

MPI: A Message-Passing Interface Standard

Version 2.1

(draft, with MPI 2.1 Ballots 1-4 and Reviews 1-22)

Message Passing Interface Forum

April 3, 2008

This work was supported in part by ARPA, NSF and DARPA under grant ASC-9310330, the National Science Foundation Science and Technology Center Cooperative Agreement No. CCR-8809615, and the NSF contract CDA-9115428, and by the Commission of the European Community through Esprit project P6643 and under project HPC Standards (21111).

1 This document describes the Message Passing Interface (MPI) standard, version 2.1.
2 The standard MPI includes point-to-point message passing, collective communications,
3 group and communicator concepts, process topologies, environmental management, pro-
4 cess creation and management, one-sided communications, extended collective operations,
5 external interfaces, I/O, some miscellaneous topics, and a profiling interface. Language
6 bindings for C, C++ and Fortran are defined.

7 Technically, this version of the standard is based on the "MPI: A Message-Passing In-
8 terface Standard, June 12, 1995" (MPI-1.1) from the MPI-1 Forum, and "MPI-2: Extensions
9 to the Message-Passing Interface, July, 1997" (MPI-1.2 and MPI-2.0) from the MPI-2 Forum,
10 and errata documents from the MPI Forum.

11 Historically, the evolution of the standards is MPI-1.0 (June 1994), MPI-1.1 (June 12,
12 1995), MPI-1.2 (July 18, 1997), with several clarifications and additions and published as
13 part of the MPI-2 document, MPI-2.0 (July 18, 1997), with new functionality, MPI-1.3
14 (date, 2008), combining for historical reasons the documents 1.1 and 1.2 and some errata
15 documnts to one combined document, and this document, MPI-2.1, combining the previous
16 documents. Additional clarifications and errata corrections to MPI-2.0 are also included.

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45 ©1993, 1994, 1995, 1996, 1997, 2008 University of Tennessee, Knoxville, Tennessee.
46 Permission to copy without fee all or part of this material is granted, provided the University
47 of Tennessee copyright notice and the title of this document appear, and notice is given that
48 copying is by permission of the University of Tennessee.

Version 2.1: XXXXXX, 2008. This document combines the previous documents MPI-1.3 (XXXXXX, 2008) and MPI-2.0 (July 18, 1997). Certain parts of MPI-2.0, such as some sections of Chapter 4, Miscellany, and Chapter 7, Extended Collective Operations have been merged into the Chapters of MPI-1.3. Additional errata and clarifications collected by the MPI Forum are also included in this document.

Version 1.3: XXXXXX, 2008. This document combines the previous documents MPI-1.1 (June 12, 1995) and the MPI-1.2 Chapter in MPI-2 (July 18, 1997). Additional errata collected by the MPI Forum referring to MPI-1.1 and MPI-1.2 are also included in this document.

Version 2.0: July 18, 1997. Beginning after the release of MPI-1.1, the MPI Forum began meeting to consider corrections and extensions. MPI-2 has been focused on process creation and management, one-sided communications, extended collective communications, external interfaces and parallel I/O. A miscellany chapter discusses items that don't fit elsewhere, in particular language interoperability.

Version 1.2: July 18, 1997. The MPI-2 Forum introduced MPI-1.2 as Chap.3 in the standard "MPI-2: Extensions to the Message-Passing Interface", July 18, 1997. This section contains clarifications and minor corrections to Version 1.1 of the MPI Standard. The only new function in MPI-1.2 is one for identifying to which version of the MPI Standard the implementation conforms. There are small differences between MPI-1 and MPI-1.1. There are very few differences between MPI-1.1 and MPI-1.2, but large differences between MPI-1.2 and MPI-2.

Version 1.1: June, 1995. Beginning in March, 1995, the Message Passing Interface Forum reconvened to correct errors and make clarifications in the MPI document of May 5, 1994, referred to below as Version 1.0. These discussions resulted in Version 1.1, which is this document. The changes from Version 1.0 are minor. A version of this document with all changes marked is available. This paragraph is an example of a change.

Version 1.0: May, 1994. The Message Passing Interface Forum (MPIF), with participation from over 40 organizations, has been meeting since January 1993 to discuss and define a set of library interface standards for message passing. MPIF is not sanctioned or supported by any official standards organization.

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.

This is the final report, Version 1.0, of the Message Passing Interface Forum. This document contains all the technical features proposed for the interface. This copy of the draft was processed by L^AT_EX on May 5, 1994.

Please send comments on MPI to mpi-comments@mpi-forum.org. Your comment will be forwarded to MPI Forum committee members who will attempt to respond.

Contents

| | |
|--|------------|
| Acknowledgments | xvi |
| 1 Introduction to MPI | 1 |
| 1.1 Overview and Goals | 1 |
| 1.2 Background of MPI-1 | 1 |
| 1.3 Background of MPI-2 | 3 |
| 1.4 Background of MPI-1.3 and MPI-2.1 | 3 |
| 1.5 Who Should Use This Standard? | 4 |
| 1.6 What Platforms Are Targets For Implementation? | 4 |
| 1.7 What Is Included In The Standard? | 5 |
| 1.8 What Is Not Included In The Standard? | 5 |
| 1.9 Organization of this Document | 6 |
| 2 MPI Terms and Conventions | 9 |
| 2.1 Document Notation | 9 |
| 2.2 Naming Conventions | 9 |
| 2.3 Procedure Specification | 10 |
| 2.4 Semantic Terms | 11 |
| 2.5 Data Types | 12 |
| 2.5.1 Opaque Objects | 12 |
| 2.5.2 Array Arguments | 14 |
| 2.5.3 State | 14 |
| 2.5.4 Named Constants | 14 |
| 2.5.5 Choice | 15 |
| 2.5.6 Addresses | 15 |
| 2.5.7 File Offsets | 15 |
| 2.6 Language Binding | 15 |
| 2.6.1 Deprecated Names and Functions | 16 |
| 2.6.2 Fortran Binding Issues | 16 |
| 2.6.3 C Binding Issues | 18 |
| 2.6.4 C++ Binding Issues | 18 |
| 2.6.5 Functions and Macros | 21 |
| 2.7 Processes | 22 |
| 2.8 Error Handling | 22 |
| 2.9 Implementation Issues | 23 |
| 2.9.1 Independence of Basic Runtime Routines | 23 |
| 2.9.2 Interaction with Signals | 24 |

| | | |
|----------|---|-----------|
| 2.10 | Examples | 24 |
| 3 | Point-to-Point Communication | 25 |
| 3.1 | Introduction | 25 |
| 3.2 | Blocking Send and Receive Operations | 26 |
| 3.2.1 | Blocking send | 26 |
| 3.2.2 | Message data | 27 |
| 3.2.3 | Message envelope | 28 |
| 3.2.4 | Blocking receive | 29 |
| 3.2.5 | Return status | 31 |
| 3.2.6 | Passing MPI_STATUS_IGNORE for Status | 33 |
| 3.3 | Data type matching and data conversion | 33 |
| 3.3.1 | Type matching rules | 33 |
| 3.3.2 | Data conversion | 36 |
| 3.4 | Communication Modes | 37 |
| 3.5 | Semantics of point-to-point communication | 41 |
| 3.6 | Buffer allocation and usage | 45 |
| 3.6.1 | Model implementation of buffered mode | 46 |
| 3.7 | Nonblocking communication | 47 |
| 3.7.1 | Communication Objects | 48 |
| 3.7.2 | Communication initiation | 49 |
| 3.7.3 | Communication Completion | 51 |
| 3.7.4 | Semantics of Nonblocking Communications | 55 |
| 3.7.5 | Multiple Completions | 56 |
| 3.7.6 | Non-destructive Test of status | 63 |
| 3.8 | Probe and Cancel | 63 |
| 3.9 | Persistent communication requests | 68 |
| 3.10 | Send-receive | 72 |
| 3.11 | Null processes | 74 |
| 3.12 | Derived datatypes | 75 |
| 3.12.1 | New Datatype Manipulation Functions | 76 |
| 3.12.2 | Type Constructors with Explicit Addresses | 77 |
| 3.12.3 | Datatype constructors | 77 |
| 3.12.4 | Subarray Datatype Constructor | 85 |
| 3.12.5 | Distributed Array Datatype Constructor | 87 |
| 3.12.6 | Address and size functions | 92 |
| 3.12.7 | Lower-bound and upper-bound markers | 94 |
| 3.12.8 | Extent and Bounds of Datatypes | 95 |
| 3.12.9 | True Extent of Datatypes | 96 |
| 3.12.10 | Commit and free | 97 |
| 3.12.11 | Duplicating a Datatype | 98 |
| 3.12.12 | Use of general datatypes in communication | 99 |
| 3.12.13 | Correct use of addresses | 101 |
| 3.12.14 | Examples | 102 |
| 3.13 | Pack and unpack | 110 |
| 3.14 | Canonical MPI_PACK and MPI_UNPACK | 116 |

| | | |
|----------|--|------------|
| 4 | Collective Communication | 119 |
| 4.1 | Introduction and Overview | 119 |
| 4.2 | Communicator argument | 122 |
| 4.3 | Extended Collective Operations | 122 |
| 4.3.1 | Introduction | 122 |
| 4.3.2 | Intercommunicator Collective Operations | 122 |
| 4.3.3 | Operations that Move Data | 125 |
| 4.4 | Barrier synchronization | 126 |
| 4.5 | Broadcast | 126 |
| 4.5.1 | Example using MPI_BCAST | 127 |
| 4.6 | Gather | 127 |
| 4.6.1 | Examples using MPI_GATHER, MPI_GATHERV | 130 |
| 4.7 | Scatter | 137 |
| 4.7.1 | Examples using MPI_SCATTER, MPI_SCATTERV | 139 |
| 4.8 | Gather-to-all | 142 |
| 4.8.1 | Examples using MPI_ALLGATHER, MPI_ALLGATHERV | 144 |
| 4.9 | All-to-All Scatter/Gather | 145 |
| 4.9.1 | Generalized All-to-all Function | 147 |
| 4.10 | Global Reduction Operations | 149 |
| 4.10.1 | Reduce | 149 |
| 4.10.2 | Predefined reduce operations | 151 |
| 4.10.3 | Signed Characters and Reductions | 153 |
| 4.10.4 | MINLOC and MAXLOC | 153 |
| 4.10.5 | User-Defined Operations | 157 |
| 4.10.6 | All-Reduce | 161 |
| 4.11 | Reduce-Scatter | 162 |
| 4.12 | Scan | 164 |
| 4.12.1 | Exclusive Scan | 164 |
| 4.12.2 | Example using MPI_SCAN | 165 |
| 4.13 | Correctness | 167 |
| 5 | Groups, Contexts, and Communicators | 171 |
| 5.1 | Introduction | 171 |
| 5.1.1 | Features Needed to Support Libraries | 171 |
| 5.1.2 | MPI's Support for Libraries | 172 |
| 5.2 | Basic Concepts | 174 |
| 5.2.1 | Groups | 174 |
| 5.2.2 | Contexts | 174 |
| 5.2.3 | Intra-Communicators | 175 |
| 5.2.4 | Predefined Intra-Communicators | 175 |
| 5.3 | Group Management | 176 |
| 5.3.1 | Group Accessors | 176 |
| 5.3.2 | Group Constructors | 177 |
| 5.3.3 | Group Destructors | 182 |
| 5.4 | Communicator Management | 183 |
| 5.4.1 | Communicator Accessors | 183 |
| 5.4.2 | Communicator Constructors | 185 |
| 5.4.3 | Communicator Destructors | 191 |

| | | |
|----------|--|------------|
| 5.5 | Motivating Examples | 192 |
| 5.5.1 | Current Practice #1 | 192 |
| 5.5.2 | Current Practice #2 | 193 |
| 5.5.3 | (Approximate) Current Practice #3 | 193 |
| 5.5.4 | Example #4 | 194 |
| 5.5.5 | Library Example #1 | 195 |
| 5.5.6 | Library Example #2 | 197 |
| 5.6 | Inter-Communication | 199 |
| 5.6.1 | Inter-communicator Accessors | 201 |
| 5.6.2 | Inter-communicator Operations | 202 |
| 5.6.3 | Inter-Communication Examples | 204 |
| 5.7 | Caching | 211 |
| 5.7.1 | New Attribute Caching Functions | 212 |
| 5.7.2 | Functionality | 212 |
| 5.7.3 | Communicators | 213 |
| 5.7.4 | Windows | 217 |
| 5.7.5 | Datatypes | 220 |
| 5.7.6 | Error Class for Invalid Keyval | 223 |
| 5.7.7 | Attributes Example | 223 |
| 5.8 | Formalizing the Loosely Synchronous Model | 225 |
| 5.8.1 | Basic Statements | 225 |
| 5.8.2 | Models of Execution | 225 |
| 6 | Process Topologies | 229 |
| 6.1 | Introduction | 229 |
| 6.2 | Virtual Topologies | 230 |
| 6.3 | Embedding in MPI | 230 |
| 6.4 | Overview of the Functions | 231 |
| 6.5 | Topology Constructors | 232 |
| 6.5.1 | Cartesian Constructor | 232 |
| 6.5.2 | Cartesian Convenience Function: <code>MPI_DIMS_CREATE</code> | 232 |
| 6.5.3 | General (Graph) Constructor | 234 |
| 6.5.4 | Topology inquiry functions | 236 |
| 6.5.5 | Cartesian Shift Coordinates | 240 |
| 6.5.6 | Partitioning of Cartesian structures | 242 |
| 6.5.7 | Low-level topology functions | 242 |
| 6.6 | An Application Example | 244 |
| 7 | MPI Environmental Management | 247 |
| 7.1 | Implementation information | 247 |
| 7.1.1 | Version Inquiries | 247 |
| 7.1.2 | Environmental Inquiries | 248 |
| 7.2 | Memory Allocation | 250 |
| 7.3 | Error handling | 252 |
| 7.3.1 | Extended Error Handling in MPI-2 | 253 |
| 7.3.2 | Error Handlers for Communicators | 254 |
| 7.3.3 | Error Handlers for Windows | 256 |
| 7.3.4 | Error Handlers for Files | 257 |

| | | |
|-----------|--|------------|
| 7.3.5 | Freeing Errorhandlers and Retrieving Error Strings | 258 |
| 7.4 | Error codes and classes | 259 |
| 7.5 | Timers and synchronization | 262 |
| 7.6 | Startup | 263 |
| 7.6.1 | Allowing User Functions at Process Termination | 268 |
| 7.6.2 | Determining Whether MPI Has Finished | 268 |
| 7.7 | Portable MPI Process Startup | 269 |
| 8 | Miscellany | 273 |
| 8.1 | The Info Object | 273 |
| 9 | Process Creation and Management | 279 |
| 9.1 | Introduction | 279 |
| 9.2 | The MPI-2 Process Model | 280 |
| 9.2.1 | Starting Processes | 280 |
| 9.2.2 | The Runtime Environment | 280 |
| 9.3 | Process Manager Interface | 282 |
| 9.3.1 | Processes in MPI | 282 |
| 9.3.2 | Starting Processes and Establishing Communication | 282 |
| 9.3.3 | Starting Multiple Executables and Establishing Communication | 287 |
| 9.3.4 | Reserved Keys | 289 |
| 9.3.5 | Spawn Example | 290 |
| 9.4 | Establishing Communication | 292 |
| 9.4.1 | Names, Addresses, Ports, and All That | 292 |
| 9.4.2 | Server Routines | 294 |
| 9.4.3 | Client Routines | 295 |
| 9.4.4 | Name Publishing | 297 |
| 9.4.5 | Reserved Key Values | 299 |
| 9.4.6 | Client/Server Examples | 299 |
| 9.5 | Other Functionality | 302 |
| 9.5.1 | Universe Size | 302 |
| 9.5.2 | Singleton MPI_INIT | 302 |
| 9.5.3 | MPI_APPNUM | 303 |
| 9.5.4 | Releasing Connections | 304 |
| 9.5.5 | Another Way to Establish MPI Communication | 305 |
| 10 | One-Sided Communications | 307 |
| 10.1 | Introduction | 307 |
| 10.2 | Initialization | 308 |
| 10.2.1 | Window Creation | 308 |
| 10.2.2 | Window Attributes | 310 |
| 10.3 | Communication Calls | 311 |
| 10.3.1 | Put | 312 |
| 10.3.2 | Get | 314 |
| 10.3.3 | Examples | 314 |
| 10.3.4 | Accumulate Functions | 317 |
| 10.4 | Synchronization Calls | 319 |
| 10.4.1 | Fence | 324 |

| | | |
|-----------|--|------------|
| 10.4.2 | General Active Target Synchronization | 325 |
| 10.4.3 | Lock | 328 |
| 10.4.4 | Assertions | 330 |
| 10.4.5 | Miscellaneous Clarifications | 332 |
| 10.5 | Examples | 332 |
| 10.6 | Error Handling | 334 |
| 10.6.1 | Error Handlers | 334 |
| 10.6.2 | Error Classes | 335 |
| 10.7 | Semantics and Correctness | 335 |
| 10.7.1 | Atomicity | 338 |
| 10.7.2 | Progress | 338 |
| 10.7.3 | Registers and Compiler Optimizations | 340 |
| 11 | External Interfaces | 343 |
| 11.1 | Introduction | 343 |
| 11.2 | Generalized Requests | 343 |
| 11.2.1 | Examples | 347 |
| 11.3 | Associating Information with Status | 349 |
| 11.4 | Naming Objects | 351 |
| 11.5 | Error Classes, Error Codes, and Error Handlers | 355 |
| 11.6 | Decoding a Datatype | 358 |
| 11.7 | MPI and Threads | 367 |
| 11.7.1 | General | 367 |
| 11.7.2 | Clarifications | 368 |
| 11.7.3 | Initialization | 370 |
| 12 | I/O | 375 |
| 12.1 | Introduction | 375 |
| 12.1.1 | Definitions | 375 |
| 12.2 | File Manipulation | 377 |
| 12.2.1 | Opening a File | 377 |
| 12.2.2 | Closing a File | 379 |
| 12.2.3 | Deleting a File | 380 |
| 12.2.4 | Resizing a File | 381 |
| 12.2.5 | Preallocating Space for a File | 381 |
| 12.2.6 | Querying the Size of a File | 382 |
| 12.2.7 | Querying File Parameters | 382 |
| 12.2.8 | File Info | 384 |
| 12.3 | File Views | 387 |
| 12.4 | Data Access | 389 |
| 12.4.1 | Data Access Routines | 389 |
| 12.4.2 | Data Access with Explicit Offsets | 392 |
| 12.4.3 | Data Access with Individual File Pointers | 396 |
| 12.4.4 | Data Access with Shared File Pointers | 401 |
| 12.4.5 | Split Collective Data Access Routines | 406 |
| 12.5 | File Interoperability | 412 |
| 12.5.1 | Datatypes for File Interoperability | 414 |
| 12.5.2 | External Data Representation: “external32” | 416 |

| | | |
|-----------|--|------------|
| 12.5.3 | User-Defined Data Representations | 418 |
| 12.5.4 | Matching Data Representations | 421 |
| 12.6 | Consistency and Semantics | 422 |
| 12.6.1 | File Consistency | 422 |
| 12.6.2 | Random Access vs. Sequential Files | 425 |
| 12.6.3 | Progress | 425 |
| 12.6.4 | Collective File Operations | 425 |
| 12.6.5 | Type Matching | 426 |
| 12.6.6 | Miscellaneous Clarifications | 426 |
| 12.6.7 | MPI_Offset Type | 426 |
| 12.6.8 | Logical vs. Physical File Layout | 426 |
| 12.6.9 | File Size | 427 |
| 12.6.10 | Examples | 427 |
| 12.7 | I/O Error Handling | 431 |
| 12.8 | I/O Error Classes | 432 |
| 12.9 | Examples | 432 |
| 12.9.1 | Double Buffering with Split Collective I/O | 432 |
| 12.9.2 | Subarray Filetype Constructor | 435 |
| 13 | Language Bindings | 437 |
| 13.1 | C++ | 437 |
| 13.1.1 | Overview | 437 |
| 13.1.2 | Design | 437 |
| 13.1.3 | C++ Classes for MPI | 438 |
| 13.1.4 | Class Member Functions for MPI | 438 |
| 13.1.5 | Semantics | 439 |
| 13.1.6 | C++ Datatypes | 441 |
| 13.1.7 | Communicators | 444 |
| 13.1.8 | Exceptions | 445 |
| 13.1.9 | Mixed-Language Operability | 446 |
| 13.1.10 | Profiling | 446 |
| 13.2 | Fortran Support | 449 |
| 13.2.1 | Overview | 449 |
| 13.2.2 | Problems With Fortran Bindings for MPI | 450 |
| 13.2.3 | Basic Fortran Support | 456 |
| 13.2.4 | Extended Fortran Support | 457 |
| 13.2.5 | Additional Support for Fortran Numeric Intrinsic Types | 458 |
| 13.3 | Language Interoperability | 466 |
| 13.3.1 | Introduction | 466 |
| 13.3.2 | Assumptions | 466 |
| 13.3.3 | Initialization | 467 |
| 13.3.4 | Transfer of Handles | 467 |
| 13.3.5 | Status | 471 |
| 13.3.6 | MPI Opaque Objects | 471 |
| 13.3.7 | Attributes | 474 |
| 13.3.8 | Extra State | 476 |
| 13.3.9 | Constants | 476 |
| 13.3.10 | Interlanguage Communication | 477 |

| | |
|--|------------|
| 14 Profiling Interface | 479 |
| 14.1 Requirements | 479 |
| 14.2 Discussion | 479 |
| 14.3 Logic of the design | 480 |
| 14.3.1 Miscellaneous control of profiling | 480 |
| 14.4 Examples | 481 |
| 14.4.1 Profiler implementation | 481 |
| 14.4.2 MPI library implementation | 482 |
| 14.4.3 Complications | 483 |
| 14.5 Multiple levels of interception | 484 |
| | |
| 15 Deprecated Functions | 485 |
| 15.1 Deprecated since MPI-2.0 | 485 |
| | |
| A Language Binding | 492 |
| A.1 Defined Values and Handles | 492 |
| A.1.1 Defined Constants | 492 |
| A.1.2 Type and prototype definitions | 500 |
| A.1.3 Info Keys | 501 |
| A.1.4 Info Values | 501 |
| A.2 C Bindings | 502 |
| A.2.1 Point-to-Point Communication C Bindings | 502 |
| A.2.2 Collective Communication C Bindings | 505 |
| A.2.3 Groups, Contexts, and Communicators C Bindings | 506 |
| A.2.4 Process Topologies C Bindings | 508 |
| A.2.5 MPI Environmenta Management C Bindings | 508 |
| A.2.6 Miscellany C Bindings | 509 |
| A.2.7 Process Creation and Management C Bindings | 510 |
| A.2.8 One-Sided Communications C Bindings | 510 |
| A.2.9 External Interfaces C Bindings | 511 |
| A.2.10 I/O C Bindings | 512 |
| A.2.11 Language Bindings C Bindings | 514 |
| A.2.12 Profiling Interface C Bindings | 515 |
| A.2.13 Deprecated C Bindings | 515 |
| A.3 Fortran Bindings | 517 |
| A.3.1 Point-to-Point Communication Fortran Bindings | 517 |
| A.3.2 Collective Communication Fortran Bindings | 521 |
| A.3.3 Groups, Contexts, and Communicators Fortran Bindings | 523 |
| A.3.4 Process Topologies Fortran Bindings | 526 |
| A.3.5 MPI Environmenta Management Fortran Bindings | 527 |
| A.3.6 Miscellany Fortran Bindings | 528 |
| A.3.7 Process Creation and Management Fortran Bindings | 529 |
| A.3.8 One-Sided Communications Fortran Bindings | 530 |
| A.3.9 External Interfaces Fortran Bindings | 531 |
| A.3.10 I/O Fortran Bindings | 532 |
| A.3.11 Language Bindings Fortran Bindings | 537 |
| A.3.12 Profiling Interface Fortran Bindings | 537 |
| A.3.13 Deprecated Fortran Bindings | 537 |

| | | |
|----------|--|------------|
| A.4 | C++ Bindings | 539 |
| A.4.1 | Point-to-Point Communication C++ Bindings | 539 |
| A.4.2 | Collective Communication C++ Bindings | 542 |
| A.4.3 | Groups, Contexts, and Communicators C++ Bindings | 544 |
| A.4.4 | Process Topologies C++ Bindings | 546 |
| A.4.5 | MPI Environmenta Management C++ Bindings | 546 |
| A.4.6 | Miscellany C++ Bindings | 548 |
| A.4.7 | Process Creation and Management C++ Bindings | 548 |
| A.4.8 | One-Sided Communications C++ Bindings | 549 |
| A.4.9 | External Interfaces C++ Bindings | 550 |
| A.4.10 | I/O C++ Bindings | 551 |
| A.4.11 | Language Bindings C++ Bindings | 554 |
| A.4.12 | Profiling Interface C++ Bindings | 554 |
| A.4.13 | Deprecated C++ Bindings | 554 |
| A.4.14 | C++ Bindings on all MPI Classes | 555 |
| A.4.15 | Construction / Destruction | 555 |
| A.4.16 | Copy / Assignment | 555 |
| A.4.17 | Comparison | 555 |
| A.4.18 | Inter-language Operability | 555 |
| A.4.19 | Function Name Cross Reference | 556 |
| B | Change-Log | 560 |
| B.1 | Changes from Version 2.0 to Version 2.1 | 560 |
| | Bibliography | 565 |
| | Index | 569 |

List of Figures

| | | |
|------|--|-----|
| 4.1 | Collective communications, an overview | 120 |
| 4.2 | Intercommunicator allgather | 124 |
| 4.3 | Intercommunicator reduce-scatter | 124 |
| 4.4 | Gather example | 131 |
| 4.5 | Gatherv example with strides | 132 |
| 4.6 | Gatherv example, 2-dimensional | 133 |
| 4.7 | Gatherv example, 2-dimensional, subarrays with different sizes | 134 |
| 4.8 | Gatherv example, 2-dimensional, subarrays with different sizes and strides | 135 |
| 4.9 | Scatter example | 140 |
| 4.10 | Scatterv example with strides | 141 |
| 4.11 | Scatterv example with different strides and counts | 142 |
| 4.12 | Race conditions with point-to-point and collective communications | 169 |
| | | |
| 5.1 | Intercommunicator create using <code>MPI_COMM_CREATE</code> | 188 |
| 5.2 | Intercommunicator construction with <code>MPI_COMM_SPLIT</code> | 190 |
| 5.3 | Three-group pipeline. | 204 |
| 5.4 | Three-group ring. | 206 |
| | | |
| 6.1 | Set-up of process structure for two-dimensional parallel Poisson solver. | 245 |
| | | |
| 10.1 | Active target communication | 321 |
| 10.2 | Active target communication, with weak synchronization | 322 |
| 10.3 | Passive target communication | 323 |
| 10.4 | Active target communication with several processes | 327 |
| 10.5 | Schematic description of window | 336 |
| 10.6 | Symmetric communication | 339 |
| 10.7 | Deadlock situation | 339 |
| 10.8 | No deadlock | 340 |
| | | |
| 12.1 | Etypes and filetypes | 376 |
| 12.2 | Partitioning a file among parallel processes | 376 |
| 12.3 | Displacements | 388 |
| 12.4 | Example array file layout | 435 |
| 12.5 | Example local array filetype for process 1 | 435 |

List of Tables

| | | |
|------|---|-----|
| 2.1 | Deprecated constructs | 17 |
| 3.1 | Predefined MPI datatypes corresponding to Fortran datatypes | 27 |
| 3.2 | Predefined MPI datatypes corresponding to C datatypes | 27 |
| 5.1 | MPI_COMM_* Function Behavior (in Inter-Communication Mode) | 201 |
| 7.1 | Error classes (Part 1) | 260 |
| 7.2 | Error classes (Part 2) | 261 |
| 10.1 | Error classes in one-sided communication routines | 335 |
| 11.1 | combiner values returned from MPI_TYPE_GET_ENVELOPE | 360 |
| 12.1 | Data access routines | 390 |
| 12.2 | “external32”-sizes of predefined datatypes | 417 |
| 12.3 | Error classes returned from MPI I/O routines. | 433 |
| 13.1 | C++ names for the MPI C and C++ predefined datatypes | 442 |
| 13.2 | C++ names for the MPI Fortran predefined datatypes | 442 |
| 13.3 | C++ names for other MPI datatypes | 443 |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Acknowledgments

This document represents the work of many people who have served on the MPI Forum. The meetings have been attended by dozens of people from many parts of the world. It is the hard and dedicated work of this group that has led to the MPI standard.

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 helped with this effort.

Those who served as primary coordinators in MPI-1.0 and MPI-1.1 are:

- 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-1.0 and MPI-1.1 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 | Steve 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.

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 that have made a contribution towards travel for the U.S. academics. Support for several European participants was provided by ESPRIT.

MPI-1.2 and MPI-2.0:

Those who served as primary coordinators in MPI-1.2 and MPI-2.0 are:

- Ewing Lusk, Convener and Meeting Chair
- Steve Huss-Lederman, Editor
- Ewing Lusk, Miscellany
- Bill Saphir, Process Creation and Management
- Marc Snir, One-Sided Communications
- Bill Gropp and Anthony Skjellum, Extended Collective Operations
- Steve Huss-Lederman, External Interfaces
- Bill Nitzberg, I/O
- Andrew Lumsdaine, Bill Saphir, and Jeff Squyres, Language Bindings
- Anthony Skjellum and Arkady Kanevsky, Real-Time

The following list includes some of the active participants who attended MPI-2 Forum meetings and are not mentioned above.

| | | | |
|----------------|------------------------|-----------------|--------------------|
| Greg Astfalk | Robert Babb | Ed Benson | Rajesh Bordawekar |
| Pete Bradley | Peter Brennan | Ron Brightwell | Maciej Brodowicz |
| Eric Brunner | Greg Burns | Margaret Cahir | Pang Chen |
| Ying Chen | Albert Cheng | Yong Cho | Joel Clark |
| Lyndon Clarke | Laurie Costello | Dennis Cottel | Jim Cownie |
| Zhenqian Cui | Suresh Damodaran-Kamal | | Raja Daoud |
| Judith Devaney | David DiNucci | Doug Doefler | Jack Dongarra |
| Terry Dontje | Nathan Doss | Anne Elster | Mark Fallon |
| Karl Feind | Sam Fineberg | Craig Fischberg | Stephen Fleischman |
| Ian Foster | Hubertus Franke | Richard Frost | Al Geist |
| Robert George | David Greenberg | John Hagedorn | Kei Harada |
| Leslie Hart | Shane Hebert | Rolf Hempel | Tom Henderson |
| Alex Ho | Hans-Christian Hoppe | Joefon Jann | Terry Jones |
| Karl Kesselman | Koichi Konishi | Susan Kraus | Steve Kubica |
| Steve Landherr | Mario Lauria | Mark Law | Juan Leon |
| Lloyd Lewins | Ziyang Lu | Bob Madahar | Peter Madams |

| | | | | |
|----|-------------------|-------------------|-------------------|------------------|
| 1 | John May | Oliver McBryan | Brian McCandless | Tyce McLarty |
| 2 | Thom McMahon | Harish Nag | Nick Nevin | Jarek Nieplocha |
| 3 | Ron Oldfield | Peter Ossadnik | Steve Otto | Peter Pacheco |
| 4 | Yoonho Park | Perry Partow | Pratap Pattnaik | Elsie Pierce |
| 5 | Paul Pierce | Heidi Poxon | Jean-Pierre Prost | Boris Protopopov |
| 6 | James Pruyve | Rolf Rabenseifner | Joe Rieken | Peter Rigsbee |
| 7 | Tom Robey | Anna Rounbehler | Nobutoshi Sagawa | Arindam Saha |
| 8 | Eric Salo | Darren Sanders | Eric Sharakan | Andrew Sherman |
| 9 | Fred Shirley | Lance Shuler | A. Gordon Smith | Ian Stockdale |
| 10 | David Taylor | Stephen Taylor | Greg Tensa | Rajeev Thakur |
| 11 | Marydell Tholburn | Dick Treumann | Simon Tsang | Manuel Ujaldon |
| 12 | David Walker | Jerrell Watts | Klaus Wolf | Parkson Wong |
| 13 | Dave Wright | | | |

14 The MPI Forum also acknowledges and appreciates the valuable input from people via
15 e-mail and in person.

17 The following institutions supported the MPI-2 effort through time and travel support
18 for the people listed above.

20 Argonne National Laboratory
21 Bolt, Beranek, and Newman
22 California Institute of Technology
23 Center for Computing Sciences
24 Convex Computer Corporation
25 Cray Research
26 Digital Equipment Corporation
27 Dolphin Interconnect Solutions, Inc.
28 Edinburgh Parallel Computing Centre
29 General Electric Company
30 German National Research Center for Information Technology
31 Hewlett-Packard
32 Hitachi
33 Hughes Aircraft Company
34 Intel Corporation
35 International Business Machines
36 Khoral Research
37 Lawrence Livermore National Laboratory
38 Los Alamos National Laboratory
39 MPI Software Technology, Inc.
40 Mississippi State University
41 NEC Corporation
42 National Aeronautics and Space Administration
43 National Energy Research Scientific Computing Center
44 National Institute of Standards and Technology
45 National Oceanic and Atmospheric Administration
46 Oak Ridge National Laboratory
47 Ohio State University
48 PALLAS GmbH

| | |
|--|----|
| Pacific Northwest National Laboratory | 1 |
| Pratt & Whitney | 2 |
| San Diego Supercomputer Center | 3 |
| Sanders, A Lockheed-Martin Company | 4 |
| Sandia National Laboratories | 5 |
| Schlumberger | 6 |
| Scientific Computing Associates, Inc. | 7 |
| Silicon Graphics Incorporated | 8 |
| Sky Computers | 9 |
| Sun Microsystems Computer Corporation | 10 |
| Syracuse University | 11 |
| The MITRE Corporation | 12 |
| Thinking Machines Corporation | 13 |
| United States Navy | 14 |
| University of Colorado | 15 |
| University of Denver | 16 |
| University of Houston | 17 |
| University of Illinois | 18 |
| University of Maryland | 19 |
| University of Notre Dame | 20 |
| University of San Fransisco | 21 |
| University of Stuttgart Computing Center | 22 |
| University of Wisconsin | 23 |

MPI-2 operated on a very tight budget (in reality, it had no budget when the first meeting was announced). Many institutions helped the MPI-2 effort by supporting the efforts and travel of the members of the MPI Forum. Direct support was given by NSF and DARPA under NSF contract CDA-9115428 for travel by U.S. academic participants and Esprit under project HPC Standards (21111) for European participants.

MPI-1.3 and MPI-2.1:

The editors and organizers of the combined documents have been:

- Richard Graham, Convener and Meeting Chair
- Jack Dongarra, Steering Committee
- Al Geist, Steering Committee
- Bill Gropp, Steering Committee
- Rainer Keller, Merge of MPI-1.3
- Andrew Lumsdaine, Steering Committee
- Ewing Lusk, Steering Committee, MPI-1.1-Errata (Oct. 12, 1998) MPI-2.1-Errata Bal-
lots 1, 2 (May 15, 2002)

- 1 ● Rolf Rabenseifner, Steering Committee, Merge of MPI-2.1 and MPI-2.1-Errata Ballots
- 2 3, 4 (2008)

3
4 All chapters have been revisited to achieve a consistent MPI-2.1 text. Those who served
5 as authors for the necessary modifications are:

- 6 ● Bill Gropp, Frontmatter, Introduction, and Bibliography
- 7
- 8 ● Richard Graham, Point-to-Point Communications
- 9
- 10 ● Adam Moody, Collective Communication
- 11
- 12 ● Richard Treumann, Groups, Contexts, and Communicators
- 13
- 14 ● Jeper Larsson Traeff, Process Topologies, Miscellany, and One-Sided Communications
- 15
- 16 ● George Bosilca, Environmental Management
- 17
- 18 ● David Solt, Process Creation and Management
- 19
- 20 ● Bronis de Supinski, External Interfaces, and Profiling
- 21
- 22 ● Rajeev Thakur, I/O
- 23
- 24 ● Jeff Squyres, Language Bindings
- 25
- 26 ● Rolf Rabenseifner, Deprecated Functions, and Annex Change-Log
- 27
- 28 ● Alexander Supalov and Denis Nagomy, Annex Language bindings

29 The following list includes some of the active participants who attended MPI-2 Forum
30 meetings and in the e-mail discussions of the errata items and are not mentioned above.

| | | | |
|--------------------|--------------------------|----------------------|---------------------|
| 31 Pavan Balaji | Purushotham V. Bangalore | Brian Barrett | Richard Barrett |
| 32 Christian Bell | Robert Blackmore | Gil Bloch | Ron Brightwell |
| 33 Jeffrey Brown | Darius Buntinas | Jonathan Carter | Nathan DeBardeleben |
| 34 Terry Dontje | Gabor Dozsa | Edric Ellis | Karl Feind |
| 35 Edgar Gabriel | Patrick Geoffray | David Gingold | Dave Goodell |
| 36 Erez Haba | Robert Harrison | Thomas Herault | Steve Hodson |
| 37 Torsten Hoefler | Joshua Hursey | Yann Kalemkarian | Matthew Koop |
| 38 Quincey Koziol | Sameer Kumar | Miron Livny | Kannan Narasimhan |
| 39 Mark Pagel | Avneesh Pant | Steve Poole | Howard Pritchard |
| 40 Craig Rasmussen | Hubert Ritzdorf | Rob Ross | Tony Skjellum |
| Brian Smith | Vinod Tipparaju | Jesper Larsson Traff | Keith Underwood |
| NAMES | from e-mails | TO BE DONE | |

41 The MPI Forum also acknowledges and appreciates the valuable input from people via
42 e-mail and in person.

43
44 The following institutions supported the MPI-2 effort through time and travel support
45 for the people listed above.

46
47 Argonne National Laboratory
48 Bull

| | |
|---|----|
| Cisco Systems, Inc. | 1 |
| Cray Incorporation | 2 |
| HDF Group | 3 |
| Hewlett-Packard | 4 |
| IBM T.J. Watson Research | 5 |
| Indiana University | 6 |
| Indiana University | 7 |
| Institut National de Recherche en Informatique et Automatique (INRIA) | 8 |
| Intel Corporation | 9 |
| Lawrence Berkeley National Laboratory | 10 |
| Lawrence Livermore National Laboratory | 11 |
| Los Alamos National Laboratory | 12 |
| Mathworks | 13 |
| Mellanox Technologies | 14 |
| Microsoft | 15 |
| Myricom | 16 |
| NEC Laboratories Europe, NEC Europe Ltd. | 17 |
| Oak Ridge National Laboratory | 18 |
| Ohio State University | 19 |
| Pacific Northwest National Laboratory | 20 |
| QLogic Corporation | 21 |
| Sandia National Laboratories | 22 |
| SiCortex | 23 |
| Silicon Graphics Incorporation | 24 |
| Sun Microsystem | 25 |
| University of Barcelona (UAB), Dept. of Computer and Information Sciences | 26 |
| University of Houston | 27 |
| University of Illinois at Urbana-Champaign | 28 |
| University of Stuttgart, High Performance Computing Center Stuttgart (HLRS) | 29 |
| University of Tennessee, Knoxville | 30 |
| University of Wisconsin | 31 |

32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 1

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 [40], Express [12], nCUBE’s Vertex [36], p4 [7, 8], and PARMACS [5, 9]. Other important contributions have come from Zipcode [42, 43], Chimp [17, 18], PVM [4, 14], Chameleon [27], and PICL [26].

1.2 Background of MPI-1

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 [51]. 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 [15]. 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

1 issues, but did not include any collective communication routines and was not thread-safe.

2 In November 1992, a meeting of the MPI working group was held in Minneapolis, at
3 which it was decided to place the standardization process on a more formal footing, and to
4 generally adopt the procedures and organization of the High Performance Fortran Forum.
5 Subcommittees were formed for the major component areas of the standard, and an email
6 discussion service established for each. In addition, the goal of producing a draft MPI
7 standard by the Fall of 1993 was set. To achieve this goal the MPI working group met every
8 6 weeks for two days throughout the first 9 months of 1993, and presented the draft MPI
9 standard at the Supercomputing 93 conference in November 1993. These meetings and the
10 email discussion together constituted the MPI Forum, membership of which has been open
11 to all members of the high performance computing community.

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

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

22 A complete list of goals follows.

- 23
- 24 ● Design an application programming interface (not necessarily for compilers or a system
25 implementation library).
- 26
- 27 ● Allow efficient communication: Avoid memory-to-memory copying and allow overlap
28 of computation and communication and offload to communication co-processor, where
29 available.
- 30
- 31 ● Allow for implementations that can be used in a heterogeneous environment.
- 32
- 33 ● Allow convenient C and Fortran 77 bindings for the interface.
- 34
- 35 ● Assume a reliable communication interface: the user need not cope with communica-
36 tion failures. Such failures are dealt with by the underlying communication subsystem.
- 37
- 38 ● Define an interface that is not too different from current practice, such as PVM, NX,
39 Express, p4, etc., and provides extensions that allow greater flexibility.
- 40
- 41 ● Define an interface that can be implemented on many vendor's platforms, with no
42 significant changes in the underlying communication and system software.
- 43
- 44 ● Semantics of the interface should be language independent.
- 45
- 46 ● The interface should be designed to allow for thread-safety.
- 47
- 48

1.3 Background of MPI-2

Beginning in March 1995, the MPI Forum began meeting to consider corrections and extensions to the original MPI Standard document [22]. The first product of these deliberations was Version 1.1 of the MPI specification, released in June of 1995 (see <http://www.mpi-forum.org> for official MPI document releases). Since that time, effort has been focused in five areas.

1. Further corrections and clarifications for the MPI-1.1 document.
2. Additions to MPI-1.1 that do not significantly change its types of functionality (new datatype constructors, language interoperability, etc.).
3. Completely new types of functionality (dynamic processes, one-sided communication, parallel I/O, etc.) that are what everyone thinks of as “MPI-2 functionality.”
4. Bindings for Fortran 90 and C++. This document specifies C++ bindings for both MPI-1 and MPI-2 functions, and extensions to the Fortran 77 binding of MPI-1 and MPI-2 to handle Fortran 90 issues.
5. Discussions of areas in which the MPI process and framework seem likely to be useful, but where more discussion and experience are needed before standardization (e.g. 0-copy semantics on shared-memory machines, real-time specifications).

Corrections and clarifications (items of type 1 in the above list) have been collected in Chapter 3 of the MPI-2 document: “Version 1.2 of MPI.” This chapter also contains the function for identifying the version number. Additions to MPI-1.1 (items of types 2, 3, and 4 in the above list) are in the remaining chapters of the MPI-2 document, and constitute the specification for MPI-2. Items of type 5 in the above list have been moved to a separate document, the “MPI Journal of Development” (JOD), and are not part of the MPI-2 Standard.

This structure makes it easy for users and implementors to understand what level of MPI compliance a given implementation has:

- MPI-1 compliance will mean compliance with MPI-1.3. This is a useful level of compliance. It means that the implementation conforms to the clarifications of MPI-1.1 function behavior given in Chapter 3 of the MPI-2 document. Some implementations may require changes to be MPI-1 compliant.
- MPI-2 compliance will mean compliance with all of MPI-2.1.
- The MPI Journal of Development is not part of the MPI Standard.

It is to be emphasized that forward compatibility is preserved. That is, a valid MPI-1.1 program is both a valid MPI-1.2 program and a valid MPI-2 program, and a valid MPI-1.2 program is a valid MPI-2 program.

1.4 Background of MPI-1.3 and MPI-2.1

After the release of MPI-2.0, the MPI Forum kept working on erratas and clarifications for both standard documents (MPI-1.1 and MPI-2.0). The short document “Errata for MPI-1.1”

1 was released October 12, 1998. July 5, 2001, a first ballot of erratas and clarifications for
2 MPI-2.0 was released, and a second ballot was voted on May 22, 2002. Both votes were done
3 electronically. Both ballots were combined into one document: “Errata for MPI-2”, May 15,
4 2002. This errata process was then interrupted, but the Forum and its e-mail reflectors
5 kept working on new requests for clarification.

6 Restarting regular work of the MPI Forum was initiated in three meetings, at Eu-
7 roPVM/MPI’06 in Bonn, at EuroPVM/MPI’07 in Paris, and at SC’07 in Reno. In Dec.
8 2007, a steering committee started the organization of new MPI Forum meetings at regular
9 8-weeks intervals. At the Jan. 14-16, 2008 meeting in Chicago, the MPI Forum decided to
10 combine the existing and future MPI documents to one single document for each version of
11 the MPI standard. For technical and historical reasons, this series was started with MPI-1.3.
12 Additional Ballots 3 and 4 are solving old questions from the errata list started in 1995 upto
13 new questions from the last years.

14 After all documents (MPI-1.1, MPI-2, Errata for MPI-1.1 (Oct. 12, 1998), and MPI-2.1
15 Ballots 1-4) were combined into one draft document, for each chapter, a chapter author and
16 review team were defined. They cleaned up the document to achieve a consistent MPI-2.1
17 document.

18 It is expected that the final MPI-2.1 standard document is finished in June 2008, and
19 finally released with a second vote in September 2008 in the meeting at Dublin, straight
20 before EuroPVM/MPI’08. The major work of the current MPI Forum is the preparation of
21 MPI-3.

22 1.5 Who Should Use This Standard?

23 This standard is intended for use by all those who want to write portable message-passing
24 programs in [Fortran](#), [C](#) and [C++](#). This includes individual application programmers, de-
25 velopers of software designed to run on parallel machines, and creators of environments
26 and tools. In order to be attractive to this wide audience, the standard must provide a
27 simple, easy-to-use interface for the basic user while not semantically precluding the high-
28 performance message-passing operations available on advanced machines.
29
30
31

32 1.6 What Platforms Are Targets For Implementation?

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

41 The interface is suitable for use by fully general MIMD programs, as well as those writ-
42 ten in the more restricted style of [SPMD](#). MPI provides many features intended to improve
43 performance on scalable parallel computers with specialized interprocessor communication
44 hardware. Thus, we expect that native, high-performance implementations of MPI will be
45 provided on such machines. At the same time, implementations of MPI on top of stan-
46 dard Unix interprocessor communication protocols will provide portability to workstation
47
48

clusters and heterogenous networks of workstations. Several proprietary, native implemen-
tations of MPI, and a public domain, portable implementation of MPI are in progress at the
time of this writing [24, 16].

1.7 What Is Included In The Standard?

The standard includes:

- Point-to-point communication
- Collective operations
- Process groups
- Communication contexts
- Process topologies
- Environmental Management and inquiry
- A Miscellany chapter
- Process creation and management
- One-sided communication
- External interfaces
- Parallel file I/O
- Language Bindings for Fortran, C and C++
- Profiling interface

1.8 What Is Not Included In The Standard?

The standard does not specify:

- 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

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.

1.9 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 Communications, 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). **With MPI-2, the semantics of collective communication was extended to include intercommunicators. It also adds more convenient methods of constructing intercommunicators and two new collective operations.**
- 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.

The following Chapters 8-13 have been added with MPI-2:

- Chapter 8, Miscellany, discusses items that don't fit elsewhere.
- Chapter 9, Process Creation and Management, defines routines that allow for creation of processes.
- Chapter 10, One-Sided Communications, defines communication routines that can be completed by a single process. These include shared-memory operations (*put/get*) and remote accumulate operations.
- Chapter 11, External Interfaces, defines routines designed to allow developers to layer on top of MPI. This includes generalized requests, routines that decode MPI opaque objects, and threads.
- Chapter 12, I/O, defines MPI-2 support for parallel I/O.
- Chapter 13, Language Bindings, describes the C++ binding and discusses Fortran-90 issues.

- Chapter 14, 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.

The Appendices are:

- Annex A, Language Bindings, gives specific syntax in Fortran, C, and C++, for all MPI functions, constants, and types.
- Annex chap:change, Change-Log, summarizes major changes since the previous version of the standard.
- The MPI Function Index is a simple index showing the location of the precise definition of each MPI function, together with C, C++, and Fortran bindings.

MPI-2 provides various interfaces to facilitate interoperability of distinct MPI implementations. Among these are the canonical data representation for MPI I/O and for MPI_PACK_EXTERNAL and MPI_UNPACK_EXTERNAL. The definition of an actual binding of these interfaces that will enable interoperability is outside the scope of this document.

A separate document consists of ideas that were discussed in the MPI Forum and deemed to have value, but are not included in the MPI Standard. They are part of the “Journal of Development” (JOD), lest good ideas be lost and in order to provide a starting point for further work. The chapters in the JOD are

- Chapter 2, Spawning Independent Processes, includes some elements of dynamic process management, in particular management of processes with which the spawning processes do not intend to communicate, that the Forum discussed at length but ultimately decided not to include in the MPI Standard.
- Chapter 3, Threads and MPI, describes some of the expected interaction between an MPI implementation and a thread library in a multi-threaded environment.
- Chapter 4, Communicator ID, describes an approach to providing identifiers for communicators.
- Chapter 5, Miscellany, discusses Miscellaneous topics in the MPI JOD, in particular single-copy routines for use in shared-memory environments and new datatype constructors.
- Chapter 6, Toward a Full Fortran 90 Interface, describes an approach to providing a more elaborate Fortran 90 interface.
- Chapter 7, Split Collective Communication, describes a specification for certain non-blocking collective operations.
- Chapter 8, Real-Time MPI, discusses MPI support for real time processing.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 2

MPI Terms and Conventions

This chapter explains notational terms and conventions used throughout the MPI-2 document, some of the choices that have been made, and the rationale behind those choices. It is similar to the MPI-1 Terms and Conventions chapter but differs in some major and minor ways. Some of the major areas of difference are the naming conventions, some semantic definitions, file objects, Fortran 90 *vs* Fortran 77, C++, processes, and interaction with signals.

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 aimed at users and that 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 Naming Conventions

MPI-1 used informal naming conventions. In many cases, MPI-1 names for C functions are of the form `Class.action_subset` and in Fortran of the form `CLASS_ACTION_SUBSET`, but this rule is not uniformly applied. In MPI-2, an attempt has been made to standardize names of new functions according to the following rules. In addition, the C++ bindings for MPI-1 functions also follow these rules (see Section 2.6.4). C and Fortran function names for MPI-1 have not been changed.

1. In C, all routines associated with a particular type of MPI object should be of the form `Class.action_subset` or, if no subset exists, of the form `Class.action`. In Fortran,

1 all routines associated with a particular type of MPI object should be of the form
 2 `CLASS_ACTION_SUBSET` or, if no subset exists, of the form `CLASS_ACTION`. For C
 3 and Fortran we use the C++ terminology to define the `Class`. In C++, the routine
 4 is a method on `Class` and is named `MPI::Class::Action_subset`. If the routine is
 5 associated with a certain class, but does not make sense as an object method, it is a
 6 static member function of the class.

- 7
- 8 2. If the routine is not associated with a class, the name should be of the form
 9 `Action_subset` in C and `ACTION_SUBSET` in Fortran, and in C++ should be scoped
 10 in the `MPI` namespace, `MPI::Action_subset`.
- 11 3. The names of certain actions have been standardized. In particular, `Create` creates
 12 a new object, `Get` retrieves information about an object, `Set` sets this information,
 13 `Delete` deletes information, `Is` asks whether or not an object has a certain property.
 14

15 C and Fortran names for MPI-1 functions violate these rules in several cases. The most
 16 common exceptions are the omission of the `Class` name from the routine and the omission
 17 of the `Action` where one can be inferred.

18 MPI identifiers are limited to 30 characters (31 with the profiling interface). This is
 19 done to avoid exceeding the limit on some compilation systems.
 20

21 2.3 Procedure Specification

22 MPI procedures are specified using a language-independent notation. The arguments of
 23 procedure calls are marked as IN, OUT or INOUT. The meanings of these are:
 24

- 25
- 26 • the call may use the input value but does not update an argument is marked IN,
- 27
- 28 • the call may update an argument but does not use its input value is marked OUT,
- 29
- 30 • the call may both use and update an argument is marked INOUT.

31 There is one special case — if an argument is a handle to an opaque object (these
 32 terms are defined in Section 2.5.1), and the object is updated by the procedure call, then
 33 the argument is marked `INOUT` or `OUT`. It is marked this way even though the handle itself
 34 is not modified — we use the `INOUT` or `OUT` attribute to denote that what the handle
 35 *references* is updated. Thus, in C++, IN arguments are either references or pointers to
 36 `const` objects.
 37

38 *Rationale.* The definition of MPI tries to avoid, to the largest possible extent, the use
 39 of INOUT arguments, because such use is error-prone, especially for scalar arguments.
 40 (*End of rationale.*)
 41

42 MPI's use of IN, OUT and INOUT is intended to indicate to the user how an argument
 43 is to be used, but does not provide a rigorous classification that can be translated directly
 44 into all language bindings (e.g., `INTENT` in Fortran 90 bindings or `const` in C bindings).
 45 For instance, the “constant” `MPI_BOTTOM` can usually be passed to OUT buffer arguments.
 46 Similarly, `MPI_STATUS_IGNORE` can be passed as the OUT status argument.

47 A common occurrence for MPI functions is an argument that is used as
 48 IN by some processes and OUT by other processes. Such an 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 on a single process.

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 an 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 ISO C version of the function is shown followed by a version of the same function in Fortran and then the C++ binding. Fortran in this document refers to Fortran 90; see Section 2.6.

2.4 Semantic Terms

When discussing MPI procedures the following semantic terms are used.

nonblocking A procedure is nonblocking if the procedure may return before the operation completes, and before the user is allowed to reuse resources (such as buffers) specified in the call. A nonblocking request is **started** by the call that initiates it, e.g., MPI_ISEND. The word complete is used with respect to operations, requests, and communications. An **operation completes** when the user is allowed to reuse resources, and any output buffers have been updated; i.e. a call to MPI_TEST will return `flag = true`. A **request is completed** by a call to wait, which returns, or a test or get status call which returns `flag = true`. This completing call has two effects: the status is extracted from the request; in the case of test and wait, if the request was nonpersistent, it is **freed**. A **communication completes** when all participating operations complete.

blocking A procedure is blocking if return from the procedure indicates the user is allowed to reuse resources specified in the call.

local A procedure is local if completion of the procedure depends only on the local executing process.

non-local A procedure is 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.

1 **collective** A procedure is collective if all processes in a process group need to invoke the
2 procedure. A collective call may or may not be synchronizing. Collective calls over
3 the same communicator must be executed in the same order by all members of the
4 process group.

5
6 **predefined** A predefined datatype is a datatype with a predefined (constant) name (such
7 as MPI_INT, MPI_FLOAT_INT, or MPI_UB) or a datatype constructed with
8 MPI_TYPE_CREATE_F90_INTEGER, MPI_TYPE_CREATE_F90_REAL, or
9 MPI_TYPE_CREATE_F90_COMPLEX. The former are **named** whereas the latter are
10 **unnamed**.

11 **derived** A derived datatype is any datatype that is not predefined.

12
13 **portable** A datatype is portable, if it is a predefined datatype, or it is derived from a
14 portable datatype using only the type constructors MPI_TYPE_CONTIGUOUS,
15 MPI_TYPE_VECTOR, MPI_TYPE_INDEXED, MPI_TYPE_INDEXED_BLOCK,
16 MPI_TYPE_CREATE_SUBARRAY, MPI_TYPE_DUP, and MPI_TYPE_CREATE_DARRAY.
17 Such a datatype is portable because all displacements in the datatype are in terms of
18 extents of one predefined datatype. Therefore, if such a datatype fits a data layout
19 in one memory, it will fit the corresponding data layout in another memory, if the
20 same declarations were used, even if the two systems have different architectures. On
21 the other hand, if a datatype was constructed using MPI_TYPE_CREATE_HINDEXED,
22 MPI_TYPE_CREATE_HVECTOR or MPI_TYPE_CREATE_STRUCT, then the datatype
23 contains explicit byte displacements (e.g., providing padding to meet alignment re-
24 strictions). These displacements are unlikely to be chosen correctly if they fit data
25 layout on one memory, but are used for data layouts on another process, running on
26 a processor with a different architecture.

27
28 **equivalent** Two datatypes are equivalent if they appear to have been created with the same
29 sequence of calls (and arguments) and thus have the same typemap. Two equivalent
30 datatypes do not necessarily have the same cached attributes or the same names.

31 2.5 Data Types

32 2.5.1 Opaque Objects

33
34 MPI manages **system memory** that is used for buffering messages and for storing internal
35 representations of various MPI objects such as groups, communicators, datatypes, etc. This
36 memory is not directly accessible to the user, and objects stored there are **opaque**: their
37 size and shape is not visible to the user. Opaque objects are accessed via **handles**, which
38 exist in user space. MPI procedures that operate on opaque objects are passed handle
39 arguments to access these objects. In addition to their use by MPI calls for object access,
40 handles can participate in assignments and comparisons.

41
42 In Fortran, all handles have type **INTEGER**. In C and C++, a different handle type is
43 defined for each category of objects. In addition, handles themselves are distinct objects
44 in C++. The C and C++ types must support the use of the assignment and equality
45 operators.

46
47 *Advice to implementors.* In Fortran, the handle can be an index into a table of
48

opaque objects in a system table; in C it can be such an index or a pointer to the object. C++ handles can simply “wrap up” a table index or pointer.

(End of advice to implementors.)

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 a deallocate routine 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. The user must not free such objects. In C++, this is enforced by declaring the handles to these predefined objects to be `static const`.

Rationale. This design hides the internal representation used for MPI data structures, thus allowing similar calls in C, 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 separation of handles in user space and objects in system space allows space-reclaiming and 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 MPI calls that allocate or deallocate such objects. *(End of advice to users.)*

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

11 2.5.2 Array Arguments

12 An MPI call may need an argument that is an array of opaque objects, or an array of
13 handles. The array-of-handles is a regular array with entries that are handles to objects
14 of the same type in consecutive locations in the array. Whenever such an array is used,
15 an additional `len` argument is required to indicate the number of valid entries (unless this
16 number can be derived otherwise). The valid entries are at the beginning of the array;
17 `len` indicates how many of them there are, and need not be the size of the entire array.
18 The same approach is followed for other array arguments. In some cases `NULL` handles are
19 considered valid entries. When a `NULL` argument is desired for an array of statuses, one
20 uses `MPI_STATUSES_IGNORE`.

22 2.5.3 State

23 MPI procedures use at various places arguments with *state* types. The values of such a data
24 type are all identified by names, and no operation is defined on them. For example, the
25 `MPI_TYPE_CREATE_SUBARRAY` routine has a state argument `order` with values
26 `MPI_ORDER_C` and `MPI_ORDER_FORTRAN`.

29 2.5.4 Named Constants

30 MPI procedures sometimes assign a special meaning to a special value of a basic type argu-
31 ment; e.g., `tag` is an integer-valued argument of point-to-point communication operations,
32 with a special wild-card value, `MPI_ANY_TAG`. Such arguments will have a range of regular
33 values, which is a proper subrange of the range of values of the corresponding basic type;
34 special values (such as `MPI_ANY_TAG`) will be outside the regular range. The range of regular
35 values, such as `tag`, can be queried using environmental inquiry functions (Chapter 7 of the
36 MPI-1 document). The range of other values, such as `source`, depends on values given by
37 other MPI routines (in the case of `source` it is the communicator size).

38 MPI also provides predefined named constant handles, such as `MPI_COMM_WORLD`.

39 All named constants, with the exceptions noted below for Fortran, can be used in
40 initialization expressions or assignments. These constants do not change values during
41 execution. Opaque objects accessed by constant handles are defined and do not change
42 value between MPI initialization (`MPI_INIT`) and MPI completion (`MPI_FINALIZE`).

43 The constants that cannot be used in initialization expressions or assignments in For-
44 tran are:

45
46 `MPI_BOTTOM`
47 `MPI_STATUS_IGNORE`
48

MPI_STATUSES_IGNORE
 MPI_ERRCODES_IGNORE
 MPI_IN_PLACE
 MPI_ARGV_NULL
 MPI_ARGVS_NULL

Advice to implementors. In Fortran the implementation of these special constants may require the use of language constructs that are outside the Fortran standard. Using special values for the constants (e.g., by defining them through `parameter` statements) is not possible because an implementation cannot distinguish these values from legal data. Typically, these constants are implemented as predefined static variables (e.g., a variable in an MPI-declared `COMMON` block), relying on the fact that the target compiler passes data by address. Inside the subroutine, this address can be extracted by some mechanism outside the Fortran standard (e.g., by Fortran extensions or by implementing the function in C). (*End of advice to implementors.*)

2.5.5 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 and C++, we use `void *`.

2.5.6 Addresses

Some MPI procedures use *address* arguments that represent an absolute address in the calling program. The datatype of such an argument is `MPI_Aint` in C, `MPI::Aint` in C++ and `INTEGER (KIND=MPI_ADDRESS_KIND)` in Fortran. There is the MPI constant `MPI_BOTTOM` to indicate the start of the address range.

2.5.7 File Offsets

For I/O there is a need to give the size, displacement, and offset into a file. These quantities can easily be larger than 32 bits which can be the default size of a Fortran integer. To overcome this, these quantities are declared to be `INTEGER (KIND=MPI_OFFSET_KIND)` in Fortran. In C one uses `MPI_Offset` whereas in C++ one uses `MPI::Offset`.

2.6 Language Binding

This section defines the rules for MPI language binding in general and for Fortran, ISO C, and C++, in particular. (Note that ANSI C has been replaced by ISO C.) Defined here are various object representations, as well as the naming conventions used for expressing this standard. The actual calling sequences are defined elsewhere.

MPI bindings are for Fortran 90, though they are designed to be usable in Fortran 77 environments.

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 and C++, however, we expect that C and C++ programmers will

1 understand the word “argument” (which has no specific meaning in C/C++), thus allowing
2 us to avoid unnecessary confusion for Fortran programmers.

3 Since Fortran is case insensitive, linkers may use either lower case or upper case when
4 resolving Fortran names. Users of case sensitive languages should avoid the “mpi_” and
5 “pmpi_” prefixes.

7 2.6.1 Deprecated Names and Functions

8 A number of chapters refer to deprecated or replaced MPI-1 constructs. These are constructs
9 that continue to be part of the MPI standard, but that users are recommended not to
10 continue using, since MPI-2 provides better solutions. For example, the Fortran binding for
11 MPI-1 functions that have address arguments uses `INTEGER`. This is not consistent with the
12 C binding, and causes problems on machines with 32 bit `INTEGERS` and 64 bit addresses.
13 In MPI-2, these functions have new names, and new bindings for the address arguments.
14 The use of the old functions is deprecated. For consistency, here and in a few other cases,
15 new C functions are also provided, even though the new functions are equivalent to the
16 old functions. The old names are deprecated. Another example is provided by the MPI-1
17 predefined datatypes `MPI_UB` and `MPI_LB`. They are deprecated, since their use is awkward
18 and error-prone, while the MPI-2 function `MPI_TYPE_CREATE_RESIZED` provides a more
19 convenient mechanism to achieve the same effect.

20 [Table 2.1](#) shows a list of all of the deprecated constructs. Note that the constants
21 `MPI_LB` and `MPI_UB` are replaced by the function `MPI_TYPE_CREATE_RESIZED`; this is
22 because their principle use was as input datatypes to `MPI_TYPE_STRUCT` to create resized
23 datatypes. Also note that some C typedefs and Fortran subroutine names are included in
24 this list; they are the types of callback functions.

27 2.6.2 Fortran Binding Issues

28 MPI-1.1 provided bindings for Fortran 77. MPI-2 retains these bindings but they are now
29 interpreted in the context of the Fortran 90 standard. MPI can still be used with most
30 Fortran 77 compilers, as noted below. When the term Fortran is used it means Fortran 90.

31 All MPI names have an `MPI_` prefix, and all characters are capitals. Programs must not
32 declare variables, parameters, or functions with names beginning with the prefix `MPI_`. To
33 avoid conflicting with the profiling interface, programs should also avoid functions with the
34 prefix `PMPI_`. This is mandated to avoid possible name collisions.

35 All MPI Fortran subroutines have a return code in the last argument. A few MPI
36 operations which are functions do not have the return code argument. The return code value
37 for successful completion is `MPI_SUCCESS`. Other error codes are implementation dependent;
38 see the error codes in Chapter 7 of the MPI-1 document and Annex A in the MPI-2 document.

39 Constants representing the maximum length of a string are one smaller in Fortran than
40 in C and C++ as discussed in Section 13.3.9.

41 Handles are represented in Fortran as `INTEGERS`. Binary-valued variables are of type
42 `LOGICAL`.

43 Array arguments are indexed from one.

44 The MPI Fortran binding is inconsistent with the Fortran 90 standard in several re-
45 spects. These inconsistencies, such as register optimization problems, have implications for
46 user codes that are discussed in detail in Section 13.2.2. They are also inconsistent with
47 Fortran 77.

| Deprecated | MPI-2 Replacement | |
|-----------------------|-------------------------------|----|
| MPI_ADDRESS | MPI_GET_ADDRESS | 1 |
| MPI_TYPE_HINDEXED | MPI_TYPE_CREATE_HINDEXED | 2 |
| MPI_TYPE_HVECTOR | MPI_TYPE_CREATE_HVECTOR | 3 |
| MPI_TYPE_STRUCT | MPI_TYPE_CREATE_STRUCT | 4 |
| MPI_TYPE_EXTENT | MPI_TYPE_GET_EXTENT | 5 |
| MPI_TYPE_UB | MPI_TYPE_GET_EXTENT | 6 |
| MPI_TYPE_LB | MPI_TYPE_GET_EXTENT | 7 |
| MPI_LB | MPI_TYPE_CREATE_RESIZED | 8 |
| MPI_UB | MPI_TYPE_CREATE_RESIZED | 9 |
| MPI_ERRHANDLER_CREATE | MPI_COMM_CREATE_ERRHANDLER | 10 |
| MPI_ERRHANDLER_GET | MPI_COMM_GET_ERRHANDLER | 11 |
| MPI_ERRHANDLER_SET | MPI_COMM_SET_ERRHANDLER | 12 |
| MPI_Handler_function | MPI_Comm_errhandler_fn | 13 |
| MPI_KEYVAL_CREATE | MPI_COMM_CREATE_KEYVAL | 14 |
| MPI_KEYVAL_FREE | MPI_COMM_FREE_KEYVAL | 15 |
| MPI_DUP_FN | MPI_COMM_DUP_FN | 16 |
| MPI_NULL_COPY_FN | MPI_COMM_NULL_COPY_FN | 17 |
| MPI_NULL_DELETE_FN | MPI_COMM_NULL_DELETE_FN | 18 |
| MPI_Copy_function | MPI_Comm_copy_attr_function | 19 |
| COPY_FUNCTION | COMM_COPY_ATTR_FN | 20 |
| MPI_Delete_function | MPI_Comm_delete_attr_function | 21 |
| DELETE_FUNCTION | COMM_DELETE_ATTR_FN | 22 |
| MPI_ATTR_DELETE | MPI_COMM_DELETE_ATTR | 23 |
| MPI_ATTR_GET | MPI_COMM_GET_ATTR | 24 |
| MPI_ATTR_PUT | MPI_COMM_SET_ATTR | 25 |

Table 2.1: [Deprecated constructs](#)

- An MPI subroutine with a choice argument may be called with different argument types.
- An MPI subroutine with an assumed-size dummy argument may be passed an actual scalar argument.
- Many MPI routines assume that actual arguments are passed by address and that arguments are not copied on entrance to or exit from the subroutine.
- An MPI implementation may read or modify user data (e.g., communication buffers used by nonblocking communications) concurrently with a user program executing outside MPI calls.
- Several named “constants,” such as `MPI_BOTTOM`, `MPI_STATUS_IGNORE`, and `MPI_ERRCODES_IGNORE`, are not ordinary Fortran constants and require a special implementation. See Section 2.5.4 on page 14 for more information.

Additionally, MPI is inconsistent with Fortran 77 in a number of ways, as noted below.

- MPI identifiers exceed 6 characters.

- 1 • MPI identifiers may contain underscores after the first character.
- 2
- 3 • MPI requires an include file, `mpif.h`. On systems that do not support include files,
- 4 the implementation should specify the values of named constants.
- 5
- 6 • Many routines in MPI-2 have KIND-parameterized integers (e.g., `MPI_ADDRESS_KIND`
- 7 and `MPI_OFFSET_KIND`) that hold address information. On systems that do not sup-
- 8 port Fortran 90-style parameterized types, `INTEGER*8` or `INTEGER` should be used
- 9 instead.
- 10 • The memory allocation routine `MPI_ALLOC_MEM` can't be usefully used in Fortran
- 11 without a language extension that allows the allocated memory to be associated with
- 12 a Fortran variable.
- 13

14 2.6.3 C Binding Issues

15 We use the ISO C declaration format. All MPI names have an `MPI_` prefix, defined constants
 16 are in all capital letters, and defined types and functions have one capital letter after the
 17 prefix. Programs must not declare variables or functions with names beginning with the
 18 prefix `MPI_`. To support the profiling interface, programs should not declare functions with
 19 names beginning with the prefix `PMPI_`.

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

22 Almost all C functions return an error code. The successful return code will be
 23 `MPI_SUCCESS`, but failure return codes are implementation dependent.

24 Type declarations are provided for handles to each category of opaque objects.

25 Array arguments are indexed from zero.

26 Logical flags are integers with value 0 meaning “false” and a non-zero value meaning
 27 “true.”

28 Choice arguments are pointers of type `void *`.

29 Address arguments are of MPI defined type `MPI_Aint`. File displacements are of type
 30 `MPI_Offset`. `MPI_Aint` is defined to be an integer of the size needed to hold any valid address
 31 on the target architecture. `MPI_Offset` is defined to be an integer of the size needed to hold
 32 any valid file size on the target architecture.

33 2.6.4 C++ Binding Issues

34 There are places in the standard that give rules for C and not for C++. In these cases,
 35 the C rule should be applied to the C++ case, as appropriate. In particular, the values of
 36 constants given in the text are the ones for C and Fortran. A cross index of these with the
 37 C++ names is given in Annex A.

38 We use the ISO C++ declaration format. All MPI names are declared within the scope
 39 of a namespace called `MPI` and therefore are referenced with an `MPI::` prefix. Defined
 40 constants are in all capital letters, and class names, defined types, and functions have only
 41 their first letter capitalized. Programs must not declare variables or functions in the `MPI`
 42 namespace. This is mandated to avoid possible name collisions.

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

Advice to implementors. The file `mpi.h` may contain both the C and C++ definitions. Usually one can simply use the defined value (generally `__cplusplus`, but not required) to see if one is using C++ to protect the C++ definitions. It is possible that a C compiler will require that the source protected this way be legal C code. In this case, all the C++ definitions can be placed in a different include file and the “`#include`” directive can be used to include the necessary C++ definitions in the `mpi.h` file. (*End of advice to implementors.*)

C++ functions that create objects or return information usually place the object or information in the return value. Since the language neutral prototypes of MPI functions include the C++ return value as an OUT parameter, semantic descriptions of MPI functions refer to the C++ return value by that parameter name (see Section A.4.19 on page 556). The remaining C++ functions return `void`.

In some circumstances, MPI permits users to indicate that they do not want a return value. For example, the user may indicate that the status is not filled in. Unlike C and Fortran where this is achieved through a special input value, in C++ this is done by having two bindings where one has the optional argument and one does not.

C++ functions do not return error codes. If the default error handler has been set to `MPI::ERRORS_THROW_EXCEPTIONS`, the C++ exception mechanism is used to signal an error by throwing an `MPI::Exception` object.

It should be noted that the default error handler (i.e., `MPI::ERRORS_ARE_FATAL`) on a given type has not changed. User error handlers are also permitted. `MPI::ERRORS_RETURN` simply returns control to the calling function; there is no provision for the user to retrieve the error code.

User callback functions that return integer error codes should not throw exceptions; the returned error will be handled by the MPI implementation by invoking the appropriate error handler.

Advice to users. C++ programmers that want to handle MPI errors on their own should use the `MPI::ERRORS_THROW_EXCEPTIONS` error handler, rather than `MPI::ERRORS_RETURN`, that is used for that purpose in C. Care should be taken using exceptions in mixed language situations. (*End of advice to users.*)

Opaque object handles must be objects in themselves, and have the assignment and equality operators overridden to perform semantically like their C and Fortran counterparts.

Array arguments are indexed from zero.

Logical flags are of type `bool`.

Choice arguments are pointers of type `void *`.

Address arguments are of MPI-defined integer type `MPI::Aint`, defined to be an integer of the size needed to hold any valid address on the target architecture. Analogously, `MPI::Offset` is an integer to hold file offsets.

Most MPI functions are methods of MPI C++ classes. MPI class names are generated from the language neutral MPI types by dropping the `MPI_` prefix and scoping the type within the MPI namespace. For example, `MPI_DATATYPE` becomes `MPI::Datatype`.

The names of MPI-2 functions generally follow the naming rules given. In some circumstances, the new MPI-2 function is related to an MPI-1 function with a name that does not follow the naming conventions. In this circumstance, the language neutral name is in analogy to the MPI-1 name even though this gives an MPI-2 name that violates the naming

conventions. The C and Fortran names are the same as the language neutral name in this case. However, the C++ names for MPI-1 do reflect the naming rules and can differ from the C and Fortran names. Thus, the analogous name in C++ to the MPI-1 name is different than the language neutral name. This results in the C++ name differing from the language neutral name. An example of this is the language neutral name of `MPI_FINALIZED` and a C++ name of `MPI::Is_finalized`.

In C++, function `typedefs` are made publicly within appropriate classes. However, these declarations then become somewhat cumbersome, as with the following:

```
typedef MPI::Grequest::Query_function();
```

would look like the following:

```
namespace MPI {
  class Request {
    // ...
  };

  class Grequest : public MPI::Request {
    // ...
    typedef Query_function(void* extra_state, MPI::Status& status);
  };
};
```

Rather than including this scaffolding when declaring C++ `typedefs`, we use an abbreviated form. In particular, we explicitly indicate the class and namespace scope for the `typedef` of the function. Thus, the example above is shown in the text as follows:

```
typedef int MPI::Grequest::Query_function(void* extra_state,
                                           MPI::Status& status)
```

The C++ bindings presented in Annex A.4 and throughout this document were generated by applying a simple set of name generation rules to the MPI function specifications. While these guidelines may be sufficient in most cases, they may not be suitable for all situations. In cases of ambiguity or where a specific semantic statement is desired, these guidelines may be superseded as the situation dictates.

1. All functions, types, and constants are declared within the scope of a `namespace` called `MPI`.
2. Arrays of MPI handles are always left in the argument list (whether they are IN or OUT arguments).
3. If the argument list of an MPI function contains a scalar IN handle, and it makes sense to define the function as a method of the object corresponding to that handle, the function is made a member function of the corresponding MPI class. The member functions are named according to the corresponding MPI function name, but without the “MPI_” prefix and without the object name prefix (if applicable). In addition:
 - (a) The scalar IN handle is dropped from the argument list, and `this` corresponds to the dropped argument.

- (b) The function is declared `const`.
- 4. MPI functions are made into class functions (static) when they belong on a class but do not have a unique scalar IN or INOUT parameter of that class.
- 5. If the argument list contains a single OUT argument that is not of type `MPI_STATUS` (or an array), that argument is dropped from the list and the function returns that value.

Example 2.1 The C++ binding for `MPI_COMM_SIZE` is `int MPI::Comm::Get_size(void) const`.

- 6. If there are multiple OUT arguments in the argument list, one is chosen as the return value and is removed from the list.
- 7. If the argument list does not contain any OUT arguments, the function returns `void`.

Example 2.2 The C++ binding for `MPI_REQUEST_FREE` is `void MPI::Request::Free(void)`

- 8. MPI functions to which the above rules do not apply are not members of any class, but are defined in the MPI namespace.

Example 2.3 The C++ binding for `MPI_BUFFER_ATTACH` is `void MPI::Attach_buffer(void* buffer, int size)`.

- 9. All class names, defined types, and function names have only their first letter capitalized. Defined constants are in all capital letters.
- 10. Any IN pointer, reference, or array argument must be declared `const`.
- 11. Handles are passed by reference.
- 12. Array arguments are denoted with square brackets (`[]`), not pointers, as this is more semantically precise.

2.6.5 Functions and Macros

An implementation is allowed to implement `MPI_WTIME`, `MPI_WTICK`, `PMPI_WTIME`, `PMPI_WTICK`, and the handle-conversion functions (`MPI_Group_f2c`, etc.) in Section 13.3.4, and no others, as macros in C.

Advice to implementors. Implementors should document which routines are implemented as macros. (*End of advice to implementors.*)

Advice to users. If these routines are implemented as macros, they will not work with the MPI profiling interface. (*End of advice to users.*)

2.7 Processes

An MPI program consists of autonomous processes, executing their own code, in a 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. The interaction of an MPI program with other possible means of communication, I/O, and process management is not specified. Unless otherwise stated in the specification of the standard, MPI places no requirements on the result of its interaction with external mechanisms that provide similar or equivalent functionality. This includes, but is not limited to, interactions with external mechanisms for process control, shared and remote memory access, file system access and control, interprocess communication, process signaling, and terminal I/O. High quality implementations should strive to make the results of such interactions intuitive to users, and attempt to document restrictions where deemed necessary.

Advice to implementors. Implementations that support such additional mechanisms for functionality supported within MPI are expected to document how these interact with MPI. (*End of advice to implementors.*)

The interaction of MPI and threads is defined in Section 11.7.

2.8 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.

Of course, MPI programs may still be erroneous. A **program error** can occur when an MPI call is made 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.

In C and Fortran, 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, except for file operations. 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 Chapter 7 of the MPI-1 document and in Section 7.3.1 of this document. The return values of C++ functions are not error codes. If the default error handler has been set to `MPI::ERRORS_THROW_EXCEPTIONS`, the C++ exception mechanism is used to signal an error by throwing an `MPI::Exception` object. See also Section 13.1.8 on page 445.

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 on the receiver 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 non-conforming behavior.

MPI-2 defines a way for users to create new error codes as defined in Section 11.5.

2.9 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 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.

2.9.1 Independence of Basic Runtime Routines

MPI programs require that library routines that are part of the basic language environment (such as `write` in Fortran and `printf` and `malloc` in ISO 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.

1 Note that this in no way prevents the creation of library routines that provide parallel
2 services whose operation is collective. However, the following program is expected to com-
3 plete in an ISO C environment regardless of the size of MPI_COMM_WORLD (assuming that
4 `printf` is available at the executing nodes).

```
5 int rank;  
6 MPI_Init((void *)0, (void *)0);  
7 MPI_Comm_rank(MPI_COMM_WORLD, &rank);  
8 if (rank == 0) printf("Starting program\n");  
9 MPI_Finalize();
```

11 The corresponding Fortran and C++ programs are also expected to complete.

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

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

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

23 2.9.2 Interaction with Signals

24 MPI does not specify the interaction of processes with signals and does not require that MPI
25 be signal safe. The implementation may reserve some signals for its own use. It is required
26 that the implementation document which signals it uses, and it is strongly recommended
27 that it not use SIGALRM, SIGFPE, or SIGIO. Implementations may also prohibit the use of
28 MPI calls from within signal handlers.

29 In multithreaded environments, users can avoid conflicts between signals and the MPI
30 library by catching signals only on threads that do not execute MPI calls. High quality
31 single-threaded implementations will be signal safe: an MPI call suspended by a signal will
32 resume and complete normally after the signal is handled.

35 2.10 Examples

36 The examples in this document are for illustration purposes only. They are not intended
37 to specify the standard. Furthermore, the examples have not been carefully checked or
38 verified.

Chapter 3

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)+1, 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 | <code>buf</code> | initial address of send buffer (choice) |
| IN | <code>count</code> | number of elements in send buffer (nonnegative integer) |
| IN | <code>datatype</code> | datatype of each send buffer element (handle) |
| IN | <code>dest</code> | rank of destination (integer) |
| IN | <code>tag</code> | message tag (integer) |
| IN | <code>comm</code> | 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
```

```
void MPI::Comm::Send(const void* buf, int count, const
                    MPI::Datatype& datatype, int dest, int tag) const
```

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

3.2.2 Message data

The send buffer specified by the `MPI_SEND` operation consists of `count` successive entries of the type indicated by `datatype`, starting with the entry at address `buf`. Note that we specify the message length in terms of number of *elements*, not number of *bytes*. The former is machine independent and closer to the application level.

The data part of the message consists of a sequence of `count` values, each of the type indicated by `datatype`. `count` may be zero, in which case the data part of the message is empty. The basic datatypes that can be specified for message data values correspond to the basic datatypes of the host language. Possible values of this argument for Fortran and the corresponding Fortran types are listed below.

| MPI datatype | Fortran datatype |
|-----------------------------------|-------------------------------|
| <code>MPI_INTEGER</code> | <code>INTEGER</code> |
| <code>MPI_REAL</code> | <code>REAL</code> |
| <code>MPI_DOUBLE_PRECISION</code> | <code>DOUBLE PRECISION</code> |
| <code>MPI_COMPLEX</code> | <code>COMPLEX</code> |
| <code>MPI_LOGICAL</code> | <code>LOGICAL</code> |
| <code>MPI_CHARACTER</code> | <code>CHARACTER(1)</code> |
| <code>MPI_BYTE</code> | |
| <code>MPI_PACKED</code> | |

Table 3.1: [Predefined MPI datatypes corresponding to Fortran datatypes](#)

Possible values for this argument for C and the corresponding C types are listed below.

| MPI datatype | C datatype |
|---|---|
| <code>MPI_CHAR</code> | <code>signed char</code> |
| <code>MPI_SHORT</code> | <code>signed short int</code> |
| <code>MPI_INT</code> | <code>signed int</code> |
| <code>MPI_LONG</code> | <code>signed long int</code> |
| <code>MPI_LONG_LONG_INT</code> | <code>signed long long int</code> |
| <code>MPI_LONG_LONG</code> (as a synonym) | <code>signed long long int</code> |
| <code>MPI_UNSIGNED_CHAR</code> | <code>unsigned char</code> |
| <code>MPI_UNSIGNED_SHORT</code> | <code>unsigned short int</code> |
| <code>MPI_UNSIGNED</code> | <code>unsigned int</code> |
| <code>MPI_UNSIGNED_LONG</code> | <code>unsigned long int</code> |
| <code>MPI_UNSIGNED_LONG_LONG</code> | <code>unsigned long long int</code> |
| <code>MPI_FLOAT</code> | <code>float</code> |
| <code>MPI_DOUBLE</code> | <code>double</code> |
| <code>MPI_LONG_DOUBLE</code> | <code>long double</code> |
| <code>MPI_WCHAR</code> | <code>wchar_t</code> (defined in <code><stddef.h></code>) |
| <code>MPI_BYTE</code> | |
| <code>MPI_PACKED</code> | |

Table 3.2: [Predefined MPI datatypes corresponding to C datatypes](#)

1 The datatypes `MPI_BYTE` and `MPI_PACKED` do not correspond to a Fortran or C
 2 datatype. A value of type `MPI_BYTE` consists of a byte (8 binary digits). A byte is
 3 uninterpreted and is different from a character. Different machines may have different
 4 representations for characters, or may use more than one byte to represent characters. On
 5 the other hand, a byte has the same binary value on all machines. The use of the type
 6 `MPI_PACKED` is explained in Section 3.13.

7 MPI requires support of the datatypes listed above, which match the basic datatypes of
 8 Fortran and ISO C. Additional MPI datatypes should be provided if the host language has
 9 additional data types: `MPI_DOUBLE_COMPLEX` for double precision complex in Fortran
 10 declared to be of type `DOUBLE COMPLEX`; `MPI_REAL2`, `MPI_REAL4` and `MPI_REAL8` for
 11 Fortran reals, declared to be of type `REAL*2`, `REAL*4` and `REAL*8`, respectively;
 12 `MPI_INTEGER1` `MPI_INTEGER2` and `MPI_INTEGER4` for Fortran integers, declared to be of
 13 type `INTEGER*1`, `INTEGER*2` and `INTEGER*4`, respectively; etc.

14
 15 *Rationale.* One goal of the design is to allow for MPI to be implemented as a
 16 library, with no need for additional preprocessing or compilation. Thus, one cannot
 17 assume that a communication call has information on the datatype of variables in the
 18 communication buffer; this information must be supplied by an explicit argument.
 19 The need for such datatype information will become clear in Section 3.3.2. (*End of*
 20 *rationale.*)

23 3.2.3 Message envelope

24 In addition to the data part, messages carry information that can be used to distinguish
 25 messages and selectively receive them. This information consists of a fixed number of fields,
 26 which we collectively call the **message envelope**. These fields are

28 source
 29 destination
 30 tag
 31 communicator

32
 33 The message source is implicitly determined by the identity of the message sender. The
 34 other fields are specified by arguments in the send operation.

35 The message destination is specified by the `dest` argument.

36 The integer-valued message tag is specified by the `tag` argument. This integer can be
 37 used by the program to distinguish different types of messages. The range of valid tag
 38 values is $0, \dots, UB$, where the value of `UB` is implementation dependent. It can be found by
 39 querying the value of the attribute `MPI_TAG_UB`, as described in Chapter 7. MPI requires
 40 that `UB` be no less than 32767.

41 The `comm` argument specifies the **communicator** that is used for the send operation.
 42 Communicators are explained in Chapter 5; below is a brief summary of their usage.

43 A communicator specifies the communication context for a communication operation.
 44 Each communication context provides a separate “communication universe:” messages are
 45 always received within the context they were sent, and messages sent in different contexts
 46 do not interfere.

47 The communicator also specifies the set of processes that share this communication
 48 context. This **process group** is ordered and processes are identified by their rank within

this group. Thus, the range of valid values for `dest` is 0, ... , n-1, where n is the number of processes in the group. (If the communicator is an inter-communicator, then destinations are identified by their rank in the remote group. See Chapter 5.)

A predefined communicator `MPI_COMM_WORLD` is provided by MPI. It allows communication with all processes that are accessible after MPI initialization and processes are identified by their rank in the group of `MPI_COMM_WORLD`.

Advice to users. Users that are comfortable with the notion of a flat name space for processes, and a single communication context, as offered by most existing communication libraries, need only use the predefined variable `MPI_COMM_WORLD` as the `comm` argument. This will allow communication with all the processes available at initialization time.

Users may define new communicators, as explained in Chapter 5. Communicators provide an important encapsulation mechanism for libraries and modules. They allow modules to have their own disjoint communication universe and their own process numbering scheme. (*End of advice to users.*)

Advice to implementors. The message envelope would normally be encoded by a fixed-length message header. However, the actual encoding is implementation dependent. Some of the information (e.g., source or destination) may be implicit, and need not be explicitly carried by messages. Also, processes may be identified by relative ranks, or absolute ids, etc. (*End of advice to implementors.*)

3.2.4 Blocking receive

The syntax of the blocking receive operation is given below.

`MPI_RECV (buf, count, datatype, source, tag, comm, status)`

| | | |
|-----|-----------------------|--|
| OUT | <code>buf</code> | initial address of receive buffer (choice) |
| IN | <code>count</code> | number of elements in receive buffer (integer) |
| IN | <code>datatype</code> | datatype of each receive buffer element (handle) |
| IN | <code>source</code> | rank of source (integer) |
| IN | <code>tag</code> | message tag (integer) |
| IN | <code>comm</code> | communicator (handle) |
| OUT | <code>status</code> | status object (Status) |

```
int MPI_Recv(void* buf, int count, MPI_Datatype datatype, int source,
            int tag, MPI_Comm comm, MPI_Status *status)
```

```
MPI_RECV(BUF, COUNT, DATATYPE, SOURCE, TAG, COMM, STATUS, IERROR)
<type> BUF(*)
INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE),
IERROR
```

```
void MPI::Comm::Recv(void* buf, int count, const MPI::Datatype& datatype,
                    int source, int tag, MPI::Status& status) const
```

```

1 void MPI::Comm::Recv(void* buf, int count, const MPI::Datatype& datatype,
2     int source, int tag) const
3

```

4 The blocking semantics of this call are described in Sec. 3.4.

5 The receive buffer consists of the storage containing `count` consecutive elements of the
6 type specified by `datatype`, starting at address `buf`. The length of the received message must
7 be less than or equal to the length of the receive buffer. An overflow error occurs if all
8 incoming data does not fit, without truncation, into the receive buffer.

9 If a message that is shorter than the receive buffer arrives, then only those locations
10 corresponding to the (shorter) message are modified.

11 *Advice to users.* The `MPI_PROBE` function described in Section 3.8 can be used to
12 receive messages of unknown length. (*End of advice to users.*)

13
14 *Advice to implementors.* Even though no specific behavior is mandated by MPI for
15 erroneous programs, the recommended handling of overflow situations is to return in
16 `status` information about the source and tag of the incoming message. The receive
17 operation will return an error code. A quality implementation will also ensure that
18 no memory that is outside the receive buffer will ever be overwritten.

19
20 In the case of a message shorter than the receive buffer, MPI is quite strict in that it
21 allows no modification of the other locations. A more lenient statement would allow
22 for some optimizations but this is not allowed. The implementation must be ready to
23 end a copy into the receiver memory exactly at the end of the receive buffer, even if
24 it is an odd address. (*End of advice to implementors.*)

25
26 The selection of a message by a receive operation is governed by the value of the
27 message envelope. A message can be received by a receive operation if its envelope matches
28 the `source`, `tag` and `comm` values specified by the receive operation. The receiver may specify
29 a wildcard `MPI_ANY_SOURCE` value for `source`, and/or a wildcard `MPI_ANY_TAG` value for
30 `tag`, indicating that any source and/or tag are acceptable. It cannot specify a wildcard value
31 for `comm`. Thus, a message can be received by a receive operation only if it is addressed to
32 the receiving process, has a matching communicator, has matching source unless `source=`
33 `MPI_ANY_SOURCE` in the pattern, and has a matching tag unless `tag=MPI_ANY_TAG` in the
34 pattern.

35 The message tag is specified by the `tag` argument of the receive operation. The
36 argument `source`, if different from `MPI_ANY_SOURCE`, is specified as a rank within the
37 process group associated with that same communicator (remote process group, for in-
38 tercommunicators). Thus, the range of valid values for the `source` argument is $\{0, \dots, n-1\} \cup \{\text{MPI_ANY_SOURCE}\}$, where n is the number of processes in this group.

39 Note the asymmetry between send and receive operations: A receive operation may
40 accept messages from an arbitrary sender, on the other hand, a send operation must specify
41 a unique receiver. This matches a “push” communication mechanism, where data transfer
42 is effected by the sender (rather than a “pull” mechanism, where data transfer is effected
43 by the receiver).

44
45 Source = destination is allowed, that is, a process can send a message to itself. (How-
46 ever, it is unsafe to do so with the blocking send and receive operations described above,
47 since this may lead to deadlock. See Sec. 3.5.)

Advice to implementors. Message context and other communicator information can be implemented as an additional tag field. It differs from the regular message tag in that wild card matching is not allowed on this field, and that value setting for this field is controlled by communicator manipulation functions. (*End of advice to implementors.*)

3.2.5 Return status

The source or tag of a received message may not be known if wildcard values were used in the receive operation. Also, if multiple requests are completed by a single MPI function (see Section 3.7.5), a distinct error code may need to be returned for each request. The information is returned by the `status` argument of `MPI_RECV`. The type of `status` is MPI-defined. Status variables need to be explicitly allocated by the user, that is, they are not system objects.

In C, `status` is a structure that contains three fields named `MPI_SOURCE`, `MPI_TAG`, and `MPI_ERROR`; the structure may contain additional fields. Thus, `status.MPI_SOURCE`, `status.MPI_TAG` and `status.MPI_ERROR` contain the source, tag, and error code, respectively, of the received message.

In Fortran, `status` is an array of `INTEGER`s of size `MPI_STATUS_SIZE`. The constants `MPI_SOURCE`, `MPI_TAG` and `MPI_ERROR` are the indices of the entries that store the source, tag and error fields. Thus, `status(MPI_SOURCE)`, `status(MPI_TAG)` and `status(MPI_ERROR)` contain, respectively, the source, tag and error code of the received message.

In C++, the `status` object is handled through the following methods:

```
int MPI::Status::Get_source() const
void MPI::Status::Set_source(int source)
int MPI::Status::Get_tag() const
void MPI::Status::Set_tag(int tag)
int MPI::Status::Get_error() const
void MPI::Status::Set_error(int error)
```

In general, message passing calls do not modify the value of the error code field of `status` variables. This field may be updated only by the functions in Section 3.7.5 which return multiple statuses. The field is updated if and only if such function returns with an error code of `MPI_ERR_IN_STATUS`.

Rationale. The error field in `status` is not needed for calls that return only one status, such as `MPI_WAIT`, since that would only duplicate the information returned by the function itself. The current design avoids the additional overhead of setting it, in such cases. The field is needed for calls that return multiple statuses, since each request may have had a different failure. (*End of rationale.*)

The `status` argument also returns information on the length of the message received. However, this information is not directly available as a field of the `status` variable and a call to `MPI_GET_COUNT` is required to “decode” this information.

```
1 MPI_GET_COUNT(status, datatype, count)
```

```
2     IN      status      return status of receive operation (Status)
3
4     IN      datatype    datatype of each receive buffer entry (handle)
5
6     OUT     count       number of received entries (integer)
```

```
7 int MPI_Get_count(MPI_Status *status, MPI_Datatype datatype, int *count)
```

```
8 MPI_GET_COUNT(STATUS, DATATYPE, COUNT, IERROR)
```

```
9     INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, COUNT, IERROR
```

```
10
11 int MPI::Status::Get_count(const MPI::Datatype& datatype) const
```

12
13 Returns the number of entries received. (Again, we count *entries*, each of type *datatype*,
14 not *bytes*.) The *datatype* argument should match the argument provided by the receive call
15 that set the status variable. (We shall later see, in Section 3.12.12, that MPI_GET_COUNT
16 may return, in certain situations, the value MPI_UNDEFINED.)

17 *Rationale.* Some message passing libraries use INOUT *count*, *tag* and
18 *source* arguments, thus using them both to specify the selection criteria for incoming
19 messages and return the actual envelope values of the received message. The use of a
20 separate status argument prevents errors that are often attached with INOUT argument
21 (e.g., using the MPI_ANY_TAG constant as the tag in a receive). Some libraries use calls
22 that refer implicitly to the “last message received.” This is not thread safe.

23 The *datatype* argument is passed to MPI_GET_COUNT so as to improve performance.
24 A message might be received without counting the number of elements it contains,
25 and the count value is often not needed. Also, this allows the same function to be
26 used after a call to MPI_PROBE or MPI_IPROBE. With a status from MPI_PROBE or
27 MPI_IPROBE, the same datatypes are allowed as in a call to MPI_RECV to receive this
28 message. (*End of rationale.*)

29
30 The value returned as the *count* argument of MPI_GET_COUNT for a datatype of length
31 zero where zero bytes have been transferred is zero. If the number of bytes transferred is
32 greater than zero, MPI_UNDEFINED is returned.

33 *Rationale.* Zero-length datatypes may be created in a number of cases. In MPI-2, an
34 important case is MPI_TYPE_CREATE_DARRAY, where the definition of the particular
35 darray results in an empty block on some MPI process. Programs written in an SPMD
36 style will not check for this special case and may want to use MPI_GET_COUNT to
37 check the status. (*End of rationale.*)

38
39 *Advice to users.* The buffer size required for the receive can be affected by data con-
40 versions and by the stride of the receive datatype. In most cases, the safest approach
41 is to use the same datatype with MPI_GET_COUNT and the receive. (*End of advice*
42 *to users.*)

43
44 All send and receive operations use the *buf*, *count*, *datatype*, *source*, *dest*, *tag*, *comm*
45 and *status* arguments in the same way as the blocking MPI_SEND and MPI_RECV operations
46 described in this section.

47 The following feature adds to, but does not change, the functionality associated with
48 MPI_STATUS.

3.2.6 Passing MPI_STATUS_IGNORE for Status

Every call to MPI_RECV includes a `status` argument, wherein the system can return details about the message received. There are also a number of other MPI calls, particularly in MPI-2, where `status` is returned. An object of type `MPI_STATUS` is not an MPI opaque object; its structure is declared in `mpi.h` and `mpif.h`, and it exists in the user's program. In many cases, application programs are constructed so that it is unnecessary for them to examine the `status` fields. In these cases, it is a waste for the user to allocate a status object, and it is particularly wasteful for the MPI implementation to fill in fields in this object.

To cope with this problem, there are two predefined constants, `MPI_STATUS_IGNORE` and `MPI_STATUSES_IGNORE`, which when passed to a receive, wait, or test function, inform the implementation that the status fields are not to be filled in. Note that `MPI_STATUS_IGNORE` is not a special type of `MPI_STATUS` object; rather, it is a special value for the argument. In C one would expect it to be `NULL`, not the address of a special `MPI_STATUS`.

`MPI_STATUS_IGNORE`, and the array version `MPI_STATUSES_IGNORE`, can be used everywhere a status argument is passed to a receive, wait, or test function. `MPI_STATUS_IGNORE` cannot be used when status is an IN argument. Note that in Fortran `MPI_STATUS_IGNORE` and `MPI_STATUSES_IGNORE` are objects like `MPI_BOTTOM` (not usable for initialization or assignment). See Section 2.5.4.

In general, this optimization can apply to all functions for which `status` or an array of `statuses` is an OUT argument. Note that this converts `status` into an INOUT argument. The functions that can be passed `MPI_STATUS_IGNORE` are all the various forms of `MPI_RECV`, `MPI_TEST`, and `MPI_WAIT`, as well as `MPI_REQUEST_GET_STATUS`. When an array is passed, as in the `ANY` and `ALL` functions, a separate constant, `MPI_STATUSES_IGNORE`, is passed for the array argument. It is possible for an MPI function to return `MPI_ERR_IN_STATUS` even when `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE` has been passed to that function.

`MPI_STATUS_IGNORE` and `MPI_STATUSES_IGNORE` are not required to have the same values in C and Fortran.

It is not allowed to have some of the statuses in an array of statuses for `_ANY` and `_ALL` functions set to `MPI_STATUS_IGNORE`; one either specifies ignoring *all* of the statuses in such a call with `MPI_STATUSES_IGNORE`, or *none* of them by passing normal statuses in all positions in the array of statuses.

There are no C++ bindings for `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE`. To allow an OUT or INOUT `MPI::Status` argument to be ignored, all MPI C++ bindings that have OUT or INOUT `MPI::Status` parameters are overloaded with a second version that omits the OUT or INOUT `MPI::Status` parameter.

Example 3.1 The C++ bindings for `MPI_PROBE` are:

```
void MPI::Comm::Probe(int source, int tag, MPI::Status& status) const
void MPI::Comm::Probe(int source, int tag) const
```

3.3 Data type matching and data conversion

3.3.1 Type matching rules

One can think of message transfer as consisting of the following three phases.

- 1 1. Data is pulled out of the send buffer and a message is assembled.
- 2
- 3 2. A message is transferred from sender to receiver.
- 4
- 5 3. Data is pulled from the incoming message and disassembled into the receive buffer.

6 Type matching has to be observed at each of these three phases: The type of each
 7 variable in the sender buffer has to match the type specified for that entry by the send
 8 operation; the type specified by the send operation has to match the type specified by the
 9 receive operation; and the type of each variable in the receive buffer has to match the type
 10 specified for that entry by the receive operation. A program that fails to observe these three
 11 rules is erroneous.

12 To define type matching more precisely, we need to deal with two issues: matching of
 13 types of the host language with types specified in communication operations; and matching
 14 of types at sender and receiver.

15 The types of a send and receive match (phase two) if both operations use identical
 16 names. That is, `MPI_INTEGER` matches `MPI_INTEGER`, `MPI_REAL` matches `MPI_REAL`,
 17 and so on. There is one exception to this rule, discussed in Sec. 3.13, the type `MPI_PACKED`
 18 can match any other type.

19 The type of a variable in a host program matches the type specified in the commu-
 20 nication operation if the datatype name used by that operation corresponds to the basic
 21 type of the host program variable. For example, an entry with type name `MPI_INTEGER`
 22 matches a Fortran variable of type `INTEGER`. A table giving this correspondence for Fortran
 23 and C appears in Sec. 3.2.2. There are two exceptions to this last rule: an entry with
 24 type name `MPI_BYTE` or `MPI_PACKED` can be used to match any byte of storage (on a
 25 byte-addressable machine), irrespective of the datatype of the variable that contains this
 26 byte. The type `MPI_PACKED` is used to send data that has been explicitly packed, or receive
 27 data that will be explicitly unpacked, see Section 3.13. The type `MPI_BYTE` allows one to
 28 transfer the binary value of a byte in memory unchanged.

29 To summarize, the type matching rules fall into the three categories below.

- 30 • Communication of typed values (e.g., with datatype different from `MPI_BYTE`), where
 31 the datatypes of the corresponding entries in the sender program, in the send call, in
 32 the receive call and in the receiver program must all match.
- 33
- 34 • Communication of untyped values (e.g., of datatype `MPI_BYTE`), where both sender
 35 and receiver use the datatype `MPI_BYTE`. In this case, there are no requirements on
 36 the types of the corresponding entries in the sender and the receiver programs, nor is
 37 it required that they be the same.
- 38
- 39 • Communication involving packed data, where `MPI_PACKED` is used.

40 The following examples illustrate the first two cases.

41 **Example 3.2** Sender and receiver specify matching types.

```

42 CALL MPI_COMM_RANK(comm, rank, ierr)
43 IF(rank.EQ.0) THEN
44     CALL MPI_SEND(a(1), 10, MPI_REAL, 1, tag, comm, ierr)
45 ELSE
46     CALL MPI_RECV(b(1), 15, MPI_REAL, 0, tag, comm, status, ierr)
47 END IF
48
```


This code is correct if both `a` and `b` are real arrays of size ≥ 10 . (In Fortran, it might be correct to use this code even if `a` or `b` have size < 10 : e.g., when `a(1)` can be equivalenced to an array with ten reals.)

Example 3.3 Sender and receiver do not specify matching types.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF(rank.EQ.0) THEN
    CALL MPI_SEND(a(1), 10, MPI_REAL, 1, tag, comm, ierr)
ELSE
    CALL MPI_RECV(b(1), 40, MPI_BYTE, 0, tag, comm, status, ierr)
END IF
```

This code is erroneous, since sender and receiver do not provide matching datatype arguments.

Example 3.4 Sender and receiver specify communication of untyped values.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF(rank.EQ.0) THEN
    CALL MPI_SEND(a(1), 40, MPI_BYTE, 1, tag, comm, ierr)
ELSE
    CALL MPI_RECV(b(1), 60, MPI_BYTE, 0, tag, comm, status, ierr)
END IF
```

This code is correct, irrespective of the type and size of `a` and `b` (unless this results in an out of bound memory access).

Advice to users. If a buffer of type `MPI_BYTE` is passed as an argument to `MPI_SEND`, then MPI will send the data stored at contiguous locations, starting from the address indicated by the `buf` argument. This may have unexpected results when the data layout is not as a casual user would expect it to be. For example, some Fortran compilers implement variables of type `CHARACTER` as a structure that contains the character length and a pointer to the actual string. In such an environment, sending and receiving a Fortran `CHARACTER` variable using the `MPI_BYTE` type will not have the anticipated result of transferring the character string. For this reason, the user is advised to use typed communications whenever possible. (*End of advice to users.*)

Type `MPI_CHARACTER`

The type `MPI_CHARACTER` matches one character of a Fortran variable of type `CHARACTER`, rather than the entire character string stored in the variable. Fortran variables of type `CHARACTER` or substrings are transferred as if they were arrays of characters. This is illustrated in the example below.

Example 3.5 Transfer of Fortran `CHARACTER`s.

```
CHARACTER*10 a
CHARACTER*10 b
```

```

1 CALL MPI_COMM_RANK(comm, rank, ierr)
2 IF(rank.EQ.0) THEN
3     CALL MPI_SEND(a, 5, MPI_CHARACTER, 1, tag, comm, ierr)
4 ELSE
5     CALL MPI_RECV(b(6:10), 5, MPI_CHARACTER, 0, tag, comm, status, ierr)
6 END IF

```

The last five characters of string `b` at process 1 are replaced by the first five characters of string `a` at process 0.

Rationale. The alternative choice would be for `MPI_CHARACTER` to match a character of arbitrary length. This runs into problems.

A Fortran character variable is a constant length string, with no special termination symbol. There is no fixed convention on how to represent characters, and how to store their length. Some compilers pass a character argument to a routine as a pair of arguments, one holding the address of the string and the other holding the length of string. Consider the case of an MPI communication call that is passed a communication buffer with type defined by a derived datatype (Section 3.12). If this communicator buffer contains variables of type `CHARACTER` then the information on their length will not be passed to the MPI routine.

This problem forces us to provide explicit information on character length with the MPI call. One could add a length parameter to the type `MPI_CHARACTER`, but this does not add much convenience and the same functionality can be achieved by defining a suitable derived datatype. (*End of rationale.*)

Advice to implementors. Some compilers pass Fortran `CHARACTER` arguments as a structure with a length and a pointer to the actual string. In such an environment, the MPI call needs to dereference the pointer in order to reach the string. (*End of advice to implementors.*)

3.3.2 Data conversion

One of the goals of MPI is to support parallel computations across heterogeneous environments. Communication in a heterogeneous environment may require data conversions. We use the following terminology.

type conversion changes the datatype of a value, e.g., by rounding a `REAL` to an `INTEGER`.

representation conversion changes the binary representation of a value, e.g., from Hex floating point to IEEE floating point.

The type matching rules imply that MPI communication never entails type conversion. On the other hand, MPI requires that a representation conversion be performed when a typed value is transferred across environments that use different representations for the datatype of this value. MPI does not specify rules for representation conversion. Such conversion is expected to preserve integer, logical or character values, and to convert a floating point value to the nearest value that can be represented on the target system.

Overflow and underflow exceptions may occur during floating point conversions. Conversion of integers or characters may also lead to exceptions when a value that can be

represented in one system cannot be represented in the other system. An exception occurring during representation conversion results in a failure of the communication. An error occurs either in the send operation, or the receive operation, or both.

If a value sent in a message is untyped (i.e., of type `MPI_BYTE`), then the binary representation of the byte stored at the receiver is identical to the binary representation of the byte loaded at the sender. This holds true, whether sender and receiver run in the same or in distinct environments. No representation conversion is required. (Note that representation conversion may occur when values of type `MPI_CHARACTER` or `MPI_CHAR` are transferred, for example, from an EBCDIC encoding to an ASCII encoding.)

No conversion need occur when an MPI program executes in a homogeneous system, where all processes run in the same environment.

Consider the three examples, 3.2–3.4. The first program is correct, assuming that `a` and `b` are `REAL` arrays of size ≥ 10 . If the sender and receiver execute in different environments, then the ten real values that are fetched from the send buffer will be converted to the representation for reals on the receiver site before they are stored in the receive buffer. While the number of real elements fetched from the send buffer equal the number of real elements stored in the receive buffer, the number of bytes stored need not equal the number of bytes loaded. For example, the sender may use a four byte representation and the receiver an eight byte representation for reals.

The second program is erroneous, and its behavior is undefined.

The third program is correct. The exact same sequence of forty bytes that were loaded from the send buffer will be stored in the receive buffer, even if sender and receiver run in a different environment. The message sent has exactly the same length (in bytes) and the same binary representation as the message received. If `a` and `b` are of different types, or if they are of the same type but different data representations are used, then the bits stored in the receive buffer may encode values that are different from the values they encoded in the send buffer.

Data representation conversion also applies to the envelope of a message: source, destination and tag are all integers that may need to be converted.

Advice to implementors. The current definition does not require messages to carry data type information. Both sender and receiver provide complete data type information. In a heterogeneous environment, one can either use a machine independent encoding such as XDR, or have the receiver convert from the sender representation to its own, or even have the sender do the conversion.

Additional type information might be added to messages in order to allow the system to detect mismatches between datatype at sender and receiver. This might be particularly useful in a slower but safer debug mode. (*End of advice to implementors.*)

MPI requires support for inter-language communication, i.e., if messages are sent by a C or C++ process and received by a Fortran process, or vice-versa. The behavior is defined in Section 13.3 on page 466.

3.4 Communication Modes

The send call described in Section 3.2.1 is **blocking**: it does not return until the message data and envelope have been safely stored away so that the sender is free to access and

1 overwrite the send buffer. The message might be copied directly into the matching receive
2 buffer, or it might be copied into a temporary system buffer.

3 Message buffering decouples the send and receive operations. A blocking send can com-
4 plete as soon as the message was buffered, even if no matching receive has been executed by
5 the receiver. On the other hand, message buffering can be expensive, as it entails additional
6 memory-to-memory copying, and it requires the allocation of memory for buffering. MPI
7 offers the choice of several communication modes that allow one to control the choice of the
8 communication protocol.

9 The send call described in Section 3.2.1 uses the **standard** communication mode. In
10 this mode, it is up to MPI to decide whether outgoing messages will be buffered. MPI may
11 buffer outgoing messages. In such a case, the send call may complete before a matching
12 receive is invoked. On the other hand, buffer space may be unavailable, or MPI may choose
13 not to buffer outgoing messages, for performance reasons. In this case, the send call will
14 not complete until a matching receive has been posted, and the data has been moved to the
15 receiver.

16 Thus, a send in standard mode can be started whether or not a matching receive has
17 been posted. It may complete before a matching receive is posted. The standard mode send
18 is **non-local**: successful completion of the send operation may depend on the occurrence
19 of a matching receive.

20
21 *Rationale.* The reluctance of MPI to mandate whether standard sends are buffering
22 or not stems from the desire to achieve portable programs. Since any system will run
23 out of buffer resources as message sizes are increased, and some implementations may
24 want to provide little buffering, MPI takes the position that correct (and therefore,
25 portable) programs do not rely on system buffering in standard mode. Buffering may
26 improve the performance of a correct program, but it doesn't affect the result of the
27 program. If the user wishes to guarantee a certain amount of buffering, the user-
28 provided buffer system of Sec. 3.6 should be used, along with the buffered-mode send.
29 (*End of rationale.*)

30
31 There are three additional communication modes.

32 A **buffered** mode send operation can be started whether or not a matching receive
33 has been posted. It may complete before a matching receive is posted. However, unlike
34 the standard send, this operation is **local**, and its completion does not depend on the
35 occurrence of a matching receive. Thus, if a send is executed and no matching receive is
36 posted, then MPI must buffer the outgoing message, so as to allow the send call to complete.
37 An error will occur if there is insufficient buffer space. The amount of available buffer space
38 is controlled by the user — see Section 3.6. Buffer allocation by the user may be required
39 for the buffered mode to be effective.

40 A send that uses the **synchronous** mode can be started whether or not a matching
41 receive was posted. However, the send will complete successfully only if a matching re-
42 ceive is posted, and the receive operation has started to receive the message sent by the
43 synchronous send. Thus, the completion of a synchronous send not only indicates that the
44 send buffer can be reused, but also indicates that the receiver has reached a certain point in
45 its execution, namely that it has started executing the matching receive. If both sends and
46 receives are blocking operations then the use of the synchronous mode provides synchronous
47 communication semantics: a communication does not complete at either end before both
48 processes rendezvous at the communication. A send executed in this mode is **non-local**.

A send that uses the **ready** communication mode may be started *only* if the matching receive is already posted. Otherwise, the operation is erroneous and its outcome is undefined. On some systems, this allows the removal of a hand-shake operation that is otherwise required and results in improved performance. The completion of the send operation does not depend on the status of a matching receive, and merely indicates that the send buffer can be reused. A send operation that uses the ready mode has the same semantics as a standard send operation, or a synchronous send operation; it is merely that the sender provides additional information to the system (namely that a matching receive is already posted), that can save some overhead. In a correct program, therefore, a ready send could be replaced by a standard send with no effect on the behavior of the program other than performance.

Three additional send functions are provided for the three additional communication modes. The communication mode is indicated by a one letter prefix: B for buffered, S for synchronous, and R for ready.

MPI_BSEND (buf, count, datatype, dest, tag, comm)

| | | |
|----|----------|---|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements in send buffer (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_Bsend(void* buf, int count, MPI_Datatype datatype, int dest,
             int tag, MPI_Comm comm)
```

```
MPI_BSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR)
```

```
<type> BUF(*)
```

```
INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR
```

```
void MPI::Comm::Bsend(const void* buf, int count, const
                     MPI::Datatype& datatype, int dest, int tag) const
```

Send in buffered mode.

MPI_SSEND (buf, count, datatype, dest, tag, comm)

| | | |
|----|----------|---|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements in send buffer (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_Ssend(void* buf, int count, MPI_Datatype datatype, int dest,
```

```

1         int tag, MPI_Comm comm)
2
3 MPI_SSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR)
4     <type> BUF(*)
5     INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR
6
7 void MPI::Comm::Ssend(const void* buf, int count, const
8     MPI::Datatype& datatype, int dest, int tag) const

```

Send in synchronous mode.

```

11 MPI_RSEND (buf, count, datatype, dest, tag, comm)
12
13     IN        buf                initial address of send buffer (choice)
14     IN        count              number of elements in send buffer (integer)
15     IN        datatype           datatype of each send buffer element (handle)
16     IN        dest               rank of destination (integer)
17     IN        tag                message tag (integer)
18     IN        comm              communicator (handle)
19
20

```

```

21 int MPI_Rsend(void* buf, int count, MPI_Datatype datatype, int dest,
22     int tag, MPI_Comm comm)
23

```

```

24 MPI_RSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR)
25     <type> BUF(*)
26     INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR

```

```

27 void MPI::Comm::Rsend(const void* buf, int count, const
28     MPI::Datatype& datatype, int dest, int tag) const
29

```

Send in ready mode.

There is only one receive operation, which can match any of the send modes. The receive operation described in the last section is **blocking**: it returns only after the receive buffer contains the newly received message. A receive can complete before the matching send has completed (of course, it can complete only after the matching send has started).

In a multi-threaded implementation of MPI, the system may de-schedule a thread that is blocked on a send or receive operation, and schedule another thread for execution in the same address space. In such a case it is the user's responsibility not to access or modify a communication buffer until the communication completes. Otherwise, the outcome of the computation is undefined.

Rationale. We prohibit read accesses to a send buffer while it is being used, even though the send operation is not supposed to alter the content of this buffer. This may seem more stringent than necessary, but the additional restriction causes little loss of functionality and allows better performance on some systems — consider the case where data transfer is done by a DMA engine that is not cache-coherent with the main processor. (*End of rationale.*)

Advice to implementors. Since a synchronous send cannot complete before a matching receive is posted, one will not normally buffer messages sent by such an operation.

It is recommended to choose buffering over blocking the sender, whenever possible, for standard sends. The programmer can signal his or her preference for blocking the sender until a matching receive occurs by using the synchronous send mode.

A possible communication protocol for the various communication modes is outlined below.

ready send: The message is sent as soon as possible.

synchronous send: The sender sends a request-to-send message. The receiver stores this request. When a matching receive is posted, the receiver sends back a permission-to-send message, and the sender now sends the message.

standard send: First protocol may be used for short messages, and second protocol for long messages.

buffered send: The sender copies the message into a buffer and then sends it with a nonblocking send (using the same protocol as for standard send).

Additional control messages might be needed for flow control and error recovery. Of course, there are many other possible protocols.

Ready send can be implemented as a standard send. In this case there will be no performance advantage (or disadvantage) for the use of ready send.

A standard send can be implemented as a synchronous send. In such a case, no data buffering is needed. However, many (most?) users expect some buffering.

In a multi-threaded environment, the execution of a blocking communication should block only the executing thread, allowing the thread scheduler to de-schedule this thread and schedule another thread for execution. (*End of advice to implementors.*)

3.5 Semantics of point-to-point communication

A valid MPI implementation guarantees certain general properties of point-to-point communication, which are described in this section.

Order Messages are *non-overtaking*: If a sender sends two messages in succession to the same destination, and both match the same receive, then this operation cannot receive the second message if the first one is still pending. If a receiver posts two receives in succession, and both match the same message, then the second receive operation cannot be satisfied by this message, if the first one is still pending. This requirement facilitates matching of sends to receives. It guarantees that message-passing code is deterministic, if processes are single-threaded and the wildcard MPI_ANY_SOURCE is not used in receives. (Some of the calls described later, such as MPI_CANCEL or MPI_WAITANY, are additional sources of nondeterminism.)

If a process has a single thread of execution, then any two communications executed by this process are ordered. On the other hand, if the process is multi-threaded, then the semantics of thread execution may not define a relative order between two send operations executed by two distinct threads. The operations are logically concurrent, even if one physically precedes the other. In such a case, the two messages sent can be received in any order. Similarly, if two receive operations that are logically concurrent receive two successively sent messages, then the two messages can match the two receives in either order.

1 **Example 3.6** An example of non-overtaking messages.

```
2
3 CALL MPI_COMM_RANK(comm, rank, ierr)
4 IF (rank.EQ.0) THEN
5     CALL MPI_BSEND(buf1, count, MPI_REAL, 1, tag, comm, ierr)
6     CALL MPI_BSEND(buf2, count, MPI_REAL, 1, tag, comm, ierr)
7 ELSE    ! rank.EQ.1
8     CALL MPI_RECV(buf1, count, MPI_REAL, 0, MPI_ANY_TAG, comm, status, ierr)
9     CALL MPI_RECV(buf2, count, MPI_REAL, 0, tag, comm, status, ierr)
10 END IF
```

11
12 The message sent by the first send must be received by the first receive, and the message
13 sent by the second send must be received by the second receive.

14
15 **Progress** If a pair of matching send and receives have been initiated on two processes, then
16 at least one of these two operations will complete, independently of other actions in the
17 system: the send operation will complete, unless the receive is satisfied by another message,
18 and completes; the receive operation will complete, unless the message sent is consumed by
19 another matching receive that was posted at the same destination process.

20
21 **Example 3.7** An example of two, intertwined matching pairs.

```
22 CALL MPI_COMM_RANK(comm, rank, ierr)
23 IF (rank.EQ.0) THEN
24     CALL MPI_BSEND(buf1, count, MPI_REAL, 1, tag1, comm, ierr)
25     CALL MPI_SSEND(buf2, count, MPI_REAL, 1, tag2, comm, ierr)
26 ELSE    ! rank.EQ.1
27     CALL MPI_RECV(buf1, count, MPI_REAL, 0, tag2, comm, status, ierr)
28     CALL MPI_RECV(buf2, count, MPI_REAL, 0, tag1, comm, status, ierr)
29 END IF
```

30
31 Both processes invoke their first communication call. Since the first send of process zero
32 uses the buffered mode, it must complete, irrespective of the state of process one. Since
33 no matching receive is posted, the message will be copied into buffer space. (If insufficient
34 buffer space is available, then the program will fail.) The second send is then invoked. At
35 that point, a matching pair of send and receive operation is enabled, and both operations
36 must complete. Process one next invokes its second receive call, which will be satisfied by
37 the buffered message. Note that process one received the messages in the reverse order they
38 were sent.

39
40 **Fairness** MPI makes no guarantee of *fairness* in the handling of communication. Suppose
41 that a send is posted. Then it is possible that the destination process repeatedly posts a
42 receive that matches this send, yet the message is never received, because it is each time
43 overtaken by another message, sent from another source. Similarly, suppose that a receive
44 was posted by a multi-threaded process. Then it is possible that messages that match this
45 receive are repeatedly received, yet the receive is never satisfied, because it is overtaken
46 by other receives posted at this node (by other executing threads). It is the programmer's
47 responsibility to prevent starvation in such situations.

48

Resource limitations Any pending communication operation consumes system resources that are limited. Errors may occur when lack of resources prevent the execution of an MPI call. A quality implementation will use a (small) fixed amount of resources for each pending send in the ready or synchronous mode and for each pending receive. However, buffer space may be consumed to store messages sent in standard mode, and must be consumed to store messages sent in buffered mode, when no matching receive is available. The amount of space available for buffering will be much smaller than program data memory on many systems. Then, it will be easy to write programs that overrun available buffer space.

MPI allows the user to provide buffer memory for messages sent in the buffered mode. Furthermore, MPI specifies a detailed operational model for the use of this buffer. An MPI implementation is required to do no worse than implied by this model. This allows users to avoid buffer overflows when they use buffered sends. Buffer allocation and use is described in Section 3.6.

A buffered send operation that cannot complete because of a lack of buffer space is erroneous. When such a situation is detected, an error is signalled that may cause the program to terminate abnormally. On the other hand, a standard send operation that cannot complete because of lack of buffer space will merely block, waiting for buffer space to become available or for a matching receive to be posted. This behavior is preferable in many situations. Consider a situation where a producer repeatedly produces new values and sends them to a consumer. Assume that the producer produces new values faster than the consumer can consume them. If buffered sends are used, then a buffer overflow will result. Additional synchronization has to be added to the program so as to prevent this from occurring. If standard sends are used, then the producer will be automatically throttled, as its send operations will block when buffer space is unavailable.

In some situations, a lack of buffer space leads to deadlock situations. This is illustrated by the examples below.

Example 3.8 An exchange of messages.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 1, tag, comm, ierr)
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 1, tag, comm, status, ierr)
ELSE    ! rank.EQ.1
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
END IF
```

This program will succeed even if no buffer space for data is available. The standard send operation can be replaced, in this example, with a synchronous send.

Example 3.9 An attempt to exchange messages.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 1, tag, comm, status, ierr)
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 1, tag, comm, ierr)
ELSE    ! rank.EQ.1
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
```

```

1     CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
2 END IF

```

The receive operation of the first process must complete before its send, and can complete only if the matching send of the second processor is executed. The receive operation of the second process must complete before its send and can complete only if the matching send of the first process is executed. This program will always deadlock. The same holds for any other send mode.

Example 3.10 An exchange that relies on buffering.

```

10 CALL MPI_COMM_RANK(comm, rank, ierr)
11 IF (rank.EQ.0) THEN
12     CALL MPI_SEND(sendbuf, count, MPI_REAL, 1, tag, comm, ierr)
13     CALL MPI_RECV(recvbuf, count, MPI_REAL, 1, tag, comm, status, ierr)
14 ELSE ! rank.EQ.1
15     CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
16     CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
17 END IF

```

The message sent by each process has to be copied out before the send operation returns and the receive operation starts. For the program to complete, it is necessary that at least one of the two messages sent be buffered. Thus, this program can succeed only if the communication system can buffer at least `count` words of data.

Advice to users. When standard send operations are used, then a deadlock situation may occur where both processes are blocked because buffer space is not available. The same will certainly happen, if the synchronous mode is used. If the buffered mode is used, and not enough buffer space is available, then the program will not complete either. However, rather than a deadlock situation, we shall have a buffer overflow error.

A program is “safe” if no message buffering is required for the program to complete. One can replace all sends in such program with synchronous sends, and the program will still run correctly. This conservative programming style provides the best portability, since program completion does not depend on the amount of buffer space available or in the communication protocol used.

Many programmers prefer to have more leeway and be able to use the “unsafe” programming style shown in example 3.10. In such cases, the use of standard sends is likely to provide the best compromise between performance and robustness: quality implementations will provide sufficient buffering so that “common practice” programs will not deadlock. The buffered send mode can be used for programs that require more buffering, or in situations where the programmer wants more control. This mode might also be used for debugging purposes, as buffer overflow conditions are easier to diagnose than deadlock conditions.

Nonblocking message-passing operations, as described in Section 3.7, can be used to avoid the need for buffering outgoing messages. This prevents deadlocks due to lack of buffer space, and improves performance, by allowing overlap of computation and communication, and avoiding the overheads of allocating buffers and copying messages into buffers. (*End of advice to users.*)

3.6 Buffer allocation and usage

A user may specify a buffer to be used for buffering messages sent in buffered mode. Buffering is done by the sender.

```
MPI_BUFFER_ATTACH( buffer, size)
```

| | | |
|----|--------|---------------------------------|
| IN | buffer | initial buffer address (choice) |
| IN | size | buffer size, in bytes (integer) |

```
int MPI_Buffer_attach( void* buffer, int size)
```

```
MPI_BUFFER_ATTACH( BUFFER, SIZE, IERROR)
```

```
<type> BUFFER(*)
```

```
INTEGER SIZE, IERROR
```

```
void MPI::Attach_buffer(void* buffer, int size)
```

Provides to MPI a buffer in the user's memory to be used for buffering outgoing messages. The buffer is used only by messages sent in buffered mode. Only one buffer can be attached to a process at a time.

```
MPI_BUFFER_DETACH( buffer_addr, size)
```

| | | |
|-----|-------------|---------------------------------|
| OUT | buffer_addr | initial buffer address (choice) |
| OUT | size | buffer size, in bytes (integer) |

```
int MPI_Buffer_detach( void* buffer_addr, int* size)
```

```
MPI_BUFFER_DETACH( BUFFER_ADDR, SIZE, IERROR)
```

```
<type> BUFFER_ADDR(*)
```

```
INTEGER SIZE, IERROR
```

```
int MPI::Detach_buffer(void*& buffer)
```

Detach the buffer currently associated with MPI. The call returns the address and the size of the detached buffer. This operation will block until all messages currently in the buffer have been transmitted. Upon return of this function, the user may reuse or deallocate the space taken by the buffer.

Example 3.11 Calls to attach and detach buffers.

```
#define BUFFSIZE 10000
```

```
int size
```

```
char *buff;
```

```
MPI_Buffer_attach( malloc(BUFFSIZE), BUFFSIZE);
```

```
/* a buffer of 10000 bytes can now be used by MPI_Bsend */
```

```
MPI_Buffer_detach( &buff, &size);
```

```
/* Buffer size reduced to zero */
```

```
MPI_Buffer_attach( buff, size);
```

```
/* Buffer of 10000 bytes available again */
```

1 *Advice to users.* Even though the C functions `MPI_Buffer_attach` and
 2 `MPI_Buffer_detach` both have a first argument of type `void*`, these arguments are used
 3 differently: A pointer to the buffer is passed to `MPI_Buffer_attach`; the address of the
 4 pointer is passed to `MPI_Buffer_detach`, so that this call can return the pointer value.
 5 (*End of advice to users.*)

6
 7 *Rationale.* Both arguments are defined to be of type `void*` (rather than `void*` and
 8 `void**`, respectively), so as to avoid complex type casts. E.g., in the last example,
 9 `&buff`, which is of type `char**`, can be passed as argument to `MPI_Buffer_detach` without
 10 type casting. If the formal parameter had type `void**` then we would need a type cast
 11 before and after the call. (*End of rationale.*)

12
 13 The statements made in this section describe the behavior of MPI for buffered-mode
 14 sends. When no buffer is currently associated, MPI behaves as if a zero-sized buffer is
 15 associated with the process.

16 MPI must provide as much buffering for outgoing messages *as if* outgoing message
 17 data were buffered by the sending process, in the specified buffer space, using a circular,
 18 contiguous-space allocation policy. We outline below a model implementation that defines
 19 this policy. MPI may provide more buffering, and may use a better buffer allocation algo-
 20 rithm than described below. On the other hand, MPI may signal an error whenever the
 21 simple buffering allocator described below would run out of space. In particular, if no buffer
 22 is explicitly associated with the process, then any buffered send may cause an error.

23 MPI does not provide mechanisms for querying or controlling buffering done by standard
 24 mode sends. It is expected that vendors will provide such information for their implemen-
 25 tations.

26
 27 *Rationale.* There is a wide spectrum of possible implementations of buffered com-
 28 munication: buffering can be done at sender, at receiver, or both; buffers can be
 29 dedicated to one sender-receiver pair, or be shared by all communications; buffering
 30 can be done in real or in virtual memory; it can use dedicated memory, or memory
 31 shared by other processes; buffer space may be allocated statically or be changed dy-
 32 namically; etc. It does not seem feasible to provide a portable mechanism for querying
 33 or controlling buffering that would be compatible with all these choices, yet provide
 34 meaningful information. (*End of rationale.*)

36 3.6.1 Model implementation of buffered mode

37 The model implementation uses the packing and unpacking functions described in Sec-
 38 tion 3.13 and the nonblocking communication functions described in Section 3.7.

39 We assume that a circular queue of pending message entries (PME) is maintained.
 40 Each entry contains a communication request handle that identifies a pending nonblocking
 41 send, a pointer to the next entry and the packed message data. The entries are stored in
 42 successive locations in the buffer. Free space is available between the queue tail and the
 43 queue head.

44 A buffered send call results in the execution of the following code.

- 45
- 46 • Traverse sequentially the PME queue from head towards the tail, deleting all entries
 47 for communications that have completed, up to the first entry with an uncompleted
 48 request; update queue head to point to that entry.

- Compute the number, n , of bytes needed to store an entry for the new message. An upper bound on n can be computed as follows: A call to the function `MPI_PACK_SIZE(count, datatype, comm, size)`, with the `count`, `datatype` and `comm` arguments used in the `MPI_BSEND` call, returns an upper bound on the amount of space needed to buffer the message data (see Section 3.13). The MPI constant `MPI_BSEND_OVERHEAD` provides an upper bound on the additional space consumed by the entry (e.g., for pointers or envelope information).
- Find the next contiguous empty space of n bytes in buffer (space following queue tail, or space at start of buffer if queue tail is too close to end of buffer). If space is not found then raise buffer overflow error.
- Append to end of PME queue in contiguous space the new entry that contains request handle, next pointer and packed message data; `MPI_PACK` is used to pack data.
- Post nonblocking send (standard mode) for packed data.
- Return

3.7 Nonblocking communication

One can improve performance on many systems by overlapping communication and computation. This is especially true on systems where communication can be executed autonomously by an intelligent communication controller. Light-weight threads are one mechanism for achieving such overlap. An alternative mechanism that often leads to better performance is to use **nonblocking communication**. A nonblocking **send start** call initiates the send operation, but does not complete it. The send start call `can` return before the message was copied out of the send buffer. A separate **send complete** call is needed to complete the communication, i.e., to verify that the data has been copied out of the send buffer. With suitable hardware, the transfer of data out of the sender memory may proceed concurrently with computations done at the sender after the send was initiated and before it completed. Similarly, a nonblocking **receive start call** initiates the receive operation, but does not complete it. The call `can` return before a message is stored into the receive buffer. A separate **receive complete** call is needed to complete the receive operation and verify that the data has been received into the receive buffer. With suitable hardware, the transfer of data into the receiver memory may proceed concurrently with computations done after the receive was initiated and before it completed. The use of nonblocking receives may also avoid system buffering and memory-to-memory copying, as information is provided early on the location of the receive buffer.

Nonblocking send start calls can use the same four modes as blocking sends: **standard**, **buffered**, **synchronous** and **ready**. These carry the same meaning. Sends of all modes, **ready** excepted, can be started whether a matching receive has been posted or not; a nonblocking **ready** send can be started only if a matching receive is posted. In all cases, the send start call is local: it returns immediately, irrespective of the status of other processes. If the call causes some system resource to be exhausted, then it will fail and return an error code. Quality implementations of MPI should ensure that this happens only in “pathological” cases. That is, an MPI implementation should be able to support a large number of pending nonblocking operations.

1 The send-complete call returns when data has been copied out of the send buffer. It
2 may carry additional meaning, depending on the send mode.

3 If the send mode is **synchronous**, then the send can complete only if a matching receive
4 has started. That is, a receive has been posted, and has been matched with the send. In
5 this case, the send-complete call is non-local. Note that a synchronous, nonblocking send
6 may complete, if matched by a nonblocking receive, before the receive complete call occurs.
7 (It can complete as soon as the sender “knows” the transfer will complete, but before the
8 receiver “knows” the transfer will complete.)

9 If the send mode is **buffered** then the message must be buffered if there is no pending
10 receive. In this case, the send-complete call is local, and must succeed irrespective of the
11 status of a matching receive.

12 If the send mode is **standard** then the send-complete call may return before a matching
13 receive occurred, if the message is buffered. On the other hand, the send-complete may not
14 complete until a matching receive occurred, and the message was copied into the receive
15 buffer.

16 Nonblocking sends can be matched with blocking receives, and vice-versa.

17
18 *Advice to users.* The completion of a send operation may be delayed, for standard
19 mode, and must be delayed, for synchronous mode, until a matching receive is posted.
20 The use of nonblocking sends in these two cases allows the sender to proceed ahead
21 of the receiver, so that the computation is more tolerant of fluctuations in the speeds
22 of the two processes.

23 Nonblocking sends in the buffered and ready modes have a more limited impact. A
24 nonblocking send will return as soon as possible, whereas a blocking send will return
25 after the data has been copied out of the sender memory. The use of nonblocking
26 sends is advantageous in these cases only if data copying can be concurrent with
27 computation.

28 The message-passing model implies that communication is initiated by the sender.
29 The communication will generally have lower overhead if a receive is already posted
30 when the sender initiates the communication (data can be moved directly to the
31 receive buffer, and there is no need to queue a pending send request). However, a
32 receive operation can complete only after the matching send has occurred. The use
33 of nonblocking receives allows one to achieve lower communication overheads without
34 blocking the receiver while it waits for the send. (*End of advice to users.*)
35

36 37 3.7.1 Communication Objects

38 Nonblocking communications use opaque **request** objects to identify communication oper-
39 ations and match the operation that initiates the communication with the operation that
40 terminates it. These are system objects that are accessed via a handle. A request object
41 identifies various properties of a communication operation, such as the send mode, the com-
42 munication buffer that is associated with it, its context, the tag and destination arguments
43 to be used for a send, or the tag and source arguments to be used for a receive. In addition,
44 this object stores information about the status of the pending communication operation.
45
46
47
48

3.7.2 Communication initiation

We use the same naming conventions as for blocking communication: a prefix of B, S, or R is used for buffered, synchronous or ready mode. In addition a prefix of I (for immediate) indicates that the call is nonblocking.

`MPI_ISEND(buf, count, datatype, dest, tag, comm, request)`

| | | |
|-----|----------|---|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements in send buffer (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) |
| OUT | request | communication request (handle) |

```
int MPI_Isend(void* buf, int count, MPI_Datatype datatype, int dest,
             int tag, MPI_Comm comm, MPI_Request *request)
```

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

```
MPI::Request MPI::Comm::Isend(const void* buf, int count, const
                               MPI::Datatype& datatype, int dest, int tag) const
```

Start a standard mode, nonblocking send.

`MPI_IBSEND(buf, count, datatype, dest, tag, comm, request)`

| | | |
|-----|----------|---|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements in send buffer (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) |
| OUT | request | communication request (handle) |

```
int MPI_Ibsend(void* buf, int count, MPI_Datatype datatype, int dest,
              int tag, MPI_Comm comm, MPI_Request *request)
```

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

```

1 MPI::Request MPI::Comm::IbSend(const void* buf, int count, const
2     MPI::Datatype& datatype, int dest, int tag) const
3
4     Start a buffered mode, nonblocking send.
5
6 MPI_ISSEND(buf, count, datatype, dest, tag, comm, request)
7
8     IN      buf          initial address of send buffer (choice)
9     IN      count       number of elements in send buffer (integer)
10    IN      datatype     datatype of each send buffer element (handle)
11    IN      dest         rank of destination (integer)
12    IN      tag          message tag (integer)
13    IN      comm         communicator (handle)
14    IN      comm         communicator (handle)
15    OUT     request      communication request (handle)
16
17
18 int MPI_Issend(void* buf, int count, MPI_Datatype datatype, int dest,
19     int tag, MPI_Comm comm, MPI_Request *request)
20
21 MPI_ISSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
22     <type> BUF(*)
23     INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR
24
25 MPI::Request MPI::Comm::Issend(const void* buf, int count, const
26     MPI::Datatype& datatype, int dest, int tag) const
27
28     Start a synchronous mode, nonblocking send.
29
30 MPI_IRSEND(buf, count, datatype, dest, tag, comm, request)
31
32    IN      buf          initial address of send buffer (choice)
33    IN      count       number of elements in send buffer (integer)
34    IN      datatype     datatype of each send buffer element (handle)
35    IN      dest         rank of destination (integer)
36    IN      tag          message tag (integer)
37    IN      comm         communicator (handle)
38    IN      comm         communicator (handle)
39    OUT     request      communication request (handle)
40
41
42 int MPI_Irsend(void* buf, int count, MPI_Datatype datatype, int dest,
43     int tag, MPI_Comm comm, MPI_Request *request)
44
45 MPI_IRSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
46     <type> BUF(*)
47     INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR
48
49 MPI::Request MPI::Comm::Irsend(const void* buf, int count, const
50     MPI::Datatype& datatype, int dest, int tag) const

```


Start a ready mode nonblocking send.

`MPI_Irecv (buf, count, datatype, source, tag, comm, request)`

| | | |
|-----|----------|--|
| OUT | buf | initial address of receive buffer (choice) |
| IN | count | number of elements in receive buffer (integer) |
| IN | datatype | datatype of each receive buffer element (handle) |
| IN | source | rank of source (integer) |
| IN | tag | message tag (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Irecv(void* buf, int count, MPI_Datatype datatype, int source,
             int tag, MPI_Comm comm, MPI_Request *request)
```

```
MPI_Irecv(BUF, COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR)
```

```
<type> BUF(*)
```

```
INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR
```

```
MPI::Request MPI::Comm::Irecv(void* buf, int count, const
                               MPI::Datatype& datatype, int source, int tag) const
```

Start a nonblocking receive.

These calls allocate a communication request object and associate it with the request handle (the argument `request`). The request can be used later to query the status of the communication or wait for its completion.

A nonblocking send call indicates that the system may start copying data out of the send buffer. The sender should not access any part of the send buffer after a nonblocking send operation is called, until the send completes.

A nonblocking receive call indicates that the system may start writing data into the receive buffer. The receiver should not access any part of the receive buffer after a nonblocking receive operation is called, until the receive completes.

Advice to users. To prevent problems with the argument copying and register optimization done by Fortran compilers, please note the hints in subsections “Problems Due to Data Copying and Sequence Association,” and “A Problem with Register Optimization” in [Section 13.2.2 on pages 451 and 454](#). (*End of advice to users.*)

3.7.3 Communication Completion

The functions `MPI_WAIT` and `MPI_TEST` are used to complete a nonblocking communication. The completion of a send operation indicates that the sender is now free to update the locations in the send buffer (the send operation itself leaves the content of the send buffer unchanged). It does not indicate that the message has been received, rather, it may have been buffered by the communication subsystem. However, if a synchronous mode send was used, the completion of the send operation indicates that a matching receive was initiated, and that the message will eventually be received by this matching receive.

The completion of a receive operation indicates that the receive buffer contains the received message, the receiver is now free to access it, and that the status object is set. It does not indicate that the matching send operation has completed (but indicates, of course, that the send was initiated).

We shall use the following terminology: A **null** handle is a handle with value `MPI_REQUEST_NULL`. A persistent request and the handle to it are **inactive** if the request is not associated with any ongoing communication (see Section 3.9). A handle is **active** if it is neither null nor inactive. An **empty** status is a status which is set to return `tag = MPI_ANY_TAG`, `source = MPI_ANY_SOURCE`, `error = MPI_SUCCESS`, and is also internally configured so that calls to `MPI_GET_COUNT` and `MPI_GET_ELEMENTS` return `count = 0` and `MPI_TEST_CANCELLED` returns false. We set a status variable to empty when the value returned by it is not significant. Status is set in this way so as to prevent errors due to accesses of stale information.

The fields in a **status** object returned by a call to `MPI_WAIT`, `MPI_TEST`, or any of the other derived functions (`MPI_{TEST, WAIT}_{ALL, SOME, ANY}`), where the **request** corresponds to a send call, are undefined, with two exceptions: The error status field will contain valid information if the wait or test call returned with `MPI_ERR_IN_STATUS`; and the returned status can be queried by the call `MPI_TEST_CANCELLED`.

Error codes belonging to the error class `MPI_ERR_IN_STATUS` should be returned only by the MPI completion functions that take arrays of `MPI_STATUS`. For the functions (`MPI_TEST`, `MPI_TESTANY`, `MPI_WAIT`, `MPI_WAITANY`) that return a single `MPI_STATUS` value, the normal MPI error return process should be used (not the `MPI_ERROR` field in the `MPI_STATUS` argument).

```
MPI_WAIT(request, status)
```

| | | |
|-------|---------|------------------------|
| INOUT | request | request (handle) |
| OUT | status | status object (Status) |

```
int MPI_Wait(MPI_Request *request, MPI_Status *status)
```

```
MPI_WAIT(REQUEST, STATUS, IERROR)
```

```
INTEGER REQUEST, STATUS(MPI_STATUS_SIZE), IERROR
```

```
void MPI::Request::Wait(MPI::Status& status)
```

```
void MPI::Request::Wait()
```

A call to `MPI_WAIT` returns when the operation identified by **request** is complete. If the communication object associated with this request was created by a nonblocking send or receive call, then the object is deallocated by the call to `MPI_WAIT` and the request handle is set to `MPI_REQUEST_NULL`. `MPI_WAIT` is a non-local operation.

The call returns, in **status**, information on the completed operation. The content of the status object for a receive operation can be accessed as described in section 3.2.5. The status object for a send operation may be queried by a call to `MPI_TEST_CANCELLED` (see Section 3.8).

One is allowed to call `MPI_WAIT` with a null or inactive **request** argument. In this case the operation returns immediately with empty **status**.

Advice to users. Successful return of `MPI_WAIT` after a `MPI_IBSEND` implies that the user send buffer can be reused — i.e., data has been sent out or copied into a buffer attached with `MPI_BUFFER_ATTACH`. Note that, at this point, we can no longer cancel the send (see Sec. 3.8). If a matching receive is never posted, then the buffer cannot be freed. This runs somewhat counter to the stated goal of `MPI_CANCEL` (always being able to free program space that was committed to the communication subsystem). (*End of advice to users.*)

Advice to implementors. In a multi-threaded environment, a call to `MPI_WAIT` should block only the calling thread, allowing the thread scheduler to schedule another thread for execution. (*End of advice to implementors.*)

`MPI_TEST(request, flag, status)`

| | | |
|-------|---------|---------------------------------------|
| INOUT | request | communication request (handle) |
| OUT | flag | true if operation completed (logical) |
| OUT | status | status object (Status) |

`int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status)`

`MPI_TEST(REQUEST, FLAG, STATUS, IERROR)`
 LOGICAL FLAG
 INTEGER REQUEST, STATUS(MPI_STATUS_SIZE), IERROR

`bool MPI::Request::Test(MPI::Status& status)`

`bool MPI::Request::Test()`

A call to `MPI_TEST` returns `flag = true` if the operation identified by `request` is complete. In such a case, the status object is set to contain information on the completed operation; if the communication object was created by a nonblocking send or receive, then it is deallocated and the request handle is set to `MPI_REQUEST_NULL`. The call returns `flag = false`, otherwise. In this case, the value of the status object is undefined. `MPI_TEST` is a local operation.

The return status object for a receive operation carries information that can be accessed as described in section 3.2.5. The status object for a send operation carries information that can be accessed by a call to `MPI_TEST_CANCELLED` (see Section 3.8).

One is allowed to call `MPI_TEST` with a null or inactive `request` argument. In such a case the operation returns with `flag = true` and empty `status`.

The functions `MPI_WAIT` and `MPI_TEST` can be used to complete both sends and receives.

Advice to users. The use of the nonblocking `MPI_TEST` call allows the user to schedule alternative activities within a single thread of execution. An event-driven thread scheduler can be emulated with periodic calls to `MPI_TEST`. (*End of advice to users.*)

Rationale. The function `MPI_TEST` returns with `flag = true` exactly in those situations where the function `MPI_WAIT` returns; both functions return in such case the

1 same value in `status`. Thus, a blocking `Wait` can be easily replaced by a nonblocking
 2 `Test`. (*End of rationale.*)

3
 4 **Example 3.12** Simple usage of nonblocking operations and `MPI_WAIT`.

```

5 CALL MPI_COMM_RANK(comm, rank, ierr)
6 IF(rank.EQ.0) THEN
7     CALL MPI_ISEND(a(1), 10, MPI_REAL, 1, tag, comm, request, ierr)
8     **** do some computation to mask latency ****
9     CALL MPI_WAIT(request, status, ierr)
10 ELSE
11     CALL MPI_IRECV(a(1), 15, MPI_REAL, 0, tag, comm, request, ierr)
12     **** do some computation to mask latency ****
13     CALL MPI_WAIT(request, status, ierr)
14 END IF

```

15
 16 A request object can be deallocated without waiting for the associated communication
 17 to complete, by using the following operation.

```

18
19 MPI_REQUEST_FREE(request)
20     INOUT      request                communication request (handle)
21
22
23 int MPI_Request_free(MPI_Request *request)
24
25 MPI_REQUEST_FREE(REQUEST, IERROR)
26     INTEGER REQUEST, IERROR
27
28 void MPI::Request::Free()

```

29 Mark the request object for deallocation and set `request` to `MPI_REQUEST_NULL`. An
 30 ongoing communication that is associated with the request will be allowed to complete.
 31 The request will be deallocated only after its completion.

32 *Rationale.* The `MPI_REQUEST_FREE` mechanism is provided for reasons of perfor-
 33 mance and convenience on the sending side. (*End of rationale.*)

34
 35 *Advice to users.* Once a request is freed by a call to `MPI_REQUEST_FREE`, it is
 36 not possible to check for the successful completion of the associated communication
 37 with calls to `MPI_WAIT` or `MPI_TEST`. Also, if an error occurs subsequently during
 38 the communication, an error code cannot be returned to the user — such an error
 39 must be treated as fatal. Questions arise as to how one knows when the operations
 40 have completed when using `MPI_REQUEST_FREE`. Depending on the program logic,
 41 there may be other ways in which the program knows that certain operations have
 42 completed and this makes usage of `MPI_REQUEST_FREE` practical. For example, an
 43 active send request could be freed when the logic of the program is such that the
 44 receiver sends a reply to the message sent — the arrival of the reply informs the
 45 sender that the send has completed and the send buffer can be reused. An active
 46 receive request should never be freed as the receiver will have no way to verify that
 47 the receive has completed and the receive buffer can be reused. (*End of advice to*
 48 *users.*)

Example 3.13 An example using MPI_REQUEST_FREE.

```

CALL MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
IF(rank.EQ.0) THEN
  DO i=1, n
    CALL MPI_ISEND(outval, 1, MPI_REAL, 1, 0, MPI_COMM_WORLD, req, ierr)
    CALL MPI_REQUEST_FREE(req, ierr)
    CALL MPI_Irecv(ival, 1, MPI_REAL, 1, 0, MPI_COMM_WORLD, req, ierr)
    CALL MPI_WAIT(req, status, ierr)
  END DO
ELSE ! rank.EQ.1
  CALL MPI_Irecv(ival, 1, MPI_REAL, 0, 0, MPI_COMM_WORLD, req, ierr)
  CALL MPI_WAIT(req, status, ierr)
  DO I=1, n-1
    CALL MPI_ISEND(outval, 1, MPI_REAL, 0, 0, MPI_COMM_WORLD, req, ierr)
    CALL MPI_REQUEST_FREE(req, ierr)
    CALL MPI_Irecv(ival, 1, MPI_REAL, 0, 0, MPI_COMM_WORLD, req, ierr)
    CALL MPI_WAIT(req, status, ierr)
  END DO
  CALL MPI_ISEND(outval, 1, MPI_REAL, 0, 0, MPI_COMM_WORLD, req, ierr)
  CALL MPI_WAIT(req, status, ierr)
END IF

```

3.7.4 Semantics of Nonblocking Communications

The semantics of nonblocking communication is defined by suitably extending the definitions in Section 3.5.

Order Nonblocking communication operations are ordered according to the execution order of the calls that initiate the communication. The non-overtaking requirement of Section 3.5 is extended to nonblocking communication, with this definition of order being used.

Example 3.14 Message ordering for nonblocking operations.

```

CALL MPI_COMM_RANK(comm, rank, ierr)
IF (RANK.EQ.0) THEN
  CALL MPI_ISEND(a, 1, MPI_REAL, 1, 0, comm, r1, ierr)
  CALL MPI_ISEND(b, 1, MPI_REAL, 1, 0, comm, r2, ierr)
ELSE ! rank.EQ.1
  CALL MPI_Irecv(a, 1, MPI_REAL, 0, MPI_ANY_TAG, comm, r1, ierr)
  CALL MPI_Irecv(b, 1, MPI_REAL, 0, 0, comm, r2, ierr)
END IF
CALL MPI_WAIT(r1, status, ierr)
CALL MPI_WAIT(r2, status, ierr)

```

The first send of process zero will match the first receive of process one, even if both messages are sent before process one executes either receive.

Progress A call to `MPI_WAIT` that completes a receive will eventually terminate and return if a matching send has been started, unless the send is satisfied by another receive. In particular, if the matching send is nonblocking, then the receive should complete even if no call is executed by the sender to complete the send. Similarly, a call to `MPI_WAIT` that completes a send will eventually return if a matching receive has been started, unless the receive is satisfied by another send, and even if no call is executed to complete the receive.

Example 3.15 An illustration of progress semantics.

```

9 CALL MPI_COMM_RANK(comm, rank, ierr)
10 IF (RANK.EQ.0) THEN
11     CALL MPI_SSEND(a, 1, MPI_REAL, 1, 0, comm, ierr)
12     CALL MPI_SEND(b, 1, MPI_REAL, 1, 1, comm, ierr)
13 ELSE ! rank.EQ.1
14     CALL MPI_IRecv(a, 1, MPI_REAL, 0, 0, comm, r, ierr)
15     CALL MPI_RECV(b, 1, MPI_REAL, 0, 1, comm, ierr)
16     CALL MPI_WAIT(r, status, ierr)
17 END IF
18

```

This code should not deadlock in a correct MPI implementation. The first synchronous send of process zero must complete after process one posts the matching (nonblocking) receive even if process one has not yet reached the completing wait call. Thus, process zero will continue and execute the second send, allowing process one to complete execution.

If an `MPI_TEST` that completes a receive is repeatedly called with the same arguments, and a matching send has been started, then the call will eventually return `flag = true`, unless the send is satisfied by another receive. If an `MPI_TEST` that completes a send is repeatedly called with the same arguments, and a matching receive has been started, then the call will eventually return `flag = true`, unless the receive is satisfied by another send.

3.7.5 Multiple Completions

It is convenient to be able to wait for the completion of any, some, or all the operations in a list, rather than having to wait for a specific message. A call to `MPI_WAITANY` or `MPI_TESTANY` can be used to wait for the completion of one out of several operations. A call to `MPI_WAITALL` or `MPI_TESTALL` can be used to wait for all pending operations in a list. A call to `MPI_WAIT SOME` or `MPI_TEST SOME` can be used to complete all enabled operations in a list.

`MPI_WAITANY` (count, array_of_requests, index, status)

| | | |
|-------|-------------------|--|
| IN | count | list length (integer) |
| INOUT | array_of_requests | array of requests (array of handles) |
| OUT | index | index of handle for operation that completed (integer) |
| OUT | status | status object (Status) |

```

47 int MPI_Waitany(int count, MPI_Request *array_of_requests, int *index,
48               MPI_Status *status)

```

```
MPI_WAITANY(COUNT, ARRAY_OF_REQUESTS, INDEX, STATUS, IERROR)
    INTEGER COUNT, ARRAY_OF_REQUESTS(*), INDEX, STATUS(MPI_STATUS_SIZE),
    IERROR
```

```
static int MPI::Request::Waitany(int count,
    MPI::Request array_of_requests[], MPI::Status& status)
```

```
static int MPI::Request::Waitany(int count,
    MPI::Request array_of_requests[])
```

Blocks until one of the operations associated with the active requests in the array has completed. If more than one operation is enabled and can terminate, one is arbitrarily chosen. Returns in `index` the index of that request in the array and returns in `status` the status of the completing communication. (The array is indexed from zero in C, and from one in Fortran.) If the request was allocated by a nonblocking communication operation, then it is deallocated and the request handle is set to `MPI_REQUEST_NULL`.

The `array_of_requests` list may contain null or inactive handles. If the list contains no active handles (list has length zero or all entries are null or inactive), then the call returns immediately with `index = MPI_UNDEFINED`, and an empty status.

The execution of `MPI_WAITANY(count, array_of_requests, index, status)` has the same effect as the execution of `MPI_WAIT(&array_of_requests[i], status)`, where `i` is the value returned by `index` (unless the value of `index` is `MPI_UNDEFINED`). `MPI_WAITANY` with an array containing one active entry is equivalent to `MPI_WAIT`.

```
MPI_TESTANY(count, array_of_requests, index, flag, status)
```

| | | |
|-------|-------------------|---|
| IN | count | list length (integer) |
| INOUT | array_of_requests | array of requests (array of handles) |
| OUT | index | index of operation that completed, or <code>MPI_UNDEFINED</code> if none completed (integer) |
| OUT | flag | true if one of the operations is complete (logical) |
| OUT | status | status object (Status) |

```
int MPI_Testany(int count, MPI_Request *array_of_requests, int *index,
    int *flag, MPI_Status *status)
```

```
MPI_TESTANY(COUNT, ARRAY_OF_REQUESTS, INDEX, FLAG, STATUS, IERROR)
    LOGICAL FLAG
    INTEGER COUNT, ARRAY_OF_REQUESTS(*), INDEX, STATUS(MPI_STATUS_SIZE),
    IERROR
```

```
static bool MPI::Request::Testany(int count,
    MPI::Request array_of_requests[], int& index,
    MPI::Status& status)
```

```
static bool MPI::Request::Testany(int count,
    MPI::Request array_of_requests[], int& index)
```

Tests for completion of either one or none of the operations associated with active handles. In the former case, it returns `flag = true`, returns in `index` the index of this request

1 in the array, and returns in `status` the status of that operation; if the request was allocated
 2 by a nonblocking communication call then the request is deallocated and the handle is set
 3 to `MPI_REQUEST_NULL`. (The array is indexed from zero in C, and from one in Fortran.)
 4 In the latter case (no operation completed), it returns `flag = false`, returns a value of
 5 `MPI_UNDEFINED` in `index` and `status` is undefined.

6 The array may contain null or inactive handles. If the array contains no active handles
 7 then the call returns immediately with `flag = true`, `index = MPI_UNDEFINED`, and an empty
 8 `status`.

9 If the array of requests contains active handles then the execution of
 10 `MPI_TESTANY(count, array_of_requests, index, status)` has the same effect as the execution
 11 of `MPI_TEST(&array_of_requests[i], flag, status)`, for $i=0, 1, \dots, \text{count}-1$, in some arbitrary
 12 order, until one call returns `flag = true`, or all fail. In the former case, `index` is set to the
 13 last value of i , and in the latter case, it is set to `MPI_UNDEFINED`. `MPI_TESTANY` with an
 14 array containing one active entry is equivalent to `MPI_TEST`.

15
 16 *Rationale.* The function `MPI_TESTANY` returns with `flag = true` exactly in those
 17 situations where the function `MPI_WAITANY` returns; both functions return in that
 18 case the same values in the remaining parameters. Thus, a blocking `MPI_WAITANY`
 19 can be easily replaced by a nonblocking `MPI_TESTANY`. The same relation holds for
 20 the other pairs of Wait and Test functions defined in this section. (*End of rationale.*)
 21
 22

23 `MPI_WAITALL(count, array_of_requests, array_of_statuses)`

| | | | |
|----|-------|-------------------|---|
| 25 | IN | count | lists length (integer) |
| 26 | INOUT | array_of_requests | array of requests (array of handles) |
| 27 | OUT | array_of_statuses | array of status objects (array of Status) |

29
 30 `int MPI_Waitall(int count, MPI_Request *array_of_requests,`
 31 `MPI_Status *array_of_statuses)`

32 `MPI_WAITALL(COUNT, ARRAY_OF_REQUESTS, ARRAY_OF_STATUSES, IERROR)`
 33 `INTEGER COUNT, ARRAY_OF_REQUESTS(*)`
 34 `INTEGER ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR`

35
 36 `static void MPI::Request::Waitall(int count,`
 37 `MPI::Request array_of_requests[],`
 38 `MPI::Status array_of_statuses[])`

39
 40 `static void MPI::Request::Waitall(int count,`
 41 `MPI::Request array_of_requests[])`

42 Blocks until all communication operations associated with active handles in the list
 43 complete, and return the status of all these operations (this includes the case where no
 44 handle in the list is active). Both arrays have the same number of valid entries. The i -th
 45 entry in `array_of_statuses` is set to the return status of the i -th operation. Requests that were
 46 created by nonblocking communication operations are deallocated and the corresponding
 47 handles in the array are set to `MPI_REQUEST_NULL`. The list may contain null or inactive
 48 handles. The call sets to empty the status of each such entry.

The error-free execution of `MPI_WAITALL(count, array_of_requests, array_of_statuses)` has the same effect as the execution of `MPI_WAIT(&array_of_request[i], &array_of_statuses[i])`, for $i=0, \dots, \text{count}-1$, in some arbitrary order. `MPI_WAITALL` with an array of length one is equivalent to `MPI_WAIT`.

When one or more of the communications completed by a call to `MPI_WAITALL` fail, it is desirable to return specific information on each communication. The function `MPI_WAITALL` will return in such case the error code `MPI_ERR_IN_STATUS` and will set the error field of each status to a specific error code. This code will be `MPI_SUCCESS`, if the specific communication completed; it will be another specific error code, if it failed; or it can be `MPI_ERR_PENDING` if it has neither failed nor completed. The function `MPI_WAITALL` will return `MPI_SUCCESS` if no request had an error, or will return another error code if it failed for other reasons (such as invalid arguments). In such cases, it will not update the error fields of the statuses.

Rationale. This design streamlines error handling in the application. The application code need only test the (single) function result to determine if an error has occurred. It needs to check each individual status only when an error occurred. (*End of rationale.*)

`MPI_TESTALL(count, array_of_requests, flag, array_of_statuses)`

| | | |
|-------|-------------------|---|
| IN | count | lists length (integer) |
| INOUT | array_of_requests | array of requests (array of handles) |
| OUT | flag | (logical) |
| OUT | array_of_statuses | array of status objects (array of Status) |

```
int MPI_Testall(int count, MPI_Request *array_of_requests, int *flag,
                MPI_Status *array_of_statuses)
```

```
MPI_TESTALL(COUNT, ARRAY_OF_REQUESTS, FLAG, ARRAY_OF_STATUSES, IERROR)
LOGICAL FLAG
INTEGER COUNT, ARRAY_OF_REQUESTS(*),
ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR
```

```
static bool MPI::Request::Testall(int count,
                                   MPI::Request array_of_requests[],
                                   MPI::Status array_of_statuses[])
```

```
static bool MPI::Request::Testall(int count,
                                   MPI::Request array_of_requests[])
```

Returns `flag = true` if all communications associated with active handles in the array have completed (this includes the case where no handle in the list is active). In this case, each status entry that corresponds to an active handle request is set to the status of the corresponding communication; if the request was allocated by a nonblocking communication call then it is deallocated, and the handle is set to `MPI_REQUEST_NULL`. Each status entry that corresponds to a null or inactive handle is set to empty.

Otherwise, `flag = false` is returned, no request is modified and the values of the status entries are undefined. This is a local operation.

Errors that occurred during the execution of `MPI_TESTALL` are handled as errors in `MPI_WAITALL`.

`MPI_WAITSSOME`(`incount`, `array_of_requests`, `outcount`, `array_of_indices`, `array_of_statuses`)

| | | |
|-------|--------------------------------|---|
| IN | <code>incount</code> | length of <code>array_of_requests</code> (integer) |
| INOUT | <code>array_of_requests</code> | array of requests (array of handles) |
| OUT | <code>outcount</code> | number of completed requests (integer) |
| OUT | <code>array_of_indices</code> | array of indices of operations that completed (array of integers) |
| OUT | <code>array_of_statuses</code> | array of status objects for operations that completed (array of Status) |

```
int MPI_Waitssome(int incount, MPI_Request *array_of_requests, int *outcount,
                 int *array_of_indices, MPI_Status *array_of_statuses)
```

```
MPI_WAITSSOME(INCOUNT, ARRAY_OF_REQUESTS, OUTCOUNT, ARRAY_OF_INDICES,
              ARRAY_OF_STATUSES, IERROR)
INTEGER INCOUNT, ARRAY_OF_REQUESTS(*), OUTCOUNT, ARRAY_OF_INDICES(*),
ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR
```

```
static int MPI::Request::Waitssome(int incount,
                                   MPI::Request array_of_requests[], int array_of_indices[],
                                   MPI::Status array_of_statuses[])
```

```
static int MPI::Request::Waitssome(int incount,
                                   MPI::Request array_of_requests[], int array_of_indices[])
```

Waits until at least one of the operations associated with active handles in the list have completed. Returns in `outcount` the number of requests from the list `array_of_requests` that have completed. Returns in the first `outcount` locations of the array `array_of_indices` the indices of these operations (index within the array `array_of_requests`; the array is indexed from zero in C and from one in Fortran). Returns in the first `outcount` locations of the array `array_of_status` the status for these completed operations. If a request that completed was allocated by a nonblocking communication call, then it is deallocated, and the associated handle is set to `MPI_REQUEST_NULL`.

If the list contains no active handles, then the call returns immediately with `outcount = MPI_UNDEFINED`.

When one or more of the communications completed by `MPI_WAITSSOME` fails, then it is desirable to return specific information on each communication. The arguments `outcount`, `array_of_indices` and `array_of_statuses` will be adjusted to indicate completion of all communications that have succeeded or failed. The call will return the error code `MPI_ERR_IN_STATUS` and the error field of each status returned will be set to indicate success or to indicate the specific error that occurred. The call will return `MPI_SUCCESS` if no request resulted in an error, and will return another error code if it failed for other reasons (such as invalid arguments). In such cases, it will not update the error fields of the statuses.

```

MPI_TESTSOME(incount, array_of_requests, outcount, array_of_indices, array_of_statuses) 1
IN      incount      length of array_of_requests (integer) 2
INOUT  array_of_requests  array of requests (array of handles) 3
OUT    outcount      number of completed requests (integer) 4
OUT    array_of_indices  array of indices of operations that completed (array of 5
                                integers) 6
OUT    array_of_statuses  array of status objects for operations that completed 7
                                (array of Status) 8
                                9
                                10
                                11
int MPI_Testsome(int incount, MPI_Request *array_of_requests, int *outcount, 12
                int *array_of_indices, MPI_Status *array_of_statuses) 13
MPI_TESTSOME(INCOUNT, ARRAY_OF_REQUESTS, OUTCOUNT, ARRAY_OF_INDICES, 14
            ARRAY_OF_STATUSES, IERROR) 15
    INTEGER INCOUNT, ARRAY_OF_REQUESTS(*), OUTCOUNT, ARRAY_OF_INDICES(*), 16
    ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR 17
static int MPI::Request::Testsome(int incount, 18
    MPI::Request array_of_requests[], int array_of_indices[], 19
    MPI::Status array_of_statuses[]) 20
static int MPI::Request::Testsome(int incount, 21
    MPI::Request array_of_requests[], int array_of_indices[]) 22

```

Behaves like `MPI_WAITSSOME`, except that it returns immediately. If no operation has completed it returns `outcount = 0`. If there is no active handle in the list it returns `outcount = MPI_UNDEFINED`.

`MPI_TESTSSOME` is a local operation, which returns immediately, whereas `MPI_WAITSSOME` will block until a communication completes, if it was passed a list that contains at least one active handle. Both calls fulfill a fairness requirement: If a request for a receive repeatedly appears in a list of requests passed to `MPI_WAITSSOME` or `MPI_TESTSSOME`, and a matching send has been posted, then the receive will eventually succeed, unless the send is satisfied by another receive; and similarly for send requests.

Errors that occur during the execution of `MPI_TESTSSOME` are handled as for `MPI_WAITSSOME`.

Advice to users. The use of `MPI_TESTSSOME` is likely to be more efficient than the use of `MPI_TESTANY`. The former returns information on all completed communications, with the latter, a new call is required for each communication that completes.

A server with multiple clients can use `MPI_WAITSSOME` so as not to starve any client. Clients send messages to the server with service requests. The server calls `MPI_WAITSSOME` with one receive request for each client, and then handles all receives that completed. If a call to `MPI_WAITANY` is used instead, then one client could starve while requests from another client always sneak in first. (*End of advice to users.*)

Advice to implementors. `MPI_TESTSSOME` should complete as many pending communications as possible. (*End of advice to implementors.*)

Example 3.16 Client-server code (starvation can occur).

```

1  CALL MPI_COMM_SIZE(comm, size, ierr)
2
3  CALL MPI_COMM_RANK(comm, rank, ierr)
4  IF(rank > 0) THEN          ! client code
5
6      DO WHILE(.TRUE.)
7          CALL MPI_ISEND(a, n, MPI_REAL, 0, tag, comm, request, ierr)
8          CALL MPI_WAIT(request, status, ierr)
9      END DO
10 ELSE
11     ! rank=0 -- server code
12     DO i=1, size-1
13         CALL MPI_Irecv(a(1,i), n, MPI_REAL, i tag,
14             comm, request_list(i), ierr)
15     END DO
16     DO WHILE(.TRUE.)
17         CALL MPI_WAITANY(size-1, request_list, index, status, ierr)
18         CALL DO_SERVICE(a(1,index)) ! handle one message
19         CALL MPI_Irecv(a(1, index), n, MPI_REAL, index, tag,
20             comm, request_list(index), ierr)
21     END DO
22 END IF

```

Example 3.17 Same code, using MPI_WAITSSOME.

```

24 CALL MPI_COMM_SIZE(comm, size, ierr)
25
26 CALL MPI_COMM_RANK(comm, rank, ierr)
27 IF(rank > 0) THEN          ! client code
28
29     DO WHILE(.TRUE.)
30         CALL MPI_ISEND(a, n, MPI_REAL, 0, tag, comm, request, ierr)
31         CALL MPI_WAIT(request, status, ierr)
32     END DO
33 ELSE
34     ! rank=0 -- server code
35     DO i=1, size-1
36         CALL MPI_Irecv(a(1,i), n, MPI_REAL, i, tag,
37             comm, request_list(i), ierr)
38     END DO
39     DO WHILE(.TRUE.)
40         CALL MPI_WAITSSOME(size, request_list, numdone,
41             indices, statuses, ierr)
42         DO i=1, numdone
43             CALL DO_SERVICE(a(1, indices(i)))
44             CALL MPI_Irecv(a(1, indices(i)), n, MPI_REAL, 0, tag,
45                 comm, request_list(indices(i)), ierr)
46         END DO
47     END DO
48 END IF

```

3.7.6 Non-destructive Test of status

This call is useful for accessing the information associated with a request, without freeing the request (in case the user is expected to access it later). It allows one to layer libraries more conveniently, since multiple layers of software may access the same completed request and extract from it the status information.

`MPI_REQUEST_GET_STATUS(request, flag, status)`

| | | |
|-----|---------|--|
| IN | request | request (handle) |
| OUT | flag | boolean flag, same as from <code>MPI_TEST</code> (logical) |
| OUT | status | <code>MPI_STATUS</code> object if flag is true (Status) |

```
int MPI_Request_get_status(MPI_Request request, int *flag,
                          MPI_Status *status)
```

```
MPI_REQUEST_GET_STATUS( REQUEST, FLAG, STATUS, IERROR)
    INTEGER REQUEST, STATUS(MPI_STATUS_SIZE), IERROR
    LOGICAL FLAG
```

```
bool MPI::Request::Get_status(MPI::Status& status) const
```

```
bool MPI::Request::Get_status() const
```

Sets `flag=true` if the operation is complete, and, if so, returns in `status` the request status. However, unlike `test` or `wait`, it does not deallocate or inactivate the request; a subsequent call to `test`, `wait` or `free` should be executed with that request. It sets `flag=false` if the operation is not complete.

3.8 Probe and Cancel

The `MPI_PROBE` and `MPI_IPROBE` operations allow incoming messages to be checked for, without actually receiving them. The user can then decide how to receive them, based on the information returned by the probe (basically, the information returned by `status`). In particular, the user may allocate memory for the receive buffer, according to the length of the probed message.

The `MPI_CANCEL` operation allows pending communications to be canceled. This is required for cleanup. Posting a send or a receive ties up user resources (send or receive buffers), and a cancel may be needed to free these resources gracefully.

```
1 MPI_IPROBE(source, tag, comm, flag, status)
```

```
2     IN     source           source rank, or MPI_ANY_SOURCE (integer)
3
4     IN     tag             tag value or MPI_ANY_TAG (integer)
5
6     IN     comm           communicator (handle)
7
8     OUT    flag           (logical)
9
10    OUT    status         status object (Status)
```

```
11 int MPI_Iprobe(int source, int tag, MPI_Comm comm, int *flag,
12               MPI_Status *status)
```

```
13 MPI_IPROBE(SOURCE, TAG, COMM, FLAG, STATUS, IERROR)
```

```
14     LOGICAL FLAG
```

```
15     INTEGER SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR
```

```
16 bool MPI::Comm::Iprobe(int source, int tag, MPI::Status& status) const
```

```
17 bool MPI::Comm::Iprobe(int source, int tag) const
```

19 MPI_IPROBE(source, tag, comm, flag, status) returns `flag = true` if there is a message
20 that can be received and that matches the pattern specified by the arguments `source`, `tag`,
21 and `comm`. The call matches the same message that would have been received by a call to
22 `MPI_RECV(..., source, tag, comm, status)` executed at the same point in the program, and
23 returns in `status` the same value that would have been returned by `MPI_RECV()`. Otherwise,
24 the call returns `flag = false`, and leaves `status` undefined.

25 If `MPI_IPROBE` returns `flag = true`, then the content of the status object can be sub-
26 sequently accessed as described in section 3.2.5 to find the source, tag and length of the
27 probed message.

28 A subsequent receive executed with the same communicator, and the source and tag
29 returned in `status` by `MPI_IPROBE` will receive the message that was matched by the probe,
30 if no other intervening receive occurs after the probe, and the send is not successfully
31 cancelled before the receive. If the receiving process is multi-threaded, it is the user's
32 responsibility to ensure that the last condition holds.

33 The `source` argument of `MPI_PROBE` can be `MPI_ANY_SOURCE`, and the `tag` argument
34 can be `MPI_ANY_TAG`, so that one can probe for messages from an arbitrary source and/or
35 with an arbitrary tag. However, a specific communication context must be provided with
36 the `comm` argument.

37 It is not necessary to receive a message immediately after it has been probed for, and
38 the same message may be probed for several times before it is received.

```
40 MPI_PROBE(source, tag, comm, status)
```

```
41
42    IN     source           source rank, or MPI_ANY_SOURCE (integer)
43
44    IN     tag             tag value, or MPI_ANY_TAG (integer)
45
46    IN     comm           communicator (handle)
47
48    OUT    status         status object (Status)
```

```
49 int MPI_Probe(int source, int tag, MPI_Comm comm, MPI_Status *status)
```

```

MPI_PROBE(SOURCE, TAG, COMM, STATUS, IERROR)
    INTEGER SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR
void MPI::Comm::Probe(int source, int tag, MPI::Status& status) const
void MPI::Comm::Probe(int source, int tag) const

```

MPI_PROBE behaves like MPI_IPROBE except that it is a blocking call that returns only after a matching message has been found.

The MPI implementation of MPI_PROBE and MPI_IPROBE needs to guarantee progress: if a call to MPI_PROBE has been issued by a process, and a send that matches the probe has been initiated by some process, then the call to MPI_PROBE will return, unless the message is received by another concurrent receive operation (that is executed by another thread at the probing process). Similarly, if a process busy waits with MPI_IPROBE and a matching message has been issued, then the call to MPI_IPROBE will eventually return `flag = true` unless the message is received by another concurrent receive operation.

Example 3.18 Use blocking probe to wait for an incoming message.

```

CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(i, 1, MPI_INTEGER, 2, 0, comm, ierr)
ELSE IF(rank.EQ.1) THEN
    CALL MPI_SEND(x, 1, MPI_REAL, 2, 0, comm, ierr)
ELSE ! rank.EQ.2
    DO i=1, 2
        CALL MPI_PROBE(MPI_ANY_SOURCE, 0,
                       comm, status, ierr)
        IF (status(MPI_SOURCE) .EQ. 0) THEN
100          CALL MPI_RECV(i, 1, MPI_INTEGER, 0, 0, comm, status, ierr)
        ELSE
200          CALL MPI_RECV(x, 1, MPI_REAL, 1, 0, comm, status, ierr)
        END IF
    END DO
END IF

```

Each message is received with the right type.

Example 3.19 A similar program to the previous example, but now it has a problem.

```

CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(i, 1, MPI_INTEGER, 2, 0, comm, ierr)
ELSE IF(rank.EQ.1) THEN
    CALL MPI_SEND(x, 1, MPI_REAL, 2, 0, comm, ierr)
ELSE
    DO i=1, 2
        CALL MPI_PROBE(MPI_ANY_SOURCE, 0,
                       comm, status, ierr)
        IF (status(MPI_SOURCE) .EQ. 0) THEN

```

```

1 100          CALL MPI_RECV(i, 1, MPI_INTEGER, MPI_ANY_SOURCE,
2                    0, comm, status, ierr)
3          ELSE
4 200          CALL MPI_RECV(x, 1, MPI_REAL, MPI_ANY_SOURCE,
5                    0, comm, status, ierr)
6          END IF
7          END DO
8          END IF
9

```

We slightly modified example 3.18, using `MPI_ANY_SOURCE` as the source argument in the two receive calls in statements labeled 100 and 200. The program is now incorrect: the receive operation may receive a message that is distinct from the message probed by the preceding call to `MPI_PROBE`.

Advice to implementors. A call to `MPI_PROBE(source, tag, comm, status)` will match the message that would have been received by a call to `MPI_RECV(..., source, tag, comm, status)` executed at the same point. Suppose that this message has source `s`, tag `t` and communicator `c`. If the tag argument in the probe call has value `MPI_ANY_TAG` then the message probed will be the earliest pending message from source `s` with communicator `c` and any tag; in any case, the message probed will be the earliest pending message from source `s` with tag `t` and communicator `c` (this is the message that would have been received, so as to preserve message order). This message continues as the earliest pending message from source `s` with tag `t` and communicator `c`, until it is received. A receive operation subsequent to the probe that uses the same communicator as the probe and uses the tag and source values returned by the probe, must receive this message, unless it has already been received by another receive operation. (*End of advice to implementors.*)

`MPI_CANCEL(request)`

IN request communication request (handle)

```
int MPI_Cancel(MPI_Request *request)
```

```
MPI_CANCEL(REQUEST, IERROR)
```

```
INTEGER REQUEST, IERROR
```

```
void MPI::Request::Cancel() const
```

A call to `MPI_CANCEL` marks for cancellation a pending, nonblocking communication operation (send or receive). The cancel call is local. It returns immediately, possibly before the communication is actually canceled. It is still necessary to complete a communication that has been marked for cancellation, using a call to `MPI_REQUEST_FREE`, `MPI_WAIT` or `MPI_TEST` (or any of the derived operations).

If a communication is marked for cancellation, then a `MPI_WAIT` call for that communication is guaranteed to return, irrespective of the activities of other processes (i.e., `MPI_WAIT` behaves as a local function); similarly if `MPI_TEST` is repeatedly called in a busy wait loop for a canceled communication, then `MPI_TEST` will eventually be successful.

`MPI_CANCEL` can be used to cancel a communication that uses a persistent request (see

Sec. 3.9), in the same way it is used for nonpersistent requests. A successful cancellation cancels the active communication, but not the request itself. After the call to `MPI_CANCEL` and the subsequent call to `MPI_WAIT` or `MPI_TEST`, the request becomes inactive and can be activated for a new communication.

The successful cancellation of a buffered send frees the buffer space occupied by the pending message.

Either the cancellation succeeds, or the communication succeeds, but not both. If a send is marked for cancellation, then it must be the case that either the send completes normally, in which case the message sent was received at the destination process, or that the send is successfully canceled, in which case no part of the message was received at the destination. Then, any matching receive has to be satisfied by another send. If a receive is marked for cancellation, then it must be the case that either the receive completes normally, or that the receive is successfully canceled, in which case no part of the receive buffer is altered. Then, any matching send has to be satisfied by another receive.

If the operation has been canceled, then information to that effect will be returned in the status argument of the operation that completes the communication.

`MPI_TEST_CANCELLED(status, flag)`

| | | |
|-----|--------|------------------------|
| IN | status | status object (Status) |
| OUT | flag | (logical) |

`int MPI_Test_cancelled(MPI_Status *status, int *flag)`

`MPI_TEST_CANCELLED(STATUS, FLAG, IERROR)`
`LOGICAL FLAG`
`INTEGER STATUS(MPI_STATUS_SIZE), IERROR`

`bool MPI::Status::Is_cancelled() const`

Returns `flag = true` if the communication associated with the status object was canceled successfully. In such a case, all other fields of `status` (such as `count` or `tag`) are undefined. Returns `flag = false`, otherwise. If a receive operation might be canceled then one should call `MPI_TEST_CANCELLED` first, to check whether the operation was canceled, before checking on the other fields of the return status.

Advice to users. Cancel can be an expensive operation that should be used only exceptionally. (*End of advice to users.*)

Advice to implementors. If a send operation uses an “eager” protocol (data is transferred to the receiver before a matching receive is posted), then the cancellation of this send may require communication with the intended receiver in order to free allocated buffers. On some systems this may require an interrupt to the intended receiver. Note that, while communication may be needed to implement `MPI_CANCEL`, this is still a local operation, since its completion does not depend on the code executed by other processes. If processing is required on another process, this should be transparent to the application (hence the need for an interrupt and an interrupt handler). (*End of advice to implementors.*)

3.9 Persistent communication requests

Often a communication with the same argument list is repeatedly executed within the inner loop of a parallel computation. In such a situation, it may be possible to optimize the communication by binding the list of communication arguments to a **persistent** communication request once and, then, repeatedly using the request to initiate and complete messages. The persistent request thus created can be thought of as a communication port or a “half-channel.” It does not provide the full functionality of a conventional channel, since there is no binding of the send port to the receive port. This construct allows reduction of the overhead for communication between the process and communication controller, but not of the overhead for communication between one communication controller and another. It is not necessary that messages sent with a persistent request be received by a receive operation using a persistent request, or vice versa.

A persistent communication request is created using one of the [five](#) following calls. These calls involve no communication.

```
MPI_SEND_INIT(buf, count, datatype, dest, tag, comm, request)
```

| | | |
|-----|----------|---|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements sent (integer) |
| IN | datatype | type of each element (handle) |
| IN | dest | rank of destination (integer) |
| IN | tag | message tag (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Send_init(void* buf, int count, MPI_Datatype datatype, int dest,
                 int tag, MPI_Comm comm, MPI_Request *request)
```

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

```
MPI::Prequest MPI::Comm::Send_init(const void* buf, int count, const
                                   MPI::Datatype& datatype, int dest, int tag) const
```

Creates a persistent communication request for a standard mode send operation, and binds to it all the arguments of a send operation.


```

1 MPI_RSEND_INIT(buf, count, datatype, dest, tag, comm, request)
2     IN      buf                initial address of send buffer (choice)
3
4     IN      count              number of elements sent (integer)
5
6     IN      datatype           type of each element (handle)
7
8     IN      dest               rank of destination (integer)
9
10    IN      tag                message tag (integer)
11
12    IN      comm               communicator (handle)
13
14    OUT     request            communication request (handle)
15
16 int MPI_Rsend_init(void* buf, int count, MPI_Datatype datatype, int dest,
17                   int tag, MPI_Comm comm, MPI_Request *request)
18
19 MPI_RSEND_INIT(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
20 <type> BUF(*)
21 INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR
22
23 MPI::Prequest MPI::Comm::Rsend_init(const void* buf, int count, const
24 MPI::Datatype& datatype, int dest, int tag) const

```

Creates a persistent communication object for a ready mode send operation.

```

25 MPI_RECV_INIT(buf, count, datatype, source, tag, comm, request)
26     OUT     buf                initial address of receive buffer (choice)
27
28     IN      count              number of elements received (integer)
29
30     IN      datatype           type of each element (handle)
31
32     IN      source             rank of source or MPI_ANY_SOURCE (integer)
33
34     IN      tag                message tag or MPI_ANY_TAG (integer)
35
36     IN      comm               communicator (handle)
37
38     OUT     request            communication request (handle)
39
40 int MPI_Recv_init(void* buf, int count, MPI_Datatype datatype, int source,
41                  int tag, MPI_Comm comm, MPI_Request *request)
42
43 MPI_RECV_INIT(BUF, COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR)
44 <type> BUF(*)
45 INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR
46
47 MPI::Prequest MPI::Comm::Recv_init(void* buf, int count, const
48 MPI::Datatype& datatype, int source, int tag) const

```

Creates a persistent communication request for a receive operation. The argument `buf` is marked as `OUT` because the user gives permission to write on the receive buffer by passing the argument to `MPI_RECV_INIT`.

A persistent communication request is inactive after it was created — no active communication is attached to the request.

A communication (send or receive) that uses a persistent request is initiated by the function `MPI_START`.

```
MPI_START(request)
```

```
    INOUT    request                communication request (handle)
```

```
int MPI_Start(MPI_Request *request)
```

```
MPI_START(REQUEST, IERROR)
```

```
    INTEGER REQUEST, IERROR
```

```
void MPI::Prequest::Start()
```

The argument, `request`, is a handle returned by one of the previous five calls. The associated request should be inactive. The request becomes active once the call is made.

If the request is for a send with ready mode, then a matching receive should be posted before the call is made. The communication buffer should not be accessed after the call, and until the operation completes.

The call is local, with similar semantics to the nonblocking communication operations described in section 3.7. That is, a call to `MPI_START` with a request created by `MPI_SEND_INIT` starts a communication in the same manner as a call to `MPI_ISEND`; a call to `MPI_START` with a request created by `MPI_BSEND_INIT` starts a communication in the same manner as a call to `MPI_IBSEND`; and so on.

```
MPI_STARTALL(count, array_of_requests)
```

```
    IN        count                list length (integer)
```

```
    INOUT    array_of_requests      array of requests (array of handle)
```

```
int MPI_Startall(int count, MPI_Request *array_of_requests)
```

```
MPI_STARTALL(COUNT, ARRAY_OF_REQUESTS, IERROR)
```

```
    INTEGER COUNT, ARRAY_OF_REQUESTS(*), IERROR
```

```
static void MPI::Prequest::Startall(int count,
                                     MPI::Prequest array_of_requests[])
```

Start all communications associated with requests in `array_of_requests`. A call to `MPI_STARTALL(count, array_of_requests)` has the same effect as calls to `MPI_START (&array_of_requests[i])`, executed for $i=0, \dots, \text{count}-1$, in some arbitrary order.

A communication started with a call to `MPI_START` or `MPI_STARTALL` is completed by a call to `MPI_WAIT`, `MPI_TEST`, or one of the derived functions described in section 3.7.5. The request becomes inactive after successful completion of such call. The request is not deallocated and it can be activated anew by an `MPI_START` or `MPI_STARTALL` call.

A persistent request is deallocated by a call to `MPI_REQUEST_FREE` (Section 3.7.3).

The call to `MPI_REQUEST_FREE` can occur at any point in the program after the persistent request was created. However, the request will be deallocated only after it becomes inactive. Active receive requests should not be freed. Otherwise, it will not be possible to check that the receive has completed. It is preferable, in general, to free requests when

1 they are inactive. If this rule is followed, then the functions described in this section will
2 be invoked in a sequence of the form,

3
4 **Create (Start Complete)* Free**

5
6 **where** * indicates zero or more repetitions. If the same communication object is used in
7 several concurrent threads, it is the user's responsibility to coordinate calls so that the
8 correct sequence is obeyed.

9 A send operation initiated with MPI_START can be matched with any receive operation
10 and, likewise, a receive operation initiated with MPI_START can receive messages generated
11 by any send operation.

12
13 *Advice to users.* To prevent problems with the argument copying and register opti-
14 mization done by Fortran compilers, please note the hints in subsections "Problems
15 Due to Data Copying and Sequence Association," and "A Problem with Register
16 Optimization" in [Section 13.2.2 on pages 451 and 454](#). (*End of advice to users.*)

17
18
19 **3.10 Send-receive**

20
21 The **send-receive** operations combine in one call the sending of a message to one desti-
22 nation and the receiving of another message, from another process. The two (source and
23 destination) are possibly the same. A send-receive operation is very useful for executing
24 a shift operation across a chain of processes. If blocking sends and receives are used for
25 such a shift, then one needs to order the sends and receives correctly (for example, even
26 processes send, then receive, odd processes receive first, then send) so as to prevent cyclic
27 dependencies that may lead to deadlock. When a send-receive operation is used, the com-
28 munication subsystem takes care of these issues. The send-receive operation can be used
29 in conjunction with the functions described in Chapter 6 in order to perform shifts on var-
30 ious logical topologies. Also, a send-receive operation is useful for implementing remote
31 procedure calls.

32 A message sent by a send-receive operation can be received by a regular receive oper-
33 ation or probed by a probe operation; a send-receive operation can receive a message sent
34 by a regular send operation.

35
36
37
38
39
40
41
42
43
44
45
46
47
48

```

MPI_SENDRECV(sendbuf, sendcount, sendtype, dest, sendtag, recvbuf, recvcount, recvtype,
source, recvtag, comm, status)
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

```

| | | |
|-----|-----------|--|
| IN | sendbuf | initial address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (integer) |
| IN | sendtype | type of elements in send buffer (handle) |
| IN | dest | rank of destination (integer) |
| IN | sendtag | send tag (integer) |
| OUT | recvbuf | initial address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (integer) |
| IN | recvtype | type of elements in receive buffer (handle) |
| IN | source | rank of source (integer) |
| IN | recvtag | receive tag (integer) |
| IN | comm | communicator (handle) |
| OUT | status | status object (Status) |

```

int MPI_Sendrecv(void *sendbuf, int sendcount, MPI_Datatype sendtype,
                int dest, int sendtag, void *recvbuf, int recvcount,
                MPI_Datatype recvtype, int source, int recvtag, MPI_Comm comm,
                MPI_Status *status)
MPI_SENDRECV(SENDBUF, SENDCOUNT, SENDTYPE, DEST, SENDTAG, RECVBUF,
                RECVCOUNT, RECVTYPE, SOURCE, RECVTAG, COMM, STATUS, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER SENDCOUNT, SENDTYPE, DEST, SENDTAG, RECVCOUNT, RECVTYPE,
SOURCE, RECVTAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR
void MPI::Comm::Sendrecv(const void *sendbuf, int sendcount, const
                        MPI::Datatype& sendtype, int dest, int sendtag, void *recvbuf,
                        int recvcount, const MPI::Datatype& recvtype, int source,
                        int recvtag, MPI::Status& status) const
void MPI::Comm::Sendrecv(const void *sendbuf, int sendcount, const
                        MPI::Datatype& sendtype, int dest, int sendtag, void *recvbuf,
                        int recvcount, const MPI::Datatype& recvtype, int source,
                        int recvtag) const

```

Execute a blocking send and receive operation. Both send and receive use the same communicator, but possibly different tags. The send buffer and receive buffers must be disjoint, and may have different lengths and datatypes.

The semantics of a send-receive operation is what would be obtained if the caller forked two concurrent threads, one to execute the send, and one to execute the receive, followed by a join of these two threads.

```

1 MPI_SENDRECV_REPLACE(buf, count, datatype, dest, sendtag, source, recvtag, comm, sta-
2 tus)
3     INOUT   buf                initial address of send and receive buffer (choice)
4     IN      count              number of elements in send and receive buffer (integer)
5     IN      datatype           type of elements in send and receive buffer (handle)
6     IN      dest               rank of destination (integer)
7     IN      sendtag            send message tag (integer)
8     IN      source             rank of source (integer)
9     IN      recvtag            receive message tag (integer)
10    IN      comm               communicator (handle)
11    IN      status              status object (Status)
12
13
14
15
16 int MPI_Sendrecv_replace(void* buf, int count, MPI_Datatype datatype,
17                          int dest, int sendtag, int source, int recvtag, MPI_Comm comm,
18                          MPI_Status *status)
19
20 MPI_SENDRECV_REPLACE(BUF, COUNT, DATATYPE, DEST, SENDTAG, SOURCE, RECVTAG,
21                      COMM, STATUS, IERROR)
22     <type> BUF(*)
23     INTEGER COUNT, DATATYPE, DEST, SENDTAG, SOURCE, RECVTAG, COMM,
24     STATUS(MPI_STATUS_SIZE), IERROR
25
26 void MPI::Comm::Sendrecv_replace(void* buf, int count, const
27     MPI::Datatype& datatype, int dest, int sendtag, int source,
28     int recvtag, MPI::Status& status) const
29
30 void MPI::Comm::Sendrecv_replace(void* buf, int count, const
31     MPI::Datatype& datatype, int dest, int sendtag, int source,
32     int recvtag) const

```

Execute a blocking send and receive. The same buffer is used both for the send and for the receive, so that the message sent is replaced by the message received.

Advice to implementors. Additional intermediate buffering is needed for the “replace” variant. (*End of advice to implementors.*)

3.11 Null processes

In many instances, it is convenient to specify a “dummy” source or destination for communication. This simplifies the code that is needed for dealing with boundaries, for example, in the case of a non-circular shift done with calls to send-receive.

The special value `MPI_PROC_NULL` can be used instead of a rank wherever a source or a destination argument is required in a call. A communication with process `MPI_PROC_NULL` has no effect. A send to `MPI_PROC_NULL` succeeds and returns as soon as possible. A receive from `MPI_PROC_NULL` succeeds and returns as soon as possible with no modifications to the receive buffer. When a receive with `source = MPI_PROC_NULL` is executed then the status object returns `source = MPI_PROC_NULL`, `tag = MPI_ANY_TAG` and `count = 0`.

3.12 Derived datatypes

Up to here, all point to point communication have involved only contiguous buffers containing a sequence of elements of the same type. This is too constraining on two accounts. One often wants to pass messages that contain values with different datatypes (e.g., an integer count, followed by a sequence of real numbers); and one often wants to send noncontiguous data (e.g., a sub-block of a matrix). One solution is to pack noncontiguous data into a contiguous buffer at the sender site and unpack it back at the receiver site. This has the disadvantage of requiring additional memory-to-memory copy operations at both sites, even when the communication subsystem has scatter-gather capabilities. Instead, MPI provides mechanisms to specify more general, mixed, and noncontiguous communication buffers. It is up to the implementation to decide whether data should be first packed in a contiguous buffer before being transmitted, or whether it can be collected directly from where it resides.

The general mechanisms provided here allow one to transfer directly, without copying, objects of various shape and size. It is not assumed that the MPI library is cognizant of the objects declared in the host language. Thus, if one wants to transfer a structure, or an array section, it will be necessary to provide in MPI a definition of a communication buffer that mimics the definition of the structure or array section in question. These facilities can be used by library designers to define communication functions that can transfer objects defined in the host language — by decoding their definitions as available in a symbol table or a dope vector. Such higher-level communication functions are not part of MPI.

More general communication buffers are specified by replacing the basic datatypes that have been used so far with derived datatypes that are constructed from basic datatypes using the constructors described in this section. These methods of constructing derived datatypes can be applied recursively.

A **general datatype** is an opaque object that specifies two things:

- A sequence of basic datatypes
- A sequence of integer (byte) displacements

The displacements are not required to be positive, distinct, or in increasing order. Therefore, the order of items need not coincide with their order in store, and an item may appear more than once. We call such a pair of sequences (or sequence of pairs) a **type map**. The sequence of basic datatypes (displacements ignored) is the **type signature** of the datatype.

Let

$$Typemap = \{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

be such a type map, where $type_i$ are basic types, and $disp_i$ are displacements. Let

$$Typesig = \{type_0, \dots, type_{n-1}\}$$

be the associated type signature. This type map, together with a base address buf , specifies a communication buffer: the communication buffer that consists of n entries, where the i -th entry is at address $buf + disp_i$ and has type $type_i$. A message assembled from such a communication buffer will consist of n values, of the types defined by $Typesig$.

Most datatype constructors have replication count or block length arguments. Allowed values are nonnegative integers. If the value is zero, no elements are generated in the type map and there is no effect on datatype bounds or extent.

We can use a handle to a general datatype as an argument in a send or receive operation, instead of a basic datatype argument. The operation `MPI_SEND(buf, 1, datatype,...)` will use the send buffer defined by the base address `buf` and the general datatype associated with `datatype`; it will generate a message with the type signature determined by the `datatype` argument. `MPI_RECV(buf, 1, datatype,...)` will use the receive buffer defined by the base address `buf` and the general datatype associated with `datatype`.

General datatypes can be used in all send and receive operations. We discuss, in Sec. 3.12.12, the case where the second argument `count` has value > 1 .

The basic datatypes presented in section 3.2.2 are particular cases of a general datatype, and are predefined. Thus, `MPI_INT` is a predefined handle to a datatype with type map $\{(int, 0)\}$, with one entry of type `int` and displacement zero. The other basic datatypes are similar.

The **extent** of a datatype is defined to be the span from the first byte to the last byte occupied by entries in this datatype, rounded up to satisfy alignment requirements. That is, if

$$Typemap = \{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

then

$$\begin{aligned} lb(Typemap) &= \min_j disp_j, \\ ub(Typemap) &= \max_j (disp_j + sizeof(type_j)) + \epsilon, \text{ and} \\ extent(Typemap) &= ub(Typemap) - lb(Typemap). \end{aligned} \tag{3.1}$$

If $type_i$ requires alignment to a byte address that is a multiple of k_i , then ϵ is the least nonnegative increment needed to round $extent(Typemap)$ to the next multiple of $\max_i k_i$. The complete definition of **extent** is given on page 94.

Example 3.20 Assume that $Type = \{(double, 0), (char, 8)\}$ (a `double` at displacement zero, followed by a `char` at displacement eight). Assume, furthermore, that doubles have to be strictly aligned at addresses that are multiples of eight. Then, the extent of this datatype is 16 (9 rounded to the next multiple of 8). A datatype that consists of a character immediately followed by a double will also have an extent of 16.

Rationale. The definition of extent is motivated by the assumption that the amount of padding added at the end of each structure in an array of structures is the least needed to fulfill alignment constraints. More explicit control of the extent is provided in section 3.12.7. Such explicit control is needed in cases where the assumption does not hold, for example, where union types are used. (*End of rationale.*)

3.12.1 New Datatype Manipulation Functions

New functions are provided to supplement the type manipulation functions that have address sized integer arguments. The new functions will use, in their Fortran binding, address-sized `INTEGERS`, thus solving problems currently encountered when the application address range is $> 2^{32}$. Also, a new, more convenient type constructor is provided to modify the lower bound and extent of a datatype. The deprecated functions replaced by the new functions here are listed in Section 2.6.1.

3.12.2 Type Constructors with Explicit Addresses

The four functions `MPI_TYPE_CREATE_HVECTOR`, `MPI_TYPE_CREATE_HINDEXED`, `MPI_TYPE_CREATE_STRUCT`, and `MPI_GET_ADDRESS` supplement the four corresponding type constructor functions from MPI-1. The new functions are synonymous with the old functions in C/C++, or on Fortran systems where default `INTEGER`s are address sized. (The old names are not available in C++.) In Fortran, these functions accept arguments of type `INTEGER(KIND=MPI_ADDRESS_KIND)`, wherever arguments of type `MPI_Aint` are used in C. On Fortran 77 systems that do not support the Fortran 90 `KIND` notation, and where addresses are 64 bits whereas default `INTEGER`s are 32 bits, these arguments will be of type `INTEGER*8`. The old functions will continue to be provided for backward compatibility. However, users are encouraged to switch to the new functions, in both Fortran and C.

The use of the old functions is deprecated.

3.12.3 Datatype constructors

Contiguous The simplest datatype constructor is `MPI_TYPE_CONTIGUOUS` which allows replication of a datatype into contiguous locations.

`MPI_TYPE_CONTIGUOUS(count, oldtype, newtype)`

| | | |
|-----|---------|---|
| IN | count | replication count (nonnegative integer) |
| IN | oldtype | old datatype (handle) |
| OUT | newtype | new datatype (handle) |

```
int MPI_Type_contiguous(int count, MPI_Datatype oldtype,
                       MPI_Datatype *newtype)
```

```
MPI_TYPE_CONTIGUOUS(COUNT, OLDTYPE, NEWTYPE, IERROR)
INTEGER COUNT, OLDTYPE, NEWTYPE, IERROR
```

```
MPI::Datatype MPI::Datatype::Create_contiguous(int count) const
```

`newtype` is the datatype obtained by concatenating `count` copies of `oldtype`. Concatenation is defined using *extent* as the size of the concatenated copies.

Example 3.21 Let `oldtype` have type map $\{(double, 0), (char, 8)\}$, with extent 16, and let `count = 3`. The type map of the datatype returned by `newtype` is

$$\{(double, 0), (char, 8), (double, 16), (char, 24), (double, 32), (char, 40)\};$$

i.e., alternating `double` and `char` elements, with displacements 0, 8, 16, 24, 32, 40.

In general, assume that the type map of `oldtype` is

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

with extent *ex*. Then `newtype` has a type map with `count · n` entries defined by:

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1}), (type_0, disp_0 + ex), \dots, (type_{n-1}, disp_{n-1} + ex),$$

1 ..., (type₀, disp₀ + ex · (count - 1)), ..., (type_{n-1}, disp_{n-1} + ex · (count - 1))}.

2
3
4
5 **Vector** The function MPI_TYPE_VECTOR is a more general constructor that allows repli-
6 cation of a datatype into locations that consist of equally spaced blocks. Each block is
7 obtained by concatenating the same number of copies of the old datatype. The spacing
8 between blocks is a multiple of the extent of the old datatype.
9

10
11 MPI_TYPE_VECTOR(count, blocklength, stride, oldtype, newtype)

| | | | |
|----|-----|-------------|--|
| 12 | IN | count | number of blocks (nonnegative integer) |
| 13 | IN | blocklength | number of elements in each block (nonnegative integer) |
| 14 | | | ger) |
| 15 | | | |
| 16 | IN | stride | number of elements between start of each block (integer) |
| 17 | | | ger) |
| 18 | | | |
| 19 | IN | oldtype | old datatype (handle) |
| 20 | OUT | newtype | new datatype (handle) |

21
22 `int MPI_Type_vector(int count, int blocklength, int stride,`
23 `MPI_Datatype oldtype, MPI_Datatype *newtype)`

24 `MPI_TYPE_VECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR)`
25 `INTEGER COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR`

26
27 `MPI::Datatype MPI::Datatype::Create_vector(int count, int blocklength,`
28 `int stride) const`

29
30 **Example 3.22** Assume, again, that oldtype has type map {(double, 0), (char, 8)}, with extent 16. A call to MPI_TYPE_VECTOR(2, 3, 4, oldtype, newtype) will create the datatype with type map,
31
32
33

34 {(double, 0), (char, 8), (double, 16), (char, 24), (double, 32), (char, 40),
35
36 (double, 64), (char, 72), (double, 80), (char, 88), (double, 96), (char, 104)}.

37
38 That is, two blocks with three copies each of the old type, with a stride of 4 elements (4 · 16
39 bytes) between the blocks.

40 **Example 3.23** A call to MPI_TYPE_VECTOR(3, 1, -2, oldtype, newtype) will create the
41 datatype,
42

43 {(double, 0), (char, 8), (double, -32), (char, -24), (double, -64), (char, -56)}.

44
45
46 In general, assume that oldtype has type map,
47

48 {(type₀, disp₀), ..., (type_{n-1}, disp_{n-1})},

This function replaces `MPI_TYPE_HVECTOR`, whose use is deprecated. See also Chapter 15.

Assume that `oldtype` has type map,

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

with extent ex . Let `bl` be the blocklength. The newly created datatype has a type map with $count \cdot bl \cdot n$ entries:

$$\begin{aligned} &\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1}), \\ &(type_0, disp_0 + ex), \dots, (type_{n-1}, disp_{n-1} + ex), \dots, \\ &(type_0, disp_0 + (bl - 1) \cdot ex), \dots, (type_{n-1}, disp_{n-1} + (bl - 1) \cdot ex), \\ &(type_0, disp_0 + stride), \dots, (type_{n-1}, disp_{n-1} + stride), \dots, \\ &(type_0, disp_0 + stride + (bl - 1) \cdot ex), \dots, \\ &(type_{n-1}, disp_{n-1} + stride + (bl - 1) \cdot ex), \dots, \\ &(type_0, disp_0 + stride \cdot (count - 1)), \dots, (type_{n-1}, disp_{n-1} + stride \cdot (count - 1)), \dots, \\ &(type_0, disp_0 + stride \cdot (count - 1) + (bl - 1) \cdot ex), \dots, \\ &(type_{n-1}, disp_{n-1} + stride \cdot (count - 1) + (bl - 1) \cdot ex)\}. \end{aligned}$$

Indexed The function `MPI_TYPE_INDEXED` allows replication of an old datatype into a sequence of blocks (each block is a concatenation of the old datatype), where each block can contain a different number of copies and have a different displacement. All block displacements are multiples of the old type extent.

`MPI_TYPE_INDEXED(count, array_of_blocklengths, array_of_displacements, oldtype, newtype)`

| | | |
|-----|------------------------|--|
| IN | count | number of blocks – also number of entries in <code>array_of_displacements</code> and <code>array_of_blocklengths</code> (non-negative integer) |
| IN | array_of_blocklengths | number of elements per block (array of nonnegative integers) |
| IN | array_of_displacements | displacement for each block, in multiples of <code>oldtype</code> extent (array of integer) |
| IN | oldtype | old datatype (handle) |
| OUT | newtype | new datatype (handle) |

```
int MPI_Type_indexed(int count, int *array_of_blocklengths,
                    int *array_of_displacements, MPI_Datatype oldtype,
                    MPI_Datatype *newtype)
```

```

MPI_TYPE_INDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
                 OLDTYPE, NEWTYPE, IERROR)
INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*),
OLDTYPE, NEWTYPE, IERROR

```

```

MPI::Datatype MPI::Datatype::Create_indexed(int count,
      const int array_of_blocklengths[],
      const int array_of_displacements[]) const

```

Example 3.24 Let `oldtype` have type map $\{(double, 0), (char, 8)\}$, with extent 16. Let $\mathbf{B} = (3, 1)$ and let $\mathbf{D} = (4, 0)$. A call to `MPI_TYPE_INDEXED(2, B, D, oldtype, newtype)` returns a datatype with type map,

$$\{(double, 64), (char, 72), (double, 80), (char, 88), (double, 96), (char, 104), (double, 0), (char, 8)\}.$$

That is, three copies of the old type starting at displacement 64, and one copy starting at displacement 0.

In general, assume that `oldtype` has type map,

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

with extent ex . Let \mathbf{B} be the `array_of_blocklength` argument and

\mathbf{D} be the `array_of_displacements` argument. The newly created datatype has $n \cdot \sum_{i=0}^{\text{count}-1} B[i]$ entries:

$$\begin{aligned} &\{(type_0, disp_0 + D[0] \cdot ex), \dots, (type_{n-1}, disp_{n-1} + D[0] \cdot ex), \dots, \\ &(type_0, disp_0 + (D[0] + B[0] - 1) \cdot ex), \dots, (type_{n-1}, disp_{n-1} + (D[0] + B[0] - 1) \cdot ex), \dots, \\ &(type_0, disp_0 + D[\text{count} - 1] \cdot ex), \dots, (type_{n-1}, disp_{n-1} + D[\text{count} - 1] \cdot ex), \dots, \\ &(type_0, disp_0 + (D[\text{count} - 1] + B[\text{count} - 1] - 1) \cdot ex), \dots, \\ &(type_{n-1}, disp_{n-1} + (D[\text{count} - 1] + B[\text{count} - 1] - 1) \cdot ex)\}. \end{aligned}$$

A call to `MPI_TYPE_VECTOR(count, blocklength, stride, oldtype, newtype)` is equivalent to a call to `MPI_TYPE_INDEXED(count, B, D, oldtype, newtype)` where

$$D[j] = j \cdot \text{stride}, \quad j = 0, \dots, \text{count} - 1,$$

and

$$B[j] = \text{blocklength}, \quad j = 0, \dots, \text{count} - 1.$$

Hindexed The function `MPI_TYPE_CREATE_HINDEXED` is identical to `MPI_TYPE_INDEXED`, except that block displacements in `array_of_displacements` are specified in bytes, rather than in multiples of the `oldtype` extent.

```
MPI_TYPE_CREATE_HINDEXED( count, array_of_blocklengths, array_of_displacements, old-
type, newtype)
```

| | | | |
|----|-----|------------------------|---|
| 8 | IN | count | number of blocks — also number of entries in array_of_displacements and array_of_blocklengths (non- negative integer) |
| 11 | IN | array_of_blocklengths | number of elements in each block (array of nonnega- tive integers) |
| 14 | IN | array_of_displacements | byte displacement of each block (array of integer) |
| 15 | IN | oldtype | old datatype (handle) |
| 16 | OUT | newtype | new datatype (handle) |

```
int MPI_Type_create_hindexed(int count, int array_of_blocklengths[],
                             MPI_Aint array_of_displacements[], MPI_Datatype oldtype,
                             MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_HINDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS,
                          ARRAY_OF_DISPLACEMENTS, OLDTYPE, NEWTYPE, IERROR)
INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), OLDTYPE, NEWTYPE, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_DISPLACEMENTS(*)
```

```
MPI::Datatype MPI::Datatype::Create_hindexed(int count,
                                               const int array_of_blocklengths[],
                                               const MPI::Aint array_of_displacements[]) const
```

This function replaces `MPI_TYPE_HINDEXED`, whose use is deprecated. See also Chapter 15.

Assume that `oldtype` has type map,

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

with extent ex . Let `B` be the `array_of_blocklength` argument and `D` be the `array_of_displacements` argument. The newly created datatype has a type map with $n \cdot \sum_{i=0}^{count-1} B[i]$ entries:

$$\{(type_0, disp_0 + D[0]), \dots, (type_{n-1}, disp_{n-1} + D[0]), \dots,$$

$$(type_0, disp_0 + D[0] + (B[0] - 1) \cdot ex), \dots,$$

$$(type_{n-1}, disp_{n-1} + D[0] + (B[0] - 1) \cdot ex), \dots,$$

$$(type_0, disp_0 + D[count - 1]), \dots, (type_{n-1}, disp_{n-1} + D[count - 1]), \dots,$$

$$(type_0, disp_0 + D[count - 1] + (B[count - 1] - 1) \cdot ex), \dots,$$

$$(type_{n-1}, disp_{n-1} + D[count - 1] + (B[count - 1] - 1) \cdot ex)\}.$$

`Indexed_block` This function is the same as `MPI_TYPE_INDEXED` except that the block-length is the same for all blocks. There are many codes using indirect addressing arising from unstructured grids where the blocksize is always 1 (gather/scatter). The following convenience function allows for constant blocksize and arbitrary displacements.

`MPI_TYPE_CREATE_INDEXED_BLOCK(count, blocklength, array_of_displacements, oldtype, newtype)`

| | | |
|-----|------------------------|--|
| IN | count | length of array of displacements (integer) |
| IN | blocklength | size of block (integer) |
| IN | array_of_displacements | array of displacements (array of integer) |
| IN | oldtype | old datatype (handle) |
| OUT | newtype | new datatype (handle) |

```
int MPI_Type_create_indexed_block(int count, int blocklength,
    int array_of_displacements[], MPI_Datatype oldtype,
    MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_INDEXED_BLOCK(COUNT, BLOCKLENGTH, ARRAY_OF_DISPLACEMENTS,
    OLDTYPE, NEWTYPE, IERROR)
    INTEGER COUNT, BLOCKLENGTH, ARRAY_OF_DISPLACEMENTS(*), OLDTYPE,
    NEWTYPE, IERROR
```

```
MPI::Datatype MPI::Datatype::Create_indexed_block( int count,
    int blocklength, const int array_of_displacements[]) const
```

Struct `MPI_TYPE_STRUCT` is the most general type constructor. It further generalizes `MPI_TYPE_CREATE_HINDEXED` in that it allows each block to consist of replications of different datatypes.

```

1 MPI_TYPE_CREATE_STRUCT(count, array_of_blocklengths, array_of_displacements,
2 array_of_types, newtype)
3     IN      count                number of blocks (nonnegative integer) — also number
4                                     of entries in arrays array_of_types,
5                                     array_of_displacements and array_of_blocklengths
6
7     IN      array_of_blocklength  number of elements in each block (array of nonnega-
8                                     tive integer)
9
10    IN      array_of_displacements byte displacement of each block (array of integer)
11
12    IN      array_of_types        type of elements in each block (array of handles to
13                                     datatype objects)
14
15    OUT     newtype               new datatype (handle)

```

```

15 int MPI_Type_create_struct(int count, int array_of_blocklengths[],
16                           MPI_Aint array_of_displacements[],
17                           MPI_Datatype array_of_types[], MPI_Datatype *newtype)
18
19 MPI_TYPE_CREATE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
20                        ARRAY_OF_TYPES, NEWTYPE, IERROR)
21     INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_TYPES(*), NEWTYPE,
22     IERROR
23     INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_DISPLACEMENTS(*)
24
25 static MPI::Datatype MPI::Datatype::Create_struct(int count,
26           const int array_of_blocklengths[], const MPI::Aint
27           array_of_displacements[], const MPI::Datatype array_of_types[])

```

This function replaces `MPI_TYPE_STRUCT`, whose use is deprecated. See also Chapter 15.

Example 3.25 Let `type1` have type map,

$$\{(\text{double}, 0), (\text{char}, 8)\},$$

with extent 16. Let $B = (2, 1, 3)$, $D = (0, 16, 26)$, and $T = (\text{MPI_FLOAT}, \text{type1}, \text{MPI_CHAR})$. Then a call to `MPI_TYPE_STRUCT(3, B, D, T, newtype)` returns a datatype with type map,

$$\{(\text{float}, 0), (\text{float}, 4), (\text{double}, 16), (\text{char}, 24), (\text{char}, 26), (\text{char}, 27), (\text{char}, 28)\}.$$

That is, two copies of `MPI_FLOAT` starting at 0, followed by one copy of `type1` starting at 16, followed by three copies of `MPI_CHAR`, starting at 26. (We assume that a float occupies four bytes.)

In general, let T be the `array_of_types` argument, where $T[i]$ is a handle to,

$$\text{typemap}_i = \{(\text{type}_0^i, \text{disp}_0^i), \dots, (\text{type}_{n_i-1}^i, \text{disp}_{n_i-1}^i)\},$$

with extent ex_i . Let B be the `array_of_blocklength` argument and D be the `array_of_displacements` argument. Let c be the `count` argument. Then the newly created datatype has a type map with $\sum_{i=0}^{c-1} B[i] \cdot n_i$ entries:

$$\{(\text{type}_0^0, \text{disp}_0^0 + D[0]), \dots, (\text{type}_{n_0}^0, \text{disp}_{n_0}^0 + D[0]), \dots,$$

$$\begin{aligned}
 & (type_0^0, disp_0^0 + D[0] + (B[0] - 1) \cdot ex_0), \dots, (type_{n_0}^0, disp_{n_0}^0 + D[0] + (B[0] - 1) \cdot ex_0), \dots, \\
 & (type_0^{c-1}, disp_0^{c-1} + D[c - 1]), \dots, (type_{n_{c-1}-1}^{c-1}, disp_{n_{c-1}-1}^{c-1} + D[c - 1]), \dots, \\
 & (type_0^{c-1}, disp_0^{c-1} + D[c - 1] + (B[c - 1] - 1) \cdot ex_{c-1}), \dots, \\
 & (type_{n_{c-1}-1}^{c-1}, disp_{n_{c-1}-1}^{c-1} + D[c - 1] + (B[c - 1] - 1) \cdot ex_{c-1}).
 \end{aligned}$$

A call to `MPI_TYPE_HINDEXED(count, B, D, oldtype, newtype)` is equivalent to a call to `MPI_TYPE_STRUCT(count, B, D, T, newtype)`, where each entry of `T` is equal to `oldtype`.

3.12.4 Subarray Datatype Constructor

`MPI_TYPE_CREATE_SUBARRAY(ndims, array_of_sizes, array_of_subsizes, array_of_starts, order, oldtype, newtype)`

| | | |
|-----|--------------------------------|--|
| IN | <code>ndims</code> | number of array dimensions (positive integer) |
| IN | <code>array_of_sizes</code> | number of elements of type <code>oldtype</code> in each dimension of the full array (array of positive integers) |
| IN | <code>array_of_subsizes</code> | number of elements of type <code>oldtype</code> in each dimension of the subarray (array of positive integers) |
| IN | <code>array_of_starts</code> | starting coordinates of the subarray in each dimension (array of nonnegative integers) |
| IN | <code>order</code> | array storage order flag (state) |
| IN | <code>oldtype</code> | array element datatype (handle) |
| OUT | <code>newtype</code> | new datatype (handle) |

```
int MPI_Type_create_subarray(int ndims, int array_of_sizes[],
                           int array_of_subsizes[], int array_of_starts[], int order,
                           MPI_Datatype oldtype, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_SUBARRAY(NDIMS, ARRAY_OF_SIZES, ARRAY_OF_SUBSIZES,
                          ARRAY_OF_STARTS, ORDER, OLDTYPE, NEWTYPE, IERROR)
INTEGER NDIMS, ARRAY_OF_SIZES(*), ARRAY_OF_SUBSIZES(*),
ARRAY_OF_STARTS(*), ORDER, OLDTYPE, NEWTYPE, IERROR
```

```
MPI::Datatype MPI::Datatype::Create_subarray(int ndims,
                                              const int array_of_sizes[], const int array_of_subsizes[],
                                              const int array_of_starts[], int order) const
```

The subarray type constructor creates an MPI datatype describing an n-dimensional subarray of an n-dimensional array. The subarray may be situated anywhere within the full array, and may be of any nonzero size up to the size of the larger array as long as it is confined within this array. This type constructor facilitates creating filetypes to access arrays distributed in blocks among processes to a single file that contains the global array.

This type constructor can handle arrays with an arbitrary number of dimensions and works for both C and Fortran ordered matrices (i.e., row-major or column-major). Note that a C program may use Fortran order and a Fortran program may use C order.

The `ndims` parameter specifies the number of dimensions in the full data array and gives the number of elements in `array_of_sizes`, `array_of_subsizes`, and `array_of_starts`.

The number of elements of type `oldtype` in each dimension of the n -dimensional array and the requested subarray are specified by `array_of_sizes` and `array_of_subsizes`, respectively. For any dimension i , it is erroneous to specify `array_of_subsizes[i] < 1` or `array_of_subsizes[i] > array_of_sizes[i]`.

The `array_of_starts` contains the starting coordinates of each dimension of the subarray. Arrays are assumed to be indexed starting from zero. For any dimension i , it is erroneous to specify `array_of_starts[i] < 0` or `array_of_starts[i] > (array_of_sizes[i] - array_of_subsizes[i])`.

Advice to users. In a Fortran program with arrays indexed starting from 1, if the starting coordinate of a particular dimension of the subarray is n , then the entry in `array_of_starts` for that dimension is $n-1$. (*End of advice to users.*)

The `order` argument specifies the storage order for the subarray as well as the full array. It must be set to one of the following:

`MPI_ORDER_C` The ordering used by C arrays, (i.e., row-major order)

`MPI_ORDER_FORTRAN` The ordering used by Fortran arrays, (i.e., column-major order)

A $ndims$ -dimensional subarray (`newtype`) with no extra padding can be defined by the function `Subarray()` as follows:

$$\begin{aligned} \text{newtype} = & \text{Subarray}(ndims, \{size_0, size_1, \dots, size_{ndims-1}\}, \\ & \{subsize_0, subsize_1, \dots, subsize_{ndims-1}\}, \\ & \{start_0, start_1, \dots, start_{ndims-1}\}, \text{oldtype}) \end{aligned}$$

Let the typemap of `oldtype` have the form:

$$\{(type_0, disp_0), (type_1, disp_1), \dots, (type_{n-1}, disp_{n-1})\}$$

where $type_i$ is a predefined MPI datatype, and let ex be the extent of `oldtype`. Then we define the `Subarray()` function recursively using the following three equations. Equation 3.2 defines the base step. Equation 3.3 defines the recursion step when `order = MPI_ORDER_FORTRAN`, and Equation 3.4 defines the recursion step when `order = MPI_ORDER_C`.

$$\begin{aligned} & \text{Subarray}(1, \{size_0\}, \{subsize_0\}, \{start_0\}, \\ & \quad \{(type_0, disp_0), (type_1, disp_1), \dots, (type_{n-1}, disp_{n-1})\}) \\ & = \{(MPI_LB, 0), \\ & \quad (type_0, disp_0 + start_0 \times ex), \dots, (type_{n-1}, disp_{n-1} + start_0 \times ex), \\ & \quad (type_0, disp_0 + (start_0 + 1) \times ex), \dots, (type_{n-1}, \\ & \quad \quad disp_{n-1} + (start_0 + 1) \times ex), \dots \\ & \quad (type_0, disp_0 + (start_0 + subsize_0 - 1) \times ex), \dots, \\ & \quad (type_{n-1}, disp_{n-1} + (start_0 + subsize_0 - 1) \times ex), \end{aligned} \tag{3.2}$$

$$(\text{MPI_UB}, \text{size}_0 \times \text{ex})\}$$

$$\begin{aligned} & \text{Subarray}(\text{ndims}, \{\text{size}_0, \text{size}_1, \dots, \text{size}_{\text{ndims}-1}\}, \\ & \quad \{\text{subsize}_0, \text{subsize}_1, \dots, \text{subsize}_{\text{ndims}-1}\}, \\ & \quad \{\text{start}_0, \text{start}_1, \dots, \text{start}_{\text{ndims}-1}\}, \text{oldtype}) \\ = & \text{Subarray}(\text{ndims} - 1, \{\text{size}_1, \text{size}_2, \dots, \text{size}_{\text{ndims}-1}\}, \\ & \quad \{\text{subsize}_1, \text{subsize}_2, \dots, \text{subsize}_{\text{ndims}-1}\}, \\ & \quad \{\text{start}_1, \text{start}_2, \dots, \text{start}_{\text{ndims}-1}\}, \\ & \quad \text{Subarray}(1, \{\text{size}_0\}, \{\text{subsize}_0\}, \{\text{start}_0\}, \text{oldtype})) \end{aligned} \tag{3.3}$$

$$\begin{aligned} & \text{Subarray}(\text{ndims}, \{\text{size}_0, \text{size}_1, \dots, \text{size}_{\text{ndims}-1}\}, \\ & \quad \{\text{subsize}_0, \text{subsize}_1, \dots, \text{subsize}_{\text{ndims}-1}\}, \\ & \quad \{\text{start}_0, \text{start}_1, \dots, \text{start}_{\text{ndims}-1}\}, \text{oldtype}) \\ = & \text{Subarray}(\text{ndims} - 1, \{\text{size}_0, \text{size}_1, \dots, \text{size}_{\text{ndims}-2}\}, \\ & \quad \{\text{subsize}_0, \text{subsize}_1, \dots, \text{subsize}_{\text{ndims}-2}\}, \\ & \quad \{\text{start}_0, \text{start}_1, \dots, \text{start}_{\text{ndims}-2}\}, \\ & \quad \text{Subarray}(1, \{\text{size}_{\text{ndims}-1}\}, \{\text{subsize}_{\text{ndims}-1}\}, \{\text{start}_{\text{ndims}-1}\}, \text{oldtype})) \end{aligned} \tag{3.4}$$

For an example use of `MPI_TYPE_CREATE_SUBARRAY` in the context of I/O see Section 12.9.2.

3.12.5 Distributed Array Datatype Constructor

The distributed array type constructor supports HPF-like [32] data distributions. However, unlike in HPF, the storage order may be specified for C arrays as well as for Fortran arrays.

Advice to users. One can create an HPF-like file view using this type constructor as follows. Complementary filetypes are created by having every process of a group call this constructor with identical arguments (with the exception of rank which should be set appropriately). These filetypes (along with identical `disp` and `etype`) are then used to define the view (via `MPI_FILE_SET_VIEW`). Using this view, a collective data access operation (with identical offsets) will yield an HPF-like distribution pattern. (*End of advice to users.*)

```

1 MPI_TYPE_CREATE_DARRAY(size, rank, ndims, array_of_gsizes, array_of_distrib,
2 array_of_dargs, array_of_psizes, order, oldtype, newtype)
3
4     IN     size                size of process group (positive integer)
5
6     IN     rank                rank in process group (nonnegative integer)
7
8     IN     ndims               number of array dimensions as well as process grid
9                                dimensions (positive integer)
10
11    IN     array_of_gsizes      number of elements of type oldtype in each dimension
12                                of global array (array of positive integers)
13
14    IN     array_of_distrib     distribution of array in each dimension (array of state)
15
16    IN     array_of_dargs       distribution argument in each dimension (array of positive
17                                integers)
18
19    IN     array_of_psizes      size of process grid in each dimension (array of positive
20                                integers)
21
22    IN     order                array storage order flag (state)
23
24    IN     oldtype              old datatype (handle)
25
26    OUT    newtype              new datatype (handle)

```

```

21
22 int MPI_Type_create_darray(int size, int rank, int ndims,
23                            int array_of_gsizes[], int array_of_distrib[], int
24                            array_of_dargs[], int array_of_psizes[], int order,
25                            MPI_Datatype oldtype, MPI_Datatype *newtype)
26
27 MPI_TYPE_CREATE_DARRAY(SIZE, RANK, NDIMS, ARRAY_OF_GSIZES, ARRAY_OF_DISTRIBS,
28                        ARRAY_OF_DARGS, ARRAY_OF_PSIZEs, ORDER, OLDTYPE, NEWTYPE,
29                        IERROR)
30
31 INTEGER SIZE, RANK, NDIMS, ARRAY_OF_GSIZES(*), ARRAY_OF_DISTRIBS(*),
32 ARRAY_OF_DARGS(*), ARRAY_OF_PSIZEs(*), ORDER, OLDTYPE, NEWTYPE, IERROR
33
34 MPI::Datatype MPI::Datatype::Create_darray(int size, int rank, int ndims,
35                            const int array_of_gsizes[], const int array_of_distrib[],
36                            const int array_of_dargs[], const int array_of_psize[],
37                            int order) const

```

MPI_TYPE_CREATE_DARRAY can be used to generate the datatypes corresponding to the distribution of an $ndims$ -dimensional array of `oldtype` elements onto an $ndims$ -dimensional grid of logical processes. Unused dimensions of `array_of_psize`s should be set to 1. (See Example 3.26, page 91.) For a call to MPI_TYPE_CREATE_DARRAY to be correct, the equation $\prod_{i=0}^{ndims-1} array_of_psizes[i] = size$ must be satisfied. The ordering of processes in the process grid is assumed to be row-major, as in the case of virtual Cartesian process topologies in MPI-1.

Advice to users. For both Fortran and C arrays, the ordering of processes in the process grid is assumed to be row-major. This is consistent with the ordering used in virtual Cartesian process topologies in MPI. To create such virtual process topologies, or to find the coordinates of a process in the process grid, etc., users may use the corresponding [process topology functions](#). (*End of advice to users.*)

Each dimension of the array can be distributed in one of three ways:

- MPI_DISTRIBUTE_BLOCK - Block distribution
- MPI_DISTRIBUTE_CYCLIC - Cyclic distribution
- MPI_DISTRIBUTE_NONE - Dimension not distributed.

The constant MPI_DISTRIBUTE_DFLT_DARG specifies a default distribution argument. The distribution argument for a dimension that is not distributed is ignored. For any dimension i in which the distribution is MPI_DISTRIBUTE_BLOCK, it is erroneous to specify $\text{array_of_dargs}[i] * \text{array_of_psizes}[i] < \text{array_of_gsizes}[i]$.

For example, the HPF layout `ARRAY(CYCLIC(15))` corresponds to MPI_DISTRIBUTE_CYCLIC with a distribution argument of 15, and the HPF layout `ARRAY(BLOCK)` corresponds to MPI_DISTRIBUTE_BLOCK with a distribution argument of MPI_DISTRIBUTE_DFLT_DARG.

The order argument is used as in MPI_TYPE_CREATE_SUBARRAY to specify the storage order. Therefore, arrays described by this type constructor may be stored in Fortran (column-major) or C (row-major) order. Valid values for order are MPI_ORDER_FORTRAN and MPI_ORDER_C.

This routine creates a new MPI datatype with a typemap defined in terms of a function called “cyclic()” (see below).

Without loss of generality, it suffices to define the typemap for the MPI_DISTRIBUTE_CYCLIC case where MPI_DISTRIBUTE_DFLT_DARG is not used.

MPI_DISTRIBUTE_BLOCK and MPI_DISTRIBUTE_NONE can be reduced to the MPI_DISTRIBUTE_CYCLIC case for dimension i as follows.

MPI_DISTRIBUTE_BLOCK with $\text{array_of_dargs}[i]$ equal to MPI_DISTRIBUTE_DFLT_DARG is equivalent to MPI_DISTRIBUTE_CYCLIC with $\text{array_of_dargs}[i]$ set to

$$(\text{array_of_gsizes}[i] + \text{array_of_psizes}[i] - 1) / \text{array_of_psizes}[i].$$

If $\text{array_of_dargs}[i]$ is not MPI_DISTRIBUTE_DFLT_DARG, then MPI_DISTRIBUTE_BLOCK and MPI_DISTRIBUTE_CYCLIC are equivalent.

MPI_DISTRIBUTE_NONE is equivalent to MPI_DISTRIBUTE_CYCLIC with $\text{array_of_dargs}[i]$ set to $\text{array_of_gsizes}[i]$.

Finally, MPI_DISTRIBUTE_CYCLIC with $\text{array_of_dargs}[i]$ equal to MPI_DISTRIBUTE_DFLT_DARG is equivalent to MPI_DISTRIBUTE_CYCLIC with $\text{array_of_dargs}[i]$ set to 1.

For MPI_ORDER_FORTRAN, an ndims -dimensional distributed array (`newtype`) is defined by the following code fragment:

```
oldtype[0] = oldtype;
for ( i = 0; i < ndims; i++ ) {
    oldtype[i+1] = cyclic(array_of_dargs[i],
                        array_of_gsizes[i],
                        r[i],
                        array_of_psizes[i],
                        oldtype[i]);
}
newtype = oldtype[ndims];
```

1 For MPI.ORDER_C, the code is:

```

2
3 oldtype[0] = oldtype;
4 for ( i = 0; i < ndims; i++ ) {
5     oldtype[i + 1] = cyclic(array_of_dargs[ndims - i - 1],
6                           array_of_gsizes[ndims - i - 1],
7                           r[ndims - i - 1],
8                           array_of_psize[ndims - i - 1],
9                           oldtype[i]);
10 }
11 newtype = oldtype[ndims];
12
13

```

14 where $r[i]$ is the position of the process (with rank rank) in the process grid at dimension i .
15 The values of $r[i]$ are given by the following code fragment:

```

16     t_rank = rank;
17     t_size = 1;
18     for (i = 0; i < ndims; i++)
19         t_size *= array_of_psize[i];
20     for (i = 0; i < ndims; i++) {
21         t_size = t_size / array_of_psize[i];
22         r[i] = t_rank / t_size;
23         t_rank = t_rank % t_size;
24     }
25

```

26 Let the typemap of oldtype have the form:

```

27     {(type0, disp0), (type1, disp1), ..., (typen-1, dispn-1)}
28

```

29 where $type_i$ is a predefined MPI datatype, and let ex be the extent of oldtype.

30 Given the above, the function cyclic() is defined as follows:

```

31 cyclic(darg, gsize, r, psize, oldtype)
32 = {(MPI_LB, 0),
33   (type0, disp0 + r × darg × ex), ...,
34   (typen-1, dispn-1 + r × darg × ex),
35   (type0, disp0 + (r × darg + 1) × ex), ...,
36   (typen-1, dispn-1 + (r × darg + 1) × ex),
37   ...
38   (type0, disp0 + ((r + 1) × darg - 1) × ex), ...,
39   (typen-1, dispn-1 + ((r + 1) × darg - 1) × ex),
40
41   (type0, disp0 + r × darg × ex + psize × darg × ex), ...,
42   (typen-1, dispn-1 + r × darg × ex + psize × darg × ex),
43
44   (type0, disp0 + (r × darg + 1) × ex + psize × darg × ex), ...,
45   (typen-1, dispn-1 + (r × darg + 1) × ex + psize × darg × ex),
46
47   (type0, disp0 + ((r + 1) × darg - 1) × ex + psize × darg × ex), ...,
48   (typen-1, dispn-1 + ((r + 1) × darg - 1) × ex + psize × darg × ex),

```



```

...
(type0, disp0 + ((r + 1) × darg - 1) × ex + psize × darg × ex), ...,
      (typen-1, dispn-1 + ((r + 1) × darg - 1) × ex + psize × darg × ex),
      ⋮
(type0, disp0 + r × darg × ex + psize × darg × ex × (count - 1)), ...,
      (typen-1, dispn-1 + r × darg × ex + psize × darg × ex × (count - 1)),
(type0, disp0 + (r × darg + 1) × ex + psize × darg × ex × (count - 1)), ...,
      (typen-1, dispn-1 + (r × darg + 1) × ex
      + psize × darg × ex × (count - 1)),
...
(type0, disp0 + (r × darg + darglast - 1) × ex
      + psize × darg × ex × (count - 1)), ...,
      (typen-1, dispn-1 + (r × darg + darglast - 1) × ex
      + psize × darg × ex × (count - 1)),
(MPI_UB, gsize * ex)

```

where *count* is defined by this code fragment:

```

nblocks = (gsize + (darg - 1)) / darg;
count = nblocks / psize;
left_over = nblocks - count * psize;
if (r < left_over)
    count = count + 1;

```

Here, *nblocks* is the number of blocks that must be distributed among the processors. Finally, *darglast* is defined by this code fragment:

```

if ((num_in_last_cyclic = gsize % (psize * darg)) == 0)
    darg_last = darg;
else
    darg_last = num_in_last_cyclic - darg * r;
    if (darg_last > darg)
        darg_last = darg;
    if (darg_last <= 0)
        darg_last = darg;

```

Example 3.26 Consider generating the filetypes corresponding to the HPF distribution:

```

<oldtype> FILEARRAY(100, 200, 300)
!HPF$ PROCESSORS PROCESSES(2, 3)
!HPF$ DISTRIBUTE FILEARRAY(CYCLIC(10), *, BLOCK) ONTO PROCESSES

```

This can be achieved by the following Fortran code, assuming there will be six processes attached to the run:

```

1     ndims = 3
2     array_of_gsizes(1) = 100
3     array_of_distribs(1) = MPI_DISTRIBUTE_CYCLIC
4     array_of_dargs(1) = 10
5     array_of_gsizes(2) = 200
6     array_of_distribs(2) = MPI_DISTRIBUTE_NONE
7     array_of_dargs(2) = 0
8     array_of_gsizes(3) = 300
9     array_of_distribs(3) = MPI_DISTRIBUTE_BLOCK
10    array_of_dargs(3) = MPI_DISTRIBUTE_DFLT_ARG
11    array_of_psizes(1) = 2
12    array_of_psizes(2) = 1
13    array_of_psizes(3) = 3
14    call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
15    call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
16    call MPI_TYPE_CREATE_DARRAY(size, rank, ndims, array_of_gsizes, &
17        array_of_distribs, array_of_dargs, array_of_psizes,      &
18        MPI_ORDER_FORTRAN, oldtype, newtype, ierr)
19
20

```

3.12.6 Address and size functions

The displacements in a general datatype are relative to some initial buffer address. **Absolute addresses** can be substituted for these displacements: we treat them as displacements relative to “address zero,” the start of the address space. This initial address zero is indicated by the constant `MPI_BOTTOM`. Thus, a datatype can specify the absolute address of the entries in the communication buffer, in which case the `buf` argument is passed the value `MPI_BOTTOM`.

The address of a location in memory can be found by invoking the function `MPI_GET_ADDRESS`.

```

32 MPI_GET_ADDRESS(location, address)

```

| | | | |
|----|-----|----------|------------------------------------|
| 33 | IN | location | location in caller memory (choice) |
| 34 | | | |
| 35 | OUT | address | address of location (integer) |

```

36
37 int MPI_Get_address(void *location, MPI_Aint *address)

```

```

38 MPI_GET_ADDRESS(LOCATION, ADDRESS, IERROR)

```

```

39 <type> LOCATION(*)

```

```

40 INTEGER IERROR

```

```

41 INTEGER(KIND=MPI_ADDRESS_KIND) ADDRESS

```

```

42
43 MPI::Aint MPI::Get_address(void* location)

```

This function replaces `MPI_ADDRESS`, whose use is deprecated. See also Chapter 15.
Returns the (byte) address of `location`.

Advice to users. Current Fortran MPI codes will run unmodified, and will port to any system. However, they may fail if addresses larger than $2^{32} - 1$ are used

in the program. New codes should be written so that they use the new functions. This provides compatibility with C/C++ and avoids errors on 64 bit architectures. However, such newly written codes may need to be (slightly) rewritten to port to old Fortran 77 environments that do not support KIND declarations. (*End of advice to users.*)

Example 3.27 Using `MPI_GET_ADDRESS` for an array.

```

REAL A(100,100)
INTEGER(KIND=MPI_ADDRESS_KIND) I1, I2, DIFF
CALL MPI_GET_ADDRESS(A(1,1), I1, IERROR)
CALL MPI_GET_ADDRESS(A(10,10), I2, IERROR)
DIFF = I2 - I1
! The value of DIFF is 909*sizeofreal; the values of I1 and I2 are
! implementation dependent.
```

Advice to users. C users may be tempted to avoid the usage of `MPI_GET_ADDRESS` and rely on the availability of the address operator `&`. Note, however, that `& cast-expression` is a pointer, not an address. ISO C does not require that the value of a pointer (or the pointer cast to int) be the absolute address of the object pointed at — although this is commonly the case. Furthermore, referencing may not have a unique definition on machines with a segmented address space. The use of `MPI_GET_ADDRESS` to “reference” C variables guarantees portability to such machines as well. (*End of advice to users.*)

Advice to users. To prevent problems with the argument copying and register optimization done by Fortran compilers, please note the hints in subsections “Problems Due to Data Copying and Sequence Association,” and “A Problem with Register Optimization” in Section 13.2.2 on pages 451 and 454. (*End of advice to users.*)

The following auxiliary function provides useful information on derived datatypes.

```

MPI_TYPE_SIZE(datatype, size)
    IN      datatype          datatype (handle)
    OUT     size              datatype size (integer)

int MPI_Type_size(MPI_Datatype datatype, int *size)
MPI_TYPE_SIZE(DATATYPE, SIZE, IERROR)
    INTEGER DATATYPE, SIZE, IERROR

int MPI::Datatype::Get_size() const
```

`MPI_TYPE_SIZE` returns the total size, in bytes, of the entries in the type signature associated with `datatype`; i.e., the total size of the data in a message that would be created with this datatype. Entries that occur multiple times in the datatype are counted with their multiplicity.

1 *Advice to users.* The MPI-1 Standard specifies that the output argument of
 2 MPI_TYPE_SIZE in C is of type `int`. The MPI Forum considered proposals to change
 3 this and decided to reiterate the original decision. (*End of advice to users.*)

6 3.12.7 Lower-bound and upper-bound markers

7 It is often convenient to define explicitly the lower bound and upper bound of a type map,
 8 and override the definition given on page 94. This allows one to define a datatype that has
 9 “holes” at its beginning or its end, or a datatype with entries that extend above the upper
 10 bound or below the lower bound. Examples of such usage are provided in Sec. 3.12.14.
 11 Also, the user may want to override the alignment rules that are used to compute upper
 12 bounds and extents. E.g., a C compiler may allow the user to override default alignment
 13 rules for some of the structures within a program. The user has to specify explicitly the
 14 bounds of the datatypes that match these structures.

15 To achieve this, we add two additional “pseudo-datatypes,” `MPI_LB` and `MPI_UB`, that
 16 can be used, respectively, to mark the lower bound or the upper bound of a datatype. These
 17 pseudo-datatypes occupy no space ($extent(MPI_LB) = extent(MPI_UB) = 0$). They do not
 18 affect the size or count of a datatype, and do not affect the content of a message created
 19 with this datatype. However, they do affect the definition of the extent of a datatype and,
 20 therefore, affect the outcome of a replication of this datatype by a datatype constructor.

21 **Example 3.28** Let $D = (-3, 0, 6)$; $T = (MPI_LB, MPI_INT, MPI_UB)$, and $B = (1, 1, 1)$.
 22 Then a call to `MPI_TYPE_STRUCT(3, B, D, T, type1)` creates a new datatype that has an
 23 extent of 9 (from -3 to 5, 5 included), and contains an integer at displacement 0. This is
 24 the datatype defined by the sequence $\{(lb, -3), (int, 0), (ub, 6)\}$. If this type is replicated
 25 twice by a call to `MPI_TYPE_CONTIGUOUS(2, type1, type2)` then the newly created type
 26 can be described by the sequence $\{(lb, -3), (int, 0), (int, 9), (ub, 15)\}$. (An entry of type `ub`
 27 can be deleted if there is another entry of type `ub` with a higher displacement; an entry of
 28 type `lb` can be deleted if there is another entry of type `lb` with a lower displacement.)

29 In general, if

$$30 Typemap = \{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

31 then the **lower bound** of *Typemap* is defined to be

$$32 lb(Typemap) = \begin{cases} \min_j disp_j & \text{if no entry has basic type lb} \\ \min_j \{disp_j \text{ such that } type_j = lb\} & \text{otherwise} \end{cases}$$

33 Similarly, the **upper bound** of *Typemap* is defined to be

$$34 ub(Typemap) = \begin{cases} \max_j disp_j + sizeof(type_j) + \epsilon & \text{if no entry has basic type ub} \\ \max_j \{disp_j \text{ such that } type_j = ub\} & \text{otherwise} \end{cases}$$

35 Then

$$36 extent(Typemap) = ub(Typemap) - lb(Typemap)$$

37 If $type_i$ requires alignment to a byte address that is a multiple of k_i , then ϵ is the least
 38 nonnegative increment needed to round $extent(Typemap)$ to the next multiple of $\max_i k_i$.

39 The formal definitions given for the various datatype constructors apply now, with the
 40 amended definition of **extent**.

Returns in `newtype` a handle to a new datatype that is identical to `oldtype`, except that the lower bound of this new datatype is set to be `lb`, and its upper bound is set to be `lb + extent`. Any previous `lb` and `ub` markers are erased, and a new pair of lower bound and upper bound markers are put in the positions indicated by the `lb` and `extent` arguments. This affects the behavior of the datatype when used in communication operations, with `count > 1`, and when used in the construction of new derived datatypes.

Advice to users. It is strongly recommended that users use these two new functions, rather than the old MPI-1 functions to set and access lower bound, upper bound and extent of datatypes. (*End of advice to users.*)

3.12.9 True Extent of Datatypes

Suppose we implement gather as a spanning tree implemented on top of point-to-point routines. Since the receive buffer is only valid on the root process, one will need to allocate some temporary space for receiving data on intermediate nodes. However, the datatype extent cannot be used as an estimate of the amount of space that needs to be allocated, if the user has modified the extent using the `MPI_UB` and `MPI_LB` values. A new function is provided which returns the true extent of the datatype.

```
MPI_TYPE_GET_TRUE_EXTENT(datatype, true_lb, true_extent)
```

| | | |
|-----|--------------------------|---|
| IN | <code>datatype</code> | datatype to get information on (handle) |
| OUT | <code>true_lb</code> | true lower bound of datatype (integer) |
| OUT | <code>true_extent</code> | true size of datatype (integer) |

```
int MPI_Type_get_true_extent(MPI_Datatype datatype, MPI_Aint *true_lb,
                             MPI_Aint *true_extent)
```

```
MPI_TYPE_GET_TRUE_EXTENT(DATATYPE, TRUE_LB, TRUE_EXTENT, IERROR)
INTEGER DATATYPE, IERROR
INTEGER(KIND = MPI_ADDRESS_KIND) TRUE_LB, TRUE_EXTENT
```

```
void MPI::Datatype::Get_true_extent(MPI::Aint& true_lb,
                                     MPI::Aint& true_extent) const
```

`true_lb` returns the offset of the lowest unit of store which is addressed by the datatype, i.e., the lower bound of the corresponding typemap, ignoring `MPI_LB` markers. `true_extent` returns the true size of the datatype, i.e., the extent of the corresponding typemap, ignoring `MPI_LB` and `MPI_UB` markers, and performing no rounding for alignment. If the typemap associated with `datatype` is

$$\text{Typemap} = \{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\}$$

Then

$$\text{true_lb}(\text{Typemap}) = \min_j \{disp_j : type_j \neq \mathbf{lb}, \mathbf{ub}\},$$

$$\text{true_ub}(\text{Typemap}) = \max_j \{disp_j + \text{sizeof}(type_j) : type_j \neq \mathbf{lb}, \mathbf{ub}\},$$

and

$$true_extent(Typemap) = true_ub(Typemap) - true_lb(typemap).$$

(Readers should compare this with the definitions in [Section 3.12.7 on page 94](#) and [Section 3.12.8 on page 95](#), which describe the function `MPI_TYPE_EXTENT`.)

The `true_extent` is the minimum number of bytes of memory necessary to hold a datatype, uncompressed.

3.12.10 Commit and free

A datatype object has to be **committed** before it can be used in a communication. [As a argument in datatype constructors, uncommitted and also committed datatypes can be used.](#) There is no need to commit basic datatypes. They are “pre-committed.”

`MPI_TYPE_COMMIT(datatype)`

INOUT datatype datatype that is committed (handle)

`int MPI_Type_commit(MPI_Datatype *datatype)`

`MPI_TYPE_COMMIT(DATATYPE, IERROR)`

INTEGER DATATYPE, IERROR

`void MPI::Datatype::Commit()`

The commit operation commits the datatype, that is, the formal description of a communication buffer, not the content of that buffer. Thus, after a datatype has been committed, it can be repeatedly reused to communicate the changing content of a buffer or, indeed, the content of different buffers, with different starting addresses.

Advice to implementors. The system may “compile” at commit time an internal representation for the datatype that facilitates communication, e.g. change from a compacted representation to a flat representation of the datatype, and select the most convenient transfer mechanism. (*End of advice to implementors.*)

`MPI_TYPE_COMMIT` will accept a committed datatype; in this case, it is equivalent to a no-op.

`MPI_TYPE_FREE(datatype)`

INOUT datatype datatype that is freed (handle)

`int MPI_Type_free(MPI_Datatype *datatype)`

`MPI_TYPE_FREE(DATATYPE, IERROR)`

INTEGER DATATYPE, IERROR

`void MPI::Datatype::Free()`

Marks the datatype object associated with `datatype` for deallocation and sets `datatype` to `MPI_DATATYPE_NULL`. Any communication that is currently using this datatype

will complete normally. Derived datatypes that were defined from the freed datatype are not affected.

Example 3.29 The following code fragment gives examples of using `MPI_TYPE_COMMIT`.

```

1  INTEGER type1, type2
2  CALL MPI_TYPE_CONTIGUOUS(5, MPI_REAL, type1, ierr)
3      ! new type object created
4  CALL MPI_TYPE_COMMIT(type1, ierr)
5      ! now type1 can be used for communication
6  type2 = type1
7      ! type2 can be used for communication
8      ! (it is a handle to same object as type1)
9  CALL MPI_TYPE_VECTOR(3, 5, 4, MPI_REAL, type1, ierr)
10     ! new uncommitted type object created
11 CALL MPI_TYPE_COMMIT(type1, ierr)
12     ! now type1 can be used anew for communication

```

Freeing a datatype does not affect any other datatype that was built from the freed datatype. The system behaves as if input datatype arguments to derived datatype constructors are passed by value.

Advice to implementors. The implementation may keep a reference count of active communications that use the datatype, in order to decide when to free it. Also, one may implement constructors of derived datatypes so that they keep pointers to their datatype arguments, rather than copying them. In this case, one needs to keep track of active datatype definition references in order to know when a datatype object can be freed. (*End of advice to implementors.*)

3.12.11 Duplicating a Datatype

`MPI_TYPE_DUP`(type, newtype)

| | | |
|-----|---------|-----------------------|
| IN | type | datatype (handle) |
| OUT | newtype | copy of type (handle) |

```
int MPI_Type_dup(MPI_Datatype type, MPI_Datatype *newtype)
```

```
MPI_TYPE_DUP(TYPE, NEWTYPE, IERROR)
INTEGER TYPE, NEWTYPE, IERROR
```

```
MPI::Datatype MPI::Datatype::Dup() const
```

`MPI_TYPE_DUP` is a [type constructor](#) which duplicates the existing type with associated key values. For each key value, the respective copy callback function determines the attribute value associated with this key in the new communicator; one particular action that a copy callback may take is to delete the attribute from the new datatype. Returns in `newtype` a new datatype with exactly the same properties as `type`

and any copied cached information, see Section 5.7.5 on page 220. The new datatype has identical upper bound and lower bound and yields the same net result when fully decoded with the functions in Section 11.6. The newtype has the same committed state as the old type.

3.12.12 Use of general datatypes in communication

Handles to derived datatypes can be passed to a communication call wherever a datatype argument is required. A call of the form `MPI_SEND(buf, count, datatype, ...)`, where `count > 1`, is interpreted as if the call was passed a new datatype which is the concatenation of `count` copies of `datatype`. Thus, `MPI_SEND(buf, count, datatype, dest, tag, comm)` is equivalent to,

```
MPI_TYPE_CONTIGUOUS(count, datatype, newtype)
MPI_TYPE_COMMIT(newtype)
MPI_SEND(buf, 1, newtype, dest, tag, comm).
```

Similar statements apply to all other communication functions that have a `count` and `datatype` argument.

Suppose that a send operation `MPI_SEND(buf, count, datatype, dest, tag, comm)` is executed, where `datatype` has type map,

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

and extent *extent*. (Empty entries of “pseudo-type” `MPI_UB` and `MPI_LB` are not listed in the type map, but they affect the value of *extent*.) The send operation sends $n \cdot \text{count}$ entries, where entry $i \cdot n + j$ is at location $addr_{i,j} = \text{buf} + \text{extent} \cdot i + disp_j$ and has type $type_j$, for $i = 0, \dots, \text{count} - 1$ and $j = 0, \dots, n - 1$. These entries need not be contiguous, nor distinct; their order can be arbitrary.

The variable stored at address $addr_{i,j}$ in the calling program should be of a type that matches $type_j$, where type matching is defined as in section 3.3.1. The message sent contains $n \cdot \text{count}$ entries, where entry $i \cdot n + j$ has type $type_j$.

Similarly, suppose that a receive operation `MPI_RECV(buf, count, datatype, source, tag, comm, status)` is executed, where `datatype` has type map,

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\},$$

with extent *extent*. (Again, empty entries of “pseudo-type” `MPI_UB` and `MPI_LB` are not listed in the type map, but they affect the value of *extent*.) This receive operation receives $n \cdot \text{count}$ entries, where entry $i \cdot n + j$ is at location $\text{buf} + \text{extent} \cdot i + disp_j$ and has type $type_j$. If the incoming message consists of k elements, then we must have $k \leq n \cdot \text{count}$; the $i \cdot n + j$ -th element of the message should have a type that matches $type_j$.

Type matching is defined according to the type signature of the corresponding datatypes, that is, the sequence of basic type components. Type matching does not depend on some aspects of the datatype definition, such as the displacements (layout in memory) or the intermediate types used.

Example 3.30 This example shows that type matching is defined in terms of the basic types that a derived type consists of.

```

1  ...
2  CALL MPI_TYPE_CONTIGUOUS( 2, MPI_REAL, type2, ...)
3  CALL MPI_TYPE_CONTIGUOUS( 4, MPI_REAL, type4, ...)
4  CALL MPI_TYPE_CONTIGUOUS( 2, type2, type22, ...)
5  ...
6  CALL MPI_SEND( a, 4, MPI_REAL, ...)
7  CALL MPI_SEND( a, 2, type2, ...)
8  CALL MPI_SEND( a, 1, type22, ...)
9  CALL MPI_SEND( a, 1, type4, ...)
10 ...
11 CALL MPI_RECV( a, 4, MPI_REAL, ...)
12 CALL MPI_RECV( a, 2, type2, ...)
13 CALL MPI_RECV( a, 1, type22, ...)
14 CALL MPI_RECV( a, 1, type4, ...)

```

Each of the sends matches any of the receives.

A datatype may specify overlapping entries. The use of such a datatype in a receive operation is erroneous. (This is erroneous even if the actual message received is short enough not to write any entry more than once.)

Suppose that `MPI_RECV(buf, count, datatype, dest, tag, comm, status)` is executed, where `datatype` has type map,

$$\{(type_0, disp_0), \dots, (type_{n-1}, disp_{n-1})\}.$$

The received message need not fill all the receive buffer, nor does it need to fill a number of locations which is a multiple of n . Any number, k , of basic elements can be received, where $0 \leq k \leq \text{count} \cdot n$. The number of basic elements received can be retrieved from `status` using the query function `MPI_GET_ELEMENTS`.

```

30 MPI_GET_ELEMENTS( status, datatype, count)

```

| | | | |
|----|-----|----------|---|
| 31 | IN | status | return status of receive operation (Status) |
| 32 | IN | datatype | datatype used by receive operation (handle) |
| 33 | OUT | count | number of received basic elements (integer) |

```

34
35
36 int MPI_Get_elements(MPI_Status *status, MPI_Datatype datatype, int *count)

```

```

37
38 MPI_GET_ELEMENTS(STATUS, DATATYPE, COUNT, IERROR)
39     INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, COUNT, IERROR

```

```

40
41 int MPI::Status::Get_elements(const MPI::Datatype& datatype) const

```

The previously defined function, `MPI_GET_COUNT` (Sec. 3.2.5), has a different behavior. It returns the number of “top-level entries” received, i.e. the number of “copies” of type `datatype`. In the previous example, `MPI_GET_COUNT` may return any integer value k , where $0 \leq k \leq \text{count}$. If `MPI_GET_COUNT` returns k , then the number of basic elements received (and the value returned by `MPI_GET_ELEMENTS`) is $n \cdot k$. If the number of basic elements received is not a multiple of n , that is, if the receive operation has not received an integral number of datatype “copies,” then `MPI_GET_COUNT` returns the value `MPI_UNDEFINED`.

The datatype argument should match the argument provided by the receive call that set the status variable.

Example 3.31 Usage of MPI_GET_COUNT and MPI_GET_ELEMENT.

```

...
CALL MPI_TYPE_CONTIGUOUS(2, MPI_REAL, Type2, ierr)
CALL MPI_TYPE_COMMIT(Type2, ierr)
...
CALL MPI_COMM_RANK(comm, rank, ierr)
IF(rank.EQ.0) THEN
    CALL MPI_SEND(a, 2, MPI_REAL, 1, 0, comm, ierr)
    CALL MPI_SEND(a, 3, MPI_REAL, 1, 0, comm, ierr)
ELSE
    CALL MPI_RECV(a, 2, Type2, 0, 0, comm, stat, ierr)
    CALL MPI_GET_COUNT(stat, Type2, i, ierr)      ! returns i=1
    CALL MPI_GET_ELEMENTS(stat, Type2, i, ierr)  ! returns i=2
    CALL MPI_RECV(a, 2, Type2, 0, 0, comm, stat, ierr)
    CALL MPI_GET_COUNT(stat, Type2, i, ierr)      ! returns i=MPI_UNDEFINED
    CALL MPI_GET_ELEMENTS(stat, Type2, i, ierr)  ! returns i=3
END IF

```

The function MPI_GET_ELEMENTS can also be used after a probe to find the number of elements in the probed message. Note that the two functions MPI_GET_COUNT and MPI_GET_ELEMENTS return the same values when they are used with basic datatypes.

Rationale. The extension given to the definition of MPI_GET_COUNT seems natural: one would expect this function to return the value of the count argument, when the receive buffer is filled. Sometimes datatype represents a basic unit of data one wants to transfer, for example, a record in an array of records (structures). One should be able to find out how many components were received without bothering to divide by the number of elements in each component. However, on other occasions, datatype is used to define a complex layout of data in the receiver memory, and does not represent a basic unit of data for transfers. In such cases, one needs to use the function MPI_GET_ELEMENTS. (*End of rationale.*)

Advice to implementors. The definition implies that a receive cannot change the value of storage outside the entries defined to compose the communication buffer. In particular, the definition implies that padding space in a structure should not be modified when such a structure is copied from one process to another. This would prevent the obvious optimization of copying the structure, together with the padding, as one contiguous block. The implementation is free to do this optimization when it does not impact the outcome of the computation. The user can “force” this optimization by explicitly including padding as part of the message. (*End of advice to implementors.*)

3.12.13 Correct use of addresses

Successively declared variables in C or Fortran are not necessarily stored at contiguous locations. Thus, care must be exercised that displacements do not cross from one variable

1 to another. Also, in machines with a segmented address space, addresses are not unique
2 and address arithmetic has some peculiar properties. Thus, the use of **addresses**, that is,
3 displacements relative to the start address `MPI_BOTTOM`, has to be restricted.

4 Variables belong to the same **sequential storage** if they belong to the same array, to
5 the same `COMMON` block in Fortran, or to the same structure in C. Valid addresses are
6 defined recursively as follows:

- 7
- 8 1. The function `MPI_GET_ADDRESS` returns a valid address, when passed as argument
9 a variable of the calling program.
- 10
- 11 2. The `buf` argument of a communication function evaluates to a valid address, when
12 passed as argument a variable of the calling program.
- 13
- 14 3. If `v` is a valid address, and `i` is an integer, then `v+i` is a valid address, provided `v` and
15 `v+i` are in the same sequential storage.
- 16
- 17 4. If `v` is a valid address then `MPI_BOTTOM + v` is a valid address.

18 A correct program uses only valid addresses to identify the locations of entries in
19 communication buffers. Furthermore, if `u` and `v` are two valid addresses, then the (integer)
20 difference `u - v` can be computed only if both `u` and `v` are in the same sequential storage.
21 No other arithmetic operations can be meaningfully executed on addresses.

22 The rules above impose no constraints on the use of derived datatypes, as long as
23 they are used to define a communication buffer that is wholly contained within the same
24 sequential storage. However, the construction of a communication buffer that contains
25 variables that are not within the same sequential storage must obey certain restrictions.
26 Basically, a communication buffer with variables that are not within the same sequential
27 storage can be used only by specifying in the communication call `buf = MPI_BOTTOM`,
28 `count = 1`, and using a `datatype` argument where all displacements are valid (absolute)
29 addresses.

30
31 *Advice to users.* It is not expected that MPI implementations will be able to detect
32 erroneous, “out of bound” displacements — unless those overflow the user address
33 space — since the MPI call may not know the extent of the arrays and records in the
34 host program. (*End of advice to users.*)

35
36 *Advice to implementors.* There is no need to distinguish (absolute) addresses and
37 (relative) displacements on a machine with contiguous address space: `MPI_BOTTOM`
38 is zero, and both addresses and displacements are integers. On machines where the
39 distinction is required, addresses are recognized as expressions that involve
40 `MPI_BOTTOM`. (*End of advice to implementors.*)

41 42 43 3.12.14 Examples

44 The following examples illustrate the use of derived datatypes.

45
46 **Example 3.32** Send and receive a section of a 3D array.

```

REAL a(100,100,100), e(9,9,9)
INTEGER oneslice, twoslice, threeslice, sizeofreal, myrank, ierr
INTEGER status(MPI_STATUS_SIZE)

C   extract the section a(1:17:2, 3:11, 2:10)
C   and store it in e(:, :, :).

CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)

CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)

C   create datatype for a 1D section
CALL MPI_TYPE_VECTOR( 9, 1, 2, MPI_REAL, oneslice, ierr)

C   create datatype for a 2D section
CALL MPI_TYPE_HVECTOR(9, 1, 100*sizeofreal, oneslice, twoslice, ierr)

C   create datatype for the entire section
CALL MPI_TYPE_HVECTOR( 9, 1, 100*100*sizeofreal, twoslice,
                      threeslice, ierr)

CALL MPI_TYPE_COMMIT( threeslice, ierr)
CALL MPI_SENDRECV(a(1,3,2), 1, threeslice, myrank, 0, e, 9*9*9,
                  MPI_REAL, myrank, 0, MPI_COMM_WORLD, status, ierr)

```

Example 3.33 Copy the (strictly) lower triangular part of a matrix.

```

REAL a(100,100), b(100,100)
INTEGER disp(100), blocklen(100), ltype, myrank, ierr
INTEGER status(MPI_STATUS_SIZE)

C   copy lower triangular part of array a
C   onto lower triangular part of array b

CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)

C   compute start and size of each column
DO i=1, 100
    disp(i) = 100*(i-1) + i
    block(i) = 100-i
END DO

C   create datatype for lower triangular part
CALL MPI_TYPE_INDEXED( 100, block, disp, MPI_REAL, ltype, ierr)

CALL MPI_TYPE_COMMIT(ltype, ierr)
CALL MPI_SENDRECV( a, 1, ltype, myrank, 0, b, 1,
                  ltype, myrank, 0, MPI_COMM_WORLD, status, ierr)

```

1 **Example 3.34** Transpose a matrix.

```

2
3     REAL a(100,100), b(100,100)
4     INTEGER row, xpose, sizeofreal, myrank, ierr
5     INTEGER status(MPI_STATUS_SIZE)
6
7 C   transpose matrix a onto b
8
9     CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)
10
11    CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)
12
13 C   create datatype for one row
14    CALL MPI_TYPE_VECTOR( 100, 1, 100, MPI_REAL, row, ierr)
15
16 C   create datatype for matrix in row-major order
17    CALL MPI_TYPE_HVECTOR( 100, 1, sizeofreal, row, xpose, ierr)
18
19    CALL MPI_TYPE_COMMIT( xpose, ierr)
20
21 C   send matrix in row-major order and receive in column major order
22    CALL MPI_SENDRECV( a, 1, xpose, myrank, 0, b, 100*100,
23                      MPI_REAL, myrank, 0, MPI_COMM_WORLD, status, ierr)
24

```

25 **Example 3.35** Another approach to the transpose problem:

```

26
27    REAL a(100,100), b(100,100)
28    INTEGER disp(2), blocklen(2), type(2), row, row1, sizeofreal
29    INTEGER myrank, ierr
30    INTEGER status(MPI_STATUS_SIZE)
31
32    CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)
33
34 C   transpose matrix a onto b
35
36    CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)
37
38 C   create datatype for one row
39    CALL MPI_TYPE_VECTOR( 100, 1, 100, MPI_REAL, row, ierr)
40
41 C   create datatype for one row, with the extent of one real number
42    disp(1) = 0
43    disp(2) = sizeofreal
44    type(1) = row
45    type(2) = MPI_UB
46    blocklen(1) = 1
47    blocklen(2) = 1
48    CALL MPI_TYPE_STRUCT( 2, blocklen, disp, type, row1, ierr)

```

```

CALL MPI_TYPE_COMMIT( row1, ierr)
C    send 100 rows and receive in column major order
CALL MPI_SENDRECV( a, 100, row1, myrank, 0, b, 100*100,
                  MPI_REAL, myrank, 0, MPI_COMM_WORLD, status, ierr)

```

Example 3.36 We manipulate an array of structures.

```

struct Partstruct
{
    int    class; /* particle class */
    double d[6]; /* particle coordinates */
    char   b[7]; /* some additional information */
};

struct Partstruct    particle[1000];

int                  i, dest, rank;
MPI_Comm             comm;

/* build datatype describing structure */

MPI_Datatype Particletype;
MPI_Datatype type[3] = {MPI_INT, MPI_DOUBLE, MPI_CHAR};
int          blocklen[3] = {1, 6, 7};
MPI_Aint     disp[3];
MPI_Aint     base;

/* compute displacements of structure components */

MPI_Address( particle, disp);
MPI_Address( particle[0].d, disp+1);
MPI_Address( particle[0].b, disp+2);
base = disp[0];
for (i=0; i <3; i++) disp[i] -= base;

MPI_Type_struct( 3, blocklen, disp, type, &Particletype);

/* If compiler does padding in mysterious ways,
the following may be safer */

MPI_Datatype type1[4] = {MPI_INT, MPI_DOUBLE, MPI_CHAR, MPI_UB};
int          blocklen1[4] = {1, 6, 7, 1};
MPI_Aint     disp1[4];

```

```

1  /* compute displacements of structure components */
2
3  MPI_Address( particle, disp1);
4  MPI_Address( particle[0].d, disp1+1);
5  MPI_Address( particle[0].b, disp1+2);
6  MPI_Address( particle+1, disp1+3);
7  base = disp1[0];
8  for (i=0; i <4; i++) disp1[i] -= base;
9
10 /* build datatype describing structure */
11
12 MPI_Type_struct( 4, blocklen1, disp1, type1, &Particletype);
13
14
15         /* 4.1:
16         send the entire array */
17
18 MPI_Type_commit( &Particletype);
19 MPI_Send( particle, 1000, Particletype, dest, tag, comm);
20
21
22         /* 4.2:
23         send only the entries of class zero particles,
24         preceded by the number of such entries */
25
26 MPI_Datatype Zparticles; /* datatype describing all particles
27                          with class zero (needs to be recomputed
28                          if classes change) */
29 MPI_Datatype Ztype;
30
31 MPI_Aint      zdisp[1000];
32 int zblock[1000], j, k;
33 int zzblock[2] = {1,1};
34 MPI_Aint      zzdisp[2];
35 MPI_Datatype  zztype[2];
36
37 /* compute displacements of class zero particles */
38 j = 0;
39 for(i=0; i < 1000; i++)
40     if (particle[i].class==0)
41         {
42             zdisp[j] = i;
43             zblock[j] = 1;
44             j++;
45         }
46
47 /* create datatype for class zero particles */
48 MPI_Type_indexed( j, zblock, zdisp, Particletype, &Zparticles);

```



```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

```

```

/* prepend particle count */
MPI_Address(&j, zzdisp);
MPI_Address(particle, zzdisp+1);
zztype[0] = MPI_INT;
zztype[1] = Zparticles;
MPI_Type_struct(2, zblock, zzdisp, zztype, &Ztype);

MPI_Type_commit( &Ztype);
MPI_Send( MPI_BOTTOM, 1, Ztype, dest, tag, comm);

/* A probably more efficient way of defining Zparticles */

/* consecutive particles with index zero are handled as one block */
j=0;
for (i=0; i < 1000; i++)
    if (particle[i].index==0)
        {
            for (k=i+1; (k < 1000)&&(particle[k].index == 0) ; k++);
            zdisp[j] = i;
            zblock[j] = k-i;
            j++;
            i = k;
        }
MPI_Type_indexed( j, zblock, zdisp, Particletype, &Zparticles);

/* 4.3:
send the first two coordinates of all entries */

MPI_Datatype Allpairs; /* datatype for all pairs of coordinates */

MPI_Aint sizeofentry;

MPI_Type_extent( Particletype, &sizeofentry);

/* sizeofentry can also be computed by subtracting the address
of particle[0] from the address of particle[1] */

MPI_Type_hvector( 1000, 2, sizeofentry, MPI_DOUBLE, &Allpairs);
MPI_Type_commit( &Allpairs);
MPI_Send( particle[0].d, 1, Allpairs, dest, tag, comm);

/* an alternative solution to 4.3 */

MPI_Datatype Onepair; /* datatype for one pair of coordinates, with
the extent of one particle entry */

```

```

1  MPI_Aint disp2[3];
2  MPI_Datatype type2[3] = {MPI_LB, MPI_DOUBLE, MPI_UB};
3  int blocklen2[3] = {1, 2, 1};
4
5  MPI_Address( particle, disp2);
6  MPI_Address( particle[0].d, disp2+1);
7  MPI_Address( particle+1, disp2+2);
8  base = disp2[0];
9  for (i=0; i<2; i++) disp2[i] -= base;
10
11 MPI_Type_struct( 3, blocklen2, disp2, type2, &Onepair);
12 MPI_Type_commit( &Onepair);
13 MPI_Send( particle[0].d, 1000, Onepair, dest, tag, comm);
14
15

```

Example 3.37 The same manipulations as in the previous example, but use absolute addresses in datatypes.

```

18
19 struct Partstruct
20 {
21     int class;
22     double d[6];
23     char b[7];
24 };
25
26 struct Partstruct particle[1000];
27
28     /* build datatype describing first array entry */
29
30 MPI_Datatype Particletype;
31 MPI_Datatype type[3] = {MPI_INT, MPI_DOUBLE, MPI_CHAR};
32 int          block[3] = {1, 6, 7};
33 MPI_Aint     disp[3];
34
35 MPI_Address( particle, disp);
36 MPI_Address( particle[0].d, disp+1);
37 MPI_Address( particle[0].b, disp+2);
38 MPI_Type_struct( 3, block, disp, type, &Particletype);
39
40 /* Particletype describes first array entry -- using absolute
41    addresses */
42
43     /* 5.1:
44    send the entire array */
45
46 MPI_Type_commit( &Particletype);
47 MPI_Send( MPI_BOTTOM, 1000, Particletype, dest, tag, comm);
48

```

```

/* 5.2:
send the entries of class zero,
preceded by the number of such entries */

MPI_Datatype Zparticles, Ztype;

MPI_Aint zdisp[1000]
int zblock[1000], i, j, k;
int zzblock[2] = {1,1};
MPI_Datatype zztype[2];
MPI_Aint    zzdisp[2];

j=0;
for (i=0; i < 1000; i++)
    if (particle[i].index==0)
        {
            for (k=i+1; (k < 1000)&&(particle[k].index = 0) ; k++);
            zdisp[j] = i;
            zblock[j] = k-i;
            j++;
            i = k;
        }
MPI_Type_indexed( j, zblock, zdisp, Particletype, &Zparticles);
/* Zparticles describe particles with class zero, using
their absolute addresses*/

/* prepend particle count */
MPI_Address(&j, zzdisp);
zzdisp[1] = MPI_BOTTOM;
zztype[0] = MPI_INT;
zztype[1] = Zparticles;
MPI_Type_struct(2, zzblock, zzdisp, zztype, &Ztype);

MPI_Type_commit( &Ztype);
MPI_Send( MPI_BOTTOM, 1, Ztype, dest, tag, comm);

Example 3.38 Handling of unions.

union {
    int    ival;
    float  fval;
} u[1000]

int    utype;

/* All entries of u have identical type; variable

```

```
1      utype keeps track of their current type */
2
3      MPI_Datatype   type[2];
4      int           blocklen[2] = {1,1};
5      MPI_Aint      disp[2];
6      MPI_Datatype   mpi_utype[2];
7      MPI_Aint      i,j;
8
9      /* compute an MPI datatype for each possible union type;
10     assume values are left-aligned in union storage. */
11
12     MPI_Address( u, &i);
13     MPI_Address( u+1, &j);
14     disp[0] = 0; disp[1] = j-i;
15     type[1] = MPI_UB;
16
17     type[0] = MPI_INT;
18     MPI_Type_struct(2, blocklen, disp, type, &mpi_utype[0]);
19
20     type[0] = MPI_FLOAT;
21     MPI_Type_struct(2, blocklen, disp, type, &mpi_utype[1]);
22
23     for(i=0; i<2; i++) MPI_Type_commit(&mpi_utype[i]);
24
25     /* actual communication */
26
27     MPI_Send(u, 1000, mpi_utype[utype], dest, tag, comm);
28
29
```

3.13 Pack and unpack

```
31
32     Some existing communication libraries provide pack/unpack functions for sending noncon-
33     tiguous data. In these, the user explicitly packs data into a contiguous buffer before sending
34     it, and unpacks it from a contiguous buffer after receiving it. Derived datatypes, which are
35     described in Section 3.12, allow one, in most cases, to avoid explicit packing and unpacking.
36     The user specifies the layout of the data to be sent or received, and the communication
37     library directly accesses a noncontiguous buffer. The pack/unpack routines are provided
38     for compatibility with previous libraries. Also, they provide some functionality that is not
39     otherwise available in MPI. For instance, a message can be received in several parts, where
40     the receive operation done on a later part may depend on the content of a former part.
41     Another use is that outgoing messages may be explicitly buffered in user supplied space,
42     thus overriding the system buffering policy. Finally, the availability of pack and unpack
43     operations facilitates the development of additional communication libraries layered on top
44     of MPI.
45
46
47
48
```

| | | | |
|---|----------|--|---|
| MPI_PACK(inbuf, incount, datatype, outbuf, outsize, position, comm) | | | 1 |
| IN | inbuf | input buffer start (choice) | 2 |
| IN | incount | number of input data items (integer) | 3 |
| IN | datatype | datatype of each input data item (handle) | 4 |
| OUT | outbuf | output buffer start (choice) | 5 |
| IN | outsize | output buffer size, in bytes (integer) | 6 |
| INOUT | position | current position in buffer, in bytes (integer) | 7 |
| IN | comm | communicator for packed message (handle) | 8 |

```
int MPI_Pack(void* inbuf, int incount, MPI_Datatype datatype, void *outbuf,
            int outsize, int *position, MPI_Comm comm) 9
```

```
MPI_PACK(INBUF, INCOUNT, DATATYPE, OUTBUF, OUTSIZE, POSITION, COMM, IERROR) 10
<type> INBUF(*), OUTBUF(*) 11
INTEGER INCOUNT, DATATYPE, OUTSIZE, POSITION, COMM, IERROR 12
```

```
void MPI::Datatype::Pack(const void* inbuf, int incount, void *outbuf,
                        int outsize, int& position, const MPI::Comm &comm) const 13
```

Packs the message in the send buffer specified by `inbuf`, `incount`, `datatype` into the buffer space specified by `outbuf` and `outsize`. The input buffer can be any communication buffer allowed in `MPI_SEND`. The output buffer is a contiguous storage area containing `outsize` bytes, starting at the address `outbuf` (length is counted in bytes, not elements, as if it were a communication buffer for a message of type `MPI_PACKED`).

The input value of `position` is the first location in the output buffer to be used for packing. `position` is incremented by the size of the packed message, and the output value of `position` is the first location in the output buffer following the locations occupied by the packed message. The `comm` argument is the communicator that will be subsequently used for sending the packed message.

| | | | |
|---|----------|--|----|
| MPI_UNPACK(inbuf, insize, position, outbuf, outcount, datatype, comm) | | | 14 |
| IN | inbuf | input buffer start (choice) | 15 |
| IN | insize | size of input buffer, in bytes (integer) | 16 |
| INOUT | position | current position in bytes (integer) | 17 |
| OUT | outbuf | output buffer start (choice) | 18 |
| IN | outcount | number of items to be unpacked (integer) | 19 |
| IN | datatype | datatype of each output data item (handle) | 20 |
| IN | comm | communicator for packed message (handle) | 21 |

```
int MPI_Unpack(void* inbuf, int insize, int *position, void *outbuf,
              int outcount, MPI_Datatype datatype, MPI_Comm comm) 22
```

```
MPI_UNPACK(INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT, DATATYPE, COMM, 23
            IERROR) 24
```

```

1     <type> INBUF(*), OUTBUF(*)
2     INTEGER INSIZE, POSITION, OUTCOUNT, DATATYPE, COMM, IERROR
3
4     void MPI::Datatype::Unpack(const void* inbuf, int insize, void *outbuf,
5                               int outcount, int& position, const MPI::Comm& comm) const

```

Unpacks a message into the receive buffer specified by `outbuf`, `outcount`, `datatype` from the buffer space specified by `inbuf` and `insize`. The output buffer can be any communication buffer allowed in `MPI_RECV`. The input buffer is a contiguous storage area containing `insize` bytes, starting at address `inbuf`. The input value of `position` is the first location in the input buffer occupied by the packed message. `position` is incremented by the size of the packed message, so that the output value of `position` is the first location in the input buffer after the locations occupied by the message that was unpacked. `comm` is the communicator used to receive the packed message.

Advice to users. Note the difference between `MPI_RECV` and `MPI_UNPACK`: in `MPI_RECV`, the `count` argument specifies the maximum number of items that can be received. The actual number of items received is determined by the length of the incoming message. In `MPI_UNPACK`, the `count` argument specifies the actual number of items that are unpacked; the “size” of the corresponding message is the increment in `position`. The reason for this change is that the “incoming message size” is not predetermined since the user decides how much to unpack; nor is it easy to determine the “message size” from the number of items to be unpacked. In fact, in a heterogeneous system, this number may not be determined *a priori*. (*End of advice to users.*)

To understand the behavior of `pack` and `unpack`, it is convenient to think of the data part of a message as being the sequence obtained by concatenating the successive values sent in that message. The `pack` operation stores this sequence in the buffer space, as if sending the message to that buffer. The `unpack` operation retrieves this sequence from buffer space, as if receiving a message from that buffer. (It is helpful to think of internal Fortran files or `sscanf` in C, for a similar function.)

Several messages can be successively packed into one **packing unit**. This is effected by several successive **related** calls to `MPI_PACK`, where the first call provides `position = 0`, and each successive call inputs the value of `position` that was output by the previous call, and the same values for `outbuf`, `outcount` and `comm`. This packing unit now contains the equivalent information that would have been stored in a message by one `send` call with a `send` buffer that is the “concatenation” of the individual `send` buffers.

A packing unit can be sent using type `MPI_PACKED`. Any point to point or collective communication function can be used to move the sequence of bytes that forms the packing unit from one process to another. This packing unit can now be received using any receive operation, with any `datatype`: the type matching rules are relaxed for messages sent with type `MPI_PACKED`.

A message sent with any type (including `MPI_PACKED`) can be received using the type `MPI_PACKED`. Such a message can then be unpacked by calls to `MPI_UNPACK`.

A packing unit (or a message created by a regular, “typed” `send`) can be unpacked into several successive messages. This is effected by several successive related calls to `MPI_UNPACK`, where the first call provides `position = 0`, and each successive call inputs the value of `position` that was output by the previous call, and the same values for `inbuf`, `insize` and `comm`.

The concatenation of two packing units is not necessarily a packing unit; nor is a substring of a packing unit necessarily a packing unit. Thus, one cannot concatenate two packing units and then unpack the result as one packing unit; nor can one unpack a substring of a packing unit as a separate packing unit. Each packing unit, that was created by a related sequence of pack calls, or by a regular send, must be unpacked as a unit, by a sequence of related unpack calls.

Rationale. The restriction on “atomic” packing and unpacking of packing units allows the implementation to add at the head of packing units additional information, such as a description of the sender architecture (to be used for type conversion, in a heterogeneous environment) (*End of rationale.*)

The following call allows the user to find out how much space is needed to pack a message and, thus, manage space allocation for buffers.

MPI_PACK_SIZE(incount, datatype, comm, size)

| | | |
|-----|----------|---|
| IN | incount | count argument to packing call (integer) |
| IN | datatype | datatype argument to packing call (handle) |
| IN | comm | communicator argument to packing call (handle) |
| OUT | size | upper bound on size of packed message, in bytes (integer) |

```
int MPI_Pack_size(int incount, MPI_Datatype datatype, MPI_Comm comm,
                 int *size)
```

```
MPI_PACK_SIZE(INCOUNT, DATATYPE, COMM, SIZE, IERROR)
INTEGER INCOUNT, DATATYPE, COMM, SIZE, IERROR
```

```
int MPI::Datatype::Pack_size(int incount, const MPI::Comm& comm) const
```

A call to MPI_PACK_SIZE(incount, datatype, comm, size) returns in size an upper bound on the increment in position that is effected by a call to MPI_PACK(inbuf, incount, datatype, outbuf, outcount, position, comm).

Rationale. The call returns an upper bound, rather than an exact bound, since the exact amount of space needed to pack the message may depend on the context (e.g., first message packed in a packing unit may take more space). (*End of rationale.*)

Example 3.39 An example using MPI_PACK.

```
int position, i, j, a[2];
char buff[1000];

....

MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank == 0)
{
```

```

1      / * SENDER CODE */
2
3      position = 0;
4      MPI_Pack(&i, 1, MPI_INT, buff, 1000, &position, MPI_COMM_WORLD);
5      MPI_Pack(&j, 1, MPI_INT, buff, 1000, &position, MPI_COMM_WORLD);
6      MPI_Send( buff, position, MPI_PACKED, 1, 0, MPI_COMM_WORLD);
7  }
8  else /* RECEIVER CODE */
9      MPI_Recv( a, 2, MPI_INT, 0, 0, MPI_COMM_WORLD)
10
11 }

```

Example 3.40 An elaborate example.

```

14 int position, i;
15 float a[1000];
16 char buff[1000]
17
18
19 .....
20
21 MPI_Comm_rank(MPI_Comm_world, &myrank);
22 if (myrank == 0)
23 {
24     / * SENDER CODE */
25
26     int len[2];
27     MPI_Aint disp[2];
28     MPI_Datatype type[2], newtype;
29
30     /* build datatype for i followed by a[0]...a[i-1] */
31
32     len[0] = 1;
33     len[1] = i;
34     MPI_Address( &i, disp);
35     MPI_Address( a, disp+1);
36     type[0] = MPI_INT;
37     type[1] = MPI_FLOAT;
38     MPI_Type_struct( 2, len, disp, type, &newtype);
39     MPI_Type_commit( &newtype);
40
41     /* Pack i followed by a[0]...a[i-1]*/
42
43     position = 0;
44     MPI_Pack( MPI_BOTTOM, 1, newtype, buff, 1000, &position, MPI_COMM_WORLD);
45
46     /* Send */
47
48     MPI_Send( buff, position, MPI_PACKED, 1, 0,

```



```

        MPI_COMM_WORLD)
1
2
/* *****
3
   One can replace the last three lines with
4
   MPI_Send( MPI_BOTTOM, 1, newtype, 1, 0, MPI_COMM_WORLD);
5
   ***** */
6
}
7
else /* myrank == 1 */
8
{
9
    /* RECEIVER CODE */
10
11
    MPI_Status status;
12
13
    /* Receive */
14
15
    MPI_Recv( buff, 1000, MPI_PACKED, 0, 0, &status);
16
17
    /* Unpack i */
18
19
    position = 0;
20
    MPI_Unpack(buff, 1000, &position, &i, 1, MPI_INT, MPI_COMM_WORLD);
21
22
    /* Unpack a[0]...a[i-1] */
23
    MPI_Unpack(buff, 1000, &position, a, i, MPI_FLOAT, MPI_COMM_WORLD);
24
}
25
26
Example 3.41 Each process sends a count, followed by count characters to the root; the
27
root concatenate all characters into one string.
28
29
int count, gsize, counts[64], totalcount, k1, k2, k,
30
    displs[64], position, concat_pos;
31
char chr[100], *lbuf, *rbuf, *cbuf;
32
...
33
MPI_Comm_size(comm, &gsize);
34
MPI_Comm_rank(comm, &myrank);
35
36
    /* allocate local pack buffer */
37
MPI_Pack_size(1, MPI_INT, comm, &k1);
38
MPI_Pack_size(count, MPI_CHAR, comm, &k2);
39
k = k1+k2;
40
lbuf = (char *)malloc(k);
41
42
    /* pack count, followed by count characters */
43
position = 0;
44
MPI_Pack(&count, 1, MPI_INT, lbuf, k, &position, comm);
45
MPI_Pack(chr, count, MPI_CHAR, lbuf, k, &position, comm);
46
47
if (myrank != root) {
48

```

```

1      /* gather at root sizes of all packed messages */
2      MPI_Gather( &position, 1, MPI_INT, NULL, NULL,
3                NULL, root, comm);
4
5      /* gather at root packed messages */
6      MPI_Gatherv( &buf, position, MPI_PACKED, NULL,
7                 NULL, NULL, NULL, root, comm);
8
9  } else { /* root code */
10     /* gather sizes of all packed messages */
11     MPI_Gather( &position, 1, MPI_INT, counts, 1,
12               MPI_INT, root, comm);
13
14     /* gather all packed messages */
15     displs[0] = 0;
16     for (i=1; i < gsize; i++)
17         displs[i] = displs[i-1] + counts[i-1];
18     totalcount = displs[gsiz-1] + counts[gsiz-1];
19     rbuf = (char *)malloc(totalcount);
20     cbuf = (char *)malloc(totalcount);
21     MPI_Gatherv( lbuf, position, MPI_PACKED, rbuf,
22                counts, displs, MPI_PACKED, root, comm);
23
24     /* unpack all messages and concatenate strings */
25     concat_pos = 0;
26     for (i=0; i < gsize; i++) {
27         position = 0;
28         MPI_Unpack( rbuf+displs[i], totalcount-displs[i],
29                   &position, &count, 1, MPI_INT, comm);
30         MPI_Unpack( rbuf+displs[i], totalcount-displs[i],
31                   &position, cbuf+concat_pos, count, MPI_CHAR, comm);
32         concat_pos += count;
33     }
34     cbuf[concat_pos] = '\0';
35 }

```

3.14 Canonical MPI_PACK and MPI_UNPACK

These functions read/write data to/from the buffer in the “external32” data format specified in Section 12.5.2, and calculate the size needed for packing. Their first arguments specify the data format, for future extensibility, but for MPI-2 the only valid value of the `datarep` argument is “external32.”

Advice to users. These functions could be used, for example, to send typed data in a portable format from one MPI implementation to another. (*End of advice to users.*)

The buffer will contain exactly the packed data, without headers. `MPI_BYTE` should be used to send and receive data that is packed using `MPI_PACK_EXTERNAL`.

Rationale. MPI_PACK_EXTERNAL specifies that there is no header on the message and further specifies the exact format of the data. Since MPI_PACK may (and is allowed to) use a header, the datatype MPI_PACKED cannot be used for data packed with MPI_PACK_EXTERNAL. (*End of rationale.*)

MPI_PACK_EXTERNAL(datarep, inbuf, incount, datatype, outbuf, outsize, position)

| | | |
|-------|----------|--|
| IN | datarep | data representation (string) |
| IN | inbuf | input buffer start (choice) |
| IN | incount | number of input data items (integer) |
| IN | datatype | datatype of each input data item (handle) |
| OUT | outbuf | output buffer start (choice) |
| IN | outsize | output buffer size, in bytes (integer) |
| INOUT | position | current position in buffer, in bytes (integer) |

```
int MPI_Pack_external(char *datarep, void *inbuf, int incount,
                    MPI_Datatype datatype, void *outbuf, MPI_Aint outsize,
                    MPI_Aint *position)
```

```
MPI_PACK_EXTERNAL(DATAREP, INBUF, INCOUNT, DATATYPE, OUTBUF, OUTSIZE,
                 POSITION, IERROR)
INTEGER INCOUNT, DATATYPE, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) OUTSIZE, POSITION
CHARACTER*(*) DATAREP
<type> INBUF(*), OUTBUF(*)
```

```
void MPI::Datatype::Pack_external(const char* datarep, const void* inbuf,
                                 int incount, void* outbuf, MPI::Aint outsize,
                                 MPI::Aint& position) const
```

MPI_UNPACK_EXTERNAL(datarep, inbuf, insize, position, outbuf, outsize, position)

| | | |
|-------|----------|--|
| IN | datarep | data representation (string) |
| IN | inbuf | input buffer start (choice) |
| IN | insize | input buffer size, in bytes (integer) |
| INOUT | position | current position in buffer, in bytes (integer) |
| OUT | outbuf | output buffer start (choice) |
| IN | outcount | number of output data items (integer) |
| IN | datatype | datatype of output data item (handle) |

```
int MPI_Unpack_external(char *datarep, void *inbuf, MPI_Aint insize,
                      MPI_Aint *position, void *outbuf, int outcount,
                      MPI_Datatype datatype)
```

```

1 MPI_UNPACK_EXTERNAL(DATAREP, INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT,
2     DATATYPE, IERROR)
3     INTEGER OUTCOUNT, DATATYPE, IERROR
4     INTEGER(KIND=MPI_ADDRESS_KIND) INSIZE, POSITION
5     CHARACTER*(*) DATAREP
6     <type> INBUF(*), OUTBUF(*)
7
8 void MPI::Datatype::Unpack_external(const char* datarep, const void* inbuf,
9     MPI::Aint insize, MPI::Aint& position, void* outbuf,
10    int outcount) const
11
12
13 MPI_PACK_EXTERNAL_SIZE( datarep, incount, datatype, size )
14     IN        datarep                data representation (string)
15     IN        incount                 number of input data items (integer)
16     IN        datatype               datatype of each input data item (handle)
17     OUT       size                   output buffer size, in bytes (integer)
18
19
20 int MPI_Pack_external_size(char *datarep, int incount,
21     MPI_Datatype datatype, MPI_Aint *size)
22
23 MPI_PACK_EXTERNAL_SIZE(DATAREP, INCOUNT, DATATYPE, SIZE, IERROR)
24     INTEGER INCOUNT, DATATYPE, IERROR
25     INTEGER(KIND=MPI_ADDRESS_KIND) SIZE
26     CHARACTER*(*) DATAREP
27
28 MPI::Aint MPI::Datatype::Pack_external_size(const char* datarep,
29     int incount) const
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

```

Chapter 4

Collective Communication

4.1 Introduction and Overview

Collective communication is defined as communication that involves a group of processes. The functions of this type provided by MPI are the following:

- Barrier synchronization across all group members (Sec. 4.4).
- Broadcast from one member to all members of a group (Sec. 4.5). This is shown in figure 4.1.
- Gather data from all group members to one member (Sec. 4.6). This is shown in figure 4.1.
- Scatter data from one member to all members of a group (Sec. 4.7). This is shown in figure 4.1.
- A variation on Gather where all members of the group receive the result (Sec. 4.8). This is shown as “allgather” in figure 4.1.
- Scatter/Gather data from all members to all members of a group (also called complete exchange or all-to-all) (Sec. 4.9). This is shown as “alltoall” in figure 4.1.
- Global reduction operations such as sum, max, min, or user-defined functions, where the result is returned to all group members and a variation where the result is returned to only one member (Sec. 4.10).
- A combined reduction and scatter operation (Sec. 4.11).
- Scan across all members of a group (also called prefix) (Sec. 4.12).

A collective operation is executed by having all processes in the group call the communication routine, with matching arguments. The syntax and semantics of the collective operations are defined to be consistent with the syntax and semantics of the point-to-point operations. Thus, general datatypes are allowed and must match between sending and receiving processes as specified in Chapter 3. One of the key arguments is a communicator that defines the group of participating processes and provides a context for the operation. Several collective routines such as broadcast and gather have a single originating or receiving process. Such processes are called the *root*. Some arguments in the collective functions

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

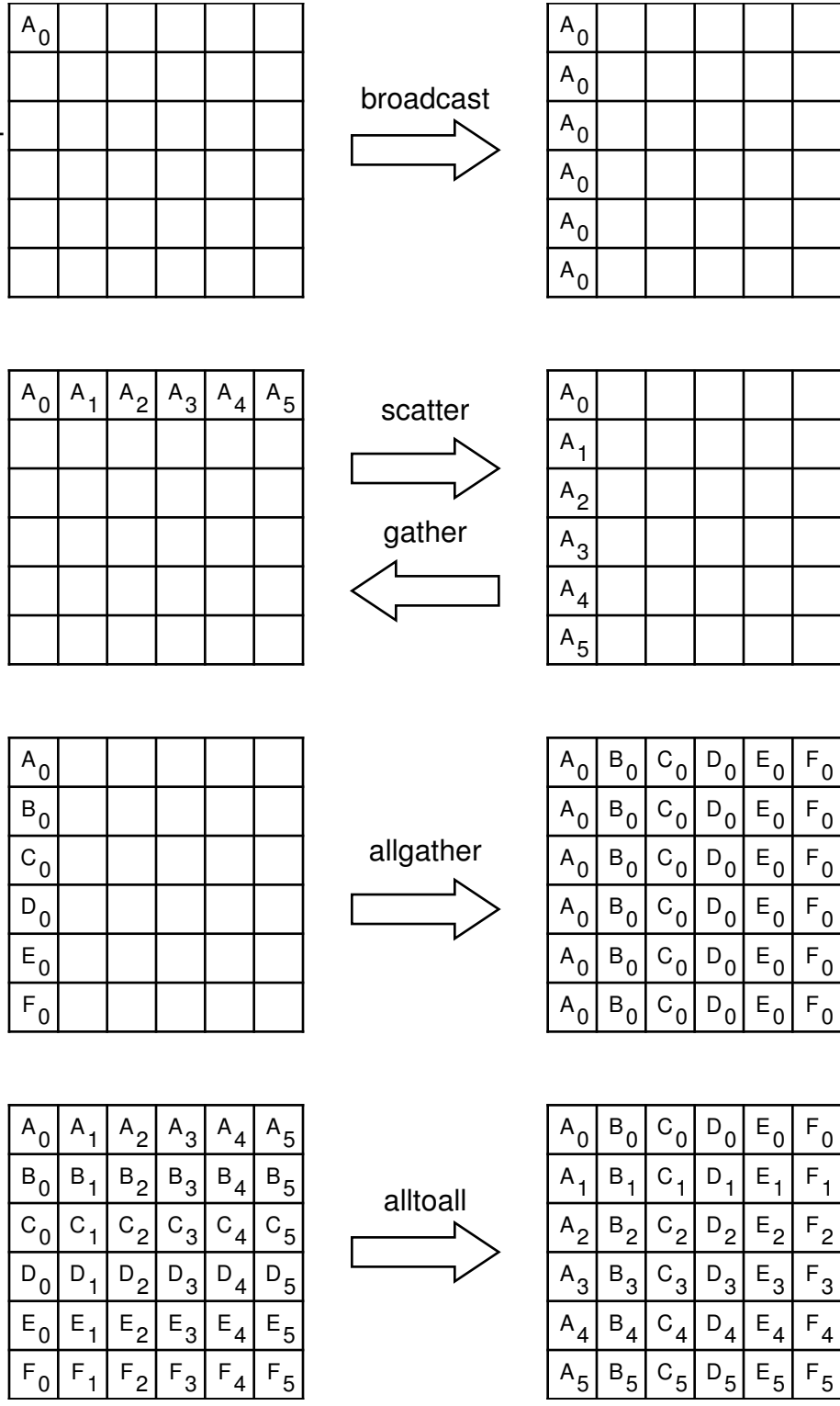


Figure 4.1: Collective move functions illustrated for a group of six processes. In each case, each row of boxes represents data locations in one process. Thus, in the broadcast, initially just the first process contains the data A_0 , but after the broadcast all processes contain it.

are specified as “significant only at root,” and are ignored for all participants except the root. The reader is referred to Chapter 3 for information concerning communication buffers, general datatypes and type matching rules, and to Chapter 5 for information on how to define groups and create communicators.

The type-matching conditions for the collective operations are more strict than the corresponding conditions between sender and receiver in point-to-point. Namely, for collective operations, the amount of data sent must exactly match the amount of data specified by the receiver. Distinct type maps (the layout in memory, see Sec. 3.12) between sender and receiver are still allowed.

Collective routine calls can (but are not required to) return as soon as their participation in the collective communication is complete. The completion of a call indicates that the caller is now free to access locations in the communication buffer. It does not indicate that other processes in the group have completed or even started the operation (unless otherwise indicated in the description of the operation). Thus, a collective communication call may, or may not, have the effect of synchronizing all calling processes. This statement excludes, of course, the barrier function.

Collective communication calls may use the same communicators as point-to-point communication; MPI guarantees that messages generated on behalf of collective communication calls will not be confused with messages generated by point-to-point communication. A more detailed discussion of correct use of collective routines is found in Sec. 4.13.

Rationale. The equal-data restriction (on type matching) was made so as to avoid the complexity of providing a facility analogous to the status argument of MPI_RECV for discovering the amount of data sent. Some of the collective routines would require an array of status values.

The statements about synchronization are made so as to allow a variety of implementations of the collective functions.

The collective operations do not accept a message tag argument. If future revisions of MPI define non-blocking collective functions, then tags (or a similar mechanism) will need to be added so as to allow the dis-ambiguation of multiple, pending, collective operations. (*End of rationale.*)

Advice to users. It is dangerous to rely on synchronization side-effects of the collective operations for program correctness. For example, even though a particular implementation may provide a broadcast routine with a side-effect of synchronization, the standard does not require this, and a program that relies on this will not be portable.

On the other hand, a correct, portable program must allow for the fact that a collective call *may* be synchronizing. Though one cannot rely on any synchronization side-effect, one must program so as to allow it. These issues are discussed further in Sec. 4.13. (*End of advice to users.*)

Advice to implementors. While vendors may write optimized collective routines matched to their architectures, a complete library of the collective communication routines can be written entirely using the MPI point-to-point communication functions and a few auxiliary functions. If implementing on top of point-to-point, a hidden,

special communicator must be created for the collective operation so as to avoid interference with any on-going point-to-point communication at the time of the collective call. This is discussed further in Sec. 4.13. (*End of advice to implementors.*)

Many of the descriptions of the collective routines provide illustrations in terms of blocking MPI point-to-point routines. These are intended solely to indicate what data is sent or received by what process. Many of these examples are *not* correct MPI programs; for purposes of simplicity, they often assume infinite buffering.

4.2 Communicator argument

The key concept of the collective functions is to have a “group” of participating processes. The routines do not have a group identifier as an explicit argument. Instead, there is a communicator argument. For the purposes of this chapter, a communicator can be thought of as a group identifier linked with a context.

4.3 Extended Collective Operations

4.3.1 Introduction

MPI-1 defined collective communication for intracommunicators and two routines, `MPI_INTERCOMM_CREATE` and `MPI_COMM_DUP`, for creating new intercommunicators. In addition, in order to avoid argument aliasing problems with Fortran, MPI-1 requires separate send and receive buffers for collective operations. MPI-2 introduces extensions of many of the MPI-1 collective routines to intercommunicators, additional routines for creating intercommunicators, and two new collective routines: a generalized all-to-all and an exclusive scan. In addition, a way to specify “in place” buffers is provided for many of the intracommunicator collective operations.

In addition, the specification of collective operations (Section 4.1 of MPI-1) requires that all collective routines are called with matching arguments. For the intercommunicator extensions, this is weakened to matching for all members of the same local group.

4.3.2 Intercommunicator Collective Operations

In the MPI-1 standard (Section 4.2), collective operations only apply to intracommunicators; however, most MPI collective operations can be generalized to intercommunicators. To understand how MPI can be extended, we can view most MPI intracommunicator collective operations as fitting one of the following categories (see, for instance, [45]):

All-To-All All processes contribute to the result. All processes receive the result.

- `MPI_Allgather`, `MPI_Allgatherv`
- `MPI_Alltoall`, `MPI_Alltoallv`
- `MPI_Allreduce`, `MPI_Reduce_scatter`

All-To-One All processes contribute to the result. One process receives the result.

- `MPI_Gather`, `MPI_Gatherv`
- `MPI_Reduce`

One-To-All One process contributes to the result. All processes receive the result.

- MPI_Bcast
- MPI_Scatter, MPI_Scatterv

Other Collective operations that do not fit into one of the above categories.

- MPI_Scan
- MPI_Barrier

The MPI_Barrier operation does not fit into this classification since no data is being moved (other than the implicit fact that a barrier has been called). The data movement pattern of MPI_Scan does not fit this taxonomy.

The extension of collective communication from intracommunicators to intercommunicators is best described in terms of the left and right groups. For example, an all-to-all MPI_Allgather operation can be described as collecting data from all members of one group with the result appearing in all members of the other group (see Figure 4.2). As another example, a one-to-all MPI_Bcast operation sends data from one member of one group to all members of the other group. Collective computation operations such as MPI_REDUCE_SCATTER have a similar interpretation (see Figure 4.3). For intracommunicators, these two groups are the same. For intercommunicators, these two groups are distinct. For the all-to-all operations, each such operation is described in two phases, so that it has a symmetric, full-duplex behavior.

For MPI-2, the following intracommunicator collective operations also apply to intercommunicators:

- MPI_BCAST,
- MPI_GATHER, MPI_GATHERV,
- MPI_SCATTER, MPI_SCATTERV,
- MPI_ALLGATHER, MPI_ALLGATHERV,
- MPI_ALLTOALL, MPI_ALLTOALLV, MPI_ALLTOALLW,
- MPI_REDUCE, MPI_ALLREDUCE,
- MPI_REDUCE_SCATTER,
- MPI_BARRIER.

(MPI_ALLTOALLW is a new function described in Section 4.9.1.)

These functions use exactly the same argument list as their MPI-1 counterparts and also work on intracommunicators, as expected. No new language bindings are consequently needed for Fortran or C. However, in C++, the bindings have been "relaxed"; these member functions have been moved from the `MPI::Intercomm` class to the `MPI::Comm` class. But since the collective operations do not make sense on a C++ `MPI::Comm` (since it is neither an intercommunicator nor an intracommunicator), the functions are all pure virtual. In an MPI-2 implementation, the bindings in this chapter supersede the corresponding bindings for MPI-1.2.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

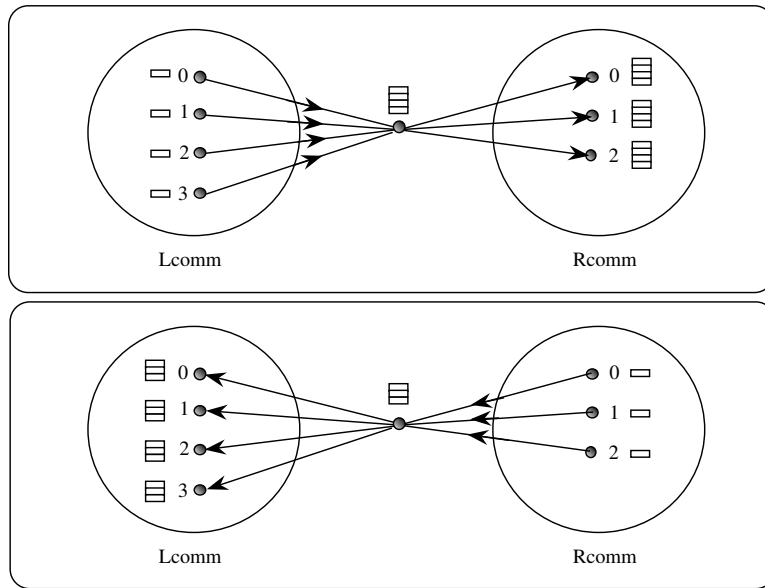


Figure 4.2: Intercommunicator allgather. The focus of data to one process is represented, not mandated by the semantics. The two phases do allgathers in both directions.

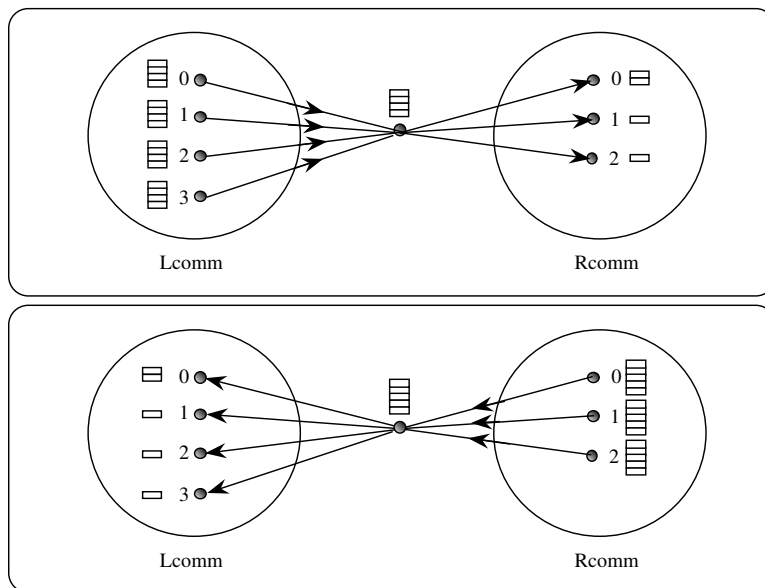


Figure 4.3: Intercommunicator reduce-scatter. The focus of data to one process is represented, not mandated by the semantics. The two phases do reduce-scatters in both directions.

4.3.3 Operations that Move Data

Two additions are made to many collective communication calls:

- Collective communication can occur “in place” for intracommunicators, with the output buffer being identical to the input buffer. This is specified by providing a special argument value, `MPI_IN_PLACE`, instead of the send buffer or the receive buffer argument.

Rationale. The “in place” operations are provided to reduce unnecessary memory motion by both the MPI implementation and by the user. Note that while the simple check of testing whether the send and receive buffers have the same address will work for some cases (e.g., `MPI_ALLREDUCE`), they are inadequate in others (e.g., `MPI_GATHER`, with root not equal to zero). Further, Fortran explicitly prohibits aliasing of arguments; the approach of using a special value to denote “in place” operation eliminates that difficulty. (*End of rationale.*)

Advice to users. By allowing the “in place” option, the receive buffer in many of the collective calls becomes a send-and-receive buffer. For this reason, a Fortran binding that includes `INTENT` must mark these as `INOUT`, not `OUT`.

Note that `MPI_IN_PLACE` is a special kind of value; it has the same restrictions on its use that `MPI_BOTTOM` has.

Some intracommunicator collective operations do not support the “in place” option (e.g., `MPI_ALLTOALLV`). (*End of advice to users.*)

- Collective communication applies to intercommunicators. If the operation is rooted (e.g., broadcast, gather, scatter), then the transfer is unidirectional. The direction of the transfer is indicated by a special value of the root argument. In this case, for the group containing the root process, all processes in the group must call the routine using a special argument for the root. The root process uses the special root value `MPI_ROOT`; all other processes in the same group as the root use `MPI_PROC_NULL`. All processes in the other group (the group that is the remote group relative to the root process) must call the collective routine and provide the rank of the root. If the operation is unrooted (e.g., `alltoall`), then the transfer is bidirectional.

Note that the “in place” option for intracommunicators does not apply to intercommunicators since in the intercommunicator case there is no communication from a process to itself.

Rationale. Rooted operations are unidirectional by nature, and there is a clear way of specifying direction. Non-rooted operations, such as all-to-all, will often occur as part of an exchange, where it makes sense to communicate in both directions at once. (*End of rationale.*)

4.4 Barrier synchronization

`MPI_BARRIER(comm)`

IN comm communicator (handle)

```
int MPI_Barrier(MPI_Comm comm )
```

```
MPI_BARRIER(COMM, IERROR)
          INTEGER COMM, IERROR
```

```
void MPI::Comm::Barrier() const = 0
```

If `comm` is an [intracommunicator](#), `MPI_BARRIER` blocks the caller until all group members have called it. The call returns at any process only after all group members have entered the call.

If `comm` is an [intercommunicator](#), the barrier is performed across all processes in the intercommunicator. In this case, all processes in the local group of the intercommunicator may exit the barrier when all of the processes in the remote group have entered the barrier.

4.5 Broadcast

`MPI_BCAST(buffer, count, datatype, root, comm)`

INOUT buffer starting address of buffer (choice)

IN count number of entries in buffer (integer)

IN datatype data type of buffer (handle)

IN root rank of broadcast root (integer)

IN comm communicator (handle)

```
int MPI_Bcast(void* buffer, int count, MPI_Datatype datatype, int root,
              MPI_Comm comm )
```

```
MPI_BCAST(BUFFER, COUNT, DATATYPE, ROOT, COMM, IERROR)
          <type> BUFFER(*)
          INTEGER COUNT, DATATYPE, ROOT, COMM, IERROR
```

```
void MPI::Comm::Bcast(void* buffer, int count,
                      const MPI::Datatype& datatype, int root) const = 0
```

`MPI_BCAST` broadcasts a message from the process with rank `root` to all processes of the group, itself included for [intracommunicators](#). It is called by all members of group using the same arguments for `comm`, `root`. On return, the contents of `root`'s communication buffer has been copied to all processes.

General, derived datatypes are allowed for `datatype`. The type signature of `count`, `datatype` on any process must be equal to the type signature of `count`, `datatype` at the root. This implies that the amount of data sent must be equal to the amount received, pairwise

between each process and the root. `MPI_BCAST` and all other data-movement collective routines make this restriction. Distinct type maps between sender and receiver are still allowed.

The “in place” option is not meaningful here.

If `comm` is an intercommunicator, then the call involves all processes in the intercommunicator, but with one group (group A) defining the root process. All processes in the other group (group B) pass the same value in argument `root`, which is the rank of the root in group A. The root passes the value `MPI_ROOT` in `root`. All other processes in group A pass the value `MPI_PROC_NULL` in `root`. Data is broadcast from the root to all processes in group B. The receive buffer arguments of the processes in group B must be consistent with the send buffer argument of the root.

4.5.1 Example using `MPI_BCAST`

Example 4.1 Broadcast 100 ints from process 0 to every process in the group.

```
MPI_Comm comm;
int array[100];
int root=0;
...
MPI_Bcast( array, 100, MPI_INT, root, comm);
```

As in many of our example code fragments, we assume that some of the variables (such as `comm` in the above) have been assigned appropriate values.

4.6 Gather

`MPI_GATHER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)`

| | | |
|-----|------------------------|---|
| IN | <code>sendbuf</code> | starting address of send buffer (choice) |
| IN | <code>sendcount</code> | number of elements in send buffer (integer) |
| IN | <code>sendtype</code> | data type of send buffer elements (handle) |
| OUT | <code>recvbuf</code> | address of receive buffer (choice, significant only at root) |
| IN | <code>recvcount</code> | number of elements for any single receive (integer, significant only at root) |
| IN | <code>recvtype</code> | data type of recv buffer elements (significant only at root) (handle) |
| IN | <code>root</code> | rank of receiving process (integer) |
| IN | <code>comm</code> | communicator (handle) |

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

```

1 MPI_GATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
2           ROOT, COMM, IERROR)
3   <type> SENDBUF(*), RECVBUF(*)
4   INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, ROOT, COMM, IERROR
5
6 void MPI::Comm::Gather(const void* sendbuf, int sendcount, const
7                       MPI::Datatype& sendtype, void* recvbuf, int recvcount,
8                       const MPI::Datatype& recvtpe, int root) const = 0

```

If `comm` is an intracommunicator, each process (root process included) sends the contents of its send buffer to the root process. The root process receives the messages and stores them in rank order. The outcome is *as if* each of the `n` processes in the group (including the root process) had executed a call to

```
MPI_Send(sendbuf, sendcount, sendtype, root, ...),
```

and the root had executed `n` calls to

```
MPI_Recv(recvbuf + i · recvcount · extent(recvtpe), recvcount, recvtpe, i, ...),
```

where `extent(recvtpe)` is the type extent obtained from a call to `MPI_Type_extent()`.

An alternative description is that the `n` messages sent by the processes in the group are concatenated in rank order, and the resulting message is received by the root as if by a call to `MPI_RECV(recvbuf, recvcount·n, recvtpe, ...)`.

The receive buffer is ignored for all non-root processes.

General, derived datatypes are allowed for both `sendtype` and `recvtpe`. The type signature of `sendcount`, `sendtype` on process `i` must be equal to the type signature of `recvcount`, `recvtpe` at the root. This implies that the amount of data sent must be equal to the amount of data received, pairwise between each process and the root. Distinct type maps between sender and receiver are still allowed.

All arguments to the function are significant on process `root`, while on other processes, only arguments `sendbuf`, `sendcount`, `sendtype`, `root`, `comm` are significant. The arguments `root` and `comm` must have identical values on all processes.

The specification of counts and types should not cause any location on the root to be written more than once. Such a call is erroneous.

Note that the `recvcount` argument at the root indicates the number of items it receives from *each* process, not the total number of items it receives.

The “in place” option for intracommunicators is specified by passing `MPI_IN_PLACE` as the value of `sendbuf` at the root. In such a case, `sendcount` and `sendtype` are ignored, and the contribution of the root to the gathered vector is assumed to be already in the correct place in the receive buffer

If `comm` is an intercommunicator, then the call involves all processes in the intercommunicator, but with one group (group A) defining the root process. All processes in the other group (group B) pass the same value in argument `root`, which is the rank of the root in group A. The root passes the value `MPI_ROOT` in `root`. All other processes in group A pass the value `MPI_PROC_NULL` in `root`. Data is gathered from all processes in group B to the root. The send buffer arguments of the processes in group B must be consistent with the receive buffer argument of the root.

| | | | | | | | | | |
|--------------|------------|------------|-----------|---|-------------|---------|-----------|-------|----|
| MPI_GATHERV(| sendbuf, | sendcount, | sendtype, | recvbuf, | recvcounts, | displs, | recvtype, | root, | 1 |
| comm) | | | | | | | | | 2 |
| IN | sendbuf | | | starting address of send buffer (choice) | | | | | 3 |
| IN | sendcount | | | number of elements in send buffer (integer) | | | | | 4 |
| IN | sendtype | | | data type of send buffer elements (handle) | | | | | 5 |
| OUT | recvbuf | | | address of receive buffer (choice, significant only at root) | | | | | 6 |
| | | | | | | | | | 7 |
| IN | recvcounts | | | integer array (of length group size) containing the number of elements that are received from each process (significant only at root) | | | | | 8 |
| | | | | | | | | | 9 |
| IN | displs | | | integer array (of length group size). Entry <i>i</i> specifies the displacement relative to <i>recvbuf</i> at which to place the incoming data from process <i>i</i> (significant only at root) | | | | | 10 |
| | | | | | | | | | 11 |
| | | | | | | | | | 12 |
| IN | recvtype | | | data type of recv buffer elements (significant only at root) (handle) | | | | | 13 |
| | | | | | | | | | 14 |
| | | | | | | | | | 15 |
| IN | root | | | rank of receiving process (integer) | | | | | 16 |
| | | | | | | | | | 17 |
| IN | comm | | | communicator (handle) | | | | | 18 |
| | | | | | | | | | 19 |
| | | | | | | | | | 20 |
| | | | | | | | | | 21 |

```

int MPI_Gatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                void* recvbuf, int *recvcounts, int *displs,
                MPI_Datatype recvtype, int root, MPI_Comm comm)
MPI_GATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
            RECVTYPE, ROOT, COMM, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPE, ROOT,
COMM, IERROR
void MPI::Comm::Gatherv(const void* sendbuf, int sendcount, const
                        MPI::Datatype& sendtype, void* recvbuf,
                        const int recvcounts[], const int displs[],
                        const MPI::Datatype& recvtype, int root) const = 0

```

MPI_GATHERV extends the functionality of MPI_GATHER by allowing a varying count of data from each process, since *recvcounts* is now an array. It also allows more flexibility as to where the data is placed on the root, by providing the new argument, *displs*.

If *comm* is an [intracommunicator](#), the outcome is *as if* each process, including the root process, sends a message to the root,

```
MPI_Send(sendbuf, sendcount, sendtype, root, ...),
```

and the root executes *n* receives,

```
MPI_Recv(recvbuf + displs[i] · extent(recvtype), recvcounts[i], recvtype, i, ...).
```

Messages are placed in the receive buffer of the root process in rank order, that is, the data sent from process *j* is placed in the *j*th portion of the receive buffer *recvbuf* on process

1 root. The j th portion of `recvbuf` begins at offset `displs[j]` elements (in terms of `recvtype`) into
 2 `recvbuf`.

3 The receive buffer is ignored for all non-root processes.

4 The type signature implied by `sendcount`, `sendtype` on process `i` must be equal to the
 5 type signature implied by `recvcounts[i]`, `recvtype` at the root. This implies that the amount
 6 of data sent must be equal to the amount of data received, pairwise between each process
 7 and the root. Distinct type maps between sender and receiver are still allowed, as illustrated
 8 in Example 4.6.

9 All arguments to the function are significant on process `root`, while on other processes,
 10 only arguments `sendbuf`, `sendcount`, `sendtype`, `root`, `comm` are significant. The arguments
 11 `root` and `comm` must have identical values on all processes.

12 The specification of counts, types, and displacements should not cause any location on
 13 the root to be written more than once. Such a call is erroneous.

14 The “in place” option for intracommunicators is specified by passing `MPI_IN_PLACE` as
 15 the value of `sendbuf` at the root. In such a case, `sendcount` and `sendtype` are ignored, and
 16 the contribution of the root to the gathered vector is assumed to be already in the correct
 17 place in the receive buffer

18 If `comm` is an intercommunicator, then the call involves all processes in the intercom-
 19 municator, but with one group (group A) defining the root process. All processes in the
 20 other group (group B) pass the same value in argument `root`, which is the rank of the root
 21 in group A. The root passes the value `MPI_ROOT` in `root`. All other processes in group A
 22 pass the value `MPI_PROC_NULL` in `root`. Data is gathered from all processes in group B to
 23 the root. The send buffer arguments of the processes in group B must be consistent with
 24 the receive buffer argument of the root.

26 4.6.1 Examples using `MPI_GATHER`, `MPI_GATHERV`

27 The examples in this section are using intracommunicators.

29 **Example 4.2** Gather 100 ints from every process in group to root. See figure 4.4.

```
31 MPI_Comm comm;
32 int gsize, sendarray[100];
33 int root, *rbuf;
34 ...
35 MPI_Comm_size( comm, &gsize);
36 rbuf = (int *)malloc(gsize*100*sizeof(int));
37 MPI_Gather( sendarray, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);
```

39 **Example 4.3** Previous example modified – only the root allocates memory for the receive
 40 buffer.

```
41 MPI_Comm comm;
42 int gsize, sendarray[100];
43 int root, myrank, *rbuf;
44 ...
45 MPI_Comm_rank( comm, myrank);
46 if ( myrank == root) {
47     MPI_Comm_size( comm, &gsize);
```

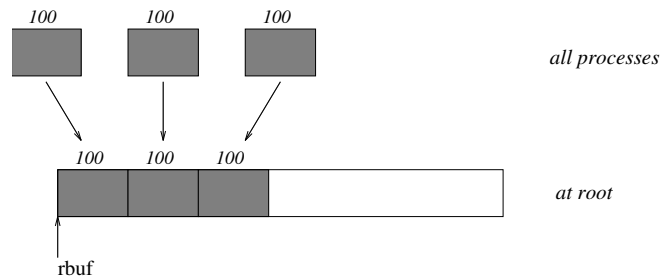



Figure 4.4: The root process gathers 100 ints from each process in the group.

```

    rbuf = (int *)malloc(gsize*100*sizeof(int));
}
MPI_Gather( sendarray, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);

```

Example 4.4 Do the same as the previous example, but use a derived datatype. Note that the type cannot be the entire set of `gsize*100` ints since type matching is defined pairwise between the root and each process in the gather.

```

MPI_Comm comm;
int gsize, sendarray[100];
int root, *rbuf;
MPI_Datatype rtype;
...
MPI_Comm_size( comm, &gsize);
MPI_Type_contiguous( 100, MPI_INT, &rtype );
MPI_Type_commit( &rtype );
rbuf = (int *)malloc(gsize*100*sizeof(int));
MPI_Gather( sendarray, 100, MPI_INT, rbuf, 1, rtype, root, comm);

```

Example 4.5 Now have each process send 100 ints to root, but place each set (of 100) *stride* ints apart at receiving end. Use `MPI_GATHERV` and the `displs` argument to achieve this effect. Assume *stride* ≥ 100 . See figure 4.5.

```

MPI_Comm comm;
int gsize, sendarray[100];
int root, *rbuf, stride;
int *displs, i, *rcounts;
...
MPI_Comm_size( comm, &gsize);
rbuf = (int *)malloc(gsize*stride*sizeof(int));
displs = (int *)malloc(gsize*sizeof(int));
rcounts = (int *)malloc(gsize*sizeof(int));
for (i=0; i<gsize; ++i) {
    displs[i] = i*stride;
    rcounts[i] = 100;
}

```

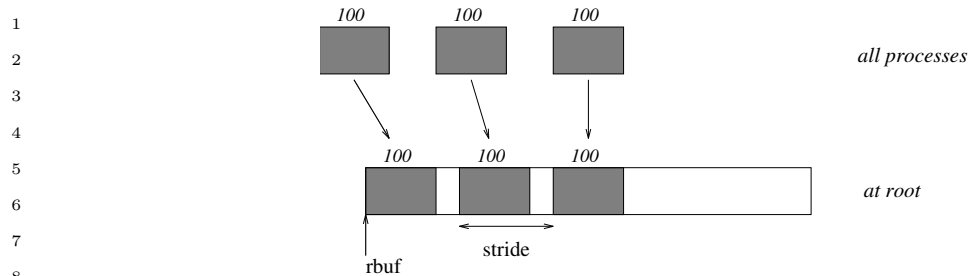


Figure 4.5: The root process gathers 100 ints from each process in the group, each set is placed `stride` ints apart.

```

13     }
14     MPI_Gatherv( sendarray, 100, MPI_INT, rbuf, rcounts, displs, MPI_INT,
15                 root, comm);

```

Note that the program is erroneous if `stride < 100`.

Example 4.6 Same as Example 4.5 on the receiving side, but send the 100 ints from the 0th column of a 100×150 int array, in C. See figure 4.6.

```

22     MPI_Comm comm;
23     int gsize, sendarray[100][150];
24     int root, *rbuf, stride;
25     MPI_Datatype stype;
26     int *displs, i, *rcounts;
27
28     ...
29
30     MPI_Comm_size( comm, &gsize);
31     rbuf = (int *)malloc(gsize*stride*sizeof(int));
32     displs = (int *)malloc(gsize*sizeof(int));
33     rcounts = (int *)malloc(gsize*sizeof(int));
34     for (i=0; i<gsize; ++i) {
35         displs[i] = i*stride;
36         rcounts[i] = 100;
37     }
38     /* Create datatype for 1 column of array
39     */
40     MPI_Type_vector( 100, 1, 150, MPI_INT, &stype);
41     MPI_Type_commit( &stype );
42     MPI_Gatherv( sendarray, 1, stype, rbuf, rcounts, displs, MPI_INT,
43                 root, comm);

```

Example 4.7 Process i sends $(100-i)$ ints from the i th column of a 100×150 int array, in C. It is received into a buffer with stride, as in the previous two examples. See figure 4.7.

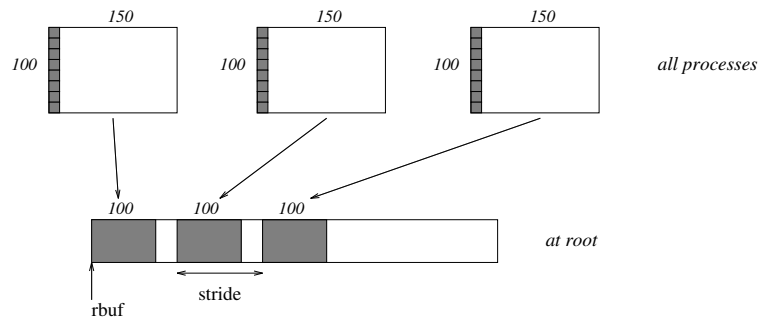


Figure 4.6: The root process gathers column 0 of a 100×150 C array, and each set is placed `stride` ints apart.

```

MPI_Comm comm;
int gsize, sendarray[100][150], *sptr;
int root, *rbuf, stride, myrank;
MPI_Datatype stype;
int *displs, i, *rcounts;
...

MPI_Comm_size( comm, &gsize);
MPI_Comm_rank( comm, &myrank );
rbuf = (int *)malloc(gsize*stride*sizeof(int));
displs = (int *)malloc(gsize*sizeof(int));
rcounts = (int *)malloc(gsize*sizeof(int));
for (i=0; i<gsize; ++i) {
    displs[i] = i*stride;
    rcounts[i] = 100-i;    /* note change from previous example */
}
/* Create datatype for the column we are sending
*/
MPI_Type_vector( 100-myrank, 1, 150, MPI_INT, &stype);
MPI_Type_commit( &stype );
/* sptr is the address of start of "myrank" column
*/
sptr = &sendarray[0][myrank];
MPI_Gatherv( sptr, 1, stype, rbuf, rcounts, displs, MPI_INT,
             root, comm);

```

Note that a different amount of data is received from each process.

Example 4.8 Same as Example 4.7, but done in a different way at the sending end. We create a datatype that causes the correct striding at the sending end so that that we read a column of a C array. A similar thing was done in Example 3.35, Section 3.12.14.

```

MPI_Comm comm;
int gsize, sendarray[100][150], *sptr;

```

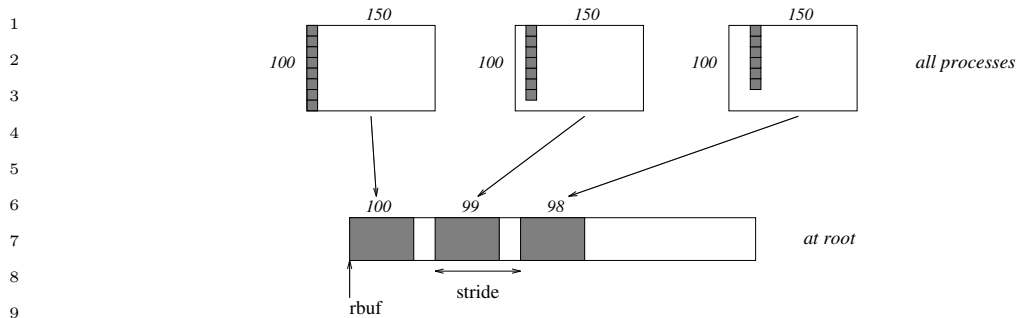


Figure 4.7: The root process gathers $100-i$ ints from column i of a 100×150 C array, and each set is placed stride ints apart.

```

14     int root, *rbuf, stride, myrank, disp[2], blocklen[2];
15     MPI_Datatype stype, type[2];
16     int *displs, i, *rcounts;
17
18     ...
19
20     MPI_Comm_size( comm, &gsize);
21     MPI_Comm_rank( comm, &myrank );
22     rbuf = (int *)malloc(gsize*stride*sizeof(int));
23     displs = (int *)malloc(gsize*sizeof(int));
24     rcounts = (int *)malloc(gsize*sizeof(int));
25     for (i=0; i<gsize; ++i) {
26         displs[i] = i*stride;
27         rcounts[i] = 100-i;
28     }
29     /* Create datatype for one int, with extent of entire row
30        */
31     disp[0] = 0;      disp[1] = 150*sizeof(int);
32     type[0] = MPI_INT; type[1] = MPI_UB;
33     blocklen[0] = 1;  blocklen[1] = 1;
34     MPI_Type_struct( 2, blocklen, disp, type, &stype );
35     MPI_Type_commit( &stype );
36     sptr = &sendarray[0][myrank];
37     MPI_Gatherv( sptr, 100-myrank, stype, rbuf, rcounts, displs, MPI_INT,
38                 root, comm);

```

Example 4.9 Same as Example 4.7 at sending side, but at receiving side we make the stride between received blocks vary from block to block. See figure 4.8.

```

43     MPI_Comm comm;
44     int gsize, sendarray[100][150], *sptr;
45     int root, *rbuf, *stride, myrank, bufsize;
46     MPI_Datatype stype;
47     int *displs, i, *rcounts, offset;

```

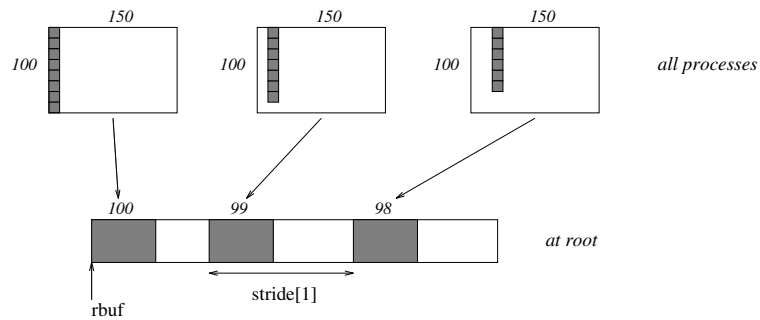


Figure 4.8: The root process gathers $100-i$ ints from column i of a 100×150 C array, and each set is placed $\text{stride}[i]$ ints apart (a varying stride).

```

...
MPI_Comm_size( comm, &gsize);
MPI_Comm_rank( comm, &myrank );

stride = (int *)malloc(gsize*sizeof(int));
...
/* stride[i] for i = 0 to gsize-1 is set somehow
*/

/* set up displs and rcounts vectors first
*/
displs = (int *)malloc(gsize*sizeof(int));
rcounts = (int *)malloc(gsize*sizeof(int));
offset = 0;
for (i=0; i<gsize; ++i) {
    displs[i] = offset;
    offset += stride[i];
    rcounts[i] = 100-i;
}
/* the required buffer size for rbuf is now easily obtained
*/
bufsize = displs[gsize-1]+rcounts[gsize-1];
rbuf = (int *)malloc(bufsize*sizeof(int));
/* Create datatype for the column we are sending
*/
MPI_Type_vector( 100-myrank, 1, 150, MPI_INT, &stype);
MPI_Type_commit( &stype );
sptr = &sendarray[0][myrank];
MPI_Gatherv( sptr, 1, stype, rbuf, rcounts, displs, MPI_INT,
             root, comm);

```

Example 4.10 Process i sends num ints from the i th column of a 100×150 int array, in C. The complicating factor is that the various values of num are not known to root , so a

1 separate gather must first be run to find these out. The data is placed contiguously at the
 2 receiving end.

```

3
4     MPI_Comm comm;
5     int gsize, sendarray[100][150], *sptr;
6     int root, *rbuf, stride, myrank, disp[2], blocklen[2];
7     MPI_Datatype stype, types[2];
8     int *displs, i, *rcounts, num;
9
10    ...
11
12    MPI_Comm_size( comm, &gsize);
13    MPI_Comm_rank( comm, &myrank );
14
15    /* First, gather nums to root
16     */
17    rcounts = (int *)malloc(gsize*sizeof(int));
18    MPI_Gather( &num, 1, MPI_INT, rcounts, 1, MPI_INT, root, comm);
19    /* root now has correct rcounts, using these we set displs[] so
20     * that data is placed contiguously (or concatenated) at receive end
21     */
22    displs = (int *)malloc(gsize*sizeof(int));
23    displs[0] = 0;
24    for (i=1; i<gsize; ++i) {
25        displs[i] = displs[i-1]+rcounts[i-1];
26    }
27    /* And, create receive buffer
28     */
29    rbuf = (int *)malloc(gsize*(displs[gsize-1]+rcounts[gsize-1])
30                                *sizeof(int));
31
32    /* Create datatype for one int, with extent of entire row
33     */
34    disp[0] = 0;      disp[1] = 150*sizeof(int);
35    type[0] = MPI_INT; type[1] = MPI_UB;
36    blocklen[0] = 1;  blocklen[1] = 1;
37    MPI_Type_struct( 2, blocklen, disp, type, &stype );
38    MPI_Type_commit( &stype );
39    sptr = &sendarray[0][myrank];
40    MPI_Gatherv( sptr, num, stype, rbuf, rcounts, displs, MPI_INT,
41                root, comm);
42
43
44
45
46
47
48
```

4.7 Scatter

| | | | |
|--|-----------|---|----|
| MPI_SCATTER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm) | | | 1 |
| | | | 2 |
| | | | 3 |
| | | | 4 |
| | | | 5 |
| IN | sendbuf | address of send buffer (choice, significant only at root) | 6 |
| IN | sendcount | number of elements sent to each process (integer, significant only at root) | 7 |
| | | | 8 |
| | | | 9 |
| IN | sendtype | data type of send buffer elements (significant only at root) (handle) | 10 |
| | | | 11 |
| OUT | recvbuf | address of receive buffer (choice) | 12 |
| IN | recvcount | number of elements in receive buffer (integer) | 13 |
| IN | recvtype | data type of receive buffer elements (handle) | 14 |
| IN | root | rank of sending process (integer) | 15 |
| IN | comm | communicator (handle) | 16 |
| | | | 17 |
| | | | 18 |

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

```
MPI_SCATTER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
            ROOT, COMM, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, ROOT, COMM, IERROR
```

```
void MPI::Comm::Scatter(const void* sendbuf, int sendcount, const
                        MPI::Datatype& sendtype, void* recvbuf, int recvcount,
                        const MPI::Datatype& recvtype, int root) const = 0
```

MPI_SCATTER is the inverse operation to MPI_GATHER.

If *comm* is an intracommunicator, the outcome is *as if* the root executed *n* send operations,

```
MPI_Send(sendbuf + i · sendcount · extent(sendtype), sendcount, sendtype, i, ...),
```

and each process executed a receive,

```
MPI_Recv(recvbuf, recvcount, recvtype, i, ...).
```

An alternative description is that the root sends a message with `MPI_Send(sendbuf, sendcount·n, sendtype, ...)`. This message is split into *n* equal segments, the *i*th segment is sent to the *i*th process in the group, and each process receives this message as above.

The send buffer is ignored for all non-root processes.

The type signature associated with `sendcount`, `sendtype` at the root must be equal to the type signature associated with `recvcount`, `recvtype` at all processes (however, the type maps may be different). This implies that the amount of data sent must be equal to the amount of data received, pairwise between each process and the root. Distinct type maps between sender and receiver are still allowed.

1 All arguments to the function are significant on process `root`, while on other processes,
 2 only arguments `recvbuf`, `recvcount`, `recvtype`, `root`, `comm` are significant. The arguments `root`
 3 and `comm` must have identical values on all processes.

4 The specification of counts and types should not cause any location on the root to be
 5 read more than once.

6
 7 *Rationale.* Though not needed, the last restriction is imposed so as to achieve
 8 symmetry with `MPI_GATHER`, where the corresponding restriction (a multiple-write
 9 restriction) is necessary. (*End of rationale.*)

10
 11 The “in place” option for intracommunicators is specified by passing `MPI_IN_PLACE` as
 12 the value of `recvbuf` at the root. In such case, `recvcount` and `recvtype` are ignored, and root
 13 “sends” no data to itself. The scattered vector is still assumed to contain n segments, where
 14 n is the group size; the $root$ -th segment, which root should “send to itself,” is not moved.

15 If `comm` is an intercommunicator, then the call involves all processes in the intercom-
 16 municator, but with one group (group A) defining the root process. All processes in the
 17 other group (group B) pass the same value in argument `root`, which is the rank of the root
 18 in group A. The root passes the value `MPI_ROOT` in `root`. All other processes in group A
 19 pass the value `MPI_PROC_NULL` in `root`. Data is scattered from the root to all processes in
 20 group B. The receive buffer arguments of the processes in group B must be consistent with
 21 the send buffer argument of the root.

22
 23 `MPI_SCATTERV(sendbuf, sendcounts, displs, sendtype, recvbuf, recvcount, recvtype, root,`
 24 `comm)`

| | | | |
|----|-----|-------------------------|--|
| 26 | IN | <code>sendbuf</code> | address of send buffer (choice, significant only at root) |
| 27 | IN | <code>sendcounts</code> | integer array (of length group size) specifying the num- ber of elements to send to each processor |
| 28 | | | |
| 29 | IN | <code>displs</code> | integer array (of length group size). Entry i specifies the displacement (relative to <code>sendbuf</code> from which to take the outgoing data to process i |
| 30 | | | |
| 31 | | | |
| 32 | | | |
| 33 | IN | <code>sendtype</code> | data type of send buffer elements (handle) |
| 34 | OUT | <code>recvbuf</code> | address of receive buffer (choice) |
| 35 | | | |
| 36 | IN | <code>recvcount</code> | number of elements in receive buffer (integer) |
| 37 | IN | <code>recvtype</code> | data type of receive buffer elements (handle) |
| 38 | IN | <code>root</code> | rank of sending process (integer) |
| 39 | IN | <code>comm</code> | communicator (handle) |
| 40 | | | |

41
 42 `int MPI_Scatterv(void* sendbuf, int *sendcounts, int *displs,`
 43 `MPI_Datatype sendtype, void* recvbuf, int recvcount,`
 44 `MPI_Datatype recvtype, int root, MPI_Comm comm)`

45 `MPI_SCATTERV(SENDBUF, SENDCOUNTS, DISPLS, SENDTYPE, RECVBUF, RECVCOUNT,`
 46 `RECVTYPE, ROOT, COMM, IERROR)`
 47 `<type> SENDBUF(*), RECVBUF(*)`
 48


```

INTEGER SENDCOUNTS(*), DISPLS(*), SENDTYPE, RECVCOUNT, RECVMODE, ROOT,
COMM, IERROR
void MPI::Comm::Scatterv(const void* sendbuf, const int sendcounts[],
                        const int displs[], const MPI::Datatype& sendtype,
                        void* recvbuf, int recvcount, const MPI::Datatype& recvmode,
                        int root) const = 0

```

MPI_SCATTERV is the inverse operation to MPI_GATHERV.

MPI_SCATTERV extends the functionality of MPI_SCATTER by allowing a varying count of data to be sent to each process, since `sendcounts` is now an array. It also allows more flexibility as to where the data is taken from on the root, by providing the new argument, `displs`.

If `comm` is an intracommunicator, the outcome is as if the root executed `n` send operations,

```
MPI_Send(sendbuf + displs[i] * extent(sendtype), sendcounts[i], sendtype, i, ...),
```

and each process executed a receive,

```
MPI_Recv(recvbuf, recvcount, recvmode, i, ...).
```

The send buffer is ignored for all non-root processes.

The type signature implied by `sendcount[i]`, `sendtype` at the root must be equal to the type signature implied by `recvcount`, `recvmode` at process `i` (however, the type maps may be different). This implies that the amount of data sent must be equal to the amount of data received, pairwise between each process and the root. Distinct type maps between sender and receiver are still allowed.

All arguments to the function are significant on process `root`, while on other processes, only arguments `recvbuf`, `recvcount`, `recvmode`, `root`, `comm` are significant. The arguments `root` and `comm` must have identical values on all processes.

The specification of counts, types, and displacements should not cause any location on the root to be read more than once.

The “in place” option for intracommunicators is specified by passing `MPI.IN_PLACE` as the value of `recvbuf` at the root. In such case, `recvcount` and `recvmode` are ignored, and root “sends” no data to itself. The scattered vector is still assumed to contain n segments, where n is the group size; the $root$ -th segment, which root should “send to itself,” is not moved.

If `comm` is an intercommunicator, then the call involves all processes in the intercommunicator, but with one group (group A) defining the root process. All processes in the other group (group B) pass the same value in argument `root`, which is the rank of the root in group A. The root passes the value `MPI_ROOT` in `root`. All other processes in group A pass the value `MPI_PROC_NULL` in `root`. Data is scattered from the root to all processes in group B. The receive buffer arguments of the processes in group B must be consistent with the send buffer argument of the root.

4.7.1 Examples using MPI_SCATTER, MPI_SCATTERV

The examples in this section are using intracommunicators.

Example 4.11 The reverse of Example 4.2. Scatter sets of 100 ints from the root to each process in the group. See figure 4.9.

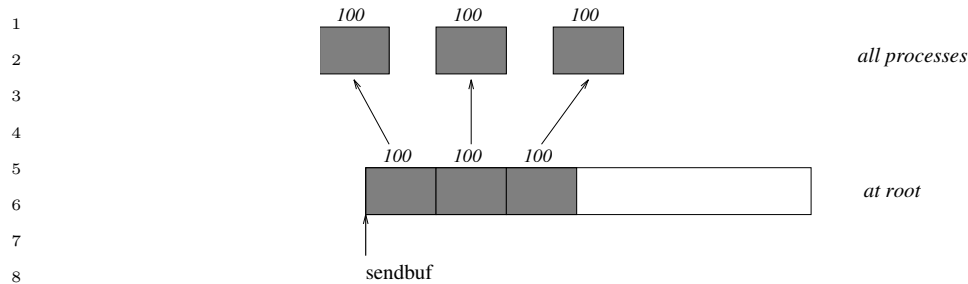


Figure 4.9: The root process scatters sets of 100 ints to each process in the group.

```

12 MPI_Comm comm;
13 int gsize,*sendbuf;
14 int root, rbuf[100];
15 ...
16 MPI_Comm_size( comm, &gsize);
17 sendbuf = (int *)malloc(gsize*100*sizeof(int));
18 ...
19 MPI_Scatter( sendbuf, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);

```

Example 4.12 The reverse of Example 4.5. The root process scatters sets of 100 ints to the other processes, but the sets of 100 are *stride* ints apart in the sending buffer. Requires use of MPI_SCATTERV. Assume *stride* \geq 100. See figure 4.10.

```

25 MPI_Comm comm;
26 int gsize,*sendbuf;
27 int root, rbuf[100], i, *displs, *scounts;
28 ...
29 ...
30 ...
31 MPI_Comm_size( comm, &gsize);
32 sendbuf = (int *)malloc(gsize*stride*sizeof(int));
33 ...
34 displs = (int *)malloc(gsize*sizeof(int));
35 scounts = (int *)malloc(gsize*sizeof(int));
36 for (i=0; i<gsize; ++i) {
37     displs[i] = i*stride;
38     scounts[i] = 100;
39 }
40 MPI_Scatterv( sendbuf, scounts, displs, MPI_INT, rbuf, 100, MPI_INT,
41             root, comm);

```

Example 4.13 The reverse of Example 4.9. We have a varying stride between blocks at sending (root) side, at the receiving side we receive into the *i*th column of a 100×150 C array. See figure 4.11.

```

48 MPI_Comm comm;

```

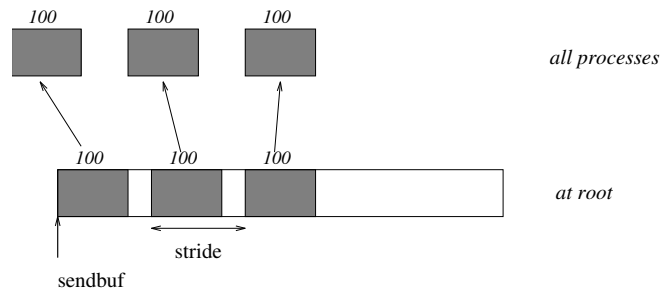


Figure 4.10: The root process scatters sets of 100 ints, moving by `stride` ints from send to send in the scatter.

```

int gsize, recvarray[100][150], *rptr;
int root, *sendbuf, myrank, bufsize, *stride;
MPI_Datatype rtype;
int i, *displs, *counts, offset;
...
MPI_Comm_size( comm, &gsize);
MPI_Comm_rank( comm, &myrank );

stride = (int *)malloc(gsize*sizeof(int));
...
/* stride[i] for i = 0 to gsize-1 is set somehow
 * sendbuf comes from elsewhere
 */
...
displs = (int *)malloc(gsize*sizeof(int));
counts = (int *)malloc(gsize*sizeof(int));
offset = 0;
for (i=0; i<gsize; ++i) {
    displs[i] = offset;
    offset += stride[i];
    counts[i] = 100 - i;
}
/* Create datatype for the column we are receiving
 */
MPI_Type_vector( 100-myrank, 1, 150, MPI_INT, &rtype);
MPI_Type_commit( &rtype );
rptr = &recvarray[0][myrank];
MPI_Scatterv( sendbuf, counts, displs, MPI_INT, rptr, 1, rtype,
              root, comm);

```

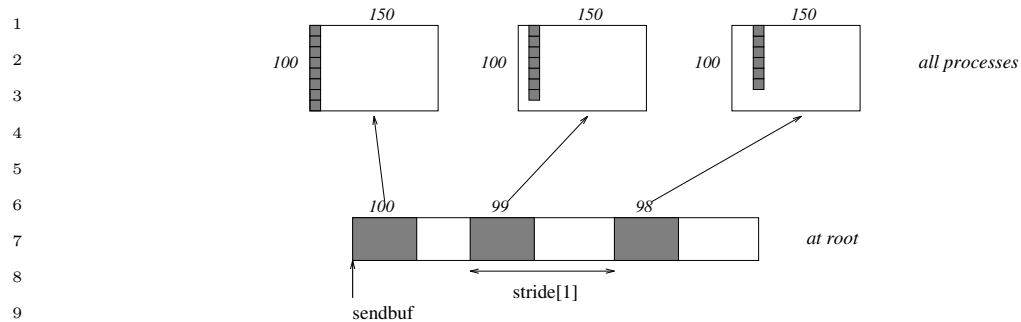


Figure 4.11: The root scatters blocks of $100-i$ ints into column i of a 100×150 C array. At the sending side, the blocks are `stride[i]` ints apart.

4.8 Gather-to-all

```

17 MPI_ALLGATHER( sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, comm)
18
19     IN      sendbuf      starting address of send buffer (choice)
20     IN      sendcount    number of elements in send buffer (integer)
21     IN      sendtype     data type of send buffer elements (handle)
22     OUT     recvbuf      address of receive buffer (choice)
23     IN      recvcount    number of elements received from any process (inte-
24                                     ger)
25     IN      recvtype     data type of receive buffer elements (handle)
26     IN      comm         communicator (handle)
27
28
29
30 int MPI_Allgather(void* sendbuf, int sendcount, MPI_Datatype sendtype,
31                 void* recvbuf, int recvcount, MPI_Datatype recvtype,
32                 MPI_Comm comm)
33
34 MPI_ALLGATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, REVCOUNT, RECVTYPE,
35               COMM, IERROR)
36     <type> SENDBUF(*), RECVBUF(*)
37     INTEGER SENDCOUNT, SENDTYPE, REVCOUNT, RECVTYPE, COMM, IERROR
38
39 void MPI::Comm::Allgather(const void* sendbuf, int sendcount, const
40                          MPI::Datatype& sendtype, void* recvbuf, int recvcount,
41                          const MPI::Datatype& recvtype) const = 0

```

`MPI_ALLGATHER` can be thought of as `MPI_GATHER`, but where all processes receive the result, instead of just the root. The block of data sent from the j th process is received by every process and placed in the j th block of the buffer `recvbuf`.

The type signature associated with `sendcount`, `sendtype`, at a process must be equal to the type signature associated with `recvcount`, `recvtype` at any other process.

If `comm` is an [intracommunicator](#), the outcome of a call to `MPI_ALLGATHER(...)` is as if all processes executed `n` calls to

```
MPI_GATHER(sendbuf, sendcount, sendtype, recvbuf, recvcount,
           recvtype, root, comm),
```

for `root = 0, ..., n-1`. The rules for correct usage of `MPI_ALLGATHER` are easily found from the corresponding rules for `MPI_GATHER`.

The “in place” option for intracommunicators is specified by passing the value `MPI_IN_PLACE` to the argument `sendbuf` at all processes. `sendcount` and `sendtype` are ignored. Then the input data of each process is assumed to be in the area where that process would receive its own contribution to the receive `buffer`.

If `comm` is an intercommunicator, then each process in group A contributes a data item; these items are concatenated and the result is stored at each process in group B. Conversely the concatenation of the contributions of the processes in group B is stored at each process in group A. The send buffer arguments in group A must be consistent with the receive buffer arguments in group B, and vice versa.

Advice to users. The communication pattern of `MPI_ALLGATHER` executed on an intercommunication domain need not be symmetric. The number of items sent by processes in group A (as specified by the arguments `sendcount`, `sendtype` in group A and the arguments `recvcount`, `recvtype` in group B), need not equal the number of items sent by processes in group B (as specified by the arguments `sendcount`, `sendtype` in group B and the arguments `recvcount`, `recvtype` in group A). In particular, one can move data in only one direction by specifying `sendcount = 0` for the communication in the reverse direction.

(End of advice to users.)

```
MPI_ALLGATHERV(sendbuf, sendcount, sendtype, recvbuf, recvcnts, displs, recvtype, comm)
```

| | | |
|-----|------------------------|---|
| IN | <code>sendbuf</code> | starting address of send buffer (choice) |
| IN | <code>sendcount</code> | number of elements in send buffer (integer) |
| IN | <code>sendtype</code> | data type of send buffer elements (handle) |
| OUT | <code>recvbuf</code> | address of receive buffer (choice) |
| IN | <code>recvcnts</code> | integer array (of length group size) containing the number of elements that are received from each process |
| IN | <code>displs</code> | integer array (of length group size). Entry <code>i</code> specifies the displacement (relative to <code>recvbuf</code>) at which to place the incoming data from process <code>i</code> |
| IN | <code>recvtype</code> | data type of receive buffer elements (handle) |
| IN | <code>comm</code> | communicator (handle) |

```
int MPI_Allgatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                  void* recvbuf, int *recvcnts, int *displs,
                  MPI_Datatype recvtype, MPI_Comm comm)
```

```
MPI_ALLGATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
               RECVTYPE, COMM, IERROR)
```

```

1     <type> SENDBUF(*), RECVBUF(*)
2     INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVMODE, COMM,
3     IERROR
4
5     void MPI::Comm::Allgatherv(const void* sendbuf, int sendcount, const
6         MPI::Datatype& sendtype, void* recvbuf,
7         const int recvcounts[], const int displs[],
8         const MPI::Datatype& recvmode) const = 0

```

MPI_ALLGATHERV can be thought of as MPI_GATHERV, but where all processes receive the result, instead of just the root. The block of data sent from the j th process is received by every process and placed in the j th block of the buffer `recvbuf`. These blocks need not all be the same size.

The type signature associated with `sendcount`, `sendtype`, at process j must be equal to the type signature associated with `recvcounts[j]`, `recvmode` at any other process.

If `comm` is an intracommunicator, the outcome is as if all processes executed calls to

```

17     MPI_GATHERV(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs,
18                 recvmode, root, comm),

```

for `root = 0, ..., n-1`. The rules for correct usage of MPI_ALLGATHERV are easily found from the corresponding rules for MPI_GATHERV.

The “in place” option for intracommunicators is specified by passing the value `MPI_IN_PLACE` to the argument `sendbuf` at all processes. `sendcount` and `sendtype` are ignored. Then the input data of each process is assumed to be in the area where that process would receive its own contribution to the receive buffer.

If `comm` is an intercommunicator, then each process in group A contributes a data item; these items are concatenated and the result is stored at each process in group B. Conversely the concatenation of the contributions of the processes in group B is stored at each process in group A. The send buffer arguments in group A must be consistent with the receive buffer arguments in group B, and vice versa.

4.8.1 Examples using MPI_ALLGATHER, MPI_ALLGATHERV

The examples in this section are using intracommunicators.

Example 4.14 The all-gather version of Example 4.2. Using MPI_ALLGATHER, we will gather 100 ints from every process in the group to every process.

```

38     MPI_Comm comm;
39     int gsize, sendarray[100];
40     int *rbuf;
41     ...
42     MPI_Comm_size( comm, &gsize);
43     rbuf = (int *)malloc(gsize*100*sizeof(int));
44     MPI_Allgather( sendarray, 100, MPI_INT, rbuf, 100, MPI_INT, comm);

```

After the call, every process has the group-wide concatenation of the sets of data.

4.9 All-to-All Scatter/Gather

`MPI_ALLTOALL(sendbuf, sendcount, sendtype, recvbuf, recvcnt, recvtype, comm)`

| | | |
|-----|-----------|--|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements sent to each process (integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcnt | number of elements received from any process (integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | comm | communicator (handle) |

```
int MPI_Alltoall(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                void* recvbuf, int recvcnt, MPI_Datatype recvtype,
                MPI_Comm comm)
```

```
MPI_ALLTOALL(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
             COMM, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, COMM, IERROR
```

```
void MPI::Comm::Alltoall(const void* sendbuf, int sendcount, const
                        MPI::Datatype& sendtype, void* recvbuf, int recvcnt,
                        const MPI::Datatype& recvtype) const = 0
```

`MPI_ALLTOALL` is an extension of `MPI_ALLGATHER` to the case where each process sends distinct data to each of the receivers. The j th block sent from process i is received by process j and is placed in the i th block of `recvbuf`.

The type signature associated with `sendcount`, `sendtype`, at a process must be equal to the type signature associated with `recvcnt`, `recvtype` at any other process. This implies that the amount of data sent must be equal to the amount of data received, pairwise between every pair of processes. As usual, however, the type maps may be different.

If `comm` is an [intracommunicator](#), the outcome is as if each process executed a send to each process (itself included) with a call to,

```
MPI_Send(sendbuf + i · sendcount · extent(sendtype), sendcount, sendtype, i, ...),
```

and a receive from every other process with a call to,

```
MPI_Recv(recvbuf + i · recvcnt · extent(recvtype), recvcnt, i, ...).
```

All arguments on all processes are significant. The argument `comm` must have identical values on all processes.

No “in place” option is supported.

If `comm` is an [intercommunicator](#), then the outcome is as if each process in group A sends a message to each process in group B, and vice versa. The j -th send buffer of process i in group A should be consistent with the i -th receive buffer of process j in group B, and vice versa.

Advice to users. When all-to-all is executed on an intercommunication domain, then the number of data items sent from processes in group A to processes in group B need not equal the number of items sent in the reverse direction. In particular, one can have unidirectional communication by specifying `sendcount = 0` in the reverse direction.

(End of advice to users.)

`MPI_ALLTOALLV(sendbuf, sendcounts, sdispls, sendtype, recvbuf, recvcoun-
ts, rdispls, recvtype, comm)`

| | | |
|-----|-----------------|--|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcounts | integer array equal to the group size specifying the number of elements to send to each processor |
| IN | sdispls | integer array (of length group size). Entry <i>j</i> specifies the displacement (relative to <code>sendbuf</code> from which to take the outgoing data destined for process <i>j</i>) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcoun- ts | integer array equal to the group size specifying the number of elements that can be received from each processor |
| IN | rdispls | integer array (of length group size). Entry <i>i</i> specifies the displacement (relative to <code>recvbuf</code> at which to place the incoming data from process <i>i</i>) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | comm | communicator (handle) |

```
int MPI_Alltoallv(void* sendbuf, int *sendcounts, int *sdispls,
                 MPI_Datatype sendtype, void* recvbuf, int *recvcoun-
                 ts, int *rdispls, MPI_Datatype recvtype, MPI_Comm comm)
MPI_ALLTOALLV(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPE, RECVBUF, RECVCOUNTS,
              RDISPLS, RECVTYPE, COMM, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPE, RECVCOUNTS(*), RDISPLS(*),
RECVTYPE, COMM, IERROR
```

```
void MPI::Comm::Alltoallv(const void* sendbuf, const int sendcounts[],
                          const int sdispls[], const MPI::Datatype& sendtype,
                          void* recvbuf, const int recvcoun-
                          ts[], const int rdispls[],
                          const MPI::Datatype& recvtype) const = 0
```

`MPI_ALLTOALLV` adds flexibility to `MPI_ALLTOALL` in that the location of data for the send is specified by `sdispls` and the location of the placement of the data on the receive side is specified by `rdispls`.

If `comm` is an intracommunicator, then the j th block sent from process i is received by process j and is placed in the i th block of `recvbuf`. These blocks need not all have the same size.

The type signature associated with `sendcount[j]`, `sendtype` at process i must be equal to the type signature associated with `recvcount[i]`, `recvtype` at process j . This implies that the amount of data sent must be equal to the amount of data received, pairwise between every pair of processes. Distinct type maps between sender and receiver are still allowed.

The outcome is as if each process sent a message to every other process with,

```
MPI_Send(sendbuf + displs[i] · extent(sendtype), sendcounts[i], sendtype, i, ...),
```

and received a message from every other process with a call to

```
MPI_Recv(recvbuf + displs[i] · extent(recvtype), recvcounts[i], recvtype, i, ...).
```

All arguments on all processes are significant. The argument `comm` must have identical values on all processes.

No “in place” option is supported.

If `comm` is an intercommunicator, then the outcome is as if each process in group A sends a message to each process in group B, and vice versa. The j -th send buffer of process i in group A should be consistent with the i -th receive buffer of process j in group B, and vice versa.

Rationale. The definitions of `MPI_ALLTOALL` and `MPI_ALLTOALLV` give as much flexibility as one would achieve by specifying n independent, point-to-point communications, with two exceptions: all messages use the same datatype, and messages are scattered from (or gathered to) sequential storage. (*End of rationale.*)

Advice to implementors. Although the discussion of collective communication in terms of point-to-point operation implies that each message is transferred directly from sender to receiver, implementations may use a tree communication pattern. Messages can be forwarded by intermediate nodes where they are split (for scatter) or concatenated (for gather), if this is more efficient. (*End of advice to implementors.*)

4.9.1 Generalized All-to-all Function

One of the basic data movement operations needed in parallel signal processing is the 2-D matrix transpose. This operation has motivated a generalization of the `MPI_ALLTOALLV` function. This new collective operation is `MPI_ALLTOALLW`; the “W” indicates that it is an extension to `MPI_ALLTOALLV`.

The following function is the most general form of All-to-all. Like `MPI_TYPE_CREATE_STRUCT`, the most general type constructor, `MPI_ALLTOALLW` allows separate specification of count, displacement and datatype. In addition, to allow maximum flexibility, the displacement of blocks within the send and receive buffers is specified in bytes.

Rationale. The `MPI_ALLTOALLW` function generalizes several MPI functions by carefully selecting the input arguments. For example, by making all but one process have `sendcounts[i] = 0`, this achieves an `MPI_SCATTERW` function. (*End of rationale.*)

```
1 MPI_ALLTOALLW(sendbuf, sendcounts, sdispls, sendtypes, recvbuf, recvcoun-
2 types, comm)
```

| | | | |
|----|-----|------------|--|
| 3 | IN | sendbuf | starting address of send buffer (choice) |
| 4 | | | |
| 5 | IN | sendcounts | integer array equal to the group size specifying the |
| 6 | | | number of elements to send to each processor (array |
| 7 | | | of integers) |
| 8 | IN | sdispls | integer array (of length group size). Entry j specifies |
| 9 | | | the displacement in bytes (relative to sendbuf) from |
| 10 | | | which to take the outgoing data destined for process j |
| 11 | | | (array of integers) |
| 12 | IN | sendtypes | array of datatypes (of length group size). Entry j spec- |
| 13 | | | ifies the type of data to send to process j (array of |
| 14 | | | handles) |
| 15 | | | |
| 16 | OUT | recvbuf | address of receive buffer (choice) |
| 17 | IN | recvcoun- | integer array equal to the group size specifying the |
| 18 | | ts | number of elements that can be received from each |
| 19 | | | processor (array of integers) |
| 20 | IN | rdispls | integer array (of length group size). Entry i specifies |
| 21 | | | the displacement in bytes (relative to recvbuf) at which |
| 22 | | | to place the incoming data from process i (array of |
| 23 | | | integers) |
| 24 | | | |
| 25 | IN | recvtypes | array of datatypes (of length group size). Entry i spec- |
| 26 | | | ifies the type of data received from process i (array of |
| 27 | | | handles) |
| 28 | IN | comm | communicator (handle) |
| 29 | | | |

```
30 int MPI_Alltoallw(void *sendbuf, int sendcounts[], int sdispls[],
31 MPI_Datatype sendtypes[], void *recvbuf, int recvcoun-
32 ts, int rdispls[], MPI_Datatype recvtypes[], MPI_Comm comm)
```

```
33
34 MPI_ALLTOALLW(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPES, RECVBUF, RECVCOUNTS,
35 RDISPLS, RECVTYPES, COMM, IERROR)
36 <type> SENDBUF(*), RECVBUF(*)
37 INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPES(*), RECVCOUNTS(*),
38 RDISPLS(*), RECVTYPES(*), COMM, IERROR
```

```
39 void MPI::Comm::Alltoallw(const void* sendbuf, const int sendcounts[],
40 const int sdispls[], const MPI::Datatype sendtypes[], void*
41 recvbuf, const int recvcoun-
42 ts, const int rdispls[], const
43 MPI::Datatype recvtypes[]) const = 0
```

44 No “in place” option is supported.

45 If *comm* is an intracommunicator, then the *j*-th block sent from process *i* is received
 46 by process *j* and is placed in the *i*-th block of *recvbuf*. These blocks need not all have the
 47 same size.

48

The type signature associated with `sendcounts[j]`, `sendtypes[j]` at process i must be equal to the type signature associated with `recvcounts[i]`, `recvtypes[i]` at process j . This implies that the amount of data sent must be equal to the amount of data received, pairwise between every pair of processes. Distinct type maps between sender and receiver are still allowed.

The outcome is as if each process sent a message to every other process with

```
MPI_Send(sendbuf + sdispls[i], sendcounts[i], sendtypes[i], i, ...),
```

and received a message from every other process with a call to

```
MPI_Recv(recvbuf + rdispls[i], recvcounts[i], recvtypes[i], i, ...).
```

All arguments on all processes are significant. The argument `comm` must describe the same communicator on all processes.

If `comm` is an intercommunicator, then the outcome is as if each process in group A sends a message to each process in group B, and vice versa. The j -th send buffer of process i in group A should be consistent with the i -th receive buffer of process j in group B, and vice versa.

4.10 Global Reduction Operations

The functions in this section perform a global reduce operation (such as sum, max, logical AND, etc.) across all the members of a group. The reduction operation can be either one of a predefined list of operations, or a user-defined operation. The global reduction functions come in several flavors: a reduce that returns the result of the reduction at one node, an all-reduce that returns this result at all nodes, and a scan (parallel prefix) operation. In addition, a reduce-scatter operation combines the functionality of a reduce and of a scatter operation.

4.10.1 Reduce

```
MPI_REDUCE( sendbuf, recvbuf, count, datatype, op, root, comm)
```

| | | |
|-----|-----------------------|--|
| IN | <code>sendbuf</code> | address of send buffer (choice) |
| OUT | <code>recvbuf</code> | address of receive buffer (choice, significant only at root) |
| IN | <code>count</code> | number of elements in send buffer (integer) |
| IN | <code>datatype</code> | data type of elements of send buffer (handle) |
| IN | <code>op</code> | reduce operation (handle) |
| IN | <code>root</code> | rank of root process (integer) |
| IN | <code>comm</code> | communicator (handle) |

```
int MPI_Reduce(void* sendbuf, void* recvbuf, int count,
               MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)
MPI_REDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, ROOT, COMM, IERROR)
```

```

1     <type> SENDBUF(*), RECVBUF(*)
2     INTEGER COUNT, DATATYPE, OP, ROOT, COMM, IERROR
3
4     void MPI::Comm::Reduce(const void* sendbuf, void* recvbuf, int count,
5                           const MPI::Datatype& datatype, const MPI::Op& op, int root)
6                           const = 0

```

MPI_REDUCE combines the elements provided in the input buffer of each process in the group, using the operation `op`, and returns the combined value in the output buffer of the process with rank `root`. The input buffer is defined by the arguments `sendbuf`, `count` and `datatype`; the output buffer is defined by the arguments `recvbuf`, `count` and `datatype`; both have the same number of elements, with the same type. The routine is called by all group members using the same arguments for `count`, `datatype`, `op`, `root` and `comm`. Thus, all processes provide input buffers and output buffers of the same length, with elements of the same type. Each process can provide one element, or a sequence of elements, in which case the combine operation is executed element-wise on each entry of the sequence. For example, if the operation is `MPI_MAX` and the send buffer contains two elements that are floating point numbers (`count = 2` and `datatype = MPI_FLOAT`), then `recvbuf(1) = global max(sendbuf(1))` and `recvbuf(2) = global max(sendbuf(2))`.

Sec. 4.10.2, lists the set of predefined operations provided by MPI. That section also enumerates the datatypes each operation can be applied to. In addition, users may define their own operations that can be overloaded to operate on several datatypes, either basic or derived. This is further explained in Sec. 4.10.5.

The operation `op` is always assumed to be associative. All predefined operations are also assumed to be commutative. Users may define operations that are assumed to be associative, but not commutative. The “canonical” evaluation order of a reduction is determined by the ranks of the processes in the group. However, the implementation can take advantage of associativity, or associativity and commutativity in order to change the order of evaluation. This may change the result of the reduction for operations that are not strictly associative and commutative, such as floating point addition.

Advice to implementors. It is strongly recommended that MPI_REDUCE be implemented so that the same result be obtained whenever the function is applied on the same arguments, appearing in the same order. Note that this may prevent optimizations that take advantage of the physical location of processors. (*End of advice to implementors.*)

The `datatype` argument of MPI_REDUCE must be compatible with `op`. Predefined operators work only with the MPI types listed in Sec. 4.10.2 and Sec. 4.10.4. Furthermore, the `datatype` and `op` given for predefined operators must be the same on all processes.

Note that it is possible for users to supply different user-defined operations to MPI_REDUCE in each process. MPI does not define which operations are used on which operands in this case. **User-defined operators may operate on general, derived datatypes.** In this case, each argument that the reduce operation is applied to is one element described by such a datatype, which may contain several basic values. This is further explained in Section 4.10.5.

Advice to users. Users should make no assumptions about how MPI_REDUCE is implemented. Safest is to ensure that the same function is passed to MPI_REDUCE

by each process. (*End of advice to users.*)

Overlapping datatypes are permitted in “send” buffers. Overlapping datatypes in “receive” buffers are erroneous and may give unpredictable results.

The “in place” option for intracommunicators is specified by passing the value `MPI_IN_PLACE` to the argument `sendbuf` at the root. In such case, the input data is taken at the root from the receive buffer, where it will be replaced by the output data.

If `comm` is an intercommunicator, then the call involves all processes in the intercommunicator, but with one group (group A) defining the root process. All processes in the other group (group B) pass the same value in argument `root`, which is the rank of the root in group A. The root passes the value `MPI_ROOT` in `root`. All other processes in group A pass the value `MPI_PROC_NULL` in `root`. Only send buffer arguments are significant in group B and only receive buffer arguments are significant at the root.

4.10.2 Predefined reduce operations

The following predefined operations are supplied for `MPI_REDUCE` and related functions `MPI_ALLREDUCE`, `MPI_REDUCE_SCATTER`, and `MPI_SCAN`. These operations are invoked by placing the following in `op`.

| Name | Meaning |
|-------------------------|------------------------|
| <code>MPI_MAX</code> | maximum |
| <code>MPI_MIN</code> | minimum |
| <code>MPI_SUM</code> | sum |
| <code>MPI_PROD</code> | product |
| <code>MPI_LAND</code> | logical and |
| <code>MPI_BAND</code> | bit-wise and |
| <code>MPI_LOR</code> | logical or |
| <code>MPI BOR</code> | bit-wise or |
| <code>MPI_LXOR</code> | logical xor |
| <code>MPI_BXOR</code> | bit-wise xor |
| <code>MPI_MAXLOC</code> | max value and location |
| <code>MPI_MINLOC</code> | min value and location |

The two operations `MPI_MINLOC` and `MPI_MAXLOC` are discussed separately in Sec. 4.10.4. For the other predefined operations, we enumerate below the allowed combinations of `op` and datatype arguments. First, define groups of MPI basic datatypes in the following way.

| | |
|------------------|--|
| C integer: | <code>MPI_INT</code> , <code>MPI_LONG</code> , <code>MPI_SHORT</code> , <code>MPI_UNSIGNED_SHORT</code> , <code>MPI_UNSIGNED</code> , <code>MPI_UNSIGNED_LONG</code> |
| Fortran integer: | <code>MPI_INTEGER</code> |
| Floating point: | <code>MPI_FLOAT</code> , <code>MPI_DOUBLE</code> , <code>MPI_REAL</code> , <code>MPI_DOUBLE_PRECISION</code> , <code>MPI_LONG_DOUBLE</code> |
| Logical: | <code>MPI_LOGICAL</code> |
| Complex: | <code>MPI_COMPLEX</code> |

1 Byte: MPI_BYTE

2 Now, the valid datatypes for each option is specified below.

| 3 | 4 | 5 | 6 |
|----|-----------------------------|---|---|
| Op | | Allowed Types | |
| 7 | MPI_MAX, MPI_MIN | C integer, Fortran integer, Floating point | |
| 8 | MPI_SUM, MPI_PROD | C integer, Fortran integer, Floating point, Complex | |
| 9 | MPI_LAND, MPI_LOR, MPI_LXOR | C integer, Logical | |
| 10 | MPI_BAND, MPI_BOR, MPI_BXOR | C integer, Fortran integer, Byte | |

11 The following examples are using intracommunicators.

12

13 **Example 4.15** A routine that computes the dot product of two vectors that are distributed
14 across a group of processes and returns the answer at node zero.

```
15
16 SUBROUTINE PAR_BLAS1(m, a, b, c, comm)
17 REAL a(m), b(m)           ! local slice of array
18 REAL c                   ! result (at node zero)
19 REAL sum
20 INTEGER m, comm, i, ierr
21
22 ! local sum
23 sum = 0.0
24 DO i = 1, m
25     sum = sum + a(i)*b(i)
26 END DO
27
28 ! global sum
29 CALL MPI_REDUCE(sum, c, 1, MPI_REAL, MPI_SUM, 0, comm, ierr)
30 RETURN
31
```

32 **Example 4.16** A routine that computes the product of a vector and an array that are
33 distributed across a group of processes and returns the answer at node zero.

```
34
35 SUBROUTINE PAR_BLAS2(m, n, a, b, c, comm)
36 REAL a(m), b(m,n)       ! local slice of array
37 REAL c(n)               ! result
38 REAL sum(n)
39 INTEGER n, comm, i, j, ierr
40
41 ! local sum
42 DO j= 1, n
43     sum(j) = 0.0
44     DO i = 1, m
45         sum(j) = sum(j) + a(i)*b(i,j)
46     END DO
47 END DO
48
```

```

! global sum
CALL MPI_REDUCE(sum, c, n, MPI_REAL, MPI_SUM, 0, comm, ierr)

! return result at node zero (and garbage at the other nodes)
RETURN

```

4.10.3 Signed Characters and Reductions

MPI-1 doesn't allow reductions on signed or unsigned `chars`. Since this restriction (formally) prevents a C programmer from performing reduction operations on such types (which could be useful, particularly in an image processing application where pixel values are often represented as "unsigned char"), we now specify a way for such reductions to be carried out.

MPI-1.2 already has the C types `MPI_CHAR` and `MPI_UNSIGNED_CHAR`. However there is a problem here in that `MPI_CHAR` is intended to represent a character, not a small integer, and therefore will be translated between machines with different character representations.

To overcome this, a new MPI predefined datatype, `MPI_SIGNED_CHAR`, is added to the predefined datatypes of MPI-2, which corresponds to the [ISO C](#) and [ISO C++](#) datatype `signed char`.

Advice to users.

The types `MPI_CHAR` and `MPI_CHARACTER` are intended for characters, and so will be translated to preserve the printable representation, rather than the bit value, if sent between machines with different character codes. The types `MPI_SIGNED_CHAR` and `MPI_UNSIGNED_CHAR` should be used in C if the integer value should be preserved.

(End of advice to users.)

The types `MPI_SIGNED_CHAR` and `MPI_UNSIGNED_CHAR` can be used in reduction operations. `MPI_CHAR` (which represents printable characters) cannot be used in reduction operations. This is an extension to MPI-1.2, since MPI-1.2 does not allow the use of `MPI_UNSIGNED_CHAR` in reduction operations (and does not have the `MPI_SIGNED_CHAR` type).

In a heterogeneous environment, `MPI_CHAR` and `MPI_WCHAR` will be translated so as to preserve the printable character, whereas `MPI_SIGNED_CHAR` and `MPI_UNSIGNED_CHAR` will be translated so as to preserve the integer value.

4.10.4 MINLOC and MAXLOC

The operator `MPI_MINLOC` is used to compute a global minimum and also an index attached to the minimum value. `MPI_MAXLOC` similarly computes a global maximum and index. One application of these is to compute a global minimum (maximum) and the rank of the process containing this value.

The operation that defines `MPI_MAXLOC` is:

$$\begin{pmatrix} u \\ i \end{pmatrix} \circ \begin{pmatrix} v \\ j \end{pmatrix} = \begin{pmatrix} w \\ k \end{pmatrix}$$

where

$$w = \max(u, v)$$

1 and

$$2 \quad k = \begin{cases} i & \text{if } u > v \\ \min(i, j) & \text{if } u = v \\ j & \text{if } u < v \end{cases}$$

6 MPI_MINLOC is defined similarly:

$$8 \quad \begin{pmatrix} u \\ i \end{pmatrix} \circ \begin{pmatrix} v \\ j \end{pmatrix} = \begin{pmatrix} w \\ k \end{pmatrix}$$

11 where

$$12 \quad w = \min(u, v)$$

14 and

$$16 \quad k = \begin{cases} i & \text{if } u < v \\ \min(i, j) & \text{if } u = v \\ j & \text{if } u > v \end{cases}$$

19 Both operations are associative and commutative. Note that if MPI_MAXLOC is applied
 20 to reduce a sequence of pairs $(u_0, 0), (u_1, 1), \dots, (u_{n-1}, n-1)$, then the value returned is
 21 (u, r) , where $u = \max_i u_i$ and r is the index of the first global maximum in the sequence.
 22 Thus, if each process supplies a value and its rank within the group, then a reduce operation
 23 with `op = MPI_MAXLOC` will return the maximum value and the rank of the first process
 24 with that value. Similarly, MPI_MINLOC can be used to return a minimum and its index.
 25 More generally, MPI_MINLOC computes a *lexicographic minimum*, where elements are ordered
 26 according to the first component of each pair, and ties are resolved according to the second
 27 component.

28 The reduce operation is defined to operate on arguments that consist of a pair: value
 29 and index. For both Fortran and C, types are provided to describe the pair. The potentially
 30 mixed-type nature of such arguments is a problem in Fortran. The problem is circumvented,
 31 for Fortran, by having the MPI-provided type consist of a pair of the same type as value,
 32 and coercing the index to this type also. In C, the MPI-provided pair type has distinct
 33 types and the index is an `int`.

34 In order to use MPI_MINLOC and MPI_MAXLOC in a reduce operation, one must provide
 35 a `datatype` argument that represents a pair (value and index). MPI provides nine such
 36 predefined datatypes. The operations MPI_MAXLOC and MPI_MINLOC can be used with each
 37 of the following datatypes.

39 Fortran:

| 40 Name | Description |
|--------------------------|------------------------------------|
| 41 MPI_2REAL | pair of REALs |
| 42 MPI_2DOUBLE_PRECISION | pair of DOUBLE PRECISION variables |
| 43 MPI_2INTEGER | pair of INTEGERS |

46 C:

| 47 Name | Description |
|------------------|---------------|
| 48 MPI_FLOAT_INT | float and int |

| | | |
|---------------------|---------------------|---|
| MPI_DOUBLE_INT | double and int | 1 |
| MPI_LONG_INT | long and int | 2 |
| MPI_2INT | pair of int | 3 |
| MPI_SHORT_INT | short and int | 4 |
| MPI_LONG_DOUBLE_INT | long double and int | 5 |

The datatype MPI_2REAL is *as if* defined by the following (see Section 3.12).

```
MPI_TYPE_CONTIGUOUS(2, MPI_REAL, MPI_2REAL)
```

Similar statements apply for MPI_2INTEGER, MPI_2DOUBLE_PRECISION, and MPI_2INT.

The datatype MPI_FLOAT_INT is *as if* defined by the following sequence of instructions.

```
type[0] = MPI_FLOAT
type[1] = MPI_INT
disp[0] = 0
disp[1] = sizeof(float)
block[0] = 1
block[1] = 1
MPI_TYPE_STRUCT(2, block, disp, type, MPI_FLOAT_INT)
```

Similar statements apply for MPI_LONG_INT and MPI_DOUBLE_INT.

The following examples are using intracommunicators.

Example 4.17 Each process has an array of 30 doubles, in C. For each of the 30 locations, compute the value and rank of the process containing the largest value.

```
...
/* each process has an array of 30 double: ain[30]
*/
double ain[30], aout[30];
int ind[30];
struct {
    double val;
    int rank;
} in[30], out[30];
int i, myrank, root;

MPI_Comm_rank(comm, &myrank);
for (i=0; i<30; ++i) {
    in[i].val = ain[i];
    in[i].rank = myrank;
}
MPI_Reduce( in, out, 30, MPI_DOUBLE_INT, MPI_MAXLOC, root, comm );
/* At this point, the answer resides on process root
*/
if (myrank == root) {
    /* read ranks out
    */
    for (i=0; i<30; ++i) {
```

```

1         aout[i] = out[i].val;
2         ind[i] = out[i].rank;
3     }
4 }
5
6

```

Example 4.18 Same example, in Fortran.

```

8     ...
9     ! each process has an array of 30 double: ain(30)
10
11    DOUBLE PRECISION ain(30), aout(30)
12    INTEGER ind(30);
13    DOUBLE PRECISION in(2,30), out(2,30)
14    INTEGER i, myrank, root, ierr;
15
16    MPI_COMM_RANK(comm, myrank, ierr);
17    DO I=1, 30
18        in(1,i) = ain(i)
19        in(2,i) = myrank    ! myrank is coerced to a double
20    END DO
21
22    MPI_REDUCE( in, out, 30, MPI_2DOUBLE_PRECISION, MPI_MAXLOC, root,
23                                     comm, ierr );
24
25    ! At this point, the answer resides on process root
26
27    IF (myrank .EQ. root) THEN
28        ! read ranks out
29        DO I= 1, 30
30            aout(i) = out(1,i)
31            ind(i) = out(2,i) ! rank is coerced back to an integer
32        END DO
33    END IF
34

```

Example 4.19 Each process has a non-empty array of values. Find the minimum global value, the rank of the process that holds it and its index on this process.

```

37 #define LEN 1000
38
39 float val[LEN];          /* local array of values */
40 int count;              /* local number of values */
41 int myrank, minrank, minindex;
42 float minval;
43
44 struct {
45     float value;
46     int index;
47 } in, out;
48

```

```

    /* local minloc */
in.value = val[0];
in.index = 0;
for (i=1; i < count; i++)
    if (in.value > val[i]) {
        in.value = val[i];
        in.index = i;
    }

    /* global minloc */
MPI_Comm_rank(comm, &myrank);
in.index = myrank*LEN + in.index;
MPI_Reduce( in, out, 1, MPI_FLOAT_INT, MPI_MINLOC, root, comm );
    /* At this point, the answer resides on process root
    */
if (myrank == root) {
    /* read answer out
    */
    minval = out.value;
    minrank = out.index / LEN;
    minindex = out.index % LEN;
}

```

Rationale. The definition of MPI_MINLOC and MPI_MAXLOC given here has the advantage that it does not require any special-case handling of these two operations: they are handled like any other reduce operation. A programmer can provide his or her own definition of MPI_MAXLOC and MPI_MINLOC, if so desired. The disadvantage is that values and indices have to be first interleaved, and that indices and values have to be coerced to the same type, in Fortran. (*End of rationale.*)

4.10.5 User-Defined Operations

MPI_OP_CREATE(function, commute, op)

| | | |
|-----|----------|---------------------------------------|
| IN | function | user defined function (function) |
| IN | commute | true if commutative; false otherwise. |
| OUT | op | operation (handle) |

int MPI_Op_create(MPI_User_function *function, int commute, MPI_Op *op)

MPI_OP_CREATE(FUNCTION, COMMUTE, OP, IERROR)

EXTERNAL FUNCTION

LOGICAL COMMUTE

INTEGER OP, IERROR

void MPI::Op::Init(MPI::User_function* function, bool commute)

1 MPI_OP_CREATE binds a user-defined global operation to an op handle that can
 2 subsequently be used in MPI_REDUCE, MPI_ALLREDUCE, MPI_REDUCE_SCATTER, and
 3 MPI_SCAN. The user-defined operation is assumed to be associative. If `commute = true`,
 4 then the operation should be both commutative and associative. If `commute = false`,
 5 then the order of operands is fixed and is defined to be in ascending, process rank order,
 6 beginning with process zero. The order of evaluation can be changed, taking advantage of
 7 the associativity of the operation. If `commute = true` then the order of evaluation can be
 8 changed, taking advantage of commutativity and associativity.

9 function is the user-defined function, which must have the following four arguments:
 10 `invec`, `inoutvec`, `len` and `datatype`.

11 The ISO C prototype for the function is the following.

```
12
13 typedef void MPI_User_function( void *invec, void *inoutvec, int *len,
14                               MPI_Datatype *datatype);
```

15 The Fortran declaration of the user-defined function appears below.

```
16
17 SUBROUTINE USER_FUNCTION( INVEC, INOUTVEC, LEN, TYPE)
18   <type> INVEC(LEN), INOUTVEC(LEN)
19   INTEGER LEN, TYPE
```

20
 21 The `datatype` argument is a handle to the data type that was passed into the call to
 22 MPI_REDUCE. The user reduce function should be written such that the following holds:
 23 Let $u[0], \dots, u[\text{len}-1]$ be the `len` elements in the communication buffer described by the
 24 arguments `invec`, `len` and `datatype` when the function is invoked; let $v[0], \dots, v[\text{len}-1]$ be `len`
 25 elements in the communication buffer described by the arguments `inoutvec`, `len` and `datatype`
 26 when the function is invoked; let $w[0], \dots, w[\text{len}-1]$ be `len` elements in the communication
 27 buffer described by the arguments `inoutvec`, `len` and `datatype` when the function returns;
 28 then $w[i] = u[i] \circ v[i]$, for $i=0, \dots, \text{len}-1$, where \circ is the reduce operation that the function
 29 computes.

30 Informally, we can think of `invec` and `inoutvec` as arrays of `len` elements that function
 31 is combining. The result of the reduction over-writes values in `inoutvec`, hence the name.
 32 Each invocation of the function results in the pointwise evaluation of the reduce operator
 33 on `len` elements: I.e, the function returns in `inoutvec[i]` the value `invec[i] \circ inoutvec[i]`, for
 34 $i = 0, \dots, \text{count} - 1$, where \circ is the combining operation computed by the function.

35
 36 *Rationale.* The `len` argument allows MPI_REDUCE to avoid calling the function for
 37 each element in the input buffer. Rather, the system can choose to apply the function
 38 to chunks of input. In C, it is passed in as a reference for reasons of compatibility
 39 with Fortran.

40 By internally comparing the value of the `datatype` argument to known, global handles,
 41 it is possible to overload the use of a single user-defined function for several, different
 42 data types. (*End of rationale.*)

43
 44 General datatypes may be passed to the user function. However, use of datatypes that
 45 are not contiguous is likely to lead to inefficiencies.

46 No MPI communication function may be called inside the user function. MPI_ABORT
 47 may be called inside the function in case of an error.

Advice to users. Suppose one defines a library of user-defined reduce functions that are overloaded: the `datatype` argument is used to select the right execution path at each invocation, according to the types of the operands. The user-defined reduce function cannot “decode” the `datatype` argument that it is passed, and cannot identify, by itself, the correspondence between the `datatype` handles and the `datatype` they represent. This correspondence was established when the `datatype`s were created. Before the library is used, a library initialization preamble must be executed. This preamble code will define the `datatype`s that are used by the library, and store handles to these `datatype`s in global, static variables that are shared by the user code and the library code.

The Fortran version of `MPI_REDUCE` will invoke a user-defined reduce function using the Fortran calling conventions and will pass a Fortran-type `datatype` argument; the C version will use C calling convention and the C representation of a `datatype` handle. Users who plan to mix languages should define their reduction functions accordingly. (*End of advice to users.*)

Advice to implementors. We outline below a naive and inefficient implementation of `MPI_REDUCE`.

```

if (rank > 0) {
    RECV(tempbuf, count, datatype, rank-1,...)
    User_reduce( tempbuf, sendbuf, count, datatype)
}
if (rank < groupsize-1) {
    SEND( sendbuf, count, datatype, rank+1, ...)
}
/* answer now resides in process groupsize-1 ... now send to root
*/
if (rank == groupsize-1) {
    SEND( sendbuf, count, datatype, root, ...)
}
if (rank == root) {
    RECV(recvbuf, count, datatype, groupsize-1,...)
}

```

The reduction computation proceeds, sequentially, from process 0 to process `group-size-1`. This order is chosen so as to respect the order of a possibly non-commutative operator defined by the function `User_reduce()`. A more efficient implementation is achieved by taking advantage of associativity and using a logarithmic tree reduction. Commutativity can be used to advantage, for those cases in which the `commute` argument to `MPI_OP_CREATE` is true. Also, the amount of temporary buffer required can be reduced, and communication can be pipelined with computation, by transferring and reducing the elements in chunks of size `len < count`.

The predefined reduce operations can be implemented as a library of user-defined operations. However, better performance might be achieved if `MPI_REDUCE` handles these functions as a special case. (*End of advice to implementors.*)

```

1 MPI_OP_FREE( op)
2     INOUT    op                operation (handle)
3

```

```

4
5 int MPI_op_free( MPI_Op *op)

```

```

6 MPI_OP_FREE( OP, IERROR)
7     INTEGER OP, IERROR
8

```

```

9 void MPI::Op::Free()

```

```

10     Marks a user-defined reduction operation for deallocation and sets op to MPI_OP_NULL.
11

```

```

12 Example of User-defined Reduce
13

```

```

14 It is time for an example of user-defined reduction. The example in this section is using an intracommunicator.
15

```

```

16 Example 4.20 Compute the product of an array of complex numbers, in C.
17

```

```

18
19 typedef struct {
20     double real,imag;
21 } Complex;
22
23 /* the user-defined function
24  */
25 void myProd( Complex *in, Complex *inout, int *len, MPI_Datatype *dptr )
26 {
27     int i;
28     Complex c;
29
30     for (i=0; i< *len; ++i) {
31         c.real = inout->real*in->real -
32             inout->imag*in->imag;
33         c.imag = inout->real*in->imag +
34             inout->imag*in->real;
35         *inout = c;
36         inout++;
37     }
38 }
39
40 /* and, to call it...
41  */
42 ...
43
44     /* each process has an array of 100 Complexes
45     */
46     Complex a[100], answer[100];
47     MPI_Op myOp;
48     MPI_Datatype ctype;

```

```

1
2  /* explain to MPI how type Complex is defined
3  */
4  MPI_Type_contiguous( 2, MPI_DOUBLE, &ctype );
5  MPI_Type_commit( &ctype );
6  /* create the complex-product user-op
7  */
8  MPI_Op_create( myProd, True, &myOp );
9
10 MPI_Reduce( a, answer, 100, ctype, myOp, root, comm );
11
12 /* At this point, the answer, which consists of 100 Complexes,
13  * resides on process root
14  */
15
16

```

4.10.6 All-Reduce

MPI includes variants of each of the reduce operations where the result is returned to all processes in the group. MPI requires that all processes participating in these operations receive identical results.

MPI_ALLREDUCE(sendbuf, recvbuf, count, datatype, op, comm)

| | | | |
|-----|----------|---|----|
| IN | sendbuf | starting address of send buffer (choice) | 25 |
| OUT | recvbuf | starting address of receive buffer (choice) | 26 |
| IN | count | number of elements in send buffer (integer) | 27 |
| IN | datatype | data type of elements of send buffer (handle) | 29 |
| IN | op | operation (handle) | 30 |
| IN | comm | communicator (handle) | 31 |

```

32
33
34 int MPI_Allreduce(void* sendbuf, void* recvbuf, int count,
35                 MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
36

```

```

37 MPI_ALLREDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, IERROR)
38 <type> SENDBUF(*), RECVBUF(*)
39 INTEGER COUNT, DATATYPE, OP, COMM, IERROR

```

```

40 void MPI::Comm::Allreduce(const void* sendbuf, void* recvbuf, int count,
41                          const MPI::Datatype& datatype, const MPI::Op& op) const = 0

```

Same as MPI_REDUCE except that the result appears in the receive buffer of all the group members.

Advice to implementors. The all-reduce operations can be implemented as a reduce, followed by a broadcast. However, a direct implementation can lead to better performance. (*End of advice to implementors.*)

The “in place” option for intracommunicators is specified by passing the value `MPI_IN_PLACE` to the argument `sendbuf` at all processes. In such case, the input data is taken at each process from the receive buffer, where it will be replaced by the output data.

If `comm` is an intercommunicator, then the result of the reduction of the data provided by processes in group A is stored at each process in group B, and vice versa. Both groups should provide `count` and `datatype` arguments that specify the same type signature.

The following example is using an intracommunicators.

Example 4.21 A routine that computes the product of a vector and an array that are distributed across a group of processes and returns the answer at all nodes (see also Example 4.16).

```

12 SUBROUTINE PAR_BLAS2(m, n, a, b, c, comm)
13 REAL a(m), b(m,n)    ! local slice of array
14 REAL c(n)           ! result
15 REAL sum(n)
16 INTEGER n, comm, i, j, ierr
17
18 ! local sum
19 DO j= 1, n
20   sum(j) = 0.0
21   DO i = 1, m
22     sum(j) = sum(j) + a(i)*b(i,j)
23   END DO
24 END DO
25
26 ! global sum
27 CALL MPI_ALLREDUCE(sum, c, n, MPI_REAL, MPI_SUM, comm, ierr)
28
29 ! return result at all nodes
30 RETURN
31

```

4.11 Reduce-Scatter

MPI includes variants of each of the reduce operations where the result is scattered to all processes in the group on return.


```

MPI_REDUCE_SCATTER( sendbuf, recvbuf, recvcnts, datatype, op, comm) 1
    IN      sendbuf      starting address of send buffer (choice) 2
    OUT     recvbuf      starting address of receive buffer (choice) 3
    IN      recvcnts     integer array specifying the number of elements in re- 4
                        sult distributed to each process. Array must be iden- 5
                        tical on all calling processes. 6
    IN      datatype     data type of elements of input buffer (handle) 7
    IN      op           operation (handle) 8
    IN      comm         communicator (handle) 9
                                10
                                11
                                12
int MPI_Reduce_scatter(void* sendbuf, void* recvbuf, int *recvcnts, 13
                      MPI_Datatype datatype, MPI_Op op, MPI_Comm comm) 14
MPI_REDUCE_SCATTER(SENDBUF, RECVBUF, RECVCOUNTS, DATATYPE, OP, COMM, 15
                  IERROR) 16
<type> SENDBUF(*), RECVBUF(*) 17
INTEGER RECVCOUNTS(*), DATATYPE, OP, COMM, IERROR 18
void MPI::Comm::Reduce_scatter(const void* sendbuf, void* recvbuf, 19
                              int recvcnts[], const MPI::Datatype& datatype, 20
                              const MPI::Op& op) const = 0 21
                                22
                                23

```

MPI_REDUCE_SCATTER first does an element-wise reduction on vector of `count = $\sum_i \text{recvcnts}[i]$` elements in the send buffer defined by `sendbuf`, `count` and `datatype`. Next, the resulting vector of results is split into `n` disjoint segments, where `n` is the number of members in the group. Segment `i` contains `recvcnts[i]` elements. The `i`th segment is sent to process `i` and stored in the receive buffer defined by `recvbuf`, `recvcnts[i]` and `datatype`.

Advice to implementors. The MPI_REDUCE_SCATTER routine is functionally equivalent to: A MPI_REDUCE operation function with `count` equal to the sum of `recvcnts[i]` followed by MPI_SCATTERV with `sendcounts` equal to `recvcnts`. However, a direct implementation may run faster. (*End of advice to implementors.*)

The “in place” option for intracommunicators is specified by passing `MPI_IN_PLACE` in the `sendbuf` argument. In this case, the input data is taken from the top of the receive buffer.

If `comm` is an intercommunicator, then the result of the reduction of the data provided by processes in group A is scattered among processes in group B, and vice versa. Within each group, all processes provide the same `recvcnts` argument, and the sum of the `recvcnts` entries should be the same for the two groups.

Rationale. The last restriction is needed so that the length of the send buffer can be determined by the sum of the local `recvcnts` entries. Otherwise, a communication is needed to figure out how many elements are reduced. (*End of rationale.*)

4.12 Scan

```
MPI_SCAN( sendbuf, recvbuf, count, datatype, op, comm )
```

| | | |
|-----|----------|--|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | count | number of elements in input buffer (integer) |
| IN | datatype | data type of elements of input buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |

```
int MPI_Scan(void* sendbuf, void* recvbuf, int count,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm )
```

```
MPI_SCAN(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER COUNT, DATATYPE, OP, COMM, IERROR
```

```
void MPI::Intracomm::Scan(const void* sendbuf, void* recvbuf, int count,
                        const MPI::Datatype& datatype, const MPI::Op& op) const
```

MPI_SCAN is used to perform a prefix reduction on data distributed across the group. The operation returns, in the receive buffer of the process with rank i , the reduction of the values in the send buffers of processes with ranks $0, \dots, i$ (inclusive). The type of operations supported, their semantics, and the constraints on send and receive buffers are as for MPI_REDUCE.

The “in place” option for intracommunicators is specified by passing MPI_IN_PLACE in the sendbuf argument. In this case, the input data is taken from the receive buffer, and replaced by the output data.

This operation is illegal for intercommunicators.

Rationale. We have defined an inclusive scan, that is, the prefix reduction on process i includes the data from process i . An alternative is to define scan in an exclusive manner, where the result on i only includes data up to $i-1$. Both definitions are useful. The latter has some advantages: the inclusive scan can always be computed from the exclusive scan with no additional communication; for non-invertible operations such as max and min, communication is required to compute the exclusive scan from the inclusive scan. There is, however, a complication with exclusive scan since one must define the “unit” element for the reduction in this case. That is, one must explicitly say what occurs for process 0. This was thought to be complex for user-defined operations and hence, the exclusive scan was dropped. (*End of rationale.*)

4.12.1 Exclusive Scan

MPI-1 provides an inclusive scan operation. The exclusive scan is described here.

| | | | |
|---|----------|--|---|
| MPI_EXSCAN(sendbuf, recvbuf, count, datatype, op, comm) | | | 1 |
| IN | sendbuf | starting address of send buffer (choice) | 2 |
| OUT | recvbuf | starting address of receive buffer (choice) | 3 |
| IN | count | number of elements in input buffer (integer) | 4 |
| IN | datatype | data type of elements of input buffer (handle) | 5 |
| IN | op | operation (handle) | 6 |
| IN | comm | intracommunicator (handle) | 7 |

```

int MPI_Exscan(void *sendbuf, void *recvbuf, int count,
               MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
MPI_EXSCAN(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, IERROR)
<type> SENDBUF(*), RECVBUF(*)
INTEGER COUNT, DATATYPE, OP, COMM, IERROR
void MPI::Intracomm::Exscan(const void* sendbuf, void* recvbuf, int count,
                           const MPI::Datatype& datatype, const MPI::Op& op) const

```

MPI_EXSCAN is used to perform a prefix reduction on data distributed across the group. The value in `recvbuf` on the process with rank 0 is undefined, and `recvbuf` is not significant on process 0. The value in `recvbuf` on the process with rank 1 is defined as the value in `sendbuf` on the process with rank 0. For processes with rank $i > 1$, the operation returns, in the receive buffer of the process with rank i , the reduction of the values in the send buffers of processes with ranks $0, \dots, i - 1$ (inclusive). The type of operations supported, their semantics, and the constraints on send and receive buffers, are as for MPI_REDUCE.

No “in place” option is supported.

Advice to users. As for MPI_SCAN, MPI does not specify which processes may call the operation, only that the result be correctly computed. In particular, note that the process with rank 1 need not call the MPI_Op, since all it needs to do is to receive the value from the process with rank 0. However, all processes, even the processes with ranks zero and one, must provide the same op. (*End of advice to users.*)

Rationale. The exclusive scan is more general than the inclusive scan provided in MPI-1 as MPI_SCAN. Any inclusive scan operation can be achieved by using the exclusive scan and then locally combining the local contribution. Note that for non-invertable operations such as MPI_MAX, the exclusive scan cannot be computed with the inclusive scan.

No in-place version is specified for MPI_EXSCAN because it is not clear what this means for the process for rank zero. The reason that MPI-1 chose the inclusive scan is that the definition of behavior on processes zero and one was thought to offer too many complexities in definition, particularly for user-defined operations. (*End of rationale.*)

4.12.2 Example using MPI_SCAN

The example in this section is using an intracommunicators.

Example 4.22 This example uses a user-defined operation to produce a *segmented scan*. A segmented scan takes, as input, a set of values and a set of logicals, and the logicals delineate the various segments of the scan. For example:

| | | | | | | | | |
|-----------------|-------|-------------|-------|-------------|-------------------|-------|-------------|-------|
| <i>values</i> | v_1 | v_2 | v_3 | v_4 | v_5 | v_6 | v_7 | v_8 |
| <i>logicals</i> | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 |
| <i>result</i> | v_1 | $v_1 + v_2$ | v_3 | $v_3 + v_4$ | $v_3 + v_4 + v_5$ | v_6 | $v_6 + v_7$ | v_8 |

The operator that produces this effect is,

$$\begin{pmatrix} u \\ i \end{pmatrix} \circ \begin{pmatrix} v \\ j \end{pmatrix} = \begin{pmatrix} w \\ j \end{pmatrix},$$

where,

$$w = \begin{cases} u + v & \text{if } i = j \\ v & \text{if } i \neq j \end{cases}.$$

Note that this is a non-commutative operator. C code that implements it is given below.

```

typedef struct {
    double val;
    int log;
} SegScanPair;

/* the user-defined function
*/
void segScan( SegScanPair *in, SegScanPair *inout, int *len,
              MPI_Datatype *dptr )
{
    int i;
    SegScanPair c;

    for (i=0; i< *len; ++i) {
        if ( in->log == inout->log )
            c.val = in->val + inout->val;
        else
            c.val = inout->val;
        c.log = inout->log;
        *inout = c;
        in++; inout++;
    }
}

```

Note that the *inout* argument to the user-defined function corresponds to the right-hand operand of the operator. When using this operator, we must be careful to specify that it is non-commutative, as in the following.

```

int i,base;
SeqScanPair a, answer;
MPI_Op      myOp;
MPI_Datatype type[2] = {MPI_DOUBLE, MPI_INT};
MPI_Aint    disp[2];
int         blocklen[2] = { 1, 1};
MPI_Datatype sspair;

/* explain to MPI how type SegScanPair is defined
*/
MPI_Address( a, disp);
MPI_Address( a.log, disp+1);
base = disp[0];
for (i=0; i<2; ++i) disp[i] -= base;
MPI_Type_struct( 2, blocklen, disp, type, &sspair );
MPI_Type_commit( &sspair );
/* create the segmented-scan user-op
*/
MPI_Op_create( segScan, False, &myOp );
...
MPI_Scan( a, answer, 1, sspair, myOp, comm );

```

4.13 Correctness

A correct, portable program must invoke collective communications so that deadlock will not occur, whether collective communications are synchronizing or not. The following examples illustrate dangerous use of collective routines [on intracommunicators](#).

Example 4.23 The following is erroneous.

```

switch(rank) {
  case 0:
    MPI_Bcast(buf1, count, type, 0, comm);
    MPI_Bcast(buf2, count, type, 1, comm);
    break;
  case 1:
    MPI_Bcast(buf2, count, type, 1, comm);
    MPI_Bcast(buf1, count, type, 0, comm);
    break;
}

```

We assume that the group of `comm` is $\{0,1\}$. Two processes execute two broadcast operations in reverse order. If the operation is synchronizing then a deadlock will occur.

Collective operations must be executed in the same order at all members of the communication group.

Example 4.24 The following is erroneous.

```

1  switch(rank) {
2      case 0:
3          MPI_Bcast(buf1, count, type, 0, comm0);
4          MPI_Bcast(buf2, count, type, 2, comm2);
5          break;
6      case 1:
7          MPI_Bcast(buf1, count, type, 1, comm1);
8          MPI_Bcast(buf2, count, type, 0, comm0);
9          break;
10     case 2:
11         MPI_Bcast(buf1, count, type, 2, comm2);
12         MPI_Bcast(buf2, count, type, 1, comm1);
13         break;
14 }

```

15

16 Assume that the group of `comm0` is $\{0,1\}$, of `comm1` is $\{1, 2\}$ and of `comm2` is $\{2,0\}$. If
17 the broadcast is a synchronizing operation, then there is a cyclic dependency: the broadcast
18 in `comm2` completes only after the broadcast in `comm0`; the broadcast in `comm0` completes
19 only after the broadcast in `comm1`; and the broadcast in `comm1` completes only after the
20 broadcast in `comm2`. Thus, the code will deadlock.

21 Collective operations must be executed in an order so that no cyclic dependences occur.

22

23 **Example 4.25** The following is erroneous.

24

```

25 switch(rank) {
26     case 0:
27         MPI_Bcast(buf1, count, type, 0, comm);
28         MPI_Send(buf2, count, type, 1, tag, comm);
29         break;
30     case 1:
31         MPI_Recv(buf2, count, type, 0, tag, comm, status);
32         MPI_Bcast(buf1, count, type, 0, comm);
33         break;
34 }

```

35

36 Process zero executes a broadcast, followed by a blocking send operation. Process one
37 first executes a blocking receive that matches the send, followed by broadcast call that
38 matches the broadcast of process zero. This program may deadlock. The broadcast call on
39 process zero *may* block until process one executes the matching broadcast call, so that the
40 send is not executed. Process one will definitely block on the receive and so, in this case,
41 never executes the broadcast.

41

42 The relative order of execution of collective operations and point-to-point operations
43 should be such, so that even if the collective operations and the point-to-point operations
44 are synchronizing, no deadlock will occur.

44

45 **Example 4.26** A correct, but non-deterministic program.

46

```

47 switch(rank) {
48     case 0:

```

48

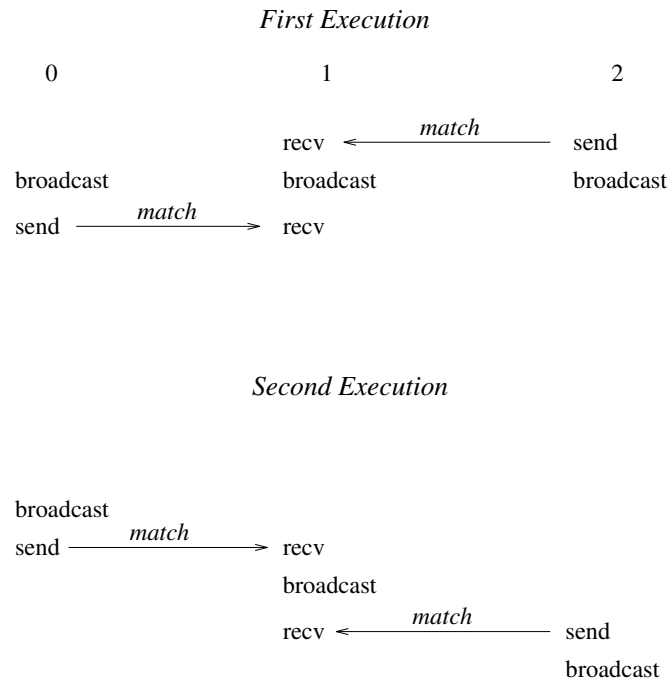


Figure 4.12: A race condition causes non-deterministic matching of sends and receives. One cannot rely on synchronization from a broadcast to make the program deterministic.

```

    MPI_Bcast(buf1, count, type, 0, comm);
    MPI_Send(buf2, count, type, 1, tag, comm);
    break;
case 1:
    MPI_Recv(buf2, count, type, MPI_ANY_SOURCE, tag, comm, status);
    MPI_Bcast(buf1, count, type, 0, comm);
    MPI_Recv(buf2, count, type, MPI_ANY_SOURCE, tag, comm, status);
    break;
case 2:
    MPI_Send(buf2, count, type, 1, tag, comm);
    MPI_Bcast(buf1, count, type, 0, comm);
    break;
}

```

All three processes participate in a broadcast. Process 0 sends a message to process 1 after the broadcast, and process 2 sends a message to process 1 before the broadcast. Process 1 receives before and after the broadcast, with a wildcard source argument.

Two possible executions of this program, with different matchings of sends and receives, are illustrated in figure 4.12. Note that the second execution has the peculiar effect that a send executed after the broadcast is received at another node before the broadcast. This example illustrates the fact that one should not rely on collective communication functions to have particular synchronization effects. A program that works correctly only when the first execution occurs (only when broadcast is synchronizing) is erroneous.

Finally, in multithreaded implementations, one can have more than one, concurrently executing, collective communication call at a process. In these situations, it is the user's re-

1 responsibility to ensure that the same communicator is not used concurrently by two different
2 collective communication calls at the same process.

3
4 *Advice to implementors.* Assume that broadcast is implemented using point-to-point
5 MPI communication. Suppose the following two rules are followed.

- 6 1. All receives specify their source explicitly (no wildcards).
- 7 2. Each process sends all messages that pertain to one collective call before sending
8 any message that pertain to a subsequent collective call.
9

10 Then, messages belonging to successive broadcasts cannot be confused, as the order
11 of point-to-point messages is preserved.
12

13 It is the implementor's responsibility to ensure that point-to-point messages are not
14 confused with collective messages. One way to accomplish this is, whenever a commu-
15 nicator is created, to also create a "hidden communicator" for collective communica-
16 tion. One could achieve a similar effect more cheaply, for example, by using a hidden
17 tag or context bit to indicate whether the communicator is used for point-to-point or
18 collective communication. (*End of advice to implementors.*)
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 5

Groups, Contexts, and Communicators

5.1 Introduction

This chapter introduces MPI features that support the development of parallel libraries. Parallel libraries are needed to encapsulate the distracting complications inherent in parallel implementations of key algorithms. They help to ensure consistent correctness of such procedures, and provide a “higher level” of portability than MPI itself can provide. As such, libraries prevent each programmer from repeating the work of defining consistent data structures, data layouts, and methods that implement key algorithms (such as matrix operations). Since the best libraries come with several variations on parallel systems (different data layouts, different strategies depending on the size of the system or problem, or type of floating point), this too needs to be hidden from the user.

We refer the reader to [44] and [3] for further information on writing libraries in MPI, using the features described in this chapter.

5.1.1 Features Needed to Support Libraries

The key features needed to support the creation of robust parallel libraries are as follows:

- Safe communication space, that guarantees that libraries can communicate as they need to, without conflicting with communication extraneous to the library,
- Group scope for collective operations, that allow libraries to avoid unnecessarily synchronizing uninvolved processes (potentially running unrelated code),
- Abstract process naming to allow libraries to describe their communication in terms suitable to their own data structures and algorithms,
- The ability to “adorn” a set of communicating processes with additional user-defined attributes, such as extra collective operations. This mechanism should provide a means for the user or library writer effectively to extend a message-passing notation.

In addition, a unified mechanism or object is needed for conveniently denoting communication context, the group of communicating processes, to house abstract process naming, and to store adornments.

5.1.2 MPI's Support for Libraries

The corresponding concepts that MPI provides, specifically to support robust libraries, are as follows:

- **Contexts** of communication,
- **Groups** of processes,
- **Virtual topologies**,
- **Attribute caching**,
- **Communicators**.

Communicators (see [20, 42, 47]) encapsulate all of these ideas in order to provide the appropriate scope for all communication operations in MPI. Communicators are divided into two kinds: intra-communicators for operations within a single group of processes, and inter-communicators, for point-to-point communication between two groups of processes.

Caching. Communicators (see below) provide a “caching” mechanism that allows one to associate new attributes with communicators, on a par with MPI built-in features. This can be used by advanced users to adorn communicators further, and by MPI to implement some communicator functions. For example, the virtual-topology functions described in Chapter 6 are likely to be supported this way.

Groups. Groups define an ordered collection of processes, each with a rank, and it is this group that defines the low-level names for inter-process communication (ranks are used for sending and receiving). Thus, groups define a scope for process names in point-to-point communication. In addition, groups define the scope of collective operations. Groups may be manipulated separately from communicators in MPI, but only communicators can be used in communication operations.

Intra-communicators. The most commonly used means for message passing in MPI is via intra-communicators. Intra-communicators contain an instance of a group, contexts of communication for both point-to-point and collective communication, and the ability to include virtual topology and other attributes. These features work as follows:

- **Contexts** provide the ability to have separate safe “universes” of message passing in MPI. A context is akin to an additional tag that differentiates messages. The system manages this differentiation process. The use of separate communication contexts by distinct libraries (or distinct library invocations) insulates communication internal to the library execution from external communication. This allows the invocation of the library even if there are pending communications on “other” communicators, and avoids the need to synchronize entry or exit into library code. Pending point-to-point communications are also guaranteed not to interfere with collective communications within a single communicator.
- **Groups** define the participants in the communication (see above) of a communicator.

- A **virtual topology** defines a special mapping of the ranks in a group to and from a topology. Special constructors for communicators are defined in chapter 6 to provide this feature. Intra-communicators as described in this chapter do not have topologies.
- **Attributes** define the local information that the user or library has added to a communicator for later reference.

Advice to users. The current practice in many communication libraries is that there is a unique, predefined communication universe that includes all processes available when the parallel program is initiated; the processes are assigned consecutive ranks. Participants in a point-to-point communication are identified by their rank; a collective communication (such as broadcast) always involves all processes. This practice can be followed in MPI by using the predefined communicator `MPI_COMM_WORLD`. *Users who are satisfied with this practice can plug in `MPI_COMM_WORLD` wherever a communicator argument is required, and can consequently disregard the rest of this chapter. (End of advice to users.)*

Inter-communicators. The discussion has dealt so far with **intra-communication**: communication within a group. MPI also supports **inter-communication**: communication between two non-overlapping groups. When an application is built by composing several parallel modules, it is convenient to allow one module to communicate with another using local ranks for addressing within the second module. This is especially convenient in a client-server computing paradigm, where either client or server are parallel. The support of inter-communication also provides a mechanism for the extension of MPI to a dynamic model where not all processes are preallocated at initialization time. In such a situation, it becomes necessary to support communication across “universes.” Inter-communication is supported by objects called **inter-communicators**. These objects bind two groups together with communication contexts shared by both groups. For inter-communicators, these features work as follows:

- **Contexts** provide the ability to have a separate safe “universe” of message passing between the two groups. A send in the local group is always a receive in the remote group, and vice versa. The system manages this differentiation process. The use of separate communication contexts by distinct libraries (or distinct library invocations) insulates communication internal to the library execution from external communication. This allows the invocation of the library even if there are pending communications on “other” communicators, and avoids the need to synchronize entry or exit into library [code](#).
- A local and remote group specify the recipients and destinations for an inter-communicator.
- Virtual topology is undefined for an inter-communicator.
- As before, attributes cache defines the local information that the user or library has added to a communicator for later reference.

MPI provides mechanisms for creating and manipulating inter-communicators. They are used for point-to-point [and collective](#) communication in an related manner to intra-communicators. Users who do not need inter-communication in their applications can safely

1 ignore this extension. **Users** who require inter-communication between overlapping groups
2 **must layer** this capability on top of MPI.
3

4 5.2 Basic Concepts 5

6 In this section, we turn to a more formal definition of the concepts introduced above.
7

8 5.2.1 Groups 9

10 A **group** is an ordered set of process identifiers (henceforth processes); processes are
11 implementation-dependent objects. Each process in a group is associated with an inte-
12 ger **rank**. Ranks are contiguous and start from zero. Groups are represented by opaque
13 **group objects**, and hence cannot be directly transferred from one process to another. A
14 group is used within a communicator to describe the participants in a communication “uni-
15 verse” and to rank such participants (thus giving them unique names within that “universe”
16 of communication).

17 There is a special pre-defined group: `MPI_GROUP_EMPTY`, which is a group with no
18 members. The predefined constant `MPI_GROUP_NULL` is the value used for invalid group
19 handles.
20

21 *Advice to users.* `MPI_GROUP_EMPTY`, which is a valid handle to an empty group,
22 should not be confused with `MPI_GROUP_NULL`, which in turn is an invalid handle.
23 The former may be used as an argument to group operations; the latter, which is
24 returned when a group is freed, is not a valid argument. (*End of advice to users.*)
25

26 *Advice to implementors.* A group may be represented by a virtual-to-real process-
27 address-translation table. Each communicator object (see below) would have a pointer
28 to such a table.

29 Simple implementations of MPI will enumerate groups, such as in a table. However,
30 more advanced data structures make sense in order to improve scalability and memory
31 usage with large numbers of processes. Such implementations are possible with MPI.
32 (*End of advice to implementors.*)
33

34 5.2.2 Contexts 35

36 A **context** is a property of communicators (defined next) that allows partitioning of the
37 communication space. A message sent in one context cannot be received in another context.
38 Furthermore, where permitted, collective operations are independent of pending point-to-
39 point operations. Contexts are not explicit MPI objects; they appear only as part of the
40 realization of communicators (below).
41

42 *Advice to implementors.* Distinct communicators in the same process have distinct
43 contexts. A context is essentially a system-managed tag (or tags) needed to make
44 a communicator safe for point-to-point and MPI-defined collective communication.
45 Safety means that collective and point-to-point communication within one commu-
46 nicator do not interfere, and that communication over distinct communicators don't
47 interfere.
48

A possible implementation for a context is as a supplemental tag attached to messages on send and matched on receive. Each intra-communicator stores the value of its two tags (one for point-to-point and one for collective communication). Communicator-generating functions use a collective communication to agree on a new group-wide unique context.

Analogously, in [inter-communication](#), two context tags are stored per communicator, one used by group A to send and group B to receive, and a second used by group B to send and for group A to receive.

Since contexts are not explicit objects, other implementations are also possible. (*End of advice to implementors.*)

5.2.3 Intra-Communicators

Intra-communicators bring together the concepts of group and context. To support implementation-specific optimizations, and application topologies (defined in the next chapter, chapter 6), communicators may also “cache” additional information (see section 5.7). MPI communication operations reference communicators to determine the scope and the “communication universe” in which a point-to-point or collective operation is to operate.

Each communicator contains a group of valid participants; this group always includes the local process. The source and destination of a message is identified by process rank within that group.

For collective communication, the intra-communicator specifies the set of processes that participate in the collective operation (and their order, when significant). Thus, the communicator restricts the “spatial” scope of communication, and provides machine-independent process addressing through ranks.

Intra-communicators are represented by opaque **intra-communicator objects**, and hence cannot be directly transferred from one process to another.

5.2.4 Predefined Intra-Communicators

An initial intra-communicator `MPI_COMM_WORLD` of all processes the local process can communicate with after initialization (itself included) is defined once `MPI_INIT` or `MPI_INIT_THREAD` has been called. In addition, the communicator `MPI_COMM_SELF` is provided, which includes only the process itself.

The predefined constant `MPI_COMM_NULL` is the value used for invalid communicator handles.

In a static-process-model implementation of MPI, all processes that participate in the computation are available after MPI is initialized. For this case, `MPI_COMM_WORLD` is a communicator of all processes available for the computation; this communicator has the same value in all processes. In an implementation of MPI where processes can dynamically join an MPI execution, it may be the case that a process starts an MPI computation without having access to all other processes. In such situations, `MPI_COMM_WORLD` is a communicator incorporating all processes with which the joining process can immediately communicate. Therefore, `MPI_COMM_WORLD` may simultaneously [represent disjoint groups](#) in different processes.

All MPI implementations are required to provide the `MPI_COMM_WORLD` communicator. It cannot be deallocated during the life of a process. The group corresponding to

1 this communicator does not appear as a pre-defined constant, but it may be accessed using
 2 `MPI_COMM_GROUP` (see below). MPI does not specify the correspondence between the
 3 process rank in `MPI_COMM_WORLD` and its (machine-dependent) absolute address. Neither
 4 does MPI specify the function of the host process, if any. Other implementation-dependent,
 5 predefined communicators may also be provided.

7 5.3 Group Management

9 This section describes the manipulation of process groups in MPI. These operations are
 10 local and their execution do not require interprocess communication.

12 5.3.1 Group Accessors

15 `MPI_GROUP_SIZE(group, size)`

| | | | |
|----|-----|-------|--|
| 17 | IN | group | group (handle) |
| 18 | OUT | size | number of processes in the group (integer) |

20 `int MPI_Group_size(MPI_Group group, int *size)`

22 `MPI_GROUP_SIZE(GROUP, SIZE, IERROR)`
 23 `INTEGER GROUP, SIZE, IERROR`

24 `int MPI::Group::Get_size() const`

27 `MPI_GROUP_RANK(group, rank)`

| | | | |
|----|-----|-------|--|
| 29 | IN | group | group (handle) |
| 30 | OUT | rank | rank of the calling process in group, or |
| 31 | | | <code>MPI_UNDEFINED</code> if the process is not a member (in- |
| 32 | | | teger) |

34 `int MPI_Group_rank(MPI_Group group, int *rank)`

36 `MPI_GROUP_RANK(GROUP, RANK, IERROR)`
 37 `INTEGER GROUP, RANK, IERROR`

38 `int MPI::Group::Get_rank() const`

40
41
42
43
44
45
46
47
48

| | | | |
|---|--------|---|---|
| MPI_GROUP_TRANSLATE_RANKS (group1, n, ranks1, group2, ranks2) | | | 1 |
| IN | group1 | group1 (handle) | 2 |
| IN | n | number of ranks in ranks1 and ranks2 arrays (integer) | 3 |
| IN | ranks1 | array of zero or more valid ranks in group1 | 4 |
| IN | group2 | group2 (handle) | 5 |
| OUT | ranks2 | array of corresponding ranks in group2, MPI_UNDEFINED when no correspondence exists. | 6 |

```
int MPI_Group_translate_ranks (MPI_Group group1, int n, int *ranks1,
                             MPI_Group group2, int *ranks2)

```

```
MPI_GROUP_TRANSLATE_RANKS(GROUP1, N, RANKS1, GROUP2, RANKS2, IERROR)
    INTEGER GROUP1, N, RANKS1(*), GROUP2, RANKS2(*), IERROR

```

```
static void MPI::Group::Translate_ranks (const MPI::Group& group1, int n,
                                         const int ranks1[], const MPI::Group& group2, int ranks2[])

```

This function is important for determining the relative numbering of the same processes in two different groups. For instance, if one knows the ranks of certain processes in the group of MPI_COMM_WORLD, one might want to know their ranks in a subset of that group.

MPI_PROC_NULL is a valid rank for input to MPI_GROUP_TRANSLATE_RANKS, which returns MPI_PROC_NULL as the translated rank.

| | | | |
|---|--------|-----------------------|----|
| MPI_GROUP_COMPARE(group1, group2, result) | | | 25 |
| IN | group1 | first group (handle) | 26 |
| IN | group2 | second group (handle) | 27 |
| OUT | result | result (integer) | 28 |

```
int MPI_Group_compare(MPI_Group group1, MPI_Group group2, int *result)

```

```
MPI_GROUP_COMPARE(GROUP1, GROUP2, RESULT, IERROR)
    INTEGER GROUP1, GROUP2, RESULT, IERROR

```

```
static int MPI::Group::Compare(const MPI::Group& group1,
                               const MPI::Group& group2)

```

MPI_IDENT results if the group members and group order is exactly the same in both groups. This happens for instance if group1 and group2 are the same handle. MPI_SIMILAR results if the group members are the same but the order is different. MPI_UNEQUAL results otherwise.

5.3.2 Group Constructors

Group constructors are used to subset and superset existing groups. These constructors construct new groups from existing groups. These are local operations, and distinct groups may be defined on different processes; a process may also define a group that does not include itself. Consistent definitions are required when groups are used as arguments in communicator-building functions. MPI does not provide a mechanism to build a group

1 from scratch, but only from other, previously defined groups. The base group, upon which
 2 all other groups are defined, is the group associated with the initial communicator
 3 MPI_COMM_WORLD (accessible through the function MPI_COMM_GROUP).

4
 5 *Rationale.* In what follows, there is no group duplication function analogous to
 6 MPI_COMM_DUP, defined later in this chapter. There is no need for a group dupli-
 7 cator. A group, once created, can have several references to it by making copies of
 8 the handle. The following constructors address the need for subsets and supersets of
 9 existing groups. (*End of rationale.*)

10
 11 *Advice to implementors.* Each group constructor behaves as if it returned a new
 12 group object. When this new group is a copy of an existing group, then one can
 13 avoid creating such new objects, using a reference-count mechanism. (*End of advice*
 14 *to implementors.*)

15
 16
 17 MPI_COMM_GROUP(comm, group)

18
 19 IN comm communicator (handle)
 20 OUT group group corresponding to comm (handle)

21
 22 int MPI_Comm_group(MPI_Comm comm, MPI_Group *group)

23
 24 MPI_COMM_GROUP(COMM, GROUP, IERROR)
 25 INTEGER COMM, GROUP, IERROR

26 MPI::Group MPI::Comm::Get_group() const

27
 28 MPI_COMM_GROUP returns in group a handle to the group of comm.

29
 30 MPI_GROUP_UNION(group1, group2, newgroup)

31
 32 IN group1 first group (handle)
 33 IN group2 second group (handle)
 34 OUT newgroup union group (handle)

35
 36
 37 int MPI_Group_union(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)

38
 39 MPI_GROUP_UNION(GROUP1, GROUP2, NEWGROUP, IERROR)
 40 INTEGER GROUP1, GROUP2, NEWGROUP, IERROR

41 static MPI::Group MPI::Group::Union(const MPI::Group& group1,
 42 const MPI::Group& group2)

43
 44
 45
 46
 47
 48


```

MPI_GROUP_INTERSECTION(group1, group2, newgroup) 1
IN      group1      first group (handle) 2
IN      group2      second group (handle) 3
OUT     newgroup    intersection group (handle) 4
                                                5
                                                6
int MPI_Group_intersection(MPI_Group group1, MPI_Group group2, 7
                          MPI_Group *newgroup) 8
                                                9
MPI_GROUP_INTERSECTION(GROUP1, GROUP2, NEWGROUP, IERROR) 10
    INTEGER GROUP1, GROUP2, NEWGROUP, IERROR 11
static MPI::Group MPI::Group::Intersect(const MPI::Group& group1, 12
                                         const MPI::Group& group2) 13

```

```

MPI_GROUP_DIFFERENCE(group1, group2, newgroup) 16
IN      group1      first group (handle) 18
IN      group2      second group (handle) 19
OUT     newgroup    difference group (handle) 20
                                                21
int MPI_Group_difference(MPI_Group group1, MPI_Group group2, 22
                          MPI_Group *newgroup) 23
                                                24
MPI_GROUP_DIFFERENCE(GROUP1, GROUP2, NEWGROUP, IERROR) 25
    INTEGER GROUP1, GROUP2, NEWGROUP, IERROR 26
static MPI::Group MPI::Group::Difference(const MPI::Group& group1, 27
                                         const MPI::Group& group2) 28

```

The set-like operations are defined as follows:

union All elements of the first group (*group1*), followed by all elements of second group (*group2*) not in first.

intersect all elements of the first group that are also in the second group, ordered as in first group.

difference all elements of the first group that are not in the second group, ordered as in the first group.

Note that for these operations the order of processes in the output group is determined primarily by order in the first group (if possible) and then, if necessary, by order in the second group. Neither union nor intersection are commutative, but both are associative.

The new group can be empty, that is, equal to `MPI_GROUP_EMPTY`.

1 MPI_GROUP_INCL(group, n, ranks, newgroup)

| | | | |
|----|-----|----------|---|
| 2 | IN | group | group (handle) |
| 3 | | | |
| 4 | IN | n | number of elements in array ranks (and size of |
| 5 | | | newgroup) (integer) |
| 6 | IN | ranks | ranks of processes in group to appear in |
| 7 | | | newgroup (array of integers) |
| 8 | | | |
| 9 | OUT | newgroup | new group derived from above, in the order defined by |
| 10 | | | ranks (handle) |

11
12 int MPI_Group_incl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup)

13 MPI_GROUP_INCL(GROUP, N, RANKS, NEWGROUP, IERROR)

14 INTEGER GROUP, N, RANKS(*), NEWGROUP, IERROR

15
16 MPI::Group MPI::Group::Incl(int n, const int ranks[]) const

17 The function MPI_GROUP_INCL creates a group newgroup that consists of the
18 n processes in group with ranks rank[0], ..., rank[n-1]; the process with rank i in newgroup
19 is the process with rank ranks[i] in group. Each of the n elements of ranks must be a valid
20 rank in group and all elements must be distinct, or else the program is erroneous. If n = 0,
21 then newgroup is MPI_GROUP_EMPTY. This function can, for instance, be used to reorder
22 the elements of a group. See also MPI_GROUP_COMPARE.
23

24
25 MPI_GROUP_EXCL(group, n, ranks, newgroup)

| | | | |
|----|-----|----------|--|
| 26 | IN | group | group (handle) |
| 27 | | | |
| 28 | IN | n | number of elements in array ranks (integer) |
| 29 | IN | ranks | array of integer ranks in group not to appear in |
| 30 | | | newgroup |
| 31 | OUT | newgroup | new group derived from above, preserving the order |
| 32 | | | defined by group (handle) |

33
34
35 int MPI_Group_excl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup)

36 MPI_GROUP_EXCL(GROUP, N, RANKS, NEWGROUP, IERROR)

37 INTEGER GROUP, N, RANKS(*), NEWGROUP, IERROR

38
39 MPI::Group MPI::Group::Excl(int n, const int ranks[]) const

40 The function MPI_GROUP_EXCL creates a group of processes newgroup that is obtained
41 by deleting from group those processes with ranks ranks[0] .. ranks[n-1]. The ordering of
42 processes in newgroup is identical to the ordering in group. Each of the n elements of ranks
43 must be a valid rank in group and all elements must be distinct; otherwise, the program is
44 erroneous. If n = 0, then newgroup is identical to group.
45
46
47
48

| | | | |
|--|----------|--|------------------|
| MPI_GROUP_RANGE_INCL(group, n, ranges, newgroup) | | | 1 |
| IN | group | group (handle) | 2 |
| IN | n | number of triplets in array ranges (integer) | 3 |
| IN | ranges | a one-dimensional array of integer triplets, of the form (first rank, last rank, stride) indicating ranks in group of processes to be included in newgroup | 4 5 6 7 |
| OUT | newgroup | new group derived from above, in the order defined by ranges (handle) | 8 9 10 |

```
int MPI_Group_range_incl(MPI_Group group, int n, int ranges[] [3],
                        MPI_Group *newgroup)

```

```
MPI_GROUP_RANGE_INCL(GROUP, N, RANGES, NEWGROUP, IERROR)
INTEGER GROUP, N, RANGES(3,*), NEWGROUP, IERROR

```

```
MPI::Group MPI::Group::Range_incl(int n, const int ranges[] [3]) const

```

If ranges consist of the triplets

$$(first_1, last_1, stride_1), \dots, (first_n, last_n, stride_n)$$

then newgroup consists of the sequence of processes in group with ranks

$$first_1, first_1 + stride_1, \dots, first_1 + \left\lfloor \frac{last_1 - first_1}{stride_1} \right\rfloor stride_1, \dots$$

$$first_n, first_n + stride_n, \dots, first_n + \left\lfloor \frac{last_n - first_n}{stride_n} \right\rfloor stride_n.$$

Each computed rank must be a valid rank in group and all computed ranks must be distinct, or else the program is erroneous. Note that we may have $first_i > last_i$, and $stride_i$ may be negative, but cannot be zero.

The functionality of this routine is specified to be equivalent to expanding the array of ranges to an array of the included ranks and passing the resulting array of ranks and other arguments to MPI_GROUP_INCL. A call to MPI_GROUP_INCL is equivalent to a call to MPI_GROUP_RANGE_INCL with each rank i in ranks replaced by the triplet $(i, i, 1)$ in the argument ranges.

```

1 MPI_GROUP_RANGE_EXCL(group, n, ranges, newgroup)
2     IN      group                group (handle)
3
4     IN      n                    number of elements in array ranges (integer)
5
6     IN      ranges               a one-dimensional array of integer triplets of the form
7                                (first rank, last rank, stride), indicating the ranks in
8                                group of processes to be excluded from the output
9                                group newgroup.
10
11     OUT     newgroup            new group derived from above, preserving the order
12                                in group (handle)

```

```

13 int MPI_Group_range_excl(MPI_Group group, int n, int ranges[][3],
14                        MPI_Group *newgroup)

```

```

15 MPI_GROUP_RANGE_EXCL(GROUP, N, RANGES, NEWGROUP, IERROR)
16     INTEGER GROUP, N, RANGES(3,*), NEWGROUP, IERROR

```

```

17 MPI::Group MPI::Group::Range_excl(int n, const int ranges[][3]) const

```

Each computed rank must be a valid rank in `group` and all computed ranks must be distinct, or else the program is erroneous.

The functionality of this routine is specified to be equivalent to expanding the array of ranges to an array of the excluded ranks and passing the resulting array of ranks and other arguments to `MPI_GROUP_EXCL`. A call to `MPI_GROUP_EXCL` is equivalent to a call to `MPI_GROUP_RANGE_EXCL` with each rank `i` in `ranges` replaced by the triplet `(i,i,1)` in the argument `ranges`.

Advice to users. The range operations do not explicitly enumerate ranks, and therefore are more scalable if implemented efficiently. Hence, we recommend MPI programmers to use them whenever possible, as high-quality implementations will take advantage of this fact. (*End of advice to users.*)

Advice to implementors. The range operations should be implemented, if possible, without enumerating the group members, in order to obtain better scalability (time and space). (*End of advice to implementors.*)

5.3.3 Group Destructors

```

39 MPI_GROUP_FREE(group)

```

```

40     INOUT   group                group (handle)

```

```

43 int MPI_Group_free(MPI_Group *group)

```

```

44 MPI_GROUP_FREE(GROUP, IERROR)
45     INTEGER GROUP, IERROR

```

```

47 void MPI::Group::Free()

```

This operation marks a group object for deallocation. The handle group is set to MPI_GROUP_NULL by the call. Any on-going operation using this group will complete normally.

Advice to implementors. One can keep a reference count that is incremented for each call to MPI_COMM_CREATE and MPI_COMM_DUP, and decremented for each call to MPI_GROUP_FREE or MPI_COMM_FREE; the group object is ultimately deallocated when the reference count drops to zero. (*End of advice to implementors.*)

5.4 Communicator Management

This section describes the manipulation of communicators in MPI. Operations that access communicators are local and their execution does not require interprocess communication. Operations that create communicators are collective and may require interprocess communication.

Advice to implementors. High-quality implementations should amortize the overheads associated with the creation of communicators (for the same group, or subsets thereof) over several calls, by allocating multiple contexts with one collective communication. (*End of advice to implementors.*)

5.4.1 Communicator Accessors

The following are all local operations.

MPI_COMM_SIZE(comm, size)

| | | |
|-----|------|--|
| IN | comm | communicator (handle) |
| OUT | size | number of processes in the group of comm (integer) |

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

```
MPI_COMM_SIZE(COMM, SIZE, IERROR)
    INTEGER COMM, SIZE, IERROR
```

```
int MPI::Comm::Get_size() const
```

Rationale. This function is equivalent to accessing the communicator's group with MPI_COMM_GROUP (see above), computing the size using MPI_GROUP_SIZE, and then freeing the temporary group via MPI_GROUP_FREE. However, this function is so commonly used, that this shortcut was introduced. (*End of rationale.*)

Advice to users. This function indicates the number of processes involved in a communicator. For MPI_COMM_WORLD, it indicates the total number of processes available (for this version of MPI, there is no standard way to change the number of processes once initialization has taken place).

This call is often used with the next call to determine the amount of concurrency available for a specific library or program. The following call, MPI_COMM_RANK

1 indicates the rank of the process that calls it in the range from 0 . . . size−1, where size
 2 is the return value of MPI_COMM_SIZE. (*End of advice to users.*)

3
 4
 5 MPI_COMM_RANK(comm, rank)

6
 7 IN comm communicator (handle)
 8 OUT rank rank of the calling process in group of comm (integer)
 9

10 int MPI_Comm_rank(MPI_Comm comm, int *rank)

11 MPI_COMM_RANK(COMM, RANK, IERROR)

12 INTEGER COMM, RANK, IERROR

13
 14 int MPI::Comm::Get_rank() const

15
 16
 17 *Rationale.* This function is equivalent to accessing the communicator’s group with
 18 MPI_COMM_GROUP (see above), computing the rank using MPI_GROUP_RANK, and
 19 then freeing the temporary group via MPI_GROUP_FREE. However, this function is so
 20 commonly used, that this shortcut was introduced. (*End of rationale.*)

21
 22 *Advice to users.* This function gives the rank of the process in the particular commu-
 23 nicator’s group. It is useful, as noted above, in conjunction with MPI_COMM_SIZE.

24 Many programs will be written with the master-slave model, where one process (such
 25 as the rank-zero process) will play a supervisory role, and the other processes will
 26 serve as compute nodes. In this framework, the two preceding calls are useful for
 27 determining the roles of the various processes of a communicator. (*End of advice to*
 28 *users.*)

29
 30
 31 MPI_COMM_COMPARE(comm1, comm2, result)

32 IN comm1 first communicator (handle)
 33 IN comm2 second communicator (handle)
 34 OUT result result (integer)
 35
 36

37 int MPI_Comm_compare(MPI_Comm comm1, MPI_Comm comm2, int *result)

38 MPI_COMM_COMPARE(COMM1, COMM2, RESULT, IERROR)

39 INTEGER COMM1, COMM2, RESULT, IERROR

40
 41 static int MPI::Comm::Compare(const MPI::Comm& comm1,
 42 const MPI::Comm& comm2)

43
 44 MPI_IDENT results if and only if comm1 and comm2 are handles for the same object (identical
 45 groups and same contexts). MPI_CONGRUENT results if the underlying groups are identical
 46 in constituents and rank order; these communicators differ only by context. MPI_SIMILAR
 47 results if the group members of both communicators are the same but the rank order differs.
 48 MPI_UNEQUAL results otherwise.

5.4.2 Communicator Constructors

The following are collective functions that are invoked by all processes in the group associated with `comm`.

Rationale. Note that there is a chicken-and-egg aspect to MPI in that a communicator is needed to create a new communicator. The base communicator for all MPI communicators is predefined outside of MPI, and is `MPI_COMM_WORLD`. This model was arrived at after considerable debate, and was chosen to increase “safety” of programs written in MPI. (*End of rationale.*)

Intercommunicator Constructors

The current MPI interface provides only two intercommunicator construction routines:

- `MPI_INTERCOMM_CREATE`, creates an intercommunicator from two intracommunicators,
- `MPI_COMM_DUP`, duplicates an existing intercommunicator (or intracommunicator).

In MPI-1, the other communicator constructors, `MPI_COMM_CREATE` and `MPI_COMM_SPLIT`, applied only to intracommunicators. These operations in fact have well-defined semantics for intercommunicators [45].

In the following discussions, the two groups in an intercommunicator are called the *left* and *right* groups. A process in an intercommunicator is a member of either the left or the right group. From the point of view of that process, the group that the process is a member of is called the *local* group; the other group (relative to that process) is the *remote* group. The left and right group labels give us a way to describe the two groups in an intercommunicator that is not relative to any particular process (as the local and remote groups are).

`MPI_COMM_DUP(comm, newcomm)`

| | | |
|-----|----------------------|------------------------------------|
| IN | <code>comm</code> | communicator (handle) |
| OUT | <code>newcomm</code> | copy of <code>comm</code> (handle) |

`int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)`

`MPI_COMM_DUP(COMM, NEWCOMM, IERROR)`
`INTEGER COMM, NEWCOMM, IERROR`

`MPI::Intracomm MPI::Intracomm::Dup() const`

`MPI::Intercomm MPI::Intercomm::Dup() const`

`MPI::Cartcomm MPI::Cartcomm::Dup() const`

`MPI::Graphcomm MPI::Graphcomm::Dup() const`

`MPI::Comm& MPI::Comm::Clone() const = 0`

`MPI::Intracomm& MPI::Intracomm::Clone() const`

```
1 MPI::Intercomm& MPI::Intercomm::Clone() const
```

```
2 MPI::Cartcomm& MPI::Cartcomm::Clone() const
```

```
3 MPI::Graphcomm& MPI::Graphcomm::Clone() const
```

5 MPI_COMM_DUP Duplicates the existing communicator `comm` with associated key values. For each key value, the respective copy callback function determines the attribute value associated with this key in the new communicator; one particular action that a copy callback may take is to delete the attribute from the new communicator. Returns in `newcomm` a new communicator with the same group, any copied cached information, but a new context (see section 5.7.2).

12 *Advice to users.* This operation is used to provide a parallel library call with a duplicate communication space that has the same properties as the original communicator. This includes any attributes (see below), and topologies (see chapter 6). This call is valid even if there are pending point-to-point communications involving the communicator `comm`. A typical call might involve a MPI_COMM_DUP at the beginning of the parallel call, and an MPI_COMM_FREE of that duplicated communicator at the end of the call. Other models of communicator management are also possible.

20 This call applies to both intra- and inter-communicators. (*End of advice to users.*)

22 *Advice to implementors.* One need not actually copy the group information, but only add a new reference and increment the reference count. Copy on write can be used for the cached information. (*End of advice to implementors.*)

```
27 MPI_COMM_CREATE(comm, group, newcomm)
```

```
29 IN      comm      communicator (handle)
```

```
30 IN      group     Group, which is a subset of the group of
31                          comm (handle)
```

```
32 OUT     newcomm   new communicator (handle)
```

```
34 int MPI_Comm_create(MPI_Comm comm, MPI_Group group, MPI_Comm *newcomm)
```

```
36 MPI_COMM_CREATE(COMM, GROUP, NEWCOMM, IERROR)
```

```
37 INTEGER COMM, GROUP, NEWCOMM, IERROR
```

```
38 MPI::Intercomm MPI::Intercomm::Create(const MPI::Group& group) const
```

```
39 MPI::Intracomm MPI::Intracomm::Create(const MPI::Group& group) const
```

41 If `comm` is an intra-communicator, this function creates a new communicator `newcomm` with communication group defined by `group` and a new context. No cached information propagates from `comm` to `newcomm`. The function returns MPI_COMM_NULL to processes that are not in `group`. The call is erroneous if not all `group` arguments have the same value, or if `group` is not a subset of the group associated with `comm`. Note that the call is to be executed by all processes in `comm`, even if they do not belong to the new group.

```
48
```


Rationale. The requirement that the entire group of `comm` participate in the call stems from the following considerations:

- It allows the implementation to layer `MPI_COMM_CREATE` on top of regular collective communications.
- It provides additional safety, in particular in the case where partially overlapping groups are used to create new communicators.
- It permits implementations sometimes to avoid communication related to context creation.

(End of rationale.)

Advice to users. `MPI_COMM_CREATE` provides a means to subset a group of processes for the purpose of separate MIMD computation, with separate communication space. `newcomm`, which emerges from `MPI_COMM_CREATE` can be used in subsequent calls to `MPI_COMM_CREATE` (or other communicator constructors) further to subdivide a computation into parallel sub-computations. A more general service is provided by `MPI_COMM_SPLIT`, below. *(End of advice to users.)*

Advice to implementors. Since all processes calling `MPI_COMM_DUP` or `MPI_COMM_CREATE` provide the same `group` argument, it is theoretically possible to agree on a group-wide unique context with no communication. However, local execution of these functions requires use of a larger context name space and reduces error checking. Implementations may strike various compromises between these conflicting goals, such as bulk allocation of multiple contexts in one collective operation.

Important: If new communicators are created without synchronizing the processes involved then the communication system should be able to cope with messages arriving in a context that has not yet been allocated at the receiving process. *(End of advice to implementors.)*

If `comm` is an intercommunicator, then the output communicator is also an intercommunicator where the local group consists only of those processes contained in `group` (see Figure 5.1). The `group` argument should only contain those processes in the local group of the input intercommunicator that are to be a part of `newcomm`. If either `group` does not specify at least one process in the local group of the intercommunicator, or if the calling process is not included in the `group`, `MPI_COMM_NULL` is returned.

Rationale. In the case where either the left or right group is empty, a null communicator is returned instead of an intercommunicator with `MPI_GROUP_EMPTY` because the side with the empty group must return `MPI_COMM_NULL`. *(End of rationale.)*

Example 5.1 The following example illustrates how the first node in the left side of an intercommunicator could be joined with all members on the right side of an intercommunicator to form a new intercommunicator.

```
MPI_Comm inter_comm, new_inter_comm;
MPI_Group local_group, group;
int      rank = 0; /* rank on left side to include in
```

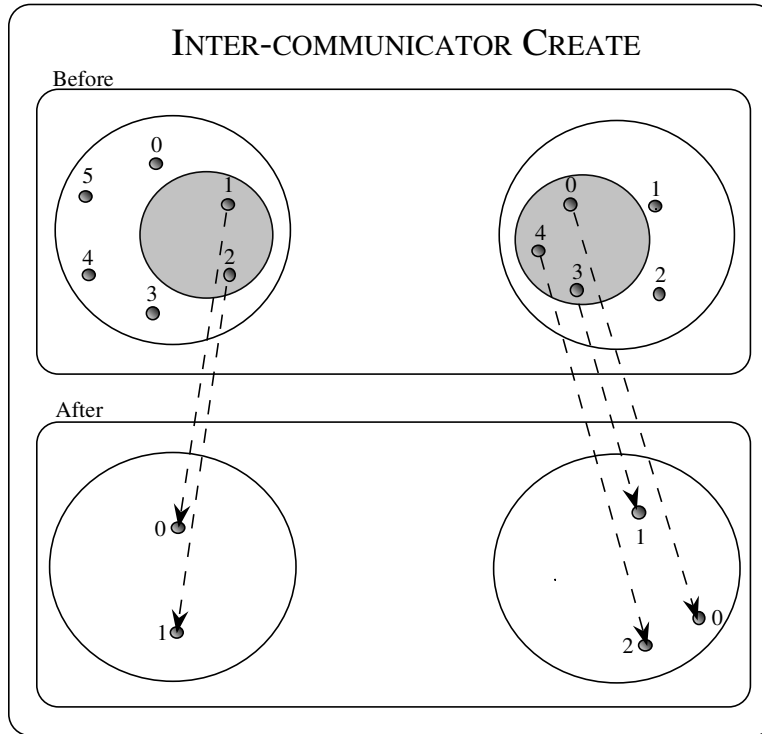


Figure 5.1: Intercommunicator create using `MPI_COMM_CREATE` extended to intercommunicators. The input groups are those in the grey circle.

```

25         new inter-comm */
26
27
28     /* Construct the original intercommunicator: "inter_comm" */
29     ...
30
31     /* Construct the group of processes to be in new
32        intercommunicator */
33     if (/* I'm on the left side of the intercommunicator */) {
34         MPI_Comm_group ( inter_comm, &local_group );
35         MPI_Group_incl ( local_group, 1, &rank, &group );
36         MPI_Group_free ( &local_group );
37     }
38     else
39         MPI_Comm_group ( inter_comm, &group );
40
41     MPI_Comm_create ( inter_comm, group, &new_inter_comm );
42     MPI_Group_free( &group );

```

| | | | |
|---|---------|--|---|
| MPI_COMM_SPLIT(comm, color, key, newcomm) | | | 1 |
| IN | comm | communicator (handle) | 2 |
| | | | 3 |
| IN | color | control of subset assignment (integer) | 4 |
| | | | 5 |
| IN | key | control of rank assignment (integer) | 6 |
| | | | 7 |
| OUT | newcomm | new communicator (handle) | 8 |

```
int MPI_Comm_split(MPI_Comm comm, int color, int key, MPI_Comm *newcomm)
```

```
MPI_COMM_SPLIT(COMM, COLOR, KEY, NEWCOMM, IERROR)
    INTEGER COMM, COLOR, KEY, NEWCOMM, IERROR
```

```
MPI::Intercomm MPI::Intercomm::Split(int color, int key) const
```

```
MPI::Intracomm MPI::Intracomm::Split(int color, int key) const
```

This function partitions the group associated with `comm` into disjoint subgroups, one for each value of `color`. Each subgroup contains all processes of the same color. Within each subgroup, the processes are ranked in the order defined by the value of the argument `key`, with ties broken according to their rank in the old group. A new communicator is created for each subgroup and returned in `newcomm`. A process may supply the color value `MPI_UNDEFINED`, in which case `newcomm` returns `MPI_COMM_NULL`. This is a collective call, but each process is permitted to provide different values for `color` and `key`.

A call to `MPI_COMM_CREATE(comm, group, newcomm)` is equivalent to a call to `MPI_COMM_SPLIT(comm, color, key, newcomm)`, where all members of `group` provide `color = 0` and `key = rank` in `group`, and all processes that are not members of `group` provide `color = MPI_UNDEFINED`. The function `MPI_COMM_SPLIT` allows more general partitioning of a group into one or more subgroups with optional reordering.

The value of `color` must be nonnegative.

Advice to users. This is an extremely powerful mechanism for dividing a single communicating group of processes into k subgroups, with k chosen implicitly by the user (by the number of colors asserted over all the processes). Each resulting communicator will be non-overlapping. Such a division could be useful for defining a hierarchy of computations, such as for multigrid, or linear algebra.

Multiple calls to `MPI_COMM_SPLIT` can be used to overcome the requirement that any call have no overlap of the resulting communicators (each process is of only one color per call). In this way, multiple overlapping communication structures can be created. Creative use of the `color` and `key` in such splitting operations is encouraged.

Note that, for a fixed `color`, the keys need not be unique. It is `MPI_COMM_SPLIT`'s responsibility to sort processes in ascending order according to this key, and to break ties in a consistent way. If all the keys are specified in the same way, then all the processes in a given color will have the relative rank order as they did in their parent group. (In general, they will have different ranks.)

Essentially, making the key value zero for all processes of a given color means that one doesn't really care about the rank-order of the processes in the new communicator. (*End of advice to users.*)

Rationale. color is restricted to be nonnegative, so as not to conflict with the value assigned to MPI_UNDEFINED. (*End of rationale.*)

The result of MPI_COMM_SPLIT on an intercommunicator is that those processes on the left with the same color as those processes on the right combine to create a new intercommunicator. The key argument describes the relative rank of processes on each side of the intercommunicator (see Figure 5.2). For those colors that are specified only on one side of the intercommunicator, MPI_COMM_NULL is returned. MPI_COMM_NULL is also returned to those processes that specify MPI_UNDEFINED as the color.

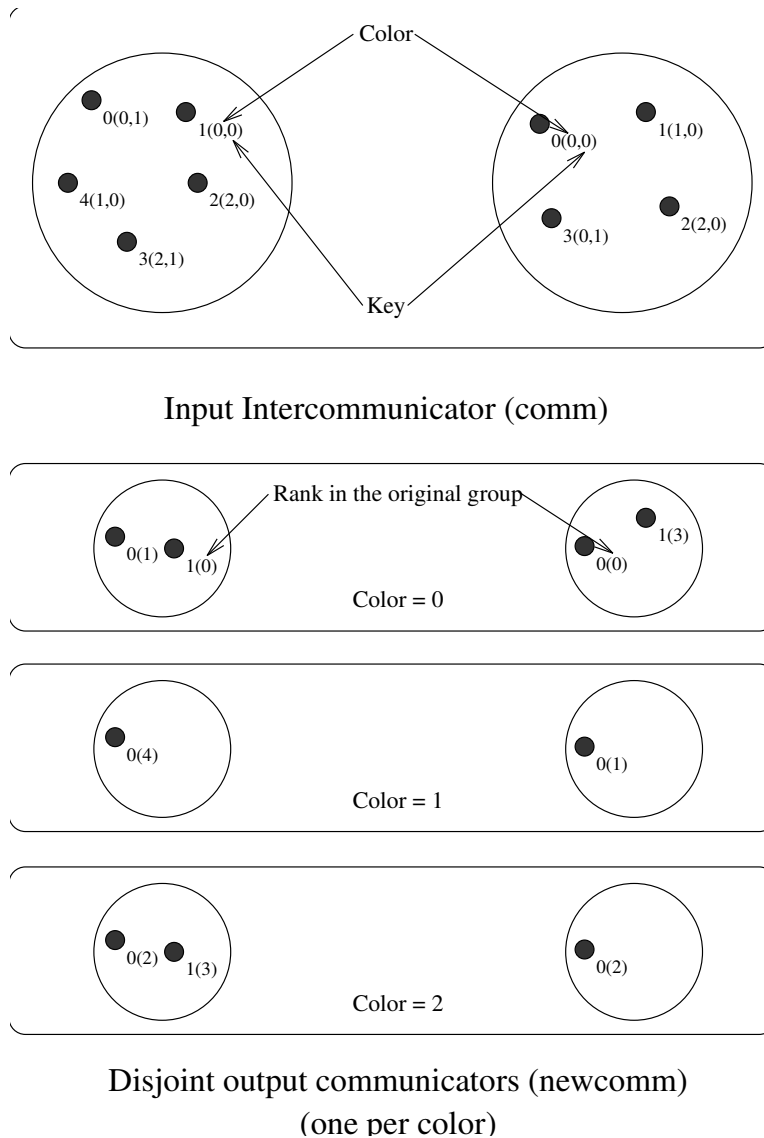


Figure 5.2: Intercommunicator construction achieved by splitting an existing intercommunicator with MPI_COMM_SPLIT extended to intercommunicators.

Example 5.2 (Parallel client-server model). The following client code illustrates how clients on the left side of an intercommunicator could be assigned to a single server from a pool of servers on the right side of an intercommunicator.

```

/* Client code */
MPI_Comm multiple_server_comm;
MPI_Comm single_server_comm;
int color, rank, num_servers;

/* Create intercommunicator with clients and servers:
   multiple_server_comm */
...

/* Find out the number of servers available */
MPI_Comm_remote_size ( multiple_server_comm, &num_servers );

/* Determine my color */
MPI_Comm_rank ( multiple_server_comm, &rank );
color = rank % num_servers;

/* Split the intercommunicator */
MPI_Comm_split ( multiple_server_comm, color, rank,
                &single_server_comm );

```

The following is the corresponding server code:

```

/* Server code */
MPI_Comm multiple_client_comm;
MPI_Comm single_server_comm;
int rank;

/* Create intercommunicator with clients and servers:
   multiple_client_comm */
...

/* Split the intercommunicator for a single server per group
   of clients */
MPI_Comm_rank ( multiple_client_comm, &rank );
MPI_Comm_split ( multiple_client_comm, rank, 0,
                &single_server_comm );

```

5.4.3 Communicator Destructors

MPI_COMM_FREE(comm)

INOUT comm communicator to be destroyed (handle)

int MPI_Comm_free(MPI_Comm *comm)

MPI_COMM_FREE(COMM, IERROR)

INTEGER COMM, IERROR

```
1 void MPI::Comm::Free()
```

2
3 This collective operation marks the communication object for deallocation. The handle
4 is set to MPI_COMM_NULL. Any pending operations that use this communicator will complete
5 normally; the object is actually deallocated only if there are no other active references to
6 it. This call applies to intra- and inter-communicators. The delete callback functions for
7 all cached attributes (see section 5.7) are called in arbitrary order.

8
9 *Advice to implementors.* A reference-count mechanism may be used: the reference
10 count is incremented by each call to MPI_COMM_DUP, and decremented by each call
11 to MPI_COMM_FREE. The object is ultimately deallocated when the count reaches
12 zero.

13 Though collective, it is anticipated that this operation will normally be implemented to
14 be local, though the debugging version of an MPI library might choose to synchronize.
15 (*End of advice to implementors.*)

17 5.5 Motivating Examples

19 5.5.1 Current Practice #1

21 Example #1a:

```
22 main(int argc, char **argv)
23 {
24     int me, size;
25     ...
26     MPI_Init ( &argc, &argv );
27     MPI_Comm_rank (MPI_COMM_WORLD, &me);
28     MPI_Comm_size (MPI_COMM_WORLD, &size);
29
30     (void)printf ("Process %d size %d\n", me, size);
31     ...
32     MPI_Finalize();
33 }
34
```

35 Example #1a is a do-nothing program that initializes itself legally, and refers to the “all”
36 communicator, and prints a message. It terminates itself legally too. This example does
37 not imply that MPI supports printf-like communication itself.

38 Example #1b (supposing that `size` is even):

```
39 main(int argc, char **argv)
40 {
41     int me, size;
42     int SOME_TAG = 0;
43     ...
44     MPI_Init(&argc, &argv);
45
46     MPI_Comm_rank(MPI_COMM_WORLD, &me); /* local */
47     MPI_Comm_size(MPI_COMM_WORLD, &size); /* local */
48
```

```

1
2   if((me % 2) == 0)
3   {
4       /* send unless highest-numbered process */
5       if((me + 1) < size)
6           MPI_Send(..., me + 1, SOME_TAG, MPI_COMM_WORLD);
7   }
8   else
9       MPI_Recv(..., me - 1, SOME_TAG, MPI_COMM_WORLD);
10
11   ...
12   MPI_Finalize();
13 }

```

Example #1b schematically illustrates message exchanges between “even” and “odd” processes in the “all” communicator.

5.5.2 Current Practice #2

```

19 main(int argc, char **argv)
20 {
21     int me, count;
22     void *data;
23     ...
24
25     MPI_Init(&argc, &argv);
26     MPI_Comm_rank(MPI_COMM_WORLD, &me);
27
28     if(me == 0)
29     {
30         /* get input, create buffer ‘‘data’’ */
31         ...
32     }
33
34     MPI_Bcast(data, count, MPI_BYTE, 0, MPI_COMM_WORLD);
35
36     ...
37
38     MPI_Finalize();
39 }
40

```

This example illustrates the use of a collective communication.

5.5.3 (Approximate) Current Practice #3

```

45 main(int argc, char **argv)
46 {
47     int me, count, count2;
48     void *send_buf, *recv_buf, *send_buf2, *recv_buf2;

```

```

1   MPI_Group MPI_GROUP_WORLD, grprem;
2   MPI_Comm commslave;
3   static int ranks[] = {0};
4   ...
5   MPI_Init(&argc, &argv);
6   MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);
7   MPI_Comm_rank(MPI_COMM_WORLD, &me); /* local */
8
9   MPI_Group_excl(MPI_GROUP_WORLD, 1, ranks, &grprem); /* local */
10  MPI_Comm_create(MPI_COMM_WORLD, grprem, &commslave);
11
12  if(me != 0)
13  {
14      /* compute on slave */
15      ...
16      MPI_Reduce(send_buf,recv_buff,count, MPI_INT, MPI_SUM, 1, commslave);
17      ...
18  }
19  /* zero falls through immediately to this reduce, others do later... */
20  MPI_Reduce(send_buf2, recv_buff2, count2,
21             MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
22
23  MPI_Comm_free(&commslave);
24  MPI_Group_free(&MPI_GROUP_WORLD);
25  MPI_Group_free(&grprem);
26  MPI_Finalize();
27  }
28

```

29 This example illustrates how a group consisting of all but the zeroth process of the “all”
30 group is created, and then how a communicator is formed (`commslave`) for that new group.
31 The new communicator is used in a collective call, and all processes execute a collective call
32 in the `MPI_COMM_WORLD` context. This example illustrates how the two communicators
33 (that inherently possess distinct contexts) protect communication. That is, communication
34 in `MPI_COMM_WORLD` is insulated from communication in `commslave`, and vice versa.

35 In summary, “group safety” is achieved via communicators because distinct contexts
36 within communicators are enforced to be unique on any process.

37 5.5.4 Example #4

38 The following example is meant to illustrate “safety” between point-to-point and collective
39 communication. MPI guarantees that a single communicator can do safe point-to-point and
40 collective communication.
41
42

```

43     #define TAG_ARBITRARY 12345
44     #define SOME_COUNT    50
45
46     main(int argc, char **argv)
47     {
48         int me;

```



```

MPI_Request request[2];
MPI_Status status[2];
MPI_Group MPI_GROUP_WORLD, subgroup;
int ranks[] = {2, 4, 6, 8};
MPI_Comm the_comm;
...
MPI_Init(&argc, &argv);
MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);

MPI_Group_incl(MPI_GROUP_WORLD, 4, ranks, &subgroup); /* local */
MPI_Group_rank(subgroup, &me); /* local */

MPI_Comm_create(MPI_COMM_WORLD, subgroup, &the_comm);

if(me != MPI_UNDEFINED)
{
    MPI_Irecv(buff1, count, MPI_DOUBLE, MPI_ANY_SOURCE, TAG_ARBITRARY,
              the_comm, request);
    MPI_Isend(buff2, count, MPI_DOUBLE, (me+1)%4, TAG_ARBITRARY,
              the_comm, request+1);
}

for(i = 0; i < SOME_COUNT, i++)
    MPI_Reduce(..., the_comm);
MPI_Waitall(2, request, status);

MPI_Comm_free(&the_comm);
MPI_Group_free(&MPI_GROUP_WORLD);
MPI_Group_free(&subgroup);
MPI_Finalize();
}

```

5.5.5 Library Example #1

The main program:

```

main(int argc, char **argv)
{
    int done = 0;
    user_lib_t *libh_a, *libh_b;
    void *dataset1, *dataset2;
    ...
    MPI_Init(&argc, &argv);
    ...
    init_user_lib(MPI_COMM_WORLD, &libh_a);
    init_user_lib(MPI_COMM_WORLD, &libh_b);
    ...
    user_start_op(libh_a, dataset1);
}

```

```

1   user_start_op(libh_b, dataset2);
2   ...
3   while(!done)
4   {
5       /* work */
6       ...
7       MPI_Reduce(..., MPI_COMM_WORLD);
8       ...
9       /* see if done */
10      ...
11  }
12  user_end_op(libh_a);
13  user_end_op(libh_b);
14
15  uninit_user_lib(libh_a);
16  uninit_user_lib(libh_b);
17  MPI_Finalize();
18  }

```

The user library initialization code:

```

21  void init_user_lib(MPI_Comm comm, user_lib_t **handle)
22  {
23      user_lib_t *save;
24
25      user_lib_initsave(&save); /* local */
26      MPI_Comm_dup(comm, &(save -> comm));
27
28      /* other inits */
29      ...
30
31      *handle = save;
32  }

```

User start-up code:

```

35  void user_start_op(user_lib_t *handle, void *data)
36  {
37      MPI_Irecv( ..., handle->comm, &(handle -> irecv_handle) );
38      MPI_Isend( ..., handle->comm, &(handle -> isend_handle) );
39  }

```

User communication clean-up code:

```

43  void user_end_op(user_lib_t *handle)
44  {
45      MPI_Status *status;
46      MPI_Wait(handle -> isend_handle, status);
47      MPI_Wait(handle -> irecv_handle, status);
48  }

```

User object clean-up code:

```
void uninit_user_lib(user_lib_t *handle)
{
    MPI_Comm_free(&(handle -> comm));
    free(handle);
}
```

5.5.6 Library Example #2

The main program:

```
main(int argc, char **argv)
{
    int ma, mb;
    MPI_Group MPI_GROUP_WORLD, group_a, group_b;
    MPI_Comm comm_a, comm_b;

    static int list_a[] = {0, 1};
    #if defined(EXAMPLE_2B) | defined(EXAMPLE_2C)
        static int list_b[] = {0, 2 ,3};
    #else/* EXAMPLE_2A */
        static int list_b[] = {0, 2};
    #endif

    int size_list_a = sizeof(list_a)/sizeof(int);
    int size_list_b = sizeof(list_b)/sizeof(int);

    ...
    MPI_Init(&argc, &argv);
    MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);

    MPI_Group_incl(MPI_GROUP_WORLD, size_list_a, list_a, &group_a);
    MPI_Group_incl(MPI_GROUP_WORLD, size_list_b, list_b, &group_b);

    MPI_Comm_create(MPI_COMM_WORLD, group_a, &comm_a);
    MPI_Comm_create(MPI_COMM_WORLD, group_b, &comm_b);

    if(comm_a != MPI_COMM_NULL)
        MPI_Comm_rank(comm_a, &ma);
    if(comm_b != MPI_COMM_NULL)
        MPI_Comm_rank(comm_b, &mb);

    if(comm_a != MPI_COMM_NULL)
        lib_call(comm_a);

    if(comm_b != MPI_COMM_NULL)
    {
        lib_call(comm_b);
        lib_call(comm_b);
    }
}
```

```

1     }
2
3     if(comm_a != MPI_COMM_NULL)
4         MPI_Comm_free(&comm_a);
5     if(comm_b != MPI_COMM_NULL)
6         MPI_Comm_free(&comm_b);
7     MPI_Group_free(&group_a);
8     MPI_Group_free(&group_b);
9     MPI_Group_free(&MPI_GROUP_WORLD);
10    MPI_Finalize();
11    }

```

The library:

```

14    void lib_call(MPI_Comm comm)
15    {
16        int me, done = 0;
17        MPI_Comm_rank(comm, &me);
18        if(me == 0)
19            while(!done)
20            {
21                MPI_Recv(..., MPI_ANY_SOURCE, MPI_ANY_TAG, comm);
22                ...
23            }
24        else
25        {
26            /* work */
27            MPI_Send(..., 0, ARBITRARY_TAG, comm);
28            ....
29        }
30    #ifdef EXAMPLE_2C
31        /* include (resp, exclude) for safety (resp, no safety): */
32        MPI_Barrier(comm);
33    #endif
34    }
35

```

The above example is really three examples, depending on whether or not one includes rank 3 in `list_b`, and whether or not a synchronize is included in `lib_call`. This example illustrates that, despite contexts, subsequent calls to `lib_call` with the same context need not be safe from one another (colloquially, “back-masking”). Safety is realized if the `MPI_Barrier` is added. What this demonstrates is that libraries have to be written carefully, even with contexts. When rank 3 is excluded, then the synchronize is not needed to get safety from back masking.

Algorithms like “reduce” and “allreduce” have strong enough source selectivity properties so that they are inherently okay (no backmasking), provided that MPI provides basic guarantees. So are multiple calls to a typical tree-broadcast algorithm with the same root or different roots (see [48]). Here we rely on two guarantees of MPI: pairwise ordering of messages between processes in the same context, and source selectivity — deleting either feature removes the guarantee that backmasking cannot be required.

Algorithms that try to do non-deterministic broadcasts or other calls that include wild-card operations will not generally have the good properties of the deterministic implementations of “reduce,” “allreduce,” and “broadcast.” Such algorithms would have to utilize the monotonically increasing tags (within a communicator scope) to keep things straight.

All of the foregoing is a supposition of “collective calls” implemented with point-to-point operations. MPI implementations may or may not implement collective calls using point-to-point operations. These algorithms are used to illustrate the issues of correctness and safety, independent of how MPI implements its collective calls. See also section 5.8.

5.6 Inter-Communication

This section introduces the concept of inter-communication and describes the portions of MPI that support it. It describes support for writing programs that contain user-level servers.

All point-to-point communication described thus far has involved communication between processes that are members of the same group. This type of communication is called “intra-communication” and the communicator used is called an “intra-communicator,” as we have noted earlier in the chapter.

In modular and multi-disciplinary applications, different process groups execute distinct modules and processes within different modules communicate with one another in a pipeline or a more general module graph. In these applications, the most natural way for a process to specify a target process is by the rank of the target process within the target group. In applications that contain internal user-level servers, each server may be a process group that provides services to one or more clients, and each client may be a process group that uses the services of one or more servers. It is again most natural to specify the target process by rank within the target group in these applications. This type of communication is called “inter-communication” and the communicator used is called an “inter-communicator,” as introduced earlier.

An inter-communication is a point-to-point communication between processes in different groups. The group containing a process that initiates an inter-communication operation is called the “local group,” that is, the sender in a send and the receiver in a receive. The group containing the target process is called the “remote group,” that is, the receiver in a send and the sender in a receive. As in intra-communication, the target process is specified using a (communicator, rank) pair. Unlike intra-communication, the rank is relative to a second, remote group.

All inter-communicator constructors are blocking and require that the local and remote groups be disjoint.

Advice to users. The groups must be disjoint for several reasons. Primarily, this is the intent of the intercommunicators — to provide a communicator for communication between disjoint groups. This is reflected in the definition of `MPI_INTERCOMM_MERGE`, which allows the user to control the ranking of the processes in the created intracommunicator; this ranking makes little sense if the groups are not disjoint. In addition, the natural extension of collective operations to intercommunicators makes the most sense when the groups are disjoint. (*End of advice to users.*)

Here is a summary of the properties of inter-communication and inter-communicators:

- 1 • The syntax of point-to-point **and collective** communication is the same for both inter-
2 and intra-communication. The same communicator can be used both for send and for
3 receive operations.
- 4 • A target process is addressed by its rank in the remote group, both for sends and for
5 receives.
- 6 • Communications using an inter-communicator are guaranteed not to conflict with any
7 communications that use a different communicator.
- 8 • A communicator will provide either intra- or inter-communication, never both.

11 The routine `MPI_COMM_TEST_INTER` may be used to determine if a communicator is an
12 inter- or intra-communicator. Inter-communicators can be used as arguments to some of the
13 other communicator access routines. Inter-communicators cannot be used as input to some
14 of the constructor routines for intra-communicators (for instance, `MPI_COMM_CREATE`).

16 *Advice to implementors.* For the purpose of point-to-point communication, commu-
17 nicators can be represented in each process by a tuple consisting of:

18 **group**
19 **send_context**
20 **receive_context**
21 **source**

24 For inter-communicators, **group** describes the remote group, and **source** is the rank of
25 the process in the local group. For intra-communicators, **group** is the communicator
26 group (remote=local), **source** is the rank of the process in this group, and **send**
27 **context** and **receive context** are identical. A group is represented by a rank-to-
28 absolute-address translation table.

29 The inter-communicator cannot be discussed sensibly without considering processes in
30 both the local and remote groups. Imagine a process **P** in group \mathcal{P} , which has an inter-
31 communicator $C_{\mathcal{P}}$, and a process **Q** in group \mathcal{Q} , which has an inter-communicator
32 $C_{\mathcal{Q}}$. Then

- 34 • $C_{\mathcal{P}}.\mathbf{group}$ describes the group \mathcal{Q} and $C_{\mathcal{Q}}.\mathbf{group}$ describes the group \mathcal{P} .
- 35 • $C_{\mathcal{P}}.\mathbf{send_context} = C_{\mathcal{Q}}.\mathbf{receive_context}$ and the context is unique in \mathcal{Q} ;
36 $C_{\mathcal{P}}.\mathbf{receive_context} = C_{\mathcal{Q}}.\mathbf{send_context}$ and this context is unique in \mathcal{P} .
- 37 • $C_{\mathcal{P}}.\mathbf{source}$ is rank of **P** in \mathcal{P} and $C_{\mathcal{Q}}.\mathbf{source}$ is rank of **Q** in \mathcal{Q} .

39 Assume that **P** sends a message to **Q** using the inter-communicator. Then **P** uses
40 the **group** table to find the absolute address of **Q**; **source** and **send_context** are
41 appended to the message.

42 Assume that **Q** posts a receive with an explicit source argument using the inter-
43 communicator. Then **Q** matches **receive_context** to the message context and source
44 argument to the message source.

45 The same algorithm is appropriate for intra-communicators as well.

46 In order to support inter-communicator accessors and constructors, it is necessary to
47 supplement this model with additional structures, that store information about the
48


```

1 MPI_COMM_REMOTE_SIZE(COMM, SIZE, IERROR)
2     INTEGER COMM, SIZE, IERROR
3
4 int MPI::Intercomm::Get_remote_size() const
5
6
7 MPI_COMM_REMOTE_GROUP(comm, group)
8     IN      comm          inter-communicator (handle)
9     OUT     group        remote group corresponding to comm (handle)
10
11
12 int MPI_Comm_remote_group(MPI_Comm comm, MPI_Group *group)
13
14 MPI_COMM_REMOTE_GROUP(COMM, GROUP, IERROR)
15     INTEGER COMM, GROUP, IERROR
16
17 MPI::Group MPI::Intercomm::Get_remote_group() const

```

Rationale. Symmetric access to both the local and remote groups of an inter-communicator is important, so this function, as well as `MPI_COMM_REMOTE_SIZE` have been provided. (*End of rationale.*)

5.6.2 Inter-communicator Operations

This section introduces four blocking inter-communicator operations.

`MPI_INTERCOMM_CREATE` is used to bind two intra-communicators into an inter-communicator; the function `MPI_INTERCOMM_MERGE` creates an intra-communicator by merging the local and remote groups of an inter-communicator. The functions `MPI_COMM_DUP` and `MPI_COMM_FREE`, introduced previously, duplicate and free an inter-communicator, respectively.

Overlap of local and remote groups that are bound into an inter-communicator is prohibited. If there is overlap, then the program is erroneous and is likely to deadlock. (If a process is multithreaded, and MPI calls block only a thread, rather than a process, then “dual membership” can be supported. It is then the user’s responsibility to make sure that calls on behalf of the two “roles” of a process are executed by two independent threads.)

The function `MPI_INTERCOMM_CREATE` can be used to create an inter-communicator from two existing intra-communicators, in the following situation: At least one selected member from each group (the “group leader”) has the ability to communicate with the selected member from the other group; that is, a “peer” communicator exists to which both leaders belong, and each leader knows the rank of the other leader in this peer communicator. Furthermore, members of each group know the rank of their leader.

Construction of an inter-communicator from two intra-communicators requires separate collective operations in the local group and in the remote group, as well as a point-to-point communication between a process in the local group and a process in the remote group.

In standard MPI implementations (with static process allocation at initialization), the `MPI_COMM_WORLD` communicator (or preferably a dedicated duplicate thereof) can be this peer communicator. [For applications that have used spawn or join, it may be necessary to first create an intracommunicator to be used as peer.](#)

The application topology functions described in chapter 6 do not apply to inter-communicators. Users that require this capability should utilize `MPI_INTERCOMM_MERGE` to build an intra-communicator, then apply the graph or cartesian topology capabilities to that intra-communicator, creating an appropriate topology-oriented intra-communicator. Alternatively, it may be reasonable to devise one’s own application topology mechanisms for this case, without loss of generality.

```
MPI_INTERCOMM_CREATE(local_comm, local_leader, peer_comm, remote_leader, tag,
newintercomm)
```

| | | |
|-----|----------------------------|---|
| IN | <code>local_comm</code> | local intra-communicator (handle) |
| IN | <code>local_leader</code> | rank of local group leader in <code>local_comm</code> (integer) |
| IN | <code>peer_comm</code> | “peer” communicator; significant only at the <code>local_leader</code> (handle) |
| IN | <code>remote_leader</code> | rank of remote group leader in <code>peer_comm</code> ; significant only at the <code>local_leader</code> (integer) |
| IN | <code>tag</code> | “safe” tag (integer) |
| OUT | <code>newintercomm</code> | new inter-communicator (handle) |

```
int MPI_Intercomm_create(MPI_Comm local_comm, int local_leader,
MPI_Comm peer_comm, int remote_leader, int tag,
MPI_Comm *newintercomm)
```

```
MPI_INTERCOMM_CREATE(LOCAL_COMM, LOCAL_LEADER, PEER_COMM, REMOTE_LEADER, TAG,
NEWINTERCOMM, IERROR)
INTEGER LOCAL_COMM, LOCAL_LEADER, PEER_COMM, REMOTE_LEADER, TAG,
NEWINTERCOMM, IERROR
```

```
MPI::Intercomm MPI::Intracomm::Create_intercomm(int local_leader, const
MPI::Comm& peer_comm, int remote_leader, int tag) const
```

This call creates an inter-communicator. It is collective over the union of the local and remote groups. Processes should provide identical `local_comm` and `local_leader` arguments within each group. Wildcards are not permitted for `remote_leader`, `local_leader`, and `tag`.

This call uses point-to-point communication with communicator `peer_comm`, and with tag `tag` between the leaders. Thus, care must be taken that there be no pending communication on `peer_comm` that could interfere with this communication.

Advice to users. We recommend using a dedicated peer communicator, such as a duplicate of `MPI_COMM_WORLD`, to avoid trouble with peer communicators. (*End of advice to users.*)

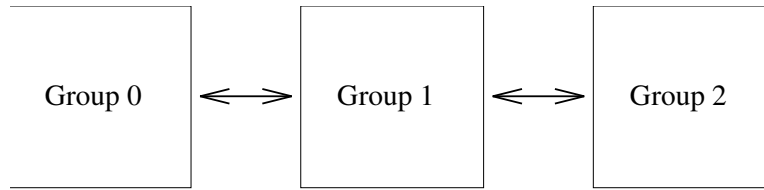


Figure 5.3: Three-group pipeline.

```

MPI_INTERCOMM_MERGE(intercomm, high, newintracomm)

```

| | | |
|-----|--------------|---------------------------------|
| IN | intercomm | Inter-Communicator (handle) |
| IN | high | (logical) |
| OUT | newintracomm | new intra-communicator (handle) |

```

int MPI_Intercomm_merge(MPI_Comm intercomm, int high,
                        MPI_Comm *newintracomm)

```

```

MPI_INTERCOMM_MERGE(INTERCOMM, HIGH, INTRACOMM, IERROR)
INTEGER INTERCOMM, INTRACOMM, IERROR
LOGICAL HIGH

```

```

MPI::Intracomm MPI::Intercomm::Merge(bool high) const

```

This function creates an intra-communicator from the union of the two groups that are associated with `intercomm`. All processes should provide the same `high` value within each of the two groups. If processes in one group provided the value `high = false` and processes in the other group provided the value `high = true` then the union orders the “low” group before the “high” group. If all processes provided the same `high` argument then the order of the union is arbitrary. This call is blocking and collective within the union of the two groups.

The error handler on the new intercommunicator in each process is inherited from the communicator that contributes the local group. Note that this can result in different processes in the same communicator having different error handlers.

Advice to implementors. The implementation of `MPI_INTERCOMM_MERGE`, `MPI_COMM_FREE` and `MPI_COMM_DUP` are similar to the implementation of `MPI_INTERCOMM_CREATE`, except that contexts private to the input inter-communicator are used for communication between group leaders rather than contexts inside a bridge communicator. (*End of advice to implementors.*)

5.6.3 Inter-Communication Examples

Example 1: Three-Group “Pipeline”

Groups 0 and 1 communicate. Groups 1 and 2 communicate. Therefore, group 0 requires one inter-communicator, group 1 requires two inter-communicators, and group 2 requires 1 inter-communicator.

```

main(int argc, char **argv)
{

```

```
MPI_Comm  myComm;          /* intra-communicator of local sub-group */  1
MPI_Comm  myFirstComm;     /* inter-communicator */  2
MPI_Comm  mySecondComm;   /* second inter-communicator (group 1 only) */  3
int  membershipKey;  4
int  rank;  5

MPI_Init(&argc, &argv);  6
MPI_Comm_rank(MPI_COMM_WORLD, &rank);  7

/* User code must generate membershipKey in the range [0, 1, 2] */  8
membershipKey = rank % 3;  9

/* Build intra-communicator for local sub-group */  10
MPI_Comm_split(MPI_COMM_WORLD, membershipKey, rank, &myComm);  11

/* Build inter-communicators. Tags are hard-coded. */  12
if (membershipKey == 0)  13
{  14
    /* Group 0 communicates with group 1. */  15
    MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,  16
                        1, &myFirstComm);  17
}  18
else if (membershipKey == 1)  19
{  20
    /* Group 1 communicates with groups 0 and 2. */  21
    MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 0,  22
                        1, &myFirstComm);  23
    MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 2,  24
                        12, &mySecondComm);  25
}  26
else if (membershipKey == 2)  27
{  28
    /* Group 2 communicates with group 1. */  29
    MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,  30
                        12, &myFirstComm);  31
}  32

/* Do work ... */  33

switch(membershipKey) /* free communicators appropriately */  34
{  35
case 1:  36
    MPI_Comm_free(&mySecondComm);  37
case 0:  38
case 2:  39
    MPI_Comm_free(&myFirstComm);  40
    break;  41
}  42

MPI_Finalize();  43
}  44
```

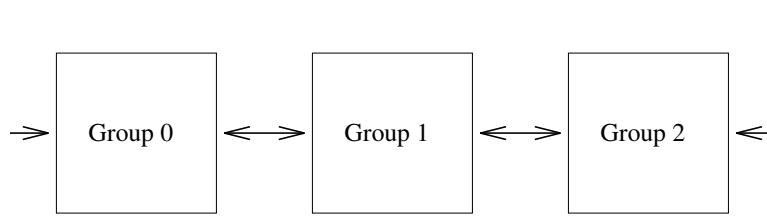


Figure 5.4: Three-group ring.

Example 2: Three-Group “Ring”

Groups 0 and 1 communicate. Groups 1 and 2 communicate. Groups 0 and 2 communicate. Therefore, each requires two inter-communicators.

```

14  main(int argc, char **argv)
15  {
16      MPI_Comm  myComm;      /* intra-communicator of local sub-group */
17      MPI_Comm  myFirstComm; /* inter-communicators */
18      MPI_Comm  mySecondComm;
19      MPI_Status status;
20      int membershipKey;
21      int rank;
22
23
24      MPI_Init(&argc, &argv);
25      MPI_Comm_rank(MPI_COMM_WORLD, &rank);
26      ...
27
28      /* User code must generate membershipKey in the range [0, 1, 2] */
29      membershipKey = rank % 3;
30
31      /* Build intra-communicator for local sub-group */
32      MPI_Comm_split(MPI_COMM_WORLD, membershipKey, rank, &myComm);
33
34      /* Build inter-communicators. Tags are hard-coded. */
35      if (membershipKey == 0)
36      {
37          /* Group 0 communicates with groups 1 and 2. */
38          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,
39                               1, &myFirstComm);
40          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 2,
41                               2, &mySecondComm);
42      }
43      else if (membershipKey == 1)
44      {
45          /* Group 1 communicates with groups 0 and 2. */
46          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 0,
47                               1, &myFirstComm);
48          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 2,
49                               12, &mySecondComm);
50      }
51  }

```

```

else if (membershipKey == 2)
{
    /* Group 2 communicates with groups 0 and 1. */
    MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 0,
                        2, &myFirstComm);
    MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,
                        12, &mySecondComm);
}

/* Do some work ... */

/* Then free communicators before terminating... */
MPI_Comm_free(&myFirstComm);
MPI_Comm_free(&mySecondComm);
MPI_Comm_free(&myComm);
MPI_Finalize();
}

```

Example 3: Building Name Service for Intercommunication

The following procedures exemplify the process by which a user could create name service for building intercommunicators via a rendezvous involving a server communicator, and a tag name selected by both groups.

After all MPI processes execute `MPI_INIT`, every process calls the example function, `Init_server()`, defined below. Then, if the `new_world` returned is `NULL`, the process getting `NULL` is required to implement a server function, in a reactive loop, `Do_server()`. Everyone else just does their prescribed computation, using `new_world` as the new effective “global” communicator. One designated process calls `Undo_Server()` to get rid of the server when it is not needed any longer.

Features of this approach include:

- Support for multiple name servers
- Ability to scope the name servers to specific processes
- Ability to make such servers come and go as desired.

```

#define INIT_SERVER_TAG_1 666
#define UNDO_SERVER_TAG_1 777

static int server_key_val;

/* for attribute management for server_comm, copy callback: */
void handle_copy_fn(MPI_Comm *oldcomm, int *keyval, void *extra_state,
void *attribute_val_in, void **attribute_val_out, int *flag)
{
    /* copy the handle */
    *attribute_val_out = attribute_val_in;
    *flag = 1; /* indicate that copy to happen */
}

```

```
1
2 int Init_server(peer_comm, rank_of_server, server_comm, new_world)
3 MPI_Comm peer_comm;
4 int rank_of_server;
5 MPI_Comm *server_comm;
6 MPI_Comm *new_world; /* new effective world, sans server */
7 {
8     MPI_Comm temp_comm, lone_comm;
9     MPI_Group peer_group, temp_group;
10    int rank_in_peer_comm, size, color, key = 0;
11    int peer_leader, peer_leader_rank_in_temp_comm;
12
13    MPI_Comm_rank(peer_comm, &rank_in_peer_comm);
14    MPI_Comm_size(peer_comm, &size);
15
16    if ((size < 2) || (0 > rank_of_server) || (rank_of_server >= size))
17        return (MPI_ERR_OTHER);
18
19    /* create two communicators, by splitting peer_comm
20       into the server process, and everyone else */
21
22    peer_leader = (rank_of_server + 1) % size; /* arbitrary choice */
23
24    if ((color = (rank_in_peer_comm == rank_of_server)))
25    {
26        MPI_Comm_split(peer_comm, color, key, &lone_comm);
27
28        MPI_Intercomm_create(lone_comm, 0, peer_comm, peer_leader,
29                            INIT_SERVER_TAG_1, server_comm);
30
31        MPI_Comm_free(&lone_comm);
32        *new_world = MPI_COMM_NULL;
33    }
34    else
35    {
36        MPI_Comm_Split(peer_comm, color, key, &temp_comm);
37
38        MPI_Comm_group(peer_comm, &peer_group);
39        MPI_Comm_group(temp_comm, &temp_group);
40        MPI_Group_translate_ranks(peer_group, 1, &peer_leader,
41 temp_group, &peer_leader_rank_in_temp_comm);
42
43        MPI_Intercomm_create(temp_comm, peer_leader_rank_in_temp_comm,
44                            peer_comm, rank_of_server,
45                            INIT_SERVER_TAG_1, server_comm);
46
47        /* attach new_world communication attribute to server_comm: */
48
```



```

1      {
2          while (de_queue(queue, MPI_ANY_TAG, &pairs_rank_in_new_world,
3                      &pairs_rank_in_server))
4              ;
5
6          MPI_Intercomm_free(&server_comm);
7          break;
8      }
9
10     if (de_queue(queue, client_tag, &pairs_rank_in_new_world,
11                &pairs_rank_in_server))
12     {
13         /* matched pair with same tag, tell them
14            about each other! */
15         buffer[0] = pairs_rank_in_new_world;
16         MPI_Send(buffer, 1, MPI_INT, client_src, client_tag,
17                 server_comm);
18
19         buffer[0] = client_rank_in_new_world;
20         MPI_Send(buffer, 1, MPI_INT, pairs_rank_in_server, client_tag,
21                 server_comm);
22     }
23     else
24         en_queue(queue, client_tag, client_source,
25                 client_rank_in_new_world);
26
27 }
28 }
29

```

A particular process would be responsible for ending the server when it is no longer needed. Its call to `Undo_server` would terminate server function.

```

32 int Undo_server(server_comm)    /* example client that ends server */
33 MPI_Comm *server_comm;
34 {
35     int buffer = 0;
36     MPI_Send(&buffer, 1, MPI_INT, 0, UNDO_SERVER_TAG_1, *server_comm);
37     MPI_Intercomm_free(server_comm);
38 }
39

```

The following is a blocking name-service for inter-communication, with same semantic restrictions as `MPI_Intercomm_create`, but simplified syntax. It uses the functionality just defined to create the name service.

```

44 int Intercomm_name_create(local_comm, server_comm, tag, comm)
45 MPI_Comm local_comm, server_comm;
46 int tag;
47 MPI_Comm *comm;
48 {

```



```

int error;
int found; /* attribute acquisition mgmt for new_world */
           /* comm in server_comm */
void *val;

MPI_Comm new_world;

int buffer[10], rank;
int local_leader = 0;

MPI_Attr_get(server_comm, server_keyval, &val, &found);
new_world = (MPI_Comm)val; /* retrieve cached handle */

MPI_Comm_rank(server_comm, &rank); /* rank in local group */

if (rank == local_leader)
{
    buffer[0] = rank;
    MPI_Send(&buffer, 1, MPI_INT, 0, tag, server_comm);
    MPI_Recv(&buffer, 1, MPI_INT, 0, tag, server_comm);
}

error = MPI_Intercomm_create(local_comm, local_leader, new_world,
                             buffer[0], tag, comm);

return(error);
}

```

5.7 Caching

MPI provides a “caching” facility that allows an application to attach arbitrary pieces of information, called **attributes**, to communicators. More precisely, the caching facility allows a portable library to do the following:

- pass information between calls by associating it with an MPI intra- or inter-communicator,
- quickly retrieve that information, and
- be guaranteed that out-of-date information is never retrieved, even if the communicator is freed and its handle subsequently reused by MPI.

The caching capabilities, in some form, are required by built-in MPI routines such as collective communication and application topology. Defining an interface to these capabilities as part of the MPI standard is valuable because it permits routines like collective communication and application topologies to be implemented as portable code, and also because it makes MPI more extensible by allowing user-written routines to use standard MPI calling sequences.

1 *Advice to users.* The communicator `MPI_COMM_SELF` is a suitable choice for posting
2 process-local attributes, via this attributing-caching mechanism. (*End of advice to*
3 *users.*)

6 5.7.1 New Attribute Caching Functions

7 Caching on communicators has been a very useful feature. In MPI-2 it is expanded to
8 include caching on windows and datatypes.

10 *Rationale.* In one extreme you can allow caching on all opaque handles. The other
11 extreme is to only allow it on communicators. Caching has a cost associated with it
12 and should only be allowed when it is clearly needed and the increased cost is modest.
13 This is the reason that windows and datatypes were added but not other handles.
14 (*End of rationale.*)

16 One difficulty in MPI-1 is the potential for size differences between Fortran integers and
17 C pointers. To overcome this problem with attribute caching on communicators, new func-
18 tions are also given for this case. The new functions to cache on datatypes and windows also
19 address this issue. For a general discussion of the address size problem, see Section 13.3.6.

20 *Advice to implementors.* High quality implementations should raise an error when
21 a keyval that was created by a call to `MPI_XXX_CREATE_KEYVAL` is used with an
22 object of the wrong type with a call to `MPI_YYY_GET_ATTR`, `MPI_YYY_SET_ATTR`,
23 `MPI_YYY_DELETE_ATTR`, or `MPI_YYY_FREE_KEYVAL`. To do so, it is necessary to
24 maintain, with each keyval, information on the type of the associated user function.
25 (*End of advice to implementors.*)

28 5.7.2 Functionality

29 Attributes are attached to communicators. Attributes are local to the process and specific
30 to the communicator to which they are attached. Attributes are not propagated by MPI
31 from one communicator to another except when the communicator is duplicated using
32 `MPI_COMM_DUP` (and even then the application must give specific permission through
33 callback functions for the attribute to be copied).

35 *Advice to users.* Attributes in C are of type `void *`. Typically, such an attribute will
36 be a pointer to a structure that contains further information, or a handle to an MPI
37 object. In Fortran, attributes are of type `INTEGER`. Such attribute can be a handle to
38 an MPI object, or just an integer-valued attribute. (*End of advice to users.*)

39 *Advice to implementors.* Attributes are scalar values, equal in size to, or larger than
40 a C-language pointer. Attributes can always hold an MPI handle. (*End of advice to*
41 *implementors.*)

43 The caching interface defined here represents that attributes be stored by MPI opaquely
44 within a communicator. Accessor functions include the following:

- 46 • obtain a key value (used to identify an attribute); the user specifies “callback” func-
47 tions by which MPI informs the application when the communicator is destroyed or
48 copied.

1 Generates a new attribute key. Keys are locally unique in a process, and opaque to
 2 user, though they are explicitly stored in integers. Once allocated, the key value can be
 3 used to associate attributes and access them on any locally defined communicator.

4 This function replaces `MPI_KEYVAL_CREATE`, whose use is deprecated. The C binding
 5 is identical. The Fortran binding differs in that `extra_state` is an address-sized integer.
 6 Also, the copy and delete callback functions have Fortran bindings that are consistent with
 7 address-sized attributes.

8 The C callback functions are:

```
9 typedef int MPI_Comm_copy_attr_function(MPI_Comm oldcomm, int comm_keyval,  
10     void *extra_state, void *attribute_val_in,  
11     void *attribute_val_out, int *flag);
```

12 and

```
13 typedef int MPI_Comm_delete_attr_function(MPI_Comm comm, int comm_keyval,  
14     void *attribute_val, void *extra_state);
```

16 which are the same as the MPI-1.1 calls but with a new name. The old names are deprecated.

17 The Fortran callback functions are:

```
18 SUBROUTINE COMM_COPY_ATTR_FN(OLDCOMM, COMM_KEYVAL, EXTRA_STATE,  
19     ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)  
20     INTEGER OLDCOMM, COMM_KEYVAL, IERROR  
21     INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,  
22     ATTRIBUTE_VAL_OUT  
23     LOGICAL FLAG
```

24 and

```
25 SUBROUTINE COMM_DELETE_ATTR_FN(COMM, COMM_KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE,  
26     IERROR)  
27     INTEGER COMM, COMM_KEYVAL, IERROR  
28     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

30 The C++ callbacks are:

```
31 typedef int MPI::Comm::Copy_attr_function(const MPI::Comm& oldcomm,  
32     int comm_keyval, void* extra_state, void* attribute_val_in,  
33     void* attribute_val_out, bool& flag);
```

34 and

```
35 typedef int MPI::Comm::Delete_attr_function(MPI::Comm& comm,  
36     int comm_keyval, void* attribute_val, void* extra_state);
```

38 The `comm_copy_attr_fn` function is invoked when a communicator is duplicated by
 39 `MPI_COMM_DUP`. `comm_copy_attr_fn` should be of type `MPI_Comm_copy_attr_function`. The
 40 copy callback function is invoked for each key value in `oldcomm` in arbitrary order. Each call
 41 to the copy callback is made with a key value and its corresponding attribute. If it returns
 42 `flag = 0`, then the attribute is deleted in the duplicated communicator. Otherwise (`flag = 1`),
 43 the new attribute value is set to the value returned in `attribute_val_out`. The function returns
 44 `MPI_SUCCESS` on success and an error code on failure (in which case `MPI_COMM_DUP` will
 45 fail).

46 The argument `comm_copy_attr_fn` may be specified as `MPI_COMM_NULL_COPY_FN` or
 47 `MPI_COMM_DUP_FN` from either C, C++, or Fortran. `MPI_COMM_NULL_COPY_FN` is a
 48

function that does nothing other than returning `flag = 0` and `MPI_SUCCESS`.
`MPI_COMM_DUP_FN` is a simple-minded copy function that sets `flag = 1`, returns the value of `attribute_val_in` in `attribute_val_out`, and returns `MPI_SUCCESS`. These replace the MPI-1 predefined callbacks `MPI_NULL_COPY_FN` and `MPI_DUP_FN`, whose use is deprecated.

Advice to users. Even though both formal arguments `attribute_val_in` and `attribute_val_out` are of type `void *`, their usage differs. The C copy function is passed by MPI in `attribute_val_in` the *value* of the attribute, and in `attribute_val_out` the *address* of the attribute, so as to allow the function to return the (new) attribute value. The use of type `void *` for both is to avoid messy type casts.

A valid copy function is one that completely duplicates the information by making a full duplicate copy of the data structures implied by an attribute; another might just make another reference to that data structure, while using a reference-count mechanism. Other types of attributes might not copy at all (they might be specific to oldcomm only). (*End of advice to users.*)

Advice to implementors. A C interface should be assumed for copy and delete functions associated with key values created in C; a Fortran calling interface should be assumed for key values created in Fortran. (*End of advice to implementors.*)

Analogous to `comm_copy_attr_fn` is a callback deletion function, defined as follows. The `comm_delete_attr_fn` function is invoked when a communicator is deleted by `MPI_COMM_FREE` or when a call is made explicitly to `MPI_ATTR_DELETE`. `comm_delete_attr_fn` should be of type `MPI_Comm_delete_attr_function`.

This function is called by `MPI_COMM_FREE`, `MPI_COMM_DELETE_ATTR`, and `MPI_COMM_SET_ATTR` to do whatever is needed to remove an attribute. The function returns `MPI_SUCCESS` on success and an error code on failure (in which case `MPI_COMM_FREE` will fail).

The argument `comm_delete_attr_fn` may be specified as `MPI_COMM_NULL_DELETE_FN` from either C, C++, or Fortran. `MPI_COMM_NULL_DELETE_FN` is a function that does nothing, other than returning `MPI_SUCCESS`. `MPI_COMM_NULL_DELETE_FN` replaces `MPI_NULL_DELETE_FN`, whose use is deprecated.

If an attribute copy function or attribute delete function returns other than `MPI_SUCCESS`, then the call that caused it to be invoked (for example, `MPI_COMM_FREE`), is erroneous.

The special key value `MPI_KEYVAL_INVALID` is never returned by `MPI_KEYVAL_CREATE`. Therefore, it can be used for static initialization of key values.

`MPI_COMM_FREE_KEYVAL(comm_keyval)`

INOUT `comm_keyval` key value (integer)

`int MPI_Comm_free_keyval(int *comm_keyval)`

`MPI_COMM_FREE_KEYVAL(COMM_KEYVAL, IERROR)`
 `INTEGER COMM_KEYVAL, IERROR`

`static void MPI::Comm::Free_keyval(int& comm_keyval)`

1 Frees an extant attribute key. This function sets the value of `keyval` to
 2 `MPI_KEYVAL_INVALID`. Note that it is not erroneous to free an attribute key that is in use,
 3 because the actual free does not transpire until after all references (in other communicators
 4 on the process) to the key have been freed. These references need to be explicitly freed by the
 5 program, either via calls to `MPI_COMM_DELETE_ATTR` that free one attribute instance,
 6 or by calls to `MPI_COMM_FREE` that free all attribute instances associated with the freed
 7 communicator.

8 This call is identical to the MPI-1 call `MPI_KEYVAL_FREE` but is needed to match the
 9 new communicator-specific creation function. The use of `MPI_KEYVAL_FREE` is deprecated.

10
 11
 12
 13 `MPI_COMM_SET_ATTR(comm, comm_keyval, attribute_val)`

| | | | |
|----|-------|----------------------------|--|
| 14 | INOUT | <code>comm</code> | communicator from which attribute will be attached (handle) |
| 15 | | | |
| 16 | IN | <code>comm_keyval</code> | key value (integer) |
| 17 | | | |
| 18 | IN | <code>attribute_val</code> | attribute value |

19
 20 `int MPI_Comm_set_attr(MPI_Comm comm, int comm_keyval, void *attribute_val)`

21 `MPI_COMM_SET_ATTR(COMM, COMM_KEYVAL, ATTRIBUTE_VAL, IERROR)`
 22 `INTEGER COMM, COMM_KEYVAL, IERROR`
 23 `INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL`

24
 25 `void MPI::Comm::Set_attr(int comm_keyval, const void* attribute_val) const`

26
 27 This function stores the stipulated attribute value `attribute_val` for subsequent retrieval
 28 by `MPI_COMM_GET_ATTR`. If the value is already present, then the outcome is as if
 29 `MPI_COMM_DELETE_ATTR` was first called to delete the previous value (and the callback
 30 function `comm_delete_attr_fn` was executed), and a new value was next stored. The call
 31 is erroneous if there is no key with value `keyval`; in particular `MPI_KEYVAL_INVALID` is an
 32 erroneous key value. The call will fail if the `comm_delete_attr_fn` function returned an error
 33 code other than `MPI_SUCCESS`.

34 This function replaces `MPI_ATTR_PUT`, whose use is deprecated. The C binding is
 35 identical. The Fortran binding differs in that `attribute_val` is an address-sized integer.

36
 37 `MPI_COMM_GET_ATTR(comm, comm_keyval, attribute_val, flag)`

| | | | |
|----|-----|----------------------------|---|
| 38 | IN | <code>comm</code> | communicator to which the attribute is attached (han- dle) |
| 39 | | | |
| 40 | | | |
| 41 | IN | <code>comm_keyval</code> | key value (integer) |
| 42 | OUT | <code>attribute_val</code> | attribute value, unless <code>flag = false</code> |
| 43 | OUT | <code>flag</code> | false if no attribute is associated with the key (logical) |

44
 45 `int MPI_Comm_get_attr(MPI_Comm comm, int comm_keyval, void *attribute_val,`
 46 `int *flag)`

47
 48 `MPI_COMM_GET_ATTR(COMM, COMM_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)`

```

INTEGER COMM, COMM_KEYVAL, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
LOGICAL FLAG

```

```
bool MPI::Comm::Get_attr(int comm_keyval, void* attribute_val) const
```

Retrieves attribute value by key. The call is erroneous if there is no key with value `keyval`. On the other hand, the call is correct if the key value exists, but no attribute is attached on `comm` for that key; in such case, the call returns `flag = false`. In particular `MPI_KEYVAL_INVALID` is an erroneous key value.

Advice to users. The call to `MPI_Comm_set_attr` passes in `attribute_val` the *value* of the attribute; the call to `MPI_Comm_get_attr` passes in `attribute_val` the *address* of the the location where the attribute value is to be returned. Thus, if the attribute value itself is a pointer of type `void*`, the actual `attribute_val` parameter to `MPI_Comm_set_attr` will be of type `void*` and the actual `attribute_val` parameter to `MPI_Comm_get_attr` will be of type `void**`. (*End of advice to users.*)

Rationale. The use of a formal parameter `attribute_val` or type `void*` (rather than `void**`) avoids the messy type casting that would be needed if the attribute value is declared with a type other than `void*`. (*End of rationale.*)

This function replaces `MPI_ATTR_GET`, whose use is deprecated. The C binding is identical. The Fortran binding differs in that `attribute_val` is an address-sized integer.

```
MPI_COMM_DELETE_ATTR(comm, comm_keyval)
```

| | | |
|-------|-------------|---|
| INOUT | comm | communicator from which the attribute is deleted (handle) |
| IN | comm_keyval | key value (integer) |

```
int MPI_Comm_delete_attr(MPI_Comm comm, int comm_keyval)
```

```

MPI_COMM_DELETE_ATTR(COMM, COMM_KEYVAL, IERROR)
INTEGER COMM, COMM_KEYVAL, IERROR

```

```
void MPI::Comm::Delete_attr(int comm_keyval)
```

Delete attribute from cache by key. This function invokes the attribute delete function `comm_delete_attr_fn` specified when the `keyval` was created. The call will fail if the `comm_delete_attr_fn` function returns an error code other than `MPI_SUCCESS`.

Whenever a communicator is replicated using the function `MPI_COMM_DUP`, all callback copy functions for attributes that are currently set are invoked (in arbitrary order). Whenever a communicator is deleted using the function `MPI_COMM_FREE` all callback delete functions for attributes that are currently set are invoked.

This function is the same as `MPI_ATTR_DELETE` but is needed to match the new communicator specific functions. The use of `MPI_ATTR_DELETE` is deprecated.

5.7.4 Windows

The new functions for caching on windows are:

```

1 MPI_WIN_CREATE_KEYVAL(win_copy_attr_fn, win_delete_attr_fn, win_keyval, extra_state)
2     IN      win_copy_attr_fn      copy callback function for win_keyval (function)
3
4     IN      win_delete_attr_fn    delete callback function for win_keyval (function)
5
6     OUT     win_keyval            key value for future access (integer)
7
8     IN      extra_state           extra state for callback functions

```

```

9 int MPI_Win_create_keyval(MPI_Win_copy_attr_function *win_copy_attr_fn,
10                          MPI_Win_delete_attr_function *win_delete_attr_fn,
11                          int *win_keyval, void *extra_state)

```

```

12 MPI_WIN_CREATE_KEYVAL(WIN_COPY_ATTR_FN, WIN_DELETE_ATTR_FN, WIN_KEYVAL,
13                      EXTRA_STATE, IERROR)
14     EXTERNAL WIN_COPY_ATTR_FN, WIN_DELETE_ATTR_FN
15     INTEGER WIN_KEYVAL, IERROR
16     INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE

```

```

17
18 static int MPI::Win::Create_keyval(MPI::Win::Copy_attr_function*
19                                   win_copy_attr_fn,
20                                   MPI::Win::Delete_attr_function* win_delete_attr_fn,
21                                   void* extra_state)

```

The argument `win_copy_attr_fn` may be specified as `MPI_WIN_NULL_COPY_FN` or `MPI_WIN_DUP_FN` from either C, C++, or Fortran. `MPI_WIN_NULL_COPY_FN` is a function that does nothing other than returning `flag = 0` and `MPI_SUCCESS`. `MPI_WIN_DUP_FN` is a simple-minded copy function that sets `flag = 1`, returns the value of `attribute_val_in` in `attribute_val_out`, and returns `MPI_SUCCESS`.

The argument `win_delete_attr_fn` may be specified as `MPI_WIN_NULL_DELETE_FN` from either C, C++, or Fortran. `MPI_WIN_NULL_DELETE_FN` is a function that does nothing, other than returning `MPI_SUCCESS`.

The C callback functions are:

```

31 typedef int MPI_Win_copy_attr_function(MPI_Win oldwin, int win_keyval,
32                                       void *extra_state, void *attribute_val_in,
33                                       void *attribute_val_out, int *flag);

```

and

```

36 typedef int MPI_Win_delete_attr_function(MPI_Win win, int win_keyval,
37                                         void *attribute_val, void *extra_state);

```

The Fortran callback functions are:

```

39 SUBROUTINE WIN_COPY_ATTR_FN(OLDWIN, WIN_KEYVAL, EXTRA_STATE,
40                            ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
41     INTEGER OLDWIN, WIN_KEYVAL, IERROR
42     INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,
43         ATTRIBUTE_VAL_OUT
44     LOGICAL FLAG

```

and

```

47 SUBROUTINE WIN_DELETE_ATTR_FN(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE,
48                               IERROR)

```



```

INTEGER WIN, WIN_KEYVAL, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE

```

The C++ callbacks are:

```

typedef int MPI::Win::Copy_attr_function(const MPI::Win& oldwin,
    int win_keyval, void* extra_state, void* attribute_val_in,
    void* attribute_val_out, bool& flag);

```

and

```

typedef int MPI::Win::Delete_attr_function(MPI::Win& win, int win_keyval,
    void* attribute_val, void* extra_state);

```

If an attribute copy function or attribute delete function returns other than MPI_SUCCESS, then the call that caused it to be invoked (for example, MPI_WIN_FREE), is erroneous.

```

MPI_WIN_FREE_KEYVAL(win_keyval)

```

```

    INOUT    win_keyval                key value (integer)

```

```

int MPI_Win_free_keyval(int *win_keyval)

```

```

MPI_WIN_FREE_KEYVAL(WIN_KEYVAL, IERROR)

```

```

    INTEGER WIN_KEYVAL, IERROR

```

```

static void MPI::Win::Free_keyval(int& win_keyval)

```

```

MPI_WIN_SET_ATTR(win, win_keyval, attribute_val)

```

```

    INOUT    win                        window to which attribute will be attached (handle)

```

```

    IN      win_keyval                  key value (integer)

```

```

    IN      attribute_val               attribute value

```

```

int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val)

```

```

MPI_WIN_SET_ATTR(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, IERROR)

```

```

    INTEGER WIN, WIN_KEYVAL, IERROR

```

```

    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL

```

```

void MPI::Win::Set_attr(int win_keyval, const void* attribute_val)

```

```

1 MPI_WIN_GET_ATTR(win, win_keyval, attribute_val, flag)
2     IN        win                window to which the attribute is attached (handle)
3
4     IN        win_keyval         key value (integer)
5
6     OUT       attribute_val      attribute value, unless flag = false
7
8     OUT       flag               false if no attribute is associated with the key (logical)
9
10
11 int MPI_Win_get_attr(MPI_Win win, int win_keyval, void *attribute_val,
12                     int *flag)
13
14 MPI_WIN_GET_ATTR(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)
15     INTEGER WIN, WIN_KEYVAL, IERROR
16     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
17     LOGICAL FLAG
18
19 bool MPI::Win::Get_attr(int win_keyval, void* attribute_val) const
20
21 MPI_WIN_DELETE_ATTR(win, win_keyval)
22     INOUT     win                window from which the attribute is deleted (handle)
23
24     IN        win_keyval         key value (integer)
25
26
27 int MPI_Win_delete_attr(MPI_Win win, int win_keyval)
28
29 MPI_WIN_DELETE_ATTR(WIN, WIN_KEYVAL, IERROR)
30     INTEGER WIN, WIN_KEYVAL, IERROR
31
32 void MPI::Win::Delete_attr(int win_keyval)
33
34
35 5.7.5 Datatypes
36
37 The new functions for caching on datatypes are:
38
39
40 MPI_TYPE_CREATE_KEYVAL(type_copy_attr_fn, type_delete_attr_fn, type_keyval, extra_state)
41
42     IN        type_copy_attr_fn   copy callback function for type_keyval (function)
43
44     IN        type_delete_attr_fn delete callback function for type_keyval (function)
45
46     OUT       type_keyval         key value for future access (integer)
47
48     IN        extra_state         extra state for callback functions
49
50
51 int MPI_Type_create_keyval(MPI_Type_copy_attr_function *type_copy_attr_fn,
52                           MPI_Type_delete_attr_function *type_delete_attr_fn,
53                           int *type_keyval, void *extra_state)
54
55 MPI_TYPE_CREATE_KEYVAL(TYPE_COPY_ATTR_FN, TYPE_DELETE_ATTR_FN, TYPE_KEYVAL,
56                       EXTRA_STATE, IERROR)

```

```

EXTERNAL TYPE_COPY_ATTR_FN, TYPE_DELETE_ATTR_FN
INTEGER TYPE_KEYVAL, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE

static int MPI::Datatype::Create_keyval(MPI::Datatype::Copy_attr_function*
    type_copy_attr_fn, MPI::Datatype::Delete_attr_function*
    type_delete_attr_fn, void* extra_state)

```

The argument `type_copy_attr_fn` may be specified as `MPI_TYPE_NULL_COPY_FN` or `MPI_TYPE_DUP_FN` from either C, C++, or Fortran. `MPI_TYPE_NULL_COPY_FN` is a function that does nothing other than returning `flag = 0` and `MPI_SUCCESS`. `MPI_TYPE_DUP_FN` is a simple-minded copy function that sets `flag = 1`, returns the value of `attribute_val_in` in `attribute_val_out`, and returns `MPI_SUCCESS`.

The argument `type_delete_attr_fn` may be specified as `MPI_TYPE_NULL_DELETE_FN` from either C, C++, or Fortran. `MPI_TYPE_NULL_DELETE_FN` is a function that does nothing, other than returning `MPI_SUCCESS`.

The C callback functions are:

```

typedef int MPI_Type_copy_attr_function(MPI_Datatype oldtype,
    int type_keyval, void *extra_state, void *attribute_val_in,
    void *attribute_val_out, int *flag);

and

typedef int MPI_Type_delete_attr_function(MPI_Datatype type, int type_keyval,
    void *attribute_val, void *extra_state);

```

The Fortran callback functions are:

```

SUBROUTINE TYPE_COPY_ATTR_FN(OLDTYPE, TYPE_KEYVAL, EXTRA_STATE,
    ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
INTEGER OLDTYPE, TYPE_KEYVAL, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE,
    ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT
LOGICAL FLAG

and

SUBROUTINE TYPE_DELETE_ATTR_FN(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE,
    IERROR)
INTEGER TYPE, TYPE_KEYVAL, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE

```

The C++ callbacks are:

```

typedef int MPI::Datatype::Copy_attr_function(const MPI::Datatype& oldtype,
    int type_keyval, void* extra_state,
    const void* attribute_val_in, void* attribute_val_out,
    bool& flag);

and

typedef int MPI::Datatype::Delete_attr_function(MPI::Datatype& type,
    int type_keyval, void* attribute_val, void* extra_state);

```

If an attribute copy function or attribute delete function returns other than `MPI_SUCCESS`, then the call that caused it to be invoked (for example, `MPI_TYPE_FREE`),

```

1  is erroneous.
2
3
4  MPI_TYPE_FREE_KEYVAL(type_keyval)
5      INOUT    type_keyval                key value (integer)
6
7
8  int MPI_Type_free_keyval(int *type_keyval)
9
10 MPI_TYPE_FREE_KEYVAL(TYPE_KEYVAL, IERROR)
11     INTEGER TYPE_KEYVAL, IERROR
12
13
14 static void MPI::Datatype::Free_keyval(int& type_keyval)
15
16 MPI_TYPE_SET_ATTR(type, type_keyval, attribute_val)
17     INOUT    type                        datatype to which attribute will be attached (handle)
18     IN       type_keyval                key value (integer)
19     IN       attribute_val              attribute value
20
21 int MPI_Type_set_attr(MPI_Datatype type, int type_keyval,
22     void *attribute_val)
23
24 MPI_TYPE_SET_ATTR(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, IERROR)
25     INTEGER TYPE, TYPE_KEYVAL, IERROR
26     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
27
28 void MPI::Datatype::Set_attr(int type_keyval, const void* attribute_val)
29
30 MPI_TYPE_GET_ATTR(type, type_keyval, attribute_val, flag)
31     IN       type                        datatype to which the attribute is attached (handle)
32     IN       type_keyval                key value (integer)
33     OUT      attribute_val              attribute value, unless flag = false
34     OUT      flag                       false if no attribute is associated with the key (logical)
35
36
37 int MPI_Type_get_attr(MPI_Datatype type, int type_keyval, void
38     *attribute_val, int *flag)
39
40 MPI_TYPE_GET_ATTR(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)
41     INTEGER TYPE, TYPE_KEYVAL, IERROR
42     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
43     LOGICAL FLAG
44
45 bool MPI::Datatype::Get_attr(int type_keyval, void* attribute_val) const
46
47
48

```

```

MPI_TYPE_DELETE_ATTR(type, type_keyval) 1
    INOUT  type  datatype from which the attribute is deleted (handle) 2
    IN     type_keyval  key value (integer) 3
                                         4
                                         5
int MPI_Type_delete_attr(MPI_Datatype type, int type_keyval) 6
MPI_TYPE_DELETE_ATTR(TYPE, TYPE_KEYVAL, IERROR) 7
    INTEGER TYPE, TYPE_KEYVAL, IERROR 8
                                         9
void MPI::Datatype::Delete_attr(int type_keyval) 10
                                         11

```

5.7.6 Error Class for Invalid Keyval

Key values for attributes are system-allocated, by `MPI_{TYPE,COMM,WIN}_CREATE_KEYVAL`. Only such values can be passed to the functions that use key values as input arguments. In order to signal that an erroneous key value has been passed to one of these functions, there is a new MPI error class: `MPI_ERR_KEYVAL`. It can be returned by `MPI_ATTR_PUT`, `MPI_ATTR_GET`, `MPI_ATTR_DELETE`, `MPI_KEYVAL_FREE`, `MPI_{TYPE,COMM,WIN}_DELETE_ATTR`, `MPI_{TYPE,COMM,WIN}_SET_ATTR`, `MPI_{TYPE,COMM,WIN}_GET_ATTR`, `MPI_{TYPE,COMM,WIN}_FREE_KEYVAL`, `MPI_COMM_DUP`, `MPI_COMM_DISCONNECT`, and `MPI_COMM_FREE`. The last three are included because `keyval` is an argument to the copy and delete functions for attributes.

5.7.7 Attributes Example

Advice to users. This example shows how to write a collective communication operation that uses caching to be more efficient after the first call. The coding style assumes that MPI function results return only error statuses. (*End of advice to users.*)

```

/* key for this module's stuff: */ 30
static int gop_key = MPI_KEYVAL_INVALID; 31
                                         32
typedef struct 33
{ 34
    int ref_count;  /* reference count */ 35
    /* other stuff, whatever else we want */ 36
} gop_stuff_type; 37
                                         38
Efficient_Collective_Op (comm, ...) 39
MPI_Comm comm; 40
{ 41
    gop_stuff_type *gop_stuff; 42
    MPI_Group group; 43
    int foundflag; 44
                                         45
    MPI_Comm_group(comm, &group); 46
                                         47

```

48

```

1   if (gop_key == MPI_KEYVAL_INVALID) /* get a key on first call ever */
2   {
3       if ( ! MPI_Comm_create_keyval( gop_stuff_copier,
4                                     gop_stuff_destructor,
5                                     &gop_key, (void *)0));
6       /* get the key while assigning its copy and delete callback
7        behavior. */
8
9       MPI_Abort (comm, 99);
10      }
11
12      MPI_Comm_get_attr (comm, gop_key, &gop_stuff, &foundflag);
13      if (foundflag)
14      { /* This module has executed in this group before.
15         We will use the cached information */
16      }
17      else
18      { /* This is a group that we have not yet cached anything in.
19         We will now do so.
20         */
21
22         /* First, allocate storage for the stuff we want,
23          and initialize the reference count */
24
25         gop_stuff = (gop_stuff_type *) malloc (sizeof(gop_stuff_type));
26         if (gop_stuff == NULL) { /* abort on out-of-memory error */ }
27
28         gop_stuff -> ref_count = 1;
29
30         /* Second, fill in *gop_stuff with whatever we want.
31          This part isn't shown here */
32
33         /* Third, store gop_stuff as the attribute value */
34         MPI_Comm_set_attr ( comm, gop_key, gop_stuff);
35      }
36      /* Then, in any case, use contents of *gop_stuff
37       to do the global op ... */
38  }
39
40  /* The following routine is called by MPI when a group is freed */
41
42  gop_stuff_destructor (comm, keyval, gop_stuff, extra)
43  MPI_Comm comm;
44  int keyval;
45  gop_stuff_type *gop_stuff;
46  void *extra;
47  {
48      if (keyval != gop_key) { /* abort -- programming error */ }

```

```

1
2  /* The group's being freed removes one reference to gop_stuff */
3  gop_stuff -> ref_count -= 1;
4
5  /* If no references remain, then free the storage */
6  if (gop_stuff -> ref_count == 0) {
7      free((void *)gop_stuff);
8  }
9
10 }
11
12 /* The following routine is called by MPI when a group is copied */
13 gop_stuff_copier (comm, keyval, extra, gop_stuff_in, gop_stuff_out, flag)
14 MPI_Comm comm;
15 int keyval;
16 gop_stuff_type *gop_stuff_in, *gop_stuff_out;
17 void *extra;
18 {
19     if (keyval != gop_key) { /* abort -- programming error */ }
20
21     /* The new group adds one reference to this gop_stuff */
22     gop_stuff -> ref_count += 1;
23     gop_stuff_out = gop_stuff_in;
24 }

```

5.8 Formalizing the Loosely Synchronous Model

In this section, we make further statements about the loosely synchronous model, with particular attention to intra-communication.

5.8.1 Basic Statements

When a caller passes a communicator (that contains a context and group) to a callee, that communicator must be free of side effects throughout execution of the subprogram: there should be no active operations on that communicator that might involve the process. This provides one model in which libraries can be written, and work “safely.” For libraries so designated, the callee has permission to do whatever communication it likes with the communicator, and under the above guarantee knows that no other communications will interfere. Since we permit good implementations to create new communicators without synchronization (such as by preallocated contexts on communicators), this does not impose a significant overhead.

This form of safety is analogous to other common computer-science usages, such as passing a descriptor of an array to a library routine. The library routine has every right to expect such a descriptor to be valid and modifiable.

5.8.2 Models of Execution

In the loosely synchronous model, transfer of control to a **parallel procedure** is effected by having each executing process invoke the procedure. The invocation is a collective operation:

1 it is executed by all processes in the execution group, and invocations are similarly ordered
2 at all processes. However, the invocation need not be synchronized.

3 We say that a parallel procedure is *active* in a process if the process belongs to a group
4 that may collectively execute the procedure, and some member of that group is currently
5 executing the procedure code. If a parallel procedure is active in a process, then this process
6 may be receiving messages pertaining to this procedure, even if it does not currently execute
7 the code of this procedure.
8

9 Static communicator allocation

10 This covers the case where, at any point in time, at most one invocation of a parallel
11 procedure can be active at any process, and the group of executing processes is fixed. For
12 example, all invocations of parallel procedures involve all processes, processes are single-
13 threaded, and there are no recursive invocations.
14

15 In such a case, a communicator can be statically allocated to each procedure. The
16 static allocation can be done in a preamble, as part of initialization code. If the parallel
17 procedures can be organized into libraries, so that only one procedure of each library can
18 be concurrently active in each processor, then it is sufficient to allocate one communicator
19 per library.
20

21 Dynamic communicator allocation

22 Calls of parallel procedures are well-nested if a new parallel procedure is always invoked in
23 a subset of a group executing the same parallel procedure. Thus, processes that execute
24 the same parallel procedure have the same execution stack.
25

26 In such a case, a new communicator needs to be dynamically allocated for each new
27 invocation of a parallel procedure. The allocation is done by the caller. A new communicator
28 can be generated by a call to `MPI_COMM_DUP`, if the callee execution group is identical to
29 the caller execution group, or by a call to `MPI_COMM_SPLIT` if the caller execution group
30 is split into several subgroups executing distinct parallel routines. The new communicator
31 is passed as an argument to the invoked routine.

32 The need for generating a new communicator at each invocation can be alleviated or
33 avoided altogether in some cases: If the execution group is not split, then one can allocate
34 a stack of communicators in a preamble, and next manage the stack in a way that mimics
35 the stack of recursive calls.

36 One can also take advantage of the well-ordering property of communication to avoid
37 confusing caller and callee communication, even if both use the same communicator. To do
38 so, one needs to abide by the following two rules:

- 39 • messages sent before a procedure call (or before a return from the procedure) are also
40 received before the matching call (or return) at the receiving end;
- 41 • messages are always selected by source (no use is made of `MPI_ANY_SOURCE`).
42

43 The General case

44 In the general case, there may be multiple concurrently active invocations of the same
45 parallel procedure within the same group; invocations may not be well-nested. A new
46 communicator needs to be created for each invocation. It is the user's responsibility to make
47
48

sure that, should two distinct parallel procedures be invoked concurrently on overlapping sets of processes, then communicator creation be properly coordinated.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 6

Process Topologies

6.1 Introduction

This chapter discusses the MPI topology mechanism. A topology is an extra, optional attribute that one can give to an intra-communicator; topologies cannot be added to inter-communicators. A topology can provide a convenient naming mechanism for the processes of a group (within a communicator), and additionally, may assist the runtime system in mapping the processes onto hardware.

As stated in chapter 5, a process group in MPI is a collection of n processes. Each process in the group is assigned a rank between 0 and $n-1$. In many parallel applications a linear ranking of processes does not adequately reflect the logical communication pattern of the processes (which is usually determined by the underlying problem geometry and the numerical algorithm used). Often the processes are arranged in topological patterns such as two- or three-dimensional grids. More generally, the logical process arrangement is described by a graph. In this chapter we will refer to this logical process arrangement as the “virtual topology.”

A clear distinction must be made between the virtual process topology and the topology of the underlying, physical hardware. The virtual topology can be exploited by the system in the assignment of processes to physical processors, if this helps to improve the communication performance on a given machine. How this mapping is done, however, is outside the scope of MPI. The description of the virtual topology, on the other hand, depends only on the application, and is machine-independent. The functions that are proposed in this chapter deal only with machine-independent mapping.

Rationale. Though physical mapping is not discussed, the existence of the virtual topology information may be used as advice by the runtime system. There are well-known techniques for mapping grid/torus structures to hardware topologies such as hypercubes or grids. For more complicated graph structures good heuristics often yield nearly optimal results [34]. On the other hand, if there is no way for the user to specify the logical process arrangement as a “virtual topology,” a random mapping is most likely to result. On some machines, this will lead to unnecessary contention in the interconnection network. Some details about predicted and measured performance improvements that result from good process-to-processor mapping on modern wormhole-routing architectures can be found in [11, 10].

Besides possible performance benefits, the virtual topology can function as a convenient, process-naming structure, with tremendous benefits for program readability

1 and notational power in message-passing programming. (*End of rationale.*)

2 3 6.2 Virtual Topologies 4

5 The communication pattern of a set of processes can be represented by a graph. The
6 nodes stand for the processes, and the edges connect processes that communicate with each
7 other. MPI provides message-passing between any pair of processes in a group. There
8 is no requirement for opening a channel explicitly. Therefore, a “missing link” in the
9 user-defined process graph does not prevent the corresponding processes from exchanging
10 messages. It means rather that this connection is neglected in the virtual topology. This
11 strategy implies that the topology gives no convenient way of naming this pathway of
12 communication. Another possible consequence is that an automatic mapping tool (if one
13 exists for the runtime environment) will not take account of this edge when mapping. Edges
14 in the communication graph are not weighted, so that processes are either simply connected
15 or not connected at all.
16

17 *Rationale.* Experience with similar techniques in PARMACS [5, 9] show that this
18 information is usually sufficient for a good mapping. Additionally, a more precise
19 specification is more difficult for the user to set up, and it would make the interface
20 functions substantially more complicated. (*End of rationale.*)
21

22 Specifying the virtual topology in terms of a graph is sufficient for all applications.
23 However, in many applications the graph structure is regular, and the detailed set-up of the
24 graph would be inconvenient for the user and might be less efficient at run time. A large frac-
25 tion of all parallel applications use process topologies like rings, two- or higher-dimensional
26 grids, or tori. These structures are completely defined by the number of dimensions and
27 the numbers of processes in each coordinate direction. Also, the mapping of grids and tori
28 is generally an easier problem than that of general graphs. Thus, it is desirable to address
29 these cases explicitly.

30 Process coordinates in a cartesian structure begin their numbering at 0. Row-major
31 numbering is always used for the processes in a cartesian structure. This means that, for
32 example, the relation between group rank and coordinates for four processes in a (2×2)
33 grid is as follows.
34

35 coord (0,0): rank 0
36 coord (0,1): rank 1
37 coord (1,0): rank 2
38 coord (1,1): rank 3
39

40 6.3 Embedding in MPI 41

42 The support for virtual topologies as defined in this chapter is consistent with other parts of
43 MPI, and, whenever possible, makes use of functions that are defined elsewhere. Topology
44 information is associated with communicators. It is added to communicators using the
45 caching mechanism described in Chapter 5.
46
47
48

6.4 Overview of the Functions

The functions `MPI_GRAPH_CREATE` and `MPI_CART_CREATE` are used to create general (graph) virtual topologies and cartesian topologies, respectively. These topology creation functions are collective. As with other collective calls, the program must be written to work correctly, whether the call synchronizes or not.

The topology creation functions take as input an existing communicator `comm_old`, which defines the set of processes on which the topology is to be mapped. All input arguments must have identical values on all processes of the group of `comm_old`. A new communicator `comm_topol` is created that carries the topological structure as cached information (see Chapter 5). In analogy to function `MPI_COMM_CREATE`, no cached information propagates from `comm_old` to `comm_topol`.

`MPI_CART_CREATE` can be used to describe cartesian structures of arbitrary dimension. For each coordinate direction one specifies whether the process structure is periodic or not. Note that an n -dimensional hypercube is an n -dimensional torus with 2 processes per coordinate direction. Thus, special support for hypercube structures is not necessary. The local auxiliary function `MPI_DIMS_CREATE` can be used to compute a balanced distribution of processes among a given number of dimensions.

Rationale. Similar functions are contained in EXPRESS [12] and PARMACS. (*End of rationale.*)

The function `MPI_TOPO_TEST` can be used to inquire about the topology associated with a communicator. The topological information can be extracted from the communicator using the functions `MPI_GRAPHDIMS_GET` and `MPI_GRAPH_GET`, for general graphs, and `MPI_CARTDIM_GET` and `MPI_CART_GET`, for cartesian topologies. Several additional functions are provided to manipulate cartesian topologies: the functions `MPI_CART_RANK` and `MPI_CART_COORDS` translate cartesian coordinates into a group rank, and vice-versa; the function `MPI_CART_SUB` can be used to extract a cartesian subspace (analogous to `MPI_COMM_SPLIT`). The function `MPI_CART_SHIFT` provides the information needed to communicate with neighbors in a cartesian dimension. The two functions `MPI_GRAPH_NEIGHBORS_COUNT` and `MPI_GRAPH_NEIGHBORS` can be used to extract the neighbors of a node in a graph. The function `MPI_CART_SUB` is collective over the input communicator's group; all other functions are local.

Two additional functions, `MPI_GRAPH_MAP` and `MPI_CART_MAP` are presented in the last section. In general these functions are not called by the user directly. However, together with the communicator manipulation functions presented in Chapter 5, they are sufficient to implement all other topology functions. Section 6.5.7 outlines such an implementation.

6.5 Topology Constructors

6.5.1 Cartesian Constructor

`MPI_CART_CREATE(comm_old, ndims, dims, periods, reorder, comm_cart)`

| | | |
|-----|------------------------|--|
| IN | <code>comm_old</code> | input communicator (handle) |
| IN | <code>ndims</code> | number of dimensions of cartesian grid (integer) |
| IN | <code>dims</code> | integer array of size <code>ndims</code> specifying the number of processes in each dimension |
| IN | <code>periods</code> | logical array of size <code>ndims</code> specifying whether the grid is periodic (<code>true</code>) or not (<code>false</code>) in each dimension |
| IN | <code>reorder</code> | ranking may be reordered (<code>true</code>) or not (<code>false</code>) (logical) |
| OUT | <code>comm_cart</code> | communicator with new cartesian topology (handle) |

```
int MPI_Cart_create(MPI_Comm comm_old, int ndims, int *dims, int *periods,
                  int reorder, MPI_Comm *comm_cart)
```

```
MPI_CART_CREATE(COMM_OLD, NDIMS, DIMS, PERIODS, REORDER, COMM_CART, IERROR)
                INTEGER COMM_OLD, NDIMS, DIMS(*), COMM_CART, IERROR
                LOGICAL PERIODS(*), REORDER
```

```
MPI::Cartcomm MPI::Intracomm::Create_cart(int ndims, const int dims[],
                                          const bool periods[], bool reorder) const
```

`MPI_CART_CREATE` returns a handle to a new communicator to which the cartesian topology information is attached. If `reorder = false` then the rank of each process in the new group is identical to its rank in the old group. Otherwise, the function may reorder the processes (possibly so as to choose a good embedding of the virtual topology onto the physical machine). If the total size of the cartesian grid is smaller than the size of the group of `comm`, then some processes are returned `MPI_COMM_NULL`, in analogy to `MPI_COMM_SPLIT`. **If `ndims` is zero then a zero-dimensional Cartesian topology is created. The call is erroneous if it specifies a grid that is larger than the group size or if `ndims` is negative.**

6.5.2 Cartesian Convenience Function: `MPI_DIMS_CREATE`

For cartesian topologies, the function `MPI_DIMS_CREATE` helps the user select a balanced distribution of processes per coordinate direction, depending on the number of processes in the group to be balanced and optional constraints that can be specified by the user. One use is to partition all the processes (the size of `MPI_COMM_WORLD`'s group) into an n -dimensional topology.

```

MPI_DIMS_CREATE(nnodes, ndims, dims)
    IN          nnodes          number of nodes in a grid (integer)
    IN          ndims          number of cartesian dimensions (integer)
    INOUT       dims           integer array of size ndims specifying the number of
                                nodes in each dimension

int MPI_Dims_create(int nnodes, int ndims, int *dims)

MPI_DIMS_CREATE(NNODES, NDIMS, DIMS, IERROR)
    INTEGER NNODES, NDIMS, DIMS(*), IERROR

void MPI::Compute_dims(int nnodes, int ndims, int dims[])

```

The entries in the array `dims` are set to describe a cartesian grid with `ndims` dimensions and a total of `nnodes` nodes. The dimensions are set to be as close to each other as possible, using an appropriate divisibility algorithm. The caller may further constrain the operation of this routine by specifying elements of array `dims`. If `dims[i]` is set to a positive number, the routine will not modify the number of nodes in dimension `i`; only those entries where `dims[i] = 0` are modified by the call.

Negative input values of `dims[i]` are erroneous. An error will occur if `nnodes` is not a multiple of $\prod_{i, \text{dims}[i] \neq 0} \text{dims}[i]$.

For `dims[i]` set by the call, `dims[i]` will be ordered in non-increasing order. Array `dims` is suitable for use as input to routine `MPI_CART_CREATE`. `MPI_DIMS_CREATE` is local.

Example 6.1

| <code>dims</code> before call | function call | <code>dims</code> on return |
|----------------------------------|--|--------------------------------|
| (0,0) | <code>MPI_DIMS_CREATE(6, 2, dims)</code> | (3,2) |
| (0,0) | <code>MPI_DIMS_CREATE(7, 2, dims)</code> | (7,1) |
| (0,3,0) | <code>MPI_DIMS_CREATE(6, 3, dims)</code> | (2,3,1) |
| (0,3,0) | <code>MPI_DIMS_CREATE(7, 3, dims)</code> | erroneous call |

6.5.3 General (Graph) Constructor

```

1 MPI_GRAPH_CREATE(comm_old, nnodes, index, edges, reorder, comm_graph)
2
3
4 MPI_GRAPH_CREATE(comm_old, nnodes, index, edges, reorder, comm_graph)
5
6     IN      comm_old      input communicator (handle)
7
8     IN      nnodes       number of nodes in graph (integer)
9
10    IN      index        array of integers describing node degrees (see below)
11
12    IN      edges        array of integers describing graph edges (see below)
13
14    IN      reorder      ranking may be reordered (true) or not (false) (logical)
15
16    OUT     comm_graph    communicator with graph topology added (handle)
17
18 int MPI_Graph_create(MPI_Comm comm_old, int nnodes, int *index, int *edges,
19                     int reorder, MPI_Comm *comm_graph)
20
21 MPI_GRAPH_CREATE(COMM_OLD, NNODES, INDEX, EDGES, REORDER, COMM_GRAPH,
22                 IERROR)
23
24     INTEGER COMM_OLD, NNODES, INDEX(*), EDGES(*), COMM_GRAPH, IERROR
25     LOGICAL REORDER
26
27 MPI::Graphcomm MPI::Intracomm::Create_graph(int nnodes, const int index[],
28                                             const int edges[], bool reorder) const

```

MPI_GRAPH_CREATE returns a handle to a new communicator to which the graph topology information is attached. If `reorder = false` then the rank of each process in the new group is identical to its rank in the old group. Otherwise, the function may reorder the processes. If the size, `nnodes`, of the graph is smaller than the size of the group of `comm`, then some processes are returned `MPI_COMM_NULL`, in analogy to `MPI_CART_CREATE` and `MPI_COMM_SPLIT`. If the graph is empty, i.e., `nnodes == 0`, then `MPI_COMM_NULL` is returned in all processes. The call is erroneous if it specifies a graph that is larger than the group size of the input communicator.

The three parameters `nnodes`, `index` and `edges` define the graph structure. `nnodes` is the number of nodes of the graph. The nodes are numbered from 0 to `nnodes-1`. The `i`th entry of array `index` stores the total number of neighbors of the first `i` graph nodes. The lists of neighbors of nodes 0, 1, ..., `nnodes-1` are stored in consecutive locations in array `edges`. The array `edges` is a flattened representation of the edge lists. The total number of entries in `index` is `nnodes` and the total number of entries in `edges` is equal to the number of graph edges.

The definitions of the arguments `nnodes`, `index`, and `edges` are illustrated with the following simple example.

Example 6.2 Assume there are four processes 0, 1, 2, 3 with the following adjacency matrix:

| process | neighbors |
|---------|-----------|
| 0 | 1, 3 |
| 1 | 0 |
| 2 | 3 |
| 3 | 0, 2 |

Then, the input arguments are:

```

nnodes = 4
index = 2, 3, 4, 6
edges = 1, 3, 0, 3, 0, 2

```

Thus, in C, `index[0]` is the degree of node zero, and `index[i] - index[i-1]` is the degree of node `i`, `i=1, ..., nnodes-1`; the list of neighbors of node zero is stored in `edges[j]`, for $0 \leq j \leq \text{index}[0] - 1$ and the list of neighbors of node `i`, `i > 0`, is stored in `edges[j]`, $\text{index}[i-1] \leq j \leq \text{index}[i] - 1$.

In Fortran, `index(1)` is the degree of node zero, and `index(i+1) - index(i)` is the degree of node `i`, `i=1, ..., nnodes-1`; the list of neighbors of node zero is stored in `edges(j)`, for $1 \leq j \leq \text{index}(1)$ and the list of neighbors of node `i`, `i > 0`, is stored in `edges(j)`, $\text{index}(i) + 1 \leq j \leq \text{index}(i + 1)$.

A single process is allowed to be defined multiple times in the list of neighbors of a process (i.e., there may be multiple edges between two processes). A process is also allowed to be a neighbor to itself (i.e., a self loop in the graph). The adjacency matrix is allowed to be non-symmetric.

Advice to users. Performance implications of using multiple edges or a non-symmetric adjacency matrix are not defined. The definition of a node-neighbor edge does not imply a direction of the communication. (*End of advice to users.*)

Advice to implementors. The following topology information is likely to be stored with a communicator:

- Type of topology (cartesian/graph),
- For a cartesian topology:
 1. `ndims` (number of dimensions),
 2. `dims` (numbers of processes per coordinate direction),
 3. `periods` (periodicity information),
 4. `own_position` (own position in grid, could also be computed from rank and `dims`)
- For a graph topology:
 1. `index`,
 2. `edges`,

which are the vectors defining the graph structure.

For a graph structure the number of nodes is equal to the number of processes in the group. Therefore, the number of nodes does not have to be stored explicitly. An additional zero entry at the start of array `index` simplifies access to the topology information. (*End of advice to implementors.*)

6.5.4 Topology inquiry functions

If a topology has been defined with one of the above functions, then the topology information can be looked up using inquiry functions. They all are local calls.

```
MPI_TOPO_TEST(comm, status)
```

| | | |
|-----|--------|--|
| IN | comm | communicator (handle) |
| OUT | status | topology type of communicator comm (state) |

```
int MPI_Topo_test(MPI_Comm comm, int *status)
```

```
MPI_TOPO_TEST(COMM, STATUS, IERROR)
INTEGER COMM, STATUS, IERROR
```

```
int MPI::Comm::Get_topology() const
```

The function `MPI_TOPO_TEST` returns the type of topology that is assigned to a communicator.

The output value `status` is one of the following:

| | |
|----------------------------|--------------------|
| <code>MPI_GRAPH</code> | graph topology |
| <code>MPI_CART</code> | cartesian topology |
| <code>MPI_UNDEFINED</code> | no topology |

```
MPI_GRAPHDIMS_GET(comm, nnodes, nedges)
```

| | | |
|-----|--------|---|
| IN | comm | communicator for group with graph structure (handle) |
| OUT | nnodes | number of nodes in graph (integer) (same as number of processes in the group) |
| OUT | nedges | number of edges in graph (integer) |

```
int MPI_Graphdims_get(MPI_Comm comm, int *nnodes, int *nedges)
```

```
MPI_GRAPHDIMS_GET(COMM, NNODES, NEDGES, IERROR)
INTEGER COMM, NNODES, NEDGES, IERROR
```

```
void MPI::Graphcomm::Get_dims(int nnodes[], int nedges[]) const
```

Functions `MPI_GRAPHDIMS_GET` and `MPI_GRAPH_GET` retrieve the graph-topology information that was associated with a communicator by `MPI_GRAPH_CREATE`.

The information provided by `MPI_GRAPHDIMS_GET` can be used to dimension the vectors `index` and `edges` correctly for the following call to `MPI_GRAPH_GET`.

```

MPI_GRAPH_GET(comm, maxindex, maxedges, index, edges) 1
  IN      comm      communicator with graph structure (handle) 2
  IN      maxindex   length of vector index in the calling program 3
                    (integer) 4
  IN      maxedges   length of vector edges in the calling program 5
                    (integer) 6
  OUT     index      array of integers containing the graph structure (for 7
                    details see the definition of MPI_GRAPH_CREATE) 8
  OUT     edges      array of integers containing the graph structure 9
                                                    10
int MPI_Graph_get(MPI_Comm comm, int maxindex, int maxedges, int *index, 11
                  int *edges) 12
MPI_GRAPH_GET(COMM, MAXINDEX, MAXEDGES, INDEX, EDGES, IERROR) 13
  INTEGER COMM, MAXINDEX, MAXEDGES, INDEX(*), EDGES(*), IERROR 14
void MPI::Graphcomm::Get_topo(int maxindex, int maxedges, int index[], 15
                              int edges[]) const 16
                                                    17
MPI_CARTDIM_GET(comm, ndims) 18
  IN      comm      communicator with cartesian structure (handle) 19
  OUT     ndims     number of dimensions of the cartesian structure (inte- 20
                    ger) 21
int MPI_Cartdim_get(MPI_Comm comm, int *ndims) 22
MPI_CARTDIM_GET(COMM, NDIMS, IERROR) 23
  INTEGER COMM, NDIMS, IERROR 24
int MPI::Cartcomm::Get_dim() const 25

```

The functions `MPI_CARTDIM_GET` and `MPI_CART_GET` return the cartesian topology information that was associated with a communicator by `MPI_CART_CREATE`. If `comm` is associated with a zero-dimensional Cartesian topology, `MPI_CARTDIM_GET` returns `ndims=0` and `MPI_CART_GET` will keep all output arguments unchanged.

```

1 MPI_CART_GET(comm, maxdims, dims, periods, coords)
2     IN      comm      communicator with cartesian structure (handle)
3
4     IN      maxdims   length of vectors dims, periods, and coords in the
5                          calling program (integer)
6
7     OUT     dims      number of processes for each cartesian dimension (ar-
8                          ray of integer)
9
10    OUT     periods   periodicity (true/false) for each cartesian dimension
11                          (array of logical)
12
13    OUT     coords   coordinates of calling process in cartesian structure
14                          (array of integer)
15
16 int MPI_Cart_get(MPI_Comm comm, int maxdims, int *dims, int *periods,
17                 int *coords)
18
19 MPI_CART_GET(COMM, MAXDIMS, DIMS, PERIODS, COORDS, IERROR)
20     INTEGER COMM, MAXDIMS, DIMS(*), COORDS(*), IERROR
21     LOGICAL PERIODS(*)
22
23 void MPI::Cartcomm::Get_topo(int maxdims, int dims[], bool periods[],
24                             int coords[]) const
25
26 MPI_CART_RANK(comm, coords, rank)
27     IN      comm      communicator with cartesian structure (handle)
28     IN      coords   integer array (of size ndims) specifying the cartesian
29                          coordinates of a process
30
31     OUT     rank      rank of specified process (integer)
32
33 int MPI_Cart_rank(MPI_Comm comm, int *coords, int *rank)
34
35 MPI_CART_RANK(COMM, COORDS, RANK, IERROR)
36     INTEGER COMM, COORDS(*), RANK, IERROR
37
38 int MPI::Cartcomm::Get_cart_rank(const int coords[]) const
39
40     For a process group with cartesian structure, the function MPI_CART_RANK translates
41     the logical process coordinates to process ranks as they are used by the point-to-point
42     routines.
43
44     For dimension i with periods(i) = true, if the coordinate, coords(i), is out of
45     range, that is, coords(i) < 0 or coords(i) ≥ dims(i), it is shifted back to the interval
46      $0 \leq \text{coords}(i) < \text{dims}(i)$  automatically. Out-of-range coordinates are erroneous for
47     non-periodic dimensions.
48
49     If comm is associated with a zero-dimensional Cartesian topology,
50     coord is not significant and 0 is returned in rank.

```

MPI_CART_COORDS(comm, rank, maxdims, coords)

| | | | |
|-----|---------|---|---|
| IN | comm | communicator with cartesian structure (handle) | 1 |
| IN | rank | rank of a process within group of comm (integer) | 2 |
| IN | maxdims | length of vector <code>coords</code> in the calling program (integer) | 3 |
| OUT | coords | integer array (of size <code>ndims</code>) containing the cartesian coordinates of specified process (array of integers) | 4 |

```
int MPI_Cart_coords(MPI_Comm comm, int rank, int maxdims, int *coords)
```

```
MPI_CART_COORDS(COMM, RANK, MAXDIMS, COORDS, IERROR)
    INTEGER COMM, RANK, MAXDIMS, COORDS(*), IERROR
```

```
void MPI::Cartcomm::Get_coords(int rank, int maxdims, int coords[]) const
```

The inverse mapping, rank-to-coordinates translation is provided by `MPI_CART_COORDS`.

If `comm` is associated with a zero-dimensional Cartesian topology, `coords` will be unchanged.

MPI_GRAPH_NEIGHBORS_COUNT(comm, rank, nneighbors)

| | | | |
|-----|------------|--|---|
| IN | comm | communicator with graph topology (handle) | 1 |
| IN | rank | rank of process in group of comm (integer) | 2 |
| OUT | nneighbors | number of neighbors of specified process (integer) | 3 |

```
int MPI_Graph_neighbors_count(MPI_Comm comm, int rank, int *nneighbors)
```

```
MPI_GRAPH_NEIGHBORS_COUNT(COMM, RANK, NNEIGHBORS, IERROR)
    INTEGER COMM, RANK, NNEIGHBORS, IERROR
```

```
int MPI::Graphcomm::Get_neighbors_count(int rank) const
```

`MPI_GRAPH_NEIGHBORS_COUNT` and `MPI_GRAPH_NEIGHBORS` provide adjacency information for a general, graph topology.

MPI_GRAPH_NEIGHBORS(comm, rank, maxneighbors, neighbors)

| | | | |
|-----|--------------|---|---|
| IN | comm | communicator with graph topology (handle) | 1 |
| IN | rank | rank of process in group of comm (integer) | 2 |
| IN | maxneighbors | size of array <code>neighbors</code> (integer) | 3 |
| OUT | neighbors | ranks of processes that are neighbors to specified process (array of integer) | 4 |

```
int MPI_Graph_neighbors(MPI_Comm comm, int rank, int maxneighbors,
    int *neighbors)
```

```

1 MPI_GRAPH_NEIGHBORS(COMM, RANK, MAXNEIGHBORS, NEIGHBORS, IERROR)
2     INTEGER COMM, RANK, MAXNEIGHBORS, NEIGHBORS(*), IERROR
3
4 void MPI::Graphcomm::Get_neighbors(int rank, int maxneighbors, int
5     neighbors[]) const
6

```

Example 6.3 Suppose that `comm` is a communicator with a shuffle-exchange topology. The group has 2^n members. Each process is labeled by a_1, \dots, a_n with $a_i \in \{0, 1\}$, and has three neighbors: $\text{exchange}(a_1, \dots, a_n) = a_1, \dots, a_{n-1}, \bar{a}_n$ ($\bar{a} = 1 - a$), $\text{shuffle}(a_1, \dots, a_n) = a_2, \dots, a_n, a_1$, and $\text{unshuffle}(a_1, \dots, a_n) = a_n, a_1, \dots, a_{n-1}$. The graph adjacency list is illustrated below for $n = 3$.

| node | exchange neighbors(1) | shuffle neighbors(2) | unshuffle neighbors(3) |
|---------|--------------------------|-------------------------|---------------------------|
| 0 (000) | 1 | 0 | 0 |
| 1 (001) | 0 | 2 | 4 |
| 2 (010) | 3 | 4 | 1 |
| 3 (011) | 2 | 6 | 5 |
| 4 (100) | 5 | 1 | 2 |
| 5 (101) | 4 | 3 | 6 |
| 6 (110) | 7 | 5 | 3 |
| 7 (111) | 6 | 7 | 7 |

Suppose that the communicator `comm` has this topology associated with it. The following code fragment cycles through the three types of neighbors and performs an appropriate permutation for each.

```

27
28 C assume: each process has stored a real number A.
29 C extract neighborhood information
30     CALL MPI_COMM_RANK(comm, myrank, ierr)
31     CALL MPI_GRAPH_NEIGHBORS(comm, myrank, 3, neighbors, ierr)
32 C perform exchange permutation
33     CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, neighbors(1), 0,
34 +     neighbors(1), 0, comm, status, ierr)
35 C perform shuffle permutation
36     CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, neighbors(2), 0,
37 +     neighbors(3), 0, comm, status, ierr)
38 C perform unshuffle permutation
39     CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, neighbors(3), 0,
40 +     neighbors(2), 0, comm, status, ierr)
41

```

6.5.5 Cartesian Shift Coordinates

If the process topology is a cartesian structure, an `MPI_SENDRECV` operation is likely to be used along a coordinate direction to perform a shift of data. As input, `MPI_SENDRECV` takes the rank of a source process for the receive, and the rank of a destination process for the send. If the function `MPI_CART_SHIFT` is called for a cartesian process group, it provides the calling process with the above identifiers, which then can be passed to `MPI_SENDRECV`.

The user specifies the coordinate direction and the size of the step (positive or negative).
The function is local.

`MPI_CART_SHIFT(comm, direction, disp, rank_source, rank_dest)`

| | | |
|-----|--------------------------|--|
| IN | <code>comm</code> | communicator with cartesian structure (handle) |
| IN | <code>direction</code> | coordinate dimension of shift (integer) |
| IN | <code>disp</code> | displacement (> 0 : upwards shift, < 0 : downwards shift) (integer) |
| OUT | <code>rank_source</code> | rank of source process (integer) |
| OUT | <code>rank_dest</code> | rank of destination process (integer) |

```
int MPI_Cart_shift(MPI_Comm comm, int direction, int disp, int *rank_source,
                  int *rank_dest)
```

```
MPI_CART_SHIFT(COMM, DIRECTION, DISP, RANK_SOURCE, RANK_DEST, IERROR)
                INTEGER COMM, DIRECTION, DISP, RANK_SOURCE, RANK_DEST, IERROR
```

```
void MPI::Cartcomm::Shift(int direction, int disp, int& rank_source,
                          int& rank_dest) const
```

The `direction` argument indicates the dimension of the shift, i.e., the coordinate which value is modified by the shift. The coordinates are numbered from 0 to `ndims-1`, when `ndims` is the number of dimensions.

Depending on the periodicity of the cartesian group in the specified coordinate direction, `MPI_CART_SHIFT` provides the identifiers for a circular or an end-off shift. In the case of an end-off shift, the value `MPI_PROC_NULL` may be returned in `rank_source` or `rank_dest`, indicating that the source or the destination for the shift is out of range.

It is erroneous to call `MPI_CART_SHIFT` with a `direction` that is either negative or greater than or equal to the number of dimensions in the Cartesian communicator. This implies that it is erroneous to call `MPI_CART_SHIFT` with a `comm` that is associated with a zero-dimensional Cartesian topology.

Example 6.4 The communicator, `comm`, has a two-dimensional, periodic, cartesian topology associated with it. A two-dimensional array of `REALs` is stored one element per process, in variable `A`. One wishes to skew this array, by shifting column `i` (vertically, i.e., along the column) by `i` steps.

```
....
C find process rank
  CALL MPI_COMM_RANK(comm, rank, ierr)
C find cartesian coordinates
  CALL MPI_CART_COORDS(comm, rank, maxdims, coords, ierr)
C compute shift source and destination
  CALL MPI_CART_SHIFT(comm, 0, coords(2), source, dest, ierr)
C skew array
  CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, dest, 0, source, 0, comm,
+                           status, ierr)
```

Advice to users. In Fortran, the dimension indicated by `DIRECTION = i` has `DIMS(i+1)` nodes, where `DIMS` is the array that was used to create the grid. In C, the dimension indicated by `direction = i` is the dimension specified by `dims[i]`. (*End of advice to users.*)

6.5.6 Partitioning of Cartesian structures

`MPI_CART_SUB(comm, remain_dims, newcomm)`

| | | |
|-----|--------------------------|--|
| IN | <code>comm</code> | communicator with cartesian structure (handle) |
| IN | <code>remain_dims</code> | the <i>i</i> th entry of <code>remain_dims</code> specifies whether the <i>i</i> th dimension is kept in the subgrid (<code>true</code>) or is dropped (<code>false</code>) (logical vector) |
| OUT | <code>newcomm</code> | communicator containing the subgrid that includes the calling process (handle) |

```
int MPI_Cart_sub(MPI_Comm comm, int *remain_dims, MPI_Comm *newcomm)
```

```
MPI_CART_SUB(COMM, REMAIN_DIMS, NEWCOMM, IERROR)
    INTEGER COMM, NEWCOMM, IERROR
    LOGICAL REMAIN_DIMS(*)
```

```
MPI::Cartcomm MPI::Cartcomm::Sub(const bool remain_dims[]) const
```

If a cartesian topology has been created with `MPI_CART_CREATE`, the function `MPI_CART_SUB` can be used to partition the communicator group into subgroups that form lower-dimensional cartesian subgrids, and to build for each subgroup a communicator with the associated subgrid cartesian topology. **If all entries in `remain_dims` are false or `comm` is already associated with a zero-dimensional Cartesian topology then `newcomm` is associated with a zero-dimensional Cartesian topology.** (This function is closely related to `MPI_COMM_SPLIT`.)

Example 6.5 Assume that `MPI_CART_CREATE(..., comm)` has defined a $(2 \times 3 \times 4)$ grid. Let `remain_dims = (true, false, true)`. Then a call to,

```
MPI_CART_SUB(comm, remain_dims, comm_new),
```

will create three communicators each with eight processes in a 2×4 cartesian topology. If `remain_dims = (false, false, true)` then the call to `MPI_CART_SUB(comm, remain_dims, comm_new)` will create six non-overlapping communicators, each with four processes, in a one-dimensional cartesian topology.

6.5.7 Low-level topology functions

The two additional functions introduced in this section can be used to implement all other topology functions. In general they will not be called by the user directly, unless he or she is creating additional virtual topology capability other than that provided by MPI.


```

MPI_CART_MAP(comm, ndims, dims, periods, newrank) 1
IN      comm      input communicator (handle) 2
IN      ndims     number of dimensions of cartesian structure (integer) 3
IN      dims      integer array of size ndims specifying the number of 4
                processes in each coordinate direction 5
IN      periods   logical array of size ndims specifying the periodicity 6
                specification in each coordinate direction 7
OUT     newrank   reordered rank of the calling process; 8
                MPI_UNDEFINED if calling process does not belong 9
                to grid (integer) 10
                11
                12
                13
int MPI_Cart_map(MPI_Comm comm, int ndims, int *dims, int *periods, 14
                int *newrank) 15
MPI_CART_MAP(COMM, NDIMS, DIMS, PERIODS, NEWRANK, IERROR) 16
INTEGER COMM, NDIMS, DIMS(*), NEWRANK, IERROR 17
LOGICAL PERIODS(*) 18
int MPI::Cartcomm::Map(int ndims, const int dims[], const bool periods[]) 19
                const 20
                21
                22

```

MPI_CART_MAP computes an “optimal” placement for the calling process on the physical machine. A possible implementation of this function is to always return the rank of the calling process, that is, not to perform any reordering.

Advice to implementors. The function MPI_CART_CREATE(comm, ndims, dims, periods, reorder, comm_cart), with reorder = true can be implemented by calling MPI_CART_MAP(comm, ndims, dims, periods, newrank), then calling MPI_COMM_SPLIT(comm, color, key, comm_cart), with color = 0 if newrank ≠ MPI_UNDEFINED, color = MPI_UNDEFINED otherwise, and key = newrank.

The function MPI_CART_SUB(comm, remain_dims, comm_new) can be implemented by a call to MPI_COMM_SPLIT(comm, color, key, comm_new), using a single number encoding of the lost dimensions as color and a single number encoding of the preserved dimensions as key.

All other cartesian topology functions can be implemented locally, using the topology information that is cached with the communicator. (*End of advice to implementors.*)

The corresponding new function for general graph structures is as follows.

```

1 MPI_GRAPH_MAP(comm, nnodes, index, edges, newrank)
2     IN      comm      input communicator (handle)
3
4     IN      nnodes    number of graph nodes (integer)
5
6     IN      index     integer array specifying the graph structure, see
7                        MPI_GRAPH_CREATE
8     IN      edges     integer array specifying the graph structure
9     OUT     newrank   reordered rank of the calling process;
10                      MPI_UNDEFINED if the calling process does not be-
11                      long to graph (integer)

```

```

12
13 int MPI_Graph_map(MPI_Comm comm, int nnodes, int *index, int *edges,
14                  int *newrank)

```

```

15 MPI_GRAPH_MAP(COMM, NNODES, INDEX, EDGES, NEWRANK, IERROR)
16     INTEGER COMM, NNODES, INDEX(*), EDGES(*), NEWRANK, IERROR

```

```

17
18 int MPI::Graphcomm::Map(int nnodes, const int index[], const int edges[])
19     const

```

21 *Advice to implementors.* The function `MPI_GRAPH_CREATE(comm, nnodes, index, edges, reorder, comm_graph)`, with `reorder = true` can be implemented by calling `MPI_GRAPH_MAP(comm, nnodes, index, edges, newrank)`, then calling `MPI_COMM_SPLIT(comm, color, key, comm_graph)`, with `color = 0` if `newrank ≠ MPI_UNDEFINED`, `color = MPI_UNDEFINED` otherwise, and `key = newrank`.

27 All other graph topology functions can be implemented locally, using the topology information that is cached with the communicator. (*End of advice to implementors.*)

30 6.6 An Application Example

32 **Example 6.6** The example in figure 6.1 shows how the grid definition and inquiry functions can be used in an application program. A partial differential equation, for instance the Poisson equation, is to be solved on a rectangular domain. First, the processes organize themselves in a two-dimensional structure. Each process then inquires about the ranks of its neighbors in the four directions (up, down, right, left). The numerical problem is solved by an iterative method, the details of which are hidden in the subroutine `relax`.

38 In each relaxation step each process computes new values for the solution grid function at all points owned by the process. Then the values at inter-process boundaries have to be exchanged with neighboring processes. For example, the exchange subroutine might contain a call like `MPI_SEND(...,neigh_rank(1),...)` to send updated values to the left-hand neighbor (`i-1,j`).

```

integer ndims, num_neigh
logical reorder
parameter (ndims=2, num_neigh=4, reorder=.true.)
integer comm, comm_cart, dims(ndims), neigh_def(ndims), ierr
integer neigh_rank(num_neigh), own_position(ndims), i, j
logical periods(ndims)
real*8 u(0:101,0:101), f(0:101,0:101)
data dims / ndims * 0 /
comm = MPI_COMM_WORLD
C   Set process grid size and periodicity
call MPI_DIMS_CREATE(comm, ndims, dims,ierr)
periods(1) = .TRUE.
periods(2) = .TRUE.
C   Create a grid structure in WORLD group and inquire about own position
call MPI_CART_CREATE (comm, ndims, dims, periods, comm_cart,ierr)
call MPI_CART_GET (comm_cart, ndims, dims, periods, own_position,ierr)
C   Look up the ranks for the neighbors. Own process coordinates are (i,j).
C   Neighbors are (i-1,j), (i+1,j), (i,j-1), (i,j+1)
i = own_position(1)
j = own_position(2)
neigh_def(1) = i-1
neigh_def(2) = j
call MPI_CART_RANK (comm_cart, neigh_def, neigh_rank(1),ierr)
neigh_def(1) = i+1
neigh_def(2) = j
call MPI_CART_RANK (comm_cart, neigh_def, neigh_rank(2),ierr)
neigh_def(1) = i
neigh_def(2) = j-1
call MPI_CART_RANK (comm_cart, neigh_def, neigh_rank(3),ierr)
neigh_def(1) = i
neigh_def(2) = j+1
call MPI_CART_RANK (comm_cart, neigh_def, neigh_rank(4),ierr)
C   Initialize the grid functions and start the iteration
call init (u, f)
do 10 it=1,100
  call relax (u, f)
C   Exchange data with neighbor processes
  call exchange (u, comm_cart, neigh_rank, num_neigh)
10 continue
call output (u)
end

```

Figure 6.1: Set-up of process structure for two-dimensional parallel Poisson solver.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 7

MPI Environmental Management

This chapter discusses routines for getting and, where appropriate, setting various parameters that relate to the MPI implementation and the execution environment (such as error handling). The procedures for entering and leaving the MPI execution environment are also described here.

7.1 Implementation information

7.1.1 Version Inquiries

In order to cope with changes to the MPI Standard, there are both compile-time and runtime ways to determine which version of the standard is in use in the environment one is using.

The “version” will be represented by two separate integers, for the version and subversion: In C and C++,

```
#define MPI_VERSION    2
#define MPI_SUBVERSION 1
```

in Fortran,

```
INTEGER MPI_VERSION, MPI_SUBVERSION
PARAMETER (MPI_VERSION    = 2)
PARAMETER (MPI_SUBVERSION = 1)
```

For runtime determination,

`MPI_GET_VERSION(version, subversion)`

| | | |
|-----|------------|-----------------------------|
| OUT | version | version number (integer) |
| OUT | subversion | subversion number (integer) |

```
int MPI_Get_version(int *version, int *subversion)
```

```
MPI_GET_VERSION(VERSION, SUBVERSION, IERROR)
    INTEGER VERSION, SUBVERSION, IERROR
```

```
void MPI::Get_version(int& version, int& subversion)
```

1 MPI_GET_VERSION is one of the few functions that can be called before MPI_INIT
2 and after MPI_FINALIZE. Valid (MPI_VERSION, MPI_SUBVERSION) pairs in this and previous
3 versions of the MPI standard are (2,1), (2,0), and (1,2).
4

5 7.1.2 Environmental Inquiries

6 A set of attributes that describe the execution environment are attached to the commu-
7 nicator MPI_COMM_WORLD when MPI is initialized. The value of these attributes can be
8 inquired by using the function MPI_ATTR_GET described in Chapter 5. It is erroneous to
9 delete these attributes, free their keys, or change their values.
10

11 The list of predefined attribute keys include

12 MPI_TAG_UB Upper bound for tag value.

13
14 MPI_HOST Host process rank, if such exists, MPI_PROC_NULL, otherwise.

15
16 MPI_IO rank of a node that has regular I/O facilities (possibly myrank). Nodes in the same
17 communicator may return different values for this parameter.

18
19 MPI_WTIME_IS_GLOBAL Boolean variable that indicates whether clocks are synchronized.

20 Vendors may add implementation specific parameters (such as node number, real mem-
21 ory size, virtual memory size, etc.)

22 These predefined attributes do not change value between MPI initialization (MPI_INIT
23 and MPI completion (MPI_FINALIZE), and cannot be updated or deleted by users.
24

25 *Advice to users.* Note that in the C binding, the value returned by these attributes
26 is a *pointer* to an `int` containing the requested value. (*End of advice to users.*)
27

28 The required parameter values are discussed in more detail below:
29

30 Tag values

31 Tag values range from 0 to the value returned for MPI_TAG_UB inclusive. These values are
32 guaranteed to be unchanging during the execution of an MPI program. In addition, the tag
33 upper bound value must be *at least* 32767. An MPI implementation is free to make the
34 value of MPI_TAG_UB larger than this; for example, the value $2^{30} - 1$ is also a legal value for
35 MPI_TAG_UB.
36

37 The attribute MPI_TAG_UB has the same value on all processes of MPI_COMM_WORLD.
38

39 Host rank

40 The value returned for MPI_HOST gets the rank of the `HOST` process in the group associated
41 with communicator MPI_COMM_WORLD, if there is such. MPI_PROC_NULL is returned if
42 there is no host. MPI does not specify what it means for a process to be a `HOST`, nor does
43 it require that a `HOST` exists.
44

45 The attribute MPI_HOST has the same value on all processes of MPI_COMM_WORLD.
46
47
48

IO rank

The value returned for `MPI_IO` is the rank of a processor that can provide language-standard I/O facilities. For Fortran, this means that all of the Fortran I/O operations are supported (e.g., `OPEN`, `REWIND`, `WRITE`). For C, this means that all of the [ISO C](#) I/O operations are supported (e.g., `fopen`, `fprintf`, `lseek`).

If every process can provide language-standard I/O, then the value `MPI_ANY_SOURCE` will be returned. Otherwise, if the calling process can provide language-standard I/O, then its rank will be returned. Otherwise, if some process can provide language-standard I/O then the rank of one such process will be returned. The same value need not be returned by all processes. If no process can provide language-standard I/O, then the value `MPI_PROC_NULL` will be returned.

Advice to users. Note that input is not collective, and this attribute does *not* indicate which process can or does provide input. (*End of advice to users.*)

Clock synchronization

The value returned for `MPI_WTIME_IS_GLOBAL` is 1 if clocks at all processes in `MPI_COMM_WORLD` are synchronized, 0 otherwise. A collection of clocks is considered synchronized if explicit effort has been taken to synchronize them. The expectation is that the variation in time, as measured by calls to `MPI_WTIME`, will be less than one half the round-trip time for an MPI message of length zero. If time is measured at a process just before a send and at another process just after a matching receive, the second time should be always higher than the first one.

The attribute `MPI_WTIME_IS_GLOBAL` need not be present when the clocks are not synchronized (however, the attribute key `MPI_WTIME_IS_GLOBAL` is always valid). This attribute may be associated with communicators other than `MPI_COMM_WORLD`.

The attribute `MPI_WTIME_IS_GLOBAL` has the same value on all processes of `MPI_COMM_WORLD`.

`MPI_GET_PROCESSOR_NAME(name, resultlen)`

| | | |
|-----|-----------|---|
| OUT | name | A unique specifier for the actual (as opposed to virtual) node. |
| OUT | resultlen | Length (in printable characters) of the result returned in name |

```
int MPI_Get_processor_name(char *name, int *resultlen)
```

```
MPI_GET_PROCESSOR_NAME( NAME, RESULTLEN, IERROR)
```

```
CHARACTER*(*) NAME
```

```
INTEGER RESULTLEN, IERROR
```

```
void MPI::Get_processor_name(char* name, int& resultlen)
```

```
void MPI::Get_processor_name(char* name, int& resultlen)
```

This routine returns the name of the processor on which it was called at the moment of the call. The name is a character string for maximum flexibility. From this value it

1 must be possible to identify a specific piece of hardware; possible values include “processor
 2 9 in rack 4 of mpp.cs.org” and “231” (where 231 is the actual processor number in the
 3 running homogeneous system). The argument `name` must represent storage that is at least
 4 `MPI_MAX_PROCESSOR_NAME` characters long. `MPI_GET_PROCESSOR_NAME` may write up
 5 to this many characters into `name`.

6 The number of characters actually written is returned in the output argument, `resultlen`.
 7 In C, a null character is additionally stored at `name[resultlen]`. The `resultlen` cannot be larger
 8 than `MPI_MAX_PROCESSOR_NAME-1`. In Fortran, `name` is padded on the right with blank
 9 characters. The `resultlen` cannot be larger than `MPI_MAX_PROCESSOR_NAME`.

10
 11 *Rationale.* This function allows MPI implementations that do process migration
 12 to return the current processor. Note that nothing in MPI *requires* or defines pro-
 13 cess migration; this definition of `MPI_GET_PROCESSOR_NAME` simply allows such
 14 an implementation. (*End of rationale.*)

15
 16 *Advice to users.* The user must provide at least `MPI_MAX_PROCESSOR_NAME` space
 17 to write the processor name — processor names can be this long. The user should
 18 examine the `output` argument, `resultlen`, to determine the actual length of the name.
 19 (*End of advice to users.*)

20 The constant `MPI_BSEND_OVERHEAD` provides an upper bound on the fixed overhead
 21 per message buffered by a call to `MPI_BSEND` (see Section 3.6.1).
 22

23 7.2 Memory Allocation

24
 25 In some systems, message-passing and remote-memory-access (RMA) operations run faster
 26 when accessing specially allocated memory (e.g., memory that is shared by the other pro-
 27 cesses in the communicating group on an SMP). MPI provides a mechanism for allocating
 28 and freeing such special memory. The use of such memory for message passing or RMA is
 29 not mandatory, and this memory can be used without restrictions as any other dynamically
 30 allocated memory. However, implementations may restrict the use of the `MPI_WIN_LOCK`
 31 and `MPI_WIN_UNLOCK` functions to windows allocated in such memory (see Section 10.4.3.)
 32

33
 34 `MPI_ALLOC_MEM(size, info, baseptr)`

| | | | |
|----|-----|----------------------|---|
| 35 | IN | <code>size</code> | size of memory segment in bytes (nonnegative integer) |
| 36 | IN | <code>info</code> | <code>info</code> argument (handle) |
| 37 | OUT | <code>baseptr</code> | pointer to beginning of memory segment allocated |

38
 39
 40 `int MPI_Alloc_mem(MPI_Aint size, MPI_Info info, void *baseptr)`

41
 42 `MPI_ALLOC_MEM(SIZE, INFO, BASEPTR, IERROR)`
 43 `INTEGER INFO, IERROR`
 44 `INTEGER(KIND=MPI_ADDRESS_KIND) SIZE, BASEPTR`

45
 46 `void* MPI::Alloc_mem(MPI::Aint size, const MPI::Info& info)`

47 The `info` argument can be used to provide directives that control the desired location
 48 of the allocated memory. Such a directive does not affect the semantics of the call. Valid

1 Since standard Fortran does not support (C-like) pointers, this code is not Fortran 77
 2 or Fortran 90 code. Some compilers (in particular, at the time of writing, g77 and Fortran
 3 compilers for Intel) do not support this code.

4
 5 **Example 7.2** Same example, in C

```
6        float (* f)[100][100] ;
7        /* no memory is allocated */
8        MPI_Alloc_mem(sizeof(float)*100*100, MPI_INFO_NULL, &f);
9        /* memory allocated */
10        ...
11        (*f)[5][3] = 2.71;
12        ...
13        MPI_Free_mem(f);
14
```

17 7.3 Error handling

18
 19 An MPI implementation cannot or may choose not to handle some errors that occur during
 20 MPI calls. These can include errors that generate exceptions or traps, such as floating point
 21 errors or access violations. The set of errors that are handled by MPI is implementation-
 22 dependent. Each such error generates an **MPI exception**.

23 The above text takes precedence over any text on error handling within this document.
 24 Specifically, text that states that errors *will* be handled should be read as *may* be handled.

25 A user can associate an error handler with a communicator. The specified error han-
 26 dling routine will be used for any MPI exception that occurs during a call to MPI for a
 27 communication with this communicator. MPI calls that are not related to any communi-
 28 cator are considered to be attached to the communicator MPI_COMM_WORLD. The attachment
 29 of error handlers to communicators is purely local: different processes may attach different
 30 error handlers to the same communicator.

31 A newly created communicator inherits the error handler that is associated with the
 32 “parent” communicator. In particular, the user can specify a “global” error handler for
 33 all communicators by associating this handler with the communicator MPI_COMM_WORLD
 34 immediately after initialization.

35 Several predefined error handlers are available in MPI:

36
 37 **MPI_ERRORS_ARE_FATAL** The handler, when called, causes the program to abort on all exe-
 38 cuting processes. This has the same effect as if MPI_ABORT was called by the process
 39 that invoked the handler.

40
 41 **MPI_ERRORS_RETURN** The handler has no effect other than returning the error code to the
 42 user.

43 Implementations may provide additional predefined error handlers and programmers
 44 can code their own error handlers.

45 The error handler MPI_ERRORS_ARE_FATAL is associated by default with MPI_COMM-
 46 WORLD after initialization. Thus, if the user chooses not to control error handling, every
 47 error that MPI handles is treated as fatal. Since (almost) all MPI calls return an error code,
 48 a user may choose to handle errors in its main code, by testing the return code of MPI calls

and executing a suitable recovery code when the call was not successful. In this case, the error handler `MPI_ERRORS_RETURN` will be used. Usually it is more convenient and more efficient not to test for errors after each MPI call, and have such error handled by a non trivial MPI error handler.

After an error is detected, the state of MPI is undefined. That is, using a user-defined error handler, or `MPI_ERRORS_RETURN`, does *not* necessarily allow the user to continue to use MPI after an error is detected. The purpose of these error handlers is to allow a user to issue user-defined error messages and to take actions unrelated to MPI (such as flushing I/O buffers) before a program exits. An MPI implementation is free to allow MPI to continue after an error but is not required to do so.

Advice to implementors. A good quality implementation will, to the greatest possible extent, circumscribe the impact of an error, so that normal processing can continue after an error handler was invoked. The implementation documentation will provide information on the possible effect of each class of errors. (*End of advice to implementors.*)

An MPI error handler is an opaque object, which is accessed by a handle. MPI calls are provided to create new error handlers, to associate error handlers with communicators, and to test which error handler is associated with a communicator.

7.3.1 Extended Error Handling in MPI-2

MPI-1 attached error handlers only to communicators. MPI-2 attaches them to three types of objects: communicators, windows, and files. The extension was done while maintaining only one type of error handler opaque object. On the other hand, there are, in C and C++, distinct typedefs for user defined error handling callback functions that accept, respectively, communicator, file, and window arguments. In Fortran there are three user routines.

An error handler object is created by a call to `MPI_XXX_CREATE_ERRHANDLER(function, errhandler)`, where XXX is, respectively, `COMM`, `WIN`, or `FILE`.

An error handler is attached to a communicator, window, or file by a call to `MPI_XXX_SET_ERRHANDLER`. The error handler must be either a predefined error handler, or an error handler that was created by a call to `MPI_XXX_CREATE_ERRHANDLER`, with matching XXX. The predefined error handlers `MPI_ERRORS_RETURN` and `MPI_ERRORS_ARE_FATAL` can be attached to communicators, windows, and files. In C++, the predefined error handler `MPI::ERRORS_THROW_EXCEPTIONS` can also be attached to communicators, windows, and files.

The error handler currently associated with a communicator, window, or file can be retrieved by a call to `MPI_XXX_GET_ERRHANDLER`.

The MPI-1 function `MPI_ERRHANDLER_FREE` can be used to free an error handler that was created by a call to `MPI_XXX_CREATE_ERRHANDLER`.

`MPI_{COMM,WIN,FILE}_GET_ERRHANDLER` behave as if a new error handler object is created. That is, once the error handler is no longer needed, `MPI_ERRHANDLER_FREE` should be called with the error handler returned from `MPI_ERRHANDLER_GET` or `MPI_{COMM,WIN,FILE}_GET_ERRHANDLER` to mark the error handler for deallocation. This provides behavior similar to that of `MPI_COMM_GROUP` and `MPI_GROUP_FREE`.

Advice to implementors. High quality implementation should raise an error when an error handler that was created by a call to `MPI_XXX_CREATE_ERRHANDLER` is

attached to an object of the wrong type with a call to `MPI_YYY_SET_ERRHANDLER`. To do so, it is necessary to maintain, with each error handler, information on the typedef of the associated user function. (*End of advice to implementors.*)

The syntax for these calls is given below.

7.3.2 Error Handlers for Communicators

`MPI_COMM_CREATE_ERRHANDLER(function, errhandler)`

| | | |
|-----|------------|--|
| IN | function | user defined error handling procedure (function) |
| OUT | errhandler | MPI error handler (handle) |

```
int MPI_Comm_create_errhandler(MPI_Comm_errhandler_fn *function,
                               MPI_Errhandler *errhandler)
```

```
MPI_COMM_CREATE_ERRHANDLER(FUNCTION, ERRHANDLER, IERROR)
EXTERNAL FUNCTION
INTEGER ERRHANDLER, IERROR
```

```
MPI::Intracomm MPI::Intracomm::Create(const MPI::Group& group) const
```

```
static MPI::Errhandler
    MPI::Comm::Create_errhandler(MPI::Comm::Errhandler_fn*
                                function)
```

Creates an error handler that can be attached to communicators. This function is identical to `MPI_ERRHANDLER_CREATE`, whose use is deprecated.

The user routine should be, in C, a function of type `MPI_Comm_errhandler_fn`, which is defined as

```
typedef void MPI_Comm_errhandler_fn(MPI_Comm *, int *, ...);
```

The first argument is the communicator in use. The second is the error code to be returned by the MPI routine that raised the error. If the routine would have returned `MPI_ERR_IN_STATUS`, it is the error code returned in the status for the request that caused the error handler to be invoked. The remaining arguments are “`stdargs`” arguments whose number and meaning is implementation-dependent. An implementation should clearly document these arguments. Addresses are used so that the handler may be written in Fortran. This typedef replaces `MPI_Handler_function`, whose use is deprecated.

In Fortran, the user routine should be of the form:

```
SUBROUTINE COMM_ERRHANDLER_FN(COMM, ERROR_CODE, ...)
INTEGER COMM, ERROR_CODE
```

Advice to users. Users are discouraged from using a Fortran `HANDLER_FUNCTION` since the routine expects a variable number of arguments. Some Fortran systems may allow this but some may fail to give the correct result or compile/link this code. Thus, it will not, in general, be possible to create portable code with a Fortran `HANDLER_FUNCTION`. (*End of advice to users.*)

In C++, the user routine should be of the form:

```
typedef void MPI::Comm::Errhandler_fn(MPI::Comm &, int *, ...);
```

Rationale. The variable argument list is provided because it provides an ISO-standard hook for providing additional information to the error handler; without this hook, ISO C prohibits additional arguments. (*End of rationale.*)

```
MPI_COMM_SET_ERRHANDLER(comm, errhandler)
```

```
INOUT  comm                communicator (handle)
```

```
IN      errhandler         new error handler for communicator (handle)
```

```
int MPI_Comm_set_errhandler(MPI_Comm comm, MPI_Errhandler errhandler)
```

```
MPI_COMM_SET_ERRHANDLER(COMM, ERRHANDLER, IERROR)
```

```
INTEGER COMM, ERRHANDLER, IERROR
```

```
void MPI::Comm::Set_errhandler(const MPI::Errhandler& errhandler)
```

Attaches a new error handler to a communicator. The error handler must be either a predefined error handler, or an error handler created by a call to `MPI_COMM_CREATE_ERRHANDLER`. This call is identical to `MPI_ERRHANDLER_SET`, whose use is deprecated.

```
MPI_COMM_GET_ERRHANDLER(comm, errhandler)
```

```
IN      comm                communicator (handle)
```

```
OUT     errhandler         error handler currently associated with communicator  
                        (handle)
```

```
int MPI_Comm_get_errhandler(MPI_Comm comm, MPI_Errhandler *errhandler)
```

```
MPI_COMM_GET_ERRHANDLER(COMM, ERRHANDLER, IERROR)
```

```
INTEGER COMM, ERRHANDLER, IERROR
```

```
MPI::Errhandler MPI::Comm::Get_errhandler() const
```

Retrieves the error handler currently associated with a communicator. This call is identical to `MPI_ERRHANDLER_GET`, whose use is deprecated.

Example: A library function may register at its entry point the current error handler for a communicator, set its own private error handler for this communicator, and restore before exiting the previous error handler.

7.3.3 Error Handlers for Windows

`MPI_WIN_CREATE_ERRHANDLER(function, errhandler)`

| | | |
|-----|------------|--|
| IN | function | user defined error handling procedure (function) |
| OUT | errhandler | MPI error handler (handle) |

```
int MPI_Win_create_errhandler(MPI_Win_errhandler_fn *function, MPI_Errhandler
                             *errhandler)
```

`MPI_WIN_CREATE_ERRHANDLER(FUNCTION, ERRHANDLER, IERROR)`

EXTERNAL FUNCTION
INTEGER ERRHANDLER, IERROR

```
static MPI::Errhandler MPI::Win::Create_errhandler(MPI::Win::Errhandler_fn*
                                                  function)
```

The user routine should be, in C, a function of type `MPI_Win_errhandler_fn`, which is defined as

```
typedef void MPI_Win_errhandler_fn(MPI_Win *, int *, ...);
```

The first argument is the window in use, the second is the error code to be returned.

In Fortran, the user routine should be of the form:

```
SUBROUTINE WIN_ERRHANDLER_FN(WIN, ERROR_CODE, ...)
  INTEGER WIN, ERROR_CODE
```

In C++, the user routine should be of the form:

```
typedef void MPI::Win::Errhandler_fn(MPI::Win &, int *, ...);
```

`MPI_WIN_SET_ERRHANDLER(win, errhandler)`

| | | |
|-------|------------|---------------------------------------|
| INOUT | win | window (handle) |
| IN | errhandler | new error handler for window (handle) |

```
int MPI_Win_set_errhandler(MPI_Win win, MPI_Errhandler errhandler)
```

`MPI_WIN_SET_ERRHANDLER(WIN, ERRHANDLER, IERROR)`

INTEGER WIN, ERRHANDLER, IERROR

```
void MPI::Win::Set_errhandler(const MPI::Errhandler& errhandler)
```

Attaches a new error handler to a window. The error handler must be either a pre-defined error handler, or an error handler created by a call to `MPI_WIN_CREATE_ERRHANDLER`.

MPI_WIN_GET_ERRHANDLER(win, errhandler) 1

IN win window (handle) 2

OUT errhandler error handler currently associated with window (handle) 3
4
5

int MPI_Win_get_errhandler(MPI_Win win, MPI_Errhandler *errhandler) 6
7

MPI_WIN_GET_ERRHANDLER(WIN, ERRHANDLER, IERROR) 8
9
INTEGER WIN, ERRHANDLER, IERROR 10

MPI::Errhandler MPI::Win::Get_errhandler() const 11

Retrieves the error handler currently associated with a window. 12
13

7.3.4 Error Handlers for Files 14

MPI_FILE_CREATE_ERRHANDLER(function, errhandler) 15
16
17

IN function user defined error handling procedure (function) 18
19

OUT errhandler MPI error handler (handle) 20
21

int MPI_File_create_errhandler(MPI_File_errhandler_fn *function, 22
MPI_Errhandler *errhandler) 23
24

MPI_FILE_CREATE_ERRHANDLER(FUNCTION, ERRHANDLER, IERROR) 25
26
EXTERNAL FUNCTION 27
INTEGER ERRHANDLER, IERROR 28

static MPI::Errhandler 29
MPI::File::Create_errhandler(MPI::File::Errhandler_fn* 30
function) 31

The user routine should be, in C, a function of type MPI_File_errhandler_fn, which is 32
defined as 33

```
typedef void MPI_File_errhandler_fn(MPI_File *, int *, ...); 34  
35
```

The first argument is the file in use, the second is the error code to be returned. 36

In Fortran, the user routine should be of the form: 37

```
SUBROUTINE FILE_ERRHANDLER_FN(FILE, ERROR_CODE, ...) 38  
INTEGER FILE, ERROR_CODE 39
```

In C++, the user routine should be of the form: 40

```
typedef void MPI::File::Errhandler_fn(MPI::File &, int *, ...); 41  
42
```

1 MPI_FILE_SET_ERRHANDLER(file, errhandler)

2 INOUT file file (handle)

3 IN errhandler new error handler for file (handle)

4 int MPI_File_set_errhandler(MPI_File file, MPI_Errhandler errhandler)

5 MPI_FILE_SET_ERRHANDLER(FILE, ERRHANDLER, IERROR)

6 INTEGER FILE, ERRHANDLER, IERROR

7 void MPI::File::Set_errhandler(const MPI::Errhandler& errhandler)

8 Attaches a new error handler to a file. The error handler must be either a predefined
9 error handler, or an error handler created by a call to MPI_FILE_CREATE_ERRHANDLER.

10 MPI_FILE_GET_ERRHANDLER(file, errhandler)

11 IN file file (handle)

12 OUT errhandler error handler currently associated with file (handle)

13 int MPI_File_get_errhandler(MPI_File file, MPI_Errhandler *errhandler)

14 MPI_FILE_GET_ERRHANDLER(FILE, ERRHANDLER, IERROR)

15 INTEGER FILE, ERRHANDLER, IERROR

16 MPI::Errhandler MPI::File::Get_errhandler() const

17 Retrieves the error handler currently associated with a file.

18 7.3.5 Freeing Errorhandlers and Retrieving Error Strings

19 MPI_ERRHANDLER_FREE(errhandler)

20 INOUT errhandler MPI error handler (handle)

21 int MPI_Errhandler_free(MPI_Errhandler *errhandler)

22 MPI_ERRHANDLER_FREE(ERRHANDLER, IERROR)

23 INTEGER ERRHANDLER, IERROR

24 void MPI::Errhandler::Free()

25 void MPI::Errhandler::Free()

26 Marks the error handler associated with `errhandler` for deallocation and sets `errhandler`
27 to `MPI_ERRHANDLER_NULL`. The error handler will be deallocated after all communicators
28 associated with it have been deallocated.

| | | | |
|--|-----------|--|---|
| MPI_ERROR_STRING(errorcode, string, resultlen) | | | 1 |
| IN | errorcode | Error code returned by an MPI routine | 2 |
| OUT | string | Text that corresponds to the errorcode | 3 |
| OUT | resultlen | Length (in printable characters) of the result returned in string | 4 |
| | | | 5 |
| | | | 6 |
| | | | 7 |

```
int MPI_Error_string(int errorcode, char *string, int *resultlen)
```

```
MPI_ERROR_STRING(ERRORCODE, STRING, RESULTLEN, IERROR)
```

```
INTEGER ERRORCODE, RESULTLEN, IERROR
```

```
CHARACTER*(*) STRING
```

```
void MPI::Get_error_string(int errorcode, char* name, int& resultlen)
```

```
void MPI::Get_error_string(int errorcode, char* name, int& resultlen)
```

Returns the error string associated with an error code or class. The argument `string` must represent storage that is at least `MPI_MAX_ERROR_STRING` characters long.

The number of characters actually written is returned in the output argument, `resultlen`.

Rationale. The form of this function was chosen to make the Fortran and C bindings similar. A version that returns a pointer to a string has two difficulties. First, the return string must be statically allocated and different for each error message (allowing the pointers returned by successive calls to `MPI_ERROR_STRING` to point to the correct message). Second, in Fortran, a function declared as returning `CHARACTER*(*)` can not be referenced in, for example, a `PRINT` statement. (*End of rationale.*)

7.4 Error codes and classes

The error codes returned by MPI are left entirely to the implementation (with the exception of `MPI_SUCCESS`). This is done to allow an implementation to provide as much information as possible in the error code (for use with `MPI_ERROR_STRING`).

To make it possible for an application to interpret an error code, the routine `MPI_ERROR_CLASS` converts any error code into one of a small set of standard error codes, called *error classes*. [Valid error classes are shown in Table 7.1 and Table 7.2.](#)

The error classes are a subset of the error codes: an MPI function may return an error class number; and the function `MPI_ERROR_STRING` can be used to compute the error string associated with an error class. An MPI error class is a valid MPI error code. Specifically, the values defined for MPI error classes are valid MPI error codes.

The error codes satisfy,

$$0 = \text{MPI_SUCCESS} < \text{MPI_ERR_...} \leq \text{MPI_ERR_LASTCODE}.$$

Rationale. The difference between `MPI_ERR_UNKNOWN` and `MPI_ERR_OTHER` is that `MPI_ERROR_STRING` can return useful information about `MPI_ERR_OTHER`.

Note that `MPI_SUCCESS = 0` is necessary to be consistent with C practice; the separation of error classes and error codes allows us to define the error classes this way. Having a known `LASTCODE` is often a nice sanity check as well. (*End of rationale.*)

| | | |
|----|----------------------|--|
| 1 | | |
| 2 | | |
| 3 | MPI_SUCCESS | No error |
| 4 | MPI_ERR_BUFFER | Invalid buffer pointer |
| 5 | MPI_ERR_COUNT | Invalid count argument |
| 6 | MPI_ERR_TYPE | Invalid datatype argument |
| 7 | MPI_ERR_TAG | Invalid tag argument |
| 8 | MPI_ERR_COMM | Invalid communicator |
| 9 | MPI_ERR_RANK | Invalid rank |
| 10 | MPI_ERR_REQUEST | Invalid request (handle) |
| 11 | MPI_ERR_ROOT | Invalid root |
| 12 | MPI_ERR_GROUP | Invalid group |
| 13 | MPI_ERR_OP | Invalid operation |
| 14 | MPI_ERR_TOPOLOGY | Invalid topology |
| 15 | MPI_ERR_DIMS | Invalid dimension argument |
| 16 | MPI_ERR_ARG | Invalid argument of some other kind |
| 17 | MPI_ERR_UNKNOWN | Unknown error |
| 18 | MPI_ERR_TRUNCATE | Message truncated on receive |
| 19 | MPI_ERR_OTHER | Known error not in this list |
| 20 | MPI_ERR_INTERN | Internal MPI (implementation) error |
| 21 | MPI_ERR_IN_STATUS | Error code is in status |
| 22 | MPI_ERR_PENDING | Pending request |
| 23 | MPI_ERR_KEYVAL | Invalid keyval has been passed |
| 24 | MPI_ERR_NO_MEM | MPI_ALLOC_MEM failed because memory is exhausted |
| 25 | MPI_ERR_BASE | Invalid base passed to MPI_FREE_MEM |
| 26 | MPI_ERR_INFO_KEY | Key longer than MPI_MAX_INFO_KEY |
| 27 | MPI_ERR_INFO_VALUE | Value longer than MPI_MAX_INFO_VAL |
| 28 | MPI_ERR_INFO_NOKEY | Invalid key passed to MPI_INFO_DELETE |
| 29 | MPI_ERR_SPAWN | Error in spawning processes |
| 30 | MPI_ERR_PORT | Invalid port name passed to MPI_COMM_CONNECT |
| 31 | | |
| 32 | MPI_ERR_SERVICE | Invalid service name passed to MPI_UNPUBLISH_NAME |
| 33 | | |
| 34 | MPI_ERR_NAME | Invalid service name passed to MPI_LOOKUP_NAME |
| 35 | | |
| 36 | MPI_ERR_WIN | invalid win argument |
| 37 | MPI_ERR_SIZE | invalid size argument |
| 38 | MPI_ERR_DISP | invalid disp argument |
| 39 | MPI_ERR_INFO | invalid info argument |
| 40 | MPI_ERR_LOCKTYPE | invalid locktype argument |
| 41 | MPI_ERR_ASSERT | invalid assert argument |
| 42 | MPI_ERR_RMA_CONFLICT | conflicting accesses to window |
| 43 | MPI_ERR_RMA_SYNC | wrong synchronization of RMA calls |
| 44 | | |
| 45 | | |
| 46 | | |
| 47 | | |
| 48 | | |

Table 7.1: Error classes (Part 1)

| | | |
|--|--|----------------------|
| <code>MPI_ERR_FILE</code> | Invalid file handle | 1 |
| <code>MPI_ERR_NOT_SAME</code> | Collective argument not identical on all processes, or collective routines called in a different order by different processes | 2 3 4 |
| <code>MPI_ERR_AMODE</code> | Error related to the <code>amode</code> passed to <code>MPI_FILE_OPEN</code> | 5 6 |
| <code>MPI_ERR_UNSUPPORTED_DATAREP</code> | Unsupported <code>datarep</code> passed to <code>MPI_FILE_SET_VIEW</code> | 7 8 |
| <code>MPI_ERR_UNSUPPORTED_OPERATION</code> | Unsupported operation, such as seeking on a file which supports sequential access only | 9 10 |
| <code>MPI_ERR_NO_SUCH_FILE</code> | File does not exist | 11 |
| <code>MPI_ERR_FILE_EXISTS</code> | File exists | 12 |
| <code>MPI_ERR_BAD_FILE</code> | Invalid file name (e.g., path name too long) | 13 |
| <code>MPI_ERR_ACCESS</code> | Permission denied | 14 |
| <code>MPI_ERR_NO_SPACE</code> | Not enough space | 15 |
| <code>MPI_ERR_QUOTA</code> | Quota exceeded | 16 |
| <code>MPI_ERR_READ_ONLY</code> | Read-only file or file system | 17 |
| <code>MPI_ERR_FILE_IN_USE</code> | File operation could not be completed, as the file is currently open by some process | 18 19 |
| <code>MPI_ERR_DUP_DATAREP</code> | Conversion functions could not be registered because a data representation identifier that was already defined was passed to <code>MPI_REGISTER_DATAREP</code> | 20 21 22 23 |
| <code>MPI_ERR_CONVERSION</code> | An error occurred in a user supplied data conversion function. | 24 25 |
| <code>MPI_ERR_IO</code> | Other I/O error | 26 |
| <code>MPI_ERR_LASTCODE</code> | Last error code | 27 |

Table 7.2: Error classes (Part 2)

| | | | |
|--|--|--|----------|
| <code>MPI_ERROR_CLASS(errorcode, errorclass)</code> | | 28 29 30 31 32 | |
| IN | <code>errorcode</code> | Error code returned by an MPI routine | 33 |
| OUT | <code>errorclass</code> | Error class associated with <code>errorcode</code> | 34 35 |
| <code>int MPI_Error_class(int errorcode, int *errorclass)</code> | | 36 37 | |
| <code>MPI_ERROR_CLASS(ERRORCODE, ERRORCLASS, IERROR)</code> | | 38 | |
| INTEGER | <code>ERRORCODE, ERRORCLASS, IERROR</code> | 39 40 | |
| <code>int MPI::Get_error_class(int errorcode)</code> | | 41 | |
| <code>int MPI::Get_error_class(int errorcode)</code> | | 42 43 | |
| The function <code>MPI_ERROR_CLASS</code> maps each standard error code (error class) onto itself. | | 44 45 46 47 48 | |

7.5 Timers and synchronization

MPI defines a timer. A timer is specified even though it is not “message-passing,” because timing parallel programs is important in “performance debugging” and because existing timers (both in POSIX 1003.1-1988 and 1003.4D 14.1 and in Fortran 90) are either inconvenient or do not provide adequate access to high-resolution timers. See also Section 2.6.5 on page 21.

MPI_WTIME()

```
double MPI_Wtime(void)
```

```
DOUBLE PRECISION MPI_WTIME()
```

```
double MPI::Wtime()
```

MPI_WTIME returns a floating-point number of seconds, representing elapsed wall-clock time since some time in the past.

The “time in the past” is guaranteed not to change during the life of the process. The user is responsible for converting large numbers of seconds to other units if they are preferred.

This function is portable (it returns seconds, not “ticks”), it allows high-resolution, and carries no unnecessary baggage. One would use it like this:

```
{
    double starttime, endtime;
    starttime = MPI_Wtime();
    .... stuff to be timed ...
    endtime = MPI_Wtime();
    printf("That took %f seconds\n",endtime-starttime);
}
```

The times returned are local to the node that called them. There is no requirement that different nodes return “the same time.” (But see also the discussion of MPI_WTIME_IS_GLOBAL).

MPI_WTICK()

```
double MPI_Wtick(void)
```

```
DOUBLE PRECISION MPI_WTICK()
```

```
double MPI::Wtick()
```

MPI_WTICK returns the resolution of MPI_WTIME in seconds. That is, it returns, as a double precision value, the number of seconds between successive clock ticks. For example, if the clock is implemented by the hardware as a counter that is incremented every millisecond, the value returned by MPI_WTICK should be 10^{-3} .

7.6 Startup

One goal of MPI is to achieve *source code portability*. By this we mean that a program written using MPI and complying with the relevant language standards is portable as written, and must not require any source code changes when moved from one system to another. This explicitly does *not* say anything about how an MPI program is started or launched from the command line, nor what the user must do to set up the environment in which an MPI program will run. However, an implementation may require some setup to be performed before other MPI routines may be called. To provide for this, MPI includes an initialization routine `MPI_INIT`.

```
MPI_INIT()
```

```
int MPI_Init(int *argc, char ***argv)
```

```
MPI_INIT(IERROR)
    INTEGER IERROR
```

```
void MPI::Init(int& argc, char**& argv)
```

```
void MPI::Init()
```

```
void MPI::Init(int& argc, char**& argv)
```

```
void MPI::Init()
```

This routine must be called before any other MPI routine. It must be called at most once; subsequent calls are erroneous (see `MPI_INITIALIZED`).

All MPI programs must contain a call to `MPI_INIT`; this routine must be called before any other MPI routine (apart from `MPI_INITIALIZED`) is called. The version for [ISO C](#) accepts the `argc` and `argv` that are provided by the arguments to `main`:

```
int main(argc, argv)
int argc;
char **argv;
{
    MPI_Init(&argc, &argv);

    /* parse arguments */
    /* main program */

    MPI_Finalize();    /* see below */
}
```

The Fortran version takes only `IERROR`.

An MPI implementation is free to require that the arguments in the C binding must be the arguments to `main`.

Rationale. The command line arguments are provided to `MPI_Init` to allow an MPI implementation to use them in initializing the MPI environment. They are passed by reference to allow an MPI implementation to *provide* them in environments where the command-line arguments are not provided to `main`. (*End of rationale.*)

1 In MPI-1.1, it is explicitly stated that an implementation is allowed to require that the
 2 arguments `argc` and `argv` passed by an application to `MPI_INIT` in C be the same arguments
 3 passed into the application as the arguments to `main`. In MPI-2 implementations are not
 4 allowed to impose this requirement. Conforming implementations of MPI are required to
 5 allow applications to pass `NULL` for both the `argc` and `argv` arguments of `main`. In C++,
 6 there is an alternative binding for `MPI::Init` that does not have these arguments at all.

7
 8 *Rationale.* In some applications, libraries may be making the call to `MPI_Init`, and
 9 may not have access to `argc` and `argv` from `main`. It is anticipated that applications
 10 requiring special information about the environment or information supplied by
 11 `mpiexec` can get that information from environment variables. (*End of rationale.*)

14 MPI_FINALIZE()

```
15 int MPI_Finalize(void)
16
17 MPI_FINALIZE(IERROR)
18     INTEGER IERROR
19
20 void MPI::Finalize()
21
22 void MPI::Finalize()
```

23 This routine cleans up all MPI state. Each process must call `MPI_FINALIZE` before
 24 it exits. Unless there has been a call to `MPI_ABORT`, each process must ensure that all
 25 pending non-blocking communications are (locally) complete before calling `MPI_FINALIZE`.
 26 Further, at the instant at which the last process calls `MPI_FINALIZE`, all pending sends
 27 must be matched by a receive, and all pending receives must be matched by a send.

28 For example, the following program is correct:

```
29
30     Process 0           Process 1
31     -----           -----
32     MPI_Init();        MPI_Init();
33     MPI_Send(dest=1);  MPI_Recv(src=0);
34     MPI_Finalize();    MPI_Finalize();
35
```

36 Without the matching receive, the program is erroneous:

```
37
38     Process 0           Process 1
39     -----           -----
40     MPI_Init();        MPI_Init();
41     MPI_Send (dest=1);
42     MPI_Finalize();    MPI_Finalize();
43
```

44 A successful return from a blocking communication operation or from `MPI_WAIT` or
 45 `MPI_TEST` tells the user that the buffer can be reused and means that the communication
 46 is completed by the user, but does not guarantee that the local process has no more work
 47 to do. A successful return from `MPI_REQUEST_FREE` with a request handle generated by
 48 an `MPI_ISEND` nullifies the handle but provides no assurance of operation completion. The
`MPI_ISEND` is complete only when it is known by some means that a matching receive has

completed. MPI_FINALIZE guarantees that all local actions required by communications the user has completed will, in fact, occur before it returns.

MPI_FINALIZE guarantees nothing about pending communications that have not been completed (completion is assured only by MPI_WAIT, MPI_TEST, or MPI_REQUEST_FREE combined with some other verification of completion).

Example 7.3 This program is correct:

```

rank 0                                rank 1
=====
...
MPI_Isend();                          MPI_Recv();
MPI_Request_free();                   MPI_Barrier();
MPI_Barrier();                        MPI_Finalize();
MPI_Finalize();                       exit();
exit();

```

Example 7.4 This program is erroneous and its behavior is undefined:

```

rank 0                                rank 1
=====
...
MPI_Isend();                          MPI_Recv();
MPI_Request_free();                   MPI_Finalize();
MPI_Finalize();                       exit();
exit();

```

If no MPI_BUFFER_DETACH occurs between an MPI_BSEND (or other buffered send) and MPI_FINALIZE, the MPI_FINALIZE implicitly supplies the MPI_BUFFER_DETACH.

Example 7.5 This program is correct, and after the MPI_Finalize, it is as if the buffer had been detached.

```

rank 0                                rank 1
=====
...
buffer = malloc(1000000);             MPI_Recv();
MPI_Buffer_attach();                  MPI_Finalize();
MPI_Bsend();                          exit();
MPI_Finalize();
free(buffer);
exit();

```

Example 7.6 In this example, MPI_lprobe() must return a FALSE flag. MPI_Test_cancelled() must return a TRUE flag, independent of the relative order of execution of MPI_Cancel() in process 0 and MPI_Finalize() in process 1.

The MPI_lprobe() call is there to make sure the implementation knows that the “tag1” message exists at the destination, without being able to claim that the user knows about it.

```

1
2 rank 0                                rank 1
3 =====
4 MPI_Init();                            MPI_Init();
5 MPI_Isend(tag1);
6 MPI_Barrier();                          MPI_Barrier();
7                                         MPI_Iprobe(tag2);
8 MPI_Barrier();                          MPI_Barrier();
9                                         MPI_Finalize();
10                                        exit();
11 MPI_Cancel();
12 MPI_Wait();
13 MPI_Test_cancelled();
14 MPI_Finalize();
15 exit();
16
17

```

Advice to implementors. An implementation may need to delay the return from MPI_FINALIZE until all potential future message cancellations have been processed. One possible solution is to place a barrier inside MPI_FINALIZE (*End of advice to implementors.*)

Once MPI_FINALIZE returns, no MPI routine (not even MPI_INIT) may be called, except for MPI_GET_VERSION, MPI_INITIALIZED, and the MPI-2 function MPI_FINALIZED. Each process must complete any pending communication it initiated before it calls MPI_FINALIZE. If the call returns, each process may continue local computations, or exit, without participating in further MPI communication with other processes. MPI_FINALIZE is collective over all connected processes. If no processes were spawned, accepted or connected then this means over MPI_COMM_WORLD; otherwise it is collective over the union of all processes that have been and continue to be connected, as explained in Section 9.5.4 on page 304.

Advice to implementors. Even though a process has completed all the communication it initiated, such communication may not yet be completed from the viewpoint of the underlying MPI system. E.g., a blocking send may have completed, even though the data is still buffered at the sender. The MPI implementation must ensure that a process has completed any involvement in MPI communication before MPI_FINALIZE returns. Thus, if a process exits after the call to MPI_FINALIZE, this will not cause an ongoing communication to fail. (*End of advice to implementors.*)

Although it is not required that all processes return from MPI_FINALIZE, it is required that at least process 0 in MPI_COMM_WORLD return, so that users can know that the MPI portion of the computation is over. In addition, in a POSIX environment, they may desire to supply an exit code for each process that returns from MPI_FINALIZE.

Example 7.7 The following illustrates the use of requiring that at least one process return and that it be known that process 0 is one of the processes that return. One wants code like the following to work no matter how many processes return.


```

...
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
...
MPI_Finalize();
if (myrank == 0) {
    resultfile = fopen("outfile","w");
    dump_results(resultfile);
    fclose(resultfile);
}
exit(0);

```

MPI_INITIALIZED(flag)

OUT flag Flag is true if MPI_INIT has been called and false otherwise.

```
int MPI_Initialized(int *flag)
```

```
MPI_INITIALIZED(FLAG, IERROR)
```

```
LOGICAL FLAG
INTEGER IERROR
```

```
bool MPI::Is_initialized()
```

```
bool MPI::Is_initialized()
```

This routine may be used to determine whether MPI_INIT has been called. MPI_INITIALIZED returns true if the calling process has called MPI_INIT. Whether MPI_FINALIZE has been called does not affect the behavior of MPI_INITIALIZED. It is one of the few routines that may be called before MPI_INIT is called.

MPI_ABORT(comm, errorcode)

IN comm communicator of tasks to abort
IN errorcode error code to return to invoking environment

```
int MPI_Abort(MPI_Comm comm, int errorcode)
```

```
MPI_ABORT(COMM, ERRORCODE, IERROR)
```

```
INTEGER COMM, ERRORCODE, IERROR
```

```
void MPI::Comm::Abort(int errorcode)
```

```
void MPI::Comm::Abort(int errorcode)
```

This routine makes a “best attempt” to abort all tasks in the group of comm. This function does not require that the invoking environment take any action with the error code. However, a Unix or POSIX environment should handle this **as a return errorcode from the main program.**

1 It may not be possible for an MPI implementation to abort only the processes repre-
2 sented by `comm` if this is a subset of the processes. In this case, the MPI implementation
3 should attempt to abort all the connected processes but should not abort any unconnected
4 processes. If no processes were spawned, accepted or connected then this has the effect of
5 aborting all the processes associated with `MPI_COMM_WORLD`.

6
7 *Rationale.* The communicator argument is provided to allow for future extensions of
8 MPI to environments with, for example, dynamic process management. In particular,
9 it allows but does not require an MPI implementation to abort a subset of
10 `MPI_COMM_WORLD`. (*End of rationale.*)

11
12
13 *Advice to users.* Whether the errorcode is returned from the executable or from the
14 MPI process startup mechanism (e.g., `mpixec`), is an aspect of quality of the MPI
15 library but not mandatory. (*End of advice to users.*)

16
17 *Advice to implementors.* Where possible, a high quality implementation will try
18 to return the errorcode from the MPI process startup mechanism (e.g. `mpixec` or
19 `singleton init`). (*End of advice to implementors.*)

20 21 22 7.6.1 Allowing User Functions at Process Termination

23 There are times in which it would be convenient to have actions happen when an MPI process
24 finishes. For example, a routine may do initializations that are useful until the MPI job (or
25 that part of the job that being terminated in the case of dynamically created processes) is
26 finished. This can be accomplished in MPI-2 by attaching an attribute to `MPI_COMM_SELF`
27 with a callback function. When `MPI_FINALIZE` is called, it will first execute the equivalent
28 of an `MPI_COMM_FREE` on `MPI_COMM_SELF`. This will cause the delete callback function
29 to be executed on all keys associated with `MPI_COMM_SELF`, in an arbitrary order. If no
30 key has been attached to `MPI_COMM_SELF`, then no callback is invoked. The “freeing” of
31 `MPI_COMM_SELF` occurs before any other parts of MPI are affected. Thus, for example,
32 calling `MPI_FINALIZED` will return `false` in any of these callback functions. Once done with
33 `MPI_COMM_SELF`, the order and rest of the actions taken by `MPI_FINALIZE` is not specified.

34
35 *Advice to implementors.* Since attributes can be added from any supported language,
36 the MPI implementation needs to remember the creating language so the correct
37 callback is made. (*End of advice to implementors.*)

38 39 40 7.6.2 Determining Whether MPI Has Finished

41 One of the goals of MPI was to allow for layered libraries. In order for a library to do this
42 cleanly, it needs to know if MPI is active. In MPI-1 the function `MPI_INITIALIZED` was
43 provided to tell if MPI had been initialized. The problem arises in knowing if MPI has been
44 finalized. Once MPI has been finalized it is no longer active and cannot be restarted. A
45 library needs to be able to determine this to act accordingly. To achieve this the following
46 function is needed:
47
48

MPI_FINALIZED(flag) 1

OUT flag true if MPI was finalized (logical) 2

int MPI_Finalized(int *flag) 4

MPI_FINALIZED(FLAG, IERROR) 6

LOGICAL FLAG 7

INTEGER IERROR 8

bool MPI::Is_finalized() 10

This routine returns true if MPI_FINALIZE has completed. It is legal to call MPI_FINALIZED before MPI_INIT and after MPI_FINALIZE. 11

Advice to users. MPI is “active” and it is thus safe to call MPI functions if MPI_INIT has completed and MPI_FINALIZE has not completed. If a library has no other way of knowing whether MPI is active or not, then it can use MPI_INITIALIZED and MPI_FINALIZED to determine this. For example, MPI is “active” in callback functions that are invoked during MPI_FINALIZE. (*End of advice to users.*) 14

7.7 Portable MPI Process Startup 21

A number of implementations of MPI-1 provide a startup command for MPI programs that is of the form 23

```
mpirun <mpirun arguments> <program> <program arguments>
```

26

Separating the command to start the program from the program itself provides flexibility, particularly for network and heterogeneous implementations. For example, the startup script need not run on one of the machines that will be executing the MPI program itself. 28

Having a standard startup mechanism also extends the portability of MPI programs one step further, to the command lines and scripts that manage them. For example, a validation suite script that runs hundreds of programs can be a portable script if it is written using such a standard startup mechanism. In order that the “standard” command not be confused with existing practice, which is not standard and not portable among implementations, instead of `mpirun` MPI specifies `mpiexec`. 31

While a standardized startup mechanism improves the usability of MPI, the range of environments is so diverse (e.g., there may not even be a command line interface) that MPI cannot mandate such a mechanism. Instead, MPI specifies an `mpiexec` startup command and recommends but does not require it, as advice to implementors. However, if an implementation does provide a command called `mpiexec`, it must be of the form described below. 37

It is suggested that 42

```
mpiexec -n <numprocs> <program>
```

44

be at least one way to start `<program>` with an initial MPI_COMM_WORLD whose group contains `<numprocs>` processes. Other arguments to `mpiexec` may be implementation-dependent. 46

Advice to implementors. Implementors, if they do provide a special startup command for MPI programs, are advised to give it the following form. The syntax is chosen in order that `mpiexec` be able to be viewed as a command-line version of `MPI_COMM_SPAWN` (See Section 9.3.4).

Analogous to `MPI_COMM_SPAWN`, we have

```

mpiexec -n    <maxprocs>
           -soft <      >
           -host <      >
           -arch <      >
           -wdir <      >
           -path <      >
           -file <      >
           ...
           <command line>

```

for the case where a single command line for the application program and its arguments will suffice. See Section 9.3.4 for the meanings of these arguments. For the case corresponding to `MPI_COMM_SPAWN_MULTIPLE` there are two possible formats:

Form A:

```

mpiexec { <above arguments> } : { ... } : { ... } : ... : { ... }

```

As with `MPI_COMM_SPAWN`, all the arguments are optional. (Even the `-n x` argument is optional; the default is implementation dependent. It might be 1, it might be taken from an environment variable, or it might be specified at compile time.) The names and meanings of the arguments are taken from the keys in the `info` argument to `MPI_COMM_SPAWN`. There may be other, implementation-dependent arguments as well.

Note that Form A, though convenient to type, prevents colons from being program arguments. Therefore an alternate, file-based form is allowed:

Form B:

```

mpiexec -configfile <filename>

```

where the lines of `<filename>` are of the form separated by the colons in Form A. Lines beginning with `#` are comments, and lines may be continued by terminating the partial line with `\`.

Example 7.8 Start 16 instances of `myprog` on the current or default machine:

```

mpiexec -n 16 myprog

```

Example 7.9 Start 10 processes on the machine called `ferrari`:

```

mpiexec -n 10 -host ferrari myprog

```

Example 7.10 Start three copies of the same program with different command-line arguments:

```
mpiexec myprog infile1 : myprog infile2 : myprog infile3
```

Example 7.11 Start the `ocean` program on five Suns and the `atmos` program on 10 RS/6000's:

```
mpiexec -n 5 -arch sun ocean : -n 10 -arch rs6000 atmos
```

It is assumed that the implementation in this case has a method for choosing hosts of the appropriate type. Their ranks are in the order specified.

Example 7.12 Start the `ocean` program on five Suns and the `atmos` program on 10 RS/6000's (Form B):

```
mpiexec -configfile myfile
```

where `myfile` contains

```
-n 5 -arch sun    ocean
-n 10 -arch rs6000 atmos
```

(End of advice to implementors.)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 8

Miscellany

This chapter contains topics that do not fit conveniently into other chapters.

8.1 The Info Object

Many of the routines in MPI-2 take an argument `info`. `info` is an opaque object with a handle of type `MPI_Info` in C, `MPI::Info` in C++, and `INTEGER` in Fortran. It stores an unordered set of (key,value) pairs (both key and value are strings). A key can have only one value. MPI reserves several keys and requires that if an implementation uses a reserved key, it must provide the specified functionality. An implementation is not required to support these keys and may support any others not reserved by MPI.

An implementation must support info objects as caches for arbitrary (key, value) pairs, regardless of whether it recognizes the key. Each function that takes hints in the form of an `MPI_Info` must be prepared to ignore any key it does not recognize. This description of info objects does not attempt to define how a particular function should react if it recognizes a key but not the associated value. `MPI_INFO_GET_NKEYS`, `MPI_INFO_GET_NTHKEY`, `MPI_INFO_GET_VALUELEN`, and `MPI_INFO_GET` must retain all (key,value) pairs so that layered functionality can also use the `Info` object.

Keys have an implementation-defined maximum length of `MPI_MAX_INFO_KEY`, which is at least 32 and at most 255. Values have an implementation-defined maximum length of `MPI_MAX_INFO_VAL`. In Fortran, leading and trailing spaces are stripped from both. Returned values will never be larger than these maximum lengths. Both key and value are case sensitive.

Rationale. Keys have a maximum length because the set of known keys will always be finite and known to the implementation and because there is no reason for keys to be complex. The small maximum size allows applications to declare keys of size `MPI_MAX_INFO_KEY`. The limitation on value sizes is so that an implementation is not forced to deal with arbitrarily long strings. (*End of rationale.*)

Advice to users. `MPI_MAX_INFO_VAL` might be very large, so it might not be wise to declare a string of that size. (*End of advice to users.*)

When it is an argument to a non-blocking routine, `info` is parsed before that routine returns, so that it may be modified or freed immediately after return.

1 When the descriptions refer to a key or value as being a boolean, an integer, or a list,
 2 they mean the string representation of these types. An implementation may define its own
 3 rules for how info value strings are converted to other types, but to ensure portability, every
 4 implementation must support the following representations. Legal values for a boolean must
 5 include the strings “true” and “false” (all lowercase). For integers, legal values must include
 6 string representations of decimal values of integers that are within the range of a standard
 7 integer type in the program. (However it is possible that not every legal integer is a legal
 8 value for a given key.) On positive numbers, + signs are optional. No space may appear
 9 between a + or – sign and the leading digit of a number. For comma separated lists, the
 10 string must contain legal elements separated by commas. Leading and trailing spaces are
 11 stripped automatically from the types of info values described above and for each element of
 12 a comma separated list. These rules apply to all info values of these types. Implementations
 13 are free to specify a different interpretation for values of other info keys.

14
 15 **MPI_INFO_CREATE**(info)

16 OUT info info object created (handle)

17
 18
 19 `int MPI_Info_create(MPI_Info *info)`

20
 21 **MPI_INFO_CREATE**(INFO, IERROR)
 22 INTEGER INFO, IERROR

23 `static MPI::Info MPI::Info::Create()`

24
 25 **MPI_INFO_CREATE** creates a new info object. The newly created object contains no
 26 key/value pairs.

27
 28 **MPI_INFO_SET**(info, key, value)

29 INOUT info info object (handle)

30 IN key key (string)

31 IN value value (string)

32
 33
 34
 35 `int MPI_Info_set(MPI_Info info, char *key, char *value)`

36
 37 **MPI_INFO_SET**(INFO, KEY, VALUE, IERROR)
 38 INTEGER INFO, IERROR
 39 CHARACTER*(*) KEY, VALUE

40 `void MPI::Info::Set(const char* key, const char* value)`

41
 42 **MPI_INFO_SET** adds the (key,value) pair to info, and overrides the value if a value for
 43 the same key was previously set. key and value are null-terminated strings in C. In Fortran,
 44 leading and trailing spaces in key and value are stripped. If either key or value are larger than
 45 the allowed maximums, the errors **MPIERR.INFO_KEY** or **MPIERR.INFO_VALUE** are raised,
 46 respectively.

MPI_INFO_DELETE(info, key) 1

INOUT info info object (handle) 2

IN key key (string) 3

int MPI_Info_delete(MPI_Info info, char *key) 4

MPI_INFO_DELETE(INFO, KEY, IERROR) 5

INTEGER INFO, IERROR 6

CHARACTER*(*) KEY 7

void MPI::Info::Delete(const char* key) 8

MPI_INFO_DELETE deletes a (key,value) pair from info. If key is not defined in info, the call raises an error of class MPI_ERR_INFO_NOKEY. 9

MPI_INFO_GET(info, key, valuelen, value, flag) 10

IN info info object (handle) 11

IN key key (string) 12

IN valuelen length of value arg (integer) 13

OUT value value (string) 14

OUT flag true if key defined, false if not (boolean) 15

int MPI_Info_get(MPI_Info info, char *key, int valuelen, char *value, int *flag) 16

MPI_INFO_GET(INFO, KEY, VALUELEN, VALUE, FLAG, IERROR) 17

INTEGER INFO, VALUELEN, IERROR 18

CHARACTER*(*) KEY, VALUE 19

LOGICAL FLAG 20

bool MPI::Info::Get(const char* key, int valuelen, char* value) const 21

This function retrieves the value associated with key in a previous call to MPI_INFO_SET. If such a key exists, it sets flag to true and returns the value in value, otherwise it sets flag to false and leaves value unchanged. valuelen is the number of characters available in value. If it is less than the actual size of the value, the value is truncated. In C, valuelen should be one less than the amount of allocated space to allow for the null terminator. 22

If key is larger than MPI_MAX_INFO_KEY, the call is erroneous. 23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

```

1 MPI_INFO_GET_VALUELEN(info, key, valuelen, flag)
2     IN      info                info object (handle)
3
4     IN      key                  key (string)
5
6     OUT     valuelen             length of value arg (integer)
7
8     OUT     flag                  true if key defined, false if not (boolean)

```

```

9 int MPI_Info_get_valuelen(MPI_Info info, char *key, int *valuelen,
10                          int *flag)

```

```

11 MPI_INFO_GET_VALUELEN(INFO, KEY, VALUELEN, FLAG, IERROR)
12     INTEGER INFO, VALUELEN, IERROR
13     LOGICAL FLAG
14     CHARACTER*(*) KEY

```

```

15 bool MPI::Info::Get_valuelen(const char* key, int& valuelen) const

```

Retrieves the length of the value associated with key. If key is defined, valuelen is set to the length of its associated value and flag is set to true. If key is not defined, valuelen is not touched and flag is set to false. The length returned in C or C++ does not include the end-of-string character.

If key is larger than MPI_MAX_INFO_KEY, the call is erroneous.

```

23
24 MPI_INFO_GET_NKEYS(info, nkeys)

```

```

25     IN      info                info object (handle)
26
27     OUT     nkeys                number of defined keys (integer)

```

```

28
29 int MPI_Info_get_nkeys(MPI_Info info, int *nkeys)

```

```

30 MPI_INFO_GET_NKEYS(INFO, NKEYS, IERROR)
31     INTEGER INFO, NKEYS, IERROR

```

```

32
33 int MPI::Info::Get_nkeys() const

```

MPI_INFO_GET_NKEYS returns the number of currently defined keys in info.

```

36
37 MPI_INFO_GET_NTHKEY(info, n, key)

```

```

38     IN      info                info object (handle)
39
40     IN      n                    key number (integer)
41
42     OUT     key                  key (string)

```

```

43
44 int MPI_Info_get_nthkey(MPI_Info info, int n, char *key)

```

```

45 MPI_INFO_GET_NTHKEY(INFO, N, KEY, IERROR)
46     INTEGER INFO, N, IERROR
47     CHARACTER*(*) KEY

```

```

48

```

```
void MPI::Info::Get_nthkey(int n, char* key) const
```

This function returns the n th defined key in `info`. Keys are numbered $0 \dots N - 1$ where N is the value returned by `MPI_INFO_GET_NKEYS`. All keys between 0 and $N - 1$ are guaranteed to be defined. The number of a given key does not change as long as `info` is not modified with `MPI_INFO_SET` or `MPI_INFO_DELETE`.

```
MPI_INFO_DUP(info, newinfo)
```

| | | |
|-----|---------|----------------------|
| IN | info | info object (handle) |
| OUT | newinfo | info object (handle) |

```
int MPI_Info_dup(MPI_Info info, MPI_Info *newinfo)
```

```
MPI_INFO_DUP(INFO, NEWINFO, IERROR)
    INTEGER INFO, NEWINFO, IERROR
```

```
MPI::Info MPI::Info::Dup() const
```

`MPI_INFO_DUP` duplicates an existing `info` object, creating a new object, with the same (key,value) pairs and the same ordering of keys.

```
MPI_INFO_FREE(info)
```

| | | |
|-------|------|----------------------|
| INOUT | info | info object (handle) |
|-------|------|----------------------|

```
int MPI_Info_free(MPI_Info *info)
```

```
MPI_INFO_FREE(INFO, IERROR)
    INTEGER INFO, IERROR
```

```
void MPI::Info::Free()
```

This function frees `info` and sets it to `MPI_INFO_NULL`. The value of an `info` argument is interpreted each time the `info` is passed to a routine. Changes to an `info` after return from a routine do not affect that interpretation.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 9

Process Creation and Management

9.1 Introduction

MPI-1 provides an interface that allows processes in a parallel program to communicate with one another. MPI-1 specifies neither how the processes are created, nor how they establish communication. Moreover, an MPI-1 application is static; that is, no processes can be added to or deleted from an application after it has been started.

MPI users have asked that the MPI-1 model be extended to allow process creation and management after an MPI application has been started. A major impetus comes from the PVM [25] research effort, which has provided a wealth of experience with process management and resource control that illustrates their benefits and potential pitfalls.

The MPI Forum decided not to address resource control in MPI-2 because it was not able to design a portable interface that would be appropriate for the broad spectrum of existing and potential resource and process controllers. Resource control can encompass a wide range of abilities, including adding and deleting nodes from a virtual parallel machine, reserving and scheduling resources, managing compute partitions of an MPP, and returning information about available resources. MPI-2 assumes that resource control is provided externally — probably by computer vendors, in the case of tightly coupled systems, or by a third party software package when the environment is a cluster of workstations.

The reasons for adding process management to MPI are both technical and practical. Important classes of message passing applications require process control. These include task farms, serial applications with parallel modules, and problems that require a run-time assessment of the number and type of processes that should be started. On the practical side, users of workstation clusters who are migrating from PVM to MPI may be accustomed to using PVM's capabilities for process and resource management. The lack of these features is a practical stumbling block to migration.

While process management is essential, adding it to MPI should not compromise the portability or performance of MPI applications. In particular:

- The MPI-2 process model must apply to the vast majority of current parallel environments. These include everything from tightly integrated MPPs to heterogeneous networks of workstations.
- MPI must not take over operating system responsibilities. It should instead provide a clean interface between an application and system software.

- 1 • MPI must continue to guarantee communication determinism, i.e., process manage-
2 ment must not introduce unavoidable race conditions.
- 3
- 4 • MPI must not contain features that compromise performance.
- 5
- 6 • MPI-1 programs must work under MPI-2, i.e., the MPI-1 static process model must be
7 a special case of the MPI-2 dynamic model.

8 The MPI-2 process management model addresses these issues in two ways. First, MPI
9 remains primarily a communication library. It does not manage the parallel environment
10 in which a parallel program executes, though it provides a minimal interface between an
11 application and external resource and process managers.

12 Second, MPI-2 does not change the concept of communicator. Once a communicator
13 is built, it behaves as specified in MPI-1. A communicator is never changed once created,
14 and it is always created using deterministic collective operations.

16 9.2 The MPI-2 Process Model

18 The MPI-2 process model allows for the creation and cooperative termination of processes
19 after an MPI application has started. It provides a mechanism to establish communication
20 between the newly created processes and the existing MPI application. It also provides a
21 mechanism to establish communication between two existing MPI applications, even when
22 one did not “start” the other.

24 9.2.1 Starting Processes

26 MPI applications may start new processes through an interface to an external process man-
27 ager, which can range from a parallel operating system (CMOST) to layered software (POE)
28 to an `rsh` command (p4).

29 `MPI_COMM_SPAWN` starts MPI processes and establishes communication with them,
30 returning an intercommunicator. `MPI_COMM_SPAWN_MULTIPLE` starts several different
31 binaries (or the same binary with different arguments), placing them in the same
32 `MPI_COMM_WORLD` and returning an intercommunicator.

33 MPI uses the existing group abstraction to represent processes. A process is identified
34 by a (group, rank) pair.

36 9.2.2 The Runtime Environment

38 The `MPI_COMM_SPAWN` and `MPI_COMM_SPAWN_MULTIPLE` routines provide an interface
39 between MPI and the *runtime environment* of an MPI application. The difficulty is that
40 there is an enormous range of runtime environments and application requirements, and MPI
41 must not be tailored to any particular one. Examples of such environments are:

- 42
- 43 • **MPP managed by a batch queueing system.** Batch queueing systems generally
44 allocate resources before an application begins, enforce limits on resource use (CPU
45 time, memory use, etc.), and do not allow a change in resource allocation after a
46 job begins. Moreover, many MPPs have special limitations or extensions, such as a
47 limit on the number of processes that may run on one processor, or the ability to
48 gang-schedule processes of a parallel application.

- **Network of workstations with PVM.** PVM (Parallel Virtual Machine) allows a user to create a “virtual machine” out of a network of workstations. An application may extend the virtual machine or manage processes (create, kill, redirect output, etc.) through the PVM library. Requests to manage the machine or processes may be intercepted and handled by an external resource manager.
- **Network of workstations managed by a load balancing system.** A load balancing system may choose the location of spawned processes based on dynamic quantities, such as load average. It may transparently migrate processes from one machine to another when a resource becomes unavailable.
- **Large SMP with Unix.** Applications are run directly by the user. They are scheduled at a low level by the operating system. Processes may have special scheduling characteristics (gang-scheduling, processor affinity, deadline scheduling, processor locking, etc.) and be subject to OS resource limits (number of processes, amount of memory, etc.).

MPI assumes, implicitly, the existence of an environment in which an application runs. It does not provide “operating system” services, such as a general ability to query what processes are running, to kill arbitrary processes, to find out properties of the runtime environment (how many processors, how much memory, etc.).

Complex interaction of an MPI application with its runtime environment should be done through an environment-specific API. An example of such an API would be the PVM task and machine management routines — `pvm_addhosts`, `pvm_config`, `pvm_tasks`, etc., possibly modified to return an MPI (group,rank) when possible. A Condor or PBS API would be another possibility.

At some low level, obviously, MPI must be able to interact with the runtime system, but the interaction is not visible at the application level and the details of the interaction are not specified by the MPI standard.

In many cases, it is impossible to keep environment-specific information out of the MPI interface without seriously compromising MPI functionality. To permit applications to take advantage of environment-specific functionality, many MPI routines take an `info` argument that allows an application to specify environment-specific information. There is a tradeoff between functionality and portability: applications that make use of `info` are not portable.

MPI does not require the existence of an underlying “virtual machine” model, in which there is a consistent global view of an MPI application and an implicit “operating system” managing resources and processes. For instance, processes spawned by one task may not be visible to another; additional hosts added to the runtime environment by one process may not be visible in another process; tasks spawned by different processes may not be automatically distributed over available resources.

Interaction between MPI and the runtime environment is limited to the following areas:

- A process may start new processes with `MPI_COMM_SPAWN` and `MPI_COMM_SPAWN_MULTIPLE`.
- When a process spawns a child process, it may optionally use an `info` argument to tell the runtime environment where or how to start the process. This extra information may be opaque to MPI.

- An attribute `MPI_UNIVERSE_SIZE` on `MPI_COMM_WORLD` tells a program how “large” the initial runtime environment is, namely how many processes can usefully be started in all. One can subtract the size of `MPI_COMM_WORLD` from this value to find out how many processes might usefully be started in addition to those already running.

9.3 Process Manager Interface

9.3.1 Processes in MPI

A process is represented in MPI by a (group, rank) pair. A (group, rank) pair specifies a unique process but a process does not determine a unique (group, rank) pair, since a process may belong to several groups.

9.3.2 Starting Processes and Establishing Communication

The following routine starts a number of MPI processes and establishes communication with them, returning an intercommunicator.

Advice to users. It is possible in MPI to start a static SPMD or MPMD application by starting first one process and having that process start its siblings with `MPI_COMM_SPAWN`. This practice is discouraged primarily for reasons of performance. If possible, it is preferable to start all processes at once, as a single MPI-1 application. (*End of advice to users.*)

`MPI_COMM_SPAWN`(command, argv, maxprocs, info, root, comm, intercomm, array_of_errcodes)

| | | |
|-----|-------------------|---|
| IN | command | name of program to be spawned (string, significant only at root) |
| IN | argv | arguments to command (array of strings, significant only at root) |
| IN | maxprocs | maximum number of processes to start (integer, significant only at root) |
| IN | info | a set of key-value pairs telling the runtime system where and how to start the processes (handle, significant only at root) |
| IN | root | rank of process in which previous arguments are examined (integer) |
| IN | comm | intracommunicator containing group of spawning processes (handle) |
| OUT | intercomm | intercommunicator between original group and the newly spawned group (handle) |
| OUT | array_of_errcodes | one code per process (array of integer) |

```
int MPI_Comm_spawn(char *command, char *argv[], int maxprocs, MPI_Info info,
                  int root, MPI_Comm comm, MPI_Comm *intercomm,
```



```

        int array_of_errcodes[])
MPI_COMM_SPAWN(COMMAND, ARGV, MAXPROCS, INFO, ROOT, COMM, INTERCOMM,
        ARRAY_OF_ERRCODES, IERROR)
CHARACTER*(*) COMMAND, ARGV(*)
INTEGER INFO, MAXPROCS, ROOT, COMM, INTERCOMM, ARRAY_OF_ERRCODES(*),
IERROR

MPI::Intercomm MPI::Intracomm::Spawn(const char* command,
        const char* argv[], int maxprocs, const MPI::Info& info,
        int root, int array_of_errcodes[]) const

MPI::Intercomm MPI::Intracomm::Spawn(const char* command,
        const char* argv[], int maxprocs, const MPI::Info& info,
        int root) const

```

MPI_COMM_SPAWN tries to start maxprocs identical copies of the MPI program specified by command, establishing communication with them and returning an intercommunicator. The spawned processes are referred to as children. The children have their own MPI_COMM_WORLD, which is separate from that of the parents. MPI_COMM_SPAWN is collective over comm, and also may not return until MPI_INIT has been called in the children. Similarly, MPI_INIT in the children may not return until all parents have called MPI_COMM_SPAWN. In this sense, MPI_COMM_SPAWN in the parents and MPI_INIT in the children form a collective operation over the union of parent and child processes. The intercommunicator returned by MPI_COMM_SPAWN contains the parent processes in the local group and the child processes in the remote group. The ordering of processes in the local and remote groups is the same as the ordering of the group of the comm in the parents and of MPI_COMM_WORLD of the children, respectively. This intercommunicator can be obtained in the children through the function MPI_COMM_GET_PARENT.

Advice to users. An implementation may automatically establish communication before MPI_INIT is called by the children. Thus, completion of MPI_COMM_SPAWN in the parent does not necessarily mean that MPI_INIT has been called in the children (although the returned intercommunicator can be used immediately). (*End of advice to users.*)

The command argument The command argument is a string containing the name of a program to be spawned. The string is null-terminated in C. In Fortran, leading and trailing spaces are stripped. MPI does not specify how to find the executable or how the working directory is determined. These rules are implementation-dependent and should be appropriate for the runtime environment.

Advice to implementors. The implementation should use a natural rule for finding executables and determining working directories. For instance, a homogeneous system with a global file system might look first in the working directory of the spawning process, or might search the directories in a PATH environment variable as do Unix shells. An implementation on top of PVM would use PVM's rules for finding executables (usually in \$HOME/pvm3/bin/\$PVM_ARCH). An MPI implementation running under POE on an IBM SP would use POE's method of finding executables. An implementation should document its rules for finding executables and determining working

1 directories, and a high-quality implementation should give the user some control over
 2 these rules. (*End of advice to implementors.*)

3
 4 If the program named in `command` does not call `MPI_INIT`, but instead forks a process
 5 that calls `MPI_INIT`, the results are undefined. Implementations may allow this case to work
 6 but are not required to.

7
 8 *Advice to users.* MPI does not say what happens if the program you start is a
 9 shell script and that shell script starts a program that calls `MPI_INIT`. Though some
 10 implementations may allow you to do this, they may also have restrictions, such as
 11 requiring that arguments supplied to the shell script be supplied to the program, or
 12 requiring that certain parts of the environment not be changed. (*End of advice to*
 13 *users.*)

14
 15 The `argv` argument `argv` is an array of strings containing arguments that are passed to
 16 the program. The first element of `argv` is the first argument passed to `command`, not, as
 17 is conventional in some contexts, the command itself. The argument list is terminated by
 18 `NULL` in C and C++ and an empty string in Fortran. In Fortran, leading and trailing spaces
 19 are always stripped, so that a string consisting of all spaces is considered an empty string.
 20 The constant `MPI_ARGV_NULL` may be used in C, C++ and Fortran to indicate an empty
 21 argument list. In C and C++, this constant is the same as `NULL`.

22 **Example 9.1** Examples of `argv` in C and Fortran

23 To run the program “ocean” with arguments “-gridfile” and “ocean1.grd” in C:

```
24
25     char command[] = "ocean";
26     char *argv[] = {"-gridfile", "ocean1.grd", NULL};
27     MPI_Comm_spawn(command, argv, ...);
28
```

29 or, if not everything is known at compile time:

```
30
31     char *command;
32     char **argv;
33     command = "ocean";
34     argv=(char **)malloc(3 * sizeof(char *));
35     argv[0] = "-gridfile";
36     argv[1] = "ocean1.grd";
37     argv[2] = NULL;
38     MPI_Comm_spawn(command, argv, ...);
39
```

40 In Fortran:

```
41     CHARACTER*25 command, argv(3)
42     command = ' ocean '
43     argv(1) = ' -gridfile '
44     argv(2) = ' ocean1.grd'
45     argv(3) = ' '
46     call MPI_COMM_SPAWN(command, argv, ...)
47
48
```

Arguments are supplied to the program if this is allowed by the operating system. In C, the `MPI_COMM_SPAWN` argument `argv` differs from the `argv` argument of `main` in two respects. First, it is shifted by one element. Specifically, `argv[0]` of `main` is provided by the implementation and conventionally contains the name of the program (`argv[1]` of `main` corresponds to `argv[0]` in `MPI_COMM_SPAWN`, `argv[2]` of `main` to `argv[1]` of `MPI_COMM_SPAWN`, etc. Second, `argv` of `MPI_COMM_SPAWN` must be null-terminated, so that its length can be determined. Passing an `argv` of `MPI_ARGV_NULL` to `MPI_COMM_SPAWN` results in `main` receiving `argc` of 1 and an `argv` whose element 0 is (conventionally) the name of the program.

If a Fortran implementation supplies routines that allow a program to obtain its arguments, the arguments may be available through that mechanism. In C, if the operating system does not support arguments appearing in `argv` of `main()`, the MPI implementation may add the arguments to the `argv` that is passed to `MPI_INIT`.

The `maxprocs` argument MPI tries to spawn `maxprocs` processes. If it is unable to spawn `maxprocs` processes, it raises an error of class `MPI_ERR_SPAWN`.

An implementation may allow the `info` argument to change the default behavior, such that if the implementation is unable to spawn all `maxprocs` processes, it may spawn a smaller number of processes instead of raising an error. In principle, the `info` argument may specify an arbitrary set $\{m_i : 0 \leq m_i \leq \text{maxprocs}\}$ of allowed values for the number of processes spawned. The set $\{m_i\}$ does not necessarily include the value `maxprocs`. If an implementation is able to spawn one of these allowed numbers of processes, `MPI_COMM_SPAWN` returns successfully and the number of spawned processes, m , is given by the size of the remote group of `intercomm`. If m is less than `maxproc`, reasons why the other processes were not spawned are given in `array_of_errcodes` as described below. If it is not possible to spawn one of the allowed numbers of processes, `MPI_COMM_SPAWN` raises an error of class `MPI_ERR_SPAWN`.

A spawn call with the default behavior is called *hard*. A spawn call for which fewer than `maxprocs` processes may be returned is called *soft*. See Section 9.3.4 on page 289 for more information on the soft key for `info`.

Advice to users. By default, requests are hard and MPI errors are fatal. This means that by default there will be a fatal error if MPI cannot spawn all the requested processes. If you want the behavior “spawn as many processes as possible, up to N ,” you should do a soft spawn, where the set of allowed values $\{m_i\}$ is $\{0 \dots N\}$. However, this is not completely portable, as implementations are not required to support soft spawning. (*End of advice to users.*)

The `info` argument The `info` argument to all of the routines in this chapter is an opaque handle of type `MPI_Info` in C, `MPI::Info` in C++ and `INTEGER` in Fortran. It is a container for a number of user-specified (key,value) pairs. `key` and `value` are strings (null-terminated `char*` in C, `character*(*)` in Fortran). Routines to create and manipulate the `info` argument are described in Section 8.1 on page 273.

For the `SPAWN` calls, `info` provides additional (and possibly implementation-dependent) instructions to MPI and the runtime system on how to start processes. An application may pass `MPI_INFO_NULL` in C or Fortran. Portable programs not requiring detailed control over process locations should use `MPI_INFO_NULL`.

MPI does not specify the content of the `info` argument, except to reserve a number of special key values (see Section 9.3.4 on page 289). The `info` argument is quite flexible and could even be used, for example, to specify the executable and its command-line arguments. In this case the `command` argument to `MPI_COMM_SPAWN` could be empty. The ability to do this follows from the fact that MPI does not specify how an executable is found, and the `info` argument can tell the runtime system where to “find” the executable “” (empty string). Of course a program that does this will not be portable across MPI implementations.

The root argument All arguments before the `root` argument are examined only on the process whose rank in `comm` is equal to `root`. The value of these arguments on other processes is ignored.

The `array_of_errcodes` argument The `array_of_errcodes` is an array of length `maxprocs` in which MPI reports the status of each process that MPI was requested to start. If all `maxprocs` processes were spawned, `array_of_errcodes` is filled in with the value `MPI_SUCCESS`. If only m ($0 \leq m < \text{maxprocs}$) processes are spawned, m of the entries will contain `MPI_SUCCESS` and the rest will contain an implementation-specific error code indicating the reason MPI could not start the process. MPI does not specify which entries correspond to failed processes. An implementation may, for instance, fill in error codes in one-to-one correspondence with a detailed specification in the `info` argument. These error codes all belong to the error class `MPI_ERR_SPAWN` if there was no error in the argument list. In C or Fortran, an application may pass `MPI_ERRCODES_IGNORE` if it is not interested in the error codes. In C++ this constant does not exist, and the `array_of_errcodes` argument may be omitted from the argument list.

Advice to implementors. `MPI_ERRCODES_IGNORE` in Fortran is a special type of constant, like `MPI_BOTTOM`. See the discussion in Section 2.5.4 on page 14. (*End of advice to implementors.*)

`MPI_COMM_GET_PARENT(parent)`

OUT parent the parent communicator (handle)

`int MPI_Comm_get_parent(MPI_Comm *parent)`

`MPI_COMM_GET_PARENT(PARENT, IERROR)`

INTEGER PARENT, IERROR

`static MPI::Intercomm MPI::Comm::Get_parent()`

If a process was started with `MPI_COMM_SPAWN` or `MPI_COMM_SPAWN_MULTIPLE`, `MPI_COMM_GET_PARENT` returns the “parent” intercommunicator of the current process. This parent intercommunicator is created implicitly inside of `MPI_INIT` and is the same intercommunicator returned by `SPAWN` in the parents.

If the process was not spawned, `MPI_COMM_GET_PARENT` returns `MPI_COMM_NULL`.

After the parent communicator is freed or disconnected, `MPI_COMM_GET_PARENT` returns `MPI_COMM_NULL`.

Advice to users. MPI_COMM_GET_PARENT returns a handle to a single intercommunicator. Calling MPI_COMM_GET_PARENT a second time returns a handle to the same intercommunicator. Freeing the handle with MPI_COMM_DISCONNECT or MPI_COMM_FREE will cause other references to the intercommunicator to become invalid (dangling). Note that calling MPI_COMM_FREE on the parent communicator is not useful. (*End of advice to users.*)

Rationale. The desire of the Forum was to create a constant MPI_COMM_PARENT similar to MPI_COMM_WORLD. Unfortunately such a constant cannot be used (syntactically) as an argument to MPI_COMM_DISCONNECT, which is explicitly allowed. (*End of rationale.*)

9.3.3 Starting Multiple Executables and Establishing Communication

While MPI_COMM_SPAWN is sufficient for most cases, it does not allow the spawning of multiple binaries, or of the same binary with multiple sets of arguments. The following routine spawns multiple binaries or the same binary with multiple sets of arguments, establishing communication with them and placing them in the same MPI_COMM_WORLD.

MPI_COMM_SPAWN_MULTIPLE(count, array_of_commands, array_of_argv, array_of_maxprocs, array_of_info, root, comm, intercomm, array_of_errcodes)

| | | |
|-----|-------------------|---|
| IN | count | number of commands (positive integer, significant to MPI only at root — see advice to users) |
| IN | array_of_commands | programs to be executed (array of strings, significant only at root) |
| IN | array_of_argv | arguments for commands (array of array of strings, significant only at root) |
| IN | array_of_maxprocs | maximum number of processes to start for each command (array of integer, significant only at root) |
| IN | array_of_info | info objects telling the runtime system where and how to start processes (array of handles, significant only at root) |
| IN | root | rank of process in which previous arguments are examined (integer) |
| IN | comm | intracommunicator containing group of spawning processes (handle) |
| OUT | intercomm | intercommunicator between original group and newly spawned group (handle) |
| OUT | array_of_errcodes | one error code per process (array of integer) |

```
int MPI_Comm_spawn_multiple(int count, char *array_of_commands[],
                           char **array_of_argv[], int array_of_maxprocs[],
                           MPI_Info array_of_info[], int root, MPI_Comm comm,
                           MPI_Comm *intercomm, int array_of_errcodes[])
```

```

1 MPI_COMM Spawn Multiple(COUNT, ARRAY_OF_COMMANDS, ARRAY_OF_ARGV,
2     ARRAY_OF_MAXPROCS, ARRAY_OF_INFO, ROOT, COMM, INTERCOMM,
3     ARRAY_OF_ERRCODES, IERROR)
4     INTEGER COUNT, ARRAY_OF_INFO(*), ARRAY_OF_MAXPROCS(*), ROOT, COMM,
5     INTERCOMM, ARRAY_OF_ERRCODES(*), IERROR
6     CHARACTER*(*) ARRAY_OF_COMMANDS(*), ARRAY_OF_ARGV(COUNT, *)
7
8 MPI::Intercomm MPI::Intracomm::Spawn_multiple(int count,
9     const char* array_of_commands[], const char** array_of_argv[],
10    const int array_of_maxprocs[], const MPI::Info array_of_info[],
11    int root, int array_of_errcodes[])
12
13 MPI::Intercomm MPI::Intracomm::Spawn_multiple(int count,
14    const char* array_of_commands[], const char** array_of_argv[],
15    const int array_of_maxprocs[], const MPI::Info array_of_info[],
16    int root)

```

MPI_COMM_SPAWN_MULTIPLE is identical to MPI_COMM_SPAWN except that there are multiple executable specifications. The first argument, `count`, gives the number of specifications. Each of the next four arguments are simply arrays of the corresponding arguments in MPI_COMM_SPAWN. For the Fortran version of `array_of_argv`, the element `array_of_argv(i,j)` is the j th argument to command number i .

Rationale. This may seem backwards to Fortran programmers who are familiar with Fortran's column-major ordering. However, it is necessary to do it this way to allow MPI_COMM_SPAWN to sort out arguments. Note that the leading dimension of `array_of_argv` must be the same as `count`. (*End of rationale.*)

Advice to users. The argument `count` is interpreted by MPI only at the root, as is `array_of_argv`. Since the leading dimension of `array_of_argv` is `count`, a non-positive value of `count` at a non-root node could theoretically cause a runtime bounds check error, even though `array_of_argv` should be ignored by the subroutine. If this happens, you should explicitly supply a reasonable value of `count` on the non-root nodes. (*End of advice to users.*)

In any language, an application may use the constant MPI_ARGVS_NULL (which is likely to be `(char ***)0` in C) to specify that no arguments should be passed to any commands. The effect of setting individual elements of `array_of_argv` to MPI_ARGV_NULL is not defined. To specify arguments for some commands but not others, the commands without arguments should have a corresponding `argv` whose first element is null (`(char *)0` in C and empty string in Fortran).

All of the spawned processes have the same MPI_COMM_WORLD. Their ranks in MPI_COMM_WORLD correspond directly to the order in which the commands are specified in MPI_COMM_SPAWN_MULTIPLE. Assume that m_1 processes are generated by the first command, m_2 by the second, etc. The processes corresponding to the first command have ranks $0, 1, \dots, m_1 - 1$. The processes in the second command have ranks $m_1, m_1 + 1, \dots, m_1 + m_2 - 1$. The processes in the third have ranks $m_1 + m_2, m_1 + m_2 + 1, \dots, m_1 + m_2 + m_3 - 1$, etc.

Advice to users. Calling MPI_COMM_SPAWN multiple times would create many sets of children with different MPI_COMM_WORLDS whereas

MPI_COMM_SPAWN_MULTIPLE creates children with a single MPI_COMM_WORLD, so the two methods are not completely equivalent. There are also two performance-related reasons why, if you need to spawn multiple executables, you may want to use MPI_COMM_SPAWN_MULTIPLE instead of calling MPI_COMM_SPAWN several times. First, spawning several things at once may be faster than spawning them sequentially. Second, in some implementations, communication between processes spawned at the same time may be faster than communication between processes spawned separately. (*End of advice to users.*)

The `array_of_errcodes` argument is 1-dimensional array of size $\sum_{i=1}^{count} n_i$, where n_i is the i th element of `array_of_maxprocs`. Command number i corresponds to the n_i contiguous slots in this array from element $\sum_{j=1}^{i-1} n_j$ to $[\sum_{j=1}^i n_j] - 1$. Error codes are treated as for MPI_COMM_SPAWN.

Example 9.2 Examples of `array_of_argv` in C and Fortran

To run the program “ocean” with arguments “-gridfile” and “ocean1.grd” and the program “atmos” with argument “atmos.grd” in C:

```
char *array_of_commands[2] = {"ocean", "atmos"};
char **array_of_argv[2];
char *argv0[] = {"-gridfile", "ocean1.grd", (char *)0};
char *argv1[] = {"atmos.grd", (char *)0};
array_of_argv[0] = argv0;
array_of_argv[1] = argv1;
MPI_Comm_spawn_multiple(2, array_of_commands, array_of_argv, ...);
```

Here’s how you do it in Fortran:

```
CHARACTER*25 commands(2), array_of_argv(2, 3)
commands(1) = ' ocean '
array_of_argv(1, 1) = ' -gridfile '
array_of_argv(1, 2) = ' ocean1.grd'
array_of_argv(1, 3) = ' '

commands(2) = ' atmos '
array_of_argv(2, 1) = ' atmos.grd '
array_of_argv(2, 2) = ' '

call MPI_COMM_SPAWN_MULTIPLE(2, commands, array_of_argv, ...)
```

9.3.4 Reserved Keys

The following keys are reserved. An implementation is not required to interpret these keys, but if it does interpret the key, it must provide the functionality described.

`host` Value is a hostname. The format of the hostname is determined by the implementation.

`arch` Value is an architecture name. Valid architecture names and what they mean are determined by the implementation.

1 wdir Value is the name of a directory on a machine on which the spawned process(es)
 2 execute(s). This directory is made the working directory of the executing process(es).
 3 The format of the directory name is determined by the implementation.

4 path Value is a directory or set of directories where the implementation should look for the
 5 executable. The format of path is determined by the implementation.

7 file Value is the name of a file in which additional information is specified. The format of
 8 the filename and internal format of the file are determined by the implementation.

10 soft Value specifies a set of numbers which are allowed values for the number of processes
 11 that MPI_COMM_SPAWN (et al.) may create. The format of the value is a comma-
 12 separated list of Fortran-90 triplets each of which specifies a set of integers and which
 13 together specify the set formed by the union of these sets. Negative values in this set
 14 and values greater than `maxprocs` are ignored. MPI will spawn the largest number of
 15 processes it can, consistent with some number in the set. The order in which triplets
 16 are given is not significant.

17 By Fortran-90 triplets, we mean:

- 18 1. `a` means a
- 19 2. `a:b` means $a, a + 1, a + 2, \dots, b$
- 20 3. `a:b:c` means $a, a + c, a + 2c, \dots, a + ck$, where for $c > 0$, k is the largest integer
 21 for which $a + ck \leq b$ and for $c < 0$, k is the largest integer for which $a + ck \geq b$.
 22 If $b > a$ then c must be positive. If $b < a$ then c must be negative.

25 Examples:

- 26 1. `a:b` gives a range between a and b
- 27 2. `0:N` gives full “soft” functionality
- 28 3. `1,2,4,8,16,32,64,128,256,512,1024,2048,4096` allows power-of-two number
 29 of processes.
- 30 4. `2:10000:2` allows even number of processes.
- 31 5. `2:10:2,7` allows 2, 4, 6, 7, 8, or 10 processes.

35 9.3.5 Spawn Example

36 Manager-worker Example, Using MPI_SPAWN.

```

38 /* manager */
39 #include "mpi.h"
40 int main(int argc, char *argv[])
41 {
42     int world_size, universe_size, *universe_sizep, flag;
43     MPI_Comm everyone;          /* intercommunicator */
44     char worker_program[100];
45
46     MPI_Init(&argc, &argv);
47     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
48 
```



```
1  /*
2  * Parallel code here.
3  * The manager is represented as the process with rank 0 in (the remote
4  * group of) the parent communicator. If the workers need to communicate
5  * among themselves, they can use MPI_COMM_WORLD.
6  */
7
8  MPI_Finalize();
9  return 0;
10 }
11
12
13
14
```

9.4 Establishing Communication

This section provides functions that establish communication between two sets of MPI processes that do not share a communicator.

Some situations in which these functions are useful are:

1. Two parts of an application that are started independently need to communicate.
2. A visualization tool wants to attach to a running process.
3. A server wants to accept connections from multiple clients. Both clients and server may be parallel programs.

In each of these situations, MPI must establish communication channels where none existed before, and there is no parent/child relationship. The routines described in this section establish communication between the two sets of processes by creating an MPI intercommunicator, where the two groups of the intercommunicator **are the original sets of processes.**

Establishing contact between two groups of processes that do not share an existing communicator is a collective but asymmetric process. One group of processes indicates its willingness to accept connections from other groups of processes. We will call this group the (parallel) *server*, even if this is not a client/server type of application. The other group connects to the server; we will call it the *client*.

Advice to users. While the names *client* and *server* are used throughout this section, MPI does not guarantee the traditional robustness of client server systems. The functionality described in this section is intended to allow two cooperating parts of the same application to communicate with one another. For instance, a client that gets a segmentation fault and dies, or one that doesn't participate in a collective operation may cause a server to crash or hang. (*End of advice to users.*)

9.4.1 Names, Addresses, Ports, and All That

Almost all of the complexity in MPI client/server routines addresses the question “how does the client find out how to contact the server?” The difficulty, of course, is that there

is no existing communication channel between them, yet they must somehow agree on a rendezvous point where they will establish communication — Catch 22.

Agreeing on a rendezvous point always involves a third party. The third party may itself provide the rendezvous point or may communicate rendezvous information from server to client. Complicating matters might be the fact that a client doesn't really care what server it contacts, only that it be able to get in touch with one that can handle its request.

Ideally, MPI can accommodate a wide variety of run-time systems while retaining the ability to write simple portable code. The following should be compatible with MPI:

- The server resides at a well-known internet address `host:port`.
- The server prints out an address to the terminal, the user gives this address to the client program.
- The server places the address information on a nameserver, where it can be retrieved with an agreed-upon name.
- The server to which the client connects is actually a broker, acting as a middleman between the client and the real server.

MPI does not require a nameserver, so not all implementations will be able to support all of the above scenarios. However, MPI provides an optional nameserver interface, and is compatible with external name servers.

A `port_name` is a *system-supplied* string that encodes a low-level network address at which a server can be contacted. Typically this is an IP address and a port number, but an implementation is free to use any protocol. The server establishes a `port_name` with the `MPI_OPEN_PORT` routine. It accepts a connection to a given port with `MPI_COMM_ACCEPT`. A client uses `port_name` to connect to the server.

By itself, the `port_name` mechanism is completely portable, but it may be clumsy to use because of the necessity to communicate `port_name` to the client. It would be more convenient if a server could specify that it be known by an *application-supplied service_name* so that the client could connect to that `service_name` without knowing the `port_name`.

An MPI implementation may allow the server to publish a (`port_name`, `service_name`) pair with `MPI_PUBLISH_NAME` and the client to retrieve the port name from the service name with `MPI_LOOKUP_NAME`. This allows three levels of portability, with increasing levels of functionality.

1. Applications that do not rely on the ability to publish names are the most portable. Typically the `port_name` must be transferred “by hand” from server to client.
2. Applications that use the `MPI_PUBLISH_NAME` mechanism are completely portable among implementations that provide this service. To be portable among all implementations, these applications should have a fall-back mechanism that can be used when names are not published.
3. Applications may ignore MPI's name publishing functionality and use their own mechanism (possibly system-supplied) to publish names. This allows arbitrary flexibility but is not portable.

9.4.2 Server Routines

A server makes itself available with two routines. First it must call `MPI_OPEN_PORT` to establish a port at which it may be contacted. Secondly it must call `MPI_COMM_ACCEPT` to accept connections from clients.

`MPI_OPEN_PORT(info, port_name)`

| | | |
|-----|-----------|---|
| IN | info | implementation-specific information on how to establish an address (handle) |
| OUT | port_name | newly established port (string) |

`int MPI_Open_port(MPI_Info info, char *port_name)`

`MPI_OPEN_PORT(INFO, PORT_NAME, IERROR)`

`CHARACTER*(*) PORT_NAME`

`INTEGER INFO, IERROR`

`void MPI::Open_port(const MPI::Info& info, char* port_name)`

This function establishes a network address, encoded in the `port_name` string, at which the server will be able to accept connections from clients. `port_name` is supplied by the system, possibly using information in the `info` argument.

MPI copies a system-supplied port name into `port_name`. `port_name` identifies the newly opened port and can be used by a client to contact the server. The maximum size string that may be supplied by the system is `MPI_MAX_PORT_NAME`.

Advice to users. The system copies the port name into `port_name`. The application must pass a buffer of sufficient size to hold this value. (*End of advice to users.*)

`port_name` is essentially a network address. It is unique within the communication universe to which it belongs (determined by the implementation), and may be used by any client within that communication universe. For instance, if it is an internet (host:port) address, it will be unique on the internet. If it is a low level switch address on an IBM SP, it will be unique to that SP.

Advice to implementors. These examples are not meant to constrain implementations. A `port_name` could, for instance, contain a user name or the name of a batch job, as long as it is unique within some well-defined communication domain. The larger the communication domain, the more useful MPI's client/server functionality will be. (*End of advice to implementors.*)

The precise form of the address is implementation-defined. For instance, an internet address may be a host name or IP address, or anything that the implementation can decode into an IP address. A port name may be reused after it is freed with `MPI_CLOSE_PORT` and released by the system.

Advice to implementors. Since the user may type in `port_name` by hand, it is useful to choose a form that is easily readable and does not have embedded spaces. (*End of advice to implementors.*)

info may be used to tell the implementation how to establish the address. It may, and usually will, be MPI_INFO_NULL in order to get the implementation defaults.

MPI_CLOSE_PORT(port_name)

IN port_name a port (string)

int MPI_Close_port(char *port_name)

MPI_CLOSE_PORT(PORT_NAME, IERROR)

CHARACTER*(*) PORT_NAME

INTEGER IERROR

void MPI::Close_port(const char* port_name)

This function releases the network address represented by port_name.

MPI_COMM_ACCEPT(port_name, info, root, comm, newcomm)

IN port_name port name (string, used only on root)

IN info implementation-dependent information (handle, used only on root)

IN root rank in comm of root node (integer)

IN comm intracommunicator over which call is collective (handle)

OUT newcomm intercommunicator with client as remote group (handle)

int MPI_Comm_accept(char *port_name, MPI_Info info, int root, MPI_Comm comm, MPI_Comm *newcomm)

MPI_COMM_ACCEPT(PORT_NAME, INFO, ROOT, COMM, NEWCOMM, IERROR)

CHARACTER*(*) PORT_NAME

INTEGER INFO, ROOT, COMM, NEWCOMM, IERROR

MPI::Intercomm MPI::Intracomm::Accept(const char* port_name, const MPI::Info& info, int root) const

MPI_COMM_ACCEPT establishes communication with a client. It is collective over the calling communicator. It returns an intercommunicator that allows communication with the client.

The port_name must have been established through a call to MPI_OPEN_PORT.

info is a implementation-defined string that may allow fine control over the ACCEPT call.

9.4.3 Client Routines

There is only one routine on the client side.

```

1 MPI_COMM_CONNECT(port_name, info, root, comm, newcomm)
2     IN      port_name      network address (string, used only on root)
3
4     IN      info           implementation-dependent information (handle, used
5                          only on root)
6
7     IN      root           rank in comm of root node (integer)
8
9     IN      comm           intracommunicator over which call is collective (han-
10                          dle)
11
12     OUT     newcomm        intercommunicator with server as remote group (han-
13                          dle)

```

```

13 int MPI_Comm_connect(char *port_name, MPI_Info info, int root,
14                    MPI_Comm comm, MPI_Comm *newcomm)

```

```

15 MPI_COMM_CONNECT(PORT_NAME, INFO, ROOT, COMM, NEWCOMM, IERROR)
16     CHARACTER*(*) PORT_NAME
17     INTEGER INFO, ROOT, COMM, NEWCOMM, IERROR
18

```

```

19 MPI::Intercomm MPI::Intracomm::Connect(const char* port_name,
20                                       const MPI::Info& info, int root) const

```

This routine establishes communication with a server specified by `port_name`. It is collective over the calling communicator and returns an intercommunicator in which the remote group participated in an `MPI_COMM_ACCEPT`.

If the named port does not exist (or has been closed), `MPI_COMM_CONNECT` raises an error of class `MPI_ERR_PORT`.

If the port exists, but does not have a pending `MPI_COMM_ACCEPT`, the connection attempt will eventually time out after an implementation-defined time, or succeed when the server calls `MPI_COMM_ACCEPT`. In the case of a time out, `MPI_COMM_CONNECT` raises an error of class `MPI_ERR_PORT`.

Advice to implementors. The time out period may be arbitrarily short or long. However, a high quality implementation will try to queue connection attempts so that a server can handle simultaneous requests from several clients. A high quality implementation may also provide a mechanism, through the `info` arguments to `MPI_OPEN_PORT`, `MPI_COMM_ACCEPT` and/or `MPI_COMM_CONNECT`, for the user to control timeout and queuing behavior. (*End of advice to implementors.*)

MPI provides no guarantee of fairness in servicing connection attempts. That is, connection attempts are not necessarily satisfied in the order they were initiated and competition from other connection attempts may prevent a particular connection attempt from being satisfied.

`port_name` is the address of the server. It must be the same as the name returned by `MPI_OPEN_PORT` on the server. Some freedom is allowed here. If there are equivalent forms of `port_name`, an implementation may accept them as well. For instance, if `port_name` is `(hostname:port)`, an implementation may accept `(ip_address:port)` as well.

9.4.4 Name Publishing

The routines in this section provide a mechanism for publishing names. A (`service_name`, `port_name`) pair is published by the server, and may be retrieved by a client using the `service_name` only. An MPI implementation defines the *scope* of the `service_name`, that is, the domain over which the `service_name` can be retrieved. If the domain is the empty set, that is, if no client can retrieve the information, then we say that name publishing is not supported. Implementations should document how the scope is determined. High quality implementations will give some control to users through the `info` arguments to name publishing functions. Examples are given in the descriptions of individual functions.

```
MPI_PUBLISH_NAME(service_name, info, port_name)
```

| | | |
|----|---------------------------|--|
| IN | <code>service_name</code> | a service name to associate with the port (string) |
| IN | <code>info</code> | implementation-specific information (handle) |
| IN | <code>port_name</code> | a port name (string) |

```
int MPI_Publish_name(char *service_name, MPI_Info info, char *port_name)
```

```
MPI_PUBLISH_NAME(SERVICE_NAME, INFO, PORT_NAME, IERROR)
```

```
INTEGER INFO, IERROR
```

```
CHARACTER*(*) SERVICE_NAME, PORT_NAME
```

```
void MPI::Publish_name(const char* service_name, const MPI::Info& info,  
                      const char* port_name)
```

This routine publishes the pair (`port_name`, `service_name`) so that an application may retrieve a system-supplied `port_name` using a well-known `service_name`.

The implementation must define the *scope* of a published service name, that is, the domain over which the service name is unique, and conversely, the domain over which the (`port name`, `service name`) pair may be retrieved. For instance, a service name may be unique to a job (where job is defined by a distributed operating system or batch scheduler), unique to a machine, or unique to a Kerberos realm. The scope may depend on the `info` argument to `MPI_PUBLISH_NAME`.

MPI permits publishing more than one `service_name` for a single `port_name`. On the other hand, if `service_name` has already been published within the scope determined by `info`, the behavior of `MPI_PUBLISH_NAME` is undefined. An MPI implementation may, through a mechanism in the `info` argument to `MPI_PUBLISH_NAME`, provide a way to allow multiple servers with the same service in the same scope. In this case, an implementation-defined policy will determine which of several port names is returned by `MPI_LOOKUP_NAME`.

Note that while `service_name` has a limited scope, determined by the implementation, `port_name` always has global scope within the communication universe used by the implementation (i.e., it is globally unique).

`port_name` should be the name of a port established by `MPI_OPEN_PORT` and not yet deleted by `MPI_CLOSE_PORT`. If it is not, the result is undefined.

Advice to implementors. In some cases, an MPI implementation may use a name service that a user can also access directly. In this case, a name published by MPI could easily conflict with a name published by a user. In order to avoid such conflicts,

1 MPI implementations should mangle service names so that they are unlikely to conflict
 2 with user code that makes use of the same service. Such name mangling will of course
 3 be completely transparent to the user.

4 The following situation is problematic but unavoidable, if we want to allow implemen-
 5 tations to use nameservers. Suppose there are multiple instances of “ocean” running
 6 on a machine. If the scope of a service name is confined to a job, then multiple
 7 oceans can coexist. If an implementation provides site-wide scope, however, multiple
 8 instances are not possible as all calls to `MPI_PUBLISH_NAME` after the first may fail.
 9 There is no universal solution to this.

10 To handle these situations, a high quality implementation should make it possible to
 11 limit the domain over which names are published. (*End of advice to implementors.*)
 12

13
 14
 15 `MPI_UNPUBLISH_NAME(service_name, info, port_name)`

16 IN service_name a service name (string)
 17 IN info implementation-specific information (handle)
 18 IN port_name a port name (string)

19
 20
 21 `int MPI_Unpublish_name(char *service_name, MPI_Info info, char *port_name)`

22 `MPI_UNPUBLISH_NAME(SERVICE_NAME, INFO, PORT_NAME, IERROR)`

23 INTEGER INFO, IERROR

24 CHARACTER*(*) SERVICE_NAME, PORT_NAME

25
 26 `void MPI::Unpublish_name(const char* service_name, const MPI::Info& info,`
 27 `const char* port_name)`

28
 29 This routine unpublishes a service name that has been previously published. Attempt-
 30 ing to unpublish a name that has not been published or has already been unpublished is
 31 erroneous and is indicated by the error class `MPI_ERR_SERVICE`.

32 All published names must be unpublished before the corresponding port is closed and
 33 before the publishing process exits. The behavior of `MPI_UNPUBLISH_NAME` is implemen-
 34 tation dependent when a process tries to unpublish a name that it did not publish.

35 If the `info` argument was used with `MPI_PUBLISH_NAME` to tell the implementation
 36 how to publish names, the implementation may require that `info` passed to
 37 `MPI_UNPUBLISH_NAME` contain information to tell the implementation how to unpublish
 38 a name.

39
 40
 41 `MPI_LOOKUP_NAME(service_name, info, port_name)`

42 IN service_name a service name (string)
 43 IN info implementation-specific information (handle)
 44 OUT port_name a port name (string)

45
 46 `int MPI_Lookup_name(char *service_name, MPI_Info info, char *port_name)`

47 `MPI_LOOKUP_NAME(SERVICE_NAME, INFO, PORT_NAME, IERROR)`
 48


```

CHARACTER*(*) SERVICE_NAME, PORT_NAME
INTEGER INFO, IERROR

void MPI::Lookup_name(const char* service_name, const MPI::Info& info,
                     char* port_name)

```

This function retrieves a `port_name` published by `MPI_PUBLISH_NAME` with `service_name`. If `service_name` has not been published, it raises an error in the error class `MPI_ERR_NAME`. The application must supply a `port_name` buffer large enough to hold the largest possible port name (see discussion above under `MPI_OPEN_PORT`).

If an implementation allows multiple entries with the same `service_name` within the same scope, a particular `port_name` is chosen in a way determined by the implementation.

If the `info` argument was used with `MPI_PUBLISH_NAME` to tell the implementation how to publish names, a similar `info` argument may be required for `MPI_LOOKUP_NAME`.

9.4.5 Reserved Key Values

The following key values are reserved. An implementation is not required to interpret these key values, but if it does interpret the key value, it must provide the functionality described.

`ip_port` Value contains IP port number at which to establish a port. (Reserved for `MPI_OPEN_PORT` only).

`ip_address` Value contains IP address at which to establish a port. If the address is not a valid IP address of the host on which the `MPI_OPEN_PORT` call is made, the results are undefined. (Reserved for `MPI_OPEN_PORT` only).

9.4.6 Client/Server Examples

Simplest Example — Completely Portable.

The following example shows the simplest way to use the client/server interface. It does not use service names at all.

On the server side:

```

char myport[MPI_MAX_PORT_NAME];
MPI_Comm intercomm;
/* ... */
MPI_Open_port(MPI_INFO_NULL, myport);
printf("port name is: %s\n", myport);

MPI_Comm_accept(myport, MPI_INFO_NULL, 0, MPI_COMM_SELF, &intercomm);
/* do something with intercomm */

```

The server prints out the port name to the terminal and the user must type it in when starting up the client (assuming the MPI implementation supports stdin such that this works). On the client side:

```

MPI_Comm intercomm;
char name[MPI_MAX_PORT_NAME];

```

```

1     printf("enter port name: ");
2     gets(name);
3     MPI_Comm_connect(name, MPI_INFO_NULL, 0, MPI_COMM_SELF, &intercomm);
4

```

Ocean/Atmosphere - Relies on Name Publishing

In this example, the “ocean” application is the “server” side of a coupled ocean-atmosphere climate model. It assumes that the MPI implementation publishes names.

```

9
10    MPI_Open_port(MPI_INFO_NULL, port_name);
11    MPI_Publish_name("ocean", MPI_INFO_NULL, port_name);
12
13
14    MPI_Comm_accept(port_name, MPI_INFO_NULL, 0, MPI_COMM_SELF, &intercomm);
15    /* do something with intercomm */
16    MPI_Unpublish_name("ocean", MPI_INFO_NULL, port_name);
17

```

On the client side:

```

18
19
20    MPI_Lookup_name("ocean", MPI_INFO_NULL, port_name);
21    MPI_Comm_connect( port_name, MPI_INFO_NULL, 0, MPI_COMM_SELF,
22                    &intercomm);
23

```

Simple Client-Server Example.

This is a simple example; the server accepts only a single connection at a time and serves that connection until the client requests to be disconnected. The server is a single process.

Here is the server. It accepts a single connection and then processes data until it receives a message with tag 1. A message with tag 0 tells the server to exit.

```

24
25
26    This is a simple example; the server accepts only a single connection at a time and serves
27    that connection until the client requests to be disconnected. The server is a single process.
28    Here is the server. It accepts a single connection and then processes data until it
29    receives a message with tag 1. A message with tag 0 tells the server to exit.
30
31    #include "mpi.h"
32    int main( int argc, char **argv )
33    {
34        MPI_Comm client;
35        MPI_Status status;
36        char port_name[MPI_MAX_PORT_NAME];
37        double buf[MAX_DATA];
38        int size, again;
39
40        MPI_Init( &argc, &argv );
41        MPI_Comm_size(MPI_COMM_WORLD, &size);
42        if (size != 1) error(FATAL, "Server too big");
43        MPI_Open_port(MPI_INFO_NULL, port_name);
44        printf("server available at %s\n",port_name);
45        while (1) {
46            MPI_Comm_accept( port_name, MPI_INFO_NULL, 0, MPI_COMM_WORLD,
47                            &client );
48            again = 1;
49            while (again) {

```

```
        MPI_Recv( buf, MAX_DATA, MPI_DOUBLE,
                MPI_ANY_SOURCE, MPI_ANY_TAG, client, &status );
switch (status.MPI_TAG) {
    case 0: MPI_Comm_free( &client );
            MPI_Close_port(port_name);
            MPI_Finalize();
            return 0;
    case 1: MPI_Comm_disconnect( &client );
            again = 0;
            break;
    case 2: /* do something */
    ...
    default:
                /* Unexpected message type */
                MPI_Abort( MPI_COMM_WORLD, 1 );
    }
    }
}
```

Here is the client.

```
#include "mpi.h"
int main( int argc, char **argv )
{
    MPI_Comm server;
    double buf[MAX_DATA];
    char port_name[MPI_MAX_PORT_NAME];

    MPI_Init( &argc, &argv );
    strcpy(port_name, argv[1] );/* assume server's name is cmd-line arg */

    MPI_Comm_connect( port_name, MPI_INFO_NULL, 0, MPI_COMM_WORLD,
                    &server );

    while (!done) {
        tag = 2; /* Action to perform */
        MPI_Send( buf, n, MPI_DOUBLE, 0, tag, server );
        /* etc */
    }
    MPI_Send( buf, 0, MPI_DOUBLE, 0, 1, server );
    MPI_Comm_disconnect( &server );
    MPI_Finalize();
    return 0;
}
```

9.5 Other Functionality

9.5.1 Universe Size

Many “dynamic” MPI applications are expected to exist in a static runtime environment, in which resources have been allocated before the application is run. When a user (or possibly a batch system) runs one of these quasi-static applications, she will usually specify a number of processes to start and a total number of processes that are expected. An application simply needs to know how many slots there are, i.e., how many processes it should spawn.

MPI provides an attribute on `MPI_COMM_WORLD`, `MPI_UNIVERSE_SIZE`, that allows the application to obtain this information in a portable manner. This attribute indicates the total number of processes that are expected. In Fortran, the attribute is the integer value. In C, the attribute is a pointer to the integer value. An application typically subtracts the size of `MPI_COMM_WORLD` from `MPI_UNIVERSE_SIZE` to find out how many processes it should spawn. `MPI_UNIVERSE_SIZE` is initialized in `MPI_INIT` and is not changed by MPI. If defined, it has the same value on all processes of `MPI_COMM_WORLD`. `MPI_UNIVERSE_SIZE` is determined by the application startup mechanism in a way not specified by MPI. (The size of `MPI_COMM_WORLD` is another example of such a parameter.)

Possibilities for how `MPI_UNIVERSE_SIZE` might be set include

- A `-universe_size` argument to a program that starts MPI processes.
- Automatic interaction with a batch scheduler to figure out how many processors have been allocated to an application.
- An environment variable set by the user.
- Extra information passed to `MPI_COMM_SPAWN` through the `info` argument.

An implementation must document how `MPI_UNIVERSE_SIZE` is set. An implementation may not support the ability to set `MPI_UNIVERSE_SIZE`, in which case the attribute `MPI_UNIVERSE_SIZE` is not set.

`MPI_UNIVERSE_SIZE` is a recommendation, not necessarily a hard limit. For instance, some implementations may allow an application to spawn 50 processes per processor, if they are requested. However, it is likely that the user only wants to spawn one process per processor.

`MPI_UNIVERSE_SIZE` is assumed to have been specified when an application was started, and is in essence a portable mechanism to allow the user to pass to the application (through the MPI process startup mechanism, such as `mpiexec`) a piece of critical runtime information. Note that no interaction with the runtime environment is required. If the runtime environment changes size while an application is running, `MPI_UNIVERSE_SIZE` is not updated, and the application must find out about the change through direct communication with the runtime system.

9.5.2 Singleton `MPI_INIT`

A high-quality implementation will allow any process (including those not started with a “parallel application” mechanism) to become an MPI process by calling `MPI_INIT`. Such a process can then connect to other MPI processes using the `MPI_COMM_ACCEPT` and

MPI_COMM_CONNECT routines, or spawn other MPI processes. MPI does not mandate this behavior, but strongly encourages it where technically feasible.

Advice to implementors. To start an MPI-1 application with more than one process requires some special coordination. The processes must be started at the “same” time, they must have a mechanism to establish communication, etc. Either the user or the operating system must take special steps beyond simply starting processes.

When an application enters MPI_INIT, clearly it must be able to determine if these special steps were taken. MPI-1 does not say what happens if these special steps were not taken — presumably this is treated as an error in starting the MPI application. MPI-2 recommends the following behavior.

If a process enters MPI_INIT and determines that no special steps were taken (i.e., it has not been given the information to form an MPI_COMM_WORLD with other processes) it succeeds and forms a singleton MPI program, that is, one in which MPI_COMM_WORLD has size 1.

In some implementations, MPI may not be able to function without an “MPI environment.” For example, MPI may require that daemons be running or MPI may not be able to work at all on the front-end of an MPP. In this case, an MPI implementation may either

1. Create the environment (e.g., start a daemon) or
2. Raise an error if it cannot create the environment and the environment has not been started independently.

A high quality implementation will try to create a singleton MPI process and not raise an error.

(End of advice to implementors.)

9.5.3 MPI_APPNUM

There is a predefined attribute MPI_APPNUM of MPI_COMM_WORLD. In Fortran, the attribute is an integer value. In C, the attribute is a pointer to an integer value. If a process was spawned with MPI_COMM_SPAWN_MULTIPLE, MPI_APPNUM is the command number that generated the current process. Numbering starts from zero. If a process was spawned with MPI_COMM_SPAWN, it will have MPI_APPNUM equal to zero.

Additionally, if the process was not started by a spawn call, but by an implementation-specific startup mechanism that can handle multiple process specifications, MPI_APPNUM should be set to the number of the corresponding process specification. In particular, if it is started with

```
mpirexec spec0 [: spec1 : spec2 : ...]
```

MPI_APPNUM should be set to the number of the corresponding specification.

If an application was not spawned with MPI_COMM_SPAWN or MPI_COMM_SPAWN_MULTIPLE, and MPI_APPNUM doesn’t make sense in the context of the implementation-specific startup mechanism, MPI_APPNUM is not set.

MPI implementations may optionally provide a mechanism to override the value of MPI_APPNUM through the info argument. MPI reserves the following key for all SPAWN calls.

1 `appnum` Value contains an integer that overrides the default value for `MPI_APPNUM` in the
 2 child.

3
 4 *Rationale.* When a single application is started, it is able to figure out how many pro-
 5 cesses there are by looking at the size of `MPI_COMM_WORLD`. An application consisting
 6 of multiple SPMD sub-applications has no way to find out how many sub-applications
 7 there are and to which sub-application the process belongs. While there are ways to
 8 figure it out in special cases, there is no general mechanism. `MPI_APPNUM` provides
 9 such a general mechanism. (*End of rationale.*)

11 9.5.4 Releasing Connections

12 Before a client and server connect, they are independent MPI applications. An error in one
 13 does not affect the other. After establishing a connection with `MPI_COMM_CONNECT` and
 14 `MPI_COMM_ACCEPT`, an error in one may affect the other. It is desirable for a client and
 15 server to be able to disconnect, so that an error in one will not affect the other. Similarly,
 16 it might be desirable for a parent and child to disconnect, so that errors in the child do not
 17 affect the parent, or vice-versa.

- 19 • Two processes are **connected** if there is a communication path (direct or indirect)
 20 between them. More precisely:
 - 22 1. Two processes are connected if
 - 23 (a) they both belong to the same communicator (inter- or intra-, including
 24 `MPI_COMM_WORLD`) *or*
 - 25 (b) they have previously belonged to a communicator that was freed with
 26 `MPI_COMM_FREE` instead of `MPI_COMM_DISCONNECT` *or*
 - 27 (c) they both belong to the group of the same window or filehandle.
 - 28 2. If A is connected to B and B to C, then A is connected to C.
- 30 • Two processes are **disconnected** (also **independent**) if they are not connected.
- 31 • By the above definitions, connectivity is a transitive property, and divides the uni-
 32 verse of MPI processes into disconnected (independent) sets (equivalence classes) of
 33 processes.
- 34 • Processes which are connected, but don't share the same `MPI_COMM_WORLD` may
 35 become disconnected (independent) if the communication path between them is bro-
 36 ken by using `MPI_COMM_DISCONNECT`.

39 The following additional rules apply to MPI-1 functions:

- 40 • `MPI_FINALIZE` is collective over a set of connected processes.
- 41 • `MPI_ABORT` does not abort independent processes. As in MPI-1, it may abort all
 42 processes in `MPI_COMM_WORLD` (ignoring its `comm` argument). Additionally, it
 43 may abort connected processes as well, though it makes a “best attempt” to abort
 44 only the processes in `comm`.
- 45 • If a process terminates without calling `MPI_FINALIZE`, independent processes are not
 46 affected but the effect on connected processes is not defined.

```
MPI_COMM_DISCONNECT(comm)
```

```
  INOUT  comm                communicator (handle)
```

```
int MPI_Comm_disconnect(MPI_Comm *comm)
```

```
MPI_COMM_DISCONNECT(COMM, IERROR)
```

```
  INTEGER COMM, IERROR
```

```
void MPI::Comm::Disconnect()
```

This function waits for all pending communication on `comm` to complete internally, deallocates the communicator object, and sets the handle to `MPI_COMM_NULL`. It is a collective operation.

It may not be called with the communicator `MPI_COMM_WORLD` or `MPI_COMM_SELF`.

`MPI_COMM_DISCONNECT` may be called only if all communication is complete and matched, so that buffered data can be delivered to its destination. This requirement is the same as for `MPI_FINALIZE`.

`MPI_COMM_DISCONNECT` has the same action as `MPI_COMM_FREE`, except that it waits for pending communication to finish internally and enables the guarantee about the behavior of disconnected processes.

Advice to users. To disconnect two processes you may need to call `MPI_COMM_DISCONNECT`, `MPI_WIN_FREE` and `MPI_FILE_CLOSE` to remove all communication paths between the two processes. Notes that it may be necessary to disconnect several communicators (or to free several windows or files) before two processes are completely independent. (*End of advice to users.*)

Rationale. It would be nice to be able to use `MPI_COMM_FREE` instead, but that function explicitly does not wait for pending communication to complete. (*End of rationale.*)

9.5.5 Another Way to Establish MPI Communication

```
MPI_COMM_JOIN(fd, intercomm)
```

```
  IN      fd                socket file descriptor
```

```
  OUT    intercomm         new intercommunicator (handle)
```

```
int MPI_Comm_join(int fd, MPI_Comm *intercomm)
```

```
MPI_COMM_JOIN(FD, INTERCOMM, IERROR)
```

```
  INTEGER FD, INTERCOMM, IERROR
```

```
static MPI::Intercomm MPI::Comm::Join(const int fd)
```

`MPI_COMM_JOIN` is intended for MPI implementations that exist in an environment supporting the Berkeley Socket interface [35, 39]. Implementations that exist in an environment not supporting Berkeley Sockets should provide the entry point for `MPI_COMM_JOIN` and should return `MPI_COMM_NULL`.

1 This call creates an intercommunicator from the union of two MPI processes which are
2 connected by a socket. `MPI_COMM_JOIN` should normally succeed if the local and remote
3 processes have access to the same implementation-defined MPI communication universe.

4
5 *Advice to users.* An MPI implementation may require a specific communication
6 medium for MPI communication, such as a shared memory segment or a special switch.
7 In this case, it may not be possible for two processes to successfully join even if there
8 is a socket connecting them and they are using the same MPI implementation. (*End*
9 *of advice to users.*)

10 *Advice to implementors.* A high quality implementation will attempt to establish
11 communication over a slow medium if its preferred one is not available. If implemen-
12 tations do not do this, they must document why they cannot do MPI communication
13 over the medium used by the socket (especially if the socket is a TCP connection).
14 (*End of advice to implementors.*)

15
16 `fd` is a file descriptor representing a socket of type `SOCK_STREAM` (a two-way reliable
17 byte-stream connection). Non-blocking I/O and asynchronous notification via `SIGIO` must
18 not be enabled for the socket. The socket must be in a connected state. The socket must
19 be quiescent when `MPI_COMM_JOIN` is called (see below). It is the responsibility of the
20 application to create the socket using standard socket API calls.

21 `MPI_COMM_JOIN` must be called by the process at each end of the socket. It does not
22 return until both processes have called `MPI_COMM_JOIN`. The two processes are referred
23 to as the local and remote processes.

24 MPI uses the socket to bootstrap creation of the intercommunicator, and for nothing
25 else. Upon return from `MPI_COMM_JOIN`, the file descriptor will be open and quiescent
26 (see below).

27 If MPI is unable to create an intercommunicator, but is able to leave the socket in its
28 original state, with no pending communication, it succeeds and sets `intercomm` to
29 `MPI_COMM_NULL`.

30 The socket must be quiescent before `MPI_COMM_JOIN` is called and after
31 `MPI_COMM_JOIN` returns. More specifically, on entry to `MPI_COMM_JOIN`, a `read` on the
32 socket will not read any data that was written to the socket before the remote process called
33 `MPI_COMM_JOIN`. On exit from `MPI_COMM_JOIN`, a `read` will not read any data that was
34 written to the socket before the remote process returned from `MPI_COMM_JOIN`. It is the
35 responsibility of the application to ensure the first condition, and the responsibility of the
36 MPI implementation to ensure the second. In a multithreaded application, the application
37 must ensure that one thread does not access the socket while another is calling
38 `MPI_COMM_JOIN`, or call `MPI_COMM_JOIN` concurrently.

39 *Advice to implementors.* MPI is free to use any available communication path(s)
40 for MPI messages in the new communicator; the socket is only used for the initial
41 handshaking. (*End of advice to implementors.*)

42
43 `MPI_COMM_JOIN` uses non-MPI communication to do its work. The interaction of
44 non-MPI communication with pending MPI communication is not defined. Therefore, the
45 result of calling `MPI_COMM_JOIN` on two connected processes (see Section 9.5.4 on page
46 304 for the definition of connected) is undefined.

47 The returned communicator may be used to establish MPI communication with addi-
48 tional processes, through the usual MPI communicator creation mechanisms.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 10

One-Sided Communications

10.1 Introduction

Remote Memory Access (RMA) extends the communication mechanisms of MPI by allowing one process to specify all communication parameters, both for the sending side and for the receiving side. This mode of communication facilitates the coding of some applications with dynamically changing data access patterns where the data distribution is fixed or slowly changing. In such a case, each process can compute what data it needs to access or update at other processes. However, processes may not know which data in their own memory need to be accessed or updated by remote processes, and may not even know the identity of these processes. Thus, the transfer parameters are all available only on one side. Regular send/receive communication requires matching operations by sender and receiver. In order to issue the matching operations, an application needs to distribute the transfer parameters. This may require all processes to participate in a time consuming global computation, or to periodically poll for potential communication requests to receive and act upon. The use of RMA communication mechanisms avoids the need for global computations or explicit polling. A generic example of this nature is the execution of an assignment of the form $A = B(\text{map})$, where map is a permutation vector, and A , B and map are distributed in the same manner.

Message-passing communication achieves two effects: *communication* of data from sender to receiver; and *synchronization* of sender with receiver. The RMA design separates these two functions. Three communication calls are provided: `MPI_PUT` (remote write), `MPI_GET` (remote read) and `MPI_ACCUMULATE` (remote update). A larger number of synchronization calls are provided that support different synchronization styles. The design is similar to that of weakly coherent memory systems: correct ordering of memory accesses has to be imposed by the user, using synchronization calls; the implementation can delay communication operations until the synchronization calls occur, for efficiency.

The design of the RMA functions allows implementors to take advantage, in many cases, of fast communication mechanisms provided by various platforms, such as coherent or noncoherent shared memory, DMA engines, hardware-supported put/get operations, communication coprocessors, etc. The most frequently used RMA communication mechanisms can be layered on top of message passing. However, support for asynchronous communication agents (handlers, threads, etc.) is needed, for certain RMA functions, in a distributed memory environment.

We shall denote by **origin** the process that performs the call, and by **target** the

process in which the memory is accessed. Thus, in a put operation, source=origin and destination=target; in a get operation, source=target and destination=origin.

10.2 Initialization

10.2.1 Window Creation

The initialization operation allows each process in an intracommunicator group to specify, in a collective operation, a “window” in its memory that is made accessible to accesses by remote processes. The call returns an opaque object that represents the group of processes that own and access the set of windows, and the attributes of each window, as specified by the initialization call.

`MPI_WIN_CREATE(base, size, disp_unit, info, comm, win)`

| | | |
|-----|-----------|--|
| IN | base | initial address of window (choice) |
| IN | size | size of window in bytes (nonnegative integer) |
| IN | disp_unit | local unit size for displacements, in bytes (positive integer) |
| IN | info | info argument (handle) |
| IN | comm | communicator (handle) |
| OUT | win | window object returned by the call (handle) |

```
int MPI_Win_create(void *base, MPI_Aint size, int disp_unit, MPI_Info info,
                  MPI_Comm comm, MPI_Win *win)
```

```
MPI_WIN_CREATE(BASE, SIZE, DISP_UNIT, INFO, COMM, WIN, IERROR)
```

```
<type> BASE(*)
INTEGER(KIND=MPI_ADDRESS_KIND) SIZE
INTEGER DISP_UNIT, INFO, COMM, WIN, IERROR
```

```
static MPI::Win MPI::Win::Create(const void* base, MPI::Aint size, int
                                disp_unit, const MPI::Info& info, const MPI::Intracomm& comm)
```

This is a collective call executed by all processes in the group of `comm`. It returns a window object that can be used by these processes to perform RMA operations. Each process specifies a window of existing memory that it exposes to RMA accesses by the processes in the group of `comm`. The window consists of `size` bytes, starting at address `base`. A process may elect to expose no memory by specifying `size = 0`.

The displacement unit argument is provided to facilitate address arithmetic in RMA operations: the target displacement argument of an RMA operation is scaled by the factor `disp_unit` specified by the target process, at window creation.

Rationale. The window size is specified using an address sized integer, so as to allow windows that span more than 4 GB of address space. (Even if the physical memory size is less than 4 GB, the address range may be larger than 4 GB, if addresses are not contiguous.) (*End of rationale.*)

Advice to users. Common choices for `disp_unit` are 1 (no scaling), and (in C syntax) `sizeof(type)`, for a window that consists of an array of elements of type `type`. The later choice will allow one to use array indices in RMA calls, and have those scaled correctly to byte displacements, even in a heterogeneous environment. (*End of advice to users.*)

The `info` argument provides optimization hints to the runtime about the expected usage pattern of the window. The following `info` key is predefined:

`no_locks` — if set to `true`, then the implementation may assume that the local window is never locked (by a call to `MPI_WIN_LOCK`). This implies that this window is not used for 3-party communication, and RMA can be implemented with no (less) asynchronous agent activity at this process.

The various processes in the group of `comm` may specify completely different target windows, in location, size, displacement units and `info` arguments. As long as all the `get`, `put` and `accumulate` accesses to a particular process fit their specific target window this should pose no problem. The same area in memory may appear in multiple windows, each associated with a different window object. However, concurrent communications to distinct, overlapping windows may lead to erroneous results.

Advice to users. A window can be created in any part of the process memory. However, on some systems, the performance of windows in memory allocated by `MPI_ALLOC_MEM` (Section 7.2, page 250) will be better. Also, on some systems, performance is improved when window boundaries are aligned at “natural” boundaries (word, double-word, cache line, page frame, etc.). (*End of advice to users.*)

Advice to implementors. In cases where RMA operations use different mechanisms in different memory areas (e.g., load/store in a shared memory segment, and an asynchronous handler in private memory), the `MPI_WIN_CREATE` call needs to figure out which type of memory is used for the window. To do so, MPI maintains, internally, the list of memory segments allocated by `MPI_ALLOC_MEM`, or by other, implementation specific, mechanisms, together with information on the type of memory segment allocated. When a call to `MPI_WIN_CREATE` occurs, then MPI checks which segment contains each window, and decides, accordingly, which mechanism to use for RMA operations.

Vendors may provide additional, implementation-specific mechanisms to allow “good” memory to be used for static variables.

Implementors should document any performance impact of window alignment. (*End of advice to implementors.*)

`MPI_WIN_FREE(win)`

INOUT win window object (handle)

`int MPI_Win_free(MPI_Win *win)`

`MPI_WIN_FREE(WIN, IERROR)`

1 INTEGER WIN, IERROR

2
3 void MPI::Win::Free()4 Frees the window object `win` and returns a null handle (equal to
5 `MPI_WIN_NULL`). This is a collective call executed by all processes in the group associated
6 with `win`. `MPI_WIN_FREE(win)` can be invoked by a process only after it has completed its
7 involvement in RMA communications on window `win`: i.e., the process has called
8 `MPI_WIN_FENCE`, or called `MPI_WIN_WAIT` to match a previous call to `MPI_WIN_POST`
9 or called `MPI_WIN_COMPLETE` to match a previous call to `MPI_WIN_START` or called
10 `MPI_WIN_UNLOCK` to match a previous call to `MPI_WIN_LOCK`. When the call returns, the
11 window memory can be freed.12
13 *Advice to implementors.* `MPI_WIN_FREE` requires a barrier synchronization: no
14 process can return from free until all processes in the group of `win` called free. This, to
15 ensure that no process will attempt to access a remote window (e.g., with lock/unlock)
16 after it was freed. (*End of advice to implementors.*)17
18 10.2.2 Window Attributes

19 The following three attributes are cached with a window, when the window is created.

20
21 MPI_WIN_BASE window base address.
22 MPI_WIN_SIZE window size, in bytes.
23 MPI_WIN_DISP_UNIT displacement unit associated with the window.
2425 In C, calls to `MPI_Win_get_attr(win, MPI_WIN_BASE, &base, &flag)`,
26 `MPI_Win_get_attr(win, MPI_WIN_SIZE, &size, &flag)` and
27 `MPI_Win_get_attr(win, MPI_WIN_DISP_UNIT, &disp_unit, &flag)` will return in `base` a pointer
28 to the start of the window `win`, and will return in `size` and `disp_unit` pointers to the size and
29 displacement unit of the window, respectively. And similarly, in C++.30 In Fortran, calls to `MPI_WIN_GET_ATTR(win, MPI_WIN_BASE, base, flag, ierror)`,
31 `MPI_WIN_GET_ATTR(win, MPI_WIN_SIZE, size, flag, ierror)` and
32 `MPI_WIN_GET_ATTR(win, MPI_WIN_DISP_UNIT, disp_unit, flag, ierror)` will return in
33 `base`, `size` and `disp_unit` the (integer representation of) the base address, the size and the
34 displacement unit of the window `win`, respectively. (The window attribute access functions
35 are defined in Section 5.7.1, page 212.)36 The other “window attribute,” namely the group of processes attached to the window,
37 can be retrieved using the call below.38
39 MPI_WIN_GET_GROUP(win, group)40
41 IN win window object (handle)
42 OUT group group of processes which share access to the window
43 (handle)
44

45 int MPI_Win_get_group(MPI_Win win, MPI_Group *group)

46 MPI_WIN_GET_GROUP(WIN, GROUP, IERROR)

47 INTEGER WIN, GROUP, IERROR
48

`MPI::Group MPI::Win::Get_group() const`

`MPI_WIN_GET_GROUP` returns a duplicate of the group of the communicator used to create the window. associated with `win`. The group is returned in `group`.

10.3 Communication Calls

MPI supports three RMA communication calls: `MPI_PUT` transfers data from the caller memory (origin) to the target memory; `MPI_GET` transfers data from the target memory to the caller memory; and `MPI_ACCUMULATE` updates locations in the target memory, e.g. by adding to these locations values sent from the caller memory. These operations are *nonblocking*: the call initiates the transfer, but the transfer may continue after the call returns. The transfer is completed, both at the origin and at the target, when a subsequent *synchronization* call is issued by the caller on the involved window object. These synchronization calls are described in Section 10.4, page 319.

The local communication buffer of an RMA call should not be updated, and the local communication buffer of a get call should not be accessed after the RMA call, until the subsequent synchronization call completes.

Rationale. The rule above is more lenient than for message passing, where we do not allow two concurrent sends, with overlapping send buffers. Here, we allow two concurrent puts with overlapping send buffers. The reasons for this relaxation are

1. Users do not like that restriction, which is not very natural (it prohibits concurrent reads).
2. Weakening the rule does not prevent efficient implementation, as far as we know.
3. Weakening the rule is important for performance of RMA: we want to associate one synchronization call with as many RMA operations is possible. If puts from overlapping buffers cannot be concurrent, then we need to needlessly add synchronization points in the code.

(End of rationale.)

It is erroneous to have concurrent conflicting accesses to the same memory location in a window; if a location is updated by a put or accumulate operation, then this location cannot be accessed by a load or another RMA operation until the updating operation has completed at the target. There is one exception to this rule; namely, the same location can be updated by several concurrent accumulate calls, the outcome being as if these updates occurred in some order. In addition, a window cannot concurrently be updated by a put or accumulate operation and by a local store operation. This, even if these two updates access different locations in the window. The last restriction enables more efficient implementations of RMA operations on many systems. These restrictions are described in more detail in Section 10.7, page 335.

The calls use general datatype arguments to specify communication buffers at the origin and at the target. Thus, a transfer operation may also gather data at the source and scatter it at the destination. However, all arguments specifying both communication buffers are provided by the caller.

For all three calls, the target process may be identical with the origin process; i.e., a process may use an RMA operation to move data in its memory.

Rationale. The choice of supporting “self-communication” is the same as for message passing. It simplifies some coding, and is very useful with accumulate operations, to allow atomic updates of local variables. (*End of rationale.*)

MPI_PROC_NULL is a valid target rank in the MPI RMA calls MPI_ACCUMULATE, MPI_GET, and MPI_PUT. The effect is the same as for MPI_PROC_NULL in MPI point-to-point communication. After any RMA operation with rank MPI_PROC_NULL, it is still necessary to finish the RMA epoch with the synchronization method that started the epoch.

10.3.1 Put

The execution of a put operation is similar to the execution of a send by the origin process and a matching receive by the target process. The obvious difference is that all arguments are provided by one call — the call executed by the origin process.

```
MPI_PUT(origin_addr, origin_count, origin_datatype, target_rank, target_disp, target_count, target_datatype, win)
```

| | | |
|----|-----------------|--|
| IN | origin_addr | initial address of origin buffer (choice) |
| IN | origin_count | number of entries in origin buffer (nonnegative integer) |
| IN | origin_datatype | datatype of each entry in origin buffer (handle) |
| IN | target_rank | rank of target (nonnegative integer) |
| IN | target_disp | displacement from start of window to target buffer (nonnegative integer) |
| IN | target_count | number of entries in target buffer (nonnegative integer) |
| IN | target_datatype | datatype of each entry in target buffer (handle) |
| IN | win | window object used for communication (handle) |

```
int MPI_Put(void *origin_addr, int origin_count, MPI_Datatype
            origin_datatype, int target_rank, MPI_Aint target_disp, int
            target_count, MPI_Datatype target_datatype, MPI_Win win)

MPI_PUT(ORIGIN_ADDR, ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_DISP,
        TARGET_COUNT, TARGET_DATATYPE, WIN, IERROR)
<type> ORIGIN_ADDR(*)
INTEGER(KIND=MPI_ADDRESS_KIND) TARGET_DISP
INTEGER ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_COUNT,
TARGET_DATATYPE, WIN, IERROR

void MPI::Win::Put(const void* origin_addr, int origin_count, const
                  MPI::Datatype& origin_datatype, int target_rank, MPI::Aint
                  target_disp, int target_count, const MPI::Datatype&
                  target_datatype) const
```

Transfers `origin_count` successive entries of the type specified by the `origin_datatype`, starting at address `origin_addr` on the origin node to the target node specified by the

win, target_rank pair. The data are written in the target buffer at address `target_addr = window_base + target_disp × disp_unit`, where `window_base` and `disp_unit` are the base address and window displacement unit specified at window initialization, by the target process.

The target buffer is specified by the arguments `target_count` and `target_datatype`.

The data transfer is the same as that which would occur if the origin process executed a send operation with arguments `origin_addr`, `origin_count`, `origin_datatype`, `target_rank`, `tag`, `comm`, and the target process executed a receive operation with arguments `target_addr`, `target_count`, `target_datatype`, `source`, `tag`, `comm`, where `target_addr` is the target buffer address computed as explained above, and `comm` is a communicator for the group of win.

The communication must satisfy the same constraints as for a similar message-passing communication. The `target_datatype` may not specify overlapping entries in the target buffer. The message sent must fit, without truncation, in the target buffer. Furthermore, the target buffer must fit in the target window.

The `target_datatype` argument is a handle to a datatype object defined at the origin process. However, this object is interpreted at the target process: the outcome is as if the target datatype object was defined at the target process, by the same sequence of calls used to define it at the origin process. The target datatype must contain only relative displacements, not absolute addresses. The same holds for `get` and `accumulate`.

Advice to users. The `target_datatype` argument is a handle to a datatype object that is defined at the origin process, even though it defines a data layout in the target process memory. This causes no problems in a homogeneous environment, or in a heterogeneous environment, if only portable datatypes are used (portable datatypes are defined in Section 2.4, page 11).

The performance of a `put` transfer can be significantly affected, on some systems, from the choice of window location and the shape and location of the origin and target buffer: transfers to a target window in memory allocated by `MPI_ALLOC_MEM` may be much faster on shared memory systems; transfers from contiguous buffers will be faster on most, if not all, systems; the alignment of the communication buffers may also impact performance. (*End of advice to users.*)

Advice to implementors. A high quality implementation will attempt to prevent remote accesses to memory outside the window that was exposed by the process. This, both for debugging purposes, and for protection with client-server codes that use RMA. I.e., a high-quality implementation will check, if possible, window bounds on each RMA call, and raise an MPI exception at the origin call if an out-of-bound situation occurred. Note that the condition can be checked at the origin. Of course, the added safety achieved by such checks has to be weighed against the added cost of such checks. (*End of advice to implementors.*)

10.3.2 Get

`MPI_GET(origin_addr, origin_count, origin_datatype, target_rank, target_disp, target_count, target_datatype, win)`

| | | |
|-----|------------------------------|--|
| OUT | <code>origin_addr</code> | initial address of origin buffer (choice) |
| IN | <code>origin_count</code> | number of entries in origin buffer (nonnegative integer) |
| IN | <code>origin_datatype</code> | datatype of each entry in origin buffer (handle) |
| IN | <code>target_rank</code> | rank of target (nonnegative integer) |
| IN | <code>target_disp</code> | displacement from window start to the beginning of the target buffer (nonnegative integer) |
| IN | <code>target_count</code> | number of entries in target buffer (nonnegative integer) |
| IN | <code>target_datatype</code> | datatype of each entry in target buffer (handle) |
| IN | <code>win</code> | window object used for communication (handle) |

```

int MPI_Get(void *origin_addr, int origin_count, MPI_Datatype
            origin_datatype, int target_rank, MPI_Aint target_disp, int
            target_count, MPI_Datatype target_datatype, MPI_Win win)
MPI_GET(ORIGIN_ADDR, ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_DISP,
        TARGET_COUNT, TARGET_DATATYPE, WIN, IERROR)
<type> ORIGIN_ADDR(*)
INTEGER(KIND=MPI_ADDRESS_KIND) TARGET_DISP
INTEGER ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_COUNT,
TARGET_DATATYPE, WIN, IERROR

void MPI::Win::Get(void *origin_addr, int origin_count, const
                  MPI::Datatype& origin_datatype, int target_rank, MPI::Aint
                  target_disp, int target_count, const MPI::Datatype&
                  target_datatype) const

```

Similar to `MPI_PUT`, except that the direction of data transfer is reversed. Data are copied from the target memory to the origin. The `origin_datatype` may not specify overlapping entries in the origin buffer. The target buffer must be contained within the target window, and the copied data must fit, without truncation, in the origin buffer.

10.3.3 Examples

Example 10.1 We show how to implement the generic indirect assignment `A = B(map)`, where `A`, `B` and `map` have the same distribution, and `map` is a permutation. To simplify, we assume a block distribution with equal size blocks.

```

SUBROUTINE MAPVALS(A, B, map, m, comm, p)
USE MPI
INTEGER m, map(m), comm, p

```



```

REAL A(m), B(m)
1
2
INTEGER otype(p), oindex(m),    & ! used to construct origin datatypes
3
   ttype(p), tindex(m),        & ! used to construct target datatypes
4
   count(p), total(p),         &
5
   sizeofreal, win, ierr
6
7
! This part does the work that depends on the locations of B.
8
! Can be reused while this does not change
9
10
CALL MPI_TYPE_EXTENT(MPI_REAL, sizeofreal, ierr)
11
CALL MPI_WIN_CREATE(B, m*sizeofreal, sizeofreal, MPI_INFO_NULL,    &
12
                   comm, win, ierr)
13
14
! This part does the work that depends on the value of map and
! the locations of the arrays.
15
! Can be reused while these do not change
16
17
! Compute number of entries to be received from each process
18
19
DO i=1,p
20
   count(i) = 0
21
END DO
22
DO i=1,m
23
   j = map(i)/m+1
24
   count(j) = count(j)+1
25
END DO
26
27
total(1) = 0
28
DO i=2,p
29
   total(i) = total(i-1) + count(i-1)
30
END DO
31
32
DO i=1,p
33
   count(i) = 0
34
END DO
35
36
! compute origin and target indices of entries.
37
! entry i at current process is received from location
38
! k at process (j-1), where map(i) = (j-1)*m + (k-1),
39
! j = 1..p and k = 1..m
40
41
DO i=1,m
42
   j = map(i)/m+1
43
   k = MOD(map(i),m)+1
44
   count(j) = count(j)+1
45
   oindex(total(j) + count(j)) = i
46
   tindex(total(j) + count(j)) = k
47
48

```

```

1  END DO
2
3  ! create origin and target datatypes for each get operation
4  DO i=1,p
5      CALL MPI_TYPE_INDEXED_BLOCK(count(i), 1, oindex(total(i)+1), &
6                                  MPI_REAL, otype(i), ierr)
7      CALL MPI_TYPE_COMMIT(otype(i), ierr)
8      CALL MPI_TYPE_INDEXED_BLOCK(count(i), 1, tindex(total(i)+1), &
9                                  MPI_REAL, ttype(i), ierr)
10     CALL MPI_TYPE_COMMIT(ttype(i), ierr)
11 END DO
12
13 ! this part does the assignment itself
14 CALL MPI_WIN_FENCE(0, win, ierr)
15 DO i=1,p
16     CALL MPI_GET(A, 1, otype(i), i-1, 0, 1, ttype(i), win, ierr)
17 END DO
18 CALL MPI_WIN_FENCE(0, win, ierr)
19
20 CALL MPI_WIN_FREE(win, ierr)
21 DO i=1,p
22     CALL MPI_TYPE_FREE(otype(i), ierr)
23     CALL MPI_TYPE_FREE(ttype(i), ierr)
24 END DO
25 RETURN
26 END

```

28 **Example 10.2** A simpler version can be written that does not require that a datatype
 29 be built for the target buffer. But, one then needs a separate get call for each entry, as
 30 illustrated below. This code is much simpler, but usually much less efficient, for large arrays.

```

31
32 SUBROUTINE MAPVALS(A, B, map, m, comm, p)
33 USE MPI
34 INTEGER m, map(m), comm, p
35 REAL A(m), B(m)
36 INTEGER sizeofreal, win, ierr
37
38 CALL MPI_TYPE_EXTENT(MPI_REAL, sizeofreal, ierr)
39 CALL MPI_WIN_CREATE(B, m*sizeofreal, sizeofreal, MPI_INFO_NULL, &
40                    comm, win, ierr)
41
42 CALL MPI_WIN_FENCE(0, win, ierr)
43 DO i=1,m
44     j = map(i)/p
45     k = MOD(map(i),p)
46     CALL MPI_GET(A(i), 1, MPI_REAL, j, k, 1, MPI_REAL, win, ierr)
47 END DO
48 CALL MPI_WIN_FENCE(0, win, ierr)

```

```
CALL MPI_WIN_FREE(win, ierr)
RETURN
END
```

10.3.4 Accumulate Functions

It is often useful in a put operation to combine the data moved to the target process with the data that resides at that process, rather than replacing the data there. This will allow, for example, the accumulation of a sum by having all involved processes add their contribution to the sum variable in the memory of one process.

```
MPI_ACCUMULATE(origin_addr, origin_count, origin_datatype, target_rank, target_disp, target_count, target_datatype, op, win)
```

| | | |
|----|-----------------|---|
| IN | origin_addr | initial address of buffer (choice) |
| IN | origin_count | number of entries in buffer (nonnegative integer) |
| IN | origin_datatype | datatype of each buffer entry (handle) |
| IN | target_rank | rank of target (nonnegative integer) |
| IN | target_disp | displacement from start of window to beginning of target buffer (nonnegative integer) |
| IN | target_count | number of entries in target buffer (nonnegative integer) |
| IN | target_datatype | datatype of each entry in target buffer (handle) |
| IN | op | reduce operation (handle) |
| IN | win | window object (handle) |

```
int MPI_Accumulate(void *origin_addr, int origin_count,
                  MPI_Datatype origin_datatype, int target_rank,
                  MPI_Aint target_disp, int target_count,
                  MPI_Datatype target_datatype, MPI_Op op, MPI_Win win)
```

```
MPI_ACCUMULATE(ORIGIN_ADDR, ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK,
               TARGET_DISP, TARGET_COUNT, TARGET_DATATYPE, OP, WIN, IERROR)
<type> ORIGIN_ADDR(*)
INTEGER(KIND=MPI_ADDRESS_KIND) TARGET_DISP
INTEGER ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_COUNT,
TARGET_DATATYPE, OP, WIN, IERROR
```

```
void MPI::Win::Accumulate(const void* origin_addr, int origin_count, const
                        MPI::Datatype& origin_datatype, int target_rank, MPI::Aint
                        target_disp, int target_count, const MPI::Datatype&
                        target_datatype, const MPI::Op& op) const
```

Accumulate the contents of the origin buffer (as defined by `origin_addr`, `origin_count` and `origin_datatype`) to the buffer specified by arguments `target_count` and `target_datatype`, at offset `target_disp`, in the target window specified by `target_rank` and `win`, using the operation

1 op. This is like MPI_PUT except that data is combined into the target area instead of
2 overwriting it.

3 Any of the predefined operations for MPI_REDUCE can be used. User-defined functions
4 cannot be used. For example, if op is MPI_SUM, each element of the origin buffer is added
5 to the corresponding element in the target, replacing the former value in the target.

6 Each datatype argument must be a predefined datatype or a derived datatype, where
7 all basic components are of the same predefined datatype. Both datatype arguments must
8 be constructed from the same predefined datatype. The operation op applies to elements of
9 that predefined type. target_datatype must not specify overlapping entries, and the target
10 buffer must fit in the target window.

11 A new predefined operation, MPI_REPLACE, is defined. It corresponds to the associative
12 function $f(a, b) = b$; i.e., the current value in the target memory is replaced by the value
13 supplied by the origin. MPI_REPLACE, like the other predefined operations, is defined only
14 for the predefined MPI datatypes.

15
16 *Rationale.* The rationale for this is that, for consistency, MPI_REPLACE should have
17 the same limitations as the other operations. Extending it to all datatypes doesn't
18 provide any real benefit. (*End of rationale.*)

19
20
21 *Advice to users.* MPI_PUT is a special case of MPI_ACCUMULATE, with the operation
22 MPI_REPLACE. Note, however, that MPI_PUT and MPI_ACCUMULATE have different
23 constraints on concurrent updates. (*End of advice to users.*)

24 **Example 10.3** We want to compute $B(j) = \sum_{\text{map}(i)=j} A(i)$. The arrays A, B and map are
25 distributed in the same manner. We write the simple version.

```
26
27 SUBROUTINE SUM(A, B, map, m, comm, p)
28 USE MPI
29 INTEGER m, map(m), comm, p, sizeofreal, win, ierr
30 REAL A(m), B(m)
31
32 CALL MPI_TYPE_EXTENT(MPI_REAL, sizeofreal, ierr)
33 CALL MPI_WIN_CREATE(B, m*sizeofreal, sizeofreal, MPI_INFO_NULL, &
34                   comm, win, ierr)
35
36 CALL MPI_WIN_FENCE(0, win, ierr)
37 DO i=1,m
38   j = map(i)/p
39   k = MOD(map(i),p)
40   CALL MPI_ACCUMULATE(A(i), 1, MPI_REAL, j, k, 1, MPI_REAL, &
41                     MPI_SUM, win, ierr)
42
43 END DO
44 CALL MPI_WIN_FENCE(0, win, ierr)
45
46 CALL MPI_WIN_FREE(win, ierr)
47 RETURN
48 END
```

This code is identical to the code in Example 10.2, page 316, except that a call to `get` has been replaced by a call to `accumulate`. (Note that, if `map` is one-to-one, then the code computes $B = A(\text{map}^{-1})$, which is the reverse assignment to the one computed in that previous example.) In a similar manner, we can replace in Example 10.1, page 314, the call to `get` by a call to `accumulate`, thus performing the computation with only one communication between any two processes.

10.4 Synchronization Calls

RMA communications fall in two categories:

- **active target** communication, where data is moved from the memory of one process to the memory of another, and both are explicitly involved in the communication. This communication pattern is similar to message passing, except that all the data transfer arguments are provided by one process, and the second process only participates in the synchronization.
- **passive target** communication, where data is moved from the memory of one process to the memory of another, and only the origin process is explicitly involved in the transfer. Thus, two origin processes may communicate by accessing the same location in a target window. The process that owns the target window may be distinct from the two communicating processes, in which case it does not participate explicitly in the communication. This communication paradigm is closest to a shared memory model, where shared data can be accessed by all processes, irrespective of location.

RMA communication calls with argument `win` must occur at a process only within an **access epoch** for `win`. Such an epoch starts with an RMA synchronization call on `win`; it proceeds with zero or more RMA communication calls (`MPI_PUT`, `MPI_GET` or `MPI_ACCUMULATE`) on `win`; it completes with another synchronization call on `win`. This allows users to amortize one synchronization with multiple data transfers and provide implementors more flexibility in the implementation of RMA operations.

Distinct access epochs for `win` at the same process must be disjoint. On the other hand, epochs pertaining to different `win` arguments may overlap. Local operations or other MPI calls may also occur during an epoch.

In active target communication, a target window can be accessed by RMA operations only within an **exposure epoch**. Such an epoch is started and completed by RMA synchronization calls executed by the target process. Distinct exposure epochs at a process on the same window must be disjoint, but such an exposure epoch may overlap with exposure epochs on other windows or with access epochs for the same or other `win` arguments. There is a one-to-one matching between access epochs at origin processes and exposure epochs on target processes: RMA operations issued by an origin process for a target window will access that target window during the same exposure epoch if and only if they were issued during the same access epoch.

In passive target communication the target process does not execute RMA synchronization calls, and there is no concept of an exposure epoch.

MPI provides three synchronization mechanisms:

1. The `MPI_WIN_FENCE` collective synchronization call supports a simple synchronization pattern that is often used in parallel computations: namely a loosely-synchronous

1 model, where global computation phases alternate with global communication phases.
 2 This mechanism is most useful for loosely synchronous algorithms where the graph
 3 of communicating processes changes very frequently, or where each process communi-
 4 cates with many others.

5 This call is used for active target communication. An access epoch at an origin
 6 process or an exposure epoch at a target process are started and completed by calls to
 7 `MPI_WIN_FENCE`. A process can access windows at all processes in the group of `win`
 8 during such an access epoch, and the local window can be accessed by all processes
 9 in the group of `win` during such an exposure epoch.

- 10
 11 2. The four functions `MPI_WIN_START`, `MPI_WIN_COMPLETE`, `MPI_WIN_POST` and
 12 `MPI_WIN_WAIT` can be used to restrict synchronization to the minimum: only pairs
 13 of communicating processes synchronize, and they do so only when a synchronization
 14 is needed to order correctly RMA accesses to a window with respect to local accesses
 15 to that same window. This mechanism may be more efficient when each process
 16 communicates with few (logical) neighbors, and the communication graph is fixed or
 17 changes infrequently.

18 These calls are used for active target communication. An access epoch is started
 19 at the origin process by a call to `MPI_WIN_START` and is terminated by a call to
 20 `MPI_WIN_COMPLETE`. The start call has a group argument that specifies the group
 21 of target processes for that epoch. An exposure epoch is started at the target process
 22 by a call to `MPI_WIN_POST` and is completed by a call to `MPI_WIN_WAIT`. The post
 23 call has a group argument that specifies the set of origin processes for that epoch.

- 24
 25 3. Finally, shared and exclusive locks are provided by the two functions `MPI_WIN_LOCK`
 26 and `MPI_WIN_UNLOCK`. Lock synchronization is useful for MPI applications that emu-
 27 late a shared memory model via MPI calls; e.g., in a “billboard” model, where
 28 processes can, at random times, access or update different parts of the billboard.

29 These two calls provide passive target communication. An access epoch is started by
 30 a call to `MPI_WIN_LOCK` and terminated by a call to `MPI_WIN_UNLOCK`. Only one
 31 target window can be accessed during that epoch with `win`.

32
 33 Figure 10.1 illustrates the general synchronization pattern for active target communi-
 34 cation. The synchronization between `post` and `start` ensures that the put call of the origin
 35 process does not start until the target process exposes the window (with the `post` call);
 36 the target process will expose the window only after preceding local accesses to the window
 37 have completed. The synchronization between `complete` and `wait` ensures that the put call
 38 of the origin process completes before the window is unexposed (with the `wait` call). The
 39 target process will execute following local accesses to the target window only after the `wait`
 40 returned.

41 Figure 10.1 shows operations occurring in the natural temporal order implied by the
 42 synchronizations: the `post` occurs before the matching `start`, and `complete` occurs before
 43 the matching `wait`. However, such **strong** synchronization is more than needed for correct
 44 ordering of window accesses. The semantics of MPI calls allow **weak** synchronization, as
 45 illustrated in Figure 10.2. The access to the target window is delayed until the window is ex-
 46 posed, after the `post`. However the `start` may complete earlier; the `put` and `complete` may
 47 also terminate earlier, if put data is buffered by the implementation. The synchronization
 48

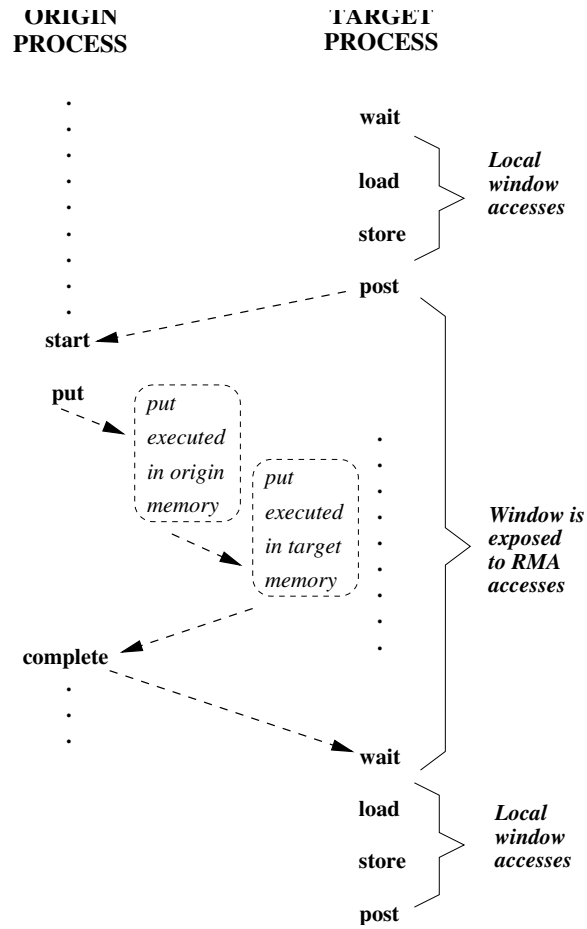


Figure 10.1: active target communication. Dashed arrows represent synchronizations (ordering of events).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

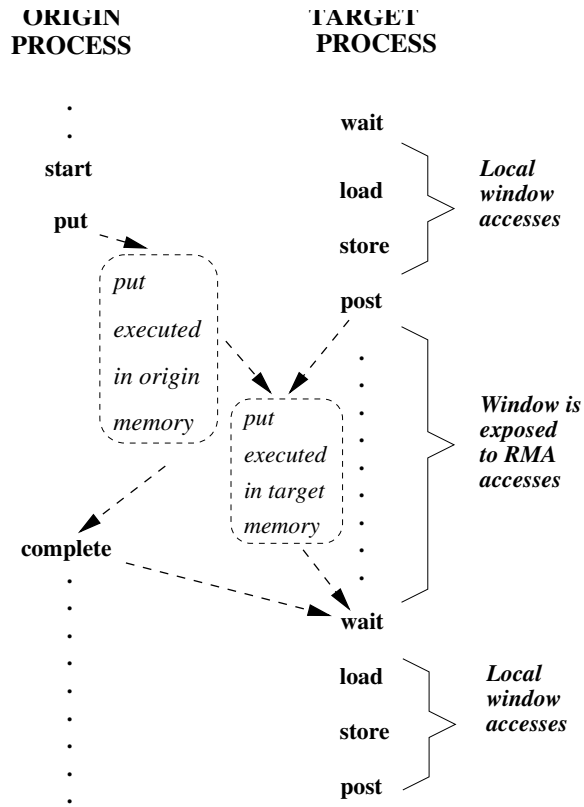


Figure 10.2: active target communication, with weak synchronization. Dashed arrows represent synchronizations (ordering of events)

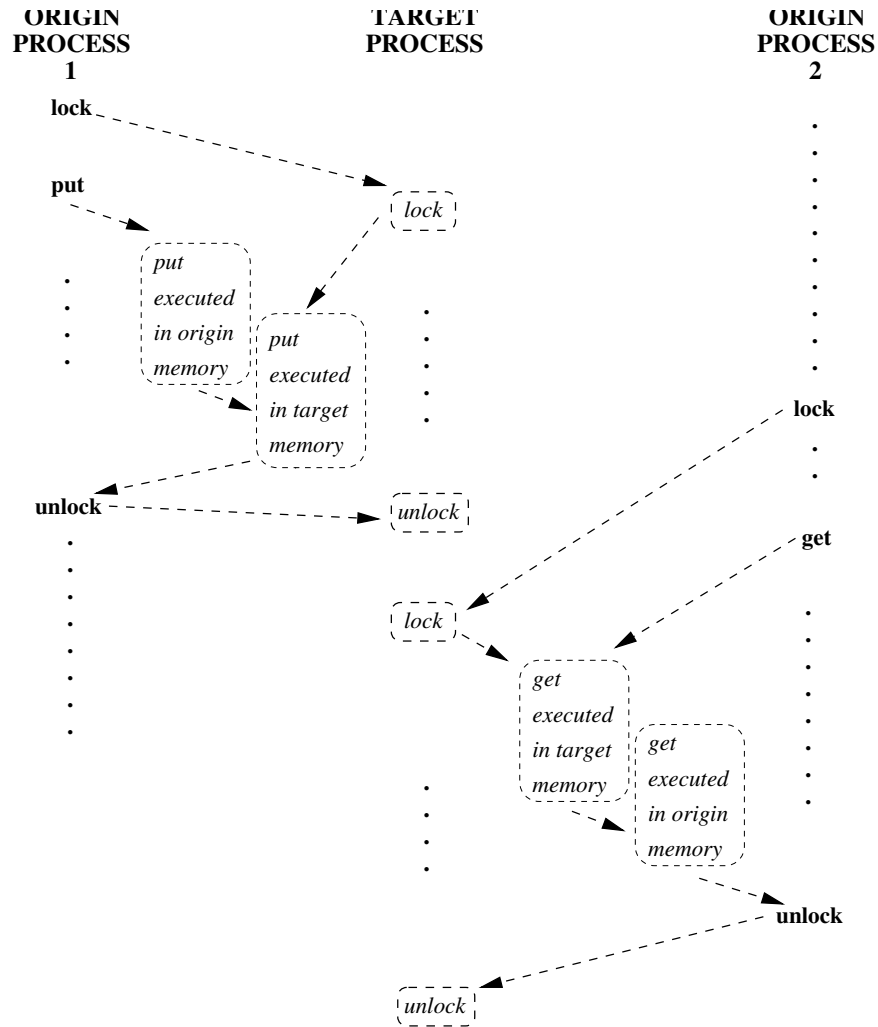


Figure 10.3: passive target communication. Dashed arrows represent synchronizations (ordering of events).

calls order correctly window accesses, but do not necessarily synchronize other operations. This weaker synchronization semantic allows for more efficient implementations.

Figure 10.3 illustrates the general synchronization pattern for passive target communication. The first origin process communicates data to the second origin process, through the memory of the target process; the target process is not explicitly involved in the communication. The **lock** and **unlock** calls ensure that the two RMA accesses do not occur concurrently. However, they do *not* ensure that the **put** by origin 1 will precede the **get** by origin 2.

10.4.1 Fence

```
MPI_WIN_FENCE(assert, win)
```

```
IN      assert          program assertion (integer)
IN      win             window object (handle)
```

```
int MPI_Win_fence(int assert, MPI_Win win)
```

```
MPI_WIN_FENCE(ASSERT, WIN, IERROR)
INTEGER ASSERT, WIN, IERROR
```

```
void MPI::Win::Fence(int assert) const
```

The MPI call `MPI_WIN_FENCE(assert, win)` synchronizes RMA calls on `win`. The call is collective on the group of `win`. All RMA operations on `win` originating at a given process and started before the fence call will complete at that process before the fence call returns. They will be completed at their target before the fence call returns at the target. RMA operations on `win` started by a process after the fence call returns will access their target window only after `MPI_WIN_FENCE` has been called by the target process.

The call completes an RMA access epoch if it was preceded by another fence call and the local process issued RMA communication calls on `win` between these two calls. The call completes an RMA exposure epoch if it was preceded by another fence call and the local window was the target of RMA accesses between these two calls. The call starts an RMA access epoch if it is followed by another fence call and by RMA communication calls issued between these two fence calls. The call starts an exposure epoch if it is followed by another fence call and the local window is the target of RMA accesses between these two fence calls. Thus, the fence call is equivalent to calls to a subset of `post`, `start`, `complete`, `wait`.

A fence call usually entails a barrier synchronization: a process completes a call to `MPI_WIN_FENCE` only after all other processes in the group entered their matching call. However, a call to `MPI_WIN_FENCE` that is known not to end any epoch (in particular, a call with `assert = MPI_MODE_NOPRECEDE`) does not necessarily act as a barrier.

The `assert` argument is used to provide assertions on the context of the call that may be used for various optimizations. This is described in Section 10.4.4. A value of `assert = 0` is always valid.

Advice to users. Calls to `MPI_WIN_FENCE` should both precede and follow calls to `put`, `get` or `accumulate` that are synchronized with fence calls. (*End of advice to users.*)

10.4.2 General Active Target Synchronization

```
MPI_WIN_START(group, assert, win)
```

| | | |
|----|--------|------------------------------------|
| IN | group | group of target processes (handle) |
| IN | assert | program assertion (integer) |
| IN | win | window object (handle) |

```
int MPI_Win_start(MPI_Group group, int assert, MPI_Win win)
```

```
MPI_WIN_START(GROUP, ASSERT, WIN, IERROR)
    INTEGER GROUP, ASSERT, WIN, IERROR
```

```
void MPI::Win::Start(const MPI::Group& group, int assert) const
```

Starts an RMA access epoch for win. RMA calls issued on win during this epoch must access only windows at processes in group. Each process in group must issue a matching call to MPI_WIN_POST. RMA accesses to each target window will be delayed, if necessary, until the target process executed the matching call to MPI_WIN_POST. MPI_WIN_START is allowed to block until the corresponding MPI_WIN_POST calls are executed, but is not required to.

The assert argument is used to provide assertions on the context of the call that may be used for various optimizations. This is described in Section 10.4.4. A value of assert = 0 is always valid.

```
MPI_WIN_COMPLETE(win)
```

| | | |
|----|-----|------------------------|
| IN | win | window object (handle) |
|----|-----|------------------------|

```
int MPI_Win_complete(MPI_Win win)
```

```
MPI_WIN_COMPLETE(WIN, IERROR)
    INTEGER WIN, IERROR
```

```
void MPI::Win::Complete() const
```

Completes an RMA access epoch on win started by a call to MPI_WIN_START. All RMA communication calls issued on win during this epoch will have completed at the origin when the call returns.

MPI_WIN_COMPLETE enforces completion of preceding RMA calls at the origin, but not at the target. A put or accumulate call may not have completed at the target when it has completed at the origin.

Consider the sequence of calls in the example below.

Example 10.4

```
MPI_Win_start(group, flag, win);
MPI_Put(...,win);
MPI_Win_complete(win);
```

1 The call to `MPI_WIN_COMPLETE` does not return until the put call has completed
 2 at the origin; and the target window will be accessed by the put operation only after
 3 the call to `MPI_WIN_START` has matched a call to `MPI_WIN_POST` by the target process.
 4 This still leaves much choice to implementors. The call to `MPI_WIN_START` can block
 5 until the matching call to `MPI_WIN_POST` occurs at all target processes. One can also
 6 have implementations where the call to `MPI_WIN_START` is nonblocking, but the call to
 7 `MPI_PUT` blocks until the matching call to `MPI_WIN_POST` occurred; or implementations
 8 where the first two calls are nonblocking, but the call to `MPI_WIN_COMPLETE` blocks
 9 until the call to `MPI_WIN_POST` occurred; or even implementations where all three calls
 10 can complete before any target process called `MPI_WIN_POST` — the data put must be
 11 buffered, in this last case, so as to allow the put to complete at the origin ahead of its
 12 completion at the target. However, once the call to `MPI_WIN_POST` is issued, the sequence
 13 above must complete, without further dependencies.

14
 15
 16 `MPI_WIN_POST(group, assert, win)`

| | | | |
|----|----|--------|------------------------------------|
| 17 | IN | group | group of origin processes (handle) |
| 18 | IN | assert | program assertion (integer) |
| 19 | IN | win | window object (handle) |

20
 21
 22 `int MPI_Win_post(MPI_Group group, int assert, MPI_Win win)`

23 `MPI_WIN_POST(GROUP, ASSERT, WIN, IERROR)`
 24 `INTEGER GROUP, ASSERT, WIN, IERROR`

25
 26 `void MPI::Win::Post(const MPI::Group& group, int assert) const`

27
 28 Starts an RMA exposure epoch for the local window associated with `win`. Only processes
 29 in `group` should access the window with RMA calls on `win` during this epoch. Each process
 30 in `group` must issue a matching call to `MPI_WIN_START`. `MPI_WIN_POST` does not block.

31
 32 `MPI_WIN_WAIT(win)`

| | | | |
|----|----|-----|------------------------|
| 33 | IN | win | window object (handle) |
|----|----|-----|------------------------|

34
 35
 36 `int MPI_Win_wait(MPI_Win win)`

37 `MPI_WIN_WAIT(WIN, IERROR)`
 38 `INTEGER WIN, IERROR`

39
 40 `void MPI::Win::Wait() const`

41 Completes an RMA exposure epoch started by a call to `MPI_WIN_POST` on `win`. This
 42 call matches calls to `MPI_WIN_COMPLETE(win)` issued by each of the origin processes that
 43 were granted access to the window during this epoch. The call to `MPI_WIN_WAIT` will block
 44 until all matching calls to `MPI_WIN_COMPLETE` have occurred. This guarantees that all
 45 these origin processes have completed their RMA accesses to the local window. When the
 46 call returns, all these RMA accesses will have completed at the target window.

47 Figure 10.4 illustrates the use of these four functions. Process 0 puts data in the
 48

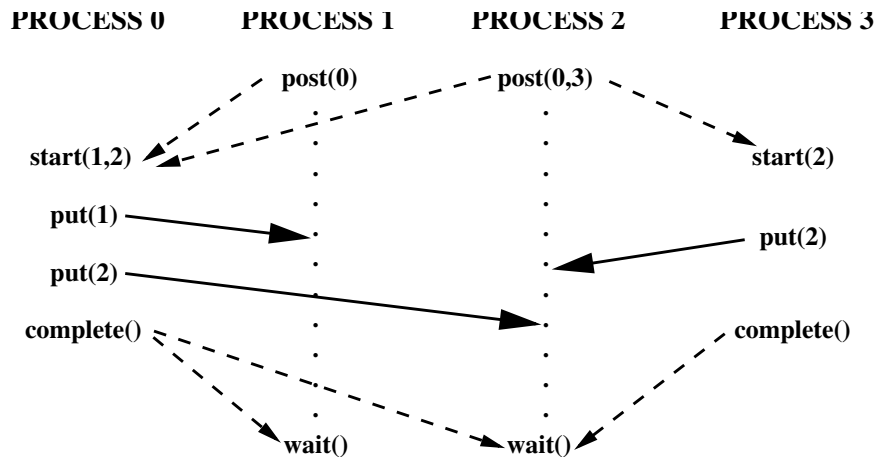


Figure 10.4: active target communication. Dashed arrows represent synchronizations and solid arrows represent data transfer.

windows of processes 1 and 2 and process 3 puts data in the window of process 2. Each start call lists the ranks of the processes whose windows will be accessed; each post call lists the ranks of the processes that access the local window. The figure illustrates a possible timing for the events, assuming strong synchronization; in a weak synchronization, the start, put or complete calls may occur ahead of the matching post calls.

`MPI_WIN_TEST(win, flag)`

| | | |
|-----|------|------------------------|
| IN | win | window object (handle) |
| OUT | flag | success flag (logical) |

```
int MPI_Win_test(MPI_Win win, int *flag)
```

```
MPI_WIN_TEST(WIN, FLAG, IERROR)
    INTEGER WIN, IERROR
    LOGICAL FLAG
```

```
bool MPI::Win::Test() const
```

This is the nonblocking version of `MPI_WIN_WAIT`. It returns `flag = true` if `MPI_WIN_WAIT` would return, `flag = false`, otherwise. The effect of return of `MPI_WIN_TEST` with `flag = true` is the same as the effect of a return of `MPI_WIN_WAIT`. If `flag = false` is returned, then the call has no visible effect.

`MPI_WIN_TEST` should be invoked only where `MPI_WIN_WAIT` can be invoked. Once the call has returned `flag = true`, it must not be invoked anew, until the window is posted anew.

Assume that window `win` is associated with a “hidden” communicator `wincomm`, used for communication by the processes of `win`. The rules for matching of post and start calls and for matching complete and wait call can be derived from the rules for matching sends and receives, by considering the following (partial) model implementation.

`MPI_WIN_POST(group,0,win)` initiate a nonblocking send with tag `tag0` to each process in

1 group, using wincomm. No need to wait for the completion of these sends.

2
3 MPI_WIN_START(group,0,win) initiate a nonblocking receive with tag
4 tag0 from each process in group, using wincomm. An RMA access to a window in
5 target process i is delayed until the receive from i is completed.

6 MPI_WIN_COMPLETE(win) initiate a nonblocking send with tag tag1 to each process in the
7 group of the preceding start call. No need to wait for the completion of these sends.

8
9 MPI_WIN_WAIT(win) initiate a nonblocking receive with tag tag1 from each process in the
10 group of the preceding post call. Wait for the completion of all receives.

11
12 No races can occur in a correct program: each of the sends matches a unique receive,
13 and vice-versa.

14 *Rationale.* The design for general active target synchronization requires the user to
15 provide complete information on the communication pattern, at each end of a com-
16 munication link: each origin specifies a list of targets, and each target specifies a list
17 of origins. This provides maximum flexibility (hence, efficiency) for the implementor:
18 each synchronization can be initiated by either side, since each “knows” the identity of
19 the other. This also provides maximum protection from possible races. On the other
20 hand, the design requires more information than RMA needs, in general: in general,
21 it is sufficient for the origin to know the rank of the target, but not vice versa. Users
22 that want more “anonymous” communication will be required to use the fence or lock
23 mechanisms. (*End of rationale.*)

24
25 *Advice to users.* Assume a communication pattern that is represented by a di-
26 rected graph $G = \langle V, E \rangle$, where $V = \{0, \dots, n - 1\}$ and $ij \in E$ if origin
27 process i accesses the window at target process j . Then each process i issues a
28 call to MPI_WIN_POST(ingroup $_i$, ...), followed by a call to
29 MPI_WIN_START(outgroup $_i$, ...), where outgroup $_i = \{j : ij \in E\}$ and ingroup $_i =$
30 $\{j : ji \in E\}$. A call is a noop, and can be skipped, if the group argument is empty.
31 After the communications calls, each process that issued a start will issue a complete.
32 Finally, each process that issued a post will issue a wait.

33 Note that each process may call with a group argument that has different members.
34 (*End of advice to users.*)

36 10.4.3 Lock

37
38
39 MPI_WIN_LOCK(lock_type, rank, assert, win)

| | | | |
|----|----|-----------|---|
| 41 | IN | lock_type | either MPI_LOCK_EXCLUSIVE or |
| 42 | | | MPI_LOCK_SHARED (state) |
| 43 | IN | rank | rank of locked window (nonnegative integer) |
| 44 | IN | assert | program assertion (integer) |
| 45 | | | |
| 46 | IN | win | window object (handle) |

47
48 int MPI_Win_lock(int lock_type, int rank, int assert, MPI_Win win)

```
MPI_WIN_LOCK(LOCK_TYPE, RANK, ASSERT, WIN, IERROR)
```

```
    INTEGER LOCK_TYPE, RANK, ASSERT, WIN, IERROR
```

```
void MPI::Win::Lock(int lock_type, int rank, int assert) const
```

Starts an RMA access epoch. Only the window at the process with rank `rank` can be accessed by RMA operations on `win` during that epoch.

```
MPI_WIN_UNLOCK(rank, win)
```

```
    IN      rank                rank of window (nonnegative integer)
```

```
    IN      win                 window object (handle)
```

```
int MPI_Win_unlock(int rank, MPI_Win win)
```

```
MPI_WIN_UNLOCK(RANK, WIN, IERROR)
```

```
    INTEGER RANK, WIN, IERROR
```

```
void MPI::Win::Unlock(int rank) const
```

Completes an RMA access epoch started by a call to `MPI_WIN_LOCK(...,win)`. RMA operations issued during this period will have completed both at the origin and at the target when the call returns.

Locks are used to protect accesses to the locked target window effected by RMA calls issued between the lock and unlock call, and to protect local load/store accesses to a locked local window executed between the lock and unlock call. Accesses that are protected by an exclusive lock will not be concurrent at the window site with other accesses to the same window that are lock protected. Accesses that are protected by a shared lock will not be concurrent at the window site with accesses protected by an exclusive lock to the same window.

It is erroneous to have a window locked and exposed (in an exposure epoch) concurrently. I.e., a process may not call `MPI_WIN_LOCK` to lock a target window if the target process has called `MPI_WIN_POST` and has not yet called `MPI_WIN_WAIT`; it is erroneous to call `MPI_WIN_POST` while the local window is locked.

Rationale. An alternative is to require MPI to enforce mutual exclusion between exposure epochs and locking periods. But this would entail additional overheads when locks or active target synchronization do not interact in support of those rare interactions between the two mechanisms. The programming style that we encourage here is that a set of windows is used with only one synchronization mechanism at a time, with shifts from one mechanism to another being rare and involving global synchronization. (*End of rationale.*)

Advice to users. Users need to use explicit synchronization code in order to enforce mutual exclusion between locking periods and exposure epochs on a window. (*End of advice to users.*)

Implementors may restrict the use of RMA communication that is synchronized by lock calls to windows in memory allocated by `MPI_ALLOC_MEM` (Section 7.2, page 250). Locks can be used portably only in such memory.

Rationale. The implementation of passive target communication when memory is not shared requires an asynchronous agent. Such an agent can be implemented more easily, and can achieve better performance, if restricted to specially allocated memory. It can be avoided altogether if shared memory is used. It seems natural to impose restrictions that allows one to use shared memory for 3-rd party communication in shared memory machines.

The downside of this decision is that passive target communication cannot be used without taking advantage of nonstandard Fortran features: namely, the availability of C-like pointers; these are not supported by some Fortran compilers (g77 and Windows/NT compilers, at the time of writing). Also, passive target communication cannot be portably targeted to COMMON blocks, or other statically declared Fortran arrays. (*End of rationale.*)

Consider the sequence of calls in the example below.

Example 10.5

```
MPI_Win_lock(MPI_LOCK_EXCLUSIVE, rank, assert, win)
MPI_Put(..., rank, ..., win)
MPI_Win_unlock(rank, win)
```

The call to MPI_WIN_UNLOCK will not return until the put transfer has completed at the origin and at the target. This still leaves much freedom to implementors. The call to MPI_WIN_LOCK may block until an exclusive lock on the window is acquired; or, the call MPI_WIN_LOCK may not block, while the call to MPI_PUT blocks until a lock is acquired; or, the first two calls may not block, while MPI_WIN_UNLOCK blocks until a lock is acquired — the update of the target window is then postponed until the call to MPI_WIN_UNLOCK occurs. However, if the call to MPI_WIN_LOCK is used to lock a local window, then the call must block until the lock is acquired, since the lock may protect local load/store accesses to the window issued after the lock call returns.

10.4.4 Assertions

The assert argument in the calls MPI_WIN_POST, MPI_WIN_START, MPI_WIN_FENCE and MPI_WIN_LOCK is used to provide assertions on the context of the call that may be used to optimize performance. The assert argument does not change program semantics if it provides correct information on the program — it is erroneous to provides incorrect information. Users may always provide assert = 0 to indicate a general case, where no guarantees are made.

Advice to users. Many implementations may not take advantage of the information in assert; some of the information is relevant only for noncoherent, shared memory machines. Users should consult their implementation manual to find which information is useful on each system. On the other hand, applications that provide correct assertions whenever applicable are portable and will take advantage of assertion specific optimizations, whenever available. (*End of advice to users.*)

Advice to implementors. Implementations can always ignore the assert argument. Implementors should document which assert values are significant on their implementation. (*End of advice to implementors.*)

assert is the bit-vector OR of zero or more of the following integer constants: MPI_MODE_NOCHECK, MPI_MODE_NOSTORE, MPI_MODE_NOPUT, MPI_MODE_NOPRECEDE and MPI_MODE_NOSUCCEED. The significant options are listed below, for each call.

Advice to users. C/C++ users can use bit vector or (|) to combine these constants; Fortran 90 users can use the bit-vector IOR intrinsic. Fortran 77 users can use (non-portably) bit vector IOR on systems that support it. Alternatively, Fortran users can portably use integer addition to OR the constants (each constant should appear at most once in the addition!). (*End of advice to users.*)

MPI_WIN_START:

MPI_MODE_NOCHECK — the matching calls to MPI_WIN_POST have already completed on all target processes when the call to MPI_WIN_START is made. The nocheck option can be specified in a start call if and only if it is specified in each matching post call. This is similar to the optimization of “ready-send” that may save a handshake when the handshake is implicit in the code. (However, ready-send is matched by a regular receive, whereas both start and post must specify the nocheck option.)

MPI_WIN_POST:

MPI_MODE_NOCHECK — the matching calls to MPI_WIN_START have not yet occurred on any origin processes when the call to MPI_WIN_POST is made. The nocheck option can be specified by a post call if and only if it is specified by each matching start call.

MPI_MODE_NOSTORE — the local window was not updated by local stores (or local get or receive calls) since last synchronization. This may avoid the need for cache synchronization at the post call.

MPI_MODE_NOPUT — the local window will not be updated by put or accumulate calls after the post call, until the ensuing (wait) synchronization. This may avoid the need for cache synchronization at the wait call.

MPI_WIN_FENCE:

MPI_MODE_NOSTORE — the local window was not updated by local stores (or local get or receive calls) since last synchronization.

MPI_MODE_NOPUT — the local window will not be updated by put or accumulate calls after the fence call, until the ensuing (fence) synchronization.

MPI_MODE_NOPRECEDE — the fence does not complete any sequence of locally issued RMA calls. If this assertion is given by any process in the window group, then it must be given by all processes in the group.

MPI_MODE_NOSUCCEED — the fence does not start any sequence of locally issued RMA calls. If the assertion is given by any process in the window group, then it must be given by all processes in the group.

MPI_WIN_LOCK:

1 MPI_MODE_NOCHECK — no other process holds, or will attempt to acquire a con-
 2 flicting lock, while the caller holds the window lock. This is useful when mutual
 3 exclusion is achieved by other means, but the coherence operations that may be
 4 attached to the lock and unlock calls are still required.

5
 6 *Advice to users.* Note that the nostore and noprecede flags provide information on
 7 what happened *before* the call; the noput and nosucceed flags provide information on
 8 what will happen *after* the call. (*End of advice to users.*)

10.4.5 Miscellaneous Clarifications

11 Once an RMA routine completes, it is safe to free any opaque objects passed as argument to
 12 that routine. For example, the datatype argument of a MPI_PUT call can be freed as soon
 13 as the call returns, even though the communication may not be complete.

14 As in message passing, datatypes must be committed before they can be used in RMA
 15 communication.

10.5 Examples

16
 17
 18 **Example 10.6** The following example shows a generic loosely synchronous, iterative code,
 19 using fence synchronization. The window at each process consists of array **A**, which contains
 20 the origin and target buffers of the put calls.

```

21       ...
22       while(!converged(A)){
23           update(A);
24           MPI_Win_fence(MPI_MODE_NOPRECEDE, win);
25           for(i=0; i < toneighbors; i++)
26               MPI_Put(&frombuf[i], 1, fromtype[i], toneighbor[i],
27                       todisp[i], 1, totype[i], win);
28           MPI_Win_fence((MPI_MODE_NOSTORE | MPI_MODE_NOSUCCEED), win);
29       }
30
31
32
33

```

34 The same code could be written with get, rather than put. Note that, during the commu-
 35 nication phase, each window is concurrently read (as origin buffer of puts) and written (as
 36 target buffer of puts). This is OK, provided that there is no overlap between the target
 37 buffer of a put and another communication buffer.

38 **Example 10.7** Same generic example, with more computation/communication overlap.
 39 We assume that the update phase is broken in two subphases: the first, where the “bound-
 40 ary,” which is involved in communication, is updated, and the second, where the “core,”
 41 which neither use nor provide communicated data, is updated.

```

42       ...
43       while(!converged(A)){
44           update_boundary(A);
45           MPI_Win_fence((MPI_MODE_NOPUT | MPI_MODE_NOPRECEDE), win);
46           for(i=0; i < fromneighbors; i++)
47               MPI_Get(&tobuf[i], 1, totype[i], fromneighbor[i],
48

```

```

        fromdisp[i], 1, fromtype[i], win);
    update_core(A);
    MPI_Win_fence(MPI_MODE_NOSUCCEED, win);
}

```

The get communication can be concurrent with the core update, since they do not access the same locations, and the local update of the origin buffer by the get call can be concurrent with the local update of the core by the `update_core` call. In order to get similar overlap with put communication we would need to use separate windows for the core and for the boundary. This is required because we do not allow local stores to be concurrent with puts on the same, or on overlapping, windows.

Example 10.8 Same code as in Example 10.6, rewritten using `post-start-complete-wait`.

```

...
while(!converged(A)){
    update(A);
    MPI_Win_post(fromgroup, 0, win);
    MPI_Win_start(togroup, 0, win);
    for(i=0; i < toneighbors; i++)
        MPI_Put(&frombuf[i], 1, fromtype[i], toneighbor[i],
               todisp[i], 1, totype[i], win);
    MPI_Win_complete(win);
    MPI_Win_wait(win);
}

```

Example 10.9 Same example, with split phases, as in Example 10.7.

```

...
while(!converged(A)){
    update_boundary(A);
    MPI_Win_post(togroup, MPI_MODE_NOPUT, win);
    MPI_Win_start(fromgroup, 0, win);
    for(i=0; i < fromneighbors; i++)
        MPI_Get(&tobuf[i], 1, totype[i], fromneighbor[i],
               fromdisp[i], 1, fromtype[i], win);
    update_core(A);
    MPI_Win_complete(win);
    MPI_Win_wait(win);
}

```

Example 10.10 A checkerboard, or double buffer communication pattern, that allows more computation/communication overlap. Array `A0` is updated using values of array `A1`, and vice versa. We assume that communication is symmetric: if process A gets data from process B, then process B gets data from process A. Window `wini` consists of array `Ai`.

```

...
if (!converged(A0,A1))
    MPI_Win_post(neighbors, (MPI_MODE_NOCHECK | MPI_MODE_NOPUT), win0);
MPI_Barrier(comm0);

```

```

1  /* the barrier is needed because the start call inside the
2  loop uses the nocheck option */
3  while(!converged(A0, A1)){
4      /* communication on A0 and computation on A1 */
5      update2(A1, A0); /* local update of A1 that depends on A0 (and A1) */
6      MPI_Win_start(neighbors, MPI_MODE_NOCHECK, win0);
7      for(i=0; i < neighbors; i++)
8          MPI_Get(&tobuf0[i], 1, totype0[i], neighbor[i],
9                 fromdisp0[i], 1, fromtype0[i], win0);
10     update1(A1); /* local update of A1 that is
11                 concurrent with communication that updates A0 */
12     MPI_Win_post(neighbors, (MPI_MODE_NOCHECK | MPI_MODE_NOPUT), win1);
13     MPI_Win_complete(win0);
14     MPI_Win_wait(win0);
15
16     /* communication on A1 and computation on A0 */
17     update2(A0, A1); /* local update of A0 that depends on A1 (and A0)*/
18     MPI_Win_start(neighbors, MPI_MODE_NOCHECK, win1);
19     for(i=0; i < neighbors; i++)
20         MPI_Get(&tobuf1[i], 1, totype1[i], neighbor[i],
21                fromdisp1[i], 1, fromtype1[i], win1);
22     update1(A0); /* local update of A0 that depends on A0 only,
23                 concurrent with communication that updates A1 */
24     if (!converged(A0,A1))
25         MPI_Win_post(neighbors, (MPI_MODE_NOCHECK | MPI_MODE_NOPUT), win0);
26     MPI_Win_complete(win1);
27     MPI_Win_wait(win1);
28 }

```

29

30 A process posts the local window associated with `win0` before it completes RMA accesses
31 to the remote windows associated with `win1`. When the `wait(win1)` call returns, then all
32 neighbors of the calling process have posted the windows associated with `win0`. Conversely,
33 when the `wait(win0)` call returns, then all neighbors of the calling process have posted the
34 windows associated with `win1`. Therefore, the `nocheck` option can be used with the calls to
35 `MPI_WIN_START`.

36 Put calls can be used, instead of get calls, if the area of array `A0` (resp. `A1`) used by
37 the `update(A1, A0)` (resp. `update(A0, A1)`) call is disjoint from the area modified by the
38 RMA communication. On some systems, a put call may be more efficient than a get call,
39 as it requires information exchange only in one direction.

40

41 10.6 Error Handling

42

43 10.6.1 Error Handlers

44

45 Errors occurring during calls to `MPI_WIN_CREATE(...,comm,...)` cause the error handler
46 currently associated with `comm` to be invoked. All other RMA calls have an input `win`
47 argument. When an error occurs during such a call, the error handler currently associated
48 with `win` is invoked.

The default error handler associated with `win` is `MPI_ERRORS_ARE_FATAL`. Users may change this default by explicitly associating a new error handler with `win` (see Section 7.3.1, page 253).

10.6.2 Error Classes

The following new error classes are defined

| | |
|-----------------------------------|--|
| <code>MPI_ERR_WIN</code> | invalid <code>win</code> argument |
| <code>MPI_ERR_BASE</code> | invalid <code>base</code> argument |
| <code>MPI_ERR_SIZE</code> | invalid <code>size</code> argument |
| <code>MPI_ERR_DISP</code> | invalid <code>disp</code> argument |
| <code>MPI_ERR_LOCKTYPE</code> | invalid <code>locktype</code> argument |
| <code>MPI_ERR_ASSERT</code> | invalid <code>assert</code> argument |
| <code>MPI_ERR_RMA_CONFLICT</code> | conflicting accesses to window |
| <code>MPI_ERR_RMA_SYNC</code> | wrong synchronization of RMA calls |

Table 10.1: Error classes in one-sided communication routines

10.7 Semantics and Correctness

The semantics of RMA operations is best understood by assuming that the system maintains a separate *public* copy of each window, in addition to the original location in process memory (the *private* window copy). There is only one instance of each variable in process memory, but a distinct *public* copy of the variable for each window that contains it. A load accesses the instance in process memory (this includes MPI sends). A store accesses and updates the instance in process memory (this includes MPI receives), but the update may affect other public copies of the same locations. A `get` on a window accesses the public copy of that window. A `put` or `accumulate` on a window accesses and updates the public copy of that window, but the update may affect the private copy of the same locations in process memory, and public copies of other overlapping windows. This is illustrated in Figure 10.5.

The following rules specify the latest time at which an operation must complete at the origin or the target. The update performed by a `get` call in the origin process memory is visible when the `get` operation is complete at the origin (or earlier); the update performed by a `put` or `accumulate` call in the public copy of the target window is visible when the `put` or `accumulate` has completed at the target (or earlier). The rules also specifies the latest time at which an update of one window copy becomes visible in another overlapping copy.

1. An RMA operation is completed at the origin by the ensuing call to `MPI_WIN_COMPLETE`, `MPI_WIN_FENCE` or `MPI_WIN_UNLOCK` that synchronizes this access at the origin.
2. If an RMA operation is completed at the origin by a call to `MPI_WIN_FENCE` then the operation is completed at the target by the matching call to `MPI_WIN_FENCE` by the target process.

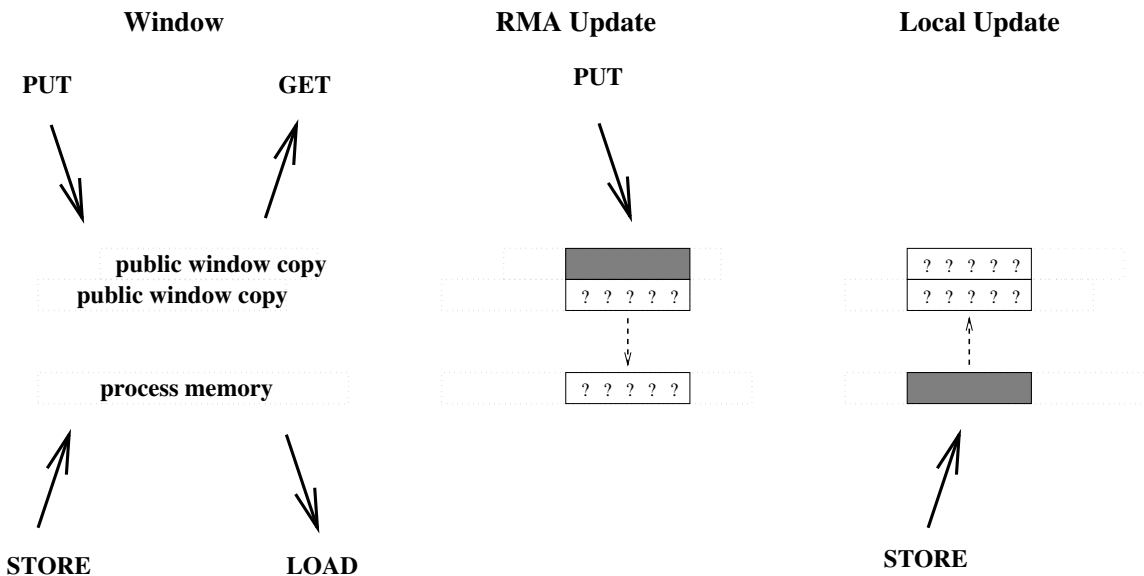


Figure 10.5: Schematic description of window

3. If an RMA operation is completed at the origin by a call to `MPI_WIN_COMPLETE` then the operation is completed at the target by the matching call to `MPI_WIN_WAIT` by the target process.
4. If an RMA operation is completed at the origin by a call to `MPI_WIN_UNLOCK` then the operation is completed at the target by that same call to `MPI_WIN_UNLOCK`.
5. An update of a location in a private window copy in process memory becomes visible in the public window copy at latest when an ensuing call to `MPI_WIN_POST`, `MPI_WIN_FENCE`, or `MPI_WIN_UNLOCK` is executed on that window by the window owner.
6. An update by a put or accumulate call to a public window copy becomes visible in the private copy in process memory at latest when an ensuing call to `MPI_WIN_WAIT`, `MPI_WIN_FENCE`, or `MPI_WIN_LOCK` is executed on that window by the window owner.

The `MPI_WIN_FENCE` or `MPI_WIN_WAIT` call that completes the transfer from public copy to private copy (6) is the same call that completes the put or accumulate operation in the window copy (2, 3). If a put or accumulate access was synchronized with a lock, then the update of the public window copy is complete as soon as the updating process executed `MPI_WIN_UNLOCK`. On the other hand, the update of private copy in the process memory may be delayed until the target process executes a synchronization call on that window (6). Thus, updates to process memory can always be delayed until the process executes a suitable synchronization call. Updates to a public window copy can also be delayed until the window owner executes a synchronization call, if fences or post-start-complete-wait synchronization is used. Only when lock synchronization is used does it becomes necessary to update the public window copy, even if the window owner does not execute any related synchronization call.

The rules above also define, by implication, when an update to a public window copy becomes visible in another overlapping public window copy. Consider, for example, two overlapping windows, win1 and win2. A call to `MPI_WIN_FENCE(0, win1)` by the window owner makes visible in the process memory previous updates to window win1 by remote processes. A subsequent call to `MPI_WIN_FENCE(0, win2)` makes these updates visible in the public copy of win2.

A correct program must obey the following rules.

1. A location in a window must not be accessed locally once an update to that location has started, until the update becomes visible in the private window copy in process memory.
2. A location in a window must not be accessed as a target of an RMA operation once an update to that location has started, until the update becomes visible in the public window copy. There is one exception to this rule, in the case where the same variable is updated by two concurrent accumulates that use the same operation, with the same predefined datatype, on the same window.
3. A put or accumulate must not access a target window once a local update or a put or accumulate update to another (overlapping) target window have started on a location in the target window, until the update becomes visible in the public copy of the window. Conversely, a local update in process memory to a location in a window must not start once a put or accumulate update to that target window has started, until the put or accumulate update becomes visible in process memory. In both cases, the restriction applies to operations even if they access disjoint locations in the window.

A program is erroneous if it violates these rules.

Rationale. The last constraint on correct RMA accesses may seem unduly restrictive, as it forbids concurrent accesses to nonoverlapping locations in a window. The reason for this constraint is that, on some architectures, explicit coherence restoring operations may be needed at synchronization points. A different operation may be needed for locations that were locally updated by stores and for locations that were remotely updated by put or accumulate operations. Without this constraint, the MPI library will have to track precisely which locations in a window were updated by a put or accumulate call. The additional overhead of maintaining such information is considered prohibitive. (*End of rationale.*)

Advice to users. A user can write correct programs by following the following rules:

fence: During each period between fence calls, each window is either updated by put or accumulate calls, or updated by local stores, but not both. Locations updated by put or accumulate calls should not be accessed during the same period (with the exception of concurrent updates to the same location by accumulate calls). Locations accessed by get calls should not be updated during the same period.

post-start-complete-wait: A window should not be updated locally while being posted, if it is being updated by put or accumulate calls. Locations updated by put or accumulate calls should not be accessed while the window is posted

(with the exception of concurrent updates to the same location by accumulate calls). Locations accessed by get calls should not be updated while the window is posted.

With the post-start synchronization, the target process can tell the origin process that its window is now ready for RMA access; with the complete-wait synchronization, the origin process can tell the target process that it has finished its RMA accesses to the window.

lock: Updates to the window are protected by exclusive locks if they may conflict. Nonconflicting accesses (such as read-only accesses or accumulate accesses) are protected by shared locks, both for local accesses and for RMA accesses.

changing window or synchronization mode: One can change synchronization mode, or change the window used to access a location that belongs to two overlapping windows, when the process memory and the window copy are guaranteed to have the same values. This is true after a local call to `MPI_WIN_FENCE`, if RMA accesses to the window are synchronized with fences; after a local call to `MPI_WIN_WAIT`, if the accesses are synchronized with post-start-complete-wait; after the call at the origin (local or remote) to `MPI_WIN_UNLOCK` if the accesses are synchronized with locks.

In addition, a process should not access the local buffer of a get operation until the operation is complete, and should not update the local buffer of a put or accumulate operation until that operation is complete. (*End of advice to users.*)

10.7.1 Atomicity

The outcome of concurrent accumulates to the same location, with the same operation and predefined datatype, is as if the accumulates were done at that location in some serial order. On the other hand, if two locations are both updated by two accumulate calls, then the updates may occur in reverse order at the two locations. Thus, there is no guarantee that the entire call to `MPI_ACCUMULATE` is executed atomically. The effect of this lack of atomicity is limited: The previous correctness conditions imply that a location updated by a call to `MPI_ACCUMULATE`, cannot be accessed by load or an RMA call other than accumulate, until the `MPI_ACCUMULATE` call has completed (at the target). Different interleavings can lead to different results only to the extent that computer arithmetics are not truly associative or commutative.

10.7.2 Progress

One-sided communication has the same progress requirements as point-to-point communication: once a communication is enabled, then it is guaranteed to complete. RMA calls must have local semantics, except when required for synchronization with other RMA calls.

There is some fuzziness in the definition of the time when a RMA communication becomes enabled. This fuzziness provides to the implementor more flexibility than with point-to-point communication. Access to a target window becomes enabled once the corresponding synchronization (such as `MPI_WIN_FENCE` or `MPI_WIN_POST`) has executed. On the origin process, an RMA communication may become enabled as soon as the corresponding put, get or accumulate call has executed, or as late as when the ensuing synchronization

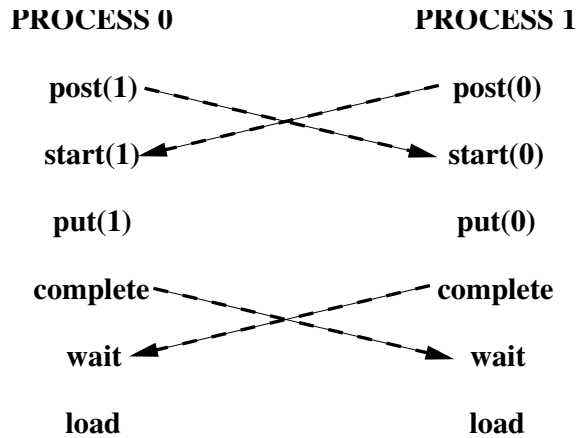


Figure 10.6: Symmetric communication

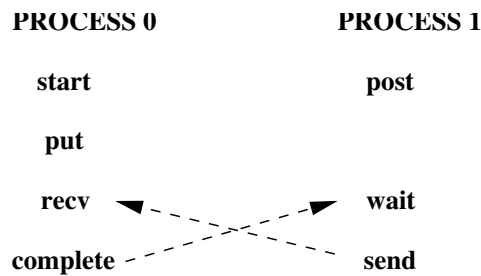


Figure 10.7: Deadlock situation

call is issued. Once the communication is enabled both at the origin and at the target, the communication must complete.

Consider the code fragment in Example 10.4, on page 325. Some of the calls may block if the target window is not posted. However, if the target window is posted, then the code fragment must complete. The data transfer may start as soon as the put call occur, but may be delayed until the ensuing complete call occurs.

Consider the code fragment in Example 10.5, on page 330. Some of the calls may block if another process holds a conflicting lock. However, if no conflicting lock is held, then the code fragment must complete.

Consider the code illustrated in Figure 10.6. Each process updates the window of the other process using a put operation, then accesses its own window. The post calls are nonblocking, and should complete. Once the post calls occur, RMA access to the windows is enabled, so that each process should complete the sequence of calls start-put-complete. Once these are done, the wait calls should complete at both processes. Thus, this communication should not deadlock, irrespective of the amount of data transferred.

Assume, in the last example, that the order of the post and start calls is reversed, at each process. Then, the code may deadlock, as each process may block on the start call, waiting for the matching post to occur. Similarly, the program will deadlock, if the order of the complete and wait calls is reversed, at each process.

The following two examples illustrate the fact that the synchronization between complete and wait is not symmetric: the wait call blocks until the complete executes, but not vice-versa. Consider the code illustrated in Figure 10.7. This code will deadlock: the wait

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

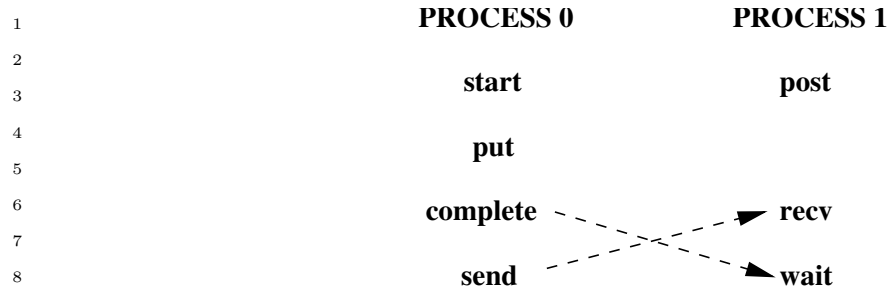


Figure 10.8: No deadlock

of process 1 blocks until process 0 calls complete, and the receive of process 0 blocks until process 1 calls send. Consider, on the other hand, the code illustrated in Figure 10.8. This code will not deadlock. Once process 1 calls post, then the sequence start, put, complete on process 0 can proceed to completion. Process 0 will reach the send call, allowing the receive call of process 1 to complete.

Rationale. MPI implementations must guarantee that a process makes progress on all enabled communications it participates in, while blocked on an MPI call. This is true for send-receive communication and applies to RMA communication as well. Thus, in the example in Figure 10.8, the put and complete calls of process 0 should complete while process 1 is blocked on the receive call. This may require the involvement of process 1, e.g., to transfer the data put, while it is blocked on the receive call.

A similar issue is whether such progress must occur while a process is busy computing, or blocked in a non-MPI call. Suppose that in the last example the send-receive pair is replaced by a write-to-socket/read-from-socket pair. Then MPI does not specify whether deadlock is avoided. Suppose that the blocking receive of process 1 is replaced by a very long compute loop. Then, according to one interpretation of the MPI standard, process 0 must return from the complete call after a bounded delay, even if process 1 does not reach any MPI call in this period of time. According to another interpretation, the complete call may block until process 1 reaches the wait call, or reaches another MPI call. The qualitative behavior is the same, under both interpretations, unless a process is caught in an infinite compute loop, in which case the difference may not matter. However, the quantitative expectations are different. Different MPI implementations reflect these different interpretations. While this ambiguity is unfortunate, it does not seem to affect many real codes. The MPI forum decided not to decide which interpretation of the standard is the correct one, since the issue is very contentious, and a decision would have much impact on implementors but less impact on users. (*End of rationale.*)

10.7.3 Registers and Compiler Optimizations

Advice to users. All the material in this section is an advice to users. (*End of advice to users.*)

A coherence problem exists between variables kept in registers and the memory value of these variables. An RMA call may access a variable in memory (or cache), while the

up-to-date value of this variable is in register. A get will not return the latest variable value, and a put may be overwritten when the register is stored back in memory.

The problem is illustrated by the following code:

| Source of Process 1 | Source of Process 2 | Executed in Process 2 |
|--------------------------------------|---------------------------------|---------------------------------------|
| <code>bbbb = 777</code> | <code>buff = 999</code> | <code>reg_A:=999</code> |
| <code>call MPI_WIN_FENCE</code> | <code>call MPI_WIN_FENCE</code> | |
| <code>call MPI_PUT(bbbb</code> | | <code>stop appl. thread</code> |
| <code>into buff of process 2)</code> | | <code>buff:=777 in PUT handler</code> |
| | | <code>continue appl. thread</code> |
| <code>call MPI_WIN_FENCE</code> | <code>call MPI_WIN_FENCE</code> | |
| | <code>ccc = buff</code> | <code>ccc:=reg_A</code> |

In this example, variable `buff` is allocated in the register `reg_A` and therefore `ccc` will have the old value of `buff` and not the new value `777`.

This problem, which also afflicts in some cases send/receive communication, is discussed more at length in Section 13.2.2.

MPI implementations will avoid this problem for standard conforming C programs. Many Fortran compilers will avoid this problem, without disabling compiler optimizations. However, in order to avoid register coherence problems in a completely portable manner, users should restrict their use of RMA windows to variables stored in `COMMON` blocks, or to variables that were declared `VOLATILE` (while `VOLATILE` is not a standard Fortran declaration, it is supported by many Fortran compilers). Details and an additional solution are discussed in Section 13.2.2, “A Problem with Register Optimization,” on page 454. See also, “Problems Due to Data Copying and Sequence Association,” on page 451, for additional Fortran problems.

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 11

External Interfaces

11.1 Introduction

This chapter begins with calls used to create **generalized requests**. The objective of this MPI-2 addition is to allow users of MPI to be able to create new nonblocking operations with an interface similar to what is present in MPI. This can be used to layer new functionality on top of MPI. Next, Section 11.3 deals with setting the information found in **status**. This is needed for generalized requests.

Section 11.4 allows users to associate names with communicators, windows, and datatypes. This will allow debuggers and profilers to identify communicators, windows, and datatypes with more useful labels. Section 11.5 allows users to add error codes, classes, and strings to MPI. With users being able to layer functionality on top of MPI, it is desirable for them to use the same error mechanisms found in MPI.

Section 11.6 deals with decoding datatypes. The opaque datatype object has found a number of uses outside MPI. Furthermore, a number of tools wish to display internal information about a datatype. To achieve this, datatype decoding functions are provided.

The chapter continues, in Section 11.7, with a discussion of how threads are to be handled in MPI-2. Although thread compliance is not required, the standard specifies how threads are to work if they are provided.

11.2 Generalized Requests

The goal of this MPI-2 extension is to allow users to define new nonblocking operations. Such an outstanding nonblocking operation is represented by a (generalized) request. A fundamental property of nonblocking operations is that progress toward the completion of this operation occurs asynchronously, i.e., concurrently with normal program execution. Typically, this requires execution of code concurrently with the execution of the user code, e.g., in a separate thread or in a signal handler. Operating systems provide a variety of mechanisms in support of concurrent execution. MPI does not attempt to standardize or replace these mechanisms: it is assumed programmers who wish to define new asynchronous operations will use the mechanisms provided by the underlying operating system. Thus, the calls in this section only provide a means for defining the effect of MPI calls such as **MPI_WAIT** or **MPI_CANCEL** when they apply to generalized requests, and for signaling to MPI the completion of a generalized operation.

Rationale. It is tempting to also define an MPI standard mechanism for achieving concurrent execution of user-defined nonblocking operations. However, it is very difficult to define such a mechanism without consideration of the specific mechanisms used in the operating system. The Forum feels that concurrency mechanisms are a proper part of the underlying operating system and should not be standardized by MPI; the MPI standard should only deal with the interaction of such mechanisms with MPI. (*End of rationale.*)

For a regular request, the operation associated with the request is performed by the MPI implementation, and the operation completes without intervention by the application. For a generalized request, the operation associated with the request is performed by the application; therefore, the application must notify MPI when the operation completes. This is done by making a call to `MPI_GREQUEST_COMPLETE`. MPI maintains the “completion” status of generalized requests. Any other request state has to be maintained by the user.

A new generalized request is started with

```
MPI_GREQUEST_START(query_fn, free_fn, cancel_fn, extra_state, request)
```

| | | |
|-----|-------------|---|
| IN | query_fn | callback function invoked when request status is queried (function) |
| IN | free_fn | callback function invoked when request is freed (function) |
| IN | cancel_fn | callback function invoked when request is cancelled (function) |
| IN | extra_state | extra state |
| OUT | request | generalized request (handle) |

```
int MPI_Grequest_start(MPI_Grequest_query_function *query_fn,
                      MPI_Grequest_free_function *free_fn,
                      MPI_Grequest_cancel_function *cancel_fn, void *extra_state,
                      MPI_Request *request)
```

```
MPI_GREQUEST_START(QUERY_FN, FREE_FN, CANCEL_FN, EXTRA_STATE, REQUEST,
                  IERROR)
    INTEGER REQUEST, IERROR
    EXTERNAL QUERY_FN, FREE_FN, CANCEL_FN
    INTEGER (KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

```
static MPI::Grequest
    MPI::Grequest::Start(const MPI::Grequest::Query_function
                        query_fn, const MPI::Grequest::Free_function free_fn,
                        const MPI::Grequest::Cancel_function cancel_fn,
                        void *extra_state)
```

Advice to users. Note that a generalized request belongs, in C++, to the class `MPI::Grequest`, which is a derived class of `MPI::Request`. It is of the same type as regular requests, in C and Fortran. (*End of advice to users.*)

The call starts a generalized request and returns a handle to it in `request`.

The syntax and meaning of the callback functions are listed below. All callback functions are passed the `extra_state` argument that was associated with the request by the starting call `MPI_GREQUEST_START`. This can be used to maintain user-defined state for the request. In C, the query function is

```
typedef int MPI_Grequest_query_function(void *extra_state,
                                       MPI_Status *status);
```

in Fortran

```
SUBROUTINE GREQUEST_QUERY_FUNCTION(EXTRA_STATE, STATUS, IERROR)
  INTEGER STATUS(MPI_STATUS_SIZE), IERROR
  INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

and in C++

```
typedef int MPI::Grequest::Query_function(void* extra_state,
                                          MPI::Status& status);
```

`query_fn` function computes the status that should be returned for the generalized request. The status also includes information about successful/unsuccesful cancellation of the request (result to be returned by `MPI_TEST_CANCELLED`).

`query_fn` callback is invoked by the `MPI_{WAIT|TEST}{ANY|SOME|ALL}` call that completed the generalized request associated with this callback. The callback function is also invoked by calls to `MPI_REQUEST_GET_STATUS`, if the request is complete when the call occurs. In both cases, the callback is passed a reference to the corresponding status variable passed by the user to the MPI call; the status set by the callback function is returned by the MPI call. If the user provided `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE` to the MPI function that causes `query_fn` to be called, then MPI will pass a valid status object to `query_fn`, and this status will be ignored upon return of the callback function. Note that `query_fn` is invoked only after `MPI_GREQUEST_COMPLETE` is called on the request; it may be invoked several times for the same generalized request, e.g., if the user calls `MPI_REQUEST_GET_STATUS` several times for this request. Note also that a call to `MPI_{WAIT|TEST}{SOME|ALL}` may cause multiple invocations of `query_fn` callback functions, one for each generalized request that is completed by the MPI call. The order of these invocations is not specified by MPI.

In C, the free function is

```
typedef int MPI_Grequest_free_function(void *extra_state);
```

and in Fortran

```
SUBROUTINE GREQUEST_FREE_FUNCTION(EXTRA_STATE, IERROR)
  INTEGER IERROR
  INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

and in C++

```
typedef int MPI::Grequest::Free_function(void* extra_state);
```

`free_fn` function is invoked to clean up user-allocated resources when the generalized request is freed.

`free_fn` callback is invoked by the `MPI_{WAIT|TEST}{ANY|SOME|ALL}` call that completed the generalized request associated with this callback. `free_fn` is invoked after the

1 call to `query_fn` for the same request. However, if the MPI call completed multiple generalized
 2 requests, the order in which `free_fn` callback functions are invoked is not specified by MPI.

3 `free_fn` callback is also invoked for generalized requests that are freed by a call to
 4 `MPI_REQUEST_FREE` (no call to `WAIT_{WAIT|TEST}{ANY|SOME|ALL}` will occur for
 5 such a request). In this case, the callback function will be called either in the MPI call
 6 `MPI_REQUEST_FREE(request)`, or in the MPI call `MPI_GREQUEST_COMPLETE(request)`,
 7 whichever happens last. I.e., in this case the actual freeing code is executed as soon as both
 8 calls `MPI_REQUEST_FREE` and `MPI_GREQUEST_COMPLETE` have occurred. The request
 9 is not deallocated until after `free_fn` completes. Note that `free_fn` will be invoked only once
 10 per request by a correct program.

11
 12 *Advice to users.* Calling `MPI_REQUEST_FREE(request)` will cause the request handle
 13 to be set to `MPI_REQUEST_NULL`. This handle to the generalized request is no longer
 14 valid. However, user copies of this handle are valid until after `free_fn` completes since
 15 MPI does not deallocate the object until then. Since `free_fn` is not called until after
 16 `MPI_GREQUEST_COMPLETE`, the user copy of the handle can be used to make this
 17 call. Users should note that MPI will deallocate the object after `free_fn` executes. At
 18 this point, user copies of the request handle no longer point to a valid request. MPI
 19 will not set user copies to `MPI_REQUEST_NULL` in this case, so it is up to the user to
 20 avoid accessing this stale handle. This is a special case where MPI defers deallocating
 21 the object until a later time that is known by the user. (*End of advice to users.*)

22
 23 In C, the cancel function is

```
24 typedef int MPI_Grequest_cancel_function(void *extra_state, int complete);
```

25
 26 in Fortran

```
27 SUBROUTINE GREQUEST_CANCEL_FUNCTION(EXTRA_STATE, COMPLETE, IERROR)
28     INTEGER IERROR
29     INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
30     LOGICAL COMPLETE
```

31
 32 and in C++

```
33 typedef int MPI::Grequest::Cancel_function(void* extra_state,
34     bool complete);
```

35
 36 `cancel_fn` function is invoked to start the cancelation of a generalized request. It is
 37 called by `MPI_CANCEL(request)`. MPI passes to the callback function `complete=true` if
 38 `MPI_GREQUEST_COMPLETE` was already called on the request, and
 39 `complete=false` otherwise.

40 All callback functions return an error code. The code is passed back and dealt with as
 41 appropriate for the error code by the MPI function that invoked the callback function. For
 42 example, if error codes are returned then the error code returned by the callback function
 43 will be returned by the MPI function that invoked the callback function. In the case of
 44 `MPI_{WAIT|TEST}{ANY}` call that invokes both `query_fn` and `free_fn`, the MPI call will
 45 return the error code returned by the last callback, namely `free_fn`. If one or more of the
 46 requests in a call to `MPI_{WAIT|TEST}{SOME|ALL}` failed, then the MPI call will return
 47 `MPI_ERR_IN_STATUS`. In such a case, if the MPI call was passed an array of statuses, then
 48 MPI will return in each of the statuses that correspond to a completed generalized request

the error code returned by the corresponding invocation of its `free_fn` callback function. However, if the MPI function was passed `MPI_STATUSES_IGNORE`, then the individual error codes returned by each callback functions will be lost.

Advice to users. `query_fn` must **not** set the error field of `status` since `query_fn` may be called by `MPI_WAIT` or `MPI_TEST`, in which case the error field of `status` should not change. The MPI library knows the “context” in which `query_fn` is invoked and can decide correctly when to put in the error field of `status` the returned error code. (*End of advice to users.*)

`MPI_GREQUEST_COMPLETE(request)`

`INOUT request` generalized request (handle)

```
int MPI_Grequest_complete(MPI_Request request)
```

```
MPI_GREQUEST_COMPLETE(REQUEST, IERROR)
```

```
INTEGER REQUEST, IERROR
```

```
void MPI::Grequest::Complete()
```

The call informs MPI that the operations represented by the generalized request `request` are complete. (See definitions in Section 2.4.) A call to `MPI_WAIT(request, status)` will return and a call to `MPI_TEST(request, flag, status)` will return `flag=true` only after a call to `MPI_GREQUEST_COMPLETE` has declared that these operations are complete.

MPI imposes no restrictions on the code executed by the callback functions. However, new nonblocking operations should be defined so that the general semantic rules about MPI calls such as `MPI_TEST`, `MPI_REQUEST_FREE`, or `MPI_CANCEL` still hold. For example, all these calls are supposed to be local and nonblocking. Therefore, the callback functions `query_fn`, `free_fn`, or `cancel_fn` should invoke blocking MPI communication calls only if the context is such that these calls are guaranteed to return in finite time. Once `MPI_CANCEL` is invoked, the cancelled operation should complete in finite time, irrespective of the state of other processes (the operation has acquired “local” semantics). It should either succeed, or fail without side-effects. The user should guarantee these same properties for newly defined operations.

Advice to implementors. A call to `MPI_GREQUEST_COMPLETE` may unblock a blocked user process/thread. The MPI library should ensure that the blocked user computation will resume. (*End of advice to implementors.*)

11.2.1 Examples

Example 11.1 This example shows the code for a user-defined reduce operation on an `int` using a binary tree: each non-root node receives two messages, sums them, and sends them up. We assume that no status is returned and that the operation cannot be cancelled.

```
typedef struct {
    MPI_Comm comm;
    int tag;
```

```
1     int root;
2     int valin;
3     int *valout;
4     MPI_Request request;
5     } ARGS;
6
7
8     int myreduce(MPI_Comm comm, int tag, int root,
9                 int valin, int *valout, MPI_Request *request)
10    {
11    ARGS *args;
12    pthread_t thread;
13
14    /* start request */
15    MPI_Grequest_start(query_fn, free_fn, cancel_fn, NULL, request);
16
17    args = (ARGS*)malloc(sizeof(ARGS));
18    args->comm = comm;
19    args->tag = tag;
20    args->root = root;
21    args->valin = valin;
22    args->valout = valout;
23    args->request = *request;
24
25    /* spawn thread to handle request */
26    /* The availability of the pthread_create call is system dependent */
27    pthread_create(&thread, NULL, reduce_thread, args);
28
29    return MPI_SUCCESS;
30    }
31
32
33    /* thread code */
34    void reduce_thread(void *ptr)
35    {
36    int lchild, rchild, parent, lval, rval, val;
37    MPI_Request req[2];
38    ARGS *args;
39
40    args = (ARGS*)ptr;
41
42    /* compute left,right child and parent in tree; set
43       to MPI_PROC_NULL if does not exist */
44    /* code not shown */
45    ...
46
47    MPI_Irecv(&lval, 1, MPI_INT, lchild, args->tag, args->comm, &req[0]);
48    MPI_Irecv(&rval, 1, MPI_INT, rchild, args->tag, args->comm, &req[1]);
```

```
MPI_Waitall(2, req, MPI_STATUSES_IGNORE);      1
val = lval + args->valin + rval;              2
MPI_Send( &val, 1, MPI_INT, parent, args->tag, args->comm );      3
if (parent == MPI_PROC_NULL) *(args->valout) = val;              4
MPI_Grequest_complete((args->request));        5
free(ptr);                                    6
return;                                       7
}                                              8
                                              9

int query_fn(void *extra_state, MPI_Status *status)      10
{                                              11
/* always send just one int */                12
MPI_Status_set_elements(status, MPI_INT, 1);      13
/* can never cancel so always true */          14
MPI_Status_set_cancelled(status, 0);           15
/* choose not to return a value for this */    16
status->MPI_SOURCE = MPI_UNDEFINED;            17
/* tag has not meaning for this generalized request */      18
status->MPI_TAG = MPI_UNDEFINED;               19
/* this generalized request never fails */     20
return MPI_SUCCESS;                            21
}                                              22
                                              23

int free_fn(void *extra_state)                  25
{                                              26
/* this generalized request does not need to do any freeing */  27
/* as a result it never fails here */         28
return MPI_SUCCESS;                            29
}                                              30
                                              31

int cancel_fn(void *extra_state, int complete)      33
{                                              34
/* This generalized request does not support cancelling.      35
   Abort if not already done. If done then treat as if cancel failed. */  36
if (!complete) {                               37
    fprintf(stderr, "Cannot cancel generalized request - aborting program\n");  38
    MPI_Abort(MPI_COMM_WORLD, 99);             39
}                                              40
return MPI_SUCCESS;                            41
}                                              42
                                              43
```

11.3 Associating Information with Status

In MPI-1, requests were associated with point-to-point operations. In MPI-2 there are several different types of requests. These range from new MPI calls for I/O to generalized requests.

1 It is desirable to allow these calls use the same request mechanism. This allows one to wait
 2 or test on different types of requests. However, `MPI_{TEST|WAIT}{ANY|SOME|ALL}`
 3 returns a status with information about the request. With the generalization of requests,
 4 one needs to define what information will be returned in the status object.

5 In MPI-2, each call fills in the appropriate fields in the status object. Any unused fields
 6 will have undefined values. A call to `MPI_{TEST|WAIT}{ANY|SOME|ALL}` can modify
 7 any of the fields in the status object. Specifically, it can modify fields that are undefined.
 8 The fields with meaningful value for a given request are defined in the sections with the
 9 new request.

10 Generalized requests raise additional considerations. Here, the user provides the func-
 11 tions to deal with the request. Unlike other MPI calls, the user needs to provide the
 12 information to be returned in status. The status argument is provided directly to the call-
 13 back function where the status needs to be set. Users can directly set the values in 3 of the
 14 5 status values. The count and cancel fields are opaque. To overcome this, new calls are
 15 provided:

16
 17 `MPI_STATUS_SET_ELEMENTS(status, datatype, count)`

| | | | |
|----|-------|-----------------------|---|
| 19 | INOUT | <code>status</code> | status to associate count with (Status) |
| 20 | IN | <code>datatype</code> | datatype associated with count (handle) |
| 21 | IN | <code>count</code> | number of elements to associate with status (integer) |

22
 23
 24 `int MPI_Status_set_elements(MPI_Status *status, MPI_Datatype datatype,`
 25 `int count)`

26 `MPI_STATUS_SET_ELEMENTS(STATUS, DATATYPE, COUNT, IERROR)`
 27 `INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, COUNT, IERROR`

28
 29 `void MPI::Status::Set_elements(const MPI::Datatype& datatype, int count)`

30
 31 This call modifies the opaque part of `status` so that a call to `MPI_GET_ELEMENTS` will
 32 return `count`. `MPI_GET_COUNT` will return a compatible value.

33 *Rationale.* The number of elements is set instead of the count because the former
 34 can deal with nonintegral number of datatypes. (*End of rationale.*)

35
 36 A subsequent call to `MPI_GET_COUNT(status, datatype, count)` or to
 37 `MPI_GET_ELEMENTS(status, datatype, count)` must use a `datatype` argument that has the
 38 same type signature as the `datatype` argument that was used in the call to
 39 `MPI_STATUS_SET_ELEMENTS`.

40
 41 *Rationale.* This is similar to the restriction that holds when when
 42 `count` is set by a receive operation: in that case, the calls to `MPI_GET_COUNT` and
 43 `MPI_GET_ELEMENTS` must use a `datatype` with the same signature as the `datatype`
 44 used in the receive call. (*End of rationale.*)

```

MPI_STATUS_SET_CANCELLED(status, flag)
    INOUT  status          status to associate cancel flag with (Status)
    IN     flag            if true indicates request was cancelled (logical)

```

```
int MPI_Status_set_cancelled(MPI_Status *status, int flag)
```

```

MPI_STATUS_SET_CANCELLED(STATUS, FLAG, IERROR)
    INTEGER STATUS(MPI_STATUS_SIZE), IERROR
    LOGICAL FLAG

```

```
void MPI::Status::Set_cancelled(bool flag)
```

If flag is set to true then a subsequent call to MPI_TEST_CANCELLED(status, flag) will also return flag = true, otherwise it will return false.

Advice to users. Users are advised not to reuse the status fields for values other than those for which they were intended. Doing so may lead to unexpected results when using the status object. For example, calling MPI_GET_ELEMENTS may cause an error if the value is out of range or it may be impossible to detect such an error. The extra_state argument provided with a generalized request can be used to return information that does not logically belong in status. Furthermore, modifying the values in a status set internally by MPI, e.g., MPI_RECV, may lead to unpredictable results and is strongly discouraged. (*End of advice to users.*)

11.4 Naming Objects

There are many occasions on which it would be useful to allow a user to associate a printable identifier with an MPI communicator, window, or datatype, for instance error reporting, debugging, and profiling. The names attached to opaque objects do not propagate when the object is duplicated or copied by MPI routines. For communicators this can be achieved using the following two functions.

```

MPI_COMM_SET_NAME(comm, comm_name)
    INOUT  comm          communicator whose identifier is to be set (handle)
    IN     comm_name     the character string which is remembered as the name
                        (string)

```

```
int MPI_Comm_set_name(MPI_Comm comm, char *comm_name)
```

```

MPI_COMM_SET_NAME(COMM, COMM_NAME, IERROR)
    INTEGER COMM, IERROR
    CHARACTER*(*) COMM_NAME

```

```
void MPI::Comm::Set_name(const char* comm_name)
```

MPI_COMM_SET_NAME allows a user to associate a name string with a communicator. The character string which is passed to MPI_COMM_SET_NAME will be saved inside the

1 MPI library (so it can be freed by the caller immediately after the call, or allocated on the
2 stack). Leading spaces in `name` are significant but trailing ones are not.

3 `MPI_COMM_SET_NAME` is a local (non-collective) operation, which only affects the
4 name of the communicator as seen in the process which made the `MPI_COMM_SET_NAME`
5 call. There is no requirement that the same (or any) name be assigned to a communicator
6 in every process where it exists.

7
8 *Advice to users.* Since `MPI_COMM_SET_NAME` is provided to help debug code, it
9 is sensible to give the same name to a communicator in all of the processes where it
10 exists, to avoid confusion. (*End of advice to users.*)

11
12 The length of the name which can be stored is limited to the value of
13 `MPI_MAX_OBJECT_NAME` in Fortran and `MPI_MAX_OBJECT_NAME-1` in C and C++ to allow
14 for the null terminator. Attempts to put names longer than this will result in truncation of
15 the name. `MPI_MAX_OBJECT_NAME` must have a value of at least 64.

16
17 *Advice to users.* Under circumstances of store exhaustion an attempt to put a name
18 of any length could fail, therefore the value of `MPI_MAX_OBJECT_NAME` should be
19 viewed only as a strict upper bound on the name length, not a guarantee that setting
20 names of less than this length will always succeed. (*End of advice to users.*)

21
22 *Advice to implementors.* Implementations which pre-allocate a fixed size space for a
23 name should use the length of that allocation as the value of `MPI_MAX_OBJECT_NAME`.
24 Implementations which allocate space for the name from the heap should still define
25 `MPI_MAX_OBJECT_NAME` to be a relatively small value, since the user has to allocate
26 space for a string of up to this size when calling `MPI_COMM_GET_NAME`. (*End of*
27 *advice to implementors.*)

28
29
30 `MPI_COMM_GET_NAME` (`comm`, `comm_name`, `resultlen`)

| | | | |
|----|-----|------------------------|--|
| 31 | IN | <code>comm</code> | communicator whose name is to be returned (handle) |
| 32 | | | |
| 33 | OUT | <code>comm_name</code> | the name previously stored on the communicator, or 34 an empty string if no such name exists (string) |
| 35 | OUT | <code>resultlen</code> | length of returned name (integer) |
| 36 | | | |

37 `int MPI_Comm_get_name(MPI_Comm comm, char *comm_name, int *resultlen)`

38 `MPI_COMM_GET_NAME(COMM, COMM_NAME, RESULTLEN, IERROR)`

39 INTEGER COMM, RESULTLEN, IERROR
40 CHARACTER*(*) COMM_NAME

41
42 `void MPI::Comm::Get_name(char* comm_name, int& resultlen) const`

43
44 `MPI_COMM_GET_NAME` returns the last name which has previously been associated
45 with the given communicator. The name may be set and got from any language. The same
46 name will be returned independent of the language used. `name` should be allocated so that
47 it can hold a resulting string of length `MPI_MAX_OBJECT_NAME` characters.

48 `MPI_COMM_GET_NAME` returns a copy of the set name in `name`. In C, a null character is

additionally stored at `name[resultlen]`. `resultlen` cannot be larger than `MPI_MAX_OBJECT-1`. In Fortran, `name` is padded on the right with blank characters. `resultlen` cannot be larger than `MPI_MAX_OBJECT`.

If the user has not associated a name with a communicator, or an error occurs, `MPI_COMM_GET_NAME` will return an empty string (all spaces in Fortran, "" in C and C++). The three predefined communicators will have predefined names associated with them. Thus, the names of `MPI_COMM_WORLD`, `MPI_COMM_SELF`, and the communicator returned by `MPI_COMM_GET_PARENT` (if not `MPI_COMM_NULL`) will have the default of `MPI_COMM_WORLD`, `MPI_COMM_SELF`, and `MPI_COMM_PARENT`. The fact that the system may have chosen to give a default name to a communicator does not prevent the user from setting a name on the same communicator; doing this removes the old name and assigns the new one.

Rationale. We provide separate functions for setting and getting the name of a communicator, rather than simply providing a predefined attribute key for the following reasons:

- It is not, in general, possible to store a string as an attribute from Fortran.
- It is not easy to set up the delete function for a string attribute unless it is known to have been allocated from the heap.
- To make the attribute key useful additional code to call `strdup` is necessary. If this is not standardized then users have to write it. This is extra unneeded work which we can easily eliminate.
- The Fortran binding is not trivial to write (it will depend on details of the Fortran compilation system), and will not be portable. Therefore it should be in the library rather than in user code.

(End of rationale.)

Advice to users. The above definition means that it is safe simply to print the string returned by `MPI_COMM_GET_NAME`, as it is always a valid string even if there was no name.

Note that associating a name with a communicator has no effect on the semantics of an MPI program, and will (necessarily) increase the store requirement of the program, since the names must be saved. Therefore there is no requirement that users use these functions to associate names with communicators. However debugging and profiling MPI applications may be made easier if names are associated with communicators, since the debugger or profiler should then be able to present information in a less cryptic manner. *(End of advice to users.)*

The following functions are used for setting and getting names of datatypes.

`MPI_Type_set_name` (`type`, `type_name`)

| | | |
|--------------------|------------------------|---|
| <code>INOUT</code> | <code>type</code> | datatype whose identifier is to be set (handle) |
| <code>IN</code> | <code>type_name</code> | the character string which is remembered as the name (string) |

`int MPI_Type_set_name(MPI_Datatype type, char *type_name)`

```

1 MPI_TYPE_SET_NAME(TYPE, TYPE_NAME, IERROR)
2     INTEGER TYPE, IERROR
3     CHARACTER*(*) TYPE_NAME
4
5 void MPI::Datatype::Set_name(const char* type_name)
6
7
8 MPI_TYPE_GET_NAME (type, type_name, resultlen)
9     IN         type           datatype whose name is to be returned (handle)
10    OUT        type_name      the name previously stored on the datatype, or a empty
11                                string if no such name exists (string)
12
13    OUT        resultlen      length of returned name (integer)
14
15 int MPI_Type_get_name(MPI_Datatype type, char *type_name, int *resultlen)
16
17 MPI_TYPE_GET_NAME(TYPE, TYPE_NAME, RESULTLEN, IERROR)
18     INTEGER TYPE, RESULTLEN, IERROR
19     CHARACTER*(*) TYPE_NAME
20
21 void MPI::Datatype::Get_name(char* type_name, int& resultlen) const
22
23     Named predefined datatypes have the default names of the datatype name. For exam-
24     ple, MPI_WCHAR has the default name of MPI_WCHAR.
25     The following functions are used for setting and getting names of windows.
26
27 MPI_WIN_SET_NAME (win, win_name)
28    INOUT     win             window whose identifier is to be set (handle)
29    IN        win_name       the character string which is remembered as the name
30                                (string)
31
32 int MPI_Win_set_name(MPI_Win win, char *win_name)
33
34 MPI_WIN_SET_NAME(WIN, WIN_NAME, IERROR)
35     INTEGER WIN, IERROR
36     CHARACTER*(*) WIN_NAME
37
38 void MPI::Win::Set_name(const char* win_name)
39
40
41 MPI_WIN_GET_NAME (win, win_name, resultlen)
42    IN        win           window whose name is to be returned (handle)
43    OUT       win_name      the name previously stored on the window, or a empty
44                                string if no such name exists (string)
45
46    OUT       resultlen     length of returned name (integer)
47
48 int MPI_Win_get_name(MPI_Win win, char *win_name, int *resultlen)

```



```

MPI_WIN_GET_NAME(WIN, WIN_NAME, RESULTLEN, IERROR)
    INTEGER WIN, RESULTLEN, IERROR
    CHARACTER*(*) WIN_NAME

void MPI::Win::Get_name(char* win_name, int& resultlen) const

```

11.5 Error Classes, Error Codes, and Error Handlers

Users may want to write a layered library on top of an existing MPI implementation, and this library may have its own set of error codes and classes. An example of such a library is an I/O library based on the I/O chapter in MPI-2. For this purpose, functions are needed to:

1. add a new error class to the ones an MPI implementation already knows.
2. associate error codes with this error class, so that `MPI_ERROR_CLASS` works.
3. associate strings with these error codes, so that `MPI_ERROR_STRING` works.
4. invoke the error handler associated with a communicator, window, or object.

Several new functions are provided to do this. They are all local. No functions are provided to free error handlers or error classes: it is not expected that an application will generate them in significant numbers.

```

MPI_ADD_ERROR_CLASS(errorclass)
    OUT      errorclass          value for the new error class (integer)

```

```

int MPI_Add_error_class(int *errorclass)

```

```

MPI_ADD_ERROR_CLASS(ERRORCLASS, IERROR)
    INTEGER ERRORCLASS, IERROR

```

```

int MPI::Add_error_class()

```

Creates a new error class and returns the value for it.

Rationale. To avoid conflicts with existing error codes and classes, the value is set by the implementation and not by the user. (*End of rationale.*)

Advice to implementors. A high quality implementation will return the value for a new `errorclass` in the same deterministic way on all processes. (*End of advice to implementors.*)

Advice to users. Since a call to `MPI_ADD_ERROR_CLASS` is local, the same `errorclass` may not be returned on all processes that make this call. Thus, it is not safe to assume that registering a new error on a set of processes at the same time will yield the same `errorclass` on all of the processes. However, if an implementation returns the new `errorclass` in a deterministic way, and they are always generated in the same

order on the same set of processes (for example, all processes), then the value will be the same. However, even if a deterministic algorithm is used, the value can vary across processes. This can happen, for example, if different but overlapping groups of processes make a series of calls. As a result of these issues, getting the “same” error on multiple processes may not cause the same value of error code to be generated. (*End of advice to users.*)

The value of `MPI_ERR_LASTCODE` is not affected by new user-defined error codes and classes. As in MPI-1, it is a constant value. Instead, a predefined attribute key `MPI_LASTUSEDCLASS` is associated with `MPI_COMM_WORLD`. The attribute value corresponding to this key is the current maximum error class including the user-defined ones. This is a local value and may be different on different processes. The value returned by this key is always greater than or equal to `MPI_ERR_LASTCODE`.

Advice to users. The value returned by the key `MPI_LASTUSEDCLASS` will not change unless the user calls a function to explicitly add an error class/code. In a multi-threaded environment, the user must take extra care in assuming this value has not changed. Note that error codes and error classes are not necessarily dense. A user may not assume that each error class below `MPI_LASTUSEDCLASS` is valid. (*End of advice to users.*)

`MPI_ADD_ERROR_CODE(errorclass, errorcode)`

| | | |
|-----|------------|--|
| IN | errorclass | error class (integer) |
| OUT | errorcode | new error code to associated with errorclass (integer) |

`int MPI_Add_error_code(int errorclass, int *errorcode)`

`MPI_ADD_ERROR_CODE(ERRORCLASS, ERRORCODE, IERROR)`

INTEGER ERRORCLASS, ERRORCODE, IERROR

`int MPI::Add_error_code(int errorclass)`

Creates new error code associated with `errorclass` and returns its value in `errorcode`.

Rationale. To avoid conflicts with existing error codes and classes, the value of the new error code is set by the implementation and not by the user. (*End of rationale.*)

Advice to implementors. A high quality implementation will return the value for a new `errorcode` in the same deterministic way on all processes. (*End of advice to implementors.*)

`MPI_ADD_ERROR_STRING(errorcode, string)`

| | | |
|----|-----------|--|
| IN | errorcode | error code or class (integer) |
| IN | string | text corresponding to errorcode (string) |

`int MPI_Add_error_string(int errorcode, char *string)`

```
MPI_ADD_ERROR_STRING(ERRORCODE, STRING, IERROR)
    INTEGER ERRORCODE, IERROR
    CHARACTER*(*) STRING
```

```
void MPI::Add_error_string(int errorcode, const char* string)
```

Associates an error string with an error code or class. The string must be no more than MPI_MAX_ERROR_STRING characters long. The length of the string is as defined in the calling language. The length of the string does not include the null terminator in C or C++. Trailing blanks will be stripped in Fortran. Calling MPI_ADD_ERROR_STRING for an errorcode that already has a string will replace the old string with the new string. It is erroneous to call MPI_ADD_ERROR_STRING for an error code or class with a value \leq MPI_ERR_LASTCODE.

If MPI_ERROR_STRING is called when no string has been set, it will return a empty string (all spaces in Fortran, "" in C and C++).

Section 7.3.1 on page 253 describes the methods for creating and associating error handlers with communicators, files, and windows.

```
MPI_COMM_CALL_ERRHANDLER (comm, errorcode)
```

```
IN      comm          communicator with error handler (handle)
IN      errorcode     error code (integer)
```

```
int MPI_Comm_call_errhandler(MPI_Comm comm, int errorcode)
```

```
MPI_COMM_CALL_ERRHANDLER(COMM, ERRORCODE, IERROR)
    INTEGER COMM, ERRORCODE, IERROR
```

```
void MPI::Comm::Call_errhandler(int errorcode) const
```

This function invokes the error handler assigned to the communicator with the error code supplied. This function returns MPI_SUCCESS in C and C++ and the same value in IERROR if the error handler was successfully called (assuming the process is not aborted and the error handler returns).

Advice to users. Users should note that the default error handler is MPI_ERRORS_ARE_FATAL. Thus, calling MPI_COMM_CALL_ERRHANDLER will abort the comm processes if the default error handler has not been changed for this communicator or on the parent before the communicator was created. (*End of advice to users.*)

```
MPI_WIN_CALL_ERRHANDLER (win, errorcode)
```

```
IN      win           window with error handler (handle)
IN      errorcode     error code (integer)
```

```
int MPI_Win_call_errhandler(MPI_Win win, int errorcode)
```

1 MPI_WIN_CALL_ERRHANDLER(WIN, ERRORCODE, IERROR)

2 INTEGER WIN, ERRORCODE, IERROR

3
4 void MPI::Win::Call_errhandler(int errorcode) const

5 This function invokes the error handler assigned to the window with the error code
6 supplied. This function returns MPI_SUCCESS in C and C++ and the same value in IERROR
7 if the error handler was successfully called (assuming the process is not aborted and the
8 error handler returns).

9
10 *Advice to users.* As with communicators, the default error handler for windows is
11 MPI_ERRORS_ARE_FATAL. (*End of advice to users.*)

12
13
14 MPI_FILE_CALL_ERRHANDLER (fh, errorcode)

15 IN fh file with error handler (handle)

16 IN errorcode error code (integer)

17
18
19 int MPI_File_call_errhandler(MPI_File fh, int errorcode)

20 MPI_FILE_CALL_ERRHANDLER(FH, ERRORCODE, IERROR)

21 INTEGER FH, ERRORCODE, IERROR

22
23 void MPI::File::Call_errhandler(int errorcode) const

24
25 This function invokes the error handler assigned to the file with the error code supplied.
26 This function returns MPI_SUCCESS in C and C++ and the same value in IERROR if the
27 error handler was successfully called (assuming the process is not aborted and the error
28 handler returns).

29
30 *Advice to users.* Unlike errors on communicators and windows, the default behavior
31 for files is to have MPI_ERRORS_RETURN (*End of advice to users.*)

32
33 *Advice to users.* Users are warned that handlers should not be called recursively
34 with MPI_COMM_CALL_ERRHANDLER, MPI_FILE_CALL_ERRHANDLER, or
35 MPI_WIN_CALL_ERRHANDLER. Doing this can create a situation where an infinite
36 recursion is created. This can occur if MPI_COMM_CALL_ERRHANDLER,
37 MPI_FILE_CALL_ERRHANDLER, or MPI_WIN_CALL_ERRHANDLER is called inside an
38 error handler.

39 Error codes and classes are associated with a process. As a result, they may be used
40 in any error handler. Error handlers should be prepared to deal with any error code
41 it is given. Furthermore, it is good practice to only call an error handler with the
42 appropriate error codes. For example, file errors would normally be sent to the file
43 error handler. (*End of advice to users.*)

44 45 11.6 Decoding a Datatype

46
47 MPI-1 provides datatype objects, which allow users to specify an arbitrary layout of data
48 in memory. The layout information, once put in a datatype, could not be decoded from

the datatype. There are several cases, however, where accessing the layout information in opaque datatype objects would be useful.

The two functions in this section are used together to decode datatypes to recreate the calling sequence used in their initial definition. These can be used to allow a user to determine the type map and type signature of a datatype.

`MPI_TYPE_GET_ENVELOPE(datatype, num_integers, num_addresses, num_datatypes, combiner)`

| | | |
|-----|----------------------------|--|
| IN | <code>datatype</code> | datatype to access (handle) |
| OUT | <code>num_integers</code> | number of input integers used in the call constructing combiner (nonnegative integer) |
| OUT | <code>num_addresses</code> | number of input addresses used in the call constructing combiner (nonnegative integer) |
| OUT | <code>num_datatypes</code> | number of input datatypes used in the call constructing combiner (nonnegative integer) |
| OUT | <code>combiner</code> | combiner (state) |

```
int MPI_Type_get_envelope(MPI_Datatype datatype, int *num_integers,
                          int *num_addresses, int *num_datatypes, int *combiner)
```

```
MPI_TYPE_GET_ENVELOPE(DATATYPE, NUM_INTEGERS, NUM_ADDRESSES, NUM_DATATYPES,
                       COMBINER, IERROR)
INTEGER DATATYPE, NUM_INTEGERS, NUM_ADDRESSES, NUM_DATATYPES, COMBINER,
IERROR
```

```
void MPI::Datatype::Get_envelope(int& num_integers, int& num_addresses,
                                 int& num_datatypes, int& combiner) const
```

For the given datatype, `MPI_TYPE_GET_ENVELOPE` returns information on the number and type of input arguments used in the call that created the datatype. The number-of-arguments values returned can be used to provide sufficiently large arrays in the decoding routine `MPI_TYPE_GET_CONTENTS`. This call and the meaning of the returned values is described below. The combiner reflects the MPI datatype constructor call that was used in creating datatype.

Rationale. By requiring that the combiner reflect the constructor used in the creation of the datatype, the decoded information can be used to effectively recreate the calling sequence used in the original creation. One call is effectively the same as another when the information obtained from `MPI_TYPE_GET_CONTENTS` may be used with either to produce the same outcome. C calls `MPI_Type_hindexed` and `MPI_Type_create_hindexed` are always effectively the same while the Fortran call `MPI_TYPE_HINDEXED` will be different than either of these in some MPI implementations. This is the most useful information and was felt to be reasonable even though it constrains implementations to remember the original constructor sequence even if the internal representation is different.

The decoded information keeps track of datatype duplications. This is important as one needs to distinguish between a predefined datatype and a dup of a predefined

datatype. The former is a constant object that cannot be freed, while the latter is a derived datatype that can be freed. (*End of rationale.*)

The list below has the values that can be returned in combiner on the left and the call associated with them on the right.

| | |
|-------------------------------|--|
| MPI_COMBINER_NAMED | a named predefined datatype |
| MPI_COMBINER_DUP | MPI_TYPE_DUP |
| MPI_COMBINER_CONTIGUOUS | MPI_TYPE_CONTIGUOUS |
| MPI_COMBINER_VECTOR | MPI_TYPE_VECTOR |
| MPI_COMBINER_HVECTOR_INTEGER | MPI_TYPE_HVECTOR from Fortran |
| MPI_COMBINER_HVECTOR | MPI_TYPE_HVECTOR from C or C++ and in some case Fortran or MPI_TYPE_CREATE_HVECTOR |
| MPI_COMBINER_INDEXED | MPI_TYPE_INDEXED |
| MPI_COMBINER_HINDEXED_INTEGER | MPI_TYPE_HINDEXED from Fortran |
| MPI_COMBINER_HINDEXED | MPI_TYPE_HINDEXED from C or C++ and in some case Fortran or MPI_TYPE_CREATE_HINDEXED |
| MPI_COMBINER_INDEXED_BLOCK | MPI_TYPE_CREATE_INDEXED_BLOCK |
| MPI_COMBINER_STRUCT_INTEGER | MPI_TYPE_STRUCT from Fortran |
| MPI_COMBINER_STRUCT | MPI_TYPE_STRUCT from C or C++ and in some case Fortran or MPI_TYPE_CREATE_STRUCT |
| MPI_COMBINER_SUBARRAY | MPI_TYPE_CREATE_SUBARRAY |
| MPI_COMBINER_DARRAY | MPI_TYPE_CREATE_DARRAY |
| MPI_COMBINER_F90_REAL | MPI_TYPE_CREATE_F90_REAL |
| MPI_COMBINER_F90_COMPLEX | MPI_TYPE_CREATE_F90_COMPLEX |
| MPI_COMBINER_F90_INTEGER | MPI_TYPE_CREATE_F90_INTEGER |
| MPI_COMBINER_RESIZED | MPI_TYPE_CREATE_RESIZED |

Table 11.1: combiner values returned from MPI_TYPE_GET_ENVELOPE

If combiner is MPI_COMBINER_NAMED then datatype is a named predefined datatype.

For calls with address arguments, we sometimes need to differentiate whether the call used an integer or an address size argument. For example, there are two combiners for hvector: MPI_COMBINER_HVECTOR_INTEGER and MPI_COMBINER_HVECTOR. The former is used if it was the MPI-1 call from Fortran, and the latter is used if it was the MPI-1 call from C or C++. However, on systems where MPI_ADDRESS_KIND = MPI_INTEGER_KIND (i.e., where integer arguments and address size arguments are the same), the combiner MPI_COMBINER_HVECTOR may be returned for a datatype constructed by a call to MPI_TYPE_HVECTOR from Fortran. Similarly, MPI_COMBINER_HINDEXED may be returned for a datatype constructed by a call to MPI_TYPE_HINDEXED from Fortran, and MPI_COMBINER_STRUCT may be returned for a datatype constructed by a call to MPI_TYPE_STRUCT from Fortran. On such systems, one need not differentiate constructors that take address size arguments from constructors that take integer arguments, since these are the same. The new MPI-2 calls all use address sized arguments.

Rationale. For recreating the original call, it is important to know if address information may have been truncated. The MPI-1 calls from Fortran for a few routines could be subject to truncation in the case where the default `INTEGER` size is smaller than the size of an address. (*End of rationale.*)

The actual arguments used in the creation call for a datatype can be obtained from the call:

```
MPI_TYPE_GET_CONTENTS(datatype, max_integers, max_addresses, max_datatypes, array_of_integers, array_of_addresses, array_of_datatypes)
```

| | | |
|-----|---------------------------------|--|
| IN | <code>datatype</code> | datatype to access (handle) |
| IN | <code>max_integers</code> | number of elements in <code>array_of_integers</code> (nonnegative integer) |
| IN | <code>max_addresses</code> | number of elements in <code>array_of_addresses</code> (nonnegative integer) |
| IN | <code>max_datatypes</code> | number of elements in <code>array_of_datatypes</code> (nonnegative integer) |
| OUT | <code>array_of_integers</code> | contains integer arguments used in constructing datatype (array of integers) |
| OUT | <code>array_of_addresses</code> | contains address arguments used in constructing datatype (array of integers) |
| OUT | <code>array_of_datatypes</code> | contains datatype arguments used in constructing datatype (array of handles) |

```
int MPI_Type_get_contents(MPI_Datatype datatype, int max_integers,
    int max_addresses, int max_datatypes, int array_of_integers[],
    MPI_Aint array_of_addresses[],
    MPI_Datatype array_of_datatypes[])
```

```
MPI_TYPE_GET_CONTENTS(DATATYPE, MAX_INTEGERS, MAX_ADDRESSES, MAX_DATATYPES,
    ARRAY_OF_INTEGERS, ARRAY_OF_ADDRESSES, ARRAY_OF_DATATYPES,
    IERROR)
    INTEGER DATATYPE, MAX_INTEGERS, MAX_ADDRESSES, MAX_DATATYPES,
    ARRAY_OF_INTEGERS(*), ARRAY_OF_DATATYPES(*), IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_ADDRESSES(*)
```

```
void MPI::Datatype::Get_contents(int max_integers, int max_addresses,
    int max_datatypes, int array_of_integers[],
    MPI::Aint array_of_addresses[],
    MPI::Datatype array_of_datatypes[]) const
```

`datatype` must be a predefined unnamed or a derived datatype; the call is erroneous if `datatype` is a predefined named datatype.

The values given for `max_integers`, `max_addresses`, and `max_datatypes` must be at least as large as the value returned in `num_integers`, `num_addresses`, and `num_datatypes`, respectively, in the call `MPI_TYPE_GET_ENVELOPE` for the same `datatype` argument.

1 *Rationale.* The arguments `max_integers`, `max_addresses`, and `max_datatypes` allow for
2 error checking in the call. This is analogous to the topology calls in MPI-1. (*End of*
3 *rationale.*)

4
5 The datatypes returned in `array_of_datatypes` are handles to datatype objects that are
6 equivalent to the datatypes used in the original construction call. If these were derived
7 datatypes, then the returned datatypes are new datatype objects, and the user is responsible
8 for freeing these datatypes with `MPI_TYPE_FREE`. If these were predefined datatypes, then
9 the returned datatype is equal to that (constant) predefined datatype and cannot be freed.

10 The committed state of returned derived datatypes is undefined, i.e., the datatypes may
11 or may not be committed. Furthermore, the content of attributes of returned datatypes is
12 undefined.

13 Note that `MPI_TYPE_GET_CONTENTS` can be invoked with a `datatype` argument that
14 was constructed using `MPI_TYPE_CREATE_F90_REAL`, `MPI_TYPE_CREATE_F90_INTEGER`,
15 or `MPI_TYPE_CREATE_F90_COMPLEX` (an unnamed predefined datatype). In such a case,
16 an empty `array_of_datatypes` is returned.

17
18 *Rationale.* The definition of datatype equivalence implies that equivalent predefined
19 datatypes are equal. By requiring the same handle for named predefined datatypes,
20 it is possible to use the `==` or `.EQ.` comparison operator to determine the datatype
21 involved. (*End of rationale.*)

22
23 *Advice to implementors.* The datatypes returned in `array_of_datatypes` must appear
24 to the user as if each is an equivalent copy of the datatype used in the type constructor
25 call. Whether this is done by creating a new datatype or via another mechanism such
26 as a reference count mechanism is up to the implementation as long as the semantics
27 are preserved. (*End of advice to implementors.*)

28
29 *Rationale.* The committed state and attributes of the returned datatype is delib-
30 erately left vague. The datatype used in the original construction may have been
31 modified since its use in the constructor call. Attributes can be added, removed, or
32 modified as well as having the datatype committed. The semantics given allow for
33 a reference count implementation without having to track these changes. (*End of*
34 *rationale.*)

35 In the MPI-1 datatype constructor calls, the address arguments in Fortran are of type
36 `INTEGER`. In the new MPI-2 calls, the address arguments are of type
37 `INTEGER(KIND=MPI_ADDRESS_KIND)`. The call `MPI_TYPE_GET_CONTENTS` returns all ad-
38 dresses in an argument of type `INTEGER(KIND=MPI_ADDRESS_KIND)`. This is true even if the
39 old MPI-1 calls were used. Thus, the location of values returned can be thought of as being
40 returned by the C bindings. It can also be determined by examining the new MPI-2 calls
41 for datatype constructors for the deprecated MPI-1 calls that involve addresses.

42
43 *Rationale.* By having all address arguments returned in the
44 `array_of_addresses` argument, the result from a C and Fortran decoding of a `datatype`
45 gives the result in the same argument. It is assumed that an integer of type
46 `INTEGER(KIND=MPI_ADDRESS_KIND)` will be at least as large as the `INTEGER` argument
47 used in datatype construction with the old MPI-1 calls so no loss of information will
48 occur. (*End of rationale.*)

The following defines what values are placed in each entry of the returned arrays depending on the datatype constructor used for `datatype`. It also specifies the size of the arrays needed which is the values returned by `MPI_TYPE_GET_ENVELOPE`. In Fortran, the following calls were made:

```

PARAMETER (LARGE = 1000)
INTEGER TYPE, NI, NA, ND, COMBINER, I(LARGE), D(LARGE), IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) A(LARGE)
! CONSTRUCT DATATYPE TYPE (NOT SHOWN)
CALL MPI_TYPE_GET_ENVELOPE(TYPE, NI, NA, ND, COMBINER, IERROR)
IF ((NI .GT. LARGE) .OR. (NA .GT. LARGE) .OR. (ND .GT. LARGE)) THEN
    WRITE (*, *) "NI, NA, OR ND = ", NI, NA, ND, &
    " RETURNED BY MPI_TYPE_GET_ENVELOPE IS LARGER THAN LARGE = ", LARGE
    CALL MPI_ABORT(MPI_COMM_WORLD, 99)
ENDIF
CALL MPI_TYPE_GET_CONTENTS(TYPE, NI, NA, ND, I, A, D, IERROR)

```

or in C the analogous calls of:

```

#define LARGE 1000
int ni, na, nd, combiner, i[LARGE];
MPI_Aint a[LARGE];
MPI_Datatype type, d[LARGE];
/* construct datatype type (not shown) */
MPI_Type_get_envelope(type, &ni, &na, &nd, &combiner);
if ((ni > LARGE) || (na > LARGE) || (nd > LARGE)) {
    fprintf(stderr, "ni, na, or nd = %d %d %d returned by ", ni, na, nd);
    fprintf(stderr, "MPI_Type_get_envelope is larger than LARGE = %d\n",
            LARGE);
    MPI_Abort(MPI_COMM_WORLD, 99);
};
MPI_Type_get_contents(type, ni, na, nd, i, a, d);

```

The C++ code is in analogy to the C code above with the same values returned.

In the descriptions that follow, the lower case name of arguments is used.

If combiner is `MPI_COMBINER_NAMED` then it is erroneous to call `MPI_TYPE_GET_CONTENTS`.

If combiner is `MPI_COMBINER_DUP` then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| oldtype | d[0] | D(1) |

and `ni = 0`, `na = 0`, `nd = 1`.

If combiner is `MPI_COMBINER_CONTIGUOUS` then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| count | i[0] | I(1) |
| oldtype | d[0] | D(1) |

and `ni = 1`, `na = 0`, `nd = 1`.

If combiner is `MPI_COMBINER_VECTOR` then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| count | i[0] | I(1) |
| blocklength | i[1] | I(2) |
| stride | i[2] | I(3) |
| oldtype | d[0] | D(1) |

and ni = 3, na = 0, nd = 1.

If combiner is MPI_COMBINER_HVECTOR_INTEGER or MPI_COMBINER_HVECTOR then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| count | i[0] | I(1) |
| blocklength | i[1] | I(2) |
| stride | a[0] | A(1) |
| oldtype | d[0] | D(1) |

and ni = 2, na = 1, nd = 1.

If combiner is MPI_COMBINER_INDEXED then

| Constructor argument | C & C++ location | Fortran location |
|------------------------|------------------------|--------------------------|
| count | i[0] | I(1) |
| array_of_blocklengths | i[1] to i[i[0]] | I(2) to I(I(1)+1) |
| array_of_displacements | i[i[0]+1] to i[2*i[0]] | I(I(1)+2) to I(2*I(1)+1) |
| oldtype | d[0] | D(1) |

and ni = 2*count+1, na = 0, nd = 1.

If combiner is MPI_COMBINER_HINDEXED_INTEGER or MPI_COMBINER_HINDEXED then

| Constructor argument | C & C++ location | Fortran location |
|------------------------|-------------------|-------------------|
| count | i[0] | I(1) |
| array_of_blocklengths | i[1] to i[i[0]] | I(2) to I(I(1)+1) |
| array_of_displacements | a[0] to a[i[0]-1] | A(1) to A(I(1)) |
| oldtype | d[0] | D(1) |

and ni = count+1, na = count, nd = 1.

If combiner is MPI_COMBINER_INDEXED_BLOCK then

| Constructor argument | C & C++ location | Fortran location |
|------------------------|-------------------|-------------------|
| count | i[0] | I(1) |
| blocklength | i[1] | I(2) |
| array_of_displacements | i[2] to i[i[0]+1] | I(3) to I(I(1)+2) |
| oldtype | d[0] | D(1) |

and ni = count+2, na = 0, nd = 1.

If combiner is MPI_COMBINER_STRUCT_INTEGER or MPI_COMBINER_STRUCT then

| Constructor argument | C & C++ location | Fortran location |
|------------------------|-------------------|-------------------|
| count | i[0] | I(1) |
| array_of_blocklengths | i[1] to i[i[0]] | I(2) to I(I(1)+1) |
| array_of_displacements | a[0] to a[i[0]-1] | A(1) to A(I(1)) |
| array_of_types | d[0] to d[i[0]-1] | D(1) to D(I(1)) |

and ni = count+1, na = count, nd = count.

If combiner is MPI_COMBINER_SUBARRAY then

46
47
48

| Constructor argument | C & C++ location | Fortran location |
|----------------------|--------------------------|----------------------------|
| ndims | i[0] | I(1) |
| array_of_sizes | i[1] to i[i[0]] | I(2) to I(I(1)+1) |
| array_of_subsizes | i[i[0]+1] to i[2*i[0]] | I(I(1)+2) to I(2*I(1)+1) |
| array_of_starts | i[2*i[0]+1] to i[3*i[0]] | I(2*I(1)+2) to I(3*I(1)+1) |
| order | i[3*i[0]+1] | I(3*I(1)+2) |
| oldtype | d[0] | D(1) |

and ni = 3*ndims+2, na = 0, nd = 1.

If combiner is MPI_COMBINER_DARRAY then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|----------------------------|----------------------------|
| size | i[0] | I(1) |
| rank | i[1] | I(2) |
| ndims | i[2] | I(3) |
| array_of_gsizes | i[3] to i[i[2]+2] | I(4) to I(I(3)+3) |
| array_of_distribs | i[i[2]+3] to i[2*i[2]+2] | I(I(3)+4) to I(2*I(3)+3) |
| array_of_dargs | i[2*i[2]+3] to i[3*i[2]+2] | I(2*I(3)+4) to I(3*I(3)+3) |
| array_of_psizes | i[3*i[2]+3] to i[4*i[2]+2] | I(3*I(3)+4) to I(4*I(3)+3) |
| order | i[4*i[2]+3] | I(4*I(3)+4) |
| oldtype | d[0] | D(1) |

and ni = 4*ndims+4, na = 0, nd = 1.

If combiner is MPI_COMBINER_F90_REAL then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| p | i[0] | I(1) |
| r | i[1] | I(2) |

and ni = 2, na = 0, nd = 0.

If combiner is MPI_COMBINER_F90_COMPLEX then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| p | i[0] | I(1) |
| r | i[1] | I(2) |

and ni = 2, na = 0, nd = 0.

If combiner is MPI_COMBINER_F90_INTEGER then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| r | i[0] | I(1) |

and ni = 1, na = 0, nd = 0.

If combiner is MPI_COMBINER_RESIZED then

| Constructor argument | C & C++ location | Fortran location |
|----------------------|------------------|------------------|
| lb | a[0] | A(1) |
| extent | a[1] | A(2) |
| oldtype | d[0] | D(1) |

and ni = 0, na = 2, nd = 1.

Example 11.2 This example shows how a datatype can be decoded. The routine `printdatatype` prints out the elements of the datatype. Note the use of `MPI_Type_free` for datatypes that are not predefined.

```

1  /*
2     Example of decoding a datatype.
3
4     Returns 0 if the datatype is predefined, 1 otherwise
5  */
6  #include <stdio.h>
7  #include <stdlib.h>
8  #include "mpi.h"
9  int printdatatype( MPI_Datatype datatype )
10 {
11     int *array_of_ints;
12     MPI_Aint *array_of_adds;
13     MPI_Datatype *array_of_dtypes;
14     int num_ints, num_adds, num_dtypes, combiner;
15     int i;
16
17     MPI_Type_get_envelope( datatype,
18                           &num_ints, &num_adds, &num_dtypes, &combiner );
19     switch (combiner) {
20     case MPI_COMBINER_NAMED:
21         printf( "Datatype is named:" );
22         /* To print the specific type, we can match against the
23            predefined forms. We can NOT use a switch statement here
24            We could also use MPI_TYPE_GET_NAME if we preferred to use
25            names that the user may have changed.
26         */
27         if (datatype == MPI_INT) printf( "MPI_INT\n" );
28         else if (datatype == MPI_DOUBLE) printf( "MPI_DOUBLE\n" );
29         ... else test for other types ...
30         return 0;
31         break;
32     case MPI_COMBINER_STRUCT:
33     case MPI_COMBINER_STRUCT_INTEGER:
34         printf( "Datatype is struct containing" );
35         array_of_ints = (int *)malloc( num_ints * sizeof(int) );
36         array_of_adds =
37             (MPI_Aint *) malloc( num_adds * sizeof(MPI_Aint) );
38         array_of_dtypes = (MPI_Datatype *)
39             malloc( num_dtypes * sizeof(MPI_Datatype) );
40         MPI_Type_get_contents( datatype, num_ints, num_adds, num_dtypes,
41                               array_of_ints, array_of_adds, array_of_dtypes );
42         printf( " %d datatypes:\n", array_of_ints[0] );
43         for (i=0; i<array_of_ints[0]; i++) {
44             printf( "blocklength %d, displacement %ld, type:\n",
45                   array_of_ints[i+1], array_of_adds[i] );
46             if (printdatatype( array_of_dtypes[i] )) {
47                 /* Note that we free the type ONLY if it
48                    is not predefined */

```

```

        MPI_Type_free( &array_of_dtypes[i] );
    }
}
free( array_of_ints );
free( array_of_adds );
free( array_of_dtypes );
break;
... other combiner values ...
default:
    printf( "Unrecognized combiner type\n" );
}
return 1;
}

```

11.7 MPI and Threads

This section specifies the interaction between MPI calls and threads. The section lists minimal requirements for **thread compliant** MPI implementations and defines functions that can be used for initializing the thread environment. MPI may be implemented in environments where threads are not supported or perform poorly. Therefore, it is not required that all MPI implementations fulfill all the requirements specified in this section.

This section generally assumes a thread package similar to POSIX threads [31], but the syntax and semantics of thread calls are not specified here — these are beyond the scope of this document.

11.7.1 General

In a thread-compliant implementation, an MPI process is a process that may be multi-threaded. Each thread can issue MPI calls; however, threads are not separately addressable: a rank in a send or receive call identifies a process, not a thread. A message sent to a process can be received by any thread in this process.

Rationale. This model corresponds to the POSIX model of interprocess communication: the fact that a process is multi-threaded, rather than single-threaded, does not affect the external interface of this process. MPI implementations where MPI ‘processes’ are POSIX threads inside a single POSIX process are not thread-compliant by this definition (indeed, their “processes” are single-threaded). (*End of rationale.*)

Advice to users. It is the user’s responsibility to prevent races when threads within the same application post conflicting communication calls. The user can make sure that two threads in the same process will not issue conflicting communication calls by using distinct communicators at each thread. (*End of advice to users.*)

The two main requirements for a thread-compliant implementation are listed below.

1. All MPI calls are *thread-safe*. I.e., two concurrently running threads may make MPI calls and the outcome will be as if the calls executed in some order, even if their execution is interleaved.

2. Blocking MPI calls will block the calling thread only, allowing another thread to execute, if available. The calling thread will be blocked until the event on which it is waiting occurs. Once the blocked communication is enabled and can proceed, then the call will complete and the thread will be marked runnable, within a finite time. A blocked thread will not prevent progress of other runnable threads on the same process, and will not prevent them from executing MPI calls.

Example 11.3 Process 0 consists of two threads. The first thread executes a blocking send call `MPI_Send(buff1, count, type, 0, 0, comm)`, whereas the second thread executes a blocking receive call `MPI_Recv(buff2, count, type, 0, 0, comm, &status)`. I.e., the first thread sends a message that is received by the second thread. This communication should always succeed. According to the first requirement, the execution will correspond to some interleaving of the two calls. According to the second requirement, a call can only block the calling thread and cannot prevent progress of the other thread. If the send call went ahead of the receive call, then the sending thread may block, but this will not prevent the receiving thread from executing. Thus, the receive call will occur. Once both calls occur, the communication is enabled and both calls will complete. On the other hand, a single-threaded process that posts a send, followed by a matching receive, may deadlock. The progress requirement for multithreaded implementations is stronger, as a blocked call cannot prevent progress in other threads.

Advice to implementors. MPI calls can be made thread-safe by executing only one at a time, e.g., by protecting MPI code with one process-global lock. However, blocked operations cannot hold the lock, as this would prevent progress of other threads in the process. The lock is held only for the duration of an atomic, locally-completing suboperation such as posting a send or completing a send, and is released in between. Finer locks can provide more concurrency, at the expense of higher locking overheads. Concurrency can also be achieved by having some of the MPI protocol executed by separate server threads. (*End of advice to implementors.*)

11.7.2 Clarifications

Initialization and Completion The call to `MPI_FINALIZE` should occur on the same thread that initialized MPI. We call this thread the **main thread**. The call should occur only after all the process threads have completed their MPI calls, and have no pending communications or I/O operations.

Rationale. This constraint simplifies implementation. (*End of rationale.*)

Multiple threads completing the same request. A program where two threads block, waiting on the same request, is erroneous. Similarly, the same request cannot appear in the array of requests of two concurrent `MPI_WAIT{ANY|SOME|ALL}` calls. In MPI, a request can only be completed once. Any combination of wait or test which violates this rule is erroneous.

Rationale. This is consistent with the view that a multithreaded execution corresponds to an interleaving of the MPI calls. In a single threaded implementation, once a wait is posted on a request the request handle will be nullified before it is possible to post a second wait on the same handle. With threads, an `MPI_WAIT{ANY|SOME|ALL}` may be blocked without having nullified its request(s)

so it becomes the user's responsibility to avoid using the same request in an `MPI_WAIT` on another thread. This constraint also simplifies implementation, as only one thread will be blocked on any communication or I/O event. (*End of rationale.*)

Probe A receive call that uses source and tag values returned by a preceding call to `MPI_PROBE` or `MPI_IPROBE` will receive the message matched by the probe call only if there was no other matching receive after the probe and before that receive. In a multithreaded environment, it is up to the user to enforce this condition using suitable mutual exclusion logic. This can be enforced by making sure that each communicator is used by only one thread on each process.

Collective calls Matching of collective calls on a communicator, window, or file handle is done according to the order in which the calls are issued at each process. If concurrent threads issue such calls on the same communicator, window or file handle, it is up to the user to make sure the calls are correctly ordered, using interthread synchronization.

Advice to users. With three concurrent threads in each MPI process of a communicator `comm`, it is allowed that thread A in each MPI process calls a collective operation on `comm`, thread B calls a file operation on an existing filehandle that was formerly opened on `comm`, and thread C invokes one-sided operations on an existing window handle that was also formerly created on `comm`. (*End of advice to users.*)

Rationale. As already specified in `MPI_FILE_OPEN` and `MPI_WIN_CREATE`, a file handle and a window handle inherit only the group of processes of the underlying communicator, but not the communicator itself. Accesses to communicators, window handles and file handles cannot affect one another. (*End of rationale.*)

Advice to implementors. Advice to implementors. If the implementation of file or window operations internally uses MPI communication then a duplicated communicator may be cached on the file or window object. (*End of advice to implementors.*)

Exception handlers An exception handler does not necessarily execute in the context of the thread that made the exception-raising MPI call; the exception handler may be executed by a thread that is distinct from the thread that will return the error code.

Rationale. The MPI implementation may be multithreaded, so that part of the communication protocol may execute on a thread that is distinct from the thread that made the MPI call. The design allows the exception handler to be executed on the thread where the exception occurred. (*End of rationale.*)

Interaction with signals and cancellations The outcome is undefined if a thread that executes an MPI call is cancelled (by another thread), or if a thread catches a signal while executing an MPI call. However, a thread of an MPI process may terminate, and may catch signals or be cancelled by another thread when not executing MPI calls.

Rationale. Few C library functions are signal safe, and many have cancellation points — points where the thread executing them may be cancelled. The above restriction simplifies implementation (no need for the MPI library to be “async-cancel-safe” or “async-signal-safe.”) (*End of rationale.*)

Advice to users. Users can catch signals in separate, non-MPI threads (e.g., by masking signals on MPI calling threads, and unmasking them in one or more non-MPI threads). A good programming practice is to have a distinct thread blocked in a call to `sigwait` for each user expected signal that may occur. Users must not catch signals used by the MPI implementation; as each MPI implementation is required to document the signals used internally, users can avoid these signals. (*End of advice to users.*)

Advice to implementors. The MPI library should not invoke library calls that are not thread safe, if multiple threads execute. (*End of advice to implementors.*)

11.7.3 Initialization

The following function may be used to initialize MPI, and initialize the MPI thread environment, instead of `MPI_INIT`.

`MPI_INIT_THREAD(required, provided)`

| | | |
|-----|----------|--|
| IN | required | desired level of thread support (integer) |
| OUT | provided | provided level of thread support (integer) |

```
int MPI_Init_thread(int *argc, char ((*argv) []), int required,
                   int *provided)
```

```
MPI_INIT_THREAD(REQUIRED, PROVIDED, IERROR)
INTEGER REQUIRED, PROVIDED, IERROR
```

```
int MPI::Init_thread(int& argc, char**& argv, int required)
```

```
int MPI::Init_thread(int required)
```

Advice to users. In C and C++, the passing of `argc` and `argv` is optional. In C, this is accomplished by passing the appropriate null pointer. In C++, this is accomplished with two separate bindings to cover these two cases. This is as with `MPI_INIT` as discussed in Section 7.6. (*End of advice to users.*)

This call initializes MPI in the same way that a call to `MPI_INIT` would. In addition, it initializes the thread environment. The argument `required` is used to specify the desired level of thread support. The possible values are listed in increasing order of thread support.

`MPI_THREAD_SINGLE` Only one thread will execute.

`MPI_THREAD_FUNNELED` The process may be multi-threaded, but **the application must ensure that only the main thread makes MPI calls (for the definition of main thread, see `MPI_IS_THREAD_MAIN` on page 372).**

`MPI_THREAD_SERIALIZED` The process may be multi-threaded, and multiple threads may make MPI calls, but only one at a time: MPI calls are not made concurrently from two distinct threads (all MPI calls are “serialized”).

`MPI_THREAD_MULTIPLE` Multiple threads may call MPI, with no restrictions.

These values are monotonic; i.e., `MPI_THREAD_SINGLE` < `MPI_THREAD_FUNNELED` < `MPI_THREAD_SERIALIZED` < `MPI_THREAD_MULTIPLE`.

Different processes in `MPI_COMM_WORLD` may require different levels of thread support.

The call returns in `provided` information about the actual level of thread support that will be provided by MPI. It can be one of the four values listed above.

The level(s) of thread support that can be provided by `MPI_INIT_THREAD` will depend on the implementation, and may depend on information provided by the user before the program started to execute (e.g., with arguments to `mpiexec`). If possible, the call will return `provided = required`. Failing this, the call will return the least supported level such that `provided > required` (thus providing a stronger level of support than required by the user). Finally, if the user requirement cannot be satisfied, then the call will return in `provided` the highest supported level.

A **thread compliant** MPI implementation will be able to return `provided = MPI_THREAD_MULTIPLE`. Such an implementation may always return `provided = MPI_THREAD_MULTIPLE`, irrespective of the value of `required`. At the other extreme, an MPI library that is not thread compliant may always return `provided = MPI_THREAD_SINGLE`, irrespective of the value of `required`.

A call to `MPI_INIT` has the same effect as a call to `MPI_INIT_THREAD` with a `required = MPI_THREAD_SINGLE`.

Vendors may provide (implementation dependent) means to specify the level(s) of thread support available when the MPI program is started, e.g., with arguments to `mpiexec`. This will affect the outcome of calls to `MPI_INIT` and `MPI_INIT_THREAD`. Suppose, for example, that an MPI program has been started so that only `MPI_THREAD_MULTIPLE` is available. Then `MPI_INIT_THREAD` will return `provided = MPI_THREAD_MULTIPLE`, irrespective of the value of `required`; a call to `MPI_INIT` will also initialize the MPI thread support level to `MPI_THREAD_MULTIPLE`. Suppose, on the other hand, that an MPI program has been started so that all four levels of thread support are available. Then, a call to `MPI_INIT_THREAD` will return `provided = required`; on the other hand, a call to `MPI_INIT` will initialize the MPI thread support level to `MPI_THREAD_SINGLE`.

Rationale. Various optimizations are possible when MPI code is executed single-threaded, or is executed on multiple threads, but not concurrently: mutual exclusion code may be omitted. Furthermore, if only one thread executes, then the MPI library can use library functions that are not thread safe, without risking conflicts with user threads. Also, the model of one communication thread, multiple computation threads fits well many applications. E.g., if the process code is a sequential Fortran/C/C++ program with MPI calls that has been parallelized by a compiler for execution on an SMP node, in a cluster of SMPs, then the process computation is multi-threaded, but MPI calls will likely execute on a single thread.

The design accommodates a static specification of the thread support level, for environments that require static binding of libraries, and for compatibility for current multi-threaded MPI codes. (*End of rationale.*)

Advice to implementors. If `provided` is not `MPI_THREAD_SINGLE` then the MPI library should not invoke C/C++/Fortran library calls that are not thread safe, e.g., in an environment where `malloc` is not thread safe, then `malloc` should not be used by the MPI library.

Some implementors may want to use different MPI libraries for different levels of thread support. They can do so using dynamic linking and selecting which library will be linked when `MPI_INIT_THREAD` is invoked. If this is not possible, then optimizations for lower levels of thread support will occur only when the level of thread support required is specified at link time. (*End of advice to implementors.*)

The following function can be used to query the current level of thread support.

```
MPI_QUERY_THREAD(provided)
```

```
OUT    provided                provided level of thread support (integer)
```

```
int MPI_Query_thread(int *provided)
```

```
MPI_QUERY_THREAD(PROVIDED, IERROR)
```

```
INTEGER PROVIDED, IERROR
```

```
int MPI::Query_thread()
```

The call returns in `provided` the current level of thread support. This will be the value returned in `provided` by `MPI_INIT_THREAD`, if MPI was initialized by a call to `MPI_INIT_THREAD()`.

```
MPI_IS_THREAD_MAIN(flag)
```

```
OUT    flag                    true if calling thread is main thread, false otherwise
                                (logical)
```

```
int MPI_Is_thread_main(int *flag)
```

```
MPI_IS_THREAD_MAIN(FLAG, IERROR)
```

```
LOGICAL FLAG
```

```
INTEGER IERROR
```

```
bool MPI::Is_thread_main()
```

This function can be called by a thread to find out whether it is the main thread (the thread that called `MPI_INIT` or `MPI_INIT_THREAD`).

All routines listed in this section must be supported by all MPI implementations.

Rationale. MPI libraries are required to provide these calls even if they do not support threads, so that portable code that contains invocations to these functions be able to link correctly. `MPI_INIT` continues to be supported so as to provide compatibility with current MPI codes. (*End of rationale.*)

Advice to users. It is possible to spawn threads before MPI is initialized, but no MPI call other than `MPI_INITIALIZED` should be executed by these threads, until `MPI_INIT_THREAD` is invoked by one thread (which, thereby, becomes the main thread). In particular, it is possible to enter the MPI execution with a multi-threaded process.

The level of thread support provided is a global property of the MPI process that can be specified only once, when MPI is initialized on that process (or before). Portable third party libraries have to be written so as to accommodate any provided level of thread support. Otherwise, their usage will be restricted to specific level(s) of thread support. If such a library can run only with specific level(s) of thread support, e.g., only with MPI_THREAD_MULTIPLE, then MPI_QUERY_THREAD can be used to check whether the user initialized MPI to the correct level of thread support and, if not, raise an exception. (*End of advice to users.*)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 12

I/O

12.1 Introduction

POSIX provides a model of a widely portable file system, but the portability and optimization needed for parallel I/O cannot be achieved with the POSIX interface.

The significant optimizations required for efficiency (e.g., grouping [37], collective buffering [6, 13, 38, 41, 49], and disk-directed I/O [33]) can only be implemented if the parallel I/O system provides a high-level interface supporting partitioning of file data among processes and a collective interface supporting complete transfers of global data structures between process memories and files. In addition, further efficiencies can be gained via support for asynchronous I/O, strided accesses, and control over physical file layout on storage devices (disks). The I/O environment described in this chapter provides these facilities.

Instead of defining I/O access modes to express the common patterns for accessing a shared file (broadcast, reduction, scatter, gather), we chose another approach in which data partitioning is expressed using derived datatypes. Compared to a limited set of predefined access patterns, this approach has the advantage of added flexibility and expressiveness.

12.1.1 Definitions

file An MPI file is an ordered collection of typed data items. MPI supports random or sequential access to any integral set of these items. A file is opened collectively by a group of processes. All collective I/O calls on a file are collective over this group.

displacement A file *displacement* is an absolute byte position relative to the beginning of a file. The displacement defines the location where a *view* begins. Note that a “file displacement” is distinct from a “typemap displacement.”

etype An *etype* (*elementary* datatype) is the unit of data access and positioning. It can be any MPI predefined or derived datatype. Derived etypes can be constructed using any of the MPI datatype constructor routines, provided all resulting typemap displacements are nonnegative and monotonically nondecreasing. Data access is performed in etype units, reading or writing whole data items of type etype. Offsets are expressed as a count of etypes; file pointers point to the beginning of etypes. Depending on context, the term “etype” is used to describe one of three aspects of an elementary datatype: a particular MPI type, a data item of that type, or the extent of that type.

filetype A *filetype* is the basis for partitioning a file among processes and defines a template for accessing the file. A filetype is either a single etype or a derived MPI datatype constructed from multiple instances of the same etype. In addition, the extent of any hole in the filetype must be a multiple of the etype’s extent. The displacements in the typemap of the filetype are not required to be distinct, but they must be nonnegative and monotonically nondecreasing.

view A *view* defines the current set of data visible and accessible from an open file as an ordered set of etypes. Each process has its own view of the file, defined by three quantities: a displacement, an etype, and a filetype. The pattern described by a filetype is repeated, beginning at the displacement, to define the view. The pattern of repetition is defined to be the same pattern that `MPI_TYPE_CONTIGUOUS` would produce if it were passed the filetype and an arbitrarily large count. Figure 12.1 shows how the tiling works; note that the filetype in this example must have explicit lower and upper bounds set in order for the initial and final holes to be repeated in the view. Views can be changed by the user during program execution. The default view is a linear byte stream (displacement is zero, etype and filetype equal to `MPI_BYTE`).

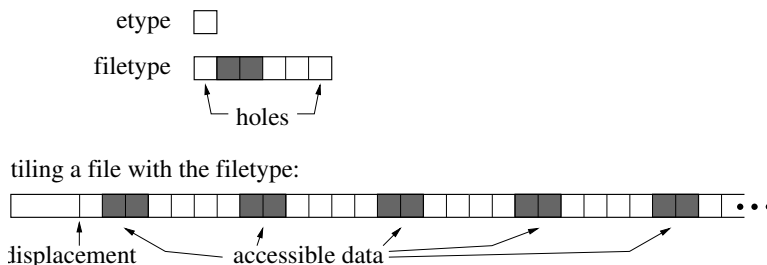


Figure 12.1: Etypes and filetypes

A group of processes can use complementary views to achieve a global data distribution such as a scatter/gather pattern (see Figure 12.2).

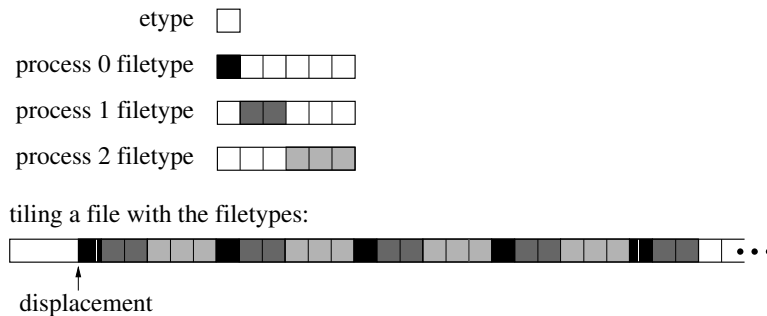


Figure 12.2: Partitioning a file among parallel processes

offset An *offset* is a position in the file relative to the current view, expressed as a count of etypes. Holes in the view’s filetype are skipped when calculating this position. Offset 0 is the location of the first etype visible in the view (after skipping the displacement and any initial holes in the view). For example, an offset of 2 for process 1 in Figure 12.2 is the position of the 8th etype in the file after the displacement. An “explicit offset” is an offset that is used as a formal parameter in explicit data access routines.

file size and end of file The *size* of an MPI file is measured in bytes from the beginning of the file. A newly created file has a size of zero bytes. Using the size as an absolute displacement gives the position of the byte immediately following the last byte in the file. For any given view, the *end of file* is the offset of the first etype accessible in the current view starting after the last byte in the file.

file pointer A *file pointer* is an implicit offset maintained by MPI. “Individual file pointers” are file pointers that are local to each process that opened the file. A “shared file pointer” is a file pointer that is shared by the group of processes that opened the file.

file handle A *file handle* is an opaque object created by `MPI_FILE_OPEN` and freed by `MPI_FILE_CLOSE`. All operations on an open file reference the file through the file handle.

12.2 File Manipulation

12.2.1 Opening a File

`MPI_FILE_OPEN(comm, filename, amode, info, fh)`

| | | |
|-----|-----------------------|-------------------------------|
| IN | <code>comm</code> | communicator (handle) |
| IN | <code>filename</code> | name of file to open (string) |
| IN | <code>amode</code> | file access mode (integer) |
| IN | <code>info</code> | info object (handle) |
| OUT | <code>fh</code> | new file handle (handle) |

```
int MPI_File_open(MPI_Comm comm, char *filename, int amode, MPI_Info info,
                 MPI_File *fh)
```

```
MPI_FILE_OPEN(COMM, FILENAME, AMODE, INFO, FH, IERROR)
CHARACTER*(*) FILENAME
INTEGER COMM, AMODE, INFO, FH, IERROR
```

```
static MPI::File MPI::File::Open(const MPI::Intracomm& comm,
                                const char* filename, int amode, const MPI::Info& info)
```

`MPI_FILE_OPEN` opens the file identified by the file name `filename` on all processes in the `comm` communicator group. `MPI_FILE_OPEN` is a collective routine: all processes must provide the same value for `amode`, and all processes must provide `filenames` that reference the same file. (Values for `info` may vary.) `comm` must be an intracommunicator; it is erroneous to pass an intercommunicator to `MPI_FILE_OPEN`. Errors in `MPI_FILE_OPEN` are raised using the default file error handler (see Section 12.7, page 431). A process can open a file independently of other processes by using the `MPI_COMM_SELF` communicator. The file handle returned, `fh`, can be subsequently used to access the file until the file is closed using `MPI_FILE_CLOSE`. Before calling `MPI_FINALIZE`, the user is required to close (via `MPI_FILE_CLOSE`) all files that were opened with `MPI_FILE_OPEN`. Note that the communicator `comm` is unaffected by `MPI_FILE_OPEN` and continues to be usable in all

1 MPI routines (e.g., `MPI_SEND`). Furthermore, the use of `comm` will not interfere with I/O
 2 behavior.

3 The format for specifying the file name in the `filename` argument is implementation
 4 dependent and must be documented by the implementation.

5
 6 *Advice to implementors.* An implementation may require that `filename` include a
 7 string or strings specifying additional information about the file. Examples include
 8 the type of filesystem (e.g., a prefix of `ufs:`), a remote hostname (e.g., a prefix of
 9 `machine.univ.edu:`), or a file password (e.g., a suffix of `/PASSWORD=SECRET`).
 10 (*End of advice to implementors.*)

11
 12 *Advice to users.* On some implementations of MPI, the file namespace may not be
 13 identical from all processes of all applications. For example, `"/tmp/foo"` may denote
 14 different files on different processes, or a single file may have many names, dependent
 15 on process location. The user is responsible for ensuring that a single file is referenced
 16 by the `filename` argument, as it may be impossible for an implementation to detect
 17 this type of namespace error. (*End of advice to users.*)

18
 19 Initially, all processes view the file as a linear byte stream, and each process views data
 20 in its own native representation (no data representation conversion is performed). (POSIX
 21 files are linear byte streams in the native representation.) The file view can be changed via
 22 the `MPI_FILE_SET_VIEW` routine.

23 The following access modes are supported (specified in `amode`, a bit vector OR of the
 24 following integer constants):

- 25 • `MPI_MODE_RDONLY` — read only,
- 26
- 27 • `MPI_MODE_RDWR` — reading and writing,
- 28
- 29 • `MPI_MODE_WRONLY` — write only,
- 30
- 31 • `MPI_MODE_CREATE` — create the file if it does not exist,
- 32
- 33 • `MPI_MODE_EXCL` — error if creating file that already exists,
- 34
- 35 • `MPI_MODE_DELETE_ON_CLOSE` — delete file on close,
- 36
- 37 • `MPI_MODE_UNIQUE_OPEN` — file will not be concurrently opened elsewhere,
- 38
- 39 • `MPI_MODE_SEQUENTIAL` — file will only be accessed sequentially,
- 40
- 41 • `MPI_MODE_APPEND` — set initial position of all file pointers to end of file.

42
 43 *Advice to users.* C/C++ users can use bit vector OR (`|`) to combine these constants;
 44 Fortran 90 users can use the bit vector `IOR` intrinsic. Fortran 77 users can use (non-
 45 portably) bit vector `IOR` on systems that support it. Alternatively, Fortran users can
 portably use integer addition to OR the constants (each constant should appear at
 most once in the addition.). (*End of advice to users.*)

46
 47 *Advice to implementors.* The values of these constants must be defined such that
 48 the bitwise OR and the sum of any distinct set of these constants is equivalent. (*End
 of advice to implementors.*)

The modes `MPI_MODE_RDONLY`, `MPI_MODE_RDWR`, `MPI_MODE_WRONLY`, `MPI_MODE_CREATE`, and `MPI_MODE_EXCL` have identical semantics to their POSIX counterparts [31]. Exactly one of `MPI_MODE_RDONLY`, `MPI_MODE_RDWR`, or `MPI_MODE_WRONLY`, must be specified. It is erroneous to specify `MPI_MODE_CREATE` or `MPI_MODE_EXCL` in conjunction with `MPI_MODE_RDONLY`; it is erroneous to specify `MPI_MODE_SEQUENTIAL` together with `MPI_MODE_RDWR`.

The `MPI_MODE_DELETE_ON_CLOSE` mode causes the file to be deleted (equivalent to performing an `MPI_FILE_DELETE`) when the file is closed.

The `MPI_MODE_UNIQUE_OPEN` mode allows an implementation to optimize access by eliminating the overhead of file locking. It is erroneous to open a file in this mode unless the file will not be concurrently opened elsewhere.

Advice to users. For `MPI_MODE_UNIQUE_OPEN`, *not opened elsewhere* includes both inside and outside the MPI environment. In particular, one needs to be aware of potential external events which may open files (e.g., automated backup facilities). When `MPI_MODE_UNIQUE_OPEN` is specified, the user is responsible for ensuring that no such external events take place. (*End of advice to users.*)

The `MPI_MODE_SEQUENTIAL` mode allows an implementation to optimize access to some sequential devices (tapes and network streams). It is erroneous to attempt nonsequential access to a file that has been opened in this mode.

Specifying `MPI_MODE_APPEND` only guarantees that all shared and individual file pointers are positioned at the initial end of file when `MPI_FILE_OPEN` returns. Subsequent positioning of file pointers is application dependent. In particular, the implementation does not ensure that all writes are appended.

Errors related to the access mode are raised in the class `MPI_ERR_AMODE`.

The `info` argument is used to provide information regarding file access patterns and file system specifics (see Section 12.2.8, page 384). The constant `MPI_INFO_NULL` can be used when no `info` needs to be specified.

Advice to users. Some file attributes are inherently implementation dependent (e.g., file permissions). These attributes must be set using either the `info` argument or facilities outside the scope of MPI. (*End of advice to users.*)

Files are opened by default using nonatomic mode file consistency semantics (see Section 12.6.1, page 422). The more stringent atomic mode consistency semantics, required for atomicity of conflicting accesses, can be set using `MPI_FILE_SET_ATOMICITY`.

12.2.2 Closing a File

`MPI_FILE_CLOSE(fh)`

INOUT `fh` file handle (handle)

```
int MPI_File_close(MPI_File *fh)
```

```
MPI_FILE_CLOSE(FH, IERROR)
```

```
INTEGER FH, IERROR
```

```
1 void MPI::File::Close()
```

```
2
3 MPI_FILE_CLOSE first synchronizes file state (equivalent to performing an
4 MPI_FILE_SYNC), then closes the file associated with fh. The file is deleted if it was opened
5 with access mode MPI_MODE_DELETE_ON_CLOSE (equivalent to performing an
6 MPI_FILE_DELETE). MPI_FILE_CLOSE is a collective routine.
```

```
7
8 Advice to users. If the file is deleted on close, and there are other processes currently
9 accessing the file, the status of the file and the behavior of future accesses by these
10 processes are implementation dependent. (End of advice to users.)
```

```
11
12 The user is responsible for ensuring that all outstanding nonblocking requests and
13 split collective operations associated with fh made by a process have completed before that
14 process calls MPI_FILE_CLOSE.
```

```
15 The MPI_FILE_CLOSE routine deallocates the file handle object and sets fh to
16 MPI_FILE_NULL.
```

17 12.2.3 Deleting a File

```
20 MPI_FILE_DELETE(filename, info)
```

```
21
22 IN filename name of file to delete (string)
23 IN info info object (handle)
```

```
24
25 int MPI_File_delete(char *filename, MPI_Info info)
```

```
26 MPI_FILE_DELETE(FILENAME, INFO, IERROR)
```

```
27 CHARACTER*(*) FILENAME
28 INTEGER INFO, IERROR
```

```
29
30 static void MPI::File::Delete(const char* filename, const MPI::Info& info)
```

```
31
32 MPI_FILE_DELETE deletes the file identified by the file name filename. If the file does
33 not exist, MPI_FILE_DELETE raises an error in the class MPI_ERR_NO_SUCH_FILE.
```

```
34 The info argument can be used to provide information regarding file system specifics
35 (see Section 12.2.8, page 384). The constant MPI_INFO_NULL refers to the null info, and can
36 be used when no info needs to be specified.
```

```
37 If a process currently has the file open, the behavior of any access to the file (as well
38 as the behavior of any outstanding accesses) is implementation dependent. In addition,
39 whether an open file is deleted or not is also implementation dependent. If the file is not
40 deleted, an error in the class MPI_ERR_FILE_IN_USE or MPI_ERR_ACCESS will be raised. Errors
41 are raised using the default error handler (see Section 12.7, page 431).
```

12.2.4 Resizing a File

`MPI_FILE_SET_SIZE(fh, size)`

| | | |
|-------|------|---|
| INOUT | fh | file handle (handle) |
| IN | size | size to truncate or expand file (integer) |

```
int MPI_File_set_size(MPI_File fh, MPI_Offset size)
```

```
MPI_FILE_SET_SIZE(FH, SIZE, IERROR)
```

```
INTEGER FH, IERROR
```

```
INTEGER(KIND=MPI_OFFSET_KIND) SIZE
```

```
void MPI::File::Set_size(MPI::Offset size)
```

`MPI_FILE_SET_SIZE` resizes the file associated with the file handle `fh`. `size` is measured in bytes from the beginning of the file. `MPI_FILE_SET_SIZE` is collective; all processes in the group must pass identical values for `size`.

If `size` is smaller than the current file size, the file is truncated at the position defined by `size`. The implementation is free to deallocate file blocks located beyond this position.

If `size` is larger than the current file size, the file size becomes `size`. Regions of the file that have been previously written are unaffected. The values of data in the new regions in the file (those locations with displacements between old file size and `size`) are undefined. It is implementation dependent whether the `MPI_FILE_SET_SIZE` routine allocates file space—use `MPI_FILE_PREALLOCATE` to force file space to be reserved.

`MPI_FILE_SET_SIZE` does not affect the individual file pointers or the shared file pointer. If `MPI.MODE_SEQUENTIAL` mode was specified when the file was opened, it is erroneous to call this routine.

Advice to users. It is possible for the file pointers to point beyond the end of file after a `MPI_FILE_SET_SIZE` operation truncates a file. This is legal, and equivalent to seeking beyond the current end of file. (*End of advice to users.*)

All nonblocking requests and split collective operations on `fh` must be completed before calling `MPI_FILE_SET_SIZE`. Otherwise, calling `MPI_FILE_SET_SIZE` is erroneous. As far as consistency semantics are concerned, `MPI_FILE_SET_SIZE` is a write operation that conflicts with operations that access bytes at displacements between the old and new file sizes (see Section 12.6.1, page 422).

12.2.5 Preallocating Space for a File

`MPI_FILE_PREALLOCATE(fh, size)`

| | | |
|-------|------|------------------------------------|
| INOUT | fh | file handle (handle) |
| IN | size | size to preallocate file (integer) |

```
int MPI_File_preallocate(MPI_File fh, MPI_Offset size)
```

```

1 MPI_FILE_PREALLOCATE(FH, SIZE, IERROR)
2     INTEGER FH, IERROR
3     INTEGER(KIND=MPI_OFFSET_KIND) SIZE
4
5 void MPI::File::Preallocate(MPI::Offset size)

```

MPI_FILE_PREALLOCATE ensures that storage space is allocated for the first `size` bytes of the file associated with `fh`. MPI_FILE_PREALLOCATE is collective; all processes in the group must pass identical values for `size`. Regions of the file that have previously been written are unaffected. For newly allocated regions of the file, MPI_FILE_PREALLOCATE has the same effect as writing undefined data. If `size` is larger than the current file size, the file size increases to `size`. If `size` is less than or equal to the current file size, the file size is unchanged.

The treatment of file pointers, pending nonblocking accesses, and file consistency is the same as with MPI_FILE_SET_SIZE. If MPI_MODE_SEQUENTIAL mode was specified when the file was opened, it is erroneous to call this routine.

Advice to users. In some implementations, file preallocation may be expensive. (*End of advice to users.*)

12.2.6 Querying the Size of a File

```

23 MPI_FILE_GET_SIZE(fh, size)
24
25     IN      fh                file handle (handle)
26     OUT    size              size of the file in bytes (integer)

```

```

28 int MPI_File_get_size(MPI_File fh, MPI_Offset *size)

```

```

30 MPI_FILE_GET_SIZE(FH, SIZE, IERROR)
31     INTEGER FH, IERROR
32     INTEGER(KIND=MPI_OFFSET_KIND) SIZE

```

```

33 MPI::Offset MPI::File::Get_size() const

```

MPI_FILE_GET_SIZE returns, in `size`, the current size in bytes of the file associated with the file handle `fh`. As far as consistency semantics are concerned, MPI_FILE_GET_SIZE is a data access operation (see Section 12.6.1, page 422).

12.2.7 Querying File Parameters

```

42 MPI_FILE_GET_GROUP(fh, group)
43
44     IN      fh                file handle (handle)
45     OUT    group              group which opened the file (handle)

```

```

47 int MPI_File_get_group(MPI_File fh, MPI_Group *group)

```

```
MPI_FILE_GET_GROUP(FH, GROUP, IERROR)
```

```
INTEGER FH, GROUP, IERROR
```

```
MPI::Group MPI::File::Get_group() const
```

MPI_FILE_GET_GROUP returns a duplicate of the group of the communicator used to open the file associated with `fh`. The group is returned in `group`. The user is responsible for freeing `group`.

```
MPI_FILE_GET_AMODE(fh, amode)
```

```
IN      fh                file handle (handle)
```

```
OUT     amode             file access mode used to open the file (integer)
```

```
int MPI_File_get_amode(MPI_File fh, int *amode)
```

```
MPI_FILE_GET_AMODE(FH, AMODE, IERROR)
```

```
INTEGER FH, AMODE, IERROR
```

```
int MPI::File::Get_amode() const
```

MPI_FILE_GET_AMODE returns, in `amode`, the access mode of the file associated with `fh`.

Example 12.1 In Fortran 77, decoding an `amode` bit vector will require a routine such as the following:

```

      SUBROUTINE BIT_QUERY(TEST_BIT, MAX_BIT, AMODE, BIT_FOUND)
!
! TEST IF THE INPUT TEST_BIT IS SET IN THE INPUT AMODE
! IF SET, RETURN 1 IN BIT_FOUND, 0 OTHERWISE
!
      INTEGER TEST_BIT, AMODE, BIT_FOUND, CP_AMODE, HIFOUND
      BIT_FOUND = 0
      CP_AMODE = AMODE
100  CONTINUE
      LBIT = 0
      HIFOUND = 0
      DO 20 L = MAX_BIT, 0, -1
          MATCHER = 2**L
          IF (CP_AMODE .GE. MATCHER .AND. HIFOUND .EQ. 0) THEN
              HIFOUND = 1
              LBIT = MATCHER
              CP_AMODE = CP_AMODE - MATCHER
          END IF
20  CONTINUE
      IF (HIFOUND .EQ. 1 .AND. LBIT .EQ. TEST_BIT) BIT_FOUND = 1
      IF (BIT_FOUND .EQ. 0 .AND. HIFOUND .EQ. 1 .AND. &
          CP_AMODE .GT. 0) GO TO 100
      END

```

This routine could be called successively to decode `amode`, one bit at a time. For example, the following code fragment would check for `MPI_MODE_RDONLY`.

```

CALL BIT_QUERY(MPI_MODE_RDONLY, 30, AMODE, BIT_FOUND)
IF (BIT_FOUND .EQ. 1) THEN
  PRINT *, ' FOUND READ-ONLY BIT IN AMODE=', AMODE
ELSE
  PRINT *, ' READ-ONLY BIT NOT FOUND IN AMODE=', AMODE
END IF

```

12.2.8 File Info

Hints specified via `info` (see Section 8.1, page 273) allow a user to provide information such as file access patterns and file system specifics to direct optimization. Providing hints may enable an implementation to deliver increased I/O performance or minimize the use of system resources. However, hints do not change the semantics of any of the I/O interfaces. In other words, an implementation is free to ignore all hints. Hints are specified on a per file basis, in `MPI_FILE_OPEN`, `MPI_FILE_DELETE`, `MPI_FILE_SET_VIEW`, and `MPI_FILE_SET_INFO`, via the opaque `info` object. *When an info object that specifies a subset of valid hints is passed to `MPI_FILE_SET_VIEW` or `MPI_FILE_SET_INFO`, there will be no effect on previously set or defaulted hints that the info does not specify.*

Advice to implementors. It may happen that a program is coded with hints for one system, and later executes on another system that does not support these hints. In general, unsupported hints should simply be ignored. Needless to say, no hint can be mandatory. However, for each hint used by a specific implementation, a default value must be provided when the user does not specify a value for this hint. (*End of advice to implementors.*)

`MPI_FILE_SET_INFO(fh, info)`

| | | |
|-------|------|----------------------|
| INOUT | fh | file handle (handle) |
| IN | info | info object (handle) |

```
int MPI_File_set_info(MPI_File fh, MPI_Info info)
```

```
MPI_FILE_SET_INFO(FH, INFO, IERROR)
```

```
INTEGER FH, INFO, IERROR
```

```
void MPI::File::Set_info(const MPI::Info& info)
```

`MPI_FILE_SET_INFO` sets new values for the hints of the file associated with `fh`. `MPI_FILE_SET_INFO` is a collective routine. The `info` object may be different on each process, but any `info` entries that an implementation requires to be the same on all processes must appear with the same value in each process's `info` object.

Advice to users. Many `info` items that an implementation can use when it creates or opens a file cannot easily be changed once the file has been created or opened. Thus, an implementation may ignore hints issued in this call that it would have accepted in an open call. (*End of advice to users.*)

```
MPI_FILE_GET_INFO(fh, info_used)
```

```
IN      fh                file handle (handle)
OUT     info_used         new info object (handle)
```

```
int MPI_File_get_info(MPI_File fh, MPI_Info *info_used)
```

```
MPI_FILE_GET_INFO(FH, INFO_USED, IERROR)
    INTEGER FH, INFO_USED, IERROR
```

```
MPI::Info MPI::File::Get_info() const
```

MPI_FILE_GET_INFO returns a new info object containing the hints of the file associated with fh. The current setting of all hints actually used by the system related to this open file is returned in info_used. *If no such hints exist, a handle to a newly created info object is returned that contains no key/value pair.* The user is responsible for freeing info_used via MPI_INFO_FREE.

Advice to users. The info object returned in info_used will contain all hints currently active for this file. This set of hints may be greater or smaller than the set of hints passed in to MPI_FILE_OPEN, MPI_FILE_SET_VIEW, and MPI_FILE_SET_INFO, as the system may not recognize some hints set by the user, and may recognize other hints that the user has not set. (*End of advice to users.*)

Reserved File Hints

Some potentially useful hints (info key values) are outlined below. The following key values are reserved. An implementation is not required to interpret these key values, but if it does interpret the key value, it must provide the functionality described. (For more details on “info,” see Section 8.1, page 273.)

These hints mainly affect access patterns and the layout of data on parallel I/O devices. For each hint name introduced, we describe the purpose of the hint, and the type of the hint value. The “[SAME]” annotation specifies that the hint values provided by all participating processes must be identical; otherwise the program is erroneous. In addition, some hints are context dependent, and are only used by an implementation at specific times (e.g., file_perm is only useful during file creation).

access_style (comma separated list of strings): This hint specifies the manner in which the file will be accessed until the file is closed or until the access_style key value is altered. The hint value is a comma separated list of the following: read_once, write_once, read_mostly, write_mostly, sequential, reverse_sequential, and random.

collective_buffering (boolean) [SAME]: This hint specifies whether the application may benefit from collective buffering. Collective buffering is an optimization performed on collective accesses. Accesses to the file are performed on behalf of all processes in the group by a number of target nodes. These target nodes coalesce small requests into large disk accesses. Legal values for this key are true and false. Collective buffering parameters are further directed via additional hints: cb_block_size, cb_buffer_size, and cb_nodes.

1 **cb_block_size (integer) [SAME]:** This hint specifies the block size to be used for collective
2 buffering file access. *Target nodes* access data in chunks of this size. The chunks are
3 distributed among target nodes in a round-robin (CYCLIC) pattern.

4 **cb_buffer_size (integer) [SAME]:** This hint specifies the total buffer space that can be used
5 for collective buffering on each target node, usually a multiple of `cb_block_size`.
6

7 **cb_nodes (integer) [SAME]:** This hint specifies the number of target nodes to be used for
8 collective buffering.
9

10 **chunked (comma separated list of integers) [SAME]:** This hint specifies that the file
11 consists of a multidimensional array that is often accessed by subarrays. The value
12 for this hint is a comma separated list of array dimensions, starting from the most
13 significant one (for an array stored in row-major order, as in C, the most significant
14 dimension is the first one; for an array stored in column-major order, as in Fortran, the
15 most significant dimension is the last one, and array dimensions should be reversed).
16

17 **chunked_item (comma separated list of integers) [SAME]:** This hint specifies the size
18 of each array entry, in bytes.

19 **chunked_size (comma separated list of integers) [SAME]:** This hint specifies the dimen-
20 sions of the subarrays. This is a comma separated list of array dimensions, starting
21 from the most significant one.
22

23 **filename (string):** This hint specifies the file name used when the file was opened. If the
24 implementation is capable of returning the file name of an open file, it will be returned
25 using this key by `MPI_FILE_GET_INFO`. This key is ignored when passed to
26 `MPI_FILE_OPEN`, `MPI_FILE_SET_VIEW`, `MPI_FILE_SET_INFO`, and
27 `MPI_FILE_DELETE`.

28 **file_perm (string) [SAME]:** This hint specifies the file permissions to use for file creation.
29 Setting this hint is only useful when passed to `MPI_FILE_OPEN` with an `amode` that
30 includes `MPI_MODE_CREATE`. The set of legal values for this key is implementation
31 dependent.
32

33 **io_node_list (comma separated list of strings) [SAME]:** This hint specifies the list of
34 I/O devices that should be used to store the file. This hint is most relevant when the
35 file is created.

36 **nb_proc (integer) [SAME]:** This hint specifies the number of parallel processes that will
37 typically be assigned to run programs that access this file. This hint is most relevant
38 when the file is created.
39

40 **num_io_nodes (integer) [SAME]:** This hint specifies the number of I/O devices in the sys-
41 tem. This hint is most relevant when the file is created.

42 **striping_factor (integer) [SAME]:** This hint specifies the number of I/O devices that the
43 file should be striped across, and is relevant only when the file is created.
44

45 **striping_unit (integer) [SAME]:** This hint specifies the suggested striping unit to be used
46 for this file. The striping unit is the amount of consecutive data assigned to one I/O
47 device before progressing to the next device, when striping across a number of devices.
48 It is expressed in bytes. This hint is relevant only when the file is created.

12.3 File Views

```
MPI_FILE_SET_VIEW(fh, disp, etype, filetype, datarep, info)
```

| | | |
|-------|----------|------------------------------|
| INOUT | fh | file handle (handle) |
| IN | disp | displacement (integer) |
| IN | etype | elementary datatype (handle) |
| IN | filetype | filetype (handle) |
| IN | datarep | data representation (string) |
| IN | info | info object (handle) |

```
int MPI_File_set_view(MPI_File fh, MPI_Offset disp, MPI_Datatype etype,
                    MPI_Datatype filetype, char *datarep, MPI_Info info)
```

```
MPI_FILE_SET_VIEW(FH, DISP, ETYPE, FILETYPE, DATAREP, INFO, IERROR)
    INTEGER FH, ETYPE, FILETYPE, INFO, IERROR
    CHARACTER*(*) DATAREP
    INTEGER(KIND=MPI_OFFSET_KIND) DISP
```

```
void MPI::File::Set_view(MPI::Offset disp, const MPI::Datatype& etype,
                        const MPI::Datatype& filetype, const char* datarep,
                        const MPI::Info& info)
```

The `MPI_FILE_SET_VIEW` routine changes the process's view of the data in the file. The start of the view is set to `disp`; the type of data is set to `etype`; the distribution of data to processes is set to `filetype`; and the representation of data in the file is set to `datarep`. In addition, `MPI_FILE_SET_VIEW` resets the individual file pointers and the shared file pointer to zero. `MPI_FILE_SET_VIEW` is collective; the values for `datarep` and the extents of `etype` in the file data representation must be identical on all processes in the group; values for `disp`, `filetype`, and `info` may vary. The datatypes passed in `etype` and `filetype` must be committed.

The `etype` always specifies the data layout in the file. If `etype` is a portable datatype (see Section 2.4, page 11), the extent of `etype` is computed by scaling any displacements in the datatype to match the file data representation. If `etype` is not a portable datatype, no scaling is done when computing the extent of `etype`. The user must be careful when using nonportable `etypes` in heterogeneous environments; see Section 12.5.1, page 414 for further details.

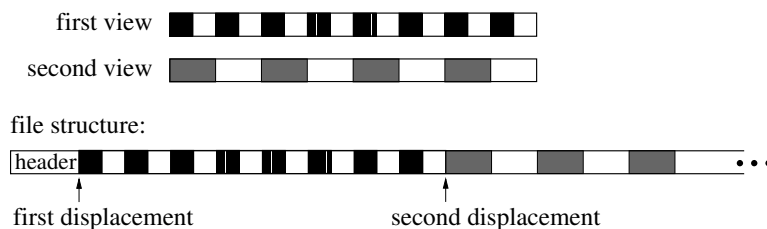
If `MPI_MODE_SEQUENTIAL` mode was specified when the file was opened, the special displacement `MPI_DISPLACEMENT_CURRENT` must be passed in `disp`. This sets the displacement to the current position of the shared file pointer. `MPI_DISPLACEMENT_CURRENT` is invalid unless the `amode` for the file has `MPI_MODE_SEQUENTIAL` set.

Rationale. For some sequential files, such as those corresponding to magnetic tapes or streaming network connections, the *displacement* may not be meaningful. `MPI_DISPLACEMENT_CURRENT` allows the view to be changed for these types of files. (*End of rationale.*)

1 *Advice to implementors.* It is expected that a call to `MPI_FILE_SET_VIEW` will immediately follow `MPI_FILE_OPEN` in numerous instances. A high quality implementation
 2 will ensure that this behavior is efficient. (*End of advice to implementors.*)
 3

4
 5 The `disp` displacement argument specifies the position (absolute offset in bytes from
 6 the beginning of the file) where the view begins.

7
 8 *Advice to users.* `disp` can be used to skip headers or when the file includes a sequence
 9 of data segments that are to be accessed in different patterns (see Figure 12.3). Sep-
 10 arate views, each using a different displacement and filetype, can be used to access
 11 each segment.



12
 13
 14
 15
 16
 17
 18
 19
 20 Figure 12.3: Displacements

21 (*End of advice to users.*)

22
 23 An *etype* (elementary datatype) is the unit of data access and positioning. It can be
 24 any MPI predefined or derived datatype. Derived etypes can be constructed by using any
 25 of the MPI datatype constructor routines, provided all resulting typemap displacements are
 26 nonnegative and monotonically nondecreasing. Data access is performed in etype units,
 27 reading or writing whole data items of type etype. Offsets are expressed as a count of
 28 etypes; file pointers point to the beginning of etypes.
 29

30 *Advice to users.* In order to ensure interoperability in a heterogeneous environment,
 31 additional restrictions must be observed when constructing the etype (see Section 12.5,
 32 page 412). (*End of advice to users.*)
 33

34 A filetype is either a single etype or a derived MPI datatype constructed from multiple
 35 instances of the same etype. In addition, the extent of any hole in the filetype must be
 36 a multiple of the etype's extent. These displacements are not required to be distinct, but
 37 they cannot be negative, and they must be monotonically nondecreasing.

38 If the file is opened for writing, neither the etype nor the filetype is permitted to contain
 39 overlapping regions. This restriction is equivalent to the "datatype used in a receive cannot
 40 specify overlapping regions" restriction for communication. Note that filetypes from different
 41 processes may still overlap each other.

42 If filetype has holes in it, then the data in the holes is inaccessible to the calling process.
 43 However, the `disp`, `etype` and `filetype` arguments can be changed via future calls to
 44 `MPI_FILE_SET_VIEW` to access a different part of the file.

45 It is erroneous to use absolute addresses in the construction of the etype and filetype.

46 The `info` argument is used to provide information regarding file access patterns and
 47 file system specifics to direct optimization (see Section 12.2.8, page 384). The constant
 48 `MPI_INFO_NULL` refers to the null info and can be used when no info needs to be specified.

The `datarep` argument is a string that specifies the representation of data in the file. See the file interoperability section (Section 12.5, page 412) for details and a discussion of valid values.

The user is responsible for ensuring that all nonblocking requests and split collective operations on `fh` have been completed before calling `MPI_FILE_SET_VIEW`—otherwise, the call to `MPI_FILE_SET_VIEW` is erroneous.

`MPI_FILE_GET_VIEW(fh, disp, etype, filetype, datarep)`

| | | |
|-----|-----------------------|------------------------------|
| IN | <code>fh</code> | file handle (handle) |
| OUT | <code>disp</code> | displacement (integer) |
| OUT | <code>etype</code> | elementary datatype (handle) |
| OUT | <code>filetype</code> | filetype (handle) |
| OUT | <code>datarep</code> | data representation (string) |

```
int MPI_File_get_view(MPI_File fh, MPI_Offset *disp, MPI_Datatype *etype,
                    MPI_Datatype *filetype, char *datarep)
```

```
MPI_FILE_GET_VIEW(FH, DISP, ETYPE, FILETYPE, DATAREP, IERROR)
    INTEGER FH, ETYPE, FILETYPE, IERROR
    CHARACTER*(*) DATAREP
    INTEGER(KIND=MPI_OFFSET_KIND) DISP
```

```
void MPI::File::Get_view(MPI::Offset& disp, MPI::Datatype& etype,
                        MPI::Datatype& filetype, char* datarep) const
```

`MPI_FILE_GET_VIEW` returns the process's view of the data in the file. The current value of the displacement is returned in `disp`. The `etype` and `filetype` are new datatypes with typemaps equal to the typemaps of the current `etype` and `filetype`, respectively.

The data representation is returned in `datarep`. The user is responsible for ensuring that `datarep` is large enough to hold the returned data representation string. The length of a data representation string is limited to the value of `MPI_MAX_DATAREP_STRING`.

In addition, if a portable datatype was used to set the current view, then the corresponding datatype returned by `MPI_FILE_GET_VIEW` is also a portable datatype. If `etype` or `filetype` are derived datatypes, the user is responsible for freeing them. The `etype` and `filetype` returned are both in a committed state.

12.4 Data Access

12.4.1 Data Access Routines

Data is moved between files and processes by issuing read and write calls. There are three orthogonal aspects to data access: positioning (explicit offset *vs.* implicit file pointer), synchronism (blocking *vs.* nonblocking and split collective), and coordination (noncollective *vs.* collective). The following combinations of these data access routines, including two types of file pointers (individual and shared) are provided in Table 12.1.

POSIX `read()/pread()` and `write()/fwrite()` are blocking, noncollective operations and use individual file pointers. The MPI equivalents are `MPI_FILE_READ` and `MPI_FILE_WRITE`.

| positioning | synchronism | coordination | |
|---------------------------------|---|---|--|
| | | <i>noncollective</i> | <i>collective</i> |
| <i>explicit offsets</i> | <i>blocking</i> | MPI_FILE_READ_AT MPI_FILE_WRITE_AT | MPI_FILE_READ_AT_ALL MPI_FILE_WRITE_AT_ALL |
| | <i>nonblocking & split collective</i> | MPI_FILE_IREAD_AT MPI_FILE_IWRITE_AT | MPI_FILE_READ_AT_ALL_BEGIN MPI_FILE_READ_AT_ALL_END MPI_FILE_WRITE_AT_ALL_BEGIN MPI_FILE_WRITE_AT_ALL_END |
| <i>individual file pointers</i> | <i>blocking</i> | MPI_FILE_READ MPI_FILE_WRITE | MPI_FILE_READ_ALL MPI_FILE_WRITE_ALL |
| | <i>nonblocking & split collective</i> | MPI_FILE_IREAD MPI_FILE_IWRITE | MPI_FILE_READ_ALL_BEGIN MPI_FILE_READ_ALL_END MPI_FILE_WRITE_ALL_BEGIN MPI_FILE_WRITE_ALL_END |
| <i>shared file pointer</i> | <i>blocking</i> | MPI_FILE_READ_SHARED MPI_FILE_WRITE_SHARED | MPI_FILE_READ_ORDERED MPI_FILE_WRITE_ORDERED |
| | <i>nonblocking & split collective</i> | MPI_FILE_IREAD_SHARED MPI_FILE_IWRITE_SHARED | MPI_FILE_READ_ORDERED_BEGIN MPI_FILE_READ_ORDERED_END MPI_FILE_WRITE_ORDERED_BEGIN MPI_FILE_WRITE_ORDERED_END |

Table 12.1: Data access routines

Implementations of data access routines may buffer data to improve performance. This does not affect reads, as the data is always available in the user’s buffer after a read operation completes. For writes, however, the `MPI_FILE_SYNC` routine provides the only guarantee that data has been transferred to the storage device.

Positioning

MPI provides three types of positioning for data access routines: explicit offsets, individual file pointers, and shared file pointers. The different positioning methods may be mixed within the same program and do not affect each other.

The data access routines that accept explicit offsets contain `_AT` in their name (e.g., `MPI_FILE_WRITE_AT`). Explicit offset operations perform data access at the file position given directly as an argument—no file pointer is used nor updated. Note that this is not equivalent to an atomic seek-and-read or seek-and-write operation, as no “seek” is issued. Operations with explicit offsets are described in Section 12.4.2, page 392.

The names of the individual file pointer routines contain no positional qualifier (e.g., `MPI_FILE_WRITE`). Operations with individual file pointers are described in Section 12.4.3, page 396. The data access routines that use shared file pointers contain `_SHARED` or `_ORDERED` in their name (e.g., `MPI_FILE_WRITE_SHARED`). Operations with shared file pointers are described in Section 12.4.4, page 401.

The main semantic issues with MPI-maintained file pointers are how and when they are updated by I/O operations. In general, each I/O operation leaves the file pointer pointing to the next data item after the last one that is accessed by the operation. In a nonblocking or split collective operation, the pointer is updated by the call that initiates the I/O, possibly before the access completes.

More formally,

$$new_file_offset = old_file_offset + \frac{elements(datatype)}{elements(etype)} \times count$$

where *count* is the number of *datatype* items to be accessed, *elements(X)* is the number of predefined datatypes in the typemap of *X*, and *old_file_offset* is the value of the implicit offset before the call. The file position, *new_file_offset*, is in terms of a count of etypes relative to the current view.

Synchronization

MPI supports blocking and nonblocking I/O routines.

A *blocking* I/O call will not return until the I/O request is completed.

A *nonblocking* I/O call initiates an I/O operation, but does not wait for it to complete. Given suitable hardware, this allows the transfer of data out/in the user's buffer to proceed concurrently with computation. A separate *request complete* call (MPI_WAIT, MPI_TEST, or any of their variants) is needed to complete the I/O request, i.e., to confirm that the data has been read or written and that it is safe for the user to reuse the buffer. The nonblocking versions of the routines are named MPI_FILE_IXXX, where the I stands for immediate.

It is erroneous to access the local buffer of a nonblocking data access operation, or to use that buffer as the source or target of other communications, between the initiation and completion of the operation.

The split collective routines support a restricted form of “nonblocking” operations for collective data access (see Section 12.4.5, page 406).

Coordination

Every noncollective data access routine MPI_FILE_XXX has a collective counterpart. For most routines, this counterpart is MPI_FILE_XXX_ALL or a pair of MPI_FILE_XXX_BEGIN and MPI_FILE_XXX_END. The counterparts to the MPI_FILE_XXX_SHARED routines are MPI_FILE_XXX_ORDERED.

The completion of a noncollective call only depends on the activity of the calling process. However, the completion of a collective call (which must be called by all members of the process group) may depend on the activity of the other processes participating in the collective call. See Section 12.6.4, page 425, for rules on semantics of collective calls.

Collective operations may perform much better than their noncollective counterparts, as global data accesses have significant potential for automatic optimization.

Data Access Conventions

Data is moved between files and processes by calling read and write routines. Read routines move data from a file into memory. Write routines move data from memory into a file. The file is designated by a file handle, *fh*. The location of the file data is specified by an offset into the current view. The data in memory is specified by a triple: *buf*, *count*, and *datatype*. Upon completion, the amount of data accessed by the calling process is returned in a *status*.

An offset designates the starting position in the file for an access. The offset is always in etype units relative to the current view. Explicit offset routines pass *offset* as an argument (negative values are erroneous). The file pointer routines use implicit offsets maintained by MPI.

1 A data access routine attempts to transfer (read or write) count data items of type
2 `datatype` between the user's buffer `buf` and the file. The `datatype` passed to the routine
3 must be a committed datatype. The layout of data in memory corresponding to `buf`, `count`,
4 `datatype` is interpreted the same way as in MPI-1 communication functions; see Section 3.12.5
5 in [23]. The data is accessed from those parts of the file specified by the current view
6 (Section 12.3, page 387). The type signature of `datatype` must match the type signature
7 of some number of contiguous copies of the `etype` of the current view. As in a receive, it
8 is erroneous to specify a `datatype` for reading that contains overlapping regions (areas of
9 memory which would be stored into more than once).

10 The nonblocking data access routines indicate that MPI can start a data access and
11 associate a request handle, `request`, with the I/O operation. Nonblocking operations are
12 completed via `MPI_TEST`, `MPI_WAIT`, or any of their variants.

13 Data access operations, when completed, return the amount of data accessed in `status`.

14
15 *Advice to users.* To prevent problems with the argument copying and register opti-
16 mization done by Fortran compilers, please note the hints in subsections “Problems
17 Due to Data Copying and Sequence Association,” and “A Problem with Register
18 Optimization” in Section 13.2.2, pages 451 and 454. (*End of advice to users.*)

19
20 For blocking routines, `status` is returned directly. For nonblocking routines and split
21 collective routines, `status` is returned when the operation is completed. The number of
22 `datatype` entries and predefined elements accessed by the calling process can be extracted
23 from `status` by using `MPI_GET_COUNT` and `MPI_GET_ELEMENTS`, respectively. The inter-
24 pretation of the `MPI_ERROR` field is the same as for other operations — normally undefined,
25 but meaningful if an MPI routine returns `MPI_ERR_IN_STATUS`. The user can pass (in C
26 and Fortran) `MPI_STATUS_IGNORE` in the `status` argument if the return value of this argu-
27 ment is not needed. In C++, the `status` argument is optional. The `status` can be passed
28 to `MPI_TEST_CANCELLED` to determine if the operation was cancelled. All other fields of
29 `status` are undefined.

30 When reading, a program can detect the end of file by noting that the amount of data
31 read is less than the amount requested. Writing past the end of file increases the file size.
32 The amount of data accessed will be the amount requested, unless an error is raised (or a
33 read reaches the end of file).

34 12.4.2 Data Access with Explicit Offsets

35
36 If `MPI_MODE_SEQUENTIAL` mode was specified when the file was opened, it is erroneous to
37 call the routines in this section.
38
39
40
41
42
43
44
45
46
47
48

```

MPI_FILE_READ_AT(fh, offset, buf, count, datatype, status) 1
    IN      fh      file handle (handle) 2
    IN      offset  file offset (integer) 3
    OUT     buf     initial address of buffer (choice) 4
    IN      count   number of elements in buffer (integer) 5
    IN      datatype  datatype of each buffer element (handle) 6
    OUT     status  status object (Status) 7
  8
  9
  10
  int MPI_File_read_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
  11
  12
  MPI_Datatype datatype, MPI_Status *status) 13
  14
MPI_FILE_READ_AT(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR) 15
  <type> BUF(*) 16
  INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR 17
  INTEGER(KIND=MPI_OFFSET_KIND) OFFSET 18
  19
void MPI::File::Read_at(MPI::Offset offset, void* buf, int count,
  20
  const MPI::Datatype& datatype, MPI::Status& status) 21
  22
void MPI::File::Read_at(MPI::Offset offset, void* buf, int count,
  23
  const MPI::Datatype& datatype) 24
  25
  MPI_FILE_READ_AT reads a file beginning at the position specified by offset.
  26
  27
MPI_FILE_READ_AT_ALL(fh, offset, buf, count, datatype, status) 28
    IN      fh      file handle (handle) 29
    IN      offset  file offset (integer) 30
    OUT     buf     initial address of buffer (choice) 31
    IN      count   number of elements in buffer (integer) 32
    IN      datatype  datatype of each buffer element (handle) 33
    OUT     status  status object (Status) 34
  35
  36
  int MPI_File_read_at_all(MPI_File fh, MPI_Offset offset, void *buf,
  37
  int count, MPI_Datatype datatype, MPI_Status *status) 38
  39
MPI_FILE_READ_AT_ALL(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR) 40
  <type> BUF(*) 41
  INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR 42
  INTEGER(KIND=MPI_OFFSET_KIND) OFFSET 43
  44
void MPI::File::Read_at_all(MPI::Offset offset, void* buf, int count,
  45
  const MPI::Datatype& datatype, MPI::Status& status) 46
  47
void MPI::File::Read_at_all(MPI::Offset offset, void* buf, int count,
  48
  const MPI::Datatype& datatype)

```

1 MPI_FILE_READ_AT_ALL is a collective version of the blocking MPI_FILE_READ_AT
2 interface.

3
4
5 MPI_FILE_WRITE_AT(fh, offset, buf, count, datatype, status)

| | | | |
|----|-------|----------|--|
| 6 | INOUT | fh | file handle (handle) |
| 7 | IN | offset | file offset (integer) |
| 8 | IN | buf | initial address of buffer (choice) |
| 9 | IN | count | number of elements in buffer (integer) |
| 10 | IN | datatype | datatype of each buffer element (handle) |
| 11 | IN | count | number of elements in buffer (integer) |
| 12 | IN | datatype | datatype of each buffer element (handle) |
| 13 | OUT | status | status object (Status) |

14
15 int MPI_File_write_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
16 MPI_Datatype datatype, MPI_Status *status)

17 MPI_FILE_WRITE_AT(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR)
18 <type> BUF(*)
19 INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
20 INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
21

22 void MPI::File::Write_at(MPI::Offset offset, const void* buf, int count,
23 const MPI::Datatype& datatype, MPI::Status& status)

24 void MPI::File::Write_at(MPI::Offset offset, const void* buf, int count,
25 const MPI::Datatype& datatype)
26

27 MPI_FILE_WRITE_AT writes a file beginning at the position specified by offset.

28
29
30 MPI_FILE_WRITE_AT_ALL(fh, offset, buf, count, datatype, status)

| | | | |
|----|-------|----------|--|
| 31 | INOUT | fh | file handle (handle) |
| 32 | IN | offset | file offset (integer) |
| 33 | IN | buf | initial address of buffer (choice) |
| 34 | IN | count | number of elements in buffer (integer) |
| 35 | IN | count | number of elements in buffer (integer) |
| 36 | IN | datatype | datatype of each buffer element (handle) |
| 37 | IN | datatype | datatype of each buffer element (handle) |
| 38 | OUT | status | status object (Status) |

39
40 int MPI_File_write_at_all(MPI_File fh, MPI_Offset offset, void *buf,
41 int count, MPI_Datatype datatype, MPI_Status *status)

42 MPI_FILE_WRITE_AT_ALL(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR)
43 <type> BUF(*)
44 INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
45 INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
46

47 void MPI::File::Write_at_all(MPI::Offset offset, const void* buf,
48 int count, const MPI::Datatype& datatype, MPI::Status& status)


```
void MPI::File::Write_at_all(MPI::Offset offset, const void* buf,
                             int count, const MPI::Datatype& datatype)
```

MPI_FILE_WRITE_AT_ALL is a collective version of the blocking MPI_FILE_WRITE_AT interface.

```
MPI_FILE_IREAD_AT(fh, offset, buf, count, datatype, request)
```

| | | |
|-----|----------|--|
| IN | fh | file handle (handle) |
| IN | offset | file offset (integer) |
| OUT | buf | initial address of buffer (choice) |
| IN | count | number of elements in buffer (integer) |
| IN | datatype | datatype of each buffer element (handle) |
| OUT | request | request object (handle) |

```
int MPI_File_iread_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
                     MPI_Datatype datatype, MPI_Request *request)
```

```
MPI_FILE_IREAD_AT(FH, OFFSET, BUF, COUNT, DATATYPE, REQUEST, IERROR)
```

```
<type> BUF(*)
```

```
INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
```

```
INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
```

```
MPI::Request MPI::File::Iread_at(MPI::Offset offset, void* buf, int count,
                                  const MPI::Datatype& datatype)
```

MPI_FILE_IREAD_AT is a nonblocking version of the MPI_FILE_READ_AT interface.

```
MPI_FILE_IWRITE_AT(fh, offset, buf, count, datatype, request)
```

| | | |
|-------|----------|--|
| INOUT | fh | file handle (handle) |
| IN | offset | file offset (integer) |
| IN | buf | initial address of buffer (choice) |
| IN | count | number of elements in buffer (integer) |
| IN | datatype | datatype of each buffer element (handle) |
| OUT | request | request object (handle) |

```
int MPI_File_ fwrite_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
                       MPI_Datatype datatype, MPI_Request *request)
```

```
MPI_FILE_IWRITE_AT(FH, OFFSET, BUF, COUNT, DATATYPE, REQUEST, IERROR)
```

```
<type> BUF(*)
```

```
INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
```

```
INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
```

```
MPI::Request MPI::File::Iwrite_at(MPI::Offset offset, const void* buf,
                                   int count, const MPI::Datatype& datatype)
```

1 MPI_FILE_IWRITE_AT is a nonblocking version of the MPI_FILE_WRITE_AT interface.

3 12.4.3 Data Access with Individual File Pointers

4 MPI maintains one individual file pointer per process per file handle. The current value
5 of this pointer implicitly specifies the offset in the data access routines described in this
6 section. These routines only use and update the individual file pointers maintained by MPI.
7 The shared file pointer is not used nor updated.

8 The individual file pointer routines have the same semantics as the data access with
9 explicit offset routines described in Section 12.4.2, page 392, with the following modification:

- 10 • the offset is defined to be the current value of the MPI-maintained individual file
11 pointer.

12 After an individual file pointer operation is initiated, the individual file pointer is updated
13 to point to the next etype after the last one that will be accessed. The file pointer is updated
14 relative to the current view of the file.

15 If MPI_MODE_SEQUENTIAL mode was specified when the file was opened, it is erroneous
16 to call the routines in this section, with the exception of MPI_FILE_GET_BYTE_OFFSET.

17 MPI_FILE_READ(fh, buf, count, datatype, status)

| | | | |
|----|-------|----------|--|
| 18 | INOUT | fh | file handle (handle) |
| 19 | OUT | buf | initial address of buffer (choice) |
| 20 | IN | count | number of elements in buffer (integer) |
| 21 | IN | datatype | datatype of each buffer element (handle) |
| 22 | OUT | status | status object (Status) |

23 int MPI_File_read(MPI_File fh, void *buf, int count, MPI_Datatype datatype,
24 MPI_Status *status)

25 MPI_FILE_READ(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)

26 <type> BUF(*)

27 INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR

28 void MPI::File::Read(void* buf, int count, const MPI::Datatype& datatype,
29 MPI::Status& status)

30 void MPI::File::Read(void* buf, int count, const MPI::Datatype& datatype)

31 MPI_FILE_READ reads a file using the individual file pointer.

32 **Example 12.2** The following Fortran code fragment is an example of reading a file until
33 the end of file is reached:

```
34       ! Read a preexisting input file until all data has been read.
35       ! Call routine "process_input" if all requested data is read.
36       ! The Fortran 90 "exit" statement exits the loop.
```

```

integer  bufsize, numread, totprocessed, status(MPI_STATUS_SIZE)  1
parameter (bufsize=100)  2
real     localbuffer(bufsize)  3
  4
call MPI_FILE_OPEN( MPI_COMM_WORLD, 'myoldfile', &  5
                    MPI_MODE_RDONLY, MPI_INFO_NULL, myfh, ierr )  6
call MPI_FILE_SET_VIEW( myfh, 0, MPI_REAL, MPI_REAL, 'native', &  7
                       MPI_INFO_NULL, ierr )  8
totprocessed = 0  9
do  10
    call MPI_FILE_READ( myfh, localbuffer, bufsize, MPI_REAL, &  11
                       status, ierr )  12
    call MPI_GET_COUNT( status, MPI_REAL, numread, ierr )  13
    call process_input( localbuffer, numread )  14
    totprocessed = totprocessed + numread  15
    if ( numread < bufsize ) exit  16
enddo  17
  18
write(6,1001) numread, bufsize, totprocessed  19
1001 format( "No more data:  read", I3, "and expected", I3, &  20
            "Processed total of", I6, "before terminating job." )  21
  22
call MPI_FILE_CLOSE( myfh, ierr )  23
  24
  25
  26

```

MPI_FILE_READ_ALL(fh, buf, count, datatype, status) 27

| | | | |
|-------|----------|--|----|
| INOUT | fh | file handle (handle) | 28 |
| OUT | buf | initial address of buffer (choice) | 29 |
| IN | count | number of elements in buffer (integer) | 30 |
| IN | datatype | datatype of each buffer element (handle) | 31 |
| OUT | status | status object (Status) | 32 |

```

int MPI_File_read_all(MPI_File fh, void *buf, int count,  33
                     MPI_Datatype datatype, MPI_Status *status)  34
  35
  36
  37

```

```

MPI_FILE_READ_ALL(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)  38
<type> BUF(*)  39
INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR  40
  41

```

```

void MPI::File::Read_all(void* buf, int count,  42
                          const MPI::Datatype& datatype, MPI::Status& status)  43
  44

```

```

void MPI::File::Read_all(void* buf, int count,  44
                          const MPI::Datatype& datatype)  45
  46

```

MPI_FILE_READ_ALL is a collective version of the blocking MPI_FILE_READ interface. 47

48

```

1 MPI_FILE_WRITE(fh, buf, count, datatype, status)
2     INOUT   fh           file handle (handle)
3
4     IN      buf         initial address of buffer (choice)
5
6     IN      count       number of elements in buffer (integer)
7
8     IN      datatype    datatype of each buffer element (handle)
9
10    OUT     status      status object (Status)

```

```

10 int MPI_File_write(MPI_File fh, void *buf, int count, MPI_Datatype datatype,
11                  MPI_Status *status)

```

```

12 MPI_FILE_WRITE(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
13     <type> BUF(*)
14     INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR

```

```

16 void MPI::File::Write(const void* buf, int count,
17                     const MPI::Datatype& datatype, MPI::Status& status)

```

```

18 void MPI::File::Write(const void* buf, int count,
19                     const MPI::Datatype& datatype)

```

21 MPI_FILE_WRITE writes a file using the individual file pointer.

```

24 MPI_FILE_WRITE_ALL(fh, buf, count, datatype, status)

```

```

25     INOUT   fh           file handle (handle)
26
27     IN      buf         initial address of buffer (choice)
28
29     IN      count       number of elements in buffer (integer)
30
31     IN      datatype    datatype of each buffer element (handle)
32
33     OUT     status      status object (Status)

```

```

32 int MPI_File_write_all(MPI_File fh, void *buf, int count,
33                      MPI_Datatype datatype, MPI_Status *status)

```

```

35 MPI_FILE_WRITE_ALL(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
36     <type> BUF(*)
37     INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR

```

```

39 void MPI::File::Write_all(const void* buf, int count,
40                          const MPI::Datatype& datatype, MPI::Status& status)

```

```

41 void MPI::File::Write_all(const void* buf, int count,
42                          const MPI::Datatype& datatype)

```

44 MPI_FILE_WRITE_ALL is a collective version of the blocking MPI_FILE_WRITE interface.

45
46
47
48

| | | | |
|---|----------|--|---|
| MPI_FILE_IREAD(fh, buf, count, datatype, request) | | | 1 |
| INOUT | fh | file handle (handle) | 2 |
| OUT | buf | initial address of buffer (choice) | 3 |
| IN | count | number of elements in buffer (integer) | 4 |
| IN | datatype | datatype of each buffer element (handle) | 5 |
| OUT | request | request object (handle) | 6 |
| | | | 7 |
| | | | 8 |
| | | | 9 |

```
int MPI_File_iread(MPI_File fh, void *buf, int count, MPI_Datatype datatype,
                  MPI_Request *request)
```

```
MPI_FILE_IREAD(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR)
<type> BUF(*)
INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
```

```
MPI::Request MPI::File::Iread(void* buf, int count,
                               const MPI::Datatype& datatype)
```

MPI_FILE_IREAD is a nonblocking version of the MPI_FILE_READ interface.

Example 12.3 The following Fortran code fragment illustrates file pointer update semantics:

```
! Read the first twenty real words in a file into two local
! buffers. Note that when the first MPI_FILE_IREAD returns,
! the file pointer has been updated to point to the
! eleventh real word in the file.

integer  bufsize, req1, req2
integer, dimension(MPI_STATUS_SIZE) :: status1, status2
parameter (bufsize=10)
real     buf1(bufsize), buf2(bufsize)

call MPI_FILE_OPEN( MPI_COMM_WORLD, 'myoldfile', &
                   MPI_MODE_RDONLY, MPI_INFO_NULL, myfh, ierr )
call MPI_FILE_SET_VIEW( myfh, 0, MPI_REAL, MPI_REAL, 'native', &
                       MPI_INFO_NULL, ierr )
call MPI_FILE_IREAD( myfh, buf1, bufsize, MPI_REAL, &
                   req1, ierr )
call MPI_FILE_IREAD( myfh, buf2, bufsize, MPI_REAL, &
                   req2, ierr )

call MPI_WAIT( req1, status1, ierr )
call MPI_WAIT( req2, status2, ierr )

call MPI_FILE_CLOSE( myfh, ierr )
```

1 MPI_FILE_IWRITE(fh, buf, count, datatype, request)

2 INOUT fh file handle (handle)
 3
 4 IN buf initial address of buffer (choice)
 5
 6 IN count number of elements in buffer (integer)
 7
 8 IN datatype datatype of each buffer element (handle)
 9
 10 OUT request request object (handle)

10 int MPI_File_irewrite(MPI_File fh, void *buf, int count,
 11 MPI_Datatype datatype, MPI_Request *request)

12 MPI_FILE_IWRITE(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR)
 13 <type> BUF(*)
 14 INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR

16 MPI::Request MPI::File::Iwrite(const void* buf, int count,
 17 const MPI::Datatype& datatype)

18 MPI_FILE_IWRITE is a nonblocking version of the MPI_FILE_WRITE interface.

21 MPI_FILE_SEEK(fh, offset, whence)

22 INOUT fh file handle (handle)
 23
 24 IN offset file offset (integer)
 25
 26 IN whence update mode (state)

27 int MPI_File_seek(MPI_File fh, MPI_Offset offset, int whence)

29 MPI_FILE_SEEK(FH, OFFSET, WHENCE, IERROR)
 30 INTEGER FH, WHENCE, IERROR
 31 INTEGER(KIND=MPI_OFFSET_KIND) OFFSET

32 void MPI::File::Seek(MPI::Offset offset, int whence)

34 MPI_FILE_SEEK updates the individual file pointer according to whence, which has the
 35 following possible values:

- 36 • MPI_SEEK_SET: the pointer is set to offset
- 37
- 38 • MPI_SEEK_CUR: the pointer is set to the current pointer position plus offset
- 39
- 40 • MPI_SEEK_END: the pointer is set to the end of file plus offset

41 The offset can be negative, which allows seeking backwards. It is erroneous to seek to
 42 a negative position in the view.

44
 45
 46
 47
 48

MPI_FILE_GET_POSITION(fh, offset) 1

IN fh file handle (handle) 2

OUT offset offset of individual pointer (integer) 3

int MPI_File_get_position(MPI_File fh, MPI_Offset *offset) 4

MPI_FILE_GET_POSITION(FH, OFFSET, IERROR) 5

INTEGER FH, IERROR 6

INTEGER(KIND=MPI_OFFSET_KIND) OFFSET 7

MPI::Offset MPI::File::Get_position() const 8

MPI_FILE_GET_POSITION returns, in `offset`, the current position of the individual file pointer in etype units relative to the current view. 9

Advice to users. The `offset` can be used in a future call to `MPI_FILE_SEEK` using `whence = MPI_SEEK_SET` to return to the current position. To set the displacement to the current file pointer position, first convert `offset` into an absolute byte position using `MPI_FILE_GET_BYTE_OFFSET`, then call `MPI_FILE_SET_VIEW` with the resulting displacement. (*End of advice to users.*) 10

MPI_FILE_GET_BYTE_OFFSET(fh, offset, disp) 11

IN fh file handle (handle) 12

IN offset offset (integer) 13

OUT disp absolute byte position of offset (integer) 14

int MPI_File_get_byte_offset(MPI_File fh, MPI_Offset offset, MPI_Offset *disp) 15

MPI_FILE_GET_BYTE_OFFSET(FH, OFFSET, DISP, IERROR) 16

INTEGER FH, IERROR 17

INTEGER(KIND=MPI_OFFSET_KIND) OFFSET, DISP 18

MPI::Offset MPI::File::Get_byte_offset(const MPI::Offset disp) const 19

MPI_FILE_GET_BYTE_OFFSET converts a view-relative offset into an absolute byte position. The absolute byte position (from the beginning of the file) of `offset` relative to the current view of `fh` is returned in `disp`. 20

12.4.4 Data Access with Shared File Pointers 21

MPI maintains exactly one shared file pointer per collective `MPI_FILE_OPEN` (shared among processes in the communicator group). The current value of this pointer implicitly specifies the offset in the data access routines described in this section. These routines only use and update the shared file pointer maintained by MPI. The individual file pointers are not used nor updated. 22

The shared file pointer routines have the same semantics as the data access with explicit offset routines described in Section 12.4.2, page 392, with the following modifications: 23

- 1 • the offset is defined to be the current value of the MPI-maintained shared file pointer,
- 2
- 3 • the effect of multiple calls to shared file pointer routines is defined to behave as if the
- 4 calls were serialized, and
- 5
- 6 • the use of shared file pointer routines is erroneous unless all processes use the same
- 7 file view.

8 For the noncollective shared file pointer routines, the serialization ordering is not determin-
9 istic. The user needs to use other synchronization means to enforce a specific order.

10 After a shared file pointer operation is initiated, the shared file pointer is updated to
11 point to the next etype after the last one that will be accessed. The file pointer is updated
12 relative to the current view of the file.

13 Noncollective Operations

14
15
16
17 MPI_FILE_READ_SHARED(fh, buf, count, datatype, status)

| | | | |
|----|-------|----------|--|
| 18 | INOUT | fh | file handle (handle) |
| 19 | OUT | buf | initial address of buffer (choice) |
| 20 | | | |
| 21 | IN | count | number of elements in buffer (integer) |
| 22 | IN | datatype | datatype of each buffer element (handle) |
| 23 | OUT | status | status object (Status) |
| 24 | | | |

25
26 int MPI_File_read_shared(MPI_File fh, void *buf, int count,
27 MPI_Datatype datatype, MPI_Status *status)

28 MPI_FILE_READ_SHARED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
29 <type> BUF(*)
30 INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR

31
32 void MPI::File::Read_shared(void* buf, int count,
33 const MPI::Datatype& datatype, MPI::Status& status)

34 void MPI::File::Read_shared(void* buf, int count,
35 const MPI::Datatype& datatype)

36
37 MPI_FILE_READ_SHARED reads a file using the shared file pointer.

38
39 MPI_FILE_WRITE_SHARED(fh, buf, count, datatype, status)

| | | | |
|----|-------|----------|--|
| 40 | INOUT | fh | file handle (handle) |
| 41 | IN | buf | initial address of buffer (choice) |
| 42 | | | |
| 43 | IN | count | number of elements in buffer (integer) |
| 44 | IN | datatype | datatype of each buffer element (handle) |
| 45 | OUT | status | status object (Status) |
| 46 | | | |

47
48 int MPI_File_write_shared(MPI_File fh, void *buf, int count,


```

        MPI_Datatype datatype, MPI_Status *status)
1
MPI_FILE_WRITE_SHARED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
2
    <type> BUF(*)
3
    INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
4
void MPI::File::Write_shared(const void* buf, int count,
5
    const MPI::Datatype& datatype, MPI::Status& status)
6
void MPI::File::Write_shared(const void* buf, int count,
7
    const MPI::Datatype& datatype)
8
    MPI_FILE_WRITE_SHARED writes a file using the shared file pointer.
9
10
11
12
13
MPI_FILE_IREAD_SHARED(fh, buf, count, datatype, request)
14
    INOUT   fh                file handle (handle)
15
    OUT     buf                initial address of buffer (choice)
16
    IN      count              number of elements in buffer (integer)
17
    IN      datatype           datatype of each buffer element (handle)
18
    OUT     request            request object (handle)
19
20
21
22
int MPI_File_iread_shared(MPI_File fh, void *buf, int count,
23
    MPI_Datatype datatype, MPI_Request *request)
24
MPI_FILE_IREAD_SHARED(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR)
25
    <type> BUF(*)
26
    INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
27
MPI::Request MPI::File::Iread_shared(void* buf, int count,
28
    const MPI::Datatype& datatype)
29
    MPI_FILE_IREAD_SHARED is a nonblocking version of the MPI_FILE_READ_SHARED
30
    interface.
31
32
33
34
MPI_FILE_IWRITE_SHARED(fh, buf, count, datatype, request)
35
    INOUT   fh                file handle (handle)
36
    IN      buf                initial address of buffer (choice)
37
    IN      count              number of elements in buffer (integer)
38
    IN      datatype           datatype of each buffer element (handle)
39
    OUT     request            request object (handle)
40
41
42
43
int MPI_File_ fwrite_shared(MPI_File fh, void *buf, int count,
44
    MPI_Datatype datatype, MPI_Request *request)
45
MPI_FILE_IWRITE_SHARED(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR)
46
    <type> BUF(*)
47
48

```

```

1     INTEGER fh, count, datatype, request, ierror
2
3     MPI::Request MPI::File::Iwrite_shared(const void* buf, int count,
4         const MPI::Datatype& datatype)

```

MPI_FILE_IWRITE_SHARED is a nonblocking version of the MPI_FILE_WRITE_SHARED interface.

Collective Operations

The semantics of a collective access using a shared file pointer is that the accesses to the file will be in the order determined by the ranks of the processes within the group. For each process, the location in the file at which data is accessed is the position at which the shared file pointer would be after all processes whose ranks within the group less than that of this process had accessed their data. In addition, in order to prevent subsequent shared offset accesses by the same processes from interfering with this collective access, the call might return only after all the processes within the group have initiated their accesses. When the call returns, the shared file pointer points to the next etype accessible, according to the file view used by all processes, after the last etype requested.

Advice to users. There may be some programs in which all processes in the group need to access the file using the shared file pointer, but the program may not *require* that data be accessed in order of process rank. In such programs, using the shared ordered routines (e.g., MPI_FILE_WRITE_ORDERED rather than MPI_FILE_WRITE_SHARED) may enable an implementation to optimize access, improving performance. (*End of advice to users.*)

Advice to implementors. Accesses to the data requested by all processes do not have to be serialized. Once all processes have issued their requests, locations within the file for all accesses can be computed, and accesses can proceed independently from each other, possibly in parallel. (*End of advice to implementors.*)

```

32 MPI_FILE_READ_ORDERED(fh, buf, count, datatype, status)

```

| | | |
|-------|----------|--|
| INOUT | fh | file handle (handle) |
| OUT | buf | initial address of buffer (choice) |
| IN | count | number of elements in buffer (integer) |
| IN | datatype | datatype of each buffer element (handle) |
| OUT | status | status object (Status) |

```

40
41 int MPI_File_read_ordered(MPI_File fh, void *buf, int count,
42     MPI_Datatype datatype, MPI_Status *status)
43
44 MPI_FILE_READ_ORDERED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
45     <type> BUF(*)
46     INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
47
48 void MPI::File::Read_ordered(void* buf, int count,
49     const MPI::Datatype& datatype, MPI::Status& status)

```

```
void MPI::File::Read_ordered(void* buf, int count,
                             const MPI::Datatype& datatype)
```

MPI_FILE_READ_ORDERED is a collective version of the MPI_FILE_READ_SHARED interface.

```
MPI_FILE_WRITE_ORDERED(fh, buf, count, datatype, status)
```

| | | |
|-------|----------|--|
| INOUT | fh | file handle (handle) |
| IN | buf | initial address of buffer (choice) |
| IN | count | number of elements in buffer (integer) |
| IN | datatype | datatype of each buffer element (handle) |
| OUT | status | status object (Status) |

```
int MPI_File_write_ordered(MPI_File fh, void *buf, int count,
                           MPI_Datatype datatype, MPI_Status *status)
```

```
MPI_FILE_WRITE_ORDERED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
    <type> BUF(*)
    INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
```

```
void MPI::File::Write_ordered(const void* buf, int count,
                              const MPI::Datatype& datatype, MPI::Status& status)
```

```
void MPI::File::Write_ordered(const void* buf, int count,
                              const MPI::Datatype& datatype)
```

MPI_FILE_WRITE_ORDERED is a collective version of the MPI_FILE_WRITE_SHARED interface.

Seek

If MPI_MODE_SEQUENTIAL mode was specified when the file was opened, it is erroneous to call the following two routines (MPI_FILE_SEEK_SHARED and MPI_FILE_GET_POSITION_SHARED).

```
MPI_FILE_SEEK_SHARED(fh, offset, whence)
```

| | | |
|-------|--------|-----------------------|
| INOUT | fh | file handle (handle) |
| IN | offset | file offset (integer) |
| IN | whence | update mode (state) |

```
int MPI_File_seek_shared(MPI_File fh, MPI_Offset offset, int whence)
```

```
MPI_FILE_SEEK_SHARED(FH, OFFSET, WHENCE, IERROR)
    INTEGER FH, WHENCE, IERROR
    INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
```

```
void MPI::File::Seek_shared(MPI::Offset offset, int whence)
```

1 `MPI_FILE_SEEK_SHARED` updates the shared file pointer according to `whence`, which
 2 has the following possible values:

- 3 • `MPI_SEEK_SET`: the pointer is set to `offset`
- 4 • `MPI_SEEK_CUR`: the pointer is set to the current pointer position plus `offset`
- 5 • `MPI_SEEK_END`: the pointer is set to the end of file plus `offset`

6 `MPI_FILE_SEEK_SHARED` is collective; all the processes in the communicator group
 7 associated with the file handle `fh` must call `MPI_FILE_SEEK_SHARED` with the same values
 8 for `offset` and `whence`.

9 The `offset` can be negative, which allows seeking backwards. It is erroneous to seek to
 10 a negative position in the view.

11 `MPI_FILE_GET_POSITION_SHARED(fh, offset)`

| | | | |
|----|-----|---------------------|------------------------------------|
| 12 | IN | <code>fh</code> | file handle (handle) |
| 13 | OUT | <code>offset</code> | offset of shared pointer (integer) |

14 `int MPI_File_get_position_shared(MPI_File fh, MPI_Offset *offset)`

15 `MPI_FILE_GET_POSITION_SHARED(FH, OFFSET, IERROR)`

16 `INTEGER FH, IERROR`

17 `INTEGER(KIND=MPI_OFFSET_KIND) OFFSET`

18 `MPI::Offset MPI::File::Get_position_shared() const`

19 `MPI_FILE_GET_POSITION_SHARED` returns, in `offset`, the current position of the shared
 20 file pointer in `etype` units relative to the current view.

21 *Advice to users.* The `offset` can be used in a future call to `MPI_FILE_SEEK_SHARED`
 22 using `whence = MPI_SEEK_SET` to return to the current position. To set the displace-
 23 ment to the current file pointer position, first convert `offset` into an absolute byte po-
 24 sition using `MPI_FILE_GET_BYTE_OFFSET`, then call `MPI_FILE_SET_VIEW` with the
 25 resulting displacement. (*End of advice to users.*)

26 12.4.5 Split Collective Data Access Routines

27 MPI provides a restricted form of “nonblocking collective” I/O operations for all data ac-
 28 cesses using split collective data access routines. These routines are referred to as “split”
 29 collective routines because a single collective operation is split in two: a begin routine and
 30 an end routine. The begin routine begins the operation, much like a nonblocking data access
 31 (e.g., `MPI_FILE_IREAD`). The end routine completes the operation, much like the matching
 32 test or wait (e.g., `MPI_WAIT`). As with nonblocking data access operations, the user must
 33 not use the buffer passed to a begin routine while the routine is outstanding; the operation
 34 must be completed with an end routine before it is safe to free buffers, etc.

35 Split collective data access operations on a file handle `fh` are subject to the semantic
 36 rules given below.

- On any MPI process, each file handle may have at most one active split collective operation at any time.
- Begin calls are collective over the group of processes that participated in the collective open and follow the ordering rules for collective calls.
- End calls are collective over the group of processes that participated in the collective open and follow the ordering rules for collective calls. Each end call matches the preceding begin call for the same collective operation. When an “end” call is made, exactly one unmatched “begin” call for the same operation must precede it.
- An implementation is free to implement any split collective data access routine using the corresponding blocking collective routine when either the begin call (e.g., `MPI_FILE_READ_ALL_BEGIN`) or the end call (e.g., `MPI_FILE_READ_ALL_END`) is issued. The begin and end calls are provided to allow the user and MPI implementation to optimize the collective operation.
- Split collective operations do not match the corresponding regular collective operation. For example, in a single collective read operation, an `MPI_FILE_READ_ALL` on one process does not match an `MPI_FILE_READ_ALL_BEGIN/MPI_FILE_READ_ALL_END` pair on another process.
- Split collective routines must specify a buffer in both the begin and end routines. By specifying the buffer that receives data in the end routine, we can avoid many (though not all) of the problems described in “A Problem with Register Optimization,” Section 13.2.2, page 454.
- No collective I/O operations are permitted on a file handle concurrently with a split collective access on that file handle (i.e., between the begin and end of the access). That is

```

MPI_File_read_all_begin(fh, ...);
...
MPI_File_read_all(fh, ...);
...
MPI_File_read_all_end(fh, ...);

```

is erroneous.

- In a multithreaded implementation, any split collective begin and end operation called by a process must be called from the same thread. This restriction is made to simplify the implementation in the multithreaded case. (Note that we have already disallowed having two threads begin a split collective operation on the same file handle since only one split collective operation can be active on a file handle at any time.)

The arguments for these routines have the same meaning as for the equivalent collective versions (e.g., the argument definitions for `MPI_FILE_READ_ALL_BEGIN` and `MPI_FILE_READ_ALL_END` are equivalent to the arguments for `MPI_FILE_READ_ALL`). The begin routine (e.g., `MPI_FILE_READ_ALL_BEGIN`) begins a split collective operation that, when completed with the matching end routine (i.e., `MPI_FILE_READ_ALL_END`) produces the result as defined for the equivalent collective routine (i.e., `MPI_FILE_READ_ALL`).

1 For the purpose of consistency semantics (Section 12.6.1, page 422), a matched pair
 2 of split collective data access operations (e.g., `MPI_FILE_READ_ALL_BEGIN` and
 3 `MPI_FILE_READ_ALL_END`) compose a single data access.

4
 5
 6 `MPI_FILE_READ_AT_ALL_BEGIN(fh, offset, buf, count, datatype)`

| | | | |
|----|-----|----------|--|
| 7 | IN | fh | file handle (handle) |
| 8 | IN | offset | file offset (integer) |
| 9 | OUT | buf | initial address of buffer (choice) |
| 10 | IN | count | number of elements in buffer (integer) |
| 11 | IN | datatype | datatype of each buffer element (handle) |

12
 13
 14
 15 `int MPI_File_read_at_all_begin(MPI_File fh, MPI_Offset offset, void *buf,`
 16 `int count, MPI_Datatype datatype)`

17 `MPI_FILE_READ_AT_ALL_BEGIN(FH, OFFSET, BUF, COUNT, DATATYPE, IERROR)`

18 `<type> BUF(*)`
 19 `INTEGER FH, COUNT, DATATYPE, IERROR`
 20 `INTEGER(KIND=MPI_OFFSET_KIND) OFFSET`

21
 22 `void MPI::File::Read_at_all_begin(MPI::Offset offset, void* buf, int count,`
 23 `const MPI::Datatype& datatype)`

24
 25
 26 `MPI_FILE_READ_AT_ALL_END(fh, buf, status)`

| | | | |
|----|-----|--------|------------------------------------|
| 27 | IN | fh | file handle (handle) |
| 28 | OUT | buf | initial address of buffer (choice) |
| 29 | OUT | status | status object (Status) |

30
 31
 32 `int MPI_File_read_at_all_end(MPI_File fh, void *buf, MPI_Status *status)`

33
 34 `MPI_FILE_READ_AT_ALL_END(FH, BUF, STATUS, IERROR)`

35 `<type> BUF(*)`
 36 `INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR`

37 `void MPI::File::Read_at_all_end(void* buf, MPI::Status& status)`

38
 39 `void MPI::File::Read_at_all_end(void* buf)`

40
 41
 42
 43
 44
 45
 46
 47
 48

```

MPI_FILE_WRITE_AT_ALL_BEGIN(fh, offset, buf, count, datatype) 1
    INOUT fh file handle (handle) 2
    IN offset file offset (integer) 3
    IN buf initial address of buffer (choice) 4
    IN count number of elements in buffer (integer) 5
    IN datatype datatype of each buffer element (handle) 6
    7
int MPI_File_write_at_all_begin(MPI_File fh, MPI_Offset offset, void *buf, 8
    int count, MPI_Datatype datatype) 9
MPI_FILE_WRITE_AT_ALL_BEGIN(FH, OFFSET, BUF, COUNT, DATATYPE, IERROR) 10
    <type> BUF(*) 11
    INTEGER FH, COUNT, DATATYPE, IERROR 12
    INTEGER(KIND=MPI_OFFSET_KIND) OFFSET 13
void MPI::File::Write_at_all_begin(MPI::Offset offset, const void* buf, 14
    int count, const MPI::Datatype& datatype) 15
    16
MPI_FILE_WRITE_AT_ALL_END(fh, buf, status) 17
    INOUT fh file handle (handle) 18
    IN buf initial address of buffer (choice) 19
    OUT status status object (Status) 20
    21
int MPI_File_write_at_all_end(MPI_File fh, void *buf, MPI_Status *status) 22
MPI_FILE_WRITE_AT_ALL_END(FH, BUF, STATUS, IERROR) 23
    <type> BUF(*) 24
    INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR 25
void MPI::File::Write_at_all_end(const void* buf, MPI::Status& status) 26
void MPI::File::Write_at_all_end(const void* buf) 27
    28
MPI_FILE_READ_ALL_BEGIN(fh, buf, count, datatype) 29
    INOUT fh file handle (handle) 30
    OUT buf initial address of buffer (choice) 31
    IN count number of elements in buffer (integer) 32
    IN datatype datatype of each buffer element (handle) 33
    34
int MPI_File_read_all_begin(MPI_File fh, void *buf, int count, 35
    MPI_Datatype datatype) 36
MPI_FILE_READ_ALL_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR) 37
    38
    39
    40
    41
    42
    43
    44
    45
    46
    47
    48

```

```

1     <type> BUF(*)
2     INTEGER FH, COUNT, DATATYPE, IERROR
3
4     void MPI::File::Read_all_begin(void* buf, int count,
5         const MPI::Datatype& datatype)
6
7
8     MPI_FILE_READ_ALL_END(fh, buf, status)
9
10    INOUT   fh           file handle (handle)
11    OUT     buf          initial address of buffer (choice)
12    OUT     status       status object (Status)
13
14    int MPI_File_read_all_end(MPI_File fh, void *buf, MPI_Status *status)
15
16    MPI_FILE_READ_ALL_END(FH, BUF, STATUS, IERROR)
17    <type> BUF(*)
18    INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
19
20    void MPI::File::Read_all_end(void* buf, MPI::Status& status)
21
22    void MPI::File::Read_all_end(void* buf)
23
24
25    MPI_FILE_WRITE_ALL_BEGIN(fh, buf, count, datatype)
26
27    INOUT   fh           file handle (handle)
28    IN      buf          initial address of buffer (choice)
29    IN      count        number of elements in buffer (integer)
30    IN      datatype     datatype of each buffer element (handle)
31
32    int MPI_File_write_all_begin(MPI_File fh, void *buf, int count,
33        MPI_Datatype datatype)
34
35    MPI_FILE_WRITE_ALL_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR)
36    <type> BUF(*)
37    INTEGER FH, COUNT, DATATYPE, IERROR
38
39    void MPI::File::Write_all_begin(const void* buf, int count,
40        const MPI::Datatype& datatype)
41
42
43    MPI_FILE_WRITE_ALL_END(fh, buf, status)
44
45    INOUT   fh           file handle (handle)
46    IN      buf          initial address of buffer (choice)
47    OUT     status       status object (Status)
48
49    int MPI_File_write_all_end(MPI_File fh, void *buf, MPI_Status *status)

```



```

MPI_FILE_WRITE_ALL_END(FH, BUF, STATUS, IERROR)
    <type> BUF(*)
    INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
void MPI::File::Write_all_end(const void* buf, MPI::Status& status)
void MPI::File::Write_all_end(const void* buf)

MPI_FILE_READ_ORDERED_BEGIN(fh, buf, count, datatype)
    INOUT  fh                file handle (handle)
    OUT    buf                initial address of buffer (choice)
    IN     count              number of elements in buffer (integer)
    IN     datatype           datatype of each buffer element (handle)
int MPI_File_read_ordered_begin(MPI_File fh, void *buf, int count,
    MPI_Datatype datatype)
MPI_FILE_READ_ORDERED_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR)
    <type> BUF(*)
    INTEGER FH, COUNT, DATATYPE, IERROR
void MPI::File::Read_ordered_begin(void* buf, int count,
    const MPI::Datatype& datatype)

MPI_FILE_READ_ORDERED_END(fh, buf, status)
    INOUT  fh                file handle (handle)
    OUT    buf                initial address of buffer (choice)
    OUT    status              status object (Status)
int MPI_File_read_ordered_end(MPI_File fh, void *buf, MPI_Status *status)
MPI_FILE_READ_ORDERED_END(FH, BUF, STATUS, IERROR)
    <type> BUF(*)
    INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
void MPI::File::Read_ordered_end(void* buf, MPI::Status& status)
void MPI::File::Read_ordered_end(void* buf)

```

```

1 MPI_FILE_WRITE_ORDERED_BEGIN(fh, buf, count, datatype)
2     INOUT    fh                file handle (handle)
3
4     IN      buf                initial address of buffer (choice)
5
6     IN      count              number of elements in buffer (integer)
7
8     IN      datatype           datatype of each buffer element (handle)
9
10 int MPI_File_write_ordered_begin(MPI_File fh, void *buf, int count,
11                                 MPI_Datatype datatype)
12
13 MPI_FILE_WRITE_ORDERED_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR)
14     <type> BUF(*)
15     INTEGER FH, COUNT, DATATYPE, IERROR
16
17 void MPI::File::Write_ordered_begin(const void* buf, int count,
18                                    const MPI::Datatype& datatype)
19
20 MPI_FILE_WRITE_ORDERED_END(fh, buf, status)
21     INOUT    fh                file handle (handle)
22
23     IN      buf                initial address of buffer (choice)
24
25     OUT     status              status object (Status)
26
27 int MPI_File_write_ordered_end(MPI_File fh, void *buf, MPI_Status *status)
28
29 MPI_FILE_WRITE_ORDERED_END(FH, BUF, STATUS, IERROR)
30     <type> BUF(*)
31     INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
32
33 void MPI::File::Write_ordered_end(const void* buf, MPI::Status& status)
34
35 void MPI::File::Write_ordered_end(const void* buf)

```

12.5 File Interoperability

At the most basic level, file interoperability is the ability to read the information previously written to a file—not just the bits of data, but the actual information the bits represent. MPI guarantees full interoperability within a single MPI environment, and supports increased interoperability outside that environment through the external data representation (Section 12.5.2, page 416) as well as the data conversion functions (Section 12.5.3, page 418).

Interoperability within a single MPI environment (which could be considered “operability”) ensures that file data written by one MPI process can be read by any other MPI process, subject to the consistency constraints (see Section 12.6.1, page 422), provided that it would have been possible to start the two processes simultaneously and have them reside in a single MPI_COMM_WORLD. Furthermore, both processes must see the same data values at every absolute byte offset in the file for which data was written.

This single environment file interoperability implies that file data is accessible regardless of the number of processes.

There are three aspects to file interoperability:

- transferring the bits,
- converting between different file structures, and
- converting between different machine representations.

The first two aspects of file interoperability are beyond the scope of this standard, as both are highly machine dependent. However, transferring the bits of a file into and out of the MPI environment (e.g., by writing a file to tape) is required to be supported by all MPI implementations. In particular, an implementation must specify how familiar operations similar to POSIX `cp`, `rm`, and `mv` can be performed on the file. Furthermore, it is expected that the facility provided maintains the correspondence between absolute byte offsets (e.g., after possible file structure conversion, the data bits at byte offset 102 in the MPI environment are at byte offset 102 outside the MPI environment). As an example, a simple off-line conversion utility that transfers and converts files between the native file system and the MPI environment would suffice, provided it maintained the offset coherence mentioned above. In a high quality implementation of MPI, users will be able to manipulate MPI files using the same or similar tools that the native file system offers for manipulating its files.

The remaining aspect of file interoperability, converting between different machine representations, is supported by the typing information specified in the `etype` and `filetype`. This facility allows the information in files to be shared between any two applications, regardless of whether they use MPI, and regardless of the machine architectures on which they run.

MPI supports multiple data representations: “native,” “internal,” and “external32.” An implementation may support additional data representations. MPI also supports user-defined data representations (see Section 12.5.3, page 418). The native and internal data representations are implementation dependent, while the external32 representation is common to all MPI implementations and facilitates file interoperability. The data representation is specified in the `datarep` argument to `MPI_FILE_SET_VIEW`.

Advice to users. MPI is not guaranteed to retain knowledge of what data representation was used when a file is written. Therefore, to correctly retrieve file data, an MPI application is responsible for specifying the same data representation as was used to create the file. (*End of advice to users.*)

“native” Data in this representation is stored in a file exactly as it is in memory. The advantage of this data representation is that data precision and I/O performance are not lost in type conversions with a purely homogeneous environment. The disadvantage is the loss of transparent interoperability within a heterogeneous MPI environment.

Advice to users. This data representation should only be used in a homogeneous MPI environment, or when the MPI application is capable of performing the data type conversions itself. (*End of advice to users.*)

1 *Advice to implementors.* When implementing read and write operations on
 2 top of MPI message passing, the message data should be typed as MPI_BYTE
 3 to ensure that the message routines do not perform any type conversions on the
 4 data. (*End of advice to implementors.*)

5
 6 **“internal”** This data representation can be used for I/O operations in a homogeneous
 7 or heterogeneous environment; the implementation will perform type conversions if
 8 necessary. The implementation is free to store data in any format of its choice, with
 9 the restriction that it will maintain constant extents for all predefined datatypes in any
 10 one file. The environment in which the resulting file can be reused is implementation-
 11 defined and must be documented by the implementation.

12
 13 *Rationale.* This data representation allows the implementation to perform I/O
 14 efficiently in a heterogeneous environment, though with implementation-defined
 15 restrictions on how the file can be reused. (*End of rationale.*)

16
 17 *Advice to implementors.* Since “external32” is a superset of the functionality
 18 provided by “internal,” an implementation may choose to implement “internal”
 19 as “external32.” (*End of advice to implementors.*)

20
 21 **“external32”** This data representation states that read and write operations convert all
 22 data from and to the “external32” representation defined in Section 12.5.2, page 416.
 23 The data conversion rules for communication also apply to these conversions (see
 24 Section 3.3.2, page 25-27, of the MPI-1 document). The data on the storage medium
 25 is always in this canonical representation, and the data in memory is always in the
 26 local process’s native representation.

27 This data representation has several advantages. First, all processes reading the file
 28 in a heterogeneous MPI environment will automatically have the data converted to
 29 their respective native representations. Second, the file can be exported from one MPI
 30 environment and imported into any other MPI environment with the guarantee that
 31 the second environment will be able to read all the data in the file.

32 The disadvantage of this data representation is that data precision and I/O perfor-
 33 mance may be lost in data type conversions.

34
 35 *Advice to implementors.* When implementing read and write operations on top
 36 of MPI message passing, the message data should be converted to and from the
 37 “external32” representation in the client, and sent as type MPI_BYTE. This will
 38 avoid possible double data type conversions and the associated further loss of
 39 precision and performance. (*End of advice to implementors.*)

40 41 12.5.1 Datatypes for File Interoperability

42
 43 If the file data representation is other than “native,” care must be taken in constructing
 44 etypes and filetypes. Any of the datatype constructor functions may be used; however,
 45 for those functions that accept displacements in bytes, the displacements must be specified
 46 in terms of their values in the file for the file data representation being used. MPI will
 47 interpret these byte displacements as is; no scaling will be done. The function
 48 MPI_FILE_GET_TYPE_EXTENT can be used to calculate the extents of datatypes in the file.

For etypes and filetypes that are portable datatypes (see Section 2.4, page 11), MPI will scale any displacements in the datatypes to match the file data representation. Datatypes passed as arguments to read/write routines specify the data layout in memory; therefore, they must always be constructed using displacements corresponding to displacements in memory.

Advice to users. One can logically think of the file as if it were stored in the memory of a file server. The `etype` and `filetype` are interpreted as if they were defined at this file server, by the same sequence of calls used to define them at the calling process. If the data representation is “native”, then this logical file server runs on the same architecture as the calling process, so that these types define the same data layout on the file as they would define in the memory of the calling process. If the `etype` and `filetype` are portable datatypes, then the data layout defined in the file is the same as would be defined in the calling process memory, up to a scaling factor. The routine `MPI_FILE_GET_FILE_EXTENT` can be used to calculate this scaling factor. Thus, two equivalent, portable datatypes will define the same data layout in the file, even in a heterogeneous environment with “internal”, “external32”, or user defined data representations. Otherwise, the `etype` and `filetype` must be constructed so that their `typemap` and `extent` are the same on any architecture. This can be achieved if they have an explicit upper bound and lower bound (defined either using `MPI_LB` and `MPI_UB` markers, or using `MPI_TYPE_CREATE_RESIZED`). This condition must also be fulfilled by any datatype that is used in the construction of the `etype` and `filetype`, if this datatype is replicated contiguously, either explicitly, by a call to `MPI_TYPE_CONTIGUOUS`, or implicitly, by a `blocklength` argument that is greater than one. If an `etype` or `filetype` is not portable, and has a `typemap` or `extent` that is architecture dependent, then the data layout specified by it on a file is implementation dependent.

File data representations other than “native” may be different from corresponding data representations in memory. Therefore, for these file data representations, it is important not to use hardwired byte offsets for file positioning, including the initial displacement that specifies the view. When a portable datatype (see Section 2.4, page 11) is used in a data access operation, any holes in the datatype are scaled to match the data representation. However, note that this technique only works when all the processes that created the file view build their etypes from the same predefined datatypes. For example, if one process uses an `etype` built from `MPI_INT` and another uses an `etype` built from `MPI_FLOAT`, the resulting views may be nonportable because the relative sizes of these types may differ from one data representation to another. (*End of advice to users.*)

`MPI_FILE_GET_TYPE_EXTENT(fh, datatype, extent)`

| | | |
|-----|-----------------------|---------------------------|
| IN | <code>fh</code> | file handle (handle) |
| IN | <code>datatype</code> | datatype (handle) |
| OUT | <code>extent</code> | datatype extent (integer) |

```
int MPI_File_get_type_extent(MPI_File fh, MPI_Datatype datatype,
                             MPI_Aint *extent)
```

```

1 MPI_FILE_GET_TYPE_EXTENT(FH, DATATYPE, EXTENT, IERROR)
2     INTEGER FH, DATATYPE, IERROR
3     INTEGER(KIND=MPI_ADDRESS_KIND) EXTENT

```

```

4 MPI::Aint MPI::File::Get_type_extent(const MPI::Datatype& datatype) const

```

Returns the extent of datatype in the file fh. This extent will be the same for all processes accessing the file fh. If the current view uses a user-defined data representation (see Section 12.5.3, page 418), MPI uses the `dtype_file_extent_fn` callback to calculate the extent.

Advice to implementors. In the case of user-defined data representations, the extent of a derived datatype can be calculated by first determining the extents of the predefined datatypes in this derived datatype using `dtype_file_extent_fn` (see Section 12.5.3, page 418). (*End of advice to implementors.*)

12.5.2 External Data Representation: “external32”

All MPI implementations are required to support the data representation defined in this section. Support of optional datatypes (e.g., `MPI_INTEGER2`) is not required.

All floating point values are in big-endian IEEE format [29] of the appropriate size. Floating point values are represented by one of three IEEE formats. These are the IEEE “Single,” “Double,” and “Double Extended” formats, requiring 4, 8 and 16 bytes of storage, respectively. For the IEEE “Double Extended” formats, MPI specifies a Format Width of 16 bytes, with 15 exponent bits, `bias = +16383`, 112 fraction bits, and an encoding analogous to the “Double” format. All integral values are in two’s complement big-endian format. Big-endian means most significant byte at lowest address byte. For Fortran `LOGICAL` and C++ `bool`, 0 implies false and nonzero implies true. Fortran `COMPLEX` and `DOUBLE COMPLEX` are represented by a pair of floating point format values for the real and imaginary components. Characters are in ISO 8859-1 format [30]. Wide characters (of type `MPI_WCHAR`) are in Unicode format [50].

All signed numerals (e.g., `MPLINT`, `MPI_REAL`) have the sign bit at the most significant bit. `MPI_COMPLEX` and `MPI_DOUBLE_COMPLEX` have the sign bit of the real and imaginary parts at the most significant bit of each part.

According to IEEE specifications [29], the “NaN” (not a number) is system dependent. It should not be interpreted within MPI as anything other than “NaN.”

Advice to implementors. The MPI treatment of “NaN” is similar to the approach used in XDR (see <ftp://ds.internic.net/rfc/rfc1832.txt>). (*End of advice to implementors.*)

All data is byte aligned, regardless of type. All data items are stored contiguously in the file (if the file view is contiguous).

Advice to implementors. All bytes of `LOGICAL` and `bool` must be checked to determine the value. (*End of advice to implementors.*)

Advice to users. The type `MPI_PACKED` is treated as bytes and is not converted. The user should be aware that `MPI_PACK` has the option of placing a header in the beginning of the pack buffer. (*End of advice to users.*)

| Type | Length | |
|------------------------|--------|----|
| ----- | ----- | |
| MPI_PACKED | 1 | 1 |
| MPI_BYTE | 1 | 2 |
| MPI_CHAR | 1 | 3 |
| MPI_UNSIGNED_CHAR | 1 | 4 |
| MPI_SIGNED_CHAR | 1 | 5 |
| MPI_WCHAR | 2 | 6 |
| MPI_SHORT | 2 | 7 |
| MPI_UNSIGNED_SHORT | 2 | 8 |
| MPI_INT | 4 | 9 |
| MPI_UNSIGNED | 4 | 10 |
| MPI_LONG | 4 | 11 |
| MPI_UNSIGNED_LONG | 4 | 12 |
| MPI_FLOAT | 4 | 13 |
| MPI_DOUBLE | 8 | 14 |
| MPI_LONG_DOUBLE | 16 | 15 |
| MPI_CHARACTER | 1 | 16 |
| MPI_LOGICAL | 4 | 17 |
| MPI_INTEGER | 4 | 18 |
| MPI_REAL | 4 | 19 |
| MPI_DOUBLE_PRECISION | 8 | 20 |
| MPI_COMPLEX | 2*4 | 21 |
| MPI_DOUBLE_COMPLEX | 2*8 | 22 |
| Optional Type | Length | 23 |
| ----- | ----- | 24 |
| MPI_INTEGER1 | 1 | 25 |
| MPI_INTEGER2 | 2 | 26 |
| MPI_INTEGER4 | 4 | 27 |
| MPI_INTEGER8 | 8 | 28 |
| MPI_LONG_LONG_INT | 8 | 29 |
| MPI_UNSIGNED_LONG_LONG | 8 | 30 |
| MPI_REAL4 | 4 | 31 |
| MPI_REAL8 | 8 | 32 |
| MPI_REAL16 | 16 | 33 |

Table 12.2: “external32”-sizes of predefined datatypes

1 The size of the predefined datatypes returned from `MPI_TYPE_CREATE_F90_REAL`,
 2 `MPI_TYPE_CREATE_F90_COMPLEX`, and `MPI_TYPE_CREATE_F90_INTEGER` are defined in
 3 Section 13.2.5, page 462.

4
 5 *Advice to implementors.* When converting a larger size integer to a smaller size
 6 integer, only the less significant bytes are moved. Care must be taken to preserve the
 7 sign bit value. This allows no conversion errors if the data range is within the range
 8 of the smaller size integer. (*End of advice to implementors.*)
 9

10 12.5.3 User-Defined Data Representations

11 There are two situations that cannot be handled by the required representations:
 12

- 13 1. a user wants to write a file in a representation unknown to the implementation, and
- 14 2. a user wants to read a file written in a representation unknown to the implementation.

15
 16 User-defined data representations allow the user to insert a third party converter into
 17 the I/O stream to do the data representation conversion.
 18

19
 20 `MPI_REGISTER_DATAREP`(`datarep`, `read_conversion_fn`, `write_conversion_fn`,
 21 `dtype_file_extent_fn`, `extra_state`)

| | | | |
|----|----|-----------------------------------|--|
| 22 | IN | <code>datarep</code> | data representation identifier (string) |
| 23 | IN | <code>read_conversion_fn</code> | function invoked to convert from file representation to native representation (function) |
| 24 | IN | <code>write_conversion_fn</code> | function invoked to convert from native representation to file representation (function) |
| 25 | IN | <code>dtype_file_extent_fn</code> | function invoked to get the extent of a datatype as represented in the file (function) |
| 26 | IN | <code>extra_state</code> | extra state |

27
 28
 29
 30
 31
 32
 33 `int MPI_Register_datarep`(`char *datarep`,
 34 `MPI_Datarep_conversion_function *read_conversion_fn`,
 35 `MPI_Datarep_conversion_function *write_conversion_fn`,
 36 `MPI_Datarep_extent_function *dtype_file_extent_fn`,
 37 `void *extra_state`)

38
 39 `MPI_REGISTER_DATAREP`(`DATAREP`, `READ_CONVERSION_FN`, `WRITE_CONVERSION_FN`,
 40 `DTYPE_FILE_EXTENT_FN`, `EXTRA_STATE`, `IERROR`)

41 `CHARACTER*(*) DATAREP`
 42 `EXTERNAL READ_CONVERSION_FN, WRITE_CONVERSION_FN, DTYPE_FILE_EXTENT_FN`
 43 `INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE`
 44 `INTEGER IERROR`

45
 46 `void MPI::Register_datarep`(`const char* datarep`,
 47 `MPI::Datarep_conversion_function* read_conversion_fn`,
 48 `MPI::Datarep_conversion_function* write_conversion_fn`,


```

MPI::Datarep_extent_function* dtype_file_extent_fn,
void* extra_state)

```

The call associates `read_conversion_fn`, `write_conversion_fn`, and `dtype_file_extent_fn` with the data representation identifier `datarep`. `datarep` can then be used as an argument to `MPI_FILE_SET_VIEW`, causing subsequent data access operations to call the conversion functions to convert all data items accessed between file data representation and native representation. `MPI_REGISTER_DATAREP` is a local operation and only registers the data representation for the calling MPI process. If `datarep` is already defined, an error in the error class `MPI_ERR_DUP_DATAREP` is raised using the default file error handler (see Section 12.7, page 431). The length of a data representation string is limited to the value of `MPI_MAX_DATAREP_STRING`. `MPI_MAX_DATAREP_STRING` must have a value of at least 64. No routines are provided to delete data representations and free the associated resources; it is not expected that an application will generate them in significant numbers.

Extent Callback

```

typedef int MPI_Datarep_extent_function(MPI_Datatype datatype,
MPI_Aint *file_extent, void *extra_state);

SUBROUTINE DATAREP_EXTENT_FUNCTION(DATATYPE, EXTENT, EXTRA_STATE, IERROR)
  INTEGER DATATYPE, IERROR
  INTEGER(KIND=MPI_ADDRESS_KIND) EXTENT, EXTRA_STATE

typedef void MPI::Datarep_extent_function(const MPI::Datatype& datatype,
MPI::Aint& file_extent, void* extra_state);

```

The function `dtype_file_extent_fn` must return, in `file_extent`, the number of bytes required to store `datatype` in the file representation. The function is passed, in `extra_state`, the argument that was passed to the `MPI_REGISTER_DATAREP` call. MPI will only call this routine with predefined datatypes employed by the user.

Datarep Conversion Functions

```

typedef int MPI_Datarep_conversion_function(void *userbuf,
MPI_Datatype datatype, int count, void *filebuf,
MPI_Offset position, void *extra_state);

SUBROUTINE DATAREP_CONVERSION_FUNCTION(USERBUF, DATATYPE, COUNT, FILEBUF,
POSITION, EXTRA_STATE, IERROR)
  <TYPE> USERBUF(*), FILEBUF(*)
  INTEGER COUNT, DATATYPE, IERROR
  INTEGER(KIND=MPI_OFFSET_KIND) POSITION
  INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE

typedef void MPI::Datarep_conversion_function(void* userbuf,
MPI::Datatype& datatype, int count, void* filebuf,
MPI::Offset position, void* extra_state);

```

The function `read_conversion_fn` must convert from file data representation to native representation. Before calling this routine, MPI allocates and fills `filebuf` with `count` contiguous data items. The type of each data item matches the corresponding entry

1 for the predefined datatype in the type signature of `datatype`. The function is passed,
2 in `extra_state`, the argument that was passed to the `MPI_REGISTER_DATAREP` call. The
3 function must copy all `count` data items from `filebuf` to `userbuf` in the distribution described
4 by `datatype`, converting each data item from file representation to native representation.
5 `datatype` will be equivalent to the datatype that the user passed to the read or write function.
6 If the size of `datatype` is less than the size of the `count` data items, the conversion function
7 must treat `datatype` as being contiguously tiled over the `userbuf`. The conversion function
8 must begin storing converted data at the location in `userbuf` specified by `position` into the
9 (tiled) `datatype`.

10
11 *Advice to users.* Although the conversion functions have similarities to `MPI_PACK`
12 and `MPI_UNPACK` in MPI-1, one should note the differences in the use of the arguments
13 `count` and `position`. In the conversion functions, `count` is a count of data items (i.e.,
14 count of typemap entries of `datatype`), and `position` is an index into this typemap. In
15 `MPI_PACK`, `incount` refers to the number of whole datatypes, and `position` is a number
16 of bytes. (*End of advice to users.*)

17
18 *Advice to implementors.* A converted read operation could be implemented as follows:

- 19 1. Get file extent of all data items
- 20 2. Allocate a `filebuf` large enough to hold all `count` data items
- 21 3. Read data from file into `filebuf`
- 22 4. Call `read_conversion_fn` to convert data and place it into `userbuf`
- 23 5. Deallocate `filebuf`

24
25
26
27 (*End of advice to implementors.*)

28
29 If MPI cannot allocate a buffer large enough to hold all the data to be converted from
30 a read operation, it may call the conversion function repeatedly using the same `datatype`
31 and `userbuf`, and reading successive chunks of data to be converted in `filebuf`. For the first
32 call (and in the case when all the data to be converted fits into `filebuf`), MPI will call the
33 function with `position` set to zero. Data converted during this call will be stored in the
34 `userbuf` according to the first `count` data items in `datatype`. Then in subsequent calls to the
35 conversion function, MPI will increment the value in `position` by the count of items converted
36 in the previous call, and the `userbuf` pointer will be unchanged.

37
38 *Rationale.* Passing the conversion function a `position` and one `datatype` for the
39 transfer allows the conversion function to decode the `datatype` only once and cache an
40 internal representation of it on the `datatype`. Then on subsequent calls, the conversion
41 function can use the `position` to quickly find its place in the `datatype` and continue
42 storing converted data where it left off at the end of the previous call. (*End of*
43 *rationale.*)

44
45 *Advice to users.* Although the conversion function may usefully cache an internal
46 representation on the `datatype`, it should not cache any state information specific to
47 an ongoing conversion operation, since it is possible for the same `datatype` to be used
48 concurrently in multiple conversion operations. (*End of advice to users.*)

The function `write_conversion_fn` must convert from native representation to file data representation. Before calling this routine, MPI allocates `filebuf` of a size large enough to hold `count` contiguous data items. The type of each data item matches the corresponding entry for the predefined datatype in the type signature of `datatype`. The function must copy `count` data items from `userbuf` in the distribution described by `datatype`, to a contiguous distribution in `filebuf`, converting each data item from native representation to file representation. If the size of `datatype` is less than the size of `count` data items, the conversion function must treat `datatype` as being contiguously tiled over the `userbuf`.

The function must begin copying at the location in `userbuf` specified by `position` into the (tiled) `datatype`. `datatype` will be equivalent to the datatype that the user passed to the read or write function. The function is passed, in `extra_state`, the argument that was passed to the `MPI_REGISTER_DATAREP` call.

The predefined constant `MPI_CONVERSION_FN_NULL` may be used as either `write_conversion_fn` or `read_conversion_fn`. In that case, MPI will not attempt to invoke `write_conversion_fn` or `read_conversion_fn`, respectively, but will perform the requested data access using the native data representation.

An MPI implementation must ensure that all data accessed is converted, either by using a `filebuf` large enough to hold all the requested data items or else by making repeated calls to the conversion function with the same `datatype` argument and appropriate values for `position`.

An implementation will only invoke the callback routines in this section (`read_conversion_fn`, `write_conversion_fn`, and `dtype_file_extent_fn`) when one of the read or write routines in Section 12.4, page 389, or `MPI_FILE_GET_TYPE_EXTENT` is called by the user. `dtype_file_extent_fn` will only be passed predefined datatypes employed by the user. The conversion functions will only be passed datatypes equivalent to those that the user has passed to one of the routines noted above.

The conversion functions must be reentrant. User defined data representations are restricted to use byte alignment for all types. Furthermore, it is erroneous for the conversion functions to call any collective routines or to free `datatype`.

The conversion functions should return an error code. If the returned error code has a value other than `MPI_SUCCESS`, the implementation will raise an error in the class `MPI_ERR_CONVERSION`.

12.5.4 Matching Data Representations

It is the user's responsibility to ensure that the data representation used to read data from a file is *compatible* with the data representation that was used to write that data to the file.

In general, using the same data representation name when writing and reading a file does not guarantee that the representation is compatible. Similarly, using different representation names on two different implementations may yield compatible representations.

Compatibility can be obtained when "external32" representation is used, although precision may be lost and the performance may be less than when "native" representation is used. Compatibility is guaranteed using "external32" provided at least one of the following conditions is met.

- The data access routines directly use types enumerated in Section 12.5.2, page 416, that are supported by all implementations participating in the I/O. The predefined

1 type used to write a data item must also be used to read a data item.

- 2
- 3 • In the case of Fortran 90 programs, the programs participating in the data accesses
- 4 obtain compatible datatypes using MPI routines that specify precision and/or range
- 5 (Section 13.2.5, page 458).
- 6
- 7 • For any given data item, the programs participating in the data accesses use compat-
- 8 ible predefined types to write and read the data item.

9 User-defined data representations may be used to provide an implementation compat-
10 ibility with another implementation’s “native” or “internal” representation.

11 *Advice to users.* Section 13.2.5, page 458, defines routines that support the use of
12 matching datatypes in heterogeneous environments and contains examples illustrating
13 their use. (*End of advice to users.*)

14 12.6 Consistency and Semantics

15 12.6.1 File Consistency

16 Consistency semantics define the outcome of multiple accesses to a single file. All file
17 accesses in MPI are relative to a specific file handle created from a collective open. MPI
18 provides three levels of consistency: sequential consistency among all accesses using a single
19 file handle, sequential consistency among all accesses using file handles created from a single
20 collective open with atomic mode enabled, and user-imposed consistency among accesses
21 other than the above. Sequential consistency means the behavior of a set of operations will
22 be as if the operations were performed in some serial order consistent with program order;
23 each access appears atomic, although the exact ordering of accesses is unspecified. User-
24 imposed consistency may be obtained using program order and calls to MPI_FILE_SYNC.

25 Let FH_1 be the set of file handles created from one particular collective open of the
26 file FOO , and FH_2 be the set of file handles created from a different collective open of
27 FOO . Note that nothing restrictive is said about FH_1 and FH_2 : the sizes of FH_1 and
28 FH_2 may be different, the groups of processes used for each open may or may not intersect,
29 the file handles in FH_1 may be destroyed before those in FH_2 are created, etc. Consider
30 the following three cases: a single file handle (e.g., $fh_1 \in FH_1$), two file handles created
31 from a single collective open (e.g., $fh_{1a} \in FH_1$ and $fh_{1b} \in FH_1$), and two file handles from
32 different collective opens (e.g., $fh_1 \in FH_1$ and $fh_2 \in FH_2$).

33 For the purpose of consistency semantics, a matched pair (Section 12.4.5, page 406)
34 of split collective data access operations (e.g., MPI_FILE_READ_ALL_BEGIN and
35 MPI_FILE_READ_ALL_END) compose a single data access operation. Similarly, a nonblock-
36 ing data access routine (e.g., MPI_FILE_IREAD) and the routine which completes the request
37 (e.g., MPI_WAIT) also compose a single data access operation. For all cases below, these data
38 access operations are subject to the same constraints as blocking data access operations.

39 *Advice to users.* For an MPI_FILE_IREAD and MPI_WAIT pair, the operation begins
40 when MPI_FILE_IREAD is called and ends when MPI_WAIT returns. (*End of advice to*
41 *users.*)

42 Assume that A_1 and A_2 are two data access operations. Let D_1 (D_2) be the set of
43 absolute byte displacements of every byte accessed in A_1 (A_2). The two data accesses
44

overlap if $D_1 \cap D_2 \neq \emptyset$. The two data accesses *conflict* if they overlap and at least one is a write access.

Let SEQ_{fh} be a sequence of file operations on a single file handle, bracketed by `MPI_FILE_SYNCs` on that file handle. (Both opening and closing a file implicitly perform an `MPI_FILE_SYNC`.) SEQ_{fh} is a “write sequence” if any of the data access operations in the sequence are writes or if any of the file manipulation operations in the sequence change the state of the file (e.g., `MPI_FILE_SET_SIZE` or `MPI_FILE_PREALLOCATE`). Given two sequences, SEQ_1 and SEQ_2 , we say they are not *concurrent* if one sequence is guaranteed to completely precede the other (temporally).

The requirements for guaranteeing sequential consistency among all accesses to a particular file are divided into the three cases given below. If any of these requirements are not met, then the value of all data in that file is implementation dependent.

Case 1: $fh_1 \in FH_1$ All operations on fh_1 are sequentially consistent if atomic mode is set. If nonatomic mode is set, then all operations on fh_1 are sequentially consistent if they are either nonconcurrent, nonconflicting, or both.

Case 2: $fh_{1a} \in FH_1$ and $fh_{1b} \in FH_1$ Assume A_1 is a data access operation using fh_{1a} , and A_2 is a data access operation using fh_{1b} . If for any access A_1 , there is no access A_2 that conflicts with A_1 , then MPI guarantees sequential consistency.

However, unlike POSIX semantics, the default MPI semantics for conflicting accesses do not guarantee sequential consistency. If A_1 and A_2 conflict, sequential consistency can be guaranteed by either enabling atomic mode via the `MPI_FILE_SET_ATOMICITY` routine, or meeting the condition described in Case 3 below.

Case 3: $fh_1 \in FH_1$ and $fh_2 \in FH_2$ Consider access to a single file using file handles from distinct collective opens. In order to guarantee sequential consistency, `MPI_FILE_SYNC` must be used (both opening and closing a file implicitly perform an `MPI_FILE_SYNC`).

Sequential consistency is guaranteed among accesses to a single file if for any write sequence SEQ_1 to the file, there is no sequence SEQ_2 to the file which is *concurrent* with SEQ_1 . To guarantee sequential consistency when there are write sequences, `MPI_FILE_SYNC` must be used together with a mechanism that guarantees nonconcurrency of the sequences.

See the examples in Section 12.6.10, page 427, for further clarification of some of these consistency semantics.

`MPI_FILE_SET_ATOMICITY(fh, flag)`

| | | |
|-------|------|--|
| INOUT | fh | file handle (handle) |
| IN | flag | true to set atomic mode, false to set nonatomic mode (logical) |

`int MPI_File_set_atomicity(MPI_File fh, int flag)`

`MPI_FILE_SET_ATOMICITY(FH, FLAG, IERROR)`

INTEGER FH, IERROR

LOGICAL FLAG

```
1 void MPI::File::Set_atomicity(bool flag)
```

2
3 Let FH be the set of file handles created by one collective open. The consistency semantics for data access operations using FH is set by collectively calling
4 `MPI_FILE_SET_ATOMICITY` on FH . `MPI_FILE_SET_ATOMICITY` is collective; all processes
5 in the group must pass identical values for fh and $flag$. If $flag$ is true, atomic mode is set; if
6 $flag$ is false, nonatomic mode is set.
7

8 Changing the consistency semantics for an open file only affects new data accesses.
9 All completed data accesses are guaranteed to abide by the consistency semantics in effect
10 during their execution. Nonblocking data accesses and split collective operations that have
11 not completed (e.g., via `MPI_WAIT`) are only guaranteed to abide by nonatomic mode
12 consistency semantics.

13 *Advice to implementors.* Since the semantics guaranteed by atomic mode are stronger
14 than those guaranteed by nonatomic mode, an implementation is free to adhere to
15 the more stringent atomic mode semantics for outstanding requests. (*End of advice*
16 *to implementors.*)
17

```
18  
19  
20 MPI_FILE_GET_ATOMICITY(fh, flag)
```

```
21     IN      fh                file handle (handle)
22     OUT     flag              true if atomic mode, false if nonatomic mode (logical)
```

```
23  
24  
25 int MPI_File_get_atomicity(MPI_File fh, int *flag)
```

```
26 MPI_FILE_GET_ATOMICITY(FH, FLAG, IERROR)
```

```
27     INTEGER FH, IERROR
```

```
28     LOGICAL FLAG
```

```
29  
30 bool MPI::File::Get_atomicity() const
```

31
32 `MPI_FILE_GET_ATOMICITY` returns the current consistency semantics for data access
33 operations on the set of file handles created by one collective open. If $flag$ is true, atomic
34 mode is enabled; if $flag$ is false, nonatomic mode is enabled.
35

```
36 MPI_FILE_SYNC(fh)
```

```
37     INOUT   fh                file handle (handle)
```

```
38  
39  
40 int MPI_File_sync(MPI_File fh)
```

```
41 MPI_FILE_SYNC(FH, IERROR)
```

```
42     INTEGER FH, IERROR
```

```
43  
44 void MPI::File::Sync()
```

45
46 Calling `MPI_FILE_SYNC` with fh causes all previous writes to fh by the calling process
47 to be transferred to the storage device. If other processes have made updates to the storage
48 device, then all such updates become visible to subsequent reads of fh by the calling process.

`MPI_FILE_SYNC` may be necessary to ensure sequential consistency in certain cases (see above).

`MPI_FILE_SYNC` is a collective operation.

The user is responsible for ensuring that all nonblocking requests and split collective operations on `fh` have been completed before calling `MPI_FILE_SYNC`—otherwise, the call to `MPI_FILE_SYNC` is erroneous.

12.6.2 Random Access vs. Sequential Files

MPI distinguishes ordinary random access files from sequential stream files, such as pipes and tape files. Sequential stream files must be opened with the `MPI_MODE_SEQUENTIAL` flag set in the `amode`. For these files, the only permitted data access operations are shared file pointer reads and writes. Filetypes and etypes with holes are erroneous. In addition, the notion of file pointer is not meaningful; therefore, calls to `MPI_FILE_SEEK_SHARED` and `MPI_FILE_GET_POSITION_SHARED` are erroneous, and the pointer update rules specified for the data access routines do not apply. The amount of data accessed by a data access operation will be the amount requested unless the end of file is reached or an error is raised.

Rationale. This implies that reading on a pipe will always wait until the requested amount of data is available or until the process writing to the pipe has issued an end of file. (*End of rationale.*)

Finally, for some sequential files, such as those corresponding to magnetic tapes or streaming network connections, writes to the file may be destructive. In other words, a write may act as a truncate (a `MPI_FILE_SET_SIZE` with `size` set to the current position) followed by the write.

12.6.3 Progress

The progress rules of MPI are both a promise to users and a set of constraints on implementors. In cases where the progress rules restrict possible implementation choices more than the interface specification alone, the progress rules take precedence.

All blocking routines must complete in finite time unless an exceptional condition (such as resource exhaustion) causes an error.

Nonblocking data access routines inherit the following progress rule from nonblocking point to point communication: a nonblocking write is equivalent to a nonblocking send for which a receive is eventually posted, and a nonblocking read is equivalent to a nonblocking receive for which a send is eventually posted.

Finally, an implementation is free to delay progress of collective routines until all processes in the group associated with the collective call have invoked the routine. Once all processes in the group have invoked the routine, the progress rule of the equivalent noncollective routine must be followed.

12.6.4 Collective File Operations

Collective file operations are subject to the same restrictions as collective communication operations. For a complete discussion, please refer to the semantics set forth in MPI-1 [23], Section 4.12.

1 Collective file operations are collective over a dup of the communicator used to open
2 the file—this duplicate communicator is implicitly specified via the file handle argument.
3 Different processes can pass different values for other arguments of a collective routine unless
4 specified otherwise.

6 12.6.5 Type Matching

7 The type matching rules for I/O mimic the type matching rules for communication with one
8 exception: if `etype` is `MPI_BYTE`, then this matches any `datatype` in a data access operation.
9 In general, the `etype` of data items written must match the `etype` used to read the items,
10 and for each data access operation, the current `etype` must also match the type declaration
11 of the data access buffer.

12
13 *Advice to users.* In most cases, use of `MPI_BYTE` as a wild card will defeat the
14 file interoperability features of MPI. File interoperability can only perform automatic
15 conversion between heterogeneous data representations when the exact datatypes ac-
16 cessed are explicitly specified. (*End of advice to users.*)

18 12.6.6 Miscellaneous Clarifications

19 Once an I/O routine completes, it is safe to free any opaque objects passed as arguments
20 to that routine. For example, the `comm` and `info` used in an `MPI_FILE_OPEN`, or the `etype`
21 and `filetype` used in an `MPI_FILE_SET_VIEW`, can be freed without affecting access to the
22 file. Note that for nonblocking routines and split collective operations, the operation must
23 be completed before it is safe to reuse data buffers passed as arguments.

24 As in communication, datatypes must be committed before they can be used in file
25 manipulation or data access operations. For example, the `etype` and
26 `filetype` must be committed before calling `MPI_FILE_SET_VIEW`, and the `datatype` must be
27 committed before calling `MPI_FILE_READ` or `MPI_FILE_WRITE`.

30 12.6.7 MPI_Offset Type

31 `MPI_Offset` is an integer type of size sufficient to represent the size (in bytes) of the largest
32 file supported by MPI. Displacements and offsets are always specified as values of type
33 `MPI_Offset`.

34 In Fortran, the corresponding integer is an integer of kind `MPI_OFFSET_KIND`, defined
35 in `mpif.h` and the `mpi` module.

36 In Fortran 77 environments that do not support `KIND` parameters,
37 `MPI_Offset` arguments should be declared as an `INTEGER` of suitable size. The language
38 interoperability implications for `MPI_Offset` are similar to those for addresses (see Section
39 13.3, page 466).

42 12.6.8 Logical vs. Physical File Layout

43 MPI specifies how the data should be laid out in a virtual file structure (the view), not
44 how that file structure is to be stored on one or more disks. Specification of the physical
45 file structure was avoided because it is expected that the mapping of files to disks will be
46 system specific, and any specific control over file layout would therefore restrict program
47 portability. However, there are still cases where some information may be necessary to
48

optimize file layout. This information can be provided as *hints* specified via *info* when a file is created (see Section 12.2.8, page 384).

12.6.9 File Size

The size of a file may be increased by writing to the file after the current end of file. The size may also be changed by calling *MPI size changing* routines, such as `MPI_FILE_SET_SIZE`. A call to a size changing routine does not necessarily change the file size. For example, calling `MPI_FILE_PREALLOCATE` with a size less than the current size does not change the size.

Consider a set of bytes that has been written to a file since the most recent call to a size changing routine, or since `MPI_FILE_OPEN` if no such routine has been called. Let the *high byte* be the byte in that set with the largest displacement. The file size is the larger of

- One plus the displacement of the high byte.
- The size immediately after the size changing routine, or `MPI_FILE_OPEN`, returned.

When applying consistency semantics, calls to `MPI_FILE_SET_SIZE` and `MPI_FILE_PREALLOCATE` are considered writes to the file (which conflict with operations that access bytes at displacements between the old and new file sizes), and `MPI_FILE_GET_SIZE` is considered a read of the file (which overlaps with all accesses to the file).

Advice to users. Any sequence of operations containing the collective routines `MPI_FILE_SET_SIZE` and `MPI_FILE_PREALLOCATE` is a write sequence. As such, sequential consistency in nonatomic mode is not guaranteed unless the conditions in Section 12.6.1, page 422, are satisfied. (*End of advice to users.*)

File pointer update semantics (i.e., file pointers are updated by the amount accessed) are only guaranteed if file size changes are sequentially consistent.

Advice to users. Consider the following example. Given two operations made by separate processes to a file containing 100 bytes: an `MPI_FILE_READ` of 10 bytes and an `MPI_FILE_SET_SIZE` to 0 bytes. If the user does not enforce sequential consistency between these two operations, the file pointer may be updated by the amount requested (10 bytes) even if the amount accessed is zero bytes. (*End of advice to users.*)

12.6.10 Examples

The examples in this section illustrate the application of the MPI consistency and semantics guarantees. These address

- conflicting accesses on file handles obtained from a single collective open, and
- all accesses on file handles obtained from two separate collective opens.

The simplest way to achieve consistency for conflicting accesses is to obtain sequential consistency by setting atomic mode. For the code below, process 1 will read either 0 or 10 integers. If the latter, every element of `b` will be 5. If nonatomic mode is set, the results of the read are undefined.

```

1  /* Process 0 */
2  int  i, a[10] ;
3  int  TRUE = 1;
4
5  for ( i=0;i<10;i++)
6      a[i] = 5 ;
7
8  MPI_File_open( MPI_COMM_WORLD, "workfile",
9                MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &fh0 ) ;
10 MPI_File_set_view( fh0, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
11 MPI_File_set_atomicity( fh0, TRUE ) ;
12 MPI_File_write_at(fh0, 0, a, 10, MPI_INT, &status) ;
13 /* MPI_Barrier( MPI_COMM_WORLD ) ; */
14
15 /* Process 1 */
16 int  b[10] ;
17 int  TRUE = 1;
18 MPI_File_open( MPI_COMM_WORLD, "workfile",
19                MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &fh1 ) ;
20 MPI_File_set_view( fh1, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
21 MPI_File_set_atomicity( fh1, TRUE ) ;
22 /* MPI_Barrier( MPI_COMM_WORLD ) ; */
23 MPI_File_read_at(fh1, 0, b, 10, MPI_INT, &status) ;
24

```

A user may guarantee that the write on process 0 precedes the read on process 1 by imposing temporal order with, for example, calls to MPI_BARRIER.

Advice to users. Routines other than MPI_BARRIER may be used to impose temporal order. In the example above, process 0 could use MPI_SEND to send a 0 byte message, received by process 1 using MPI_RECV. (*End of advice to users.*)

Alternatively, a user can impose consistency with nonatomic mode set:

```

33 /* Process 0 */
34 int  i, a[10] ;
35 for ( i=0;i<10;i++)
36     a[i] = 5 ;
37
38 MPI_File_open( MPI_COMM_WORLD, "workfile",
39                MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &fh0 ) ;
40 MPI_File_set_view( fh0, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
41 MPI_File_write_at(fh0, 0, a, 10, MPI_INT, &status) ;
42 MPI_File_sync( fh0 ) ;
43 MPI_Barrier( MPI_COMM_WORLD ) ;
44 MPI_File_sync( fh0 ) ;
45
46 /* Process 1 */
47 int  b[10] ;
48 MPI_File_open( MPI_COMM_WORLD, "workfile",

```

```

        MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &fh1 ) ;
MPI_File_set_view( fh1, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
MPI_File_sync( fh1 ) ;
MPI_Barrier( MPI_COMM_WORLD ) ;
MPI_File_sync( fh1 ) ;
MPI_File_read_at(fh1, 0, b, 10, MPI_INT, &status ) ;

```

The “sync-barrier-sync” construct is required because:

- The barrier ensures that the write on process 0 occurs before the read on process 1.
- The first sync guarantees that the data written by all processes is transferred to the storage device.
- The second sync guarantees that all data which has been transferred to the storage device is visible to all processes. (This does not affect process 0 in this example.)

The following program represents an erroneous attempt to achieve consistency by eliminating the apparently superfluous second “sync” call for each process.

```

/* ----- THIS EXAMPLE IS ERRONEOUS ----- */
/* Process 0 */
int i, a[10] ;
for ( i=0;i<10;i++)
    a[i] = 5 ;

MPI_File_open( MPI_COMM_WORLD, "workfile",
               MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &fh0 ) ;
MPI_File_set_view( fh0, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
MPI_File_write_at(fh0, 0, a, 10, MPI_INT, &status ) ;
MPI_File_sync( fh0 ) ;
MPI_Barrier( MPI_COMM_WORLD ) ;

/* Process 1 */
int b[10] ;
MPI_File_open( MPI_COMM_WORLD, "workfile",
               MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &fh1 ) ;
MPI_File_set_view( fh1, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
MPI_Barrier( MPI_COMM_WORLD ) ;
MPI_File_sync( fh1 ) ;
MPI_File_read_at(fh1, 0, b, 10, MPI_INT, &status ) ;

/* ----- THIS EXAMPLE IS ERRONEOUS ----- */

```

The above program also violates the MPI rule against out-of-order collective operations and will deadlock for implementations in which MPI_FILE_SYNC blocks.

Advice to users. Some implementations may choose to implement MPI_FILE_SYNC as a temporally synchronizing function. When using such an implementation, the “sync-barrier-sync” construct above can be replaced by a single “sync.” The results of using such code with an implementation for which MPI_FILE_SYNC is not temporally synchronizing is undefined. (*End of advice to users.*)

1 Asynchronous I/O

2 The behavior of asynchronous I/O operations is determined by applying the rules specified
3 above for synchronous I/O operations.

4 The following examples all access a preexisting file “myfile.” Word 10 in myfile initially
5 contains the integer 2. Each example writes and reads word 10.

6 First consider the following code fragment:

```
7
8 int a = 4, b, TRUE=1;
9 MPI_File_open( MPI_COMM_WORLD, "myfile",
10               MPI_MODE_RDWR, MPI_INFO_NULL, &fh ) ;
11 MPI_File_set_view( fh, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
12 /* MPI_File_set_atomicsity( fh, TRUE ) ; Use this to set atomic mode. */
13 MPI_File_iread_at(fh, 10, &a, 1, MPI_INT, &reqs[0]) ;
14 MPI_File_iread_at(fh, 10, &b, 1, MPI_INT, &reqs[1]) ;
15 MPI_Waitall(2, reqs, statuses) ;
16
```

17 For asynchronous data access operations, MPI specifies that the access occurs at any time
18 between the call to the asynchronous data access routine and the return from the corre-
19 sponding request complete routine. Thus, executing either the read before the write, or the
20 write before the read is consistent with program order. If atomic mode is set, then MPI
21 guarantees sequential consistency, and the program will read either 2 or 4 into b. If atomic
22 mode is not set, then sequential consistency is not guaranteed and the program may read
23 something other than 2 or 4 due to the conflicting data access.

24 Similarly, the following code fragment does not order file accesses:

```
25
26 int a = 4, b;
27 MPI_File_open( MPI_COMM_WORLD, "myfile",
28               MPI_MODE_RDWR, MPI_INFO_NULL, &fh ) ;
29 MPI_File_set_view( fh, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
30 /* MPI_File_set_atomicsity( fh, TRUE ) ; Use this to set atomic mode. */
31 MPI_File_iread_at(fh, 10, &a, 1, MPI_INT, &reqs[0]) ;
32 MPI_File_iread_at(fh, 10, &b, 1, MPI_INT, &reqs[1]) ;
33 MPI_Wait(&reqs[0], &status) ;
34 MPI_Wait(&reqs[1], &status) ;
35
```

36 If atomic mode is set, either 2 or 4 will be read into b. Again, MPI does not guarantee
37 sequential consistency in nonatomic mode.

38 On the other hand, the following code fragment:

```
39
40 int a = 4, b;
41 MPI_File_open( MPI_COMM_WORLD, "myfile",
42               MPI_MODE_RDWR, MPI_INFO_NULL, &fh ) ;
43 MPI_File_set_view( fh, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
44 MPI_File_iread_at(fh, 10, &a, 1, MPI_INT, &reqs[0]) ;
45 MPI_Wait(&reqs[0], &status) ;
46 MPI_File_iread_at(fh, 10, &b, 1, MPI_INT, &reqs[1]) ;
47 MPI_Wait(&reqs[1], &status) ;
48
```

defines the same ordering as:

```

int a = 4, b;
MPI_File_open( MPI_COMM_WORLD, "myfile",
               MPI_MODE_RDWR, MPI_INFO_NULL, &fh ) ;
MPI_File_set_view( fh, 0, MPI_INT, MPI_INT, "native", MPI_INFO_NULL ) ;
MPI_File_write_at(fh, 10, &a, 1, MPI_INT, &status ) ;
MPI_File_read_at(fh, 10, &b, 1, MPI_INT, &status ) ;

```

Since

- nonconcurrent operations on a single file handle are sequentially consistent, and
- the program fragments specify an order for the operations,

MPI guarantees that both program fragments will read the value 4 into `b`. There is no need to set atomic mode for this example.

Similar considerations apply to conflicting accesses of the form:

```

MPI_File_write_all_begin(fh,...) ;
MPI_File_iread(fh,...) ;
MPI_Wait(fh,...) ;
MPI_File_write_all_end(fh,...) ;

```

Recall that constraints governing consistency and semantics are not relevant to the following:

```

MPI_File_write_all_begin(fh,...) ;
MPI_File_read_all_begin(fh,...) ;
MPI_File_read_all_end(fh,...) ;
MPI_File_write_all_end(fh,...) ;

```

since split collective operations on the same file handle may not overlap (see Section 12.4.5, page 406).

12.7 I/O Error Handling

By default, communication errors are fatal—`MPI_ERRORS_ARE_FATAL` is the default error handler associated with `MPI_COMM_WORLD`. I/O errors are usually less catastrophic (e.g., “file not found”) than communication errors, and common practice is to catch these errors and continue executing. For this reason, MPI provides additional error facilities for I/O.

Advice to users. MPI does not specify the state of a computation after an erroneous MPI call has occurred. A high quality implementation will support the I/O error handling facilities, allowing users to write programs using common practice for I/O. (*End of advice to users.*)

Like communicators, each file handle has an error handler associated with it. The MPI-2 I/O error handling routines are defined in Section 7.3.1, page 253.

When MPI calls a user-defined error handler resulting from an error on a particular file handle, the first two arguments passed to the file error handler are the file handle and the error code. For I/O errors that are not associated with a valid file handle (e.g., in

1 MPI_FILE_OPEN or MPI_FILE_DELETE), the first argument passed to the error handler is
 2 MPI_FILE_NULL,

3 I/O error handling differs from communication error handling in another important
 4 aspect. By default, the predefined error handler for file handles is MPI_ERRORS_RETURN.
 5 The default file error handler has two purposes: when a new file handle is created (by
 6 MPI_FILE_OPEN), the error handler for the new file handle is initially set to the default
 7 error handler, and I/O routines that have no valid file handle on which to raise an error
 8 (e.g., MPI_FILE_OPEN or MPI_FILE_DELETE) use the default file error handler. The default
 9 file error handler can be changed by specifying MPI_FILE_NULL as the fh argument to
 10 MPI_FILE_SET_ERRHANDLER. The current value of the default file error handler can be
 11 determined by passing MPI_FILE_NULL as the fh argument to MPI_FILE_GET_ERRHANDLER.

12
 13 *Rationale.* For communication, the default error handler is inherited from
 14 MPI_COMM_WORLD. In I/O, there is no analogous “root” file handle from which default
 15 properties can be inherited. Rather than invent a new global file handle, the default
 16 file error handler is manipulated as if it were attached to MPI_FILE_NULL. (*End of*
 17 *rationale.*)

19 12.8 I/O Error Classes

21 The implementation dependent error codes returned by the I/O routines can be converted
 22 into the following error classes. In addition, calls to routines in this chapter may raise errors
 23 in other MPI classes, such as MPI_ERR_TYPE.

26 12.9 Examples

28 12.9.1 Double Buffering with Split Collective I/O

29 This example shows how to overlap computation and output. The computation is performed
 30 by the function `compute_buffer()`.

```

32 /*=====
33 *
34 * Function:          double_buffer
35 *
36 * Synopsis:
37 *   void double_buffer(
38 *       MPI_File fh,                ** IN
39 *       MPI_Datatype buftype,      ** IN
40 *       int bufcount               ** IN
41 *   )
42 *
43 * Description:
44 *   Performs the steps to overlap computation with a collective write
45 *   by using a double-buffering technique.
46 *
47 * Parameters:
48 *   fh                previously opened MPI file handle

```

| | | |
|-------------------------------|--|----|
| | | 1 |
| | | 2 |
| | | 3 |
| | | 4 |
| | | 5 |
| | | 6 |
| | | 7 |
| | | 8 |
| | | 9 |
| | | 10 |
| MPI_ERR_FILE | Invalid file handle | 11 |
| MPI_ERR_NOT_SAME | Collective argument not identical on all processes, or collective routines called in a different order by different processes | 12 |
| | | 13 |
| | | 14 |
| MPI_ERR_AMODE | Error related to the <code>amode</code> passed to <code>MPI_FILE_OPEN</code> | 15 |
| | | 16 |
| MPI_ERR_UNSUPPORTED_DATAREP | Unsupported <code>datarep</code> passed to <code>MPI_FILE_SET_VIEW</code> | 17 |
| | | 18 |
| MPI_ERR_UNSUPPORTED_OPERATION | Unsupported operation, such as seeking on a file which supports sequential access only | 19 |
| | | 20 |
| MPI_ERR_NO_SUCH_FILE | File does not exist | 21 |
| MPI_ERR_FILE_EXISTS | File exists | 22 |
| MPI_ERR_BAD_FILE | Invalid file name (e.g., path name too long) | 23 |
| MPI_ERR_ACCESS | Permission denied | 24 |
| MPI_ERR_NO_SPACE | Not enough space | 25 |
| MPI_ERR_QUOTA | Quota exceeded | 26 |
| MPI_ERR_READ_ONLY | Read-only file or file system | 27 |
| MPI_ERR_FILE_IN_USE | File operation could not be completed, as the file is currently open by some process | 28 |
| | | 29 |
| MPI_ERR_DUP_DATAREP | Conversion functions could not be registered because a data representation identifier that was already defined was passed to <code>MPI_REGISTER_DATAREP</code> | 30 |
| | | 31 |
| | | 32 |
| | | 33 |
| MPI_ERR_CONVERSION | An error occurred in a user supplied data conversion function. | 34 |
| | | 35 |
| MPI_ERR_IO | Other I/O error | 36 |
| | | 37 |

Table 12.3: Error classes returned from MPI I/O routines.

38
39
40
41
42
43
44
45
46
47
48

```

1      *      buftype          MPI datatype for memory layout
2      *
3      *      bufcount        # buftype elements to transfer
4      *-----*/
5
6      /* this macro switches which buffer "x" is pointing to */
7      #define TOGGLE_PTR(x) (((x)==(buffer1)) ? (x=buffer2) : (x=buffer1))
8
9      void double_buffer( MPI_File fh, MPI_Datatype buftype, int bufcount)
10     {
11
12         MPI_Status status;          /* status for MPI calls */
13         float *buffer1, *buffer2;  /* buffers to hold results */
14         float *compute_buf_ptr;    /* destination buffer */
15                                     /* for computing */
16         float *write_buf_ptr;      /* source for writing */
17         int done;                  /* determines when to quit */
18
19         /* buffer initialization */
20         buffer1 = (float *)
21                 malloc(bufcount*sizeof(float)) ;
22         buffer2 = (float *)
23                 malloc(bufcount*sizeof(float)) ;
24         compute_buf_ptr = buffer1 ; /* initially point to buffer1 */
25         write_buf_ptr  = buffer1 ; /* initially point to buffer1 */
26
27
28         /* DOUBLE-BUFFER prolog:
29          *   compute buffer1; then initiate writing buffer1 to disk
30          */
31         compute_buffer(compute_buf_ptr, bufcount, &done);
32         MPI_File_write_all_begin(fh, write_buf_ptr, bufcount, buftype);
33
34         /* DOUBLE-BUFFER steady state:
35          *   Overlap writing old results from buffer pointed to by write_buf_ptr
36          *   with computing new results into buffer pointed to by compute_buf_ptr.
37          *
38          *   There is always one write-buffer and one compute-buffer in use
39          *   during steady state.
40          */
41         while (!done) {
42             TOGGLE_PTR(compute_buf_ptr);
43             compute_buffer(compute_buf_ptr, bufcount, &done);
44             MPI_File_write_all_end(fh, write_buf_ptr, &status);
45             TOGGLE_PTR(write_buf_ptr);
46             MPI_File_write_all_begin(fh, write_buf_ptr, bufcount, buftype);
47         }
48

```



```

/* DOUBLE-BUFFER epilog:
 *   wait for final write to complete.
 */
MPI_File_write_all_end(fh, write_buf_ptr, &status);

/* buffer cleanup */
free(buffer1);
free(buffer2);
}

```

12.9.2 Subarray Filetype Constructor

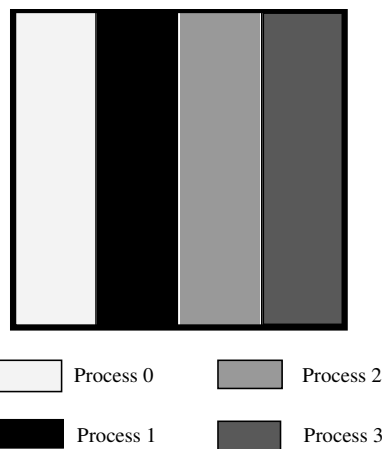


Figure 12.4: Example array file layout

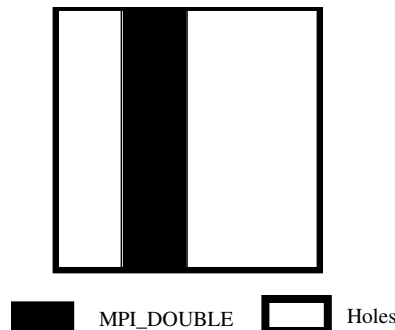


Figure 12.5: Example local array filetype for process 1

Assume we are writing out a 100x100 2D array of double precision floating point numbers that is distributed among 4 processes such that each process has a block of 25 columns (e.g., process 0 has columns 0-24, process 1 has columns 25-49, etc.; see Figure 12.4). To create the filetypes for each process one could use the following C program:

```

double subarray[100][25];
MPI_Datatype filetype;

```

```

1   int sizes[2], subsizes[2], starts[2];
2   int rank;
3
4   MPI_Comm_rank(MPI_COMM_WORLD, &rank);
5   sizes[0]=100; sizes[1]=100;
6   subsizes[0]=100; subsizes[1]=25;
7   starts[0]=0; starts[1]=rank*subsizes[1];
8
9   MPI_Type_create_subarray(2, sizes, subsizes, starts, MPI_ORDER_C,
10                          MPI_DOUBLE, &filetype);
11

```

Or, equivalently in Fortran:

```

12
13
14     double precision subarray(100,25)
15     integer filetype, rank, ierror
16     integer sizes(2), subsizes(2), starts(2)
17
18     call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierror)
19     sizes(1)=100
20     sizes(2)=100
21     subsizes(1)=100
22     subsizes(2)=25
23     starts(1)=0
24     starts(2)=rank*subsizes(2)
25
26     call MPI_TYPE_CREATE_SUBARRAY(2, sizes, subsizes, starts, &
27                                  MPI_ORDER_FORTRAN, MPI_DOUBLE_PRECISION, &
28                                  filetype, ierror)
29

```

The generated filetype will then describe the portion of the file contained within the process's subarray with holes for the space taken by the other processes. Figure 12.5 shows the filetype created for process 1.

```

30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 13

Language Bindings

13.1 C++

13.1.1 Overview

This section presents a complete C++ language interface for MPI. There are some issues specific to C++ that must be considered in the design of this interface that go beyond the simple description of language bindings. In particular, in C++, we must be concerned with the design of objects and their interfaces, rather than just the design of a language-specific functional interface to MPI. Fortunately, the original design of MPI was based on the notion of objects, so a natural set of classes is already part of MPI.

Since the original design of MPI-1 did not include a C++ language interface, a complete list of C++ bindings for MPI-1 functions is provided in Annex A.4. MPI-2 includes C++ bindings as part of its function specifications. In some cases, MPI-2 provides new names for the C bindings of MPI-1 functions. In this case, the C++ binding matches the new C name — there is no binding for the deprecated name. As such, the C++ binding for the new name appears in Annex A, not Annex A.4.

13.1.2 Design

The C++ language interface for MPI is designed according to the following criteria:

1. The C++ language interface consists of a small set of classes with a lightweight functional interface to MPI. The classes are based upon the fundamental MPI object types (e.g., communicator, group, etc.).
2. The MPI C++ language bindings provide a semantically correct interface to MPI.
3. To the greatest extent possible, the C++ bindings for MPI functions are member functions of MPI classes.

Rationale. Providing a lightweight set of MPI objects that correspond to the basic MPI types is the best fit to MPI's implicit object-based design; methods can be supplied for these objects to realize MPI functionality. The existing C bindings can be used in C++ programs, but much of the expressive power of the C++ language is forfeited. On the other hand, while a comprehensive class library would make user programming more elegant, such a library it is not suitable as a language binding for MPI since a

1 binding must provide a direct and unambiguous mapping to the specified functionality
 2 of MPI. (*End of rationale.*)
 3
 4 .
 5

6 13.1.3 C++ Classes for MPI

7 All MPI classes, constants, and functions are declared within the scope of an MPI namespace.
 8 Thus, instead of the MPI_ prefix that is used in C and Fortran, MPI functions essentially
 9 have an MPI:: prefix.
 10

11 The members of the MPI namespace are those classes corresponding to objects implicitly
 12 used by MPI. An abbreviated definition of the MPI namespace for MPI-1 and its member
 13 classes is as follows:

```
14 namespace MPI {
15     class Comm                {...};
16     class Intracomm : public Comm    {...};
17     class Graphcomm : public Intracomm    {...};
18     class Cartcomm  : public Intracomm    {...};
19     class Intercomm : public Comm        {...};
20     class Datatype                {...};
21     class Errhandler              {...};
22     class Exception               {...};
23     class Group                   {...};
24     class Op                       {...};
25     class Request                 {...};
26     class Prerequest  : public Request    {...};
27     class Status                  {...};
28 };
29
```

30 Additionally, the following classes defined for MPI-2:

```
31 namespace MPI {
32     class File                  {...};
33     class Grequest  : public Request    {...};
34     class Info                  {...};
35     class Win                   {...};
36 };
37
```

38 Note that there are a small number of derived classes, and that virtual inheritance is
 39 *not* used.
 40

41 13.1.4 Class Member Functions for MPI

42 Besides the member functions which constitute the C++ language bindings for MPI, the
 43 C++ language interface has additional functions (as required by the C++ language). In
 44 particular, the C++ language interface must provide a constructor and destructor, an
 45 assignment operator, and comparison operators.
 46

47 The complete set of C++ language bindings for MPI-1 is presented in Annex A.4. The
 48 bindings take advantage of some important C++ features, such as references and `const`.

Declarations (which apply to all MPI member classes) for construction, destruction, copying, assignment, comparison, and mixed-language operability are also provided. To maintain consistency with what has gone before, the binding definitions are given in the same order as given for the C bindings in [23].

Except where indicated, all non-static member functions (except for constructors and the assignment operator) of MPI member classes are virtual functions.

Rationale. Providing virtual member functions is an important part of design for inheritance. Virtual functions can be bound at run-time, which allows users of libraries to re-define the behavior of objects already contained in a library. There is a small performance penalty that must be paid (the virtual function must be looked up before it can be called). However, users concerned about this performance penalty can force compile-time function binding. (*End of rationale.*)

Example 13.1 Example showing a derived MPI class.

```
class foo_comm : public MPI::Intracomm {
public:
    void Send(const void* buf, int count, const MPI::Datatype& type,
              int dest, int tag) const
    {
        // Class library functionality
        MPI::Intracomm::Send(buf, count, type, dest, tag);
        // More class library functionality
    }
};
```

Advice to implementors. Implementors must be careful to avoid unintended side effects from class libraries that use inheritance, especially in layered implementations. For example, if MPI_BCAST is implemented by repeated calls to MPI_SEND or MPI_RECV, the behavior of MPI_BCAST cannot be changed by derived communicator classes that might redefine MPI_SEND or MPI_RECV. The implementation of MPI_BCAST must explicitly use the MPI_SEND (or MPI_RECV) of the base MPI::Comm class. (*End of advice to implementors.*)

13.1.5 Semantics

The semantics of the member functions constituting the C++ language binding for MPI are specified by the MPI function description itself. Here, we specify the semantics for those portions of the C++ language interface that are not part of the language binding. In this subsection, functions are prototyped using the type MPI::<CLASS> rather than listing each function for every MPI class; the word <CLASS> can be replaced with any valid MPI class name (e.g., Group), except as noted.

Construction / Destruction The default constructor and destructor are prototyped as follows:

```
MPI::<CLASS>()
```

```
~MPI::<CLASS>()
```

In terms of construction and destruction, opaque MPI user level objects behave like handles. Default constructors for all MPI objects except `MPI::Status` create corresponding `MPI::*_NULL` handles. That is, when an MPI object is instantiated, comparing it with its corresponding `MPI::*_NULL` object will return `true`. The default constructors do not create new MPI opaque objects. Some classes have a member function `Create()` for this purpose.

Example 13.2 In the following code fragment, the test will return `true` and the message will be sent to `cout`.

```

1 void foo()
2 {
3     MPI::Intracomm bar;
4
5     if (bar == MPI::COMM_NULL)
6         cout << "bar is MPI::COMM_NULL" << endl;
7 }

```

The destructor for each MPI user level object does *not* invoke the corresponding `MPI*_FREE` function (if it exists).

Rationale. `MPI*_FREE` functions are not automatically invoked for the following reasons:

1. Automatic destruction contradicts the shallow-copy semantics of the MPI classes.
2. The model put forth in MPI makes memory allocation and deallocation the responsibility of the user, not the implementation.
3. Calling `MPI*_FREE` upon destruction could have unintended side effects, including triggering collective operations (this also affects the copy, assignment, and construction semantics). In the following example, we would want neither `foo_comm` nor `bar_comm` to automatically invoke `MPI*_FREE` upon exit from the function.

```

32 void example_function()
33 {
34     MPI::Intracomm foo_comm(MPI::COMM_WORLD), bar_comm;
35     bar_comm = MPI::COMM_WORLD.Dup();
36     // rest of function
37 }

```

(End of rationale.)

Copy / Assignment The copy constructor and assignment operator are prototyped as follows:

```

43 MPI::<CLASS>(const MPI::<CLASS>& data)
44
45 MPI::<CLASS>& MPI::<CLASS>::operator=(const MPI::<CLASS>& data)

```

In terms of copying and assignment, opaque MPI user level objects behave like handles. Copy constructors perform handle-based (shallow) copies. `MPI::Status` objects are exceptions to this rule. These objects perform deep copies for assignment and copy construction.

Advice to implementors. Each MPI user level object is likely to contain, by value or by reference, implementation-dependent state information. The assignment and copying of MPI object handles may simply copy this value (or reference). (*End of advice to implementors.*)

Example 13.3 Example using assignment operator. In this example, `MPI::Intracomm::Dup()` is *not* called for `foo_comm`. The object `foo_comm` is simply an alias for `MPI::COMM_WORLD`. But `bar_comm` is created with a call to `MPI::Intracomm::Dup()` and is therefore a different communicator than `foo_comm` (and thus different from `MPI::COMM_WORLD`). `baz_comm` becomes an alias for `bar_comm`. If one of `bar_comm` or `baz_comm` is freed with `MPI_COMM_FREE` it will be set to `MPI::COMM_NULL`. The state of the other handle will be undefined — it will be invalid, but not necessarily set to `MPI::COMM_NULL`.

```
MPI::Intracomm foo_comm, bar_comm, baz_comm;

foo_comm = MPI::COMM_WORLD;
bar_comm = MPI::COMM_WORLD.Dup();
baz_comm = bar_comm;
```

Comparison The comparison operators are prototyped as follows:

```
bool MPI::<CLASS>::operator==(const MPI::<CLASS>& data) const
bool MPI::<CLASS>::operator!=(const MPI::<CLASS>& data) const
```

The member function `operator==()` returns `true` only when the handles reference the same internal MPI object, `false` otherwise. `operator!=()` returns the boolean complement of `operator==()`. However, since the `Status` class is not a handle to an underlying MPI object, it does not make sense to compare `Status` instances. Therefore, the `operator==()` and `operator!=()` functions are not defined on the `Status` class.

Constants Constants are singleton objects and are declared `const`. Note that not all globally defined MPI objects are constant. For example, `MPI::COMM_WORLD` and `MPI::COMM_SELF` are not `const`.

13.1.6 C++ Datatypes

Table 13.1 lists all of the C++ predefined MPI datatypes and their corresponding C and C++ datatypes, Table 13.2 lists all of the Fortran predefined MPI datatypes and their corresponding Fortran 77 datatypes. Table 13.3 lists the C++ names for all other MPI datatypes.

`MPI::BYTE` and `MPI::PACKED` conform to the same restrictions as `MPI_BYTE` and `MPI_PACKED`, listed in Sections 3.2.2 and 3.13 of MPI-1, respectively.

The following table defines groups of MPI predefined datatypes:

| | |
|------------|--|
| C integer: | <code>MPI::INT</code> , <code>MPI::LONG</code> , <code>MPI::SHORT</code> , |
| | <code>MPI::UNSIGNED_SHORT</code> , <code>MPI::UNSIGNED</code> , |
| | <code>MPI::UNSIGNED_LONG</code> , <code>MPI::SIGNED_CHAR</code> , |
| | <code>MPI::UNSIGNED_CHAR</code> |

| MPI datatype | C datatype | C++ datatype |
|--------------------------|----------------|----------------------|
| MPI::CHAR | char | char |
| MPI::WCHAR | wchar_t | wchar_t |
| MPI::SHORT | signed short | signed short |
| MPI::INT | signed int | signed int |
| MPI::LONG | signed long | signed long |
| MPI::SIGNED_CHAR | signed char | signed char |
| MPI::UNSIGNED_CHAR | unsigned char | unsigned char |
| MPI::UNSIGNED_SHORT | unsigned short | unsigned short |
| MPI::UNSIGNED | unsigned int | unsigned int |
| MPI::UNSIGNED_LONG | unsigned long | unsigned long int |
| MPI::FLOAT | float | float |
| MPI::DOUBLE | double | double |
| MPI::LONG_DOUBLE | long double | long double |
| MPI::BOOL | | bool |
| MPI::COMPLEX | | Complex<float> |
| MPI::DOUBLE_COMPLEX | | Complex<double> |
| MPI::LONG_DOUBLE_COMPLEX | | Complex<long double> |
| MPI::BYTE | | |
| MPI::PACKED | | |

Table 13.1: C++ names for the MPI C and C++ predefined datatypes, and their corresponding C/C++ datatypes.

| MPI datatype | Fortran datatype |
|-----------------------|------------------|
| MPI::CHARACTER | CHARACTER(1) |
| MPI::INTEGER | INTEGER |
| MPI::REAL | REAL |
| MPI::DOUBLE_PRECISION | DOUBLE PRECISION |
| MPI::LOGICAL | LOGICAL |
| MPI::F_COMPLEX | COMPLEX |
| MPI::BYTE | |
| MPI::PACKED | |

Table 13.2: C++ names for the MPI Fortran predefined datatypes, and their corresponding Fortran 77 datatypes.

| MPI datatype | Description |
|--------------------------|------------------------|
| MPI::FLOAT_INT | C/C++ reduction type |
| MPI::DOUBLE_INT | C/C++ reduction type |
| MPI::LONG_INT | C/C++ reduction type |
| MPI::TWOINT | C/C++ reduction type |
| MPI::SHORT_INT | C/C++ reduction type |
| MPI::LONG_DOUBLE_INT | C/C++ reduction type |
| MPI::LONG_LONG | Optional C/C++ type |
| MPI::UNSIGNED_LONG_LONG | Optional C/C++ type |
| MPI::TWOREAL | Fortran reduction type |
| MPI::TWODOUBLE_PRECISION | Fortran reduction type |
| MPI::TWOINTEGER | Fortran reduction type |
| MPI::F_DOUBLE_COMPLEX | Optional Fortran type |
| MPI::INTEGER1 | Explicit size type |
| MPI::INTEGER2 | Explicit size type |
| MPI::INTEGER4 | Explicit size type |
| MPI::INTEGER8 | Explicit size type |
| MPI::REAL4 | Explicit size type |
| MPI::REAL8 | Explicit size type |
| MPI::REAL16 | Explicit size type |

Table 13.3: C++ names for other MPI datatypes. Implementations may also define other optional types (e.g., `MPI::INTEGER8`).

| | |
|------------------|--|
| Fortran integer: | <code>MPI::INTEGER</code> |
| Floating point: | <code>MPI::FLOAT</code> , <code>MPI::DOUBLE</code> , <code>MPI::REAL</code> , <code>MPI::DOUBLE_PRECISION</code> , <code>MPI::LONG_DOUBLE</code> |
| Logical: | <code>MPI::LOGICAL</code> , <code>MPI::BOOL</code> |
| Complex: | <code>MPI::F_COMPLEX</code> , <code>MPI::COMPLEX</code> , <code>MPI::F_DOUBLE_COMPLEX</code> , <code>MPI::DOUBLE_COMPLEX</code> , <code>MPI::LONG_DOUBLE_COMPLEX</code> |
| Byte: | <code>MPI::BYTE</code> |

Valid datatypes for each reduction operation are specified below in terms of the groups defined above.

| Op | Allowed Types |
|---|---|
| <code>MPI::MAX</code> , <code>MPI::MIN</code> | C integer, Fortran integer, Floating point |
| <code>MPI::SUM</code> , <code>MPI::PROD</code> | C integer, Fortran integer, Floating point, Complex |
| <code>MPI::LAND</code> , <code>MPI::LOR</code> , <code>MPI::LXOR</code> | C integer, Logical |
| <code>MPI::BAND</code> , <code>MPI::BOR</code> , <code>MPI::BXOR</code> | C integer, Fortran integer, Byte |

`MPI::MINLOC` and `MPI::MAXLOC` perform just as their C and Fortran counterparts; see Section 4.9.3 in MPI-1.

13.1.7 Communicators

The `MPI::Comm` class hierarchy makes explicit the different kinds of communicators implicitly defined by MPI and allows them to be strongly typed. Since the original design of MPI defined only one type of handle for all types of communicators, the following clarifications are provided for the C++ design.

Types of communicators There are five different types of communicators: `MPI::Comm`, `MPI::Intercomm`, `MPI::Intracomm`, `MPI::Cartcomm`, and `MPI::Graphcomm`. `MPI::Comm` is the abstract base communicator class, encapsulating the functionality common to all MPI communicators. `MPI::Intercomm` and `MPI::Intracomm` are derived from `MPI::Comm`. `MPI::Cartcomm` and `MPI::Graphcomm` are derived from `MPI::Intracomm`.

Advice to users. Initializing a derived class with an instance of a base class is not legal in C++. For instance, it is not legal to initialize a `Cartcomm` from an `Intracomm`. Moreover, because `MPI::Comm` is an abstract base class, it is non-instantiable, so that it is not possible to have an object of class `MPI::Comm`. However, it is possible to have a reference or a pointer to an `MPI::Comm`.

Example 13.4 The following code is erroneous.

```
Intracomm intra = MPI::COMM_WORLD.Dup();
Cartcomm cart(intra);           // This is erroneous
```

(End of advice to users.)

MPI::COMM_NULL The specific type of `MPI::COMM_NULL` is implementation dependent. `MPI::COMM_NULL` must be able to be used in comparisons and initializations with all types of communicators. `MPI::COMM_NULL` must also be able to be passed to a function that expects a communicator argument in the parameter list (provided that `MPI::COMM_NULL` is an allowed value for the communicator argument).

Rationale. There are several possibilities for implementation of `MPI::COMM_NULL`. Specifying its required behavior, rather than its realization, provides maximum flexibility to implementors. *(End of rationale.)*

Example 13.5 The following example demonstrates the behavior of assignment and comparison using `MPI::COMM_NULL`.

```
MPI::Intercomm comm;
comm = MPI::COMM_NULL;           // assign with COMM_NULL
if (comm == MPI::COMM_NULL)     // true
    cout << "comm is NULL" << endl;
if (MPI::COMM_NULL == comm)     // note -- a different function!
    cout << "comm is still NULL" << endl;
```

`Dup()` is not defined as a member function of `MPI::Comm`, but it is defined for the derived classes of `MPI::Comm`. `Dup()` is not virtual and it returns its `OUT/` parameter by value.

`MPI::Comm::Clone()` The C++ language interface for MPI includes a new function `Clone()`. `MPI::Comm::Clone()` is a pure virtual function. For the derived communicator classes, `Clone()` behaves like `Dup()` except that it returns a new object by reference. The `Clone()` functions are prototyped as follows:

```
Comm& Comm::Clone() const = 0
Intracomm& Intracomm::Clone() const
Intercomm& Intercomm::Clone() const
Cartcomm& Cartcomm::Clone() const
Graphcomm& Graphcomm::Clone() const
```

Rationale. `Clone()` provides the “virtual dup” functionality that is expected by C++ programmers and library writers. Since `Clone()` returns a new object by reference, users are responsible for eventually deleting the object. A new name is introduced rather than changing the functionality of `Dup()`. (*End of rationale.*)

Advice to implementors. Within their class declarations, prototypes for `Clone()` and `Dup()` would look like the following:

```
namespace MPI {
  class Comm {
    virtual Comm& Clone() const = 0;
  };
  class Intracomm : public Comm {
    Intracomm Dup() const { ... };
    virtual Intracomm& Clone() const { ... };
  };
  class Intercomm : public Comm {
    Intercomm Dup() const { ... };
    virtual Intercomm& Clone() const { ... };
  };
  // Cartcomm and Graphcomm are similarly defined
};
```

Compilers that do not support the variable return type feature of virtual functions may return a reference to `Comm`. Users can cast to the appropriate type as necessary. (*End of advice to implementors.*)

13.1.8 Exceptions

The C++ language interface for MPI includes the predefined error handler `MPI::ERRORS_THROW_EXCEPTIONS` for use with the `Set_errhandler()` member functions. `MPI::ERRORS_THROW_EXCEPTIONS` can only be set or retrieved by C++ functions. If a non-C++ program causes an error that invokes the `MPI::ERRORS_THROW_EXCEPTIONS` error handler, the exception will pass up the calling stack until C++ code can catch it. If there is no C++ code to catch it, the behavior is undefined. In a multi-threaded environment or if

1 a non-blocking MPI call throws an exception while making progress in the background, the
 2 behavior is implementation dependent.

3 The error handler `MPI::ERRORS_THROW_EXCEPTIONS` causes an `MPI::Exception` to be
 4 thrown for any MPI result code other than `MPI::SUCCESS`. The public interface to
 5 `MPI::Exception` class is defined as follows:

```
6 namespace MPI {
7     class Exception {
8     public:
9
10
11         Exception(int error_code);
12
13         int Get_error_code() const;
14         int Get_error_class() const;
15         const char *Get_error_string() const;
16     };
17 };
```

18 *Advice to implementors.*

19 The exception will be thrown within the body of `MPI::ERRORS_THROW_EXCEPTIONS`. It
 20 is expected that control will be returned to the user when the exception is thrown.
 21 Some MPI functions specify certain return information in their parameters in the case
 22 of an error and `MPI_ERRORS_RETURN` is specified. The same type of return information
 23 must be provided when exceptions are thrown.

24 For example, `MPI_WAITALL` puts an error code for each request in the corresponding
 25 entry in the status array and returns `MPI_ERR_IN_STATUS`. When using
 26 `MPI::ERRORS_THROW_EXCEPTIONS`, it is expected that the error codes in the status
 27 array will be set appropriately before the exception is thrown.

28 *(End of advice to implementors.)*

31 13.1.9 Mixed-Language Operability

32 The C++ language interface provides functions listed below for mixed-language operability.
 33 These functions provide for a seamless transition between C and C++. For the case where
 34 the C++ class corresponding to `<CLASS>` has derived classes, functions are also provided
 35 for converting between the derived classes and the C `MPI_<CLASS>`.

```
36 MPI_<CLASS>& MPI_<CLASS>::operator=(const MPI_<CLASS>& data)
```

```
37 MPI_<CLASS>(const MPI_<CLASS>& data)
```

```
38 MPI_<CLASS>::operator MPI_<CLASS>() const
```

39 These functions are discussed in Section 13.3.4.

40 13.1.10 Profiling

41 This section specifies the requirements of a C++ profiling interface to MPI.

Advice to implementors. Since the main goal of profiling is to intercept function calls from user code, it is the implementor’s decision how to layer the underlying implementation to allow function calls to be intercepted and profiled. If an implementation of the MPI C++ bindings is layered on top of MPI bindings in another language (such as C), or if the C++ bindings are layered on top of a profiling interface in another language, no extra profiling interface is necessary because the underlying MPI implementation already meets the MPI profiling interface requirements.

Native C++ MPI implementations that do not have access to other profiling interfaces must implement an interface that meets the requirements outlined in this section.

High quality implementations can implement the interface outlined in this section in order to promote portable C++ profiling libraries. Implementors may wish to provide an option whether to build the C++ profiling interface or not; C++ implementations that are already layered on top of bindings in another language or another profiling interface will have to insert a third layer to implement the C++ profiling interface. (*End of advice to implementors.*)

To meet the requirements of the C++ MPI profiling interface, an implementation of the MPI functions *must*:

1. Provide a mechanism through which all of the MPI defined functions may be accessed with a name shift. Thus all of the MPI functions (which normally start with the prefix “MPI:.”) should also be accessible with the prefix “PMPI:.”
2. Ensure that those MPI functions which are not replaced may still be linked into an executable image without causing name clashes.
3. Document the implementation of different language bindings of the MPI interface if they are layered on top of each other, so that profiler developer knows whether they must implement the profile interface for each binding, or can economize by implementing it only for the lowest level routines.
4. Where the implementation of different language bindings is done through a layered approach (e.g., the C++ binding is a set of “wrapper” functions which call the C implementation), ensure that these wrapper functions are separable from the rest of the library.

This is necessary to allow a separate profiling library to be correctly implemented, since (at least with Unix linker semantics) the profiling library must contain these wrapper functions if it is to perform as expected. This requirement allows the author of the profiling library to extract these functions from the original MPI library and add them into the profiling library without bringing along any other unnecessary code.

5. Provide a no-op routine `MPI:Pcontrol` in the MPI library.

Advice to implementors. There are (at least) two apparent options for implementing the C++ profiling interface: inheritance or caching. An inheritance-based approach may not be attractive because it may require a virtual inheritance implementation of the communicator classes. Thus, it is most likely that implementors will cache `PMPI` objects on their corresponding `MPI` objects. The caching scheme is outlined below.

The “real” entry points to each routine can be provided within a `namespace PMPI`. The non-profiling version can then be provided within a `namespace MPI`.

Caching instances of PMPI objects in the MPI handles provides the “has a” relationship that is necessary to implement the profiling scheme.

Each instance of an MPI object simply “wraps up” an instance of a PMPI object. MPI objects can then perform profiling actions before invoking the corresponding function in their internal PMPI object.

The key to making the profiling work by simply re-linking programs is by having a header file that *declares* all the MPI functions. The functions must be *defined* elsewhere, and compiled into a library. MPI constants should be declared `extern` in the MPI namespace. For example, the following is an excerpt from a sample `mpi.h` file:

Example 13.6 Sample `mpi.h` file.

```

namespace PMPI {
    class Comm {
    public:
        int Get_size() const;
    };
    // etc.
};

namespace MPI {
public:
    class Comm {
    public:
        int Get_size() const;

    private:
        PMPI::Comm pmpi_comm;
    };
};

```

Note that all constructors, the assignment operator, and the destructor in the MPI class will need to initialize/destroy the internal PMPI object as appropriate.

The definitions of the functions must be in separate object files; the PMPI class member functions and the non-profiling versions of the MPI class member functions can be compiled into `libmpi.a`, while the profiling versions can be compiled into `libpmpi.a`. Note that the PMPI class member functions and the MPI constants must be in different object files than the non-profiling MPI class member functions in the `libmpi.a` library to prevent multiple definitions of MPI class member function names when linking both `libmpi.a` and `libpmpi.a`. For example:

Example 13.7 `pmpi.cc`, to be compiled into `libmpi.a`.

```

int PMPI::Comm::Get_size() const
{
    // Implementation of MPI_COMM_SIZE
}

```

Example 13.8 constants.cc, to be compiled into libmpi.a.

```
const MPI::Intracomm MPI::COMM_WORLD;
```

Example 13.9 mpi_no_profile.cc, to be compiled into libmpi.a.

```

int MPI::Comm::Get_size() const
{
    return pmpi_comm.Get_size();
}

```

Example 13.10 mpi_profile.cc, to be compiled into libpmpi.a.

```

int MPI::Comm::Get_size() const
{
    // Do profiling stuff
    int ret = pmpi_comm.Get_size();
    // More profiling stuff
    return ret;
}

```

(End of advice to implementors.)

13.2 Fortran Support

13.2.1 Overview

Fortran 90 is the current international Fortran standard. MPI-2 Fortran bindings are Fortran 90 bindings that in most cases are “Fortran 77 friendly.” That is, with few exceptions (e.g., KIND-parameterized types, and the `mpi` module, both of which can be avoided) Fortran 77 compilers should be able to compile MPI programs.

Rationale. Fortran 90 contains numerous features designed to make it a more “modern” language than Fortran 77. It seems natural that MPI should be able to take advantage of these new features with a set of bindings tailored to Fortran 90. MPI does not (yet) use many of these features because of a number of technical difficulties.
(End of rationale.)

MPI defines two levels of Fortran support, described in Sections 13.2.3 and 13.2.4. A third level of Fortran support is envisioned, but is deferred to future standardization efforts. In the rest of this section, “Fortran” shall refer to Fortran 90 (or its successor) unless qualified.

- 1 1. **Basic Fortran Support** An implementation with this level of Fortran support pro-
2 vides the original Fortran bindings specified in MPI-1, with small additional require-
3 ments specified in Section 13.2.3.
- 4
5 2. **Extended Fortran Support** An implementation with this level of Fortran sup-
6 port provides Basic Fortran Support plus additional features that specifically support
7 Fortran 90, as described in Section 13.2.4.

8
9 A compliant MPI-2 implementation providing a Fortran interface must provide Ex-
10 tended Fortran Support unless the target compiler does not support modules or KIND-
11 parameterized types.

12 13.2.2 Problems With Fortran Bindings for MPI

13
14 This section discusses a number of problems that may arise when using MPI in a Fortran
15 program. It is intended as advice to users, and clarifies how MPI interacts with Fortran. It
16 does not add to the standard, but is intended to clarify the standard.

17 As noted in the original MPI specification, the interface violates the Fortran standard
18 in several ways. While these cause few problems for Fortran 77 programs, they become
19 more significant for Fortran 90 programs, so that users must exercise care when using new
20 Fortran 90 features. The violations were originally adopted and have been retained because
21 they are important for the usability of MPI. The rest of this section describes the potential
22 problems in detail. It supersedes and replaces the discussion of Fortran bindings in the
23 original MPI specification (for Fortran 90, not Fortran 77).

24 The following MPI features are inconsistent with Fortran 90.

- 25
26 1. An MPI subroutine with a choice argument may be called with different argument
27 types.
- 28
29 2. An MPI subroutine with an assumed-size dummy argument may be passed an actual
30 scalar argument.
- 31
32 3. Many MPI routines assume that actual arguments are passed by address and that
33 arguments are not copied on entrance to or exit from the subroutine.
- 34
35 4. An MPI implementation may read or modify user data (e.g., communication buffers
36 used by nonblocking communications) concurrently with a user program that is exe-
37 cuting outside of MPI calls.
- 38
39 5. Several named “constants,” such as MPI_BOTTOM, MPI_IN_PLACE,
40 MPI_STATUS_IGNORE, MPI_STATUSES_IGNORE, MPI_ERRCODES_IGNORE,
41 MPI_ARGV_NULL, and MPI_ARGVS_NULL are not ordinary Fortran constants and require
42 a special implementation. See Section 2.5.4 on page 14 for more information.
- 43
44 6. The memory allocation routine MPI_ALLOC_MEM can’t be usefully used in Fortran
45 without a language extension that allows the allocated memory to be associated with
46 a Fortran variable.

46 MPI-1 contained several routines that take address-sized information as input or return
47 address-sized information as output. In C such arguments were of type MPI_Aint and in
48 Fortran of type INTEGER. On machines where integers are smaller than addresses, these

routines can lose information. In MPI-2 the use of these functions has been deprecated and they have been replaced by routines taking `INTEGER` arguments of `KIND=MPI_ADDRESS_KIND`. A number of new MPI-2 functions also take `INTEGER` arguments of non-default `KIND`. See Section 2.6 on page 15 and Section 3.12.1 on page 76 for more information.

Problems Due to Strong Typing

All MPI functions with choice arguments associate actual arguments of different Fortran datatypes with the same dummy argument. This is not allowed by Fortran 77, and in Fortran 90 is technically only allowed if the function is overloaded with a different function for each type. In C, the use of `void*` formal arguments avoids these problems.

The following code fragment is technically illegal and may generate a compile-time error.

```
integer i(5)
real    x(5)
...
call mpi_send(x, 5, MPI_REAL, ...)
call mpi_send(i, 5, MPI_INTEGER, ...)
```

In practice, it is rare for compilers to do more than issue a warning, though there is concern that Fortran 90 compilers are more likely to return errors.

It is also technically illegal in Fortran to pass a scalar actual argument to an array dummy argument. Thus the following code fragment may generate an error since the `buf` argument to `MPI_SEND` is declared as an assumed-size array `<type> buf(*)`.

```
integer a
call mpi_send(a, 1, MPI_INTEGER, ...)
```

Advice to users. In the event that you run into one of the problems related to type checking, you may be able to work around it by using a compiler flag, by compiling separately, or by using an MPI implementation with Extended Fortran Support as described in Section 13.2.4. An alternative that will usually work with variables local to a routine but not with arguments to a function or subroutine is to use the `EQUIVALENCE` statement to create another variable with a type accepted by the compiler. (*End of advice to users.*)

Problems Due to Data Copying and Sequence Association

Implicit in MPI is the idea of a contiguous chunk of memory accessible through a linear address space. MPI copies data to and from this memory. An MPI program specifies the location of data by providing memory addresses and offsets. In the C language, sequence association rules plus pointers provide all the necessary low-level structure.

In Fortran 90, user data is not necessarily stored contiguously. For example, the array section `A(1:N:2)` involves only the elements of `A` with indices 1, 3, 5, The same is true for a pointer array whose target is such a section. Most compilers ensure that an array that is a dummy argument is held in contiguous memory if it is declared with an explicit shape (e.g., `B(N)`) or is of assumed size (e.g., `B(*)`). If necessary, they do this by making a copy of the array into contiguous memory. Both Fortran 77 and Fortran 90 are carefully worded to allow such copying to occur, but few Fortran 77 compilers do it.¹

¹Technically, the Fortran standards are worded to allow non-contiguous storage of any array data.

1 Because MPI dummy buffer arguments are assumed-size arrays, this leads to a serious
 2 problem for a non-blocking call: the compiler copies the temporary array back on return
 3 but MPI continues to copy data to the memory that held it. For example, consider the
 4 following code fragment:

```
5
6       real a(100)
7       call MPI_IRECV(a(1:100:2), MPI_REAL, 50, ...)
```

8 Since the first dummy argument to MPI_IRECV is an assumed-size array (<type> buf(*)),
 9 the array section a(1:100:2) is copied to a temporary before being passed to MPI_IRECV,
 10 so that it is contiguous in memory. MPI_IRECV returns immediately, and data is copied
 11 from the temporary back into the array a. Sometime later, MPI may write to the address of
 12 the deallocated temporary. Copying is also a problem for MPI_ISEND since the temporary
 13 array may be deallocated before the data has all been sent from it.

14 Most Fortran 90 compilers do not make a copy if the actual argument is the whole of
 15 an explicit-shape or assumed-size array or is a ‘simple’ section such as A(1:N) of such an
 16 array. (We define ‘simple’ more fully in the next paragraph.) Also, many compilers treat
 17 allocatable arrays the same as they treat explicit-shape arrays in this regard (though we
 18 know of one that does not). However, the same is not true for assumed-shape and pointer
 19 arrays; since they may be discontinuous, copying is often done. It is this copying that causes
 20 problems for MPI as described in the previous paragraph.

21 Our formal definition of a ‘simple’ array section is

```
22
23       name ( [:,]... [<subscript>]:<subscript> [,<subscript>]... )
```

24 That is, there are zero or more dimensions that are selected in full, then one dimension
 25 selected without a stride, then zero or more dimensions that are selected with a simple
 26 subscript. Examples are

```
27
28       A(1:N), A(:,N), A(:,1:N,1), A(1:6,N), A(:, :, 1:N)
```

29 Because of Fortran’s column-major ordering, where the first index varies fastest, a simple
 30 section of a contiguous array will also be contiguous.²

31 The same problem can occur with a scalar argument. Some compilers, even for Fortran
 32 77, make a copy of some scalar dummy arguments within a called procedure. That this can
 33 cause a problem is illustrated by the example

```
34
35
36       call user1(a,rq)
37       call MPI_WAIT(rq,status,ierr)
38       write (*,*) a
39
40       subroutine user1(buf,request)
41       call MPI_IRECV(buf,...,request,...)
42       end
```

43 ²To keep the definition of ‘simple’ simple, we have chosen to require all but one of the section subscripts
 44 to be without bounds. A colon without bounds makes it obvious both to the compiler and to the reader
 45 that the whole of the dimension is selected. It would have been possible to allow cases where the whole
 46 dimension is selected with one or two bounds, but this means for the reader that the array declaration or
 47 most recent allocation has to be consulted and for the compiler that a run-time check may be required.

If `a` is copied, `MPI_RECV` will alter the copy when it completes the communication and will not alter `a` itself.

Note that copying will almost certainly occur for an argument that is a non-trivial expression (one with at least one operator or function call), a section that does not select a contiguous part of its parent (e.g., `A(1:n:2)`), a pointer whose target is such a section, or an assumed-shape array that is (directly or indirectly) associated with such a section.

If there is a compiler option that inhibits copying of arguments, in either the calling or called procedure, this should be employed.

If a compiler makes copies in the calling procedure of arguments that are explicit-shape or assumed-size arrays, simple array sections of such arrays, or scalars, and if there is no compiler option to inhibit this, then the compiler cannot be used for applications that use `MPI_GET_ADDRESS`, or any non-blocking MPI routine. If a compiler copies scalar arguments in the called procedure and there is no compiler option to inhibit this, then this compiler cannot be used for applications that use memory references across subroutine calls as in the example above.

Special Constants

MPI requires a number of special “constants” that cannot be implemented as normal Fortran constants, including `MPI_BOTTOM`, `MPI_STATUS_IGNORE`, `MPI_IN_PLACE`, `MPI_STATUSES_IGNORE` and `MPI_ERRCODES_IGNORE`. In C, these are implemented as constant pointers, usually as `NULL` and are used where the function prototype calls for a pointer to a variable, not the variable itself.

In Fortran the implementation of these special constants may require the use of language constructs that are outside the Fortran standard. Using special values for the constants (e.g., by defining them through `parameter` statements) is not possible because an implementation cannot distinguish these values from legal data. Typically these constants are implemented as predefined static variables (e.g., a variable in an MPI-declared `COMMON` block), relying on the fact that the target compiler passes data by address. Inside the subroutine, this address can be extracted by some mechanism outside the Fortran standard (e.g., by Fortran extensions or by implementing the function in C).

Fortran 90 Derived Types

MPI does not explicitly support passing Fortran 90 derived types to choice dummy arguments. Indeed, for MPI implementations that provide explicit interfaces through the `mpi` module a compiler will reject derived type actual arguments at compile time. Even when no explicit interfaces are given, users should be aware that Fortran 90 provides no guarantee of sequence association for derived types or arrays of derived types. For instance, an array of a derived type consisting of two elements may be implemented as an array of the first elements followed by an array of the second. Use of the `SEQUENCE` attribute may help here, somewhat.

The following code fragment shows one possible way to send a derived type in Fortran. The example assumes that all data is passed by address.

```

type mytype
  integer i
  real x
  double precision d

```

```

1      end type mytype
2
3      type(mytype) foo
4      integer blocklen(3), type(3)
5      integer(MPI_ADDRESS_KIND) disp(3), base
6
7      call MPI_GET_ADDRESS(foo%i, disp(1), ierr)
8      call MPI_GET_ADDRESS(foo%x, disp(2), ierr)
9      call MPI_GET_ADDRESS(foo%d, disp(3), ierr)
10
11     base = disp(1)
12     disp(1) = disp(1) - base
13     disp(2) = disp(2) - base
14     disp(3) = disp(3) - base
15
16     blocklen(1) = 1
17     blocklen(2) = 1
18     blocklen(3) = 1
19
20     type(1) = MPI_INTEGER
21     type(2) = MPI_REAL
22     type(3) = MPI_DOUBLE_PRECISION
23
24     call MPI_TYPE_CREATE_STRUCT(3, blocklen, disp, type, newtype, ierr)
25     call MPI_TYPE_COMMIT(newtype, ierr)
26
27     ! unpleasant to send foo%i instead of foo, but it works for scalar
28     ! entities of type mytype
29     call MPI_SEND(foo%i, 1, newtype, ...)
30
31

```

32 A Problem with Register Optimization

33 MPI provides operations that may be hidden from the user code and run concurrently with
34 it, accessing the same memory as user code. Examples include the data transfer for an
35 MPI_RECV. The optimizer of a compiler will assume that it can recognize periods when a
36 copy of a variable can be kept in a register without reloading from or storing to memory.
37 When the user code is working with a register copy of some variable while the hidden
38 operation reads or writes the memory copy, problems occur. This section discusses register
39 optimization pitfalls.

40 When a variable is local to a Fortran subroutine (i.e., not in a module or COMMON
41 block), the compiler will assume that it cannot be modified by a called subroutine unless it
42 is an actual argument of the call. In the most common linkage convention, the subroutine
43 is expected to save and restore certain registers. Thus, the optimizer will assume that a
44 register which held a valid copy of such a variable before the call will still hold a valid copy
45 on return.

46 Normally users are not afflicted with this. But the user should pay attention to this
47 section if in his/her program a buffer argument to an MPI_SEND, MPI_RECV etc., uses
48

a name which hides the actual variables involved. `MPI_BOTTOM` with an `MPI_Datatype` containing absolute addresses is one example. Creating a datatype which uses one variable as an anchor and brings along others by using `MPI_GET_ADDRESS` to determine their offsets from the anchor is another. The anchor variable would be the only one mentioned in the call. Also attention must be paid if MPI operations are used that run in parallel with the user's application.

The following example shows what Fortran compilers are allowed to do.

| | |
|--|---|
| This source ... | can be compiled as: |
| <code>call MPI_GET_ADDRESS(buf,bufaddr, ierror)</code> | <code>call MPI_GET_ADDRESS(buf,...)</code> |
| <code>call MPI_TYPE_CREATE_STRUCT(1,1, bufaddr, MPI_REAL,type,ierror)</code> | <code>call MPI_TYPE_CREATE_STRUCT(...)</code> |
| <code>call MPI_TYPE_COMMIT(type,ierror)</code> | <code>call MPI_TYPE_COMMIT(...)</code> |
| <code>val_old = buf</code> | <code>register = buf</code> |
| | <code>val_old = register</code> |
| <code>call MPI_RECV(MPI_BOTTOM,1,type,...)</code> | <code>call MPI_RECV(MPI_BOTTOM,...)</code> |
| <code>val_new = buf</code> | <code>val_new = register</code> |

The compiler does not invalidate the register because it cannot see that `MPI_RECV` changes the value of `buf`. The access of `buf` is hidden by the use of `MPI_GET_ADDRESS` and `MPI_BOTTOM`.

The next example shows extreme, but allowed, possibilities.

| | | |
|--|--|--|
| Source | compiled as | or compiled as |
| <code>call MPI_IRecv(buf,..req)</code> | <code>call MPI_IRecv(buf,..req)</code> | <code>call MPI_IRecv(buf,..req)</code> |
| | <code>register = buf</code> | <code>b1 = buf</code> |
| <code>call MPI_WAIT(req,..)</code> | <code>call MPI_WAIT(req,..)</code> | <code>call MPI_WAIT(req,..)</code> |
| <code>b1 = buf</code> | <code>b1 := register</code> | |

`MPI_WAIT` on a concurrent thread modifies `buf` between the invocation of `MPI_IRecv` and the finish of `MPI_WAIT`. But the compiler cannot see any possibility that `buf` can be changed after `MPI_IRecv` has returned, and may schedule the load of `buf` earlier than typed in the source. It has no reason to avoid using a register to hold `buf` across the call to `MPI_WAIT`. It also may reorder the instructions as in the case on the right.

To prevent instruction reordering or the allocation of a buffer in a register there are two possibilities in portable Fortran code:

- The compiler may be prevented from moving a reference to a buffer across a call to an MPI subroutine by surrounding the call by calls to an external subroutine with the buffer as an actual argument. Note that if the intent is declared in the external subroutine, it must be `OUT` or `INOUT`. The subroutine itself may have an empty body, but the compiler does not know this and has to assume that the buffer may be altered. For example, the above call of `MPI_RECV` might be replaced by

```
call DD(buf)
call MPI_RECV(MPI_BOTTOM,...)
```

```
1         call DD(buf)
```

2
3 with the separately compiled

```
4  
5         subroutine DD(buf)
6             integer buf
7         end
8
```

9 (assuming that `buf` has type `INTEGER`). The compiler may be similarly prevented from
10 moving a reference to a variable across a call to an MPI subroutine.

11
12 In the case of a non-blocking call, as in the above call of `MPI_WAIT`, no reference to
13 the buffer is permitted until it has been verified that the transfer has been completed.
14 Therefore, in this case, the extra call ahead of the MPI call is not necessary, i.e., the
15 call of `MPI_WAIT` in the example might be replaced by

```
16  
17         call MPI_WAIT(req,..)
18         call DD(buf)
19
```

- 20 • An alternative is to put the buffer or variable into a module or a common block and
21 access it through a `USE` or `COMMON` statement in each scope where it is referenced,
22 defined or appears as an actual argument in a call to an MPI routine. The compiler
23 will then have to assume that the MPI procedure (`MPI_RECV` in the above example)
24 may alter the buffer or variable, provided that the compiler cannot analyze that the
25 MPI procedure does not reference the module or common block.
26

27 In the longer term, the attribute `VOLATILE` is under consideration for Fortran 2000 and
28 would give the buffer or variable the properties needed, but it would inhibit optimization
29 of any code containing the buffer or variable.

30 In C, subroutines which modify variables that are not in the argument list will not cause
31 register optimization problems. This is because taking pointers to storage objects by using
32 the `&` operator and later referencing the objects by way of the pointer is an integral part of
33 the language. A C compiler understands the implications, so that the problem should not
34 occur, in general. However, some compilers do offer optional aggressive optimization levels
35 which may not be safe.
36

37 13.2.3 Basic Fortran Support

38 Because Fortran 90 is (for all practical purposes) a superset of Fortran 77, Fortran 90
39 (and future) programs can use the original Fortran interface. The following additional
40 requirements are added:
41

- 42 1. Implementations are required to provide the file `mpif.h`, as described in the original
43 MPI-1 specification.
- 44 2. `mpif.h` must be valid and equivalent for both fixed- and free- source form.
45
46
47
48

Advice to implementors. To make `mpif.h` compatible with both fixed- and free-source forms, to allow automatic inclusion by preprocessors, and to allow extended fixed-form line length, it is recommended that requirement two be met by constructing `mpif.h` without any continuation lines. This should be possible because `mpif.h` contains only declarations, and because common block declarations can be split among several lines. To support Fortran 77 as well as Fortran 90, it may be necessary to eliminate all comments from `mpif.h`. (*End of advice to implementors.*)

13.2.4 Extended Fortran Support

Implementations with Extended Fortran support must provide:

1. An `mpi` module
2. A new set of functions to provide additional support for Fortran intrinsic numeric types, including parameterized types: `MPI_SIZEOF`, `MPI_TYPE_MATCH_SIZE`, `MPI_TYPE_CREATE_F90_INTEGER`, `MPI_TYPE_CREATE_F90_REAL` and `MPI_TYPE_CREATE_F90_COMPLEX`. Parameterized types are Fortran intrinsic types which are specified using `KIND` type parameters. These routines are described in detail in Section 13.2.5.

Additionally, high quality implementations should provide a mechanism to prevent fatal type mismatch errors for MPI routines with choice arguments.

The `mpi` Module

An MPI implementation must provide a module named `mpi` that can be `USED` in a Fortran 90 program. This module must:

- Define all named MPI constants
- Declare MPI functions that return a value.

An MPI implementation may provide in the `mpi` module other features that enhance the usability of MPI while maintaining adherence to the standard. For example, it may:

- Provide interfaces for all or for a subset of MPI routines.
- Provide `INTENT` information in these interface blocks.

Advice to implementors. The appropriate `INTENT` may be different from what is given in the MPI generic interface. Implementations must choose `INTENT` so that the function adheres to the MPI standard. (*End of advice to implementors.*)

Rationale. The intent given by the MPI generic interface is not precisely defined and does not in all cases correspond to the correct Fortran `INTENT`. For instance, receiving into a buffer specified by a datatype with absolute addresses may require associating `MPI_BOTTOM` with a dummy `OUT` argument. Moreover, “constants” such as `MPI_BOTTOM` and `MPI_STATUS_IGNORE` are not constants as defined by Fortran, but “special addresses” used in a nonstandard way. Finally, the MPI-1 generic intent is changed in several places by MPI-2. For instance, `MPI_IN_PLACE` changes the sense of an `OUT` argument to be `INOUT`. (*End of rationale.*)

Applications may use either the `mpi` module or the `mpif.h` include file. An implementation may require use of the module to prevent type mismatch errors (see below).

Advice to users. It is recommended to use the `mpi` module even if it is not necessary to use it to avoid type mismatch errors on a particular system. Using a module provides several potential advantages over using an include file. (*End of advice to users.*)

It must be possible to link together routines some of which USE `mpi` and others of which INCLUDE `mpif.h`.

No Type Mismatch Problems for Subroutines with Choice Arguments

A high quality MPI implementation should provide a mechanism to ensure that MPI choice arguments do not cause fatal compile-time or run-time errors due to type mismatch. An MPI implementation may require applications to use the `mpi` module, or require that it be compiled with a particular compiler flag, in order to avoid type mismatch problems.

Advice to implementors. In the case where the compiler does not generate errors, nothing needs to be done to the existing interface. In the case where the compiler may generate errors, a set of overloaded functions may be used. See the paper of M. Hennecke [28]. Even if the compiler does not generate errors, explicit interfaces for all routines would be useful for detecting errors in the argument list. Also, explicit interfaces which give INTENT information can reduce the amount of copying for BUF(*) arguments. (*End of advice to implementors.*)

13.2.5 Additional Support for Fortran Numeric Intrinsic Types

The routines in this section are part of Extended Fortran Support described in Section 13.2.4.

MPI-1 provides a small number of named datatypes that correspond to named intrinsic types supported by C and Fortran. These include `MPI_INTEGER`, `MPI_REAL`, `MPI_INT`, `MPI_DOUBLE`, etc., as well as the optional types `MPI_REAL4`, `MPI_REAL8`, etc. There is a one-to-one correspondence between language declarations and MPI types.

Fortran (starting with Fortran 90) provides so-called KIND-parameterized types. These types are declared using an intrinsic type (one of `INTEGER`, `REAL`, `COMPLEX`, `LOGICAL` and `CHARACTER`) with an optional integer `KIND` parameter that selects from among one or more variants. The specific meaning of different `KIND` values themselves are implementation dependent and not specified by the language. Fortran provides the `KIND` selection functions `selected_real_kind` for `REAL` and `COMPLEX` types, and `selected_int_kind` for `INTEGER` types that allow users to declare variables with a minimum precision or number of digits. These functions provide a portable way to declare `KIND`-parameterized `REAL`, `COMPLEX` and `INTEGER` variables in Fortran. This scheme is backward compatible with Fortran 77. `REAL` and `INTEGER` Fortran variables have a default `KIND` if none is specified. Fortran `DOUBLE PRECISION` variables are of intrinsic type `REAL` with a non-default `KIND`. The following two declarations are equivalent:

```
double precision x
real(KIND(0.0d0)) x
```


MPI provides two orthogonal methods to communicate using numeric intrinsic types. The first method can be used when variables have been declared in a portable way — using default `KIND` or using `KIND` parameters obtained with the `selected_int_kind` or `selected_real_kind` functions. With this method, MPI automatically selects the correct data size (e.g., 4 or 8 bytes) and provides representation conversion in heterogeneous environments. The second method gives the user complete control over communication by exposing machine representations.

Parameterized Datatypes with Specified Precision and Exponent Range

MPI-1 provides named datatypes corresponding to standard Fortran 77 numeric types — `MPI_INTEGER`, `MPI_COMPLEX`, `MPI_REAL`, `MPI_DOUBLE_PRECISION` and `MPI_DOUBLE_COMPLEX`. MPI automatically selects the correct data size and provides representation conversion in heterogeneous environments. The mechanism described in this section extends this MPI-1 model to support portable parameterized numeric types.

The model for supporting portable parameterized types is as follows. Real variables are declared (perhaps indirectly) using `selected_real_kind(p, r)` to determine the `KIND` parameter, where `p` is decimal digits of precision and `r` is an exponent range. Implicitly MPI maintains a two-dimensional array of predefined MPI datatypes `D(p, r)`. `D(p, r)` is defined for each value of `(p, r)` supported by the compiler, including pairs for which one value is unspecified. Attempting to access an element of the array with an index `(p, r)` not supported by the compiler is erroneous. MPI implicitly maintains a similar array of `COMPLEX` datatypes. For integers, there is a similar implicit array related to `selected_int_kind` and indexed by the requested number of digits `r`. Note that the predefined datatypes contained in these implicit arrays are not the same as the named MPI datatypes `MPI_REAL`, etc., but a new set.

Advice to implementors. The above description is for explanatory purposes only. It is not expected that implementations will have such internal arrays. (*End of advice to implementors.*)

Advice to users. `selected_real_kind()` maps a large number of `(p,r)` pairs to a much smaller number of `KIND` parameters supported by the compiler. `KIND` parameters are not specified by the language and are not portable. From the language point of view intrinsic types of the same base type and `KIND` parameter are of the same type. In order to allow interoperability in a heterogeneous environment, MPI is more stringent. The corresponding MPI datatypes match if and only if they have the same `(p,r)` value (`REAL` and `COMPLEX`) or `r` value (`INTEGER`). Thus MPI has many more datatypes than there are fundamental language types. (*End of advice to users.*)

`MPI_TYPE_CREATE_F90_REAL(p, r, newtype)`

| | | |
|-----|----------------------|--|
| IN | <code>p</code> | precision, in decimal digits (integer) |
| IN | <code>r</code> | decimal exponent range (integer) |
| OUT | <code>newtype</code> | the requested MPI datatype (handle) |

`int MPI_Type_create_f90_real(int p, int r, MPI_Datatype *newtype)`

```
1 MPI_TYPE_CREATE_F90_REAL(P, R, NEWTYPE, IERROR)
```

```
2     INTEGER P, R, NEWTYPE, IERROR
```

```
3
4 static MPI::Datatype MPI::Datatype::Create_f90_real(int p, int r)
```

5 This function returns a predefined MPI datatype that matches a REAL variable of KIND
6 `selected_real_kind(p, r)`. In the model described above it returns a handle for the
7 element `D(p, r)`. Either `p` or `r` may be omitted from calls to `selected_real_kind(p, r)`
8 (but not both). Analogously, either `p` or `r` may be set to `MPI_UNDEFINED`. In communication,
9 an MPI datatype `A` returned by `MPI_TYPE_CREATE_F90_REAL` matches a datatype `B` if and
10 only if `B` was returned by `MPI_TYPE_CREATE_F90_REAL` called with the same values for `p`
11 and `r` or `B` is a duplicate of such a datatype. Restrictions on using the returned datatype
12 with the “external32” data representation are given on page 462.

13 It is erroneous to supply values for `p` and `r` not supported by the compiler.

```
14
15
16 MPI_TYPE_CREATE_F90_COMPLEX(p, r, newtype)
```

```
17     IN         p                precision, in decimal digits (integer)
```

```
18     IN         r                decimal exponent range (integer)
```

```
19     OUT        newtype          the requested MPI datatype (handle)
```

```
20
21
22 int MPI_Type_create_f90_complex(int p, int r, MPI_Datatype *newtype)
```

```
23 MPI_TYPE_CREATE_F90_COMPLEX(P, R, NEWTYPE, IERROR)
```

```
24     INTEGER P, R, NEWTYPE, IERROR
```

```
25
26 static MPI::Datatype MPI::Datatype::Create_f90_complex(int p, int r)
```

27 This function returns a predefined MPI datatype that matches a
28 `COMPLEX` variable of KIND `selected_real_kind(p, r)`. Either `p` or `r` may be omitted from
29 calls to `selected_real_kind(p, r)` (but not both). Analogously, either `p` or `r` may be set
30 to `MPI_UNDEFINED`. Matching rules for datatypes created by this function are analogous
31 to the matching rules for datatypes created by `MPI_TYPE_CREATE_F90_REAL`. Restrictions
32 on using the returned datatype with the “external32” data representation are given on page
33 462.

34 It is erroneous to supply values for `p` and `r` not supported by the compiler.

```
35
36
37 MPI_TYPE_CREATE_F90_INTEGER(r, newtype)
```

```
38     IN         r                decimal exponent range, i.e., number of decimal digits  
39                                     (integer)
```

```
40     OUT        newtype          the requested MPI datatype (handle)
```

```
41
42
43 int MPI_Type_create_f90_integer(int r, MPI_Datatype *newtype)
```

```
44 MPI_TYPE_CREATE_F90_INTEGER(R, NEWTYPE, IERROR)
```

```
45     INTEGER R, NEWTYPE, IERROR
```

```
46
47 static MPI::Datatype MPI::Datatype::Create_f90_integer(int r)
```

```
48
```

This function returns a predefined MPI datatype that matches a `INTEGER` variable of `KIND selected_int_kind(r)`. Matching rules for datatypes created by this function are analogous to the matching rules for datatypes created by `MPI_TYPE_CREATE_F90_REAL`. Restrictions on using the returned datatype with the “external32” data representation are given on page 462.

It is erroneous to supply a value for `r` that is not supported by the compiler.

Example:

```
integer      longtype, quadtype
integer, parameter :: long = selected_int_kind(15)
integer(long) ii(10)
real(selected_real_kind(30)) x(10)
call MPI_TYPE_CREATE_F90_INTEGER(15, longtype, ierror)
call MPI_TYPE_CREATE_F90_REAL(30, MPI_UNDEFINED, quadtype, ierror)
...

call MPI_SEND(ii, 10, longtype, ...)
call MPI_SEND(x, 10, quadtype, ...)
```

Advice to users. The datatypes returned by the above functions are predefined datatypes. They cannot be freed; they do not need to be committed; they can be used with predefined reduction operations. There are two situations in which they behave differently syntactically, but not semantically, from the MPI named predefined datatypes.

1. `MPI_TYPE_GET_ENVELOPE` returns special combinators that allow a program to retrieve the values of `p` and `r`.
2. Because the datatypes are not named, they cannot be used as compile-time initializers or otherwise accessed before a call to one of the `MPI_TYPE_CREATE_F90_` routines.

If a variable was declared specifying a non-default `KIND` value that was not obtained with `selected_real_kind()` or `selected_int_kind()`, the only way to obtain a matching MPI datatype is to use the size-based mechanism described in the next section.

(End of advice to users.)

Advice to implementors. An application may often repeat a call to `MPI_TYPE_CREATE_F90_XXXX` with the same combination of `(XXXX,p,r)`. The application is not allowed to free the returned predefined, unnamed datatype handles. To prevent the creation of a potentially huge amount of handles, the MPI implementation should return the same datatype handle for the same `(REAL/COMPLEX/INTEGER,p,r)` combination. Checking for the combination `(p,r)` in the preceding call to `MPI_TYPE_CREATE_F90_XXXX` and using a hash-table to find formerly generated handles should limit the overhead of finding a previously generated datatype with same combination of `(XXXX,p,r)`. *(End of advice to implementors.)*

Rationale. The `MPI_TYPE_CREATE_F90_REAL/COMPLEX/INTEGER` interface needs as input the original range and precision values to be able to define useful and compiler-independent external (Section 12.5.2 on page 416) or user-defined (Section 12.5.3 on page 418) data representations, and in order to be able to perform automatic and efficient data conversions in a heterogeneous environment. (*End of rationale.*)

We now specify how the datatypes described in this section behave when used with the “external32” external data representation described in Section 12.5.2 on page 416.

The external32 representation specifies data formats for integer and floating point values. Integer values are represented in two’s complement big-endian format. Floating point values are represented by one of three IEEE formats. These are the IEEE “Single,” “Double” and “Double Extended” formats, requiring 4, 8 and 16 bytes of storage, respectively. For the IEEE “Double Extended” formats, MPI specifies a Format Width of 16 bytes, with 15 exponent bits, bias = +10383, 112 fraction bits, and an encoding analogous to the “Double” format.

The external32 representations of the datatypes returned by `MPI_TYPE_CREATE_F90_REAL/COMPLEX/INTEGER` are given by the following rules.

For `MPI_TYPE_CREATE_F90_REAL`:

```

if      (p > 33) or (r > 4931) then  external32 representation
                                     is undefined
else if (p > 15) or (r > 307) then  external32_size = 16
else if (p > 6) or (r > 37) then   external32_size = 8
else                                     external32_size = 4

```

For `MPI_TYPE_CREATE_F90_COMPLEX`: twice the size as for `MPI_TYPE_CREATE_F90_REAL`.

For `MPI_TYPE_CREATE_F90_INTEGER`:

```

if      (r > 38) then  external32 representation is undefined
else if (r > 18) then  external32_size = 16
else if (r > 9) then   external32_size = 8
else if (r > 4) then   external32_size = 4
else if (r > 2) then   external32_size = 2
else                                     external32_size = 1

```

If the external32 representation of a datatype is undefined, the result of using the datatype directly or indirectly (i.e., as part of another datatype or through a duplicated datatype) in operations that require the external32 representation is undefined. These operations include `MPI_PACK_EXTERNAL`, `MPI_UNPACK_EXTERNAL` and many `MPI_FILE` functions, when the “external32” data representation is used. The ranges for which the external32 representation is undefined are reserved for future standardization.

Support for Size-specific MPI Datatypes

MPI-1 provides named datatypes corresponding to optional Fortran 77 numeric types that contain explicit byte lengths — `MPI_REAL4`, `MPI_INTEGER8`, etc. This section describes a mechanism that generalizes this model to support all Fortran numeric intrinsic types.

We assume that for each **typeclass** (integer, real, complex) and each word size there is a unique machine representation. For every pair (**typeclass**, **n**) supported by a compiler, MPI must provide a named size-specific datatype. The name of this datatype is of the form

MPI_<TYPE>n in C and Fortran and of the form MPI::<TYPE>n in C++ where <TYPE> is one of REAL, INTEGER and COMPLEX, and **n** is the length in bytes of the machine representation. This datatype locally matches all variables of type (**typeclass, n**). The list of names for such types includes:

```
MPI_REAL4
MPI_REAL8
MPI_REAL16
MPI_COMPLEX8
MPI_COMPLEX16
MPI_COMPLEX32
MPI_INTEGER1
MPI_INTEGER2
MPI_INTEGER4
MPI_INTEGER8
MPI_INTEGER16
```

In MPI-1 these datatypes are all optional and correspond to the optional, nonstandard declarations supported by many Fortran compilers. In MPI-2, one datatype is required for each representation supported by the compiler. To be backward compatible with the interpretation of these types in MPI-1, we assume that the nonstandard declarations **REAL*n**, **INTEGER*n**, always create a variable whose representation is of size **n**. All these datatypes are predefined.

The following functions allow a user to obtain a size-specific MPI datatype for any intrinsic Fortran type.

```
MPI_SIZEOF(x, size)
```

| | | |
|-----|------|---|
| IN | x | a Fortran variable of numeric intrinsic type (choice) |
| OUT | size | size of machine representation of that type (integer) |

```
MPI_SIZEOF(X, SIZE, IERROR)
```

```
<type> X
INTEGER SIZE, IERROR
```

This function returns the size in bytes of the machine representation of the given variable. It is a generic Fortran routine and has a Fortran binding only.

Advice to users. This function is similar to the C and C++ *sizeof* operator but behaves slightly differently. If given an array argument, it returns the size of the base element, not the size of the whole array. (*End of advice to users.*)

Rationale. This function is not available in other languages because it would not be useful. (*End of rationale.*)

```

1 MPI_TYPE_MATCH_SIZE(typeclass, size, type)
2     IN     typeclass           generic type specifier (integer)
3
4     IN     size                size, in bytes, of representation (integer)
5
6     OUT    type                datatype with correct type, size (handle)

```

```

7 int MPI_Type_match_size(int typeclass, int size, MPI_Datatype *type)

```

```

8 MPI_TYPE_MATCH_SIZE(TYPECLASS, SIZE, TYPE, IERROR)
9     INTEGER TYPECLASS, SIZE, TYPE, IERROR

```

```

11 static MPI::Datatype MPI::Datatype::Match_size(int typeclass, int size)

```

12 typeclass is one of MPI_TYPECLASS_REAL, MPI_TYPECLASS_INTEGER and
13 MPI_TYPECLASS_COMPLEX, corresponding to the desired **typeclass**. The function returns
14 an MPI datatype matching a local variable of type (**typeclass**, **size**).

15 This function returns a reference (handle) to one of the predefined named datatypes, not
16 a duplicate. This type cannot be freed. MPI_TYPE_MATCH_SIZE can be used to obtain a
17 size-specific type that matches a Fortran numeric intrinsic type by first calling MPI_SIZEOF
18 in order to compute the variable size, and then calling MPI_TYPE_MATCH_SIZE to find a
19 suitable datatype. In C and C++, one can use the C function sizeof(), instead of
20 MPI_SIZEOF. In addition, for variables of default kind the variable's size can be computed
21 by a call to MPI_TYPE_GET_EXTENT, if the typeclass is known. It is erroneous to specify
22 a size not supported by the compiler.

23 *Rationale.* This is a convenience function. Without it, it can be tedious to find the
24 correct named type. See note to implementors below. (*End of rationale.*)

25 *Advice to implementors.* This function could be implemented as a series of tests.

```

27
28
29 int MPI_Type_match_size(int typeclass, int size, MPI_Datatype *rtype)
30 {
31     switch(typeclass) {
32         case MPI_TYPECLASS_REAL: switch(size) {
33             case 4: *rtype = MPI_REAL4; return MPI_SUCCESS;
34             case 8: *rtype = MPI_REAL8; return MPI_SUCCESS;
35             default: error(...);
36         }
37         case MPI_TYPECLASS_INTEGER: switch(size) {
38             case 4: *rtype = MPI_INTEGER4; return MPI_SUCCESS;
39             case 8: *rtype = MPI_INTEGER8; return MPI_SUCCESS;
40             default: error(...); }
41         ... etc ...
42     }
43 }

```

44 (*End of advice to implementors.*)

Communication With Size-specific Types

The usual type matching rules apply to size-specific datatypes: a value sent with datatype `MPI_<TYPE>n` can be received with this same datatype on another process. Most modern computers use 2's complement for integers and IEEE format for floating point. Thus, communication using these size-specific datatypes will not entail loss of precision or truncation errors.

Advice to users. Care is required when communicating in a heterogeneous environment. Consider the following code:

```

real(selected_real_kind(5)) x(100)
call MPI_SIZEOF(x, size, ierror)
call MPI_TYPE_MATCH_SIZE(MPI_TYPECLASS_REAL, size, xtype, ierror)
if (myrank .eq. 0) then
    ... initialize x ...
    call MPI_SEND(x, xtype, 100, 1, ...)
else if (myrank .eq. 1) then
    call MPI_RECV(x, xtype, 100, 0, ...)
endif

```

This may not work in a heterogeneous environment if the value of `size` is not the same on process 1 and process 0. There should be no problem in a homogeneous environment. To communicate in a heterogeneous environment, there are at least four options. The first is to declare variables of default type and use the MPI datatypes for these types, e.g., declare a variable of type `REAL` and use `MPI_REAL`. The second is to use `selected_real_kind` or `selected_int_kind` and with the functions of the previous section. The third is to declare a variable that is known to be the same size on all architectures (e.g., `selected_real_kind(12)` on almost all compilers will result in an 8-byte representation). The fourth is to carefully check representation size before communication. This may require explicit conversion to a variable of size that can be communicated and handshaking between sender and receiver to agree on a size.

Note finally that using the “external32” representation for I/O requires explicit attention to the representation sizes. Consider the following code:

```

real(selected_real_kind(5)) x(100)
call MPI_SIZEOF(x, size, ierror)
call MPI_TYPE_MATCH_SIZE(MPI_TYPECLASS_REAL, size, xtype, ierror)

if (myrank .eq. 0) then
    call MPI_FILE_OPEN(MPI_COMM_SELF, 'foo',
                      MPI_MODE_CREATE+MPI_MODE_WRONLY,
                      MPI_INFO_NULL, fh, ierror)
    call MPI_FILE_SET_VIEW(fh, 0, xtype, xtype, 'external32',
                          MPI_INFO_NULL, ierror)
    call MPI_FILE_WRITE(fh, x, 100, xtype, status, ierror)
    call MPI_FILE_CLOSE(fh, ierror)

```

```

1      endif
2
3      call MPI_BARRIER(MPI_COMM_WORLD, ierror)
4
5      if (myrank .eq. 1) then
6          call MPI_FILE_OPEN(MPI_COMM_SELF, 'foo', MPI_MODE_RDONLY, &
7                          MPI_INFO_NULL, fh, ierror)
8          call MPI_FILE_SET_VIEW(fh, 0, xtype, xtype, 'external32', &
9                          MPI_INFO_NULL, ierror)
10         call MPI_FILE_WRITE(fh, x, 100, xtype, status, ierror)
11         call MPI_FILE_CLOSE(fh, ierror)
12     endif

```

If processes 0 and 1 are on different machines, this code may not work as expected if the size is different on the two machines. (*End of advice to users.*)

13.3 Language Interoperability

13.3.1 Introduction

It is not uncommon for library developers to use one language to develop an applications library that may be called by an application program written in a different language. MPI currently supports ISO (previously ANSI) C, C++, and Fortran bindings. It should be possible for applications in any of the supported languages to call MPI-related functions in another language.

Moreover, MPI allows the development of client-server code, with MPI communication used between a parallel client and a parallel server. It should be possible to code the server in one language and the clients in another language. To do so, communications should be possible between applications written in different languages.

There are several issues that need to be addressed in order to achieve interoperability.

Initialization We need to specify how the MPI environment is initialized for all languages.

Interlanguage passing of MPI opaque objects We need to specify how MPI object handles are passed between languages. We also need to specify what happens when an MPI object is accessed in one language, to retrieve information (e.g., attributes) set in another language.

Interlanguage communication We need to specify how messages sent in one language can be received in another language.

It is highly desirable that the solution for interlanguage interoperability be extendable to new languages, should MPI bindings be defined for such languages.

13.3.2 Assumptions

We assume that conventions exist for programs written in one language to call **functions written in another language**. These conventions specify how to link routines in different

languages into one program, how to call functions in a different language, how to pass arguments between languages, and the correspondence between basic data types in different languages. In general, these conventions will be implementation dependent. Furthermore, not every basic datatype may have a matching type in other languages. For example, C/C++ character strings may not be compatible with Fortran CHARACTER variables. However, we assume that a Fortran INTEGER, as well as a (sequence associated) Fortran array of INTEGERS, can be passed to a C or C++ program. We also assume that Fortran, C, and C++ have address-sized integers. This does not mean that the default-size integers are the same size as default-sized pointers, but only that there is some way to hold (and pass) a C address in a Fortran integer. It is also assumed that INTEGER(KIND=MPI_OFFSET_KIND) can be passed from Fortran to C as MPI_Offset.

13.3.3 Initialization

A call to MPI_INIT or MPI_THREAD_INIT, from any language, initializes MPI for execution in all languages.

Advice to users. Certain implementations use the (inout) argc, argv arguments of the C/C++ version of MPI_INIT in order to propagate values for argc and argv to all executing processes. Use of the Fortran version of MPI_INIT to initialize MPI may result in a loss of this ability. (*End of advice to users.*)

The function MPI_INITIALIZED returns the same answer in all languages.

The function MPI_FINALIZE finalizes the MPI environments for all languages.

The function MPI_FINALIZED returns the same answer in all languages.

The function MPI_ABORT kills processes, irrespective of the language used by the caller or by the processes killed.

The MPI environment is initialized in the same manner for all languages by MPI_INIT. E.g., MPI_COMM_WORLD carries the same information regardless of language: same processes, same environmental attributes, same error handlers.

Information can be added to info objects in one language and retrieved in another.

Advice to users. The use of several languages in one MPI program may require the use of special options at compile and/or link time. (*End of advice to users.*)

Advice to implementors. Implementations may selectively link language specific MPI libraries only to codes that need them, so as not to increase the size of binaries for codes that use only one language. The MPI initialization code need perform initialization for a language only if that language library is loaded. (*End of advice to implementors.*)

13.3.4 Transfer of Handles

Handles are passed between Fortran and C or C++ by using an explicit C wrapper to convert Fortran handles to C handles. There is no direct access to C or C++ handles in Fortran. Handles are passed between C and C++ using overloaded C++ operators called from C++ code. There is no direct access to C++ objects from C.

The type definition MPI_Fint is provided in C/C++ for an integer of the size that matches a Fortran INTEGER; usually, MPI_Fint will be equivalent to int.

The following functions are provided in C to convert from a Fortran communicator handle (which is an integer) to a C communicator handle, and vice versa. See also Section 2.6.5 on page 21.

```
MPI_Comm MPI_Comm_f2c(MPI_Fint comm)
```

If `comm` is a valid Fortran handle to a communicator, then `MPI_Comm_f2c` returns a valid C handle to that same communicator; if `comm = MPI_COMM_NULL` (Fortran value), then `MPI_Comm_f2c` returns a null C handle; if `comm` is an invalid Fortran handle, then `MPI_Comm_f2c` returns an invalid C handle.

```
MPI_Fint MPI_Comm_c2f(MPI_Comm comm)
```

The function `MPI_Comm_c2f` translates a C communicator handle into a Fortran handle to the same communicator; it maps a null handle into a null handle and an invalid handle into an invalid handle.

Similar functions are provided for the other types of opaque objects.

```
MPI_Datatype MPI_Type_f2c(MPI_Fint datatype)
```

```
MPI_Fint MPI_Type_c2f(MPI_Datatype datatype)
```

```
MPI_Group MPI_Group_f2c(MPI_Fint group)
```

```
MPI_Fint MPI_Group_c2f(MPI_Group group)
```

```
MPI_Request MPI_Request_f2c(MPI_Fint request)
```

```
MPI_Fint MPI_Request_c2f(MPI_Request request)
```

```
MPI_File MPI_File_f2c(MPI_Fint file)
```

```
MPI_Fint MPI_File_c2f(MPI_File file)
```

```
MPI_Win MPI_Win_f2c(MPI_Fint win)
```

```
MPI_Fint MPI_Win_c2f(MPI_Win win)
```

```
MPI_Op MPI_Op_f2c(MPI_Fint op)
```

```
MPI_Fint MPI_Op_c2f(MPI_Op op)
```

```
MPI_Info MPI_Info_f2c(MPI_Fint info)
```

```
MPI_Fint MPI_Info_c2f(MPI_Info info)
```

```
MPI_Errhandler MPI_Errhandler_f2c(MPI_Fint errhandler)
```

```
MPI_Fint MPI_Errhandler_c2f(MPI_Errhandler errhandler)
```

Example 13.11 The example below illustrates how the Fortran MPI function `MPI_TYPE_COMMIT` can be implemented by wrapping the C MPI function `MPI_Type_commit` with a C wrapper to do handle conversions. In this example a Fortran-C interface is assumed where a Fortran function is all upper case when referred to from C and arguments are passed by addresses.

```
! FORTRAN PROCEDURE
```

```

SUBROUTINE MPI_TYPE_COMMIT( DATATYPE, IERR)
INTEGER DATATYPE, IERR
CALL MPI_X_TYPE_COMMIT(DATATYPE, IERR)
RETURN
END

/* C wrapper */

void MPI_X_TYPE_COMMIT( MPI_Fint *f_handle, MPI_Fint *ierr)
{
MPI_Datatype datatype;

datatype = MPI_Type_f2c( *f_handle);
*ierr = (MPI_Fint)MPI_Type_commit( &datatype);
*f_handle = MPI_Type_c2f(datatype);
return;
}

```

The same approach can be used for all other MPI functions. The call to `MPI_xxx_f2c` (resp. `MPI_xxx_c2f`) can be omitted when the handle is an OUT (resp. IN) argument, rather than INOUT.

Rationale. The design here provides a convenient solution for the prevalent case, where a C wrapper is used to allow Fortran code to call a C library, or C code to call a Fortran library. The use of C wrappers is much more likely than the use of Fortran wrappers, because it is much more likely that a variable of type INTEGER can be passed to C, than a C handle can be passed to Fortran.

Returning the converted value as a function value rather than through the argument list allows the generation of efficient inlined code when these functions are simple (e.g., the identity). The conversion function in the wrapper does not catch an invalid handle argument. Instead, an invalid handle is passed below to the library function, which, presumably, checks its input arguments. (*End of rationale.*)

C and C++ The C++ language interface provides the functions listed below for mixed-language interoperability. The token `<CLASS>` is used below to indicate any valid MPI opaque handle name (e.g., `Group`), except where noted. For the case where the C++ class corresponding to `<CLASS>` has derived classes, functions are also provided for converting between the derived classes and the C MPI `<CLASS>`.

The following function allows assignment from a C MPI handle to a C++ MPI handle.

```
MPI::<CLASS>& MPI::<CLASS>::operator=(const MPI.<CLASS>& data)
```

The constructor below creates a C++ MPI object from a C MPI handle. This allows the automatic promotion of a C MPI handle to a C++ MPI handle.

```
MPI::<CLASS>::<CLASS>(const MPI.<CLASS>& data)
```

Example 13.12 In order for a C program to use a C++ library, the C++ library must export a C interface that provides appropriate conversions before invoking the underlying

1 C++ library call. This example shows a C interface function that invokes a C++ library
 2 call with a C communicator; the communicator is automatically promoted to a C++ handle
 3 when the underlying C++ function is invoked.

```

4 // C++ library function prototype
5 void cpp_lib_call(MPI::Comm cpp_comm);
6
7 // Exported C function prototype
8 extern "C" {
9 void c_interface(MPI_Comm c_comm);
10 }
11
12 void c_interface(MPI_Comm c_comm)
13 {
14 // the MPI_Comm (c_comm) is automatically promoted to MPI::Comm
15 cpp_lib_call(c_comm);
16 }
17

```

18 The following function allows conversion from C++ objects to C MPI handles. In this
 19 case, the casting operator is overloaded to provide the functionality.

```

20 MPI::<CLASS>::operator MPI_<CLASS>() const
21

```

22 **Example 13.13** A C library routine is called from a C++ program. The C library routine
 23 is prototyped to take an MPI_Comm as an argument.

```

24 // C function prototype
25 extern "C" {
26 void c_lib_call(MPI_Comm c_comm);
27 }
28
29 void cpp_function()
30 {
31 // Create a C++ communicator, and initialize it with a dup of
32 // MPI::COMM_WORLD
33 MPI::Intracomm cpp_comm(MPI::COMM_WORLD.Dup());
34 c_lib_call(cpp_comm);
35 }
36

```

37 *Rationale.* Providing conversion from C to C++ via constructors and from C++
 38 to C via casting allows the compiler to make automatic conversions. Calling C from
 39 C++ becomes trivial, as does the provision of a C or Fortran interface to a C++
 40 library. (*End of rationale.*)

41 *Advice to users.* Note that the casting and promotion operators return new handles
 42 by value. Using these new handles as INOUT parameters will affect the internal MPI
 43 object, but will *not* affect the original handle from which it was cast. (*End of advice*
 44 *to users.*)

45
 46 It is important to note that all C++ objects and their corresponding C handles can be
 47 used interchangeably by an application. For example, an application can cache an attribute
 48 on MPI_COMM_WORLD and later retrieve it from MPI::COMM_WORLD.

13.3.5 Status

The following two procedures are provided in C to convert from a Fortran status (which is an array of integers) to a C status (which is a structure), and vice versa. The conversion occurs on all the information in status, including that which is hidden. That is, no status information is lost in the conversion.

```
int MPI_Status_f2c(MPI_Fint *f_status, MPI_Status *c_status)
```

If `f_status` is a valid Fortran status, but not the Fortran value of `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE`, then `MPI_Status_f2c` returns in `c_status` a valid C status with the same content. If `f_status` is the Fortran value of `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE`, or if `f_status` is not a valid Fortran status, then the call is erroneous.

The C status has the same source, tag and error code values as the Fortran status, and returns the same answers when queried for count, elements, and cancellation. The conversion function may be called with a Fortran status argument that has an undefined error field, in which case the value of the error field in the C status argument is undefined.

Two global variables of type `MPI_Fint*`, `MPI_F_STATUS_IGNORE` and `MPI_F_STATUSES_IGNORE` are declared in `mpi.h`. They can be used to test, in C, whether `f_status` is the Fortran value of `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE`, respectively. These are global variables, not C constant expressions and cannot be used in places where C requires constant expressions. Their value is defined only between the calls to `MPI_INIT` and `MPI_FINALIZE` and should not be changed by user code.

To do the conversion in the other direction, we have the following:

```
int MPI_Status_c2f(MPI_Status *c_status, MPI_Fint *f_status)
```

This call converts a C status into a Fortran status, and has a behavior similar to `MPI_Status_f2c`. That is, the value of `c_status` must not be either `MPI_STATUS_IGNORE` or `MPI_STATUSES_IGNORE`.

Advice to users. There is not a separate conversion function for arrays of statuses, since one can simply loop through the array, converting each status. (*End of advice to users.*)

Rationale. The handling of `MPI_STATUS_IGNORE` is required in order to layer libraries with only a C wrapper: if the Fortran call has passed `MPI_STATUS_IGNORE`, then the C wrapper must handle this correctly. Note that this constant need not have the same value in Fortran and C. If `MPI_Status_f2c` were to handle `MPI_STATUS_IGNORE`, then the type of its result would have to be `MPI_Status**`, which was considered an inferior solution. (*End of rationale.*)

13.3.6 MPI Opaque Objects

Unless said otherwise, opaque objects are “the same” in all languages: they carry the same information, and have the same meaning in both languages. The mechanism described in the previous section can be used to pass references to MPI objects from language to language. An object created in one language can be accessed, modified or freed in another language.

We examine below in more detail, issues that arise for each type of MPI object.

Datatypes

Datatypes encode the same information in all languages. E.g., a datatype accessor like `MPI_TYPE_GET_EXTENT` will return the same information in all languages. If a datatype defined in one language is used for a communication call in another language, then the message sent will be identical to the message that would be sent from the first language: the same communication buffer is accessed, and the same representation conversion is performed, if needed. All predefined datatypes can be used in datatype constructors in any language. If a datatype is committed, it can be used for communication in any language.

The function `MPI_GET_ADDRESS` returns the same value in all languages. Note that we do not require that the constant `MPI_BOTTOM` have the same value in all languages (see 13.3.9, page 476).

Example 13.14

```

15 ! FORTRAN CODE
16 REAL R(5)
17 INTEGER TYPE, IERR, AOBLEN(1), AOTYPE(1)
18 INTEGER (KIND=MPI_ADDRESS_KIND) AODISP(1)
19
20 ! create an absolute datatype for array R
21 AOBLEN(1) = 5
22 CALL MPI_GET_ADDRESS( R, AODISP(1), IERR)
23 AOTYPE(1) = MPI_REAL
24 CALL MPI_TYPE_CREATE_STRUCT(1, AOBLEN,AODISP,AOTYPE, TYPE, IERR)
25 CALL C_ROUTINE(TYPE)
26
27
28 /* C code */
29
30 void C_ROUTINE(MPI_Fint *ftype)
31 {
32     int count = 5;
33     int lens[2] = {1,1};
34     MPI_Aint displs[2];
35     MPI_Datatype types[2], newtype;
36
37     /* create an absolute datatype for buffer that consists */
38     /* of count, followed by R(5) */
39
40     MPI_Get_address(&count, &displs[0]);
41     displs[1] = 0;
42     types[0] = MPI_INT;
43     types[1] = MPI_Type_f2c(*ftype);
44     MPI_Type_create_struct(2, lens, displs, types, &newtype);
45     MPI_Type_commit(&newtype);
46
47     MPI_Send(MPI_BOTTOM, 1, newtype, 1, 0, MPI_COMM_WORLD);
48     /* the message sent contains an int count of 5, followed */

```

```
/* by the 5 REAL entries of the Fortran array R.          */
}
```

Advice to implementors. The following implementation can be used: MPI addresses, as returned by `MPI_GET_ADDRESS`, will have the same value in all languages. One obvious choice is that MPI addresses be identical to regular addresses. The address is stored in the datatype, when datatypes with absolute addresses are constructed. When a send or receive operation is performed, then addresses stored in a datatype are interpreted as displacements that are all augmented by a base address. This base address is (the address of) `buf`, or zero, if `buf = MPI_BOTTOM`. Thus, if `MPI_BOTTOM` is zero then a send or receive call with `buf = MPI_BOTTOM` is implemented exactly as a call with a regular buffer argument: in both cases the base address is `buf`. On the other hand, if `MPI_BOTTOM` is not zero, then the implementation has to be slightly different. A test is performed to check whether `buf = MPI_BOTTOM`. If true, then the base address is zero, otherwise it is `buf`. In particular, if `MPI_BOTTOM` does not have the same value in Fortran and C/C++, then an additional test for `buf = MPI_BOTTOM` is needed in at least one of the languages.

It may be desirable to use a value other than zero for `MPI_BOTTOM` even in C/C++, so as to distinguish it from a `NULL` pointer. If `MPI_BOTTOM = c` then one can still avoid the test `buf = MPI_BOTTOM`, by using the displacement from `MPI_BOTTOM`, i.e., the regular address - `c`, as the MPI address returned by `MPI_GET_ADDRESS` and stored in absolute datatypes. (*End of advice to implementors.*)

Callback Functions

MPI calls may associate callback functions with MPI objects: error handlers are associated with communicators and files, attribute copy and delete functions are associated with attribute keys, reduce operations are associated with operation objects, etc. In a multilanguage environment, a function passed in an MPI call in one language may be invoked by an MPI call in another language. MPI implementations must make sure that such invocation will use the calling convention of the language the function is bound to.

Advice to implementors. Callback functions need to have a language tag. This tag is set when the callback function is passed in by the library function (which is presumably different for each language), and is used to generate the right calling sequence when the callback function is invoked. (*End of advice to implementors.*)

Error Handlers

Advice to implementors. Error handlers, have, in C and C++, a “`stdargs`” argument list. It might be useful to provide to the handler information on the language environment where the error occurred. (*End of advice to implementors.*)

Reduce Operations

Advice to users. Reduce operations receive as one of their arguments the datatype of the operands. Thus, one can define “polymorphic” reduce operations that work for C, C++, and Fortran datatypes. (*End of advice to users.*)

Addresses

Some of the datatype accessors and constructors have arguments of type `MPI_Aint` (in C) or `MPI::Aint` in C++, to hold addresses. The corresponding arguments, in Fortran, have type `INTEGER`. This causes Fortran and C/C++ to be incompatible, in an environment where addresses have 64 bits, but Fortran `INTEGER`s have 32 bits.

This is a problem, irrespective of interlanguage issues. Suppose that a Fortran process has an address space of ≥ 4 GB. What should be the value returned in Fortran by `MPI_ADDRESS`, for a variable with an address above 2^{32} ? The design described here addresses this issue, while maintaining compatibility with current Fortran codes.

The constant `MPI_ADDRESS_KIND` is defined so that, in Fortran 90, `INTEGER(KIND=MPI_ADDRESS_KIND)` is an address sized integer type (typically, but not necessarily, the size of an `INTEGER(KIND=MPI_ADDRESS_KIND)` is 4 on 32 bit address machines and 8 on 64 bit address machines). Similarly, the constant `MPI_INTEGER_KIND` is defined so that `INTEGER(KIND=MPI_INTEGER_KIND)` is a default size `INTEGER`.

There are seven functions that have address arguments: `MPI_TYPE_HVECTOR`, `MPI_TYPE_HINDEXED`, `MPI_TYPE_STRUCT`, `MPI_ADDRESS`, `MPI_TYPE_EXTENT`, `MPI_TYPE_LB` and `MPI_TYPE_UB`.

Four new functions are provided to supplement the first four functions in this list. These functions are described in Section 3.12.1, page 76. The remaining three functions are supplemented by the new function `MPI_TYPE_GET_EXTENT`, described in that same section. The new functions have the same functionality as the old functions in C/C++, or on Fortran systems where default `INTEGER`s are address sized. In Fortran, they accept arguments of type `INTEGER(KIND=MPI_ADDRESS_KIND)`, wherever arguments of type `MPI_Aint` are used in C. On Fortran 77 systems that do not support the Fortran 90 `KIND` notation, and where addresses are 64 bits whereas default `INTEGER`s are 32 bits, these arguments will be of an appropriate integer type. The old functions will continue to be provided, for backward compatibility. However, users are encouraged to switch to the new functions, in Fortran, so as to avoid problems on systems with an address range $> 2^{32}$, and to provide compatibility across languages.

13.3.7 Attributes

Attribute keys can be allocated in one language and freed in another. Similarly, attribute values can be set in one language and accessed in another. To achieve this, attribute keys will be allocated in an integer range that is valid all languages. The same holds true for system-defined attribute values (such as `MPI_TAG_UB`, `MPI_WTIME_IS_GLOBAL`, etc.)

Attribute keys declared in one language are associated with copy and delete functions in that language (the functions provided by the `MPI_{TYPE,COMM,WIN}_CREATE_KEYVAL` call). When a communicator is duplicated, for each attribute, the corresponding copy function is called, using the right calling convention for the language of that function; and similarly, for the delete callback function.

Advice to implementors. This requires that attributes be tagged either as “C,” “C++” or “Fortran,” and that the language tag be checked in order to use the right calling convention for the callback function. (*End of advice to implementors.*)

The attribute manipulation functions described in Section 5.7 of the MPI-1 standard define attributes arguments to be of type `void*` in C, and of type `INTEGER`, in Fortran. On

some systems, INTEGERS will have 32 bits, while C/C++ pointers will have 64 bits. This is a problem if communicator attributes are used to move information from a Fortran caller to a C/C++ callee, or vice-versa.

MPI will store, internally, address sized attributes. If Fortran INTEGERS are smaller, then the Fortran function `MPI_ATTR_GET` will return the least significant part of the attribute word; the Fortran function `MPI_ATTR_PUT` will set the least significant part of the attribute word, which will be sign extended to the entire word. (These two functions may be invoked explicitly by user code, or implicitly, by attribute copying callback functions.)

As for addresses, new functions are provided that manipulate Fortran address sized attributes, and have the same functionality as the old functions in C/C++. These functions are described in Section 5.7.1, page 212. Users are encouraged to use these new functions.

MPI supports two types of attributes: address-valued (pointer) attributes, and integer valued attributes. C and C++ attribute functions put and get address valued attributes. Fortran attribute functions put and get integer valued attributes. When an integer valued attribute is accessed from C or C++, then `MPI_xxx_get_attr` will return the address of (a pointer to) the integer valued attribute. When an address valued attribute is accessed from Fortran, then `MPI_xxx_GET_ATTR` will convert the address into an integer and return the result of this conversion. This conversion is lossless if new style (MPI-2) attribute functions are used, and an integer of kind `MPI_ADDRESS_KIND` is returned. The conversion may cause truncation if old style (MPI-1) attribute functions are used.

Example 13.15 A. C to Fortran

C code

```
static int i = 5;
void *p;
p = &i;
MPI_Comm_put_attr(..., p);
....
```

Fortran code

```
INTEGER(kind = MPI_ADDRESS_KIND) val
CALL MPI_COMM_GET_ATTR(...,val,...)
IF(val.NE.address_of_i) THEN CALL ERROR
```

B. Fortran to C

Fortran code

```
INTEGER(kind=MPI_ADDRESS_KIND) val
val = 55555
CALL MPI_COMM_PUT_ATTR(...,val,ierr)
```

C code

```

1  int *p;
2  MPI_Comm_get_attr(...,&p, ...);
3  if (*p != 55555) error();
4
5

```

The predefined MPI attributes can be integer valued or address valued. Predefined integer valued attributes, such as MPI_TAG_UB, behave as if they were put by a Fortran call. I.e., in Fortran, MPI_COMM_GET_ATTR(MPI_COMM_WORLD, MPI_TAG_UB, val, flag, ierr) will return in val the upper bound for tag value; in C, MPI_Comm_get_attr(MPI_COMM_WORLD, MPI_TAG_UB, &p, &flag) will return in p a pointer to an int containing the upper bound for tag value.

Address valued predefined attributes, such as MPI_WIN_BASE behave as if they were put by a C call. I.e., in Fortran, MPI_WIN_GET_ATTR(win, MPI_WIN_BASE, val, flag, ierror) will return in val the base address of the window, converted to an integer. In C, MPI_Win_get_attr(win, MPI_WIN_BASE, &p, &flag) will return in p a pointer to the window base, cast to (void *).

Rationale. The design is consistent with the behavior specified in MPI-1 for predefined attributes, and ensures that no information is lost when attributes are passed from language to language. (*End of rationale.*)

Advice to implementors. Implementations should tag attributes either as address attributes or as integer attributes, according to whether they were set in C or in Fortran. Thus, the right choice can be made when the attribute is retrieved. (*End of advice to implementors.*)

13.3.8 Extra State

Extra-state should not be modified by the copy or delete callback functions. (This is obvious from the C binding, but not obvious from the Fortran binding). However, these functions may update state that is indirectly accessed via extra-state. E.g., in C, extra-state can be a pointer to a data structure that is modified by the copy or callback functions; in Fortran, extra-state can be an index into an entry in a COMMON array that is modified by the copy or callback functions. In a multithreaded environment, users should be aware that distinct threads may invoke the same callback function concurrently: if this function modifies state associated with extra-state, then mutual exclusion code must be used to protect updates and accesses to the shared state.

13.3.9 Constants

MPI constants have the same value in all languages, unless specified otherwise. This does not apply to constant handles (MPI_INT, MPI_COMM_WORLD, MPI_ERRORS_RETURN, MPI_SUM, etc.) These handles need to be converted, as explained in Section 13.3.4. Constants that specify maximum lengths of strings (see Section A.1.1 for a listing) have a value one less in Fortran than C/C++ since in C/C++ the length includes the null terminating character. Thus, these constants represent the amount of space which must be allocated to hold the largest possible such string, rather than the maximum number of printable characters the string could contain.

Advice to users. This definition means that it is safe in C/C++ to allocate a buffer to receive a string using a declaration like

```
char name [MPI_MAX_OBJECT_NAME];
```

(End of advice to users.)

Also constant “addresses,” i.e., special values for reference arguments that are not handles, such as MPI_BOTTOM or MPI_STATUS_IGNORE may have different values in different languages.

Rationale. The current MPI standard specifies that MPI_BOTTOM can be used in initialization expressions in C, but not in Fortran. Since Fortran does not normally support call by value, then MPI_BOTTOM must be in Fortran the name of a predefined static variable, e.g., a variable in an MPI declared COMMON block. On the other hand, in C, it is natural to take MPI_BOTTOM = 0 (Caveat: Defining MPI_BOTTOM = 0 implies that NULL pointer cannot be distinguished from MPI_BOTTOM; it may be that MPI_BOTTOM = 1 is better ...) Requiring that the Fortran and C values be the same will complicate the initialization process. *(End of rationale.)*

13.3.10 Interlanguage Communication

The type matching rules for communications in MPI are not changed: the datatype specification for each item sent should match, in type signature, the datatype specification used to receive this item (unless one of the types is MPI_PACKED). Also, the type of a message item should match the type declaration for the corresponding communication buffer location, unless the type is MPI_BYTE or MPI_PACKED. Interlanguage communication is allowed if it complies with these rules.

Example 13.16 In the example below, a Fortran array is sent from Fortran and received in C.

```
! FORTRAN CODE
REAL R(5)
INTEGER TYPE, IERR, MYRANK, AOBLN(1), AOTYPE(1)
INTEGER (KIND=MPI_ADDRESS_KIND) AODISP(1)

! create an absolute datatype for array R
AOBLN(1) = 5
CALL MPI_GET_ADDRESS( R, AODISP(1), IERR)
AOTYPE(1) = MPI_REAL
CALL MPI_TYPE_CREATE_STRUCT(1, AOBLN,AODISP,AOTYPE, TYPE, IERR)
CALL MPI_TYPE_COMMIT(TYPE, IERR)

CALL MPI_COMM_RANK( MPI_COMM_WORLD, MYRANK, IERR)
IF (MYRANK.EQ.0) THEN
    CALL MPI_SEND( MPI_BOTTOM, 1, TYPE, 1, 0, MPI_COMM_WORLD, IERR)
ELSE
    CALL C_ROUTINE(TYPE)
END IF
```

```
1
2  /* C code */
3
4  void C_ROUTINE(MPI_Fint *fhandle)
5  {
6  MPI_Datatype type;
7  MPI_Status status;
8
9  type = MPI_Type_f2c(*fhandle);
10
11 MPI_Recv( MPI_BOTTOM, 1, type, 0, 0, MPI_COMM_WORLD, &status);
12 }
13
14
```

15 MPI implementors may weaken these type matching rules, and allow messages to be sent
16 with Fortran types and received with C types, and vice versa, when those types match. I.e.,
17 if the Fortran type INTEGER is identical to the C type int, then an MPI implementation may
18 allow data to be sent with datatype MPI_INTEGER and be received with datatype MPLINT.
19 However, such code is not portable.

```
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
```

Chapter 14

Profiling Interface

14.1 Requirements

To meet the MPI profiling interface, an implementation of the MPI functions *must*

1. provide a mechanism through which all of the MPI defined functions **except** those allowed as macros (See Section 2.6.5). This requires, in C and Fortran, an alternate entry point name, with the prefix `PMPI_` for each MPI function. The profiling interface in C++ is described in Section 13.1.10. For routines implemented as macros, it is still required that the `PMPI_` version be supplied and work as expected, but it is not possible to replace at link time the `MPI_` version with a user-defined version.
2. ensure that those MPI functions **that** are not replaced may still be linked into an executable image without causing name clashes.
3. document the implementation of different language bindings of the MPI interface if they are layered on top of each other, so that the profiler developer knows whether she must implement the profile interface for each binding, or can economise by implementing it only for the lowest level routines.
4. where the implementation of different language bindings is done through a layered approach (e.g. the Fortran binding is a set of “wrapper” functions **that** call the C implementation), ensure that these wrapper functions are separable from the rest of the library.

This **separability** is necessary to allow a separate profiling library to be correctly implemented, since (at least with Unix linker semantics) the profiling library must contain these wrapper functions if it is to perform as expected. This requirement allows the person who builds the profiling library to extract these functions from the original MPI library and add them into the profiling library without bringing along any other unnecessary code.

5. provide a no-op routine `MPI_PCONTROL` in the MPI library.

14.2 Discussion

The objective of the MPI profiling interface is to ensure that it is relatively easy for authors of profiling (and other similar) tools to interface their codes to MPI implementations on

1 different machines.

2 Since MPI is a machine independent standard with many different implementations,
3 it is unreasonable to expect that the authors of profiling tools for MPI will have access to
4 the source code **that** implements MPI on any particular machine. It is therefore necessary
5 to provide a mechanism by which the implementors of such tools can collect whatever
6 performance information they wish *without* access to the underlying implementation.

7 We believe that having such an interface is important if MPI is to be attractive to end
8 users, since the availability of many different tools will be a significant factor in attracting
9 users to the MPI standard.

10 The profiling interface is just that, an interface. It says *nothing* about the way in which
11 it is used. There is therefore no attempt to lay down what information is collected through
12 the interface, or how the collected information is saved, filtered, or displayed.

13 While the initial impetus for the development of this interface arose from the desire to
14 permit the implementation of profiling tools, it is clear that an interface like that specified
15 may also prove useful for other purposes, such as “internetworking” multiple MPI imple-
16 mentations. Since all that is defined is an interface, there is no objection to its being used
17 wherever it is useful.

18 As the issues being addressed here are intimately tied up with the way in which ex-
19 ecutable images are built, which may differ greatly on different machines, the examples
20 given below should be treated solely as one way of implementing the objective of the MPI
21 profiling interface. The actual requirements made of an implementation are those detailed
22 in the Requirements section above, the whole of the rest of this chapter is only present as
23 justification and discussion of the logic for those requirements.

24 The examples below show one way in which an implementation could be constructed to
25 meet the requirements on a Unix system (there are doubtless others **that** would be equally
26 valid).

28 14.3 Logic of the design

30 Provided that an MPI implementation meets the requirements above, it is possible for the
31 implementor of the profiling system to intercept all of the MPI calls **that** are made by
32 the user program. She can then collect whatever information she requires before calling
33 the underlying MPI implementation (through its name shifted entry points) to achieve the
34 desired effects.

36 14.3.1 Miscellaneous control of profiling

38 There is a clear requirement for the user code to be able to control the profiler dynamically
39 at run time. This is normally used for (at least) the purposes of

- 40 • Enabling and disabling profiling depending on the state of the calculation.
- 42 • Flushing trace buffers at non-critical points in the calculation
- 44 • Adding user events to a trace file.

45 These requirements are met by use of the MPI_PCONTROL.

```
MPI_PCONTROL(level, ...)
```

```
IN          level                      Profiling level
```

```
int MPI_Pcontrol(const int level, ...)
```

```
MPI_PCONTROL(LEVEL)
    INTEGER LEVEL, ...
```

```
void MPI::Pcontrol(const int level, ...)
```

MPI libraries themselves make no use of this routine, and simply return immediately to the user code. However the presence of calls to this routine allows a profiling package to be explicitly called by the user.

Since MPI has no control of the implementation of the profiling code, we are unable to specify precisely the semantics **that** will be provided by calls to `MPI_PCONTROL`. This vagueness extends to the number of arguments to the function, and their datatypes.

However to provide some level of portability of user codes to different profiling libraries, we request the following meanings for certain values of `level`.

- `level==0` Profiling is disabled.
- `level==1` Profiling is enabled at a normal default level of detail.
- `level==2` Profile buffers are flushed. (This may be a no-op in some profilers).
- All other values of `level` have profile library defined effects and additional arguments.

We also request that the default state after `MPI_INIT` has been called is for profiling to be enabled at the normal default level. (i.e. as if `MPI_PCONTROL` had just been called with the argument 1). This allows users to link with a profiling library and obtain profile output without having to modify their source code at all.

The provision of `MPI_PCONTROL` as a no-op in the standard MPI library allows them to modify their source code to obtain more detailed profiling information, but still be able to link exactly the same code against the standard MPI library.

14.4 Examples

14.4.1 Profiler implementation

Suppose that the profiler wishes to accumulate the total amount of data sent by the `MPI_SEND` function, along with the total elapsed time spent in the function. This could trivially be achieved thus

```
static int totalBytes;
static double totalTime;

int MPI_SEND(void * buffer, const int count, MPI_Datatype datatype,
             int dest, int tag, MPI_comm comm)
{
    double tstart = MPI_Wtime();      /* Pass on all the arguments */
    int extent;
```

```

1   int result    = PMPI_Send(buffer,count,datatype,dest,tag,comm);
2
3   MPI_Type_size(datatype, &extent); /* Compute size */
4   totalBytes += count*extent;
5
6   totalTime += MPI_Wtime() - tstart;      /* and time      */
7
8   return result;
9 }

```

14.4.2 MPI library implementation

On a Unix system, in which the MPI library is implemented in C, then there are various possible options, of which two of the most obvious are presented here. Which is better depends on whether the linker and compiler support weak symbols.

Systems with weak symbols

If the compiler and linker support weak external symbols (e.g. Solaris 2.x, other system V.4 machines), then only a single library is required through the use of `#pragma weak` thus

```

21 #pragma weak MPI_Example = PMPI_Example
22
23 int PMPI_Example(/* appropriate args */)
24 {
25     /* Useful content */
26 }

```

The effect of this `#pragma` is to define the external symbol `MPI_Example` as a weak definition. This means that the linker will not complain if there is another definition of the symbol (for instance in the profiling library), however if no other definition exists, then the linker will use the weak definition.

Systems without weak symbols

In the absence of weak symbols then one possible solution would be to use the C macro pre-processor thus

```

37
38 #ifndef PROFILELIB
39 #   ifdef __STDC__
40 #       define FUNCTION(name) P##name
41 #   else
42 #       define FUNCTION(name) P/**/name
43 #   endif
44 #else
45 #   define FUNCTION(name) name
46 #endif

```

Each of the user visible functions in the library would then be declared thus


```

int FUNCTION(MPI_Example)(/* appropriate args */)
{
    /* Useful content */
}

```

The same source file can then be compiled to produce both versions of the library, depending on the state of the `PROFILELIB` macro symbol.

It is required that the standard MPI library be built in such a way that the inclusion of MPI functions can be achieved one at a time. This is a somewhat unpleasant requirement, since it may mean that each external function has to be compiled from a separate file. However this is necessary so that the author of the profiling library need only define those MPI functions *that* she wishes to intercept, references to any others being fulfilled by the normal MPI library. Therefore the link step can look something like this

```
% cc ... -lmyprof -lmpmi -lmpi
```

Here `libmyprof.a` contains the profiler functions *that* intercept some of the MPI functions. `libmpmi.a` contains the “name shifted” MPI functions, and `libmpi.a` contains the normal definitions of the MPI functions.

14.4.3 Complications

Multiple counting

Since parts of the MPI library may themselves be implemented using more basic MPI functions (e.g. a portable implementation of the collective operations implemented using point to point communications), there is potential for profiling functions to be called from within an MPI function *that* was called from a profiling function. This could lead to “double counting” of the time spent in the inner routine. Since this effect could actually be useful under some circumstances (e.g. it might allow one to answer the question “How much time is spent in the point to point routines when they’re called from collective functions?”), we have decided not to enforce any restrictions on the author of the MPI library *that* would overcome this. Therefore the author of the profiling library should be aware of this problem, and guard against it herself. In a single threaded world this is easily achieved through use of a static variable in the profiling code *that* remembers if you are already inside a profiling routine. It becomes more complex in a multi-threaded environment (as does the meaning of the times recorded !)

Linker oddities

The Unix linker traditionally operates in one pass : the effect of this is that functions from libraries are only included in the image if they are needed at the time the library is scanned. When combined with weak symbols, or multiple definitions of the same function, this can cause odd (and unexpected) effects.

Consider, for instance, an implementation of MPI in which the Fortran binding is achieved by using wrapper functions on top of the C implementation. The author of the profile library then assumes that it is reasonable only to provide profile functions for the C binding, since Fortran will eventually call these, and the cost of the wrappers is assumed to be small. However, if the wrapper functions are not in the profiling library, then none

1 of the profiled entry points will be undefined when the profiling library is called. Therefore
2 none of the profiling code will be included in the image. When the standard MPI library
3 is scanned, the Fortran wrappers will be resolved, and will also pull in the base versions of
4 the MPI functions. The overall effect is that the code will link successfully, but will not be
5 profiled.

6 To overcome this we must ensure that the Fortran wrapper functions are included in
7 the profiling version of the library. We ensure that this is possible by requiring that these
8 be separable from the rest of the base MPI library. This allows them to be aared out of the
9 base library and into the profiling one.

11 14.5 Multiple levels of interception

12
13 The scheme given here does not directly support the nesting of profiling functions, since it
14 provides only a single alternative name for each MPI function. Consideration was given to
15 an implementation **that** would allow multiple levels of call interception, however we were
16 unable to construct an implementation of this **that** did not have the following disadvantages
17

- 18 • assuming a particular implementation language.
- 19
- 20 • imposing a run time cost even when no profiling was taking place.

21 Since one of the objectives of MPI is to permit efficient, low latency implementations, and
22 it is not the business of a standard to require a particular implementation language, we
23 decided to accept the scheme outlined above.

24 Note, however, that it is possible to use the scheme above to implement a multi-level
25 system, since the function called by the user may call many different profiling functions
26 before calling the underlying MPI function.

27 Unfortunately such an implementation may require more cooperation between the dif-
28 ferent profiling libraries than is required for the single level implementation detailed above.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Chapter 15

Deprecated Functions

15.1 Deprecated since MPI-2.0

The following function is deprecated and is superseded by `MPI_TYPE_CREATE_HVECTOR` in MPI-2.0. The language independent definition and the C binding of the deprecated function is the same as of the new function, except of the function name. Only the Fortran language binding is different.

`MPI_TYPE_HVECTOR(count, blocklength, stride, oldtype, newtype)`

| | | |
|-----|-------------|--|
| IN | count | number of blocks (nonnegative integer) |
| IN | blocklength | number of elements in each block (nonnegative integer) |
| IN | stride | number of bytes between start of each block (integer) |
| IN | oldtype | old datatype (handle) |
| OUT | newtype | new datatype (handle) |

```
int MPI_Type_hvector(int count, int blocklength, MPI_Aint stride,  
                    MPI_Datatype oldtype, MPI_Datatype *newtype)
```

```
MPI_TYPE_HVECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR)  
INTEGER COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR
```

The following function is deprecated and is superseded by `MPI_TYPE_CREATE_HINDEXED` in MPI-2.0. The language independent definition and the C binding of the deprecated function is the same as of the new function, except of the function name. Only the Fortran language binding is different.

1 **MPI_TYPE_HINDEXED**(count, array_of_blocklengths, array_of_displacements, oldtype, new-
2 type)

| | | | |
|----|-----|------------------------|---|
| 3 | IN | count | number of blocks – also number of entries in 4 array_of_displacements and array_of_blocklengths (non- 5 negative integer) |
| 6 | | | |
| 7 | IN | array_of_blocklengths | number of elements in each block (array of nonnega- 8 tive integers) |
| 9 | IN | array_of_displacements | byte displacement of each block (array of integer) |
| 10 | IN | oldtype | old datatype (handle) |
| 11 | | | |
| 12 | OUT | newtype | new datatype (handle) |

```
13
14 int MPI_Type_hindexed(int count, int *array_of_blocklengths,
15 MPI_Aint *array_of_displacements, MPI_Datatype oldtype,
16 MPI_Datatype *newtype)
```

```
17 MPI_TYPE_HINDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
18 OLDTYPE, NEWTYPE, IERROR)
19 INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*),
20 OLDTYPE, NEWTYPE, IERROR
21
```

22 The following function is deprecated and is superseded by
23 **MPI_TYPE_CREATE_STRUCT** in MPI-2.0. The language independent definition and the C
24 binding of the deprecated function is the same as of the new function, except of the function
25 name. Only the Fortran language binding is different.

26
27
28 **MPI_TYPE_STRUCT**(count, array_of_blocklengths, array_of_displacements, array_of_types, new-
29 type)

| | | | |
|----|-----|------------------------|---|
| 30 | IN | count | number of blocks (integer) (nonnegative integer) – also 31 number of entries in arrays array_of_types, 32 array_of_displacements and array_of_blocklengths |
| 33 | IN | array_of_blocklength | number of elements in each block (array of nonnega- 34 tive integer) |
| 35 | | | |
| 36 | IN | array_of_displacements | byte displacement of each block (array of integer) |
| 37 | IN | array_of_types | type of elements in each block (array of handles to 38 datatype objects) |
| 39 | OUT | newtype | new datatype (handle) |

```
40
41
42 int MPI_Type_struct(int count, int *array_of_blocklengths,
43 MPI_Aint *array_of_displacements, MPI_Datatype *array_of_types,
44 MPI_Datatype *newtype)
```

```
45 MPI_TYPE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
46 ARRAY_OF_TYPES, NEWTYPE, IERROR)
47 INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*),
48 ARRAY_OF_TYPES(*), NEWTYPE, IERROR
```

The following function is deprecated and is superseded by `MPI_GET_ADDRESS` in MPI-2.0. The language independent definition and the C binding of the deprecated function is the same as of the new function, except of the function name. Only the Fortran language binding is different.

`MPI_ADDRESS(location, address)`

| | | |
|-----|----------|------------------------------------|
| IN | location | location in caller memory (choice) |
| OUT | address | address of location (integer) |

`int MPI_Address(void* location, MPI_Aint *address)`

`MPI_ADDRESS(LOCATION, ADDRESS, IERROR)`

`<type> LOCATION(*)`
`INTEGER ADDRESS, IERROR`

The following functions are deprecated and are superseded by `MPI_TYPE_GET_EXTENT` in MPI-2.0.

`MPI_TYPE_EXTENT(datatype, extent)`

| | | |
|-----|----------|---------------------------|
| IN | datatype | datatype (handle) |
| OUT | extent | datatype extent (integer) |

`int MPI_Type_extent(MPI_Datatype datatype, MPI_Aint *extent)`

`MPI_TYPE_EXTENT(DATATYPE, EXTENT, IERROR)`

`INTEGER DATATYPE, EXTENT, IERROR`

Returns the extent of a datatype, where extent is as defined on page 94.

The two functions below can be used for finding the lower bound and the upper bound of a datatype.

`MPI_TYPE_LB(datatype, displacement)`

| | | |
|-----|--------------|---|
| IN | datatype | datatype (handle) |
| OUT | displacement | displacement of lower bound from origin, in bytes (integer) |

`int MPI_Type_lb(MPI_Datatype datatype, MPI_Aint* displacement)`

`MPI_TYPE_LB(DATATYPE, DISPLACEMENT, IERROR)`

`INTEGER DATATYPE, DISPLACEMENT, IERROR`

```

1 MPI_TYPE_UB( datatype, displacement)
2     IN      datatype          datatype (handle)
3
4     OUT     displacement      displacement of upper bound from origin, in bytes (in-
5                                     teger)
6

```

```

7 int MPI_Type_ub(MPI_Datatype datatype, MPI_Aint* displacement)
8

```

```

9 MPI_TYPE_UB( DATATYPE, DISPLACEMENT, IERROR)
10     INTEGER DATATYPE, DISPLACEMENT, IERROR
11

```

The following function is deprecated and is superseded by `MPI_COMM_CREATE_KEYVAL` in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

```

12 MPI_KEYVAL_CREATE(copy_fn, delete_fn, keyval, extra_state)
13
14     IN      copy_fn           Copy callback function for keyval
15
16     IN      delete_fn        Delete callback function for keyval
17
18     OUT     keyval           key value for future access (integer)
19
20     IN      extra_state      Extra state for callback functions
21
22

```

```

23
24 int MPI_Keyval_create(MPI_Copy_function *copy_fn, MPI_Delete_function
25     *delete_fn, int *keyval, void* extra_state)
26

```

```

27 MPI_KEYVAL_CREATE(COPY_FN, DELETE_FN, KEYVAL, EXTRA_STATE, IERROR)
28     EXTERNAL COPY_FN, DELETE_FN
29     INTEGER KEYVAL, EXTRA_STATE, IERROR
30

```

The `copy_fn` function is invoked when a communicator is duplicated by `MPI_COMM_DUP`. `copy_fn` should be of type `MPI_Copy_function`, which is defined as follows:

```

31
32 typedef int MPI_Copy_function(MPI_Comm oldcomm, int keyval,
33     void *extra_state, void *attribute_val_in,
34     void *attribute_val_out, int *flag)
35

```

A Fortran declaration for such a function is as follows:

```

36
37 SUBROUTINE COPY_FUNCTION(OLDCOMM, KEYVAL, EXTRA_STATE, ATTRIBUTE_VAL_IN,
38     ATTRIBUTE_VAL_OUT, FLAG, IERR)
39     INTEGER OLDCOMM, KEYVAL, EXTRA_STATE, ATTRIBUTE_VAL_IN,
40     ATTRIBUTE_VAL_OUT, IERR
41     LOGICAL FLAG
42

```

`copy_fn` may be specified as `MPI_NULL_COPY_FN` or `MPI_DUP_FN` from either C or FORTRAN; `MPI_NULL_COPY_FN` is a function that does nothing other than returning `flag = 0` and `MPI_SUCCESS`. `MPI_DUP_FN` is a simple-minded copy function that sets `flag = 1`, returns the value of `attribute_val_in` in `attribute_val_out`, and returns `MPI_SUCCESS`.

Analogous to `copy_fn` is a callback deletion function, defined as follows. The `delete_fn` function is invoked when a communicator is deleted by `MPI_COMM_FREE` or when a call

is made explicitly to MPI_ATTR_DELETE. delete_fn should be of type MPI_Delete_function, which is defined as follows:

```
typedef int MPI_Delete_function(MPI_Comm comm, int keyval,
void *attribute_val, void *extra_state);
```

A Fortran declaration for such a function is as follows:

```
SUBROUTINE DELETE_FUNCTION(COMM, KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE, IERR)
  INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE, IERR
```

delete_fn may be specified as MPI_NULL_DELETE_FN from either C or FORTRAN; MPI_NULL_DELETE_FN is a function that does nothing, other than returning MPI_SUCCESS.

The following function is deprecated and is superseded by MPI_COMM_FREE_KEYVAL in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

MPI_KEYVAL_FREE(keyval)

| | | |
|-------|--------|---------------------------------------|
| INOUT | keyval | Frees the integer key value (integer) |
|-------|--------|---------------------------------------|

```
int MPI_Keyval_free(int *keyval)
```

```
MPI_KEYVAL_FREE(KEYVAL, IERROR)
```

```
  INTEGER KEYVAL, IERROR
```

The following function is deprecated and is superseded by MPI_COMM_SET_ATTR in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

MPI_ATTR_PUT(comm, keyval, attribute_val)

| | | |
|-------|------|---|
| INOUT | comm | communicator to which attribute will be attached (handle) |
|-------|------|---|

| | | |
|----|--------|---|
| IN | keyval | key value, as returned by MPI_KEYVAL_CREATE (integer) |
|----|--------|---|

| | | |
|----|---------------|-----------------|
| IN | attribute_val | attribute value |
|----|---------------|-----------------|

```
int MPI_Attr_put(MPI_Comm comm, int keyval, void* attribute_val)
```

```
MPI_ATTR_PUT(COMM, KEYVAL, ATTRIBUTE_VAL, IERROR)
```

```
  INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, IERROR
```

The following function is deprecated and is superseded by MPI_COMM_GET_ATTR in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

```

1 MPI_ATTR_GET(comm, keyval, attribute_val, flag)
2     IN      comm      communicator to which attribute is attached (handle)
3
4     IN      keyval    key value (integer)
5
6     OUT     attribute_val  attribute value, unless flag = false
7
8     OUT     flag      true if an attribute value was extracted; false if no
                       attribute is associated with the key

```

```

9
10 int MPI_Attr_get(MPI_Comm comm, int keyval, void *attribute_val, int *flag)

```

```

11 MPI_ATTR_GET(COMM, KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)
12     INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, IERROR
13     LOGICAL FLAG

```

The following function is deprecated and is superseded by `MPI_COMM_DELETE_ATTR` in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

```

19 MPI_ATTR_DELETE(comm, keyval)
20
21     INOUT   comm      communicator to which attribute is attached (handle)
22
23     IN      keyval    The key value of the deleted attribute (integer)

```

```

25 int MPI_Attr_delete(MPI_Comm comm, int keyval)

```

```

26
27 MPI_ATTR_DELETE(COMM, KEYVAL, IERROR)
28     INTEGER COMM, KEYVAL, IERROR

```

The following function is deprecated and is superseded by `MPI_COMM_CREATE_ERRHANDLER` in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

```

35 MPI_ERRHANDLER_CREATE( function, errhandler )
36     IN      function  user defined error handling procedure
37
38     OUT     errhandler MPI error handler (handle)

```

```

39
40 int MPI_Errhandler_create(MPI_Handler_function *function,
41                          MPI_Errhandler *errhandler)

```

```

42 MPI_ERRHANDLER_CREATE(FUNCTION, ERRHANDLER, IERROR)
43     EXTERNAL FUNCTION
44     INTEGER ERRHANDLER, IERROR

```

Register the user routine `function` for use as an MPI exception handler. Returns in `errhandler` a handle to the registered exception handler.

In the C language, the user routine should be a C function of type `MPI_Handler_function`, which is defined as:

```
typedef void (MPI_Handler_function)(MPI_Comm *, int *, ...);
```

The first argument is the communicator in use, the second is the error code to be returned.

In the Fortran language, the user routine should be of the form:

```
SUBROUTINE HANDLER_FUNCTION(COMM, ERROR_CODE, .....)  
  INTEGER COMM, ERROR_CODE
```

The following function is deprecated and is superseded by `MPI_COMM_SET_ERRHANDLER` in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

`MPI_ERRHANDLER_SET(comm, errhandler)`

| | | |
|--------------|-------------------|--|
| INOUT | comm | communicator to set the error handler for (handle) |
| IN | errhandler | new MPI error handler for communicator (handle) |

```
int MPI_Errhandler_set(MPI_Comm comm, MPI_Errhandler errhandler)
```

```
MPI_ERRHANDLER_SET(COMM, ERRHANDLER, IERROR)  
  INTEGER COMM, ERRHANDLER, IERROR
```

Associates the new error handler `errorhandler` with communicator `comm` at the calling process. Note that an error handler is always associated with the communicator.

The following function is deprecated and is superseded by `MPI_COMM_GET_ERRHANDLER` in MPI-2.0. The language independent definition of the deprecated function is the same as of the new function, except of the function name. The language bindings are modified.

`MPI_ERRHANDLER_GET(comm, errhandler)`

| | | |
|------------|-------------------|---|
| IN | comm | communicator to get the error handler from (handle) |
| OUT | errhandler | MPI error handler currently associated with communicator (handle) |

```
int MPI_Errhandler_get(MPI_Comm comm, MPI_Errhandler *errhandler)
```

```
MPI_ERRHANDLER_GET(COMM, ERRHANDLER, IERROR)  
  INTEGER COMM, ERRHANDLER, IERROR
```

Returns in `errhandler` (a handle to) the error handler that is currently associated with communicator `comm`.

Annex A

Language Binding

In this section we summarize the specific bindings for C, Fortran, and C++. First, all constants are summarized. Then, for each binding, all function bindings are presented. Listings are alphabetical within these chapter.

A.1 Defined Values and Handles

A.1.1 Defined Constants

The C and Fortran name is listed in the left column and the C++ name is listed in the right column.

Return Codes

| | |
|-------------------|--------------------|
| MPI_SUCCESS | MPI::SUCCESS |
| MPI_ERR_BUFFER | MPI::ERR_BUFFER |
| MPI_ERR_COUNT | MPI::ERR_COUNT |
| MPI_ERR_TYPE | MPI::ERR_TYPE |
| MPI_ERR_TAG | MPI::ERR_TAG |
| MPI_ERR_COMM | MPI::ERR_COMM |
| MPI_ERR_RANK | MPI::ERR_RANK |
| MPI_ERR_REQUEST | MPI::ERR_REQUEST |
| MPI_ERR_ROOT | MPI::ERR_ROOT |
| MPI_ERR_GROUP | MPI::ERR_GROUP |
| MPI_ERR_OP | MPI::ERR_OP |
| MPI_ERR_TOPOLOGY | MPI::ERR_TOPOLOGY |
| MPI_ERR_DIMS | MPI::ERR_DIMS |
| MPI_ERR_ARG | MPI::ERR_ARG |
| MPI_ERR_UNKNOWN | MPI::ERR_UNKNOWN |
| MPI_ERR_TRUNCATE | MPI::ERR_TRUNCATE |
| MPI_ERR_OTHER | MPI::ERR_OTHER |
| MPI_ERR_INTERN | MPI::ERR_INTERN |
| MPI_ERR_PENDING | MPI::ERR_PENDING |
| MPI_ERR_IN_STATUS | MPI::ERR_IN_STATUS |

(Continued on next page)

Return Codes (continued)

| | | |
|-------------------------------|--------------------------------|----|
| MPI_ERR_ACCESS | MPI::ERR_ACCESS | 1 |
| MPI_ERR_AMODE | MPI::ERR_AMODE | 2 |
| MPI_ERR_ASSERT | MPI::ERR_ASSERT | 3 |
| MPI_ERR_BAD_FILE | MPI::ERR_BAD_FILE | 4 |
| MPI_ERR_BASE | MPI::ERR_BASE | 5 |
| MPI_ERR_CONVERSION | MPI::ERR_CONVERSION | 6 |
| MPI_ERR_DISP | MPI::ERR_DISP | 7 |
| MPI_ERR_DUP_DATAREP | MPI::ERR_DUP_DATAREP | 8 |
| MPI_ERR_FILE_EXISTS | MPI::ERR_FILE_EXISTS | 9 |
| MPI_ERR_FILE_IN_USE | MPI::ERR_FILE_IN_USE | 10 |
| MPI_ERR_FILE | MPI::ERR_FILE | 11 |
| MPI_ERR_INFO_KEY | MPI::ERR_INFO_VALUE | 12 |
| MPI_ERR_INFO_NOKEY | MPI::ERR_INFO_NOKEY | 13 |
| MPI_ERR_INFO_VALUE | MPI::ERR_INFO_KEY | 14 |
| MPI_ERR_INFO | MPI::ERR_INFO | 15 |
| MPI_ERR_IO | MPI::ERR_IO | 16 |
| MPI_ERR_KEYVAL | MPI::ERR_KEYVAL | 17 |
| MPI_ERR_LOCKTYPE | MPI::ERR_LOCKTYPE | 18 |
| MPI_ERR_NAME | MPI::ERR_NAME | 19 |
| MPI_ERR_NO_MEM | MPI::ERR_NO_MEM | 20 |
| MPI_ERR_NOT_SAME | MPI::ERR_NOT_SAME | 21 |
| MPI_ERR_NO_SPACE | MPI::ERR_NO_SPACE | 22 |
| MPI_ERR_NO_SUCH_FILE | MPI::ERR_NO_SUCH_FILE | 23 |
| MPI_ERR_PORT | MPI::ERR_PORT | 24 |
| MPI_ERR_QUOTA | MPI::ERR_QUOTA | 25 |
| MPI_ERR_READ_ONLY | MPI::ERR_READ_ONLY | 26 |
| MPI_ERR_RMA_CONFLICT | MPI::ERR_RMA_CONFLICT | 27 |
| MPI_ERR_RMA_SYNC | MPI::ERR_RMA_SYNC | 28 |
| MPI_ERR_SERVICE | MPI::ERR_SERVICE | 29 |
| MPI_ERR_SIZE | MPI::ERR_SIZE | 30 |
| MPI_ERR_SPAWN | MPI::ERR_SPAWN | 31 |
| MPI_ERR_UNSUPPORTED_DATAREP | MPI::ERR_UNSUPPORTED_DATAREP | 32 |
| MPI_ERR_UNSUPPORTED_OPERATION | MPI::ERR_UNSUPPORTED_OPERATION | 33 |
| MPI_ERR_WIN | MPI::ERR_WIN | 34 |
| MPI_ERR_LASTCODE | MPI::ERR_LASTCODE | 35 |

36

37

38

39

40

41

42

43

44

45

46

47

48

Assorted Constants

| | |
|--------------------|---------------------|
| MPI_BOTTOM | MPI::BOTTOM |
| MPI_PROC_NULL | MPI::PROC_NULL |
| MPI_ANY_SOURCE | MPI::ANY_SOURCE |
| MPI_ANY_TAG | MPI::ANY_TAG |
| MPI_UNDEFINED | MPI::UNDEFINED |
| MPI_BSEND_OVERHEAD | MPI::BSEND_OVERHEAD |
| MPI_KEYVAL_INVALID | MPI::KEYVAL_INVALID |
| MPI_IN_PLACE | MPI::IN_PLACE |
| MPI_LOCK_EXCLUSIVE | MPI::LOCK_EXCLUSIVE |
| MPI_LOCK_SHARED | MPI::LOCK_SHARED |
| MPI_ROOT | MPI::ROOT |

Status size and reserved index values (Fortran)

| | |
|-----------------|---------------------|
| MPI_STATUS_SIZE | Not defined for C++ |
| MPI_SOURCE | Not defined for C++ |
| MPI_TAG | Not defined for C++ |
| MPI_ERROR | Not defined for C++ |

Variable Address Size (Fortran only)

| | |
|------------------|---------------------|
| MPI_ADDRESS_KIND | Not defined for C++ |
| MPI_INTEGER_KIND | Not defined for C++ |
| MPI_OFFSET_KIND | Not defined for C++ |

Error-handling specifiers

| | |
|----------------------|-----------------------|
| MPI_ERRORS_ARE_FATAL | MPI::ERRORS_ARE_FATAL |
| MPI_ERRORS_RETURN | MPI::ERRORS_RETURN |

Maximum Sizes for Strings

| | |
|------------------------|-------------------------|
| MPI_MAX_PROCESSOR_NAME | MPI::MAX_PROCESSOR_NAME |
| MPI_MAX_ERROR_STRING | MPI::MAX_ERROR_STRING |
| MPI_MAX_DATAREP_STRING | MPI::MAX_DATAREP_STRING |
| MPI_MAX_INFO_KEY | MPI::MAX_INFO_KEY |
| MPI_MAX_INFO_VAL | MPI::MAX_INFO_VAL |
| MPI_MAX_OBJECT_NAME | MPI::MAX_OBJECT_NAME |
| MPI_MAX_PORT_NAME | MPI::MAX_PORT_NAME |

| Named Predefined Datatypes | | C/C++ types | |
|----------------------------|-------------------------|-------------------------|----|
| MPI_CHAR | MPI::CHAR | signed char | 1 |
| MPI_SHORT | MPI::SHORT | signed short int | 2 |
| MPI_INT | MPI::INT | signed int | 3 |
| MPI_LONG | MPI::LONG | signed long | 4 |
| MPI_LONG_LONG_INT | MPI::LONG_LONG_INT | signed long long | 5 |
| MPI_LONG_LONG | MPI::LONG_LONG | long long (synonym) | 6 |
| MPI_UNSIGNED_CHAR | MPI::UNSIGNED_CHAR | unsigned char | 7 |
| MPI_UNSIGNED_SHORT | MPI::UNSIGNED_SHORT | unsigned short | 8 |
| MPI_UNSIGNED | MPI::UNSIGNED | unsigned int | 9 |
| MPI_UNSIGNED_LONG | MPI::UNSIGNED_LONG | unsigned long | 10 |
| MPI_UNSIGNED_LONG_LONG | MPI::UNSIGNED_LONG_LONG | unsigned long long | 11 |
| MPI_FLOAT | MPI::FLOAT | float | 12 |
| MPI_DOUBLE | MPI::DOUBLE | double | 13 |
| MPI_LONG_DOUBLE | MPI::LONG_DOUBLE | long double | 14 |
| MPI_WCHAR | MPI::WCHAR | wchar_t | 15 |
| | | (defined in <stddef.h>) | 16 |
| MPI_BYTE | MPI::BYTE | | 17 |
| MPI_PACKED | MPI::PACKED | | 18 |

C and C++ (no Fortran) Named Predefined Datatypes

| | |
|----------|-----------|
| MPI_Fint | MPI::Fint |
|----------|-----------|

| Named Predefined Datatypes | | Fortran types | |
|----------------------------|-----------------------|------------------|----|
| MPI_INTEGER | MPI::INTEGER | INTEGER | 21 |
| MPI_REAL | MPI::REAL | REAL | 22 |
| MPI_DOUBLE_PRECISION | MPI::DOUBLE_PRECISION | DOUBLE PRECISION | 23 |
| MPI_COMPLEX | MPI::COMPLEX | COMPLEX | 24 |
| MPI_LOGICAL | MPI::LOGICAL | LOGICAL | 25 |
| MPI_CHARACTER | MPI::CHARACTER | CHARACTER(1) | 26 |
| MPI_BYTE | MPI::BYTE | | 27 |
| MPI_PACKED | MPI::PACKED | | 28 |

Optional C and C++ (no Fortran) Named Predefined Datatypes

| | | |
|-----------------|------------------|-------------|
| MPI_SIGNED_CHAR | MPI::SIGNED_CHAR | signed char |
|-----------------|------------------|-------------|

| Optional datatypes (Fortran) | | Fortran types | |
|------------------------------|---------------------|----------------|----|
| MPI_DOUBLE_COMPLEX | MPI::DOUBLE_COMPLEX | DOUBLE COMPLEX | 29 |
| MPI_INTEGER1 | MPI::INTEGER1 | INTEGER*1 | 30 |
| MPI_INTEGER2 | MPI::INTEGER2 | INTEGER*2 | 31 |
| MPI_INTEGER4 | MPI::INTEGER4 | INTEGER*4 | 32 |
| MPI_REAL2 | MPI::REAL2 | REAL*2 | 33 |
| MPI_REAL4 | MPI::REAL4 | REAL*4 | 34 |
| MPI_REAL8 | MPI::REAL8 | REAL*8 | 35 |

Datatypes for reduction functions (C and C++)

| | |
|---------------------|----------------------|
| MPI_FLOAT_INT | MPI::FLOAT_INT |
| MPI_DOUBLE_INT | MPI::DOUBLE_INT |
| MPI_LONG_INT | MPI::LONG_INT |
| MPI_2INT | MPI::2INT |
| MPI_SHORT_INT | MPI::SHORT_INT |
| MPI_LONG_DOUBLE_INT | MPI::LONG_DOUBLE_INT |

Datatypes for reduction functions (Fortran)

| | |
|-----------------------|------------------------|
| MPI_2REAL | MPI::2REAL |
| MPI_2DOUBLE_PRECISION | MPI::2DOUBLE_PRECISION |
| MPI_2INTEGER | MPI::2INTEGER |

Special datatypes for constructing derived datatypes

| | |
|--------|---------|
| MPI_UB | MPI::UB |
| MPI_LB | MPI::LB |

Reserved communicators

| | |
|----------------|-----------------|
| MPI_COMM_WORLD | MPI::COMM_WORLD |
| MPI_COMM_SELF | MPI::COMM_SELF |

Results of communicator and group comparisons

| | |
|---------------|----------------|
| MPI_IDENT | MPI::IDENT |
| MPI_CONGRUENT | MPI::CONGRUENT |
| MPI_SIMILAR | MPI::SIMILAR |
| MPI_UNEQUAL | MPI::UNEQUAL |

Environmental inquiry keys

| | |
|---------------------|----------------------|
| MPI_TAG_UB | MPI::TAG_UB |
| MPI_IO | MPI::IO |
| MPI_HOST | MPI::HOST |
| MPI_WTIME_IS_GLOBAL | MPI::WTIME_IS_GLOBAL |

Collective Operations

| | | |
|--------------------|---------------------|----|
| MPI_MAX | MPI::MAX | 1 |
| MPI_MIN | MPI::MIN | 2 |
| MPI_SUM | MPI::SUM | 3 |
| MPI_PROD | MPI::PROD | 4 |
| MPI_MAXLOC | MPI::MAXLOC | 5 |
| MPI_MINLOC | MPI::MINLOC | 6 |
| MPI_BAND | MPI::BAND | 7 |
| MPI_BOR | MPI::BOR | 8 |
| MPI_BXOR | MPI::BXOR | 9 |
| MPI_LAND | MPI::LAND | 10 |
| MPI_LOR | MPI::LOR | 11 |
| MPI_LXOR | MPI::LXOR | 12 |
| <u>MPI_REPLACE</u> | <u>MPI::REPLACE</u> | 13 |

Null Handles

| | | |
|----------------------|-----------------------|----|
| MPI_GROUP_NULL | MPI::GROUP_NULL | 14 |
| MPI_COMM_NULL | MPI::COMM_NULL | 15 |
| MPI_DATATYPE_NULL | MPI::DATATYPE_NULL | 16 |
| MPI_REQUEST_NULL | MPI::REQUEST_NULL | 17 |
| MPI_OP_NULL | MPI::OP_NULL | 18 |
| MPI_ERRHANDLER_NULL | MPI::ERRHANDLER_NULL | 19 |
| <u>MPI_FILE_NULL</u> | <u>MPI::FILE_NULL</u> | 20 |
| <u>MPI_INFO_NULL</u> | <u>MPI::INFO_NULL</u> | 21 |
| <u>MPI_WIN_NULL</u> | <u>MPI::WIN_NULL</u> | 22 |

Empty group

| | | |
|------------------------|-------------------------|----|
| <u>MPI_GROUP_EMPTY</u> | <u>MPI::GROUP_EMPTY</u> | 23 |
|------------------------|-------------------------|----|

Topologies

| | | |
|-----------|------------|----|
| MPI_GRAPH | MPI::GRAPH | 24 |
| MPI_CART | MPI::CART | 25 |

Predefined functions

| | | |
|--------------------|---------------------|----|
| MPI_NULL_COPY_FN | MPI::NULL_COPY_FN | 26 |
| MPI_NULL_DELETE_FN | MPI::NULL_DELETE_FN | 27 |
| MPI_DUP_FN | MPI::DUP_FN | 28 |

Predefined Attribute Keys

| | |
|-------------------|--------------------|
| MPI_APPNUM | MPI::APPNUM |
| MPI_LASTUSEDCODE | MPI::LASTUSEDCODE |
| MPI_UNIVERSE_SIZE | MPI::UNIVERSE_SIZE |
| MPI_WIN_BASE | MPI::WIN_BASE |
| MPI_WIN_DISP_UNIT | MPI::WIN_DISP_UNIT |
| MPI_WIN_SIZE | MPI::WIN_SIZE |

Mode Constants

| | |
|--------------------------|---------------------------|
| MPI_MODE_APPEND | MPI::MODE_APPEND |
| MPI_MODE_CREATE | MPI::MODE_CREATE |
| MPI_MODE_DELETE_ON_CLOSE | MPI::MODE_DELETE_ON_CLOSE |
| MPI_MODE_EXCL | MPI::MODE_EXCL |
| MPI_MODE_NOCHECK | MPI::MODE_NOCHECK |
| MPI_MODE_NOPRECEDE | MPI::MODE_NOPRECEDE |
| MPI_MODE_NOPUT | MPI::MODE_NOPUT |
| MPI_MODE_NOSTORE | MPI::MODE_NOSTORE |
| MPI_MODE_NOSUCCEED | MPI::MODE_NOSUCCEED |
| MPI_MODE_RDONLY | MPI::MODE_RDONLY |
| MPI_MODE_RDWR | MPI::MODE_RDWR |
| MPI_MODE_SEQUENTIAL | MPI::MODE_SEQUENTIAL |
| MPI_MODE_UNIQUE_OPEN | MPI::MODE_UNIQUE_OPEN |
| MPI_MODE_WRONLY | MPI::MODE_WRONLY |

Datatype Decoding Constants

| | |
|-------------------------------|--------------------------------|
| MPI_COMBINER_CONTIGUOUS | MPI::COMBINER_CONTIGUOUS |
| MPI_COMBINER_DARRAY | MPI::COMBINER_DARRAY |
| MPI_COMBINER_DUP | MPI::COMBINER_DUP |
| MPI_COMBINER_F90_COMPLEX | MPI::COMBINER_F90_COMPLEX |
| MPI_COMBINER_F90_INTEGER | MPI::COMBINER_F90_INTEGER |
| MPI_COMBINER_F90_REAL | MPI::COMBINER_F90_REAL |
| MPI_COMBINER_HINDEXED_INTEGER | MPI::COMBINER_HINDEXED_INTEGER |
| MPI_COMBINER_HINDEXED | MPI::COMBINER_HINDEXED |
| MPI_COMBINER_HVECTOR_INTEGER | MPI::COMBINER_HVECTOR_INTEGER |
| MPI_COMBINER_HVECTOR | MPI::COMBINER_HVECTOR |
| MPI_COMBINER_INDEXED_BLOCK | MPI::COMBINER_INDEXED_BLOCK |
| MPI_COMBINER_INDEXED | MPI::COMBINER_INDEXED |
| MPI_COMBINER_NAMED | MPI::COMBINER_NAMED |
| MPI_COMBINER_RESIZED | MPI::COMBINER_RESIZED |
| MPI_COMBINER_STRUCT_INTEGER | MPI::COMBINER_STRUCT_INTEGER |
| MPI_COMBINER_STRUCT | MPI::COMBINER_STRUCT |
| MPI_COMBINER_SUBARRAY | MPI::COMBINER_SUBARRAY |
| MPI_COMBINER_VECTOR | MPI::COMBINER_VECTOR |

Threads Constants

| | |
|-----------------------|------------------------|
| MPI_THREAD_FUNNELED | MPI::THREAD_FUNNELED |
| MPI_THREAD_MULTIPLE | MPI::THREAD_MULTIPLE |
| MPI_THREAD_SERIALIZED | MPI::THREAD_SERIALIZED |
| MPI_THREAD_SINGLE | MPI::THREAD_SINGLE |

File Operation Constants

| | |
|--------------------------|---------------------------|
| MPI_DISPLACEMENT_CURRENT | MPI::DISPLACEMENT_CURRENT |
| MPI_DISTRIBUTE_BLOCK | MPI::DISTRIBUTE_BLOCK |
| MPI_DISTRIBUTE_CYCLIC | MPI::DISTRIBUTE_CYCLIC |
| MPI_DISTRIBUTE_DFLT_DARG | MPI::DISTRIBUTE_DFLT_DARG |
| MPI_DISTRIBUTE_NONE | MPI::DISTRIBUTE_NONE |
| MPI_ORDER_C | MPI::ORDER_C |
| MPI_ORDER_FORTRAN | MPI::ORDER_FORTRAN |
| MPI_SEEK_CUR | MPI::SEEK_CUR |
| MPI_SEEK_END | MPI::SEEK_END |
| MPI_SEEK_SET | MPI::SEEK_SET |

F90 Datatype Matching Constants

| | |
|-----------------------|------------------------|
| MPI_TYPECLASS_COMPLEX | MPI::TYPECLASS_COMPLEX |
| MPI_TYPECLASS_INTEGER | MPI::TYPECLASS_INTEGER |
| MPI_TYPECLASS_REAL | MPI::TYPECLASS_REAL |

Handles to Assorted Structures in C and C++ (no Fortran)

| | |
|----------|-----------|
| MPI_File | MPI::File |
| MPI_Info | MPI::Info |
| MPI_Win | MPI::Win |

Constants Specifying Empty or Ignored Input

| | |
|---------------------|---------------------|
| MPI_ARGVS_NULL | MPI::ARGVS_NULL |
| MPI_ARGV_NULL | MPI::ARGV_NULL |
| MPI_ERRCODES_IGNORE | Not defined for C++ |
| MPI_STATUSES_IGNORE | Not defined for C++ |
| MPI_STATUS_IGNORE | Not defined for C++ |

C Constants Specifying Ignored Input (no C++ or Fortran)

| | |
|-----------------------|---------------------|
| MPI_F_STATUSES_IGNORE | Not defined for C++ |
| MPI_F_STATUS_IGNORE | Not defined for C++ |

C and C++ preprocessor Constants and Fortran Parameters

| |
|----------------|
| MPI_SUBVERSION |
| MPI_VERSION |

A.1.2 Type and prototype definitions

The following are defined C type definitions, also included in the file `mpi.h`.

```

1  /* opaque types (C) */
2  MPI_Aint
3  MPI_Status
4
5  /* handles to assorted structures (C) */
6  MPI_Group
7  MPI_Comm
8  MPI_Datatype
9  MPI_Request
10 MPI_Op
11 MPI_Errhandler
12
13 /* prototypes for user-defined functions (C) */
14 typedef int MPI_Copy_function(MPI_Comm oldcomm, int keyval,
15                               void *extra_state, void *attribute_val_in,
16                               void *attribute_val_out, int *flag);
17 typedef int MPI_Delete_function(MPI_Comm comm, int keyval,
18                                 void *attribute_val, void *extra_state)
19 typedef void MPI_Handler_function(MPI_Comm *, int *, ...);
20 typedef void MPI_User_function( void *invec, void *inoutvec, int *len,
21                                 MPI_Datatype *datatype);
22
23
24
25

```

For Fortran, here are examples of how each of the user-defined functions should be declared.

The user-function argument to `MPI_OP_CREATE` should be declared like this:

```

26 SUBROUTINE USER_FUNCTION( INVEC, INOUTVEC, LEN, TYPE)
27   <type> INVEC(LEN), INOUTVEC(LEN)
28   INTEGER LEN, TYPE
29

```

The copy-function argument to `MPI_KEYVAL_CREATE` should be declared like this:

```

30 SUBROUTINE COPY_FUNCTION(OLDCOMM, KEYVAL, EXTRA_STATE,
31                          ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERR)
32   INTEGER OLDCOMM, KEYVAL, EXTRA_STATE, ATTRIBUTE_VAL_IN,
33     ATTRIBUTE_VAL_OUT, IERR
34   LOGICAL FLAG
35

```

The delete-function argument to `MPI_KEYVAL_CREATE` should be declared like this:

```

36 SUBROUTINE DELETE_FUNCTION(COMM, KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE, IERR)
37   INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE, IERR
38

```

The handler-function for error handlers should be declared like this:

```

39 SUBROUTINE HANDLER_FUNCTION(COMM, ERROR_CODE, ..... )
40   INTEGER COMM, ERROR_CODE
41

```

A.1.3 Info Keys

| | |
|----------------------|----|
| | 1 |
| | 2 |
| access_style | 3 |
| appnum | 4 |
| arch | 5 |
| cb_block_size | 6 |
| cb_buffer_size | 7 |
| cb_nodes | 8 |
| chunked_item | 9 |
| chunked_size | 10 |
| chunked | 11 |
| collective_buffering | 12 |
| file_perm | 13 |
| filename | 14 |
| file | 15 |
| host | 16 |
| io_node_list | 17 |
| ip_address | 18 |
| ip_port | 19 |
| nb_proc | 20 |
| no_locks | 21 |
| num_io_nodes | 22 |
| path | 23 |
| soft | 24 |
| striping_factor | 25 |
| striping_unit | 26 |
| wdir | 27 |

A.1.4 Info Values

| | |
|--------------------|----|
| | 28 |
| | 29 |
| | 30 |
| false | 31 |
| random | 32 |
| read_mostly | 33 |
| read_once | 34 |
| reverse_sequential | 35 |
| sequential | 36 |
| true | 37 |
| write_mostly | 38 |
| write_once | 39 |

40

41

42

43

44

45

46

47

48

A.2 C Bindings

A.2.1 Point-to-Point Communication C Bindings

```
1  int MPI_Bsend(void* buf, int count, MPI_Datatype datatype, int dest,
2      int tag, MPI_Comm comm)
3
4  int MPI_Bsend_init(void* buf, int count, MPI_Datatype datatype, int dest,
5      int tag, MPI_Comm comm, MPI_Request *request)
6
7  int MPI_Buffer_attach( void* buffer, int size)
8
9  int MPI_Buffer_detach( void* buffer_addr, int* size)
10
11 int MPI_Cancel(MPI_Request *request)
12
13 int MPI_Get_address(void *location, MPI_Aint *address)
14
15 int MPI_Get_count(MPI_Status *status, MPI_Datatype datatype, int *count)
16
17 int MPI_Get_elements(MPI_Status *status, MPI_Datatype datatype, int *count)
18
19 int MPI_Ibsend(void* buf, int count, MPI_Datatype datatype, int dest,
20     int tag, MPI_Comm comm, MPI_Request *request)
21
22 int MPI_Iprobe(int source, int tag, MPI_Comm comm, int *flag,
23     MPI_Status *status)
24
25 int MPI_Irecv(void* buf, int count, MPI_Datatype datatype, int source,
26     int tag, MPI_Comm comm, MPI_Request *request)
27
28 int MPI_Irsend(void* buf, int count, MPI_Datatype datatype, int dest,
29     int tag, MPI_Comm comm, MPI_Request *request)
30
31 int MPI_Isend(void* buf, int count, MPI_Datatype datatype, int dest,
32     int tag, MPI_Comm comm, MPI_Request *request)
33
34 int MPI_issend(void* buf, int count, MPI_Datatype datatype, int dest,
35     int tag, MPI_Comm comm, MPI_Request *request)
36
37 int MPI_Pack_external(char *datarep, void *inbuf, int incount,
38     MPI_Datatype datatype, void *outbuf, MPI_Aint outsize,
39     MPI_Aint *position)
40
41 int MPI_Pack_external_size(char *datarep, int incount,
42     MPI_Datatype datatype, MPI_Aint *size)
43
44 int MPI_Pack_size(int incount, MPI_Datatype datatype, MPI_Comm comm,
45     int *size)
46
47 int MPI_Probe(int source, int tag, MPI_Comm comm, MPI_Status *status)
48
49 int MPI_Recv(void* buf, int count, MPI_Datatype datatype, int source,
50     int tag, MPI_Comm comm, MPI_Status *status)
```

```
int MPI_Recv_init(void* buf, int count, MPI_Datatype datatype, int source, 1
                  int tag, MPI_Comm comm, MPI_Request *request) 2
int MPI_Request_free(MPI_Request *request) 3
int MPI_Request_get_status(MPI_Request request, int *flag, 4
                           MPI_Status *status) 5
int MPI_Rsend(void* buf, int count, MPI_Datatype datatype, int dest, 6
              int tag, MPI_Comm comm) 7
int MPI_Rsend_init(void* buf, int count, MPI_Datatype datatype, int dest, 8
                  int tag, MPI_Comm comm, MPI_Request *request) 9
int MPI_Send(void* buf, int count, MPI_Datatype datatype, int dest, 10
             int tag, MPI_Comm comm) 11
int MPI_Send_init(void* buf, int count, MPI_Datatype datatype, int dest, 12
                 int tag, MPI_Comm comm, MPI_Request *request) 13
int MPI_Sendrecv(void *sendbuf, int sendcount, MPI_Datatype sendtype, 14
                int dest, int sendtag, void *recvbuf, int recvcount, 15
                MPI_Datatype recvtype, int source, int recvtag, MPI_Comm comm, 16
                MPI_Status *status) 17
int MPI_Sendrecv_replace(void* buf, int count, MPI_Datatype datatype, 18
                        int dest, int sendtag, int source, int recvtag, MPI_Comm comm, 19
                        MPI_Status *status) 20
int MPI_Ssend(void* buf, int count, MPI_Datatype datatype, int dest, 21
              int tag, MPI_Comm comm) 22
int MPI_Ssend_init(void* buf, int count, MPI_Datatype datatype, int dest, 23
                  int tag, MPI_Comm comm, MPI_Request *request) 24
int MPI_Start(MPI_Request *request) 25
int MPI_Startall(int count, MPI_Request *array_of_requests) 26
int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status) 27
int MPI_Test_cancelled(MPI_Status *status, int *flag) 28
int MPI_Testall(int count, MPI_Request *array_of_requests, int *flag, 29
               MPI_Status *array_of_statuses) 30
int MPI_Testany(int count, MPI_Request *array_of_requests, int *index, 31
               int *flag, MPI_Status *status) 32
int MPI_Testsome(int incount, MPI_Request *array_of_requests, int *outcount, 33
                int *array_of_indices, MPI_Status *array_of_statuses) 34
int MPI_Type_commit(MPI_Datatype *datatype) 35
int MPI_Type_contiguous(int count, MPI_Datatype oldtype, 36
                       MPI_Datatype *newtype) 37

```

```
1 int MPI_Type_create_darray(int size, int rank, int ndims,
2     int array_of_gsizes[], int array_of_distribs[], int
3     array_of_dargs[], int array_of_psizes[], int order,
4     MPI_Datatype oldtype, MPI_Datatype *newtype)
5
6 int MPI_Type_create_hindexed(int count, int array_of_blocklengths[],
7     MPI_Aint array_of_displacements[], MPI_Datatype oldtype,
8     MPI_Datatype *newtype)
9
10 int MPI_Type_create_hvector(int count, int blocklength, MPI_Aint stride,
11     MPI_Datatype oldtype, MPI_Datatype *newtype)
12
13 int MPI_Type_create_indexed_block(int count, int blocklength,
14     int array_of_displacements[], MPI_Datatype oldtype,
15     MPI_Datatype *newtype)
16
17 int MPI_Type_create_resized(MPI_Datatype oldtype, MPI_Aint lb, MPI_Aint
18     extent, MPI_Datatype *newtype)
19
20 int MPI_Type_create_struct(int count, int array_of_blocklengths[],
21     MPI_Aint array_of_displacements[],
22     MPI_Datatype array_of_types[], MPI_Datatype *newtype)
23
24 int MPI_Type_create_subarray(int ndims, int array_of_sizes[],
25     int array_of_subsizes[], int array_of_starts[], int order,
26     MPI_Datatype oldtype, MPI_Datatype *newtype)
27
28 int MPI_Type_dup(MPI_Datatype type, MPI_Datatype *newtype)
29
30 int MPI_Type_free(MPI_Datatype *datatype)
31
32 int MPI_Type_get_extent(MPI_Datatype datatype, MPI_Aint *lb,
33     MPI_Aint *extent)
34
35 int MPI_Type_get_true_extent(MPI_Datatype datatype, MPI_Aint *true_lb,
36     MPI_Aint *true_extent)
37
38 int MPI_Type_indexed(int count, int *array_of_blocklengths,
39     int *array_of_displacements, MPI_Datatype oldtype,
40     MPI_Datatype *newtype)
41
42 int MPI_Type_size(MPI_Datatype datatype, int *size)
43
44 int MPI_Type_vector(int count, int blocklength, int stride,
45     MPI_Datatype oldtype, MPI_Datatype *newtype)
46
47 int MPI_Unpack(void* inbuf, int insize, int *position, void *outbuf,
48     int outcount, MPI_Datatype datatype, MPI_Comm comm)
49
50 int MPI_Unpack_external(char *datarep, void *inbuf, MPI_Aint insize,
51     MPI_Aint *position, void *outbuf, int outcount,
52     MPI_Datatype datatype)
53
54 int MPI_Wait(MPI_Request *request, MPI_Status *status)
```

```
int MPI_Waitall(int count, MPI_Request *array_of_requests,
               MPI_Status *array_of_statuses)
int MPI_Waitany(int count, MPI_Request *array_of_requests, int *index,
               MPI_Status *status)
int MPI_Waitsome(int incount, MPI_Request *array_of_requests, int *outcount,
                int *array_of_indices, MPI_Status *array_of_statuses)

A.2.2 Collective Communication C Bindings

int MPI_Allgather(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                 void* recvbuf, int recvcount, MPI_Datatype recvttype,
                 MPI_Comm comm)
int MPI_Allgatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                  void* recvbuf, int *recvcounts, int *displs,
                  MPI_Datatype recvttype, MPI_Comm comm)
int MPI_Allreduce(void* sendbuf, void* recvbuf, int count,
                 MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
int MPI_Alltoall(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                 void* recvbuf, int recvcount, MPI_Datatype recvttype,
                 MPI_Comm comm)
int MPI_Alltoallv(void* sendbuf, int *sendcounts, int *sdispls,
                 MPI_Datatype sendtype, void* recvbuf, int *recvcounts,
                 int *rdispls, MPI_Datatype recvttype, MPI_Comm comm)
int MPI_Alltoallw(void *sendbuf, int sendcounts[], int sdispls[],
                 MPI_Datatype sendtypes[], void *recvbuf, int recvcounts[],
                 int rdispls[], MPI_Datatype recvtypes[], MPI_Comm comm)
int MPI_Barrier(MPI_Comm comm )
int MPI_Bcast(void* buffer, int count, MPI_Datatype datatype, int root,
             MPI_Comm comm )
int MPI_Exscan(void *sendbuf, void *recvbuf, int count,
              MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
int MPI_Gather(void* sendbuf, int sendcount, MPI_Datatype sendtype,
              void* recvbuf, int recvcount, MPI_Datatype recvttype, int root,
              MPI_Comm comm)
int MPI_Gatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
               void* recvbuf, int *recvcounts, int *displs,
               MPI_Datatype recvttype, int root, MPI_Comm comm)
int MPI_Op_create(MPI_User_function *function, int commute, MPI_Op *op)
int MPI_Reduce(void* sendbuf, void* recvbuf, int count,
              MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)
```

```
1 int MPI_Reduce_scatter(void* sendbuf, void* recvbuf, int *recvcnts,
2     MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
3
4 int MPI_Scan(void* sendbuf, void* recvbuf, int count,
5     MPI_Datatype datatype, MPI_Op op, MPI_Comm comm )
6
7 int MPI_Scatter(void* sendbuf, int sendcount, MPI_Datatype sendtype,
8     void* recvbuf, int recvcnt, MPI_Datatype recvtpe, int root,
9     MPI_Comm comm)
10
11 int MPI_Scatterv(void* sendbuf, int *sendcounts, int *displs,
12     MPI_Datatype sendtype, void* recvbuf, int recvcnt,
13     MPI_Datatype recvtpe, int root, MPI_Comm comm)
14
15
```

A.2.3 Groups, Contexts, and Communicators C Bindings

```
17 int MPI_Comm_compare(MPI_Comm comm1, MPI_Comm comm2, int *result)
18
19 int MPI_Comm_create(MPI_Comm comm, MPI_Group group, MPI_Comm *newcomm)
20
21 int MPI_Comm_create_keyval(MPI_Comm_copy_attr_function *comm_copy_attr_fn,
22     MPI_Comm_delete_attr_function *comm_delete_attr_fn,
23     int *comm_keyval, void *extra_state)
24
25 int MPI_Comm_delete_attr(MPI_Comm comm, int comm_keyval)
26
27 int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)
28
29 int MPI_Comm_free(MPI_Comm *comm)
30
31 int MPI_Comm_free_keyval(int *comm_keyval)
32
33 int MPI_Comm_get_attr(MPI_Comm comm, int comm_keyval, void *attribute_val,
34     int *flag)
35
36 int MPI_Comm_group(MPI_Comm comm, MPI_Group *group)
37
38 int MPI_Comm_rank(MPI_Comm comm, int *rank)
39
40 int MPI_Comm_remote_group(MPI_Comm comm, MPI_Group *group)
41
42 int MPI_Comm_remote_size(MPI_Comm comm, int *size)
43
44 int MPI_Comm_set_attr(MPI_Comm comm, int comm_keyval, void *attribute_val)
45
46 int MPI_Comm_size(MPI_Comm comm, int *size)
47
48 int MPI_Comm_split(MPI_Comm comm, int color, int key, MPI_Comm *newcomm)
49
50 int MPI_Comm_test_inter(MPI_Comm comm, int *flag)
51
52 int MPI_Group_compare(MPI_Group group1, MPI_Group group2, int *result)
53
54 int MPI_Group_difference(MPI_Group group1, MPI_Group group2,
55     MPI_Group *newgroup)
```



```
int MPI_Group_excl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup) 1
int MPI_Group_free(MPI_Group *group) 2
int MPI_Group_incl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup) 3
int MPI_Group_intersection(MPI_Group group1, MPI_Group group2, 4
    MPI_Group *newgroup) 5
int MPI_Group_range_excl(MPI_Group group, int n, int ranges[][3], 6
    MPI_Group *newgroup) 7
int MPI_Group_range_incl(MPI_Group group, int n, int ranges[][3], 8
    MPI_Group *newgroup) 9
int MPI_Group_rank(MPI_Group group, int *rank) 10
int MPI_Group_size(MPI_Group group, int *size) 11
int MPI_Group_translate_ranks (MPI_Group group1, int n, int *ranks1, 12
    MPI_Group group2, int *ranks2) 13
int MPI_Group_union(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup) 14
int MPI_Intercomm_create(MPI_Comm local_comm, int local_leader, 15
    MPI_Comm peer_comm, int remote_leader, int tag, 16
    MPI_Comm *newintercomm) 17
int MPI_Intercomm_merge(MPI_Comm intercomm, int high, 18
    MPI_Comm *newintracomm) 19
int MPI_Type_create_keyval(MPI_Type_copy_attr_function *type_copy_attr_fn, 20
    MPI_Type_delete_attr_function *type_delete_attr_fn, 21
    int *type_keyval, void *extra_state) 22
int MPI_Type_delete_attr(MPI_Datatype type, int type_keyval) 23
int MPI_Type_free_keyval(int *type_keyval) 24
int MPI_Type_get_attr(MPI_Datatype type, int type_keyval, void 25
    *attribute_val, int *flag) 26
int MPI_Type_set_attr(MPI_Datatype type, int type_keyval, 27
    void *attribute_val) 28
int MPI_Win_create_keyval(MPI_Win_copy_attr_function *win_copy_attr_fn, 29
    MPI_Win_delete_attr_function *win_delete_attr_fn, 30
    int *win_keyval, void *extra_state) 31
int MPI_Win_delete_attr(MPI_Win win, int win_keyval) 32
int MPI_Win_free_keyval(int *win_keyval) 33
int MPI_Win_get_attr(MPI_Win win, int win_keyval, void *attribute_val, 34
    int *flag) 35
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 36
int MPI_Win_get_attr(MPI_Win win, int win_keyval, void *attribute_val, 37
    int *flag) 38
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 39
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 40
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 41
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 42
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 43
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 44
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 45
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 46
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 47
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 48
```

A.2.4 Process Topologies C Bindings

```
1 int MPI_Cart_coords(MPI_Comm comm, int rank, int maxdims, int *coords)
2
3 int MPI_Cart_create(MPI_Comm comm_old, int ndims, int *dims, int *periods,
4                     int reorder, MPI_Comm *comm_cart)
5
6 int MPI_Cart_get(MPI_Comm comm, int maxdims, int *dims, int *periods,
7                 int *coords)
8
9 int MPI_Cart_map(MPI_Comm comm, int ndims, int *dims, int *periods,
10                int *newrank)
11
12 int MPI_Cart_rank(MPI_Comm comm, int *coords, int *rank)
13
14 int MPI_Cart_shift(MPI_Comm comm, int direction, int disp, int *rank_source,
15                  int *rank_dest)
16
17 int MPI_Cart_sub(MPI_Comm comm, int *remain_dims, MPI_Comm *newcomm)
18
19 int MPI_Cartdim_get(MPI_Comm comm, int *ndims)
20
21 int MPI_Dims_create(int nnodes, int ndims, int *dims)
22
23 int MPI_Graph_create(MPI_Comm comm_old, int nnodes, int *index, int *edges,
24                    int reorder, MPI_Comm *comm_graph)
25
26 int MPI_Graph_get(MPI_Comm comm, int maxindex, int maxedges, int *index,
27                 int *edges)
28
29 int MPI_Graph_map(MPI_Comm comm, int nnodes, int *index, int *edges,
30                 int *newrank)
31
32 int MPI_Graph_neighbors(MPI_Comm comm, int rank, int maxneighbors,
33                        int *neighbors)
34
35 int MPI_Graph_neighbors_count(MPI_Comm comm, int rank, int *nneighbors)
36
37 int MPI_Graphdims_get(MPI_Comm comm, int *nnodes, int *nedges)
38
39 int MPI_Topo_test(MPI_Comm comm, int *status)
```

A.2.5 MPI Environmenta Management C Bindings

```
1 int MPI_Abort(MPI_Comm comm, int errorcode)
2
3 int MPI_Alloc_mem(MPI_Aint size, MPI_Info info, void *baseptr)
4
5 int MPI_Comm_create_errhandler(MPI_Comm_errhandler_fn *function,
6                               MPI_Errhandler *errhandler)
7
8 int MPI_Comm_get_errhandler(MPI_Comm comm, MPI_Errhandler *errhandler)
9
10 int MPI_Comm_set_errhandler(MPI_Comm comm, MPI_Errhandler errhandler)
11
12 int MPI_Errhandler_free(MPI_Errhandler *errhandler)
13
14 int MPI_Error_class(int errorcode, int *errorclass)
```

```
int MPI_Error_string(int errorcode, char *string, int *resultlen) 1
int MPI_File_create_errhandler(MPI_File_errhandler_fn *function, 2
                             MPI_Errhandler *errhandler) 3
int MPI_File_get_errhandler(MPI_File file, MPI_Errhandler *errhandler) 5
int MPI_File_set_errhandler(MPI_File file, MPI_Errhandler errhandler) 7
int MPI_Finalize(void) 8
int MPI_Finalized(int *flag) 9
int MPI_Free_mem(void *base) 11
int MPI_Get_processor_name(char *name, int *resultlen) 13
int MPI_Get_version(int *version, int *subversion) 14
int MPI_Init(int *argc, char ***argv) 16
int MPI_Initialized(int *flag) 17
int MPI_Win_create_errhandler(MPI_Win_errhandler_fn *function, MPI_Errhandler 19
                             *errhandler) 20
int MPI_Win_get_errhandler(MPI_Win win, MPI_Errhandler *errhandler) 21
int MPI_Win_set_errhandler(MPI_Win win, MPI_Errhandler errhandler) 23
double MPI_Wtick(void) 24
double MPI_Wtime(void) 26
```

A.2.6 Miscellany C Bindings

```
int MPI_Info_create(MPI_Info *info) 29
int MPI_Info_delete(MPI_Info info, char *key) 31
int MPI_Info_dup(MPI_Info info, MPI_Info *newinfo) 33
int MPI_Info_free(MPI_Info *info) 34
int MPI_Info_get(MPI_Info info, char *key, int valuelen, char *value, 36
                int *flag) 37
int MPI_Info_get_nkeys(MPI_Info info, int *nkeys) 38
int MPI_Info_get_nthkey(MPI_Info info, int n, char *key) 39
int MPI_Info_get_valuelen(MPI_Info info, char *key, int *valuelen, 40
                          int *flag) 41
int MPI_Info_set(MPI_Info info, char *key, char *value) 42
```

A.2.7 Process Creation and Management C Bindings

```
1 int MPI_Close_port(char *port_name)
2
3 int MPI_Comm_accept(char *port_name, MPI_Info info, int root, MPI_Comm comm,
4 MPI_Comm *newcomm)
5
6 int MPI_Comm_connect(char *port_name, MPI_Info info, int root,
7 MPI_Comm comm, MPI_Comm *newcomm)
8
9 int MPI_Comm_disconnect(MPI_Comm *comm)
10
11 int MPI_Comm_get_parent(MPI_Comm *parent)
12
13 int MPI_Comm_join(int fd, MPI_Comm *intercomm)
14
15 int MPI_Comm_spawn(char *command, char *argv[], int maxprocs, MPI_Info info,
16 int root, MPI_Comm comm, MPI_Comm *intercomm,
17 int array_of_errcodes[])
18
19 int MPI_Comm_spawn_multiple(int count, char *array_of_commands[],
20 char **array_of_argv[], int array_of_maxprocs[],
21 MPI_Info array_of_info[], int root, MPI_Comm comm,
22 MPI_Comm *intercomm, int array_of_errcodes[])
23
24 int MPI_Lookup_name(char *service_name, MPI_Info info, char *port_name)
25
26 int MPI_Open_port(MPI_Info info, char *port_name)
27
28 int MPI_Publish_name(char *service_name, MPI_Info info, char *port_name)
29
30 int MPI_Unpublish_name(char *service_name, MPI_Info info, char *port_name)
```

A.2.8 One-Sided Communications C Bindings

```
31 int MPI_Accumulate(void *origin_addr, int origin_count,
32 MPI_Datatype origin_datatype, int target_rank,
33 MPI_Aint target_disp, int target_count,
34 MPI_Datatype target_datatype, MPI_Op op, MPI_Win win)
35
36 int MPI_Get(void *origin_addr, int origin_count, MPI_Datatype
37 origin_datatype, int target_rank, MPI_Aint target_disp, int
38 target_count, MPI_Datatype target_datatype, MPI_Win win)
39
40 int MPI_Put(void *origin_addr, int origin_count, MPI_Datatype
41 origin_datatype, int target_rank, MPI_Aint target_disp, int
42 target_count, MPI_Datatype target_datatype, MPI_Win win)
43
44 int MPI_Win_complete(MPI_Win win)
45
46 int MPI_Win_create(void *base, MPI_Aint size, int disp_unit, MPI_Info info,
47 MPI_Comm comm, MPI_Win *win)
48
49 int MPI_Win_fence(int assert, MPI_Win win)
50
51 int MPI_Win_free(MPI_Win *win)
```

```
int MPI_Win_get_group(MPI_Win win, MPI_Group *group) 1
int MPI_Win_lock(int lock_type, int rank, int assert, MPI_Win win) 2
int MPI_Win_post(MPI_Group group, int assert, MPI_Win win) 3
int MPI_Win_start(MPI_Group group, int assert, MPI_Win win) 4
int MPI_Win_test(MPI_Win win, int *flag) 5
int MPI_Win_unlock(int rank, MPI_Win win) 6
int MPI_Win_wait(MPI_Win win) 7
```

A.2.9 External Interfaces C Bindings

```
int MPI_Add_error_class(int *errorclass) 8
int MPI_Add_error_code(int errorclass, int *errorcode) 9
int MPI_Add_error_string(int errorcode, char *string) 10
int MPI_Comm_call_errhandler(MPI_Comm comm, int errorcode) 11
int MPI_Comm_get_name(MPI_Comm comm, char *comm_name, int *resultlen) 12
int MPI_Comm_set_name(MPI_Comm comm, char *comm_name) 13
int MPI_File_call_errhandler(MPI_File fh, int errorcode) 14
int MPI_Grequest_complete(MPI_Request request) 15
int MPI_Grequest_start(MPI_Grequest_query_function *query_fn,
    MPI_Grequest_free_function *free_fn,
    MPI_Grequest_cancel_function *cancel_fn, void *extra_state,
    MPI_Request *request) 16
int MPI_Init_thread(int *argc, char ((*argv)[]), int required,
    int *provided) 17
int MPI_Is_thread_main(int *flag) 18
int MPI_Query_thread(int *provided) 19
int MPI_Status_set_cancelled(MPI_Status *status, int flag) 20
int MPI_Status_set_elements(MPI_Status *status, MPI_Datatype datatype,
    int count) 21
int MPI_Type_get_contents(MPI_Datatype datatype, int max_integers,
    int max_addresses, int max_datatypes, int array_of_integers[],
    MPI_Aint array_of_addresses[],
    MPI_Datatype array_of_datatypes[]) 22
int MPI_Type_get_envelope(MPI_Datatype datatype, int *num_integers,
    int *num_addresses, int *num_datatypes, int *combiner) 23
int MPI_Type_get_name(MPI_Datatype type, char *type_name, int *resultlen) 24
```

```
1 int MPI_Type_set_name(MPI_Datatype type, char *type_name)
2
3 int MPI_Win_call_errhandler(MPI_Win win, int errorcode)
4
5 int MPI_Win_get_name(MPI_Win win, char *win_name, int *resultlen)
6
7 int MPI_Win_set_name(MPI_Win win, char *win_name)
```

8 A.2.10 I/O C Bindings

```
9
10 int MPI_File_close(MPI_File *fh)
11
12 int MPI_File_delete(char *filename, MPI_Info info)
13
14 int MPI_File_get_amode(MPI_File fh, int *amode)
15
16 int MPI_File_get_atomicity(MPI_File fh, int *flag)
17
18 int MPI_File_get_byte_offset(MPI_File fh, MPI_Offset offset,
19                             MPI_Offset *disp)
20
21 int MPI_File_get_group(MPI_File fh, MPI_Group *group)
22
23 int MPI_File_get_info(MPI_File fh, MPI_Info *info_used)
24
25 int MPI_File_get_position(MPI_File fh, MPI_Offset *offset)
26
27 int MPI_File_get_position_shared(MPI_File fh, MPI_Offset *offset)
28
29 int MPI_File_get_size(MPI_File fh, MPI_Offset *size)
30
31 int MPI_File_get_type_extent(MPI_File fh, MPI_Datatype datatype,
32                             MPI_Aint *extent)
33
34 int MPI_File_get_view(MPI_File fh, MPI_Offset *disp, MPI_Datatype *etype,
35                     MPI_Datatype *filetype, char *datarep)
36
37 int MPI_File_iread(MPI_File fh, void *buf, int count, MPI_Datatype datatype,
38                 MPI_Request *request)
39
40 int MPI_File_iread_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
41                     MPI_Datatype datatype, MPI_Request *request)
42
43 int MPI_File_iread_shared(MPI_File fh, void *buf, int count,
44                         MPI_Datatype datatype, MPI_Request *request)
45
46 int MPI_File_iwrite(MPI_File fh, void *buf, int count,
47                   MPI_Datatype datatype, MPI_Request *request)
48
49 int MPI_File_iwrite_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
50                      MPI_Datatype datatype, MPI_Request *request)
51
52 int MPI_File_iwrite_shared(MPI_File fh, void *buf, int count,
53                          MPI_Datatype datatype, MPI_Request *request)
54
55 int MPI_File_open(MPI_Comm comm, char *filename, int amode, MPI_Info info,
56                 MPI_File *fh)
```

```
int MPI_File_preallocate(MPI_File fh, MPI_Offset size) 1
int MPI_File_read(MPI_File fh, void *buf, int count, MPI_Datatype datatype, 2
    MPI_Status *status) 3
int MPI_File_read_all(MPI_File fh, void *buf, int count, 4
    MPI_Datatype datatype, MPI_Status *status) 5
int MPI_File_read_all_begin(MPI_File fh, void *buf, int count, 6
    MPI_Datatype datatype) 7
int MPI_File_read_all_end(MPI_File fh, void *buf, MPI_Status *status) 8
int MPI_File_read_at(MPI_File fh, MPI_Offset offset, void *buf, int count, 9
    MPI_Datatype datatype, MPI_Status *status) 10
int MPI_File_read_at_all(MPI_File fh, MPI_Offset offset, void *buf, 11
    int count, MPI_Datatype datatype, MPI_Status *status) 12
int MPI_File_read_at_all_begin(MPI_File fh, MPI_Offset offset, void *buf, 13
    int count, MPI_Datatype datatype) 14
int MPI_File_read_at_all_end(MPI_File fh, void *buf, MPI_Status *status) 15
int MPI_File_read_ordered(MPI_File fh, void *buf, int count, 16
    MPI_Datatype datatype, MPI_Status *status) 17
int MPI_File_read_ordered_begin(MPI_File fh, void *buf, int count, 18
    MPI_Datatype datatype) 19
int MPI_File_read_ordered_end(MPI_File fh, void *buf, MPI_Status *status) 20
int MPI_File_read_shared(MPI_File fh, void *buf, int count, 21
    MPI_Datatype datatype, MPI_Status *status) 22
int MPI_File_seek(MPI_File fh, MPI_Offset offset, int whence) 23
int MPI_File_seek_shared(MPI_File fh, MPI_Offset offset, int whence) 24
int MPI_File_set_atomicity(MPI_File fh, int flag) 25
int MPI_File_set_info(MPI_File fh, MPI_Info info) 26
int MPI_File_set_size(MPI_File fh, MPI_Offset size) 27
int MPI_File_set_view(MPI_File fh, MPI_Offset disp, MPI_Datatype etype, 28
    MPI_Datatype filetype, char *datarep, MPI_Info info) 29
int MPI_File_sync(MPI_File fh) 30
int MPI_File_write(MPI_File fh, void *buf, int count, MPI_Datatype datatype, 31
    MPI_Status *status) 32
int MPI_File_write_all(MPI_File fh, void *buf, int count, 33
    MPI_Datatype datatype, MPI_Status *status) 34
int MPI_File_write_all_begin(MPI_File fh, void *buf, int count, 35
    MPI_Datatype datatype) 36
int MPI_File_write_all_end(MPI_File fh, void *buf, MPI_Status *status) 37
int MPI_File_write_at(MPI_File fh, MPI_Offset offset, void *buf, int count, 38
    MPI_Datatype datatype, MPI_Status *status) 39
int MPI_File_write_at_all(MPI_File fh, MPI_Offset offset, void *buf, int count, 40
    MPI_Datatype datatype, MPI_Status *status) 41
int MPI_File_write_at_all_begin(MPI_File fh, MPI_Offset offset, void *buf, int count, 42
    MPI_Datatype datatype) 43
int MPI_File_write_at_all_end(MPI_File fh, void *buf, MPI_Status *status) 44
int MPI_File_write_shared(MPI_File fh, void *buf, int count, MPI_Datatype datatype, 45
    MPI_Status *status) 46
int MPI_File_write_shared_begin(MPI_File fh, void *buf, int count, MPI_Datatype datatype, 47
    MPI_Status *status) 48
int MPI_File_write_shared_end(MPI_File fh, void *buf, MPI_Status *status) 49
```

```

1  int MPI_File_write_all_end(MPI_File fh, void *buf, MPI_Status *status)
2
3  int MPI_File_write_at(MPI_File fh, MPI_Offset offset, void *buf, int count,
4                        MPI_Datatype datatype, MPI_Status *status)
5
6  int MPI_File_write_at_all(MPI_File fh, MPI_Offset offset, void *buf,
7                            int count, MPI_Datatype datatype, MPI_Status *status)
8
9  int MPI_File_write_at_all_begin(MPI_File fh, MPI_Offset offset, void *buf,
10                                 int count, MPI_Datatype datatype)
11
12 int MPI_File_write_at_all_end(MPI_File fh, void *buf, MPI_Status *status)
13
14 int MPI_File_write_ordered(MPI_File fh, void *buf, int count,
15                            MPI_Datatype datatype, MPI_Status *status)
16
17 int MPI_File_write_ordered_begin(MPI_File fh, void *buf, int count,
18                                 MPI_Datatype datatype)
19
20 int MPI_File_write_ordered_end(MPI_File fh, void *buf, MPI_Status *status)
21
22 int MPI_File_write_shared(MPI_File fh, void *buf, int count,
23                           MPI_Datatype datatype, MPI_Status *status)
24
25 int MPI_Register_datarep(char *datarep,
26                          MPI_Datarep_conversion_function *read_conversion_fn,
27                          MPI_Datarep_conversion_function *write_conversion_fn,
28                          MPI_Datarep_extent_function *dtype_file_extent_fn,
29                          void *extra_state)

```

A.2.11 Language Bindings C Bindings

```

29 MPI_Fint MPI_Comm_c2f(MPI_Comm comm)
30
31 MPI_Comm MPI_Comm_f2c(MPI_Fint comm)
32
33 MPI_Fint MPI_Errhandler_c2f(MPI_Errhandler errhandler)
34
35 MPI_Errhandler MPI_Errhandler_f2c(MPI_Fint errhandler)
36
37 MPI_Fint MPI_File_c2f(MPI_File file)
38
39 MPI_File MPI_File_f2c(MPI_Fint file)
40
41 MPI_Fint MPI_Group_c2f(MPI_Group group)
42
43 MPI_Group MPI_Group_f2c(MPI_Fint group)
44
45 MPI_Fint MPI_Info_c2f(MPI_Info info)
46
47 MPI_Info MPI_Info_f2c(MPI_Fint info)
48
49 MPI_Fint MPI_Op_c2f(MPI_Op op)
50
51 MPI_Op MPI_Op_f2c(MPI_Fint op)
52
53 MPI_Fint MPI_Request_c2f(MPI_Request request)

```



```

MPI_Request MPI_Request_f2c(MPI_Fint request) 1
int MPI_Status_c2f(MPI_Status *c_status, MPI_Fint *f_status) 2
int MPI_Status_f2c(MPI_Fint *f_status, MPI_Status *c_status) 3
MPI_Fint MPI_Type_c2f(MPI_Datatype datatype) 4
int MPI_Type_create_f90_complex(int p, int r, MPI_Datatype *newtype) 5
int MPI_Type_create_f90_integer(int r, MPI_Datatype *newtype) 6
int MPI_Type_create_f90_real(int p, int r, MPI_Datatype *newtype) 7
MPI_Datatype MPI_Type_f2c(MPI_Fint datatype) 8
int MPI_Type_match_size(int typeclass, int size, MPI_Datatype *type) 9
MPI_Fint MPI_Win_c2f(MPI_Win win) 10
MPI_Win MPI_Win_f2c(MPI_Fint win) 11

```

A.2.12 Profiling Interface C Bindings

```

int MPI_Pcontrol(const int level, ...) 12

```

A.2.13 Deprecated C Bindings

```

int MPI_Address(void* location, MPI_Aint *address) 13
int MPI_Attr_delete(MPI_Comm comm, int keyval) 14
int MPI_Attr_get(MPI_Comm comm, int keyval, void *attribute_val, int *flag) 15
int MPI_Attr_put(MPI_Comm comm, int keyval, void* attribute_val) 16
int MPI_Errhandler_create(MPI_Handler_function *function,
    MPI_Errhandler *errhandler) 17
int MPI_Errhandler_get(MPI_Comm comm, MPI_Errhandler *errhandler) 18
int MPI_Errhandler_set(MPI_Comm comm, MPI_Errhandler errhandler) 19
int MPI_Keyval_create(MPI_Copy_function *copy_fn, MPI_Delete_function
    *delete_fn, int *keyval, void* extra_state) 20
int MPI_Keyval_free(int *keyval) 21
int MPI_Type_extent(MPI_Datatype datatype, MPI_Aint *extent) 22
int MPI_Type_hindexed(int count, int *array_of_blocklengths,
    MPI_Aint *array_of_displacements, MPI_Datatype oldtype,
    MPI_Datatype *newtype) 23
int MPI_Type_hvector(int count, int blocklength, MPI_Aint stride,
    MPI_Datatype oldtype, MPI_Datatype *newtype) 24

```

```
1 int MPI_Type_lb(MPI_Datatype datatype, MPI_Aint* displacement)
2
3 int MPI_Type_struct(int count, int *array_of_blocklengths,
4                     MPI_Aint *array_of_displacements, MPI_Datatype *array_of_types,
5                     MPI_Datatype *newtype)
6
7 int MPI_Type_ub(MPI_Datatype datatype, MPI_Aint* displacement)
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
```

A.3 Fortran Bindings

A.3.1 Point-to-Point Communication Fortran Bindings

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

MPI_BSEND_INIT(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
  <type> BUF(*)
  INTEGER REQUEST, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR

MPI_BUFFER_ATTACH( BUFFER, SIZE, IERROR)
  <type> BUFFER(*)
  INTEGER SIZE, IERROR

MPI_BUFFER_DETACH( BUFFER_ADDR, SIZE, IERROR)
  <type> BUFFER_ADDR(*)
  INTEGER SIZE, IERROR

MPI_CANCEL(REQUEST, IERROR)
  INTEGER REQUEST, IERROR

MPI_GET_ADDRESS(LOCATION, ADDRESS, IERROR)
  <type> LOCATION(*)
  INTEGER IERROR
  INTEGER(KIND=MPI_ADDRESS_KIND) ADDRESS

MPI_GET_COUNT(STATUS, DATATYPE, COUNT, IERROR)
  INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, COUNT, IERROR

MPI_GET_ELEMENTS(STATUS, DATATYPE, COUNT, IERROR)
  INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, COUNT, IERROR

MPI_IBSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
  <type> BUF(*)
  INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR

MPI_IProbe(SOURCE, TAG, COMM, FLAG, STATUS, IERROR)
  LOGICAL FLAG
  INTEGER SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR

MPI_IRECV(BUF, COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR)
  <type> BUF(*)
  INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR

MPI_IRSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
  <type> BUF(*)
  INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR

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

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

```
1 MPI_ISSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
2   <type> BUF(*)
3   INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR
4
5 MPI_PACK(INBUF, INCOUNT, DATATYPE, OUTBUF, OUTSIZE, POSITION, COMM, IERROR)
6   <type> INBUF(*), OUTBUF(*)
7   INTEGER INCOUNT, DATATYPE, OUTSIZE, POSITION, COMM, IERROR
8
9 MPI_PACK_EXTERNAL(DATAREP, INBUF, INCOUNT, DATATYPE, OUTBUF, OUTSIZE,
10   POSITION, IERROR)
11   INTEGER INCOUNT, DATATYPE, IERROR
12   INTEGER(KIND=MPI_ADDRESS_KIND) OUTSIZE, POSITION
13   CHARACTER*(*) DATAREP
14   <type> INBUF(*), OUTBUF(*)
15
16 MPI_PACK_EXTERNAL_SIZE(DATAREP, INCOUNT, DATATYPE, SIZE, IERROR)
17   INTEGER INCOUNT, DATATYPE, IERROR
18   INTEGER(KIND=MPI_ADDRESS_KIND) SIZE
19   CHARACTER*(*) DATAREP
20
21 MPI_PACK_SIZE(INCOUNT, DATATYPE, COMM, SIZE, IERROR)
22   INTEGER INCOUNT, DATATYPE, COMM, SIZE, IERROR
23
24 MPI_PROBE(SOURCE, TAG, COMM, STATUS, IERROR)
25   INTEGER SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR
26
27 MPI_RECV(BUF, COUNT, DATATYPE, SOURCE, TAG, COMM, STATUS, IERROR)
28   <type> BUF(*)
29   INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE),
30   IERROR
31
32 MPI_RECV_INIT(BUF, COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR)
33   <type> BUF(*)
34   INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM, REQUEST, IERROR
35
36 MPI_REQUEST_FREE(REQUEST, IERROR)
37   INTEGER REQUEST, IERROR
38
39 MPI_REQUEST_GET_STATUS(REQUEST, FLAG, STATUS, IERROR)
40   INTEGER REQUEST, STATUS(MPI_STATUS_SIZE), IERROR
41   LOGICAL FLAG
42
43 MPI_RSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR)
44   <type> BUF(*)
45   INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR
46
47 MPI_RSEND_INIT(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR)
48   <type> BUF(*)
49   INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR
50
51 MPI_SEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR)
52   <type> BUF(*)
53   INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR
```

```
MPI_SENDRECV(SENDBUF, SENDCOUNT, SENDTYPE, DEST, SENDTAG, RECVBUFF, 1
              RECVCOUNT, RECVMYPE, SOURCE, RECVTAG, COMM, STATUS, IERROR) 2
    <type> SENDBUF(*), RECVBUFF(*) 3
    INTEGER SENDCOUNT, SENDTYPE, DEST, SENDTAG, RECVCOUNT, RECVMYPE, 4
    SOURCE, RECVTAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR 5
    6
MPI_SENDRECV_REPLACE(BUF, COUNT, DATATYPE, DEST, SENDTAG, SOURCE, RECVTAG, 7
                    COMM, STATUS, IERROR) 8
    <type> BUF(*) 9
    INTEGER COUNT, DATATYPE, DEST, SENDTAG, SOURCE, RECVTAG, COMM, 10
    STATUS(MPI_STATUS_SIZE), IERROR 11
    12
MPI_SEND_INIT(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR) 13
    <type> BUF(*) 14
    INTEGER REQUEST, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR 15
    16
MPI_SSEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, IERROR) 17
    <type> BUF(*) 18
    INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR 19
    20
MPI_SSEND_INIT(BUF, COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR) 21
    <type> BUF(*) 22
    INTEGER COUNT, DATATYPE, DEST, TAG, COMM, REQUEST, IERROR 23
    24
MPI_START(REQUEST, IERROR) 25
    INTEGER REQUEST, IERROR 26
    27
MPI_STARTALL(COUNT, ARRAY_OF_REQUESTS, IERROR) 28
    INTEGER COUNT, ARRAY_OF_REQUESTS(*), IERROR 29
    30
MPI_TEST(REQUEST, FLAG, STATUS, IERROR) 31
    LOGICAL FLAG 32
    INTEGER REQUEST, STATUS(MPI_STATUS_SIZE), IERROR 33
    34
MPI_TESTALL(COUNT, ARRAY_OF_REQUESTS, FLAG, ARRAY_OF_STATUSES, IERROR) 35
    LOGICAL FLAG 36
    INTEGER COUNT, ARRAY_OF_REQUESTS(*), 37
    ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR 38
    39
MPI_TESTANY(COUNT, ARRAY_OF_REQUESTS, INDEX, FLAG, STATUS, IERROR) 40
    LOGICAL FLAG 41
    INTEGER COUNT, ARRAY_OF_REQUESTS(*), INDEX, STATUS(MPI_STATUS_SIZE), 42
    IERROR 43
    44
MPI_TESTSOME(INCOUNT, ARRAY_OF_REQUESTS, OUTCOUNT, ARRAY_OF_INDICES, 45
             ARRAY_OF_STATUSES, IERROR) 46
    INTEGER INCOUNT, ARRAY_OF_REQUESTS(*), OUTCOUNT, ARRAY_OF_INDICES(*), 47
    ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR 48
    49
MPI_TEST_CANCELLED(STATUS, FLAG, IERROR) 50
    LOGICAL FLAG 51
    INTEGER STATUS(MPI_STATUS_SIZE), IERROR 52
    53
```

```
1 MPI_TYPE_COMMIT(DATATYPE, IERROR)
2     INTEGER DATATYPE, IERROR
3
4 MPI_TYPE_CONTIGUOUS(COUNT, OLDTYPE, NEWTYPE, IERROR)
5     INTEGER COUNT, OLDTYPE, NEWTYPE, IERROR
6
7 MPI_TYPE_CREATE_DARRAY(SIZE, RANK, NDIMS, ARRAY_OF_GSIZES, ARRAY_OF_DISTRIBS,
8     ARRAY_OF_DARGS, ARRAY_OF_PSIZEs, ORDER, OLDTYPE, NEWTYPE,
9     IERROR)
10    INTEGER SIZE, RANK, NDIMS, ARRAY_OF_GSIZES(*), ARRAY_OF_DISTRIBS(*),
11    ARRAY_OF_DARGS(*), ARRAY_OF_PSIZEs(*), ORDER, OLDTYPE, NEWTYPE, IERROR
12
13 MPI_TYPE_CREATE_HINDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS,
14     ARRAY_OF_DISPLACEMENTS, OLDTYPE, NEWTYPE, IERROR)
15    INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), OLDTYPE, NEWTYPE, IERROR
16    INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_DISPLACEMENTS(*)
17
18 MPI_TYPE_CREATE_HVECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE,
19     IERROR)
20    INTEGER COUNT, BLOCKLENGTH, OLDTYPE, NEWTYPE, IERROR
21    INTEGER(KIND=MPI_ADDRESS_KIND) STRIDE
22
23 MPI_TYPE_CREATE_INDEXED_BLOCK(COUNT, BLOCKLENGTH, ARRAY_OF_DISPLACEMENTS,
24     OLDTYPE, NEWTYPE, IERROR)
25    INTEGER COUNT, BLOCKLENGTH, ARRAY_OF_DISPLACEMENTS(*), OLDTYPE,
26    NEWTYPE, IERROR
27
28 MPI_TYPE_CREATE_RESIZED(OLDTYPE, LB, EXTENT, NEWTYPE, IERROR)
29    INTEGER OLDTYPE, NEWTYPE, IERROR
30    INTEGER(KIND=MPI_ADDRESS_KIND) LB, EXTENT
31
32 MPI_TYPE_CREATE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
33     ARRAY_OF_TYPES, NEWTYPE, IERROR)
34    INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_TYPES(*), NEWTYPE,
35    IERROR
36    INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_DISPLACEMENTS(*)
37
38 MPI_TYPE_CREATE_SUBARRAY(NDIMS, ARRAY_OF_SIZES, ARRAY_OF_SUBSIZES,
39     ARRAY_OF_STARTS, ORDER, OLDTYPE, NEWTYPE, IERROR)
40    INTEGER NDIMS, ARRAY_OF_SIZES(*), ARRAY_OF_SUBSIZES(*),
41    ARRAY_OF_STARTS(*), ORDER, OLDTYPE, NEWTYPE, IERROR
42
43 MPI_TYPE_DUP(TYPE, NEWTYPE, IERROR)
44    INTEGER TYPE, NEWTYPE, IERROR
45
46 MPI_TYPE_FREE(DATATYPE, IERROR)
47    INTEGER DATATYPE, IERROR
48
49 MPI_TYPE_GET_EXTENT(DATATYPE, LB, EXTENT, IERROR)
50    INTEGER DATATYPE, IERROR
51    INTEGER(KIND = MPI_ADDRESS_KIND) LB, EXTENT
52
53 MPI_TYPE_GET_TRUE_EXTENT(DATATYPE, TRUE_LB, TRUE_EXTENT, IERROR)
```

```

INTEGER DATATYPE, IERROR 1
INTEGER(KIND = MPI_ADDRESS_KIND) TRUE_LB, TRUE_EXTENT 2
3
MPI_TYPE_INDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS, 4
                OLDTYPE, NEWTYPE, IERROR) 5
INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*), 6
                OLDTYPE, NEWTYPE, IERROR 7
MPI_TYPE_SIZE(DATATYPE, SIZE, IERROR) 8
INTEGER DATATYPE, SIZE, IERROR 9
10
MPI_TYPE_VECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR) 11
INTEGER COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR 12
MPI_UNPACK(INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT, DATATYPE, COMM, 13
           IERROR) 14
<type> INBUF(*), OUTBUF(*) 15
INTEGER INSIZE, POSITION, OUTCOUNT, DATATYPE, COMM, IERROR 16
17
MPI_UNPACK_EXTERNAL(DATAREP, INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT, 18
                   DATATYPE, IERROR) 19
INTEGER OUTCOUNT, DATATYPE, IERROR 20
INTEGER(KIND=MPI_ADDRESS_KIND) INSIZE, POSITION 21
CHARACTER*(*) DATAREP 22
<type> INBUF(*), OUTBUF(*) 23
24
MPI_WAIT(REQUEST, STATUS, IERROR) 25
INTEGER REQUEST, STATUS(MPI_STATUS_SIZE), IERROR 26
MPI_WAITALL(COUNT, ARRAY_OF_REQUESTS, ARRAY_OF_STATUSES, IERROR) 27
INTEGER COUNT, ARRAY_OF_REQUESTS(*) 28
INTEGER ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR 29
MPI_WAITANY(COUNT, ARRAY_OF_REQUESTS, INDEX, STATUS, IERROR) 30
INTEGER COUNT, ARRAY_OF_REQUESTS(*), INDEX, STATUS(MPI_STATUS_SIZE), 31
IERROR 32
33
MPI_WAITSOME(INCOUNT, ARRAY_OF_REQUESTS, OUTCOUNT, ARRAY_OF_INDICES, 34
            ARRAY_OF_STATUSES, IERROR) 35
INTEGER INCOUNT, ARRAY_OF_REQUESTS(*), OUTCOUNT, ARRAY_OF_INDICES(*), 36
ARRAY_OF_STATUSES(MPI_STATUS_SIZE,*), IERROR 37
38

```

A.3.2 Collective Communication Fortran Bindings

```

MPI_ALLGATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, 41
             COMM, IERROR) 42
<type> SENDBUF(*), RECVCOUNT(*) 43
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, COMM, IERROR 44
45
MPI_ALLGATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, DISPLS, 46
              RECVTYPE, COMM, IERROR) 47
<type> SENDBUF(*), RECVCOUNT(*) 48

```

```

1     INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVMODE, COMM,
2     IERROR
3
4     MPI_ALLREDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, IERROR)
5     <type> SENDBUF(*), RECVBUF(*)
6     INTEGER COUNT, DATATYPE, OP, COMM, IERROR
7
8     MPI_ALLTOALL(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVMODE,
9     COMM, IERROR)
10    <type> SENDBUF(*), RECVBUF(*)
11    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVMODE, COMM, IERROR
12
13    MPI_ALLTOALLV(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPE, RECVBUF, RECVCOUNTS,
14    RDISPLS, RECVMODE, COMM, IERROR)
15    <type> SENDBUF(*), RECVBUF(*)
16    INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPE, RECVCOUNTS(*), RDISPLS(*),
17    RECVMODE, COMM, IERROR
18
19    MPI_ALLTOALLW(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPES, RECVBUF, RECVCOUNTS,
20    RDISPLS, RECVMODES, COMM, IERROR)
21    <type> SENDBUF(*), RECVBUF(*)
22    INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPES(*), RECVCOUNTS(*),
23    RDISPLS(*), RECVMODES(*), COMM, IERROR
24
25    MPI_BARRIER(COMM, IERROR)
26    INTEGER COMM, IERROR
27
28    MPI_BCAST(BUFFER, COUNT, DATATYPE, ROOT, COMM, IERROR)
29    <type> BUFFER(*)
30    INTEGER COUNT, DATATYPE, ROOT, COMM, IERROR
31
32    MPI_EXSCAN(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, IERROR)
33    <type> SENDBUF(*), RECVBUF(*)
34    INTEGER COUNT, DATATYPE, OP, COMM, IERROR
35
36    MPI_GATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVMODE,
37    ROOT, COMM, IERROR)
38    <type> SENDBUF(*), RECVBUF(*)
39    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVMODE, ROOT, COMM, IERROR
40
41    MPI_GATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
42    RECVMODE, ROOT, COMM, IERROR)
43    <type> SENDBUF(*), RECVBUF(*)
44    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVMODE, ROOT,
45    COMM, IERROR
46
47    MPI_OP_CREATE( FUNCTION, COMMUTE, OP, IERROR)
48    EXTERNAL FUNCTION
49    LOGICAL COMMUTE
50    INTEGER OP, IERROR
51
52    MPI_OP_FREE( OP, IERROR)
53    INTEGER OP, IERROR

```



```
1     LOGICAL FLAG
2
3     MPI_COMM_GROUP(COMM, GROUP, IERROR)
4         INTEGER COMM, GROUP, IERROR
5
6     MPI_COMM_RANK(COMM, RANK, IERROR)
7         INTEGER COMM, RANK, IERROR
8
9     MPI_COMM_REMOTE_GROUP(COMM, GROUP, IERROR)
10        INTEGER COMM, GROUP, IERROR
11
12    MPI_COMM_REMOTE_SIZE(COMM, SIZE, IERROR)
13        INTEGER COMM, SIZE, IERROR
14
15    MPI_COMM_SET_ATTR(COMM, COMM_KEYVAL, ATTRIBUTE_VAL, IERROR)
16        INTEGER COMM, COMM_KEYVAL, IERROR
17        INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
18
19    MPI_COMM_SIZE(COMM, SIZE, IERROR)
20        INTEGER COMM, SIZE, IERROR
21
22    MPI_COMM_SPLIT(COMM, COLOR, KEY, NEWCOMM, IERROR)
23        INTEGER COMM, COLOR, KEY, NEWCOMM, IERROR
24
25    MPI_COMM_TEST_INTER(COMM, FLAG, IERROR)
26        INTEGER COMM, IERROR
27        LOGICAL FLAG
28
29    MPI_GROUP_COMPARE(GROUP1, GROUP2, RESULT, IERROR)
30        INTEGER GROUP1, GROUP2, RESULT, IERROR
31
32    MPI_GROUP_DIFFERENCE(GROUP1, GROUP2, NEWGROUP, IERROR)
33        INTEGER GROUP1, GROUP2, NEWGROUP, IERROR
34
35    MPI_GROUP_EXCL(GROUP, N, RANKS, NEWGROUP, IERROR)
36        INTEGER GROUP, N, RANKS(*), NEWGROUP, IERROR
37
38    MPI_GROUP_FREE(GROUP, IERROR)
39        INTEGER GROUP, IERROR
40
41    MPI_GROUP_INCL(GROUP, N, RANKS, NEWGROUP, IERROR)
42        INTEGER GROUP, N, RANKS(*), NEWGROUP, IERROR
43
44    MPI_GROUP_INTERSECTION(GROUP1, GROUP2, NEWGROUP, IERROR)
45        INTEGER GROUP1, GROUP2, NEWGROUP, IERROR
46
47    MPI_GROUP_RANGE_EXCL(GROUP, N, RANGES, NEWGROUP, IERROR)
48        INTEGER GROUP, N, RANGES(3,*), NEWGROUP, IERROR
49
50    MPI_GROUP_RANGE_INCL(GROUP, N, RANGES, NEWGROUP, IERROR)
51        INTEGER GROUP, N, RANGES(3,*), NEWGROUP, IERROR
52
53    MPI_GROUP_RANK(GROUP, RANK, IERROR)
54        INTEGER GROUP, RANK, IERROR
55
56    MPI_GROUP_SIZE(GROUP, SIZE, IERROR)
```

```
INTEGER GROUP, SIZE, IERROR 1
MPI_GROUP_TRANSLATE_RANKS(GROUP1, N, RANKS1, GROUP2, RANKS2, IERROR) 2
  INTEGER GROUP1, N, RANKS1(*), GROUP2, RANKS2(*), IERROR 3
MPI_GROUP_UNION(GROUP1, GROUP2, NEWGROUP, IERROR) 4
  INTEGER GROUP1, GROUP2, NEWGROUP, IERROR 5
MPI_INTERCOMM_CREATE(LOCAL_COMM, LOCAL_LEADER, PEER_COMM, REMOTE_LEADER, TAG, 6
  NEWINTERCOMM, IERROR) 7
  INTEGER LOCAL_COMM, LOCAL_LEADER, PEER_COMM, REMOTE_LEADER, TAG, 8
  NEWINTERCOMM, IERROR 9
MPI_INTERCOMM_MERGE(INTERCOMM, HIGH, INTRACOMM, IERROR) 10
  INTEGER INTERCOMM, INTRACOMM, IERROR 11
  LOGICAL HIGH 12
MPI_TYPE_CREATE_KEYVAL(TYPE_COPY_ATTR_FN, TYPE_DELETE_ATTR_FN, TYPE_KEYVAL, 13
  EXTRA_STATE, IERROR) 14
  EXTERNAL TYPE_COPY_ATTR_FN, TYPE_DELETE_ATTR_FN 15
  INTEGER TYPE_KEYVAL, IERROR 16
  INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE 17
MPI_TYPE_DELETE_ATTR(TYPE, TYPE_KEYVAL, IERROR) 18
  INTEGER TYPE, TYPE_KEYVAL, IERROR 19
MPI_TYPE_FREE_KEYVAL(TYPE_KEYVAL, IERROR) 20
  INTEGER TYPE_KEYVAL, IERROR 21
MPI_TYPE_GET_ATTR(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR) 22
  INTEGER TYPE, TYPE_KEYVAL, IERROR 23
  INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 24
  LOGICAL FLAG 25
MPI_TYPE_SET_ATTR(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, IERROR) 26
  INTEGER TYPE, TYPE_KEYVAL, IERROR 27
  INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 28
MPI_WIN_CREATE_KEYVAL(WIN_COPY_ATTR_FN, WIN_DELETE_ATTR_FN, WIN_KEYVAL, 29
  EXTRA_STATE, IERROR) 30
  EXTERNAL WIN_COPY_ATTR_FN, WIN_DELETE_ATTR_FN 31
  INTEGER WIN_KEYVAL, IERROR 32
  INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE 33
MPI_WIN_DELETE_ATTR(WIN, WIN_KEYVAL, IERROR) 34
  INTEGER WIN, WIN_KEYVAL, IERROR 35
MPI_WIN_FREE_KEYVAL(WIN_KEYVAL, IERROR) 36
  INTEGER WIN_KEYVAL, IERROR 37
MPI_WIN_GET_ATTR(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR) 38
  INTEGER WIN, WIN_KEYVAL, IERROR 39
  INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 40
  LOGICAL FLAG 41
```

```
1 MPI_WIN_SET_ATTR(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, IERROR)
2     INTEGER WIN, WIN_KEYVAL, IERROR
3     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
4
5
```

6 A.3.4 Process Topologies Fortran Bindings

```
7 MPI_CARTDIM_GET(COMM, NDIMS, IERROR)
8     INTEGER COMM, NDIMS, IERROR
9
10 MPI_CART_COORDS(COMM, RANK, MAXDIMS, COORDS, IERROR)
11     INTEGER COMM, RANK, MAXDIMS, COORDS(*), IERROR
12
13 MPI_CART_CREATE(COMM_OLD, NDIMS, DIMS, PERIODS, REORDER, COMM_CART, IERROR)
14     INTEGER COMM_OLD, NDIMS, DIMS(*), COMM_CART, IERROR
15     LOGICAL PERIODS(*), REORDER
16
17 MPI_CART_GET(COMM, MAXDIMS, DIMS, PERIODS, COORDS, IERROR)
18     INTEGER COMM, MAXDIMS, DIMS(*), COORDS(*), IERROR
19     LOGICAL PERIODS(*)
20
21 MPI_CART_MAP(COMM, NDIMS, DIMS, PERIODS, NEWRANK, IERROR)
22     INTEGER COMM, NDIMS, DIMS(*), NEWRANK, IERROR
23     LOGICAL PERIODS(*)
24
25 MPI_CART_RANK(COMM, COORDS, RANK, IERROR)
26     INTEGER COMM, COORDS(*), RANK, IERROR
27
28 MPI_CART_SHIFT(COMM, DIRECTION, DISP, RANK_SOURCE, RANK_DEST, IERROR)
29     INTEGER COMM, DIRECTION, DISP, RANK_SOURCE, RANK_DEST, IERROR
30
31 MPI_CART_SUB(COMM, REMAIN_DIMS, NEWCOMM, IERROR)
32     INTEGER COMM, NEWCOMM, IERROR
33     LOGICAL REMAIN_DIMS(*)
34
35 MPI_DIMS_CREATE(NNODES, NDIMS, DIMS, IERROR)
36     INTEGER NNODES, NDIMS, DIMS(*), IERROR
37
38 MPI_GRAPHDIMS_GET(COMM, NNODES, NEDGES, IERROR)
39     INTEGER COMM, NNODES, NEDGES, IERROR
40
41 MPI_GRAPH_CREATE(COMM_OLD, NNODES, INDEX, EDGES, REORDER, COMM_GRAPH,
42     IERROR)
43     INTEGER COMM_OLD, NNODES, INDEX(*), EDGES(*), COMM_GRAPH, IERROR
44     LOGICAL REORDER
45
46 MPI_GRAPH_GET(COMM, MAXINDEX, MAXEDGES, INDEX, EDGES, IERROR)
47     INTEGER COMM, MAXINDEX, MAXEDGES, INDEX(*), EDGES(*), IERROR
48
49 MPI_GRAPH_MAP(COMM, NNODES, INDEX, EDGES, NEWRANK, IERROR)
50     INTEGER COMM, NNODES, INDEX(*), EDGES(*), NEWRANK, IERROR
51
52 MPI_GRAPH_NEIGHBORS(COMM, RANK, MAXNEIGHBORS, NEIGHBORS, IERROR)
53     INTEGER COMM, RANK, MAXNEIGHBORS, NEIGHBORS(*), IERROR
54
```

| | |
|---|----|
| MPI_GRAPH_NEIGHBORS_COUNT(COMM, RANK, NNEIGHBORS, IERROR) | 1 |
| INTEGER COMM, RANK, NNEIGHBORS, IERROR | 2 |
| | 3 |
| MPI_TOPO_TEST(COMM, STATUS, IERROR) | 4 |
| INTEGER COMM, STATUS, IERROR | 5 |
| | 6 |
| A.3.5 MPI Environmenta Management Fortran Bindings | 7 |
| | 8 |
| DOUBLE PRECISION MPI_WTICK() | 9 |
| | 10 |
| DOUBLE PRECISION MPI_WTIME() | 11 |
| | 12 |
| MPI_ABORT(COMM, ERRORCODE, IERROR) | 13 |
| INTEGER COMM, ERRORCODE, IERROR | 14 |
| | 15 |
| MPI_ALLOC_MEM(SIZE, INFO, BASEPTR, IERROR) | 16 |
| INTEGER INFO, IERROR | 17 |
| INTEGER(KIND=MPI_ADDRESS_KIND) SIZE, BASEPTR | 18 |
| | 19 |
| MPI_COMM_CREATE_ERRHANDLER(FUNCTION, ERRHANDLER, IERROR) | 20 |
| EXTERNAL FUNCTION | 21 |
| INTEGER ERRHANDLER, IERROR | 22 |
| | 23 |
| MPI_COMM_GET_ERRHANDLER(COMM, ERRHANDLER, IERROR) | 24 |
| INTEGER COMM, ERRHANDLER, IERROR | 25 |
| | 26 |
| MPI_COMM_SET_ERRHANDLER(COMM, ERRHANDLER, IERROR) | 27 |
| INTEGER COMM, ERRHANDLER, IERROR | 28 |
| | 29 |
| MPI_ERRHANDLER_FREE(ERRHANDLER, IERROR) | 30 |
| INTEGER ERRHANDLER, IERROR | 31 |
| | 32 |
| MPI_ERROR_CLASS(ERRORCODE, ERRORCLASS, IERROR) | 33 |
| INTEGER ERRORCODE, ERRORCLASS, IERROR | 34 |
| | 35 |
| MPI_ERROR_STRING(ERRORCODE, STRING, RESULTLEN, IERROR) | 36 |
| INTEGER ERRORCODE, RESULTLEN, IERROR | 37 |
| CHARACTER*(*) STRING | 38 |
| | 39 |
| MPI_FILE_CREATE_ERRHANDLER(FUNCTION, ERRHANDLER, IERROR) | 40 |
| EXTERNAL FUNCTION | 41 |
| INTEGER ERRHANDLER, IERROR | 42 |
| | 43 |
| MPI_FILE_GET_ERRHANDLER(FILE, ERRHANDLER, IERROR) | 44 |
| INTEGER FILE, ERRHANDLER, IERROR | 45 |
| | 46 |
| MPI_FILE_SET_ERRHANDLER(FILE, ERRHANDLER, IERROR) | 47 |
| INTEGER FILE, ERRHANDLER, IERROR | 48 |
| | |
| MPI_FINALIZE(IERROR) | |
| INTEGER IERROR | |
| | |
| MPI_FINALIZED(FLAG, IERROR) | |
| LOGICAL FLAG | |

```
1      INTEGER IERROR
2
3      MPI_FREE_MEM(BASE, IERROR)
4          <type> BASE(*)
5          INTEGER IERROR
6
7      MPI_GET_PROCESSOR_NAME( NAME, RESULTLEN, IERROR)
8          CHARACTER*(*) NAME
9          INTEGER RESULTLEN, IERROR
10
11     MPI_GET_VERSION(VERSION, SUBVERSION, IERROR)
12         INTEGER VERSION, SUBVERSION, IERROR
13
14     MPI_INIT(IERROR)
15         INTEGER IERROR
16
17     MPI_INITIALIZED(FLAG, IERROR)
18         LOGICAL FLAG
19         INTEGER IERROR
20
21     MPI_WIN_CREATE_ERRHANDLER(FUNCTION, ERRHANDLER, IERROR)
22         EXTERNAL FUNCTION
23         INTEGER ERRHANDLER, IERROR
24
25     MPI_WIN_GET_ERRHANDLER(WIN, ERRHANDLER, IERROR)
26         INTEGER WIN, ERRHANDLER, IERROR
27
28     MPI_WIN_SET_ERRHANDLER(WIN, ERRHANDLER, IERROR)
29         INTEGER WIN, ERRHANDLER, IERROR
```

A.3.6 Miscellany Fortran Bindings

```
29     MPI_INFO_CREATE(INFO, IERROR)
30         INTEGER INFO, IERROR
31
32     MPI_INFO_DELETE(INFO, KEY, IERROR)
33         INTEGER INFO, IERROR
34         CHARACTER*(*) KEY
35
36     MPI_INFO_DUP(INFO, NEWINFO, IERROR)
37         INTEGER INFO, NEWINFO, IERROR
38
39     MPI_INFO_FREE(INFO, IERROR)
40         INTEGER INFO, IERROR
41
42     MPI_INFO_GET(INFO, KEY, VALUELEN, VALUE, FLAG, IERROR)
43         INTEGER INFO, VALUELEN, IERROR
44         CHARACTER*(*) KEY, VALUE
45         LOGICAL FLAG
46
47     MPI_INFO_GET_NKEYS(INFO, NKEYS, IERROR)
48         INTEGER INFO, NKEYS, IERROR
49
50     MPI_INFO_GET_NTHKEY(INFO, N, KEY, IERROR)
```

| | |
|--|----|
| INTEGER INFO, N, IERROR | 1 |
| CHARACTER*(*) KEY | 2 |
| MPI_INFO_GET_VALUELEN(INFO, KEY, VALUELEN, FLAG, IERROR) | 3 |
| INTEGER INFO, VALUELEN, IERROR | 4 |
| LOGICAL FLAG | 5 |
| CHARACTER*(*) KEY | 6 |
| MPI_INFO_SET(INFO, KEY, VALUE, IERROR) | 7 |
| INTEGER INFO, IERROR | 8 |
| CHARACTER*(*) KEY, VALUE | 9 |
| | 10 |
| | 11 |
| | 12 |
| A.3.7 Process Creation and Management Fortran Bindings | 13 |
| MPI_CLOSE_PORT(PORT_NAME, IERROR) | 14 |
| CHARACTER*(*) PORT_NAME | 15 |
| INTEGER IERROR | 16 |
| | 17 |
| MPI_COMM_ACCEPT(PORT_NAME, INFO, ROOT, COMM, NEWCOMM, IERROR) | 18 |
| CHARACTER*(*) PORT_NAME | 19 |
| INTEGER INFO, ROOT, COMM, NEWCOMM, IERROR | 20 |
| | 21 |
| MPI_COMM_CONNECT(PORT_NAME, INFO, ROOT, COMM, NEWCOMM, IERROR) | 22 |
| CHARACTER*(*) PORT_NAME | 23 |
| INTEGER INFO, ROOT, COMM, NEWCOMM, IERROR | 24 |
| | 25 |
| MPI_COMM_DISCONNECT(COMM, IERROR) | 26 |
| INTEGER COMM, IERROR | 27 |
| | 28 |
| MPI_COMM_GET_PARENT(PARENT, IERROR) | 29 |
| INTEGER PARENT, IERROR | 30 |
| | 31 |
| MPI_COMM_JOIN(FD, INTERCOMM, IERROR) | 32 |
| INTEGER FD, INTERCOMM, IERROR | 33 |
| | 34 |
| MPI_COMM_SPAWN(COMMAND, ARGV, MAXPROCS, INFO, ROOT, COMM, INTERCOMM, ARRAY_OF_ERRCODES, IERROR) | 35 |
| CHARACTER*(*) COMMAND, ARGV(*) | 36 |
| INTEGER INFO, MAXPROCS, ROOT, COMM, INTERCOMM, ARRAY_OF_ERRCODES(*), IERROR | 37 |
| | 38 |
| MPI_COMM_SPAWN_MULTIPLE(COUNT, ARRAY_OF_COMMANDS, ARRAY_OF_ARGV, ARRAY_OF_MAXPROCS, ARRAY_OF_INFO, ROOT, COMM, INTERCOMM, ARRAY_OF_ERRCODES, IERROR) | 39 |
| INTEGER COUNT, ARRAY_OF_INFO(*), ARRAY_OF_MAXPROCS(*), ROOT, COMM, INTERCOMM, ARRAY_OF_ERRCODES(*), IERROR | 40 |
| CHARACTER*(*) ARRAY_OF_COMMANDS(*), ARRAY_OF_ARGV(COUNT, *) | 41 |
| | 42 |
| MPI_LOOKUP_NAME(SERVICE_NAME, INFO, PORT_NAME, IERROR) | 43 |
| CHARACTER*(*) SERVICE_NAME, PORT_NAME | 44 |
| INTEGER INFO, IERROR | 45 |
| | 46 |
| MPI_OPEN_PORT(INFO, PORT_NAME, IERROR) | 47 |
| | 48 |

```

1     CHARACTER*(*) PORT_NAME
2     INTEGER INFO, IERROR
3
4     MPI_PUBLISH_NAME(SERVICE_NAME, INFO, PORT_NAME, IERROR)
5     INTEGER INFO, IERROR
6     CHARACTER*(*) SERVICE_NAME, PORT_NAME
7
8     MPI_UNPUBLISH_NAME(SERVICE_NAME, INFO, PORT_NAME, IERROR)
9     INTEGER INFO, IERROR
10    CHARACTER*(*) SERVICE_NAME, PORT_NAME

```

11

12 A.3.8 One-Sided Communications Fortran Bindings

```

13    MPI_ACCUMULATE(ORIGIN_ADDR, ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK,
14                  TARGET_DISP, TARGET_COUNT, TARGET_DATATYPE, OP, WIN, IERROR)
15    <type> ORIGIN_ADDR(*)
16    INTEGER(KIND=MPI_ADDRESS_KIND) TARGET_DISP
17    INTEGER ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_COUNT,
18    TARGET_DATATYPE, OP, WIN, IERROR
19
20    MPI_GET(ORIGIN_ADDR, ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_DISP,
21           TARGET_COUNT, TARGET_DATATYPE, WIN, IERROR)
22    <type> ORIGIN_ADDR(*)
23    INTEGER(KIND=MPI_ADDRESS_KIND) TARGET_DISP
24    INTEGER ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_COUNT,
25    TARGET_DATATYPE, WIN, IERROR
26
27    MPI_PUT(ORIGIN_ADDR, ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_DISP,
28           TARGET_COUNT, TARGET_DATATYPE, WIN, IERROR)
29    <type> ORIGIN_ADDR(*)
30    INTEGER(KIND=MPI_ADDRESS_KIND) TARGET_DISP
31    INTEGER ORIGIN_COUNT, ORIGIN_DATATYPE, TARGET_RANK, TARGET_COUNT,
32    TARGET_DATATYPE, WIN, IERROR
33
34    MPI_WIN_COMPLETE(WIN, IERROR)
35    INTEGER WIN, IERROR
36
37    MPI_WIN_CREATE(BASE, SIZE, DISP_UNIT, INFO, COMM, WIN, IERROR)
38    <type> BASE(*)
39    INTEGER(KIND=MPI_ADDRESS_KIND) SIZE
40    INTEGER DISP_UNIT, INFO, COMM, WIN, IERROR
41
42    MPI_WIN_FENCE(ASSERT, WIN, IERROR)
43    INTEGER ASSERT, WIN, IERROR
44
45    MPI_WIN_FREE(WIN, IERROR)
46    INTEGER WIN, IERROR
47
48    MPI_WIN_GET_GROUP(WIN, GROUP, IERROR)
49    INTEGER WIN, GROUP, IERROR
50
51    MPI_WIN_LOCK(LOCK_TYPE, RANK, ASSERT, WIN, IERROR)

```


| | |
|--|----|
| INTEGER LOCK_TYPE, RANK, ASSERT, WIN, IERROR | 1 |
| MPI_WIN_POST(GROUP, ASSERT, WIN, IERROR) | 2 |
| INTEGER GROUP, ASSERT, WIN, IERROR | 3 |
| MPI_WIN_START(GROUP, ASSERT, WIN, IERROR) | 4 |
| INTEGER GROUP, ASSERT, WIN, IERROR | 5 |
| MPI_WIN_TEST(WIN, FLAG, IERROR) | 6 |
| INTEGER WIN, IERROR | 7 |
| LOGICAL FLAG | 8 |
| MPI_WIN_UNLOCK(RANK, WIN, IERROR) | 9 |
| INTEGER RANK, WIN, IERROR | 10 |
| MPI_WIN_WAIT(WIN, IERROR) | 11 |
| INTEGER WIN, IERROR | 12 |

A.3.9 External Interfaces Fortran Bindings

| | |
|--|----|
| MPI_ADD_ERROR_CLASS(ERRORCLASS, IERROR) | 13 |
| INTEGER ERRORCLASS, IERROR | 14 |
| MPI_ADD_ERROR_CODE(ERRORCLASS, ERRORCODE, IERROR) | 15 |
| INTEGER ERRORCLASS, ERRORCODE, IERROR | 16 |
| MPI_ADD_ERROR_STRING(ERRORCODE, STRING, IERROR) | 17 |
| INTEGER ERRORCODE, IERROR | 18 |
| CHARACTER*(*) STRING | 19 |
| MPI_COMM_CALL_ERRHANDLER(COMM, ERRORCODE, IERROR) | 20 |
| INTEGER COMM, ERRORCODE, IERROR | 21 |
| MPI_COMM_GET_NAME(COMM, COMM_NAME, RESULTLEN, IERROR) | 22 |
| INTEGER COMM, RESULTLEN, IERROR | 23 |
| CHARACTER*(*) COMM_NAME | 24 |
| MPI_COMM_SET_NAME(COMM, COMM_NAME, IERROR) | 25 |
| INTEGER COMM, IERROR | 26 |
| CHARACTER*(*) COMM_NAME | 27 |
| MPI_FILE_CALL_ERRHANDLER(FH, ERRORCODE, IERROR) | 28 |
| INTEGER FH, ERRORCODE, IERROR | 29 |
| MPI_GREQUEST_COMPLETE(REQUEST, IERROR) | 30 |
| INTEGER REQUEST, IERROR | 31 |
| MPI_GREQUEST_START(QUERY_FN, FREE_FN, CANCEL_FN, EXTRA_STATE, REQUEST, | 32 |
| IERROR) | 33 |
| INTEGER REQUEST, IERROR | 34 |
| EXTERNAL QUERY_FN, FREE_FN, CANCEL_FN | 35 |
| INTEGER (KIND=MPI_ADDRESS_KIND) EXTRA_STATE | 36 |
| MPI_INIT_THREAD(REQUIRED, PROVIDED, IERROR) | 37 |

```

1      INTEGER REQUIRED, PROVIDED, IERROR
2
3      MPI_IS_THREAD_MAIN(FLAG, IERROR)
4          LOGICAL FLAG
5          INTEGER IERROR
6
7      MPI_QUERY_THREAD(PROVIDED, IERROR)
8          INTEGER PROVIDED, IERROR
9
10     MPI_STATUS_SET_CANCELLED(STATUS, FLAG, IERROR)
11         INTEGER STATUS(MPI_STATUS_SIZE), IERROR
12         LOGICAL FLAG
13
14     MPI_STATUS_SET_ELEMENTS(STATUS, DATATYPE, COUNT, IERROR)
15         INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, COUNT, IERROR
16
17     MPI_TYPE_GET_CONTENTS(DATATYPE, MAX_INTEGERS, MAX_ADDRESSES, MAX_DATATYPES,
18         ARRAY_OF_INTEGERS, ARRAY_OF_ADDRESSES, ARRAY_OF_DATATYPES,
19         IERROR)
20         INTEGER DATATYPE, MAX_INTEGERS, MAX_ADDRESSES, MAX_DATATYPES,
21         ARRAY_OF_INTEGERS(*), ARRAY_OF_DATATYPES(*), IERROR
22         INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_ADDRESSES(*)
23
24     MPI_TYPE_GET_ENVELOPE(DATATYPE, NUM_INTEGERS, NUM_ADDRESSES, NUM_DATATYPES,
25         COMBINER, IERROR)
26         INTEGER DATATYPE, NUM_INTEGERS, NUM_ADDRESSES, NUM_DATATYPES, COMBINER,
27         IERROR
28
29     MPI_TYPE_GET_NAME(TYPE, TYPE_NAME, RESULTLEN, IERROR)
30         INTEGER TYPE, RESULTLEN, IERROR
31         CHARACTER*(*) TYPE_NAME
32
33     MPI_TYPE_SET_NAME(TYPE, TYPE_NAME, IERROR)
34         INTEGER TYPE, IERROR
35         CHARACTER*(*) TYPE_NAME
36
37     MPI_WIN_CALL_ERRHANDLER(WIN, ERRORCODE, IERROR)
38         INTEGER WIN, ERRORCODE, IERROR
39
40     MPI_WIN_GET_NAME(WIN, WIN_NAME, RESULTLEN, IERROR)
41         INTEGER WIN, RESULTLEN, IERROR
42         CHARACTER*(*) WIN_NAME
43
44     MPI_WIN_SET_NAME(WIN, WIN_NAME, IERROR)
45         INTEGER WIN, IERROR
46         CHARACTER*(*) WIN_NAME
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

```

A.3.10 I/O Fortran Bindings

```

1      MPI_FILE_CLOSE(FH, IERROR)
2          INTEGER FH, IERROR
3
4      MPI_FILE_DELETE(FILENAME, INFO, IERROR)

```

| | |
|--|----|
| CHARACTER*(*) FILENAME | 1 |
| INTEGER INFO, IERROR | 2 |
| | 3 |
| MPI_FILE_GET_AMODE(FH, AMODE, IERROR) | 4 |
| INTEGER FH, AMODE, IERROR | 5 |
| | 6 |
| MPI_FILE_GET_ATOMICITY(FH, FLAG, IERROR) | 7 |
| INTEGER FH, IERROR | 8 |
| LOGICAL FLAG | 9 |
| | 10 |
| MPI_FILE_GET_BYTE_OFFSET(FH, OFFSET, DISP, IERROR) | 11 |
| INTEGER FH, IERROR | 12 |
| INTEGER(KIND=MPI_OFFSET_KIND) OFFSET, DISP | 13 |
| | 14 |
| MPI_FILE_GET_GROUP(FH, GROUP, IERROR) | 15 |
| INTEGER FH, GROUP, IERROR | 16 |
| | 17 |
| MPI_FILE_GET_INFO(FH, INFO_USED, IERROR) | 18 |
| INTEGER FH, INFO_USED, IERROR | 19 |
| | 20 |
| MPI_FILE_GET_POSITION(FH, OFFSET, IERROR) | 21 |
| INTEGER FH, IERROR | 22 |
| INTEGER(KIND=MPI_OFFSET_KIND) OFFSET | 23 |
| | 24 |
| MPI_FILE_GET_POSITION_SHARED(FH, OFFSET, IERROR) | 25 |
| INTEGER FH, IERROR | 26 |
| INTEGER(KIND=MPI_OFFSET_KIND) OFFSET | 27 |
| | 28 |
| MPI_FILE_GET_SIZE(FH, SIZE, IERROR) | 29 |
| INTEGER FH, IERROR | 30 |
| INTEGER(KIND=MPI_OFFSET_KIND) SIZE | 31 |
| | 32 |
| MPI_FILE_GET_TYPE_EXTENT(FH, DATATYPE, EXTENT, IERROR) | 33 |
| INTEGER FH, DATATYPE, IERROR | 34 |
| INTEGER(KIND=MPI_ADDRESS_KIND) EXTENT | 35 |
| | 36 |
| MPI_FILE_GET_VIEW(FH, DISP, ETYPE, FILETYPE, DATAREP, IERROR) | 37 |
| INTEGER FH, ETYPE, FILETYPE, IERROR | 38 |
| CHARACTER*(*) DATAREP | 39 |
| INTEGER(KIND=MPI_OFFSET_KIND) DISP | 40 |
| | 41 |
| MPI_FILE_IREAD(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR) | 42 |
| <type> BUF(*) | 43 |
| INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR | 44 |
| | 45 |
| MPI_FILE_IREAD_AT(FH, OFFSET, BUF, COUNT, DATATYPE, REQUEST, IERROR) | 46 |
| <type> BUF(*) | 47 |
| INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR | 48 |
| INTEGER(KIND=MPI_OFFSET_KIND) OFFSET | |
| | |
| MPI_FILE_IREAD_SHARED(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR) | |
| <type> BUF(*) | |
| INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR | |

```
1 MPI_FILE_IWRITE(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR)
2   <type> BUF(*)
3   INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
4
5 MPI_FILE_IWRITE_AT(FH, OFFSET, BUF, COUNT, DATATYPE, REQUEST, IERROR)
6   <type> BUF(*)
7   INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
8   INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
9
10 MPI_FILE_IWRITE_SHARED(FH, BUF, COUNT, DATATYPE, REQUEST, IERROR)
11   <type> BUF(*)
12   INTEGER FH, COUNT, DATATYPE, REQUEST, IERROR
13
14 MPI_FILE_OPEN(COMM, FILENAME, AMODE, INFO, FH, IERROR)
15   CHARACTER*(*) FILENAME
16   INTEGER COMM, AMODE, INFO, FH, IERROR
17
18 MPI_FILE_PREALLOCATE(FH, SIZE, IERROR)
19   INTEGER FH, IERROR
20   INTEGER(KIND=MPI_OFFSET_KIND) SIZE
21
22 MPI_FILE_READ(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
23   <type> BUF(*)
24   INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
25
26 MPI_FILE_READ_ALL(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
27   <type> BUF(*)
28   INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
29
30 MPI_FILE_READ_ALL_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR)
31   <type> BUF(*)
32   INTEGER FH, COUNT, DATATYPE, IERROR
33
34 MPI_FILE_READ_ALL_END(FH, BUF, STATUS, IERROR)
35   <type> BUF(*)
36   INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
37
38 MPI_FILE_READ_AT(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR)
39   <type> BUF(*)
40   INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
41   INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
42
43 MPI_FILE_READ_AT_ALL(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR)
44   <type> BUF(*)
45   INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
46   INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
47
48 MPI_FILE_READ_AT_ALL_BEGIN(FH, OFFSET, BUF, COUNT, DATATYPE, IERROR)
49   <type> BUF(*)
50   INTEGER FH, COUNT, DATATYPE, IERROR
51   INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
```

| | |
|---|----|
| <type> BUF(*) | 1 |
| INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR | 2 |
| MPI_FILE_READ_ORDERED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR) | 3 |
| <type> BUF(*) | 4 |
| INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR | 5 |
| MPI_FILE_READ_ORDERED_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR) | 6 |
| <type> BUF(*) | 7 |
| INTEGER FH, COUNT, DATATYPE, IERROR | 8 |
| MPI_FILE_READ_ORDERED_END(FH, BUF, STATUS, IERROR) | 9 |
| <type> BUF(*) | 10 |
| INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR | 11 |
| MPI_FILE_READ_SHARED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR) | 12 |
| <type> BUF(*) | 13 |
| INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR | 14 |
| MPI_FILE_SEEK(FH, OFFSET, WHENCE, IERROR) | 15 |
| INTEGER FH, WHENCE, IERROR | 16 |
| INTEGER(KIND=MPI_OFFSET_KIND) OFFSET | 17 |
| MPI_FILE_SEEK_SHARED(FH, OFFSET, WHENCE, IERROR) | 18 |
| INTEGER FH, WHENCE, IERROR | 19 |
| INTEGER(KIND=MPI_OFFSET_KIND) OFFSET | 20 |
| MPI_FILE_SET_ATOMICITY(FH, FLAG, IERROR) | 21 |
| INTEGER FH, IERROR | 22 |
| LOGICAL FLAG | 23 |
| MPI_FILE_SET_INFO(FH, INFO, IERROR) | 24 |
| INTEGER FH, INFO, IERROR | 25 |
| MPI_FILE_SET_SIZE(FH, SIZE, IERROR) | 26 |
| INTEGER FH, IERROR | 27 |
| INTEGER(KIND=MPI_OFFSET_KIND) SIZE | 28 |
| MPI_FILE_SET_VIEW(FH, DISP, ETYPE, FILETYPE, DATAREP, INFO, IERROR) | 29 |
| INTEGER FH, ETYPE, FILETYPE, INFO, IERROR | 30 |
| CHARACTER*(*) DATAREP | 31 |
| INTEGER(KIND=MPI_OFFSET_KIND) DISP | 32 |
| MPI_FILE_SYNC(FH, IERROR) | 33 |
| INTEGER FH, IERROR | 34 |
| MPI_FILE_WRITE(FH, BUF, COUNT, DATATYPE, STATUS, IERROR) | 35 |
| <type> BUF(*) | 36 |
| INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR | 37 |
| MPI_FILE_WRITE_ALL(FH, BUF, COUNT, DATATYPE, STATUS, IERROR) | 38 |
| <type> BUF(*) | 39 |
| INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR | 40 |
| | 41 |
| | 42 |
| | 43 |
| | 44 |
| | 45 |
| | 46 |
| | 47 |
| | 48 |

```
1 MPI_FILE_WRITE_ALL_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR)
2     <type> BUF(*)
3     INTEGER FH, COUNT, DATATYPE, IERROR
4
5 MPI_FILE_WRITE_ALL_END(FH, BUF, STATUS, IERROR)
6     <type> BUF(*)
7     INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
8
9 MPI_FILE_WRITE_AT(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR)
10    <type> BUF(*)
11    INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
12    INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
13
14 MPI_FILE_WRITE_AT_ALL(FH, OFFSET, BUF, COUNT, DATATYPE, STATUS, IERROR)
15    <type> BUF(*)
16    INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
17    INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
18
19 MPI_FILE_WRITE_AT_ALL_BEGIN(FH, OFFSET, BUF, COUNT, DATATYPE, IERROR)
20    <type> BUF(*)
21    INTEGER FH, COUNT, DATATYPE, IERROR
22    INTEGER(KIND=MPI_OFFSET_KIND) OFFSET
23
24 MPI_FILE_WRITE_AT_ALL_END(FH, BUF, STATUS, IERROR)
25    <type> BUF(*)
26    INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
27
28 MPI_FILE_WRITE_ORDERED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
29    <type> BUF(*)
30    INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
31
32 MPI_FILE_WRITE_ORDERED_BEGIN(FH, BUF, COUNT, DATATYPE, IERROR)
33    <type> BUF(*)
34    INTEGER FH, COUNT, DATATYPE, IERROR
35
36 MPI_FILE_WRITE_ORDERED_END(FH, BUF, STATUS, IERROR)
37    <type> BUF(*)
38    INTEGER FH, STATUS(MPI_STATUS_SIZE), IERROR
39
40 MPI_FILE_WRITE_SHARED(FH, BUF, COUNT, DATATYPE, STATUS, IERROR)
41    <type> BUF(*)
42    INTEGER FH, COUNT, DATATYPE, STATUS(MPI_STATUS_SIZE), IERROR
43
44 MPI_REGISTER_DATAREP(DATAREP, READ_CONVERSION_FN, WRITE_CONVERSION_FN,
45     DTYPE_FILE_EXTENT_FN, EXTRA_STATE, IERROR)
46     CHARACTER*(*) DATAREP
47     EXTERNAL READ_CONVERSION_FN, WRITE_CONVERSION_FN, DTYPE_FILE_EXTENT_FN
48     INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
49     INTEGER IERROR
```

A.3.11 Language Bindings Fortran Bindings

```

MPI_SIZEOF(X, SIZE, IERROR)
  <type> X
  INTEGER SIZE, IERROR

MPI_TYPE_CREATE_F90_COMPLEX(P, R, NEWTYPE, IERROR)
  INTEGER P, R, NEWTYPE, IERROR

MPI_TYPE_CREATE_F90_INTEGER(R, NEWTYPE, IERROR)
  INTEGER R, NEWTYPE, IERROR

MPI_TYPE_CREATE_F90_REAL(P, R, NEWTYPE, IERROR)
  INTEGER P, R, NEWTYPE, IERROR

MPI_TYPE_MATCH_SIZE(TYPECLASS, SIZE, TYPE, IERROR)
  INTEGER TYPECLASS, SIZE, TYPE, IERROR

```

A.3.12 Profiling Interface Fortran Bindings

```

MPI_PCONTROL(LEVEL)
  INTEGER LEVEL, ...

```

A.3.13 Deprecated Fortran Bindings

```

MPI_ADDRESS(LOCATION, ADDRESS, IERROR)
  <type> LOCATION(*)
  INTEGER ADDRESS, IERROR

MPI_ATTR_DELETE(COMM, KEYVAL, IERROR)
  INTEGER COMM, KEYVAL, IERROR

MPI_ATTR_GET(COMM, KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)
  INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, IERROR
  LOGICAL FLAG

MPI_ATTR_PUT(COMM, KEYVAL, ATTRIBUTE_VAL, IERROR)
  INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, IERROR

MPI_ERRHANDLER_CREATE(FUNCTION, ERRHANDLER, IERROR)
  EXTERNAL FUNCTION
  INTEGER ERRHANDLER, IERROR

MPI_ERRHANDLER_GET(COMM, ERRHANDLER, IERROR)
  INTEGER COMM, ERRHANDLER, IERROR

MPI_ERRHANDLER_SET(COMM, ERRHANDLER, IERROR)
  INTEGER COMM, ERRHANDLER, IERROR

MPI_KEYVAL_CREATE(COPY_FN, DELETE_FN, KEYVAL, EXTRA_STATE, IERROR)
  EXTERNAL COPY_FN, DELETE_FN
  INTEGER KEYVAL, EXTRA_STATE, IERROR

```

```
1 MPI_KEYVAL_FREE(KEYVAL, IERROR)
2     INTEGER KEYVAL, IERROR
3
4 MPI_TYPE_EXTENT(DATATYPE, EXTENT, IERROR)
5     INTEGER DATATYPE, EXTENT, IERROR
6
7 MPI_TYPE_HINDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
8     OLDTYPE, NEWTYPE, IERROR)
9     INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*),
10    OLDTYPE, NEWTYPE, IERROR
11
12 MPI_TYPE_HVECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR)
13     INTEGER COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR
14
15 MPI_TYPE_LB( DATATYPE, DISPLACEMENT, IERROR)
16     INTEGER DATATYPE, DISPLACEMENT, IERROR
17
18 MPI_TYPE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS,
19     ARRAY_OF_TYPES, NEWTYPE, IERROR)
20     INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*),
21     ARRAY_OF_TYPES(*), NEWTYPE, IERROR
22
23 MPI_TYPE_UB( DATATYPE, DISPLACEMENT, IERROR)
24     INTEGER DATATYPE, DISPLACEMENT, IERROR
25
26 SUBROUTINE COPY_FUNCTION(OLDCOMM, KEYVAL, EXTRA_STATE, ATTRIBUTE_VAL_IN,
27     ATTRIBUTE_VAL_OUT, FLAG, IERR)
28     INTEGER OLDCOMM, KEYVAL, EXTRA_STATE, ATTRIBUTE_VAL_IN,
29     ATTRIBUTE_VAL_OUT, IERR
30     LOGICAL FLAG
31
32 SUBROUTINE DELETE_FUNCTION(COMM, KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE, IERR)
33     INTEGER COMM, KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE, IERR
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
```


A.4 C++ Bindings

A.4.1 Point-to-Point Communication C++ Bindings

```
namespace MPI {  
  
    void Attach_buffer(void* buffer, int size)  
  
    void Comm::Bsend(const void* buf, int count, const Datatype& datatype,  
                    int dest, int tag) const  
  
    Prequest Comm::Bsend_init(const void* buf, int count, const  
                              Datatype& datatype, int dest, int tag) const  
  
    Request Comm::Ibsend(const void* buf, int count, const  
                        Datatype& datatype, int dest, int tag) const  
  
    bool Comm::Iprobe(int source, int tag) const  
  
    bool Comm::Iprobe(int source, int tag, Status& status) const  
  
    Request Comm::Irecv(void* buf, int count, const Datatype& datatype,  
                       int source, int tag) const  
  
    Request Comm::Irsend(const void* buf, int count, const  
                        Datatype& datatype, int dest, int tag) const  
  
    Request Comm::Isend(const void* buf, int count, const Datatype& datatype,  
                       int dest, int tag) const  
  
    Request Comm::Issend(const void* buf, int count, const  
                        Datatype& datatype, int dest, int tag) const  
  
    void Comm::Probe(int source, int tag) const  
  
    void Comm::Probe(int source, int tag, Status& status) const  
  
    void Comm::Recv(void* buf, int count, const Datatype& datatype,  
                   int source, int tag) const  
  
    void Comm::Recv(void* buf, int count, const Datatype& datatype,  
                   int source, int tag, Status& status) const  
  
    Prequest Comm::Recv_init(void* buf, int count, const Datatype& datatype,  
                             int source, int tag) const  
  
    void Comm::Rsend(const void* buf, int count, const Datatype& datatype,  
                    int dest, int tag) const  
  
    Prequest Comm::Rsend_init(const void* buf, int count, const  
                              Datatype& datatype, int dest, int tag) const  
  
    void Comm::Send(const void* buf, int count, const Datatype& datatype,  
                   int dest, int tag) const  
  
    Prequest Comm::Send_init(const void* buf, int count, const  
                             Datatype& datatype, int dest, int tag) const
```

```

1 void Comm::Sendrecv(const void *sendbuf, int sendcount, const
2     Datatype& sendtype, int dest, int sendtag, void *recvbuf,
3     int recvcount, const Datatype& recvtype, int source,
4     int recvtag) const
5
6 void Comm::Sendrecv(const void *sendbuf, int sendcount, const
7     Datatype& sendtype, int dest, int sendtag, void *recvbuf,
8     int recvcount, const Datatype& recvtype, int source,
9     int recvtag, Status& status) const
10
11 void Comm::Sendrecv_replace(void* buf, int count, const
12     Datatype& datatype, int dest, int sendtag, int source,
13     int recvtag) const
14
15 void Comm::Sendrecv_replace(void* buf, int count, const
16     Datatype& datatype, int dest, int sendtag, int source,
17     int recvtag, Status& status) const
18
19 void Comm::Ssend(const void* buf, int count, const Datatype& datatype,
20     int dest, int tag) const
21
22 Prequest Comm::Ssend_init(const void* buf, int count, const
23     Datatype& datatype, int dest, int tag) const
24
25 void Datatype::Commit()
26
27 Datatype Datatype::Create_contiguous(int count) const
28
29 Datatype Datatype::Create_darray(int size, int rank, int ndims,
30     const int array_of_gsizes[], const int array_of_distrib[],
31     const int array_of_dargs[], const int array_of_psizes[],
32     int order) const
33
34 Datatype Datatype::Create_hindexed(int count,
35     const int array_of_blocklengths[],
36     const Aint array_of_displacements[]) const
37
38 Datatype Datatype::Create_hvector(int count, int blocklength, Aint
39     stride) const
40
41 Datatype Datatype::Create_indexed(int count,
42     const int array_of_blocklengths[],
43     const int array_of_displacements[]) const
44
45 Datatype Datatype::Create_indexed_block(int count, int blocklength,
46     const int array_of_displacements[]) const
47
48 Datatype Datatype::Create_resized(const Aint lb, const Aint extent) const
49
50 static Datatype Datatype::Create_struct(int count,
51     const int array_of_blocklengths[], const Aint
52     array_of_displacements[], const Datatype array_of_types[])
53
54 Datatype Datatype::Create_subarray(int ndims, const int array_of_sizes[],
55     const int array_of_subsizes[], const int array_of_starts[],

```

```

        int order) const 1
Datatype Datatype::Create_vector(int count, int blocklength, int stride) 2
        const 3
Datatype Datatype::Dup() const 4
void Datatype::Free() 5
void Datatype::Get_extent(Aint& lb, Aint& extent) const 6
int Datatype::Get_size() const 7
void Datatype::Get_true_extent(Aint& true_lb, Aint& true_extent) const 8
void Datatype::Pack(const void* inbuf, int incount, void *outbuf, 9
        int outsize, int& position, const Comm &comm) const 10
void Datatype::Pack_external(const char* datarep, const void* inbuf, 11
        int incount, void* outbuf, Aint outsize, Aint& position) const 12
Aint Datatype::Pack_external_size(const char* datarep, int incount) const 13
int Datatype::Pack_size(int incount, const Comm& comm) const 14
void Datatype::Unpack(const void* inbuf, int insize, void *outbuf, 15
        int outcount, int& position, const Comm& comm) const 16
void Datatype::Unpack_external(const char* datarep, const void* inbuf, 17
        Aint insize, Aint& position, void* outbuf, int outcount) const 18
int Detach_buffer(void*& buffer) 19
Aint Get_address(void* location) 20
void Prequest::Start() 21
static void Prequest::Startall(int count, Prequest array_of_requests[]) 22
void Request::Cancel() const 23
void Request::Free() 24
bool Request::Get_status() const 25
bool Request::Get_status(Status& status) const 26
bool Request::Test() 27
bool Request::Test(Status& status) 28
static bool Request::Testall(int count, Request array_of_requests[]) 29
static bool Request::Testall(int count, Request array_of_requests[], 30
        Status array_of_statuses[]) 31
static bool Request::Testany(int count, Request array_of_requests[], 32
        int& index) 33

```

```

1   static bool Request::Testany(int count, Request array_of_requests[],
2       int& index, Status& status)
3
4   static int Request::Testsome(int incount, Request array_of_requests[],
5       int array_of_indices[])
6
7   static int Request::Testsome(int incount, Request array_of_requests[],
8       int array_of_indices[], Status array_of_statuses[])
9
10  void Request::Wait()
11
12  void Request::Wait(Status& status)
13
14  static void Request::Waitall(int count, Request array_of_requests[])
15
16  static void Request::Waitall(int count, Request array_of_requests[],
17      Status array_of_statuses[])
18
19  static int Request::Waitany(int count, Request array_of_requests[])
20
21  static int Request::Waitany(int count, Request array_of_requests[],
22      Status& status)
23
24  static int Request::Waitsome(int incount, Request array_of_requests[],
25      int array_of_indices[])
26
27  static int Request::Waitsome(int incount, Request array_of_requests[],
28      int array_of_indices[], Status array_of_statuses[])
29
30  int Status::Get_count(const Datatype& datatype) const
31
32  int Status::Get_elements(const Datatype& datatype) const
33
34  int Status::Get_error() const
35
36  int Status::Get_source() const
37
38  int Status::Get_tag() const
39
40  bool Status::Is_cancelled() const
41
42  void Status::Set_error(int error)
43
44  void Status::Set_source(int source)
45
46  void Status::Set_tag(int tag)
47
48  };

```

A.4.2 Collective Communication C++ Bindings

```

43 namespace MPI {
44
45     void Comm::Allgather(const void* sendbuf, int sendcount, const
46         Datatype& sendtype, void* recvbuf, int recvcount,
47         const Datatype& recvtype) const = 0
48

```



```

1   void Intracomm::Scan(const void* sendbuf, void* recvbuf, int count, const
2       Datatype& datatype, const Op& op) const
3
4   void Op::Free()
5
6   void Op::Init(User_function* function, bool commute)
7
8 };

```

A.4.3 Groups, Contexts, and Communicators C++ Bindings

```

11 namespace MPI {
12
13     Cartcomm& Cartcomm::Clone() const
14
15     Cartcomm Cartcomm::Dup() const
16
17     Comm& Comm::Clone() const = 0
18
19     static int Comm::Compare(const Comm& comm1, const Comm& comm2)
20
21     static int Comm::Create_keyval(Comm::Copy_attr_function* comm_copy_attr_fn,
22         Comm::Delete_attr_function* comm_delete_attr_fn,
23         void* extra_state)
24
25     void Comm::Delete_attr(int comm_keyval)
26
27     void Comm::Free()
28
29     static void Comm::Free_keyval(int& comm_keyval)
30
31     bool Comm::Get_attr(int comm_keyval, void* attribute_val) const
32
33     Group Comm::Get_group() const
34
35     int Comm::Get_rank() const
36
37     int Comm::Get_size() const
38
39     bool Comm::Is_inter() const
40
41     void Comm::Set_attr(int comm_keyval, const void* attribute_val) const
42
43     static int Datatype::Create_keyval(Datatype::Copy_attr_function*
44         type_copy_attr_fn, Datatype::Delete_attr_function*
45         type_delete_attr_fn, void* extra_state)
46
47     void Datatype::Delete_attr(int type_keyval)
48
49     static void Datatype::Free_keyval(int& type_keyval)
50
51     bool Datatype::Get_attr(int type_keyval, void* attribute_val) const
52
53     void Datatype::Set_attr(int type_keyval, const void* attribute_val)
54
55     Graphcomm& Graphcomm::Clone() const
56
57     Graphcomm Graphcomm::Dup() const

```

```

static int Group::Compare(const Group& group1, const Group& group2)      1
static Group Group::Difference(const Group& group1, const Group& group2)  2
Group Group::Excl(int n, const int ranks[]) const                        3
void Group::Free()                                                       4
int Group::Get_rank() const                                              5
int Group::Get_size() const                                              6
Group Group::Incl(int n, const int ranks[]) const                        7
static Group Group::Intersect(const Group& group1, const Group& group2)  8
Group Group::Range_excl(int n, const int ranges[][3]) const            9
Group Group::Range_incl(int n, const int ranges[][3]) const           10
static void Group::Translate_ranks (const Group& group1, int n,
                                   const int ranks1[], const Group& group2, int ranks2[]) 11
static Group Group::Union(const Group& group1, const Group& group2)     12
Intercomm& Intercomm::Clone() const                                     13
Intercomm Intercomm::Create(const Group& group) const                  14
Intercomm Intercomm::Dup() const                                       15
Group Intercomm::Get_remote_group() const                               16
int Intercomm::Get_remote_size() const                                  17
Intracomm Intercomm::Merge(bool high) const                             18
Intercomm Intercomm::Split(int color, int key) const                   19
Intracomm& Intracomm::Clone() const                                    20
Intracomm Intracomm::Create(const Group& group) const                  21
Intercomm Intracomm::Create_intercomm(int local_leader, const
                                       Comm& peer_comm, int remote_leader, int tag) const 22
Intracomm Intracomm::Dup() const                                       23
Intracomm Intracomm::Split(int color, int key) const                   24
static int Win::Create_keyval(Win::Copy_attr_function* win_copy_attr_fn, 25
                              Win::Delete_attr_function* win_delete_attr_fn,
                              void* extra_state)                        26
void Win::Delete_attr(int win_keyval)                                    27
static void Win::Free_keyval(int& win_keyval)                            28
bool Win::Get_attr(int win_keyval, void* attribute_val) const           29
void Win::Set_attr(int win_keyval, const void* attribute_val)           30

```

```

1  };
2
3  A.4.4 Process Topologies C++ Bindings
4
5  namespace MPI {
6
7      int Cartcomm::Get_cart_rank(const int coords[]) const
8
9      void Cartcomm::Get_coords(int rank, int maxdims, int coords[]) const
10
11     int Cartcomm::Get_dim() const
12
13     void Cartcomm::Get_topo(int maxdims, int dims[], bool periods[],
14                             int coords[]) const
15
16     int Cartcomm::Map(int ndims, const int dims[], const bool periods[])
17         const
18
19     void Cartcomm::Shift(int direction, int disp, int& rank_source,
20                           int& rank_dest) const
21
22     Cartcomm Cartcomm::Sub(const bool remain_dims[]) const
23
24     int Comm::Get_topology() const
25
26     void Compute_dims(int nnodes, int ndims, int dims[])
27
28     void Graphcomm::Get_dims(int nnodes[], int nedges[]) const
29
30     void Graphcomm::Get_neighbors(int rank, int maxneighbors, int
31                                   neighbors[]) const
32
33     int Graphcomm::Get_neighbors_count(int rank) const
34
35     void Graphcomm::Get_topo(int maxindex, int maxedges, int index[],
36                               int edges[]) const
37
38     int Graphcomm::Map(int nnodes, const int index[], const int edges[])
39         const
40
41     Cartcomm Intracomm::Create_cart(int ndims, const int dims[],
42                                     const bool periods[], bool reorder) const
43
44     Graphcomm Intracomm::Create_graph(int nnodes, const int index[],
45                                       const int edges[], bool reorder) const
46
47 };
48
49 A.4.5 MPI Environmenta Management C++ Bindings
50
51 namespace MPI {
52
53     void* Alloc_mem(Aint size, const Info& info)
54
55     void Comm::Abort(int errorcode)
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

```



```
void Comm::Abort(int errorcode) 1
static Errhandler Comm::Create_errhandler(Comm::Errhandler_fn* function) 2
Errhandler Comm::Get_errhandler() const 3
void Comm::Set_errhandler(const Errhandler& errhandler) 4
void Errhandler::Free() 5
void Errhandler::Free() 6
static Errhandler File::Create_errhandler(File::Errhandler_fn* function) 7
Errhandler File::Get_errhandler() const 8
void File::Set_errhandler(const Errhandler& errhandler) 9
void Finalize() 10
void Finalize() 11
void Free_mem(void *base) 12
int Get_error_class(int errorcode) 13
int Get_error_class(int errorcode) 14
void Get_error_string(int errorcode, char* name, int& resultlen) 15
void Get_error_string(int errorcode, char* name, int& resultlen) 16
void Get_processor_name(char* name, int& resultlen) 17
void Get_processor_name(char* name, int& resultlen) 18
void Get_version(int& version, int& subversion) 19
void Init() 20
void Init() 21
void Init(int& argc, char**& argv) 22
void Init(int& argc, char**& argv) 23
Intracomm Intracomm::Create(const Group& group) const 24
bool Is_finalized() 25
bool Is_initialized() 26
bool Is_initialized() 27
static Errhandler Win::Create_errhandler(Win::Errhandler_fn* function) 28
Errhandler Win::Get_errhandler() const 29
void Win::Set_errhandler(const Errhandler& errhandler) 30
double Wtick() 31
```

```

1     double Wtime()
2
3 };
4

```

A.4.6 Miscellany C++ Bindings

```

7 namespace MPI {
8
9     static Info Info::Create()
10    void Info::Delete(const char* key)
11
12    Info Info::Dup() const
13
14    void Info::Free()
15
16    bool Info::Get(const char* key, int valuelen, char* value) const
17
18    int Info::Get_nkeys() const
19
20    void Info::Get_nthkey(int n, char* key) const
21
22    bool Info::Get_valuelen(const char* key, int& valuelen) const
23
24    void Info::Set(const char* key, const char* value)
25
26 };

```

A.4.7 Process Creation and Management C++ Bindings

```

27 namespace MPI {
28
29    void Close_port(const char* port_name)
30
31    void Comm::Disconnect()
32
33    static Intercomm Comm::Get_parent()
34
35    static Intercomm Comm::Join(const int fd)
36
37    Intercomm Intracomm::Accept(const char* port_name, const Info& info,
38                                int root) const
39
40    Intercomm Intracomm::Connect(const char* port_name, const Info& info,
41                                int root) const
42
43    Intercomm Intracomm::Spawn(const char* command, const char* argv[],
44                                int maxprocs, const Info& info, int root) const
45
46    Intercomm Intracomm::Spawn(const char* command, const char* argv[],
47                                int maxprocs, const Info& info, int root,
48                                int array_of_errcodes[]) const
49
50    Intercomm Intracomm::Spawn_multiple(int count,
51                                        const char* array_of_commands[], const char** array_of_argv[],

```

```

        const int array_of_maxprocs[], const Info array_of_info[],
        int root)
1
2
3
Intercomm Intracomm::Spawn_multiple(int count,
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
void Lookup_name(const char* service_name, const Info& info,
char* port_name)
void Open_port(const Info& info, char* port_name)
void Publish_name(const char* service_name, const Info& info,
const char* port_name)
void Unpublish_name(const char* service_name, const Info& info,
const char* port_name)
};

```

A.4.8 One-Sided Communications C++ Bindings

```

namespace MPI {
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
void Win::Accumulate(const void* origin_addr, int origin_count, const
Datatype& origin_datatype, int target_rank, Aint target_disp,
int target_count, const Datatype& target_datatype, const Op&
op) const
void Win::Complete() const
static Win Win::Create(const void* base, Aint size, int disp_unit, const
Info& info, const Intracomm& comm)
void Win::Fence(int assert) const
void Win::Free()
void Win::Get(void *origin_addr, int origin_count, const Datatype&
origin_datatype, int target_rank, Aint target_disp, int
target_count, const Datatype& target_datatype) const
Group Win::Get_group() const
void Win::Lock(int lock_type, int rank, int assert) const
void Win::Post(const Group& group, int assert) const
void Win::Put(const void* origin_addr, int origin_count, const Datatype&
origin_datatype, int target_rank, Aint target_disp, int
target_count, const Datatype& target_datatype) const
void Win::Start(const Group& group, int assert) const

```

```

1   bool Win::Test() const
2
3   void Win::Unlock(int rank) const
4
5   void Win::Wait() const
6
7 };

```

A.4.9 External Interfaces C++ Bindings

```

10 namespace MPI {
11
12     int Add_error_class()
13
14     int Add_error_code(int errorclass)
15
16     void Add_error_string(int errorcode, const char* string)
17
18     void Comm::Call_errhandler(int errorcode) const
19
20     void Comm::Get_name(char* comm_name, int& resultlen) const
21
22     void Comm::Set_name(const char* comm_name)
23
24     void Datatype::Get_contents(int max_integers, int max_addresses,
25                                int max_datatypes, int array_of_integers[],
26                                Aint array_of_addresses[], Datatype array_of_datatypes[]) const
27
28     void Datatype::Get_envelope(int& num_integers, int& num_addresses,
29                                int& num_datatypes, int& combiner) const
30
31     void Datatype::Get_name(char* type_name, int& resultlen) const
32
33     void Datatype::Set_name(const char* type_name)
34
35     void File::Call_errhandler(int errorcode) const
36
37     void Grequest::Complete()
38
39     static Grequest Grequest::Start(const Grequest::Query_function query_fn,
40                                    const Grequest::Free_function free_fn,
41                                    const Grequest::Cancel_function cancel_fn, void *extra_state)
42
43     int Init_thread(int required)
44
45     int Init_thread(int& argc, char**& argv, int required)
46
47     bool Is_thread_main()
48
49     int Query_thread()
50
51     void Status::Set_cancelled(bool flag)
52
53     void Status::Set_elements(const Datatype& datatype, int count)
54
55     void Win::Call_errhandler(int errorcode) const
56
57     void Win::Get_name(char* win_name, int& resultlen) const

```

```
void Win::Set_name(const char* win_name) 1  
2  
}; 3  
4
```

A.4.10 I/O C++ Bindings

```
namespace MPI { 5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48
```

```
void File::Close()  
static void File::Delete(const char* filename, const Info& info)  
int File::Get_amode() const  
bool File::Get_atomicity() const  
Offset File::Get_byte_offset(const Offset disp) const  
Group File::Get_group() const  
Info File::Get_info() const  
Offset File::Get_position() const  
Offset File::Get_position_shared() const  
Offset File::Get_size() const  
Aint File::Get_type_extent(const Datatype& datatype) const  
void File::Get_view(Offset& disp, Datatype& etype, Datatype& filetype,  
char* datarep) const  
Request File::Iread(void* buf, int count, const Datatype& datatype)  
Request File::Iread_at(Offset offset, void* buf, int count,  
const Datatype& datatype)  
Request File::Iread_shared(void* buf, int count,  
const Datatype& datatype)  
Request File::Iwrite(const void* buf, int count,  
const Datatype& datatype)  
Request File::Iwrite_at(Offset offset, const void* buf, int count,  
const Datatype& datatype)  
Request File::Iwrite_shared(const void* buf, int count,  
const Datatype& datatype)  
static File File::Open(const Intracomm& comm, const char* filename,  
int amode, const Info& info)  
void File::Preallocate(Offset size)  
void File::Read(void* buf, int count, const Datatype& datatype)
```

```
1 void File::Read(void* buf, int count, const Datatype& datatype, Status&
2     status)
3
4 void File::Read_all(void* buf, int count, const Datatype& datatype)
5
6 void File::Read_all(void* buf, int count, const Datatype& datatype,
7     Status& status)
8
9 void File::Read_all_begin(void* buf, int count, const Datatype& datatype)
10
11 void File::Read_all_end(void* buf)
12
13 void File::Read_all_end(void* buf, Status& status)
14
15 void File::Read_at(Offset offset, void* buf, int count,
16     const Datatype& datatype)
17
18 void File::Read_at(Offset offset, void* buf, int count,
19     const Datatype& datatype, Status& status)
20
21 void File::Read_at_all(Offset offset, void* buf, int count,
22     const Datatype& datatype)
23
24 void File::Read_at_all(Offset offset, void* buf, int count,
25     const Datatype& datatype, Status& status)
26
27 void File::Read_at_all_begin(Offset offset, void* buf, int count,
28     const Datatype& datatype)
29
30 void File::Read_at_all_end(void* buf)
31
32 void File::Read_at_all_end(void* buf, Status& status)
33
34 void File::Read_ordered(void* buf, int count, const Datatype& datatype)
35
36 void File::Read_ordered(void* buf, int count, const Datatype& datatype,
37     Status& status)
38
39 void File::Read_ordered_begin(void* buf, int count,
40     const Datatype& datatype)
41
42 void File::Read_ordered_end(void* buf)
43
44 void File::Read_ordered_end(void* buf, Status& status)
45
46 void File::Read_shared(void* buf, int count, const Datatype& datatype)
47
48 void File::Read_shared(void* buf, int count, const Datatype& datatype,
49     Status& status)
50
51 void File::Seek(Offset offset, int whence)
52
53 void File::Seek_shared(Offset offset, int whence)
54
55 void File::Set_atomicity(bool flag)
56
57 void File::Set_info(const Info& info)
58
59 void File::Set_size(Offset size)
```

```
void File::Set_view(Offset disp, const Datatype& etype,      1
                    const Datatype& filetype, const char* datarep,  2
                    const Info& info)                          3
void File::Sync()                                           4
void File::Write(const void* buf, int count, const Datatype& datatype)  6
void File::Write(const void* buf, int count, const Datatype& datatype,  8
                 Status& status)                             9
void File::Write_all(const void* buf, int count,             10
                    const Datatype& datatype)               11
void File::Write_all(const void* buf, int count,             12
                    const Datatype& datatype, Status& status)  14
void File::Write_all_begin(const void* buf, int count,       15
                           const Datatype& datatype)         16
void File::Write_all_end(const void* buf)                    18
void File::Write_all_end(const void* buf, Status& status)    19
void File::Write_at(Offset offset, const void* buf, int count,  21
                   const Datatype& datatype)                 22
void File::Write_at(Offset offset, const void* buf, int count,  23
                   const Datatype& datatype, Status& status)  24
void File::Write_at.all(Offset offset, const void* buf, int count,  26
                       const Datatype& datatype)              27
void File::Write_at.all(Offset offset, const void* buf, int count,  28
                       const Datatype& datatype, Status& status)  29
void File::Write_at.all_begin(Offset offset, const void* buf, int count,  31
                              const Datatype& datatype)        32
void File::Write_at.all_end(const void* buf)                 34
void File::Write_at.all_end(const void* buf, Status& status)  35
void File::Write_ordered(const void* buf, int count,         37
                        const Datatype& datatype)            38
void File::Write_ordered(const void* buf, int count,         39
                        const Datatype& datatype, Status& status)  40
void File::Write_ordered_begin(const void* buf, int count,   42
                               const Datatype& datatype)      43
void File::Write_ordered_end(const void* buf)                44
void File::Write_ordered_end(const void* buf, Status& status)  46
```

```

1  void File::Write_shared(const void* buf, int count,
2  const Datatype& datatype)
3
4  void File::Write_shared(const void* buf, int count,
5  const Datatype& datatype, Status& status)
6
7  void Register_datarep(const char* datarep,
8  Datarep_conversion_function* read_conversion_fn,
9  Datarep_conversion_function* write_conversion_fn,
10 Datarep_extent_function* dtype_file_extent_fn,
11 void* extra_state)
12
13 };

```

A.4.11 Language Bindings C++ Bindings

```

14 namespace MPI {
15
16     static Datatype Datatype::Create_f90_complex(int p, int r)
17
18     static Datatype Datatype::Create_f90_integer(int r)
19
20     static Datatype Datatype::Create_f90_real(int p, int r)
21
22     static Datatype Datatype::Match_size(int typeclass, int size)
23
24     Exception::Exception(int error_code)
25
26     int Exception::Get_error_class() const
27
28     int Exception::Get_error_code() const
29
30     const char* Exception::Get_error_string() const
31
32 };

```

A.4.12 Profiling Interface C++ Bindings

```

33 namespace MPI {
34
35     void Pcontrol(const int level, ...)
36
37 };
38
39 };
40

```

A.4.13 Deprecated C++ Bindings

```

41 namespace MPI {
42
43 };
44
45 };
46
47
48

```


A.4.14 C++ Bindings on all MPI Classes

The C++ language requires all classes to have four special functions: a default constructor, a copy constructor, a destructor, and an assignment operator. The bindings for these functions are listed below; their semantics are discussed in Section 13.1.5. The two constructors are *not virtual*. The bindings prototype functions using the type `<CLASS>` rather than listing each function for every MPI class; the token `<CLASS>` can be replaced with valid MPI-2 class names, such as `Group`, `Datatype`, etc., except when noted. In addition, bindings are provided for comparison and inter-language operability from Sections 13.1.5 and 13.1.9.

A.4.15 Construction / Destruction

```
namespace MPI {
    <CLASS>::<CLASS>()
    <CLASS>::~~<CLASS>()
};
```

A.4.16 Copy / Assignment

```
namespace MPI {
    <CLASS>::<CLASS>(const <CLASS>& data)
    <CLASS>& <CLASS>::operator=(const <CLASS>& data)
};
```

A.4.17 Comparison

Since `Status` instances are not handles to underlying MPI objects, the `operator==()` and `operator!=()` functions are not defined on the `Status` class.

```
namespace MPI {
    bool <CLASS>::operator==(const <CLASS>& data) const
    bool <CLASS>::operator!=(const <CLASS>& data) const
};
```

A.4.18 Inter-language Operability

Since there are no C++ `MPI::STATUS_IGNORE` and `MPI::STATUSES_IGNORE` objects, the results of promoting the C or Fortran handles (`MPI_STATUS_IGNORE` and `MPI_STATUSES_IGNORE`) to C++ is undefined.

```
namespace MPI {
```

```

1     <CLASS>& <CLASS>::operator=(const MPI_<CLASS>& data)
2
3     <CLASS>::<CLASS>(const MPI_<CLASS>& data)
4
5     <CLASS>::operator MPI_<CLASS>() const
6
7 };
8

```

A.4.19 Function Name Cross Reference

Since some of the C++ bindings have slightly different names than their C and Fortran counterparts, this section maps each language neutral MPI-1 name to its corresponding C++ binding.

For brevity, the “MPI: :” prefix is assumed for all C++ class names.

Where MPI-1 names have been deprecated, the `<none>` keyword is used in the “Member name” column to indicate that this function is supported with a new name (see Annex A).

Where non-void values are listed in the “Return value” column, the given name is that of the corresponding parameter in the language neutral specification.

| MPI Function | C++ class | Member name | Return value | |
|-----------------------|-----------|------------------|--------------------|----|
| MPI_ABORT | Comm | Abort | void | 1 |
| MPI_ADDRESS | | <none> | | 2 |
| MPI_ALLGATHERV | Intracomm | Allgatherv | void | 3 |
| MPI_ALLGATHER | Intracomm | Allgather | void | 4 |
| MPI_ALLREDUCE | Intracomm | Allreduce | void | 5 |
| MPI_ALLTOALLV | Intracomm | Alltoallv | void | 6 |
| MPI_ALLTOALL | Intracomm | Alltoall | void | 7 |
| MPI_ATTR_DELETE | | <none> | | 8 |
| MPI_ATTR_GET | | <none> | | 9 |
| MPI_ATTR_PUT | | <none> | | 10 |
| MPI_BARRIER | Intracomm | Barrier | void | 11 |
| MPI_BCAST | Intracomm | Bcast | void | 12 |
| MPI_BSEND_INIT | Comm | Bsend_init | Prequest request | 13 |
| MPI_BSEND | Comm | Bsend | void | 14 |
| MPI_BUFFER_ATTACH | | Attach_buffer | void | 15 |
| MPI_BUFFER_DETACH | | Detach_buffer | void* buffer | 16 |
| MPI_CANCEL | Request | Cancel | void | 17 |
| MPI_CARTDIM_GET | Cartcomm | Get_dim | int ndims | 18 |
| MPI_CART_COORDS | Cartcomm | Get_coords | void | 19 |
| MPI_CART_CREATE | Intracomm | Create_cart | Cartcomm newcomm | 20 |
| MPI_CART_GET | Cartcomm | Get_topo | void | 21 |
| MPI_CART_MAP | Cartcomm | Map | int newrank | 22 |
| MPI_CART_RANK | Cartcomm | Get_rank | int rank | 23 |
| MPI_CART_SHIFT | Cartcomm | Shift | void | 24 |
| MPI_CART_SUB | Cartcomm | Sub | Cartcomm newcomm | 25 |
| MPI_COMM_COMPARE | Comm | static Compare | int result | 26 |
| MPI_COMM_CREATE | Intracomm | Create | Intracomm newcomm | 27 |
| MPI_COMM_DUP | Intracomm | Dup | Intracomm newcomm | 28 |
| | Cartcomm | Dup | Cartcomm newcomm | 29 |
| | Graphcomm | Dup | Graphcomm newcomm | 30 |
| | Intercomm | Dup | Intercomm newcomm | 31 |
| | Comm | Clone | Comm& newcomm | 32 |
| | Intracomm | Clone | Intracomm& newcomm | 33 |
| | Cartcomm | Clone | Cartcomm& newcomm | 34 |
| | Graphcomm | Clone | Graphcomm& newcomm | 35 |
| | Intercomm | Clone | Intercomm& newcomm | 36 |
| MPI_COMM_FREE | Comm | Free | void | 37 |
| MPI_COMM_GROUP | Comm | Get_group | Group group | 38 |
| MPI_COMM_RANK | Comm | Get_rank | int rank | 39 |
| MPI_COMM_REMOTE_GROUP | Intercomm | Get_remote_group | Group group | 40 |
| MPI_COMM_REMOTE_SIZE | Intercomm | Get_remote_size | int size | 41 |
| MPI_COMM_SIZE | Comm | Get_size | int size | 42 |
| MPI_COMM_SPLIT | Intracomm | Split | Intracomm newcomm | 43 |
| MPI_COMM_TEST_INTER | Comm | Is_inter | bool flag | 44 |
| MPI_DIMS_CREATE | | Compute_dims | void | 45 |

47

48

| 1 | MPI Function | C++ class | Member name | Return value |
|----|---------------------------|------------|------------------------|-------------------|
| 2 | MPI_ERRHANDLER_CREATE | | <none> | |
| 3 | MPI_ERRHANDLER_FREE | Errhandler | Free | void |
| 4 | MPI_ERRHANDLER_GET | | <none> | |
| 5 | MPI_ERRHANDLER_SET | | <none> | |
| 6 | MPI_ERROR_CLASS | | Get_error_class | int errorclass |
| 7 | MPI_ERROR_STRING | | Get_error_string | void |
| 8 | MPI_FINALIZE | | Finalize | void |
| 9 | MPI_GATHERV | Intracomm | Gatherv | void |
| 10 | MPI_GATHER | Intracomm | Gather | void |
| 11 | MPI_GET_COUNT | Status | Get_count | int count |
| 12 | MPI_GET_ELEMENTS | Status | Get_elements | int count |
| 13 | MPI_GET_PROCESSOR_NAME | | Get_processor_name | void |
| 14 | MPI_GRAPHDIMS_GET | Graphcomm | Get_dims | void |
| 15 | MPI_GRAPH_CREATE | Intracomm | Create_graph | Graphcomm newcomm |
| 16 | MPI_GRAPH_GET | Graphcomm | Get_topo | void |
| 17 | MPI_GRAPH_MAP | Graphcomm | Map | int newrank |
| 18 | MPI_GRAPH_NEIGHBORS_COUNT | Graphcomm | Get_neighbors_count | int nneighbors |
| 19 | MPI_GRAPH_NEIGHBORS | Graphcomm | Get_neighbors | void |
| 20 | MPI_GROUP_COMPARE | Group | static Compare | int result |
| 21 | MPI_GROUP_DIFFERENCE | Group | static Difference | Group newgroup |
| 22 | MPI_GROUP_EXCL | Group | Excl | Group newgroup |
| 23 | MPI_GROUP_FREE | Group | Free | void |
| 24 | MPI_GROUP_INCL | Group | Incl | Group newgroup |
| 25 | MPI_GROUP_INTERSECTION | Group | static Intersect | Group newgroup |
| 26 | MPI_GROUP_RANGE_EXCL | Group | Range_excl | Group newgroup |
| 27 | MPI_GROUP_RANGE_INCL | Group | Range_incl | Group newgroup |
| 28 | MPI_GROUP_RANK | Group | Get_rank | int rank |
| 29 | MPI_GROUP_SIZE | Group | Get_size | int size |
| 30 | MPI_GROUP_TRANSLATE_RANKS | Group | static Translate_ranks | void |
| 31 | MPI_GROUP_UNION | Group | static Union | Group newgroup |
| 32 | MPI_IBSEND | Comm | Ibsend | Request request |
| 33 | MPI_INITIALIZED | | Is_initialized | bool flag |
| 34 | MPI_INIT | | Init | void |
| 35 | MPI_INTERCOMM_CREATE | Intracomm | Create_intercomm | Intercomm newcomm |
| 36 | MPI_INTERCOMM_MERGE | Intercomm | Merge | Intracomm newcomm |
| 37 | MPI_IPROBE | Comm | Iprobe | bool flag |
| 38 | MPI_IRECV | Comm | Irecv | Request request |
| 39 | MPI_IRSEND | Comm | Irsend | Request request |
| 40 | MPI_ISEND | Comm | Isend | Request request |
| 41 | MPI_ISSEND | Comm | Issend | Request request |
| 42 | MPI_KEYVAL_CREATE | | <none> | |
| 43 | MPI_KEYVAL_FREE | | <none> | |
| 44 | MPI_OP_CREATE | Op | Init | void |
| 45 | MPI_OP_FREE | Op | Free | void |
| 46 | MPI_PACK_SIZE | Datatype | Pack_size | int size |
| 47 | MPI_PACK | Datatype | Pack | void |

48

| MPI Function | C++ class | Member name | Return value | |
|----------------------|-----------|-------------------|------------------|----|
| MPI_PCONTROL | | Pcontrol | void | 1 |
| MPI_PROBE | Comm | Probe | void | 2 |
| MPI_RECV_INIT | Comm | Recv_init | Prequest request | 3 |
| MPI_RECV | Comm | Recv | void | 4 |
| MPI_REDUCE_SCATTER | Intracomm | Reduce_scatter | void | 5 |
| MPI_REDUCE | Intracomm | Reduce | void | 6 |
| MPI_REQUEST_FREE | Request | Free | void | 7 |
| MPI_RSEND_INIT | Comm | Rsend_init | Prequest request | 8 |
| MPI_RSEND | Comm | Rsend | void | 9 |
| MPI_SCAN | Intracomm | Scan | void | 10 |
| MPI_SCATTERV | Intracomm | Scatterv | void | 11 |
| MPI_SCATTER | Intracomm | Scatter | void | 12 |
| MPI_SENDRECV_REPLACE | Comm | Sendrecv_replace | void | 13 |
| MPI_SENDRECV | Comm | Sendrecv | void | 14 |
| MPI_SEND_INIT | Comm | Send_init | Prequest request | 15 |
| MPI_SEND | Comm | Send | void | 16 |
| MPI_SSEND_INIT | Comm | Ssend_init | Prequest request | 17 |
| MPI_SSEND | Comm | Ssend | void | 18 |
| MPI_STARTALL | Prequest | static Startall | void | 19 |
| MPI_START | Prequest | Start | void | 20 |
| MPI_TESTALL | Request | static Testall | bool flag | 21 |
| MPI_TESTANY | Request | static Testany | bool flag | 22 |
| MPI_TESTSOME | Request | static Testsome | int outcount | 23 |
| MPI_TEST_CANCELLED | Status | Is_cancelled | bool flag | 24 |
| MPI_TEST | Request | Test | bool flag | 25 |
| MPI_TOPO_TEST | Comm | Get_topo | int status | 26 |
| MPI_TYPE_COMMIT | Datatype | Commit | void | 27 |
| MPI_TYPE_CONTIGUOUS | Datatype | Create_contiguous | Datatype | 28 |
| MPI_TYPE_EXTENT | | <none> | | 29 |
| MPI_TYPE_FREE | Datatype | Free | void | 30 |
| MPI_TYPE_HINDEXED | | <none> | | 31 |
| MPI_TYPE_HVECTOR | | <none> | | 32 |
| MPI_TYPE_INDEXED | Datatype | Create_indexed | Datatype | 33 |
| MPI_TYPE_LB | | <none> | | 34 |
| MPI_TYPE_SIZE | Datatype | Get_size | int | 35 |
| MPI_TYPE_STRUCT | | <none> | | 36 |
| MPI_TYPE_UB | | <none> | | 37 |
| MPI_TYPE_VECTOR | Datatype | Create_vector | Datatype | 38 |
| MPI_UNPACK | Datatype | Unpack | void | 39 |
| MPI_WAITALL | Request | static Waitall | void | 40 |
| MPI_WAITANY | Request | static Waitany | int index | 41 |
| MPI_WAITSOME | Request | static Waitsome | int outcount | 42 |
| MPI_WAIT | Request | Wait | void | 43 |
| MPI_WTICK | | Wtick | double wtick | 44 |
| MPI_WTIME | | Wtime | double wtime | 45 |

47

48

Annex B

Change-Log

This annex summarizes changes from the previous version of the MPI standard to the version presented by this document. Only changes (i.e., clarifications and new features) are presented that may cause implementation effort in the MPI libraries. Editorial modifications, formatting, typo corrections and minor clarifications are not shown.

B.1 Changes from Version 2.0 to Version 2.1

1. `MPI_FINALIZE` is collective over all connected processes. If no processes were spawned, accepted or connected then this means over `MPI_COMM_WORLD`; otherwise it is collective over the union of all processes that have been and continue to be connected, as explained in Section 9.5.4 on page 304. (MPI-2.1 Ballot 1, Item 1)
2. In Section 12.5.2 on page 416, the bias of 16 byte doubles was defined with 10383. The correct value is 16383. (MPI-2.1 Ballot 1, Item 7)
3. In the example in Section 13.1.4 on page 438, the buffer should be declared as `const void* buf`. (MPI-2.1 Ballot 1, Item 8)
4. In addition, the `MPI_LONG_LONG` should be added as an optional type; it is a synonym for `MPI_LONG_LONG_INT`. (MPI-2.1 Ballot 1, Item 16)
5. `MPI_GET_COUNT` with zero-length datatypes: The value returned as the count argument of `MPI_GET_COUNT` for a datatype of length zero where zero bytes have been transferred is zero. If the number of bytes transferred is greater than zero, `MPI_UNDEFINED` is returned. (MPI-2.1 Ballot 2, Item 1)
6. `MPI_GROUP_TRANSLATE_RANKS` and `MPI_PROC_NULL`: `MPI_PROC_NULL` is a valid rank for input to `MPI_GROUP_TRANSLATE_RANKS`, which returns `MPI_PROC_NULL` as the translated rank. (MPI-2.1 Ballot 2, Item 2)
7. `MPI_{COMM,WIN,FILE}_GET_ERRHANDLER` behave as if a new error handler object is created. That is, once the error handler is no longer needed, `MPI_ERRHANDLER_FREE` should be called with the error handler returned from `MPI_ERRHANDLER_GET` or `MPI_{COMM,WIN,FILE}_GET_ERRHANDLER` to mark the error handler for deallocation. This provides behavior similar to that of `MPI_COMM_GROUP` and `MPI_GROUP_FREE`. (MPI-2.1 Ballot 2, Item 3)

8. `MPI_BYTE` should be used to send and receive data that is packed using `MPI_PACK_EXTERNAL`. (MPI-2.1 Ballot 2, Item 4)
9. `MPI_PROC_NULL` is a valid target rank in the MPI RMA calls `MPI_ACCUMULATE`, `MPI_GET`, and `MPI_PUT`. The effect is the same as for `MPI_PROC_NULL` in MPI point-to-point communication. See also item 16 in this list. (MPI-2.1 Ballot 2, Item 6)
10. `MPI_REPLACE` in `MPI_ACCUMULATE`, like the other predefined operations, is defined only for the predefined MPI datatypes. (MPI-2.1 Ballot 2, Item 7)
11. If `comm` is an intercommunicator in `MPI_ALLREDUCE`, then both groups should provide `count` and `datatype` arguments that specify the same type signature (i.e., it is not necessary that both groups provide the same count value). (MPI-2.1 Ballot 2, Item 8)
12. About the attribute caching functions:

Advice to implementors. High quality implementations should raise an error when a keyval that was created by a call to `MPI_XXX_CREATE_KEYVAL` is used with an object of the wrong type with a call to `MPI_YYY_GET_ATTR`, `MPI_YYY_SET_ATTR`, `MPI_YYY_DELETE_ATTR`, or `MPI_YYY_FREE_KEYVAL`. To do so, it is necessary to maintain, with each keyval, information on the type of the associated user function. (*End of advice to implementors.*)

(MPI-2.1 Ballot 2, Item 9)
13. If a file does not have the mode `MPI_MODE_SEQUENTIAL`, then `MPI_DISPLACEMENT_CURRENT` is invalid as `disp` in `MPI_FILE_SET_VIEW`. (MPI-2.1 Ballot 2, Item 10)
14. `MPI_BOTTOM` is defined as `void * const MPI::BOTTOM`. (MPI-2.1 Ballot 3, Item 12)
15. An implementation must support info objects as caches for arbitrary (key, value) pairs, regardless of whether it recognizes the key. Each function that takes hints in the form of an `MPI_Info` must be prepared to ignore any key it does not recognize. This description of info objects does not attempt to define how a particular function should react if it recognizes a key but not the associated value. `MPI_INFO_GET_NKEYS`, `MPI_INFO_GET_NTHKEY`, `MPI_INFO_GET_VALUELEN`, and `MPI_INFO_GET` must retain all (key,value) pairs so that layered functionality can also use the Info object. (MPI-2.1 Ballot 4, Item 4)
16. After any RMA operation with rank `MPI_PROC_NULL`, it is still necessary to finish the RMA epoch with the synchronization method that started the epoch. See also item 9 in this list. (MPI-2.1 Ballot 4, Item 6)
17. In `MPI_CART_SUB`: If all entries in `remain_dims` are false or `comm` is already associated with a zero-dimensional Cartesian topology then `newcomm` is associated with a zero-dimensional Cartesian topology. (MPI-2.1 Ballot 4, Item 10.a)

- 1 18. In `MPI_CARTDIM_GET` and `MPI_CART_GET`: If `comm` is associated with a zero-
2 dimensional Cartesian topology, `MPI_CARTDIM_GET` returns `ndims=0` and
3 `MPI_CART_GET` will keep all output arguments unchanged. (MPI-2.1 Ballot 4, Item
4 10.b.i)
- 5 19. In `MPI_CART_RANK`: If `comm` is associated with a zero-dimensional Cartesian topol-
6 ogy, `coord` is not significant and 0 is returned in `rank`. (MPI-2.1 Ballot 4, Item 10.b.ii)
- 7 20. In `MPI_CART_COORDS`: If `comm` is associated with a zero-dimensional Cartesian
8 topology, `coords` will be unchanged. (MPI-2.1 Ballot 4, Item 10.b.iii)
- 9 21. In `MPI_CART_SHIFT`: It is erroneous to call `MPI_CART_SHIFT` with a direction that is
10 either negative or greater than or equal to the number of dimensions in the Cartesian
11 communicator. This implies that it is erroneous to call `MPI_CART_SHIFT` with a
12 `comm` that is associated with a zero-dimensional Cartesian topology. (MPI-2.1 Ballot
13 4, Item 10.c)
- 14 22. In `MPI_CART_CREATE`: If `ndims` is zero then a zero-dimensional Cartesian topology
15 is created. The call is erroneous if it specifies a grid that is larger than the group size
16 or if `ndims` is negative. (MPI-2.1 Ballot 4, Item 10.d)
- 17 23. In `MPI_GRAPH_CREATE`: If the graph is empty, i.e., `nnodes == 0`, then
18 `MPI_COMM_NULL` is returned in all processes. (MPI-2.1 Ballot 4, Item 11.a)
- 19 24. In `MPI_GRAPH_CREATE`: A single process is allowed to be defined multiple times
20 in the list of neighbors of a process (i.e., there may be multiple edges between two
21 processes). A process is also allowed to be a neighbor to itself (i.e., a self loop in the
22 graph). The adjacency matrix is allowed to be non-symmetric.
23
24 *Advice to users.* Performance implications of using multiple edges or a non-
25 symmetric adjacency matrix are not defined. The definition of a node-neighbor
26 edge does not imply a direction of the communication. (*End of advice to users.*)
27
28 (MPI-2.1 Ballot 4, Item 11.b)
- 29 25. About `MPI_GRAPH_CREATE` and `MPI_CART_CREATE`: All input arguments must
30 have identical values on all processes of the group of `comm_old`. (MPI-2.1 Ballot
31 4, Item 11.c)
- 32 26. In `MPI_GET_PROCESSOR_NAME`: In C, a null character is additionally stored at
33 `name[resultlen]`. `resultlen` cannot be larger than `MPI_MAX_PROCESSOR_NAME-1`. In
34 Fortran, `name` is padded on the right with blank characters. `resultlen` cannot be larger
35 than `MPI_MAX_PROCESSOR_NAME`. (MPI-2.1 Ballot 4, Item 13.i)
- 36 27. In `MPI_COMM_GET_NAME`: In C, a null character is additionally stored at
37 `name[resultlen]`. `resultlen` cannot be larger than `MPI_MAX_OBJECT-1`. In Fortran, `name`
38 is padded on the right with blank characters. `resultlen` cannot be larger than
39 `MPI_MAX_OBJECT`. (MPI-2.1 Ballot 4, Item 13.iii)
- 40 28. About `MPI_ABORT`:
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48

Advice to users. Whether the errorcode is returned from the executable or from the MPI process startup mechanism (e.g., mpiexec), is an aspect of quality of the MPI library but not mandatory. (*End of advice to users.*)

Advice to implementors. Where possible, a high quality implementation will try to return the errorcode from the MPI process startup mechanism (e.g. mpiexec or singleton init). (*End of advice to implementors.*)

(MPI-2.1 Ballot 4, Item 15.b)

29. About MPI_TYPE_CREATE_F90_xxxx:

Advice to implementors. An application may often repeat a call to MPI_TYPE_CREATE_F90_xxxx with the same combination of (xxxx,p,r). The application is not allowed to free the returned predefined, unnamed datatype handles. To prevent the creation of a potentially huge amount of handles, the MPI implementation should return the same datatype handle for the same (REAL/COMPLEX/INTEGER,p,r) combination. Checking for the combination (p,r) in the preceding call to MPI_TYPE_CREATE_F90_xxxx and using a hash-table to find formerly generated handles should limit the overhead of finding a previously generated datatype with same combination of (xxxx,p,r). (*End of advice to implementors.*)

(MPI-2.1 Ballot 4, Item 16)

30. About MPI_FILE_GET_INFO: If no hint exists for the file associated with fh, a handle to a newly created info object is returned that contains no key/value pair. (MPI-2.1 Ballot 4, Item 17)

31. About MPI_FILE_SET_VIEW and MPI_FILE_SET_INFO: When an info object that specifies a subset of valid hints is passed to MPI_FILE_SET_VIEW or MPI_FILE_SET_INFO, there will be no effect on previously set or defaulted hints that the info does not specify. (MPI-2.1 Ballot 4, Item 18)

32. General rule about derived datatypes: Most datatype constructors have replication count or block length arguments. Allowed values are nonnegative integers. If the value is zero, no elements are generated in the type map and there is no effect on datatype bounds or extent. (MPI-2.1 Ballot 4, Item 20)

33. MPI_LONG_LONG_INT, MPI_LONG_LONG (as synonym), MPI_UNSIGNED_LONG_LONG, and MPI_WCHAR are moved from optional to official and they are therefore defined for all three language bindings, see also Annex A.1 on page 492. (MPI-2.1 Review Item 15.h')

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Bibliography

- [1] V. Bala and S. Kipnis. Process groups: a mechanism for the coordination of and communication among processes in the Venus collective communication library. Technical report, IBM T. J. Watson Research Center, October 1992. Preprint. [1.1](#)
- [2] V. Bala, S. Kipnis, L. Rudolph, and Marc Snir. Designing efficient, scalable, and portable collective communication libraries. Technical report, IBM T. J. Watson Research Center, October 1992. Preprint. [1.1](#)
- [3] Purushotham V. Bangalore, Nathan E. Doss, and Anthony Skjellum. MPI++: Issues and Features. In *OOO-SKI '94*, page in press, 1994. [5.1](#)
- [4] A. Beguelin, J. Dongarra, A. Geist, R. Manchek, and V. Sunderam. Visualization and debugging in a heterogeneous environment. *IEEE Computer*, 26(6):88–95, June 1993. [1.1](#)
- [5] Luc Bomans and Rolf Hempel. The Argonne/GMD macros in FORTRAN for portable parallel programming and their implementation on the Intel iPSC/2. *Parallel Computing*, 15:119–132, 1990. [1.1](#), [6.2](#)
- [6] Rajesh Bordawekar, Juan Miguel del Rosario, and Alok Choudhary. Design and evaluation of primitives for parallel I/O. In *Proceedings of Supercomputing '93*, pages 452–461, 1993. [12.1](#)
- [7] R. Butler and E. Lusk. User's guide to the p4 programming system. Technical Report TM-ANL-92/17, Argonne National Laboratory, 1992. [1.1](#)
- [8] Ralph Butler and Ewing Lusk. Monitors, messages, and clusters: the p4 parallel programming system. *Parallel Computing, Special issue on message-passing interfaces*, 20:547–564, 4 (April) 1994. (Also Argonne National Laboratory Mathematics and Computer Science Division preprint P362-0493). [1.1](#)
- [9] Robin Calkin, Rolf Hempel, Hans-Christian Hoppe, and Peter Wypior. Portable programming with the parmacs message-passing library. *Parallel Computing, Special issue on message-passing interfaces*, 20:615–632, 4 (April) 1994. [1.1](#), [6.2](#)
- [10] S. Chittor and R. J. Enbody. Performance evaluation of mesh-connected wormhole-routed networks for interprocessor communication in multicomputers. In *Proceedings of the 1990 Supercomputing Conference*, pages 647–656, 1990. [6.1](#)
- [11] S. Chittor and R. J. Enbody. Predicting the effect of mapping on the communication performance of large multicomputers. In *Proceedings of the 1991 International Conference on Parallel Processing, vol. II (Software)*, pages II-1 – II-4, 1991. [6.1](#)

- 1 [12] Parasoftware Corporation. Express version 1.0: A communication environment for parallel
2 computers, 1988. [1.1](#), [6.4](#)
- 3
4 [13] Juan Miguel del Rosario, Rajesh Bordawekar, and Alok Choudhary. Improved parallel
5 I/O via a two-phase run-time access strategy. In *IPPS '93 Workshop on Input/Output*
6 *in Parallel Computer Systems*, pages 56–70, 1993. Also published in *Computer Archi-*
7 *itecture News* 21(5), December 1993, pages 31–38. [12.1](#)
- 8 [14] J. Dongarra, A. Geist, R. Manchek, and V. Sunderam. Integrated PVM framework
9 supports heterogeneous network computing. *Computers in Physics*, 7(2):166–75, April
10 1993. [1.1](#)
- 11
12 [15] J. J. Dongarra, R. Hempel, A. J. G. Hey, and D. W. Walker. A proposal for a user-
13 level, message passing interface in a distributed memory environment. Technical Report
14 TM-12231, Oak Ridge National Laboratory, February 1993. [1.2](#)
- 15 [16] Nathan Doss, William Gropp, Ewing Lusk, and Anthony Skjellum. A model imple-
16 mentation of MPI. Technical report, Argonne National Laboratory, 1993. [1.6](#)
- 17
18 [17] Edinburgh Parallel Computing Centre, University of Edinburgh. *CHIMP Concepts*,
19 June 1991. [1.1](#)
- 20 [18] Edinburgh Parallel Computing Centre, University of Edinburgh. *CHIMP Version 1.0*
21 *Interface*, May 1992. [1.1](#)
- 22
23 [19] Margaret A. Ellis and Bjarne Stroustrup. *The Annotated C++ Reference Manual*.
24 Addison Wesley, 1990.
- 25 [20] D. Feitelson. Communicators: Object-based multiparty interactions for parallel pro-
26 gramming. Technical Report 91-12, Dept. Computer Science, The Hebrew University
27 of Jerusalem, November 1991. [5.1.2](#)
- 28
29 [21] C++ Forum. Working paper for draft proposed international standard for informa-
30 tion systems — programming language C++. Technical report, American National
31 Standards Institute, 1995.
- 32 [22] Message Passing Interface Forum. MPI: A Message-Passing Interface standard. *The In-*
33 *ternational Journal of Supercomputer Applications and High Performance Computing*,
34 8, 1994. [1.3](#)
- 35
36 [23] Message Passing Interface Forum. MPI: A Message-Passing Interface standard (version
37 1.1). Technical report, 1995. <http://www.mpi-forum.org>. [12.4.1](#), [12.6.4](#), [13.1.4](#)
- 38
39 [24] Hubertus Franke, Peter Hochschild, Pratap Pattnaik, and Marc Snir. An efficient
40 implementation of MPI. In *1994 International Conference on Parallel Processing*,
41 1994. [1.6](#)
- 42 [25] Al Geist, Adam Beguelin, Jack Dongarra, Weicheng Jiang, Bob Manchek, and Vaidy
43 Sunderam. *PVM: Parallel Virtual Machine—A User's Guide and Tutorial for Network*
44 *Parallel Computing*. MIT Press, 1994. [9.1](#)
- 45
46 [26] G. A. Geist, M. T. Heath, B. W. Peyton, and P. H. Worley. PICL: A portable in-
47 strumented communications library, C reference manual. Technical Report TM-11130,
48 Oak Ridge National Laboratory, Oak Ridge, TN, July 1990. [1.1](#)

- [27] William D. Gropp and Barry Smith. Chameleon parallel programming tools users manual. Technical Report ANL-93/23, Argonne National Laboratory, March 1993. [1.1](#)
- [28] Michael Hennecke. A Fortran 90 interface to MPI version 1.1. Technical Report Internal Report 63/96, Rechenzentrum, Universität Karlsruhe, D-76128 Karlsruhe, Germany, June 1996. Available via world wide web from http://www.uni-karlsruhe.de/~Michael.Hennecke/Publications/#MPI_F90. [13.2.4](#)
- [29] Institute of Electrical and Electronics Engineers, New York. *IEEE Standard for Binary Floating-Point Arithmetic, ANSI/IEEE Standard 754-1985*, 1985. [12.5.2](#)
- [30] International Organization for Standardization, Geneva. *Information processing — 8-bit single-byte coded graphic character sets — Part 1: Latin alphabet No. 1*, 1987. [12.5.2](#)
- [31] International Organization for Standardization, Geneva. *Information technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language]*, December 1996. [11.7](#), [12.2.1](#)
- [32] Charles H. Koelbel, David B. Loveman, Robert S. Schreiber, Guy L. Steele Jr., and Mary E. Zosel. *The High Performance Fortran Handbook*. MIT Press, 1993. [3.12.5](#)
- [33] David Kotz. Disk-directed I/O for MIMD multiprocessors. In *Proceedings of the 1994 Symposium on Operating Systems Design and Implementation*, pages 61–74, November 1994. Updated as Dartmouth TR PCS-TR94-226 on November 8, 1994. [12.1](#)
- [34] O. Krämer and H. Mühlenbein. Mapping strategies in message-based multiprocessor systems. *Parallel Computing*, 9:213–225, 1989. [6.1](#)
- [35] S. J. Lefflet, R. S. Fabry, W. N. Joy, P. Lapsley, S. Miller, and C. Torek. An advanced 4.4BSD interprocess communication tutorial, Unix programmer’s supplementary documents (PSD) 21. Technical report, Computer Systems Research Group, Department of Electrical Engineering and Computer Science, University of California, Berkeley, 1993. Also available at <http://www.netbsd.org/Documentation/lite2/psd/>. [9.5.5](#)
- [36] nCUBE Corporation. *nCUBE 2 Programmers Guide, r2.0*, December 1990. [1.1](#)
- [37] Bill Nitzberg. Performance of the iPSC/860 Concurrent File System. Technical Report RND-92-020, NAS Systems Division, NASA Ames, December 1992. [12.1](#)
- [38] William J. Nitzberg. *Collective Parallel I/O*. PhD thesis, Department of Computer and Information Science, University of Oregon, December 1995. [12.1](#)
- [39] *4.4BSD Programmer’s Supplementary Documents (PSD)*. O’Reilly and Associates, 1994. [9.5.5](#)
- [40] Paul Pierce. The NX/2 operating system. In *Proceedings of the Third Conference on Hypercube Concurrent Computers and Applications*, pages 384–390. ACM Press, 1988. [1.1](#)

- 1 [41] K. E. Seamons, Y. Chen, P. Jones, J. Jozwiak, and M. Winslett. Server-directed
2 collective I/O in Panda. In *Proceedings of Supercomputing '95*, December 1995. [12.1](#)
3
- 4 [42] A. Skjellum and A. Leung. Zipcode: a portable multicomputer communication library
5 atop the reactive kernel. In D. W. Walker and Q. F. Stout, editors, *Proceedings of the*
6 *Fifth Distributed Memory Concurrent Computing Conference*, pages 767–776. IEEE
7 Press, 1990. [1.1](#), [5.1.2](#)
- 8 [43] A. Skjellum, S. Smith, C. Still, A. Leung, and M. Morari. The Zipcode message passing
9 system. Technical report, Lawrence Livermore National Laboratory, September 1992.
10 [1.1](#)
11
- 12 [44] Anthony Skjellum, Nathan E. Doss, and Purushotham V. Bangalore. Writing Libraries
13 in MPI. In Anthony Skjellum and Donna S. Reese, editors, *Proceedings of the Scalable*
14 *Parallel Libraries Conference*, pages 166–173. IEEE Computer Society Press, October
15 1993. [5.1](#)
16
- 17 [45] Anthony Skjellum, Nathan E. Doss, and Kishore Viswanathan. Inter-communicator
18 extensions to MPI in the MPIX (MPI eXtension) Library. Technical Report MSU-
19 940722, Mississippi State University — Dept. of Computer Science, April 1994.
20 <http://www.erc.msstate.edu/mpi/mpix.html>. [4.3.2](#), [5.4.2](#)
- 21 [46] Anthony Skjellum, Ziyang Lu, Purushotham V. Bangalore, and Nathan E. Doss. Ex-
22 plicit parallel programming in C++ based on the message-passing interface (MPI). In
23 Gregory V. Wilson, editor, *Parallel Programming Using C++*, Engineering Computa-
24 tion Series. MIT Press, July 1996. ISBN 0-262-73118-5.
25
- 26 [47] Anthony Skjellum, Steven G. Smith, Nathan E. Doss, Alvin P. Leung, and Manfred
27 Morari. The Design and Evolution of Zipcode. *Parallel Computing*, 20(4):565–596,
28 April 1994. [5.1.2](#)
29
- 30 [48] Anthony Skjellum, Steven G. Smith, Nathan E. Doss, Charles H. Still, Alvin P. Le-
31 ung, and Manfred Morari. Zipcode: A Portable Communication Layer for High Per-
32 formance Multicomputing. Technical Report UCRL-JC-106725 (revised 9/92, 12/93,
33 4/94), Lawrence Livermore National Laboratory, March 1991. To appear in *Concur-*
34 *rency: Practice & Experience*. [5.5.6](#)
- 35 [49] Rajeev Thakur and Alok Choudhary. An Extended Two-Phase Method for Accessing
36 Sections of Out-of-Core Arrays. *Scientific Programming*, 5(4):301–317, Winter 1996.
37 [12.1](#)
38
- 39 [50] *The Unicode Standard, Version 2.0*. Addison-Wesley, 1996. ISBN 0-201-48345-9. [12.5.2](#)
40
- 41 [51] D. Walker. Standards for message passing in a distributed memory environment. Tech-
42 nical Report TM-12147, Oak Ridge National Laboratory, August 1992. [1.2](#)
43
44
45
46
47
48

Index

- ACTION_SUBSET, 10
- Action_subset, 10

- cancel_fn, 346, 347
- Class, 10
- CLASS_ACTION, 10
- Class_action, 9
- CLASS_ACTION_SUBSET, 9, 10
- Class_action_subset, 9
- COMM_COPY_ATTR_FN, 17
- comm_copy_attr_fn, 214, 215
- COMM_DELETE_ATTR_FN, 17
- comm_delete_attr_fn, 215, 216
- CONST:, 151, 152
- CONST: MPL_ANY_SOURCE, 64
- CONST: MPL_ANY_TAG, 30, 64, 74
- CONST: MPL_COMM_WORLD, 194
- CONST: MPI_LB, 99
- CONST: MPI_PROC_NULL, 74
- CONST: MPI_UB, 99
- CONST: MPI_UNDEFINED, 176
- CONST:&, 93
- CONST:_WORLD, 252
- CONST:a, 36
- CONST:access_style, 385
- CONST:appnum, 304
- CONST:arch, 289
- CONST:argc, 263
- CONST:argv, 263
- CONST:b, 36
- CONST:Byte:, 152
- CONST:C integer, Fortran integer, Byte, 152
- CONST:C integer, Fortran integer, Floating point, 152
- CONST:C integer, Fortran integer, Floating point, Complex, 152
- CONST:C integer, Logical, 152
- CONST:C integer:, 151
- CONST:C:, 154
- CONST:cb_block_size, 385, 386
- CONST:cb_buffer_size, 385, 386
- CONST:cb_nodes, 385, 386
- CONST:char**, 46
- CONST:CHARACTER, 35, 36
- CONST:CHARACTER*, 259
- CONST:chunked, 386
- CONST:chunked_item, 386
- CONST:chunked_size, 386
- CONST:collective_buffering, 385
- CONST:Complex:, 151
- CONST:DIMS, 242
- CONST:dims[i], 242
- CONST:DIRECTION = i, 242
- CONST:direction = i, 242
- CONST:errorcode, 259, 261
- CONST:external32, 414
- CONST:false, 52, 201, 232, 234, 238, 385
- CONST:file, 290
- CONST:file_perm, 385, 386
- CONST:filename, 386
- CONST:flag = 0, 214, 488
- CONST:flag = 1, 214, 488
- CONST:Floating point:, 151
- CONST:Fortran integer:, 151
- CONST:Fortran:, 154
- CONST:host, 289
- CONST:IERROR, 263
- CONST:int, 93
- CONST:INTEGER, 212
- CONST:internal, 414
- CONST:io_node_list, 386
- CONST:ip_address, 299
- CONST:ip_port, 299
- CONST:LASTCODE, 259
- CONST:Logical:, 151
- CONST:MPI::ERRORS_THROW_EXCEPTIONS, 19, 23
- CONST:MPI_2DOUBLE_PRECISION, 154, 155

- 1 CONST:MPL2INT, 155
2 CONST:MPL2INTEGER, 154, 155
3 CONST:MPL2REAL, 154, 155
4 CONST:MPLADDRESS_KIND, 18, 360,
5 474
6 CONST:MPLAint, 18
7 CONST:MPLANY_SOURCE, 30, 41, 64,
8 66, 226, 249
9 CONST:MPLANY_TAG, 14, 30, 32, 64,
10 66
11 CONST:MPLAPPNUM, 303, 304
12 CONST:MPLARGV_NULL, 284, 285, 450
13 CONST:MPLARGVS_NULL, 288, 450
14 CONST:MPLBAND, 151, 152
15 CONST:MPLBOR, 151, 152
16 CONST:MPLBOTTOM, 17, 92, 102, 561
17 CONST:MPLBOTTOM + v, 102
18 CONST:MPLBSEND_OVERHEAD, 47, 250
19 CONST:MPLBXOR, 151, 152
20 CONST:MPLBYTE, 35, 116, 561
21 CONST:MPLCART, 236
22 CONST:MPLCHARACTER, 35
23 CONST:MPLCOMBINER_CONTIGUOUS,
24 360, 363
25 CONST:MPLCOMBINER_DARRAY, 360,
26 365
27 CONST:MPLCOMBINER_DUP, 360, 363
28 CONST:MPLCOMBINER_F90_COMPLEX,
29 360, 365
30 CONST:MPLCOMBINER_F90_INTEGER,
31 360, 365
32 CONST:MPLCOMBINER_F90_REAL, 360,
33 365
34 CONST:MPLCOMBINER_HINDEXED, 360,
35 364
36 CONST:MPLCOMBINER_HINDEXED_INTEGER,
37 360, 364
38 CONST:MPLCOMBINER_HVECTOR, 360,
39 364
40 CONST:MPLCOMBINER_HVECTOR_INTEGER,
41 360, 364
42 CONST:MPLCOMBINER_INDEXED, 360,
43 364
44 CONST:MPLCOMBINER_INDEXED_BLOCK,
45 360, 364
46 CONST:MPLCOMBINER_NAMED, 360,
47 363
48 CONST:MPLCOMBINER_RESIZED, 360,
 365
 CONST:MPLCOMBINER_STRUCT, 360,
 364
 CONST:MPLCOMBINER_STRUCT_INTEGER,
 360, 364
 CONST:MPLCOMBINER_SUBARRAY, 360,
 364
 CONST:MPLCOMBINER_VECTOR, 360,
 363
 CONST:MPLCOMM-, 252
 CONST:MPLComm_copy_attr_function, 214
 CONST:MPLComm_delete_attr_function, 215
 CONST:MPLCOMM_NULL, 175, 186, 189,
 192, 232, 234, 353, 562
 CONST:MPLCOMM_SELF, 175, 212
 CONST:MPLCOMM_WORLD, 29, 175–
 177, 183, 185, 232, 248, 249, 252,
 268
 CONST:MPLCONGRUENT, 184, 201
 CONST:MPLCONVERSION_FN_NULL, 421
 CONST:MPLCopy_function, 488
 CONST:MPLDATATYPE_NULL, 97
 CONST:MPLDelete_function, 489
 CONST:MPLDISPLACEMENT_CURRENT,
 387, 561
 CONST:MPLDISTRIBUTE_BLOCK, 89
 CONST:MPLDISTRIBUTE_CYCLIC, 89
 CONST:MPLDISTRIBUTE_DFLT_DARG,
 89
 CONST:MPLDISTRIBUTE_NONE, 89
 CONST:MPLDOUBLE_INT, 155
 CONST:MPLERR..., 259
 CONST:MPLERR_ACCESS, 261, 380, 433
 CONST:MPLERR_AMODE, 261, 379, 433
 CONST:MPLERR_ARG, 260
 CONST:MPLERR_ASSERT, 260, 335
 CONST:MPLERR_BAD_FILE, 261, 433
 CONST:MPLERR_BASE, 251, 260, 335
 CONST:MPLERR_BUFFER, 260
 CONST:MPLERR_COMM, 260
 CONST:MPLERR_CONVERSION, 261, 421,
 433
 CONST:MPLERR_COUNT, 260
 CONST:MPLERR_DIMS, 260
 CONST:MPLERR_DISP, 260, 335
 CONST:MPLERR_DUP_DATAREP, 261,
 419, 433

- CONST:MPIERR_FILE, [261](#), [433](#)
- CONST:MPIERR_FILE_EXISTS, [261](#), [433](#)
- CONST:MPIERR_FILE_IN_USE, [261](#), [380](#), [433](#)
- CONST:MPIERR_GROUP, [260](#)
- CONST:MPIERR_IN_STATUS, [31](#), [52](#), [59](#), [60](#), [254](#), [260](#)
- CONST:MPIERR_INFO, [260](#)
- CONST:MPIERR_INFO_KEY, [260](#), [274](#)
- CONST:MPIERR_INFO_NOKEY, [260](#), [275](#)
- CONST:MPIERR_INFO_VALUE, [260](#), [274](#)
- CONST:MPIERR_INTERN, [260](#)
- CONST:MPIERR_IO, [261](#), [433](#)
- CONST:MPIERR_KEYVAL, [223](#), [260](#)
- CONST:MPIERR_LASTCODE, [259](#), [261](#)
- CONST:MPIERR_LOCKTYPE, [260](#), [335](#)
- CONST:MPIERR_NAME, [260](#), [299](#)
- CONST:MPIERR_NO_MEM, [251](#), [260](#)
- CONST:MPIERR_NO_SPACE, [261](#), [433](#)
- CONST:MPIERR_NO_SUCH_FILE, [261](#), [380](#), [433](#)
- CONST:MPIERR_NOT_SAME, [261](#), [433](#)
- CONST:MPIERR_OP, [260](#)
- CONST:MPIERR_OTHER, [259](#), [260](#)
- CONST:MPIERR_PENDING, [59](#), [260](#)
- CONST:MPIERR_PORT, [260](#), [296](#)
- CONST:MPIERR_QUOTA, [261](#), [433](#)
- CONST:MPIERR_RANK, [260](#)
- CONST:MPIERR_READ_ONLY, [261](#), [433](#)
- CONST:MPIERR_REQUEST, [260](#)
- CONST:MPIERR_RMA_CONFLICT, [260](#), [335](#)
- CONST:MPIERR_RMA_SYNC, [260](#), [335](#)
- CONST:MPIERR_ROOT, [260](#)
- CONST:MPIERR_SERVICE, [260](#), [298](#)
- CONST:MPIERR_SIZE, [260](#), [335](#)
- CONST:MPIERR_SPAWN, [260](#), [285](#), [286](#)
- CONST:MPIERR_TAG, [260](#)
- CONST:MPIERR_TOPOLOGY, [260](#)
- CONST:MPIERR_TRUNCATE, [260](#)
- CONST:MPIERR_TYPE, [260](#)
- CONST:MPIERR_UNKNOWN, [259](#), [260](#)
- CONST:MPIERR_UNSUPPORTED_DATAREP, [261](#), [433](#)
- CONST:MPIERR_UNSUPPORTED_OPERATION, [261](#), [433](#)
- CONST:MPIERR_WIN, [260](#), [335](#)
- CONST:MPIERRCODES_IGNORE, [17](#), [286](#), [450](#), [453](#)
- CONST:MPIERRHANDLER_NULL, [258](#)
- CONST:MPIERROR, [31](#), [52](#)
- CONST:MPIERROR_STRING, [259](#)
- CONST:MPIERRORS_ARE_FATAL, [252](#)
- CONST:MPIERRORS_RETURN, [252](#), [253](#)
- CONST:MPIF_STATUS_IGNORE, [471](#)
- CONST:MPIF_STATUSES_IGNORE, [471](#)
- CONST:MPIFILE_NULL, [380](#), [432](#)
- CONST:MPIFLOAT_INT, [154](#), [155](#)
- CONST:MPIGRAPH, [236](#)
- CONST:MPIGROUP_EMPTY, [174](#), [179](#), [180](#)
- CONST:MPIGROUP_NULL, [174](#), [183](#)
- CONST:MPIHandler_function, [491](#)
- CONST:MPIHOST, [248](#)
- CONST:MPIIDENT, [177](#), [184](#)
- CONST:MPIIN_PLACE, [125](#), [450](#), [453](#), [457](#)
- CONST:MPIINFO_NULL, [277](#), [285](#), [295](#), [379](#), [380](#), [388](#)
- CONST:MPIINTEGER_KIND, [360](#), [474](#)
- CONST:MPIIO, [248](#), [249](#)
- CONST:MPIKEYVAL_INVALID, [215](#)–[217](#)
- CONST:MPILAND, [151](#), [152](#)
- CONST:MPILASTUSED_CODE, [356](#)
- CONST:MPILOCK_EXCLUSIVE, [328](#)
- CONST:MPILOCK_SHARED, [328](#)
- CONST:MPILONG_DOUBLE_INT, [155](#)
- CONST:MPILONG_INT, [155](#)
- CONST:MPI_LOR, [151](#), [152](#)
- CONST:MPI_LXOR, [151](#), [152](#)
- CONST:MPI_MAX, [150](#)–[152](#)
- CONST:MPI_MAX_DATAREP_STRING, [389](#), [419](#)
- CONST:MPI_MAX_ERROR_STRING, [259](#)
- CONST:MPI_MAX_INFO_KEY, [260](#), [273](#), [275](#), [276](#)
- CONST:MPI_MAX_INFO_VAL, [260](#), [273](#)
- CONST:MPI_MAX_OBJECT, [353](#), [562](#)
- CONST:MPI_MAX_OBJECT_NAME, [352](#)
- CONST:MPI_MAX_PORT_NAME, [294](#)
- CONST:MPI_MAX_PROCESSOR_NAME, [250](#), [562](#)
- CONST:MPI_MAXLOC, [151](#), [153](#), [154](#), [157](#)
- CONST:MPI_MIN, [151](#), [152](#)
- CONST:MPI_MINLOC, [151](#), [153](#), [154](#), [157](#)

- 1 CONST:MPLMODE_APPEND, 378, 379
2 CONST:MPLMODE_CREATE, 378, 379,
3 386
4 CONST:MPLMODE_DELETE_ON_CLOSE,
5 378–380
6 CONST:MPLMODE_EXCL, 378, 379
7 CONST:MPLMODE_NOCHECK, 331, 332
8 CONST:MPLMODE_NOPRECEDE, 331
9 CONST:MPLMODE_NOPUT, 331
10 CONST:MPLMODE_NOSTORE, 331
11 CONST:MPLMODE_NOSUCCEED, 331
12 CONST:MPLMODE_RDONLY, 378, 379,
13 384
14 CONST:MPLMODE_RDWR, 378, 379
15 CONST:MPLMODE_SEQUENTIAL, 378,
16 379, 381, 382, 387, 392, 396, 405,
17 425, 561
18 CONST:MPLMODE_UNIQUE_OPEN, 378,
19 379
20 CONST:MPLMODE_WRONGLY, 378, 379
21 CONST:MPLOFFSET_KIND, 18, 426
22 CONST:MPL_OP_NULL, 160
23 CONST:MPL_ORDER_C, 14, 86, 89, 90
24 CONST:MPL_ORDER_FORTRAN, 14, 86,
25 89
26 CONST:MPL_PACKED, 117
27 CONST:MPL_PROC_NULL, 177, 241, 248,
28 249, 312, 560, 561
29 CONST:MPL_PROD, 151, 152
30 CONST:MPL_REPLACE, 318, 561
31 CONST:MPL_REQUEST_NULL, 52–54, 57–
32 60
33 CONST:MPL_ROOT, 125
34 CONST:MPL_SEEK_CUR, 400, 406
35 CONST:MPL_SEEK_END, 400, 406
36 CONST:MPL_SEEK_SET, 400, 401, 406
37 CONST:MPL_SHORT_INT, 155
38 CONST:MPL_SIMILAR, 177, 184, 201
39 CONST:MPL_SOURCE, 31
40 CONST:MPL_STATUS_IGNORE, 10, 17,
41 33, 345, 392, 450, 453, 457, 471,
42 477
43 CONST:MPL_STATUS_SIZE, 31
44 CONST:MPL_STATUSES_IGNORE, 14, 33,
45 345, 347, 450, 453, 471
46 CONST:MPL_SUBVERSION, 248
47 CONST:MPL_SUCCESS, 59, 60, 214–217,
48 259, 260, 488, 489
- CONST:MPLSUM, 151, 152
CONST:MPLTAG, 31
CONST:MPLTAG_UB, 28, 248
CONST:MPLTEST, 52
CONST:MPLTESTANY, 52
CONST:MPLTHREAD_FUNNELED, 370,
371
CONST:MPLTHREAD_MULTIPLE, 370,
371, 373
CONST:MPLTHREAD_SERIALIZED, 370,
371
CONST:MPLTHREAD_SINGLE, 370, 371
CONST:MPLTYPECLASS_COMPLEX, 464
CONST:MPLTYPECLASS_INTEGER, 464
CONST:MPLTYPECLASS_REAL, 464
CONST:MPL_UNDEFINED, 32, 57, 58, 100,
177, 190, 236, 243, 244, 560
CONST:MPLUNEQUAL, 177, 184, 201
CONST:MPLUNIVERSE_SIZE, 282, 302
CONST:MPL_VERSION, 248
CONST:MPL_WAIT, 52
CONST:MPL_WAITANY, 52
CONST:MPL_WCHAR, 416
CONST:MPL_WIN_BASE, 310
CONST:MPL_WIN_DISP_UNIT, 310
CONST:MPL_WIN_NULL, 310
CONST:MPL_WIN_SIZE, 310
CONST:MPL_WTIME_IS_GLOBAL, 248,
249, 262
CONST:Name, 151, 154
CONST:native, 413
CONST:nb_proc, 386
CONST:no_locks, 309
CONST:num_io_nodes, 386
CONST:Op, 152
CONST:path, 290
CONST:PRINT, 259
CONST:random, 385
CONST:read_mostly, 385
CONST:read_once, 385
CONST:reverse_sequential, 385
CONST:sequential, 385
CONST:soft, 285, 290
CONST:string, 259
CONST:striping_factor, 386
CONST:striping_unit, 386
CONST:true, 57, 201, 232, 234, 238, 385
CONST:void *, 212, 215

- CONST:void*, 46, 217
- CONST:void**, 46, 217
- CONST:wdir, 290
- CONST:write_mostly, 385
- CONST:write_once, 385
- copy_fn, 488
- COPY_FUNCTION, 17

- delete_fn, 488, 489
- DELETE_FUNCTION, 17

- EXAMPLES:Client-server code, 62
- EXAMPLES:Client-server code with blocking probe, 65
- EXAMPLES:Client-server code with blocking probe, wrong, 65
- EXAMPLES:Datatype - 3D array, 102
- EXAMPLES:Datatype - absolute addresses, 108
- EXAMPLES:Datatype - array of structures, 105
- EXAMPLES:Datatype - elaborate example, 114, 115
- EXAMPLES:Datatype - matching type, 99
- EXAMPLES:Datatype - matrix transpose, 104
- EXAMPLES:Datatype - union, 109
- EXAMPLES:Datatypes - matching, 34
- EXAMPLES:Datatypes - not matching, 35
- EXAMPLES:Datatypes - untyped, 35
- EXAMPLES:Deadlock with MPI_Bcast, 167, 168
- EXAMPLES:Deadlock, if not buffered, 44
- EXAMPLES:Deadlock, wrong message exchange, 43
- EXAMPLES:Intercommunicator, 187, 190
- EXAMPLES:Intertwined matching pairs, 42
- EXAMPLES:Message exchange, 43
- EXAMPLES:MPI::Comm::Probe, 33
- EXAMPLES:MPI_ACCUMULATE, 318
- EXAMPLES:MPI_ADDRESS, 93
- EXAMPLES:MPI_Address, 105, 108, 109, 114
- EXAMPLES:MPI_Allgather, 144
- EXAMPLES:MPI_ALLOC_MEM, 251
- EXAMPLES:MPI_Alloc_mem, 252
- EXAMPLES:MPI_ALLREDUCE, 162
- EXAMPLES:MPI_Barrier, 265, 333
- EXAMPLES:MPI_Bcast, 127, 167, 168
- EXAMPLES:MPI_BSEND, 42
- EXAMPLES:MPI_Buffer_attach, 45, 265
- EXAMPLES:MPI_Buffer_detach, 45
- EXAMPLES:MPI_BYTE, 35
- EXAMPLES:MPI_Cancel, 265
- EXAMPLES:MPI_CART_COORDS, 241
- EXAMPLES:MPI_CART_GET, 244
- EXAMPLES:MPI_CART_RANK, 244
- EXAMPLES:MPI_CART_SHIFT, 241
- EXAMPLES:MPI_CART_SUB, 242
- EXAMPLES:MPI_CHARACTER, 35
- EXAMPLES:MPI_Comm_create, 187
- EXAMPLES:MPI_Comm_group, 187
- EXAMPLES:MPI_Comm_remote_size, 190
- EXAMPLES:MPI_COMM_SPAWN, 284
- EXAMPLES:MPI_Comm_spawn, 284
- EXAMPLES:MPI_COMM_SPAWN_MULTIPLE, 289
- EXAMPLES:MPI_Comm_spawn_multiple, 289
- EXAMPLES:MPI_Comm_split, 190
- EXAMPLES:MPI_DIMS_CREATE, 233, 244
- EXAMPLES:MPI_FILE_CLOSE, 396, 399
- EXAMPLES:MPI_FILE_GET_AMODE, 383
- EXAMPLES:MPI_FILE_IREAD, 399
- EXAMPLES:MPI_FILE_OPEN, 396, 399
- EXAMPLES:MPI_FILE_READ, 396
- EXAMPLES:MPI_FILE_SET_ATOMICITY, 427
- EXAMPLES:MPI_FILE_SET_VIEW, 396, 399
- EXAMPLES:MPI_FILE_SYNC, 428
- EXAMPLES:MPI_Finalize, 265, 266
- EXAMPLES:MPI_FREE_MEM, 251
- EXAMPLES:MPI_Gather, 115, 130, 131, 135
- EXAMPLES:MPI_Gatherv, 115, 132–135
- EXAMPLES:MPI_GET, 314, 316
- EXAMPLES:MPI_Get, 332, 333
- EXAMPLES:MPI_GET_ADDRESS, 93
- EXAMPLES:MPI_Get_address, 105, 108, 109, 114
- EXAMPLES:MPI_GET_COUNT, 101
- EXAMPLES:MPI_GET_ELEMENT, 101
- EXAMPLES:MPI_GRAPH_CREATE, 234, 240

- 1 EXAMPLES:MPI_Grequest_complete, 347
 2 EXAMPLES:MPI_Grequest_start, 347
 3 EXAMPLES:MPI_Group_free, 187
 4 EXAMPLES:MPI_Group_incl, 187
 5 EXAMPLES:MPI_Iprobe, 265
 6 EXAMPLES:MPI_Irecv, 54–56, 62
 7 EXAMPLES:MPI_Isend, 54, 55, 62
 8 EXAMPLES:MPI_Op_create, 160, 166
 9 EXAMPLES:MPI_Pack, 113–115
 10 EXAMPLES:MPI_Pack_size, 115
 11 EXAMPLES:MPI_PROBE, 65
 12 EXAMPLES:MPI_Put, 325, 330, 332, 333
 13 EXAMPLES:MPI_Recv, 34, 35, 42–44,
 14 56, 65, 99
 15 EXAMPLES:MPI_Reduce, 152, 156
 16 EXAMPLES:MPI_Reduce, 155, 156, 160
 17 EXAMPLES:MPI_REQUEST_FREE, 55
 18 EXAMPLES:MPI_Request_free, 265
 19 EXAMPLES:MPI_Scan, 166
 20 EXAMPLES:MPI_Scatter, 139
 21 EXAMPLES:MPI_Scatterv, 140
 22 EXAMPLES:MPI_Send, 34, 35, 43, 44,
 23 56, 65, 99
 24 EXAMPLES:MPI_Send, 105, 108, 109, 114
 25 EXAMPLES:MPI_SENDRECV, 102–104
 26 EXAMPLES:MPI_SENDRECV_REPLACE,
 27 241
 28 EXAMPLES:MPI_Ssend, 42, 56
 29 EXAMPLES:MPI_Test_cancelled, 265
 30 EXAMPLES:MPI_TYPE_COMMIT, 98, 102–
 31 104, 314
 32 EXAMPLES:MPI_Type_commit, 105, 108,
 33 109, 114, 131–135, 140, 166
 34 EXAMPLES:MPI_TYPE_CONTIGUOUS,
 35 77, 94, 99, 101
 36 EXAMPLES:MPI_Type_contiguous, 131
 37 EXAMPLES:MPI_TYPE_CREATE_DARRAY,
 38 91
 39 EXAMPLES:MPI_TYPE_CREATE_HVECTOR,
 40 102, 104
 41 EXAMPLES:MPI_Type_create_hvector, 105,
 42 108
 43 EXAMPLES:MPI_TYPE_CREATE_STRUCT,
 44 84, 94, 104
 45 EXAMPLES:MPI_Type_create_struct, 105,
 46 108, 109, 114, 133, 135, 166
 47 EXAMPLES:MPI_TYPE_CREATE_SUBARRAY,
 48 435
- EXAMPLES:MPI_TYPE_EXTENT, 102,
 104, 314, 316, 318
 EXAMPLES:MPI_Type_extent, 105
 EXAMPLES:MPI_TYPE_FREE, 314
 EXAMPLES:MPI_Type_get_contents, 365
 EXAMPLES:MPI_Type_get_envelope, 365
 EXAMPLES:MPI_TYPE_HVECTOR, 102,
 104
 EXAMPLES:MPI_Type_hvector, 105, 108
 EXAMPLES:MPI_TYPE_INDEXED, 81,
 103
 EXAMPLES:MPI_Type_indexed, 105, 108
 EXAMPLES:MPI_TYPE_INDEXED_BLOCK,
 314
 EXAMPLES:MPI_TYPE_STRUCT, 84, 94,
 104
 EXAMPLES:MPI_Type_struct, 105, 108,
 109, 114, 133, 135, 166
 EXAMPLES:MPI_TYPE_VECTOR, 78, 102,
 104
 EXAMPLES:MPI_Type_vector, 132, 134,
 140
 EXAMPLES:MPI_Unpack, 114, 115
 EXAMPLES:MPI_WAIT, 54–56, 62, 399
 EXAMPLES:MPI_WAITANY, 62
 EXAMPLES:MPI_WAIT SOME, 62
 EXAMPLES:MPI_Win_complete, 325, 333
 EXAMPLES:MPI_WIN_CREATE, 314, 316,
 318
 EXAMPLES:MPI_WIN_FENCE, 314, 316,
 318
 EXAMPLES:MPI_Win_fence, 332
 EXAMPLES:MPI_Win_lock, 330
 EXAMPLES:MPI_Win_post, 333
 EXAMPLES:MPI_Win_start, 325, 333
 EXAMPLES:MPI_Win_unlock, 330
 EXAMPLES:MPI_Win_wait, 333
 EXAMPLES:mpiexec, 270, 271
 EXAMPLES:Non-deterministic program with
 MPI_Bcast, 168
 EXAMPLES:Non-overtaking messages, 42
 EXAMPLES:Nonblocking operations, 54,
 55
 EXAMPLES:Nonblocking operations - mes-
 sage ordering, 55
 EXAMPLES:Nonblocking operations - progress,
 56
 EXAMPLES:Threads and MPI, 368

- EXAMPLES:Typemap, [76–78](#), [81](#), [84](#), [91](#)
- free_fn, [345–347](#)
- HANDLER_FUNCTION, [254](#)
- MPI::COMM_NULL, [441](#)
- MPI::ERRORS_ARE_FATAL, [19](#)
- MPI::ERRORS_RETURN, [19](#)
- MPI::ERRORS_THROW_EXCEPTIONS, [19](#), [253](#)
- MPI::Info, [273](#), [285](#)
- MPI::Is_finalized, [20](#)
- MPI_, [479](#)
- MPI*_FREE, [440](#)
- MPI_{TYPE,COMM,WIN}_CREATE_KEYVAL, [474](#)
- MPI_{WAIT,\$TEST}{SOME}346
- MPIABORT, [158](#), [252](#), [264](#), [267](#), [304](#), [467](#), [562](#)
- MPIACCUMULATE, [307](#), [311](#), [312](#), [317](#), [318](#), [319](#), [338](#), [561](#)
- MPIADD_ERROR_CLASS, [355](#), [355](#)
- MPIADD_ERROR_CODE, [356](#)
- MPIADD_ERROR_STRING, [356](#), [357](#)
- MPIADDRESS, [17](#), [92](#), [474](#), [487](#)
- MPIADDRESS_KIND, [475](#)
- MPIALLGATHER, [123](#), [142](#), [142](#), [143–145](#)
- MPIAllgather, [122](#), [123](#)
- MPIALLGATHERV, [123](#), [143](#), [144](#)
- MPIAllgatherv, [122](#)
- MPIALLOC_MEM, [18](#), [250](#), [251](#), [260](#), [309](#), [313](#), [329](#), [450](#)
- MPIAlloc_mem, [251](#)
- MPIALLREDUCE, [123](#), [125](#), [151](#), [158](#), [161](#), [561](#)
- MPIAllreduce, [122](#)
- MPIALLTOALL, [123](#), [145](#), [145](#), [146](#), [147](#)
- MPIAlltoall, [122](#)
- MPIALLTOALLV, [123](#), [125](#), [146](#), [146](#), [147](#)
- MPIAlltoallv, [122](#)
- MPIALLTOALLW, [123](#), [147](#), [148](#)
- MPIATTR_DELETE, [17](#), [215](#), [217](#), [223](#), [489](#), [490](#)
- MPIATTR_GET, [17](#), [217](#), [223](#), [248](#), [475](#), [490](#)
- MPIATTR_PUT, [17](#), [216](#), [223](#), [475](#), [489](#)
- MPIBARRIER, [123](#), [126](#), [126](#), [428](#)
- MPIBarrier, [123](#), [198](#)
- MPIBCAST, [123](#), [126](#), [126](#), [127](#), [439](#)
- MPIBcast, [123](#)
- MPIBOTTOM, [455](#)
- MPIBSEND, [39](#), [47](#), [250](#), [265](#)
- MPIBSEND_INIT, [69](#), [71](#)
- MPIBUFFER_ATTACH, [21](#), [45](#), [53](#)
- MPIBUFFER_DETACH, [45](#), [265](#)
- MPIBuffer_detach, [46](#)
- MPICANCEL, [41](#), [53](#), [63](#), [66](#), [66](#), [67](#), [343](#), [346](#), [347](#)
- MPICancel, [265](#)
- MPI_CART_COORDS, [231](#), [239](#), [239](#), [562](#)
- MPI_CART_CREATE, [231](#), [232](#), [232](#), [233](#), [234](#), [237](#), [242](#), [243](#), [562](#)
- MPI_CART_GET, [231](#), [237](#), [238](#), [562](#)
- MPI_CART_MAP, [231](#), [243](#), [243](#)
- MPI_CART_RANK, [231](#), [238](#), [238](#), [562](#)
- MPI_CART_SHIFT, [231](#), [240](#), [241](#), [241](#), [562](#)
- MPI_CART_SUB, [231](#), [242](#), [242](#), [243](#), [561](#)
- MPI_CARTDIM_GET, [231](#), [237](#), [237](#), [562](#)
- MPICHAR, [153](#)
- MPICHARACTER, [153](#)
- MPICLOSE_PORT, [294](#), [295](#), [297](#)
- MPICOMM_ACCEPT, [293–295](#), [295](#), [296](#), [302](#), [304](#)
- MPICOMM_C2F, [468](#)
- MPIComm_c2f, [468](#)
- MPICOMM_CALL_ERRHANDLER, [357](#), [358](#)
- MPICOMM_CLONE, [445](#)
- MPICOMM_COMPARE, [184](#), [201](#)
- MPICOMM_CONNECT, [260](#), [296](#), [296](#), [303](#), [304](#)
- MPIComm_copy_attr_function, [17](#)
- MPICOMM_CREATE, [183](#), [185](#), [186](#), [187–189](#), [200](#), [231](#)
- MPICOMM_CREATE_ERRHANDLER, [17](#), [253](#), [254](#), [255](#), [490](#)
- MPICOMM_CREATE_KEYVAL, [17](#), [212](#), [213](#), [223](#), [474](#), [488](#), [561](#)
- MPICOMM_DELETE_ATTR, [17](#), [215](#), [216](#), [217](#), [223](#), [490](#)
- MPIComm_delete_attr_function, [17](#)
- MPICOMM_DISCONNECT, [223](#), [287](#), [304](#), [305](#), [305](#)

- 1 MPI_COMM_DUP, [122](#), [178](#), [183](#), [185](#), [185](#),
2 [186](#), [187](#), [192](#), [202](#), [204](#), [212](#), [214](#),
3 [217](#), [223](#), [226](#), [488](#)
4 MPI_COMM_DUP_FN, [17](#), [214](#), [215](#)
5 MPI_Comm_errhandler_fn, [17](#)
6 MPI_COMM_F2C, [468](#)
7 MPI_COMM_FREE, [183](#), [186](#), [191](#), [192](#),
8 [202](#), [204](#), [215–217](#), [223](#), [268](#), [287](#),
9 [304](#), [305](#), [441](#), [488](#)
10 MPI_COMM_FREE_KEYVAL, [17](#), [215](#), [223](#),
11 [489](#)
12 MPI_COMM_GET_ATTR, [17](#), [216](#), [216](#),
13 [223](#), [476](#), [489](#)
14 MPI_Comm_get_attr, [217](#), [476](#)
15 MPI_COMM_GET_ERRHANDLER, [17](#), [253](#),
16 [255](#), [491](#), [560](#)
17 MPI_COMM_GET_NAME, [352](#), [352](#), [353](#),
18 [562](#)
19 MPI_COMM_GET_PARENT, [283](#), [286](#), [286](#),
20 [287](#), [353](#)
21 MPI_COMM_GROUP, [14](#), [176](#), [178](#), [178](#),
22 [183](#), [184](#), [201](#), [253](#), [560](#)
23 MPI_COMM_JOIN, [305](#), [305](#), [306](#)
24 MPI_COMM_NULL, [305](#)
25 MPI_COMM_NULL_COPY_FN, [17](#), [214](#),
26 [214](#)
27 MPI_COMM_NULL_DELETE_FN, [17](#), [215](#),
28 [215](#)
29 MPI_COMM_RANK, [183](#), [184](#), [201](#)
30 MPI_COMM_REMOTE_GROUP, [202](#)
31 MPI_COMM_REMOTE_SIZE, [201](#), [202](#)
32 MPI_COMM_SELF, [268](#), [305](#)
33 MPI_COMM_SET_ATTR, [17](#), [215](#), [216](#), [223](#),
34 [489](#)
35 MPI_Comm_set_attr, [217](#)
36 MPI_COMM_SET_ERRHANDLER, [17](#), [253](#),
37 [255](#), [491](#)
38 MPI_COMM_SET_NAME, [351](#), [351](#), [352](#)
39 MPI_COMM_SIZE, [21](#), [183](#), [184](#), [201](#)
40 MPI_COMM_SPAWN, [270](#), [280–282](#), [282](#),
41 [283](#), [285–290](#), [302](#), [303](#)
42 MPI_COMM_SPAWN_MULTIPLE, [270](#), [280](#),
43 [281](#), [286](#), [287](#), [288](#), [289](#), [303](#)
44 MPI_COMM_SPLIT, [185](#), [187](#), [189](#), [189](#),
45 [190](#), [226](#), [231](#), [232](#), [234](#), [242–244](#)
46 MPI_COMM_TEST_INTER, [200](#), [201](#)
47 MPI_COMM_WORLD, [202](#), [266](#), [280](#), [283](#),
48 [287–289](#), [304](#), [305](#), [560](#)
- MPI_Copy_function, [17](#)
MPI_Datatype, [455](#)
MPI_Delete_function, [17](#)
MPI_DIMS_CREATE, [231–233](#), [233](#)
MPI_DUP_FN, [17](#), [215](#), [488](#)
MPI_ERR_IN_STATUS, [392](#)
MPI_ERRHANDLER_C2F, [468](#)
MPI_ERRHANDLER_CREATE, [17](#), [254](#),
[490](#)
MPI_ERRHANDLER_F2C, [468](#)
MPI_ERRHANDLER_FREE, [253](#), [258](#), [560](#)
MPI_ERRHANDLER_GET, [17](#), [253](#), [255](#),
[491](#), [560](#)
MPI_ERRHANDLER_SET, [17](#), [491](#)
MPI_ERROR, [392](#)
MPI_ERROR_CLASS, [259](#), [261](#), [261](#), [355](#)
MPI_ERROR_STRING, [259](#), [259](#), [355](#), [357](#)
MPI_ERRORS_RETURN, [358](#)
MPI_EXSCAN, [165](#), [165](#)
MPI_FILE, [462](#)
MPI_FILE_C2F, [468](#)
MPI_FILE_CALL_ERRHANDLER, [358](#), [358](#)
MPI_FILE_CLOSE, [305](#), [377](#), [379](#), [380](#)
MPI_FILE_CREATE_ERRHANDLER, [253](#),
[257](#), [258](#)
MPI_FILE_DELETE, [379](#), [380](#), [380](#), [384](#),
[386](#), [432](#)
MPI_FILE_F2C, [468](#)
MPI_FILE_GET_AMODE, [383](#), [383](#)
MPI_FILE_GET_ATOMICITY, [424](#), [424](#)
MPI_FILE_GET_BYTE_OFFSET, [396](#), [401](#),
[401](#), [406](#)
MPI_FILE_GET_ERRHANDLER, [253](#), [258](#),
[432](#), [560](#)
MPI_FILE_GET_FILE_EXTENT, [415](#)
MPI_FILE_GET_GROUP, [382](#), [383](#)
MPI_FILE_GET_INFO, [385](#), [385](#), [386](#), [563](#)
MPI_FILE_GET_POSITION, [401](#), [401](#)
MPI_FILE_GET_POSITION_SHARED, [405](#),
[406](#), [406](#), [425](#)
MPI_FILE_GET_SIZE, [382](#), [382](#), [427](#)
MPI_FILE_GET_TYPE_EXTENT, [414](#), [415](#),
[421](#)
MPI_FILE_GET_VIEW, [389](#), [389](#)
MPI_FILE_IREAD, [390](#), [399](#), [399](#), [406](#), [422](#)
MPI_FILE_IREAD_AT, [390](#), [395](#), [395](#)
MPI_FILE_IREAD_SHARED, [390](#), [403](#), [403](#)
MPI_FILE_IWRITE, [390](#), [400](#), [400](#)

- MPI_FILE_IWRITE_AT, [390](#), [395](#), [396](#)
- MPI_FILE_IWRITE_SHARED, [390](#), [403](#), [404](#)
- MPI_FILE_NULL, [432](#)
- MPI_FILE_OPEN, [261](#), [369](#), [377](#), [377](#), [379](#), [384–386](#), [388](#), [401](#), [426](#), [427](#), [432](#), [433](#)
- MPI_FILE_PREALLOCATE, [381](#), [381](#), [382](#), [423](#), [427](#)
- MPI_FILE_READ, [389](#), [390](#), [396](#), [396](#), [397](#), [399](#), [426](#), [427](#)
- MPI_FILE_READ_ALL, [390](#), [397](#), [397](#), [407](#)
- MPI_FILE_READ_ALL_BEGIN, [390](#), [407](#), [408](#), [409](#), [422](#)
- MPI_FILE_READ_ALL_END, [390](#), [407](#), [408](#), [410](#), [422](#)
- MPI_FILE_READ_AT, [390](#), [393](#), [393](#), [394](#), [395](#)
- MPI_FILE_READ_AT_ALL, [390](#), [393](#), [394](#)
- MPI_FILE_READ_AT_ALL_BEGIN, [390](#), [408](#)
- MPI_FILE_READ_AT_ALL_END, [390](#), [408](#)
- MPI_FILE_READ_ORDERED, [390](#), [404](#), [405](#)
- MPI_FILE_READ_ORDERED_BEGIN, [390](#), [411](#)
- MPI_FILE_READ_ORDERED_END, [390](#), [411](#)
- MPI_FILE_READ_SHARED, [390](#), [402](#), [402](#), [403](#), [405](#)
- MPI_FILE_SEEK, [400](#), [400](#), [401](#)
- MPI_FILE_SEEK_SHARED, [405](#), [405](#), [406](#), [425](#)
- MPI_FILE_SET_ATOMICITY, [379](#), [423](#), [423](#), [424](#)
- MPI_FILE_SET_ERRHANDLER, [253](#), [258](#), [432](#)
- MPI_FILE_SET_INFO, [384](#), [384](#), [385](#), [386](#), [563](#)
- MPI_FILE_SET_SIZE, [381](#), [381](#), [382](#), [423](#), [425](#), [427](#)
- MPI_FILE_SET_VIEW, [87](#), [261](#), [378](#), [384–387](#), [387](#), [388](#), [389](#), [401](#), [406](#), [413](#), [419](#), [426](#), [433](#), [561](#), [563](#)
- MPI_FILE_SYNC, [380](#), [390](#), [422–424](#), [424](#), [425](#), [429](#)
- MPI_FILE_WRITE, [389](#), [390](#), [398](#), [398](#), [400](#), [426](#)
- MPI_FILE_WRITE_ALL, [390](#), [398](#), [398](#)
- MPI_FILE_WRITE_ALL_BEGIN, [390](#), [410](#)
- MPI_FILE_WRITE_ALL_END, [390](#), [410](#)
- MPI_FILE_WRITE_AT, [390](#), [394](#), [394](#), [395](#), [396](#)
- MPI_FILE_WRITE_AT_ALL, [390](#), [394](#), [395](#)
- MPI_FILE_WRITE_AT_ALL_BEGIN, [390](#), [409](#)
- MPI_FILE_WRITE_AT_ALL_END, [390](#), [409](#)
- MPI_FILE_WRITE_ORDERED, [390](#), [404](#), [405](#), [405](#)
- MPI_FILE_WRITE_ORDERED_BEGIN, [390](#), [412](#)
- MPI_FILE_WRITE_ORDERED_END, [390](#), [412](#)
- MPI_FILE_WRITE_SHARED, [390](#), [402](#), [403–405](#)
- MPI_FINALIZE, [14](#), [23](#), [248](#), [264](#), [264](#), [265–269](#), [304](#), [305](#), [368](#), [377](#), [467](#), [471](#), [560](#)
- MPI_Finalize, [265](#)
- MPI_FINALIZED, [20](#), [266](#), [268](#), [269](#), [269](#), [467](#)
- MPI_FREE_MEM, [251](#), [251](#), [260](#)
- MPI_Free_mem, [251](#)
- MPI_GATHER, [123](#), [125](#), [127](#), [129](#), [130](#), [137](#), [138](#), [142](#), [143](#)
- MPI_Gather, [122](#)
- MPI_GATHERV, [123](#), [129](#), [129](#), [130](#), [131](#), [139](#), [144](#)
- MPI_Gatherv, [122](#)
- MPI_GET, [307](#), [311](#), [312](#), [314](#), [319](#), [561](#)
- MPI_GET_ADDRESS, [17](#), [77](#), [92](#), [92](#), [93](#), [102](#), [453](#), [455](#), [472](#), [473](#), [487](#)
- MPI_GET_COUNT, [31](#), [32](#), [32](#), [52](#), [100](#), [101](#), [350](#), [392](#), [560](#)
- MPI_GET_ELEMENT, [101](#)
- MPI_GET_ELEMENTS, [52](#), [100](#), [100](#), [101](#), [350](#), [351](#), [392](#)
- MPI_GET_PROCESSOR_NAME, [249](#), [250](#), [562](#)
- MPI_GET_VERSION, [247](#), [248](#), [266](#)
- MPI_GRAPH_CREATE, [231](#), [234](#), [234](#), [236](#), [237](#), [244](#), [562](#)
- MPI_GRAPH_GET, [231](#), [236](#), [237](#)
- MPI_GRAPH_MAP, [231](#), [244](#), [244](#)
- MPI_GRAPH_NEIGHBORS, [231](#), [239](#), [239](#)

- 1 MPI_GRAPH_NEIGHBORS_COUNT, [231](#),
2 [239](#), [239](#)
- 3 MPI_GRAPHDIMS_GET, [231](#), [236](#), [236](#)
- 4 MPI_GREQUEST_COMPLETE, [344–347](#),
5 [347](#)
- 6 MPI_GREQUEST_START, [344](#), [345](#)
- 7 MPI_GROUP_C2F, [468](#)
- 8 MPI_GROUP_COMPARE, [177](#), [180](#)
- 9 MPI_GROUP_DIFFERENCE, [179](#)
- 10 MPI_GROUP_EXCL, [180](#), [180](#), [182](#)
- 11 MPI_GROUP_F2C, [468](#)
- 12 MPI_Group_f2c, [21](#)
- 13 MPI_GROUP_FREE, [182](#), [183](#), [184](#), [253](#),
14 [560](#)
- 15 MPI_GROUP_INCL, [180](#), [180](#), [181](#)
- 16 MPI_GROUP_INTERSECTION, [179](#)
- 17 MPI_GROUP_RANGE_EXCL, [182](#), [182](#)
- 18 MPI_GROUP_RANGE_INCL, [181](#), [181](#)
- 19 MPI_GROUP_RANK, [176](#), [184](#)
- 20 MPI_GROUP_SIZE, [176](#), [183](#)
- 21 MPI_GROUP_TRANSLATE_RANKS, [177](#),
22 [177](#), [560](#)
- 23 MPI_GROUP_UNION, [178](#)
- 24 MPI_Handler_function, [17](#), [254](#)
- 25 MPI_IBSEND, [49](#), [53](#), [71](#)
- 26 MPI_Info, [273](#), [285](#), [561](#)
- 27 MPI_INFO_C2F, [468](#)
- 28 MPI_INFO_CREATE, [274](#), [274](#)
- 29 MPI_INFO_DELETE, [260](#), [275](#), [275](#), [277](#)
- 30 MPI_INFO_DUP, [277](#), [277](#)
- 31 MPI_INFO_F2C, [468](#)
- 32 MPI_INFO_FREE, [277](#), [385](#)
- 33 MPI_INFO_GET, [273](#), [275](#), [561](#)
- 34 MPI_INFO_GET_NKEYS, [273](#), [276](#), [276](#),
35 [277](#), [561](#)
- 36 MPI_INFO_GET_NTHKEY, [273](#), [276](#), [561](#)
- 37 MPI_INFO_GET_VALUELEN, [273](#), [276](#),
38 [561](#)
- 39 MPI_INFO_SET, [274](#), [274](#), [275](#), [277](#)
- 40 MPI_INIT, [14](#), [23](#), [175](#), [207](#), [248](#), [263](#), [263](#),
41 [264](#), [266](#), [267](#), [269](#), [283–286](#), [302](#),
42 [303](#), [370–372](#), [467](#), [471](#), [481](#)
- 43 MPI_Init, [263](#)
- 44 MPI_INIT_THREAD, [175](#), [370](#), [371](#), [372](#)
- 45 MPI_INITIALIZED, [263](#), [266](#), [267](#), [267](#),
46 [268](#), [269](#), [372](#), [467](#)
- 47 MPI_INTERCOMM_CREATE, [122](#), [185](#),
48 [202](#), [203](#), [204](#)
- MPI_Intercomm_create, [210](#)
- MPI_INTERCOMM_MERGE, [199](#), [202–204](#),
[204](#)
- MPI_IPROBE, [32](#), [63](#), [64](#), [64](#), [65](#), [369](#)
- MPI_Iprobe, [265](#)
- MPI_Irecv, [51](#), [452–455](#)
- MPI_IRSEND, [50](#)
- MPI_IS_THREAD_MAIN, [370](#), [372](#)
- MPI_ISEND, [11](#), [49](#), [71](#), [264](#), [452](#)
- MPI_ISSEND, [50](#)
- MPI_KEYVAL_CREATE, [17](#), [214](#), [215](#), [488](#),
[489](#), [500](#)
- MPI_KEYVAL_FREE, [17](#), [216](#), [223](#), [489](#)
- MPI_LB, [17](#)
- MPI_LOOKUP_NAME, [260](#), [293](#), [297](#), [298](#),
[299](#)
- MPI_MAX, [165](#)
- MPI_NULL_COPY_FN, [17](#), [215](#), [488](#)
- MPI_NULL_DELETE_FN, [17](#), [215](#), [489](#)
- MPI_OP_C2F, [468](#)
- MPI_OP_CREATE, [157](#), [158](#), [159](#), [500](#)
- MPI_OP_F2C, [468](#)
- MPI_OP_FREE, [160](#)
- MPI_OPEN_PORT, [293](#), [294](#), [294](#), [295–](#)
[297](#), [299](#)
- MPI_PACK, [47](#), [111](#), [113](#), [116](#), [117](#), [416](#),
[420](#)
- MPI_PACK_EXTERNAL, [7](#), [116](#), [117](#), [117](#),
[462](#), [561](#)
- MPI_PACK_EXTERNAL_SIZE, [118](#)
- MPI_PACK_SIZE, [47](#), [113](#), [113](#)
- MPI_PCONTROL, [479–481](#), [481](#)
- MPI_PROBE, [30](#), [32](#), [33](#), [63](#), [64](#), [64](#), [65](#),
[66](#), [369](#)
- MPI_PUBLISH_NAME, [293](#), [297](#), [297](#), [298](#),
[299](#)
- MPI_PUT, [307](#), [311](#), [312](#), [312](#), [314](#), [318](#),
[319](#), [326](#), [330](#), [332](#), [561](#)
- MPI_QUERY_THREAD, [372](#), [373](#)
- MPI_RECV, [26](#), [29](#), [31–33](#), [64](#), [66](#), [76](#), [99](#),
[100](#), [112](#), [121](#), [128](#), [351](#), [428](#), [439](#),
[454–456](#)
- MPI_Recv, [368](#)
- MPI_RECV_INIT, [70](#), [70](#)
- MPI_REDUCE, [123](#), [149](#), [150](#), [151](#), [158](#),
[159](#), [161](#), [163–165](#), [318](#)
- MPI_Reduce, [122](#)

- MPIREDUCE_SCATTER, [123](#), [151](#), [158](#),
 [163](#), [163](#)
 MPI.Reduce_scatter, [122](#)
 MPI.REGISTER_DATAREP, [261](#), [418](#), [419](#)–
 [421](#), [433](#)
 MPI.REQUEST_C2F, [468](#)
 MPI.REQUEST_F2C, [468](#)
 MPI.REQUEST_FREE, [21](#), [54](#), [54](#), [55](#), [66](#),
 [71](#), [264](#), [265](#), [346](#), [347](#)
 MPI.REQUEST_GET_STATUS, [33](#), [63](#), [345](#)
 MPI.RSEND, [40](#)
 MPI.RSEND_INIT, [70](#)
 MPI.SCAN, [151](#), [158](#), [164](#), [164](#), [165](#)
 MPI.Scan, [123](#)
 MPI.SCATTER, [123](#), [137](#), [137](#), [139](#)
 MPI.Scatter, [123](#)
 MPI.SCATTERV, [123](#), [138](#), [139](#), [140](#), [163](#)
 MPI.Scatterv, [123](#)
 MPI.SEND, [25](#), [26](#), [27](#), [32](#), [35](#), [76](#), [99](#), [111](#),
 [244](#), [378](#), [428](#), [439](#), [451](#), [454](#)
 MPI.Send, [137](#), [368](#)
 MPI.SEND_INIT, [68](#), [71](#)
 MPI.SENDRECV, [73](#), [240](#)
 MPI.SENDRECV_REPLACE, [74](#)
 MPI.SIGNED_CHAR, [153](#)
 MPI.SIZEOF, [457](#), [463](#), [464](#)
 MPI.SPAWN, [290](#)
 MPI.SSEND, [39](#)
 MPI.SSEND_INIT, [69](#)
 MPI.START, [71](#), [71](#), [72](#)
 MPI.STARTALL, [71](#), [71](#)
 MPI.STATUS, [32](#), [33](#)
 MPI.STATUS_C2F, [471](#)
 MPI.STATUS_F2C, [471](#)
 MPI.Status_f2c, [471](#)
 MPI.STATUS_IGNORE, [33](#), [555](#)
 MPI.STATUS_SET_CANCELLED, [351](#)
 MPI.STATUS_SET_ELEMENTS, [350](#), [350](#)
 MPI.STATUS_IGNORE, [33](#), [555](#)
 MPI.SUCCESS, [59](#), [286](#)
 MPI.TEST, [11](#), [33](#), [51](#)–[53](#), [53](#), [54](#), [56](#), [58](#),
 [66](#), [67](#), [71](#), [264](#), [265](#), [347](#), [391](#), [392](#)
 MPI.TEST_CANCELLED, [52](#), [53](#), [67](#), [67](#),
 [345](#), [351](#), [392](#)
 MPI.Test_cancelled, [265](#)
 MPI.TESTALL, [56](#), [59](#), [60](#), [345](#), [346](#), [350](#)
 MPI.TESTANY, [56](#), [57](#), [58](#), [61](#), [345](#), [346](#),
 [350](#)
- MPI.TESTSOME, [56](#), [61](#), [61](#), [345](#), [346](#),
 [350](#)
 MPI.THREAD_INIT, [467](#)
 MPI.TOPO_TEST, [231](#), [236](#), [236](#)
 MPI.TYPE_C2F, [468](#)
 MPI.TYPE_COMMIT, [97](#), [97](#), [98](#), [468](#)
 MPI.Type_commit, [468](#)
 MPI.TYPE_CONTIGUOUS, [12](#), [77](#), [77](#),
 [79](#), [94](#), [360](#), [376](#), [415](#)
 MPI.TYPE_CREATE_DARRAY, [12](#), [32](#),
 [88](#), [88](#), [360](#)
 MPI.TYPE_CREATE_F90_COMPLEX, [12](#),
 [360](#), [362](#), [418](#), [457](#), [460](#), [462](#)
 MPI.TYPE_CREATE_F90_INTEGER, [12](#),
 [360](#), [362](#), [418](#), [457](#), [460](#), [462](#)
 MPI.TYPE_CREATE_F90_REAL, [12](#), [360](#),
 [362](#), [418](#), [457](#), [459](#), [460](#)–[462](#)
 MPI.TYPE_CREATE_F90_REAL/COMPLEX/INTEGER,
 [462](#)
 MPI.TYPE_CREATE_F90_xxxx, [461](#), [563](#)
 MPI.TYPE_CREATE_HINDEXED, [12](#), [17](#),
 [77](#), [82](#), [82](#), [83](#), [360](#), [485](#)
 MPI.Type_create_hindexed, [359](#)
 MPI.TYPE_CREATE_HVECTOR, [12](#), [17](#),
 [77](#), [79](#), [79](#), [360](#), [485](#)
 MPI.TYPE_CREATE_INDEXED_BLOCK,
 [83](#), [360](#)
 MPI.TYPE_CREATE_KEYVAL, [212](#), [220](#),
 [223](#), [474](#), [561](#)
 MPI.TYPE_CREATE_RESIZED, [16](#), [17](#),
 [95](#), [360](#), [415](#)
 MPI.TYPE_CREATE_STRUCT, [12](#), [17](#),
 [77](#), [84](#), [147](#), [360](#), [486](#)
 MPI.TYPE_CREATE_SUBARRAY, [12](#), [14](#),
 [85](#), [87](#), [89](#), [360](#)
 MPI.TYPE_DELETE_ATTR, [223](#), [223](#)
 MPI.TYPE_DUP, [12](#), [98](#), [98](#), [360](#)
 MPI.TYPE_DUP_FN, [221](#), [221](#)
 MPI.TYPE_EXTENT, [17](#), [95](#), [97](#), [474](#), [487](#)
 MPI.TYPE_F2C, [468](#)
 MPI.TYPE_FREE, [97](#), [221](#), [362](#)
 MPI.Type_free, [365](#)
 MPI.TYPE_FREE_KEYVAL, [222](#), [223](#)
 MPI.TYPE_GET_ATTR, [222](#), [223](#)
 MPI.TYPE_GET_CONTENTS, [359](#), [361](#),
 [362](#), [363](#)
 MPI.TYPE_GET_ENVELOPE, [359](#), [359](#),
 [361](#), [363](#), [461](#)

- 1 MPI.TYPE_GET_EXTENT, [17](#), [95](#), [464](#),
2 [472](#), [474](#), [487](#)
- 3 MPI.TYPE_GET_NAME, [354](#)
- 4 MPI.TYPE_GET_TRUE_EXTENT, [96](#)
- 5 MPI.TYPE_HINDEXED, [17](#), [82](#), [85](#), [359](#),
6 [360](#), [474](#), [486](#)
- 7 MPI.Type_hindexed, [359](#)
- 8 MPI.TYPE_HVECTOR, [17](#), [80](#), [360](#), [474](#),
9 [485](#)
- 10 MPI.TYPE_INDEXED, [12](#), [80](#), [80](#), [81–83](#),
11 [360](#)
- 12 MPI.TYPE_INDEXED_BLOCK, [12](#)
- 13 MPI.TYPE_LB, [17](#), [95](#), [487](#)
- 14 MPI.TYPE_MATCH_SIZE, [457](#), [464](#), [464](#)
- 15 MPI.TYPE_NULL_COPY_FN, [221](#), [221](#)
- 16 MPI.TYPE_NULL_DELETE_FN, [221](#)
- 17 MPI.TYPE_SET_ATTR, [222](#), [223](#)
- 18 MPI.TYPE_SET_NAME, [353](#)
- 19 MPI.TYPE_SIZE, [93](#), [93](#), [94](#)
- 20 MPI.TYPE_STRUCT, [16](#), [17](#), [84](#), [85](#), [94](#),
21 [360](#), [474](#), [486](#)
- 22 MPI.TYPE_UB, [17](#), [95](#), [474](#), [488](#)
- 23 MPI.TYPE_UB, MPI.TYPE_LB, [95](#)
- 24 MPI.TYPE_VECTOR, [12](#), [78](#), [78](#), [79](#), [81](#),
25 [360](#)
- 26 MPI_UB, [17](#)
- 27 MPI.UNPACK, [111](#), [112](#), [116](#), [420](#)
- 28 MPI.UNPACK_EXTERNAL, [7](#), [117](#), [462](#)
- 29 MPI.UNPUBLISH_NAME, [260](#), [298](#), [298](#)
- 30 MPI.UNSIGNED_CHAR, [153](#)
- 31 MPI.WAIT, [31](#), [33](#), [51](#), [52](#), [52](#), [53](#), [54](#), [56](#),
32 [57](#), [59](#), [66](#), [67](#), [71](#), [264](#), [265](#), [343](#),
33 [347](#), [369](#), [391](#), [392](#), [406](#), [422](#), [424](#),
34 [455](#), [456](#)
- 35 MPI.WAITALL, [56](#), [58](#), [59](#), [60](#), [345](#), [346](#),
36 [350](#), [368](#), [446](#)
- 37 MPI.WAITANY, [41](#), [56](#), [56](#), [57](#), [58](#), [61](#),
38 [345](#), [346](#), [350](#), [368](#)
- 39 MPI.WAITSOME, [56](#), [60](#), [60](#), [61](#), [62](#), [345](#),
40 [346](#), [350](#), [368](#)
- 41 MPI.WCHAR, [153](#)
- 42 MPI.WIN_BASE, [476](#)
- 43 MPI.WIN_C2F, [468](#)
- 44 MPI.WIN_CALL_ERRHANDLER, [357](#), [358](#)
- 45 MPI.WIN_COMPLETE, [310](#), [320](#), [325](#), [325](#),
46 [326](#), [328](#), [335](#), [336](#)
- 47 MPI.WIN_CREATE, [308](#), [309](#), [334](#), [369](#)
- 48 MPI.WIN_CREATE_ERRHANDLER, [253](#),
[256](#), [256](#)
- MPI.WIN_CREATE_KEYVAL, [212](#), [218](#),
[223](#), [474](#), [561](#)
- MPI.WIN_DELETE_ATTR, [220](#), [223](#)
- MPI.WIN_DUP_FN, [218](#), [218](#)
- MPI.WIN_F2C, [468](#)
- MPI.WIN_FENCE, [310](#), [319](#), [324](#), [324](#), [330](#),
[335–338](#)
- MPI.WIN_FREE, [219](#), [305](#), [309](#), [310](#)
- MPI.WIN_FREE_KEYVAL, [219](#), [223](#)
- MPI.WIN_GET_ATTR, [220](#), [223](#), [310](#), [476](#)
- MPI.Win_get_attr, [310](#), [476](#)
- MPI.WIN_GET_ERRHANDLER, [253](#), [257](#),
[560](#)
- MPI.WIN_GET_GROUP, [310](#), [311](#)
- MPI.WIN_GET_NAME, [354](#)
- MPI.WIN_LOCK, [250](#), [309](#), [320](#), [328](#), [329](#),
[330](#), [336](#)
- MPI.WIN_NULL_COPY_FN, [218](#), [218](#)
- MPI.WIN_NULL_DELETE_FN, [218](#)
- MPI.WIN_POST, [310](#), [320](#), [325](#), [326](#), [326](#),
[327–331](#), [338](#)
- MPI.WIN_POST, MPI.WIN_FENCE, [336](#)
- MPI.WIN_SET_ATTR, [219](#), [223](#)
- MPI.WIN_SET_ERRHANDLER, [253](#), [256](#)
- MPI.WIN_SET_NAME, [354](#)
- MPI.WIN_START, [320](#), [325](#), [325](#), [326](#), [328](#),
[330](#), [331](#), [334](#)
- MPI.WIN_TEST, [327](#), [327](#)
- MPI.WIN_UNLOCK, [320](#), [329](#), [330](#), [335](#),
[336](#), [338](#)
- MPI.WIN_WAIT, [310](#), [320](#), [326](#), [326](#), [327–](#)
[329](#), [336](#), [338](#)
- MPI.WTICK, [21](#), [262](#), [262](#)
- MPI.WTIME, [21](#), [249](#), [262](#), [262](#)
- MPI_xxx_c2f, [469](#)
- MPI_xxx_f2c, [469](#)
- MPI_xxx_GET_ATTR, [475](#)
- MPI_xxx_get_attr, [475](#)
- MPI_XXX_GET_ERRHANDLER, [253](#)
- MPI.YYY_DELETE_ATTR, [212](#), [561](#)
- MPI.YYY_FREE_KEYVAL, [212](#), [561](#)
- MPI.YYY_GET_ATTR, [212](#), [561](#)
- MPI.YYY_SET_ATTR, [212](#), [561](#)
- mpiexec, [371](#)
- PMPI_, [479](#)

| | |
|---|----|
| PMPI.WTICK, 21 | 1 |
| PMPI.WTIME, 21 | 2 |
| | 3 |
| query_fn, 345–347 | 4 |
| | 5 |
| read_conversion_fn, 419–421 | 6 |
| | 7 |
| sizeof, 464 | 8 |
| | 9 |
| TYPEDEF:MPI_Comm_copy_attr_function, 214 | 10 |
| TYPEDEF:MPI_Comm_delete_attr_function, 214 | 11 |
| | 12 |
| TYPEDEF:MPI_Comm_errhandler_fn, 254 | 13 |
| TYPEDEF:MPI_Datarep_conversion_function, 419 | 14 |
| | 15 |
| TYPEDEF:MPI_Datarep_extent_function, 419 | 16 |
| | 17 |
| TYPEDEF:MPI_File_errhandler_fn, 257 | 18 |
| TYPEDEF:MPI_Grequest_cancel_function, 346 | 19 |
| | 20 |
| TYPEDEF:MPI_Grequest_free_function, 345 | 21 |
| TYPEDEF:MPI_Grequest_query_function, 345 | 22 |
| | 23 |
| TYPEDEF:MPI_Type_copy_attr_function, 221 | 24 |
| | 25 |
| TYPEDEF:MPI_Type_delete_attr_function, 221 | 26 |
| | 27 |
| TYPEDEF:MPI_Win_copy_attr_function, 218 | 28 |
| TYPEDEF:MPI_Win_delete_attr_function, 218 | 29 |
| | 30 |
| TYPEDEF:MPI_Win_errhandler_fn, 256 | 31 |
| | 32 |
| write_conversion_fn, 421 | 33 |
| | 34 |
| | 35 |
| | 36 |
| | 37 |
| | 38 |
| | 39 |
| | 40 |
| | 41 |
| | 42 |
| | 43 |
| | 44 |
| | 45 |
| | 46 |
| | 47 |
| | 48 |