

**The International Journal of Supercomputer Applications and High Performance Computing**

**Editors-in-Chief**  
Joanne L. Martin  
Jack Dongarra

**Editorial Assistant**  
Jan Jones

**The International Journal of Supercomputer Applications and High Performance Computing** is published quarterly (spring, summer, fall, and winter) by The MIT Press, Cambridge, MA 02142. Subscriptions and address changes should be addressed to MIT Press Journals, 55 Hayward Street, Cambridge, MA 02142; (617) 253-2889; E-mail: journals-orders@mit.edu. Subscriptions are on a volume year basis. Rates are: Individuals \$75.00, Institutions \$210.00, Students/Retired \$45.00. Outside U.S. add \$16.00 for postage and handling. Canadians add additional 7% GST. Current issues are \$53.00. Back issues are: Individuals \$26.50, Institutions \$53.00. Outside U.S. add \$5.00 per issue for postage and handling. Canadians add additional 7% GST. To be honored free, claims for missing issues must be made immediately upon receipt of the next published issue. Prices subject to change without notice. *The International Journal of Supercomputer Applications and High Performance Computing* is abstracted or indexed in the following: *Applied Mechanics Reviews*, *Artificial Intelligence Abstracts*, *Cambridge Scientific Abstracts*, *CompuMath Citation Index*, *Current Contents/Engineering, Technology & Computer Science*, *Applied Sciences*, *Energy Data Base*, *Information Science Abstracts*, *Research Alert*, and *Science Citation Index/SciSearch*.

**Editorial Addresses**

for individual manuscripts,  
Jack Dongarra

*The International Journal of Supercomputer Applications and High Performance Computing* Dept. of Computer Science, 104 Ayres Hall University of Tennessee, Knoxville, TN 37996-1301

for special issue proposals,

Joanne L. Martin  
*The International Journal of Supercomputer Applications and High Performance Computing* IBM Corporation, Mailstop 601; Bldg 201 Neighborhood Rd. Kingston, NY 12401

**Business Offices**

Subscriptions, address changes, and mailing list correspondence should be addressed to MIT Press Journals, 55 Hayward Street, Cambridge, MA 02142.

**Copyright Information**

Permission to photocopy articles for internal or personal use, or the internal or personal use of specific clients, is granted by the copyright owner for users registered with the Copyright Clearance Center (CCC) Transactional Reporting Service, provided that the fee of \$5.00 per article-copy is paid directly to CCC, 222 Rosewood Drive, Danvers, MA 01923. The fee code for users of the Transactional Reporting Service is 0890-2720/94 \$5.00. For those organizations that have been granted a photocopy license with CCC, a separate system of payment has been arranged. Address all other

inquiries to the Subsidiary Rights Manager, MIT Press Journals, 55 Hayward St., Cambridge, MA 02142

Copyright © 1994 by The Massachusetts Institute of Technology

**Postmaster**

Send address changes to *The International Journal of Supercomputer Applications and High Performance Computing*, 55 Hayward Street, Cambridge, MA 02142. Second Class postage paid at Boston, MA, and at additional post offices.

**Marketing and Mailing List Rental**

Address inquiries to the Advertising Manager, MIT Press Journals, 55 Hayward Street, Cambridge, MA 02142. (617) 253-2866, e-mail: journals-info@mit.edu.

**Editorial Board**

Petter E. Bjonstad  
Institutt for Informatikk  
Bergen, Norway

James L. Black  
IBM Corporation  
Dallas, TX

Ernest R. Davidson  
Indiana University  
Bloomington, IN

David Dixon  
E. I. du Pont de Nemours & Co.  
Wilmington, DE

Iain Duff  
Rutherford Appleton Laboratory  
Didcot, Oxfordshire  
United Kingdom

Charbel Farhat  
The University of Colorado  
Boulder, CO

William Gropp  
Argonne National Laboratory  
Argonne, IL

A. J. G. Hey  
University of Southampton  
England

Michael Heath  
Beckman Institute  
University of Illinois  
Urbana, IL

Rolf Hempel  
GMD  
Germany

William Jalby  
University of Rennes,  
Rennes, France

Anthony Jameson  
Princeton University  
Princeton, NJ

Dennis Jespersen  
NASA/Ames Research Center  
Moffett Field, CA

Olin Johnson  
University of Houston  
Houston, TX

Lennart Johnsson  
Thinking Machines Corporation  
and Harvard University  
Cambridge, MA

David K. Kahaner  
US Office of Naval Research Asia  
Tokyo, Japan

David E. Keyes  
Old Dominion University and  
ICASE/NASA Langley Research Center  
Hampton, VA

Thomas A. Kitchens  
Office of Energy Research  
Washington, DC

Monica Lam  
Stanford University  
Stanford, CA

Jill P. Mesirov  
Thinking Machines Corporation  
Cambridge, MA

Gerard Meurant  
Centres d'Etudes de Limeil-Valenton  
Villeneuve-St. Georges, France

Anna Nagurny  
University of Massachusetts  
Amherst, MA

Kenneth W. Neves  
Boeing Computer Services  
Seattle, WA

Steve W. Otto  
Oregon Graduate Institute  
of Science & Technology  
Portland, OR

Yoshio Oyanagi  
University of Tokyo  
Tokyo, Japan

John P. Riganati  
David Sarnoff Research Center  
Princeton, NJ

Yves Robert  
Ecole Normale Supérieure de Lyon  
Lyon, France

Garry H. Rodrigue  
Lawrence Livermore  
National Laboratory  
Livermore, CA

John Rust  
University of Wisconsin-Madison  
Madison, WI

Robert S. Schreiber  
RIACS/NASA Ames Research Center  
Moffett Field, CA

James A. Sethian  
University of California, Berkeley  
Berkeley, CA

Horst Simon  
Computer Sciences Corporation/  
NASA Ames Research Center  
Moffett Field, CA

Anthony Skjellum  
Mississippi State University  
Mississippi State, MS

Malcolm Stocks  
Oak Ridge National Laboratory  
Oak Ridge, TN

Robert G. Voigt  
National Science Foundation  
Washington, DC

THE INTERNATIONAL JOURNAL OF  
**SUPERCOMPUTER APPLICATIONS  
AND HIGH PERFORMANCE COMPUTING**

**SPECIAL ISSUE**

**MPI: A MESSAGE-PASSING INTERFACE STANDARD**

The standard, called the Message Passing Interface (MPI), provides a common interface for distributed memory concurrent computers and networks of workstations. MPI functionality includes point-to-point and collective communication routines, as well as support for process groups, communication contexts, and application topologies. While making use of new ideas, the MPI standard is based largely on current practice, such as Express, PVM, NX/2 Vertex, and P4.

The main advantages of establishing a message passing interface are portability and ease of use; a standard for message passing is a key component in building a concurrent computing environment in which applications, software libraries, and tools can be transparently ported between different machines. Furthermore, the definition of a standard provides vendors with a clearly defined set of routines that they can implement efficiently, or in some cases provide hardware or low-level system support for, thereby enhancing scalability.

The MPI standardization effort involved about 60 people from 40 organizations, mainly from the United States and Europe. Most of the major vendors of concurrent computers have been involved in MPI, along with researchers from universities, government laboratories and industry. MPI is intended to be a standard message passing interface for applications running on MIMD distributed memory concurrent computers and workstation networks.

**VOLUME 8 NUMBER 3/4**

**FALL/WINTER 1994**

I S S N 0 8 9 0 - 2 7 2 0

# CONTENTS

<b>Acknowledgments</b>	<b>165</b>
<b>Preface</b>	<b>167</b>
<b>1 Introduction to MPI</b>	<b>169</b>
1.1 Overview and Goals	169
1.2 Who Should Use This Standard?	171
1.3 What Platforms Are Targets For Implementation?	171
1.4 What Is Included in the Standard?	172
1.5 What Is Not Included in the Standard?	172
1.6 Organization of This Document	172
<b>2 MPI Terms and Conventions</b>	<b>175</b>
2.1 Document Notation	175
2.2 Procedure Specification	175
2.3 Semantic Terms	176
2.4 Data Types	177
2.4.1 Opaque Objects	177
2.4.2 Array Arguments	178
2.4.3 State	179
2.4.4 Named Constants	179
2.4.5 Choice	179
2.4.6 Addresses	179
2.5 Language Binding	179
2.5.1 Fortran 77 Binding Issues	180
2.5.2 C Binding Issues	181
2.6 Processes	181
2.7 Error Handling	182
2.8 Implementation Issues	183
2.8.1 Independence of Basic Runtime Routines	183
2.8.2 Interaction with Signals in POSIX	184

<b>3</b>	<b>Point-to-Point Communication</b>	<b>185</b>
3.1	Introduction	185
3.2	Blocking Send and Receive Operations	186
3.2.1	Blocking Send	186
3.2.2	Message Data	187
3.2.3	Message Envelope	188
3.2.4	Blocking Receive	189
3.2.5	Return Status	191
3.3	Data Type Matching and Data Conversion	192
3.3.1	Type Matching Rules	192
3.3.2	Data Conversion	195
3.4	Communication Modes	197
3.5	Semantics of Point-to-Point Communication	201
3.6	Buffer Allocation and Usage	205
3.6.1	Model Implementation of Buffered Mode	206
3.7	Nonblocking Communication	206
3.7.1	Communication Objects	208
3.7.2	Communication Initiation	208
3.7.3	Communication Completion	211
3.7.4	Semantics of Nonblocking Communications	214
3.7.5	Multiple Completions	215
3.8	Probe and Cancel	221
3.9	Persistent Communication Requests	225
3.10	Send-receive	229
3.11	Null Processes	231
3.12	Derived Datatypes	231
3.12.1	Datatype Constructors	233
3.12.2	Address and Extent Functions	241
3.12.3	Lower-bound and Upper-bound Markers	243
3.12.4	Commit and Free	245
3.12.5	Use of General Datatypes in Communication	246
3.12.6	Correct Use of Addresses	249
3.12.7	Examples	250
3.13	Pack and Unpack	258
<b>4</b>	<b>Collective Communication</b>	<b>267</b>
4.1	Introduction and Overview	267
4.2	Communicator Argument	270
4.3	Barrier Synchronization	270
4.4	Broadcast	270
4.4.1	Example using MPI.BCAST	271
4.5	Gather	271
4.5.1	Examples using MPI.GATHER, MPI.GATHERV	274
4.6	Scatter	280



4.6.1	Examples using MPI_SCATTER, MPI_SCATTERV	283
4.7	Gather-to-All	285
4.7.1	Examples using MPI_ALLGATHER, MPI_ALLGATHERV	287
4.8	All-to-All Scatter/Gather	287
4.9	Global Reduction Operations	290
4.9.1	Reduce	290
4.9.2	Predefined Reduce Operations	291
4.9.3	MINLOC and MAXLOC	293
4.9.4	User-Defined Operations	297
4.9.5	All-Reduce	301
4.10	Reduce-Scatter	302
4.11	Scan	303
4.11.1	Example using MPI_SCAN	304
4.12	Correctness	305
<b>5</b>	<b>Groups, Contexts, and Communicators</b>	<b>311</b>
5.1	Introduction	311
5.1.1	Features Needed to Support Libraries	311
5.1.2	MPI's Support for Libraries	312
5.2	Basic Concepts	314
5.2.1	Groups	314
5.2.2	Contexts	315
5.2.3	Intra-Communicators	315
5.2.4	Predefined Intra-Communicators	316
5.3	Group Management	316
5.3.1	Group Accessors	316
5.3.2	Group Constructors	318
5.3.3	Group Destructors	322
5.4	Communicator Management	322
5.4.1	Communicator Accessors	322
5.4.2	Communicator Constructors	324
5.4.3	Communicator Destructors	327
5.5	Motivating Examples	328
5.5.1	Current Practice #1	328
5.5.2	Current Practice #2	329
5.5.3	(Approximate) Current Practice #3	329
5.5.4	Example #4	330
5.5.5	Library Example #1	331
5.5.6	Library Example #2	333
5.6	Inter-Communication	335
5.6.1	Inter-Communicator Accessors	337
5.6.2	Inter-Communicator Operations	338
5.6.3	Inter-Communication Examples	341
5.7	Caching	348

5.7.1	Functionality	349
5.7.2	Attributes Example	353
5.8	Formalizing the Loosely Synchronous Model	355
5.8.1	Basic Statements	355
5.8.2	Models of Execution	355
<b>6</b>	<b>Process Topologies</b>	<b>357</b>
6.1	Introduction	357
6.2	Virtual Topologies	358
6.3	Embedding in MPI	359
6.4	Overview of the Functions	359
6.5	Topology Constructors	360
6.5.1	Cartesian Constructor	360
6.5.2	Cartesian Convenience Function: <code>MPI_DIMS_CREATE</code>	361
6.5.3	General (Graph) Constructor	361
6.5.4	Topology Inquiry Functions	363
6.5.5	Cartesian Shift Coordinates	367
6.5.6	Partitioning of Cartesian structures	368
6.5.7	Low-level Topology Functions	369
6.6	An Application Example	372
<b>7</b>	<b>MPI Environmental Management</b>	<b>373</b>
7.1	Implementation Information	373
7.1.1	Environmental Inquiries	373
7.2	Error Handling	375
7.3	Error Codes and Classes	378
7.4	Timers	379
7.5	Startup	380
<b>8</b>	<b>Profiling Interface</b>	<b>383</b>
8.1	Requirements	383
8.2	Discussion	383
8.3	Logic of the Design	384
8.3.1	Miscellaneous Control of Profiling	384
8.4	Examples	385
8.4.1	Profiler Implementation	385
8.4.2	MPI Library Implementation	386
8.4.3	Complications	387
8.5	Multiple Levels of Interception	388
	<b>Bibliography</b>	<b>389</b>
<b>A</b>	<b>Language Binding</b>	<b>393</b>
A.1	Introduction	393
A.2	Defined Constants for C and Fortran	393

A.3 C Bindings for Point-to-Point Communication	397
A.4 C Bindings for Collective Communication	400
A.5 C Bindings for Groups, Contexts, and Communicators	401
A.6 C Bindings for Process Topologies	403
A.7 C Bindings for Environmental Inquiry	403
A.8 C Bindings for Profiling	404
A.9 Fortran Bindings for Point-to-Point Communication	404
A.10 Fortran Bindings for Collective Communication	408
A.11 Fortran Bindings for Groups, Contexts, etc.	410
A.12 Fortran Bindings for Process Topologies	412
A.13 Fortran Bindings for Environmental Inquiry	413
A.14 Fortran Bindings for Profiling	414

The technical development was carried out by subgroups, whose work was reviewed by the full committee. During the period of development of the Message Passing Interface (MPI), many people served in positions of responsibility and are listed below.

- Jack Dongarra, David Walker, Conveners and Meeting Chairs
- Ewing Lusk, Bob Knighten, Minutes
- Marc Snir, William Gropp, Ewing Lusk, Point-to-Point Communications
- Al Geist, Marc Snir, Steve Otto, Collective Communications
- Steve Otto, Editor
- Rolf Hempel, Process Topologies
- Ewing Lusk, Language Binding
- William Gropp, Environmental Management
- James Cowrie, Profiling
- Tony Skjellum, Lyndon Clarke, Marc Snir, Richard Littlefield, Mark Sears, Groups, Contexts, and Communicators
- Steven Huss-Lederman, Initial Implementation Subset

The following list includes some of the active participants in the MPI process not mentioned above.

Ed Anderson	Robert Balch	Joe Eron	Eric Barnack
Scott Berryman	Rob Bjornson	Nathan Dow	Anne Elver
Jim Fennedy	Vince Fernando	Sam Fineberg	Jon Flower
Daniel Frye	Ian Glendinning	Adam Greenberg	Robert Harrison
Ledie Hart	Tom Hump	Don Heller	Tom Henderson
Alex Ho	C.T. Howard Ho	Gary Howell	John Karpaga
James Kohl	Susan Roman	Bob Leary	Arthur MacCabe
Peter Madhus	Alan Mainwaring	Olivier McBryan	Phil McKinley
Charles Mosher	Dan Newett	Peter Padroco	Howard Palmer
Paul Pierce	Sanjay Ranka	Peter Rignbee	Arch Robison
Erich Scholz	Amby Singh	Alan Sussman	Robert Tomlinson
Robert G. Voigt	Dennis Weeks	Stephen Wurst	Steven Zanolli

## ACKNOWLEDGMENTS

The technical development was carried out by subgroups, whose work was reviewed by the full committee. During the period of development of the Message Passing Interface (MPI), many people served in positions of responsibility and are listed below.

- Jack Dongarra, David Walker, Conveners and Meeting Chairs
- Ewing Lusk, Bob Knighten, Minutes
- Marc Snir, William Gropp, Ewing Lusk, Point-to-Point Communications
- Al Geist, Marc Snir, Steve Otto, Collective Communications
- Steve Otto, Editor
- Rolf Hempel, Process Topologies
- Ewing Lusk, Language Binding
- William Gropp, Environmental Management
- James Cownie, Profiling
- Tony Skjellum, Lyndon Clarke, Marc Snir, Richard Littlefield, Mark Sears, Groups, Contexts, and Communicators
- Steven Huss-Lederman, Initial Implementation Subset

The following list includes some of the active participants in the MPI process not mentioned above.

Ed Anderson	Robert Babb	Joe Baron	Eric Barszcz
Scott Berryman	Rob Bjornson	Nathan Doss	Anne Elster
Jim Feeney	Vince Fernando	Sam Fineberg	Jon Flower
Daniel Frye	Ian Glendinning	Adam Greenberg	Robert Harrison
Leslie Hart	Tom Haupt	Don Heller	Tom Henderson
Alex Ho	C.T. Howard Ho	Gary Howell	John Kapenga
James Kohl	Susan Krauss	Bob Leary	Arthur Maccabe
Peter Madams	Alan Mainwaring	Oliver McBryan	Phil McKinley
Charles Mosher	Dan Nessett	Peter Pacheco	Howard Palmer
Paul Pierce	Sanjay Ranka	Peter Rigsbee	Arch Robison
Erich Schikuta	Ambuj Singh	Alan Sussman	Robert Tomlinson
Robert G. Voigt	Dennis Weeks	Stephen Wheat	Steven Zenith



The University of Tennessee and Oak Ridge National Laboratory made the draft available by anonymous FTP mail servers and were instrumental in distributing the document.

This work was supported in part by ARPA and NSF under grant ASC-9310330, the National Science Foundation Science and Technology Center Cooperative Agreement No. CCR-8809615, and by the Commission of the European Community through Esprit project P6643 (PPPE).

The technical development was carried out by subgroups whose work was reviewed by the full committee. During the period of development of the Message Passing Interface (MPI), many people served in positions of responsibility and are listed below.

- Jack Dongarra, David Walker, Convenors and Meeting Chairs
- Ewing Lusk, Bob Knighten, Minutes
- Marc Snit, William Gropp, Ewing Lusk, Point-to-Point Communications
- Al Geist, Marc Snit, Steve Otto, Collective Communications
- Steve Otto, Editor
- Rolf Hempel, Process Topologies
- Ewing Lusk, Language Binding
- William Gropp, Environmental Management
- James Cowme, Prototyping
- Tony Skjellum, Lyndon Clark, Marc Snit, Richard Linglefield, Mark Sear, Groups, Conferences, and Communications
- Steven Huss, Initial Implementation Subgroup

The following list includes some of the active participants in the MPI process not mentioned above.

Ed Anderson	Robert Bobb	Joe Brown	Eric Bunker
Scott Brannen	Rob Brannen	Nathan Dow	Anne Elser
Jim Evans	Franc Fernandez	Sam Fluckey	Jon Fowler
David Fry	Jan Glesne	Adam Greenberg	Robert Harrison
Lash Hsu	Tom Hays	Don Heller	Tom Henderson
Alfred Ho	C.J. Howard Ho	Ray Howell	John Hwang
James Koll	Sean Kline	Bob Levy	Arthur MacKay
Peter Mathis	Alan McInerney	Oliver McNamara	Phil McKinney
Charles Mohr	Don Norwest	Peter Pacheco	Howard Palmer
Paul Price	George Rucka	Peter Rucka	Arch Robinson
Erich Scholz	Angus Smith	Alan Swann	Robert Tammann
Robert C. Vogt	Dennis Webb	Stephen Wilson	Steven Zorn

## PREFACE

**Jack Dongarra**

The Message Passing Interface effort began in the summer of 1991 when a small group of researchers started discussions at a mountain retreat in Austria. Out of that discussion came a Workshop on Standards for Message Passing in a Distributed Memory Environment held on April 29–30, 1992, in Williamsburg, Virginia. At this workshop the basic features essential to a standard message-passing interface were discussed, and a working group established to continue the standardization process. More formal meetings and discussions began in January 1993 and continued with meetings every six weeks with discussions via e-mail. The MPI Standard was completed in March of 1994. The MPI effort involved about 60 people from 40 organizations, mainly from the United States and Europe. Most of the major vendors of concurrent computers were involved in MPI, along with researchers from universities, government laboratories, and industry.

This effort defines the syntax and semantics of a core of library routines useful to a wide range of users writing portable message-passing programs in Fortran or C. The MPI effort was conducted in a similar spirit to the High-Performance Fortran Forum (HPFF).

MPI provides parallel hardware vendors with a clearly defined base set of routines that can be efficiently implemented. As a result, hardware vendors can build upon this collection of standard low-level routines to create higher-level routines for the distributed-memory communication environments supplied with their parallel machines. MPI provides a simple-to-use portable interface for the basic user, yet is powerful enough to allow programmers to use the high-performance message-passing operations available on advanced machines.

As an effort to create a "true" standard for message passing, researchers incorporated into MPI the most useful features of several systems, rather than choosing one system to adopt as a standard. Features were used from systems by IBM, Intel, nCUBE, PVM, Express, P4, and PARMACS. The message paradigm, we believe, will be attractive because of its wide portability and can be used in

communications for distributed-memory and shared-memory multiprocessors, networks of workstations, and any combination of these elements. The paradigm will not be made obsolete by increases in network speeds or by architectures combining shared and distributed-memory components. As this standard is printed, we have implementations of MPI on various platforms with more expected in the months to come.

MPI operated on a very tight budget (in reality, it had no budget when the first meeting was announced). ARPA and NSF have supported research at various institutions and have made a contribution toward travel for the U.S. academics. Support for several European participants was provided by ESPRIT.

This issue of the Journal is also available in Postscript and HTML forms over the Internet. To retrieve the postscript file you can anonymous ftp to netlib2.cs.utk.edu; cd mpi; get mpi-report.ps. The HTML form can be found by the URL: <http://www.mcs.anl.gov/mpi/mpi-report/mpi-report.html>. An up-to-date list of errata for this document is being maintained. To receive a copy, send an e-mail message to [netlib@ornl.gov](mailto:netlib@ornl.gov) with contents: send mpi.errata.ps from mpi.

This effort defines the syntax and semantics of a core of library routines useful to a wide range of users writing portable message-passing programs in Fortran or C. The MPI effort was conducted in a similar spirit to the High-Performance Fortran Forum (HPFF).

MPI provides parallel hardware vendors with a clearly defined base set of routines that can be efficiently implemented. As a result, hardware vendors can build upon this collection of standard low-level routines to create higher-level routines for the distributed-memory communication environments supported with their parallel machines. MPI provides a simple-to-use portable interface for the basic user, yet is powerful enough to allow programmers to use the high-performance message-passing operations available on advanced machines.

As an effort to create a "true" standard for message passing, researchers incorporated into MPI the most useful features of several systems, rather than choosing one system to adopt as a standard. Features were used from systems by IBM, Intel, NCUBE, EVM, Express, P4, and PARMAK. The message paradigm, we believe, will be attractive because of its wide portability and can be used in

## INTRODUCTION TO MPI

### 1.1 Overview and Goals

Message passing is a paradigm used widely on certain classes of parallel machines, especially those with distributed memory. Although there are many variations, the basic concept of processes communicating through messages is well understood. Over the last ten years, substantial progress has been made in casting significant applications in this paradigm. Each vendor has implemented its own variant. More recently, several systems have demonstrated that a message-passing system can be efficiently and portably implemented. It is thus an appropriate time to try to define both the syntax and semantics of a core of library routines that will be useful to a wide range of users and efficiently implementable on a wide range of computers.

In designing MPI we have sought to make use of the most attractive features of a number of existing message-passing systems, rather than selecting one of them and adopting it as the standard. Thus, MPI has been strongly influenced by work at the IBM T.J. Watson Research Center [1, 2], Intel's NX/2 [23], Express [22], nCUBE's Vertex [21], p4 [7, 6], and PARMACS [5, 8]. Other important contributions have come from Zipcode [24, 25], Chimp [14, 15], PVM [4, 11], Chameleon [19], and PICL [18].

The MPI standardization effort involved about 60 people from 40 organizations mainly from the United States and Europe. Most of the major vendors of concurrent computers were involved in MPI, along with researchers from universities, government laboratories, and industry. The standardization process began with the Workshop on Standards for Message Passing in a Distributed Memory Environment, sponsored by the Center for Research on Parallel Computing, held April 29–30, 1992, in Williamsburg, Virginia [29]. At this workshop the basic features essential to a standard message-passing interface were discussed, and a working group established to continue the standardization process.



A preliminary draft proposal, known as MPI1, was put forward by Dongarra, Hempel, Hey, and Walker in November 1992, and a revised version was completed in February 1993 [12]. MPI1 embodied the main features that were identified at the Williamsburg workshop as being necessary in a message-passing standard. Since MPI1 was primarily intended to promote discussion and "get the ball rolling," it focused mainly on point-to-point communications. MPI1 brought to the forefront a number of important standardization issues, but did not include any collective communication routines and was not thread-safe.

In November 1992, a meeting of the MPI working group was held in Minneapolis, at which it was decided to place the standardization process on a more formal footing, and to generally adopt the procedures and organization of the High Performance Fortran Forum. Subcommittees were formed for the major component areas of the standard, and an e-mail discussion service established for each. In addition, the goal of producing a draft MPI standard by the Fall of 1993 was set. To achieve this goal the MPI working group met every 6 weeks for two days throughout the first 9 months of 1993, and presented the draft MPI standard at the Supercomputing 93 conference in November 1993. These meetings and the e-mail discussion together constituted the MPI Forum, membership of which has been open to all members of the high performance computing community.

The main advantages of establishing a message-passing standard are portability and ease-of-use. In a distributed memory communication environment in which the higher level routines and/or abstractions are built upon lower level message-passing routines the benefits of standardization are particularly apparent. Furthermore, the definition of a message-passing standard, such as that proposed here, provides vendors with a clearly defined base set of routines that they can implement efficiently, or in some cases provide hardware support for, thereby enhancing scalability.

The goal of the Message Passing Interface simply stated is to develop a widely used standard for writing message-passing programs. As such the interface should establish a practical, portable, efficient, and flexible standard for message passing.

A complete list of goals follows.

- Design an application programming interface (not necessarily for compilers or a system implementation library).
- Allow efficient communication: Avoid memory-to-memory copying and allow overlap of computation and communication and offload to communication co-processor, where available.
- Allow for implementations that can be used in a heterogeneous environment.
- Allow convenient C and Fortran 77 bindings for the interface.
- Assume a reliable communication interface: the user need not cope with communication failures. Such failures are dealt with by the underlying communication subsystem.

- Define an interface that is not too different from current practice, such as PVM, NX, Express, p4, etc., and provides extensions that allow greater flexibility.
- Define an interface that can be implemented on many vendors' platforms, with no significant changes in the underlying communication and system software.
- Semantics of the interface should be language independent.
- The interface should be designed to allow for thread-safety.

## 1.2 Who Should Use This Standard?

This standard is intended for use by all those who want to write portable message-passing programs in Fortran 77 and C. This includes individual application programmers, developers of software designed to run on parallel machines, and creators of environments and tools. In order to be attractive to this wide audience, the standard must provide a simple, easy-to-use interface for the basic user while not semantically precluding the high-performance message-passing operations available on advanced machines.

## 1.3 What Platforms Are Targets For Implementation?

The attractiveness of the message-passing paradigm at least partially stems from its wide portability. Programs expressed this way may run on distributed-memory multiprocessors, networks of workstations, and combinations of all of these. In addition, shared-memory implementations are possible. The paradigm will not be made obsolete by architectures combining the shared- and distributed-memory views, or by increases in network speeds. It thus should be both possible and useful to implement this standard on a great variety of machines, including those "machines" consisting of collections of other machines, parallel or not, connected by a communication network.

The interface is suitable for use by fully general MIMD programs, as well as those written in the more restricted style of SPMD. Although no explicit support for threads is provided, the interface has been designed so as not to prejudice their use. With this version of MPI no support is provided for dynamic spawning of tasks.

MPI provides many features intended to improve performance on scalable parallel computers with specialized interprocessor communication hardware. Thus, we expect that native, high-performance implementations of MPI will be provided on such machines. At the same time, implementations of MPI on top of standard Unix interprocessor communication protocols will provide portability to workstation clusters and heterogeneous networks of workstations. Several proprietary, native implementations of MPI, and a public domain, portable implementation of MPI are in progress at the time of this writing [17, 13].

## 1.4 What Is Included in the Standard?

The standard includes:

- Point-to-point communication
- Collective operations
- Process groups
- Communication contexts
- Process topologies
- Bindings for Fortran 77 and C
- Environmental Management and inquiry
- Profiling interface

## 1.5 What Is Not Included in the Standard?

The standard does not specify:

- Explicit shared-memory operations
- Operations that require more operating system support than is currently standard; for example, interrupt-driven receives, remote execution, or active messages
- Program construction tools
- Debugging facilities
- Explicit support for threads
- Support for task management
- I/O functions

There are many features that have been considered and not included in this standard. This happened for a number of reasons, one of which is the time constraint that was self-imposed in finishing the standard. Features that are not included can always be offered as extensions by specific implementations. Perhaps future versions of MPI will address some of these issues.

## 1.6 Organization of This Document

The following is a list of the remaining chapters in this document, along with a brief description of each.

- Chapter 2, MPI Terms and Conventions, explains notational terms and conventions used throughout the MPI document.
- Chapter 3, Point-to-Point Communication, defines the basic, pairwise communication subset of MPI. *send* and *receive* are found here, along with many associated functions designed to make basic communication powerful and efficient.
- Chapter 4, Collective Communication, defines process-group collective communication operations. Well-known examples of this are barrier and broadcast over a group of processes (not necessarily all the processes).

- Chapter 5, Groups, Contexts, and Communicators, shows how groups of processes are formed and manipulated, how unique communication contexts are obtained, and how the two are bound together into a *communicator*.
- Chapter 6, Process Topologies, explains a set of utility functions meant to assist in the mapping of process groups (a linearly ordered set) to richer topological structures such as multi-dimensional grids.
- Chapter 7, MPI Environmental Management, explains how the programmer can manage and make inquiries of the current MPI environment. These functions are needed for the writing of correct, robust programs, and are especially important for the construction of highly portable message-passing programs.
- Chapter 8, Profiling Interface, explains a simple name-shifting convention that any MPI implementation must support. One motivation for this is the ability to put performance profiling calls into MPI without the need for access to the MPI source code. The name shift is merely an interface; it says nothing about how the actual profiling should be done and in fact, the name shift can be useful for other purposes.
- Annex A, Language Binding, gives specific syntax in Fortran 77 and C, for all MPI functions, constants, and types.
- The MPI Function Index is a simple index showing the location of the precise definition of each MPI function, together with both C and Fortran bindings.

*You should read them carefully. (End of rationale.)*

*Advice to users.* Throughout this document, material that speaks to users and illustrates usage is set off in this format. Some readers may wish to skip these sections, while readers interested in programming in MPI may want to read them carefully. *(End of advice to users.)*

*Advice to implementors.* Throughout this document, material that is primarily commentary to implementors is set off in this format. Some readers may wish to skip these sections, while readers interested in MPI implementations may want to read them carefully. *(End of advice to implementors.)*

## 2.2 Procedure Specification

MPI procedures are specified using a language-independent notation. The arguments of procedure calls are marked as IN, OUT or INOUT. The meanings of these are:

- the call uses but does not update an argument marked IN,
- the call may update an argument marked OUT,
- the call both uses and updates an argument marked INOUT.

There is one special case: if an argument is a handle to an opaque object (these terms are defined in Section 2.1.1), and the object is updated by the procedure call, then the argument is marked OUT. It is marked this way even



## MPI TERMS AND CONVENTIONS

This chapter explains notational terms and conventions used throughout the MPI document, some of the choices that have been made, and the rationale behind those choices.

### 2.1 Document Notation

*Rationale.* Throughout this document, the rationale for the design choices made in the interface specification is set off in this format. Some readers may wish to skip these sections, while readers interested in interface design may want to read them carefully. (*End of rationale.*)

*Advice to users.* Throughout this document, material that speaks to users and illustrates usage is set off in this format. Some readers may wish to skip these sections, while readers interested in programming in MPI may want to read them carefully. (*End of advice to users.*)

*Advice to implementors.* Throughout this document, material that is primarily commentary to implementors is set off in this format. Some readers may wish to skip these sections, while readers interested in MPI implementations may want to read them carefully. (*End of advice to implementors.*)

### 2.2 Procedure Specification

MPI procedures are specified using a language-independent notation. The arguments of procedure calls are marked as IN, OUT or INOUT. The meanings of these are:

- the call uses but does not update an argument marked IN,
- the call may update an argument marked OUT,
- the call both uses and updates an argument marked INOUT.

There is one special case—if an argument is a handle to an opaque object (these terms are defined in Section 2.4.1), and the object is updated by the procedure call, then the argument is marked OUT. It is marked this way even

though the handle itself is not modified—we use the OUT attribute to denote that what the handle *references* is updated.

The definition of MPI tries to avoid, to the largest possible extent, the use of INOUT arguments, because such use is error-prone, especially for scalar arguments.

A common occurrence for MPI functions is an argument that is used as IN by some processes and OUT by other processes. Such argument is, syntactically, an INOUT argument and is marked as such, although, semantically, it is not used in one call both for input and for output.

Another frequent situation arises when an argument value is needed only by a subset of the processes. When an argument is not significant at a process then an arbitrary value can be passed as argument.

Unless specified otherwise, an argument of type OUT or type INOUT cannot be aliased with any other argument passed to an MPI procedure. An example of argument aliasing in C appears below. If we define a C procedure like this,

```
void copyIntBuffer( int *pin, int *pout, int len )
{
    int i;
    for (i=0; i<len; ++i) *pout++ = *pin++;
}
```

then a call to it in the following code fragment has aliased arguments.

```
int a[10];
copyIntBuffer( a, a+3, 7);
```

Although the C language allows this, such usage of MPI procedures is forbidden unless otherwise specified. Note that Fortran prohibits aliasing of arguments.

All MPI functions are first specified in the language-independent notation. Immediately below this, the ANSI C version of the function is shown, and below this, a version of the same function in Fortran 77.

## 2.3 Semantic Terms

When discussing MPI procedures the following semantic terms are used. The first two are usually applied to communication operations.

**nonblocking** If the procedure may return before the operation completes, and before the user is allowed to re-use resources (such as buffers) specified in the call.

**blocking** If return from the procedure indicates the user is allowed to re-use resources specified in the call.

**local** If completion of the procedure depends only on the local executing process. Such an operation does not require communication with another user process.

**non-local** If completion of the operation may require the execution of some

MPI procedure on another process. Such an operation may require communication occurring with another user process.

**collective** If all processes in a process group need to invoke the procedure.

## 2.4 Data Types

### 2.4.1 OPAQUE OBJECTS

MPI manages **system memory** that is used for buffering messages and for storing internal representations of various MPI objects such as groups, communicators, datatypes, etc. This memory is not directly accessible to the user, and objects stored there are **opaque**: their size and shape is not visible to the user. Opaque objects are accessed via **handles**, which exist in user space. MPI procedures that operate on opaque objects are passed handle arguments to access these objects. In addition to their use by MPI calls for object access, handles can participate in assignment and comparisons.

In Fortran, all handles have type `INTEGER`. In C, a different handle type is defined for each category of objects. These should be types that support assignment and equality operators.

In Fortran, the handle can be an index to a table of opaque objects in system table; in C it can be such index or a pointer to the object. More bizarre possibilities exist.

Opaque objects are allocated and deallocated by calls that are specific to each object type. These are listed in the sections where the objects are described. The calls accept a handle argument of matching type. In an allocate call this is an OUT argument that returns a valid reference to the object. In a call to deallocate this is an INOUT argument which returns with an "invalid handle" value. MPI provides an "invalid handle" constant for each object type. Comparisons to this constant are used to test for validity of the handle.

A call to deallocate invalidates the handle and marks the object for deallocation. The object is not accessible to the user after the call. However, MPI need not deallocate the object immediately. Any operation pending (at the time of the deallocate) that involves this object will complete normally; the object will be deallocated afterwards.

An opaque object and its handle are significant only at the process where the object was created, and cannot be transferred to another process.

MPI provides certain predefined opaque objects and predefined, static handles to these objects. Such objects may not be destroyed.

*Rationale.* This design hides the internal representation used for MPI data structures, thus allowing similar calls in C and Fortran. It also avoids conflicts with the typing rules in these languages, and easily allows future extensions of functionality. The mechanism for opaque objects used here loosely follows the POSIX Fortran binding standard.

The explicit separating of handles in user space, objects in system space, allows space-reclaiming, deallocation calls to be made at appropriate points

in the user program. If the opaque objects were in user space, one would have to be very careful not to go out of scope before any pending operation requiring that object completed. The specified design allows an object to be marked for deallocation, the user program can then go out of scope, and the object itself still persists until any pending operations are complete.

The requirement that handles support assignment/comparison is made since such operations are common. This restricts the domain of possible implementations. The alternative would have been to allow handles to have been an arbitrary, opaque type. This would force the introduction of routines to do assignment and comparison, adding complexity, and was therefore ruled out. (*End of rationale.*)

*Advice to users.* A user may accidentally create a dangling reference by assigning to a handle the value of another handle, and then deallocating the object associated with these handles. Conversely, if a handle variable is deallocated before the associated object is freed, then the object becomes inaccessible (this may occur, for example, if the handle is a local variable within a subroutine, and the subroutine is exited before the associated object is deallocated). It is the user's responsibility to avoid adding or deleting references to opaque objects, except as a result of calls that allocate or deallocate such objects. (*End of advice to users.*)

*Advice to implementors.* The intended semantics of opaque objects is that each opaque object is separate from each other; each call to allocate such an object copies all the information required for the object. Implementations may avoid excessive copying by substituting referencing for copying. For example, a derived datatype may contain references to its components, rather than copies of its components; a call to `MPI_COMM_GROUP` may return a reference to the group associated with the communicator, rather than a copy of this group. In such cases, the implementation must maintain reference counts, and allocate and deallocate objects such that the visible effect is as if the objects were copied. (*End of advice to implementors.*)

## 2.4.2 ARRAY ARGUMENTS

An MPI call may need an argument that is an array of opaque objects, or an array of handles. The array-of-handles is a regular array with entries that are handles to objects of the same type in consecutive locations in the array. Whenever such an array is used, an additional `len` argument is required to indicate the number of valid entries (unless this number can be derived otherwise). The valid entries are at the beginning of the array; `len` indicates how many of them there are, and need not be the entire size of the array. The same approach is followed for other array arguments.



### 2.4.3 STATE

MPI procedures use at various places arguments with *state* types. The values of such data type are all identified by names, and no operation is defined on them. For example, the `MPI_ERRHANDLER_SET` routine has a state type argument with values `MPI_ERRORS_RETURN`, `MPI_ERRORS_FATAL`, etc.

### 2.4.4 NAMED CONSTANTS

MPI procedures sometimes assign a special meaning to a special value of a basic type argument; e.g. `tag` is an integer-valued argument of point-to-point communication operations, with a special wild-card value, `MPI_ANY_TAG`. Such arguments will have a range of regular values, which is a proper subrange of the range of values of the corresponding basic type; special values (such as `MPI_ANY_TAG`) will be outside the regular range. The range of regular values can be queried using environmental inquiry functions (Chapter 7).

### 2.4.5 CHOICE

MPI functions sometimes use arguments with a *choice* (or union) data type. Distinct calls to the same routine may pass by reference actual arguments of different types. The mechanism for providing such arguments will differ from language to language. For Fortran, the document uses `<type>` to represent a choice variable, for C, we use `(void *)`.

### 2.4.6 ADDRESSES

Some MPI procedures use *address* arguments that represent an absolute address in the calling program. The datatype of such an argument is an integer of the size needed to hold any valid address in the execution environment.

## 2.5 Language Binding

This section defines the rules for MPI language binding in general and for Fortran 77 and ANSI C in particular. Defined here are various object representations, as well as the naming conventions used for expressing this standard. The actual calling sequences are defined elsewhere.

It is expected that any Fortran 90 and C++ implementations use the Fortran 77 and ANSI C bindings, respectively. Although we consider it premature to define other bindings to Fortran 90 and C++, the current bindings are designed to encourage, rather than discourage, experimentation with better bindings that might be adopted later.

Since the word `PARAMETER` is a keyword in the Fortran language, we use the word "argument" to denote the arguments to a subroutine. These are normally referred to as parameters in C, however, we expect that C programmers will understand the word "argument" (which has no specific meaning in C), thus allowing us to avoid unnecessary confusion for Fortran programmers.

```

double precision a
integer b
...
call MPI_send(a,...)
call MPI_send(b,...)

```

**Fig. 2.1** An example of calling a routine with mismatched formal and actual arguments.

There are several important language binding issues not addressed by this standard. This standard does not discuss the interoperability of message passing between languages. It is fully expected that many implementations will have such features, and that such features are a sign of the quality of the implementation.

### 2.5.1 FORTRAN 77 BINDING ISSUES

All MPI names have an `MPI_` prefix, and all characters are capitals. Programs must not declare variables or functions with names beginning with the prefix, `MPI_`. This is mandated to avoid possible name collisions.

All MPI Fortran subroutines have a return code in the last argument. A few MPI operations are functions, which do not have the return code argument. The return code value for successful completion is `MPI_SUCCESS`. Other error codes are implementation dependent; see Chapter 7.

Handles are represented in Fortran as `INTEGER`s. Binary-valued variables are of type `LOGICAL`.

Array arguments are indexed from one.

Unless explicitly stated, the MPI F77 binding is consistent with ANSI standard Fortran 77. There are several points where this standard diverges from the ANSI Fortran 77 standard. These exceptions are consistent with common practice in the Fortran community. In particular:

- MPI identifiers are limited to thirty, not six, significant characters.
- MPI identifiers may contain underscores after the first character.
- An MPI subroutine with a choice argument may be called with different argument types. An example is shown in Figure 2.1. This violates the letter of the Fortran standard, but such a violation is common practice. An alternative would be to have a separate version of `MPI_SEND` for each data type.
- Although not required, it is strongly suggested that named MPI constants (`PARAMETERS`) be provided in an include file, called `mpif.h`. On systems that do not support include files, the implementation should specify the values of named constants.
- Vendors are encouraged to provide type declarations in the `mpif.h` file on Fortran systems that support user-defined types. One should define, if possible, the type `MPI_ADDRESS`, which is an `INTEGER` of the size needed

to hold an address in the execution environment. On systems where type definition is not supported, it is up to the user to use an `INTEGER` of the right kind to represent addresses (i.e., `INTEGER*4` on a 32 bit machine, `INTEGER*8` on a 64 bit machine, etc.).

## 2.5.2 C BINDING ISSUES

We use the ANSI C declaration format. All MPI names have an `MPI_` prefix, defined constants are in all capital letters, and defined types and functions have one capital letter after the prefix. Programs must not declare variables or functions with names beginning with the prefix, `MPI_`. This is mandated to avoid possible name collisions.

The definition of named constants, function prototypes, and type definitions must be supplied in an include file `mpi.h`.

Almost all C functions return an error code. The successful return code will be `MPI_SUCCESS`, but failure return codes are implementation dependent. A few C functions do not return values, so that they can be implemented as macros.

Type declarations are provided for handles to each category of opaque objects. Either a pointer or an integer type is used.

Array arguments are indexed from zero.

Logical flags are integers with value 0 meaning "false" and a non-zero value meaning "true."

Choice arguments are pointers of type `void*`.

Address arguments are of MPI defined type `MPLAint`. This is defined to be an int of the size needed to hold any valid address on the target architecture.

## 2.6 Processes

An MPI program consists of autonomous processes, executing their own code, in an MIMD style. The codes executed by each process need not be identical. The processes communicate via calls to MPI communication primitives. Typically, each process executes in its own address space, although shared-memory implementations of MPI are possible. This document specifies the behavior of a parallel program assuming that only MPI calls are used for communication. The interaction of an MPI program with other possible means of communication (e.g., shared memory) is not specified.

MPI does not specify the execution model for each process. A process can be sequential, or can be multi-threaded, with threads possibly executing concurrently. Care has been taken to make MPI "thread-safe," by avoiding the use of implicit state. The desired interaction of MPI with threads is that concurrent threads be all allowed to execute MPI calls, and calls be reentrant; a blocking MPI call blocks only the invoking thread, allowing the scheduling of another thread.

MPI does not provide mechanisms to specify the initial allocation of processes to an MPI computation and their binding to physical processors. It is expected

that vendors will provide mechanisms to do so either at load time or at run time. Such mechanisms will allow the specification of the initial number of required processes, the code to be executed by each initial process, and the allocation of processes to processors. Also, the current proposal does not provide for dynamic creation or deletion of processes during program execution (the total number of processes is fixed), although it is intended to be consistent with such extensions. Finally, we always identify processes according to their relative rank in a group, that is, consecutive integers in the range  $0..groupsize-1$ .

## 2.7 Error Handling

MPI provides the user with reliable message transmission. A message sent is always received correctly, and the user does not need to check for transmission errors, time-outs, or other error conditions. In other words, MPI does not provide mechanisms for dealing with failures in the communication system. If the MPI implementation is built on an unreliable underlying mechanism, then it is the job of the implementor of the MPI subsystem to insulate the user from this unreliability, or to reflect unrecoverable errors as failures. Whenever possible, such failures will be reflected as errors in the relevant communication call. Similarly, MPI itself provides no mechanisms for handling processor failures. The error handling facilities described in Section 7.2 can be used to restrict the scope of an unrecoverable error, or design error recovery at the application level.

Of course, MPI programs may still be erroneous. A **program error** can occur when an MPI call is called with an incorrect argument (non-existing destination in a send operation, buffer too small in a receive operation, etc.). This type of error would occur in any implementation. In addition, a **resource error** may occur when a program exceeds the amount of available system resources (number of pending messages, system buffers, etc.). The occurrence of this type of error depends on the amount of available resources in the system and the resource allocation mechanism used; this may differ from system to system. A high-quality implementation will provide generous limits on the important resources so as to alleviate the portability problem this represents.

Almost all MPI calls return a code that indicates successful completion of the operation. Whenever possible, MPI calls return an error code, if an error occurred during the call. By default, an error detected during the execution of the MPI library causes the parallel computation to abort. However, MPI provides mechanisms for users to change this default and to handle recoverable errors. The user may specify that no error is fatal, and handle error codes returned by MPI calls by himself or herself. Also, the user may provide his or her own error-handling routines, which will be invoked whenever an MPI call returns abnormally. The MPI error handling facilities are described in Section 7.2.

Several factors limit the ability of MPI calls to return with meaningful error codes when an error occurs. MPI may not be able to detect some errors; other errors may be too expensive to detect in normal execution mode; finally some



errors may be “catastrophic” and may prevent MPI from returning control to the caller in a consistent state.

Another subtle issue arises because of the nature of asynchronous communications: MPI calls may initiate operations that continue asynchronously after the call returned. Thus, the operation may return with a code indicating successful completion, yet later cause an error exception to be raised. If there is a subsequent call that relates to the same operation (e.g., a call that verifies that an asynchronous operation has completed) then the error argument associated with this call will be used to indicate the nature of the error. In a few cases, the error may occur after all calls that relate to the operation have completed, so that no error value can be used to indicate the nature of the error (e.g., an error in a send with the ready mode). Such an error must be treated as fatal, since information cannot be returned for the user to recover from it.

This document does not specify the state of a computation after an erroneous MPI call has occurred. The desired behavior is that a relevant error code be returned, and the effect of the error be localized to the greatest possible extent. E.g., it is highly desirable that an erroneous receive call will not cause any part of the receiver’s memory to be overwritten, beyond the area specified for receiving the message.

Implementations may go beyond this document in supporting in a meaningful manner MPI calls that are defined here to be erroneous. For example, MPI specifies strict type matching rules between matching send and receive operations: it is erroneous to send a floating point variable and receive an integer. Implementations may go beyond these type matching rules, and provide automatic type conversion in such situations. It will be helpful to generate warnings for such nonconforming behavior.

## 2.8 Implementation Issues

There are a number of areas where an MPI implementation may interact with the operating environment and system. While MPI does not mandate that any services (such as I/O or signal handling) be provided, it does strongly suggest the behavior to be provided if those services are available. This is an important point in achieving portability across platforms that provide the same set of services.

### 2.8.1 INDEPENDENCE OF BASIC RUNTIME ROUTINES

MPI programs require that library routines that are part of the basic language environment (such as `date` and `write` in Fortran and `printf` and `malloc` in ANSI C) and are executed after `MPI_INIT` and before `MPI_FINALIZE` operate independently and that their *completion* is independent of the action of other processes in an MPI program.

Note that this in no way prevents the creation of library routines that provide parallel services whose operation is collective. However, the following program

is expected to complete in an ANSI C environment regardless of the size of `MPI_COMM_WORLD` (assuming that I/O is available at the executing nodes).

```
int rank;
MPI_Init( argc, argv );
MPI_Comm_rank( MPI_COMM_WORLD, &rank );
if (rank == 0) printf( "Starting program\n" );
MPI_Finalize();
```

The corresponding Fortran 77 program is also expected to complete.

An example of what is *not* required is any particular ordering of the action of these routines when called by several tasks. For example, MPI makes neither requirements nor recommendations for the output from the following program (again assuming that I/O is available at the executing nodes).

```
MPI_Comm_rank( MPI_COMM_WORLD, &rank );
printf( "Output from task rank %d\n", rank );
```

In addition, calls that fail because of resource exhaustion or other error are not considered a violation of the requirements here (however, they are required to complete, just not to complete successfully).

### 2.8.2 INTERACTION WITH SIGNALS IN POSIX

MPI does not specify either the interaction of processes with signals, in a UNIX environment, or with other events that do not relate to MPI communication. That is, signals are not significant from the viewpoint of MPI, and implementors should attempt to implement MPI so that signals are transparent: an MPI call suspended by a signal should resume and complete after the signal is handled. Generally, the state of a computation that is visible or significant from the viewpoint of MPI should only be affected by MPI calls.

The intent of MPI to be thread and signal safe has a number of subtle effects. For example, on Unix systems, a catchable signal such as `SIGALRM` (an alarm signal) must not cause an MPI routine to behave differently than it would have in the absence of the signal. Of course, if the signal handler issues MPI calls or changes the environment in which the MPI routine is operating (for example, consuming all available memory space), the MPI routine should behave as appropriate for that situation (in particular, in this case, the behavior should be the same as for a multi-threaded MPI implementation).

A second effect is that a signal handler that performs MPI calls must not interfere with the operation of MPI. For example, an MPI receive of any type that occurs within a signal handler must not cause erroneous behavior by the MPI implementation. Note that an implementation is permitted to prohibit the use of MPI calls from within a signal handler, and is not required to detect such use.

It is highly desirable that MPI not use `SIGALRM`, `SIGFPE`, or `SIGIO`. An implementation is *required* to clearly document all of the signals that the MPI implementation uses; a good place for this information is a Unix 'man' page on MPI.

## POINT-TO-POINT COMMUNICATION

## 3.1 Introduction

Sending and receiving of messages by processes is the basic MPI communication mechanism. The basic point-to-point communication operations are **send** and **receive**. Their use is illustrated in the example below.

```
#include "mpi.h"
main( argc, argv )
int argc;
char **argv;
{
    char message[20];
    int myrank;
    MPI_Status status;
    MPI_Init( &argc, &argv );
    MPI_Comm_rank( MPI_COMM_WORLD, &myrank );
    if (myrank == 0) /* code for process zero */
    {
        strcpy(message, "Hello, there");
        MPI_Send(message, strlen(message), MPI_CHAR, 1, 99, MPI_COMM_WORLD);
    }
    else /* code for process one */
    {
        MPI_Recv(message, 20, MPI_CHAR, 0, 99, MPI_COMM_WORLD, &status);
        printf("received :%s:\n", message);
    }
    MPI_Finalize();
}
```

In this example, process zero (`myrank = 0`) sends a message to process one using the **send** operation `MPI_Send`. The operation specifies a **send buffer** in the sender memory from which the message data is taken. In the example above, the send buffer consists of the storage containing the variable `message` in the

memory of process zero. The location, size, and type of the send buffer are specified by the first three parameters of the send operation. The message sent will contain the 13 characters of this variable. In addition, the send operation associates an **envelope** with the message. This envelope specifies the message destination and contains distinguishing information that can be used by the **receive** operation to select a particular message. The last three parameters of the send operation specify the envelope for the message sent.

Process one (`myrank = 1`) receives this message with the **receive** operation `MPI_RECV`. The message to be received is selected according to the value of its envelope, and the message data is stored into the **receive buffer**. In the example above, the receive buffer consists of the storage containing the string `message` in the memory of process one. The first three parameters of the receive operation specify the location, size, and type of the receive buffer. The next three parameters are used for selecting the incoming message. The last parameter is used to return information on the message just received.

The next sections describe the blocking send and receive operations. We discuss send, receive, blocking communication semantics, type matching requirements, type conversion in heterogeneous environments, and more general communication modes. Nonblocking communication is addressed next, followed by channel-like constructs and send-receive operations. We then consider general datatypes that allow one to transfer efficiently heterogeneous and noncontiguous data. We conclude with the description of calls for explicit packing and unpacking of messages.

## 3.2 Blocking Send and Receive Operations

### 3.2.1 BLOCKING SEND

The syntax of the blocking send operation is given below.

`MPI_SEND(buf, count, datatype, dest, tag, comm)`

IN	buf	initial address of send buffer (choice)
IN	count	number of elements in send buffer (nonnegative integer)
IN	datatype	datatype of each send buffer element (handle)
IN	dest	rank of destination (integer)
IN	tag	message tag (integer)
IN	comm	communicator (handle)

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

```
MPI_SEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR)
<type> BUF(*)
INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR
```

The blocking semantics of this call are described in Section 3.4.