# Chapter 1

# Introduction to MPI

## 1.1 Overview and Goals

MPI (Message-Passing Interface) is a *message-passing library interface specification*. All parts of this definition are significant. MPI addresses primarily the message-passing parallel programming model, in which data is moved from the address space of one process to that of another process through cooperative operations on each process. [(]Extensions to the "classical" message-passing model are provided in collective operations, remote-memory access operations, dynamic process creation, and parallel I/O.[)] MPI is a *specification*, not an implementation; there are multiple implementations of MPI. This specification is for a *library interface*; MPI is not a language, and all MPI operations are expressed as functions, subroutines, or methods, according to the appropriate language bindings, which for C, C++, Fortran-77, and Fortran-95, are part of the MPI standard. The standard has been defined through an open process by a community of parallel computing vendors, computer scientists, and application developers. The next few sections provide an overview of the history of MPI's development.

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

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

A complete list of goals follows.

- Design an application programming interface (not necessarily for compilers or a system implementation library).

- Allow efficient communication: Avoid memory-to-memory copying, allow overlap of computation and communication, and offload to communication co-processor, where available.

- Allow for implementations that can be used in a heterogeneous environment.

ticket0.

ticket0.

ticket0.

- Allow convenient C, C++, Fortran-77, and Fortran-95 bindings for the interface.

- Assume a reliable communication interface: the user need not cope with communication failures. Such failures are dealt with by the underlying communication subsystem.

- Define an interface that can be implemented on many vendor's platforms, with no significant changes in the underlying communication and system software.

- Semantics of the interface should be language independent.

- The interface should be designed to allow for thread safety.

## 1.2   Background of MPI-1.0

MPI 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 was strongly influenced by work at the IBM T. J. Watson Research Center [1, 2], Intel's NX/2 [42], Express [12], nCUBE's Vertex [38], p4 [7, 8], and PARMACS [5, 9]. Other important contributions have come from Zipcode [44, 45], Chimp [16, 17], PVM [4, 14], Chameleon [25], and PICL [24].

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

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

## 1.3 Background of MPI-1.1, MPI-1.2, and MPI-2.0

Beginning in March 1995, the MPI Forum began meeting to consider corrections and extensions to the original MPI Standard document [21]. The first product of these deliberations was Version 1.1 of the MPI specification, released in June of 1995 [22] (see `http://www.mpi-forum.org` for official MPI document releases). At that time, effort 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++. MPI-2 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. zero-copy semantics on shared-memory machines, real-time specifications).

Corrections and clarifications (items of type 1 in the above list) were collected in Chapter 3 of the MPI-2 document: "Version 1.2 of MPI." That 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-1.3 program and a valid MPI-2.1 program, and a valid MPI-1.3 program is a valid MPI-2.1 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 errata and clarifications for both standard documents (MPI-1.1 and MPI-2.0). The short document "Errata for MPI-1.1" was released October 12, 1998. On July 5, 2001, a first ballot of errata and clarifications for

MPI-2.0 was released, and a second ballot was voted on May 22, 2002. Both votes were done electronically. Both ballots were combined into one document: "Errata for MPI-2", May 15, 2002. This errata process was then interrupted, but the Forum and its e-mail reflectors kept working on new requests for clarification.

Restarting regular work of the MPI Forum was initiated in three meetings, at EuroPVM/MPI'06 in Bonn, at EuroPVM/MPI'07 in Paris, and at SC'07 in Reno. In December 2007, a steering committee started the organization of new MPI Forum meetings at regular 8-weeks intervals. At the January 14-16, 2008 meeting in Chicago, the MPI Forum decided to combine the existing and future MPI documents to one [single] document for each version of the MPI standard. For technical and historical reasons, this series was started with MPI-1.3. Additional Ballots 3 and 4 solved old questions from the errata list started in 1995 up to new questions from the last years. After all documents (MPI-1.1, MPI-2, Errata for MPI-1.1 (Oct. 12, 1998), and MPI-2.1 Ballots 1-4) were combined into one draft document, for each chapter, a chapter author and review team were defined. They cleaned up the document to achieve a consistent MPI-2.1 document. The final MPI-2.1 standard document was finished in June 2008, and finally released with a second vote in September 2008 in the meeting at Dublin, just before EuroPVM/MPI'08. The major work of the current MPI Forum is the preparation of MPI-3.

## 1.5   Background of MPI-2.2

MPI-2.2 is a minor update to the MPI-2.1 standard. This version addresses additional errors and ambiguities that were not corrected in the MPI-2.1 standard as well as a small number of extensions to MPI-2.1 that met the following criteria:

- Any correct MPI-2.1 program is a correct MPI-2.2 program.

- Any extension must have significant benefit for users.

- Any extension must not require significant implementation effort. To that end, all such changes are accompanied by an open source implementation.

The discussions of MPI-2.2 proceeded concurrently with the MPI-3 discussions; in some cases, extensions were proposed for MPI-2.2 but were later moved to MPI-3.

## 1.6   Background of MPI-3.0

MPI-3.0 is a major update to the MPI standard. Areas of particular interest are the extension of collective operations to include nonblocking, with other areas under consideration. This *draft* contains the MPI Forum's current draft of nonblocking collective routines.

## 1.7   Who Should Use This Standard?

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

ticket0.

ticket0.

## 1.8 What Platforms Are Targets For Implementation?

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

The interface is suitable for use by fully general MIMD programs, as well as those written in the more restricted style of SPMD. MPI provides many features intended to improve performance on scalable parallel computers with specialized interprocessor communication hardware. Thus, we expect that native, high-performance implementations of MPI will be provided on such machines. At the same time, implementations of MPI on top of standard Unix interprocessor communication protocols will provide portability to workstation clusters and heterogenous networks of workstations.

## 1.9 What Is Included In The Standard?

The standard includes:

- Point-to-point communication,

- Datatypes,

- Collective operations,

- Process groups,

- Communication contexts,

- Process topologies,

- Environmental [M]management and inquiry,

- The [i]Info object,

- Process creation and management,

- One-sided communication,

- External interfaces,

- Parallel file I/O,

- Language [B]bindings for Fortran, C and C++,

- Profiling interface.

**Unofficial Draft for Comment Only**

## 1.10   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.

## 1.11   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, Datatypes, defines a method to describe any data layout, e.g., an array of structures in the memory, which can be used as message send or receive buffer.

- Chapter 5, 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 two new collective operations. MPI-3 adds nonblocking collective operations.

- Chapter 6, Groups, Contexts, Communicators, and Caching, 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 7, 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 8, 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.

ticket0.

ticket0.

- Chapter 9, The Info Object, defines an opaque object, that is used as input [of]in several MPI routines.

- Chapter 10, Process Creation and Management, defines routines that allow for creation of processes.

- Chapter 11, 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 12, 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 13, I/O, defines MPI support for parallel I/O.

- 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.

- Chapter 15, Deprecated Functions, describes routines that are kept for reference. However usage of these functions is discouraged, as they may be deleted in future versions of the standard.

- Chapter 16, Language Bindings, describes the C++ binding, discusses Fortran issues, and describes language interoperability aspects between C, C++, and Fortran.

The Appendices are:

- Annex A, Language Bindings Summary, gives specific syntax in C, C++, and Fortran, for all MPI functions, constants, and types.

- Annex B, Change-Log, summarizes major changes since the previous version of the standard.

- Several Index pages [are showing]show the locations of examples, constants and pre-defined handles, callback routine[s'] prototypes, and all MPI functions. ticket0. ticket0.

MPI 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.

**Unofficial Draft for Comment Only**

- 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.

# Chapter 2

# MPI Terms and Conventions

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

In many cases MPI names for C functions are of the form MPI_Class_action_subset. This convention originated with MPI-1. Since MPI-2 an attempt has been made to standardize the names of MPI functions according to the following rules. The C++ bindings in particular follow these rules (see Section 2.6.4 on page 18).

1. In C, all routines associated with a particular type of MPI object should be of the form MPI_Class_action_subset or, if no subset exists, of the form MPI_Class_action. In Fortran, all routines associated with a particular type of MPI object should be of the form MPI_CLASS_ACTION_SUBSET or, if no subset exists, of the form

MPI_CLASS_ACTION. For C and Fortran we use the C++ terminology to define the Class. In C++, the routine is a method on **Class** and is named MPI::Class::Action_subset. If the routine is associated with a certain class, but does not make sense as an object method, it is a static member function of the class.

2. If the routine is not associated with a class, the name should be of the form MPI_Action_subset in C and MPI_ACTION_SUBSET in Fortran, and in C++ should be scoped in the MPI namespace, MPI::Action_subset.

3. The names of certain actions have been standardized. In particular, **Create** creates a new object, **Get** retrieves information about an object, **Set** sets this information, **Delete** deletes information, **Is** asks whether or not an object has a certain property.

   C and Fortran names for some MPI functions (that were defined during the MPI-1 process) violate these rules in several cases. The most common exceptions are the omission of the **Class** name from the routine and the omission of the **Action** where one can be inferred.

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

## 2.3   Procedure Specification

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

   • IN: the call may use the input value but does not update the argument,

   • OUT: the call may update the argument but does not use its input value,

   • INOUT: the call may both use and update the argument.

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

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

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

   A common occurrence for MPI functions is an argument that is used as 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 opera-tion 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 ef-fects: the status is extracted from the request; in the case of test and wait, if the request was nonpersistent, it is **freed**, and becomes **inactive** if it was persistent. 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 exe-cution of some MPI procedure on another process. Such an operation may require communication occurring with another user process.

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

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

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

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

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

## 2.5   Data Types

### 2.5.1   Opaque Objects

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

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

*Advice to implementors.* In Fortran, the handle can be an index into a table of 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 accidently 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.*)

**Unofficial Draft for Comment Only**

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

### 2.5.2   Array Arguments

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

### 2.5.3   State

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

### 2.5.4   Named Constants

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

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

All named constants, with the exceptions noted below for Fortran, can be used in
initialization expressions or assignments, but not necessarily in array declarations or as
labels in C/C++ switch or Fortran select/case statements. This implies named constants
to be link-time but not necessarily compile-time constants. The named constants listed
below are required to be compile-time constants in both C/C++ and Fortran.  These
constants do not change values during execution. Opaque objects accessed by constant
handles are defined and do not change value between MPI initialization (MPI_INIT) and
MPI completion (MPI_FINALIZE). The handles themselves are constants and can be also
used in initialization expressions or assignments.

The constants that are required to be compile-time constants (and can thus be used for array length declarations and labels in C/C++ `switch` and Fortran `case`/`select` statements) are:

MPI_MAX_PROCESSOR_NAME
MPI_MAX_ERROR_STRING
MPI_MAX_DATAREP_STRING
MPI_MAX_INFO_KEY
MPI_MAX_INFO_VAL
MPI_MAX_OBJECT_NAME
MPI_MAX_PORT_NAME
MPI_STATUS_SIZE (Fortran only)
MPI_ADDRESS_KIND (Fortran only)
MPI_INTEGER_KIND (Fortran only)
MPI_OFFSET_KIND (Fortran only)

and their C++ counterparts where appropriate.

The constants that cannot be used in initialization expressions or assignments in Fortran are:

MPI_BOTTOM
MPI_STATUS_IGNORE
MPI_STATUSES_IGNORE
MPI_ERRCODES_IGNORE
MPI_IN_PLACE
MPI_ARGV_NULL
MPI_ARGVS_NULL
MPI_UNWEIGHTED

> *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. These types must have the same

width and encode address values in the same manner such that address values in one language may be passed directly to another language without conversion. 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. These types must have the same width and encode address values in the same manner such that offset values in one language may be passed directly to another language without conversion.

### 2.5.8   Counts

Derived datatypes can be created representing more elements than can be encoded in a C int or Fortran INTEGER. MPI_GET_COUNT, MPI_GET_ELEMENTS , and associated functions cannot properly express these quantities. To overcome this, these quantities are declared to be INTEGER (KIND=MPI_COUNT_KIND) in Fortran.  In C one uses MPI_Count. These types must have the same width and encode values in the same manner such that count values in one language may be passed directly to another language without conversion. The size of the MPI_Count type is determined by the MPI implementation with the restriction that it must be minimally capable of encoding a C int and Fortran INTEGER.

## 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.) The C++ language bindings have been deprecated. 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 understand the word "argument" (which has no specific meaning in C/C++), thus allowing us to avoid unnecessary confusion for Fortran programmers.

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

### 2.6.1   Deprecated Names and Functions

A number of chapters refer to deprecated or replaced MPI-1 constructs. These are constructs that continue to be part of the MPI standard, as documented in Chapter 15, but that users are recommended not to continue using, since better solutions were provided with MPI-2. For example, the Fortran binding for MPI-1 functions that have address arguments uses

ticket265.

INTEGER. This is not consistent with the C binding, and causes problems on machines with 32 bit INTEGERs and 64 bit addresses. In MPI-2, these functions were given new names with new bindings for the address arguments. The use of the old functions is deprecated. For consistency, here and in a few other cases, new C functions are also provided, even though the new functions are equivalent to the old functions. The old names are deprecated. Another example is provided by the MPI-1 predefined datatypes MPI_UB and MPI_LB. They are deprecated, since their use is awkward and error-prone. The MPI-2 function MPI_TYPE_CREATE_RESIZED provides a more convenient mechanism to achieve the same effect.

Table 2.1 shows a list of all of the deprecated constructs. Note that the constants MPI_LB and MPI_UB are replaced by the function MPI_TYPE_CREATE_RESIZED; this is because their principal use was as input datatypes to MPI_TYPE_STRUCT to create resized datatypes. Also note that some C typedefs and Fortran subroutine names are included in this list; they are the types of callback functions.

| Deprecated | MPI-2 Replacement |
|---|---|
| MPI_ADDRESS | MPI_GET_ADDRESS |
| MPI_TYPE_HINDEXED | MPI_TYPE_CREATE_HINDEXED |
| MPI_TYPE_HVECTOR | MPI_TYPE_CREATE_HVECTOR |
| MPI_TYPE_STRUCT | MPI_TYPE_CREATE_STRUCT |
| MPI_TYPE_EXTENT | MPI_TYPE_GET_EXTENT |
| MPI_TYPE_UB | MPI_TYPE_GET_EXTENT |
| MPI_TYPE_LB | MPI_TYPE_GET_EXTENT |
| MPI_LB | MPI_TYPE_CREATE_RESIZED |
| MPI_UB | MPI_TYPE_CREATE_RESIZED |
| MPI_ERRHANDLER_CREATE | MPI_COMM_CREATE_ERRHANDLER |
| MPI_ERRHANDLER_GET | MPI_COMM_GET_ERRHANDLER |
| MPI_ERRHANDLER_SET | MPI_COMM_SET_ERRHANDLER |
| MPI_Handler_function | MPI_Comm_errhandler_function |
| MPI_KEYVAL_CREATE | MPI_COMM_CREATE_KEYVAL |
| MPI_KEYVAL_FREE | MPI_COMM_FREE_KEYVAL |
| MPI_DUP_FN | MPI_COMM_DUP_FN |
| MPI_NULL_COPY_FN | MPI_COMM_NULL_COPY_FN |
| MPI_NULL_DELETE_FN | MPI_COMM_NULL_DELETE_FN |
| MPI_Copy_function | MPI_Comm_copy_attr_function |
| COPY_FUNCTION | COMM_COPY_ATTR_FN |
| MPI_Delete_function | MPI_Comm_delete_attr_function |
| DELETE_FUNCTION | COMM_DELETE_ATTR_FN |
| MPI_ATTR_DELETE | MPI_COMM_DELETE_ATTR |
| MPI_ATTR_GET | MPI_COMM_GET_ATTR |
| MPI_ATTR_PUT | MPI_COMM_SET_ATTR |

Table 2.1: Deprecated constructs

### 2.6.2   Fortran Binding Issues

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

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

All MPI Fortran subroutines have a return code in the last argument. A few MPI operations which are functions do not have the return code argument. The return code value for successful completion is `MPI_SUCCESS`. Other error codes are implementation dependent; see the error codes in Chapter 8 and Annex A.

Constants representing the maximum length of a string are one smaller in Fortran than in C and C++ as discussed in Section 16.3.9.

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

Array arguments are indexed from one.

The MPI Fortran binding is inconsistent with the Fortran 90 standard in several respects. These inconsistencies, such as register optimization problems, have implications for user codes that are discussed in detail in Section 16.2.2. They are also inconsistent with Fortran 77.

### 2.6.3   C Binding Issues

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

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

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

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

Array arguments are indexed from zero.

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

Choice arguments are pointers of type `void *`.

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

### 2.6.4   C++ Binding Issues

The C++ language bindings have been deprecated. There are places in the standard that give rules for C and not for C++. In these cases, the C rule should be applied to the C++

case, as appropriate. In particular, the values of constants given in the text are the ones
for C and Fortran. A cross index of these with the C++ names is given in Annex A.

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

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

> *Advice to implementors.* The file mpi.h may contain both the C and C++ defini-
> tions. 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. 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 *.

**Unofficial Draft for Comment Only**

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 functions generally follow the naming rules given. In some circumstances, the MPI function is related to a function defined already for MPI-1 with a name that does not follow the naming conventions. In this circumstance, the language neutral name is in analogy to the MPI 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 do reflect the naming rules and can differ from the C and Fortran names. Thus, the analogous name in C++ to the MPI name may be 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 `typedef`s are made publicly within appropriate classes. However, these declarations then become somewhat cumbersome, as with the following:

`{typedef MPI::Grequest::Query_function();` *(binding deprecated, see Section 15.2)*`}`

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++ `typedef`s, 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 16.3.4, 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 an MIMD style. The codes executed by each process need not be identical. The processes communicate via calls to MPI communication primitives. Typically, each process executes in its own address space, although shared-memory implementations of MPI are possible.

This document specifies the behavior of a parallel program assuming that only MPI calls are used. 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 12.4.

## 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 Section 8.3. 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 16.1.8 on page 506.

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 defines a way for users to create new error codes as defined in Section 8.5.

## 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.

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

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

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

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

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

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

### 2.9.2   Interaction with Signals

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

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

## 2.10   Examples

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

```c
#include "mpi.h"
int main( 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 if (myrank == 1)  /* 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, along with the rank of the sender, 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, Nonblocking communication is addressed next, followed by channel-like constructs and send-receive operations, ending with a description of the "dummy" process, MPI_PROC_NULL.

## 3.2   Blocking Send and Receive Operations

### 3.2.1   Blocking Send

The syntax of the blocking send operation is given below.

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

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

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

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

{void MPI::Comm::Send(const void* buf, int count, const
            MPI::Datatype& datatype, int dest, int tag) const*(binding
            deprecated, see Section 15.2)* }

The blocking semantics of this call are described in Section 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 in Table 3.1.

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

Table 3.1: Predefined MPI datatypes corresponding to Fortran datatypes

Possible values for this argument for C and the corresponding C types are listed in Table 3.2.

The datatypes MPI_BYTE and MPI_PACKED do not correspond to a Fortran or C datatype. A value of type MPI_BYTE consists of a byte (8 binary digits). A byte is uninterpreted and is different from a character. Different machines may have different representations for characters, or may use more than one byte to represent characters. On the other hand, a byte has the same binary value on all machines. The use of the type MPI_PACKED is explained in Section 4.2.

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

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

> *Rationale.* The datatypes MPI_C_BOOL, MPI_INT8_T, MPI_INT16_T, MPI_INT32_T, MPI_UINT8_T, MPI_UINT16_T, MPI_UINT32_T, MPI_C_COMPLEX,

| MPI datatype | C datatype |
|---|---|
| MPI_CHAR | char |
|  | (treated as printable character) |
| MPI_SHORT | signed short int |
| MPI_INT | signed int |
| MPI_LONG | signed long int |
| MPI_LONG_LONG_INT | signed long long int |
| MPI_LONG_LONG (as a synonym) | signed long long int |
| MPI_SIGNED_CHAR | signed char |
|  | (treated as integral value) |
| MPI_UNSIGNED_CHAR | unsigned char |
|  | (treated as integral value) |
| MPI_UNSIGNED_SHORT | unsigned short int |
| MPI_UNSIGNED | unsigned int |
| MPI_UNSIGNED_LONG | unsigned long int |
| MPI_UNSIGNED_LONG_LONG | unsigned long long int |
| MPI_FLOAT | float |
| MPI_DOUBLE | double |
| MPI_LONG_DOUBLE | long double |
| MPI_WCHAR | wchar_t |
|  | (defined in <stddef.h>) |
|  | (treated as printable character) |
| MPI_C_BOOL | _Bool |
| MPI_INT8_T | int8_t |
| MPI_INT16_T | int16_t |
| MPI_INT32_T | int32_t |
| MPI_INT64_T | int64_t |
| MPI_UINT8_T | uint8_t |
| MPI_UINT16_T | uint16_t |
| MPI_UINT32_T | uint32_t |
| MPI_UINT64_T | uint64_t |
| MPI_C_COMPLEX | float _Complex |
| MPI_C_FLOAT_COMPLEX (as a synonym) | float _Complex |
| MPI_C_DOUBLE_COMPLEX | double _Complex |
| MPI_C_LONG_DOUBLE_COMPLEX | long double _Complex |
| MPI_BYTE |  |
| MPI_PACKED |  |

Table 3.2: Predefined MPI datatypes corresponding to C datatypes

MPI_C_FLOAT_COMPLEX, MPI_C_DOUBLE_COMPLEX, and
MPI_C_LONG_DOUBLE_COMPLEX have no corresponding C++ bindings. This was
intentionally done to avoid potential collisions with the C preprocessor and names-
paced C++ names. C++ applications can use the C bindings with no loss of func-
tionality. (*End of rationale.*)

The datatypes MPI_AINT [and ], MPI_OFFSET , and MPI_COUNT correspond to the

ticket265.

ticket265.

| MPI datatype | C datatype | Fortran datatype | |
|---|---|---|---|
| MPI_AINT | MPI_Aint | INTEGER (KIND=MPI_ADDRESS_KIND) | |
| MPI_OFFSET | MPI_Offset | INTEGER (KIND=MPI_OFFSET_KIND) | |
| [ticket265.] MPI_COUNT | [ticket265.]MPI_Count | [ticket265.]INTEGER (KIND=MPI_COUNT_KIND) | |

Table 3.3: Predefined MPI datatypes corresponding to both C and Fortran datatypes

MPI-defined C types MPI_Aint [and ], MPI_Offset , and MPI_COUNT and their Fortran equivalents INTEGER (KIND=MPI_ADDRESS_KIND) [and ], INTEGER (KIND= MPI_OFFSET_KIND) , and INTEGER (KIND=MPI_COUNT_KIND) . This is described in Table 3.3. See Section 16.3.10 for information on interlanguage communication with these types.

### 3.2.3  Message Envelope

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

<div align="center">

source

destination

tag

communicator

</div>

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

The message destination is specified by the dest argument.

The integer-valued message tag is specified by the tag argument. This integer can be used by the program to distinguish different types of messages. The range of valid tag values is 0,...,UB, where the value of UB is implementation dependent. It can be found by querying the value of the attribute MPI_TAG_UB, as described in Chapter 8. MPI requires that UB be no less than 32767.

The comm argument specifies the **communicator** that is used for the send operation. Communicators are explained in Chapter 6; below is a brief summary of their usage.

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

The communicator also specifies the set of processes that share this communication 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 6.)

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 6. 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  | buf      | initial address of receive buffer (choice)            |
| IN   | count    | number of elements in receive buffer (non-negative integer) |
| IN   | datatype | datatype of each receive buffer element (handle)      |
| IN   | source   | rank of source or MPI_ANY_SOURCE (integer)            |
| IN   | tag      | message tag or MPI_ANY_TAG (integer)                  |
| IN   | comm     | communicator (handle)                                 |
| OUT  | status   | 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*(binding deprecated, see Section 15.2)* }

{void MPI::Comm::Recv(void* buf, int count, const MPI::Datatype& datatype, int source, int tag) const*(binding deprecated, see Section 15.2)* }

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

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

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

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

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

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

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

The message tag is specified by the tag argument of the receive operation. The argument source, if different from MPI_ANY_SOURCE, is specified as a rank within the process group associated with that same communicator (remote process group, for intercommunicators). Thus, the range of valid values for the source argument is {0,...,n-1}∪{MPI_ANY_SOURCE}, where n is the number of processes in this group.

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

Source = destination is allowed, that is, a process can send a message to itself. (However, it is unsafe to do so with the blocking send and receive operations described above, since this may lead to deadlock. See Section 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 INTEGERs 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 *(binding deprecated, see Section 15.2)* }

{void MPI::Status::Set_source(int source) *(binding deprecated, see Section 15.2)* }

{int MPI::Status::Get_tag() const *(binding deprecated, see Section 15.2)* }

{void MPI::Status::Set_tag(int tag) *(binding deprecated, see Section 15.2)* }

{int MPI::Status::Get_error() const *(binding deprecated, see Section 15.2)* }

{void MPI::Status::Set_error(int error) *(binding deprecated, see Section 15.2)* }

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.

MPI_GET_COUNT(status, datatype, count)

| IN | status | return status of receive operation (Status) |
|---|---|---|
| IN | datatype | datatype of each receive buffer entry (handle) |
| OUT | count | number of received entries (integer) |

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

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

{int MPI::Status::Get_count(const MPI::Datatype& datatype) const *(binding deprecated, see Section 15.2)* }

Returns the number of entries received. (Again, we count *entries*, each of type *datatype*, not *bytes*.) The datatype argument should match the argument provided by the receive call that set the status variable. (We shall later see, in Section 4.1.11, that MPI_GET_COUNT may return, in certain situations, the value MPI_UNDEFINED.)

> *Rationale.*     Some message-passing libraries use INOUT count, tag and source arguments, thus using them both to specify the selection criteria for incoming messages and return the actual envelope values of the received message. The use of a separate status argument prevents errors that are often attached with INOUT argument (e.g., using the MPI_ANY_TAG constant as the tag in a receive). Some libraries use calls that refer implicitly to the "last message received." This is not thread safe.
>
> The datatype argument is passed to MPI_GET_COUNT so as to improve performance. A message might be received without counting the number of elements it contains, and the count value is often not needed. Also, this allows the same function to be used after a call to MPI_PROBE or MPI_IPROBE. With a status from MPI_PROBE or MPI_IPROBE, the same datatypes are allowed as in a call to MPI_RECV to receive this message. (*End of rationale.*)

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 transfered is greater than zero, MPI_UNDEFINED is returned.

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

> *Advice to users.*   The buffer size required for the receive can be affected by data conversions and by the stride of the receive datatype. In most cases, the safest approach is to use the same datatype with MPI_GET_COUNT and the receive. (*End of advice to users.*)

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

### 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 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 `status`es 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 MPI_{TEST|WAIT}{ALL|SOME} 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 MPI_{TEST|WAIT}{ALL|SOME} 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. Data is pulled out of the send buffer and a message is assembled.

2. A message is transferred from sender to receiver.

3. Data is pulled from the incoming message and disassembled into the receive buffer.

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

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

The types of a send and receive match (phase two) if both operations use identical names. That is, MPI_INTEGER matches MPI_INTEGER, MPI_REAL matches MPI_REAL, and so on. There is one exception to this rule, discussed in Section 4.2, the type MPI_PACKED can match any other type.

The type of a variable in a host program matches the type specified in the communication operation if the datatype name used by that operation corresponds to the basic type of the host program variable. For example, an entry with type name MPI_INTEGER matches a Fortran variable of type INTEGER. A table giving this correspondence for Fortran and C appears in Section 3.2.2. There are two exceptions to this last rule: an entry with type name MPI_BYTE or MPI_PACKED can be used to match any byte of storage (on a byte-addressable machine), irrespective of the datatype of the variable that contains this byte. The type MPI_PACKED is used to send data that has been explicitly packed, or receive data that will be explicitly unpacked, see Section 4.2. The type MPI_BYTE allows one to transfer the binary value of a byte in memory unchanged.

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

- Communication of typed values (e.g., with datatype different from MPI_BYTE), where the datatypes of the corresponding entries in the sender program, in the send call, in the receive call and in the receiver program must all match.

- Communication of untyped values (e.g., of datatype MPI_BYTE), where both sender and receiver use the datatype MPI_BYTE. In this case, there are no requirements on the types of the corresponding entries in the sender and the receiver programs, nor is it required that they be the same.

- Communication involving packed data, where MPI_PACKED is used.

The following examples illustrate the first two cases.

**Example 3.2**   Sender and receiver 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 IF (rank.EQ.1) THEN
    CALL MPI_RECV(b(1), 15, MPI_REAL, 0, tag, comm, status, ierr)
END IF
```

This code is correct if both a and b are real arrays of size $\geq 10$. (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.)

**Unofficial Draft for Comment Only**

**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 IF (rank.EQ.1) THEN
    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 IF (rank.EQ.1) THEN
    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 then 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 CHARACTERs.

```
CHARACTER*10 a
CHARACTER*10 b

CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(a, 5, MPI_CHARACTER, 1, tag, comm, ierr)
ELSE IF (rank.EQ.1) THEN
    CALL MPI_RECV(b(6:10), 5, MPI_CHARACTER, 0, tag, comm, status, ierr)
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 4.1). 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 16.3 on page 527.

## 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 modify the send buffer. The message might be copied directly into the matching receive buffer, or it might be copied into a temporary system buffer.

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

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

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

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

There are three additional communication modes.

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

A send that uses the **synchronous** mode can be started whether or not a matching receive was posted. However, the send will complete successfully only if a matching receive is posted, and the receive operation has started to receive the message sent by the synchronous send. Thus, the completion of a synchronous send not only indicates that the send buffer can be reused, but it also indicates that the receiver has reached a certain point in its

execution, namely that it has started executing the matching receive. If both sends and receives are blocking operations then the use of the synchronous mode provides synchronous communication semantics: a communication does not complete at either end before both 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 (non-negative 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*(binding
              deprecated, see Section 15.2)* }

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 (non-negative 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,
              int tag, MPI_Comm comm)
```

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

{void MPI::Comm::Ssend(const void* buf, int count, const
            MPI::Datatype& datatype, int dest, int tag) const*(binding
            deprecated, see Section 15.2)* }

Send in synchronous mode.

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

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

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

{void MPI::Comm::Rsend(const void* buf, int count, const
            MPI::Datatype& datatype, int dest, int tag) const*(binding
            deprecated, see Section 15.2)* }

Send in ready mode.

There is only one receive operation, but it matches 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 modify a communication buffer until the communication completes. Otherwise, the outcome of the computation is undefined.

> *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, users may 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.

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

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_BSEND(buf1, count, MPI_REAL, 1, tag, comm, ierr)
    CALL MPI_BSEND(buf2, count, MPI_REAL, 1, tag, comm, ierr)
ELSE IF (rank.EQ.1) THEN
    CALL MPI_RECV(buf1, count, MPI_REAL, 0, MPI_ANY_TAG, comm, status, ierr)
    CALL MPI_RECV(buf2, count, MPI_REAL, 0, tag, comm, status, ierr)
END IF
```

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

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

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

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_BSEND(buf1, count, MPI_REAL, 1, tag1, comm, ierr)
    CALL MPI_SSEND(buf2, count, MPI_REAL, 1, tag2, comm, ierr)
ELSE IF (rank.EQ.1) THEN
    CALL MPI_RECV(buf1, count, MPI_REAL, 0, tag2, comm, status, ierr)
    CALL MPI_RECV(buf2, count, MPI_REAL, 0, tag1, comm, status, ierr)
END IF
```

Both processes invoke their first communication call. Since the first send of process zero uses the buffered mode, it must complete, irrespective of the state of process one. Since no matching receive is posted, the message will be copied into buffer space. (If insufficient buffer space is available, then the program will fail.) The second send is then invoked. At that point, a matching pair of send and receive operation is enabled, and both operations must complete. Process one next invokes its second receive call, which will be satisfied by

the buffered message. Note that process one received the messages in the reverse order they were sent.

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

**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 IF (rank.EQ.1) THEN
    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 errant 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 IF (rank.EQ.1) THEN
    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
```

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.

```
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 IF (rank.EQ.1) THEN
    CALL MPI_SEND(sendbuf, count, MPI_REAL, 0, tag, comm, ierr)
    CALL MPI_RECV(recvbuf, count, MPI_REAL, 0, tag, comm, status, ierr)
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 on the communication protocol used.

Many programmers prefer to have more leeway and opt 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 (non-negative 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)*(binding deprecated, see Section 15.2)* }

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 (non-negative 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)*(binding deprecated, see Section 15.2)* }

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 */
```

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

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

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

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

MPI does not provide mechanisms for querying or controlling buffering done by standard mode sends. It is expected that vendors will provide such information for their implementations.

> *Rationale.*   There is a wide spectrum of possible implementations of buffered communication: buffering can be done at sender, at receiver, or both; buffers can be dedicated to one sender-receiver pair, or be shared by all communications; buffering can be done in real or in virtual memory; it can use dedicated memory, or memory

shared by other processes; buffer space may be allocated statically or be changed dynamically; etc. It does not seem feasible to provide a portable mechanism for querying or controlling buffering that would be compatible with all these choices, yet provide meaningful information. (*End of rationale.*)

### 3.6.1  Model Implementation of Buffered Mode

The model implementation uses the packing and unpacking functions described in Section 4.2 and the nonblocking communication functions described in Section 3.7.

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

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

- Traverse sequentially the PME queue from head towards the tail, deleting all entries for communications that have completed, up to the first entry with an uncompleted 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 4.2). 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.

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

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

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

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

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

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

Nonblocking sends in the buffered and ready modes have a more limited impact, e.g., the blocking version of buffered send is capable of completing regardless of when a matching receive call is made. However, separating the start from the completion of these sends still gives some opportunity for optimization within the MPI library. For example, starting a buffered send gives an implementation more flexibility in determining if and how the message is buffered. There are also advantages for both nonblocking buffered and ready modes when data copying can be done concurrently with computation.

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

### 3.7.1  Communication Request Objects

Nonblocking communications use opaque request objects to identify communication operations and match the operation that initiates the communication with the operation that terminates it. These are system objects that are accessed via a handle. A request object identifies various properties of a communication operation, such as the send mode, the communication buffer that is associated with it, its context, the tag and destination arguments to be used for a send, or the tag and source arguments to be used for a receive. In addition, this object stores information about the status of the pending communication operation.

### 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 (non-negative 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*(binding
            deprecated, see Section 15.2)* }

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 (non-negative 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
```

{MPI::Request MPI::Comm::Ibsend(const void* buf, int count, const
          MPI::Datatype& datatype, int dest, int tag) const*(binding
          deprecated, see Section 15.2)* }

Start a buffered mode, nonblocking send.

MPI_ISSEND(buf, count, datatype, dest, tag, comm, request)

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

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

{MPI::Request MPI::Comm::Issend(const void* buf, int count, const
          MPI::Datatype& datatype, int dest, int tag) const*(binding
          deprecated, see Section 15.2)* }

Start a synchronous mode, nonblocking send.

MPI_IRSEND(buf, count, datatype, dest, tag, comm, request)

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

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

{MPI::Request MPI::Comm::Irsend(const void* buf, int count, const
              MPI::Datatype& datatype, int dest, int tag) const*(binding
              deprecated, see Section 15.2)* }

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 (non-negative integer) |
| IN | datatype | datatype of each receive buffer element (handle) |
| IN | source | rank of source or MPI_ANY_SOURCE (integer) |
| IN | tag | message tag or MPI_ANY_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*(binding
              deprecated, see Section 15.2)* }

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 modify 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 16.2.2 on pages 512 and 515. (*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, and MPI_WAITANY, which 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) *(binding deprecated, see Section 15.2)* }

{void MPI::Request::Wait() *(binding deprecated, see Section 15.2)* }

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 Section 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) *(binding deprecated, see Section 15.2)* }

{bool MPI::Request::Test() *(binding deprecated, see Section 15.2)* }

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.*)

**Example 3.12**  Simple usage of nonblocking operations and MPI_WAIT.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_ISEND(a(1), 10, MPI_REAL, 1, tag, comm, request, ierr)
    **** do some computation to mask latency ****
    CALL MPI_WAIT(request, status, ierr)
ELSE IF (rank.EQ.1) THEN
    CALL MPI_IRECV(a(1), 15, MPI_REAL, 0, tag, comm, request, ierr)
    **** do some computation to mask latency ****
    CALL MPI_WAIT(request, status, ierr)
END IF
```

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

MPI_REQUEST_FREE(request)

| | | |
|---|---|---|
| INOUT | request | communication request (handle) |

```
int MPI_Request_free(MPI_Request *request)
```

```
MPI_REQUEST_FREE(REQUEST, IERROR)
    INTEGER REQUEST, IERROR
```

{void MPI::Request::Free() *(binding deprecated, see Section 15.2)* }

**Unofficial Draft for Comment Only**

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

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

> *Advice to users.*   Once a request is freed by a call to MPI_REQUEST_FREE, it is not possible to check for the successful completion of the associated communication with calls to MPI_WAIT or MPI_TEST. Also, if an error occurs subsequently during the communication, an error code cannot be returned to the user — such an error must be treated as fatal. An active receive request should never be freed as the receiver will have no way to verify that the receive has completed and the receive buffer can be reused. (*End of advice to 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(inval, 1, MPI_REAL, 1, 0, MPI_COMM_WORLD, req, ierr)
      CALL MPI_WAIT(req, status, ierr)
    END DO
ELSE IF (rank.EQ.1) THEN
    CALL MPI_IRECV(inval, 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(inval, 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 IF (rank.EQ.1) THEN
      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.

```
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (RANK.EQ.0) THEN
      CALL MPI_SSEND(a, 1, MPI_REAL, 1, 0, comm, ierr)
      CALL MPI_SEND(b, 1, MPI_REAL, 1, 1, comm, ierr)
ELSE IF (rank.EQ.1) THEN
      CALL MPI_IRECV(a, 1, MPI_REAL, 0, 0, comm, r, ierr)
      CALL MPI_RECV(b, 1, MPI_REAL, 0, 1, comm, status, ierr)
      CALL MPI_WAIT(r, status, ierr)
END IF
```

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_WAITSOME or MPI_TESTSOME can be used to complete all enabled operations in a list.

MPI_WAITANY (count, array_of_requests, index, status)

| IN | count | list length (non-negative 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) |

```
int MPI_Waitany(int count, MPI_Request *array_of_requests, int *index,
            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)*(binding deprecated, see Section 15.2)* }

{static int MPI::Request::Waitany(int count,
            MPI::Request array_of_requests[])*(binding deprecated, see Section 15.2)* }

Blocks until one of the operations associated with the active requests in the array has completed. If more then 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 a 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 (non-negative integer) |
| INOUT | array_of_requests | array of requests (array of handles) |
| OUT | index | index of operation that completed, or MPI_UNDEFINED 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) *(binding deprecated, see Section 15.2)* }

{static bool MPI::Request::Testany(int count,
            MPI::Request array_of_requests[], int& index) *(binding deprecated, see Section 15.2)* }

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

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

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

MPI_WAITALL( count, array_of_requests, array_of_statuses)

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

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

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

```
{static void MPI::Request::Waitall(int count,
                MPI::Request array_of_requests[],
                MPI::Status array_of_statuses[])(binding deprecated, see
                Section 15.2) }
```

```
{static void MPI::Request::Waitall(int count,
                MPI::Request array_of_requests[])(binding deprecated, see
                Section 15.2) }
```

Blocks until all communication operations associated with active handles in the list complete, and return the status of all these operations (this includes the case where no handle in the list is active). Both arrays have the same number of valid entries. The i-th entry in array_of_statuses is set to the return status of the i-th operation. Requests that were created by nonblocking communication operations are deallocated and the corresponding handles in the array are set to MPI_REQUEST_NULL. The list may contain null or inactive 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 ,..., 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 desireable 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 (non-negative 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[])*(binding deprecated, see
            Section 15.2)* }

{static bool MPI::Request::Testall(int count,
            MPI::Request array_of_requests[])*(binding deprecated, see
            Section 15.2)* }

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_WAITSOME(incount, array_of_requests, outcount, array_of_indices, array_of_statuses)

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

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

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

```
{static int MPI::Request::Waitsome(int incount,
            MPI::Request array_of_requests[], int array_of_indices[],
            MPI::Status array_of_statuses[])(binding deprecated, see
            Section 15.2) }

{static int MPI::Request::Waitsome(int incount,
            MPI::Request array_of_requests[],
            int array_of_indices[])(binding deprecated, see Section 15.2) }
```

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_WAITSOME 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)

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

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

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

**Unofficial Draft for Comment Only**

```
{static int MPI::Request::Testsome(int incount,
            MPI::Request array_of_requests[], int array_of_indices[],
            MPI::Status array_of_statuses[])(binding deprecated, see
            Section 15.2) }

{static int MPI::Request::Testsome(int incount,
            MPI::Request array_of_requests[],
            int array_of_indices[])(binding deprecated, see Section 15.2) }
```

Behaves like MPI_WAITSOME, 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_TESTSOME is a local operation, which returns immediately, whereas MPI_WAITSOME 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_WAITSOME or MPI_TESTSOME, 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_TESTSOME are handled as for MPI_WAITSOME.

> *Advice to users.* The use of MPI_TESTSOME 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_WAITSOME so as not to starve any client. Clients send messages to the server with service requests. The server calls MPI_WAITSOME 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_TESTSOME should complete as many pending communications as possible. (*End of advice to implementors.*)

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

```
CALL MPI_COMM_SIZE(comm, size, ierr)
CALL MPI_COMM_RANK(comm, rank, ierr)
IF(rank .GT. 0) THEN          ! client code
    DO WHILE(.TRUE.)
        CALL MPI_ISEND(a, n, MPI_REAL, 0, tag, comm, request, ierr)
        CALL MPI_WAIT(request, status, ierr)
    END DO
ELSE          ! rank=0 -- server code
    DO i=1, size-1
        CALL MPI_IRECV(a(1,i), n, MPI_REAL, i, tag,
                comm, request_list(i), ierr)
    END DO
    DO WHILE(.TRUE.)
```

```
        CALL MPI_WAITANY(size-1, request_list, index, status, ierr)
        CALL DO_SERVICE(a(1,index))  ! handle one message
        CALL MPI_IRECV(a(1, index), n, MPI_REAL, index, tag,
                    comm, request_list(index), ierr)
    END DO
END IF
```

**Example 3.17**    Same code, using MPI_WAITSOME.

```
CALL MPI_COMM_SIZE(comm, size, ierr)
CALL MPI_COMM_RANK(comm, rank, ierr)
IF(rank .GT. 0) THEN          ! client code
    DO WHILE(.TRUE.)
        CALL MPI_ISEND(a, n, MPI_REAL, 0, tag, comm, request, ierr)
        CALL MPI_WAIT(request, status, ierr)
    END DO
ELSE          ! rank=0 -- server code
    DO i=1, size-1
        CALL MPI_IRECV(a(1,i), n, MPI_REAL, i, tag,
                    comm, request_list(i), ierr)
    END DO
    DO WHILE(.TRUE.)
        CALL MPI_WAITSOME(size, request_list, numdone,
                        indices, statuses, ierr)
        DO i=1, numdone
            CALL DO_SERVICE(a(1, indices(i)))
            CALL MPI_IRECV(a(1, indices(i)), n, MPI_REAL, 0, tag,
                        comm, request_list(indices(i)), ierr)
        END DO
    END DO
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 MPI_TEST (logical) |
| OUT | status | MPI_STATUS 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*(binding deprecated, see Section 15.2)* }

{bool MPI::Request::Get_status() const*(binding deprecated, see Section 15.2)* }

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.

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

## 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.

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

| | | |
|-----|--------|-------------------------------------------|
| IN | source | rank of source or MPI_ANY_SOURCE (integer) |
| IN | tag | message tag or MPI_ANY_TAG (integer) |
| IN | comm | communicator (handle) |
| OUT | flag | (logical) |
| OUT | status | status object (Status) |

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

```
MPI_IPROBE(SOURCE, TAG, COMM, FLAG, STATUS, IERROR)
    LOGICAL FLAG
    INTEGER SOURCE, TAG, COMM, STATUS(MPI_STATUS_SIZE), IERROR
```

{bool MPI::Comm::Iprobe(int source, int tag, MPI::Status& status) const*(binding deprecated, see Section 15.2)* }

{bool MPI::Comm::Iprobe(int source, int tag) const*(binding deprecated, see Section 15.2)* }

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

If MPI_IPROBE returns flag = true, then the content of the status object can be subsequently accessed as described in Section 3.2.5 to find the source, tag and length of the probed message.

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

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

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

MPI_PROBE(source, tag, comm, status)

| | | |
|---|---|---|
| IN | source | rank of source or MPI_ANY_SOURCE (integer) |
| IN | tag | message tag or MPI_ANY_TAG (integer) |
| IN | comm | communicator (handle) |
| OUT | status | status object (Status) |

```
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*(binding deprecated, see Section 15.2)* }

{void MPI::Comm::Probe(int source, int tag) const*(binding deprecated, see Section 15.2)* }

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 IF (rank.EQ.2) THEN
          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 IF (rank.EQ.2) THEN
          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, MPI_ANY_SOURCE,
                                 0, comm, status, ierr)
             ELSE
200              CALL MPI_RECV(x, 1, MPI_REAL, MPI_ANY_SOURCE,
                                 0, comm, status, ierr)
             END IF
          END DO
      END IF
```

    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*(binding deprecated, see Section 15.2)* }

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 Section 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.

> *Rationale.* Although the IN request handle parameter should not need to be passed by reference, the C binding has listed the argument type as MPI_Request* since MPI-1.0. This function signature therefore cannot be changed without breaking existing MPI applications. (*End of rationale.*)

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*(binding deprecated, see Section 15.2)* }

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 (non-negative 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*(binding
            deprecated, see Section 15.2)* }

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

MPI_BSEND_INIT(buf, count, datatype, dest, tag, comm, request)

| | | |
|---|---|---|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements sent (non-negative 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_Bsend_init(void* buf, int count, MPI_Datatype datatype, int dest,
            int tag, MPI_Comm comm, MPI_Request *request)
```

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

{MPI::Prequest MPI::Comm::Bsend_init(const void* buf, int count, const
            MPI::Datatype& datatype, int dest, int tag) const*(binding
            deprecated, see Section 15.2)* }

Creates a persistent communication request for a buffered mode send.


MPI_SSEND_INIT(buf, count, datatype, dest, tag, comm, request)

| | | |
|----|----------|--------------------------------------------|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements sent (non-negative 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_Ssend_init(void* buf, int count, MPI_Datatype datatype, int dest,
            int tag, MPI_Comm comm, MPI_Request *request)
```

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

{MPI::Prequest MPI::Comm::Ssend_init(const void* buf, int count, const
            MPI::Datatype& datatype, int dest, int tag) const*(binding
            deprecated, see Section 15.2)* }

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


MPI_RSEND_INIT(buf, count, datatype, dest, tag, comm, request)

| | | |
|----|----------|--------------------------------------------|
| IN | buf | initial address of send buffer (choice) |
| IN | count | number of elements sent (non-negative 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_Rsend_init(void* buf, int count, MPI_Datatype datatype, int dest,
                 int tag, MPI_Comm comm, MPI_Request *request)
```

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

{MPI::Prequest MPI::Comm::Rsend_init(const void* buf, int count, const
          MPI::Datatype& datatype, int dest, int tag) const*(binding
          deprecated, see Section 15.2)* }

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


MPI_RECV_INIT(buf, count, datatype, source, tag, comm, request)

| | | |
|---|---|---|
| OUT | buf | initial address of receive buffer (choice) |
| IN | count | number of elements received (non-negative integer) |
| IN | datatype | type of each element (handle) |
| IN | source | rank of source or MPI_ANY_SOURCE (integer) |
| IN | tag | message tag or MPI_ANY_TAG (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

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

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

{MPI::Prequest MPI::Comm::Recv_init(void* buf, int count, const
          MPI::Datatype& datatype, int source, int tag) const*(binding
          deprecated, see Section 15.2)* }

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 com-
munication 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() *(binding deprecated, see Section 15.2)* }

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 modified 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 (non-negative 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[]) *(binding deprecated, see Section 15.2)* }

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

**Create (Start Complete)* Free**

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

**Unofficial Draft for Comment Only**

correct sequence is obeyed.

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

> *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 16.2.2 on pages 512 and 515. (*End of advice to users.*)

## 3.10   Send-Receive

The **send-receive** operations combine in one call the sending of a message to one destination and the receiving of another message, from another process. The two (source and destination) are possibly the same.  A send-receive operation is very useful for executing a shift operation across a chain of processes.  If blocking sends and receives are used for such a shift, then one needs to order the sends and receives correctly (for example, even processes send, then receive, odd processes receive first, then send) so as to prevent cyclic dependencies that may lead to deadlock.  When a send-receive operation is used, the communication subsystem takes care of these issues.  The send-receive operation can be used in conjunction with the functions described in Chapter 7 in order to perform shifts on various logical topologies.  Also, a send-receive operation is useful for implementing remote procedure calls.

A message sent by a send-receive operation can be received by a regular receive operation or probed by a probe operation; a send-receive operation can receive a message sent by a regular send operation.

MPI_SENDRECV(sendbuf, sendcount, sendtype, dest, sendtag, recvbuf, recvcount, recvtype, source, recvtag, comm, status)

| | | |
|------|----------|----------------------------------------------------|
| IN   | sendbuf  | initial address of send buffer (choice)            |
| IN   | sendcount| number of elements in send buffer (non-negative 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 (non-negative integer) |
| IN   | recvtype | type of elements in receive buffer (handle)        |
| IN   | source   | rank of source or MPI_ANY_SOURCE (integer)         |
| IN   | recvtag  | receive tag or MPI_ANY_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*(binding deprecated, see
            Section 15.2)* }

{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*(binding deprecated, see Section 15.2)* }

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.

MPI_SENDRECV_REPLACE(buf, count, datatype, dest, sendtag, source, recvtag, comm, status)

| | | |
|---|---|---|
| INOUT | buf | initial address of send and receive buffer (choice) |
| IN | count | number of elements in send and receive buffer (non-negative integer) |
| IN | datatype | type of elements in send and receive buffer (handle) |
| IN | dest | rank of destination (integer) |
| IN | sendtag | send message tag (integer) |
| IN | source | rank of source or MPI_ANY_SOURCE (integer) |
| IN | recvtag | receive message tag or MPI_ANY_TAG (integer) |
| IN | comm | communicator (handle) |
| OUT | status | status object (Status) |

```
int MPI_Sendrecv_replace(void* buf, int count, MPI_Datatype datatype,
            int dest, int sendtag, int source, int recvtag, MPI_Comm comm,
            MPI_Status *status)
```

```
MPI_SENDRECV_REPLACE(BUF, COUNT, DATATYPE, DEST, SENDTAG, SOURCE, RECVTAG,
                COMM, STATUS, IERROR)
    <type> BUF(*)
    INTEGER COUNT, DATATYPE, DEST, SENDTAG, SOURCE, RECVTAG, COMM,
    STATUS(MPI_STATUS_SIZE), IERROR
```

{void MPI::Comm::Sendrecv_replace(void* buf, int count, const
                MPI::Datatype& datatype, int dest, int sendtag, int source,
                int recvtag, MPI::Status& status) const*(binding deprecated, see
                Section 15.2)* }

{void MPI::Comm::Sendrecv_replace(void* buf, int count, const
                MPI::Datatype& datatype, int dest, int sendtag, int source,
                int recvtag) const*(binding deprecated, see Section 15.2)* }

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.

# Chapter 4

# Datatypes

Basic datatypes were introduced in Section 3.2.2 Message Data on page 27 and in Section 3.3 Data Type Matching and Data Conversion on page 34. In this chapter, this model is extended to describe any data layout. We 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.

## 4.1 Derived Datatypes

Up to here, all point to point communication have involved only buffers containing a sequence of identical basic datatypes. 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 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

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

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

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

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 non-negative 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 Section 4.1.11, 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

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

then

$$
\begin{aligned}
lb(Typemap) &= \min_j disp_j, \\
ub(Typemap) &= \max_j(disp_j + sizeof(type_j)) + \epsilon, \text{ and} \\
extent(Typemap) &= ub(Typemap) - lb(Typemap).
\end{aligned}
\tag{4.1}
$$

If $type_i$ requires alignment to a byte address that is a multiple of $k_i$, then $\epsilon$ is the least non-negative increment needed to round $extent(Typemap)$ to the next multiple of $\max_i k_i$. The complete definition of **extent** is given on page 97.

**Unofficial Draft for Comment Only**

**Example 4.1** Assume that $Type = \{(\text{double}, 0), (\text{char}, 8)\}$ (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 4.1.6. Such explicit control is needed in cases where the assumption does not hold, for example, where union types are used. (*End of rationale.*)

### 4.1.1 Type Constructors with Explicit Addresses

In Fortran, the functions MPI_TYPE_CREATE_HVECTOR, MPI_TYPE_CREATE_HINDEXED, MPI_TYPE_CREATE_STRUCT, and MPI_GET_ADDRESS accept arguments of type INTEGER(KIND=MPI_ADDRESS_KIND), wherever arguments of type MPI_Aint and MPI::Aint are used in C and C++. On Fortran 77 systems that do not support the Fortran 90 KIND notation, and where addresses are 64 bits whereas default INTEGERs are 32 bits, these arguments will be of type INTEGER*8.

### 4.1.2 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 (non-negative 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*(binding deprecated, see Section 15.2)* }

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

**Example 4.2** Let oldtype have type map $\{(\text{double}, 0), (\text{char}, 8)\}$, with extent 16, and let count $= 3$. The type map of the datatype returned by newtype is

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

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

**Unofficial Draft for Comment Only**

In general, assume that the type map of oldtype is

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

with extent $ex$. Then newtype has a type map with count $\cdot$ n entries defined by:

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

$$..., (type_0, disp_0 + ex \cdot (\text{count} - 1)), ..., (type_{n-1}, disp_{n-1} + ex \cdot (\text{count} - 1))\}.$$

Vector    The function MPI_TYPE_VECTOR is a more general constructor that allows replication of a datatype into locations that consist of equally spaced blocks. Each block is obtained by concatenating the same number of copies of the old datatype. The spacing between blocks is a multiple of the extent of the old datatype.

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

| | | |
|---|---|---|
| IN | count | number of blocks (non-negative integer) |
| IN | blocklength | number of elements in each block (non-negative integer) |
| IN | stride | number of elements between start of each block (integer) |
| IN | oldtype | old datatype (handle) |
| OUT | newtype | new datatype (handle) |

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

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

```
{MPI::Datatype MPI::Datatype::Create_vector(int count, int blocklength,
              int stride) const(binding deprecated, see Section 15.2) }
```

**Example 4.3** Assume, again, that oldtype has type map $\{(\text{double}, 0), (\text{char}, 8)\}$, with extent 16. A call to MPI_TYPE_VECTOR( 2, 3, 4, oldtype, newtype) will create the datatype with type map,

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

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

That is, two blocks with three copies each of the old type, with a stride of 4 elements ($4 \cdot 16$ bytes) between the blocks.

**Example 4.4** A call to MPI_TYPE_VECTOR(3, 1, -2, oldtype, newtype) will create the datatype,

$$\{(\mathsf{double}, 0), (\mathsf{char}, 8), (\mathsf{double}, -32), (\mathsf{char}, -24), (\mathsf{double}, -64), (\mathsf{char}, -56)\}.$$

In general, assume that oldtype has type map,

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

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

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

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

$$(type_0, disp_0 + (\mathsf{bl} - 1) \cdot ex), ..., (type_{n-1}, disp_{n-1} + (\mathsf{bl} - 1) \cdot ex),$$

$$(type_0, disp_0 + \mathsf{stride} \cdot ex), ..., (type_{n-1}, disp_{n-1} + \mathsf{stride} \cdot ex), ...,$$

$$(type_0, disp_0 + (\mathsf{stride} + \mathsf{bl} - 1) \cdot ex), ..., (type_{n-1}, disp_{n-1} + (\mathsf{stride} + \mathsf{bl} - 1) \cdot ex), ....,$$

$$(type_0, disp_0 + \mathsf{stride} \cdot (\mathsf{count} - 1) \cdot ex), ...,$$

$$(type_{n-1}, disp_{n-1} + \mathsf{stride} \cdot (\mathsf{count} - 1) \cdot ex), ...,$$

$$(type_0, disp_0 + (\mathsf{stride} \cdot (\mathsf{count} - 1) + \mathsf{bl} - 1) \cdot ex), ...,$$

$$(type_{n-1}, disp_{n-1} + (\mathsf{stride} \cdot (\mathsf{count} - 1) + \mathsf{bl} - 1) \cdot ex)\}.$$

A call to MPI_TYPE_CONTIGUOUS(count, oldtype, newtype) is equivalent to a call to MPI_TYPE_VECTOR(count, 1, 1, oldtype, newtype), or to a call to MPI_TYPE_VECTOR(1, count, n, oldtype, newtype), n arbitrary.

Hvector    The function MPI_TYPE_CREATE_HVECTOR is identical to MPI_TYPE_VECTOR, except that stride is given in bytes, rather than in elements. The use for both types of vector constructors is illustrated in Section 4.1.14. (H stands for "heterogeneous").

MPI_TYPE_CREATE_HVECTOR( count, blocklength, stride, oldtype, newtype)

| | | |
|---|---|---|
| IN | count | number of blocks (non-negative integer) |
| IN | blocklength | number of elements in each block (non-negative 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_create_hvector(int count, int blocklength, MPI_Aint stride,
               MPI_Datatype oldtype, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_HVECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE,
               IERROR)
    INTEGER COUNT, BLOCKLENGTH, OLDTYPE, NEWTYPE, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) STRIDE
```

{MPI::Datatype MPI::Datatype::Create_hvector(int count, int blocklength,
               MPI::Aint stride) const *(binding deprecated, see Section 15.2)* }

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

Assume that oldtype has type map,

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

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

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

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

$$(type_0, disp_0 + (\text{bl} - 1) \cdot ex), ..., (type_{n-1}, disp_{n-1} + (\text{bl} - 1) \cdot ex),$$

$$(type_0, disp_0 + \text{stride}), ..., (type_{n-1}, disp_{n-1} + \text{stride}), ...,$$

$$(type_0, disp_0 + \text{stride} + (\text{bl} - 1) \cdot ex), ...,$$

$$(type_{n-1}, disp_{n-1} + \text{stride} + (\text{bl} - 1) \cdot ex), ....,$$

$$(type_0, disp_0 + \text{stride} \cdot (\text{count} - 1)), ..., (type_{n-1}, disp_{n-1} + \text{stride} \cdot (\text{count} - 1)), ...,$$

$$(type_0, disp_0 + \text{stride} \cdot (\text{count} - 1) + (\text{bl} - 1) \cdot ex), ...,$$

$$(type_{n-1}, disp_{n-1} + \text{stride} \cdot (\text{count} - 1) + (\text{bl} - 1) \cdot ex)\}.$$

**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, new-
type)

| | | |
|---|---|---|
| IN | count | number of blocks – also number of entries in array_of_displacements and array_of_blocklengths (non-negative integer) |
| IN | array_of_blocklengths | number of elements per block (array of non-negative integers) |
| IN | array_of_displacements | displacement for each block, in multiples of oldtype 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(binding deprecated, see
            Section 15.2) }
```

**Example 4.5**
Let oldtype have type map $\{(\text{double}, 0), (\text{char}, 8)\}$, with extent 16. Let B = (3, 1) and let D = (4, 0). A call to MPI_TYPE_INDEXED(2, B, D, oldtype, newtype) returns a datatype with type map,

$$\{(\text{double}, 64), (\text{char}, 72), (\text{double}, 80), (\text{char}, 88), (\text{double}, 96), (\text{char}, 104),$$

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

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,

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

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

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

$$(type_0, disp_0 + (D[0] + B[0] - 1) \cdot ex), ..., (type_{n-1}, disp_{n-1} + (D[0] + B[0] - 1) \cdot ex), ...,$$

$$(type_0, disp_0 + \text{D[count-1]} \cdot ex), ..., (type_{n-1}, disp_{n-1} + \text{D[count-1]} \cdot ex), ...,$$

$$(type_0, disp_0 + (\text{D[count-1]} + \text{B[count-1]} - 1) \cdot ex), ...,$$

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

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

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

and

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

**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, oldtype, newtype)

| | | |
|---|---|---|
| IN | count | number of blocks — also number of entries in array_of_displacements and array_of_blocklengths (non-negative integer) |
| IN | array_of_blocklengths | number of elements in each block (array of non-negative integers) |
| IN | array_of_displacements | byte displacement of each block (array of integer) |
| IN | oldtype | old datatype (handle) |
| 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 *(binding deprecated, see Section 15.2)* }

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

Assume that oldtype has type map,

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

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 $n \cdot \sum_{i=0}^{\text{count}-1} B[i]$ entries:

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

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

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

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

$$(type_0, disp_0 + D[\text{count-1}] + (B[\text{count-1}] - 1) \cdot ex), ...,$$

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

Indexed_block   This function is the same as MPI_TYPE_INDEXED except that the blocklength 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 (non-negative integer) |
| IN | blocklength | size of block (non-negative 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*(binding deprecated, see Section 15.2)* }

**Unofficial Draft for Comment Only**

**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.

MPI_TYPE_CREATE_STRUCT(count, array_of_blocklengths, array_of_displacements,
                 array_of_types, newtype)

| | | |
|---|---|---|
| IN | count | number of blocks (non-negative integer) — also number of entries in arrays array_of_types, array_of_displacements and array_of_blocklengths |
| IN | array_of_blocklength | number of elements in each block (array of non-negative integer) |
| IN | array_of_displacements | byte displacement of each block (array of integer) |
| IN | array_of_types | type of elements in each block (array of handles to datatype objects) |
| OUT | newtype | new datatype (handle) |

```
int MPI_Type_create_struct(int count, int array_of_blocklengths[],
            MPI_Aint array_of_displacements[],
            MPI_Datatype array_of_types[], MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS,
            ARRAY_OF_DISPLACEMENTS, ARRAY_OF_TYPES, NEWTYPE, IERROR)
    INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_TYPES(*), NEWTYPE,
    IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_DISPLACEMENTS(*)
```

```
{static MPI::Datatype MPI::Datatype::Create_struct(int count,
            const int array_of_blocklengths[], const MPI::Aint
            array_of_displacements[],
            const MPI::Datatype array_of_types[])
```
*(binding deprecated, see Section 15.2)* }

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

**Example 4.6** 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 = (MPI_FLOAT, type1, 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,

$$typemap_i = \{(type_0^i, disp_0^i), ..., (type_{n_i-1}^i, disp_{n_i-1}^i)\},$$

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:

$$\{(type_0^0, disp_0^0 + D[0]), ..., (type_{n_0}^0, disp_{n_0}^0 + D[0]), ...,$$

$$(type_0^0, disp_0^0 + D[0] + (B[0]-1) \cdot ex_0), ..., (type_{n_0}^0, disp_{n_0}^0 + D[0] + (B[0]\text{-}1) \cdot ex_0), ...,$$

$$(type_0^{c-1}, disp_0^{c-1} + D[c\text{-}1]), ..., (type_{n_{c-1}-1}^{c-1}, disp_{n_{c-1}-1}^{c-1} + D[c\text{-}1]), ...,$$

$$(type_0^{c-1}, disp_0^{c-1} + D[c\text{-}1] + (B[c\text{-}1]-1) \cdot ex_{c-1}), ...,$$

$$(type_{n_{c-1}-1}^{c-1}, disp_{n_{c-1}-1}^{c-1} + D[c\text{-}1] + (B[c\text{-}1]\text{-}1) \cdot ex_{c-1})\}.$$

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

## 4.1.3 Subarray Datatype Constructor

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

| | | |
|----|----|----|
| IN | ndims | number of array dimensions (positive integer) |
| IN | array_of_sizes | number of elements of type oldtype in each dimension of the full array (array of positive integers) |
| IN | array_of_subsizes | number of elements of type oldtype in each dimension of the subarray (array of positive integers) |
| IN | array_of_starts | starting coordinates of the subarray in each dimension (array of non-negative integers) |
| IN | order | array storage order flag (state) |
| IN | oldtype | array element datatype (handle) |
| OUT | newtype | 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(binding deprecated,
            see Section 15.2) }
```

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, see MPI I/O, especially Section 13.1.1 on page 419.

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:

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

Let the typemap of oldtype have the form:

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

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 4.2 defines the base step. Equation 4.3 defines the recursion step when order = MPI_ORDER_FORTRAN, and Equation 4.4 defines the recursion step when order = MPI_ORDER_C.

$$\text{Subarray}(1, \{size_0\}, \{subsize_0\}, \{start_0\}, \tag{4.2}$$
$$\{(type_0, disp_0), (type_1, disp_1), \ldots, (type_{n-1}, disp_{n-1})\})$$
$$= \{(\mathsf{MPI\_LB}, 0),$$
$$(type_0, disp_0 + start_0 \times ex), \ldots, (type_{n-1}, disp_{n-1} + start_0 \times ex),$$
$$(type_0, disp_0 + (start_0 + 1) \times ex), \ldots, (type_{n-1},$$
$$disp_{n-1} + (start_0 + 1) \times ex), \ldots$$
$$(type_0, disp_0 + (start_0 + subsize_0 - 1) \times ex), \ldots,$$
$$(type_{n-1}, disp_{n-1} + (start_0 + subsize_0 - 1) \times ex),$$
$$(\mathsf{MPI\_UB}, size_0 \times ex)\}$$

$$\text{Subarray}(ndims, \{size_0, size_1, \ldots, size_{ndims-1}\}, \tag{4.3}$$
$$\{subsize_0, subsize_1, \ldots, subsize_{ndims-1}\},$$
$$\{start_0, start_1, \ldots, start_{ndims-1}\}, \mathsf{oldtype})$$
$$= \text{Subarray}(ndims - 1, \{size_1, size_2, \ldots, size_{ndims-1}\},$$
$$\{subsize_1, subsize_2, \ldots, subsize_{ndims-1}\},$$
$$\{start_1, start_2, \ldots, start_{ndims-1}\},$$
$$\text{Subarray}(1, \{size_0\}, \{subsize_0\}, \{start_0\}, \mathsf{oldtype}))$$

$$\text{Subarray}(ndims, \{size_0, size_1, \ldots, size_{ndims-1}\}, \tag{4.4}$$
$$\{subsize_0, subsize_1, \ldots, subsize_{ndims-1}\},$$
$$\{start_0, start_1, \ldots, start_{ndims-1}\}, \mathsf{oldtype})$$
$$= \text{Subarray}(ndims - 1, \{size_0, size_1, \ldots, size_{ndims-2}\},$$
$$\{subsize_0, subsize_1, \ldots, subsize_{ndims-2}\},$$
$$\{start_0, start_1, \ldots, start_{ndims-2}\},$$
$$\text{Subarray}(1, \{size_{ndims-1}\}, \{subsize_{ndims-1}\}, \{start_{ndims-1}\}, \mathsf{oldtype}))$$

For an example use of MPI_TYPE_CREATE_SUBARRAY in the context of I/O see Section 13.9.2.

## 4.1.4 Distributed Array Datatype Constructor

The distributed array type constructor supports HPF-like [34] 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), see MPI I/O, especially Section 13.1.1 on page 419 and Section 13.3 on page 431. Using this view, a collective data access operation (with identical offsets) will yield an HPF-like distribution pattern. (*End of advice to users.*)

**Unofficial Draft for Comment Only**

MPI_TYPE_CREATE_DARRAY(size, rank, ndims, array_of_gsizes, array_of_distribs,
                   array_of_dargs, array_of_psizes, order, oldtype, newtype)

| | | |
|---|---|---|
| IN | size | size of process group (positive integer) |
| IN | rank | rank in process group (non-negative integer) |
| IN | ndims | number of array dimensions as well as process grid dimensions (positive integer) |
| IN | array_of_gsizes | number of elements of type oldtype in each dimension of global array (array of positive integers) |
| IN | array_of_distribs | distribution of array in each dimension (array of state) |
| IN | array_of_dargs | distribution argument in each dimension (array of positive integers) |
| IN | array_of_psizes | size of process grid in each dimension (array of positive integers) |
| IN | order | array storage order flag (state) |
| IN | oldtype | old datatype (handle) |
| OUT | newtype | new datatype (handle) |

```
int MPI_Type_create_darray(int size, int rank, int ndims,
              int array_of_gsizes[], int array_of_distribs[], int
              array_of_dargs[], int array_of_psizes[], int order,
              MPI_Datatype oldtype, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_DARRAY(SIZE, RANK, NDIMS, ARRAY_OF_GSIZES,
              ARRAY_OF_DISTRIBS, ARRAY_OF_DARGS, ARRAY_OF_PSIZES, ORDER,
              OLDTYPE, NEWTYPE, IERROR)
    INTEGER SIZE, RANK, NDIMS, ARRAY_OF_GSIZES(*), ARRAY_OF_DISTRIBS(*),
    ARRAY_OF_DARGS(*), ARRAY_OF_PSIZES(*), ORDER, OLDTYPE, NEWTYPE, IERROR
```

```
{MPI::Datatype MPI::Datatype::Create_darray(int size, int rank, int ndims,
              const int array_of_gsizes[], const int array_of_distribs[],
              const int array_of_dargs[], const int array_of_psizes[],
              int order) const(binding deprecated, see Section 15.2) }
```

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_psizes should be set to 1. (See Example 4.7, page 93.) 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 .

> *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, see Chapter 7 on page 275. (*End of advice to users.*)

**Unofficial Draft for Comment Only**

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 array_of_dargs[i] * array_of_psizes[i] < 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 AR-RAY(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 array_of_dargs[i] equal to MPI_DISTRIBUTE_DFLT_DARG is equivalent to MPI_DISTRIBUTE_CYCLIC with array_of_dargs[i] set to

$$(\text{array\_of\_gsizes}[i] + \text{array\_of\_psizes}[i] - 1)/\text{array\_of\_psizes}[i].$$

If 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 array_of_dargs[i] set to array_of_gsizes[i].

Finally, MPI_DISTRIBUTE_CYCLIC with array_of_dargs[i] equal to MPI_DISTRIBUTE_DFLT_DARG is equivalent to MPI_DISTRIBUTE_CYCLIC with 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];
```

**Unofficial Draft for Comment Only**

For `MPI_ORDER_C`, the code is:

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

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

```
t_rank = rank;
t_size = 1;
for (i = 0; i < ndims; i++)
        t_size *= array_of_psizes[i];
for (i = 0; i < ndims; i++) {
    t_size = t_size / array_of_psizes[i];
    r[i] = t_rank / t_size;
    t_rank = t_rank % t_size;
}
```

Let the typemap of oldtype have the form:

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

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

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

$$
\begin{aligned}
\text{cyclic}(darg&, gsize, r, psize, \text{oldtype}) \\
= \quad &\{(\text{MPI\_LB}, 0), \\
&(type_0, disp_0 + r \times darg \times ex), \ldots, \\
&\qquad (type_{n-1}, disp_{n-1} + r \times darg \times ex), \\
&(type_0, disp_0 + (r \times darg + 1) \times ex), \ldots, \\
&\qquad (type_{n-1}, disp_{n-1} + (r \times darg + 1) \times ex), \\
&\ldots \\
&(type_0, disp_0 + ((r + 1) \times darg - 1) \times ex), \ldots, \\
&\qquad (type_{n-1}, disp_{n-1} + ((r + 1) \times darg - 1) \times ex), \\
\\
&(type_0, disp_0 + r \times darg \times ex + psize \times darg \times ex), \ldots, \\
&\qquad (type_{n-1}, disp_{n-1} + r \times darg \times ex + psize \times darg \times ex), \\
&(type_0, disp_0 + (r \times darg + 1) \times ex + psize \times darg \times ex), \ldots, \\
&\qquad (type_{n-1}, disp_{n-1} + (r \times darg + 1) \times ex + psize \times darg \times ex),
\end{aligned}
$$

$$\ldots$$

$$(type_0, disp_0 + ((r+1) \times darg - 1) \times ex + psize \times darg \times ex), \ldots,$$

$$(type_{n-1}, disp_{n-1} + ((r+1) \times darg - 1) \times ex + psize \times darg \times ex),$$

$$\vdots$$

$$(type_0, disp_0 + r \times darg \times ex + psize \times darg \times ex \times (count - 1)), \ldots,$$

$$(type_{n-1}, disp_{n-1} + r \times darg \times ex + psize \times darg \times ex \times (count - 1)),$$

$$(type_0, disp_0 + (r \times darg + 1) \times ex + psize \times darg \times ex \times (count - 1)), \ldots,$$

$$(type_{n-1}, disp_{n-1} + (r \times darg + 1) \times ex$$

$$+ psize \times darg \times ex \times (count - 1)),$$

$$\ldots$$

$$(type_0, disp_0 + (r \times darg + darg_{last} - 1) \times ex$$

$$+ psize \times darg \times ex \times (count - 1)), \ldots,$$

$$(type_{n-1}, disp_{n-1} + (r \times darg + darg_{last} - 1) \times ex$$

$$+ psize \times darg \times ex \times (count - 1)),$$

$$(\mathsf{MPI\_UB}, gsize * ex)\}$$

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, $darg_{last}$ 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 4.7** 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:

```
ndims = 3
array_of_gsizes(1) = 100
array_of_distribs(1) = MPI_DISTRIBUTE_CYCLIC
array_of_dargs(1) = 10
array_of_gsizes(2) = 200
array_of_distribs(2) = MPI_DISTRIBUTE_NONE
array_of_dargs(2) = 0
array_of_gsizes(3) = 300
array_of_distribs(3) = MPI_DISTRIBUTE_BLOCK
array_of_dargs(3) = MPI_DISTRIBUTE_DFLT_DARG
array_of_psizes(1) = 2
array_of_psizes(2) = 1
array_of_psizes(3) = 3
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_TYPE_CREATE_DARRAY(size, rank, ndims, array_of_gsizes, &
     array_of_distribs, array_of_dargs, array_of_psizes,        &
     MPI_ORDER_FORTRAN, oldtype, newtype, ierr)
```

### 4.1.5   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.

---

MPI_GET_ADDRESS(location, address)

| | | |
|------|----------|------------------------------------|
| IN   | location | location in caller memory (choice) |
| OUT  | address  | address of location (integer)      |

int MPI_Get_address(void *location, MPI_Aint *address)

```
MPI_GET_ADDRESS(LOCATION, ADDRESS, IERROR)
    <type> LOCATION(*)
    INTEGER IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ADDRESS
```

{MPI::Aint MPI::Get_address(void* location) *(binding deprecated, see Section 15.2)* }

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 4.8** 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 16.2.2 on pages 512 and 515. (*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*(binding deprecated, see Section 15.2)* }

ticket265.

MPI_TYPE_SIZE_X(datatype, size)

  IN        datatype                           datatype (handle)

  OUT     size                              datatype size (integer)

```
int MPI_Type_size_x(MPI_Datatype datatype, MPI_Count *size)
```

```
MPI_TYPE_SIZE_X(DATATYPE, SIZE, IERROR)
    INTEGER DATATYPE, IERROR
    INTEGER (KIND=MPI_COUNT_KIND) size
```

MPI_TYPE_SIZE and MPI_TYPE_SIZE_X [returns]return 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.

If the total size of the datatype can not be expressed by the size parameter then MPI_TYPE_SIZE and MPI_TYPE_SIZE_X return the value MPI_UNDEFINED.

### 4.1.6   Lower-Bound and Upper-Bound Markers

It is often convenient to define explicitly the lower bound and upper bound of a type map, and override the definition given on page 97. This allows one to define a datatype that has "holes" at its beginning or its end, or a datatype with entries that extend above the upper bound or below the lower bound. Examples of such usage are provided in Section 4.1.14. Also, the user may want to overide the alignment rules that are used to compute upper bounds and extents. E.g., a C compiler may allow the user to overide default alignment rules for some of the structures within a program. The user has to specify explicitly the bounds of the datatypes that match these structures.

To achieve this, we add two additional "pseudo-datatypes," MPI_LB and MPI_UB, that can be used, respectively, to mark the lower bound or the upper bound of a datatype. These pseudo-datatypes occupy no space ($extent(\mathsf{MPI\_LB}) = extent(\mathsf{MPI\_UB}) = 0$). They do not affect the size or count of a datatype, and do not affect the content of a message created with this datatype. However, they do affect the definition of the extent of a datatype and, therefore, affect the outcome of a replication of this datatype by a datatype constructor.

**Example 4.9**  Let D = (-3, 0, 6); T = (MPI_LB, MPI_INT, MPI_UB), and B = (1, 1, 1). Then a call to MPI_TYPE_STRUCT(3, B, D, T, type1) creates a new datatype that has an extent of 9 (from -3 to 5, 5 included), and contains an integer at displacement 0. This is the datatype defined by the sequence {(lb, -3), (int, 0), (ub, 6)} . If this type is replicated twice by a call to MPI_TYPE_CONTIGUOUS(2, type1, type2) then the newly created type can be described by the sequence {(lb, -3), (int, 0), (int,9), (ub, 15)} . (An entry of type ub can be deleted if there is another entry of type ub with a higher displacement; an entry of type lb can be deleted if there is another entry of type lb with a lower displacement.)

In general, if

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

then the **lower bound** of $Typemap$ is defined to be

$$lb(Typemap) = \begin{cases} \min_j disp_j & \text{if no entry has basic type lb} \\ \min_j\{disp_j \text{ such that } type_j = \mathsf{lb}\} & \text{otherwise} \end{cases}$$

Similarly, the **upper bound** of *Typemap* is defined to be

$$ub(Typemap) = \begin{cases} \max_j disp_j + sizeof(type_j) + \epsilon & \text{if no entry has basic type } \mathsf{ub} \\ \max_j\{disp_j \text{ such that } type_j = \mathsf{ub}\} & \text{otherwise} \end{cases}$$

Then

$$extent(Typemap) = ub(Typemap) - lb(Typemap)$$

If $type_i$ requires alignment to a byte address that is a multiple of $k_i$, then $\epsilon$ is the least non-negative increment needed to round $extent(Typemap)$ to the next multiple of $\max_i k_i$.

The formal definitions given for the various datatype constructors apply now, with the amended definition of **extent**.

### 4.1.7 Extent and Bounds of Datatypes

The following function replaces the three functions MPI_TYPE_UB, MPI_TYPE_LB and MPI_TYPE_EXTENT. It also returns address sized integers, in the Fortran binding. The use of MPI_TYPE_UB, MPI_TYPE_LB and MPI_TYPE_EXTENT is deprecated.

MPI_TYPE_GET_EXTENT(datatype, lb, extent)

| | | |
|---|---|---|
| IN | datatype | datatype to get information on (handle) |
| OUT | lb | lower bound of datatype (integer) |
| OUT | extent | extent of datatype (integer) |

```
int MPI_Type_get_extent(MPI_Datatype datatype, MPI_Aint *lb,
            MPI_Aint *extent)
```

```
MPI_TYPE_GET_EXTENT(DATATYPE, LB, EXTENT, IERROR)
    INTEGER DATATYPE, IERROR
    INTEGER(KIND = MPI_ADDRESS_KIND) LB, EXTENT
```

{void MPI::Datatype::Get_extent(MPI::Aint& lb, MPI::Aint& extent)
        const *(binding deprecated, see Section 15.2)* }

Returns the lower bound and the extent of datatype (as defined in Section 4.1.6 on page 96).

MPI allows one to change the extent of a datatype, using lower bound and upper bound markers (MPI_LB and MPI_UB). This is useful, as it allows to control the stride of successive datatypes that are replicated by datatype constructors, or are replicated by the count argument in a send or receive call. However, the current mechanism for achieving it is painful; also it is restrictive. MPI_LB and MPI_UB are "sticky": once present in a datatype, they cannot be overridden (e.g., the upper bound can be moved up, by adding a new MPI_UB marker, but cannot be moved down below an existing MPI_UB marker). A new type constructor is provided to facilitate these changes. The use of MPI_LB and MPI_UB is deprecated.

MPI_TYPE_CREATE_RESIZED(oldtype, lb, extent, newtype)

IN          oldtype                                  input datatype (handle)

IN          lb                                       new lower bound of datatype (integer)

IN          extent                                   new extent of datatype (integer)

OUT         newtype                                  output datatype (handle)

```
int MPI_Type_create_resized(MPI_Datatype oldtype, MPI_Aint lb, MPI_Aint
              extent, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_RESIZED(OLDTYPE, LB, EXTENT, NEWTYPE, IERROR)
    INTEGER OLDTYPE, NEWTYPE, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) LB, EXTENT
```

```
{MPI::Datatype MPI::Datatype::Create_resized(const MPI::Aint lb,
              const MPI::Aint extent) const(binding deprecated, see Section 15.2) }
```

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.*)

### 4.1.8   True Extent of Datatypes

Suppose we implement gather (see also Section 5.5 on page 140) 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 function is provided which returns the true extent of the datatype.

MPI_TYPE_GET_TRUE_EXTENT(datatype, true_lb, true_extent)

IN          datatype                                 datatype to get information on (handle)

OUT         true_lb                                  true lower bound of datatype (integer)

OUT         true_extent                              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*(binding deprecated, see Section 15.2)* }

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

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

Then

$$true\_lb(Typemap) = min_j\{disp_j \; : \; type_j \neq \mathbf{lb}, \mathbf{ub}\},$$

$$true\_ub(Typemap) = max_j\{disp_j + sizeof(type_j) \; : \; type_j \neq \mathbf{lb}, \mathbf{ub}\},$$

and

$$true\_extent(Typemap) = true\_ub(Typemap) - true\_lb(typemap).$$

(Readers should compare this with the definitions in Section 4.1.6 on page 96 and Section 4.1.7 on page 97, which describe the function MPI_TYPE_GET_EXTENT.)

The true_extent is the minimum number of bytes of memory necessary to hold a datatype, uncompressed.

## 4.1.9 Commit and Free

A datatype object has to be **committed** before it can be used in a communication. As an 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()*(binding deprecated, see Section 15.2)* }

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.

**Example 4.10** The following code fragment gives examples of using MPI_TYPE_COMMIT.

```
INTEGER type1, type2
CALL MPI_TYPE_CONTIGUOUS(5, MPI_REAL, type1, ierr)
               ! new type object created
CALL MPI_TYPE_COMMIT(type1, ierr)
               ! now type1 can be used for communication
type2 = type1
               ! type2 can be used for communication
               ! (it is a handle to same object as type1)
CALL MPI_TYPE_VECTOR(3, 5, 4, MPI_REAL, type1, ierr)
               ! new uncommitted type object created
CALL MPI_TYPE_COMMIT(type1, ierr)
               ! now type1 can be used anew for communication
```

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()*(binding deprecated, see Section 15.2)* }

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. 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 then 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.*)

### 4.1.10  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*(binding deprecated, see Section 15.2)* }

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 6.7.4 on page 263. The new datatype has identical upper bound and lower bound and yields the same net result when fully decoded with the functions in Section 4.1.13. The newtype has the same committed state as the old type.

### 4.1.11  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,

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

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 $n \cdot$ count entries, where entry $i \cdot n + j$ is at location $addr_{i,j} = $ buf $+ extent \cdot i + disp_j$ and has type $type_j$, for $i = 0, ...,$ count $- 1$ and $j = 0, ..., n - 1$. 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 $n \cdot$ count entries, where entry $i \cdot n + j$ has type $type_j$.

**Unofficial Draft for Comment Only**

Similarly, suppose that a receive operation MPI_RECV(buf, count, datatype, source, tag, comm, status) is executed, where datatype has type map,

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

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 $n \cdot$ count entries, where entry $i \cdot n + j$ is at location buf $+ extent \cdot i + disp_j$ and has type $type_j$. If the incoming message consists of $k$ elements, then we must have $k \leq n \cdot$ 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 4.11**   This example shows that type matching is defined in terms of the basic types that a derived type consists of.

```
...
CALL MPI_TYPE_CONTIGUOUS( 2, MPI_REAL, type2, ...)
CALL MPI_TYPE_CONTIGUOUS( 4, MPI_REAL, type4, ...)
CALL MPI_TYPE_CONTIGUOUS( 2, type2, type22, ...)
...
CALL MPI_SEND( a, 4, MPI_REAL, ...)
CALL MPI_SEND( a, 2, type2, ...)
CALL MPI_SEND( a, 1, type22, ...)
CALL MPI_SEND( a, 1, type4, ...)
...
CALL MPI_RECV( a, 4, MPI_REAL, ...)
CALL MPI_RECV( a, 2, type2, ...)
CALL MPI_RECV( a, 1, type22, ...)
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,

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

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 $0 \leq k \leq$ count $\cdot n$. The number of basic elements received can be retrieved from status using the query function MPI_GET_ELEMENTS.

MPI_GET_ELEMENTS( status, datatype, count)

| IN | status | return status of receive operation (Status) |
|---|---|---|
| IN | datatype | datatype used by receive operation (handle) |
| OUT | count | number of received basic elements (integer) |

```
int MPI_Get_elements(MPI_Status *status, MPI_Datatype datatype, int *count)
```

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

{int MPI::Status::Get_elements(const MPI::Datatype& datatype) const*(binding deprecated, see Section 15.2)* }

MPI_GET_ELEMENTS_X( status, datatype, count)

| IN | status | return status of receive operation (Status) |
|---|---|---|
| IN | datatype | datatype used by receive operation (handle) |
| OUT | count | number of received basic elements (integer) |

```
int MPI_Get_elements_x(MPI_Status *status, MPI_Datatype datatype,
            MPI_Count *count)
```

```
MPI_GET_ELEMENTS_X(STATUS, DATATYPE, COUNT, IERROR)
    INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, IERROR
    INTEGER (KIND=MPI_COUNT_KIND) COUNT
```

The previously defined [function]functions, MPI_GET_COUNT MPI_GET_COUNT_X (Section 3.2.5), [has]have a different behavior. [It returns]They return 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 $0 \leq k \leq$ 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 and MPI_GET_COUNT_X [returns]return the value MPI_UNDEFINED. The datatype argument should match the argument provided by the receive call that set the status variable.

If the number of basic elements received can not be expressed by the count parameter then MPI_GET_ELEMENTS and MPI_GET_ELEMENTS_X return the value MPI_UNDEFINED.

**Example 4.12** Usage of MPI_GET_COUNT and MPI_GET_ELEMENTS.

```
...
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
```

1
2
3
4
5
6
7
8
9
10
11
12
13 ticket265.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28 ticket265.
29 ticket265.
30 ticket265.
31 ticket265.
32
33
34
35 ticket265.
36 ticket265.
37
38
39
40
41
42
43
44
45
46
47
48

```
1          CALL MPI_SEND(a, 2, MPI_REAL, 1, 0, comm, ierr)
2          CALL MPI_SEND(a, 3, MPI_REAL, 1, 0, comm, ierr)
3    ELSE IF (rank.EQ.1) THEN
4          CALL MPI_RECV(a, 2, Type2, 0, 0, comm, stat, ierr)
5          CALL MPI_GET_COUNT(stat, Type2, i, ierr)     ! returns i=1
6          CALL MPI_GET_ELEMENTS(stat, Type2, i, ierr)  ! returns i=2
7          CALL MPI_RECV(a, 2, Type2, 0, 0, comm, stat, ierr)
8          CALL MPI_GET_COUNT(stat, Type2, i, ierr)       ! returns i=MPI_UNDEFINED
9          CALL MPI_GET_ELEMENTS(stat, Type2, i, ierr)  ! returns i=3
10   END IF
```

The [function]functions MPI_GET_ELEMENTS and MPI_GET_ELEMENTS_X 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 ( or MPI_GET_ELEMENTS_X) 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 or MPI_GET_ELEMENTS_X . (*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.*)

### 4.1.12  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 to another. Also, in machines with a segmented address space, addresses are not unique and address arithmetic has some peculiar properties. Thus, the use of **addresses**, that is, displacements relative to the start address MPI_BOTTOM, has to be restricted.

Variables belong to the same **sequential storage** if they belong to the same array, to the same COMMON block in Fortran, or to the same structure in C. Valid addresses are defined recursively as follows:

1. The function MPI_GET_ADDRESS returns a valid address, when passed as argument a variable of the calling program.

2. The buf argument of a communication function evaluates to a valid address, when passed as argument a variable of the calling program.

3. If v is a valid address, and i is an integer, then v+i is a valid address, provided v and v+i are in the same sequential storage.

4. If v is a valid address then MPI_BOTTOM + v is a valid address.

A correct program uses only valid addresses to identify the locations of entries in communication buffers. Furthermore, if u and v are two valid addresses, then the (integer) difference u - v can be computed only if both u and v are in the same sequential storage. No other arithmetic operations can be meaningfully executed on addresses.

The rules above impose no constraints on the use of derived datatypes, as long as they are used to define a communication buffer that is wholly contained within the same sequential storage. However, the construction of a communication buffer that contains variables that are not within the same sequential storage must obey certain restrictions. Basically, a communication buffer with variables that are not within the same sequential storage can be used only by specifying in the communication call buf = MPI_BOTTOM, count = 1, and using a datatype argument where all displacements are valid (absolute) addresses.

*Advice to users.* It is not expected that MPI implementations will be able to detect erroneous, "out of bound" displacements — unless those overflow the user address space — since the MPI call may not know the extent of the arrays and records in the host program. (*End of advice to users.*)

*Advice to implementors.* There is no need to distinguish (absolute) addresses and (relative) displacements on a machine with contiguous address space: MPI_BOTTOM is zero, and both addresses and displacements are integers. On machines where the distinction is required, addresses are recognized as expressions that involve MPI_BOTTOM. (*End of advice to implementors.*)

### 4.1.13 Decoding a Datatype

MPI datatype objects allow users to specify an arbitrary layout of data in memory. There are several cases where accessing the layout information in opaque datatype objects would be useful. 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 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, com-
                biner)

| IN | datatype | datatype to access (handle) |
| OUT | num_integers | number of input integers used in the call constructing combiner (non-negative integer) |
| OUT | num_addresses | number of input addresses used in the call constructing combiner (non-negative integer) |
| OUT | num_datatypes | number of input datatypes used in the call constructing combiner (non-negative integer) |
| OUT | combiner | 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 *(binding deprecated, see
              Section 15.2)* }

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 4.1: combiner values returned from MPI_TYPE_GET_ENVELOPE

If combiner is MPI_COMBINER_NAMED then datatype is a named predefined datatype.

For deprecated 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 preferred calls all use address sized arguments so two combiners are not required for them.

*Rationale.* For recreating the original call, it is important to know if address information may have been truncated. The deprecated calls from Fortran for a few routines could be subject to truncation in the case where the default INTEGER size is smaller

**Unofficial Draft for Comment Only**

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,  ar-
                      ray_of_integers, array_of_addresses, array_of_datatypes)

| IN  | datatype | datatype to access (handle) |
|-----|----------|------------------------------|
| IN  | max_integers | number of elements in array_of_integers (non-negative integer) |
| IN  | max_addresses | number of elements in array_of_addresses (non-negative integer) |
| IN  | max_datatypes | number of elements in array_of_datatypes (non-negative integer) |
| OUT | array_of_integers | contains integer arguments used in constructing datatype (array of integers) |
| OUT | array_of_addresses | contains address arguments used in constructing datatype (array of integers) |
| OUT | array_of_datatypes | 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
```
*(binding deprecated, see Section 15.2)* }

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.

> *Rationale.*   The arguments max_integers, max_addresses, and max_datatypes allow for error checking in the call. (*End of rationale.*)

The datatypes returned in `array_of_datatypes` are handles to datatype objects that are equivalent to the datatypes used in the original construction call. If these were derived datatypes, then the returned datatypes are new datatype objects, and the user is responsible for freeing these datatypes with MPI_TYPE_FREE. If these were predefined datatypes, then the returned datatype is equal to that (constant) predefined datatype and cannot be freed.

The committed state of returned derived datatypes is undefined, i.e., the datatypes may or may not be committed. Furthermore, the content of attributes of returned datatypes is undefined.

Note that MPI_TYPE_GET_CONTENTS can be invoked with a `datatype` argument that was constructed using MPI_TYPE_CREATE_F90_REAL, MPI_TYPE_CREATE_F90_INTEGER, or MPI_TYPE_CREATE_F90_COMPLEX (an unnamed predefined datatype). In such a case, an empty `array_of_datatypes` is returned.

> *Rationale.* The definition of datatype equivalence implies that equivalent predefined datatypes are equal. By requiring the same handle for named predefined datatypes, it is possible to use the `==` or `.EQ.` comparison operator to determine the datatype involved. (*End of rationale.*)

> *Advice to implementors.* The datatypes returned in `array_of_datatypes` must appear to the user as if each is an equivalent copy of the datatype used in the type constructor call. Whether this is done by creating a new datatype or via another mechanism such as a reference count mechanism is up to the implementation as long as the semantics are preserved. (*End of advice to implementors.*)

> *Rationale.* The committed state and attributes of the returned datatype is deliberately left vague. The datatype used in the original construction may have been modified since its use in the constructor call. Attributes can be added, removed, or modified as well as having the datatype committed. The semantics given allow for a reference count implementation without having to track these changes. (*End of rationale.*)

In the deprecated datatype constructor calls, the address arguments in Fortran are of type `INTEGER`. In the preferred calls, the address arguments are of type `INTEGER(KIND=MPI_ADDRESS_KIND)`. The call MPI_TYPE_GET_CONTENTS returns all addresses in an argument of type `INTEGER(KIND=MPI_ADDRESS_KIND)`. This is true even if the deprecated calls were used. Thus, the location of values returned can be thought of as being returned by the C bindings. It can also be determined by examining the preferred calls for datatype constructors for the deprecated calls that involve addresses.

> *Rationale.* By having all address arguments returned in the `array_of_addresses` argument, the result from a C and Fortran decoding of a `datatype` gives the result in the same argument. It is assumed that an integer of type `INTEGER(KIND=MPI_ADDRESS_KIND)` will be at least as large as the `INTEGER` argument used in datatype construction with the old MPI-1 calls so no loss of information will 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, IERROR)
       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

| 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) |

**Unofficial Draft for Comment Only**

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.

### 4.1.14   Examples

The following examples illustrate the use of derived datatypes.

**Example 4.13**  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 4.14** 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
        blocklen(i) = 100-i
      END DO

C     create datatype for lower triangular part
      CALL MPI_TYPE_INDEXED( 100, blocklen, 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
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

**Example 4.15** Transpose a matrix.

```
        REAL a(100,100), b(100,100)
        INTEGER row, xpose, sizeofreal, myrank, ierr
        INTEGER status(MPI_STATUS_SIZE)

C       transpose matrix a onto b

        CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)

        CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)

C       create datatype for one row
        CALL MPI_TYPE_VECTOR( 100, 1, 100, MPI_REAL, row, ierr)

C       create datatype for matrix in row-major order
        CALL MPI_TYPE_HVECTOR( 100, 1, sizeofreal, row, xpose, ierr)

        CALL MPI_TYPE_COMMIT( xpose, ierr)

C       send matrix in row-major order and receive in column major order
        CALL MPI_SENDRECV( a, 1, xpose, myrank, 0, b, 100*100,
                  MPI_REAL, myrank, 0, MPI_COMM_WORLD, status, ierr)
```

**Example 4.16** Another approach to the transpose problem:

```
        REAL a(100,100), b(100,100)
        INTEGER  disp(2), blocklen(2), type(2), row, row1, sizeofreal
        INTEGER  myrank, ierr
        INTEGER status(MPI_STATUS_SIZE)

        CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)

C       transpose matrix a onto b

        CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)

C       create datatype for one row
        CALL MPI_TYPE_VECTOR( 100, 1, 100, MPI_REAL, row, ierr)

C       create datatype for one row, with the extent of one real number
        disp(1) = 0
        disp(2) = sizeofreal
        type(1)  = row
        type(2)  = MPI_UB
        blocklen(1)  = 1
        blocklen(2)  = 1
        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 4.17** 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, tag;
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
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    /* 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);
```

```
/* prepend particle count */
MPI_Address(&j, zzdisp);
MPI_Address(particle, zzdisp+1);
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);


        /* 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
16
```

**Example 4.18** The same manipulations as in the previous example, but use absolute addresses in datatypes.

```
struct Partstruct
    {
       int class;
       double d[6];
       char b[7];
    };

struct Partstruct particle[1000];

            /* build datatype describing first array entry */

MPI_Datatype Particletype;
MPI_Datatype type[3] = {MPI_INT, MPI_DOUBLE, MPI_CHAR};
int         block[3] = {1, 6, 7};
MPI_Aint    disp[3];

MPI_Address( particle, disp);
MPI_Address( particle[0].d, disp+1);
MPI_Address( particle[0].b, disp+2);
MPI_Type_struct( 3, block, disp, type, &Particletype);

/* Particletype describes first array entry -- using absolute
    addresses */

                /* 5.1:
            send the entire array */

MPI_Type_commit( &Particletype);
MPI_Send( MPI_BOTTOM, 1000, Particletype, dest, tag, comm);
```

```
              /* 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 4.19** Handling of unions.

```
union {
   int     ival;
   float   fval;
     } u[1000];

int     utype;
```

```
/* All entries of u have identical type; variable
   utype keeps track of their current type */

MPI_Datatype    type[2];
int             blocklen[2] = {1,1};
MPI_Aint        disp[2];
MPI_Datatype    mpi_utype[2];
MPI_Aint        i,j;

/* compute an MPI datatype for each possible union type;
   assume values are left-aligned in union storage. */

MPI_Address( u, &i);
MPI_Address( u+1, &j);
disp[0] = 0; disp[1] = j-i;
type[1] = MPI_UB;

type[0] = MPI_INT;
MPI_Type_struct(2, blocklen, disp, type, &mpi_utype[0]);

type[0] = MPI_FLOAT;
MPI_Type_struct(2, blocklen, disp, type, &mpi_utype[1]);

for(i=0; i<2; i++) MPI_Type_commit(&mpi_utype[i]);

/* actual communication */

MPI_Send(u, 1000, mpi_utype[utype], dest, tag, comm);
```

**Example 4.20**  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.

```
/*
  Example of decoding a datatype.

  Returns 0 if the datatype is predefined, 1 otherwise
 */
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"
int printdatatype( MPI_Datatype datatype )
{
    int *array_of_ints;
    MPI_Aint *array_of_adds;
    MPI_Datatype *array_of_dtypes;
    int num_ints, num_adds, num_dtypes, combiner;
```

```
int i;                                                                          1
                                                                                2
MPI_Type_get_envelope( datatype,                                                3
                      &num_ints, &num_adds, &num_dtypes, &combiner );           4
switch (combiner) {                                                             5
case MPI_COMBINER_NAMED:                                                        6
    printf( "Datatype is named:" );                                            7
    /* To print the specific type, we can match against the                    8
       predefined forms. We can NOT use a switch statement here                 9
       We could also use MPI_TYPE_GET_NAME if we prefered to use               10
       names that the user may have changed.                                   11
     */                                                                        12
    if      (datatype == MPI_INT)    printf( "MPI_INT\n" );                    13
    else if (datatype == MPI_DOUBLE) printf( "MPI_DOUBLE\n" );                 14
    ... else test for other types ...                                         15
    return 0;                                                                  16
    break;                                                                     17
case MPI_COMBINER_STRUCT:                                                      18
case MPI_COMBINER_STRUCT_INTEGER:                                             19
    printf( "Datatype is struct containing" );                               20
    array_of_ints  = (int *)malloc( num_ints * sizeof(int) );                21
    array_of_adds  =                                                         22
            (MPI_Aint *) malloc( num_adds * sizeof(MPI_Aint) );              23
    array_of_dtypes = (MPI_Datatype *)                                       24
        malloc( num_dtypes * sizeof(MPI_Datatype) );                         25
    MPI_Type_get_contents( datatype, num_ints, num_adds, num_dtypes,         26
                  array_of_ints, array_of_adds, array_of_dtypes );           27
    printf( " %d datatypes:\n", array_of_ints[0] );                          28
    for (i=0; i<array_of_ints[0]; i++) {                                     29
        printf( "blocklength %d, displacement %ld, type:\n",                 30
                array_of_ints[i+1], array_of_adds[i] );                      31
        if (printdatatype( array_of_dtypes[i] )) {                           32
            /* Note that we free the type ONLY if it                         33
               is not predefined */                                          34
            MPI_Type_free( &array_of_dtypes[i] );                            35
        }                                                                    36
    }                                                                        37
    free( array_of_ints );                                                   38
    free( array_of_adds );                                                   39
    free( array_of_dtypes );                                                 40
    break;                                                                   41
    ... other combiner values ...                                           42
default:                                                                    43
    printf( "Unrecognized combiner type\n" );                               44
}                                                                            45
return 1;                                                                    46
}                                                                            47
                                                                            48
```

## 4.2   Pack and Unpack

Some existing communication libraries provide pack/unpack functions for sending noncontiguous data. In these, the user explicitly packs data into a contiguous buffer before sending it, and unpacks it from a contiguous buffer after receiving it. Derived datatypes, which are described in Section 4.1, allow one, in most cases, to avoid explicit packing and unpacking. The user specifies the layout of the data to be sent or received, and the communication library directly accesses a noncontiguous buffer. The pack/unpack routines are provided for compatibility with previous libraries. Also, they provide some functionality that is not otherwise available in MPI. For instance, a message can be received in several parts, where the receive operation done on a later part may depend on the content of a former part. Another use is that outgoing messages may be explicitly buffered in user supplied space, thus overriding the system buffering policy. Finally, the availability of pack and unpack operations facilitates the development of additional communication libraries layered on top of MPI.

MPI_PACK(inbuf, incount, datatype, outbuf, outsize, position, comm)

| | | |
|---|---|---|
| IN | inbuf | input buffer start (choice) |
| IN | incount | number of input data items (non-negative integer) |
| IN | datatype | datatype of each input data item (handle) |
| OUT | outbuf | output buffer start (choice) |
| IN | outsize | output buffer size, in bytes (non-negative integer) |
| INOUT | position | current position in buffer, in bytes (integer) |
| IN | comm | communicator for packed message (handle) |

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

```
MPI_PACK(INBUF, INCOUNT, DATATYPE, OUTBUF, OUTSIZE, POSITION, COMM, IERROR)
    <type> INBUF(*), OUTBUF(*)
    INTEGER INCOUNT, DATATYPE, OUTSIZE, POSITION, COMM, IERROR
```

```
{void MPI::Datatype::Pack(const void* inbuf, int incount, void *outbuf,
             int outsize, int& position, const MPI::Comm &comm)
             const(binding deprecated, see Section 15.2) }
```

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)

| IN | inbuf | input buffer start (choice) |
|---|---|---|
| IN | insize | size of input buffer, in bytes (non-negative integer) |
| INOUT | position | current position in bytes (integer) |
| OUT | outbuf | output buffer start (choice) |
| IN | outcount | number of items to be unpacked (integer) |
| IN | datatype | datatype of each output data item (handle) |
| IN | comm | communicator for packed message (handle) |

```
int MPI_Unpack(void* inbuf, int insize, int *position, void *outbuf,
           int outcount, MPI_Datatype datatype, MPI_Comm comm)
```

```
MPI_UNPACK(INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT, DATATYPE, COMM,
           IERROR)
    <type> INBUF(*), OUTBUF(*)
    INTEGER INSIZE, POSITION, OUTCOUNT, DATATYPE, COMM, IERROR
```

{void MPI::Datatype::Unpack(const void* inbuf, int insize, void *outbuf,
           int outcount, int& position, const MPI::Comm& comm)
           const *(binding deprecated, see Section 15.2)* }

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.)

**Unofficial Draft for Comment Only**

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 (non-negative 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 (non-negative 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 *(binding deprecated, see Section 15.2)* }

A call to MPI_PACK_SIZE(incount, datatype, comm, size) returns in size an upper bound [ on the increment in position that is effected] on the size of the output buffer, in bytes, required  by a call to MPI_PACK(inbuf, incount, datatype, outbuf, outcount, position, comm). The value returned as the size argument of MPI_PACK_SIZE for a datatype larger than what can be represented by a C integer of Fortran INTEGER is MPI_UNDEFINED.

> *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 4.21**  An example using MPI_PACK.

```
int        position, i, j, a[2];
char       buff[1000];

MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank == 0)
{
   /* SENDER CODE */

   position = 0;
   MPI_Pack(&i, 1, MPI_INT, buff, 1000, &position, MPI_COMM_WORLD);
   MPI_Pack(&j, 1, MPI_INT, buff, 1000, &position, MPI_COMM_WORLD);
   MPI_Send( buff, position, MPI_PACKED, 1, 0, MPI_COMM_WORLD);
}
else  /* RECEIVER CODE */
   MPI_Recv( a, 2, MPI_INT, 0, 0, MPI_COMM_WORLD);
```

**Example 4.22**  An elaborate example.

```
int   position, i;
float a[1000];
char  buff[1000];

MPI_Comm_rank(MPI_Comm_world, &myrank);
if (myrank == 0)
{
  /* SENDER CODE */

  int len[2];
  MPI_Aint disp[2];
  MPI_Datatype type[2], newtype;

  /* build datatype for i followed by a[0]...a[i-1] */
```

ticket265.

ticket265.

```
1      len[0] = 1;
2      len[1] = i;
3      MPI_Address( &i, disp);
4      MPI_Address( a, disp+1);
5      type[0] = MPI_INT;
6      type[1] = MPI_FLOAT;
7      MPI_Type_struct( 2, len, disp, type, &newtype);
8      MPI_Type_commit( &newtype);
9
10     /* Pack i followed by a[0]...a[i-1]*/
11
12     position = 0;
13     MPI_Pack( MPI_BOTTOM, 1, newtype, buff, 1000, &position, MPI_COMM_WORLD);
14
15     /* Send */
16
17     MPI_Send( buff, position, MPI_PACKED, 1, 0,
18               MPI_COMM_WORLD);
19
20   /* *****
21      One can replace the last three lines with
22      MPI_Send( MPI_BOTTOM, 1, newtype, 1, 0, MPI_COMM_WORLD);
23      ***** */
24   }
25   else if (myrank == 1)
26   {
27      /* RECEIVER CODE */
28
29     MPI_Status status;
30
31     /* Receive */
32
33     MPI_Recv( buff, 1000, MPI_PACKED, 0, 0, MPI_COMM_WORLD, &status);
34
35     /* Unpack i */
36
37     position = 0;
38     MPI_Unpack(buff, 1000, &position, &i, 1, MPI_INT, MPI_COMM_WORLD);
39
40     /* Unpack a[0]...a[i-1] */
41     MPI_Unpack(buff, 1000, &position, a, i, MPI_FLOAT, MPI_COMM_WORLD);
42   }
```

**Example 4.23** Each process sends a count, followed by count characters to the root; the root concatenates all characters into one string.

```
int  count, gsize, counts[64], totalcount, k1, k2, k,
     displs[64], position, concat_pos;
```

```
char chr[100], *lbuf, *rbuf, *cbuf;                              1
                                                                2
MPI_Comm_size(comm, &gsize);                                    3
MPI_Comm_rank(comm, &myrank);                                   4
                                                                5
    /* allocate local pack buffer */                           6
MPI_Pack_size(1, MPI_INT, comm, &k1);                          7
MPI_Pack_size(count, MPI_CHAR, comm, &k2);                     8
k = k1+k2;                                                      9
lbuf = (char *)malloc(k);                                      10
                                                               11
    /* pack count, followed by count characters */            12
position = 0;                                                  13
MPI_Pack(&count, 1, MPI_INT, lbuf, k, &position, comm);       14
MPI_Pack(chr, count, MPI_CHAR, lbuf, k, &position, comm);     15
                                                               16
if (myrank != root) {                                          17
    /* gather at root sizes of all packed messages */         18
  MPI_Gather( &position, 1, MPI_INT, NULL, 0,                  19
          MPI_DATATYPE_NULL, root, comm);                      20
                                                               21
    /* gather at root packed messages */                      22
  MPI_Gatherv( lbuf, position, MPI_PACKED, NULL,              23
          NULL, NULL, NULL, root, comm);                       24
                                                               25
} else {   /* root code */                                    26
    /* gather sizes of all packed messages */                 27
  MPI_Gather( &position, 1, MPI_INT, counts, 1,               28
          MPI_INT, root, comm);                                29
                                                               30
    /* gather all packed messages */                          31
  displs[0] = 0;                                               32
  for (i=1; i < gsize; i++)                                    33
    displs[i] = displs[i-1] + counts[i-1];                    34
  totalcount = displs[gsize-1] + counts[gsize-1];            35
  rbuf = (char *)malloc(totalcount);                          36
  cbuf = (char *)malloc(totalcount);                          37
  MPI_Gatherv( lbuf, position, MPI_PACKED, rbuf,             38
          counts, displs, MPI_PACKED, root, comm);           39
                                                               40
     /* unpack all messages and concatenate strings */       41
  concat_pos = 0;                                             42
  for (i=0; i < gsize; i++) {                                 43
     position = 0;                                            44
     MPI_Unpack( rbuf+displs[i], totalcount-displs[i],       45
          &position, &count, 1, MPI_INT, comm);              46
     MPI_Unpack( rbuf+displs[i], totalcount-displs[i],       47
          &position, cbuf+concat_pos, count, MPI_CHAR, comm); 48
```

**Unofficial Draft for Comment Only**

```
        concat_pos += count;
    }
    cbuf[concat_pos] = '\0';
}
```

## 4.3  Canonical MPI_PACK and MPI_UNPACK

These functions read/write data to/from the buffer in the "external32" data format specified
in Section 13.5.2, and calculate the size needed for packing. Their first arguments specify
the data format, for future extensibility, but currently 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(binding deprecated, see Section 15.2) }
```

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)
```

```
MPI_UNPACK_EXTERNAL(DATAREP, INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT,
            DATATYPE, IERROR)
    INTEGER OUTCOUNT, DATATYPE, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) INSIZE, POSITION
    CHARACTER*(*) DATAREP
    <type> INBUF(*), OUTBUF(*)
```

{void MPI::Datatype::Unpack_external(const char* datarep,
            const void* inbuf, MPI::Aint insize, MPI::Aint& position,
            void* outbuf, int outcount) const *(binding deprecated, see
            Section 15.2)* }

MPI_PACK_EXTERNAL_SIZE( datarep, incount, datatype, size )

| | | |
|---|---|---|
| IN | datarep | data representation (string) |
| IN | incount | number of input data items (integer) |
| IN | datatype | datatype of each input data item (handle) |
| OUT | size | output buffer size, in bytes (integer) |

```
int MPI_Pack_external_size(char *datarep, int incount,
            MPI_Datatype datatype, MPI_Aint *size)
```

```
MPI_PACK_EXTERNAL_SIZE(DATAREP, INCOUNT, DATATYPE, SIZE, IERROR)
    INTEGER INCOUNT, DATATYPE, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) SIZE
    CHARACTER*(*) DATAREP
```

{MPI::Aint MPI::Datatype::Pack_external_size(const char* datarep,
            int incount) const*(binding deprecated, see Section 15.2)* }

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

# Chapter 5

# Collective Communication

## 5.1 Introduction and Overview

Collective communication is defined as communication that involves a group or groups of processes. The functions of this type provided by MPI are the following:

- MPI_BARRIER, MPI_IBARRIER: Barrier synchronization across all members of a group (Section 5.3  and Section 5.12.1  ).

- MPI_BCAST  , MPI_IBCAST  : Broadcast from one member to all members of a group (Section 5.4 and Section 5.12.2). This is shown as "broadcast" in Figure 5.1.

- MPI_GATHER, MPI_IGATHER, MPI_GATHERV, MPI_IGATHERV , MPI_GATHERW, MPI_IGATHERW: Gather data from all members of a group to one member (Section 5.5 and Section 5.12.3). This is shown as "gather" in Figure 5.1.

- MPI_SCATTER, MPI_ISCATTER , MPI_SCATTERV, MPI_ISCATTERV  , MPI_SCATTERW, MPI_ISCATTERW: Scatter data from one member to all members of a group (Section 5.6 and Section 5.12.4). This is shown as "scatter" in Figure 5.1.

- MPI_ALLGATHER, MPI_IALLGATHER , MPI_ALLGATHERV, MPI_IALLGATHERV , MPI_ALLGATHERW, MPI_IALLGATHERW: A variation on Gather where all members of a group receive the result (Section 5.7 and Section 5.12.5). This is shown as "allgather" in Figure 5.1.

- MPI_ALLTOALL, MPI_IALLTOALL , MPI_ALLTOALLV, MPI_IALLTOALLV , MPI_ALLTOALLW, MPI_IALLTOALLW : Scatter/Gather data from all members to all members of a group (also called complete exchange) (Section 5.8 and Section 5.12.6). This is shown as "complete exchange" in Figure 5.1.

- MPI_ALLREDUCE, MPI_IALLREDUCE , MPI_REDUCE, MPI_IREDUCE : Global reduction operations such as sum, max, min, or user-defined functions, where the result is returned to all members of a group (Section 5.9.6 and Section 5.12.8)  and a variation where the result is returned to only one member (Section 5.9 and Section 5.12.7).

-  MPI_REDUCE_SCATTER_BLOCK, MPI_IREDUCE_SCATTER_BLOCK, MPI_REDUCE_SCATTER , MPI_IREDUCE_SCATTER : A combined reduction and scatter operation (Section 5.10, Section 5.12.9, and Section 5.12.10).

ticket109.
ticket109.

- MPI_SCAN, MPI_ISCAN, MPI_EXSCAN, MPI_IEXSCAN: Scan across all members of a group (also called prefix) (Section 5.11, Section 5.11.2, Section 5.12.11, and Section 5.12.12).

One of the key arguments in a call to a collective routine is a communicator that defines the group or groups of participating processes and provides a context for the operation. This is discussed further in Section 5.2. 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 4. Several collective routines such as broadcast and gather have a single originating or receiving process. Such a process is called the *root*. Some arguments in the collective functions are specified as "significant only at root," and are ignored for all participants except the root. The reader is referred to Chapter 4 for information concerning communication buffers, general datatypes and type matching rules, and to Chapter 6 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. Different type maps (the layout in memory, see Section 4.1) between sender and receiver are still allowed.

ticket109.
ticket109.
ticket109.
ticket109.
ticket109.
ticket109.
ticket109.
ticket109.

Collective [routine calls]operations can (but are not required to) [return]complete as soon as [their]the caller's participation in the collective communication is [complete]finished. A blocking operation is complete as soon as the call returns. A nonblocking (immediate) call requires a separate completion call (cf. Section 3.7). The completion of a [call]collective operation indicates that the caller is [now] free to modify locations in the communication buffer. It does not indicate that other processes in the group have completed or even started the operation (unless otherwise implied by 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]Thus, a collective communication operation may, or may not, have the effect of synchronizing all calling processes. This statement excludes, of course, the barrier operation.

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.

ticket109.

The collective operations do not have a message tag argument. A more detailed discussion of correct use of collective routines is found in Section 5.13.

> *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.

ticket109.

> [The collective operations do not accept a message tag argument. If future revisions of MPI define nonblocking collective functions, then tags (or a similar mechanism) might need to be added so as to allow the dis-ambiguation of multiple, pending, collective operations.] (*End of rationale.*)
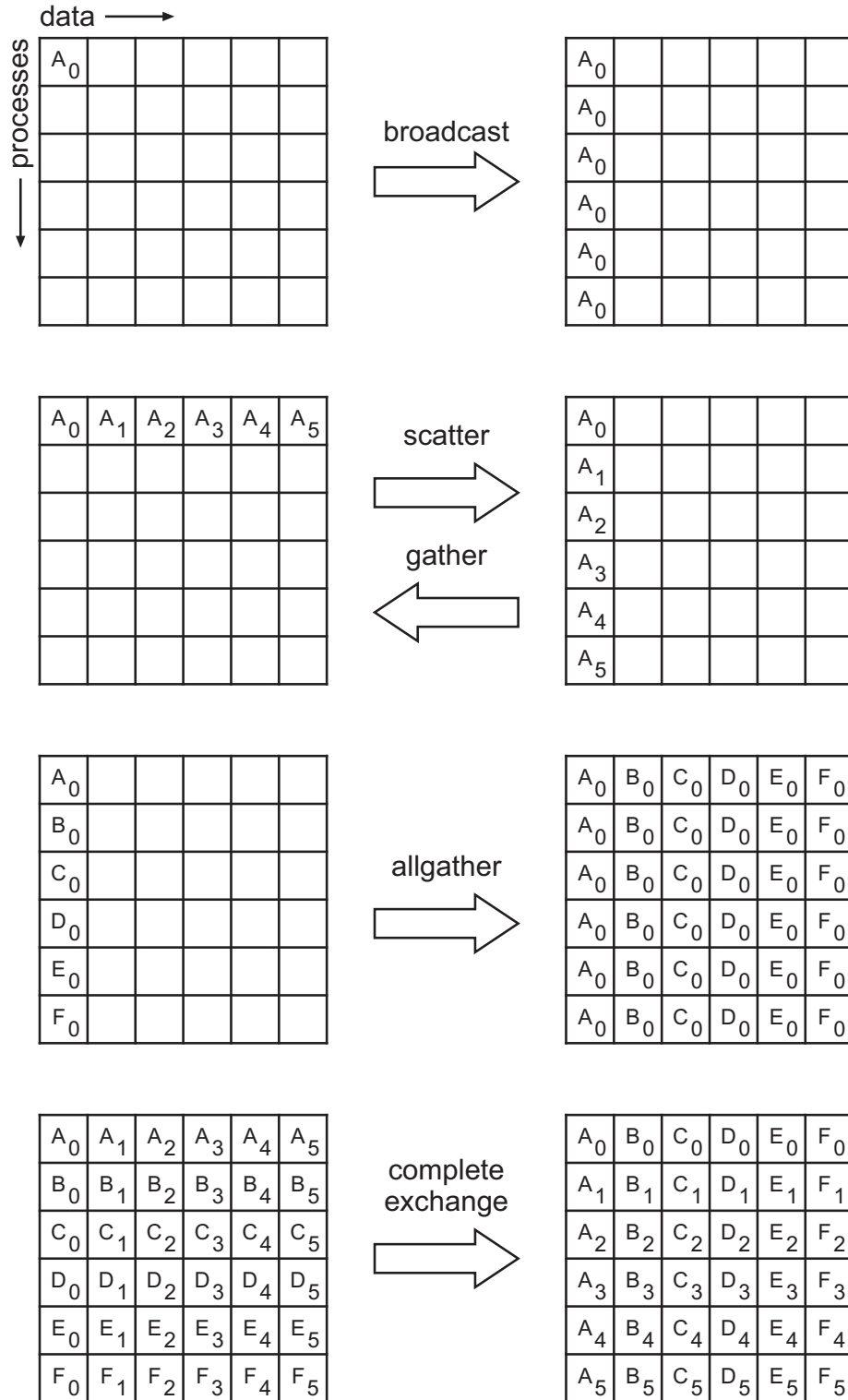
Figure 5.1: Collective move functions illustrated for a group of six processes. In each case, each row of boxes represents data locations in one process. Thus, in the broadcast, initially just the first process contains the data $A_0$, but after the broadcast all processes contain it.

*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 Section 5.13. (*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 might 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 Section 5.13. (*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.

## 5.2   Communicator Argument

The key concept of the collective functions is to have a group or groups of participating processes. The routines do not have group identifiers as explicit arguments. Instead, there is a communicator argument. Groups and communicators are discussed in full detail in Chapter 6. For the purposes of this chapter, it is sufficient to know that there are two types of communicators: *intra-communicators* and *inter-communicators*. An intracommunicator can be thought of as an i[n]dentifier for a single group of processes linked with a context. An intercommunicator identifies two distinct groups of processes linked with a context.

ticket109.

### 5.2.1   Specifics for Intracommunicator Collective Operations

All processes in the group identified by the intracommunicator must call the collective routine.

In many cases, 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, depending on the operation performed.

*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.*)

## 5.2.2 Applying Collective Operations to Intercommunicators

To understand how collective operations apply to intercommunicators, we can view most MPI intracommunicator collective operations as fitting one of the following categories (see, for instance, [47]):

**All-To-All** All processes contribute to the result. All processes receive the result.

- MPI_ALLGATHER, MPI_IALLGATHER, MPI_ALLGATHERV, MPI_IALLGATHERV , MPI_ALLGATHERW, MPI_IALLGATHERW
- MPI_ALLTOALL, MPI_IALLTOALL, MPI_ALLTOALLV, MPI_IALLTOALLV, MPI_ALLTOALLW, MPI_IALLTOALLW
- MPI_ALLREDUCE, MPI_IALLREDUCE, MPI_REDUCE_SCATTER_BLOCK, MPI_IREDUCE_SCATTER_BLOCK, MPI_REDUCE_SCATTER, MPI_IREDUCE_SCATTER
- MPI_BARRIER, MPI_IBARRIER

**All-To-One** All processes contribute to the result. One process receives the result.

- MPI_GATHER, MPI_IGATHER, MPI_GATHERV, MPI_IGATHERV , MPI_GATHERW, MPI_IGATHERW
- MPI_REDUCE, MPI_IREDUCE

**One-To-All** One process contributes to the result. All processes receive the result.

- MPI_BCAST, MPI_IBCAST
- MPI_SCATTER, MPI_ISCATTER, MPI_SCATTERV, MPI_ISCATTERV , MPI_SCATTERW, MPI_ISCATTERW

**Other** Collective operations that do not fit into one of the above categories.

- MPI_SCAN, MPI_ISCAN, MPI_EXSCAN, MPI_IEXSCAN

The data movement patterns of MPI_SCAN, MPI_ISCAN [and], MPI_EXSCAN, and MPI_IEXSCAN do not fit this taxonomy.

The application of collective communication to intercommunicators is best described in terms of two 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 5.2). 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 5.3). 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.

The following collective operations also apply to intercommunicators:

- MPI_BARRIER, MPI_IBARRIER

- MPI_BCAST, MPI_IBCAST

- MPI_GATHER, MPI_IGATHER, MPI_GATHERV, MPI_IGATHERV,

- MPI_SCATTER, MPI_ISCATTER, MPI_SCATTERV, MPI_ISCATTERV,

- MPI_ALLGATHER, MPI_IALLGATHER, MPI_ALLGATHERV, MPI_IALLGATHERV,

- MPI_ALLTOALL, MPI_IALLTOALL, MPI_ALLTOALLV, MPI_IALLTOALLV, MPI_ALLTOALLW, MPI_IALLTOALLW,

- MPI_ALLREDUCE, MPI_IALLREDUCE, MPI_REDUCE, MPI_IREDUCE,

- MPI_REDUCE_SCATTER_BLOCK, MPI_IREDUCE_SCATTER_BLOCK, MPI_REDUCE_SCATTER , MPI_IREDUCE_SCATTER.

In C++, the bindings for these functions are in the `MPI::Comm` class. However, since the collective operations do not make sense on a C++ `MPI::Comm` (as it is neither an intercommunicator nor an intracommunicator), the functions are all pure virtual.



Figure 5.2: 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 5.3: 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.

### 5.2.3  Specifics for Intercommunicator Collective Operations

All processes in both groups identified by the intercommunicator must call the collective routine.

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.

For intercommunicator collective communication, if the operation is in the All-To-One or One-To-All categories, 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. For this, 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 in the All-To-All category, then the transfer is bidirectional.

> *Rationale.* Operations in the All-To-One and One-To-All categories are unidirectional by nature, and there is a clear way of specifying direction. Operations in the All-To-All category will often occur as part of an exchange, where it makes sense to communicate in both directions at once. (*End of rationale.*)

## 5.3   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*(binding deprecated, see Section 15.2)* }

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, MPI_BARRIER involves two groups. The call returns at processes in one group (group A) of the intercommunicator only after all members of the other group (group B) have entered the call (and vice versa). A process may return from the call before all processes in its own group have entered the call.

## 5.4   Broadcast

MPI_BCAST(buffer, count, datatype, root, comm)

  INOUT    buffer                          starting address of buffer (choice)

  IN         count                         number of entries in buffer (non-negative 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*(binding deprecated, see Section 15.2)* }

If comm is an intracommunicator, MPI_BCAST broadcasts a message from the process with rank root to all processes of the group, itself included. It is called by all members of the group using the same arguments for comm and root. On return, the content of root's buffer is copied to all other 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 buffer arguments of the processes in group B must be consistent with the buffer argument of the root.

### 5.4.1  Example using MPI_BCAST

The examples in this section use intracommunicators.

**Example 5.1**

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.

## 5.5  Gather

MPI_GATHER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | recvcount | number of elements for any single receive (non-negative integer, significant only at root) |
| IN | recvtype | data type of recv buffer elements (significant only at root) (handle) |
| IN | root | rank of receiving process (integer) |
| IN | comm | communicator (handle) |

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

```
MPI_GATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
               ROOT, COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, ROOT, COMM, IERROR
```

```
{void MPI::Comm::Gather(const void* sendbuf, int sendcount, const
               MPI::Datatype& sendtype, void* recvbuf, int recvcount,
               const MPI::Datatype& recvtype, int root) const = 0(binding
               deprecated, see Section 15.2) }
```

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

$$MPI\_Recv(recvbuf + i \cdot recvcount \cdot extent(recvtype), recvcount, recvtype, i, ...),$$

where extent(recvtype) is the type extent obtained from a call to MPI_Type_get_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, recvtype, ...).
The receive buffer is ignored for all non-root processes.

**Unofficial Draft for Comment Only**

General, derived datatypes are allowed for both sendtype and recvtype. The type signature of sendcount, sendtype on each process must be equal to the type signature of recvcount, recvtype 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, and 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, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process (significant only at root) |
| IN | displs | integer array (of length group size). Entry i specifies the displacement relative to recvbuf at which to place the incoming data from process i (significant only at root) |
| IN | recvtype | data type of recv buffer elements (significant only at root) (handle) |
| IN | root | rank of receiving process (integer) |
| IN | comm | communicator (handle) |

```
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(binding
               deprecated, see Section 15.2) }
```

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,

$$\text{MPI\_Send}(\text{sendbuf}, \text{sendcount}, \text{sendtype}, \text{root}, ...),$$

and the root executes n receives,

$$\text{MPI\_Recv}(\text{recvbuf} + \text{displs}[j] \cdot \text{extent}(\text{recvtype}), \text{recvcounts}[j], \text{recvtype}, i, ...).$$

The data received from process j is placed into recvbuf of the root process beginning at offset displs[j] elements (in terms of the recvtype).

The receive buffer is ignored for all non-root processes.

The type signature implied by sendcount, sendtype on process i must be equal to the type signature implied by recvcounts[i], recvtype 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, as illustrated in Example 5.6.

All arguments to the function are significant on process root, while on other processes, only arguments sendbuf, sendcount, sendtype, root, and 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 written more than once. Such a call is erroneous.

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.

ticket265.

MPI_GATHERW(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtypes, root, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process (significant only at root) |
| IN | displs | integer array (of length group size). Entry i specifies the displacement relative to recvbuf at which to place the incoming data from process i (significant only at root) |
| IN | recvtypes | array of datatypes (of length group size). Entry i specifies the type of data received from process i (array of handles) |
| IN | root | rank of receiving process (integer) |
| IN | comm | communicator (handle) |

```
int MPI_Gatherw(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int *recvcounts, int *displs,
            MPI_Datatype *recvtype, int root, MPI_Comm comm)
```

```
MPI_GATHERW(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
            RECVTYPES, ROOT, COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPES(*),
    ROOT, COMM, IERROR
```

MPI_GATHERW extends the functionality of MPI_GATHERV by allowing varying datatype specifications for data received from each process, since datatypes is now an array. If comm is an intracommunicator, the outcome is as if each process, including the root process, sends a message to the root,

$$\mathrm{MPI\_Send}(\texttt{sendbuf}, \texttt{sendcount}, \texttt{sendtype}, \texttt{root}, ...),$$

and the root executes n receives,

$$\mathrm{MPI\_Recv}(\texttt{recvbuf} + \texttt{displs}[j] \cdot \texttt{extent}(\texttt{recvtypes}[j]), \texttt{recvcounts}[j], \texttt{recvtypes}[j], i, ...).$$

The data received from process j is placed into recvbuf of the root process beginning at offset displs[j] elements (in terms of recvtypes[j]).

The receive buffer is ignored for all non-root processes.

The type signature implied by sendcount, sendtype on process i must be equal to the type signature implied by recvcounts[i], recvtypes[i] 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.

All arguments to the function are significant on process root, while on other processes, only arguments sendbuf, sendcount, sendtype, root, and 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 written more than once. Such a call is erroneous.

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 arguments of the root.

## 5.5.1  Examples using MPI_GATHER, MPI_GATHERV

The examples in this section use intracommunicators.

**Example 5.2**
    Gather 100 `ints` from every process in group to root. See [f]Figure 5.4.

```
MPI_Comm comm;
int gsize,sendarray[100];
int root, *rbuf;
...
MPI_Comm_size(comm, &gsize);
rbuf = (int *)malloc(gsize*100*sizeof(int));
MPI_Gather(sendarray, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);
```
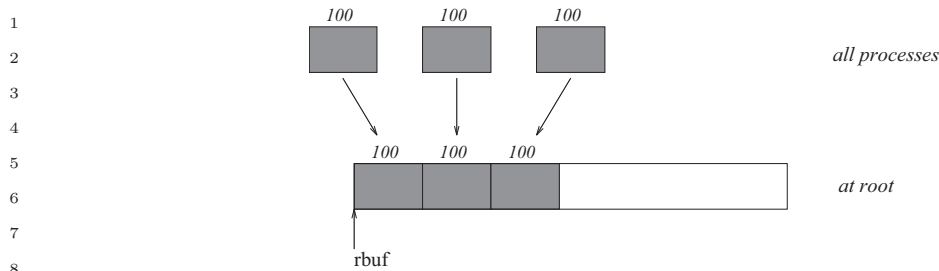
**Example 5.3**
    Previous example modified – only the root allocates memory for the receive buffer.

```
MPI_Comm comm;
int gsize,sendarray[100];
int root, myrank, *rbuf;
...
MPI_Comm_rank(comm, &myrank);
if (myrank == root) {
   MPI_Comm_size(comm, &gsize);
   rbuf = (int *)malloc(gsize*100*sizeof(int));
}
MPI_Gather(sendarray, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);
```

ticket0.

Figure 5.4: The root process gathers 100 `ints` from each process in the group.

**Example 5.4**

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 5.5**

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 $stride \geq 100$. See Figure 5.5.

```
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;
}
MPI_Gatherv(sendarray, 100, MPI_INT, rbuf, rcounts, displs, MPI_INT,
                                        root, comm);
```
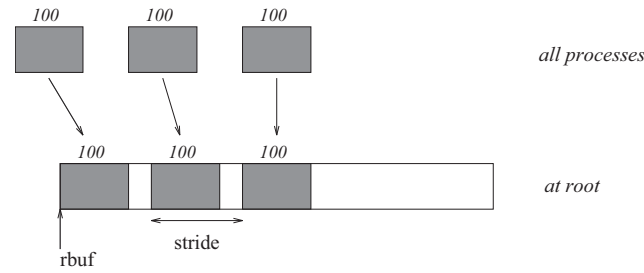
Figure 5.5: The root process gathers 100 `int`s from each process in the group, each set is placed `stride int`s apart.

Note that the program is erroneous if $stride < 100$.

**Example 5.6**
Same as Example 5.5 on the receiving side, but send the 100 `int`s from the 0th column of a 100×150 `int` array, in C. See Figure 5.6.

```
MPI_Comm comm;
int gsize,sendarray[100][150];
int root, *rbuf, stride;
MPI_Datatype stype;
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;
}
/* Create datatype for 1 column of array
 */
MPI_Type_vector(100, 1, 150, MPI_INT, &stype);
MPI_Type_commit(&stype);
MPI_Gatherv(sendarray, 1, stype, rbuf, rcounts, displs, MPI_INT,
                                            root, comm);
```

**Example 5.7**
Process `i` sends (100-i) `int`s from the `i`-th column of a $100 \times 150$ `int` array, in C. It is received into a buffer with stride, as in the previous two examples. See Figure 5.7.
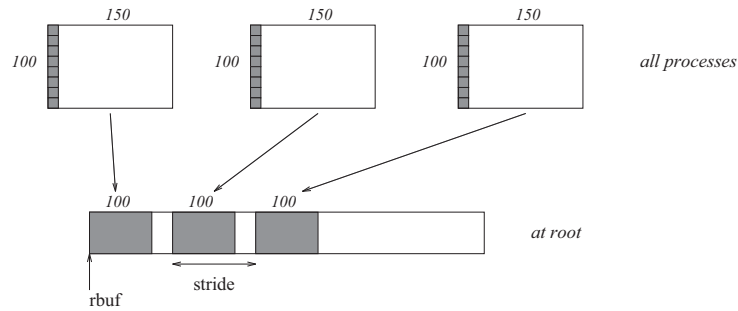
Figure 5.6: The root process gathers column 0 of a 100×150 C array, and each set is placed `stride int`s apart.

```
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 5.8**
    Same as Example 5.7, 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 we read a column of a C array. A similar thing was done in Example 4.16, Section 4.1.14.

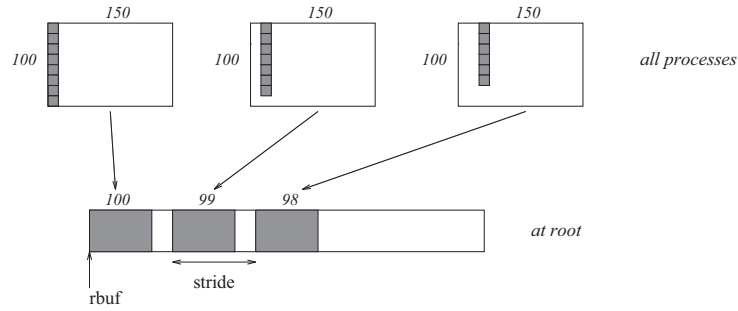```
MPI_Comm comm;
int gsize,sendarray[100][150],*sptr;
```

Figure 5.7: The root process gathers `100-i ints` from column `i` of a $100\times150$ C array, and each set is placed `stride ints` apart.

```
int root, *rbuf, stride, myrank, disp[2], blocklen[2];
MPI_Datatype stype,type[2];
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;
}
/* Create datatype for one int, with extent of entire row
 */
disp[0] = 0;        disp[1] = 150*sizeof(int);
type[0] = MPI_INT; type[1] = MPI_UB;
blocklen[0] = 1;   blocklen[1] = 1;
MPI_Type_create_struct(2, blocklen, disp, type, &stype);
MPI_Type_commit(&stype);
sptr = &sendarray[0][myrank];
MPI_Gatherv(sptr, 100-myrank, stype, rbuf, rcounts, displs, MPI_INT,
                                                    root, comm);
```

**Example 5.9**

Same as Example 5.7 at sending side, but at receiving side we make the stride between received blocks vary from block to block. See Figure 5.8.

```
MPI_Comm comm;
int gsize,sendarray[100][150],*sptr;
int root, *rbuf, *stride, myrank, bufsize;
MPI_Datatype stype;
int *displs,i,*rcounts,offset;
```
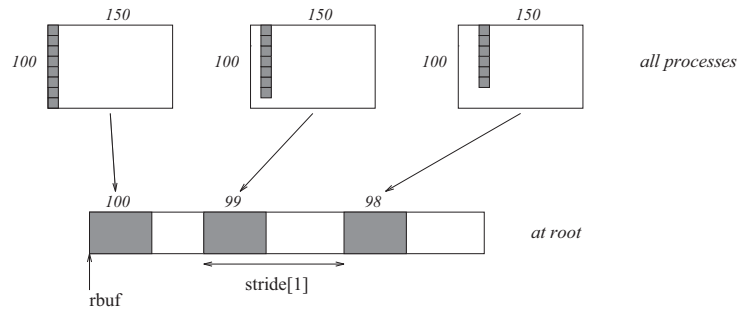
Figure 5.8: The root process gathers `100-i ints` from column `i` of a $100 \times 150$ C array, and each set is placed `stride[i] ints` apart (a varying stride).

```
    ...

    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 5.10**

Process `i` sends `num` `ints` from the `i`-th column of a $100 \times 150$ `int` array, in C. The complicating factor is that the various values of `num` are not known to `root`, so a separate gather must first be run to find these out. The data is placed contiguously at the receiving end.

```
MPI_Comm comm;
int gsize,sendarray[100][150],*sptr;
int root, *rbuf, myrank, disp[2], blocklen[2];
MPI_Datatype stype,type[2];
int *displs,i,*rcounts,num;

...

MPI_Comm_size(comm, &gsize);
MPI_Comm_rank(comm, &myrank);

/* First, gather nums to root
 */
rcounts = (int *)malloc(gsize*sizeof(int));
MPI_Gather(&num, 1, MPI_INT, rcounts, 1, MPI_INT, root, comm);
/* root now has correct rcounts, using these we set displs[] so
 * that data is placed contiguously (or concatenated) at receive end
 */
displs = (int *)malloc(gsize*sizeof(int));
displs[0] = 0;
for (i=1; i<gsize; ++i) {
    displs[i] = displs[i-1]+rcounts[i-1];
}
/* And, create receive buffer
 */
rbuf = (int *)malloc(gsize*(displs[gsize-1]+rcounts[gsize-1])
                                              *sizeof(int));
/* Create datatype for one int, with extent of entire row
 */
disp[0] = 0;      disp[1] = 150*sizeof(int);
type[0] = MPI_INT; type[1] = MPI_UB;
blocklen[0] = 1;   blocklen[1] = 1;
MPI_Type_create_struct( 2, blocklen, disp, type, &stype );
MPI_Type_commit(&stype);
sptr = &sendarray[0][myrank];
MPI_Gatherv(sptr, num, stype, rbuf, rcounts, displs, MPI_INT,
                                              root, comm);
```

## 5.6   Scatter

MPI_SCATTER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)

| | | |
|---|---|---|
| IN | sendbuf | address of send buffer (choice, significant only at root) |
| IN | sendcount | number of elements sent to each process (non-negative integer, significant only at root) |
| IN | sendtype | data type of send buffer elements (significant only at root) (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |

```
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(binding
                deprecated, see Section 15.2) }
```

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,

$$\text{MPI\_Send}(\text{sendbuf} + i \cdot \text{sendcount} \cdot \text{extent}(\text{sendtype}), \text{sendcount}, \text{sendtype}, i, ...),$$

and each process executed a receive,

$$\text{MPI\_Recv}(\text{recvbuf}, \text{recvcount}, \text{recvtype}, i, ...).$$

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.

All arguments to the function are significant on process root, while on other processes, only arguments recvbuf, recvcount, recvtype, root, and 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 read more than once.

> *Rationale.* Though not needed, the last restriction is imposed so as to achieve symmetry with MPI_GATHER, where the corresponding restriction (a multiple-write restriction) is necessary. (*End of rationale.*)

The "in place" option for intracommunicators is specified by passing MPI_IN_PLACE as the value of recvbuf at the root. In such a case, recvcount and recvtype 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.

MPI_SCATTERV(sendbuf, sendcounts, displs, sendtype, recvbuf, recvcount, recvtype, root, comm)

| | | |
|------|-----------|---|
| IN | sendbuf | address of send buffer (choice, significant only at root) |
| IN | sendcounts | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN | displs | integer array (of length group size). Entry i specifies the displacement (relative to sendbuf) from which to take the outgoing data to process i |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |

ticket109.

```
int MPI_Scatterv(void* sendbuf, int *sendcounts, int *displs,
            MPI_Datatype sendtype, void* recvbuf, int recvcount,
            MPI_Datatype recvtype, int root, MPI_Comm comm)
```

**Unofficial Draft for Comment Only**

```
MPI_SCATTERV(SENDBUF, SENDCOUNTS, DISPLS, SENDTYPE, RECVBUF, RECVCOUNT,
             RECVTYPE, ROOT, COMM, IERROR)
   <type> SENDBUF(*), RECVBUF(*)
   INTEGER SENDCOUNTS(*), DISPLS(*), SENDTYPE, RECVCOUNT, RECVTYPE, 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& recvtype,
             int root) const = 0(binding deprecated, see Section 15.2) }
```

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 an additional argument, displs.

If comm is an intracommunicator, the outcome is as if the root executed n send operations,

$$MPI\_Send(sendbuf + displs[i] \cdot extent(sendtype), sendcounts[i], sendtype, i, ...),$$

and each process executed a receive,

$$MPI\_Recv(recvbuf, recvcount, recvtype, i, ...).$$

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, recvtype 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, recvtype, root, and 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 a case, recvcount and recvtype 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.

ticket265.

MPI_SCATTERW(sendbuf, sendcount, displs, sendtypes, recvbuf, recvcount, recvtype, root, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice, significant only at root) |
| IN | sendcount | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN | displs | integer array (of length group size). Entry i specifies the displacement relative to sendbuf from which to take the outgoing data to process i |
| IN | sendtypes | array of datatypes (of length group size). Entry t j specifies the type of data to send to process t j (array of handles) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |

```
int MPI_Scatterw(void* sendbuf, int sendcounts[], int displs[],
            MPI_Datatype sendtypes[], void* recvbuf, int *recvcount,
            MPI_Datatype *recvtype, int root, MPI_Comm comm)
```

```
MPI_SCATTERW(SENDBUF, SENDCOUNTS, DISPLS, SENDTYPES, RECVBUF, RECVCOUNT,
            RECVTYPE, ROOT, COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER DISPLS(*), SENDCOUNTS(*), SENDTYPES, RECVCOUNT, RECVTYPE, ROOT,
    COMM, IERROR
```

MPI_SCATTERW is the inverse operation to MPI_GATHERW.

MPI_SCATTERW extends the functionality of MPI_SCATTERV by allowing varying datatype specifications for data sent to each process, since datatypes is now an array.

If comm is an intracommunicator, the outcome is as if the root executed n send operations,

```
MPI_Send(sendbuf + displs[i] extent (sendtypes[i]), sendcounts[i], sendtype, i, ...),
```

and each process executed a receive,

```
MPI_Recv(recvbuf, recvcount, recvtype, i, ...).
```

The send buffer is ignored for all non-root processes.

The type signature implied by sendcount[i], sendtypes[i] at the root must be equal to the type signature implied by recvcount, recvtype 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.
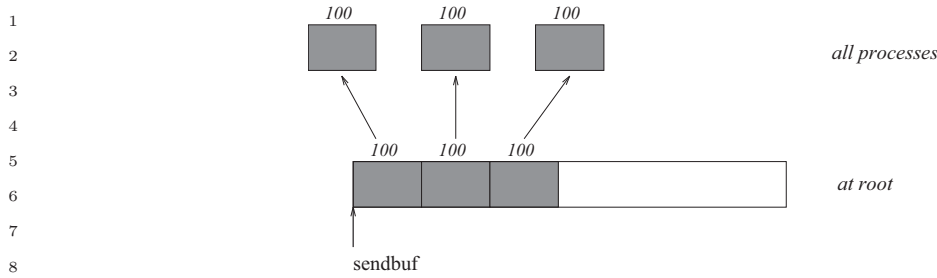
Figure 5.9: The root process scatters sets of 100 `int`s to each process in the group.

All arguments to the function are significant on process **root**, while on other processes, only arguments **recvbuf**, **recvcount**, **recvtype**, **root**, and **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 a case, **recvcount** and **recvtype** 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 arguments of the root.

## 5.6.1   Examples using MPI_SCATTER, MPI_SCATTERV

The examples in this section use intracommunicators.

**Example 5.11**
The reverse of Example 5.2. Scatter sets of 100 `int`s from the root to each process in the group. See Figure 5.9.

```
MPI_Comm comm;
int gsize,*sendbuf;
int root, rbuf[100];
...
MPI_Comm_size(comm, &gsize);
sendbuf = (int *)malloc(gsize*100*sizeof(int));
...
MPI_Scatter(sendbuf, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);
```
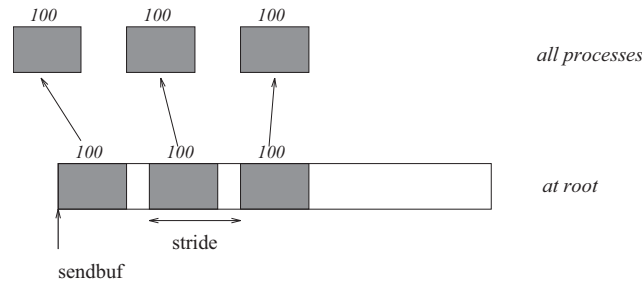
**Example 5.12**

Figure 5.10: The root process scatters sets of 100 `int`s, moving by `stride int`s from send to send in the scatter.

The reverse of Example 5.5. The root process scatters sets of 100 `int`s to the other processes, but the sets of 100 are *stride int*s apart in the sending buffer. Requires use of MPI_SCATTERV. Assume $stride \geq 100$. See Figure 5.10.

```
MPI_Comm comm;
int gsize,*sendbuf;
int root, rbuf[100], i, *displs, *scounts;

...

MPI_Comm_size(comm, &gsize);
sendbuf = (int *)malloc(gsize*stride*sizeof(int));
...
displs = (int *)malloc(gsize*sizeof(int));
scounts = (int *)malloc(gsize*sizeof(int));
for (i=0; i<gsize; ++i) {
    displs[i] = i*stride;
    scounts[i] = 100;
}
MPI_Scatterv(sendbuf, scounts, displs, MPI_INT, rbuf, 100, MPI_INT,
                                              root, comm);
```

**Example 5.13**

The reverse of Example 5.9. 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 5.11.

```
MPI_Comm comm;
int gsize,recvarray[100][150],*rptr;
int root, *sendbuf, myrank, *stride;
MPI_Datatype rtype;
int i, *displs, *scounts, offset;
...
MPI_Comm_size(comm, &gsize);
MPI_Comm_rank(comm, &myrank);
```
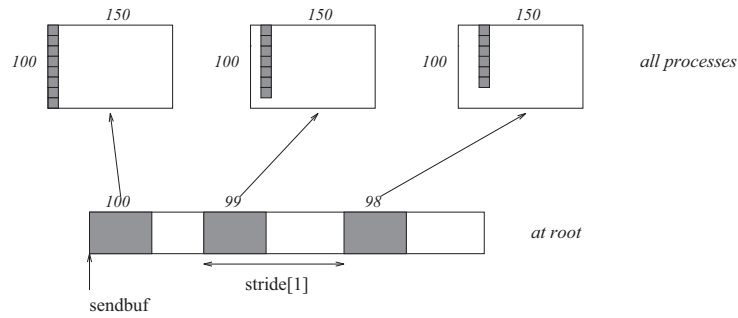
Figure 5.11: The root scatters blocks of `100-i` `int`s into column `i` of a $100{\times}150$ C array. At the sending side, the blocks are `stride[i]` `int`s apart.

```
    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));
    scounts = (int *)malloc(gsize*sizeof(int));
    offset = 0;
    for (i=0; i<gsize; ++i) {
        displs[i] = offset;
        offset += stride[i];
        scounts[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, scounts, displs, MPI_INT, rptr, 1, rtype,
                                              root, comm);
```

## 5.7 Gather-to-all

MPI_ALLGATHER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, comm)

| IN | sendbuf | starting address of send buffer (choice) |
|----|---------|------------------------------------------|
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements received from any process (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | comm | communicator (handle) |

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

```
MPI_ALLGATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
            COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, COMM, IERROR
```

```
{void MPI::Comm::Allgather(const void* sendbuf, int sendcount, const
            MPI::Datatype& sendtype, void* recvbuf, int recvcount,
            const MPI::Datatype& recvtype) const = 0(binding deprecated, see
            Section 15.2) }
```

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 of one group (group A) contributes sendcount data items; these data are concatenated and the result is stored at each process

in the other group (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, recvcounts, displs, recvtype, comm)

| | | |
|------|------------|-----------------------------------------------------|
| IN   | sendbuf    | starting address of send buffer (choice)            |
| IN   | sendcount  | number of elements in send buffer (non-negative integer) |
| IN   | sendtype   | data type of send buffer elements (handle)          |
| OUT  | recvbuf    | address of receive buffer (choice)                  |
| IN   | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process |
| IN   | displs     | integer array (of length group size). Entry i specifies the displacement (relative to recvbuf) at which to place the incoming data from process i |
| IN   | recvtype   | data type of receive buffer elements (handle)       |
| IN   | comm       | communicator (handle)                               |

```
int MPI_Allgatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
              void* recvbuf, int *recvcounts, int *displs,
              MPI_Datatype recvtype, MPI_Comm comm)
```

```
MPI_ALLGATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
              RECVTYPE, COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPE, COMM,
    IERROR
```

```
{void MPI::Comm::Allgatherv(const void* sendbuf, int sendcount, const
              MPI::Datatype& sendtype, void* recvbuf,
              const int recvcounts[], const int displs[],
              const MPI::Datatype& recvtype) const = 0(binding deprecated, see
              Section 15.2) }
```

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], recvtype at any other process.

If comm is an intracommunicator, the outcome is as if all processes executed calls to

MPI_GATHERV($sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtype, root, comm$),

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. In such a case, sendcount and sendtype are ignored, and 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 of one group (group A) contributes sendcount data items; these data are concatenated and the result is stored at each process in the other group (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.                                ticket265.

MPI_ALLGATHERW(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtypes, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process |
| IN | displs | integer array (of length group size). Entry i specifies the displacement (relative to recvbuf) at which to place the incoming data from process i |
| IN | recvtypes | array of datatypes (of length group size). Entry i specifies the type of data received from process i (array of handles) |
| IN | comm | communicator (handle) |

```
int MPI_Allgatherw(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int *recvcounts, int *displs,
            MPI_Datatype *recvtypes, MPI_Comm comm)
```

```
MPI_ALLGATHERW(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
            RECVTYPES, COMM, IERROR)
```

```
<type> SENDBUF(*), RECVBUF(*)
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPES(*),
COMM, IERROR
```

MPI_ALLGATHERW can be thought of as MPI_GATHERW, 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], recvtypes[j] at any other process.

If comm is an intracommunicator, the outcome is as if all processes executed calls to

MPI_GATHERW(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtypes, root, comm)

for root = 0 , ..., n-1. The rules for correct usage of MPI_ALLGATHERW are easily found from the corresponding rules for MPI_GATHERW.

The "in place" option for intracommunicators is specified by passing the value MPI_IN_PLACE to the argument sendbuf at all processes. In such a case, sendcount and sendtype are ignored, and 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 of one group (group A) contributes sendcount data items; these data are concatenated and the result is stored at each process in the other group (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.

## 5.7.1   Example using MPI_ALLGATHER

The example in this section uses intracommunicators.

**Example 5.14**
    The all-gather version of Example 5.2. Using MPI_ALLGATHER, we will gather 100 ints from every process in the group to every process.

```
MPI_Comm comm;
int gsize,sendarray[100];
int *rbuf;
...
MPI_Comm_size(comm, &gsize);
rbuf = (int *)malloc(gsize*100*sizeof(int));
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.

## 5.8 All-to-All Scatter/Gather

MPI_ALLTOALL(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements sent to each process (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements received from any process (non-negative 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 recvcount, 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 recvcount,
            const MPI::Datatype& recvtype) const = 0*(binding deprecated, see
            Section 15.2)* }

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 recvcount, 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,

$$\text{MPI\_Send}(\text{sendbuf} + i \cdot \text{sendcount} \cdot \text{extent(sendtype)}, \text{sendcount}, \text{sendtype}, i, ...),$$

and a receive from every other process with a call to,

$$\text{MPI\_Recv}(\text{recvbuf} + i \cdot \text{recvcount} \cdot \text{extent(recvtype)}, \text{recvcount}, \text{recvtype}, i, ...).$$

All arguments on all processes are significant. The argument comm must have identical values on all processes.

The "in place" option for intracommunicators is specified by passing MPI_IN_PLACE to the argument sendbuf at *all* processes. In such a case, sendcount and sendtype are ignored.

The data to be sent is taken from the recvbuf and replaced by the received data. Data sent and received must have the same type map as specified by recvcount and recvtype.

> *Rationale.*    For large MPI_ALLTOALL instances, allocating both send and receive buffers may consume too much memory. The "in place" option effectively halves the application memory consumption and is useful in situations where the data to be sent will not be used by the sending process after the MPI_ALLTOALL exchange (e.g., in parallel Fast Fourier Transforms). (*End of rationale.*)

> *Advice to implementors.*    Users may opt to use the "in place" option in order to conserve memory. Quality MPI implementations should thus strive to minimize system buffering. (*End of advice to implementors.*)

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 a complete exchange 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,   recvcounts,   rdispls, recvtype, comm)

| | | |
|----|----------|----------------------------------------------|
| IN  | sendbuf    | starting address of send buffer (choice) |
| IN  | sendcounts | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN  | sdispls    | integer array (of length group size). Entry j specifies the displacement (relative to sendbuf) from which to take the outgoing data destined for process j |
| IN  | sendtype   | data type of send buffer elements (handle) |
| OUT | recvbuf    | address of receive buffer (choice) |
| IN  | recvcounts | non-negative integer array (of length group size) specifying the number of elements that can be received from each processor |
| IN  | rdispls    | integer array (of length group size). Entry i specifies the displacement (relative to recvbuf) at which to place the incoming data from process i |
| IN  | recvtype   | data type of receive buffer elements (handle) |
| IN  | comm       | communicator (handle) |

ticket109.

ticket109.

```
int MPI_Alltoallv(void* sendbuf, int *sendcounts, int *sdispls,
             MPI_Datatype sendtype, void* recvbuf, int *recvcounts,
```

```
                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 recvcounts[], const int rdispls[],
            const MPI::Datatype& recvtype) const = 0*(binding deprecated, see
            Section 15.2) }*

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 sendcounts[j], sendtype at process i must be equal to the type signature associated with recvcounts[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,

$$\text{MPI\_Send}(\text{sendbuf} + \text{sdispls}[i] \cdot \text{extent}(\text{sendtype}), \text{sendcounts}[i], \text{sendtype}, i, ...),$$

and received a message from every other process with a call to

$$\text{MPI\_Recv}(\text{recvbuf} + \text{rdispls}[i] \cdot \text{extent}(\text{recvtype}), \text{recvcounts}[i], \text{recvtype}, i, ...).$$

All arguments on all processes are significant. The argument comm must have identical values on all processes.

The "in place" option for intracommunicators is specified by passing MPI_IN_PLACE to the argument sendbuf at *all* processes. In such a case, sendcounts, sdispls and sendtype are ignored. The data to be sent is taken from the recvbuf and replaced by the received data. Data sent and received must have the same type map as specified by the recvcounts array and the recvtype, and is taken from the locations of the receive buffer specified by rdispls.

> *Advice to users.* Specifying the "in place" option (which must be given on all processes) implies that the same amount and type of data is sent and received between any two processes in the group of the communicator. Different pairs of processes can exchange different amounts of data. Users must ensure that recvcounts[j] and recvtype on process i match recvcounts[i] and recvtype on process j. This symmetric exchange can be useful in applications where the data to be sent will not be used by the sending process after the MPI_ALLTOALLV exchange. (*End of advice to users.*)

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.*)

MPI_ALLTOALLW(sendbuf, sendcounts, sdispls, sendtypes, recvbuf, recvcounts, rdispls, recvtypes, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcounts | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN | sdispls | integer array (of length group size). Entry j specifies the displacement in bytes (relative to sendbuf) from which to take the outgoing data destined for process j (array of integers) |
| IN | sendtypes | array of datatypes (of length group size). Entry j specifies the type of data to send to process j (array of handles) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) specifying the number of elements that can be received from each processor |
| IN | rdispls | integer array (of length group size). Entry i specifies the displacement in bytes (relative to recvbuf) at which to place the incoming data from process i (array of integers) |
| IN | recvtypes | array of datatypes (of length group size). Entry i specifies the type of data received from process i (array of handles) |
| IN | comm | communicator (handle) |

```
int MPI_Alltoallw(void* sendbuf, int sendcounts[], int sdispls[],
            MPI_Datatype sendtypes[], void* recvbuf, int recvcounts[],
            int rdispls[], MPI_Datatype recvtypes[], MPI_Comm comm)
```

```
MPI_ALLTOALLW(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPES, RECVBUF, RECVCOUNTS,
            RDISPLS, RECVTYPES, COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
```

```
        INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPES(*), RECVCOUNTS(*),
        RDISPLS(*), RECVTYPES(*), COMM, IERROR
```

```
{void MPI::Comm::Alltoallw(const void* sendbuf, const int sendcounts[],
            const int sdispls[], const MPI::Datatype sendtypes[], void*
            recvbuf, const int recvcounts[], const int rdispls[], const
            MPI::Datatype recvtypes[]) const = 0(binding deprecated, see
            Section 15.2) }
```

MPI_ALLTOALLW is the most general form of complete exchange. 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.

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 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

$$\text{MPI\_Send}(\text{sendbuf} + \text{sdispls}[i], \text{sendcounts}[i], \text{sendtypes}[i], i, ...),$$

and received a message from every other process with a call to

$$\text{MPI\_Recv}(\text{recvbuf} + \text{rdispls}[i], \text{recvcounts}[i], \text{recvtypes}[i], i, ...).$$

All arguments on all processes are significant. The argument comm must describe the same communicator on all processes.

Like for MPI_ALLTOALLV, the "in place" option for intracommunicators is specified by passing MPI_IN_PLACE to the argument sendbuf at *all* processes. In such a case, sendcounts, sdispls and sendtypes are ignored. The data to be sent is taken from the recvbuf and replaced by the received data. Data sent and received must have the same type map as specified by the recvcounts and recvtypes arrays, and is taken from the locations of the receive buffer specified by rdispls.

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 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.*)

## 5.9   Global Reduction Operations

The functions in this section perform a global reduce operation (for example sum, maximum, and logical and) across all 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 to one member of a group, an all-reduce that returns this result to all members of a group, and two scan (parallel prefix) operations. In addition, a reduce-scatter operation combines the functionality of a reduce and of a scatter operation.

### 5.9.1 Reduce

MPI_REDUCE(sendbuf, recvbuf, count, datatype, op, root, comm)

| | | |
|---|---|---|
| IN | sendbuf | address of send buffer (choice) |
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | count | number of elements in send buffer (non-negative integer) |
| IN | datatype | data type of elements of send buffer (handle) |
| IN | op | reduce operation (handle) |
| IN | root | rank of root process (integer) |
| IN | comm | 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)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER COUNT, DATATYPE, OP, ROOT, COMM, IERROR
```

```
{void MPI::Comm::Reduce(const void* sendbuf, void* recvbuf, int count,
               const MPI::Datatype& datatype, const MPI::Op& op, int root)
               const = 0(binding deprecated, see Section 15.2) }
```

If comm is an intracommunicator, 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)).

Section 5.9.2, lists the set of predefined operations provided by MPI. That section also enumerates the datatypes to which each operation can be applied.

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 Section 5.9.5.

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.*)

> *Advice to users.* Some applications may not be able to ignore the non-associative nature of floating-point operations or may use user-defined operations (see Section 5.9.5) that require a special reduction order and cannot be treated as associative. Such applications should enforce the order of evaluation explicitly. For example, in the case of operations that require a strict left-to-right (or right-to-left) evaluation order, this could be done by gathering all operands at a single process (e.g., with MPI_GATHER), applying the reduction operation in the desired order (e.g., with MPI_REDUCE_LOCAL), and if needed, broadcast or scatter the result to the other processes (e.g., with MPI_BCAST). (*End of advice to users.*)

The datatype argument of MPI_REDUCE must be compatible with op. Predefined operators work only with the MPI types listed in Section 5.9.2 and Section 5.9.4. 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 5.9.5.

> *Advice to users.* Users should make no assumptions about how MPI_REDUCE is implemented. It is safest 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 a 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.

**Unofficial Draft for Comment Only**

### 5.9.2   Predefined Reduction Operations

The following predefined operations are supplied for MPI_REDUCE and related functions MPI_ALLREDUCE, MPI_REDUCE_SCATTER, MPI_SCAN, and MPI_EXSCAN. These operations are invoked by placing the following in op.

| Name | Meaning |
|------|---------|
| MPI_MAX | maximum |
| MPI_MIN | minimum |
| MPI_SUM | sum |
| MPI_PROD | product |
| MPI_LAND | logical and |
| MPI_BAND | bit-wise and |
| MPI_LOR | logical or |
| MPI_BOR | bit-wise or |
| MPI_LXOR | logical exclusive or (xor) |
| MPI_BXOR | bit-wise exclusive or (xor) |
| MPI_MAXLOC | max value and location |
| MPI_MINLOC | min value and location |

The two operations MPI_MINLOC and MPI_MAXLOC are discussed separately in Section 5.9.4. 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: | MPI_INT, MPI_LONG, MPI_SHORT, MPI_UNSIGNED_SHORT, MPI_UNSIGNED, MPI_UNSIGNED_LONG, MPI_LONG_LONG_INT, MPI_LONG_LONG (as synonym), MPI_UNSIGNED_LONG_LONG, MPI_SIGNED_CHAR, MPI_UNSIGNED_CHAR, MPI_INT8_T, MPI_INT16_T, MPI_INT32_T, MPI_INT64_T, MPI_UINT8_T, MPI_UINT16_T, MPI_UINT32_T, MPI_UINT64_T |
| Fortran integer: | MPI_INTEGER, MPI_AINT, MPI_COUNT, MPI_OFFSET, and handles returned from MPI_TYPE_CREATE_F90_INTEGER, and if available: MPI_INTEGER1, MPI_INTEGER2, MPI_INTEGER4, MPI_INTEGER8, MPI_INTEGER16 |
| Floating point: | MPI_FLOAT, MPI_DOUBLE, MPI_REAL, MPI_DOUBLE_PRECISION MPI_LONG_DOUBLE and handles returned from |

ticket265.

**Unofficial Draft for Comment Only**

|  | MPI_TYPE_CREATE_F90_REAL, |
|---|---|
|  | and if available: MPI_REAL2, |
|  | MPI_REAL4, MPI_REAL8, MPI_REAL16 |
| Logical: | MPI_LOGICAL, MPI_C_BOOL |
| Complex: | MPI_COMPLEX, |
|  | MPI_C_FLOAT_COMPLEX, |
|  | MPI_C_DOUBLE_COMPLEX, |
|  | MPI_C_LONG_DOUBLE_COMPLEX, |
|  | and handles returned from |
|  | MPI_TYPE_CREATE_F90_COMPLEX, |
|  | and if available: MPI_DOUBLE_COMPLEX, |
|  | MPI_COMPLEX4, MPI_COMPLEX8, |
|  | MPI_COMPLEX16, MPI_COMPLEX32 |
| Byte: | MPI_BYTE |

Now, the valid datatypes for each option is specified below.

| Op | Allowed Types |
|---|---|
| MPI_MAX, MPI_MIN | C integer, Fortran integer, Floating point |
| MPI_SUM, MPI_PROD | C integer, Fortran integer, Floating point, Complex |
| MPI_LAND, MPI_LOR, MPI_LXOR | C integer, Logical |
| MPI_BAND, MPI_BOR, MPI_BXOR | C integer, Fortran integer, Byte |

The following examples use intracommunicators.

**Example 5.15**
A routine that computes the dot product of two vectors that are distributed across a group of processes and returns the answer at node zero.

```
SUBROUTINE PAR_BLAS1(m, a, b, c, comm)
REAL a(m), b(m)        ! local slice of array
REAL c                 ! result (at node zero)
REAL sum
INTEGER m, comm, i, ierr

! local sum
sum = 0.0
DO i = 1, m
   sum = sum + a(i)*b(i)
END DO

! global sum
CALL MPI_REDUCE(sum, c, 1, MPI_REAL, MPI_SUM, 0, comm, ierr)
RETURN
```

**Example 5.16**
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 node zero.

**Unofficial Draft for Comment Only**

```
1    SUBROUTINE PAR_BLAS2(m, n, a, b, c, comm)
2    REAL a(m), b(m,n)    ! local slice of array
3    REAL c(n)                ! result
4    REAL sum(n)
5    INTEGER n, comm, i, j, ierr
6
7    ! local sum
8    DO j= 1, n
9      sum(j) = 0.0
10     DO i = 1, m
11       sum(j) = sum(j) + a(i)*b(i,j)
12     END DO
13   END DO
14
15   ! global sum
16   CALL MPI_REDUCE(sum, c, n, MPI_REAL, MPI_SUM, 0, comm, ierr)
17
18   ! return result at node zero (and garbage at the other nodes)
19   RETURN
20
```

### 5.9.3  Signed Characters and Reductions

The types MPI_SIGNED_CHAR and MPI_UNSIGNED_CHAR can be used in reduction operations. MPI_CHAR, MPI_WCHAR, and MPI_CHARACTER (which represent printable characters) cannot be used in reduction operations. In a heterogeneous environment, MPI_CHAR, MPI_WCHAR, and MPI_CHARACTER 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.

> *Advice to users.*    The types MPI_CHAR, MPI_WCHAR, and MPI_CHARACTER are intended for characters, and so will be translated to preserve the printable representation, rather than the integer 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.*)

### 5.9.4  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

$$w = \max(u, v)$$

and

$$k = \begin{cases} i & \text{if } u > v \\ \min(i, j) & \text{if } u = v \\ j & \text{if } u < v \end{cases}$$

MPI_MINLOC is defined similarly:

$$\begin{pmatrix} u \\ i \end{pmatrix} \circ \begin{pmatrix} v \\ j \end{pmatrix} = \begin{pmatrix} w \\ k \end{pmatrix}$$

where

$$w = \min(u, v)$$

and

$$k = \begin{cases} i & \text{if } u < v \\ \min(i, j) & \text{if } u = v \\ j & \text{if } u > v \end{cases}$$

Both operations are associative and commutative. Note that if MPI_MAXLOC is applied to reduce a sequence of pairs $(u_0, 0), (u_1, 1), \ldots, (u_{n-1}, n-1)$, then the value returned is $(u, r)$, where $u = \max_i u_i$ and $r$ is the index of the first global maximum in the sequence. Thus, if each process supplies a value and its rank within the group, then a reduce operation with op = MPI_MAXLOC will return the maximum value and the rank of the first process with that value. Similarly, MPI_MINLOC can be used to return a minimum and its index. More generally, MPI_MINLOC computes a *lexicographic minimum*, where elements are ordered according to the first component of each pair, and ties are resolved according to the second component.

The reduce operation is defined to operate on arguments that consist of a pair: value and index. For both Fortran and C, types are provided to describe the pair. The potentially mixed-type nature of such arguments is a problem in Fortran. The problem is circumvented, for Fortran, by having the MPI-provided type consist of a pair of the same type as value, and coercing the index to this type also. In C, the MPI-provided pair type has distinct types and the index is an int.

In order to use MPI_MINLOC and MPI_MAXLOC in a reduce operation, one must provide a datatype argument that represents a pair (value and index). MPI provides nine such predefined datatypes. The operations MPI_MAXLOC and MPI_MINLOC can be used with each of the following datatypes.

Fortran:

| Name | Description |
|---|---|
| MPI_2REAL | pair of REALs |
| MPI_2DOUBLE_PRECISION | pair of DOUBLE PRECISION variables |
| MPI_2INTEGER | pair of INTEGERs |

C:

| Name | Description |
|---|---|
| MPI_FLOAT_INT | float and int |

**Unofficial Draft for Comment Only**

```
MPI_DOUBLE_INT                       double and int
MPI_LONG_INT                         long and int
MPI_2INT                             pair of int
MPI_SHORT_INT                        short and int
MPI_LONG_DOUBLE_INT                  long double and int
```

The datatype MPI_2REAL is *as if* defined by the following (see Section 4.1).

```
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_CREATE_STRUCT(2, block, disp, type, MPI_FLOAT_INT)
```

Similar statements apply for MPI_LONG_INT and MPI_DOUBLE_INT.
The following examples use intracommunicators.

**Example 5.17**
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
         */
```

**Unofficial Draft for Comment Only**

```
    for (i=0; i<30; ++i) {                                              1
        aout[i] = out[i].val;                                          2
        ind[i] = out[i].rank;                                          3
    }                                                                  4
}                                                                      5
                                                                       6
                                                                       7
```

**Example 5.18**

Same example, in Fortran.

```
...
! each process has an array of 30 double: ain(30)

DOUBLE PRECISION ain(30), aout(30)
INTEGER ind(30)
DOUBLE PRECISION in(2,30), out(2,30)
INTEGER i, myrank, root, ierr

CALL MPI_COMM_RANK(comm, myrank, ierr)
DO I=1, 30
    in(1,i) = ain(i)
    in(2,i) = myrank    ! myrank is coerced to a double
END DO

CALL MPI_REDUCE(in, out, 30, MPI_2DOUBLE_PRECISION, MPI_MAXLOC, root,
                                                    comm, ierr)
! At this point, the answer resides on process root

IF (myrank .EQ. root) THEN
    ! read ranks out
    DO I= 1, 30
        aout(i) = out(1,i)
        ind(i) = out(2,i)  ! rank is coerced back to an integer
    END DO
END IF
```

**Example 5.19**

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.

```
#define  LEN    1000

float val[LEN];       /* local array of values */
int count;            /* local number of values */
int myrank, minrank, minindex;
float minval;

struct {
```

```
    float value;
    int   index;
} in, out;

    /* 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.*)

### 5.9.5   User-Defined Reduction 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) *(binding deprecated, see Section 15.2)* }

MPI_OP_CREATE binds a user-defined reduction operation to an op handle that can subsequently be used in MPI_REDUCE, MPI_ALLREDUCE, MPI_REDUCE_SCATTER, MPI_SCAN, and MPI_EXSCAN. The user-defined operation is assumed to be associative. If commute = true, then the operation should be both commutative and associative. If commute = false, then the order of operands is fixed and is defined to be in ascending, process rank order, beginning with process zero. The order of evaluation can be changed, talking advantage of the associativity of the operation. If commute = true then the order of evaluation can be changed, taking advantage of commutativity and associativity.

The argument function is the user-defined function, which must have the following four arguments: invec, inoutvec, len and datatype.

The ISO C prototype for the function is the following.

```
typedef void MPI_User_function(void* invec, void* inoutvec, int *len,
          MPI_Datatype *datatype);
```

The Fortran declaration of the user-defined function appears below.

```
SUBROUTINE USER_FUNCTION(INVEC, INOUTVEC, LEN, TYPE)
    <type> INVEC(LEN), INOUTVEC(LEN)
    INTEGER LEN, TYPE
```

The C++ declaration of the user-defined function appears below.

{typedef void MPI::User_function(const void* invec, void* inoutvec, int len, const Datatype& datatype); *(binding deprecated, see Section 15.2)*}

The datatype argument is a handle to the data type that was passed into the call to MPI_REDUCE. The user reduce function should be written such that the following holds: Let u[0], ... , u[len-1] be the len elements in the communication buffer described by the arguments invec, len and datatype when the function is invoked; let v[0], ... , v[len-1] be len elements in the communication buffer described by the arguments inoutvec, len and datatype when the function is invoked; let w[0], ... , w[len-1] be len elements in the communication buffer described by the arguments inoutvec, len and datatype when the function returns; then w[i] = u[i]∘v[i], for i=0 , ... , len-1, where ∘ is the reduce operation that the function computes.

Informally, we can think of invec and inoutvec as arrays of len elements that function is combining. The result of the reduction over-writes values in inoutvec, hence the name. Each invocation of the function results in the pointwise evaluation of the reduce operator on len elements: i.e., the function returns in inoutvec[i] the value invec[i] ∘ inoutvec[i], for $i = 0, \dots, \text{count} - 1$, where ∘ is the combining operation computed by the function.

> *Rationale.* The len argument allows MPI_REDUCE to avoid calling the function for each element in the input buffer. Rather, the system can choose to apply the function to chunks of input. In C, it is passed in as a reference for reasons of compatibility with Fortran.

> By internally comparing the value of the datatype argument to known, global handles, it is possible to overload the use of a single user-defined function for several, different

data types. (*End of rationale.*)

General datatypes may be passed to the user function. However, use of datatypes that are not contiguous is likely to lead to inefficiencies.

No MPI communication function may be called inside the user function. MPI_ABORT 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 datatypes were created. Before the library is used, a library initialization preamble must be executed. This preamble code will define the datatypes that are used by the library, and store handles to these datatypes 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 not supporting the "in place" option.

```
MPI_Comm_size(comm, &groupsize);
MPI_Comm_rank(comm, &rank);
if (rank > 0) {
    MPI_Recv(tempbuf, count, datatype, rank-1,...);
    User_reduce(tempbuf, sendbuf, count, datatype);
}
if (rank < groupsize-1) {
    MPI_Send(sendbuf, count, datatype, rank+1, ...);
}
/* answer now resides in process groupsize-1 ... now send to root
 */
if (rank == root) {
    MPI_Irecv(recvbuf, count, datatype, groupsize-1,..., &req);
}
if (rank == groupsize-1) {
    MPI_Send(sendbuf, count, datatype, root, ...);
}
if (rank == root) {
    MPI_Wait(&req, &status);
}
```

The reduction computation proceeds, sequentially, from process 0 to process
`groupsize-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 im-
plementation 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 computa-
tion, 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.*)

MPI_OP_FREE(op)

  INOUT    op                                operation (handle)

```
int MPI_Op_free(MPI_Op *op)
```

```
MPI_OP_FREE(OP, IERROR)
    INTEGER OP, IERROR
```

{void MPI::Op::Free()*(binding deprecated, see Section 15.2)* }

Marks a user-defined reduction operation for deallocation and sets op to MPI_OP_NULL.

## Example of User-defined Reduce

It is time for an example of user-defined reduction. The example in this section uses an
intracommunicator.

**Example 5.20** Compute the product of an array of complex numbers, in C.

```
typedef struct {
    double real,imag;
} Complex;

/* the user-defined function
 */
void myProd(Complex *in, Complex *inout, int *len, MPI_Datatype *dptr)
{
    int i;
    Complex c;

    for (i=0; i< *len; ++i) {
        c.real = inout->real*in->real -
                   inout->imag*in->imag;
        c.imag = inout->real*in->imag +
                   inout->imag*in->real;
        *inout = c;
```

```
        in++; inout++;
    }
}

/* and, to call it...
 */
...

    /* each process has an array of 100 Complexes
     */
    Complex a[100], answer[100];
    MPI_Op myOp;
    MPI_Datatype ctype;

    /* explain to MPI how type Complex is defined
     */
    MPI_Type_contiguous(2, MPI_DOUBLE, &ctype);
    MPI_Type_commit(&ctype);
    /* create the complex-product user-op
     */
    MPI_Op_create( myProd, 1, &myOp );

    MPI_Reduce(a, answer, 100, ctype, myOp, root, comm);

    /* At this point, the answer, which consists of 100 Complexes,
     * resides on process root
     */
```

### 5.9.6  All-Reduce

MPI includes a variant of the reduce operations where the result is returned to all processes in a group. MPI requires that all processes from the same group participating in these operations receive identical results.

MPI_ALLREDUCE(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 send buffer (non-negative integer) |
| IN | datatype | data type of elements of send buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |

```
int MPI_Allreduce(void* sendbuf, void* recvbuf, int count,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

```
MPI_ALLREDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER COUNT, DATATYPE, OP, COMM, IERROR
```

{void MPI::Comm::Allreduce(const void* sendbuf, void* recvbuf, int count,
            const MPI::Datatype& datatype, const MPI::Op& op)
            const = 0*(binding deprecated, see Section 15.2)* }

If comm is an intracommunicator, MPI_ALLREDUCE behaves the 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 this 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 uses an intracommunicator.

**Example 5.21**

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 5.16).

```
SUBROUTINE PAR_BLAS2(m, n, a, b, c, comm)
REAL a(m), b(m,n)    ! local slice of array
REAL c(n)            ! result
REAL sum(n)
INTEGER n, comm, i, j, ierr

! local sum
DO j= 1, n
  sum(j) = 0.0
  DO i = 1, m
    sum(j) = sum(j) + a(i)*b(i,j)
  END DO
END DO

! global sum
CALL MPI_ALLREDUCE(sum, c, n, MPI_REAL, MPI_SUM, comm, ierr)

! return result at all nodes
RETURN
```

### 5.9.7  Process-[l]Local [r]Reduction

The functions in this section are of importance to library implementors who may want to
implement special reduction patterns that are otherwise not easily covered by the standard

MPI operations.

The following function applies a reduction operator to local arguments.

MPI_REDUCE_LOCAL( inbuf, inoutbuf, count, datatype, op)

| IN | inbuf | input buffer (choice) |
|---|---|---|
| INOUT | inoutbuf | combined input and output buffer (choice) |
| IN | count | number of elements in inbuf and inoutbuf buffers (non-negative integer) |
| IN | datatype | data type of elements of inbuf and inoutbuf buffers (handle) |
| IN | op | operation (handle) |

```
int MPI_Reduce_local(void* inbuf, void* inoutbuf, int count,
              MPI_Datatype datatype, MPI_Op op)
```

```
MPI_REDUCE_LOCAL(INBUF, INOUBUF, COUNT, DATATYPE, OP, IERROR)
    <type> INBUF(*), INOUTBUF(*)
    INTEGER COUNT, DATATYPE, OP, IERROR
```

`{void MPI::Op::Reduce_local(const void* inbuf, void* inoutbuf, int count, const MPI::Datatype& datatype) const`*(binding deprecated, see Section 15.2)* `}`

The function applies the operation given by op element-wise to the elements of inbuf and inoutbuf with the result stored element-wise in inoutbuf, as explained for user-defined operations in Section 5.9.5. Both inbuf and inoutbuf (input as well as result) have the same number of elements given by count and the same datatype given by datatype. The MPI_IN_PLACE option is not allowed.

Reduction operations can be queried for their commutativity.

MPI_OP_COMMUTATIVE( op, commute)

| IN | op | operation (handle) |
|---|---|---|
| OUT | commute | true if op is commutative, false otherwise (logical) |

```
int MPI_Op_commutative(MPI_Op op, int *commute)
```

```
MPI_OP_COMMUTATIVE(OP, COMMUTE, IERROR)
    LOGICAL COMMUTE
    INTEGER OP, IERROR
```

`{bool MPI::Op::Is_commutative() const`*(binding deprecated, see Section 15.2)* `}`

## 5.10 Reduce-Scatter

MPI includes variants of the reduce operations where the result is scattered to all processes in a group on return. One variant scatters equal-sized blocks to all processes, while another variant scatters blocks that may vary in size for each process.

### 5.10.1 MPI_REDUCE_SCATTER_BLOCK

MPI_REDUCE_SCATTER_BLOCK( sendbuf, recvbuf, recvcount, datatype, op, comm)

| | | |
|----|----|----|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | recvcount | element count per block (non-negative integer) |
| IN | datatype | data type of elements of send and receive buffers (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |

```
int MPI_Reduce_scatter_block(void* sendbuf, void* recvbuf, int recvcount,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

```
MPI_REDUCE_SCATTER_BLOCK(SENDBUF, RECVBUF, RECVCOUNT, DATATYPE, OP, COMM,
            IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER RECVCOUNT, DATATYPE, OP, COMM, IERROR
```

{void MPI::Comm::Reduce_scatter_block(const void* sendbuf, void* recvbuf,
            int recvcount, const MPI::Datatype& datatype,
            const MPI::Op& op) const = 0*(binding deprecated, see Section 15.2)* }

If comm is an intracommunicator, MPI_REDUCE_SCATTER_BLOCK first performs a global, element-wise reduction on vectors of count = n*recvcount elements in the send buffers defined by sendbuf, count and datatype, using the operation op, where n is the number of processes in the group of comm. The routine is called by all group members using the same arguments for recvcount, datatype, op and comm. The resulting vector is treated as n consecutive blocks of recvcount elements that are scattered to the processes of the group. The i-th block is sent to process i and stored in the receive buffer defined by recvbuf, recvcount, and datatype.

> *Advice to implementors.* The MPI_REDUCE_SCATTER_BLOCK routine is functionally equivalent to: an MPI_REDUCE collective operation with count equal to recvcount*n, followed by an MPI_SCATTER with sendcount equal to recvcount. However, a direct implementation may run faster. (*End of advice to implementors.*)

The "in place" option for intracommunictors is specified by passing MPI_IN_PLACE in the sendbuf argument on *all* processes. In this case, the input data is taken from the receive buffer.

**Unofficial Draft for Comment Only**

If comm is an intercommunicator, then the result of the reduction of the data provided by processes in one group (group A) is scattered among processes in the other group (group B) and vice versa. Within each group, all processes provide the same value for the recvcount argument, and provide input vectors of count = n*recvcount elements stored in the send buffers, where n is the size of the group. The number of elements count must be the same for the two groups. The resulting vector from the other group is scattered in blocks of recvcount elements among the processes in the group.

> *Rationale.*    The last restriction is needed so that the length of the send buffer of one group can be determined by the local recvcount argument of the other group. Otherwise, a communication is needed to figure out how many elements are reduced. (*End of rationale.*)

## 5.10.2   MPI_REDUCE_SCATTER

MPI_REDUCE_SCATTER extends the functionality of MPI_REDUCE_SCATTER_BLOCK such that the scattered blocks can vary in size. Block sizes are determined by the recvcounts array, such that the i-th block contains recvcounts[i] elements.

MPI_REDUCE_SCATTER( sendbuf, recvbuf, recvcounts, datatype, op, comm)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) specifying the number of elements of the result distributed to each process. |
| IN | datatype | data type of elements of send and receive buffers (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |

```
int MPI_Reduce_scatter(void* sendbuf, void* recvbuf, int *recvcounts,
              MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

```
MPI_REDUCE_SCATTER(SENDBUF, RECVBUF, RECVCOUNTS, DATATYPE, OP, COMM,
              IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER RECVCOUNTS(*), DATATYPE, OP, COMM, IERROR
```

```
{void MPI::Comm::Reduce_scatter(const void* sendbuf, void* recvbuf,
              int recvcounts[], const MPI::Datatype& datatype,
              const MPI::Op& op) const = 0(binding deprecated, see Section 15.2) }
```

If comm is an intracommunicator, MPI_REDUCE_SCATTER first performs a global, element-wise reduction on vectors of count = $\sum_{i=0}^{n-1}$ recvcounts[i] elements in the send buffers defined by sendbuf, count and datatype, using the operation op, where n is the number of processes in the group of comm. The routine is called by all group members using the same arguments for recvcounts, datatype, op and comm. The resulting vector is treated as

n consecutive blocks where the number of elements of the i-th block is recvcounts[i]. The
blocks are scattered to the processes of the group. The i-th block is sent to process i and
stored in the receive buffer defined by recvbuf, recvcounts[i] and datatype.

> *Advice to implementors.* The MPI_REDUCE_SCATTER routine is functionally equiv-
> alent to: an MPI_REDUCE collective operation with count equal to the sum of
> recvcounts[i] followed by MPI_SCATTERV with sendcounts equal to recvcounts. How-
> ever, 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 receive buffer. It is
not required to specify the "in place" option on all processes, since the processes for which
recvcounts[i]==0 may not have allocated a receive buffer.

If comm is an intercommunicator, then the result of the reduction of the data provided
by processes in one group (group A) is scattered among processes in the other group (group
B), and vice versa. Within each group, all processes provide the same recvcounts argument,
and provide input vectors of count = $\sum_{i=0}^{n-1}$ recvcounts[i] elements stored in the send buffers,
where n is the size of the group. The resulting vector from the other group is scattered in
blocks of recvcounts[i] elements among the processes in the group. The number of elements
count must 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 recvcounts entries. Otherwise, a communication
> is needed to figure out how many elements are reduced. (*End of rationale.*)

## 5.11   Scan

### 5.11.1   Inclusive 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 (non-negative 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
```

**Unofficial Draft for Comment Only**

```
{void MPI::Intracomm::Scan(const void* sendbuf, void* recvbuf, int count,
            const MPI::Datatype& datatype, const MPI::Op& op) const(binding
            deprecated, see Section 15.2) }
```

If comm is an intracommunicator, 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,...,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 invalid for intercommunicators.

## 5.11.2   Exclusive Scan

MPI_EXSCAN(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 (non-negative integer) |
| IN | datatype | data type of elements of input buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | intracommunicator (handle) |

```
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(binding
            deprecated, see Section 15.2) }
```

If comm is an intracommunicator, 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 signficant 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 $0, \ldots, i - 1$ (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. The receive buffer on rank 0 is not changed by this operation. This operation is invalid for intercommunicators.

> *Rationale.* The exclusive scan is more general than the inclusive 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. (*End of rationale.*)

### 5.11.3 Example using MPI_SCAN

The example in this section uses an intracommunicator.

**Example 5.22**

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:

| values | $v_1$ | $v_2$ | $v_3$ | $v_4$ | $v_5$ | $v_6$ | $v_7$ | $v_8$ |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| logicals | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 |
| result | $v_1$ | $v_1 + v_2$ | $v_3$ | $v_3 + v_4$ | $v_3 + v_4 + v_5$ | $v_6$ | $v_6 + v_7$ | $v_8$ |

The operator that produces this effect is,

$$\left( \begin{array}{c} u \\ i \end{array} \right) \circ \left( \begin{array}{c} v \\ j \end{array} \right) = \left( \begin{array}{c} w \\ j \end{array} \right),$$

where,

$$w = \left\{ \begin{array}{ll} u + v & \text{if } i = j \\ v & \text{if } i \neq j \end{array} \right. .$$

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;
    SegScanPair  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_Get_address( a, disp);
    MPI_Get_address( a.log, disp+1);
    base = disp[0];
    for (i=0; i<2; ++i) disp[i] -= base;
    MPI_Type_create_struct( 2, blocklen, disp, type, &sspair );
    MPI_Type_commit( &sspair );
    /* create the segmented-scan user-op
     */
    MPI_Op_create(segScan, 0, &myOp);
    ...
    MPI_Scan( &a, &answer, 1, sspair, myOp, comm );
```

ticket109.

## 5.12   Nonblocking Collective Operations

As described in Section 3.7, performance of many applications can be improved by overlapping communication and computation, and many systems enable this. Nonblocking collective operations combine the potential benefits of nonblocking point-to-point operations, to exploit overlap and to avoid synchronization, with the optimized implementation and message scheduling provided by collective operations [27, 30]. One way of doing this would be to perform a blocking collective operation in a separate thread. An alternative mechanism that often leads to better performance (e.g., avoids context switching, scheduler overheads, and thread management) is to use nonblocking collective communication [28].

The nonblocking collective communication model is similar to the model used for nonblocking point-to-point communication. A nonblocking call initiates a collective operation, which must be completed in a separate completion call. Once initiated, the operation

may progress independently of any computation or other communication at participating processes. In this manner, nonblocking collective operations can mitigate possible synchronizing effects of collective operations by running them in the "background." In addition to enabling communication-computation overlap, nonblocking collective operations can perform collective operations on overlapping communicators, which would lead to deadlocks with blocking operations. Their semantic advantages can also be useful in combination with point-to-point communication.

As in the nonblocking point-to-point case, all calls are local and return immediately, irrespective of the status of other processes. The call initiates the operation, which indicates that the system may start to copy data out of the send buffer and into the receive buffer. Once initiated, all associated send buffers and buffers associated with input arguments (such as arrays of counts, displacements, or datatypes in the vector versions of the collectives) should not be modified, and all associated receive buffers should not be accessed, until the collective operation completes. The call returns a request handle, which must be passed to a completion call.

All completion calls (e.g., MPI_WAIT) described in Section 3.7.3 are supported for nonblocking collective operations. Similarly to the blocking case, nonblocking collective operations are considered to be complete when the local part of the operation is finished, i.e., for the caller, the semantics of the operation are guaranteed and all buffers can be safely accessed and modified. Completion does not indicate that other processes have completed or even started the operation (unless otherwise implied by the description of the operation). Completion of a particular nonblocking collective operation also does not indicate completion of any other posted nonblocking collective (or send-receive) operations, whether they are posted before or after the completed operation.

> *Advice to users.*    Users should be aware that implementations are allowed, but not required (with exception of MPI_IBARRIER), to synchronize processes during the completion of a nonblocking collective operation. (*End of advice to users.*)

Upon returning from a completion call in which a nonblocking collective operation completes, the MPI_ERROR field in the associated status object is set appropriately, see Section 3.2.5 on page 32. The values of the MPI_SOURCE and MPI_TAG fields are undefined. It is valid to mix different request types (i.e., any combination of collective requests, I/O requests, generalized requests, or point-to-point requests) in functions that enable multiple completions (e.g., MPI_WAITALL). It is erroneous to call MPI_REQUEST_FREE or MPI_CANCEL for a request associated with a nonblocking collective operation. Nonblocking collective requests are not persistent.

> *Rationale.*    Freeing an active nonblocking collective request could cause similar problems as discussed for point-to-point requests (see Section 3.7.3). Cancelling a request is not supported because the semantics of this operation are not well-defined. (*End of rationale.*)

Multiple nonblocking collective operations can be outstanding on a single communicator. If the nonblocking call causes some system resource to be exhausted, then it may fail and generate an MPI exception. 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.

Unlike point-to-point operations, nonblocking collective operations do not match with blocking collective operations, and collective operations do not have a tag argument. All processes must call collective operations (blocking and nonblocking) in the same order per communicator. In particular, once a process calls a collective operation, all other processes in the communicator must eventually call the same collective operation, and no other collective operation with the same communicator in between. This is consistent with the ordering rules for blocking collective operations in threaded environments.

> *Rationale.*   Matching blocking and nonblocking collective operations is not allowed because the implementation might use different communication algorithms for the two cases. Blocking collective operations may be optimized for minimal time to completion, while nonblocking collective operations may balance time to completion with CPU overhead and asynchronous progression.
>
> The use of tags for collective operations can prevent certain hardware optimizations. (*End of rationale.*)

> *Advice to users.*   If program semantics require matching blocking and nonblocking collective operations, then a nonblocking collective operation can be initiated and immediately completed with a blocking wait to emulate blocking behavior. (*End of advice to users.*)

In terms of data movements, each nonblocking collective operation has the same effect as its blocking counterpart for intracommunicators and intercommunicators after completion. Likewise, upon completion, nonblocking collective reduction operations have the same effect as their blocking counterparts, and the same restrictions and recommendations on reduction orders apply.

The use of the "in place" option is allowed exactly as described for the corresponding blocking collective operations. When using the "in place" option, message buffers function as both send and receive buffers. Such buffers should not be modified or accessed until the operation completes.

Progression rules for nonblocking collective operations are similar to progression of nonblocking point-to-point operations, refer to Section 3.7.4.

> *Advice to implementors.*   Nonblocking collective operations can be implemented with local execution schedules [29] using nonblocking point-to-point communication and a reserved tag-space. (*End of advice to implementors.*)

### 5.12.1   Nonblocking Barrier Synchronization

MPI_IBARRIER(comm , request)

| | | |
|---|---|---|
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ibarrier(MPI_Comm comm, MPI_Request *request)
```

```
MPI_IBARRIER(COMM, REQUEST, IERROR)
```

```
    INTEGER COMM, REQUEST, IERROR
```

{MPI::Request MPI::Comm::Ibarrier() const = 0*(binding deprecated, see
            Section 15.2)* }

MPI_IBARRIER is a nonblocking version of MPI_BARRIER. By calling MPI_IBARRIER,
a process notifies that it has reached the barrier. The call returns immediately, independent of whether other processes have called MPI_IBARRIER. The usual barrier semantics
are enforced at the corresponding completion operation (test or wait), which in the intra-
communicator case will complete only after all other processes in the communicator have
called MPI_IBARRIER. In the intercommunicator case, it will complete when all processes
in the remote group have called MPI_IBARRIER.

> *Advice to users.* A nonblocking barrier can be used to hide latency. Moving independent computations between the MPI_IBARRIER and the subsequent completion call
> can overlap the barrier latency and therefore shorten possible waiting times. The se-
> mantic properties are also useful when mixing collective operations and point-to-point
> messages. (*End of advice to users.*)

### 5.12.2 Nonblocking Broadcast

MPI_IBCAST(buffer, count, datatype, root, comm, request)

| | | |
|---|---|---|
| INOUT | buffer | starting address of buffer (choice) |
| IN | count | number of entries in buffer (non-negative integer) |
| IN | datatype | data type of buffer (handle) |
| IN | root | rank of broadcast root (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ibcast(void* buffer, int count, MPI_Datatype datatype, int root,
            MPI_Comm comm, MPI_Request *request)
```

```
MPI_IBCAST(BUFFER, COUNT, DATATYPE, ROOT, COMM, REQUEST, IERROR)
    <type> BUFFER(*)
    INTEGER COUNT, DATATYPE, ROOT, COMM, REQUEST, IERROR
```

{MPI::Request MPI::Comm::Ibcast(void* buffer, int count,
            const MPI::Datatype& datatype, int root) const = 0*(binding
            deprecated, see Section 15.2)* }

This call starts a nonblocking variant of MPI_BCAST (see Section 5.4).

### Example using MPI_IBCAST

The example in this section uses an intracommunicator.

**Example 5.23**

Start a broadcast of 100 `int`s from process 0 to every process in the group, perform some computation on independent data, and then complete the outstanding broadcast operation.

```
    MPI_Comm comm;
    int array1[100], array2[100];
    int root=0;
    MPI_Request req;
    ...
    MPI_Ibcast(array1, 100, MPI_INT, root, comm, &req);
    compute(array2, 100);
    MPI_Wait(&req, MPI_STATUS_IGNORE);
```

### 5.12.3   Nonblocking Gather

MPI_IGATHER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | recvcount | number of elements for any single receive (non-negative integer, significant only at root) |
| IN | recvtype | data type of recv buffer elements (significant only at root) (handle) |
| IN | root | rank of receiving process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Igather(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int recvcount, MPI_Datatype recvtype, int root,
            MPI_Comm comm, MPI_Request *request)
```

```
MPI_IGATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
            ROOT, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, ROOT, COMM, REQUEST,
    IERROR
```

```
{MPI::Request MPI::Comm::Igather(const void* sendbuf, int sendcount, const
            MPI::Datatype& sendtype, void* recvbuf, int recvcount,
            const MPI::Datatype& recvtype, int root) const = 0(binding
            deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_GATHER (see Section 5.5).

MPI_IGATHERV(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtype, root, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process (significant only at root) |
| IN | displs | integer array (of length group size). Entry i specifies the displacement relative to recvbuf at which to place the incoming data from process i (significant only at root) |
| IN | recvtype | data type of recv buffer elements (significant only at root) (handle) |
| IN | root | rank of receiving process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Igatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int *recvcounts, int *displs,
            MPI_Datatype recvtype, int root, MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_IGATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
            RECVTYPE, ROOT, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPE, ROOT,
    COMM, REQUEST, IERROR
```

```
{MPI::Request MPI::Comm::Igatherv(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(binding
            deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_GATHERV (see Section 5.5).

ticket265.

MPI_IGATHERW(sendbuf, sendcount, sendtype, recvbuf, displs, recvcounts, recvtype, root,
             comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process (signifiant only at root) |
| IN | displs | integer array (of length group size). Entry i specifies the displacement relative to recvbuf from which to take the incoming data from process i (significant only at root) |
| IN | recvtypes | array of datatypes (of length group size). Entry i specifies the type of data received from process i (significant only at root) (array of handles) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Igatherw(void* sendbuf, int sendcount, MPI_Datatype sendtype,
             void* recvbuf, int *recvcounts[], int displs[],
             MPI_Datatype *recvtypes[], int root, MPI_Comm comm,
             MPI_Request *request)
```

```
MPI_IGATHERW(SENDBUF, SENDCOUNT, SENDTYPES, RECVBUF, RECVCOUNTS, DISPLS,
             RECVTYPES, ROOT, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPES(*),
    ROOT, COMM, REQUEST, IERROR
```

ticket109.      This call starts a nonblocking variant of MPI_GATHERW (see Section 5.5).

### 5.12.4 Nonblocking Scatter

MPI_ISCATTER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | address of send buffer (choice, significant only at root) |
| IN | sendcount | number of elements sent to each process (non-negative integer, significant only at root) |
| IN | sendtype | data type of send buffer elements (significant only at root) (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iscatter(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int recvcount, MPI_Datatype recvtype, int root,
            MPI_Comm comm, MPI_Request *request)
```

```
MPI_ISCATTER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
            ROOT, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, ROOT, COMM, REQUEST,
    IERROR
```

```
{MPI::Request MPI::Comm::Iscatter(const void* sendbuf, int sendcount, const
            MPI::Datatype& sendtype, void* recvbuf, int recvcount,
            const MPI::Datatype& recvtype, int root) const = 0(binding
            deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_SCATTER (see Section 5.6).

MPI_ISCATTERV(sendbuf, sendcounts, displs, sendtype, recvbuf, recvcount, recvtype, root, comm, request)

| IN | sendbuf | address of send buffer (choice, significant only at root) |
|----|---------|-----------------------------------------------------------|
| IN | sendcounts | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN | displs | integer array (of length group size). Entry i specifies the displacement (relative to sendbuf) from which to take the outgoing data to process i |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iscatterv(void* sendbuf, int *sendcounts, int *displs,
            MPI_Datatype sendtype, void* recvbuf, int recvcount,
            MPI_Datatype recvtype, int root, MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_ISCATTERV(SENDBUF, SENDCOUNTS, DISPLS, SENDTYPE, RECVBUF, RECVCOUNT,
            RECVTYPE, ROOT, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNTS(*), DISPLS(*), SENDTYPE, RECVCOUNT, RECVTYPE, ROOT,
    COMM, REQUEST, IERROR
```

```
{MPI::Request MPI::Comm::Iscatterv(const void* sendbuf,
            const int sendcounts[], const int displs[],
            const MPI::Datatype& sendtype, void* recvbuf, int recvcount,
            const MPI::Datatype& recvtype, int root) const = 0
```
*(binding deprecated, see Section 15.2)* }

ticket265.     This call starts a nonblocking variant of MPI_SCATTERV (see Section 5.6).

MPI_ISCATTERW(sendbuf, sendcount, displs, sendtypes, recvbuf, recvcount, recvtype, root, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice, significant only at root) |
| IN | sendcount | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN | displs | integer array (of length group size). Entry i specifies the displacement relative to sendbuf from which to take the outgoing data to process i |
| IN | sendtypes | array of datatypes (of length group size). Entry t j specifies the type of data to send to process t j (array of handles) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements in receive buffer (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | root | rank of sending process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iscatterw(void* sendbuf, int sendcounts[], int displs[],
            MPI_Datatype sendtypes[], void* recvbuf, int *recvcount,
            MPI_Datatype *recvtype, int root, MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_ISCATTERW(SENDBUF, SENDCOUNTS, DISPLS, SENDTYPES, RECVBUF, RECVCOUNT,
            RECVTYPE, ROOT, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER DISPLS(*), SENDCOUNTS(*), SENDTYPES, RECVCOUNT, RECVTYPE, ROOT,
    COMM, REQUEST, IERROR
```

This call starts a nonblocking variant of MPI_SCATTERW (see Section 5.6).     ticket109.

## 5.12.5   Nonblocking Gather-to-all

MPI_IALLGATHER(sendbuf,  sendcount,  sendtype,  recvbuf,  recvcount,  recvtype,  comm,
              request)

| IN  | sendbuf   | starting address of send buffer (choice) |
|-----|-----------|------------------------------------------|
| IN  | sendcount | number of elements in send buffer (non-negative integer) |
| IN  | sendtype  | data type of send buffer elements (handle) |
| OUT | recvbuf   | address of receive buffer (choice) |
| IN  | recvcount | number of elements received from any process (non-negative integer) |
| IN  | recvtype  | data type of receive buffer elements (handle) |
| IN  | comm      | communicator (handle) |
| OUT | request   | communication request (handle) |

```
int MPI_Iallgather(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int recvcount, MPI_Datatype recvtype,
            MPI_Comm comm, MPI_Request *request)
```

```
MPI_IALLGATHER(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
            COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, COMM, REQUEST, IERROR
```

```
{MPI::Request MPI::Comm::Iallgather(const void* sendbuf, int sendcount,
            const MPI::Datatype& sendtype, void* recvbuf, int recvcount,
            const MPI::Datatype& recvtype) const = 0(binding deprecated, see
            Section 15.2) }
```

This call starts a nonblocking variant of MPI_ALLGATHER (see Section 5.7).

MPI_IALLGATHERV(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtype, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process |
| IN | displs | integer array (of length group size). Entry i specifies the displacement (relative to recvbuf) at which to place the incoming data from process i |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iallgatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int *recvcounts, int *displs,
            MPI_Datatype recvtype, MPI_Comm comm, MPI_Request* request)
```

```
MPI_IALLGATHERV(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
            RECVTYPE, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPE, COMM,
    REQUEST, IERROR
```

{MPI::Request MPI::Comm::Iallgatherv(const void* sendbuf, int sendcount,
            const MPI::Datatype& sendtype, void* recvbuf,
            const int recvcounts[], const int displs[],
            const MPI::Datatype& recvtype) const = 0*(binding deprecated, see Section 15.2)* }

This call starts a nonblocking variant of MPI_ALLGATHERV (see Section 5.7). ticket265.

MPI_IALLGATHERW(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtypes, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements in send buffer (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) containing the number of elements that are received from each process |
| IN | displs | integer array (of length group size). Entry i specifies the displacement (relative to recvbuf) at which to place the incoming data from process i |
| IN | recvtypes | array of datatypes (of length group size). Entry i specifies the type of data received from process i (array of handles) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iallgatherw(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int *recvcounts, int *displs,
            MPI_Datatype *recvtypes, MPI_Comm comm, MPI_Request *request)
```

```
MPI_IALLGATHERW(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNTS, DISPLS,
            RECVTYPES, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNTS(*), DISPLS(*), RECVTYPES(*),
    COMM, REQUEST, IERROR
```

ticket109.    This call starts a nonblocking variant of MPI_ALLGATHERW (see Section 5.7).

### 5.12.6   Nonblocking All-to-All Scatter/Gather

**MPI_IALLTOALL(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, comm, request)**

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcount | number of elements sent to each process (non-negative integer) |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcount | number of elements received from any process (non-negative integer) |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ialltoall(void* sendbuf, int sendcount, MPI_Datatype sendtype,
            void* recvbuf, int recvcount, MPI_Datatype recvtype,
            MPI_Comm comm, MPI_Request *request)
```

```
MPI_IALLTOALL(SENDBUF, SENDCOUNT, SENDTYPE, RECVBUF, RECVCOUNT, RECVTYPE,
            COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT, RECVTYPE, COMM, REQUEST, IERROR
```

```
{MPI::Request MPI::Comm::Ialltoall(const void* sendbuf, int sendcount,
            const MPI::Datatype& sendtype, void* recvbuf, int recvcount,
            const MPI::Datatype& recvtype) const = 0
```
*(binding deprecated, see Section 15.2)* }

This call starts a nonblocking variant of MPI_ALLTOALL (see Section 5.8).

MPI_IALLTOALLV(sendbuf,  sendcounts,  sdispls,  sendtype,  recvbuf,  recvcounts,  rdispls,
                recvtype, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcounts | non-negative integer array (of length group size) specifying the number of elements to send to each processor |
| IN | sdispls | integer array (of length group size). Entry j specifies the displacement (relative to sendbuf) from which to take the outgoing data destined for process j |
| IN | sendtype | data type of send buffer elements (handle) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array (of length group size) specifying the number of elements that can be received from each processor |
| IN | rdispls | integer array (of length group size). Entry i specifies the displacement (relative to recvbuf) at which to place the incoming data from process i |
| IN | recvtype | data type of receive buffer elements (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ialltoallv(void* sendbuf, int *sendcounts, int *sdispls,
              MPI_Datatype sendtype, void* recvbuf, int *recvcounts,
              int *rdispls, MPI_Datatype recvtype, MPI_Comm comm,
              MPI_Request *request)

MPI_IALLTOALLV(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPE, RECVBUF, RECVCOUNTS,
              RDISPLS, RECVTYPE, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPE, RECVCOUNTS(*), RDISPLS(*),
    RECVTYPE, COMM, REQUEST, IERROR

{MPI::Request MPI::Comm::Ialltoallv(const void* sendbuf,
              const int sendcounts[], const int sdispls[],
              const MPI::Datatype& sendtype, void* recvbuf,
              const int recvcounts[], const int rdispls[],
              const MPI::Datatype& recvtype) const = 0(binding deprecated, see
              Section 15.2) }
```

This call starts a nonblocking variant of MPI_ALLTOALLV (see Section 5.8).

MPI_IALLTOALLW(sendbuf, sendcounts, sdispls, sendtypes, recvbuf, recvcounts, rdispls, recvtypes, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| IN | sendcounts | integer array (of length group size) specifying the number of elements to send to each processor (array of non-negative integers) |
| IN | sdispls | integer array (of length group size). Entry j specifies the displacement in bytes (relative to sendbuf) from which to take the outgoing data destined for process j (array of integers) |
| IN | sendtypes | array of datatypes (of length group size). Entry j specifies the type of data to send to process j (array of handles) |
| OUT | recvbuf | address of receive buffer (choice) |
| IN | recvcounts | integer array (of length group size) specifying the number of elements that can be received from each processor (array of non-negative integers) |
| IN | rdispls | integer array (of length group size). Entry i specifies the displacement in bytes (relative to recvbuf) at which to place the incoming data from process i (array of integers) |
| IN | recvtypes | array of datatypes (of length group size). Entry i specifies the type of data received from process i (array of handles) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ialltoallw(void* sendbuf, int sendcounts[], int sdispls[],
            MPI_Datatype sendtypes[], void* recvbuf, int recvcounts[],
            int rdispls[], MPI_Datatype recvtypes[], MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_IALLTOALLW(SENDBUF, SENDCOUNTS, SDISPLS, SENDTYPES, RECVBUF,
            RECVCOUNTS, RDISPLS, RECVTYPES, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER SENDCOUNTS(*), SDISPLS(*), SENDTYPES(*), RECVCOUNTS(*),
    RDISPLS(*), RECVTYPES(*), COMM, REQUEST, IERROR
```

{MPI::Request MPI::Comm::Ialltoallw(const void* sendbuf, const int sendcounts[], const int sdispls[], const MPI::Datatype sendtypes[], void* recvbuf, const int recvcounts[], const int rdispls[], const MPI::Datatype recvtypes[]) const = 0 *(binding deprecated, see Section 15.2)* }

This call starts a nonblocking variant of MPI_ALLTOALLW (see Section 5.8).

## 5.12.7   Nonblocking Reduce

MPI_IREDUCE(sendbuf, recvbuf, count, datatype, op, root, comm, request)

| IN | sendbuf | address of send buffer (choice) |
|---|---|---|
| OUT | recvbuf | address of receive buffer (choice, significant only at root) |
| IN | count | number of elements in send buffer (non-negative integer) |
| IN | datatype | data type of elements of send buffer (handle) |
| IN | op | reduce operation (handle) |
| IN | root | rank of root process (integer) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ireduce(void* sendbuf, void* recvbuf, int count,
              MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm,
              MPI_Request *request)
```

```
MPI_IREDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, ROOT, COMM, REQUEST,
              IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER COUNT, DATATYPE, OP, ROOT, COMM, REQUEST, IERROR
```

{MPI::Request MPI::Comm::Ireduce(const void* sendbuf, void* recvbuf,
              int count, const MPI::Datatype& datatype, const MPI::Op& op,
              int root) const = 0*(binding deprecated, see Section 15.2)* }

This call starts a nonblocking variant of MPI_REDUCE (see Section 5.9.1).

*Advice to implementors.* The implementation is explicitly allowed to use different algorithms for blocking and nonblocking reduction operations that might change the order of evaluation of the operations. However, as for MPI_REDUCE, it is strongly recommended that MPI_IREDUCE 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 processes. (*End of advice to implementors.*)

*Advice to users.* For operations which are not truly associative, the result delivered upon completion of the nonblocking reduction may not exactly equal the result delivered by the blocking reduction, even when specifying the same arguments in the same order. (*End of advice to users.*)

### 5.12.8   Nonblocking All-Reduce

MPI_IALLREDUCE(sendbuf, recvbuf, count, datatype, op, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | count | number of elements in send buffer (non-negative integer) |
| IN | datatype | data type of elements of send buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iallreduce(void* sendbuf, void* recvbuf, int count,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_IALLREDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, REQUEST,
            IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER COUNT, DATATYPE, OP, COMM, REQUEST, IERROR
```

```
{MPI::Request MPI::Comm::Iallreduce(const void* sendbuf, void* recvbuf,
            int count, const MPI::Datatype& datatype, const MPI::Op& op)
            const = 0(binding deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_ALLREDUCE (see Section 5.9.6).

### 5.12.9   Nonblocking Reduce-Scatter with Equal Blocks

MPI_IREDUCE_SCATTER_BLOCK(sendbuf, recvbuf, recvcount, datatype, op, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | recvcount | element count per block (non-negative integer) |
| IN | datatype | data type of elements of send and receive buffers (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ireduce_scatter_block(void* sendbuf, void* recvbuf, int recvcount,
             MPI_Datatype datatype, MPI_Op op, MPI_Comm comm,
             MPI_Request *request)

MPI_IREDUCE_SCATTER_BLOCK(SENDBUF, RECVBUF, RECVCOUNT, DATATYPE, OP, COMM,
             REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER RECVCOUNT, DATATYPE, OP, COMM, REQUEST, IERROR

{MPI::Request MPI::Comm::Ireduce_scatter_block(const void* sendbuf,
             void* recvbuf, int recvcount, const MPI::Datatype& datatype,
             const MPI::Op& op) const = 0(binding deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_REDUCE_SCATTER_BLOCK (see Section 5.10.1).

### 5.12.10   Nonblocking Reduce-Scatter

MPI_IREDUCE_SCATTER(sendbuf, recvbuf, recvcounts, datatype, op, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | recvcounts | non-negative integer array specifying the number of elements in result distributed to each process. Array must be identical on all calling processes. |
| IN | datatype | data type of elements of input buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Ireduce_scatter(void* sendbuf, void* recvbuf, int *recvcounts,
             MPI_Datatype datatype, MPI_Op op, MPI_Comm comm,
             MPI_Request *request)

MPI_IREDUCE_SCATTER(SENDBUF, RECVBUF, RECVCOUNTS, DATATYPE, OP, COMM,
             REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER RECVCOUNTS(*), DATATYPE, OP, COMM, REQUEST, IERROR

{MPI::Request MPI::Comm::Ireduce_scatter(const void* sendbuf,
             void* recvbuf, int recvcounts[],
             const MPI::Datatype& datatype, const MPI::Op& op)
             const = 0(binding deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_REDUCE_SCATTER (see Section 5.10.2).

### 5.12.11 Nonblocking Inclusive Scan

MPI_ISCAN(sendbuf, recvbuf, count, datatype, op, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | count | number of elements in input buffer (non-negative integer) |
| IN | datatype | data type of elements of input buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | communicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iscan(void* sendbuf, void* recvbuf, int count,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_ISCAN(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER COUNT, DATATYPE, OP, COMM, REQUEST, IERROR
```

```
{MPI::Request MPI::Intracomm::Iscan(const void* sendbuf, void* recvbuf,
            int count, const MPI::Datatype& datatype, const MPI::Op& op)
            const(binding deprecated, see Section 15.2) }
```

This call starts a nonblocking variant of MPI_SCAN (see Section 5.11).

### 5.12.12 Nonblocking Exclusive Scan

MPI_IEXSCAN(sendbuf, recvbuf, count, datatype, op, comm, request)

| | | |
|---|---|---|
| IN | sendbuf | starting address of send buffer (choice) |
| OUT | recvbuf | starting address of receive buffer (choice) |
| IN | count | number of elements in input buffer (non-negative integer) |
| IN | datatype | data type of elements of input buffer (handle) |
| IN | op | operation (handle) |
| IN | comm | intracommunicator (handle) |
| OUT | request | communication request (handle) |

```
int MPI_Iexscan(void* sendbuf, void* recvbuf, int count,
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm,
            MPI_Request *request)
```

```
MPI_IEXSCAN(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, REQUEST, IERROR)
    <type> SENDBUF(*), RECVBUF(*)
    INTEGER COUNT, DATATYPE, OP, COMM, REQUEST, IERROR
```

{MPI::Request MPI::Intracomm::Iexscan(const void* sendbuf, void* recvbuf,
              int count, const MPI::Datatype& datatype, const MPI::Op& op)
              const*(binding deprecated, see Section 15.2) }*

This call starts a nonblocking variant of MPI_EXSCAN (see Section 5.11.2).

## 5.13   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 5.24**
    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 com-
munication group.

**Example 5.25**
    The following is erroneous.

```
switch(rank) {
    case 0:
        MPI_Bcast(buf1, count, type, 0, comm0);
        MPI_Bcast(buf2, count, type, 2, comm2);
        break;
    case 1:
        MPI_Bcast(buf1, count, type, 1, comm1);
        MPI_Bcast(buf2, count, type, 0, comm0);
        break;
    case 2:
        MPI_Bcast(buf1, count, type, 2, comm2);
```

```
        MPI_Bcast(buf2, count, type, 1, comm1);
        break;
}
```

Assume that the group of comm0 is {0,1}, of comm1 is {1, 2} and of comm2 is {2,0}. If the broadcast is a synchronizing operation, then there is a cyclic dependency: the broadcast in comm2 completes only after the broadcast in comm0; the broadcast in comm0 completes only after the broadcast in comm1; and the broadcast in comm1 completes only after the broadcast in comm2. Thus, the code will deadlock.

Collective operations must be executed in an order so that no cyclic dependencies occur. Nonblocking collective operations can alleviate this issue.

**Example 5.26**
The following is erroneous.

```
switch(rank) {
    case 0:
        MPI_Bcast(buf1, count, type, 0, comm);
        MPI_Send(buf2, count, type, 1, tag, comm);
        break;
    case 1:
        MPI_Recv(buf2, count, type, 0, tag, comm, status);
        MPI_Bcast(buf1, count, type, 0, comm);
        break;
}
```

Process zero executes a broadcast, followed by a blocking send operation. Process one first executes a blocking receive that matches the send, followed by broadcast call that matches the broadcast of process zero. This program may deadlock. The broadcast call on process zero *may* block until process one executes the matching broadcast call, so that the send is not executed. Process one will definitely block on the receive and so, in this case, never executes the broadcast.

The relative order of execution of collective operations and point-to-point operations should be such, so that even if the collective operations and the point-to-point operations are synchronizing, no deadlock will occur.

**Example 5.27**
An unsafe, non-deterministic program.

```
switch(rank) {
    case 0:
        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;
```

*First Execution*

```
        0                    1                    2

                      recv  ◄──── match ──── send

   broadcast        broadcast            broadcast

   send ──── match ────►  recv
```

*Second Execution*

```
   broadcast
   send ──── match ────►  recv

                      broadcast

                      recv  ◄──── match ──── send

                                           broadcast
```
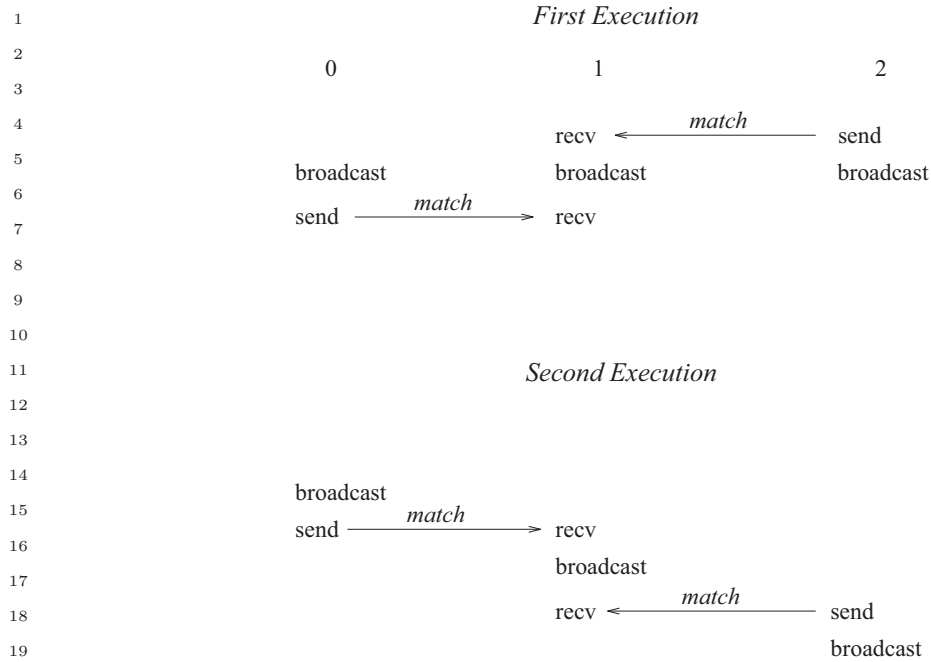
Figure 5.12: A race condition causes non-deterministic matching of sends and receives. One cannot rely on synchronization from a broadcast to make the program deterministic.

```
    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 5.12. 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 responsibility to ensure that the same communicator is not used concurrently by two different collective communication calls at the same process.

> *Advice to implementors.*   Assume that broadcast is implemented using point-to-point MPI communication. Suppose the following two rules are followed.
>
>   1. All receives specify their source explicitly (no wildcards).
>
>   2. Each process sends all messages that pertain to one collective call before sending any message that pertain to a subsequent collective call.

**Unofficial Draft for Comment Only**

Then, messages belonging to successive broadcasts cannot be confused, as the order of point-to-point messages is preserved.

It is the implementor's responsibility to ensure that point-to-point messages are not confused with collective messages. One way to accomplish this is, whenever a communicator is created, to also create a "hidden communicator" for collective communication. One could achieve a similar effect more cheaply, for example, by using a hidden tag or context bit to indicate whether the communicator is used for point-to-point or collective communication. (*End of advice to implementors.*)

**Example 5.28**
Blocking and nonblocking collective operations can be interleaved, i.e., a blocking collective operation can be posted even if there is a nonblocking collective operation outstanding.

```
MPI_Request req;

MPI_Ibarrier(comm, &req);
MPI_Bcast(buf1, count, type, 0, comm);
MPI_Wait(&req, MPI_STATUS_IGNORE);
```

Each process starts a nonblocking barrier operation, participates in a blocking broadcast and then waits until every other process started the barrier operation. This effectively turns the broadcast into a synchronizing broadcast with possible communication/communication overlap (MPI_Bcast is allowed, but not required to synchronize).

**Example 5.29**
The starting order of collective operations on a particular communicator defines their matching. The following example shows an erroneous matching of different collective operations on the same communicator.

```
MPI_Request req;
switch(rank) {
    case 0:
        /* erroneous matching */
        MPI_Ibarrier(comm, &req);
        MPI_Bcast(buf1, count, type, 0, comm);
        MPI_Wait(&req, MPI_STATUS_IGNORE);
        break;
    case 1:
        /* erroneous matching */
        MPI_Bcast(buf1, count, type, 0, comm);
        MPI_Ibarrier(comm, &req);
        MPI_Wait(&req, MPI_STATUS_IGNORE);
        break;
}
```

This ordering would match MPI_Ibarrier on rank 0 with MPI_Bcast on rank 1 which is erroneous and the program behavior is undefined. However, if such an order is required, the user must create different duplicate communicators and perform the operations on them. If started with two processes, the following program would be correct:

ticket109.

**Unofficial Draft for Comment Only**

```
MPI_Request req;
MPI_Comm dupcomm;
MPI_Comm_dup(comm, &dupcomm);
switch(rank) {
    case 0:
        MPI_Ibarrier(comm, &req);
        MPI_Bcast(buf1, count, type, 0, dupcomm);
        MPI_Wait(&req, MPI_STATUS_IGNORE);
        break;
    case 1:
        MPI_Bcast(buf1, count, type, 0, dupcomm);
        MPI_Ibarrier(comm, &req);
        MPI_Wait(&req, MPI_STATUS_IGNORE);
        break;
}
```

> *Advice to users.*   The use of different communicators offers some flexibility regarding the matching of nonblocking collective operations. In this sense, communicators could be used as an equivalent to tags. However, communicator construction might induce overheads so that this should be used carefully. (*End of advice to users.*)

**Example 5.30**
Nonblocking collective operations can rely on the same progression rules as nonblocking point-to-point messages. Thus, if started with two processes, the following program is a valid MPI program and is guaranteed to terminate:

```
MPI_Request req;

switch(rank) {
    case 0:
      MPI_Ibarrier(comm, &req);
      MPI_Wait(&req, MPI_STATUS_IGNORE);
      MPI_Send(buf, count, dtype, 1, tag, comm);
      break;
    case 1:
      MPI_Ibarrier(comm, &req);
      MPI_Recv(buf, count, dtype, 0, tag, comm, MPI_STATUS_IGNORE);
      MPI_Wait(&req, MPI_STATUS_IGNORE);
      break;
}
```

The MPI library must progress the barrier in the MPI_Recv call. Thus, the MPI_Wait call in rank 0 will eventually complete, which enables the matching MPI_Send so all calls eventually return.

**Example 5.31**
Blocking and nonblocking collective operations do not match. The following example is erroneous.

```
MPI_Request req;

switch(rank) {
    case 0:
      /* erroneous false matching of Alltoall and Ialltoall */
      MPI_Ialltoall(sbuf, scnt, stype, rbuf, rcnt, rtype, comm, &req);
      MPI_Wait(&req, MPI_STATUS_IGNORE);
      break;
    case 1:
      /* erroneous false matching of Alltoall and Ialltoall */
      MPI_Alltoall(sbuf, scnt, stype, rbuf, rcnt, rtype, comm);
      break;
}
```

**Example 5.32**
    Collective and point-to-point requests can be mixed in functions that enable multiple completions. If started with two processes, the following program is valid.

```
MPI_Request reqs[2];

switch(rank) {
    case 0:
      MPI_Ibarrier(comm, &reqs[0]);
      MPI_Send(buf, count, dtype, 1, tag, comm);
      MPI_Wait(&reqs[0], MPI_STATUS_IGNORE);
      break;
    case 1:
      MPI_Irecv(buf, count, dtype, 0, tag, comm, &reqs[0]);
      MPI_Ibarrier(comm, &reqs[1]);
      MPI_Waitall(2, reqs, MPI_STATUSES_IGNORE);
      break;
}
```

    The Waitall call returns only after the barrier and the receive completed.

**Example 5.33**
    Multiple nonblocking collective operations can be outstanding on a single communicator and match in order.

```
MPI_Request reqs[3];

compute(buf1);
MPI_Ibcast(buf1, count, type, 0, comm, &reqs[0]);
compute(buf2);
MPI_Ibcast(buf2, count, type, 0, comm, &reqs[1]);
compute(buf3);
MPI_Ibcast(buf3, count, type, 0, comm, &reqs[2]);
MPI_Waitall(3, reqs, MPI_STATUSES_IGNORE);
```

*Advice to users.*   Pipelining and double-buffering techniques can efficiently be used to overlap computation and communication. However, having too many outstanding requests might have a negative impact on performance. (*End of advice to users.*)

*Advice to implementors.*   The use of pipelining may generate many outstanding requests. A high-quality hardware-supported implementation with limited resources should be able to fall back to a software implementation if its resources are exhausted. In this way, the implementation could limit the number of outstanding requests only by the available memory. (*End of advice to implementors.*)

**Example 5.34**
    Nonblocking collective operations can also be used to enable simultaneous collective operations on multiple overlapping communicators (see Figure 5.13). The following example is started with three processes and three communicators. The first communicator `comm1` includes ranks 0 and 1, `comm2` includes ranks 1 and 2 and `comm3` spans ranks 0 and 2. It is not possible to perform a blocking collective operation on all communicators because there exists no deadlock-free order to invoke them. However, nonblocking collective operations can easily be used to achieve this task.

```
MPI_Request reqs[2];

switch(rank) {
    case 0:
      MPI_Iallreduce(sbuf1, rbuf1, count, dtype, MPI_SUM, comm1, &reqs[0]);
      MPI_Iallreduce(sbuf3, rbuf3, count, dtype, MPI_SUM, comm3, &reqs[1]);
      break;
    case 1:
      MPI_Iallreduce(sbuf1, rbuf1, count, dtype, MPI_SUM, comm1, &reqs[0]);
      MPI_Iallreduce(sbuf2, rbuf2, count, dtype, MPI_SUM, comm2, &reqs[1]);
      break;
    case 2:
      MPI_Iallreduce(sbuf2, rbuf2, count, dtype, MPI_SUM, comm2, &reqs[0]);
      MPI_Iallreduce(sbuf3, rbuf3, count, dtype, MPI_SUM, comm3, &reqs[1]);
      break;
}
MPI_Waitall(2, reqs, MPI_STATUSES_IGNORE);
```

*Advice to users.*   This method can be useful if overlapping neighboring regions (halo or ghost zones) are used in collective operations. The sequence of the two calls in each process is irrelevant because the two nonblocking operations are performed on different communicators. (*End of advice to users.*)

**Example 5.35**
    The progress of multiple outstanding nonblocking collective operations is completely independent.
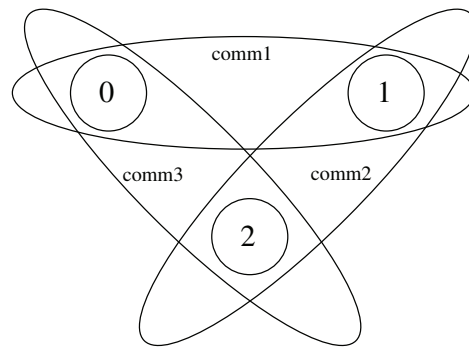
Figure 5.13: Example with overlapping communicators.

```
MPI_Request reqs[2];

compute(buf1);
MPI_Ibcast(buf1, count, type, 0, comm, &reqs[0]);
compute(buf2);
MPI_Ibcast(buf2, count, type, 0, comm, &reqs[1]);
MPI_Wait(&reqs[1], MPI_STATUS_IGNORE);
/* nothing is known about the status of the first bcast here */
MPI_Wait(&reqs[0], MPI_STATUS_IGNORE);
```

Finishing the second MPI_IBCAST is completely independent of the first one. This means that it is not guaranteed that the first broadcast operation is finished or even started after the second one is completed via `reqs[1]`.

# Chapter 12

# External Interfaces

## 12.1 Introduction

This chapter begins with calls used to create **generalized requests**, which allow users 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 12.3 deals with setting the information found in status. This is needed for generalized requests.

The chapter continues, in Section 12.4, with a discussion of how threads are to be handled in MPI. Although thread compliance is not required, the standard specifies how threads are to work if they are provided.

## 12.2 Generalized Requests

The goal of generalized requests 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) (binding deprecated, see Section 15.2) }
```

> *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); (binding deprecated, see Section 15.2)}
```

query_fn function computes the status that should be returned for the generalized request. The status also includes information about successful/unsuccessful 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); (binding
            deprecated, see Section 15.2)}
```

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 call to query_fn for the same request. However, if the MPI call completed multiple generalized requests, the order in which free_fn callback functions are invoked is not specified by MPI.

free_fn callback is also invoked for generalized requests that are freed by a call to MPI_REQUEST_FREE (no call to WAIT_{WAIT|TEST}{ANY|SOME|ALL} will occur for such a request). In this case, the callback function will be called either in the MPI call MPI_REQUEST_FREE(request), or in the MPI call MPI_GREQUEST_COMPLETE(request),

whichever happens last, i.e., in this case the actual freeing code is executed as soon as both calls MPI_REQUEST_FREE and MPI_GREQUEST_COMPLETE have occurred. The request is not deallocated until after free_fn completes. Note that free_fn will be invoked only once per request by a correct program.

> *Advice to users.*  Calling MPI_REQUEST_FREE(request) will cause the request handle to be set to MPI_REQUEST_NULL. This handle to the generalized request is no longer valid. However, user copies of this handle are valid until after free_fn completes since MPI does not deallocate the object until then. Since free_fn is not called until after MPI_GREQUEST_COMPLETE, the user copy of the handle can be used to make this call. Users should note that MPI will deallocate the object after free_fn executes. At this point, user copies of the request handle no longer point to a valid request. MPI will not set user copies to MPI_REQUEST_NULL in this case, so it is up to the user to avoid accessing this stale handle. This is a special case where MPI defers deallocating the object until a later time that is known by the user. (*End of advice to users.*)

In C, the cancel function is

```
typedef int MPI_Grequest_cancel_function(void *extra_state, int complete);
```

in Fortran

```
SUBROUTINE GREQUEST_CANCEL_FUNCTION(EXTRA_STATE, COMPLETE, IERROR)
    INTEGER IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
    LOGICAL COMPLETE
```

and in C++

```
{typedef int MPI::Grequest::Cancel_function(void* extra_state,
             bool complete); (binding deprecated, see Section 15.2)}
```

cancel_fn function is invoked to start the cancelation of a generalized request. It is called by MPI_CANCEL(request). MPI passes to the callback function complete=true if MPI_GREQUEST_COMPLETE was already called on the request, and complete=false otherwise.

All callback functions return an error code. The code is passed back and dealt with as appropriate for the error code by the MPI function that invoked the callback function. For example, if error codes are returned then the error code returned by the callback function will be returned by the MPI function that invoked the callback function. In the case of an MPI_{WAIT|TEST}{ANY} call that invokes both query_fn and free_fn, the MPI call will return the error code returned by the last callback, namely free_fn. If one or more of the requests in a call to MPI_{WAIT|TEST}{SOME|ALL} failed, then the MPI call will return MPI_ERR_IN_STATUS. In such a case, if the MPI call was passed an array of statuses, then 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() *(binding deprecated, see Section 15.2)* }

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.*)

### 12.2.1 Examples

**Example 12.1** 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;
  int root;
  int valin;
  int *valout;
  MPI_Request request;
  } ARGS;
```

```
int myreduce(MPI_Comm comm, int tag, int root,
             int valin, int *valout, MPI_Request *request)
{
   ARGS *args;
   pthread_t thread;

   /* start request */
   MPI_Grequest_start(query_fn, free_fn, cancel_fn, NULL, request);

   args = (ARGS*)malloc(sizeof(ARGS));
   args->comm = comm;
   args->tag = tag;
   args->root = root;
   args->valin = valin;
   args->valout = valout;
   args->request = *request;

   /* spawn thread to handle request */
   /* The availability of the pthread_create call is system dependent */
   pthread_create(&thread, NULL, reduce_thread, args);

   return MPI_SUCCESS;
}

/* thread code */
void* reduce_thread(void *ptr)
{
   int lchild, rchild, parent, lval, rval, val;
   MPI_Request req[2];
   ARGS *args;

   args = (ARGS*)ptr;

   /* compute left,right child and parent in tree; set
      to MPI_PROC_NULL if does not exist  */
   /* code not shown */
   ...

   MPI_Irecv(&lval, 1, MPI_INT, lchild, args->tag, args->comm, &req[0]);
   MPI_Irecv(&rval, 1, MPI_INT, rchild, args->tag, args->comm, &req[1]);
   MPI_Waitall(2, req, MPI_STATUSES_IGNORE);
   val = lval + args->valin + rval;
   MPI_Send( &val, 1, MPI_INT, parent, args->tag, args->comm );
   if (parent == MPI_PROC_NULL) *(args->valout) = val;
   MPI_Grequest_complete((args->request));
   free(ptr);
```

```
   return(NULL);
}

int query_fn(void *extra_state, MPI_Status *status)
{
   /* always send just one int */
   MPI_Status_set_elements(status, MPI_INT, 1);
   /* can never cancel so always true */
   MPI_Status_set_cancelled(status, 0);
   /* choose not to return a value for this */
   status->MPI_SOURCE = MPI_UNDEFINED;
   /* tag has no meaning for this generalized request */
   status->MPI_TAG = MPI_UNDEFINED;
   /* this generalized request never fails */
   return MPI_SUCCESS;
}


int free_fn(void *extra_state)
{
   /* this generalized request does not need to do any freeing */
   /* as a result it never fails here */
   return MPI_SUCCESS;
}


int cancel_fn(void *extra_state, int complete)
{
   /* This generalized request does not support cancelling.
      Abort if not already done.  If done then treat as if cancel failed.*/
   if (!complete) {
     fprintf(stderr,
             "Cannot cancel generalized request - aborting program\n");
     MPI_Abort(MPI_COMM_WORLD, 99);
     }
   return MPI_SUCCESS;
}
```

## 12.3   Associating Information with Status

MPI supports several different types of requests besides those for point-to-point operations. These range from MPI calls for I/O to generalized requests. It is desirable to allow these calls use the same request mechanism. This allows one to wait or test on different types of requests. However, MPI_{TEST|WAIT}{ANY|SOME|ALL} returns a status with information about the request. With the generalization of requests, one needs to define what information will be returned in the status object.

Each MPI call fills in the appropriate fields in the status object. Any unused fields will

have undefined values. A call to MPI_{TEST|WAIT}{ANY|SOME|ALL} can modify any of the fields in the status object. Specifically, it can modify fields that are undefined. The fields with meaningful value for a given request are defined in the sections with the new request.

Generalized requests raise additional considerations. Here, the user provides the functions to deal with the request. Unlike other MPI calls, the user needs to provide the information to be returned in status. The status argument is provided directly to the callback function where the status needs to be set. Users can directly set the values in 3 of the 5 status values. The count and cancel fields are opaque. To overcome this, these calls are provided:

MPI_STATUS_SET_ELEMENTS(status, datatype, count)

| | | |
|------|----------|-----------------------------------------------|
| INOUT | status | status with which to associate count (Status) |
| IN | datatype | datatype associated with count (handle) |
| IN | count | number of elements to associate with status (integer) |

```
int MPI_Status_set_elements(MPI_Status *status, MPI_Datatype datatype,
              int count)
```

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

```
{void MPI::Status::Set_elements(const MPI::Datatype& datatype, int
              count)(binding deprecated, see Section 15.2) }
```

MPI_STATUS_SET_ELEMENTS_X(status, datatype, count)

| | | |
|------|----------|-----------------------------------------------|
| INOUT | status | status with which to associate count (Status) |
| IN | datatype | datatype associated with count (handle) |
| IN | count | number of elements to associate with status (integer) |

```
int MPI_Status_set_elements_x(MPI_Status *status, MPI_Datatype datatype,
              MPI_Count count)
```

```
MPI_STATUS_SET_ELEMENTS_X(STATUS, DATATYPE, COUNT, IERROR)
    INTEGER STATUS(MPI_STATUS_SIZE), DATATYPE, IERROR
    INTEGER (KIND=MPI_COUNT_KIND) COUNT
```

This call modifies the opaque part of status so that a call to MPI_GET_ELEMENTS or MPI_GET_ELEMENTS_X will return count. MPI_GET_COUNT will return a compatible value.

> *Rationale.* The number of elements is set instead of the count because the former can deal with a nonintegral number of datatypes. (*End of rationale.*)

A subsequent call to MPI_GET_COUNT(status, datatype, count) [ or to], MPI_GET_ELEMENTS(status, datatype, count) , or MPI_GET_ELEMENTS_X(status, datatype,

ticket265.

ticket265.

ticket265.
ticket265.

count) must use a datatype argument that has the same type signature as the datatype argument that was used in the call to MPI_STATUS_SET_ELEMENTS   or
MPI_STATUS_SET_ELEMENTS_X.

> *Rationale.*   This is similar to the restriction that holds when count is set by a receive operation: in that case, the calls to MPI_GET_COUNT and MPI_GET_ELEMENTS must use a datatype with the same signature as the datatype used in the receive call. (*End of rationale.*)

MPI_STATUS_SET_CANCELLED(status, flag)

| INOUT | status | status with which to associate cancel flag (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) *(binding deprecated, see Section 15.2)*
          }

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.*)

## 12.4   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 [33], but the syntax and semantics of thread calls are not specified here — these are beyond the scope of this document.

ticket265.

### 12.4.1   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 12.2**  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.*)

### 12.4.2 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|TEST}{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 multi-threaded 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.*)

### 12.4.3   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)*(binding deprecated, see Section 15.2)* }

{int MPI::Init_thread(int required)*(binding deprecated, see Section 15.2)* }

> *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 8.7. (*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 417).

**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.

**Unofficial Draft for Comment Only**

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 many applications well, 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() *(binding deprecated, see Section 15.2)* }

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() *(binding deprecated, see Section 15.2)* }

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

# Chapter 16

# Language Bindings

## 16.1 C++

### 16.1.1 Overview

The C++ language bindings have been deprecated.

There are some issues specific to C++ that must be considered in the design of an 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 design of MPI was based on the notion of objects, so a natural set of classes is already part of MPI.

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.

### 16.1.2 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 binding must provide a direct and unambiguous mapping to the specified functionality of MPI. (*End of rationale.*)

   .

### 16.1.3   C++ Classes for MPI

All MPI classes, constants, and functions are declared within the scope of an `MPI namespace`. Thus, instead of the `MPI_` prefix that is used in C and Fortran, MPI functions essentially have an `MPI::` prefix.

The members of the `MPI` namespace are those classes corresponding to objects implicitly used by MPI. An abbreviated definition of the `MPI` namespace and its member classes is as follows:

```
namespace MPI {
  class Comm                              {...};
  class Intracomm : public Comm           {...};
  class Graphcomm : public Intracomm      {...};
  class Distgraphcomm : public Intracomm {...};
  class Cartcomm  : public Intracomm      {...};
  class Intercomm : public Comm           {...};
  class Datatype                          {...};
  class Errhandler                        {...};
  class Exception                         {...};
  class File                              {...};
  class Group                             {...};
  class Info                              {...};
  class Op                                {...};
  class Request                           {...};
  class Prequest  : public Request        {...};
  class Grequest  : public Request        {...};
  class Status                            {...};
  class Win                               {...};
};
```

Note that there are a small number of derived classes, and that virtual inheritance is *not* used.

### 16.1.4   Class Member Functions for MPI

Besides the member functions which constitute the C++ language bindings for MPI, the C++ language interface has additional functions (as required by the C++ language). In particular, the C++ language interface must provide a constructor and destructor, an assignment operator, and comparison operators.

The complete set of C++ language bindings for MPI is presented in Annex A.4. The 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.

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 16.1** 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.*)

### 16.1.5 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>()` *(binding deprecated, see Section 15.2)* }

{ `~MPI::<CLASS>()` *(binding deprecated, see Section 15.2)* }

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 16.2** In the following code fragment, the test will return `true` and the message will be sent to `cout`.

```
void foo()
{
  MPI::Intracomm bar;

  if (bar == MPI::COMM_NULL)
    cout << "bar is MPI::COMM_NULL" << endl;
}
```

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.
>
>    ```
>    void example_function()
>    {
>      MPI::Intracomm foo_comm(MPI::COMM_WORLD), bar_comm;
>      bar_comm = MPI::COMM_WORLD.Dup();
>      // rest of function
>    }
>    ```
>
> (*End of rationale.*)

Copy / Assignment   The copy constructor and assignment operator are prototyped as follows:

{ MPI::<CLASS>(const MPI::<CLASS>& data)*(binding deprecated, see Section 15.2)* }

{ MPI::<CLASS>& MPI::<CLASS>::operator=(const MPI::<CLASS>& data)*(binding deprecated, see Section 15.2)* }

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 16.3** 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*(binding deprecated, see Section 15.2)* }

{bool MPI::<CLASS>::operator!=(const MPI::<CLASS>& data) const*(binding deprecated, see Section 15.2)* }

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`.

## 16.1.6   C++ Datatypes

Table 16.1 lists all of the C++ predefined MPI datatypes and their corresponding C and C++ datatypes, Table 16.2 lists all of the Fortran predefined MPI datatypes and their corresponding Fortran 77 datatypes.  Table 16.3 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 on page 27 and Sections 4.2 on page 122, respectively.

The following table defines groups of MPI predefined datatypes:

| | |
|---|---|
| C integer: | MPI::INT, MPI::LONG, MPI::SHORT, MPI::UNSIGNED_SHORT, MPI::UNSIGNED, MPI::UNSIGNED_LONG, MPI::_LONG_LONG, MPI::UNSIGNED_LONG_LONG, MPI::SIGNED_CHAR, MPI::UNSIGNED_CHAR |
| Fortran integer: | MPI::INTEGER and handles returned from |

| MPI datatype | C datatype | C++ datatype |
|---|---|---|
| MPI::CHAR | char | char |
| MPI::SHORT | signed short | signed short |
| MPI::INT | signed int | signed int |
| MPI::LONG | signed long | signed long |
| MPI::LONG_LONG | signed long long | signed long 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::UNSIGNED_LONG_LONG | unsigned long long | unsigned long long |
| 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::WCHAR | wchar_t | wchar_t |
| MPI::BYTE | | |
| MPI::PACKED | | |

Table 16.1:  C++ names for the MPI C and C++ predefined datatypes, and their corresponding C/C++ datatypes.

| MPI datatype | Fortran datatype |
|---|---|
| MPI::INTEGER | INTEGER |
| MPI::REAL | REAL |
| MPI::DOUBLE_PRECISION | DOUBLE PRECISION |
| MPI::F_COMPLEX | COMPLEX |
| MPI::LOGICAL | LOGICAL |
| MPI::CHARACTER | CHARACTER(1) |
| MPI::BYTE | |
| MPI::PACKED | |

Table 16.2:  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::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::INTEGER16 | Explicit size type |
| MPI::REAL2 | Explicit size type |
| MPI::REAL4 | Explicit size type |
| MPI::REAL8 | Explicit size type |
| MPI::REAL16 | Explicit size type |
| MPI::F_COMPLEX4 | Explicit size type |
| MPI::F_COMPLEX8 | Explicit size type |
| MPI::F_COMPLEX16 | Explicit size type |
| MPI::F_COMPLEX32 | Explicit size type |

Table 16.3: C++ names for other MPI datatypes. Implementations may also define other optional types (e.g., `MPI::INTEGER8`).

|  | MPI::Datatype::Create_f90_integer, |
|  | and if available: MPI::INTEGER1, |
|  | MPI::INTEGER2, MPI::INTEGER4, |
|  | MPI::INTEGER8, MPI::INTEGER16 |
| Floating point: | MPI::FLOAT, MPI::DOUBLE, MPI::REAL, |
|  | MPI::DOUBLE_PRECISION, |
|  | MPI::LONG_DOUBLE |
|  | and handles returned from |
|  | MPI::Datatype::Create_f90_real, |
|  | and if available: MPI::REAL2, |
|  | MPI::REAL4, MPI::REAL8, MPI::REAL16 |
| Logical: | MPI::LOGICAL, MPI::BOOL |
| Complex: | MPI::F_COMPLEX, MPI::COMPLEX, |
|  | MPI::F_DOUBLE_COMPLEX, |
|  | MPI::DOUBLE_COMPLEX, |
|  | MPI::LONG_DOUBLE_COMPLEX |
|  | and handles returned from |
|  | MPI::Datatype::Create_f90_complex, |
|  | and if available: MPI::F_DOUBLE_COMPLEX, |
|  | MPI::F_COMPLEX4, MPI::F_COMPLEX8, |
|  | MPI::F_COMPLEX16, MPI::F_COMPLEX32 |
| Byte: | MPI::BYTE |

Valid datatypes for each reduction operation are specified below in terms of the groups defined above.

| Op | Allowed Types |
|---|---|
| MPI::MAX, MPI::MIN | C integer, Fortran integer, Floating point |
| MPI::SUM, MPI::PROD | C integer, Fortran integer, Floating point, Complex |
| MPI::LAND, MPI::LOR, MPI::LXOR | C integer, Logical |
| MPI::BAND, MPI::BOR, MPI::BXOR | C integer, Fortran integer, Byte |

MPI::MINLOC and MPI::MAXLOC perform just as their C and Fortran counterparts; see Section 5.9.4 on page 172.

### 16.1.7 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 six different types of communicators: `MPI::Comm`, `MPI::Intercomm`, `MPI::Intracomm`, `MPI::Cartcomm`, `MPI::Graphcomm`, and `MPI::Distgraphcomm`. `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`, `MPI::Graphcomm`, and `MPI::Distgraphcomm` 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 16.4** 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 16.5** 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
```

```
Distgraphcomm& Distgraphcomm::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, Graphcomm,
  // and Distgraphcomm are similarly defined
};
```

> (*End of advice to implementors.*)

## 16.1.8   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 a nonblocking MPI call throws an exception while making progress in the background,
the behavior is implementation dependent.

The error handler `MPI::ERRORS_THROW_EXCEPTIONS` causes an `MPI::Exception` to be
thrown for any MPI result code other than `MPI::SUCCESS`. The public interface to
`MPI::Exception` class is defined as follows:

```
namespace MPI {
  class Exception {
  public:

     Exception(int error_code);

     int Get_error_code() const;
```

```
    int Get_error_class() const;
    const char *Get_error_string() const;
  };
};
```

> *Advice to implementors.*
>
> The exception will be thrown within the body of `MPI::ERRORS_THROW_EXCEPTIONS`. It is expected that control will be returned to the user when the exception is thrown. Some `MPI` functions specify certain return information in their parameters in the case of an error and `MPI_ERRORS_RETURN` is specified. The same type of return information must be provided when exceptions are thrown.
>
> For example, `MPI_WAITALL` puts an error code for each request in the corresponding entry in the status array and returns `MPI_ERR_IN_STATUS`. When using `MPI::ERRORS_THROW_EXCEPTIONS`, it is expected that the error codes in the status array will be set appropriately before the exception is thrown.
>
> (*End of advice to implementors.*)

### 16.1.9  Mixed-Language Operability

The C++ language interface provides functions listed below for mixed-language operability. These functions provide for a seamless transition between C and C++. 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>`.

```
MPI::<CLASS>& MPI::<CLASS>::operator=(const MPI_<CLASS>& data)
```

```
MPI::<CLASS>(const MPI_<CLASS>& data)
```

```
MPI::<CLASS>::operator MPI_<CLASS>() const
```

These functions are discussed in Section 16.3.4.

### 16.1.10  Profiling

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 imple-
   menting 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 16.6** 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 16.7** pmpi.cc, to be compiled into libmpi.a.

```
int PMPI::Comm::Get_size() const
{
  // Implementation of MPI_COMM_SIZE
}
```

**Example 16.8** constants.cc, to be compiled into libmpi.a.

```
const MPI::Intracomm MPI::COMM_WORLD;
```

**Example 16.9** mpi_no_profile.cc, to be compiled into libmpi.a.

```
int MPI::Comm::Get_size() const
{
  return pmpi_comm.Get_size();
}
```

**Example 16.10** `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.*)

## 16.2  Fortran Support

### 16.2.1  Overview

The Fortran MPI-2 language bindings have been designed to be compatible with the Fortran 90 standard (and later). These bindings are in most cases compatible with Fortran 77, implicit-style interfaces.

> *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 16.2.3 and 16.2.4. In the rest of this section, "Fortran" and "Fortran 90" shall refer to "Fortran 90" and its successors, unless qualified.

1. **Basic Fortran Support** An implementation with this level of Fortran support provides the original Fortran bindings specified in MPI-1, with small additional requirements specified in Section 16.2.3.

2. **Extended Fortran Support** An implementation with this level of Fortran support provides Basic Fortran Support plus additional features that specifically support Fortran 90, as described in Section 16.2.4.

A compliant MPI-2 implementation providing a Fortran interface must provide Extended Fortran Support unless the target compiler does not support modules or `KIND`-parameterized types.

### 16.2.2  Problems With Fortran Bindings for MPI

This section discusses a number of problems that may arise when using MPI in a Fortran program. It is intended as advice to users, and clarifies how MPI interacts with Fortran. It does not add to the standard, but is intended to clarify the standard.

As noted in the original MPI specification, the interface violates the Fortran standard in several ways. While these cause few problems for Fortran 77 programs, they become

more significant for Fortran 90 programs, so that users must exercise care when using new
Fortran 90 features. The violations were originally adopted and have been retained because
they are important for the usability of MPI. The rest of this section describes the potential
problems in detail. It supersedes and replaces the discussion of Fortran bindings in the
original MPI specification (for Fortran 90, not Fortran 77).

The following MPI features are inconsistent with Fortran 90.

1. An MPI subroutine with a choice argument may be called with different argument
   types.

2. An MPI subroutine with an assumed-size dummy argument may be passed an actual
   scalar argument.

3. 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.

4. An MPI implementation may read or modify user data (e.g., communication buffers
   used by nonblocking communications) concurrently with a user program that is exe-
   cuting outside of MPI calls.

5. Several named "constants," such as MPI_BOTTOM, MPI_IN_PLACE,
   MPI_STATUS_IGNORE, MPI_STATUSES_IGNORE, MPI_ERRCODES_IGNORE,
   MPI_UNWEIGHTED, MPI_ARGV_NULL, and MPI_ARGVS_NULL are not ordinary Fortran
   constants and require a special implementation. See Section 2.5.4 on page 14 for more
   information.

6. The memory allocation routine MPI_ALLOC_MEM can't be usefully used in Fortran
   without a language extension that allows the allocated memory to be associated with
   a Fortran variable.

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

- MPI identifiers exceed 6 characters.

- MPI identifiers may contain underscores after the first character.

- MPI requires an include file, `mpif.h`. On systems that do not support include files,
  the implementation should specify the values of named constants.

- Many routines in MPI have KIND-parameterized integers (e.g., MPI_ADDRESS_KIND
  and MPI_OFFSET_KIND) that hold address information. On systems that do not sup-
  port Fortran 90-style parameterized types, `INTEGER*8` or `INTEGER` should be used
  instead.

MPI-1 contained several routines that take address-sized information as input or return
address-sized information as output. In C such arguments were of type MPI_Aint and in
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 16 and Section 4.1.1 on page 79 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 16.2.4. 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]

Because MPI dummy buffer arguments are assumed-size arrays, this leads to a serious problem for a nonblocking call: the compiler copies the temporary array back on return but MPI continues to copy data to the memory that held it. For example, consider the following code fragment:

---

[1]Technically, the Fortran standards are worded to allow non-contiguous storage of any array data.

```
real a(100)
call MPI_IRECV(a(1:100:2), MPI_REAL, 50, ...)
```

Since the first dummy argument to MPI_IRECV is an assumed-size array (`<type> buf(*)`), the array section `a(1:100:2)` is copied to a temporary before being passed to MPI_IRECV, so that it is contiguous in memory. MPI_IRECV returns immediately, and data is copied from the temporary back into the array `a`. Sometime later, MPI may write to the address of the deallocated temporary. Copying is also a problem for MPI_ISEND since the temporary array may be deallocated before the data has all been sent from it.

Most Fortran 90 compilers do not make a copy if the actual argument is the whole of an explicit-shape or assumed-size array or is a 'simple' section such as `A(1:N)` of such an array. (We define 'simple' more fully in the next paragraph.) Also, many compilers treat allocatable arrays the same as they treat explicit-shape arrays in this regard (though we know of one that does not). However, the same is not true for assumed-shape and pointer arrays; since they may be discontiguous, copying is often done. It is this copying that causes problems for MPI as described in the previous paragraph.

Our formal definition of a 'simple' array section is

```
name ( [:,]... [<subscript>]:[<subscript>] [,<subscript>]... )
```

That is, there are zero or more dimensions that are selected in full, then one dimension selected without a stride, then zero or more dimensions that are selected with a simple subscript. Examples are

```
A(1:N), A(:,N), A(:,1:N,1), A(1:6,N), A(:,:,1:N)
```

Because of Fortran's column-major ordering, where the first index varies fastest, a simple section of a contiguous array will also be contiguous.[2]

The same problem can occur with a scalar argument. Some compilers, even for Fortran 77, make a copy of some scalar dummy arguments within a called procedure. That this can cause a problem is illustrated by the example

```
call user1(a,rq)
call MPI_WAIT(rq,status,ierr)
write (*,*) a

subroutine user1(buf,request)
call MPI_IRECV(buf,...,request,...)
end
```

If `a` is copied, MPI_IRECV 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.

---

[2]To keep the definition of 'simple' simple, we have chosen to require all but one of the section subscripts to be without bounds. A colon without bounds makes it obvious both to the compiler and to the reader that the whole of the dimension is selected. It would have been possible to allow cases where the whole dimension is selected with one or two bounds, but this means for the reader that the array declaration or most recent allocation has to be consulted and for the compiler that a run-time check may be required.

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 nonblocking 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, e.g., MPI_BOTTOM. The complete list can be found in Section 2.5.4 on page 14. 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
end type mytype

type(mytype) foo
integer blocklen(3), type(3)
integer(MPI_ADDRESS_KIND) disp(3), base

call MPI_GET_ADDRESS(foo%i, disp(1), ierr)
```

```
    call MPI_GET_ADDRESS(foo%x, disp(2), ierr)
    call MPI_GET_ADDRESS(foo%d, disp(3), ierr)

    base = disp(1)
    disp(1) = disp(1) - base
    disp(2) = disp(2) - base
    disp(3) = disp(3) - base

    blocklen(1) = 1
    blocklen(2) = 1
    blocklen(3) = 1

    type(1) = MPI_INTEGER
    type(2) = MPI_REAL
    type(3) = MPI_DOUBLE_PRECISION

    call MPI_TYPE_CREATE_STRUCT(3, blocklen, disp, type, newtype, ierr)
    call MPI_TYPE_COMMIT(newtype, ierr)

! unpleasant to send foo%i instead of foo, but it works for scalar
! entities of type mytype
    call MPI_SEND(foo%i, 1, newtype, ...)
```

## A Problem with Register Optimization

MPI provides operations that may be hidden from the user code and run concurrently with it, accessing the same memory as user code. Examples include the data transfer for an MPI_IRECV. The optimizer of a compiler will assume that it can recognize periods when a copy of a variable can be kept in a register without reloading from or storing to memory. When the user code is working with a register copy of some variable while the hidden operation reads or writes the memory copy, problems occur. This section discusses register optimization pitfalls.

When a variable is local to a Fortran subroutine (i.e., not in a module or `COMMON block`), the compiler will assume that it cannot be modified by a called subroutine unless it is an actual argument of the call. In the most common linkage convention, the subroutine is expected to save and restore certain registers. Thus, the optimizer will assume that a register which held a valid copy of such a variable before the call will still hold a valid copy on return.

Normally users are not afflicted with this. But the user should pay attention to this section if in his/her program a buffer argument to an MPI_SEND, MPI_RECV etc., uses 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.

Example 16.11 shows what Fortran compilers are allowed to do.

**Unofficial Draft for Comment Only**

**Example 16.11** Fortran 90 register optimization.

This source ...                                  can be compiled as:

```
call MPI_GET_ADDRESS(buf,bufaddr,        call MPI_GET_ADDRESS(buf,...)
            ierror)
call MPI_TYPE_CREATE_STRUCT(1,1,         call MPI_TYPE_CREATE_STRUCT(...)
            bufaddr,
            MPI_REAL,type,ierror)
call MPI_TYPE_COMMIT(type,ierror)        call MPI_TYPE_COMMIT(...)
val_old = buf                            register = buf
                                         val_old = register
call MPI_RECV(MPI_BOTTOM,1,type,...)     call MPI_RECV(MPI_BOTTOM,...)
val_new = buf                            val_new = register
```

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.

Example 16.12 shows extreme, but allowed, possibilities.

**Example 16.12** Fortran 90 register optimization – extreme.

```
Source                     compiled as                or compiled as
call MPI_IRECV(buf,..req)  call MPI_IRECV(buf,..req)  call MPI_IRECV(buf,..req)
                           register = buf             b1 = buf
call MPI_WAIT(req,..)      call MPI_WAIT(req,..)      call MPI_WAIT(req,..)
b1 = buf                   b1 := register
```

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,...)
    call DD(buf)
    ```

    with the separately compiled

```
subroutine DD(buf)
  integer buf
end
```

(assuming that `buf` has type `INTEGER`). The compiler may be similarly prevented from moving a reference to a variable across a call to an MPI subroutine.

In the case of a nonblocking call, as in the above call of MPI_WAIT, no reference to the buffer is permitted until it has been verified that the transfer has been completed. Therefore, in this case, the extra call ahead of the MPI call is not necessary, i.e., the call of MPI_WAIT in the example might be replaced by

```
call MPI_WAIT(req,..)
call DD(buf)
```

- An alternative is to put the buffer or variable into a module or a common block and access it through a `USE` or `COMMON` statement in each scope where it is referenced, defined or appears as an actual argument in a call to an MPI routine. The compiler will then have to assume that the MPI procedure (MPI_RECV in the above example) may alter the buffer or variable, provided that the compiler cannot analyze that the MPI procedure does not reference the module or common block.

The `VOLATILE` attribute, available in later versions of Fortran, gives the buffer or variable the properties needed, but it may inhibit optimization of any code containing the buffer or variable.

In C, subroutines which modify variables that are not in the argument list will not cause register optimization problems. This is because taking pointers to storage objects by using the & operator and later referencing the objects by way of the pointer is an integral part of the language. A C compiler understands the implications, so that the problem should not occur, in general. However, some compilers do offer optional aggressive optimization levels which may not be safe.

### 16.2.3  Basic Fortran Support

Because Fortran 90 is (for all practical purposes) a superset of Fortran 77, Fortran 90 (and future) programs can use the original Fortran interface. The following additional requirements are added:

1. Implementations are required to provide the file `mpif.h`, as described in the original MPI-1 specification.

2. `mpif.h` must be valid and equivalent for both fixed- and free- source form.

   *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.*)

**Unofficial Draft for Comment Only**

### 16.2.4   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 16.2.5.

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 `use`d 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 [26]. 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.*)

### 16.2.5 Additional Support for Fortran Numeric Intrinsic Types

The routines in this section are part of Extended Fortran Support described in Section 16.2.4.

MPI 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 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 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 | p | precision, in decimal digits (integer) |
| IN | r | decimal exponent range (integer) |
| OUT | newtype | the requested MPI datatype (handle) |

```
int MPI_Type_create_f90_real(int p, int r, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_F90_REAL(P, R, NEWTYPE, IERROR)
    INTEGER P, R, NEWTYPE, IERROR
```

{static MPI::Datatype MPI::Datatype::Create_f90_real(int p, int r) *(binding deprecated, see Section 15.2)* }

This function returns a predefined MPI datatype that matches a REAL variable of KIND `selected_real_kind(p, r)`. In the model described above it returns a handle for the element D(p, r). Either p or r may be omitted from calls to `selected_real_kind(p, r)`

(but not both). Analogously, either p or r may be set to MPI_UNDEFINED. In communication, an MPI datatype A returned by MPI_TYPE_CREATE_F90_REAL matches a datatype B if and only if B was returned by MPI_TYPE_CREATE_F90_REAL called with the same values for p and r or B is a duplicate of such a datatype. Restrictions on using the returned datatype with the "external32" data representation are given on page 523.

It is erroneous to supply values for p and r not supported by the compiler.

MPI_TYPE_CREATE_F90_COMPLEX(p, r, newtype)

| IN | p | precision, in decimal digits (integer) |
|-----|---------|-----|
| IN | r | decimal exponent range (integer) |
| OUT | newtype | the requested MPI datatype (handle) |

```
int MPI_Type_create_f90_complex(int p, int r, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_F90_COMPLEX(P, R, NEWTYPE, IERROR)
    INTEGER P, R, NEWTYPE, IERROR
```

{static MPI::Datatype MPI::Datatype::Create_f90_complex(int p,
            int r)*(binding deprecated, see Section 15.2)* }

This function returns a predefined MPI datatype that matches a COMPLEX variable of KIND selected_real_kind(p, r). Either p or r may be omitted from calls to selected_real_kind(p, r) (but not both). Analogously, either p or r may be set to MPI_UNDEFINED. 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 523.

It is erroneous to supply values for p and r not supported by the compiler.

MPI_TYPE_CREATE_F90_INTEGER(r, newtype)

| IN | r | decimal exponent range, i.e., number of decimal digits (integer) |
|-----|---------|-----|
| OUT | newtype | the requested MPI datatype (handle) |

```
int MPI_Type_create_f90_integer(int r, MPI_Datatype *newtype)
```

```
MPI_TYPE_CREATE_F90_INTEGER(R, NEWTYPE, IERROR)
    INTEGER R, NEWTYPE, IERROR
```

{static MPI::Datatype MPI::Datatype::Create_f90_integer(int r)*(binding deprecated, see Section 15.2)* }

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 523.

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 combiners 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 appli-
> cation is not allowed to free the returned predefined, unnamed datatype handles. To
> prevent the creation of a potentially huge amount of handles, a high quality MPI imple-
> mentation 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 13.5.2 on page 461) or user-defined (Sec-
> tion 13.5.3 on page 462) 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 13.5.2 on page 461.

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 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
```

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.*)


MPI_TYPE_MATCH_SIZE(typeclass, size, type)

  IN        typeclass                      generic type specifier (integer)

  IN        size                           size, in bytes, of representation (integer)

  OUT       type                           datatype with correct type, size (handle)


```
int MPI_Type_match_size(int typeclass, int size, MPI_Datatype *type)
```

```
MPI_TYPE_MATCH_SIZE(TYPECLASS, SIZE, TYPE, IERROR)
    INTEGER TYPECLASS, SIZE, TYPE, IERROR
```

{static MPI::Datatype MPI::Datatype::Match_size(int typeclass,
        int size)*(binding deprecated, see Section 15.2)* }

typeclass is one of MPI_TYPECLASS_REAL, MPI_TYPECLASS_INTEGER and
MPI_TYPECLASS_COMPLEX, corresponding to the desired **typeclass**. The function returns
an MPI datatype matching a local variable of type (**typeclass**, **size**).

This function returns a reference (handle) to one of the predefined named datatypes, not
a duplicate. This type cannot be freed. MPI_TYPE_MATCH_SIZE can be used to obtain a
size-specific type that matches a Fortran numeric intrinsic type by first calling MPI_SIZEOF
in order to compute the variable size, and then calling MPI_TYPE_MATCH_SIZE to find a
suitable datatype. In C and C++, one can use the C function sizeof(), instead of
MPI_SIZEOF. In addition, for variables of default kind the variable's size can be computed
by a call to MPI_TYPE_GET_EXTENT, if the typeclass is known. It is erroneous to specify
a size not supported by the compiler.

> *Rationale.* This is a convenience function. Without it, it can be tedious to find the
> correct named type. See note to implementors below. (*End of rationale.*)

> *Advice to implementors.* This function could be implemented as a series of tests.

```
int MPI_Type_match_size(int typeclass, int size, MPI_Datatype *rtype)
{
  switch(typeclass) {
      case MPI_TYPECLASS_REAL: switch(size) {
        case 4: *rtype = MPI_REAL4; return MPI_SUCCESS;
        case 8: *rtype = MPI_REAL8; return MPI_SUCCESS;
        default: error(...);
      }
      case MPI_TYPECLASS_INTEGER: switch(size) {
         case 4: *rtype = MPI_INTEGER4; return MPI_SUCCESS;
         case 8: *rtype = MPI_INTEGER8; return MPI_SUCCESS;
         default: error(...);        }
      ... etc. ...
  }
}
```

> (*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, com-
munication 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 environ-
> ment. 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)
endif

call MPI_BARRIER(MPI_COMM_WORLD, ierror)

if (myrank .eq. 1) then
   call MPI_FILE_OPEN(MPI_COMM_SELF, 'foo', MPI_MODE_RDONLY,  &
                  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)
```

```
    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.*)

## 16.3 Language Interoperability

### 16.3.1 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.

### 16.3.2 Assumptions

We assume that conventions exist for programs written in one language to call routines 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 `INTEGER`s, 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.

### 16.3.3   Initialization

A call to MPI_INIT or MPI_INIT_THREAD, 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.*)

### 16.3.4   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 16.13** 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 16.14**  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 C++ library call. This example shows a C interface function that invokes a C++ library call with a C communicator; the communicator is automatically promoted to a C++ handle when the underlying C++ function is invoked.

```
// C++ library function prototype
void cpp_lib_call(MPI::Comm cpp_comm);

// Exported C function prototype
extern "C" {
    void c_interface(MPI_Comm c_comm);
}

void c_interface(MPI_Comm c_comm)
{
```

```
   // the MPI_Comm (c_comm) is automatically promoted to MPI::Comm
   cpp_lib_call(c_comm);
}
```

The following function allows conversion from C++ objects to C MPI handles. In this case, the casting operator is overloaded to provide the functionality.

```
MPI::<CLASS>::operator MPI_<CLASS>() const
```

**Example 16.15** A C library routine is called from a C++ program. The C library routine is prototyped to take an `MPI_Comm` as an argument.

```
// C function prototype
extern "C" {
   void c_lib_call(MPI_Comm c_comm);
}

void cpp_function()
{
   // Create a C++ communicator, and initialize it with a dup of
   //   MPI::COMM_WORLD
   MPI::Intracomm cpp_comm(MPI::COMM_WORLD.Dup());
   c_lib_call(cpp_comm);
}
```

> *Rationale.* Providing conversion from C to C++ via constructors and from C++ to C via casting allows the compiler to make automatic conversions. Calling C from C++ becomes trivial, as does the provision of a C or Fortran interface to a C++ library. (*End of rationale.*)

> *Advice to users.* Note that the casting and promotion operators return new handles by value. Using these new handles as INOUT parameters will affect the internal MPI object, but will *not* affect the original handle from which it was cast. (*End of advice to users.*)

It is important to note that all C++ objects with corresponding C handles can be used interchangeably by an application. For example, an application can cache an attribute on `MPI_COMM_WORLD` and later retrieve it from `MPI::COMM_WORLD`.

### 16.3.5 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.*)

### 16.3.6   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
16.3.9, page 539).

**Example 16.16**

```
! FORTRAN CODE
REAL R(5)
INTEGER TYPE, IERR, AOBLEN(1), AOTYPE(1)
INTEGER (KIND=MPI_ADDRESS_KIND) AODISP(1)

! create an absolute datatype for array R
AOBLEN(1) = 5
CALL MPI_GET_ADDRESS( R, AODISP(1), IERR)
AOTYPE(1) = MPI_REAL
CALL MPI_TYPE_CREATE_STRUCT(1, AOBLEN,AODISP,AOTYPE, TYPE, IERR)
CALL C_ROUTINE(TYPE)

/* C code */

void C_ROUTINE(MPI_Fint *ftype)
{
   int count = 5;
   int lens[2] = {1,1};
   MPI_Aint displs[2];
   MPI_Datatype types[2], newtype;

   /* create an absolute datatype for buffer that consists   */
   /*  of count, followed by R(5)                            */

   MPI_Get_address(&count, &displs[0]);
   displs[1] = 0;
   types[0] = MPI_INT;
   types[1] = MPI_Type_f2c(*ftype);
   MPI_Type_create_struct(2, lens, displs, types, &newtype);
   MPI_Type_commit(&newtype);

   MPI_Send(MPI_BOTTOM, 1, newtype, 1, 0, MPI_COMM_WORLD);
   /* 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 INTEGERs have 32 bits.

This is a problem, irrespective of interlanguage issues. Suppose that a Fortran process has an address space of $\geq$ 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.

**Unofficial Draft for Comment Only**

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 4.1.1 on page 79. 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 and MPI::Aint are used in C and 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.

## 16.3.7  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 6.7 on page 254 define attributes arguments to be of type void* in C, and of type `INTEGER`, in Fortran. On some systems, `INTEGER`s 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 behaves as if it stores, internally, address sized attributes. If Fortran `INTEGER`s 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 6.7, page 254. 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, which is a pointer to MPI_Aint if the attribute was stored with Fortran MPI_xxx_SET_ATTR, and a pointer to int if it was stored with the deprecated Fortran MPI_ATTR_PUT. 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 attribute functions are used, and an integer of kind MPI_ADDRESS_KIND is returned. The conversion may cause truncation if deprecated attribute functions are used. In C, the deprecated routines MPI_Attr_put and MPI_Attr_get behave identical to MPI_Comm_set_attr and MPI_Comm_get_attr.

**Example 16.17**
   A. Setting an attribute value in C

```
int set_val = 3;
struct foo set_struct;

/* Set a value that is a pointer to an int */

MPI_Comm_set_attr(MPI_COMM_WORLD, keyval1, &set_val);
/* Set a value that is a pointer to a struct */
MPI_Comm_set_attr(MPI_COMM_WORLD, keyval2, &set_struct);
/* Set an integer value */
MPI_Comm_set_attr(MPI_COMM_WORLD, keyval3, (void *) 17);
```

   B. Reading the attribute value in C

```
int flag, *get_val;
struct foo *get_struct;

/* Upon successful return, get_val == &set_val
   (and therefore *get_val == 3) */
MPI_Comm_get_attr(MPI_COMM_WORLD, keyval1, &get_val, &flag);
/* Upon successful return, get_struct == &set_struct */
MPI_Comm_get_attr(MPI_COMM_WORLD, keyval2, &get_struct, &flag);
/* Upon successful return, get_val == (void*) 17 */
/*         i.e., (MPI_Aint) get_val == 17 */
MPI_Comm_get_attr(MPI_COMM_WORLD, keyval3, &get_val, &flag);
```

   C. Reading the attribute value with (deprecated) Fortran MPI-1 calls

```
LOGICAL FLAG
INTEGER IERR, GET_VAL, GET_STRUCT

! Upon successful return, GET_VAL == &set_val, possibly truncated
CALL MPI_ATTR_GET(MPI_COMM_WORLD, KEYVAL1, GET_VAL, FLAG, IERR)
! Upon successful return, GET_STRUCT == &set_struct, possibly truncated
CALL MPI_ATTR_GET(MPI_COMM_WORLD, KEYVAL2, GET_STRUCT, FLAG, IERR)
! Upon successful return, GET_VAL == 17
CALL MPI_ATTR_GET(MPI_COMM_WORLD, KEYVAL3, GET_VAL, FLAG, IERR)
```

   D. Reading the attribute value with Fortran MPI-2 calls

```
LOGICAL FLAG
INTEGER IERR
INTEGER (KIND=MPI_ADDRESS_KIND) GET_VAL, GET_STRUCT

! Upon successful return, GET_VAL == &set_val
CALL MPI_COMM_GET_ATTR(MPI_COMM_WORLD, KEYVAL1, GET_VAL, FLAG, IERR)
! Upon successful return, GET_STRUCT == &set_struct
CALL MPI_COMM_GET_ATTR(MPI_COMM_WORLD, KEYVAL2, GET_STRUCT, FLAG, IERR)
! Upon successful return, GET_VAL == 17
CALL MPI_COMM_GET_ATTR(MPI_COMM_WORLD, KEYVAL3, GET_VAL, FLAG, IERR)
```

**Example 16.18**
   A. Setting an attribute value with the (deprecated) Fortran MPI-1 call

```
INTEGER IERR, VAL
VAL = 7
CALL MPI_ATTR_PUT(MPI_COMM_WORLD, KEYVAL, VAL, IERR)
```

   B. Reading the attribute value in C

```
int flag;
int *value;

/* Upon successful return, value points to internal MPI storage and
   *value == (int) 7 */
MPI_Comm_get_attr(MPI_COMM_WORLD, keyval, &value, &flag);
```

   C. Reading the attribute value with (deprecated) Fortran MPI-1 calls

```
LOGICAL FLAG
INTEGER IERR, VALUE

! Upon successful return, VALUE == 7
CALL MPI_ATTR_GET(MPI_COMM_WORLD, KEYVAL, VALUE, FLAG, IERR)
```

   D. Reading the attribute value with Fortran MPI-2 calls

**Unofficial Draft for Comment Only**

```
1    LOGICAL FLAG
2    INTEGER IERR
3    INTEGER (KIND=MPI_ADDRESS_KIND) VALUE
4
5    ! Upon successful return, VALUE == 7 (sign extended)
6    CALL MPI_COMM_GET_ATTR(MPI_COMM_WORLD, KEYVAL, VALUE, FLAG, IERR)
7
```

**Example 16.19**  A. Setting an attribute value via a Fortran MPI-2 call

```
10   INTEGER IERR
11   INTEGER(KIND=MPI_ADDRESS_KIND) VALUE1
12   INTEGER(KIND=MPI_ADDRESS_KIND) VALUE2
13   VALUE1 = 42
14   VALUE2 = INT(2, KIND=MPI_ADDRESS_KIND) ** 40
15
16   CALL MPI_COMM_SET_ATTR(MPI_COMM_WORLD, KEYVAL1, VALUE1, IERR)
17   CALL MPI_COMM_SET_ATTR(MPI_COMM_WORLD, KEYVAL2, VALUE2, IERR)
```

   B. Reading the attribute value in C

```
20   int flag;
21   MPI_Aint *value1, *value2;
22
23   /* Upon successful return, value1 points to internal MPI storage and
24      *value1 == 42 */
25   MPI_Comm_get_attr(MPI_COMM_WORLD, keyval1, &value1, &flag);
26   /* Upon successful return, value2 points to internal MPI storage and
27      *value2 == 2^40 */
28   MPI_Comm_get_attr(MPI_COMM_WORLD, keyval2, &value2, &flag);
```

   C. Reading the attribute value with (deprecated) Fortran MPI-1 calls

```
31   LOGICAL FLAG
32   INTEGER IERR, VALUE1, VALUE2
33
34   ! Upon successful return, VALUE1 == 42
35   CALL MPI_ATTR_GET(MPI_COMM_WORLD, KEYVAL1, VALUE1, FLAG, IERR)
36   ! Upon successful return, VALUE2 == 2^40, or 0 if truncation
37   ! needed (i.e., the least significant part of the attribute word)
38   CALL MPI_ATTR_GET(MPI_COMM_WORLD, KEYVAL2, VALUE2, FLAG, IERR)
```

   D. Reading the attribute value with Fortran MPI-2 calls

```
41   LOGICAL FLAG
42   INTEGER IERR
43   INTEGER (KIND=MPI_ADDRESS_KIND) VALUE1, VALUE2
44
45   ! Upon successful return, VALUE1 == 42
46   CALL MPI_COMM_GET_ATTR(MPI_COMM_WORLD, KEYVAL1, VALUE1, FLAG, IERR)
47   ! Upon successful return, VALUE2 == 2^40
48   CALL MPI_COMM_GET_ATTR(MPI_COMM_WORLD, KEYVAL2, VALUE2, FLAG, IERR)
```

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 call to the deprecated Fortran routine MPI_ATTR_PUT, 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 for predefined attributes, and ensures that no information is lost when attributes are passed from language to language. Because the language interoperability for predefined attributes was defined based on MPI_ATTR_PUT, this definition is kept for compatibility reasons although the routine itself is now deprecated. (*End of rationale.*)

> *Advice to implementors.* Implementations should tag attributes either as (1) address attributes, (2) as INTEGER(KIND=MPI_ADDRESS_KIND) attributes or (3) as INTEGER attributes, according to whether they were set in (1) C (with MPI_Attr_put or MPI_Xxx_set_attr), (2) in Fortran with MPI_XXX_SET_ATTR or (3) with the deprecated Fortran routine MPI_ATTR_PUT. Thus, the right choice can be made when the attribute is retrieved. (*End of advice to implementors.*)

### 16.3.8 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.

### 16.3.9 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 16.3.4. 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.*)

## 16.3.10   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 16.20** In the example below, a Fortran array is sent from Fortran and received in C.

```
! FORTRAN CODE
REAL R(5)
INTEGER TYPE, IERR, MYRANK, AOBLEN(1), AOTYPE(1)
INTEGER (KIND=MPI_ADDRESS_KIND) AODISP(1)

! create an absolute datatype for array R
AOBLEN(1) = 5
CALL MPI_GET_ADDRESS( R, AODISP(1), IERR)
AOTYPE(1) = MPI_REAL
CALL MPI_TYPE_CREATE_STRUCT(1, AOBLEN,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
```

```
/* C code */

void C_ROUTINE(MPI_Fint *fhandle)
{
   MPI_Datatype type;
   MPI_Status status;

   type = MPI_Type_f2c(*fhandle);

   MPI_Recv( MPI_BOTTOM, 1, type, 0, 0, MPI_COMM_WORLD, &status);
}
```

MPI implementors may weaken these type matching rules, and allow messages to be sent with Fortran types and received with C types, and vice versa, when those types match. I.e., if the Fortran type INTEGER is identical to the C type int, then an MPI implementation may allow data to be sent with datatype MPI_INTEGER and be received with datatype MPI_INT. However, such code is not portable.

# Annex A

# Language Bindings Summary

In this section we summarize the specific bindings for C, Fortran, and C++. First we present the constants, type definitions, info values and keys. Then we present the routine prototypes separately for each binding. Listings are alphabetical within 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 middle or right column. Constants with the type `const int` may also be implemented as literal integer constants substituted by the preprocessor.

<div align="center">

**Return Codes**

| C type: `const int` (or unnamed `enum`) Fortran type: `INTEGER` | C++ type: `const int` (or unnamed `enum`) |
|---|---|
| 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 |

**(Continued on next page)**

</div>

**Return Codes (continued)**

| | |
|---|---|
| MPI_ERR_IN_STATUS | MPI::ERR_IN_STATUS |
| MPI_ERR_ACCESS | MPI::ERR_ACCESS |
| MPI_ERR_AMODE | MPI::ERR_AMODE |
| MPI_ERR_ASSERT | MPI::ERR_ASSERT |
| MPI_ERR_BAD_FILE | MPI::ERR_BAD_FILE |
| MPI_ERR_BASE | MPI::ERR_BASE |
| MPI_ERR_CONVERSION | MPI::ERR_CONVERSION |
| MPI_ERR_DISP | MPI::ERR_DISP |
| MPI_ERR_DUP_DATAREP | MPI::ERR_DUP_DATAREP |
| MPI_ERR_FILE_EXISTS | MPI::ERR_FILE_EXISTS |
| MPI_ERR_FILE_IN_USE | MPI::ERR_FILE_IN_USE |
| MPI_ERR_FILE | MPI::ERR_FILE |
| MPI_ERR_INFO_KEY | MPI::ERR_INFO_VALUE |
| MPI_ERR_INFO_NOKEY | MPI::ERR_INFO_NOKEY |
| MPI_ERR_INFO_VALUE | MPI::ERR_INFO_KEY |
| MPI_ERR_INFO | MPI::ERR_INFO |
| MPI_ERR_IO | MPI::ERR_IO |
| MPI_ERR_KEYVAL | MPI::ERR_KEYVAL |
| MPI_ERR_LOCKTYPE | MPI::ERR_LOCKTYPE |
| MPI_ERR_NAME | MPI::ERR_NAME |
| MPI_ERR_NO_MEM | MPI::ERR_NO_MEM |
| MPI_ERR_NOT_SAME | MPI::ERR_NOT_SAME |
| MPI_ERR_NO_SPACE | MPI::ERR_NO_SPACE |
| MPI_ERR_NO_SUCH_FILE | MPI::ERR_NO_SUCH_FILE |
| MPI_ERR_PORT | MPI::ERR_PORT |
| MPI_ERR_QUOTA | MPI::ERR_QUOTA |
| MPI_ERR_READ_ONLY | MPI::ERR_READ_ONLY |
| MPI_ERR_RMA_CONFLICT | MPI::ERR_RMA_CONFLICT |
| MPI_ERR_RMA_SYNC | MPI::ERR_RMA_SYNC |
| MPI_ERR_SERVICE | MPI::ERR_SERVICE |
| MPI_ERR_SIZE | MPI::ERR_SIZE |
| MPI_ERR_SPAWN | MPI::ERR_SPAWN |
| MPI_ERR_UNSUPPORTED_DATAREP | MPI::ERR_UNSUPPORTED_DATAREP |
| MPI_ERR_UNSUPPORTED_OPERATION | MPI::ERR_UNSUPPORTED_OPERATION |
| MPI_ERR_WIN | MPI::ERR_WIN |
| MPI_ERR_LASTCODE | MPI::ERR_LASTCODE |

**Buffer Address Constants**

| C type: `void * const` | C++ type: |
|---|---|
| Fortran type: (predefined memory location) | `void * const` |
| MPI_BOTTOM | MPI::BOTTOM |
| MPI_IN_PLACE | MPI::IN_PLACE |

### Assorted Constants

| C type: `const int` (or unnamed `enum`) | C++ type: |
| Fortran type: `INTEGER` | `const int` (or unnamed `enum`) |
| --- | --- |
| 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_LOCK_EXCLUSIVE | MPI::LOCK_EXCLUSIVE |
| MPI_LOCK_SHARED | MPI::LOCK_SHARED |
| MPI_ROOT | MPI::ROOT |

### Status size and reserved index values (Fortran only)

| Fortran type: `INTEGER` | |
| --- | --- |
| 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)

| Fortran type: `INTEGER` | |
| --- | --- |
| MPI_ADDRESS_KIND | Not defined for C++ |
| [ticket265.] MPI_COUNT_KIND | Not defined for C++ |
| MPI_INTEGER_KIND | Not defined for C++ |
| MPI_OFFSET_KIND | Not defined for C++ |

### Error-handling specifiers

| C type: `MPI_Errhandler` | C++ type: `MPI::Errhandler` |
| Fortran type: `INTEGER` | |
| --- | --- |
| MPI_ERRORS_ARE_FATAL | MPI::ERRORS_ARE_FATAL |
| MPI_ERRORS_RETURN | MPI::ERRORS_RETURN |
| | MPI::ERRORS_THROW_EXCEPTIONS |

### Maximum Sizes for Strings

| C type: `const int` (or unnamed `enum`) | C++ type: |
| Fortran type: `INTEGER` | `const int` (or unnamed `enum`) |
| --- | --- |
| 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 |
|---|---|---|
| C type: `MPI_Datatype` | C++ type: `MPI::Datatype` | |
| Fortran type: `INTEGER` | | |
| MPI_CHAR | MPI::CHAR | `char` (treated as printable character) |
| MPI_SHORT | MPI::SHORT | `signed short int` |
| MPI_INT | MPI::INT | `signed int` |
| MPI_LONG | MPI::LONG | `signed long` |
| MPI_LONG_LONG_INT | MPI::LONG_LONG_INT | `signed long long` |
| MPI_LONG_LONG | MPI::LONG_LONG | `long long` (synonym) |
| MPI_SIGNED_CHAR | MPI::SIGNED_CHAR | `signed char` (treated as integral value) |
| MPI_UNSIGNED_CHAR | MPI::UNSIGNED_CHAR | `unsigned char` (treated as integral value) |
| MPI_UNSIGNED_SHORT | MPI::UNSIGNED_SHORT | `unsigned short` |
| MPI_UNSIGNED | MPI::UNSIGNED | `unsigned int` |
| MPI_UNSIGNED_LONG | MPI::UNSIGNED_LONG | `unsigned long` |
| MPI_UNSIGNED_LONG_LONG | MPI::UNSIGNED_LONG_LONG | `unsigned long long` |
| MPI_FLOAT | MPI::FLOAT | `float` |
| MPI_DOUBLE | MPI::DOUBLE | `double` |
| MPI_LONG_DOUBLE | MPI::LONG_DOUBLE | `long double` |
| MPI_WCHAR | MPI::WCHAR | `wchar_t` (defined in `<stddef.h>`) (treated as printable character) |
| MPI_C_BOOL | (use C datatype handle) | `_Bool` |
| MPI_INT8_T | (use C datatype handle) | `int8_t` |
| MPI_INT16_T | (use C datatype handle) | `int16_t` |
| MPI_INT32_T | (use C datatype handle) | `int32_t` |
| MPI_INT64_T | (use C datatype handle) | `int64_t` |
| MPI_UINT8_T | (use C datatype handle) | `uint8_t` |
| MPI_UINT16_T | (use C datatype handle) | `uint16_t` |
| MPI_UINT32_T | (use C datatype handle) | `uint32_t` |
| MPI_UINT64_T | (use C datatype handle) | `uint64_t` |
| MPI_AINT | (use C datatype handle) | `MPI_Aint` |
| [ticket265.] MPI_COUNT | (use C datatype handle) | `MPI_Count` |
| MPI_OFFSET | (use C datatype handle) | `MPI_Offset` |
| MPI_C_COMPLEX | (use C datatype handle) | `float _Complex` |
| MPI_C_FLOAT_COMPLEX | (use C datatype handle) | `float _Complex` |
| MPI_C_DOUBLE_COMPLEX | (use C datatype handle) | `double _Complex` |
| MPI_C_LONG_DOUBLE_COMPLEX | (use C datatype handle) | `long double _Complex` |
| MPI_BYTE | MPI::BYTE | (any C/C++ type) |
| MPI_PACKED | MPI::PACKED | (any C/C++ type) |

**Unofficial Draft for Comment Only**

| Named Predefined Datatypes | | Fortran types |
|---|---|---|
| C type: `MPI_Datatype` | C++ type: `MPI::Datatype` | |
| Fortran type: `INTEGER` | | |
| `MPI_INTEGER` | `MPI::INTEGER` | `INTEGER` |
| `MPI_REAL` | `MPI::REAL` | `REAL` |
| `MPI_DOUBLE_PRECISION` | `MPI::DOUBLE_PRECISION` | `DOUBLE PRECISION` |
| `MPI_COMPLEX` | `MPI::F_COMPLEX` | `COMPLEX` |
| `MPI_LOGICAL` | `MPI::LOGICAL` | `LOGICAL` |
| `MPI_CHARACTER` | `MPI::CHARACTER` | `CHARACTER(1)` |
| `MPI_AINT` | (use C datatype handle) | `INTEGER (KIND=MPI_ADDRESS_KIND)` |
| `MPI_OFFSET` | (use C datatype handle) | `INTEGER (KIND=MPI_OFFSET_KIND)` |
| [ticket265.] MPI_COUNT | (use C datatype handle) | `INTEGER (KIND=MPI_COUNT_KIND)` |
| 265 `MPI_BYTE` | `MPI::BYTE` | (any Fortran type) |
| `MPI_PACKED` | `MPI::PACKED` | (any Fortran type) |

| C++-Only Named Predefined Datatypes | C++ types |
|---|---|
| C++ type: `MPI::Datatype` | |
| `MPI::BOOL` | `bool` |
| `MPI::COMPLEX` | `Complex<float>` |
| `MPI::DOUBLE_COMPLEX` | `Complex<double>` |
| `MPI::LONG_DOUBLE_COMPLEX` | `Complex<long double>` |

| Optional datatypes (Fortran) | | Fortran types |
|---|---|---|
| C type: `MPI_Datatype` | C++ type: `MPI::Datatype` | |
| Fortran type: `INTEGER` | | |
| `MPI_DOUBLE_COMPLEX` | `MPI::F_DOUBLE_COMPLEX` | `DOUBLE COMPLEX` |
| `MPI_INTEGER1` | `MPI::INTEGER1` | `INTEGER*1` |
| `MPI_INTEGER2` | `MPI::INTEGER2` | `INTEGER*8` |
| `MPI_INTEGER4` | `MPI::INTEGER4` | `INTEGER*4` |
| `MPI_INTEGER8` | `MPI::INTEGER8` | `INTEGER*8` |
| `MPI_INTEGER16` | | `INTEGER*16` |
| `MPI_REAL2` | `MPI::REAL2` | `REAL*2` |
| `MPI_REAL4` | `MPI::REAL4` | `REAL*4` |
| `MPI_REAL8` | `MPI::REAL8` | `REAL*8` |
| `MPI_REAL16` | | `REAL*16` |
| `MPI_COMPLEX4` | | `COMPLEX*4` |
| `MPI_COMPLEX8` | | `COMPLEX*8` |
| `MPI_COMPLEX16` | | `COMPLEX*16` |
| `MPI_COMPLEX32` | | `COMPLEX*32` |

### Datatypes for reduction functions (C and C++)

| C type: `MPI_Datatype` | C++ type: `MPI::Datatype` |
|---|---|
| Fortran type: `INTEGER` | |
| MPI_FLOAT_INT | MPI::FLOAT_INT |
| MPI_DOUBLE_INT | MPI::DOUBLE_INT |
| MPI_LONG_INT | MPI::LONG_INT |
| MPI_2INT | MPI::TWOINT |
| MPI_SHORT_INT | MPI::SHORT_INT |
| MPI_LONG_DOUBLE_INT | MPI::LONG_DOUBLE_INT |

### Datatypes for reduction functions (Fortran)

| C type: `MPI_Datatype` | C++ type: `MPI::Datatype` |
|---|---|
| Fortran type: `INTEGER` | |
| MPI_2REAL | MPI::TWOREAL |
| MPI_2DOUBLE_PRECISION | MPI::TWODOUBLE_PRECISION |
| MPI_2INTEGER | MPI::TWOINTEGER |

### Special datatypes for constructing derived datatypes

| C type: `MPI_Datatype` | C++ type: `MPI::Datatype` |
|---|---|
| Fortran type: `INTEGER` | |
| MPI_UB | MPI::UB |
| MPI_LB | MPI::LB |

### Reserved communicators

| C type: `MPI_Comm` | C++ type: `MPI::Intracomm` |
|---|---|
| Fortran type: `INTEGER` | |
| MPI_COMM_WORLD | MPI::COMM_WORLD |
| MPI_COMM_SELF | MPI::COMM_SELF |

### Results of communicator and group comparisons

| C type: `const int` (or unnamed `enum`) | C++ type: `const int` |
|---|---|
| Fortran type: `INTEGER` | (or unnamed `enum`) |
| MPI_IDENT | MPI::IDENT |
| MPI_CONGRUENT | MPI::CONGRUENT |
| MPI_SIMILAR | MPI::SIMILAR |
| MPI_UNEQUAL | MPI::UNEQUAL |

### Environmental inquiry keys

| C type: `const int` (or unnamed `enum`) | C++ type: `const int` |
|---|---|
| Fortran type: `INTEGER` | (or unnamed `enum`) |
| MPI_TAG_UB | MPI::TAG_UB |
| MPI_IO | MPI::IO |
| MPI_HOST | MPI::HOST |
| MPI_WTIME_IS_GLOBAL | MPI::WTIME_IS_GLOBAL |

**Unofficial Draft for Comment Only**