of file types of a specified data type (starting at the specified offset) into a specified user buffer.

The offset is measured in elementary data type units relative to the current view; moreover, holes are not counted when locating an offset. The data is read from (in the case of a read) or written into (in the case of a write) those parts of the file specified by the current view. These routines store the number of buffer elements of a particular data type actually read (or written) in the status object, and all the other fields associated with the status object are undefined. The number of elements that are read or written can be accessed using `MPI_Get_count()`.

`MPI_File_read_at()` attempts to read from the file by the associated file handle returned from a successful `MPI_File_open()`. Similarly, `MPI_File_write_at()` attempts to write data from a user buffer to a file. `MPI_File_iread_at()` and `MPI_File_iwrite_at()` are the nonblocking versions of `MPI_File_read_at()` and `MPI_File_write_at()`, respectively.

`MPI_File_read_at_all()` and `MPI_File_write_at_all()` are collective versions of `MPI_File_read_at()` and `MPI_File_write_at()`, in which each process provides an explicit offset. The split collective versions of these nonblocking routines are listed in the preceding table.

## Data Access With Individual File Pointers

Synchronism	Noncollective coordination	Collective coordination
Blocking	<code>MPI_File_read()</code> <code>MPI_File_write()</code>	<code>MPI_File_read_all()</code> <code>MPI_File_write_all()</code>
Nonblocking or split collective	<code>MPI_File_iread()</code> <code>MPI_File_iwrite()</code>	<code>MPI_File_read_all_begin()</code> <code>MPI_File_read_all_end()</code> <code>MPI_File_write_all_begin()</code> <code>MPI_File_write_all_end()</code>

For each open file, Sun MPI I/O maintains one individual file pointer per process per collective `MPI_File_open()`. For these data-access routines, MPI I/O implicitly uses the value of the individual file pointer. These routines use and update only the individual file pointers maintained by MPI I/O by pointing to the next elementary data type after the one that most recently has been accessed. The individual file pointer is updated relative to the current view of the file. The shared file pointer is neither used nor updated. (For data access with shared file pointers, see the next section.)

These routines have similar semantics to the explicit-offset data-access routines, except that the offset is defined here to be the current value of the individual file pointer.

`MPI_File_read_all()` and `MPI_File_write_all()` are collective versions of `MPI_File_read()` and `MPI_File_write()`, with each process using its individual file pointer.

`MPI_File_iread()` and `MPI_File_iwrite()` are the nonblocking versions of `MPI_File_read()` and `MPI_File_write()`, respectively. The split collective versions of `MPI_File_read_all()` and `MPI_File_write_all()` are listed in the preceding table.

## Pointer Manipulation

`MPI_File_seek`

`MPI_File_get_position`

`MPI_File_get_byte_offset`

Each process can call the routine `MPI_File_seek()` to update its individual file pointer according to the update mode. The update mode has the following possible values:

- `MPI_SEEK_SET` – The pointer is set to the offset.
- `MPI_SEEK_CUR` – The pointer is set to the current pointer position plus the offset.
- `MPI_SEEK_END` – The pointer is set to the end of the file plus the offset.

The offset can be negative for seeking backwards, but you cannot seek to a negative position in the file. The current position is defined as the elementary data item immediately following the last-accessed data item.

`MPI_File_get_position()` returns the current position of the individual file pointer relative to the current displacement and file type.

`MPI_File_get_byte_offset()` converts the offset specified for the current view to the displacement value, or absolute byte position, for the file.

## Data Access With Shared File Pointers

Synchronism	Noncollective coordination	Collective coordination
Blocking	<code>MPI_File_read_shared()</code>	<code>MPI_File_read_ordered()</code>
	<code>MPI_File_write_shared()</code>	<code>MPI_File_write_ordered()</code>
		<code>MPI_File_seek_shared()</code>
		<code>MPI_File_get_position_shared()</code>
Nonblocking or split collective	<code>MPI_File_iread_shared()</code>	<code>MPI_File_read_ordered_begin()</code>
	<code>MPI_File_iwrite_shared()</code>	<code>MPI_File_read_ordered_end()</code>
		<code>MPI_File_write_ordered_begin()</code>
		<code>MPI_File_write_ordered_end()</code>

Sun MPI I/O maintains one shared file pointer per collective `MPI_File_open()` (shared among processes in the communicator group that opened the file). As with the routines for data access with individual file pointers, you can also use the current value of the shared file pointer to specify the offset of data accesses implicitly. These routines use and update only the shared file pointer; the individual file pointers are neither used nor updated by any of these routines.

These routines have similar semantics to the explicit-offset data-access routines, except:

- The offset is defined here to be the current value of the shared file pointer.
- Multiple calls (one for each process in the communicator group) affect the shared file pointer routines as if the calls were serialized.
- All processes must use the same file view.

After a shared file pointer operation is initiated, it is updated, relative to the current view of the file, to point to the elementary data item immediately following the last one requested, regardless of the number of items actually accessed.

`MPI_File_read_shared()` and `MPI_File_write_shared()` are blocking routines that use the shared file pointer to read and write files, respectively. The order of serialization is not deterministic for these noncollective routines, so you need to use other methods of synchronization if you want to impose a particular order.

`MPI_File_iread_shared()` and `MPI_File_iwrite_shared()` are the nonblocking versions of `MPI_File_read_shared()` and `MPI_File_write_shared()`, respectively.

`MPI_File_read_ordered()` and `MPI_File_write_ordered()` are the collective versions of `MPI_File_read_shared()` and `MPI_File_write_shared()`. They must be called by all processes in the communicator group associated with the file handle, and the accesses to the file occur in the order determined by the ranks of the processes within the group. After all the processes in the group have issued their respective calls, for each process in the group, these routines determine the position of the shared file pointer after all processes with ranks lower than this process's rank had accessed their data. Then data is accessed (read or written) at that position. The shared file pointer is then updated by the amount of data requested by all processes of the group.

The split collective versions of `MPI_File_read_ordered()` and `MPI_File_write_ordered()` are listed in the preceding table.

`MPI_File_seek_shared()` is a collective routine, and all processes in the communicator group associated with the particular file handler must call `MPI_File_seek_shared()` with the same file offset and the same update mode. All the processes are synchronized with a barrier before the shared file pointer is updated.

The offset can be negative for seeking backwards, but you cannot seek to a negative position in the file. The current position is defined as the elementary data item immediately following the last-accessed data item, even if that location is a hole.

`MPI_File_get_position_shared()` returns the current position of the shared file pointer relative to the current displacement and file type.

## File Interoperability

`MPI_Register_datarep()`

`MPI_File_get_type_extent()`

Sun MPI I/O supports the basic data representations described in Section 9.5 of the *MPI-2 Standard*:

- *native* – With native representation, data is stored exactly as in memory, in other words, in Solaris/UltraSPARC data representation. This format offers the highest performance and no loss of arithmetic precision. It should be used only in a

homogeneous environment, that is, on Solaris/UltraSPARC nodes running Sun HPC ClusterTools software. It also can be used when the MPI application will perform the data type conversions itself.

- *internal* – With internal representation, data is stored in an implementation-dependent format, such as for Sun MPI.
- *external32* – With external32 representation, data is stored in a portable format, prescribed by the MPI-2 and IEEE standards.

These data representations, as well as any user-defined representations, are specified as an argument to `MPI_File_set_view()`.

You can create user-defined data representations with `MPI_Register_datarep()`. Once a data representation has been defined with this routine, you can specify it as an argument to `MPI_File_set_view()`, so that subsequent data-access operations can call the conversion functions specified with `MPI_Register_datarep()`.

If the file data representation is anything but native, you must be careful when constructing elementary data types and file types. For those functions that accept displacements in bytes, the displacements must be specified in terms of their values in the file for the file data representation being used.

You can use `MPI_File_get_type_extent()` to calculate the extents of data types in the file. The extent is the same for all processes accessing the specified file. If the current view uses a user-defined data representation, `MPI_File_get_type_extent()` uses one of the functions specified in setting the data representation to calculate the extent.

## File Consistency and Semantics

Noncollective coordination	Collective coordination
<code>MPI_File_get_atomicsity()</code>	<code>MPI_File_set_atomicsity()</code>
	<code>MPI_File_sync()</code>

The routines ending in `_atomicity` enable you either to set a file's mode as atomic or nonatomic, or to query which mode it is in. In *atomic mode*, all operations within the communicator group that opens a file are completed as if sequentialized into a serial order. In *nonatomic mode*, no such guarantee is made. In nonatomic mode, `MPI_File_sync()` can be used to ensure weak consistency.

The default mode varies with the number of nodes you are using. If you are running a job on a single node, a file is in *nonatomic* mode by default when it is opened. If you are running a job on more than one node, a file is in *atomic* mode by default.

`MPI_File_set_atomicsity()` is a collective call that sets the consistency semantics for data-access operations. All the processes in the group must pass identical values for both the file handle and the Boolean flag that indicates whether atomic mode is set.

`MPI_File_get_atomicsity()` returns the current consistency semantics for data-access operations. Again, a Boolean flag indicates whether the atomic mode is set.

---

**Note** – In some cases, setting `atomicity` to `false` can provide better performance. The default `atomicity` value on a cluster is `true`. The lack of synchronization among the distributed caches on a cluster can prevent your data from completing in the desired state. In these circumstances, you might suffer performance disadvantages with `atomicity` set to `true`, especially when the data accesses overlap.

---

## Sample Code

This section provides sample code to get you started with programming your I/O using Sun MPI. The first example shows how a parallel job can partition file data among its processes. That example is then adapted to use a broad range of other I/O programming styles supported by Sun MPI I/O. The last code sample illustrates the use of the nonblocking MPI I/O routines.

Remember that MPI I/O is part of MPI, so be sure to call `MPI_Init()` before calling any MPI I/O routines, and call `MPI_Finalize()` at the end of your program, even if you use only MPI I/O routines.

## Partitioned Writing and Reading in a Parallel Job

MPI I/O was designed to enable processes in a parallel job to request multiple data items that are noncontiguous within a file. Typically, a parallel job partitions file data among the processes.

One method of partitioning a file is to derive the offset at which to access data from the rank of the process. The rich set of MPI derived types also makes it easy to partition file data. For example, you could create an MPI vector type as the filetype passed into `MPI_File_set_view()`. Because vector types do not end with a hole, you would make a call to either `MPI_Type_create_resized()` or `MPI_Type_ub()` to complete the partition. This call would lengthen the extent to include holes at the end of the type for processes with higher ranks. You can create a partitioned file by passing various displacements to `MPI_File_set_view()`. Each

of these displacements would be derived from the process's rank. Consequently, offsets would not need to be derived from the ranks, because only the data in the portion of the partition belonging to the process would be visible to the process.

The following example uses the first method that derives the file offsets directly from the rank of the process. Each process writes and reads `NUM_INTS` integers starting at the offset `rank * NUM_INTS`. It passes an explicit offset to the MPI I/O data-access routines `MPI_File_write_at()` and `MPI_File_read_at()`. It calls `MPI_Get_elements()` to find out how many elements were written or read. To verify that the write was successful, it compares the data written and read as well as set up an `MPI_Barrier()` before calling `MPI_File_get_size()` to verify that the file is the size expected upon completion of all the writes of the process.

Note that `MPI_File_set_view()` was called to set the view of the file as essentially an array of integers instead of the UNIX-like view of the file as an array of bytes. Thus, the offsets that are passed to `MPI_File_write_at()` and `MPI_File_read_at()` are indices into an array of integers and not a byte offset.

In CODE EXAMPLE 4-1, each process writes and reads `NUM_INTS` integers to a file using `MPI_File_write_at()` and `MPI_File_read_at()`, respectively.

#### CODE EXAMPLE 4-1 Writing and Reading Integers to a File

```
/* wr_at.c
 *
 * Example to demonstrate use of MPI_File_write_at and MPI_File_read_at
 *
 */

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

#define NUM_INTS 100

void sample_error(int error, char *string)
{
    fprintf(stderr, "Error %d in %s\n", error, string);
    MPI_Finalize();
    exit(-1);
}

void
main( int argc, char **argv )
{
    char filename[128];
    int i, rank, comm_size;
    int *buff1, *buff2;
```



**CODE EXAMPLE 4-1 Writing and Reading Integers to a File (Continued)**

```
MPI_File fh;
MPI_Offset disp, offset, file_size;
MPI_Datatype etype, ftype, buftype;
MPI_Info info;
MPI_Status status;
int result, count, differs;

if(argc < 2) {
    fprintf(stdout, "Missing argument: filename\n");
    exit(-1);
}
strcpy(filename, argv[1]);

MPI_Init(&argc, &argv);

/* get this processor's rank */
result = MPI_Comm_rank(MPI_COMM_WORLD, &rank);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_Comm_rank");

result = MPI_Comm_size(MPI_COMM_WORLD, &comm_size);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_Comm_size");

/* communicator group MPI_COMM_WORLD opens file "foo"
   for reading and writing (and creating, if necessary) */
result = MPI_File_open(MPI_COMM_WORLD, filename,
                      MPI_MODE_RDWR | MPI_MODE_CREATE, (int)NULL, &fh);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_File_open");

/* Set the file view which tiles the file type MPI_INT, starting
   at displacement 0. In this example, the etype is also MPI_INT. */
disp = 0;
etype = MPI_INT;
ftype = MPI_INT;
info = (MPI_Info)NULL;
result = MPI_File_set_view(fh, disp, etype, ftype, (char *)NULL, info);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_File_set_view");

/* Allocate and initialize a buffer (buff1) containing NUM_INTS integers,
   where the integer in location i is set to i. */
buff1 = (int *)malloc(NUM_INTS*sizeof(int));
for(i=0;i<NUM_INTS;i++) buff1[i] = i;

/* Set the buffer type to also be MPI_INT, then write the buffer (buff1)
```

**CODE EXAMPLE 4-1 Writing and Reading Integers to a File (Continued)**

```
    starting at offset 0, i.e., the first etype in the file. */
buftype = MPI_INT;
offset = rank * NUM_INTS;
result = MPI_File_write_at(fh, offset, buff1, NUM_INTS, buftype, &status);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_File_write_at");

result = MPI_Get_elements(&status, MPI_BYTE, &count);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_Get_elements");
if(count != NUM_INTS*sizeof(int))
    fprintf(stderr, "Did not write the same number of bytes as requested\n");
else
    fprintf(stdout, "Wrote %d bytes\n", count);

/* Allocate another buffer (buff2) to read into, then read NUM_INTS
   integers into this buffer. */
buff2 = (int *)malloc(NUM_INTS*sizeof(int));
result = MPI_File_read_at(fh, offset, buff2, NUM_INTS, buftype, &status);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_File_read_at");

/* Find out how many bytes were read and compare to how many
   we expected */
result = MPI_Get_elements(&status, MPI_BYTE, &count);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_Get_elements");
if(count != NUM_INTS*sizeof(int))
    fprintf(stderr, "Did not read the same number of bytes as requested\n");
else
    fprintf(stdout, "Read %d bytes\n", count);

/* Check to see that each integer read from each location is
   the same as the integer written to that location. */
differs = 0;
for(i=0; i<NUM_INTS; i++) {
    if(buff1[i] != buff2[i]) {
        fprintf(stderr, "Integer number %d differs\n", i);
        differs = 1;
    }
}
if(!differs)
    fprintf(stdout, "Wrote and read the same data\n");

MPI_Barrier(MPI_COMM_WORLD);

result = MPI_File_get_size(fh, &file_size);
```

#### CODE EXAMPLE 4-1 Writing and Reading Integers to a File (Continued)

```
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_File_get_size");

/* Compare the file size with what we expect */
/* To see a negative response, make the file preexist with a larger
   size than what is written by this program */
if(file_size != (comm_size * NUM_INTS * sizeof(int)))
    fprintf(stderr, "File size is not equal to the write size\n");

result = MPI_File_close(&fh);
if(result != MPI_SUCCESS)
    sample_error(result, "MPI_File_close");

MPI_Finalize();

free(buff1);
free(buff2);
}
```

## Data Access Styles

You can adapt CODE EXAMPLE 4-1 to support the I/O programming style that best suits your application. Essentially, there are three dimensions on which to choose an appropriate data access routine for a particular task: file pointer type, collective or noncollective, and blocking or nonblocking.

You need to choose which file pointer type to use: explicit, individual, or shared. The CODE EXAMPLE 4-1 used an explicit pointer and passed it directly as the offset parameter to the `MPI_File_write_at()` and `MPI_File_read_at()` routines. Using an explicit pointer is equivalent to calling `MPI_File_seek()` to set the individual file pointer to offset, and then calling `MPI_File_write()` or `MPI_File_read()`, which is directly analogous to calling UNIX `lseek()` and `write()` or `read()`. If each process accesses the file sequentially, individual file pointers save you the effort of recalculating offset for each data access. A different shared file pointer could be used in situations where all the processes needed to cooperatively access a file in a sequential way, such as to write log files.

Collective data-access routines enable you to enforce some implicit coordination among the processes in a parallel job when making data accesses. For example, if a parallel job alternately reads in a matrix and performs computation on it, but cannot progress to the next stage of computation until all processes have completed the last stage, then a coordinated effort between processes when accessing data might be more efficient. In CODE EXAMPLE 4-1, you could easily append the suffix `_all` to `MPI_File_write_at()` and `MPI_File_read_at()` to make the accesses collective. By coordinating the processes, you could achieve greater efficiency in the

MPI library or at the file system level in buffering or caching the next matrix. In contrast, noncollective accesses are used when it is not evident that any benefit would be gained by coordinating disparate accesses by each process. UNIX file accesses are noncollective.

## Overlapping I/O With Computation and Communication

MPI I/O also supports nonblocking versions of each of the data-access routines—that is, the data-access routines that have the letter *i* before *write* or *read* in the routine name (*i* stands for *immediate*). By definition, nonblocking I/O routines return immediately after the I/O request has been issued and do not wait until the I/O request has completed. This functionality enables you to perform computation and communication at the same time as the I/O. Because large I/O requests can take a long time to complete, this provides a way to more efficiently utilize your program's waiting time.

As in the previous example, parallel jobs often partition large matrices stored in files. These parallel jobs can use many large matrices, or matrices that are too large to fit into memory at once. Thus, each process can access the multiple and/or large matrices in stages. During each stage, a process reads in a chunk of data, and then performs a computation on it (which can involve communicating with the other processes in the parallel job). While performing the computation and communication, the process could issue a nonblocking I/O read request for the next chunk of data. Similarly, once the computation on a particular chunk has completed, a nonblocking write request could be issued before performing computation and communication on the next chunk.

The following example code illustrates the use of a nonblocking data-access routine. Note that like nonblocking communication routines, the nonblocking I/O routines require a call to `MPI_Wait()` to wait for the nonblocking request to complete, or repeated calls to `MPI_Test()` to determine when the nonblocking data access has completed. Once complete, the write or read buffer is available for use again by the program.

In CODE EXAMPLE 4-2, each process reads and writes `NUM_BYTES` bytes to a file using the nonblocking MPI I/O routines `MPI_File_iread_at()` and `MPI_File_iwrite_at()`, respectively. Note the use of `MPI_Wait()` and `MPI_Test()` to determine when the nonblocking requests have completed.

**CODE EXAMPLE 4-2** Reading and Writing Bytes to a File

```
/* iwr_at.c
 *
 * Example to demonstrate use of MPI_File_iwrite_at and MPI_File_iread_at
 *
 */

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

#define NUM_BYTES 100

void sample_error(int error, char *string)
{
    fprintf(stderr, "Error %d in %s\n", error, string);
    MPI_Finalize();
    exit(-1);
}

void
main( int argc, char **argv )
{
    char filename[128];
    char *buff;
    MPI_File fh;
    MPI_Offset offset;
    MPI_Request request;
    MPI_Status status;
    int i, rank, flag, result;

    if(argc < 2) {
        fprintf(stdout, "Missing argument: filename\n");
        exit(-1);
    }
    strcpy(filename, argv[1]);

    MPI_Init(&argc, &argv);

    result = MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    if(result != MPI_SUCCESS)
        sample_error(result, "MPI_Comm_rank");

    result = MPI_File_open(MPI_COMM_WORLD, filename,
                          MPI_MODE_RDWR | MPI_MODE_CREATE,
                          (MPI_Info)NULL, &fh);
    if(result != MPI_SUCCESS)
```

**CODE EXAMPLE 4-2** Reading and Writing Bytes to a File *(Continued)*

```
    sample_error(result, "MPI_File_open");

    buff = (char *)malloc(NUM_BYTES*sizeof(char));
    for(i=0;i<NUM_BYTES;i++) buff[i] = i;

    offset = rank * NUM_BYTES;
    result = MPI_File_iread_at(fh, offset, buff, NUM_BYTES,
                              MPI_BYTE, &request);
    if(result != MPI_SUCCESS)
        sample_error(result, "MPI_File_iread_at");

    /* Perform some useful computation and/or communication */

    result = MPI_Wait(&request, &status);

    buff = (char *)malloc(NUM_BYTES*sizeof(char));
    for(i=0;i<NUM_BYTES;i++) buff[i] = i;
    result = MPI_File_iwrite_at(fh, offset, buff, NUM_BYTES,
                                MPI_BYTE, &request);
    if(result != MPI_SUCCESS)
        sample_error(result, "MPI_File_iwrite_at");

    /* Perform some useful computation and/or communication */

    flag = 0;
    i = 0;
    while(!flag) {
        result = MPI_Test(&request, &flag, &status);
        i++;
        /* Perform some more computation or communication, if possible */
    }

    result = MPI_File_close(&fh);
    if(result != MPI_SUCCESS)
        sample_error(result, "MPI_File_close");

    MPI_Finalize();

    fprintf(stdout, "Successful completion\n");

    free(buff);
}
```

## Sun MPI and Sun MPI I/O Routines

---

The tables in this appendix list the Sun MPI and Sun MPI I/O routines, along with the C syntax of the routines and a brief description of each. The routines are found in two sections:

- “Sun MPI Routines” on page 67
- “Sun MPI I/O Routines” on page 95

For more information about the routines, see their online man pages, usually found in

`/opt/SUNWhpc/man/`

Your system administrator can tell you where they are installed at your site.

---

## Sun MPI Routines

TABLE A-1 lists the Sun MPI routines in alphabetical order. The following sections list the routines by functional category.

# Point-to-Point Communication

## Blocking Routines

```
MPI_Send()  
MPI_Bsend()  
MPI_Ssend()  
MPI_Rsend()  
MPI_Recv()  
MPI_Sendrecv()  
MPI_Sendrecv_replace()
```

## Nonblocking Routines

```
MPI_Isend()  
MPI_Ibsend()  
MPI_Issend()  
MPI_Irsend()  
MPI_Irecv()
```

## Communication Buffer Allocation

```
MPI_Buffer_attach()  
MPI_Buffer_detach()
```

## Status Data Structure

```
MPI_Get_count()  
MPI_Get_elements()
```



## Persistent (Half-Channel) Communication

```
MPI_Send_init()  
MPI_Bsend_init()  
MPI_Rsend_init()  
MPI_Ssend_init()  
MPI_Recv_init()  
MPI_Start()  
MPI_Startall()
```

## Completion Tests

```
MPI_Wait()  
MPI_Waitany()  
MPI_Waitsome()  
MPI_Waitall()  
MPI_Test()  
MPI_Testany()  
MPI_Testsome()  
MPI_Testall()  
MPI_Request_free()  
MPI_Cancel()  
MPI_Test_cancelled()
```

## Probing for Messages (Blocking and Nonblocking)

```
MPI_Probe()  
MPI_Iprobe()
```

## Packing and Unpacking Functions

```
MPI_Pack()  
MPI_Unpack()  
MPI_Pack_size()  
MPI_Pack_external()  
MPI_Unpack_external()  
MPI_Pack_external_size()
```

## Derived Data Type Constructors and Functions

```
MPI_Address(): Deprecated - Use MPI_Get_address()  
MPI_Type_commit()  
MPI_Type_contiguous()  
MPI_Type_create_f90_complex()  
MPI_Type_create_f90_integer()  
MPI_Type_create_F90_real()  
MPI_Type_match_size()  
MPI_Type_Sizeof()  
MPI_Type_create_indexed_block()  
MPI_Type_create_keyval()  
MPI_Type_create_resized()  
MPI_Type_delete_attr()  
MPI_Type_dup()  
MPI_Type_free()  
MPI_Type_free_keyval()  
MPI_Type_get_attr()  
MPI_Type_set_attr()  
MPI_Type_get_contents()  
MPI_Type_get_envelope()  
MPI_Type_get_name()  
MPI_Type_get_true_extent()  
MPI_Type_set_name()  
MPI_Type_hvector(): Deprecated - Use MPI_Type_create_hvector()  
MPI_Type_indexed()  
MPI_Type_hindexed(): Deprecated - Use MPI_Type_create_hindexed()  
MPI_Type_struct(): Deprecated - Use MPI_Type_create_struct()  
MPI_Type_lb(): Deprecated - Use MPI_Type_get_extent()  
MPI_Type_ub(): Deprecated - Use MPI_Type_get_extent()  
MPI_Type_vector()  
MPI_Type_extent(): Deprecated - Use MPI_Type_get_extent()  
MPI_Type_size()
```

## One-Sided Communication

### Initialization

```
MPI_Win_create()  
MPI_Win_free()  
MPI_Win_get_group()
```

## Communication Calls

```
MPI_Put()  
MPI_Get()  
MPI_Accumulate()
```

## Synchronization Calls

```
MPI_Win_fence()  
MPI_Win_lock()  
MPI_Win_unlock()  
MPI_Win_start()  
MPI_Win_complete()  
MPI_Win_post()  
MPI_Win_wait()  
MPI_Win_test()
```

## Collective Communication

### Barrier

```
MPI_Barrier()
```

### Broadcast

```
MPI_Bcast()
```

## Processor Gather and Scatter

```
MPI_Gather()  
MPI_Gatherv()  
MPI_Allgather()  
MPI_Allgatherv()  
MPI_Scatter()  
MPI_Scatterv()  
MPI_Alltoall()  
MPI_Alltoallv()  
MPI_Alltoallw()
```

## Global Reduction and Scan Operations

```
MPI_Reduce()  
MPI_Allreduce()  
MPI_Reduce_scatter()  
MPI_Scan()  
MPI_Exscan()  
MPI_Op_create()  
MPI_Op_free()
```

## Groups and Communicators

### Group Management

#### *Group Accessors*

```
MPI_Group_size()  
MPI_Group_rank()  
MPI_Group_translate_ranks()  
MPI_Group_compare()
```

## *Group Constructors*

```
MPI_Comm_group()  
MPI_Group_union()  
MPI_Group_intersection()  
MPI_Group_difference()  
MPI_Group_incl()  
MPI_Group_excl()  
MPI_Group_range_incl()  
MPI_Group_range_excl()  
MPI_Group_free()
```

## Communicator Management

### *Communicator Accessors*

```
MPI_Comm_size()  
MPI_Comm_rank()  
MPI_Comm_compare()
```

### *Communicator Constructors*

```
MPI_Comm_dup()  
MPI_Comm_create()  
MPI_Comm_split()  
MPI_Comm_free()
```

### *Intercommunicators*

```
MPI_Comm_test_inter()  
MPI_Comm_remote_group()  
MPI_Comm_remote_size()  
MPI_Intercomm_create()  
MPI_Intercomm_merge()
```

## *Communicator Attributes*

`MPI_Keyval_create()`: *Deprecated* – Use `MPI_Comm_create_keyval()`  
`MPI_Keyval_free()`: *Deprecated* – Use `MPI_Comm_free_keyval()`  
`MPI_Attr_put()`: *Deprecated* – Use `MPI_Comm_set_attr()`  
`MPI_Attr_get()`: *Deprecated* – Use `MPI_Comm_get_attr()`  
`MPI_Attr_delete()`: *Deprecated* – Use `MPI_Comm_delete_attr()`

## Process Topologies

`MPI_Cart_create()`  
`MPI_Dims_create()`  
`MPI_Graph_create()`  
`MPI_Topo_test()`  
`MPI_Graphdims_get()`  
`MPI_Graph_get()`  
`MPI_Cartdim_get()`  
`MPI_Cart_get()`  
`MPI_Cart_rank()`  
`MPI_Cart_coords()`  
`MPI_Graph_neighbors()`  
`MPI_Graph_neighbors_count()`  
`MPI_Cart_shift()`  
`MPI_Cart_sub()`  
`MPI_Cart_map()`  
`MPI_Graph_map()`

## Process Creation and Management

### Establishing Communication

`MPI_Close_port()`  
`MPI_Comm_accept()`  
`MPI_Comm_connect()`  
`MPI_Comm_disconnect()`  
`MPI_Open_port()`  
`MPI_Comm_join()`

## Name Publishing

```
MPI_Publish_name()  
MPI_Unpublish_name()  
MPI_Lookup_name()
```

## Process Manager Interface

```
MPI_Comm_get_parent()  
MPI_Comm_spawn()  
MPI_Comm_spawn_multiple()
```

## Environmental Inquiry Functions and Profiling

### Startup and Shutdown

```
MPI_Init()  
MPI_Finalize()  
MPI_Finalized()  
MPI_Initialized()  
MPI_Abort()  
MPI_Get_processor_name()  
MPI_Get_version()
```

## Error Handler Functions

```
MPI_Add_error_class()  
MPI_Add_error_code()  
MPI_Add_error_string()  
MPI_Comm_call_errhandler()  
MPI_File_call_errhandler()  
MPI_Win_call_errhandler()  
MPI_Errhandler_create(): Deprecated - Use MPI_Comm_create_errhandler()  
MPI_Errhandler_set(): Deprecated - Use MPI_Comm_set_errhandler()  
MPI_Errhandler_get(): Deprecated - Use MPI_Comm_get_errhandler()  
MPI_Errhandler_free()  
MPI_Error_string()  
MPI_Error_class()
```

## Info Objects

```
MPI_Info_create()  
MPI_Info_delete()  
MPI_Info_dup()  
MPI_Info_free()  
MPI_Info_get()  
MPI_Info_get_nkeys()  
MPI_Info_get_nthkey()  
MPI_Info_get_valuelen()  
MPI_Info_set()
```

## Timers

```
MPI_Wtime()  
MPI_Wtick()
```

## Profiling

```
MPI_Pcontrol()
```

## Miscellaneous



## Associating Information With Status

```
MPI_Status_set_cancelled()  
MPI_Status_set_elements()
```

## Generalized Requests

```
MPI_Grequest_complete()  
MPI_Grequest_start()
```

## Naming Objects

```
MPI_Comm_get_name()  
MPI_Comm_set_name()  
MPI_Type_get_name()  
MPI_Type_set_name()
```

## Threads

```
MPI_Query_thread()  
MPI_Init_thread()  
MPI_Is_thread_main()
```

## Handle Translation

MPI\_Comm\_c2f()  
MPI\_Comm\_f2c()  
MPI\_Group\_c2f()  
MPI\_Group\_f2c()  
MPI\_Info\_c2f()  
MPI\_Info\_f2c()  
MPI\_Op\_c2f()  
MPI\_Op\_f2c()  
MPI\_Request\_c2f()  
MPI\_Request\_f2c()  
MPI\_Type\_c2f()  
MPI\_Type\_f2c()

## Status Conversion

MPI\_Status\_c2f()  
MPI\_Status\_f2c()

## MPI Routines: Alphabetical Listing

TABLE A-1 Sun MPI Routines

Routine and C Syntax	Description
<b>MPI_Abort</b> (MPI_Comm <i>comm</i> , int <i>errorcode</i> )	Terminates MPI execution environment.
<b>MPI_Accumulate</b> (void * <i>origin_addr</i> , int <i>origin_count</i> , MPI_Datatype <i>origin_datatype</i> , int <i>target_rank</i> , MPI_Aint <i>target_disp</i> , int <i>target_count</i> , MPI_Datatype <i>target_datatype</i> , MPI_Op <i>op</i> , MPI_Win <i>win</i> )	Combines the contents of the origin buffer with that of a target buffer.
<b>MPI_Add_error_class</b> (int * <i>errorclass</i> )	Creates a new, local error class.
<b>MPI_Add_error_code</b> (int <i>errorclass</i> , int * <i>errorcode</i> )	Creates a new error code and associates it with an error class.
<b>MPI_Add_error_string</b> (int * <i>errorcode</i> , char * <i>string</i> )	Creates an error string and associates it with an error code or an error class.

TABLE A-1 Sun MPI Routines (Continued)

Routine and C Syntax	Description
<b>MPI_Address</b> (void *location, MPI_Aint *address)	<i>Deprecated:</i> Use instead <code>MPI_Get_address()</code> . Gets the address of a location in memory.
<b>MPI_Allgather</b> (void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcnt, MPI_Datatype recvtpe, MPI_Comm comm)	Gathers data from all processes and distributes it to all.
<b>MPI_Allgatherv</b> (void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int *recvcnt, int *displs, MPI_Datatype recvtpe, MPI_Comm comm)	Gathers data from all processes and delivers it to all. Each process can contribute a different amount of data.
<b>MPI_Alloc_mem</b> (MPI_Aint size, MPI_Info info, void *baseptr)	Allocates a specified memory segment.
<b>MPI_Allreduce</b> (void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)	Combines values from all processes and distributes the result back to all processes.
<b>MPI_Alltoall</b> (void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcnt, MPI_Datatype recvtpe, MPI_Comm comm)	All processes send data to, and receive data from, all processes.
<b>MPI_Alltoallv</b> (void *sendbuf, int *sendcounts, int *sdispls, MPI_Datatype sendtype, void *recvbuf, int *recvcnts, int *rdispls, MPI_Datatype recvtpe, MPI_Comm comm)	All processes send data to, and receive data from, all processes, but user can specify different amounts of data to send and receive.
<b>MPI_Alltoallw</b> (void *sendbuf, int *sendcounts, int *sdispls, MPI_Datatype sendtypes, void *recvbuf, int *recvcnts, int *rdispls, MPI_Datatype recvtypes, MPI_Comm comm)	All processes send data to, and receive data from, all other processes, and user can specify database of individual datablocks of different types.
<b>MPI_Attr_delete</b> (MPI_Comm comm, int keyval)	<i>Deprecated:</i> Use instead <code>MPI_Comm_delete_attr()</code> . Deletes attribute value associated with a key.
<b>MPI_Attr_get</b> (MPI_Comm comm, int keyval, void *attribute_val, int *flag)	<i>Deprecated:</i> Use instead <code>MPI_Comm_get_attr()</code> . Retrieves attribute value by key.
<b>MPI_Attr_put</b> (MPI_Comm comm, int keyval, void *attribute_val)	<i>Deprecated:</i> Use instead <code>MPI_Comm_set_attr()</code> . Stores attribute value associated with a key.
<b>MPI_Barrier</b> (MPI_Comm comm)	Blocks until all processes have reached this routine.

**TABLE A-1** Sun MPI Routines (*Continued*)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Bcast</b> (void * <i>buffer</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>root</i> , MPI_Comm <i>comm</i> )	Broadcasts a message from the process with rank <code>root</code> to all other processes of the group.
<b>MPI_Bsend</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> )	Basic send with user-specified buffering.
<b>MPI_Bsend_init</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Builds a handle for a buffered send.
<b>MPI_Buffer_attach</b> (void * <i>buf</i> , int <i>size</i> )	Attaches a user-defined buffer for sending.
<b>MPI_Buffer_detach</b> (void * <i>buf</i> , int * <i>size</i> )	Removes an existing buffer (for use in <code>MPI_Bsend()</code> , etc.).
<b>MPI_Cancel</b> (MPI_Request * <i>request</i> )	Cancels a communication request.
<b>MPI_Cart_coords</b> (MPI_Comm <i>comm</i> , int <i>rank</i> , int <i>maxdims</i> , int * <i>coords</i> )	Determines process coordinates in Cartesian topology given rank in group.
<b>MPI_Cart_create</b> (MPI_Comm <i>comm_old</i> , int <i>ndims</i> , int * <i>dims</i> , int * <i>periods</i> , int <i>reorder</i> , MPI_Comm * <i>comm_cart</i> )	Makes a new communicator to which Cartesian topology information has been attached.
<b>MPI_Cart_get</b> (MPI_Comm <i>comm</i> , int <i>maxdims</i> , int * <i>dims</i> , int * <i>periods</i> , int * <i>coords</i> )	Retrieves Cartesian topology information associated with a communicator.
<b>MPI_Cart_map</b> (MPI_Comm <i>comm</i> , int <i>ndims</i> , int * <i>dims</i> , int * <i>periods</i> , int * <i>newrank</i> )	Maps process to Cartesian topology information.
<b>MPI_Cart_rank</b> (MPI_Comm <i>comm</i> , int * <i>coords</i> , int * <i>rank</i> )	Determines process rank in communicator given Cartesian location.
<b>MPI_Cart_shift</b> (MPI_Comm <i>comm</i> , int <i>direction</i> , int <i>disp</i> , int * <i>rank_source</i> , int * <i>rank_dest</i> )	Returns the shifted source and destination ranks, given a shift direction and amount.
<b>MPI_Cart_sub</b> (MPI_Comm <i>comm</i> , int * <i>remain_dims</i> , MPI_Comm * <i>comm_new</i> )	Partitions a communicator into subcommunicators that form lower-dimensional Cartesian subgrids.
<b>MPI_Cartdim_get</b> (MPI_Comm <i>comm</i> , int * <i>ndims</i> )	Retrieves Cartesian topology information associated with a communicator.
<b>MPI_Close_port</b> (char * <i>port_name</i> )	Releases the specified network address.

TABLE A-1 Sun MPI Routines (Continued)

Routine and C Syntax	Description
<b>MPI_Comm_accept</b> (char * <i>port_name</i> , MPI_Info <i>info</i> , int <i>root</i> , MPI_Comm <i>comm</i> , MPI_Comm * <i>newcomm</i> )	Establishes communication with a client (collective).
<b>MPI_Comm_c2f</b> (MPI_Comm <i>comm</i> )	Translates a C handle into a Fortran handle.
<b>MPI_Comm_compare</b> (MPI_Comm <i>comm1</i> , MPI_Comm <i>comm2</i> , int * <i>result</i> )	Compares two communicators.
<b>MPI_Comm_connect</b> (char * <i>port_name</i> , MPI_Info <i>info</i> , int <i>root</i> , MPI_Comm <i>comm</i> , MPI_Comm * <i>newcomm</i> )	Establishes communication with a server (collective).
<b>MPI_Comm_create</b> (MPI_Comm <i>comm</i> , MPI_Group <i>group</i> , MPI_Comm * <i>newcomm</i> )	Creates a new communicator from a group.
<b>MPI_Comm_create_errhandler</b> (MPI_Comm_errhandler_fn * <i>function</i> , MPI_Errhandler * <i>errhandler</i> )	Creates an error handler that can be attached to communicators.
<b>MPI_Comm_create_keyval</b> (MPI_Comm_copy_attr_function * <i>comm_copy_attr_fn</i> , MPI_Comm_delete_attr_function * <i>comm_delete_attr_fn</i> , int * <i>comm_keyval</i> , void * <i>extra_state</i> )	Generates a new attribute key.
<b>MPI_Comm_delete_attr</b> (MPI_Comm <i>comm</i> , int <i>comm_keyval</i> )	Deletes attribute value associated with a key.
<b>MPI_Comm_disconnect</b> (MPI_Comm * <i>comm</i> )	De-allocates communicator object and sets handle to MPI_COMM_NULL (collective).
<b>MPI_Comm_dup</b> (MPI_Comm <i>comm</i> , MPI_Comm * <i>newcomm</i> )	Duplicates an existing communicator with all its cached information.
<b>MPI_Comm_f2c</b> (MPI_Fint <i>comm</i> )	Translates a Fortran handle into a C handle.
<b>MPI_Comm_free</b> (MPI_Comm * <i>comm</i> )	Marks the communicator object for de-allocation.
<b>MPI_Comm_free_keyval</b> (int * <i>comm_keyval</i> )	Frees attribute key for communicator cache attribute.
<b>MPI_Comm_get_attr</b> (MPI_Comm <i>comm</i> , int <i>comm_keyval</i> , void * <i>attribute_val</i> , int * <i>flag</i> )	Retrieves attribute value by key.
<b>MPI_Comm_get_errhandler</b> (MPI_Comm <i>comm</i> , MPI_Errhandler * <i>errhandler</i> )	Retrieves error handler associated with a communicator.

TABLE A-1 Sun MPI Routines (Continued)

Routine and C Syntax	Description
<b>MPI_Comm_get_name</b> (MPI_Comm <i>comm</i> , char * <i>comm_name</i> , int * <i>resultlen</i> )	Returns the name that was most recently associated with a communicator.
<b>MPI_Comm_get_parent</b> (MPI_Comm * <i>parent</i> )	Returns the parent intercommunicator of current spawned process.
<b>MPI_Comm_group</b> (MPI_Comm <i>comm</i> , MPI_Group * <i>group</i> )	Accesses the group associated with a communicator.
<b>MPI_Comm_join</b> (int <i>fd</i> , MPI_Comm * <i>intercomm</i> )	Creates an intercommunicator from the union of two MPI processes connected by a socket.
<b>MPI_Comm_rank</b> (MPI_Comm <i>comm</i> , int * <i>rank</i> )	Determines the rank of the calling process in a communicator.
<b>MPI_Comm_remote_group</b> (MPI_Comm <i>comm</i> , MPI_Group * <i>group</i> )	Accesses the remote group associated with an intercommunicator.
<b>MPI_Comm_remote_size</b> (MPI_Comm <i>comm</i> , int * <i>size</i> )	Determines the size of the remote group associated with an intercommunicator.
<b>MPI_Comm_set_attr</b> (MPI_Comm <i>comm</i> , int <i>comm_keyval</i> , void * <i>attribute_val</i> )	Stores attribute value associated with a key.
<b>MPI_Comm_set_errhandler</b> (MPI_Comm <i>comm</i> , MPI_Errhandler * <i>errhandler</i> )	Attaches a new error handler to a communicator.
<b>MPI_Comm_set_name</b> (MPI_Comm <i>comm</i> , char * <i>comm_name</i> )	Associates a name with a communicator.
<b>MPI_Comm_size</b> (MPI_Comm <i>comm</i> , int * <i>size</i> )	Determines the size of the group associated with a communicator.
<b>MPI_Comm_spawn</b> (char * <i>command</i> , char * <i>argv</i> [], int <i>maxprocs</i> , MPI_Info <i>info</i> , int <i>root</i> , MPI_Comm <i>comm</i> , MPI_Comm * <i>intercomm</i> , int <i>array_of_errcodes</i> [])	Spawns a number of identical binaries.
<b>MPI_Comm_spawn_multiple</b> (int <i>count</i> , char * <i>array_of_commands</i> [], char ** <i>array_of_argv</i> [], int <i>array_of_maxprocs</i> [], MPI_Info <i>array_of_info</i> [], int <i>root</i> , MPI_Comm <i>comm</i> , MPI_Comm * <i>intercomm</i> , int <i>array_of_errcodes</i> [])	Spawns multiple binaries, or the same binary with multiple sets of arguments.
<b>MPI_Comm_split</b> (MPI_Comm <i>comm</i> , int <i>color</i> , int <i>key</i> , MPI_Comm * <i>newcomm</i> )	Creates new communicators based on colors and keys.
<b>MPI_Comm_test_inter</b> (MPI_Comm <i>comm</i> , int * <i>flag</i> )	Tests whether a communicator is an intercommunicator.

TABLE A-1 Sun MPI Routines (Continued)

Routine and C Syntax	Description
<b>MPI_Dims_create</b> (int <i>nnodes</i> , int <i>ndims</i> , int <i>*dims</i> )	Creates a division of processors in a Cartesian grid.
<b>MPI_Errhandler_create</b> ( MPI_Handler_function <i>*function</i> , MPI_Errhandler <i>*errhandler</i> )	<i>Deprecated:</i> Use instead <code>MPI_Comm_create_errhandler()</code> . Creates an MPI error handler.
<b>MPI_Errhandler_free</b> (MPI_Errhandler <i>*errhandler</i> )	Frees an MPI error handler.
<b>MPI_Errhandler_get</b> (MPI_Comm <i>comm</i> , MPI_Errhandler <i>*errhandler</i> )	<i>Deprecated:</i> Use instead <code>MPI_Comm_get_errhandler()</code> . Gets the error handler for a communicator.
<b>MPI_Errhandler_set</b> (MPI_Comm <i>comm</i> , MPI_Errhandler <i>errhandler</i> )	<i>Deprecated:</i> Use instead <code>MPI_Comm_set_errhandler()</code> . Sets the error handler for a communicator.
<b>MPI_Error_class</b> (int <i>errorcode</i> , int <i>*errorclass</i> )	Converts an error code into an error class.
<b>MPI_Error_string</b> (int <i>errorcode</i> , char <i>*string</i> , int <i>*resultlen</i> )	Returns a string for a given error code.
<b>MPI_Exscan</b> (void <i>*sendbuf</i> , void <i>*recvbuf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Op <i>op</i> , MPI_Comm, <i>comm</i> )	Performs an exclusive prefix reduction on data distributed across the calling processes.
<b>MPI_Finalize</b> ( )	Terminates MPI execution environment.
<b>MPI_Finalized</b> (int <i>*flag</i> )	Checks whether <code>MPI_Finalize()</code> has completed.
<b>MPI_Free_mem</b> (void <i>*base</i> )	Frees memory that has been allocated using <code>MPI_Alloc_mem</code> .
<b>MPI_Gather</b> (void <i>*sendbuf</i> , int <i>*sendcount</i> , MPI_Datatype <i>sendtype</i> , void <i>*recvbuf</i> , int <i>recvcount</i> , MPI_Datatype <i>recvtype</i> , int <i>root</i> , MPI_Comm <i>comm</i> )	Gathers values from a group of processes.
<b>MPI_Gatherv</b> (void <i>*sendbuf</i> , int <i>sendcount</i> , MPI_Datatype <i>sendtype</i> , void <i>*recvbuf</i> , int <i>*recvcounts</i> , int <i>*displs</i> , MPI_Datatype <i>recvtype</i> , int <i>root</i> , MPI_Comm <i>comm</i> )	Gathers into specified locations from all processes in a group. Each process can contribute a different amount of data.
<b>MPI_Get</b> (void <i>*origin_addr</i> , int <i>origin_count</i> , MPI_Datatype <i>origin_datatype</i> , int <i>target_rank</i> , MPI_Aint <i>target_disp</i> , int <i>target_count</i> , MPI_Datatype <i>target_datatype</i> , MPI_Win <i>win</i> )	Copies data from the target memory to the origin.

**TABLE A-1** Sun MPI Routines (Continued)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Get_address</b> (void *location, MPI_Aint *address)	Gets the address of a location in memory.
<b>MPI_Get_count</b> (MPI_Status *status, MPI_Datatype datatype, int *count)	Gets the number of top-level elements received.
<b>MPI_Get_elements</b> (MPI_Status *status, MPI_Datatype datatype, int *count)	Returns the number of basic elements in a data type.
<b>MPI_Get_processor_name</b> (char *name, int *resultlen)	Gets the name of the processor.
<b>MPI_Get_version</b> (int *version, int *subversion)	Returns the version of the standard corresponding to the current implementation.
<b>MPI_Graph_create</b> (MPI_Comm comm_old, int nnodes, int *index, int *edges, int reorder, MPI_Comm *comm_graph)	Makes a new communicator to which graph topology information has been attached.
<b>MPI_Graph_get</b> (MPI_Comm comm, int maxindex, int maxedges, int *index, int *edges)	Retrieves graph topology information associated with a communicator.
<b>MPI_Graph_map</b> (MPI_Comm comm, int nnodes, int *index, int *edges, int *newrank)	Maps process to graph topology information.
<b>MPI_Graph_neighbors</b> (MPI_Comm comm, int rank, int maxneighbors, int *neighbors)	Returns the neighbors of a node associated with a graph topology.
<b>MPI_Graph_neighbors_count</b> (MPI_Comm comm, int rank, int *nneighbors)	Returns the number of neighbors of a node associated with a graph topology.
<b>MPI_Graphdims_get</b> (MPI_Comm comm, int *nnodes, int *nedges)	Retrieves graph topology information associated with a communicator.
<b>MPI_Grequest_complete</b> (MPI_Request request)	Reports that a generalized request is complete.
<b>MPI_Grequest_start</b> (MPI_Grequest_query_function *query_fn, MPI_Grequest_free_function *free_fn, MPI_Grequest_cancel_function *cancel_fn, void *extra_state, MPI_Request *request)	Starts a generalized request and returns a handle to it.
<b>MPI_Group_c2f</b> (MPI_Group group)	Translates a C handle into a Fortran handle.
<b>MPI_Group_compare</b> (MPI_Group group1, MPI_Group group2, int *result)	Compares two groups.



**TABLE A-1 Sun MPI Routines (Continued)**

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Group_difference</b> (MPI_Group <i>group1</i> , MPI_Group <i>group2</i> , MPI_Group * <i>group_out</i> )	Makes a group from the difference of two groups.
<b>MPI_Group_excl</b> (MPI_Group <i>group</i> , int <i>n</i> , int * <i>ranks</i> , MPI_Group * <i>newgroup</i> )	Produces a group by reordering an existing group and taking only unlisted members.
<b>MPI_Group_f2c</b> (MPI_Fint <i>group</i> )	Translates a Fortran handle into a C handle.
<b>MPI_Group_free</b> (MPI_Group * <i>group</i> )	Frees a group.
<b>MPI_Group_incl</b> (MPI_Group <i>group</i> , int <i>n</i> , int * <i>ranks</i> , MPI_Group * <i>group_out</i> )	Produces a group by reordering an existing group and taking only listed members.
<b>MPI_Group_intersection</b> (MPI_Group <i>group1</i> , MPI_Group <i>group2</i> , MPI_Group * <i>group_out</i> )	Produces a group at the intersection of two existing groups.
<b>MPI_Group_range_excl</b> (MPI_Group <i>group</i> , int <i>n</i> , int <i>ranges</i> [ ][3], MPI_Group * <i>newgroup</i> )	Produces a group by excluding ranges of processes from an existing group.
<b>MPI_Group_range_incl</b> (MPI_Group <i>group</i> , int <i>n</i> , int <i>ranges</i> [ ][3], MPI_Group * <i>newgroup</i> )	Creates a new group from ranges of ranks in an existing group.
<b>MPI_Group_rank</b> (MPI_Group <i>group</i> , int * <i>rank</i> )	Returns the rank of this process in the given group.
<b>MPI_Group_size</b> (MPI_Group <i>group</i> , int * <i>size</i> )	Returns the size of a group.
<b>MPI_Group_translate_ranks</b> (MPI_Group <i>group1</i> , int <i>n</i> , int * <i>ranks1</i> , MPI_Group <i>group2</i> , int * <i>ranks2</i> )	Translates the ranks of processes in one group to those in another group.
<b>MPI_Group_union</b> (MPI_Group <i>group1</i> , MPI_Group <i>group2</i> , MPI_Group * <i>group_out</i> )	Produces a group by combining two groups.
<b>MPI_IbSEND</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Starts a nonblocking buffered send.
<b>MPI_Info_c2f</b> (MPI_Info <i>info</i> )	Translates a C handle into a Fortran handle.
<b>MPI_Info_create</b> (MPI_Info * <i>info</i> )	Creates a new info object.
<b>MPI_Info_delete</b> (MPI_Info * <i>info</i> , char * <i>key</i> , char * <i>value</i> )	Deletes a key/value pair from <i>info</i> .
<b>MPI_Info_dup</b> (MPI_Info <i>info</i> , MPI_Info * <i>newinfo</i> )	Duplicates an info object.
<b>MPI_Info_f2c</b> (MPI_Fint <i>info</i> )	Translates a Fortran handle into a C handle.

**TABLE A-1** Sun MPI Routines (Continued)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Info_free</b> (MPI_Info * <i>info</i> )	Frees <i>info</i> and sets it to MPI_INFO_NULL.
<b>MPI_Info_get</b> (MPI_Info * <i>info</i> , char * <i>key</i> , char * <i>value</i> )	Retrieves key value for an <i>info</i> object.
<b>MPI_Info_get_nkeys</b> (MPI_Info <i>info</i> , int * <i>nkeys</i> )	Returns the number of currently defined keys in <i>info</i> .
<b>MPI_Info_get_nthkey</b> (MPI_Info <i>info</i> , int <i>n</i> , char * <i>key</i> )	Returns the <i>nth</i> defined key in <i>info</i> .
<b>MPI_Info_get_valuelen</b> (MPI_Info <i>info</i> , char * <i>key</i> , int * <i>valuelen</i> , int * <i>flag</i> )	Retrieves the length of the key value associated with an <i>info</i> object.
<b>MPI_Info_set</b> (MPI_Info * <i>info</i> , char * <i>key</i> , char * <i>value</i> )	Adds a key/value pair to <i>info</i> .
<b>MPI_Init</b> (int * <i>argc</i> , char *** <i>argv</i> )	Initializes the MPI execution environment.
<b>MPI_Initialized</b> (int * <i>flag</i> )	Indicates whether <code>MPI_Init()</code> has been called.
<b>MPI_Init_thread</b> (int * <i>argc</i> , char *** <i>argv</i> , int <i>required</i> , int * <i>provided</i> )	Initializes the MPI execution environment with a predetermined level of thread support. Thread that calls this function becomes the main thread. Call instead of <code>MPI_Init()</code> , not in addition to.
<b>MPI_Intercomm_create</b> (MPI_Comm <i>local_comm</i> , int <i>local_leader</i> , MPI_Comm <i>peer_comm</i> , int <i>remote_leader</i> , int <i>tag</i> , MPI_Comm * <i>newintercomm</i> )	Creates an intercommunicator.
<b>MPI_Intercomm_merge</b> (MPI_Comm <i>intercomm</i> , int <i>high</i> , MPI_Comm * <i>newintracomm</i> )	Creates an intracommunicator from an intercommunicator.
<b>MPI_Iprobe</b> (int <i>source</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , int * <i>flag</i> , MPI_Status * <i>status</i> )	Nonblocking test for a message.
<b>MPI_Irecv</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>source</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Begins a nonblocking receive.
<b>MPI_Irsend</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Begins a nonblocking ready send.
<b>MPI_Isend</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Begins a nonblocking send.

TABLE A-1 Sun MPI Routines (Continued)

Routine and C Syntax	Description
<b>MPI_Issend</b> (void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Begins a nonblocking synchronous send.
<b>MPI_Is_thread_main</b> (int *flag)	Called by a thread to determine whether it is the main thread. See MPI_Init_thread().
<b>MPI_Keyval_create</b> (MPI_Copy_function *copy_fn, MPI_Delete_function *delete_fn, int *keyval, void *extra_state)	<i>Deprecated:</i> Use instead MPI_Comm_create_keyval(). Generates a new attribute key.
<b>MPI_Keyval_free</b> (int *keyval)	<i>Deprecated:</i> Use instead MPI_Comm_free_keyval(). Frees attribute key for communicator cache attribute.
<b>MPI_Lookup_name</b> (char *service-name, MPI_Info info, char *port-name)	Retrieves the port name associated with a service-name published by MPI_Publish_name().
<b>MPI_Op_c2f</b> (MPI_Op op)	Translates a C handle into a Fortran handle.
<b>MPI_Op_create</b> (MPI_User_function *function, int commute, MPI_Op *op)	Creates a user-defined combination function handle.
<b>MPI_Op_f2c</b> (MPI_Fint op)	Translates a Fortran handle into a C handle.
<b>MPI_Op_free</b> (MPI_Op *op)	Frees a user-defined combination function handle.
<b>MPI_Open_port</b> (MPI_Info info, char *port_name)	Establishes a network address for a server to accept connections from clients.
<b>MPI_Pack</b> (void *inbuf, int incount, MPI_Datatype datatype, void *outbuf, int outsize, int *position, MPI_Comm comm)	Packs data of a given data type into contiguous memory.
<b>MPI_Pack_external</b> (char *datarep, void *inbuf, int incount, MPI_Datatype datatype, void *outbuf, MPI_Aint outsize, MPI_Aint *position)	Packs data into the external32 format, used to exchange data between MPI implementations, or when writing data to a file.
<b>MPI_Pack_external_size</b> (char *datarep, int incount, MPI_Datatype datatype, MPI_Aint *size)	Returns the upper bound on the amount of space necessary to pack a message in the external32 format.
<b>MPI_Pack_size</b> (int incount, MPI_Datatype datatype, MPI_Comm comm, int *size)	Returns the upper bound on the amount of space necessary to pack a message.

**TABLE A-1** Sun MPI Routines (*Continued*)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Pcontrol</b> (int <i>level</i> , ...)	Controls profiling.
<b>MPI_Probe</b> (int <i>source</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Status * <i>status</i> )	Blocking test for a message.
<b>MPI_Publish_name</b> (char * <i>service-name</i> , MPI_Info <i>info</i> , char * <i>port-name</i> )	Publishes the pair ( <i>service-name</i> , <i>port-name</i> ) so that applications can use MPI_Lookup_name( ) to find <i>port-name</i> .
<b>MPI_Put</b> (void * <i>origin_addr</i> , int <i>origin_count</i> , MPI_Datatype <i>origin_datatype</i> , int <i>target_rank</i> , MPI_Aint <i>target_disp</i> , int <i>target_count</i> , MPI_Datatype <i>target_datatype</i> , MPI_Win <i>win</i> )	Copies data from the origin memory to the target.
<b>MPI_Query_thread</b> (int * <i>provided</i> )	Returns the current level of thread support.
<b>MPI_Recv</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>source</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Status * <i>status</i> )	Performs a standard receive.
<b>MPI_Recv_init</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>source</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Builds a persistent receive request handle.
<b>MPI_Reduce</b> (void * <i>sendbuf</i> , void * <i>recvbuf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Op <i>op</i> , int <i>root</i> , MPI_Comm <i>comm</i> )	Reduces values on all processes to a single value.
<b>MPI_Reduce_scatter</b> (void * <i>sendbuf</i> , void * <i>recvbuf</i> , int * <i>recvcounts</i> , MPI_Datatype <i>datatype</i> , MPI_Op <i>op</i> , MPI_Comm <i>comm</i> )	Combines values and scatters the results.
<b>MPI_Request_c2f</b> (MPI_Request <i>request</i> )	Translates a C handle into a Fortran handle.
<b>MPI_Request_f2c</b> (MPI_Fint <i>request</i> )	Translates a Fortran handle into a C handle.
<b>MPI_Request_free</b> (MPI_Request * <i>request</i> )	Frees a communication request object.
<b>MPI_Request_get_status</b> (MPI_Request <i>request</i> , int * <i>flag</i> , MPI_Status * <i>status</i> )	Accesses information associated with a request without freeing the request.
<b>MPI_Rsend</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> )	Performs a ready send.
<b>MPI_Rsend_init</b> (void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , int <i>dest</i> , int <i>tag</i> , MPI_Comm <i>comm</i> , MPI_Request * <i>request</i> )	Builds a persistent ready send request handle.

**TABLE A-1 Sun MPI Routines (Continued)**

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Scan</b> (void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)	Computes the scan (partial reductions) of data on a collection of processes.
<b>MPI_Scatter</b> (void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtpe, int root, MPI_Comm comm)	Sends data from one job to all other processes in a group.
<b>MPI_Scatterv</b> (void *sendbuf, int *sendcounts, int *displs, MPI_Datatype sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtpe, int root, MPI_Comm comm)	Scatters a buffer in parts to all processes in a group.
<b>MPI_Send</b> (int *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Performs a standard send.
<b>MPI_Send_init</b> (void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Builds a persistent send request handle.
<b>MPI_Sendrecv</b> (void *sendbuf, int sendcount, MPI_Datatype sendtype, int dest, int sendtag, void *recvbuf, int recvcount, MPI_Datatype recvtpe, int source, int recvtag, MPI_Comm comm, MPI_Status *status)	Sends and receives two messages at the same time.
<b>MPI_Sendrecv_replace</b> (void *buf, int count, MPI_Datatype datatype, int dest, int sendtag, int source, int recvtag, MPI_Comm comm, MPI_Status *status)	Sends and receives using a single buffer.
<b>MPI_Size_of</b> (x, size, ierror)	Fortran only. Returns the size, in bytes, of the datatype of variable x.
<b>MPI_Ssend</b> (void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)	Performs a synchronous send.
<b>MPI_Ssend_init</b> (void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)	Builds a persistent synchronous send request handle.
<b>MPI_Start</b> (MPI_Request *request)	Initiates a communication using a persistent request handle.
<b>MPI_Startall</b> (int count, MPI_Request array_of_requests[ ])	Starts a collection of requests.
<b>MPI_Status_c2f</b> (MPI_Status *c_status, MPI_Fint *f_status)	Translates a C status into a Fortran status.
<b>MPI_Status_f2c</b> (MPI_Fint *f_status, MPI_Status *c_status)	Translates a Fortran status into a C status.

**TABLE A-1** Sun MPI Routines (*Continued*)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Status_set_cancelled</b> (MPI_Status * <i>status</i> , int <i>flag</i> )	Sets <i>status</i> to indicate that a request has been cancelled.
<b>MPI_Status_set_elements</b> (MPI_Status * <i>status</i> , MPI_Datatype <i>datatype</i> , int <i>count</i> )	Modifies opaque part of <i>status</i> to enable MPI_Get_elements() to return <i>count</i> .
<b>MPI_Test</b> (MPI_Request * <i>request</i> , int * <i>flag</i> , MPI_Status * <i>status</i> )	Tests for the completion of a send or receive.
<b>MPI_Test_cancelled</b> (MPI_Status * <i>status</i> , int * <i>flag</i> )	Tests whether a request was canceled.
<b>MPI_Testall</b> (int <i>count</i> , MPI_Request <i>array_of_requests</i> , int * <i>flag</i> , MPI_Status * <i>array_of_statuses</i> )	Tests for the completion of all the given communications.
<b>MPI_Testany</b> (int <i>count</i> , MPI_Request <i>array_of_requests</i> [ ], int * <i>index</i> , int * <i>flag</i> , MPI_Status <i>status</i> )	Tests for completion of any of the given communications.
<b>MPI_Testsome</b> (int <i>incount</i> , MPI_Request <i>array_of_requests</i> [ ], int * <i>outcount</i> , int * <i>array_of_indices</i> , MPI_Status * <i>array_of_statuses</i> )	Tests for some given communications to complete.
<b>MPI_Topo_test</b> (MPI_Comm <i>comm</i> , int * <i>top_type</i> )	Determines the type of topology (if any) associated with a communicator.
<b>MPI_Type_c2f</b> (MPI_Datatype <i>datatype</i> )	Translates a C handle into a Fortran handle.
<b>MPI_Type_commit</b> (MPI_Datatype * <i>datatype</i> )	Commits a data type.
<b>MPI_Type_contiguous</b> (int <i>count</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype * <i>newtype</i> )	Creates a contiguous data type.
<b>MPI_Type_create_darray</b> (int <i>size</i> , int <i>rank</i> , int <i>ndims</i> , int <i>array_of_gsizes</i> [ ], int <i>array_of_distrib</i> [ ], int <i>array_of_dargs</i> [ ], int <i>array_of_psize</i> [ ], int <i>order</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype * <i>newtype</i> )	Creates an array of data types.
<b>MPI_Type_create_f90_complex</b> (int <i>p</i> , int <i>r</i> , MPI_Datatype * <i>newtype</i> )	Returns a bounded MPI complex datatype.
<b>MPI_Type_create_f90_integer</b> (int <i>r</i> , MPI_Datatype * <i>newtype</i> )	Returns a bounded MPI integer datatype.
<b>MPI_Type_create_f90_real</b> (int <i>p</i> , int <i>r</i> , MPI_Datatype * <i>newtype</i> )	Returns a bounded MPI real datatype.

**TABLE A-1 Sun MPI Routines (Continued)**

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Type_create_hindexed</b> (int <i>count</i> , int <i>array_of_blocklengths</i> , MPI_Aint <i>array_of_displacements</i> [], MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	Creates an indexed data type with offsets in bytes.
<b>MPI_Type_create_hvector</b> (int <i>count</i> , int <i>blocklength</i> , MPI_Aint <i>stride</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	Creates a vector (strided) data type with offset in bytes.
<b>MPI_Type_create_indexed_block</b> (int <i>count</i> , int <i>blocklength</i> , int <i>array_of_displacements</i> [], MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	Creates an indexed block.
<b>MPI_Type_create_keyval</b> (MPI_Type_copy_attr_function <i>*type_copy_attr_fn</i> , MPI_Type_delete_attr_function <i>*type_delete_attr_fn</i> , int <i>*type_keyval</i> , void <i>*extra_state</i> )	Generates a new attribute key.
<b>MPI_Type_create_resized</b> (MPI_Datatype <i>oldtype</i> , MPI_Aint <i>lb</i> , MPI_Aint <i>extent</i> , MPI_Datatype <i>*newtype</i> )	Returns a new data type with new extent and upper and lower bounds.
<b>MPI_Type_create_struct</b> (int <i>count</i> , int <i>array_of_blocklengths</i> [], MPI_Aint <i>array_of_displacements</i> [], MPI_Datatype <i>array_of_types</i> [], MPI_Datatype <i>*newtype</i> )	Creates a struct data type.
<b>MPI_Type_create_subarray</b> (int <i>ndims</i> , int <i>array_of_sizes</i> [], int <i>array_of_subsizes</i> [], int <i>array_of_starts</i> [], int <i>order</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	Creates a data type describing a subarray of an array.
<b>MPI_Type_delete_attr</b> (MPI_Datatype <i>type</i> , int <i>type_keyval</i> )	Deletes attribute value associated with a key.
<b>MPI_Type_dup</b> (MPI_Datatype <i>type</i> , MPI_Datatype <i>*newtype</i> )	Duplicates a data type with associated key values.
<b>MPI_Type_extent</b> (MPI_Datatype <i>datatype</i> , MPI_Aint <i>*extent</i> )	<i>Deprecated:</i> Use instead <code>MPI_Type_get_extent()</code> . Returns the extent of a data type, the difference between the upper and lower bounds of the data type.
<b>MPI_Type_f2c</b> (MPI_Fint <i>datatype</i> )	Translates a Fortran handle into a C handle.
<b>MPI_Type_free</b> (MPI_Datatype <i>*datatype</i> )	Frees a data type.
<b>MPI_Type_free_keyval</b> (int <i>*type_keyval</i> )	Frees an attribute key.
<b>MPI_Type_get_attr</b> (MPI_Datatype <i>type</i> , int <i>type_keyval</i> , void <i>*attribute_val</i> , int <i>*flag</i> )	Returns the attribute associated with a data type.

**TABLE A-1** Sun MPI Routines (*Continued*)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Type_get_contents</b> (MPI_Datatype <i>datatype</i> , int <i>max_integers</i> , int <i>max_addresses</i> , int <i>max_datatypes</i> , int <i>array_of_integers</i> [ ], MPI_Aint <i>array_of_addresses</i> [ ], MPI_Datatype <i>array_of_datatypes</i> [ ])	Returns information about arguments used in creation of a data type.
<b>MPI_Type_get_envelope</b> (MPI_Datatype <i>datatype</i> , int <i>*num_integers</i> , int <i>*num_addresses</i> , int <i>*num_datatypes</i> , int <i>*combiner</i> )	Returns information about input arguments associated with a data type.
<b>MPI_Type_get_extent</b> (MPI_Datatype <i>datatype</i> , MPI_Aint <i>*lb</i> , MPI_Aint <i>*extent</i> )	Returns the lower bound and extent of a data type.
<b>MPI_Type_get_name</b> (MPI_Datatype <i>type</i> , char <i>*type_name</i> , int <i>*resultlen</i> )	Gets the name of a data type.
<b>MPI_Type_get_true_extent</b> (MPI_Datatype <i>datatype</i> , MPI_Aint <i>*true_lb</i> , MPI_Aint <i>*true_extent</i> )	Returns the true lower bound and extent of a data type's corresponding type map, ignoring MPI_UB and MPI_LB markers.
<b>MPI_Type_hindexed</b> (int <i>count</i> , int <i>*array_of_blocklengths</i> , MPI_Aint <i>*array_of_displacements</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	<i>Deprecated:</i> Use instead MPI_Type_create_hindexed(). Creates an indexed data type with offsets in bytes.
<b>MPI_Type_hvector</b> (int <i>count</i> , int <i>blocklength</i> , MPI_Aint <i>stride</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	<i>Deprecated:</i> Use instead MPI_Type_create_hvector(). Creates a vector (strided) data type with offset in bytes.
<b>MPI_Type_indexed</b> (int <i>count</i> , int <i>*array_of_blocklengths</i> , int <i>*array_of_displacements</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype <i>*newtype</i> )	Creates an indexed data type.
<b>MPI_Type_lb</b> (MPI_Datatype <i>datatype</i> , MPI_Aint <i>*displacement</i> )	<i>Deprecated:</i> Use instead MPI_Type_get_extent(). Returns the lower bound of a data type.
<b>MPI_Type_match_size</b> (int <i>typeclass</i> , int <i>size</i> , MPI_Datatype <i>*type</i> )	Returns an MPI data type of a given type and size.
<b>MPI_Type_set_attr</b> (MPI_Datatype <i>type</i> , int <i>type_keyval</i> , void <i>*attribute_val</i> )	Stores attribute value associated with a key.
<b>MPI_Type_set_name</b> (MPI_Comm <i>comm</i> , char <i>*type_name</i> )	Sets the name of a data type.
<b>MPI_Type_size</b> (MPI_Datatype <i>datatype</i> , int <i>*size</i> )	Returns the number of bytes occupied by entries in the data type.



TABLE A-1 Sun MPI Routines (Continued)

Routine and C Syntax	Description
<b>MPI_Type_struct</b> (int <i>count</i> , int * <i>array_of_blocklengths</i> , MPI_Aint * <i>array_of_displacements</i> , MPI_Datatype * <i>array_of_types</i> , MPI_Datatype * <i>newtype</i> )	<i>Deprecated:</i> Use instead <code>MPI_Type_create_struct()</code> . Creates a struct data type.
<b>MPI_Type_ub</b> (MPI_Datatype <i>datatype</i> , MPI_Aint * <i>displacement</i> )	<i>Deprecated:</i> Use instead <code>MPI_Type_get_extent()</code> . Returns the upper bound of a data type.
<b>MPI_Type_vector</b> (int <i>count</i> , int <i>blocklength</i> , int <i>stride</i> , MPI_Datatype <i>oldtype</i> , MPI_Datatype * <i>newtype</i> )	Creates a vector (strided) data type.
<b>MPI_Unpack</b> (void * <i>inbuf</i> , int <i>insize</i> , int * <i>position</i> , void * <i>outbuf</i> , int <i>outcount</i> , MPI_Datatype <i>datatype</i> , MPI_Comm <i>comm</i> )	Unpacks a data type into contiguous memory.
<b>MPI_Unpack_external</b> (void * <i>inbuf</i> , int <i>insize</i> , int * <i>position</i> , void * <i>outbuf</i> , int <i>outcount</i> , MPI_Datatype <i>datatype</i> , MPI_Comm <i>comm</i> )	Unpacks into contiguous memory a data type packed in the external32 format.
<b>MPI_Unpublish_name</b> (char * <i>service-name</i> , MPI_Info <i>info</i> , char * <i>port-name</i> )	Removes the pair ( <i>service-name</i> , <i>port-name</i> ) published by <code>MPI_Publish_name()</code> , so that applications can no longer use <code>MPI_Lookup_name()</code> to find <i>port-name</i> .
<b>MPI_Wait</b> (MPI_Request * <i>request</i> , MPI_Status * <i>status</i> )	Waits for an MPI send or receive to complete.
<b>MPI_Waitall</b> (int <i>count</i> , MPI_Request <i>array_of_requests</i> [ ], MPI_Status <i>array_of_statuses</i> [ ])	Waits for all of the given communications to complete.
<b>MPI_Waitany</b> (int <i>count</i> , MPI_Request <i>array_of_requests</i> [ ], int * <i>index</i> , MPI_Status * <i>status</i> )	Waits for any of the given communications to complete.
<b>MPI_Waitsome</b> (int <i>incount</i> , MPI_Request <i>array_of_requests</i> [ ], int * <i>outcount</i> , int <i>array_of_indices</i> [ ], MPI_Status <i>array_of_statuses</i> [ ])	Waits for some given communications to complete.
<b>MPI_Win_c2f</b> (MPI_Win <i>win</i> )	Translates a C handle into a Fortran handle.
<b>MPI_Win_create</b> (void * <i>base</i> , MPI_Aint <i>size</i> , int <i>disp_unit</i> , MPI_Info <i>info</i> , MPI_Comm <i>comm</i> , MPI_Win * <i>win</i> )	Opens a communication window in memory.
<b>MPI_Win_create_errhandler</b> (MPI_Win_errhandler_fn * <i>function</i> , MPI_Errhandler * <i>errhandler</i> )	Creates an error handler that can be attached to windows.

**TABLE A-1** Sun MPI Routines (Continued)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Win_create_keyval</b> (MPI_Win_copy_attr_function *win_copy_attr_fn, MPI_Win_delete_attr_function *win_delete_attr_fn, int *win_keyval, void *extra_state)	Creates a caching attribute that can be associated with a window.
<b>MPI_Win_delete_attr</b> (MPI_Win win, int win_keyval)	Deletes the attribute created with MPI_Win_create_keyval.
<b>MPI_Win_f2c</b> (MPI_Fint win)	Translates a Fortran handle into a C handle.
<b>MPI_Win_fence</b> (int assert, MPI_Win win)	Synchronizes RMA calls on a window.
<b>MPI_Win_free</b> (MPI_Win *win)	Frees the window object and returns a null handle.
<b>MPI_Win_free_keyval</b> (int *win_keyval)	Releases a window attribute.
<b>MPI_Win_get_attr</b> (MPI_Win win, int win_keyval, void *attribute_val, int *flag)	Obtains the value of a window attribute.
<b>MPI_Win_get_errhandler</b> (MPI_Win win, MPI_Errhandler *errhandler)	Retrieves the error handler currently associated with a window.
<b>MPI_Win_get_group</b> (MPI_Win win, MPI_Group *group)	Returns a duplicate of the group of the communicator used to create the window.
<b>MPI_Win_get_name</b> (MPI_Win win, char *win_name, int *resultlen)	Returns the last name associated with a window object.
<b>MPI_Win_lock</b> (int lock_type, int rank, int assert, MPI_Win win)	Starts an RMA access epoch, during which only the window at the process with the specified rank can be accessed.
<b>MPI_Win_post</b> (MPI_Group group, int assert, MPI_Win win)	Starts an RMA exposure epoch for the local window associated with win.
<b>MPI_Win_set_attr</b> (MPI_Win win, int win_keyval, void *attribute_val)	Associates an attribute with a window.
<b>MPI_Win_test</b> (MPI_Win win, int *flag)	Attempts to complete an RMA exposure epoch; a nonblocking version of MPI_Win_wait.
<b>MPI_Win_set_errhandler</b> (MPI_Win win, MPI_Errhandler errhandler)	Attaches a new error handler to a window.
<b>MPI_Win_set_name</b> (MPI_Win win, char *win_name)	Assigns a name to a window.

**TABLE A-1** Sun MPI Routines (*Continued*)

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_Win_start</b> (MPI_Group <i>group</i> , int <i>assert</i> , MPI_Win <i>win</i> )	Starts an RMA access epoch for <i>win</i> .
<b>MPI_Win_unlock</b> (int <i>rank</i> , MPI_Win <i>win</i> )	Completes an RMA access epoch started by a call to <code>MPI_Win_lock()</code> .
double <b>MPI_Wtick</b> ()	Returns the resolution of <code>MPI_Wtime()</code> .
double <b>MPI_Wtime</b> ()	Returns an elapsed time on the calling processor.

---

## Sun MPI I/O Routines

TABLE A-2 lists the Sun MPI I/O routines in alphabetical order. The following sections list the routines by functional category.

# File Manipulation

---

<b>Collective coordination</b>	<b>Noncollective coordination</b>
MPI_File_open()	MPI_File_delete()
MPI_File_close()	MPI_File_get_size()
MPI_File_set_size()	MPI_File_get_group()
MPI_File_preallocate()	MPI_File_get_amode()

---

## File Info

---

<b>Noncollective coordination</b>	<b>Collective coordination</b>
MPI_File_get_info()	MPI_File_set_info()

---

## Data access

### Data Access With Explicit Offsets

---

<b>Synchronism</b>	<b>Noncollective coordination</b>	<b>Collective coordination</b>
<b>Blocking</b>	MPI_File_read_at()	MPI_File_read_at_all()
	MPI_File_write_at()	MPI_File_write_at_all()
<b>Nonblocking or split collective</b>	MPI_File_iread_at()	MPI_File_read_at_all_begin()
		MPI_File_read_at_all_end()
	MPI_File_iwrite_at()	MPI_File_write_at_all_begin()
		MPI_File_write_at_all_end()

---

## Data Access With Individual File Pointers

Synchronism	Noncollective coordination	Collective coordination
Blocking	<code>MPI_File_read()</code>	<code>MPI_File_read_all()</code>
	<code>MPI_File_write()</code>	<code>MPI_File_write_all()</code>
Nonblocking or split collective	<code>MPI_File_iread()</code>	<code>MPI_File_read_all_begin()</code>
		<code>MPI_File_read_all_end()</code>
	<code>MPI_File_iwrite()</code>	<code>MPI_File_write_all_begin()</code>
		<code>MPI_File_write_all_end()</code>

## Data Access With Shared File Pointers

Synchronism	Noncollective coordination	Collective coordination
Blocking	<code>MPI_File_read_shared()</code>	<code>MPI_File_read_ordered()</code>
	<code>MPI_File_write_shared()</code>	<code>MPI_File_write_ordered()</code>
		<code>MPI_File_seek_shared()</code>
		<code>MPI_File_get_position_shared()</code>
Nonblocking or split collective	<code>MPI_File_iread_shared()</code>	<code>MPI_File_read_ordered_begin()</code>
		<code>MPI_File_read_ordered_end()</code>
	<code>MPI_File_iwrite_shared()</code>	<code>MPI_File_write_ordered_begin()</code>
		<code>MPI_File_write_ordered_end()</code>

## Pointer Manipulation

```
MPI_File_seek()  
MPI_File_get_position()  
MPI_File_get_byte_offset()
```

## File Interoperability

```
MPI_Register_datarep()  
MPI_File_get_type_extent()
```

# File Consistency and Semantics

```
MPI_File_set_atomicity()  
MPI_File_get_atomicity()  
MPI_File_sync()
```

## Handle Translation

```
MPI_File_f2c()  
MPI_File_c2f()
```

## MPI I/O Routines: Alphabetical Listing

TABLE A-2 Sun MPI I/O Routines

Routine and C Syntax	Description
<b>MPI_File_c2f</b> (MPI_File <i>file</i> )	Translates a C handle into a Fortran handle.
<b>MPI_File_close</b> (MPI_File * <i>fh</i> )	Closes a file (collective).
<b>MPI_File_create_errhandler</b> ( MPI_File_errhandler_fn * <i>function</i> , MPI_Errhandler * <i>errhandler</i> )	Creates an MPI-style error handler that can be attached to a file.
<b>MPI_File_delete</b> (char * <i>filename</i> , MPI_Info <i>info</i> )	Deletes a file.
<b>MPI_File_f2c</b> (MPI_File <i>file</i> )	Translates a Fortran handle into a C handle.
<b>MPI_File_get_amode</b> (MPI_File <i>fh</i> , int * <i>amode</i> )	Returns mode associated with open file.
<b>MPI_File_get_atomicity</b> (MPI_File <i>fh</i> , int * <i>flag</i> )	Returns current consistency semantics for data-access operations.
<b>MPI_File_get_byte_offset</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , MPI_Offset * <i>disp</i> )	Converts a view-relative offset into an absolute byte position.
<b>MPI_File_get_errhandler</b> (MPI_Comm <i>file</i> , MPI_Errhandler * <i>errhandler</i> )	Gets the error handler for a file.
<b>MPI_File_get_group</b> (MPI_File <i>fh</i> , MPI_Group * <i>group</i> )	Returns the process group of file.

**TABLE A-2 Sun MPI I/O Routines (Continued)**

Routine and C Syntax	Description
<b>MPI_File_get_info</b> (MPI_File <i>fh</i> , MPI_Info * <i>info_used</i> )	Returns a new info object containing hints.
<b>MPI_File_get_position</b> (MPI_File <i>fh</i> , MPI_Offset * <i>offset</i> )	Returns current position of individual file pointer.
<b>MPI_File_get_position_shared</b> (MPI_File <i>fh</i> , MPI_Offset * <i>offset</i> )	Returns current position of the shared file pointer (collective).
<b>MPI_File_get_size</b> (MPI_File <i>fh</i> , MPI_Offset * <i>size</i> )	Returns current size of file.
<b>MPI_File_get_type_extent</b> (MPI_File <i>fh</i> , MPI_Datatype <i>datatype</i> , MPI_Aint * <i>extent</i> )	Returns the extent of the data type in a file.
<b>MPI_File_get_view</b> (MPI_File <i>fh</i> , MPI_Offset * <i>disp</i> , MPI_Datatype * <i>etype</i> , MPI_Datatype * <i>filetype</i> , char * <i>datarep</i> )	Returns process's view of data in file.
<b>MPI_File_iread</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Request * <i>request</i> )	Reads a file starting at the location specified by the individual file pointer (nonblocking, noncollective).
<b>MPI_File_iread_at</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Request * <i>request</i> )	Reads a file at an explicitly specified offset (nonblocking, noncollective).
<b>MPI_File_iread_shared</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Request * <i>request</i> )	Reads a file using the shared file pointer (nonblocking, noncollective).
<b>MPI_File_iwrite</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Request * <i>request</i> )	Writes a file starting at the location specified by the individual file pointer (nonblocking, noncollective).
<b>MPI_File_iwrite_at</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Request * <i>request</i> )	Writes a file at an explicitly specified offset (nonblocking, noncollective).
<b>MPI_File_iwrite_shared</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Request * <i>request</i> )	Writes a file using the shared file pointer (nonblocking, noncollective).
<b>MPI_File_open</b> (MPI_Comm <i>comm</i> , char * <i>filename</i> , int <i>amode</i> , MPI_Info <i>info</i> , MPI_File * <i>fh</i> )	Opens a file (collective).
<b>MPI_File_preallocate</b> (MPI_File <i>fh</i> , MPI_Offset <i>size</i> )	Preallocates storage space for a portion of a file (collective).

**TABLE A-2** Sun MPI I/O Routines (Continued)

Routine and C Syntax	Description
<b>MPI_File_read</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Reads a file starting at the location specified by the individual file pointer.
<b>MPI_File_read_all</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Reads a file starting at the locations specified by individual file pointers (collective).
<b>MPI_File_read_all_begin</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> )	Reads a file starting at the locations specified by individual file pointers; the beginning part of a split collective routine (nonblocking).
<b>MPI_File_read_all_end</b> (MPI_File <i>fh</i> , void * <i>buf</i> , MPI_Status * <i>status</i> )	Reads a file starting at the locations specified by individual file pointers; the ending part of a split collective routine (blocking).
<b>MPI_File_read_at</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Reads a file at an explicitly specified offset.
<b>MPI_File_read_at_all</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Reads a file at explicitly specified offsets (collective).
<b>MPI_File_read_at_all_begin</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> )	Reads a file at explicitly specified offsets; the beginning part of a split collective routine (nonblocking).
<b>MPI_File_read_at_all_end</b> (MPI_File <i>fh</i> , void * <i>buf</i> , MPI_Status * <i>status</i> )	Reads a file at explicitly specified offsets; the ending part of a split collective routine (blocking).
<b>MPI_File_read_ordered</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Reads a file at a location specified by a shared file pointer (collective).
<b>MPI_File_read_ordered_begin</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> )	Reads a file at a location specified by a shared file pointer; the beginning part of a split collective routine (nonblocking).
<b>MPI_File_read_ordered_end</b> (MPI_File <i>fh</i> , void * <i>buf</i> , MPI_Status * <i>status</i> )	Reads a file at a location specified by a shared file pointer; the ending part of a split collective routine (blocking).
<b>MPI_File_read_shared</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Reads a file using the shared file pointer (blocking, noncollective).
<b>MPI_File_seek</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , int <i>whence</i> )	Updates individual file pointers.



**TABLE A-2 Sun MPI I/O Routines (Continued)**

<b>Routine and C Syntax</b>	<b>Description</b>
<b>MPI_File_seek_shared</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , int <i>whence</i> )	Updates the global shared file pointer (collective).
<b>MPI_File_set_atomicity</b> (MPI_File <i>fh</i> , int <i>flag</i> )	Sets consistency semantics for data-access operations (collective).
<b>MPI_File_set_errhandler</b> (MPI_File <i>file</i> , MPI_Errhandler <i>errhandler</i> )	Sets the error handler for a file.
<b>MPI_File_set_info</b> (MPI_File <i>fh</i> , MPI_Info <i>info</i> )	Sets new values for hints (collective).
<b>MPI_File_set_size</b> (MPI_File <i>fh</i> , MPI_Offset <i>size</i> )	Resizes a file (collective).
<b>MPI_File_set_view</b> (MPI_File <i>fh</i> , MPI_Offset <i>disp</i> , MPI_Datatype <i>etype</i> , MPI_Datatype <i>filetype</i> , char <i>*datarep</i> , MPI_Info <i>info</i> )	Changes the process's view of data in file (collective).
<b>MPI_File_sync</b> (MPI_File <i>fh</i> )	Makes semantics consistent for data-access operations (collective).
<b>MPI_File_write</b> (MPI_File <i>fh</i> , void <i>*buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status <i>*status</i> )	Writes a file starting at the location specified by the individual file pointer.
<b>MPI_File_write_all</b> (MPI_File <i>fh</i> , void <i>*buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status <i>*status</i> )	Writes a file starting at the locations specified by individual file pointers (collective).
<b>MPI_File_write_all_begin</b> (MPI_File <i>fh</i> , void <i>*buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> )	Writes a file starting at the locations specified by individual file pointers; the beginning part of a split collective routine (nonblocking).
<b>MPI_File_write_all_end</b> (MPI_File <i>fh</i> , void <i>*buf</i> , MPI_Status <i>*status</i> )	Writes a file starting at the locations specified by individual file pointers; the ending part of a split collective routine (blocking).
<b>MPI_File_write_at</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void <i>*buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status <i>*status</i> )	Writes a file at an explicitly specified offset.
<b>MPI_File_write_at_all</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void <i>*buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status <i>*status</i> )	Writes a file at explicitly specified offsets (collective).
<b>MPI_File_write_at_all_begin</b> (MPI_File <i>fh</i> , MPI_Offset <i>offset</i> , void <i>*buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> )	Writes a file at explicitly specified offsets; the beginning part of a split collective routine (nonblocking).

**TABLE A-2** Sun MPI I/O Routines (Continued)

Routine and C Syntax	Description
<b>MPI_File_write_at_all_end</b> (MPI_File <i>fh</i> , void * <i>buf</i> , MPI_Status * <i>status</i> )	Writes a file at explicitly specified offsets; the ending part of a split collective routine (blocking).
<b>MPI_File_write_ordered</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Writes a file at a location specified by a shared file pointer (collective).
<b>MPI_File_write_ordered_begin</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> )	Writes a file at a location specified by a shared file pointer; the beginning part of a split collective routine (nonblocking).
<b>MPI_File_write_ordered_end</b> (MPI_File <i>fh</i> , void * <i>buf</i> , MPI_Status * <i>status</i> )	Writes a file at a location specified by a shared file pointer; the ending part of a split collective routine (blocking).
<b>MPI_File_write_shared</b> (MPI_File <i>fh</i> , void * <i>buf</i> , int <i>count</i> , MPI_Datatype <i>datatype</i> , MPI_Status * <i>status</i> )	Writes a file using the shared file pointer (blocking, noncollective).
<b>MPI_Register_datarep</b> (char * <i>datarep</i> , MPI_Datarep_conversion_function * <i>read_conversion_fn</i> , MPI_Datarep_conversion_function * <i>write_conversion_fn</i> , MPI_Datarep_extent_function * <i>dtype_file_extent_fn</i> , void * <i>extra_state</i> )	Defines data representation.

## Environment Variables

---

Many environment variables are available for fine-tuning your Sun MPI environment. All Sun MPI environment variables are listed here with brief descriptions. The same descriptions are also available on the `MPI` man page. If you want to return to the default setting after setting a variable, simply unset it (using `unsetenv`). The effects of some of the variables are explained in more detail in the *Sun HPC ClusterTools Software Performance Guide*.

The environment variables are listed here in six groups:

- “Informational Variables” on page 103
- “General Performance Tuning” on page 104
- “Tuning Memory for Point-to-Point Performance” on page 106
- “Numerics” on page 109
- “Synchronization of One-Sided Communications” on page 109
- “MPPProf” on page 112
- “Miscellaneous” on page 113

---

## Informational Variables

### `MPI_PRINTENV`

When set to 1, this variable causes other environment variables and the `hpc.conf` parameters associated with the MPI job to be printed. The default value is 0.

## MPI\_QUIET

If set to 1, this variable suppresses Sun MPI warning messages. The default value is 0.

## MPI\_SHOW\_ERRORS

If set to 1, the `MPI_ERRORS_RETURN` error handler prints the error message and returns the error. The default value is 0.

## MPI\_SHOW\_INTERFACES

When set to 1, 2, or 3, information regarding which interfaces are being used by an MPI application is printed to `stdout`. Set `MPI_SHOW_INTERFACES` to 1 to print the selected internode interface. Set it to 2 to print all the interfaces and their rankings. Set it to 3 for verbose output. The default value, 0, does not print information to `stdout`.

---

# General Performance Tuning

## MPI\_POLLALL

When this variable is set to 1, the default value, all connections are polled for receives, also known as *full polling*. When set to 0, only those connections are polled where receives are posted. Full polling helps drain system buffers, lessening the chance of deadlock for “unsafe” codes. Well-written codes should set `MPI_POLLALL` to 0 for best performance.

---

**Note** – If `MPI_POLLALL` is set to 0 (zero) and your program performs an `MPI_Send/MPI_Cancel` without a corresponding MPI call on the receiving process, the `MPI_Cancel` may not succeed. Your program may hang.

---

## MPI\_PROCBIND

Binds each MPI process to its own processor. The system administrator can allow or disable processor binding by setting the `allow_pbind` parameter in the `CREOptions` section of the `hpc.conf` file. If this parameter is set, the `MPI_PROCBIND` environment variable is disabled. Performance can be enhanced with processor binding, but very poor performance will result if processor binding is used for multithreaded jobs or for more than one job at a time.

By default, `MPI_PROCBIND` is set to 0, which turns off processor binding. To turn on processor binding, set the value to 1. With processor binding turned on, the processes in a job are assigned to all CPUs that are not already bound to other processes.

You can further control CPU binding by using these values for `MPI_PROCBIND`:

```
MPI_PROCBIND {P|L|T} [list | range]
```

- **L** – Only the first thread in each process is bound to a CPU. Any additional threads created by the job are not bound.
- **P** – Every thread created by the process is bound to a CPU.
- **T** – Every thread in every process that is part of the same task is bound to a CPU. (See the `mprun -P` option.)

To specify the CPUs to which the threads are bound, you use either *list* or *range*.

- *list* – Provides explicit binding and can be a single CPU ID or a list of CPU IDs separated by commas. For example, to assign a thread to the fifth CPU in a 24-CPU node, use:

```
% setenv MPI_PROCBIND L4
```

If a node had fewer CPUs than the number you specified, CRE would assign the thread to the CPU that was the modulus of the number you specified, divided by the number of CPUs. For example, a list of L12 in a node with only 8 CPUs would result in the process being assigned to CPU number 3.

When you use a list, the CRE environment does not check to see whether those CPU's are already bound. As a result, you could have two threads bound to the same CPU.

- *range* – Provides automatic binding and assigns threads to all available (unbound) CPUs within the range. The range is expressed as:

*N-MxI*

In other words, specify the range with a starting number (*N*), an ending number (*M*), and a counting interval (*I*). The counting interval (*I*) is optional, and its default value is 1.

For example:

```
% setenv MPI_PROCBIND L0-11
```

The preceding setting would search through the first 12 CPUs, and assign processes to those that are unbound. If the number of available CPUs is less than the number of processes, the extra processes would remain unbound.

The `mprun -v` option (verbose) prints the CPU assignments that have been made. You can also run `pbind(1M)` on each node to verify their CPU bindings.

## MPI\_SPIN

This variable sets the spin policy. The default value is 0, which causes MPI processes to spin nonaggressively, allowing best performance when the load is at least as great as the number of CPUs. A value of 1 causes MPI processes to spin aggressively, leading to best performance if extra CPUs are available on each node to handle system daemons and other background activities.

---

# Tuning Memory for Point-to-Point Performance

## MPI\_RSM\_CPOOLSIZE

This is the requested size, in bytes, to be allocated per stripe for buffers for each RSM connection. This value can be overridden when connections are established on the basis of the size of the segment allocated. The default value is 256 Kbytes.

## MPI\_RSM\_NUMPOSTBOX

This is the number of postboxes per stripe per RSM connection. The maximum number of postboxes depends on the value of `rsm_maxsegsz`. The default is 128 postboxes.

## MPI\_RSM\_PIPE\_SIZE

This limits the size (in bytes) of a message that can be sent over remote shared memory through the buffer list of one postbox per stripe. The default is 64 Kbytes. This size also depends on the block size used for sending data. The maximum size is equal to  $\min (cpoolsize/2, (10 * \max(blk1sz, blk2sz)))$ .

## MPI\_RSM\_SBPOOL\_SIZE

If set, this variable specifies the requested size in bytes of each RSM send buffer pool. An RSM send buffer pool is the pool of buffers on a node that a remote process would use to send to processes on the node. A multiple of 1024 must be used. If unset, the size of buffer pool is equal to *cpoolsize* times processes per node. The *max* value allowed is *maxsegsz* minus the memory used for postboxes.

## MPI\_RSM\_SHORTMSG\_SIZE

This variable specifies the maximum size, in bytes, of a message that will be sent by means of remote shared memory without using buffers. The default value is 3918 bytes. The upper limit is determined by the number of postboxes available.

## MPI\_RSM\_STRONGPARTITION

If set to 1, the RSM protocol module uses strong partition to manage memory. Every connection has a set of blocks in the buffer pool reserved for communication. Otherwise, the pool is shared by all the receivers. The default value is 0.

## MPI\_SHM\_CPOOL\_SIZE

This variable represents the amount of memory, in bytes, that can be allocated to each connection pool. When *MPI\_SHM\_SBPOOL\_SIZE* is not set, the default value is 24,576 bytes. Otherwise, the default value is *MPI\_SHM\_SBPOOL\_SIZE*.

## MPI\_SHM\_CYCLESIZE

This variable represents the limit, in bytes, on the portion of a shared-memory message that will be sent via the buffer list of a single postbox during a cyclic transfer. The default value is 8192 bytes. A multiple of 1024 that is at most  $\text{MPI\_SHM\_CPOOLSIZE}/2$  must be used.

## MPI\_SHM\_CYCLESTART

Shared-memory transfers that are larger than  $\text{MPI\_SHM\_CYCLESTART}$  bytes are cyclic. The default value is 24,576 bytes.

## MPI\_SHM\_NUMPOSTBOX

This variable represents the number of postboxes dedicated to each shared-memory connection. The default value is 16.

## MPI\_SHM\_PIPESIZE

This variable represents the limit, in bytes, on the portion of a shared-memory message that will be sent via the buffer list of a single postbox during a pipeline transfer. The default value is 8192 bytes. The value must be a multiple of 1024.

## MPI\_SHM\_PIPESTART

This variable represents the size, in bytes, at which shared-memory transfers starts to be pipelined. The default value is 2048. Multiples of 1024 must be used.

## MPI\_SHM\_SBPOOLSIZE

If set, this variable represents the size, in bytes, of the pool of shared-memory buffers dedicated to each sender. A multiple of 1024 must be used. If unset, then pools of shared-memory buffers are dedicated to connections rather than to senders.



## MPI\_SHM\_SHORTMSGSIZE

This variable represents the size (in bytes) of the section of a postbox that contains either data or a buffer list. The default value is 256 bytes.

---

**Note** – If `MPI_SHM_PIPESTART`, `MPI_SHM_PIPE_SIZE`, or `MPI_SHM_CYCLE_SIZE` is increased to a size larger than 31,744 bytes, then `MPI_SHM_SHORTMSGSIZE` might also have to be increased. See the *Sun HPC ClusterTools Software Performance Guide* for more information.

---

---

## Numerics

### MPI\_CANONREDUCE

Prevents reduction operations from using any optimizations that take advantage of the physical location of processors. This can provide more consistent results in the case of floating-point addition, for example. However, the operation can take longer to complete. The default value is 0, meaning that optimizations are allowed. To prevent optimizations, set the value to 1.

---

## Synchronization of One-Sided Communications

### MPI\_USE\_AGENT\_THREAD

If the environmental variable `MPI_USE_AGENT_THREAD` is set to 1, upon the first call to `MPI_Win_create` the Sun MPI library creates one agent thread for processes that need such a thread. (If `MPI_USE_AGENT_THREAD` is not set, or it is set to 0 [zero], no such thread is created.)

The two purposes of `MPI_USE_AGENT_THREAD` are to ensure progress in passive target RMA synchronization and to perform MPI RMA operations on local window memory on behalf of other processes when those processes do not have direct (shared-memory or RSM) access to window memory.

The agent thread does not run user code. Thread-safety in the non-thread-safe MPI library is achieved by a monitor around MPI communication calls. If no windows requiring the use of an agent thread are active, the agent thread is suspended. If `MPI_USE_AGENT_THREAD` is not set, one-sided MPI operations can be delayed till the next synchronization call.

## `MPI_RSM_PUTSIZE`

Use the environmental variable `MPI_RSM_PUTSIZE` to control the method used to put small buffers over the RSM PM. Buffers passed to `MPI_Put` that are smaller than `MPI_RSM_PUTSIZE` are transferred to the remote side using a message protocol that has low latency but requires an implicit synchronization with the remote side, which may delay the sending side and reduce performance. The delay is usually small if both sides perform frequent small transfers. Buffers larger than `MPI_RSM_PUTSIZE` are transferred using a direct one-sided protocol which does not require synchronization, but has higher latency. The cost of the higher latency is not significant for large enough messages. Choose a value of 0 (zero) to always use the direct protocol, and a large value to always use the message protocol. The default value is 16 Kbytes.

## `MPI_RSM_GETSIZE`

Use the environmental variable `MPI_RSM_GETSIZE` to control the method used to get small buffers over the RSM PM. Buffers passed to `MPI_Get` that are larger than `MPI_RSM_GETSIZE` are transferred using a message protocol that achieves high bandwidth but requires an implicit synchronization with the remote side, which may delay the sending side and reduce performance. For large enough messages, the higher bandwidth offsets the cost of the synchronization delay. Buffers passed to `MPI_Get` that are smaller than `MPI_RSM_GETSIZE` are transferred using a direct one-sided protocol which does not require synchronization, but has lower bandwidth. For small enough messages, the lower bandwidth has little impact on performance. Choose a value of 0 (zero) to always use the message protocol, and a large value to always use the direct protocol. The default value is 16Kbytes.

---

# Tuning Rendezvous

## `MPI_EAGERONLY`

When this variable is set to 1, the default, only the eager protocol is used. When it is set to 0, both eager and rendezvous protocols are used.

## `MPI_RSM_RENDVSIZE`

Messages communicated by remote shared memory that are greater than this size use the rendezvous protocol unless the environment variable `MPI_EAGERONLY` is set to 1. The default value is 256 Kbytes.

## `MPI_SHM_RENDVSIZE`

Messages communicated by shared memory that are greater than this size use the rendezvous protocol unless the environment variable `MPI_EAGERONLY` is set. The default value is 24,576 bytes.

## `MPI_TCP_RENDVSIZE`

Messages communicated by TCP that contain data of this size and greater use the rendezvous protocol unless the environment variable `MPI_EAGERONLY` is set. The default value is 49,152 bytes.

---

# MPPProf

## MPI\_PROFILE

Setting this variable to 1 enables a profiling session. When profiling is enabled, profiling data for the MPI process ranks are written to a set of intermediate files, one file per process rank. MPPProf also creates an index file of the form: `mpprof.index.rm.jid` (where *rm* is the resource manager and *jid* is the job ID) that contains pointers to the intermediate files of the form `mpprof.n.rm.jid` (where *n* is the process rank, *rm* is the resource manager, and *jid* is the job ID). If `MPI_PROFILE` is not set, program execution proceeds without generating profiling data.

## MPI\_PROFDATADIR

By default, the temporary files generated during MPPProf profiling are located in `/usr/tmp/`. Set an alternative location as a value for environment variable `MPI_PROFDATADIR`.

## MPI\_PROFINDEXDIR

By default, the index file for MPPProf profiling is located in the current directory. Set an alternative, nondefault location as a value for `MPI_PROFINDEXDIR`.

## MPI\_PROFINTERVAL

The variable `MPI_PROFINTERVAL` can be used to specify a time interval for controlling when snapshots of the profiling data will be written to the intermediate files.

Setting `MPI_PROFINTERVAL` to 0 forces a snapshot for every MPI call that is made. Setting `MPI_PROFINTERVAL` to `Inf` causes only one snapshot to be recorded at `MPI_Finalize` time. If `MPI_PROFINTERVAL` is unset or has no value, the default value of 60 seconds will be used.

If time intervals are used and an MPI program terminates before the `MPI_Finalize` call, any snapshots that were recorded can be used by `mpprof` to generate a profile of program operations up to the point of termination.

## MPI\_PROFMAXFILESIZE

This variable can be used to specify the maximum size, in Kbytes, that can be written to the intermediate files.

The default intermediate file size limit files is 51,200 Kbytes (50 Mbytes). If a process records data that exceeds the file size limit, that write operation completes, but it cannot record additional profiling data. Other intermediate files that have not reached the limit can continue to receive data. The file size limit can be removed altogether by setting `MPI_PROFMAXFILESIZE` to unlimited.

---

## Miscellaneous

### MPI\_COSCHED

This variable specifies the user's preference regarding use of the `spind` daemon for coscheduling. The value can be 0 (prefer no use) or 1 (prefer use). This preference can be overridden by the system administrator's policy. This policy is set in the `hpc.conf` file and can be 0 (forbid use), 1 (require use), or 2 (no policy). If no policy is set and no user preference is specified, coscheduling is not used.

---

**Note** – If no user preference is specified, the value 2 is displayed when environment variables are printed with `MPI_PRINTENV`.

---

### MPI\_CHECK\_ARGS

When this variable is set to 1, argument checking is performed on MPI calls, and errors are printed when they occur. The default is 0.

### MPI\_FLOWCONTROL

This variable limits the number of unexpected messages that can be queued from a particular connection. Once this quantity of unexpected messages has been received, polling the connection for incoming messages stops. The default value, 0, indicates that no limit is set. To limit flow, set the value to an integer greater than 0.

## MPI\_FULLCONNINIT

This variable ensures that all connections are established during initialization. By default, connections are established lazily. However, you can override this default by setting the environment variable `MPI_FULLCONNINIT` to 1, forcing full-connection initialization mode. The default value is 0.

## MPI\_MAXFHANDLES

This variable represents the maximum number of Fortran handles for objects other than requests. `MPI_MAXFHANDLES` specifies the upper limit on the number of concurrently allocated Fortran handles for MPI objects other than requests. This variable is ignored in the default 32-bit library. The default value is 1024. Users should take care to free MPI objects that are no longer in use. There is no limit on handle allocation for C codes.

## MPI\_MAXPROCS

This variable overrides the value specified by `maxprocs_default` in `hpc.conf`; it cannot exceed the value specified by `maxprocs_limit` in `hpc.conf`. If the value does exceed the `maxprocs_limit` value, the job aborts with an error when the program calls `MPI_Init`. Note that the upper limit of support for RSM communication is 2048 processes.

## MPI\_MAXREQHANDLES

This variable represents the maximum number of Fortran request handles. `MPI_MAXREQHANDLES` specifies the upper limit on the number of concurrently allocated MPI request handles. Users must take care to free up request handles by properly completing requests. The default value is 1024. This variable is ignored in the default 32-bit library.

## MPI\_OPTCOLL

The MPI collectives are implemented using a variety of optimizations. Some of these optimizations can inhibit performance of point-to-point messages for “unsafe” programs. The default value of this variable, 1, indicates that optimized collectives are used. The optimizations can be turned off by setting the value to 0.

## MPI\_RSM\_MAXSTRIPE

This variable specifies the maximum number of interfaces that can be used for striping data during communication over RSM. The value cannot be higher than 64 and the number of installed interfaces. The default is 4.

## MPI\_SHM\_BCASTSIZE

On SMPs, `MPI_Bcast()` is implemented for large messages using a double-buffering scheme. The size of each buffer (in bytes) is settable by using this environment variable. The default value is 32,768 bytes.

## MPI\_SHM\_GBPOOLSIZE

This variable represents the amount of memory available, in bytes, to the general buffer pool for use by collective operations. The default value is 20,971,520 bytes.

## MPI\_SHM\_REDUCE SIZE

On SMPs, calling `MPI_Reduce()` causes all processors to participate in the reduce. Each processor works on a piece of data equal to the `MPI_SHM_REDUCE SIZE` setting. The default value is 256 bytes.

You must take care when setting this variable because the system reserves  $MPI\_SHM\_REDUCE SIZE * np * np$  memory to execute the reduce.

## MPI\_SPINDTIMEOUT

When coscheduling is enabled, this variable limits the length of time (in milliseconds) a message remains in the poll waiting for the `spind` daemon to return. If the timeout occurs before the daemon finds any messages, the process reenters the polling loop. The default value is 1000 milliseconds. A default can also be set by a system administrator in the `hpc.conf` file.

## MPI\_TCP\_CONNLOOP

This variable sets the number of times `MPI_TCP_CONNTIMEOUT` occurs before signaling an error. The default value for this variable is 0, meaning that the program aborts on the first occurrence of `MPI_TCP_CONNTIMEOUT`.

## MPI\_TCP\_CONNTIMEOUT

This variable sets the timeout value in seconds that is used for an `accept()` call. The default value for this variable is 600 seconds (10 minutes). This timeout can be triggered in both full- and lazy-connection initialization. After the timeout is reached, a warning message is printed. If `MPI_TCP_CONNLOOP` is set to 0, then the first timeout causes the program to abort.

## MPI\_TCP\_SAFEGATHER

This variable allows use of a congestion-avoidance algorithm for `MPI_Gather()` and `MPI_Gatherv()` over TCP. By default, `MPI_TCP_SAFEGATHER` is set to 1, which means that use of this algorithm is on. If you know that your underlying network can handle gathering large amounts of data on a single node, you might want to override this algorithm by setting `MPI_TCP_SAFEGATHER` to 0.



# Troubleshooting

---

This appendix describes some common problem situations, resulting error messages, and suggestions for fixing the problems. It includes the following topics:

- “MPI Messages” on page 117
- “MPI I/O Error Handling” on page 120
- “Adjusting System V Shared-Memory Limits” on page 122

Sun MPI error reporting, including I/O, follows the *MPI-2 Standard*. By default, errors are reported in the form of standard error classes. These classes and their meanings are listed in TABLE C-1 (for non-I/O MPI) and TABLE C-2 (for MPI I/O) and are also available on the MPI man page.

Three predefined error handlers are available in Sun MPI:

- `MPI_ERRORS_RETURN` – The default; returns an error code if an error occurs.
- `MPI_ERRORS_ARE_FATAL` – I/O errors are fatal, and no error code is returned.
- `MPI_THROW_EXCEPTION` – A special error handler to be used only with C++.

---

## MPI Messages

You can make changes to and get information about the error handler by using any of the following routines:

- `MPI_Comm_call_errhandler`
- `MPI_File_call_errhandler`
- `MPI_Win_call_errhandler`
- `MPI_Comm_create_errhandler`
- `MPI_Comm_get_errhandler`
- `MPI_Comm_set_errhandler`
- `MPI_Add_error_class`
- `MPI_Add_error_code`
- `MPI_Add-error-string`

Messages resulting from an MPI program fall into two categories:

- *Error messages* – Error messages stem from within MPI. Usually an error message explains why your program cannot complete, and the program aborts.
- *Warning messages* – Warnings stem from the environment in which you are running your MPI program and are usually sent by `MPI_Init()`. They are not associated with an aborted program; that is, programs continue to run despite warning messages.

## Error Messages

Sun MPI error messages use a standard format:

```
[xyz] Error in function_name: errclass_string:intern(a): description: unixerrstring
```

Where

- `[xyz]` is the *process communication identifier*, which is present in every error message, and:
  - `x` is the job ID (or jid).
  - `y` is the name of the communicator if a name exists; otherwise it is the address of the opaque object.
  - `z` is the rank of the process.
- *function\_name* is the name of the associated MPI function. It is present in every error message.
- *errclass\_string* is the string associated with the MPI error class. It is present in every error message.
- *intern* is an internal function. It is optional.
- `a` is a system call if one is the cause of the error. It is optional.
- *description* is a description of the error. It is optional.
- *unixerrstring* is the UNIX error string that describes system call `a`. It is optional.

## Warning Messages

Sun MPI warning messages also use a standard format:

```
[xyz] Warning message
```

Where *message* is a description of the error.

# Standard Error Classes

TABLE C-1 lists the error return classes you can encounter in your MPI programs. Error values may also be found in `mpi.h` (for C), `mpif.h` (for Fortran), and `mpi++.h` (for C++).

MPI I/O messages are listed separately, in TABLE C-2.

**TABLE C-1** Sun MPI Standard Error Classes

Error Code	Value	Meaning
<code>MPI_SUCCESS</code>	0	Successful return code.
<code>MPI_ERR_BUFFER</code>	1	Invalid buffer pointer.
<code>MPI_ERR_COUNT</code>	2	Invalid count argument.
<code>MPI_ERR_TYPE</code>	3	Invalid datatype argument.
<code>MPI_ERR_TAG</code>	4	Invalid tag argument.
<code>MPI_ERR_COMM</code>	5	Invalid communicator.
<code>MPI_ERR_RANK</code>	6	Invalid rank.
<code>MPI_ERR_ROOT</code>	7	Invalid root.
<code>MPI_ERR_GROUP</code>	8	Null group passed to function.
<code>MPI_ERR_OP</code>	9	Invalid operation.
<code>MPI_ERR_TOPOLOGY</code>	10	Invalid topology.
<code>MPI_ERR_DIMS</code>	11	Illegal dimension argument.
<code>MPI_ERR_ARG</code>	12	Invalid argument.
<code>MPI_ERR_UNKNOWN</code>	13	Unknown error.
<code>MPI_ERR_TRUNCATE</code>	14	Message truncated on receive.
<code>MPI_ERR_OTHER</code>	15	Other error; use <code>Error_string</code> .
<code>MPI_ERR_INTERN</code>	16	Internal error code.
<code>MPI_ERR_IN_STATUS</code>	17	Look in status for error value.
<code>MPI_ERR_PENDING</code>	18	Pending request.
<code>MPI_ERR_REQUEST</code>	19	Illegal <code>MPI_Request()</code> handle.
<code>MPI_ERR_KEYVAL</code>	36	Illegal key value.

**TABLE C-1** Sun MPI Standard Error Classes (Continued)

Error Code	Value	Meaning
MPI_ERR_INFO	37	Invalid info object.
MPI_ERR_INFO_KEY	38	Illegal info key.
MPI_ERR_INFO_NOKEY	39	No such key.
MPI_ERR_INFO_VALUE	40	Illegal info value.
MPI_ERR_TIMEOUT	41	Timed out.
MPI_ERR_RESOURCES	42	Out of resources.
MPI_ERR_TRANSPORT	43	Transport layer error.
MPI_ERR_HANDSHAKE	44	Error accepting/connecting.
MPI_ERR_SPAWN	45	Error spawning.
MPI_ERR_WIN	46	Invalid window.
MPI_ERR_BASE	47	Invalid base.
MPI_ERR_SIZE	48	Invalid size.
MPI_ERR_DISP	49	Invalid displacement.
MPI_ERR_LOCKTYPE	50	Invalid lock type.
MPI_ERR_ASSERT	51	Invalid assert.
MPI_ERR_RMA_CONFLICT	52	Conflicting accesses to window.
MPI_ERR_RMA_SYNC	53	Erroneous RMA synchronization.
MPI_ERR_NO_MEM	54	Memory exhausted.
MPI_ERR_LASTCODE	55	Last error code.

---

## MPI I/O Error Handling

Sun MPI I/O error reporting follows the *MPI-2 Standard*. By default, errors are reported in the form of standard error codes (found in `/opt/SUNWhpc/include/mpi.h`). Error classes and their meanings are listed in TABLE C-2. You can also find them in `mpif.h` (for Fortran) and `mpi++.h` (for C++).

You can change the default error handler by specifying `MPI_FILE_NULL` as the file handle with the routine `MPI_File_set_errhandler()`, even if no file is currently open. Or, you can use the same routine to change a specific file's error handler.

**TABLE C-2** Sun MPI I/O Error Classes

<b>Error Class</b>	<b>Value</b>	<b>Meaning</b>
<code>MPI_ERR_FILE</code>	20	Bad file handle.
<code>MPI_ERR_NOT_SAME</code>	21	Collective argument not identical on all processes.
<code>MPI_ERR_AMODE</code>	22	Unsupported <code>amode</code> passed to <code>open</code> .
<code>MPI_ERR_UNSUPPORTED_DATAREP</code>	23	Unsupported <code>datarep</code> passed to <code>MPI_File_set_view()</code> .
<code>MPI_ERR_UNSUPPORTED_OPERATION</code>	24	Unsupported operation, such as seeking on a file that supports only sequential access.
<code>MPI_ERR_NO_SUCH_FILE</code>	25	File (or directory) does not exist.
<code>MPI_ERR_FILE_EXISTS</code>	26	File exists.
<code>MPI_ERR_BAD_FILE</code>	27	Invalid file name (for example, path name too long).
<code>MPI_ERR_ACCESS</code>	28	Permission denied.
<code>MPI_ERR_NO_SPACE</code>	29	Not enough space.
<code>MPI_ERR_QUOTA</code>	30	Quota exceeded.
<code>MPI_ERR_READ_ONLY</code>	31	Read-only file system.
<code>MPI_ERR_FILE_IN_USE</code>	32	File operation could not be completed, because the file is currently open by a process.
<code>MPI_ERR_DUP_DATAREP</code>	33	Conversion functions could not be registered because a data representation identifier that was already defined was passed to <code>MPI_REGISTER_DATAREP</code> .
<code>MPI_ERR_CONVERSION</code>	34	An error occurred in a user-supplied data-conversion function.
<code>MPI_ERR_IO</code>	35	I/O error.
<code>MPI_ERR_INFO</code>	37	Invalid info object.

**TABLE C-2** Sun MPI I/O Error Classes (Continued)

Error Class	Value	Meaning
MPI_ERR_INFO_KEY	38	Illegal info key.
MPI_ERR_INFO_NOKEY	39	No such key.
MPI_ERR_INFO_VALUE	40	Illegal info value.
MPI_ERR_LASTCODE	55	Last error code.

## Adjusting System V Shared-Memory Limits

When using the Sun MPI library on a cluster equipped with RSM links, memory allocated from System V shared memory is used for process communication.

Given that the maximum number of such shared-memory segments is set to a conservative value by default, calls to allocate new memory segments can fail due to these system-imposed limits. When such calls fail, check the values specified for the following variables in `/etc/system`:

- `shmsys:shminfo_shmseg`
- `shmsys:shminfo_shmmni`
- `shmsys:shminfo_shmmax`

If the values are too low, increase them by editing `/etc/system` and rebooting the system.

# Index

---

## A

Argonne National Laboratory  
and MPE, 47  
array sections, 17  
attributes, with communicators, 16

## B

blocking routines. *See* routines, blocking.

## C

caching, with communicators, 16  
Cartesian topology. *See* topology, Cartesian.  
code samples. *See* sample programs.  
collective communication. *See* communication,  
collective.  
communication  
buffers, 17  
collective, 12, 16, 21  
in multithreaded programs, 25  
restrictions, 13  
half-channel, 14  
interprocess, 16  
persistent request, defined, 14  
point-to-point, 16  
port, 14  
communicator  
default, 16  
defined, 16

and MPI I/O, 49  
and multithreaded programming, 25  
and process topologies, 14  
compiling, 34 to 37  
with profiling library, 27  
*See also* include syntax.  
context, defined, 16

## D

-dalign option, 36  
data type  
possible values for C, 19  
derived (user-defined), 17, 49  
possible values for Fortran, 17, 18  
primitive, 17  
debugging, 39 to 47  
with mpe, 47  
*See also* Prism.  
displacement (*disp*), 50, 53

## E

elementary data type (*etype*), 50  
environmental inquiry functions, 15  
error handling, 15  
and MPE, 30  
and multithreaded programming, 25  
error messages  
standard error classes (Sun MPI I/O), 120  
standard error values (Sun MPI), 119

**F**

- features, 1
- file type (`filetype`), 50
- Fortran
  - compiling with `-dalign` option, 36
  - compiling with `-xalias` option, 35

**G**

- graph topology. *See* topology, graph.
- grid topology. *See* topology, Cartesian.
- group, defined, 16

**H**

- header files, 31
- “holes” (in an MPI I/O file type), 50, 54

**I**

- I/O. *See* Sun MPI I/O, MPI I/O.
- include syntax, 31
- intercommunicator, defined, 16
- intracommunicator, 14
  - defined, 16

**L**

- libraries
  - `libthread.so`, 37
  - linking, 34 to 37
  - linking with profiling library, 27
- linking, 34 to 37

**M**

- man pages
  - Solaris, location, 37
  - Sun MPI, location, xi
- modes for point-to-point communication, 4
- MPI
  - Forum, URL, x
  - Mississippi State University URL, x
  - Standards
    - profiling, 27
    - URL, x
  - University of New Mexico URL, x
- MPI I/O, 49
  - Sun MPI implementation.
    - See* Sun MPI I/O.
- `MPI_COMM_GROUP`, 16
- `MPI_COMM_WORLD`
  - as default communicator, 16
- multiprocessing environment (MPE), 29 to 30
  - and debugging, 47
- multithreaded programming, 23 to 26
  - stubbing thread calls, 37

**N**

- nonblocking routines. *See* routines, nonblocking.

**O**

- offset, 50
- one-sided communication, 5
- options
  - `-dalign`, 36
  - `-xalias`, 35

**P**

- persistent communication request. *See* communication, persistent request.
- point-to-point
  - communication. *See* communication, point-to-point.
  - routines. *See* routines, point-to-point.
- Prism, 39 to 41
  - compilers to use, 36
- process
  - relation to group, 16
- process topologies, 14
- profiling, 27 to 28



## R

rank, of a process, 14, 16  
ready mode. *See* modes for point-to-point communication.  
receive. *See* routines, receive.  
routines  
  all-gather, 12  
  all-to-all, 12  
  blocking, 4, 12  
  broadcast, 12  
  collective, 12, 16  
    in multithreaded programs, 25  
  for constructing communicators, 16  
  data access (MPI I/O), 53 to 57  
    pointer manipulation, 55 to 56  
    with explicit offsets, 54  
    with individual file pointers, 55  
    with shared file pointers, 56 to 57  
  error-handling, 15  
  file consistency (MPI I/O), 58  
  file manipulation (MPI I/O), 52  
  gather, 12  
  for constructing groups, 16  
  local, 16  
  nonblocking, 4  
  point-to-point, 4  
  receive, 4, 21  
  reduction, 12  
  scan, 12  
  scatter, 12  
  semantics (MPI I/O), 58  
  send, 4, 21  
  Sun MPI  
    listed alphabetically, 78 to 95  
    listed by functional category, 67 to 78  
  Sun MPI I/O  
    listed alphabetically, 98 to 102  
    listed by functional category, 95 to 98

## S

sample programs  
  Sun MPI, 32 to 34  
  Sun MPI I/O, 59 to 66  
send. *See* routines, send.  
shutting down, 15  
SPMD programs

  defined, 39  
  standard mode. *See* modes for point-to-point communication.  
  starting up, 15  
  static libraries, and relinking, 37  
  Sun MPI I/O, 49 to 66  
  synchronous mode. *See* modes for point-to-point communication.

## T

thread safety. *See* multithreaded programming.  
timers, 15  
topology  
  Cartesian, 14  
  graph, 14  
  virtual, defined, 14  
  *See also* process topologies.

## V

view, 50

## X

-xaliasoption, 35

