# MPI: A Message-Passing Interface Standard
## Version 3.0

(Draft, with MPI 3 Nonblocking Collectives

and new Fortran 2008 Interface)

Unofficial, for comment only

Message Passing Interface Forum

May 2, 2011

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  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 and sparse-group routines and more flexible and powerful one-sided operations. This *draft* contains the MPI Forum's current draft of nonblocking collective routines.

A new Fortran `mpi_f08` module is introduced to provide extended compile-time argument checking and buffer handling in nonblocking routines. The existing `mpi` module provides compile-time argument checking for the existing MPI-2.2 routine definitions. The use of `mpif.h` is strongly discouraged.

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, language dependent bindings follow: ticket230-B.

- The ISO C version of the function.

- The Fortran version of the same function used with USE mpi or INCLUDE 'mpif.h'

- The Fortran version used with USE mpi_f08.

- The C++ binding (which is deprecated).

Fortran in this document refers to Fortran 90 and higher; see Section 2.6. ticket230-B.

## 2.4   Semantic Terms

When discussing MPI procedures the following semantic terms are used.

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

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 with USE mpi or INCLUDE 'mpif.h', all handles have type INTEGER. In Fortran with USE mpi_f08, and in C and C++, a different handle type is defined for each category of objects.  With Fortran USE mpi_f08, the handles are defined as Fortran sequenced derived types that consist of only one element INTEGER :: MPI_VAL. The internal handle value is identical to the Fortran INTEGER value used in the mpi module and mpif.h. The type names are identical to the names in C, except that they are not case sensitive. For example:

```
TYPE MPI_Comm
  SEQUENCE
  INTEGER    :: MPI_VAL
END TYPE MPI_Comm
```

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.*)

> *Rationale.*   Due to the sequence attribute in the definition of handles in the mpi_f08 module, the new Fortran handles are associated with one numerical storage unit; i.e., they can be stored in every application data structure, including common blocks and sequence derived types. They may have the same C binding as the INTEGER handles of the mpi module. Since the integer values are equivalent, applications can easily convert MPI handles between all three supported Fortran methods. For example, an integer communicator handle COMM can be converted directly into an exactly equivalent mpi_f08 communicator handle named comm_f08 by comm_f08%MPI_VAL=COMM, and vice versa. (*End of rationale.*)

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.

ticket231-C.
ticket231-C.
ticket231-C.

ticket231-C.

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

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

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

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

*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. With the mpi_f08 module and

ticket244-P.

the C++ bindings, optional arguments through function overloading are used instead of MPI_STATUS_IGNORE, MPI_STATUSES_IGNORE, MPI_ERRCODES_IGNORE, and MPI_UNWEIGHTED. This is done by having two (or four in the case of two such independent arguments) bindings where one has the optional argument and one does not. The constants MPI_ARGV_NULL and MPI_ARGVS_NULL are not substituted by function overloading.

### 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 argument; e.g., tag is an integer-valued argument of point-to-point communication operations, with a special wild-card value, MPI_ANY_TAG. Such arguments will have a range of regular values, which is a proper subrange of the range of values of the corresponding basic type; special values (such as MPI_ANY_TAG) will be outside the regular range. The range of regular values, 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)

ticket234-F.

1  MPI_SUBARRAYS (Fortran only)
2  MPI_SUBARRAYS_SUPPORTED (Fortran only)
3  MPI_SUBARRAYS_NOT_SUPPORTED (Fortran only)
4  and their C++ counterparts where appropriate.
5  The constants that cannot be used in initialization expressions or assignments in For-
6  tran are:
7  MPI_BOTTOM
8  MPI_STATUS_IGNORE
9  MPI_STATUSES_IGNORE
10  MPI_ERRCODES_IGNORE
11  MPI_IN_PLACE
12  MPI_ARGV_NULL
13  MPI_ARGVS_NULL
14  MPI_UNWEIGHTED
15

16  *Advice to implementors.* In Fortran the implementation of these special constants
17  may require the use of language constructs that are outside the Fortran standard.
18  Using special values for the constants (e.g., by defining them through PARAMETER
19  statements) is not possible because an implementation cannot distinguish these val-
20  ues from legal data. Typically, these constants are implemented as predefined static
21  variables (e.g., a variable in an MPI-declared COMMON block), relying on the fact that
22  the target compiler passes data by address. Inside the subroutine, this address can
23  be extracted by some mechanism outside the Fortran standard (e.g., by Fortran ex-
24  tensions or by implementing the function in C). (*End of advice to implementors.*)
25

26  ### 2.5.5  Choice
27
28  MPI functions sometimes use arguments with a *choice* (or union) data type. Distinct calls to
29  the same routine may pass by reference actual arguments of different types. The mechanism
30  for providing such arguments will differ from language to language. For Fortran with the
31  include file `mpif.h` or the `mpi` module, the document uses <type> to represent a choice
32  variable; with the Fortran `mpi_f08` module, such arguments are declared with the new
33  Fortran syntax TYPE(*), DIMENSION(..); for C and C++, we use void *.
34
35  *Advice to implementors.* The implementor can freely choose how to implement choice
36  arguments in the `mpi` module, e.g., with a non-standard compiler-dependent method
37  that has the quality of the call mechanism in the implicit Fortran interfaces, or with
38  the method defined for the `mpi_f08` module. See details in Section 16.2.1 on page 544.
39  (*End of advice to implementors.*)
40
41
42  ### 2.5.6  Addresses
43
44  Some MPI procedures use *address* arguments that represent an absolute address in the
45  calling program. The datatype of such an argument is MPI_Aint in C, MPI::Aint in C++
46  and INTEGER (KIND=MPI_ADDRESS_KIND) in Fortran. These types must have the same
47  width and encode address values in the same manner such that address values in one
48  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.

ticket234-F. (line 30)
ticket234-F. (line 32)
ticket234-F. (line 33)

| 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_[ticket250-V.]FUNCTION |
| MPI_Delete_function | MPI_Comm_delete_attr_function |
| DELETE_FUNCTION | COMM_DELETE_ATTR_[ticket250-V.]FUNCTION |
| 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 or later; it means Fortran 2008 + TR 29113 and later if the `mpi_f08` module is used.

ticket234-F.

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. With `USE`
`mpi_f08`, this last argument is declared as `OPTIONAL`, except for user-defined callback func-
tions (e.g., `COMM_COPY_ATTR_FUNCTION`) and their predefined callbacks (e.g.,
`MPI_NULL_COPY_FN`). 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.

ticket239-K.
ticket250-V.
ticket239-K.
ticket239-K.

> *Advice to implementors.*   With the `mpi_f08` module, an MPI library may implement
> the `ierror` argument through function overloading instead of using the

OPTIONAL attribute. With function overloading, the branch is implemented at link-time; with the the OPTIONAL attribute, it is implemented at run-time. (*End of advice to implementors.*)

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 INTEGERs, or with the mpi_f08 module as a sequence derived type, see Section 2.5.1 on page 12. Binary-valued variables are of type LOGICAL.

Array arguments are indexed from one.

The MPI Fortran bindings are inconsistent with the Fortran standard in several respects. These inconsistencies, such as register optimization problems, have implications for user codes that are discussed in detail in Section 16.2.8.

### 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.

ticket231-C.

ticket230-B.
ticket230-B.
ticket230-B.

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

**Unofficial Draft for Comment Only**

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 ($\mathsf{myrank} = 1$) 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
```

ticket-248T.

```
MPI_Send(buf, count, datatype, dest, tag, comm, ierror)
    TYPE(*), DIMENSION(..)  ::  buf
    INTEGER, INTENT(IN) ::  count, dest, tag
    TYPE(MPI_Datatype), INTENT(IN) ::  datatype
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

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.

ticket243-O.   In Fortran with USE mpi or INCLUDE 'mpif.h', 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

ticket243-O.   code of the received message.

With Fortran USE mpi_f08, status is defined as the Fortran sequence derived type TYPE(MPI_Status) containing three public fields named MPI_SOURCE, MPI_TAG, and MPI_ERROR. TYPE(MPI_Status) may contain additional, implementation-specific fields. Thus, status%MPI_SOURCE, status%MPI_TAG and status%MPI_ERROR contain the source, tag, and error code of a received message respectively. Additionally, within both the mpi and the mpi_f08 modules, the constants MPI_STATUS_SIZE, MPI_SOURCE, MPI_TAG, MPI_ERROR, and TYPE(MPI_Status) are defined to allow conversion between both status representations.

> *Rationale.*   The Fortran TYPE(MPI_Status) is defined as a sequence derived type so that it can be used at any location where the status integer array representation can be used, e.g., in user defined sequence derived types or common blocks. (*End of rationale.*)

> *Rationale.*   It is allowed to have the same name (e.g., MPI_SOURCE) defined as a constant (e.g., Fortran parameter) and as a field of a derived type. (*End of rationale.*)

In C++, the status object is handled through the following methods:
{int MPI::Status::Get_source() const *(binding deprecated, see Section 15.2)* }

cannot be used when status is an IN argument. Note that in Fortran MPI_STATUS_IGNORE and MPI_STATUSES_IGNORE are objects like MPI_BOTTOM (not usable for initialization or assignment). See Section 2.5.4.

In general, this optimization can apply to all functions for which status or an array of statuses is an OUT argument. Note that this converts status into an INOUT argument. The functions that can be passed MPI_STATUS_IGNORE are all the various forms of MPI_RECV, MPI_TEST, and MPI_WAIT, as well as MPI_REQUEST_GET_STATUS. When an array is passed, as in the 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.

With the Fortran bindings through the `mpi_f08` module and the C++ bindings, MPI_STATUS_IGNORE and MPI_STATUSES_IGNORE do not exist. To allow an OUT or INOUT `TYPE(MPI_Status)` or MPI::Status argument to be ignored, all MPI `mpi_f08` or C++ bindings that have OUT or INOUT `TYPE(MPI_Status)` or MPI::Status parameters are overloaded with a second version that omits the OUT or INOUT `TYPE(MPI_Status)` or MPI::Status parameter.

**Example 3.1** The `mpi_f08` bindings for MPI_PROBE (the specific routine names, `...with/without_status`, are implementation dependent):

```
INTERFACE MPI_Probe
    SUBROUTINE MPI_Probe_with_status(source, tag, comm, status, ierror)
        INTEGER, INTENT(IN) ::  source, tag
        TYPE(MPI_Comm), INTENT(IN) ::  comm
        TYPE(MPI_Status), INTENT(OUT) ::  status
        INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
    END SUBROUTINE
    SUBROUTINE MPI_Probe_without_status(source, tag, comm, ierror)
        INTEGER, INTENT(IN) ::  source, tag
        TYPE(MPI_Comm), INTENT(IN) ::  comm
        INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
    END SUBROUTINE
END INTERFACE MPI_Probe
```

**Example 3.2** 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
```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17 ticket244-P.
18
19
20 ticket244-P.
21 ticket244-P.
22 ticket244-P.
23 ticket244-P.
24 ticket244-P.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

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

ticket-248T.

```
MPI_Irecv(buf, count, datatype, source, tag, comm, request, ierror)
    TYPE(*), DIMENSION(..)  ::  buf
    INTEGER, INTENT(IN) ::  count, source, tag
    TYPE(MPI_Datatype), INTENT(IN) ::  datatype
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Request), INTENT(OUT) ::  request
    INTEGER, OPTIONAL, INTENT(OUT) ::  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 Sections 16.2.2-16.2.12, especially in Sections 16.2.4 and 16.2.5 on pages 547-550 about *problems due to data copying and sequence association with subscript triplets* and *vector subscripts*, and in Sections 16.2.8 to 16.2.11 on pages 551 to 560 about *optimization problems, code movements and register optimization*, and *temporary* and *permanent data movements*. (*End of advice to users.*)

ticket238-J.
ticket238-J.
ticket236-H.
ticket238-J.

### 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

### 4.1.10 Duplicating a Datatype

ticket252-W.

MPI_TYPE_DUP(oldtype, newtype)

| IN | type | datatype (handle) |
|---|---|---|
| OUT | newtype | copy of oldtype (handle) |

ticket252-W.

```
int MPI_Type_dup(MPI_Datatype oldtype, MPI_Datatype *newtype)
```

```
MPI_TYPE_DUP(OLDTYPE, NEWTYPE, IERROR)
    INTEGER OLDTYPE, NEWTYPE, IERROR
```

```
MPI_Type_dup(oldtype, newtype, ierror)
    TYPE(MPI_Datatype), INTENT(IN) ::  oldtype
    TYPE(MPI_Datatype), INTENT(OUT) ::  newtype
    INTEGER, OPTIONAL, INTENT(OUT) ::  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 oldtype and any copied cached information, see Section 6.7.4 on page 276. 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 oldtype.

### 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

*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(user_fn, commute, op)

| IN | [ticket252-W.]user_fn | user defined function (function) |
|----|----|----|
| IN | commute | true if commutative; false otherwise. |
| OUT | op | operation (handle) |

```
int MPI_Op_create(MPI_User_function* user_fn, int commute, MPI_Op* op)
```

```
MPI_OP_CREATE( USER_FN, COMMUTE, OP, IERROR)
    EXTERNAL USER_FN
    LOGICAL COMMUTE
    INTEGER OP, IERROR
```

```
MPI_Op_create(user_fn, commute, op, ierror)
    EXTERNAL ::  user_fn
    LOGICAL, INTENT(IN) ::  commute
    TYPE(MPI_Op), INTENT(OUT) ::  op
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{void MPI::Op::Init(MPI::User_function* user_fn, 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 user_fn 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 user_fn appears below. To reduce problems with the choice buffer arguments, implicit interfaces are used with all three Fortran

support methods; the first one with the `mpi` module or `mpif.h`, the second one with
`mpi_f08`.

```
SUBROUTINE USER_FUNCTION(INVEC, INOUTVEC, LEN, DATATYPE)
    <type> INVEC(LEN), INOUTVEC(LEN)
    INTEGER LEN, DATATYPE
```

```
SUBROUTINE USER_FUNCTION(invec, inoutvec, len, datatype)
    <type> ::  invec(len), inoutvec(len)
    INTEGER ::  len
    TYPE(MPI_Datatype) ::  datatype
```

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 `user_fn`
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
$i = 0, \ldots, count - 1$, 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.

1
2
3 ticket252-W.
4
5 ticket252-W.
6 ticket230-B.
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25 ticket252-W.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

- obtain a key value (used to identify an attribute); the user specifies "callback" functions by which MPI informs the application when the communicator is destroyed or copied.

- store and retrieve the value of an attribute;

   *Advice to implementors.*   Caching and callback functions are only called synchronously, in response to explicit application requests. This avoid problems that result from repeated crossings between user and system space. (This synchronous calling rule is a general property of MPI.)

   The choice of key values is under control of MPI. This allows MPI to optimize its implementation of attribute sets. It also avoids conflict between independent modules caching information on the same communicators.

   A much smaller interface, consisting of just a callback facility, would allow the entire caching facility to be implemented by portable code. However, with the minimal callback interface, some form of table searching is implied by the need to handle arbitrary communicators. In contrast, the more complete interface defined here permits rapid access to attributes through the use of pointers in communicators (to find the attribute table) and cleverly chosen key values (to retrieve individual attributes). In light of the efficiency "hit" inherent in the minimal interface, the more complete interface defined here is seen to be superior. (*End of advice to implementors.*)

MPI provides the following services related to caching. They are all process local.

### 6.7.2   Communicators

Functions for caching on communicators are:

MPI_COMM_CREATE_KEYVAL(comm_copy_attr_fn, comm_delete_attr_fn, comm_keyval, extra_state)

| | | |
|---|---|---|
| IN | comm_copy_attr_fn | copy callback function for comm_keyval (function) |
| IN | comm_delete_attr_fn | delete callback function for comm_keyval (function) |
| OUT | comm_keyval | key value for future access (integer) |
| IN | extra_state | extra state for callback functions |

```
int MPI_Comm_create_keyval(MPI_Comm_copy_attr_function *comm_copy_attr_fn,
            MPI_Comm_delete_attr_function *comm_delete_attr_fn,
            int *comm_keyval, void *extra_state)
```

```
MPI_COMM_CREATE_KEYVAL(COMM_COPY_ATTR_FN, COMM_DELETE_ATTR_FN, COMM_KEYVAL,
            EXTRA_STATE, IERROR)
    EXTERNAL COMM_COPY_ATTR_FN, COMM_DELETE_ATTR_FN
    INTEGER COMM_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

ticket-248T.

```
MPI_Comm_create_keyval(comm_copy_attr_fn, comm_delete_attr_fn, comm_keyval,
            extra_state, ierror)
```

```
   PROCEDURE(MPI_Comm_copy_attr_function) ::  comm_copy_attr_fn
   PROCEDURE(MPI_Comm_delete_attr_function) ::  comm_delete_attr_fn
   INTEGER, INTENT(OUT) ::  comm_keyval
   INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state
   INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{static int MPI::Comm::Create_keyval(MPI::Comm::Copy_attr_function*
            comm_copy_attr_fn,
            MPI::Comm::Delete_attr_function* comm_delete_attr_fn,
            void* extra_state) *(binding deprecated, see Section 15.2)* }

Generates a new attribute key. Keys are locally unique in a process, and opaque to user, though they are explicitly stored in integers. Once allocated, the key value can be used to associate attributes and access them on any locally defined communicator.

This function replaces MPI_KEYVAL_CREATE, whose use is deprecated. The C binding is identical. The Fortran binding differs in that extra_state is an address-sized integer. Also, the copy and delete callback functions have Fortran bindings that are consistent with address-sized attributes.

The C callback functions are:

```
typedef int MPI_Comm_copy_attr_function(MPI_Comm oldcomm, int comm_keyval,
            void *extra_state, void *attribute_val_in,
            void *attribute_val_out, int *flag);
```

and

```
typedef int MPI_Comm_delete_attr_function(MPI_Comm comm, int comm_keyval,
            void *attribute_val, void *extra_state);
```

which are the same as the MPI-1.1 calls but with a new name. The old names are deprecated.

With the mpi module and mpif.h, the Fortran callback functions are:

```
SUBROUTINE COMM_COPY_ATTR_FUNCTION(OLDCOMM, COMM_KEYVAL, EXTRA_STATE,
            ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
   INTEGER OLDCOMM, COMM_KEYVAL, IERROR
   INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,
      ATTRIBUTE_VAL_OUT
   LOGICAL FLAG
```

and

```
SUBROUTINE COMM_DELETE_ATTR_FUNCTION(COMM, COMM_KEYVAL, ATTRIBUTE_VAL,
            EXTRA_STATE, IERROR)
   INTEGER COMM, COMM_KEYVAL, IERROR
   INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

With the mpi_f08 module, the Fortran callback functions are:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Comm_copy_attr_function(oldcomm, comm_keyval, extra_state,
  attribute_val_in, attribute_val_out, flag, ierror)
     TYPE(MPI_Comm) ::  oldcomm
     INTEGER ::  comm_keyval, ierror
     INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state, attribute_val_in,
     attribute_val_out
```

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 ticket230-B.
27
28 ticket250-V.
29
30
31
32
33
34
35
36 ticket250-V.
37
38
39
40 ticket230-B.
ticket-248T.
41
42
43
44
45
46
47
48

```
          LOGICAL flag
```

      and
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Comm_delete_attr_function(comm, comm_keyval,
  attribute_val, extra_state, ierror)
        TYPE(MPI_Comm) ::  comm
        INTEGER ::  comm_keyval, ierror
        INTEGER(KIND=MPI_ADDRESS_KIND) ::  attribute_val, extra_state
```

The C++ callbacks are:
```
{typedef int MPI::Comm::Copy_attr_function(const MPI::Comm& oldcomm,
             int comm_keyval, void* extra_state, void* attribute_val_in,
             void* attribute_val_out, bool& flag); (binding deprecated, see
             Section 15.2)}
```

and
```
{typedef int MPI::Comm::Delete_attr_function(MPI::Comm& comm,
             int comm_keyval, void* attribute_val, void* extra_state);
             (binding deprecated, see Section 15.2)}
```

The comm_copy_attr_fn function is invoked when a communicator is duplicated by MPI_COMM_DUP. comm_copy_attr_fn should be of type MPI_Comm_copy_attr_function. The copy callback function is invoked for each key value in oldcomm in arbitrary order. Each call to the copy callback is made with a key value and its corresponding attribute. If it returns flag = 0, then the attribute is deleted in the duplicated communicator. Otherwise (flag = 1), the new attribute value is set to the value returned in attribute_val_out. The function returns MPI_SUCCESS on success and an error code on failure (in which case MPI_COMM_DUP will fail).

The argument comm_copy_attr_fn may be specified as MPI_COMM_NULL_COPY_FN or MPI_COMM_DUP_FN from either C, C++, or Fortran. MPI_COMM_NULL_COPY_FN is a function that does nothing other than returning flag = 0 and MPI_SUCCESS. MPI_COMM_DUP_FN is a simple-minded copy function that sets flag = 1, returns the value of attribute_val_in in attribute_val_out, and returns MPI_SUCCESS. These replace the MPI-1 predefined callbacks MPI_NULL_COPY_FN and MPI_DUP_FN, whose use is deprecated.

> *Advice to users.*   Even though both formal arguments attribute_val_in and attribute_val_out are of type void *, their usage differs. The C copy function is passed by MPI in attribute_val_in the *value* of the attribute, and in attribute_val_out the *address* of the attribute, so as to allow the function to return the (new) attribute value. The use of type void * for both is to avoid messy type casts.
>
> A valid copy function is one that completely duplicates the information by making a full duplicate copy of the data structures implied by an attribute; another might just make another reference to that data structure, while using a reference-count mechanism. Other types of attributes might not copy at all (they might be specific to oldcomm only). (*End of advice to users.*)

> *Advice to implementors.*   A C interface should be assumed for copy and delete functions associated with key values created in C; a Fortran calling interface should be assumed for key values created in Fortran. (*End of advice to implementors.*)

Analogous to `comm_copy_attr_fn` is a callback deletion function, defined as follows. The `comm_delete_attr_fn` function is invoked when a communicator is deleted by MPI_COMM_FREE or when a call is made explicitly to MPI_COMM_DELETE_ATTR. `comm_delete_attr_fn` should be of type MPI_Comm_delete_attr_function.

This function is called by MPI_COMM_FREE, MPI_COMM_DELETE_ATTR, and MPI_COMM_SET_ATTR to do whatever is needed to remove an attribute. The function returns MPI_SUCCESS on success and an error code on failure (in which case MPI_COMM_FREE will fail).

The argument `comm_delete_attr_fn` may be specified as MPI_COMM_NULL_DELETE_FN from either C, C++, or Fortran. MPI_COMM_NULL_DELETE_FN is a function that does nothing, other than returning MPI_SUCCESS. MPI_COMM_NULL_DELETE_FN replaces MPI_NULL_DELETE_FN, whose use is deprecated.

If an attribute copy function or attribute delete function returns other than MPI_SUCCESS, then the call that caused it to be invoked (for example, MPI_COMM_FREE), is erroneous.

The special key value MPI_KEYVAL_INVALID is never returned by MPI_KEYVAL_CREATE. Therefore, it can be used for static initialization of key values.

> *Advice to implementors.* To be able to use the predefined C functions MPI_COMM_NULL_COPY_FN or MPI_COMM_DUP_FN as `comm_copy_attr_fn` argument and/or MPI_COMM_NULL_DELETE_FN as the `comm_delete_attr_fn` argument in a call to the C++ routine MPI::Comm::Create_keyval, this routine may be overloaded with 3 additional routines that accept the C functions as the first, the second, or both input arguments (instead of an argument that matches the C++ prototype). (*End of advice to implementors.*)

> *Advice to users.* If a user wants to write a "wrapper" routine that internally calls MPI::Comm::Create_keyval and `comm_copy_attr_fn` and/or `comm_delete_attr_fn` are arguments of this wrapper routine, and if this wrapper routine should be callable with both user-defined C++ copy and delete functions and with the predefined C functions, then the same overloading as described above in the advice to implementors may be necessary. (*End of advice to users.*)

ticket230-B.

> *Advice to implementors.* The predefined Fortran functions MPI_COMM_NULL_COPY_FN, MPI_COMM_DUP_FN, and MPI_COMM_NULL_DELETE_FN are defined in the `mpi` module (and `mpif.h`) and the `mpi_f08` module with a different interface. Each function can coexist twice with the same name in the same MPI library, one routine as an implicit interface outside of the `mpi` module, i.e., declared as `EXTERNAL`, and the other routine within `mpi_f08` declared with `CONTAINS`. (*End of advice to implementors.*)

> *Advice to users.* One should not pass the predefined Fortran functions MPI_COMM_NULL_COPY_FN, MPI_COMM_DUP_FN, and MPI_COMM_NULL_DELETE_FN from one application routine that uses the `mpi_f08` module to another application routine that uses the `mpi` module or `mpif.h`, and vice versa. (*End of advice to users.*)

```
MPI_WIN_CREATE_KEYVAL(WIN_COPY_ATTR_FN, WIN_DELETE_ATTR_FN, WIN_KEYVAL,
            EXTRA_STATE, IERROR)
    EXTERNAL WIN_COPY_ATTR_FN, WIN_DELETE_ATTR_FN
    INTEGER WIN_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

ticket-248T.

```
MPI_Win_create_keyval(win_copy_attr_fn, win_delete_attr_fn, win_keyval,
            extra_state, ierror)
    PROCEDURE(MPI_Win_copy_attr_function) ::  win_copy_attr_fn
    PROCEDURE(MPI_Win_delete_attr_function) ::  win_delete_attr_fn
    INTEGER, INTENT(OUT) ::  win_keyval
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{static int MPI::Win::Create_keyval(MPI::Win::Copy_attr_function*
            win_copy_attr_fn,
            MPI::Win::Delete_attr_function* win_delete_attr_fn,
            void* extra_state)
```
*(binding deprecated, see Section 15.2)* }

The argument win_copy_attr_fn may be specified as MPI_WIN_NULL_COPY_FN or MPI_WIN_DUP_FN from either C, C++, or Fortran.  MPI_WIN_NULL_COPY_FN is a function that does nothing other than returning flag = 0 and MPI_SUCCESS. MPI_WIN_DUP_FN is a simple-minded copy function that sets flag = 1, returns the value of attribute_val_in in attribute_val_out, and returns MPI_SUCCESS.

The argument win_delete_attr_fn may be specified as MPI_WIN_NULL_DELETE_FN from either C, C++, or Fortran.  MPI_WIN_NULL_DELETE_FN is a function that does nothing, other than returning MPI_SUCCESS.

The C callback functions are:

```
typedef int MPI_Win_copy_attr_function(MPI_Win oldwin, int win_keyval,
            void *extra_state, void *attribute_val_in,
            void *attribute_val_out, int *flag);
```

and

```
typedef int MPI_Win_delete_attr_function(MPI_Win win, int win_keyval,
            void *attribute_val, void *extra_state);
```

ticket230-B.

With the mpi module and mpif.h, the Fortran callback functions are:

ticket250-V.

```
SUBROUTINE WIN_COPY_ATTR_FUNCTION(OLDWIN, WIN_KEYVAL, EXTRA_STATE,
            ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
    INTEGER OLDWIN, WIN_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,
        ATTRIBUTE_VAL_OUT
    LOGICAL FLAG
```

and

ticket250-V.

```
SUBROUTINE WIN_DELETE_ATTR_FUNCTION(WIN, WIN_KEYVAL, ATTRIBUTE_VAL,
            EXTRA_STATE, IERROR)
    INTEGER WIN, WIN_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

ticket230-B.

With the mpi_f08 module, the Fortran callback functions are:

ticket-248T.

**Unofficial Draft for Comment Only**

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Win_copy_attr_function(oldwin, win_keyval, extra_state,
  attribute_val_in, attribute_val_out, flag, ierror)
      TYPE(MPI_Win) ::  oldwin
      INTEGER ::  win_keyval, ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state, attribute_val_in,
      attribute_val_out
      LOGICAL flag
```

ticket230-B.

ticket-248T.

    and

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Win_delete_attr_function(win, win_keyval, attribute_val,
  extra_state, ierror)
      TYPE(MPI_Win) ::  win
      INTEGER ::  win_keyval, ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  attribute_val, extra_state
```

The C++ callbacks are:

{typedef int MPI::Win::Copy_attr_function(const MPI::Win& oldwin,
                  int win_keyval, void* extra_state, void* attribute_val_in,
                  void* attribute_val_out, bool& flag); *(binding deprecated, see
                  Section 15.2)*}

    and

{typedef int MPI::Win::Delete_attr_function(MPI::Win& win, int win_keyval,
                  void* attribute_val, void* extra_state); *(binding deprecated, see
                  Section 15.2)*}

If an attribute copy function or attribute delete function returns other than
MPI_SUCCESS, then the call that caused it to be invoked (for example, MPI_WIN_FREE), is
erroneous.


MPI_WIN_FREE_KEYVAL(win_keyval)

  INOUT    win_keyval                          key value (integer)


```
int MPI_Win_free_keyval(int *win_keyval)
```

```
MPI_WIN_FREE_KEYVAL(WIN_KEYVAL, IERROR)
    INTEGER WIN_KEYVAL, IERROR
```

ticket-248T.

```
MPI_Win_free_keyval(win_keyval, ierror)
    INTEGER, INTENT(INOUT) ::  win_keyval
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{static void MPI::Win::Free_keyval(int& win_keyval)*(binding deprecated, see
                  Section 15.2)* }

MPI_WIN_DELETE_ATTR(win, win_keyval)

  INOUT     win                                    window from which the attribute is deleted (handle)

  IN        win_keyval                             key value (integer)

```
int MPI_Win_delete_attr(MPI_Win win, int win_keyval)
```

```
MPI_WIN_DELETE_ATTR(WIN, WIN_KEYVAL, IERROR)
    INTEGER WIN, WIN_KEYVAL, IERROR
```

ticket-248T.
```
MPI_Win_delete_attr(win, win_keyval, ierror)
    TYPE(MPI_Win), INTENT(IN) ::  win
    INTEGER, INTENT(IN) ::  win_keyval
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{void MPI::Win::Delete_attr(int win_keyval) *(binding deprecated, see Section 15.2)*
               }


### 6.7.4   Datatypes

The new functions for caching on datatypes are:


MPI_TYPE_CREATE_KEYVAL(type_copy_attr_fn, type_delete_attr_fn, type_keyval, extra_state)

  IN        type_copy_attr_fn          copy callback function for type_keyval (function)

  IN        type_delete_attr_fn        delete callback function for type_keyval (function)

  OUT       type_keyval                key value for future access (integer)

  IN        extra_state                extra state for callback functions

```
int MPI_Type_create_keyval(MPI_Type_copy_attr_function *type_copy_attr_fn,
             MPI_Type_delete_attr_function *type_delete_attr_fn,
             int *type_keyval, void *extra_state)
```

```
MPI_TYPE_CREATE_KEYVAL(TYPE_COPY_ATTR_FN, TYPE_DELETE_ATTR_FN, TYPE_KEYVAL,
             EXTRA_STATE, IERROR)
    EXTERNAL TYPE_COPY_ATTR_FN, TYPE_DELETE_ATTR_FN
    INTEGER TYPE_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

ticket-248T.
```
MPI_Type_create_keyval(type_copy_attr_fn, type_delete_attr_fn, type_keyval,
             extra_state, ierror)
    PROCEDURE(MPI_Type_copy_attr_function) ::  type_copy_attr_fn
    PROCEDURE(MPI_Type_delete_attr_function) ::  type_delete_attr_fn
    INTEGER, INTENT(OUT) ::  type_keyval
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{static int MPI::Datatype::Create_keyval(MPI::Datatype::Copy_attr_function*
           type_copy_attr_fn, MPI::Datatype::Delete_attr_function*
           type_delete_attr_fn, void* extra_state)
```
*(binding deprecated, see Section 15.2)* }

The argument type_copy_attr_fn may be specified as MPI_TYPE_NULL_COPY_FN or MPI_TYPE_DUP_FN from either C, C++, or Fortran. MPI_TYPE_NULL_COPY_FN is a function that does nothing other than returning flag = 0 and MPI_SUCCESS. MPI_TYPE_DUP_FN is a simple-minded copy function that sets flag = 1, returns the value of attribute_val_in in attribute_val_out, and returns MPI_SUCCESS.

The argument type_delete_attr_fn may be specified as MPI_TYPE_NULL_DELETE_FN from either C, C++, or Fortran. MPI_TYPE_NULL_DELETE_FN is a function that does nothing, other than returning MPI_SUCCESS.

The C callback functions are:
```
typedef int MPI_Type_copy_attr_function(MPI_Datatype oldtype,
           int type_keyval, void *extra_state, void *attribute_val_in,
           void *attribute_val_out, int *flag);
```

and
```
typedef int MPI_Type_delete_attr_function(MPI_Datatype datatype,
           int type_keyval, void *attribute_val, void *extra_state);
```

With the mpi module and mpif.h, the Fortran callback functions are:
```
SUBROUTINE TYPE_COPY_ATTR_FUNCTION(OLDTYPE, TYPE_KEYVAL, EXTRA_STATE,
           ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
    INTEGER OLDTYPE, TYPE_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE,
        ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT
    LOGICAL FLAG
```

and
```
SUBROUTINE TYPE_DELETE_ATTR_FUNCTION(DATATYPE, TYPE_KEYVAL, ATTRIBUTE_VAL,
           EXTRA_STATE, IERROR)
    INTEGER DATATYPE, TYPE_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

With the mpi_f08 module, the Fortran callback functions are:
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Type_copy_attr_function(oldtype, type_keyval, extra_state,
  attribute_val_in, attribute_val_out, flag, ierror)
     TYPE(MPI_Datatype) ::  oldtype
     INTEGER ::  type_keyval, ierror
     INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state, attribute_val_in,
     attribute_val_out
     LOGICAL flag
```

and
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Type_delete_attr_function(datatype, type_keyval,
  attribute_val, extra_state, ierror)
```

```
TYPE(MPI_Datatype) ::  datatype
INTEGER ::  type_keyval, ierror
INTEGER(KIND=MPI_ADDRESS_KIND) ::  attribute_val, extra_state
```

The C++ callbacks are:

```
{typedef int
            MPI::Datatype::Copy_attr_function(const MPI::Datatype& oldtype,
            int type_keyval, void* extra_state,
            const void* attribute_val_in, void* attribute_val_out,
            bool& flag); (binding deprecated, see Section 15.2)}
```

and

ticket252-W.
```
{typedef int MPI::Datatype::Delete_attr_function(MPI::Datatype& datatype,
            int type_keyval, void* attribute_val, void* extra_state);
            (binding deprecated, see Section 15.2)}
```

If an attribute copy function or attribute delete function returns other than MPI_SUCCESS, then the call that caused it to be invoked (for example, MPI_TYPE_FREE), is erroneous.


MPI_TYPE_FREE_KEYVAL(type_keyval)

  INOUT     type_keyval                    key value (integer)


```
int MPI_Type_free_keyval(int *type_keyval)
```

```
MPI_TYPE_FREE_KEYVAL(TYPE_KEYVAL, IERROR)
    INTEGER TYPE_KEYVAL, IERROR
```

ticket-248T.
```
MPI_Type_free_keyval(type_keyval, ierror)
    INTEGER, INTENT(INOUT) ::  type_keyval
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{static void MPI::Datatype::Free_keyval(int& type_keyval) (binding deprecated,
            see Section 15.2) }
```


ticket252-W.
MPI_TYPE_SET_ATTR(datatype, type_keyval, attribute_val)

  INOUT     [ticket252-W.]datatype         datatype to which attribute will be attached (handle)

  IN        type_keyval                    key value (integer)

  IN        attribute_val                  attribute value


ticket252-W.
```
int MPI_Type_set_attr(MPI_Datatype datatype, int type_keyval,
            void *attribute_val)
```

ticket252-W.
ticket252-W.
```
MPI_TYPE_SET_ATTR(DATATYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, IERROR)
    INTEGER DATATYPE, TYPE_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
```

ticket-248T.
```
MPI_Type_set_attr(datatype, type_keyval, attribute_val, ierror)
```

MPI_DIST_GRAPH_CREATE_ADJACENT(comm_old, indegree, sources, sourceweights, out-
degree, destinations, destweights, info, reorder, comm_dist_graph)

| | | |
|---|---|---|
| IN | comm_old | input communicator (handle) |
| IN | indegree | size of sources and sourceweights arrays (non-negative integer) |
| IN | sources | ranks of processes for which the calling process is a destination (array of non-negative integers) |
| IN | sourceweights | weights of the edges into the calling process (array of non-negative integers) |
| IN | outdegree | size of destinations and destweights arrays (non-negative integer) |
| IN | destinations | ranks of processes for which the calling process is a source (array of non-negative integers) |
| IN | destweights | weights of the edges out of the calling process (array of non-negative integers) |
| IN | info | hints on optimization and interpretation of weights (handle) |
| IN | reorder | the ranks may be reordered (true) or not (false) (logical) |
| OUT | comm_dist_graph | communicator with distributed graph topology (handle) |

```
int MPI_Dist_graph_create_adjacent(MPI_Comm comm_old, int indegree,
            int sources[], int sourceweights[], int outdegree,
            int destinations[], int destweights[], MPI_Info info,
            int reorder, MPI_Comm *comm_dist_graph)
```

```
MPI_DIST_GRAPH_CREATE_ADJACENT(COMM_OLD, INDEGREE, SOURCES, SOURCEWEIGHTS,
            OUTDEGREE, DESTINATIONS, DESTWEIGHTS, INFO, REORDER,
            COMM_DIST_GRAPH, IERROR)
    INTEGER COMM_OLD, INDEGREE, SOURCES(*), SOURCEWEIGHTS(*), OUTDEGREE,
        DESTINATIONS(*), DESTWEIGHTS(*), INFO, COMM_DIST_GRAPH, IERROR
    LOGICAL REORDER
```

ticket-248T.

```
MPI_Dist_graph_create_adjacent(comm_old, indegree, sources, sourceweights,
            outdegree, destinations, destweights, info, reorder,
            comm_dist_graph, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm_old
    INTEGER, INTENT(IN) ::  indegree, sources(*), outdegree,
    destinations(*)
    INTEGER, INTENT(IN) ::  sourceweights(*) !  optional by overloading
    INTEGER, INTENT(IN) ::  destweights(*) !  optional by overloading
    TYPE(MPI_Info), INTENT(IN) ::  info
    LOGICAL, INTENT(IN) ::  reorder
    TYPE(MPI_Comm), INTENT(OUT) ::  comm_dist_graph
```

**Unofficial Draft for Comment Only**

```
      INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{MPI::Distgraphcomm MPI::Intracomm::Dist_graph_create_adjacent(int
            indegree, const int sources[], const int sourceweights[],
            int outdegree, const int destinations[],
            const int destweights[], const MPI::Info& info, bool reorder)
            const(binding deprecated, see Section 15.2) }
```

```
{MPI::Distgraphcomm
            MPI::Intracomm::Dist_graph_create_adjacent(int indegree,
            const int sources[], int outdegree, const int destinations[],
            const MPI::Info& info, bool reorder) const(binding deprecated, see
            Section 15.2) }
```

MPI_DIST_GRAPH_CREATE_ADJACENT returns a handle to a new communicator to which the distributed graph topology information is attached. Each process passes all information about the edges to its neighbors in the virtual distributed graph topology. The calling processes must ensure that each edge of the graph is described in the source and in the destination process with the same weights. If there are multiple edges for a given (source,dest) pair, then the sequence of the weights of these edges does not matter. The complete communication topology is the combination of all edges shown in the sources arrays of all processes in comm_old, which must be identical to the combination of all edges shown in the destinations arrays. Source and destination ranks must be process ranks of comm_old. This allows a fully distributed specification of the communication graph. Isolated processes (i.e., processes with no outgoing or incoming edges, that is, processes that have specified indegree and outdegree as zero and that thus do not occur as source or destination rank in the graph specification) are allowed.

The call creates a new communicator comm_dist_graph of distributed graph topology type to which topology information has been attached. The number of processes in comm_dist_graph is identical to the number of processes in comm_old. The call to MPI_DIST_GRAPH_CREATE_ADJACENT is collective.

Weights are specified as non-negative integers and can be used to influence the process remapping strategy and other internal MPI optimizations. For instance, approximate count arguments of later communication calls along specific edges could be used as their edge weights. Multiplicity of edges can likewise indicate more intense communication between pairs of processes. However, the exact meaning of edge weights is not specified by the MPI standard and is left to the implementation. In C or in the Fortran mpi module or mpif.h include file, an application can supply the special value MPI_UNWEIGHTED for the weight array to indicate that all edges have the same (effectively no) weight. In the Fortran mpi_f08 module or in C++, this constant does not exist and the weight arguments may be omitted from the argument list.[1] It is erroneous to supply MPI_UNWEIGHTED, or in mpi_f08 or C++ to omit the weight arrays, for some but not all processes of comm_old. Note that MPI_UNWEIGHTED is not a special weight value; rather it is a special value for the total array argument. In C, one would expect it to be NULL. In Fortran, MPI_UNWEIGHTED is an object like MPI_BOTTOM (not usable for initialization or assignment). See Section 2.5.4.

The meaning of the info and reorder arguments is defined in the description of the following routine.

ticket244-P.
ticket244-P.
ticket244-P.
ticket244-P.
ticket244-P.
ticket244-P.

---

[1] The Fortran generic interface must contain four specific routines with none, both, or only the sourceweight or destweight array is omitted, although in the C++ binding only two of them are provided.

| node | | exchange neighbors(1) | shuffle neighbors(2) | unshuffle neighbors(3) |
|---|---|---|---|---|
| 0 | (000) | 1 | 0 | 0 |
| 1 | (001) | 0 | 2 | 4 |
| 2 | (010) | 3 | 4 | 1 |
| 3 | (011) | 2 | 6 | 5 |
| 4 | (100) | 5 | 1 | 2 |
| 5 | (101) | 4 | 3 | 6 |
| 6 | (110) | 7 | 5 | 3 |
| 7 | (111) | 6 | 7 | 7 |

Suppose that the communicator `comm` has this topology associated with it. The following code fragment cycles through the three types of neighbors and performs an appropriate permutation for each.

```
C  assume: each process has stored a real number A.
C  extract neighborhood information
      CALL MPI_COMM_RANK(comm, myrank, ierr)
      CALL MPI_GRAPH_NEIGHBORS(comm, myrank, 3, neighbors, ierr)
C  perform exchange permutation
      CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, neighbors(1), 0,
     +     neighbors(1), 0, comm, status, ierr)
C  perform shuffle permutation
      CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, neighbors(2), 0,
     +     neighbors(3), 0, comm, status, ierr)
C  perform unshuffle permutation
      CALL MPI_SENDRECV_REPLACE(A, 1, MPI_REAL, neighbors(3), 0,
     +     neighbors(2), 0, comm, status, ierr)
```

MPI_DIST_GRAPH_NEIGHBORS_COUNT and MPI_DIST_GRAPH_NEIGHBORS provide adjacency information for a distributed graph topology.

MPI_DIST_GRAPH_NEIGHBORS_COUNT(comm, indegree, outdegree, weighted)

| IN | comm | communicator with distributed graph topology (handle) |
|---|---|---|
| OUT | indegree | number of edges into this process (non-negative integer) |
| OUT | outdegree | number of edges out of this process (non-negative integer) |
| OUT | weighted | false if MPI_UNWEIGHTED was supplied or the weights were omitted during creation, true otherwise (logical) |

ticket244-P.

```
int MPI_Dist_graph_neighbors_count(MPI_Comm comm, int *indegree,
            int *outdegree, int *weighted)
```

```
MPI_DIST_GRAPH_NEIGHBORS_COUNT(COMM, INDEGREE, OUTDEGREE, WEIGHTED, IERROR)
```

```
      INTEGER COMM, INDEGREE, OUTDEGREE, IERROR
      LOGICAL WEIGHTED

MPI_Dist_graph_neighbors_count(comm, indegree, outdegree, weighted, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    INTEGER, INTENT(OUT) ::  indegree, outdegree
    LOGICAL, INTENT(OUT) ::  weighted
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{void MPI::Distgraphcomm::Get_dist_neighbors_count(int rank,
            int indegree[], int outdegree[], bool& weighted) const*(binding
            deprecated, see Section 15.2)* }


MPI_DIST_GRAPH_NEIGHBORS(comm, maxindegree, sources, sourceweights, maxoutdegree, destinations, destweights)

| | | |
|---|---|---|
| IN | comm | communicator with distributed graph topology (handle) |
| IN | maxindegree | size of sources and sourceweights arrays (non-negative integer) |
| OUT | sources | processes for which the calling process is a destination (array of non-negative integers) |
| OUT | sourceweights | weights of the edges into the calling process (array of non-negative integers) |
| IN | maxoutdegree | size of destinations and destweights arrays (non-negative integer) |
| OUT | destinations | processes for which the calling process is a source (array of non-negative integers) |
| OUT | destweights | weights of the edges out of the calling process (array of non-negative integers) |

```
int MPI_Dist_graph_neighbors(MPI_Comm comm, int maxindegree, int sources[],
            int sourceweights[], int maxoutdegree, int destinations[],
            int destweights[])

MPI_DIST_GRAPH_NEIGHBORS(COMM, MAXINDEGREE, SOURCES, SOURCEWEIGHTS,
            MAXOUTDEGREE, DESTINATIONS, DESTWEIGHTS, IERROR)
    INTEGER COMM, MAXINDEGREE, SOURCES(*), SOURCEWEIGHTS(*), MAXOUTDEGREE,
        DESTINATIONS(*), DESTWEIGHTS(*), IERROR
```

```
MPI_Dist_graph_neighbors(comm, maxindegree, sources, sourceweights,
            maxoutdegree, destinations, destweights, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    INTEGER, INTENT(IN) ::  maxindegree, maxoutdegree
    INTEGER, INTENT(OUT) ::  sources(*), destinations(*)
    INTEGER, INTENT(OUT) ::  sourceweights(*), destweights(*) !  optional
    by overloading
```

<span style="color:red">INTEGER, OPTIONAL, INTENT(OUT) :: ierror</span>

{void MPI::Distgraphcomm::Get_dist_neighbors(int maxindegree,
            int sources[], int sourceweights[], int maxoutdegree,
            int destinations[], int destweights[])*(binding deprecated, see
            Section 15.2)* }

These calls are local. The number of edges into and out of the process returned by MPI_DIST_GRAPH_NEIGHBORS_COUNT are the total number of such edges given in the call to MPI_DIST_GRAPH_CREATE_ADJACENT or MPI_DIST_GRAPH_CREATE (potentially by processes other than the calling process in the case of MPI_DIST_GRAPH_CREATE). Multiply defined edges are all counted and returned by MPI_DIST_GRAPH_NEIGHBORS in some order. If MPI_UNWEIGHTED is supplied <span style="color:red">(in C or the `mpi` module or `mpif.h`) or the argument is omitted (in `mpi_f08`)</span> for sourceweights or destweights or <span style="color:red">both,</span>[2] or if MPI_UNWEIGHTED was supplied <span style="color:red">or the weights were omitted</span> during the construction of the graph then no weight information is returned in that array or those arrays. The only requirement on the order of values in sources and destinations is that two calls to the routine with same input argument comm will return the same sequence of edges. If maxindegree or maxoutdegree is smaller than the numbers returned by MPI_DIST_GRAPH_NEIGHBOR_COUNT, then only the first part of the full list is returned. Note, that the order of returned edges does need not to be identical to the order that was provided in the creation of comm for the case that MPI_DIST_GRAPH_CREATE_ADJACENT was used.

> *Advice to implementors.* Since the query calls are defined to be local, each process needs to store the list of its neighbors with incoming and outgoing edges. Communication is required at the collective MPI_DIST_GRAPH_CREATE call in order to compute the neighbor lists for each process from the distributed graph specification. (*End of advice to implementors.*)

### 7.5.6 Cartesian Shift Coordinates

If the process topology is a Cartesian structure, an MPI_SENDRECV operation is likely to be used along a coordinate direction to perform a shift of data. As input, MPI_SENDRECV takes the rank of a source process for the receive, and the rank of a destination process for the send. If the function MPI_CART_SHIFT is called for a Cartesian process group, it provides the calling process with the above identifiers, which then can be passed to MPI_SENDRECV. The user specifies the coordinate direction and the size of the step (positive or negative). The function is local.

---

[2]The generic Fortran interface must contain 4 specific routines.

<span style="color:red">ticket244-P.</span>
ticket244-P.
ticket244-P.

# Chapter 8

# **MPI** Environmental Management

This chapter discusses routines for getting and, where appropriate, setting various parameters that relate to the MPI implementation and the execution environment (such as error handling). The procedures for entering and leaving the MPI execution environment are also described here.

## 8.1  Implementation Information

### 8.1.1  Version Inquiries

In order to cope with changes to the MPI Standard, there are both compile-time and run-time ways to determine which version of the standard is in use in the environment one is using.

The "version" will be represented by two separate integers, for the version and subversion: In C and C++,

```
#define MPI_VERSION    2
#define MPI_SUBVERSION 2
```

in Fortran,

```
INTEGER [ticket240-L.]:: MPI_VERSION, MPI_SUBVERSION
PARAMETER (MPI_VERSION    = 2)
PARAMETER (MPI_SUBVERSION = 2)
```

For runtime determination,

---

MPI_GET_VERSION( version, subversion )

| | | |
|---|---|---|
| OUT | version | version number (integer) |
| OUT | subversion | subversion number (integer) |

---

```
int MPI_Get_version(int *version, int *subversion)
```

```
MPI_GET_VERSION(VERSION, SUBVERSION, IERROR)
    INTEGER VERSION, SUBVERSION, IERROR
```

```
MPI_Get_version(version, subversion, ierror)
```

{void MPI::Get_processor_name(char* name, int& resultlen)*(binding deprecated, see Section 15.2)* }

This routine returns the name of the processor on which it was called at the moment of the call. The name is a character string for maximum flexibility. From this value it must be possible to identify a specific piece of hardware; possible values include "processor 9 in rack 4 of mpp.cs.org" and "231" (where 231 is the actual processor number in the running homogeneous system). The argument name must represent storage that is at least MPI_MAX_PROCESSOR_NAME characters long. MPI_GET_PROCESSOR_NAME may write up to this many characters into name.

The number of characters actually written is returned in the output argument, resultlen. In C, a null character is additionally stored at name[resultlen]. The resultlen cannot be larger then MPI_MAX_PROCESSOR_NAME-1. In Fortran, name is padded on the right with blank characters. The resultlen cannot be larger then MPI_MAX_PROCESSOR_NAME.

> *Rationale.* This function allows MPI implementations that do process migration to return the current processor. Note that nothing in MPI *requires* or defines process migration; this definition of MPI_GET_PROCESSOR_NAME simply allows such an implementation. (*End of rationale.*)

> *Advice to users.* The user must provide at least MPI_MAX_PROCESSOR_NAME space to write the processor name — processor names can be this long. The user should examine the output argument, resultlen, to determine the actual length of the name. (*End of advice to users.*)

The constant MPI_BSEND_OVERHEAD provides an upper bound on the fixed overhead per message buffered by a call to MPI_BSEND (see Section 3.6.1).

## 8.2   Memory Allocation

In some systems, message-passing and remote-memory-access (RMA) operations run faster when accessing specially allocated memory (e.g., memory that is shared by the other processes in the communicating group on an SMP). MPI provides a mechanism for allocating and freeing such special memory. The use of such memory for message-passing or RMA is not mandatory, and this memory can be used without restrictions as any other dynamically allocated memory. However, implementations may restrict the use of the MPI_WIN_LOCK and MPI_WIN_UNLOCK functions to windows allocated in such memory (see Section 11.4.3.)

MPI_ALLOC_MEM(size, info, baseptr)

| | | |
|---|---|---|
| IN | size | size of memory segment in bytes (non-negative integer) |
| IN | info | info argument (handle) |
| OUT | baseptr | pointer to beginning of memory segment allocated |

```
int MPI_Alloc_mem(MPI_Aint size, MPI_Info info, void *baseptr)
```

MPI_ALLOC_MEM(SIZE, INFO, BASEPTR, IERROR)

```
    INTEGER INFO, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) SIZE, BASEPTR
```

ticket-248T.

```
MPI_Alloc_mem(size, info, baseptr, ierror)
    USE, INTRINSIC ::  ISO_C_BINDING
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  size
    TYPE(MPI_Info), INTENT(IN) ::  info
    TYPE(C_PTR), INTENT(OUT) ::  baseptr
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{void* MPI::Alloc_mem(MPI::Aint size, const MPI::Info& info)*(binding deprecated, see Section 15.2)* }

The info argument can be used to provide directives that control the desired location of the allocated memory. Such a directive does not affect the semantics of the call. Valid info values are implementation-dependent; a null directive value of info = MPI_INFO_NULL is always valid.

The function MPI_ALLOC_MEM may return an error code of class MPI_ERR_NO_MEM to indicate it failed because memory is exhausted.

MPI_FREE_MEM(base)

| IN | base | initial address of memory segment allocated by MPI_ALLOC_MEM (choice) |
|---|---|---|

```
int MPI_Free_mem(void *base)
```

```
MPI_FREE_MEM(BASE, IERROR)
    <type> BASE(*)
    INTEGER IERROR
```

ticket-248T.

```
MPI_Free_mem(base, ierror)
    USE, INTRINSIC ::  ISO_C_BINDING
    TYPE(C_PTR), INTENT(IN) ::  base
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{void MPI::Free_mem(void *base)*(binding deprecated, see Section 15.2)* }

The function MPI_FREE_MEM may return an error code of class MPI_ERR_BASE to indicate an invalid base argument.

> *Rationale.* The C and C++ bindings of MPI_ALLOC_MEM and MPI_FREE_MEM are similar to the bindings for the `malloc` and `free` C library calls: a call to MPI_Alloc_mem(..., &base) should be paired with a call to MPI_Free_mem(base) (one less level of indirection). Both arguments are declared to be of same type void* so as to facilitate type casting. The Fortran binding is consistent with the C and C++ bindings: the Fortran MPI_ALLOC_MEM call returns in baseptr the (integer valued) address of the allocated memory. The base argument of MPI_FREE_MEM is a choice argument, which passes (a reference to) the variable stored at that location. (*End of rationale.*)

> *Advice to implementors.*     If MPI_ALLOC_MEM allocates special memory, then a
> design similar to the design of C `malloc` and `free` functions has to be used, in order
> to find out the size of a memory segment, when the segment is freed. If no special
> memory is used, MPI_ALLOC_MEM simply invokes `malloc`, and MPI_FREE_MEM
> invokes `free`.
>
> A call to MPI_ALLOC_MEM can be used in shared memory systems to allocate mem-
> ory in a shared memory segment. (*End of advice to implementors.*)

**Example 8.1**
Example of use of MPI_ALLOC_MEM, in Fortran with pointer support. We assume
4-byte REALs, and assume that pointers are address-sized.

```
REAL A
POINTER (P, A(100,100))   ! no memory is allocated
CALL MPI_ALLOC_MEM(4*100*100, MPI_INFO_NULL, P, IERR)
! memory is allocated
...
A(3,5) = 2.71;
...
CALL MPI_FREE_MEM(A, IERR) ! memory is freed
```

Since standard Fortran does not support (C-like) pointers, this code is not Fortran 77
or Fortran 90 code. Some compilers (in particular, at the time of writing, g77 and Fortran
compilers for Intel) do not support this code.

**Example 8.2**  Same example, in C

```
float  (* f)[100][100] ;
/* no memory is allocated */
MPI_Alloc_mem(sizeof(float)*100*100, MPI_INFO_NULL, &f);
/* memory allocated */
...
(*f)[5][3] = 2.71;
...
MPI_Free_mem(f);
```

## 8.3   Error Handling

An MPI implementation cannot or may choose not to handle some errors that occur during
MPI calls. These can include errors that generate exceptions or traps, such as floating point
errors or access violations. The set of errors that are handled by MPI is implementation-
dependent. Each such error generates an **MPI exception**.

The above text takes precedence over any text on error handling within this document.
Specifically, text that states that errors *will* be handled should be read as *may* be handled.

A user can associate error handlers to three types of objects: communicators, windows,
and files. The specified error handling routine will be used for any MPI exception that occurs
during a call to MPI for the respective object. MPI calls that are not related to any objects
are considered to be attached to the communicator MPI_COMM_WORLD. The attachment

The MPI function MPI_ERRHANDLER_FREE can be used to free an error handler that was created by a call to MPI_XXX_CREATE_ERRHANDLER.

MPI_{COMM,WIN,FILE}_GET_ERRHANDLER behave as if a new error handler object is created. That is, once the error handler is no longer needed, MPI_ERRHANDLER_FREE should be called with the error handler returned from MPI_ERRHANDLER_GET or MPI_{COMM,WIN,FILE}_GET_ERRHANDLER to mark the error handler for deallocation. This provides behavior similar to that of MPI_COMM_GROUP and MPI_GROUP_FREE.

> *Advice to implementors.*   High-quality implementation should raise an error when an error handler that was created by a call to MPI_XXX_CREATE_ERRHANDLER is attached to an object of the wrong type with a call to MPI_YYY_SET_ERRHANDLER. To do so, it is necessary to maintain, with each error handler, information on the typedef of the associated user function. (*End of advice to implementors.*)

The syntax for these calls is given below.

## 8.3.1   Error Handlers for Communicators

<sub>ticket252-W.</sub> MPI_COMM_CREATE_ERRHANDLER(comm_errhandler_fn, errhandler)

| IN | [ticket252-W.]comm_errhandler_fn | user defined error handling procedure (function) |
| OUT | errhandler | MPI error handler (handle) |

```
int MPI_Comm_create_errhandler(MPI_Comm_errhandler_function
               *comm_errhandler_fn, MPI_Errhandler *errhandler)
```

ticket252-W.
```
MPI_COMM_CREATE_ERRHANDLER(COMM_ERRHANDLER_FN, ERRHANDLER, IERROR)
```
ticket252-W.
```
    EXTERNAL COMM_ERRHANDLER_FN
    INTEGER ERRHANDLER, IERROR
```
ticket-248T.
```
MPI_Comm_create_errhandler(comm_errhandler_fn, errhandler, ierror)
    PROCEDURE(MPI_Comm_errhandler_function) ::  comm_errhandler_fn
    TYPE(MPI_Errhandler), INTENT(OUT) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{static MPI::Errhandler
          MPI::Comm::Create_errhandler(MPI::Comm::Errhandler_function*
```
ticket252-W.
```
          comm_errhandler_fn)
```
*(binding deprecated, see Section 15.2)* }

Creates an error handler that can be attached to communicators. This function is identical to MPI_ERRHANDLER_CREATE, whose use is deprecated.

The user routine should be, in C, a function of type MPI_Comm_errhandler_function, which is defined as
```
typedef void MPI_Comm_errhandler_function(MPI_Comm *, int *, ...);
```

The first argument is the communicator in use. The second is the error code to be returned by the MPI routine that raised the error. If the routine would have returned MPI_ERR_IN_STATUS, it is the error code returned in the status for the request that caused the error handler to be invoked. The remaining arguments are "stdargs" arguments whose

number and meaning is implementation-dependent. An implementation should clearly document these arguments. Addresses are used so that the handler may be written in Fortran. This typedef replaces MPI_Handler_function, whose use is deprecated.

With the Fortran mpi module and mpif.h, the user routine should be of the form:

```
SUBROUTINE COMM_ERRHANDLER_FUNCTION(COMM, ERROR_CODE)
    INTEGER COMM, ERROR_CODE
```

With the Fortran mpi_f08 module, the user routine should be of the form:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Comm_errhandler_function(comm, error_code)
      TYPE(MPI_Comm) ::  comm
      INTEGER ::  error_code
```

In C++, the user routine should be of the form:

{typedef void MPI::Comm::Errhandler_function(MPI::Comm &, int *, ...);
*(binding deprecated, see Section 15.2)*}

*Rationale.* The variable argument list is provided because it provides an ISO-standard hook for providing additional information to the error handler; without this hook, ISO C prohibits additional arguments. (*End of rationale.*)

*Advice to users.* A newly created communicator inherits the error handler that is associated with the "parent" communicator. In particular, the user can specify a "global" error handler for all communicators by associating this handler with the communicator MPI_COMM_WORLD immediately after initialization. (*End of advice to users.*)

MPI_COMM_SET_ERRHANDLER(comm, errhandler)

| | | |
|---|---|---|
| INOUT | comm | communicator (handle) |
| IN | errhandler | new error handler for communicator (handle) |

```
int MPI_Comm_set_errhandler(MPI_Comm comm, MPI_Errhandler errhandler)
```

```
MPI_COMM_SET_ERRHANDLER(COMM, ERRHANDLER, IERROR)
    INTEGER COMM, ERRHANDLER, IERROR
```

```
MPI_Comm_set_errhandler(comm, errhandler, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Errhandler), INTENT(IN) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{void MPI::Comm::Set_errhandler(const MPI::Errhandler& errhandler)*(binding deprecated, see Section 15.2)* }

Attaches a new error handler to a communicator. The error handler must be either a predefined error handler, or an error handler created by a call to MPI_COMM_CREATE_ERRHANDLER. This call is identical to MPI_ERRHANDLER_SET, whose use is deprecated.

ticket230-B.

ticket230-B.
ticket-248T.

ticket-248T.

MPI_COMM_GET_ERRHANDLER(comm, errhandler)

| | | |
|---|---|---|
| IN | comm | communicator (handle) |
| OUT | errhandler | error handler currently associated with communicator (handle) |

```
int MPI_Comm_get_errhandler(MPI_Comm comm, MPI_Errhandler *errhandler)
```

```
MPI_COMM_GET_ERRHANDLER(COMM, ERRHANDLER, IERROR)
    INTEGER COMM, ERRHANDLER, IERROR
```

ticket-248T.

```
MPI_Comm_get_errhandler(comm, errhandler, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Errhandler), INTENT(OUT) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{MPI::Errhandler MPI::Comm::Get_errhandler() const*(binding deprecated, see
          Section 15.2) }

Retrieves the error handler currently associated with a communicator.  This call is identical to MPI_ERRHANDLER_GET, whose use is deprecated.

Example: A library function may register at its entry point the current error handler for a communicator, set its own private error handler for this communicator, and restore before exiting the previous error handler.

### 8.3.2   Error Handlers for Windows

ticket252-W.

MPI_WIN_CREATE_ERRHANDLER(win_errhandler_fn, errhandler)

| | | |
|---|---|---|
| IN | [ticket252-W.]win_errhandler_fn | user defined error handling procedure (function) |
| OUT | errhandler | MPI error handler (handle) |

ticket252-W.

```
int MPI_Win_create_errhandler(MPI_Win_errhandler_function
            *win_errhandler_fn, MPI_Errhandler *errhandler)
```

ticket252-W.

ticket252-W.

```
MPI_WIN_CREATE_ERRHANDLER(WIN_ERRHANDLER_FN, ERRHANDLER, IERROR)
    EXTERNAL WIN_ERRHANDLER_FN
    INTEGER ERRHANDLER, IERROR
```

ticket-248T.

```
MPI_Win_create_errhandler(win_errhandler_fn, errhandler, ierror)
    PROCEDURE(MPI_Win_errhandler_function) ::  win_errhandler_fn
    TYPE(MPI_Errhandler), INTENT(OUT) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{static MPI::Errhandler
            MPI::Win::Create_errhandler(MPI::Win::Errhandler_function*

ticket252-W.

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

Creates an error handler that can be attached to a window object. The user routine should be, in C, a function of type MPI_Win_errhandler_function which is defined as

```
typedef void MPI_Win_errhandler_function(MPI_Win *, int *, ...);
```

The first argument is the window in use, the second is the error code to be returned.
With the Fortran `mpi` module and `mpif.h`, the user routine should be of the form:

```
SUBROUTINE WIN_ERRHANDLER_FUNCTION(WIN, ERROR_CODE)
    INTEGER WIN, ERROR_CODE
```

With the Fortran `mpi_f08` module, the user routine should be of the form:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Win_errhandler_function(win, error_code)
    TYPE(MPI_Win) ::  win
    INTEGER ::  error_code
```

In C++, the user routine should be of the form:

```
{typedef void MPI::Win::Errhandler_function(MPI::Win &, int *, ...);
```
*(binding deprecated, see Section 15.2)*}


MPI_WIN_SET_ERRHANDLER(win, errhandler)

| | | |
|---|---|---|
| INOUT | win | window (handle) |
| IN | errhandler | new error handler for window (handle) |

```
int MPI_Win_set_errhandler(MPI_Win win, MPI_Errhandler errhandler)
```

```
MPI_WIN_SET_ERRHANDLER(WIN, ERRHANDLER, IERROR)
    INTEGER WIN, ERRHANDLER, IERROR
```

```
MPI_Win_set_errhandler(win, errhandler, ierror)
    TYPE(MPI_Win), INTENT(IN) ::  win
    TYPE(MPI_Errhandler), INTENT(IN) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{void MPI::Win::Set_errhandler(const MPI::Errhandler& errhandler)(binding
```
*deprecated, see Section 15.2)* }

Attaches a new error handler to a window. The error handler must be either a pre-
defined error handler, or an error handler created by a call to
MPI_WIN_CREATE_ERRHANDLER.


MPI_WIN_GET_ERRHANDLER(win, errhandler)

| | | |
|---|---|---|
| IN | win | window (handle) |
| OUT | errhandler | error handler currently associated with window (handle) |

```
int MPI_Win_get_errhandler(MPI_Win win, MPI_Errhandler *errhandler)
```

```
MPI_WIN_GET_ERRHANDLER(WIN, ERRHANDLER, IERROR)
    INTEGER WIN, ERRHANDLER, IERROR
```

ticket230-B.

ticket230-B.
ticket-248T.

ticket-248T.

ticket-248T.

```
MPI_Win_get_errhandler(win, errhandler, ierror)
    TYPE(MPI_Win), INTENT(IN) ::  win
    TYPE(MPI_Errhandler), INTENT(OUT) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{MPI::Errhandler MPI::Win::Get_errhandler() const *(binding deprecated, see Section 15.2)* }

Retrieves the error handler currently associated with a window.

### 8.3.3 Error Handlers for Files

ticket252-W.

MPI_FILE_CREATE_ERRHANDLER(file_errhandler_fn, errhandler)

| | | |
|---|---|---|
| IN | [ticket252-W.]file_errhandler_fn | user defined error handling procedure (function) |
| OUT | errhandler | MPI error handler (handle) |

ticket252-W.
```
int MPI_File_create_errhandler(MPI_File_errhandler_function
            *file_errhandler_fn, MPI_Errhandler *errhandler)
```

ticket252-W.
```
MPI_FILE_CREATE_ERRHANDLER(FILE_ERRHANDLER_FN, ERRHANDLER, IERROR)
```
ticket252-W.
```
    EXTERNAL FILE_ERRHANDLER_FN
    INTEGER ERRHANDLER, IERROR
```
ticket-248T.
```
MPI_File_create_errhandler(file_errhandler_fn, errhandler, ierror)
    PROCEDURE(MPI_File_errhandler_function) ::  file_errhandler_fn
    TYPE(MPI_Errhandler), INTENT(OUT) ::  errhandler
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

{static MPI::Errhandler
            MPI::File::Create_errhandler(MPI::File::Errhandler_function*
ticket252-W.            file_errhandler_fn) *(binding deprecated, see Section 15.2)* }

Creates an error handler that can be attached to a file object. The user routine should be, in C, a function of type MPI_File_errhandler_function, which is defined as
```
typedef void MPI_File_errhandler_function(MPI_File *, int *, ...);
```

ticket230-B.
The first argument is the file in use, the second is the error code to be returned.
With the Fortran mpi module and mpif.h, the user routine should be of the form:
```
SUBROUTINE FILE_ERRHANDLER_FUNCTION(FILE, ERROR_CODE)
    INTEGER FILE, ERROR_CODE
```
ticket230-B.
ticket-248T.
With the Fortran mpi_f08 module, the user routine should be of the form:
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_File_errhandler_function(file, error_code)
    TYPE(MPI_File) ::  file
    INTEGER ::  error_code
```

In C++, the user routine should be of the form:
```
{typedef void MPI::File::Errhandler_function(MPI::File &, int *, ...);
            (binding deprecated, see Section 15.2)}
```

MPI standard, process 0 must return from the complete call after a bounded delay, even if process 1 does not reach any MPI call in this period of time. According to another interpretation, the complete call may block until process 1 reaches the wait call, or reaches another MPI call. The qualitative behavior is the same, under both interpretations, unless a process is caught in an infinite compute loop, in which case the difference may not matter. However, the quantitative expectations are different. Different MPI implementations reflect these different interpretations. While this ambiguity is unfortunate, it does not seem to affect many real codes. The MPI forum decided not to decide which interpretation of the standard is the correct one, since the issue is very contentious, and a decision would have much impact on implementors but less impact on users. (*End of rationale.*)

### 11.7.3 Registers and Compiler Optimizations

*Advice to users.* All the material in this section is an advice to users. (*End of advice to users.*)

A coherence problem exists between variables kept in registers and the memory value of these variables. An RMA call may access a variable in memory (or cache), while the up-to-date value of this variable is in register. A get will not return the latest variable value, and a put may be overwritten when the register is stored back in memory.

The problem is illustrated by the following code:

| **Source of Process 1** | **Source of Process 2** | **Executed in Process 2** |
|---|---|---|
| bbbb = 777 | buff = 999 | reg_A:=999 |
| call MPI_WIN_FENCE | call MPI_WIN_FENCE | |
| call MPI_PUT(bbbb | | stop appl.thread |
| into buff of process 2) | | buff:=777 in PUT handler |
| | | continue appl.thread |
| call MPI_WIN_FENCE | call MPI_WIN_FENCE | |
| | ccc = buff | ccc:=reg_A |

In this example, variable `buff` is allocated in the register `reg_A` and therefore `ccc` will have the old value of `buff` and not the new value 777.

This problem, which also afflicts in some cases send/receive communication, is discussed more at length in Section 16.2.8.

MPI implementations will avoid this problem for standard conforming C programs. Many Fortran compilers will avoid this problem, without disabling compiler optimizations. However, in order to avoid register coherence problems in a completely portable manner, users should restrict their use of RMA windows to variables stored in in modules or COMMON blocks. To prevent problems with the argument copying and register optimization done by Fortran compilers, please note the hints in Sections 16.2.2-16.2.12, especially in Sections 16.2.4 and 16.2.5 on pages 547-550 about *problems due to data copying and sequence association with subscript triplets* and *vector subscripts*, and in Sections 16.2.8 to 16.2.11 on pages 551 to 560 about *optimization problems*, *code movements and register optimization*, and *temporary* and *permanent data movements*. Sections "Solutions" to "VOLATILE" on pages 555-558 discuss several solutions for the problem in this example.

ticket238-J.
ticket238-J.

ticket236-H.
ticket238-J.

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
```

ticket-248T.

```
MPI_Grequest_start(query_fn, free_fn, cancel_fn, extra_state, request,
              ierror)
    PROCEDURE(MPI_Grequest_query_function) ::  query_fn
    PROCEDURE(MPI_Grequest_free_function) ::  free_fn
    PROCEDURE(MPI_Grequest_cancel_function) ::  cancel_fn
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state
    TYPE(MPI_Request), INTENT(OUT) ::  request
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{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.

**Unofficial Draft for Comment Only**

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; extra_state 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 with the mpi module and mpif.h

```
SUBROUTINE GREQUEST_QUERY_FUNCTION(EXTRA_STATE, STATUS, IERROR)
    INTEGER STATUS(MPI_STATUS_SIZE), IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

in Fortran with the mpi_f08 module

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Grequest_query_function(extra_state, status, ierror)
     TYPE(MPI_Status) ::  status
     INTEGER ::  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)}
```

The `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).

The `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 or has omitted the status argument (with the mpi_f08 Fortran module or C++), 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);
```

in Fortran with the mpi module and mpif.h

```
SUBROUTINE GREQUEST_FREE_FUNCTION(EXTRA_STATE, IERROR)
    INTEGER IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

**Unofficial Draft for Comment Only**

ticket230-B. ¹    in Fortran with the `mpi_f08` module

ticket-248T. ²

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_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)*}

ticket0.    The free_fn function is invoked to clean up user-allocated resources when the generalized

ticket0.    request is freed.

The 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

ticket0.    by MPI.

The 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

ticket0.    accessing this stale handle. This is a special case in which 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);
```

ticket230-B.    in Fortran with the `mpi` module and `mpif.h`

```
SUBROUTINE GREQUEST_CANCEL_FUNCTION(EXTRA_STATE, COMPLETE, IERROR)
    INTEGER IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
    LOGICAL COMPLETE
```

ticket230-B.    in Fortran with the `mpi_f08` module

ticket-248T.

```
ABSTRACT INTERFACE
```

```
SUBROUTINE MPI_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)}
```

The cancel_fn function is invoked to start the cancelation of a generalized request. It is called by MPI_CANCEL(request). MPI passes complete=true to the callback function 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 or the status argument was omitted, 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
```

```
MPI_Grequest_complete(request, ierror)
    TYPE(MPI_Request), INTENT(IN) ::  request
    INTEGER, OPTIONAL, INTENT(OUT) ::  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

ticket0.

ticket0.

ticket244-P.

ticket-248T.

- On any MPI process, each file handle may have at most one active split collective operation at any time.

- Begin calls are collective over the group of processes that participated in the collective open and follow the ordering rules for collective calls.

- End calls are collective over the group of processes that participated in the collective open and follow the ordering rules for collective calls. Each end call matches the preceding begin call for the same collective operation. When an "end" call is made, exactly one unmatched "begin" call for the same operation must precede it.

- An implementation is free to implement any split collective data access routine using the corresponding blocking collective routine when either the begin call (e.g., MPI_FILE_READ_ALL_BEGIN) or the end call (e.g., MPI_FILE_READ_ALL_END) is issued. The begin and end calls are provided to allow the user and MPI implementation to optimize the collective operation.

- Split collective operations do not match the corresponding regular collective operation. For example, in a single collective read operation, an MPI_FILE_READ_ALL on one process does not match an MPI_FILE_READ_ALL_BEGIN/ MPI_FILE_READ_ALL_END pair on another process.

- Split collective routines must specify a buffer in both the begin and end routines. By specifying the buffer that receives data in the end routine, we can avoid the problems described in "A Problem with Code Movements and Register Optimization," Section 16.2.9 on page 552, but not all of the problems described in Section 16.2.8 on page 551.

- No collective I/O operations are permitted on a file handle concurrently with a split collective access on that file handle (i.e., between the begin and end of the access). That is

```
MPI_File_read_all_begin(fh, ...);
...
MPI_File_read_all(fh, ...);
...
MPI_File_read_all_end(fh, ...);
```

is erroneous.

- In a multithreaded implementation, any split collective begin and end operation called by a process must be called from the same thread. This restriction is made to simplify the implementation in the multithreaded case. (Note that we have already disallowed having two threads begin a split collective operation on the same file handle since only one split collective operation can be active on a file handle at any time.)

The arguments for these routines have the same meaning as for the equivalent collective versions (e.g., the argument definitions for MPI_FILE_READ_ALL_BEGIN and MPI_FILE_READ_ALL_END are equivalent to the arguments for MPI_FILE_READ_ALL). The begin routine (e.g., MPI_FILE_READ_ALL_BEGIN) begins a split collective operation

ticket238-J.
ticket238-J.
ticket238-J.

*Advice to implementors.*     When converting a larger size integer to a smaller size integer, only the less significant bytes are moved. Care must be taken to preserve the sign bit value. This allows no conversion errors if the data range is within the range of the smaller size integer. (*End of advice to implementors.*)

Table 13.2 specifies the sizes of predefined datatypes in "external32" format.

### 13.5.3  User-Defined Data Representations

There are two situations that cannot be handled by the required representations:

1. a user wants to write a file in a representation unknown to the implementation, and

2. a user wants to read a file written in a representation unknown to the implementation.

User-defined data representations allow the user to insert a third party converter into the I/O stream to do the data representation conversion.

MPI_REGISTER_DATAREP(datarep, read_conversion_fn, write_conversion_fn,
                dtype_file_extent_fn, extra_state)

| | | |
|---|---|---|
| IN | datarep | data representation identifier (string) |
| IN | read_conversion_fn | function invoked to convert from file representation to native representation (function) |
| IN | write_conversion_fn | function invoked to convert from native representation to file representation (function) |
| IN | dtype_file_extent_fn | function invoked to get the extent of a datatype as represented in the file (function) |
| IN | extra_state | extra state |

```
int MPI_Register_datarep(char *datarep,
            MPI_Datarep_conversion_function *read_conversion_fn,
            MPI_Datarep_conversion_function *write_conversion_fn,
            MPI_Datarep_extent_function *dtype_file_extent_fn,
            void *extra_state)
```

```
MPI_REGISTER_DATAREP(DATAREP, READ_CONVERSION_FN, WRITE_CONVERSION_FN,
            DTYPE_FILE_EXTENT_FN, EXTRA_STATE, IERROR)
    CHARACTER*(*) DATAREP
    EXTERNAL READ_CONVERSION_FN, WRITE_CONVERSION_FN, DTYPE_FILE_EXTENT_FN
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
    INTEGER IERROR
```

ticket-248T.

```
MPI_Register_datarep(datarep, read_conversion_fn, write_conversion_fn,
            dtype_file_extent_fn, extra_state, ierror)
    CHARACTER(LEN=*), INTENT(IN) ::  datarep
    EXTERNAL ::  read_conversion_fn, write_conversion_fn
    PROCEDURE(MPI_Datarep_extent_function) ::  dtype_file_extent_fn
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state
```

```
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

```
{void MPI::Register_datarep(const char* datarep,
            MPI::Datarep_conversion_function* read_conversion_fn,
            MPI::Datarep_conversion_function* write_conversion_fn,
            MPI::Datarep_extent_function* dtype_file_extent_fn,
            void* extra_state)(binding deprecated, see Section 15.2) }
```

The call associates read_conversion_fn, write_conversion_fn, and dtype_file_extent_fn with the data representation identifier datarep. datarep can then be used as an argument to MPI_FILE_SET_VIEW, causing subsequent data access operations to call the conversion functions to convert all data items accessed between file data representation and native representation. MPI_REGISTER_DATAREP is a local operation and only registers the data representation for the calling MPI process. If datarep is already defined, an error in the error class MPI_ERR_DUP_DATAREP is raised using the default file error handler (see Section 13.7, page 510). The length of a data representation string is limited to the value of MPI_MAX_DATAREP_STRING. MPI_MAX_DATAREP_STRING must have a value of at least 64. No routines are provided to delete data representations and free the associated resources; it is not expected that an application will generate them in significant numbers.

Extent Callback

```
typedef int MPI_Datarep_extent_function(MPI_Datatype datatype,
            MPI_Aint *file_extent, void *extra_state);
```

```
SUBROUTINE DATAREP_EXTENT_FUNCTION(DATATYPE, EXTENT, EXTRA_STATE, IERROR)
    INTEGER DATATYPE, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTENT, EXTRA_STATE
```

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Datarep_extent_function(datatype, extent, extra_state,
  ierror)
      TYPE(MPI_Datatype) ::  datatype
      INTEGER ::  ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extent, extra_state
```

```
{typedef void MPI::Datarep_extent_function(const MPI::Datatype& datatype,
            MPI::Aint& file_extent, void* extra_state); (binding deprecated,
            see Section 15.2)}
```

The function dtype_file_extent_fn must return, in file_extent, the number of bytes required to store datatype in the file representation. The function is passed, in extra_state, the argument that was passed to the MPI_REGISTER_DATAREP call. MPI will only call this routine with predefined datatypes employed by the user.

Datarep Conversion Functions

```
typedef int MPI_Datarep_conversion_function(void *userbuf,
            MPI_Datatype datatype, int count, void *filebuf,
            MPI_Offset position, void *extra_state);
```

ticket-248T.

ticket230-B.

To reduce problems with the choice buffer arguments, implicit interfaces are used with all three Fortran support methods; the first one with the `mpi` module or `mpif.h`, the second one with `mpi_f08`.

```
SUBROUTINE DATAREP_CONVERSION_FUNCTION(USERBUF, DATATYPE, COUNT, FILEBUF,
             POSITION, EXTRA_STATE, IERROR)
    <TYPE> USERBUF(*), FILEBUF(*)
    INTEGER COUNT, DATATYPE, IERROR
    INTEGER(KIND=MPI_OFFSET_KIND) POSITION
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

ticket230-B.

```
SUBROUTINE DATAREP_CONVERSION_FUNCTION(userbuf, datatype, count, filebuf,
             position, extra_state, ierror)
    <TYPE> userbuf(*), filebuf(*)
    TYPE(MPI_Datatype datatype
    INTEGER count, ierror
    INTEGER(KIND=MPI_OFFSET_KIND) position
    INTEGER(KIND=MPI_ADDRESS_KIND) extra_state
```

```
{typedef void MPI::Datarep_conversion_function(void* userbuf,
             MPI::Datatype& datatype, int count, void* filebuf,
             MPI::Offset position, void* extra_state);
```
*(binding deprecated, see Section 15.2)*}

The function read_conversion_fn must convert from file data representation to native representation. Before calling this routine, MPI allocates and fills filebuf with count contiguous data items. The type of each data item matches the corresponding entry for the predefined datatype in the type signature of datatype. The function is passed, in extra_state, the argument that was passed to the MPI_REGISTER_DATAREP call. The function must copy all count data items from filebuf to userbuf in the distribution described by datatype, converting each data item from file representation to native representation. datatype will be equivalent to the datatype that the user passed to the read function. If the size of datatype is less than the size of the count data items, the conversion function must treat datatype as being contiguously tiled over the userbuf. The conversion function must begin storing converted data at the location in userbuf specified by position into the (tiled) datatype.

> *Advice to users.* Although the conversion functions have similarities to MPI_PACK and MPI_UNPACK, one should note the differences in the use of the arguments count and position. In the conversion functions, count is a count of data items (i.e., count of typemap entries of datatype), and position is an index into this typemap. In MPI_PACK, incount refers to the number of whole datatypes, and position is a number of bytes. (*End of advice to users.*)

> *Advice to implementors.* A converted read operation could be implemented as follows:
>
>   1. Get file extent of all data items
>
>   2. Allocate a filebuf large enough to hold all count data items
>
>   3. Read data from file into filebuf
>
>   4. Call read_conversion_fn to convert data and place it into userbuf

**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 language bindings have been designed to be compatible with the Fortran 90 standard with additional features from Fortran 2003 and Fortran 2008 [34] + TR 29113 [35].

> *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. In Fortran 2008 + TR 29113, the major new language features used are assumed-type and assumed-rank dummy arguments. They were defined to allow the definition of choice arguments as part of the Fortran language. Further requirements for compiler support are listed in Section 16.2.16 on page 566. (*End of rationale.*)

MPI defines three methods of Fortran support:

1. **INCULDE 'mpif.h'** This method is described in Section 16.2.13. The use of the include file `mpif.h` is strongly discouraged starting with MPI-3.0.

2. **USE mpi** This method is described in Section 16.2.14 and requires compile-time argument checking.

3. **USE mpi_f08** This method is described in Section 16.2.15 and requires compile-time argument checking that also includes unique handle types.

Compliant MPI-3 implementations providing a Fortran interface must provide all three Fortran support methods.

Application subroutines and functions may use either one of the modules or the `mpif.h` include file. An implementation may require the use of modules to prevent type mismatch errors.

> *Advice to users.*   Users are advised to utilize one of the MPI modules even if `mpif.h` enforces type checking on a particular system. Using a module provides several potential advantages over using an include file. (*End of advice to users.*)

ticket230-B.
ticket230-B.
ticket230-B.
ticket0.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket233-E.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.
ticket230-B.

In a single application, it must be possible to link together routines which USE `mpi_f08`, USE `mpi` and INCLUDE `mpif.h`.

The `INTEGER` compile-time constant `MPI_SUBARRAYS` is set to `MPI_SUBARRAYS_SUPPORTED` if all buffer choice arguments are defined in explicit interfaces with assumed-type and assumed-rank [35]; otherwise it is set to `MPI_SUBARRAYS_NOT_SUPPORTED`. This constant exists with each Fortran support method, but not in the C/C++ header files. The value may be different for each Fortran support method.

Section 16.2.2 to 16.2.11 gives an overview and details on known problems when using Fortran together with MPI; Section 16.2.12 compares the Fortran problems with C. Section 16.2.13 to 16.2.15 define the Fortran support methods. Section 16.2.16 summarizes major requirements for valid MPI-3.0 implementations with Fortran support. Section 16.2.17 and Section 16.2.18 describe additional functionality that is part of the Fortran support. `MPI_F_SYNC_REG` is needed for one of the methods to prevent register optimization problems. A set of functions provides 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. In the context of MPI, parameterized types are Fortran intrinsic types which are specified using `KIND` type parameters.

### 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.

The following MPI features are inconsistent with Fortran 90.

1. An MPI subroutine with a choice argument may be called with different argument types. When using the `mpi_f08` module together with a compiler that supports Fortran 2008 + TR 29113, this problem is resolved.

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 executing 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 15 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. TODO: THIS ITEM IS WRONG AND MUST BE CORRECTED

ticket245-Q.

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 support 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 17 and Section 4.1.1 on page 87 for more information.

ticket230-B.

Sections 16.2.3 to 16.2.11 describe in detail several problems between MPI and Fortran and their solutions. Some of these solutions require special capabilities from the compilers. Major requirements are summarized in Section 16.2.16 on page 566.

ticket230-B.

### 16.2.3   Problems Due to Strong Typing

ticket230-B.

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. Similar to C, with Fortran 2008 + TR 29113 (and later) together with the `mpi_f08` module, the problem is avoided by declaring choice arguments with `TYPE(*), DIMENSION(..)`, i.e., as assumed-type and assumed-rank dummy arguments.

ticket235-G.

Using `INCLUDE mpif.h`, the following code fragment might technically be invalid and may generate a compile-time error.

ticket235-G.
ticket235-G.

```
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. Using the `mpi_f08` or `mpi` module, the problem is usually resolved through the assumed-type and assumed-rank declarations of the dummy arguments, or with a compiler-dependent mechanism that overrides type checking for choice arguments.

It is also technically invalid in Fortran to pass a scalar actual argument to an array dummy argument. Thus, when using the `mpi_f08` or `mpi` module, the following code fragment usually generates an error since the `dims` and `periods` arguments to `MPI_CART_CREATE` are declared as assumed size arrays `INTEGER :: DIMS(*)` and `LOGICAL :: PERIODS(*)`.

```
USE mpi_f08
INTEGER size
CALL MPI_Cart_create( comm_old,1,size,.TRUE.,.TRUE.,comm_cart,ierror )
```

Using `INCLUDE 'mpif.h'`, compiler warnings are not expected unless this include file also uses Fortran explicit interfaces.

### 16.2.4 Problems Due to Data Copying and Sequence Association with Subscript Triplets

Arrays with subscript **triplets** describe in Fortran subarrays with or without strides, e.g.,

```
REAL a(100,100,100)
CALL MPI_Send( a(11:17, 12:99:3, 1:100), 7*30*100, MPI_REAL, ...)
```

The handling of subscript triplets depends on the value of the constant `MPI_SUBARRAYS`:

- If `MPI_SUBARRAYS` equals `MPI_SUBARRAYS_SUPPORTED`:

  Choice buffer arguments are declared as `TYPE(*)`, `DIMENSION(..)`. For example, considering the following code fragment:

  ```
  REAL s(100), r(100)
  CALL MPI_Isend(s(1:100:5), 3, MPI_REAL, ..., rq, ierror)
  CALL MPI_Wait(rq, status, ierror)
  CALL MPI_Irecv(r(1:100:5), 3, MPI_REAL, ..., rq, ierror)
  CALL MPI_Wait(rq, status, ierror)
  ```

  In this case, the individual elements `s(1)`, `s(6)`, and `s(11)` are sent between the start of MPI_ISEND and the end of MPI_WAIT even though the compiled code must not copy `s(1:100:5)` to a contiguous temporary scratch buffer. Instead, the compiled code will pass a descriptor to MPI_ISEND that allows MPI to operate directly on `s(1)`, `s(6)`, `s(11)`, ..., `s(96)`.

  All nonblocking MPI functions (e.g., MPI_ISEND, MPI_PUT, MPI_FILE_WRITE_ALL_BEGIN) behave as if the user-specified elements of choice buffers are copied to a contiguous scratch buffer in the MPI runtime environment. All datatype descriptions (in the example above, "`3, MPI_REAL`") read and store data from and to this virtual contiguous scratch buffer. Displacements in MPI derived datatypes are relative to the beginning of this virtual contiguous scratch buffer. Upon completion of a nonblocking receive operation (e.g., when MPI_WAIT on a corresponding MPI_Request returns), it is as if the received data has been copied from

the virtual contiguous scratch buffer back to the non-contiguous application buffer. In the example above, `r(1)`, `r(6)`, and `r(11)` will be filled with the received data when MPI_WAIT returns.

> *Advice to implementors.* The Fortran descriptor for `TYPE(*), DIMENSION(..)` arguments contains enough information that the MPI library can make a real contiguous copy of non-contiguous user buffers if desired. Efficient implementations may avoid such additional memory-to-memory data copying. (*End of advice to implementors.*)

> *Rationale.* If MPI_SUBARRAYS equals MPI_SUBARRAYS_SUPPORTED, non-contiguous buffers are handled inside of the MPI library instead of by the compiler through argument association conventions. Therefore the scope of MPI library scratch buffers can be from the beginning of a nonblocking operation until the completion of the operation although beginning and completion are implemented in different routines. If MPI_SUBARRAYS equals MPI_SUBARRAYS_NOT_SUPPORTED, temporary copies made by the compiler will not exist throughout the duration of the entire nonblocking MPI operation, which is too short for implementing the entire MPI operation. (*End of rationale.*)

• If MPI_SUBARRAYS equals MPI_SUBARRAYS_NOT_SUPPORTED:

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.[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:

```
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

ticket236-H.

---

[1]Technically, the Fortran standards are worded to allow non-contiguous storage of any array data, unless the dummy argument has the CONTIGUOUS attribute.

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 'simply contiguous' section such as A(1:N) of such an array. (We define 'simply contiguous' 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 'simply contiguous' 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 'simply contiguous' section of a contiguous array will also be contiguous.[2]

The same problem can occur with a scalar argument. Some compilers make a copy of some scalar dummy arguments within a called procedure when passed as an actual argument to a choice buffer routine. 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 'simply contiguous' 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.

ticket236-H.

If a compiler makes copies in the calling procedure of arguments that are explicit-shape or assumed-size arrays, simply contiguous 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.

ticket236-H.

ticket236-H.
ticket236-H.

### 16.2.5   Problems Due to Data Copying and Sequence Association with Vector Subscripts

Arrays with **vector** subscripts describe in Fortran subarrays containing an possibly irregular set of elements

```
REAL a(100)
CALL MPI_Send( A((/7,9,23,81,82/)), 5, MPI_REAL, ...)
```

Arrays with a vector subscript must not be used as actual choice buffer arguments in any nonblocking or split collective MPI operations; they may be used in blocking MPI operations.

ticket230-B.

### 16.2.6   Special Constants

ticket230-B.

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 15. 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.

ticket250-V.
ticket250-V.

In Fortran, 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, the address of the actual choice buffer argument can be compared with the address of such a predefined static variable.

ticket250-V.
ticket242-N.

With USE mpi_f08, the attributes INTENT(IN), INTENT(OUT), and INTENT(INOUT) are used in the Fortran interface. In most cases, INTENT(IN) is used if the C interface uses call-by-value. For all buffer arguments and for OUT and INOUT dummy arguments that allow one of these special constants as input, an INTENT(...) is not specified.

ticket230-B.
ticket237-I.

### 16.2.7   Fortran Derived Types

ticket230-B.
ticket237-I.
ticket237-I.
ticket237-I.
ticket237-I.
ticket237-I.
ticket237-I.
ticket237-I.

MPI supports passing Fortran sequence and BIND(C) derived types to choice dummy arguments, but does not support Fortran derived types that neither have the SEQUENCE nor the BIND(C) attribute.

The following code fragment shows one possible way to send a sequence derived type in Fortran. The example assumes that all data is passed by address.

```
type mytype
    [ticket237-I.]SEQUENCE
```

```
      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)

[ticket237-I.]
[ticket237-I.]
   call MPI_SEND(foo%i, 1, newtype, ...)
```

### 16.2.8 Optimization Problems, an Overview

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. These problems are independent of the Fortran support method, i.e., they occur with the `mpi_f08` module, the `mpi` module, and the `mpif.h` include file.

This section shows four problematic usage areas (the abbrevations in parentheses are used in the table below):

- Usage of nonblocking routines *(Nonbl.)*.

- Usage of one-sided routines *(1-sided)*.

- Usage of MPI parallel file I/O split collective operations *(Split)*.

- Use of `MPI_BOTTOM` together with absolute displacements in MPI datatypes, or relative displacements between two variables in such datatypes *(Bottom)*.

The following compiler optimization strategies (valid for serial code) may cause problems in MPI applications:

- Code movements and register optimization problems; see Section 16.2.9 on page 552.

- Temporary data movements and temporary memory modifications; see Section 16.2.10 on page 558.

- Permanent data movements (e.g., through garbage collection); see Section 16.2.11 on page 560.

Table 16.4 shows in which usage areas the optimization problems may only occur.

| Optimization ... | ... may cause a problem in following usage areas | | | |
|---|---|---|---|---|
| | Nonbl. | 1-sided | Split | Bottom |
| Code movements and register optimization | yes | yes | no | yes |
| Temporary data movements | yes | yes | yes | no |
| Permanent data movements | yes | yes | yes | yes |

Table 16.4: Occurrence of Fortran optimization problems in several usage areas

The solutions in the following sections are based on compromises

- to minimize the burden for the application programmer, e.g., as shown in Sections "Solutions" to "VOLATILE" on pages 555-558,

- to minimize the additional needs in the Fortran standard, e.g., in the Fortran 2008 TR 29113 [35],

- to minimize the drawbacks on compiler based optimization, and

- to minimize the requirements defined in Section 16.2.16 on page 566.

### 16.2.9   Problems with Code Movements and Register Optimization

#### Nonblocking operations

If 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.

ticket238-J.
ticket238-J.

**Example 16.11** Fortran 90 register optimization – extreme.

```
Source                    compiled as              or compiled as
[ticket238-J.]REAL :: buf, b1             REAL :: buf, b1            REAL :: buf, b1
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
```

Example 16.11 shows extreme, but allowed, possibilities. 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 column.

[ticket238-J.]

**Example 16.12** Similar example with MPI_ISEND

```
Source                    compiled as              with a possible MPI-internal
                                                   execution sequence
REAL :: buf, copy         REAL :: buf, copy        REAL :: buf, copy
buf = val                 buf = val                buf = val
call MPI_ISEND(buf,..req)  call MPI_ISEND(buf,..req)  addr = &buf
copy = buf                copy= buf                copy= buf
                          buf = val_overwrite      buf = val_overwrite
call MPI_WAIT(req,..)     call MPI_WAIT(req,..)     send(*addr)
buf = val_overwrite
```

Due to valid compiler code movement optimizations in Example 16.12, the content of buf may already be overwritten when the content of buf is sent. The code movement is permitted because the compiler cannot detect a possible access to buf in MPI_WAIT (or in a second thread between the start of MPI_ISEND and the end of MPI_WAIT).

Such register optimization is based on moving code; here, the access to buf was moved from after MPI_Wait to before MPI_Wait. Note that code movement may also occur across subroutine boundaries when subroutines or functions are inlined.

This register optimization / code movement problem for nonblocking operations does not occur with MPI parallel file I/O split collective operations, because in the ..._BEGIN and ..._END calls, the same buffer has to be provided as an actual argument. The register optimization / code movement problem for MPI_BOTTOM and derived MPI datatypes may occur in each blocking and nonblocking communication or parallel file I/O operation.

### One-sided communication

An example with instruction reordering due to register optimization can be found in Section 11.7.3 on page 427.

### MPI_BOTTOM and combining independent variables in datatypes

This section is only relevant if the MPI program uses a buffer argument to an
MPI_SEND, MPI_RECV etc., 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 referenced in the call. Also attention must be paid if MPI operations are used that run
in parallel with the user's application.

Example 16.13 shows what Fortran compilers are allowed to do.

**Example 16.13** 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
```

In Example 16.13, the compiler does not invalidate the register because it cannot
see that MPI_RECV changes the value of buf. The access to buf is hidden by the use of
MPI_GET_ADDRESS and MPI_BOTTOM.

[ticket238-J.]

**Example 16.14** Similar example with MPI_SEND

This source ...                                        can be compiled as:

```
! buf contains val_old                   ! buf contains val_old
buf = val_new
call MPI_SEND(MPI_BOTTOM,1,type,...)      call MPI_SEND(...)
! with buf as a displacement in type     ! i.e. val_old is sent
                                          !
                                          ! buf=val_new  is moved to here
                                          ! and detected as dead code
                                          ! and therefore removed
                                          !
buf = val_overwrite                       buf = val_overwrite
```

In Example 16.14, several successive assigments to the same variable
buf can be combined in the way, that only the last assignment is executed. Successive
means that no interfering load access to this variable occurs between the assignments. The
compiler cannot detect that the call to MPI_SEND statement is interfering, because the
load access to buf is hidden by the usage of MPI_BOTTOM.

## Solutions

The following sections show in detail how the problems with code movements and register optimizations can be solved in a portable way. Application writers can partially or fully avoid these compiler optimization problems by using one or more of the special Fortran declarations with the send and receive buffers used in nonblocking operations, or in all operation if MPI_BOTTOM is used, or datatype handles that combine several variables:

- Usage of the Fortran TARGET attribute.

- Usage of the helper routine MPI_F_SYNC_REG, or an equivalent user-written dummy routine.

- Declaring the buffer as a Fortran module variable or within a Fortran common block.

- Usage of the Fortran VOLATILE attribute.

Each of these methods can solve the problems of code movement and register optimization, however, may involve different degree of performance impact, and may not be usable in every application context. The VOLATILE attribute may have the most negative impact on performance. There is one attribute that cannot be used for this purpose:

- The Fortran ASYNCHRONOUS attribute may not solve code movement problems in MPI applications.

Table 16.5 shows the usability of each method. Each problem category is descibed in detail below.

[ticket238-J.]

| | Nonbl. | 1-sided | Bottom | Overhead may be |
|---|---|---|---|---|
| Examples | 16.11, 16.12 | Section 11.7.3 | 16.13, 16.14 | |
| TARGET | solved | solved | solved | low - medium |
| MPI_F_SYNC_REG | solved | solved | solved | low |
| Module Data | solved | solved | solved | low - medium |
| VOLATILE | solved | solved | solved | high - huge |
| ASYNCHRONOUS | NOT solved | NOT solved | NOT solved | medium - high |

Table 16.5: Usability of methods to prevent Fortran optimization problems

## The Fortran TARGET attribute

Declaring a buffer with the Fortran TARGET attribute in a scoping unit (or BLOCK) tells the compiler that any statement of the scoping unit may be executed while some pointer is

associated with the buffer. Calling a library routine (e.g., an MPI routine) may imply that such a pointer is used to modify the buffer.

- The TARGET attribute solves problems of instruction reordering, code movement, and register optimization related to nonblocking and one-sided communication, or related to the usage of MPI_BOTTOM and derived datatype handles. Declaring REAL, TARGET :: buf solves the register optimization problem in Examples 16.11, 16.12, 16.13, and 16.14

- Unfortunately, the TARGET attribute does **not** solve problems caused by asynchronous accesses between the start and end of a nonblocking or one-sided communication specifically, problems caused by temporary memory modifications are not solved. Example 16.15 on page 559 can **not** be solved with the TARGET attribute.

### Calling MPI_F_SYNC_REG

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. The MPI library provides the MPI_F_SYNC_REG routine for this purpose; see Section 16.2.17 on page 568.

- The problems illustrated by the Examples 16.11 and 16.12 can be solved by calling MPI_F_SYNC_REG(buf) once immediately after MPI_WAIT.

| Example 16.11 | Example 16.12 |
|---|---|
| can be solved with | can be solved with |
| `call MPI_IRECV(buf,..req)` | `buf = val` |
| | `call MPI_ISEND(buf,..req)` |
| | `copy = buf` |
| `call MPI_WAIT(req,..)` | `call MPI_WAIT(req,..)` |
| `call MPI_F_SYNC_REG(buf)` | `call MPI_F_SYNC_REG(buf)` |
| `b1 = buf` | `buf = val_overwrite` |

  The call to MPI_F_SYNC_REG(buf) prevents moving the last line before the MPI_WAIT call. Further calls to MPI_F_SYNC_REG(buf) are not needed because it is still correct if the additional read access `copy=buf` is moved below MPI_WAIT and before `buf=val_overwrite`.

- The problems illustrated by the Examples 16.13 and 16.14 can be solved with two additional MPI_F_SYNC_REG(buf) statements; one directly before MPI_RECV/ MPI_SEND, and one directly after this communication operation.

| Example 16.13 | Example 16.14 |
|---|---|
| can be solved with | can be solved with |
| `call MPI_F_SYNC_REG(buf)` | `call MPI_F_SYNC_REG(buf)` |
| `call MPI_RECV(MPI_BOTTOM,...)` | `call MPI_SEND(MPI_BOTTOM,...)` |
| `call MPI_F_SYNC_REG(buf)` | `call MPI_F_SYNC_REG(buf)` |

ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.

> The first call to MPI_F_SYNC_REG(buf) is needed to finish all load and store references to buf prior to MPI_RECV/MPI_SEND; the second call is needed to assure that the subsequent access to buf are not moved before MPI_RECV/SEND.

- In the example in Section 11.7.3 on page 427, two asynchronous accesses must be protected: in Process 1, the access to bbbb must be protected similar to Example 16.11, i.e., a call to MPI_F_SYNC_REG(bbbb) is needed after the second MPI_WIN_FENCE to guarantee that further accesses to bbbb are not moved ahead of the call to MPI_WIN_FENCE. In Process 2, both calls to MPI_WIN_FENCE together act as a communication call with MPI_BOTTOM as the buffer, i.e., before the first fence and after the second fence, a call to MPI_F_SYNC_REG(buff) is needed to guarantee that accesses to buff are not moved after or ahead of the calls to MPI_WIN_FENCE. Using MPI_GET instead of MPI_PUT, the same calls to MPI_F_SYNC_REG are necessary.

```
Source of Process 1            Source of Process 2
bbbb = 777                     buff = 999
                               call MPI_F_SYNC_REG(buff)
call MPI_WIN_FENCE             call MPI_WIN_FENCE
call MPI_PUT(bbbb
into buff of process 2)


call MPI_WIN_FENCE             call MPI_WIN_FENCE
call MPI_F_SYNC_REG(bbbb)      call MPI_F_SYNC_REG(buff)
                               ccc = buff
```

- The temporary memory modification problem, i.e., Example 16.15 on page 559, can **not** be solved with this method.

A user defined routine instead of MPI_F_SYNC_REG

Instead of MPI_F_SYNC_REG, one can also use a user defined external subroutine, which is separately compiled:

```
subroutine DD(buf)
  integer buf
end
```

Note that if the intent is declared in an explicit interface for 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, a call to MPI_RECV with MPI_BOTTOM as buffer might be replaced by

```
call DD(buf)
call MPI_RECV(MPI_BOTTOM,...)
call DD(buf)
```

Such a user-defined routine was introduced in MPI-2.0 and is still included here to document such usage in existing application programs although new applications should prefer MPI_F_SYNC_REG or one of the other posibilities.

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 ticket238-J.
29 ticket238-J.
30
31
32
33 ticket238-J.
34
35
36
37
38 ticket238-J.
39
40 ticket238-J.
41 ticket238-J.
42 ticket238-J.
43
44
45 ticket238-J.
46
47
48 ticket238-J.
ticket238-J.
ticket238-J.

### Module variables and COMMON blocks

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 may alter the buffer or variable, provided that the compiler cannot infer that the MPI procedure does not reference the module or common block.

- This method solves problems of instruction reordering, code movement, and register optimization related to nonblocking and one-sided communication, or related to the usage of MPI_BOTTOM and derived datatype handles.

- Unfortunately, this method does **not** solve problems caused by asynchronous accesses between the start and end of a nonblocking or one-sided communication specifically, problems caused by temporary memory modifications are not solved.

### The (evil) Fortran VOLATILE attribute

The VOLATILE attribute, gives the buffer or variable the properties needed, but it may inhibit optimization of any code containing references or definitions of the buffer or variable.

### The Fortran ASYNCHRONOUS attribute

Declaring a buffer with the ASYNCHRONOUS Fortran attribute in a scoping unit (or BLOCK) tells the compiler that any statement of the scoping unit may be executed while the buffer is affected by a pending asynchronous Fortran input/output operation. Because a Fortran compiler may implement asynchronous Fortran input/output operations with blocking I/O, the ASYNCHRONOUS attribute may be totally ignored by such a compiler. Therefore, the ASYNCHRONOUS attribute can **not** be used to solve the code movement and register optimization problem in a portable MPI program.

### 16.2.10    Temporary Data Movements and Temporary Memory Modifications

The compiler is allowed to modify temporarily data in memory. Normally, this problem may occur only when overlapping communication and computation. Example 16.15 shows a possibility that could be problematic.

In the compiler-generated, possible optimization in Example 16.16, buf(100,100) from Example 16.15 is equivalenced with the 1-dimensional array buf_1dim(10000). The nonblocking receive may asynchronously receive the data in the boundary buf(1,1:100) while the fused loop is temporarily using this part of the buffer. When the tmp data is written back to buf, the previous data of buf(1,1:100) is restored and the received data is lost. The principle behind this optimization is that the receive buffer data buf(1,1:100) was temporarily moved to tmp.

Example 16.17 shows a second possible optimization. The whole array is temporarily moved to local_buf. When storing local_buf back to the original location buf, then this includes also an overwriting of the receive buffer part buf(1,1:100), i.e., this storing back may overwrite the asynchronously received data.

Note, that this problem may also occurs:

ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.

**Unofficial Draft for Comment Only**

[ticket238-J.]

**Example 16.15** Overlapping Communication and Computation.

```
USE mpi_f08
REAL :: buf(100,100)
CALL MPI_Irecv(buf(1,1:100),...req,...)
DO j=1,100
  DO i=2,100
    buf(i,j)=....
  END DO
END DO
CALL MPI_Wait(req,...)
```

[ticket238-J.]

**Example 16.16** The compiler may substitute the nested loops through loop fusion.

```
REAL :: buf(100,100),  buf_1dim(10000)
EQUIVALENCE (buf(1,1), buf_1dim(1))
CALL MPI_Irecv(buf(1,1:100),...req,...)
DO h=1,100
  tmp(h)=buf(1,h)
END DO
DO j=1,10000
  buf_1dim(h)=...
END DO
DO h=1,100
  buf(1,h)=tmp(h)
END DO
CALL MPI_Wait(req,...)
```

- With the local buffer at the origin process, between an RMA call and the ensuing synchronization call.

- With the window buffer at the target process between two ensuing RMA synchronization calls,

- With the local buffer in MPI parallel file I/O split collective operations with between the ..._BEGIN and ..._END calls.

Note that such compiler optimization with temporary data movements can **not** be prevented when buf is declared with the ASYNCHRONOUS Fortran attribute.

*Rationale.*   Using the ASYNCHRONOUS attribute for a buffer while asynchronous operations are pending, the access is restricted as if the pending operation involves the whole buffer (see the *pending I/O storage sequence affector*, Section 5.3.4 and Section 9.6.4.1 paragraphs 5 and 6 of the Fortran 2008 standard [34]). Said this, it is invalid to use parts of an array in pending nonblocking operations and other parts in numerical

[ticket238-J.]

**Example 16.17** Another optimization is based on the usage of a local memory, e.g., in a GPU.

```
REAL :: buf(100,100), local_buf(100:100)
CALL MPI_Irecv(buf(1,1:100),...req,...)
local_buf = buf
DO j=1,100
  DO i=2,100
    local_buf(i,j)=....
  END DO
END DO
buf = local_buf ! may overwrite asynchronously received
                ! data in buf(1,1:100)
CALL MPI_Wait(req,...)
```

computation if the array is declared as `ASYNCHRONOUS`. Due to this reason and also due to the reason in Section *The Fortran `ASYNCHRONOUS` attribute* on page 558, the `ASYNCHRONOUS` attribute can **not** be used to solve temporary data movement problems even if the meaning of `ASYNCHRONOUS` would be extended within the Fortran standard from Fortran input/output to general nonblocking operation. (*End of rationale.*)

Note also that compiler optimization with temporary data movements should **not** be prevented by declaring `buf` as `VOLATILE` because: The `VOLATILE` implies that all accesses to any storage unit (word) of `buf` must be directly done in the main memory exactly in the sequence defined by the application program. The `VOLATILE` attribute prevents every register or cache optimization. Therefore, `VOLATILE` may cause a huge performance degradation.

Instead of solving the problem, it is better to **prevent** the problem, i.e., when overlapping communication and computation, the nonblocking communication (or nonblocking or split collective IO) and the computation should be executed **on different sets of variables**. In this case, the temporary memory modifications are done only on the variables used in the computation and cannot have any side effect on the data used in the nonblocking MPI operations.

*Rationale.* This is a strong restriction for application programs. To weaken this restriction, a new or modified asynchronous feature in the Fortran language would be necessary, e.g., an asynchronous attribute that can be used on parts of an array and together with asynchronous operations outside the scope of Fortran. If this would be available in a later version of the Fortran standard then this restriction also may be weakened in a later version of the MPI standard. (*End of rationale.*)

ticket238-J.
ticket238-J.

## 16.2.11 Permanent Data Movements

A Fortran compiler may implement permanent data movements during the execution of a Fortran program. This would require that pointers to such data are appropriately updated. An automatic garbage collection is a use case. Such permanent data movements are in conflict with MPI in several areas:

- MPI datatype handles with absolute addresses in combination with MPI_BOTTOM.

- Nonblocking MPI operations (communication, one-sided, I/O) if the internally used pointers to the buffers are not updated by the Fortran runtime, or if within an MPI process, the data movements are executed in parallel to the MPI operation.

This MPI standard requires that the problems with permanent data movements are solved by the MPI library together with the used compiler; see Section 16.2.16 on page 566.

### 16.2.12  Comparison with C

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. Problems due to temporary memory modifications can also occur in C. Also here, the best advice is to avoid the problem specifically, to use different variables for buffers in nonblocking MPI operations and computation that is executed while the nonblocking operations are pending.

### 16.2.13  Fortran Support Through the `mpif.h` Include File

The use of the `mpif.h` include file is strongly discouraged.

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 Fortran bindings are compatible with Fortran 77 implicit-style interfaces in most cases. The include file `mpif.h` must:

- Define all named MPI constants.

- Declare MPI functions that return a value.

- Define all handles as `INTEGER`. This is reflected in the first of the two Fortran interfaces in each MPI function definition throughout this document.

- Be valid and equivalent for both fixed- and free- source form.

For each MPI routine, an implementation can choose to use an implicit or explicit interface. With each MPI routine definition, two Fortran bindings are provided (in deprecated routines, the second one may be omitted). With explicit interfaces, the first of the two must be implemented.

> *Advice to users.*  Instead of using `mpif.h`, the use of the `mpi_f08` or `mpi` module is strongly encouraged. Reasons are:
>
> - Most `mpif.h` include files do not implement compile-time argument checking.
> - Therefore, too many bugs in MPI applications are still undetected:
>   – Missing `ierror` as last additional argument in most Fortran bindings.
>   – Declaration of a `status` as an `INTEGER` variable instead of an `INTEGER` array with size MPI_STATUS_SIZE.

1
2
3
4
5
6
7  ticket238-J.
8
9
10
11
12
13
14
15
16  ticket238-J.
17
18
19  ticket238-J.
20
21  ticket230-B.
22  ticket230-B.
23  ticket233-E.
24
25  ticket230-B.
26  ticket230-B.
27  ticket230-B.
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

– Wrong argument positions, e.g., interchanging the `count` and `datatype` arguments.

– Passing wrong MPI handles.

• The migration from `mpif.h` to the `mpi` module should be without problems (i.e., substituting `include 'mpif.h'` below an `implicit none` statement by `use mpi` before such `implicit` statement) as long as the application syntax is correct.

• Migrating portable applications to the `mpi` module, it is not expected to experience any compile or runtime problems because an `mpif.h` include file was always allowed to provide explicit Fortran interfaces.

(*End of advice to users.*)

*Rationale.*     With MPI-3.0, the `mpif.h` include file was not deprecated because of backward compatiblity reasons, and internally `mpif.h` and the `mpi` module can be implemented so that the same library implementation of the MPI routines can be used. (*End of rationale.*)

*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.*)

## 16.2.14   Fortran Support Through the mpi Module

An MPI implementation must provide a module named `mpi` that can be `use`d in a Fortran program. This module must:

• Define all named MPI constants

• Declare MPI functions that return a value.

• Provide explicit interfaces based on the first of the two Fortran bindings that are provided for all MPI routines, i.e., this module guarantees compile-time argument checking, and allows positional and keyword-based argument lists.

• Define all handles as `INTEGER`. This is reflected in the first of the two Fortran interfaces in each MPI function definition throughout this document.

• Define all named handle types  and `MPI_Status`  that are used in the `mpi_f08` module. They are needed only when the application converts old-style `INTEGER` handles into a new-style handles with a named type.

   An MPI implementation may provide other features in the `mpi` module that enhance the usability of MPI while maintaining adherence to the standard.  For example, it may provide `INTENT` information in these interface blocks.

ticket242-N.

*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, e.g., by defining the `INTENT` as provided in the `mpi_f08` bindings. (*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 intent of an `OUT` argument to be `INOUT`. (*End of rationale.*)

*Advice to implementors.* Some compilers allow to implement a choice buffer argument in the `mpi` module with the following explicit interface:

```
INTERFACE
  SUBROUTINE MPI_...(buf, ...)
    !DEC$ ATTRIBUTES NO_ARG_CHECK :: buf
    !$PRAGMA IGNORE_TKR buf
    !DIR$ IGNORE_TKR buf
    !IBM* IGNORE_TKR buf
    REAL, DIMENSION(*) :: buf
```

In this case, the `INTEGER` compile-time constant MPI_SUBARRAYS must be set to MPI_SUBARRAYS_NOT_SUPPORTED. Note, however, that it is explicitly allowed that the choice arguments can be implemented in the same way as with the `mpi_f08` module. In the case where the compiler does not provide such functionality, a set of overloaded functions may be used. See the paper of M. Hennecke [26]. (*End of advice to implementors.*)

### 16.2.15 Fortran Support Through the `mpi_f08` Module

An MPI implementation must provide a module named `mpi_f08` that can be used in a Fortran program. With this module, a second Fortran definition is added for each MPI routine, except for routines that are deprecated since MPI-3.0 or earlier. This module must:

- Define all named MPI constants.

- Declare MPI functions that return a value.

- Provide explicit interfaces for the second Fortran interface of all MPI routines, i.e., this module guarantees compile-time argument checking for all arguments which are not `TYPE(*)`.

- Provide explicit interfaces for the first Fortran interface if the second one is omitted; this is only the case in functions that are deprecated since MPI-3.0 or earlier.

1
2
3
4
5
6
7
8
9
10
11
12 ticket250-V.
13
14 ticket230-B.
15 ticket232-D.
16 ticket232-D.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33 ticket234-F.
34
35 ticket230-B.
36 ticket230-B.
37 ticket247-S.
38 ticket241-M.
39 ticket230-B.
40 ticket230-B.
41
42
43
44
45
46
47
48

ticket231-C.

*Advice to users.* It is strongly recommended to substitute calls to deprecated routines when changing from `mpif.h` or the `mpi` module to the `mpi_f08` module. (*End of advice to users.*)

ticket234-F.

- Define all handles with uniquely named handle types (instead of INTEGER handles in the `mpi` module). This is reflected in the second of the two Fortran interfaces in each MPI function definition throughout this document.

- Set the `INTEGER` compile-time constant `MPI_SUBARRAYS` to `MPI_SUBARRAYS_SUPPORTED` and declare choice buffers with the Fortran 2008 TR 29113 feature assumed-type and assumed-rank `TYPE(*), DIMENSION(..)` if the underlying Fortran compiler supports it. With this, the use of non-contiguous sub-arrays is valid also in nonblocking routines.

- Set the MPI_SUBARRAYS compile-time constant to MPI_SUBARRAYS_NOT_SUPPORTED and declare choice buffers with a compiler-dependent mechanism that overrides type checking if the underlying Fortran compiler does not support the Fortran 2008 TR 29113 assumed-type and assumed-rank notation. In this case, the use of non-contiguous sub-arrays in nonblocking calls may be restricted as with the `mpi` module.

ticket234-F.

*Advice to implementors.* In the MPI_SUBARRAYS_NOT_SUPPORTED case, the choice argument may be implemented with an explicit interface with compiler directives, for example:

```
INTERFACE
  SUBROUTINE MPI_...(buf, ...)
    !DEC$ ATTRIBUTES NO_ARG_CHECK :: buf
    !$PRAGMA IGNORE_TKR buf
    !DIR$ IGNORE_TKR buf
    !IBM* IGNORE_TKR buf
    REAL, DIMENSION(*) :: buf
```

(*End of advice to implementors.*)

ticket242-N.

- Declare each argument with an `INTENT=IN`, `OUT`, or `INOUT` as defined in this standard.

ticket242-N.

*Rationale.* For these definitions in the `mpi_f08` bindings, in most cases, `INTENT(IN)` is used if the C interface uses call-by-value. For all buffer arguments and for `OUT` and `INOUT` dummy arguments that allow one of the non-ordinary Fortran constants (see `MPI_BOTTOM`, etc. in Section 2.5.4 on page 15) as input, an `INTENT(...)` is not specified. (*End of rationale.*)

ticket244-P.

- Declare all `status` and `array_of_statuses` output arguments as optional through function overloading, instead of using `MPI_STATUS_IGNORE` and `MPI_STATUSES_IGNORE`. For this, two specific bindings are provided in a generic interface where one has the optional argument and one does not.

ticket244-P.

ticket244-P.

- Declare all `array_of_errcodes` output arguments as optional through function overloading, instead of using `MPI_ERRCODES_IGNORE`.

- Declare all `sourceweights`, `destweights`, and `weights` arguments as optional through function overloading, instead of using `MPI_UNWEIGHTED`.

  *Rationale.* To assure that the correspondence of the optional `ierror` argument can always be established by position, it is necessary that other arguments are not optional through using the `OPTIONAL` attribute. Therefore, function overloading is used for such other arguments. (*End of rationale.*)

- Declare all `ierror` output arguments as `OPTIONAL`, except for user-defined callback functions (e.g., `COMM_COPY_ATTR_FUNCTION`) and their predefined callbacks (e.g., `MPI_NULL_COPY_FN`).

  *Rationale.* For user-defined callback functions (e.g., `COMM_COPY_ATTR_FUNCTION`) and their predefined callbacks (e.g., `MPI_NULL_COPY_FN`), the ierror argument is not optional. The MPI library must always call these routines with an actual `ierror` argument. Therefore, these user-defined functions need not to check whether the MPI library calls these routines with or without an actual `ierror` output argument. (*End of rationale.*)

The MPI Fortran bindings in the `mpi_f08` module are designed based on the Fortran 2008 standard [34] together with the Technical Report (TR 29113) on Further Interoperability with C [35] of the ISO/IEC JTC1/SC22/WG5 (Fortran) working group.

  *Rationale.* The features in TR 29113 on further interoperability with C were decided on by ISO/IEC JTC1/SC22/WG5 and designed by PL22.3 (formerly J3) to support a higher level of integration between Fortran-specific features and C than provided in the Fortran 2008 standard; part of this design is based on requirements from the MPI Forum to support MPI-3.0. "It is the intention of ISO/IEC JTC1/SC22/WG5 that the semantics and syntax specified by this technical report be included in the next revision of the Fortran International Standard without change unless experience in the implementation and use of this feature identifies errors that need to be corrected, or changes are needed to achieve proper integration, in which case every reasonable effort will be made to minimize the impact of such changes on existing implementations."[3]

  This TR 29113 contains language features that are needed for the MPI bindings in the `mpi_f08` module: assumed-type and assumed-rank. It is important that any possible actual argument can be used for such dummy arguments, e.g., scalars, arrays, assumed-shape arrays, assumed-size arrays, allocatable arrays, and with any element type, e.g., `REAL`, `CHARACTER*5`, `CHARACTER*(*)`, sequence derived types, BIND(C) derived types. Especially for backward compatibility reasons, it is important that any possible actual argument in an implicit interface implementation of a choice buffer dummy argument can be used in an implementation with assumed-type and assumed-rank in an explicit interface.

_____
[3][36] page iv, paragraph 7.

1
2
3
4
5
6 ticket244-P.
7
8
9
10
11 ticket239-K.
12
13
14 ticket250-V.
15 ticket239-K.
16 ticket239-K.
17
18
19
20
21
22
23 ticket230-B.
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

Furthermore, the implementors of the MPI Fortran bindings can freely choose whether the interfaces are declared within a Fortran `INTERFACE` or `CONTAINS` construct, and whether all bindings are defined as native Fortran or a `BIND(C)` interface. The `INTERFACE` construct in combination with `BIND(C)` allows to implement the Fortran `mpi_f08` interface with a single set of portable wrapper routines written in C; to support all desired features in the `mpi_f08` interface, TR 29113 also has a provision for `OPTIONAL` arguments in `BIND(C)` interfaces.

The MPI Forum hereby wishes to acknowledge this important effort by the Fortran PL22.3 and WG5 committee. (*End of rationale.*)

### 16.2.16   Requirements on Fortran Compilers

The compliance to MPI-3.0 (and later) Fortran bindings is not only a property of the MPI library itself, but is always a property of an MPI library together with the Fortran compiler it is compiled for.

*Advice to users.* Many MPI libraries are shipped together with special compilation scripts (e.g., `mpif90`, `mpicc`). These scripts start the compiler probably together with special options to guarantee this compliance. (*End of advice to users.*)

An MPI library is only compliant with MPI-3.0 (and later), as referred by MPI_GET_VERSION, if all the solutions described in Sections 16.2.3 to 16.2.11 work correctly. Based on this rule, major requirements for all three Fortran support methods (i.e., the `mpi_f08` and `mpi` modules, and `mpif.h`) are:

- The language features assumed-type and assumed-rank from Fortran 2008 TR 29113 [35] are available. This is required only for `mpi_f08`. As long as this requirement is not supported by the compiler, it is valid to build a preliminary MPI-3.0 (and not later) library, which implements the `mpi_f08` module with MPI_SUBARRAYS set to MPI_SUBARRAYS_NOT_SUPPORTED.

- Simply contiguous arrays and scalars must be passed to choice buffer dummy arguments with call by reference.

- `SEQUENCE` and `BIND(C)` derived types are valid as actual arguments passed to choice buffer dummy arguments and they are passed with call by reference.

- The `TARGET` attribute (as described in Section *The Fortran TARGET attribute* on page 555) solves the problems described in Section 16.2.9 on page 552 independent of whether the MPI library internally uses Fortran or C pointers to memorize the location of a buffer between start and completion of a nonblocking operation, and to handle absolute and relative addresses in MPI derived datatype handles.

- A separately compiled empty routine (as MPI_F_SYNC_REG on page 556 and Section 16.2.17 on page 568, and DD on page 557) solves the problems described in Section 16.2.9 on page 552.

- The problems with temporary data movements (as described in Section 16.2.10 on page 558) are solved as long as the application uses different sets of variables for the nonblocking communication (or nonblocking or split collective IO) and the computation when overlapping communication and computation.

ticket230-B.
ticket230-B.
ticket234-F.
ticket230-B.
ticket237-I.
ticket238-J.
ticket238-J.
ticket238-J.
ticket238-J.

- Problems caused by automatic and permanent data movements (e.g., within a garbage collection, see Section 16.2.11 on page 560) are resolved **without** any further requirements on the application program, neither on the usage of the buffers, nor on the declaration of application routines that are involved in calling MPI operations.

- All actual arguments that are allowed for a dummy argument in an implicitly defined and separately compiled Fortran routine with the given compiler (e.g., `CHARACTER*(*)` strings and array of strings) must also be valid for choice buffer dummy arguments with all Fortran support methods.

- The handle and status types in `mpi_f08` (i.e., sequence derived types with `INTEGER` elements) are (handle) or can be (status) identical to one numerical storage unit or a sequence of those. These types must be valid at every location where an `INTEGER` and a fixed-size array of `INTEGER`s (i.e., handle and status in the `mpi` module and `mpif.h`) is valid, especially also within `BIND(C)` derived types defined by the application.

  *Rationale.* This is not yet part of the draft N1845 of TR 29113 [36], but may be part of the final version of this TR 29113 [35]. It is already implemented in some of the available Fortran compilers (e.g., ifort and pgi).

  A contrary definition of handles and statuses would have been `BIND(C)` derived types. With such definition, it would have been necessary that such `BIND(C)` derived types could be part of application-defined sequence derived types which is contratictory to the current rules about `BIND(C)` and Fortran storage units. (*End of rationale.*)

All of these rules are valid independently of whether the MPI routine interfaces in the `mpi_f08` and `mpi` modules are internally defined with an `INTERFACE` or `CONTAINS` construct, and with or without `BIND(C)`, and also when `mpif.h` uses explicit interfaces.

  *Advice to implementors.* Some of these rules are already part of the Fortran 2003 standard if the MPI interfaces are defined without `BIND(C)`. Additional compiler support may be necessary if `BIND(C)` is used. Some of these additional requirements are defined in the Fortran 2008 TR 29113 [35]. Some of these requirements for MPI-3.0 are beyond of TR 29113. (*End of advice to implementors.*)

Further requirements apply if the MPI library internally uses `BIND(C)` routine interfaces:

- Non-buffer arguments are `INTEGER`, `INTEGER(KIND=...)`, `LOGICAL`, `CHARACTER*(*)`, and sequence derived types (handles and status in `mpi_f08`) variables and arrays, and `EXTERNAL` routines; function results are `DOUBLE PRECISION`. All these types must be valid as dummy arguments in the `BIND(C)` MPI routine interfaces. When compiling an MPI application, the compiler should not issue warnings because these types may not be interoperable with an existing type in C. Some of these types are already valid in `BIND(C)` interfaces since Fortran 2003, some may be valid through TR 29113 (e.g., `CHARACTER*(*)` and sequence derived types (not yet part of draft N1845)).

- `OPTIONAL` dummy arguments are also valid within `BIND(C)` interfaces. This requirements is part of the TR 29113.

1
2
3
4 ticket232-D.
5 ticket234-F.
6
7
8
9 ticket231-C.
10 ticket243-O.
11
12
13
14
15
16
17
18
19
20
21
22
23
24 ticket230-B.
25
26
27
28
29
30
31
32
33
34 ticket230-B.
35
36
37 ticket231-C.
38 ticket230-B.
39
40
41
42
43 ticket231-C.
44 ticket230-B.
45 ticket239-K.
46
47
48

### 16.2.17   Additional Support for Fortran Register-Memory-Synchronization

As described in Section 16.2.9 on page 552, a dummy call is needed to tell the compiler that registers are to be flushed for a given buffer or that accesses to a buffer may not be moved across a given point in the sequence of execution. It has only a Fortran binding.

MPI_F_SYNC_REG(buf)

  INOUT     buf                                                        initial address of buffer (choice)


```
MPI_F_SYNC_REG(buf)
    <type> buf(*)
```

```
MPI_F_sync_reg(buf)
    TYPE(*), DIMENSION(..)  ::  buf
```

This routine is a no-operation. It must be compiled in the MPI library so that a Fortran compiler cannot detect in the module that the routine has an empty body. It is used only to tell the compiler that a cached register value of a variable or buffer should be flushed, i.e., stored back to the memory (when necessary) or invalidated.

> *Rationale.*   This function is not available in other languages because it would not be useful. This routine has no ierror return argument because there is no operation that can fail. (*End of rationale.*)

> *Advice to implementors.*    It is recommended to bind this routine to a C routine to minimize the risk that the Fortran compiler can learn that this routine is empty (and that the call to this routine can be removed as part of an optimization). It is explicitly allowed to implement this routine within the `mpi_f08` module according to the definition for the `mpi` module or `mpif.h` to circumvent the overhead of building the internal dope vector to handle the assumed-type, assumed-rank argument. (*End of advice to implementors.*)

> *Advice to users.*    If only a part of an array (e.g., defined by a subscript triplet) is used in a nonblocking routine, it is recommended to pass the whole array to MPI_F_SYNC_REG to minimize the overhead of this no-operation call. (*End of advice to users.*)

### 16.2.18   Additional Support for Fortran Numeric Intrinsic Types

The routines in this section are part of Extended Fortran Support described in Section 16.2.14.
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

*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.

ticket231-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. With the Fortran mpi module or the `mpif.h` include file, a Fortran handle is a Fortran INTEGER value that can be used in the following conversion functions. With the Fortran `mpi_f08` module, a Fortran handle is a sequence derived type that contains an INTEGER field named MPI_VAL. This INTEGER value can be used in the following conversion functions.

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 22.

```
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.

```
   // 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 16.20** 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 (with the mpi module or `mpif.h`) 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, respec-
tively.  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.

ticket250-V.
> *Advice to users.*   There exists no 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.*)

ticket243-O.

Using the mpi_f08 Fortran module, a status is declared as TYPE(MPI_Status). The C
type MPI_F_status can be used to pass a Fortran TYPE(MPI_Status) argument into a C
routine. Figure 16.1 illustrates all status conversion routines. Some are only available in C,
some in both C and Fortran.

```
int MPI_Status_f082c(MPI_F_status *f08_status, MPI_Status *c_status)
```

This C routine converts a Fortran mpi_f08 TYPE(MPI_Status) into a C MPI_Status.

```
int MPI_Status_c2f08(MPI_Status *c_status, MPI_F_status *f08_status)
```

This C routine converts a C MPI_Status into a Fortran mpi_f08 TYPE(MPI_Status).
Conversion between the two Fortran versions of a status can be done with:

MPI_STATUS_F2F08(f_status, f08_status)

| | | |
|---|---|---|
| IN | f_status | status object declared as array |
| OUT | f08_status | status object declared as named type |

```
int MPI_Status_f2f08(MPI_Fint *f_status, MPI_F_status *f08_status)
```
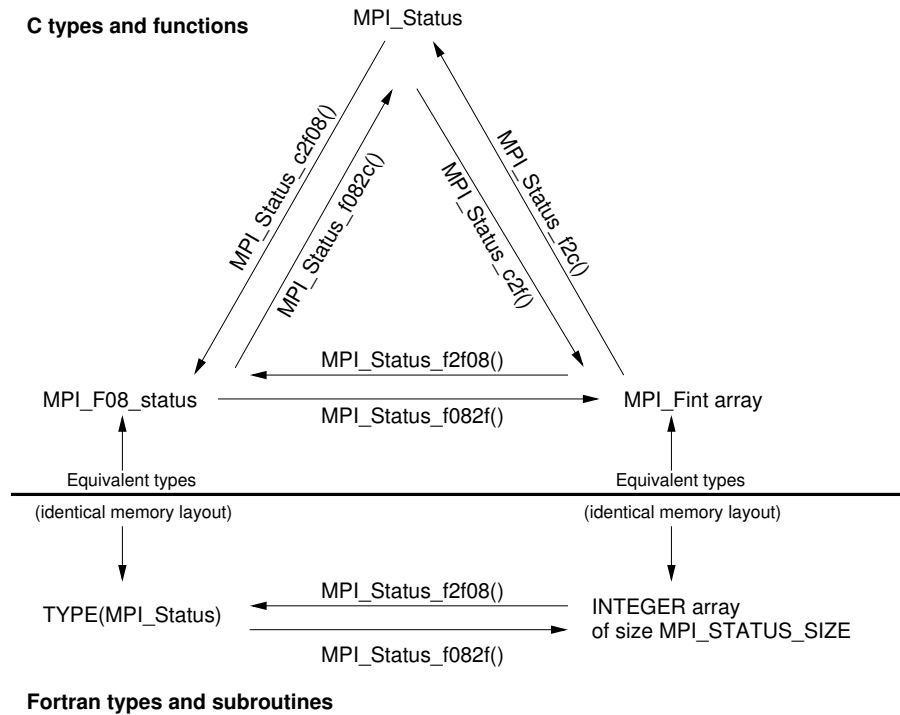
Figure 16.1: Status conversion routines

```
MPI_STATUS_F2F08(F_STATUS, F08_STATUS, IERROR)
    INTEGER ::  F_STATUS(MPI_STATUS_SIZE)
    TYPE(MPI_Status) ::  F08_STATUS
    INTEGER IERROR

MPI_Status_f2f08(f_status, f08_status, ierror)
    INTEGER, INTENT(IN) ::  f_status(MPI_STATUS_SIZE)
    TYPE(MPI_Status), INTENT(OUT) ::  f08_status
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

This routine converts a Fortran `INTEGER, DIMENSION(MPI_STATUS_SIZE)` status array into a Fortran `mpi_f08 TYPE(MPI_Status)`.

```
MPI_STATUS_F082F(f08_status, f_status)
```

|    |            |                                        |
|----|------------|----------------------------------------|
| IN | f08_status | status object declared as named type   |
| OUT| f_status   | status object declared as array        |

```
int MPI_Status_f082f(MPI_F_status *f08_status, MPI_Fint *f_status)

MPI_STATUS_F082F(F08_STATUS, F_STATUS, IERROR)
    TYPE(MPI_Status) ::  F08_STATUS
    INTEGER ::  F_STATUS(MPI_STATUS_SIZE)
    INTEGER IERROR

MPI_Status_f082f(f08_status, f_status, ierror)
```

*(margin line numbers 1–48, with "ticket-248T." noted at lines 28 and 47)*

```
       TYPE(MPI_Status), INTENT(IN) ::  f08_status
       INTEGER, INTENT(OUT) ::  f_status(MPI_STATUS_SIZE)
       INTEGER, OPTIONAL, INTENT(OUT) ::  ierror
```

This routine converts a Fortran `mpi_f08` TYPE(MPI_Status) into a Fortran INTEGER, DIMENSION(MPI_STATUS_SIZE) status array.

### 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 591).

**Example 16.21**

```
! 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];
```

**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 |
| [ticket234-F.]MPI_SUBARRAYS (Fortran only) | |
| [ticket234-F.]MPI_SUBARRAYS_SUPPORTED (Fortran only) | |
| [ticket234-F.]MPI_SUBARRAYS_NOT_SUPPORTED (Fortran only) | |

**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++ |
| 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` | |
| [ticket231-C.]or `TYPE(MPI_Errhandler)` | |
| MPI_ERRORS_ARE_FATAL | MPI::ERRORS_ARE_FATAL |
| MPI_ERRORS_RETURN | MPI::ERRORS_RETURN |
| | MPI::ERRORS_THROW_EXCEPTIONS |

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

**Null Handles**

| C/Fortran name | C++ name |
|---|---|
| C type / Fortran type | C++ type |
| MPI_GROUP_NULL | MPI::GROUP_NULL |
| MPI_Group / INTEGER | const MPI::Group |
| [ticket231-C.]   or TYPE(MPI_Group) | |
| MPI_COMM_NULL | MPI::COMM_NULL |
| MPI_Comm / INTEGER | [1]) |
| [ticket231-C.]   or TYPE(MPI_Comm) | |
| MPI_DATATYPE_NULL | MPI::DATATYPE_NULL |
| MPI_Datatype / INTEGER | const MPI::Datatype |
| [ticket231-C.]   or TYPE(MPI_Datatype) | |
| MPI_REQUEST_NULL | MPI::REQUEST_NULL |
| MPI_Request / INTEGER | const MPI::Request |
| [ticket231-C.]   or TYPE(MPI_Request) | |
| MPI_OP_NULL | MPI::OP_NULL |
| MPI_Op / INTEGER | const MPI::Op |
| [ticket231-C.]   or TYPE(MPI_Op) | |
| MPI_ERRHANDLER_NULL | MPI::ERRHANDLER_NULL |
| MPI_Errhandler / INTEGER | const MPI::Errhandler |
| [ticket231-C.]   or TYPE(MPI_Errhandler) | |
| MPI_FILE_NULL | MPI::FILE_NULL |
| MPI_File / INTEGER | |
| [ticket231-C.]   or TYPE(MPI_File) | |
| MPI_INFO_NULL | MPI::INFO_NULL |
| MPI_Info / INTEGER | const MPI::Info |
| [ticket231-C.]   or TYPE(MPI_Info) | |
| MPI_WIN_NULL | MPI::WIN_NULL |
| MPI_Win / INTEGER | |
| [ticket231-C.]   or TYPE(MPI_Win) | |

[1]) C++ type: See Section 16.1.7 on page 538 regarding
class hierarchy and the specific type of MPI::COMM_NULL

**Empty group**

| C type: MPI_Group | C++ type: const MPI::Group |
|---|---|
| Fortran type: INTEGER | |
| [ticket231-C.]or TYPE(MPI_Group) | |
| MPI_GROUP_EMPTY | MPI::GROUP_EMPTY |

**Topologies**

| C type: const int (or unnamed enum) | C++ type: const int |
|---|---|
| Fortran type: INTEGER | (or unnamed enum) |
| MPI_GRAPH | MPI::GRAPH |
| MPI_CART | MPI::CART |
| MPI_DIST_GRAPH | MPI::DIST_GRAPH |

**Unofficial Draft for Comment Only**

**Predefined functions**

| C/Fortran name | C++ name |
|---|---|
| C type / Fortran type [ticket230-B.]with `mpi` module | C++ type |
| [ticket230-B.]/ Fortran type with `mpi_f08` module | |
| MPI_COMM_NULL_COPY_FN | MPI_COMM_NULL_COPY_FN |
|    `MPI_Comm_copy_attr_function` | same as in C [1]) |
|    / `COMM_COPY_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Comm_copy_attr_function)` [2]) | |
| MPI_COMM_DUP_FN | MPI_COMM_DUP_FN |
|    `MPI_Comm_copy_attr_function` | same as in C [1]) |
|    / `COMM_COPY_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Comm_copy_attr_function)` [2]) | |
| MPI_COMM_NULL_DELETE_FN | MPI_COMM_NULL_DELETE_FN |
|    `MPI_Comm_delete_attr_function` | same as in C [1]) |
|    / `COMM_DELETE_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Comm_delete_attr_function)` [2]) | |
| MPI_WIN_NULL_COPY_FN | MPI_WIN_NULL_COPY_FN |
|    `MPI_Win_copy_attr_function` | same as in C [1]) |
|    / `WIN_COPY_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Win_copy_attr_function)` [2]) | |
| MPI_WIN_DUP_FN | MPI_WIN_DUP_FN |
|    `MPI_Win_copy_attr_function` | same as in C [1]) |
|    / `WIN_COPY_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Win_copy_attr_function)` [2]) | |
| MPI_WIN_NULL_DELETE_FN | MPI_WIN_NULL_DELETE_FN |
|    `MPI_Win_delete_attr_function` | same as in C [1]) |
|    / `WIN_DELETE_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Win_delete_attr_function)` [2]) | |
| MPI_TYPE_NULL_COPY_FN | MPI_TYPE_NULL_COPY_FN |
|    `MPI_Type_copy_attr_function` | same as in C [1]) |
|    / `TYPE_COPY_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Type_copy_attr_function)` [2]) | |
| MPI_TYPE_DUP_FN | MPI_TYPE_DUP_FN |
|    `MPI_Type_copy_attr_function` | same as in C [1]) |
|    / `TYPE_COPY_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Type_copy_attr_function)` [2]) | |
| MPI_TYPE_NULL_DELETE_FN | MPI_TYPE_NULL_DELETE_FN |
|    `MPI_Type_delete_attr_function` | same as in C [1]) |
|    / `TYPE_DELETE_ATTR_`[ticket250-V.]`FUNCTION` | |
|    / [ticket230-B.]`PROCEDURE(MPI_Type_delete_attr_function)` [2]) | |

[1] See the advice to implementors [ticket230-B.](on page 269) and advice to users (on page 269) on [ticket230-B.]the predefined C functions MPI_COMM_NULL_COPY_FN, ... in Section 6.7.2 on page 266

[ticket230-B.][2] See the advice to implementors (on page 269) and advice to users (on page 269) [ticket230-B.]    on the predefined Fortran functions MPI_COMM_NULL_COPY_FN, ... in [ticket230-B.]    Section 6.7.2 on page 266

**File Operation Constants, Part 2**

| C type: const int (or unnamed enum)<br>Fortran type: INTEGER | C++ type:<br>const int (or unnamed enum) |
|---|---|
| MPI_DISTRIBUTE_BLOCK | MPI::DISTRIBUTE_BLOCK |
| MPI_DISTRIBUTE_CYCLIC | MPI::DISTRIBUTE_CYCLIC |
| MPI_DISTRIBUTE_DFLT_DARG | MPI::DISTRIBUTE_DFLT_DARG |
| MPI_DISTRIBUTE_NONE | MPI::DISTRIBUTE_NONE |
| MPI_ORDER_C | MPI::ORDER_C |
| MPI_ORDER_FORTRAN | MPI::ORDER_FORTRAN |
| MPI_SEEK_CUR | MPI::SEEK_CUR |
| MPI_SEEK_END | MPI::SEEK_END |
| MPI_SEEK_SET | MPI::SEEK_SET |

**F90 Datatype Matching Constants**

| C type: const int (or unnamed enum)<br>Fortran type: INTEGER | C++ type:<br>const int (or unnamed enum) |
|---|---|
| MPI_TYPECLASS_COMPLEX | MPI::TYPECLASS_COMPLEX |
| MPI_TYPECLASS_INTEGER | MPI::TYPECLASS_INTEGER |
| MPI_TYPECLASS_REAL | MPI::TYPECLASS_REAL |

**Constants Specifying Empty or Ignored Input**

| C/Fortran name<br>C type / Fortran type [ticket244-P.]with mpi module<br>[ticket244-P.]/ Fortran type with mpi_f08 module | C++ name<br>C++ type |
|---|---|
| MPI_ARGVS_NULL<br>char*** / 2-dim. array of CHARACTER*(*)<br>[ticket244-P.]/ 2-dim. array of CHARACTER*(*) | MPI::ARGVS_NULL<br>const char *** |
| MPI_ARGV_NULL<br>char** / array of CHARACTER*(*)<br>[ticket244-P.]/ array of CHARACTER*(*) | MPI::ARGV_NULL<br>const char ** |
| MPI_ERRCODES_IGNORE<br>int* / INTEGER array<br>[ticket244-P.]/ not defined | Not defined for C++ |
| MPI_STATUSES_IGNORE<br>MPI_Status* / INTEGER, DIMENSION(MPI_STATUS_SIZE,*)<br>[ticket244-P.]/ not defined | Not defined for C++ |
| MPI_STATUS_IGNORE<br>MPI_Status* / INTEGER, DIMENSION(MPI_STATUS_SIZE)<br>[ticket244-P.]/ not defined | Not defined for C++ |
| MPI_UNWEIGHTED<br>int* / INTEGER array<br>[ticket244-P.]/ not defined | Not defined for C++ |

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

MPI::File

MPI::Group

MPI::Info

MPI::Op

MPI::Request

MPI::Prequest

MPI::Grequest

MPI::Win

The following are defined Fortran type definitions, included in the `mpi_f08` and `mpi` module.

```
!  Fortran opaque types in the mpi_f08 and mpi module
TYPE(MPI_Status)

!  Fortran handles in the mpi_f08 and mpi module
TYPE(MPI_Comm)
TYPE(MPI_Datatype)
TYPE(MPI_Errhandler)
TYPE(MPI_File)
TYPE(MPI_Group)
TYPE(MPI_Info)
TYPE(MPI_Op)
TYPE(MPI_Request)
TYPE(MPI_Win)
```

## A.1.3 Prototype Definitions

### C Bindings

The following are defined C typedefs for user-defined functions, also included in the file `mpi.h`.

```
/* prototypes for user-defined functions */
typedef void MPI_User_function(void *invec, void *inoutvec, int *len,
            MPI_Datatype *datatype);

typedef int MPI_Comm_copy_attr_function(MPI_Comm oldcomm,
            int comm_keyval, void *extra_state, void *attribute_val_in,
            void *attribute_val_out, int*flag);
typedef int MPI_Comm_delete_attr_function(MPI_Comm comm,
            int comm_keyval, void *attribute_val, void *extra_state);

typedef int MPI_Win_copy_attr_function(MPI_Win oldwin, int win_keyval,
            void *extra_state, void *attribute_val_in,
            void *attribute_val_out, int *flag);
typedef int MPI_Win_delete_attr_function(MPI_Win win, int win_keyval,
            void *attribute_val, void *extra_state);
```

ticket243-O.

ticket231-C.

ticket0.

ticket230-B.

```
typedef int MPI_Type_copy_attr_function(MPI_Datatype oldtype,
              int type_keyval, void *extra_state,
              void *attribute_val_in, void *attribute_val_out, int *flag);
typedef int MPI_Type_delete_attr_function(MPI_Datatype [ticket252-W.]datatype,
              int type_keyval, void *attribute_val, void *extra_state);

typedef void MPI_Comm_errhandler_function(MPI_Comm *, int *, ...);
typedef void MPI_Win_errhandler_function(MPI_Win *, int *, ...);
typedef void MPI_File_errhandler_function(MPI_File *, int *, ...);

typedef int MPI_Grequest_query_function(void *extra_state,
           MPI_Status *status);
typedef int MPI_Grequest_free_function(void *extra_state);
typedef int MPI_Grequest_cancel_function(void *extra_state, int complete);

typedef int MPI_Datarep_extent_function(MPI_Datatype datatype,
           MPI_Aint *file_extent, void *extra_state);
typedef int MPI_Datarep_conversion_function(void *userbuf,
           MPI_Datatype datatype, int count, void *filebuf,
           MPI_Offset position, void *extra_state);
```

ticket230-B.

### Fortran Bindings with mpif.h or the mpi Module

With the Fortran `mpi` module or `mpif.h`, here are examples of how each of the user-defined subroutines should be declared.

ticket230-B.

The user-function argument to MPI_OP_CREATE should be declared like this:

```
SUBROUTINE USER_FUNCTION(INVEC, INOUTVEC, LEN, [ticket252-W.]DATATYPE)
    <type> INVEC(LEN), INOUTVEC(LEN)
    INTEGER LEN, [ticket252-W.]DATATYPE
```

The copy and delete function arguments to MPI_COMM_CREATE_KEYVAL should be declared like these:

```
SUBROUTINE COMM_COPY_ATTR_[ticket250-V.]FUNCTION(OLDCOMM, COMM_KEYVAL, EXTRA_STATE,
              ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
    INTEGER OLDCOMM, COMM_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,
              ATTRIBUTE_VAL_OUT
    LOGICAL FLAG

SUBROUTINE COMM_DELETE_ATTR_[ticket250-V.]FUNCTION(COMM, COMM_KEYVAL, ATTRIBUTE_VAL,
              EXTRA_STATE, IERROR)
    INTEGER COMM, COMM_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

The copy and delete function arguments to MPI_WIN_CREATE_KEYVAL should be declared like these:

```
SUBROUTINE WIN_COPY_ATTR_[ticket250-V.]FUNCTION(OLDWIN, WIN_KEYVAL, EXTRA_STATE,
            ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
    INTEGER OLDWIN, WIN_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,
            ATTRIBUTE_VAL_OUT
    LOGICAL FLAG


SUBROUTINE WIN_DELETE_ATTR_[ticket250-V.]FUNCTION(WIN, WIN_KEYVAL, ATTRIBUTE_VAL,
            EXTRA_STATE, IERROR)
    INTEGER WIN, WIN_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

The copy and delete function arguments to MPI_TYPE_CREATE_KEYVAL should be declared like these:

```
SUBROUTINE TYPE_COPY_ATTR_[ticket250-V.]FUNCTION(OLDTYPE, TYPE_KEYVAL, EXTRA_STATE,
            ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
    INTEGER OLDTYPE, TYPE_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE,
            ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT
    LOGICAL FLAG

SUBROUTINE TYPE_DELETE_ATTR_[ticket250-V.]FUNCTION([ticket252-W.]DATATYPE, TYPE_KEYVAL, ATT
            EXTRA_STATE, IERROR)
    INTEGER [ticket252-W.]DATATYPE, TYPE_KEYVAL, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
```

The handler-function argument to MPI_COMM_CREATE_ERRHANDLER should be declared like this:

```
SUBROUTINE COMM_ERRHANDLER_FUNCTION(COMM, ERROR_CODE)
    INTEGER COMM, ERROR_CODE
```

The handler-function argument to MPI_WIN_CREATE_ERRHANDLER should be declared like this:

```
SUBROUTINE WIN_ERRHANDLER_FUNCTION(WIN, ERROR_CODE)
    INTEGER WIN, ERROR_CODE
```

The handler-function argument to MPI_FILE_CREATE_ERRHANDLER should be declared like this:

```
SUBROUTINE FILE_ERRHANDLER_FUNCTION(FILE, ERROR_CODE)
    INTEGER FILE, ERROR_CODE
```

The query, free, and cancel function arguments to MPI_GREQUEST_START should be declared like these:

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   SUBROUTINE GREQUEST_QUERY_FUNCTION(EXTRA_STATE, STATUS, IERROR)
²       INTEGER STATUS(MPI_STATUS_SIZE), IERROR
³       INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
⁴
⁵   SUBROUTINE GREQUEST_FREE_FUNCTION(EXTRA_STATE, IERROR)
⁶       INTEGER IERROR
⁷       INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
⁸
⁹   SUBROUTINE GREQUEST_CANCEL_FUNCTION(EXTRA_STATE, COMPLETE, IERROR)
¹⁰      INTEGER IERROR
¹¹      INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
¹²      LOGICAL COMPLETE
```

The extend and conversion function arguments to MPI_REGISTER_DATAREP should be declared like these:

```
SUBROUTINE DATAREP_EXTENT_FUNCTION(DATATYPE, EXTENT, EXTRA_STATE, IERROR)
    INTEGER DATATYPE, IERROR
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTENT, EXTRA_STATE

SUBROUTINE DATAREP_CONVERSION_FUNCTION(USERBUF, DATATYPE, COUNT, FILEBUF,
            POSITION, EXTRA_STATE, IERROR)
    <TYPE> USERBUF(*), FILEBUF(*)
    INTEGER COUNT, DATATYPE, IERROR
    INTEGER(KIND=MPI_OFFSET_KIND) POSITION
    INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE
```

ticket230-B.

### Fortran 2008 Bindings with the mpi_f08 Module

ticket230-B.

With the Fortran `mpi_f08` module, the callback prototypes are:

The user-function argument to MPI_Op_create should be declared like this:

```
            SUBROUTINE USER_FUNCTION(invec, inoutvec, len, datatype)
<type> ::  invec(len), inoutvec(len)
INTEGER ::  len
TYPE(MPI_Datatype) ::  datatype
```

ticket-248T.

The copy and delete function arguments to MPI_Comm_create_keyval should be declared according to:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Comm_copy_attr_function(oldcomm, comm_keyval, extra_state,
  attribute_val_in, attribute_val_out, flag, ierror)
      TYPE(MPI_Comm) ::  oldcomm
      INTEGER ::  comm_keyval, ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state, attribute_val_in,
      attribute_val_out
      LOGICAL flag
```

ticket-248T.

```
ABSTRACT INTERFACE
```

```
SUBROUTINE MPI_Comm_delete_attr_function(comm, comm_keyval,
attribute_val, extra_state, ierror)
    TYPE(MPI_Comm) ::  comm
    INTEGER ::  comm_keyval, ierror
    INTEGER(KIND=MPI_ADDRESS_KIND) ::  attribute_val, extra_state
```

The copy and delete function arguments to MPI_Win_create_keyval should be declared according to:
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Win_copy_attr_function(oldwin, win_keyval, extra_state,
  attribute_val_in, attribute_val_out, flag, ierror)
    TYPE(MPI_Win) ::  oldwin
    INTEGER ::  win_keyval, ierror
    INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state, attribute_val_in,
    attribute_val_out
    LOGICAL flag

ABSTRACT INTERFACE
  SUBROUTINE MPI_Win_delete_attr_function(win, win_keyval, attribute_val,
  extra_state, ierror)
    TYPE(MPI_Win) ::  win
    INTEGER ::  win_keyval, ierror
    INTEGER(KIND=MPI_ADDRESS_KIND) ::  attribute_val, extra_state
```

The copy and delete function arguments to MPI_Type_create_keyval should be declared according to:
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Type_copy_attr_function(oldtype, type_keyval, extra_state,
  attribute_val_in, attribute_val_out, flag, ierror)
    TYPE(MPI_Datatype) ::  oldtype
    INTEGER ::  type_keyval, ierror
    INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state, attribute_val_in,
    attribute_val_out
    LOGICAL flag

ABSTRACT INTERFACE
  SUBROUTINE MPI_Type_delete_attr_function(datatype, type_keyval,
  attribute_val, extra_state, ierror)
    TYPE(MPI_Datatype) ::  datatype
    INTEGER ::  type_keyval, ierror
    INTEGER(KIND=MPI_ADDRESS_KIND) ::  attribute_val, extra_state
```

The handler-function argument to MPI_Comm_create_errhandler should be declared like this:
```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Comm_errhandler_function(comm, error_code)
    TYPE(MPI_Comm) ::  comm
    INTEGER ::  error_code
```

1
2
3
4
5
6 ticket230-B.
7
8 ticket-248T.
9
10
11
12
13
14
15
16 ticket-248T.
17
18
19
20
21
22
23 ticket230-B.
24
25 ticket-248T.
26
27
28
29
30
31
32
33 ticket-248T.
34
35
36
37
38
39
40 ticket230-B.
41
42 ticket-248T.
43
44
45
46 ticket230-B.
47
48

The handler-function argument to MPI_Win_create_errhandler should be declared like

this:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Win_errhandler_function(win, error_code)
      TYPE(MPI_Win) ::  win
      INTEGER ::  error_code
```

The handler-function argument to MPI_File_create_errhandler should be declared like

this:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_File_errhandler_function(file, error_code)
      TYPE(MPI_File) ::  file
      INTEGER ::  error_code
```

The query, free, and cancel function arguments to MPI_Grequest_start should be de-

clared according to:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Grequest_query_function(extra_state, status, ierror)
      TYPE(MPI_Status) ::  status
      INTEGER ::  ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state
```

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Grequest_free_function(extra_state, ierror)
      INTEGER ::  ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state
```

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Grequest_cancel_function(extra_state, complete, ierror)
      INTEGER ::  ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extra_state
      LOGICAL complete
```

The extend and conversion function arguments to MPI_Register_datarep should be de-

clared according to:

```
ABSTRACT INTERFACE
  SUBROUTINE MPI_Datarep_extent_function(datatype, extent, extra_state,
  ierror)
      TYPE(MPI_Datatype) ::  datatype
      INTEGER ::  ierror
      INTEGER(KIND=MPI_ADDRESS_KIND) ::  extent, extra_state
```

and like these:

```
SUBROUTINE DATAREP_CONVERSION_FUNCTION(userbuf, datatype, count, filebuf,
              position, extra_state, ierror)
    <TYPE> userbuf(*), filebuf(*)
    TYPE(MPI_Datatype datatype
    INTEGER count, ierror
    INTEGER(KIND=MPI_OFFSET_KIND) position
    INTEGER(KIND=MPI_ADDRESS_KIND) extra_state
```

## A.4   Fortran 2008 Bindings with the mpi_f08 Module

### A.4.1   Point-to-Point Communication Fortran 2008 Bindings

```
MPI_Bsend(buf, count, datatype, dest, tag, comm, ierror)
    TYPE(*), DIMENSION(..)  ::  buf
    INTEGER, INTENT(IN) ::  count, dest, tag
    TYPE(MPI_Datatype), INTENT(IN) ::  datatype
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Bsend_init(buf, count, datatype, dest, tag, comm, request, ierror)
    TYPE(*), DIMENSION(..)  ::  buf
    INTEGER, INTENT(IN) ::  count, dest, tag
    TYPE(MPI_Datatype), INTENT(IN) ::  datatype
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Request), INTENT(OUT) ::  request
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Buffer_attach(buffer, size, ierror)
    TYPE(*), DIMENSION(..)  ::  buffer
    INTEGER, INTENT(IN) ::  size
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Buffer_detach(buffer_addr, size, ierror)
    TYPE(*), DIMENSION(..)  ::  buffer_addr
    INTEGER, INTENT(OUT) ::  size
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Cancel(request, ierror)
    TYPE(MPI_Request), INTENT(IN) ::  request
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Get_count(status, datatype, count, ierror)
    TYPE(MPI_Status), INTENT(IN) ::  status
    TYPE(MPI_Datatype), INTENT(IN) ::  datatype
    INTEGER, INTENT(OUT) ::  count
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Ibsend(buf, count, datatype, dest, tag, comm, request, ierror)
    TYPE(*), DIMENSION(..)  ::  buf
    INTEGER, INTENT(IN) ::  count, dest, tag
    TYPE(MPI_Datatype), INTENT(IN) ::  datatype
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Request), INTENT(OUT) ::  request
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Iprobe(source, tag, comm, flag, status, ierror)
    INTEGER, INTENT(IN) ::  source, tag
```

```fortran
      INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_delete_attr(comm, comm_keyval, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    INTEGER, INTENT(IN) ::  comm_keyval
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_dup(comm, newcomm, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Comm), INTENT(OUT) ::  newcomm
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_COMM_DUP_FN(oldcomm, comm_keyval, extra_state, attribute_val_in,
              attribute_val_out, flag, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  oldcomm
    INTEGER, INTENT(IN) ::  comm_keyval
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state,
    attribute_val_in
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(OUT) ::  attribute_val_out
    LOGICAL, INTENT(OUT) ::  flag
    INTEGER, INTENT(OUT) ::  ierror

MPI_Comm_free(comm, ierror)
    TYPE(MPI_Comm), INTENT(INOUT) ::  comm
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_free_keyval(comm_keyval, ierror)
    INTEGER, INTENT(INOUT) ::  comm_keyval
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_get_attr(comm, comm_keyval, attribute_val, flag, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    INTEGER, INTENT(IN) ::  comm_keyval
    INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(OUT) ::  attribute_val
    LOGICAL, INTENT(OUT) ::  flag
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_get_name(comm, comm_name, resultlen, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    CHARACTER(LEN=*), INTENT(OUT) ::  comm_name
    INTEGER, INTENT(OUT) ::  resultlen
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_group(comm, group, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  comm
    TYPE(MPI_Group), INTENT(OUT) ::  group
    INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_COMM_NULL_COPY_FN(oldcomm, comm_keyval, extra_state, attribute_val_in,
              attribute_val_out, flag, ierror)
    TYPE(MPI_Comm), INTENT(IN) ::  oldcomm
    INTEGER, INTENT(IN) ::  comm_keyval
```

```
     INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  extra_state,
     attribute_val_in
     INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(OUT) ::  attribute_val_out
     LOGICAL, INTENT(OUT) ::  flag
     INTEGER, INTENT(OUT) ::  ierror

MPI_COMM_NULL_DELETE_FN(comm, comm_keyval, attribute_val, extra_state,
            ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     INTEGER, INTENT(IN) ::  comm_keyval
     INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  attribute_val,
     extra_state
     INTEGER, INTENT(OUT) ::  ierror

MPI_Comm_rank(comm, rank, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     INTEGER, INTENT(OUT) ::  rank
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_remote_group(comm, group, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     TYPE(MPI_Group), INTENT(OUT) ::  group
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_remote_size(comm, size, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     INTEGER, INTENT(OUT) ::  size
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_set_attr(comm, comm_keyval, attribute_val, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     INTEGER, INTENT(IN) ::  comm_keyval
     INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) ::  attribute_val
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_set_name(comm, comm_name, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     CHARACTER(LEN=*), INTENT(IN) ::  comm_name
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_size(comm, size, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     INTEGER, INTENT(OUT) ::  size
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_split(comm, color, key, newcomm, ierror)
     TYPE(MPI_Comm), INTENT(IN) ::  comm
     INTEGER, INTENT(IN) ::  color, key
     TYPE(MPI_Comm), INTENT(OUT) ::  newcomm
     INTEGER, OPTIONAL, INTENT(OUT) ::  ierror

MPI_Comm_test_inter(comm, flag, ierror)
```

1
2
3
4
5
6 ticket-248T.
7
8
9
10
11
12
13 ticket-248T.
14
15
16
17
18 ticket-248T.
19
20
21
22 ticket-248T.
23
24
25
26
27 ticket-248T.
28
29
30
31
32 ticket-248T.
33
34
35
36
37 ticket-248T.
38
39
40
41 ticket-248T.
42
43
44
45
46
47 ticket-248T.
48

# Annex B

# Change-Log

This annex summarizes changes from the previous version of the MPI standard to the version presented by this document. Only significant changes (i.e., clarifications and new features) that might either require implementation effort in the MPI libraries or change the understanding of MPI from a user's perspective are presented. Editorial modifications, formatting, typo corrections and minor clarifications are not shown.

## B.1 Changes from Version 2.2 to Version 3.0

1. Section 2.3 on page 10, and Sections 16.2.1, 16.2.15, 16.2.16 on pages 544, 563, and 566.
   The new mpi_08 Fortran module is introduced.

2. Section 2.5.1 on page 12, Section 16.2.14 on page 562, Section 16.2.15 on page 563, and Section 16.2.16 on page 566.
   Handles to opaque objects are defined as named types within the mpi_08 Fortran module. The handle types are also available through the mpi Fortran module.

3. Section 2.5.2 on page 14, Section 3.2.6 on page 36, Section 7.5.4 on page 296, Section 7.5.5 on page 303, Section 10.3 on page 360, Section 10.3.3 on page 365, Section 12.2 on page 429, and Section 16.2.15 on page 563.
   With the mpi_f08 module, optional arguments through function overloading are used instead of MPI_STATUS_IGNORE, MPI_STATUSES_IGNORE, MPI_ERRCODES_IGNORE, and MPI_UNWEIGHTED.

4. Sections 2.5.4, 2.5.5 on pages 15, 16, Sections 16.2.1, 16.2.2, 16.2.3, 16.2.4, 16.2.5 on pages 544, 545, 546, 547, 550, and Sections 16.2.14, 16.2.15, 16.2.16 on pages 562, 563, 566.
   Within the mpi_08 Fortran module, choice buffers are defined as assumed-type and assumed-rank according to Fortran 2008, TR 29113 [35], and the compile-time constant MPI_SUBARRAYS is set to MPI_SUBARRAYS_SUPPORTED. With this, Fortran subscript triplets can be used in nonblocking MPI operations; vector subscripts are not supported in nonblocking operations. If the compiler does not support this Fortran TR 29113 feature, the constant is set to MPI_SUBARRAYS_NOT_SUPPORTED.

5. Section 2.6.2 on page 18, Section 16.2.15 on page 563, and Section 16.2.16 on page 566.
   The ierror dummy arguments are OPTIONAL within the mpi_08 Fortran module.

6. Section 3.2.5 on page 34, Section 16.2.14 on page 562, Section 16.2.15 on page 563, Section 16.2.16 on page 566, and Section 16.3.5 on page 581.
Within the `mpi_08` Fortran module, the status is defined as `TYPE(MPI_Status)`. New conversion routines are added: MPI_STATUS_F2F08, MPI_STATUS_F082F, MPI_Status_c2f08, and MPI_Status_f082c,

7. Sections 4.1.10, 5.9.5, 5.9.7, 6.7.4, 6.8, 8.3.1, 8.3.2, 8.3.3, 15.1, 16.2.18 on pages 111, 186, 191, 276, 282, 326, 328, 330, 523, and 568. In some routines, the dummy argument names were changed, because they were identical to the Fortran keywords `TYPE` and `FUNCTION`. The new dummy argument names must be used because the `mpi` and `mpi_08` modules guarantee keyword-based actual argument lists. The argument name type was changed into oldtype in MPI_TYPE_DUP, and into datatype in the Fortran USER_FUNCTION of MPI_OP_CREATE, and in MPI_TYPE_SET_ATTR, MPI_TYPE_GET_ATTR, MPI_TYPE_DELETE_ATTR, MPI_TYPE_SET_NAME, MPI_TYPE_GET_NAME, MPI_TYPE_MATCH_SIZE, in the callback prototype definition MPI_Type_delete_attr_function, and the predefined callback function MPI_TYPE_NULL_DELETE_FN; function was changed into user_fn in MPI_OP_CREATE, into comm_errhandler_fn in MPI_COMM_CREATE_ERRHANDLER, into win_errhandler_fn in MPI_WIN_CREATE_ERRHANDLER, into file_errhandler_fn in MPI_FILE_CREATE_ERRHANDLER, into handler_fn in MPI_ERRHANDLER_CREATE. For consistency reasons, INOUBUF was changed into INOUTBUF in MPI_REDUCE_LOCAL, and intracomm into newintracomm in MPI_INTERCOMM_MERGE.

8. Section 8.2 on page 322.
TODO: MPI_ALLOC_MEM and C pointer (C_PTR) instead of only with Cray pointer.

9. Section 16.2.7 on page 550, and Section 16.2.16 on page 566.
Fortran `SEQUENCE` and `BIND(C)` derived application types can be used as buffers in MPI operations.

10. Section 16.2.8 on page 551 to Section 16.2.11 on page 560, Section 16.2.16 on page 566, and Section 16.2.17 on page 568.
The sections about Fortran optimization problems and their solution is partially rewritten and new methods are added, e.g., the use of the `TARGET` attribute. The Fortran routine MPI_F_SYNC_REG is added. To achieve a secure and portable programming interfaces, in Section 16.2.16, several requirements are defined for the combination of an MPI library and a Fortran compiler to be MPI-3.0 compliant.

11. Section 16.2.13 on page 561.
The use of the `mpif.h` Fortran include file is strongly discouraged.

12. Section 16.2.14 on page 562, and Section 16.2.16 on page 566.
The existing `mpi` Fortran module must implement compile-time argument checking.

13. Section 16.2.15 on page 563.
Within the `mpi_08` Fortran module, dummy arguments are declared with `INTENT=IN`, `OUT`, or `INOUT` as defined in the `mpi_08` interfaces.

ticket230-B.

14. Section A.1.1, Table "*Predefined functions*" on page 602, Section A.1.3 on page 607, and Section A.4.4 on page 671.
Within the new `mpi_f08` module, all callback prototype definitions are defined with explicit interfaces `PROCEDURE(MPI_...)`, except for the `user_fn` argument in `MPI_OP_CREATE` and for the `read_conversion_fn` and `write_conversion_fn` arguments in `MPI_REGISTER_DATAREP`.

ticket250-V.

15. Section A.1.3 on page 607.
In some routines, the Fortran callback prototype names were changed from ..._`FN` to ..._`FUNCTION` to be consistent with the other language bindings.

## B.2   Changes from Version 2.1 to Version 2.2

1. Section 2.5.4 on page 15.
It is now guaranteed that predefined named constant handles (as other constants) can be used in initialization expressions or assignments, i.e., also before the call to `MPI_INIT`.

2. Section 2.6 on page 17, Section 2.6.4 on page 19, and Section 16.1 on page 531.
The C++ language bindings have been deprecated and may be removed in a future version of the MPI specification.

3. Section 3.2.2 on page 29.
`MPI_CHAR` for printable characters is now defined for C type char (instead of signed char). This change should not have any impact on applications nor on MPI libraries (except some comment lines), because printable characters could and can be stored in any of the C types char, signed char, and unsigned char, and `MPI_CHAR` is not allowed for predefined reduction operations.

4. Section 3.2.2 on page 29.
`MPI_(U)INT{8,16,32,64}_T`, `MPI_AINT`, `MPI_OFFSET`, `MPI_C_BOOL`, `MPI_C_COMPLEX`, `MPI_C_FLOAT_COMPLEX`, `MPI_C_DOUBLE_COMPLEX`, and `MPI_C_LONG_DOUBLE_COMPLEX` are now valid predefined MPI datatypes.

5. Section 3.4 on page 42, Section 3.7.2 on page 53, Section 3.9 on page 76, and Section 5.1 on page 143.
The read access restriction on the send buffer for blocking, non blocking and collective API has been lifted. It is permitted to access for read the send buffer while the operation is in progress.

6. Section 3.7 on page 52.
The Advice to users for IBSEND and IRSEND was slightly changed.

7. Section 3.7.3 on page 57.
The advice to free an active request was removed in the Advice to users for `MPI_REQUEST_FREE`.

8. Section 3.7.6 on page 70.
`MPI_REQUEST_GET_STATUS` changed to permit inactive or null requests as input.

**Unofficial Draft for Comment Only**

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

[26] Michael Hennecke. A Fortran 90 interface to MPI version 1.1. Technical Report Internal Report 63/96, Rechenzentrum, Universität Karlsruhe, D-76128 Karlsruhe, Germany, June 1996. Available via world wide web from http://www.uni-karlsruhe.de/~Michael.Hennecke/Publications/#MPI_F90. 16.2.14

[27] T. Hoefler, P. Gottschling, A. Lumsdaine, and W. Rehm. Optimizing a Conjugate Gradient Solver with Non-Blocking Collective Operations. *Elsevier Journal of Parallel Computing (PARCO)*, 33(9):624–633, Sep. 2007. 5.12

[28] T. Hoefler and A. Lumsdaine. Message Progression in Parallel Computing - To Thread or not to Thread? In *Proceedings of the 2008 IEEE International Conference on Cluster Computing*. IEEE Computer Society, Oct. 2008. 5.12

[29] T. Hoefler, A. Lumsdaine, and W. Rehm. Implementation and Performance Analysis of Non-Blocking Collective Operations for MPI. In *Proceedings of the 2007 International Conference on High Performance Computing, Networking, Storage and Analysis, SC07*. IEEE Computer Society/ACM, Nov. 2007. 5.12

[30] T. Hoefler, M. Schellmann, S. Gorlatch, and A. Lumsdaine. Communication Optimization for Medical Image Reconstruction Algorithms. In *Recent Advances in Parallel Virtual Machine and Message Passing Interface, 15th European PVM/MPI Users' Group Meeting*, volume LNCS 5205, pages 75–83. Springer, Sep. 2008. 5.12

[31] Institute of Electrical and Electronics Engineers, New York. *IEEE Standard for Binary Floating-Point Arithmetic, ANSI/IEEE Standard 754-1985*, 1985. 13.5.2

[32] International Organization for Standardization, Geneva, ISO 8859-1:1987. *Information processing — 8-bit single-byte coded graphic character sets — Part 1: Latin alphabet No. 1*, 1987. 13.5.2

[33] International Organization for Standardization, Geneva, ISO/IEC 9945-1:1996(E). *Information technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language]*, December 1996. 12.4, 13.2.1

[34] International Organization for Standardization, Geneva, ISO/IEC 1539-1:2010. *Information technology – Programming languages – Fortran – Part 1: Base language*, November 2010. 16.2.1, 16.2.10, 16.2.15

[35] International Organization for Standardization, ISO/IEC/SC22/WG5 (Fortran), Geneva, TR 29113. *TR on further interoperability with C*, MONTH NOT YET DEFINED 2011. http://www.nag.co.uk/sc22wg5/. 16.2.1, 16.2.1, 16.2.5, 16.2.15, 16.2.16, 4

[36] International Organization for Standardization, ISO/IEC/SC22/WG5 (Fortran), Geneva, TR 29113, Draft N1845. *TR on further interoperability with C*, March, 3 2011. http://www.nag.co.uk/sc22wg5/ and ftp://ftp.nag.co.uk/sc22wg5/N1801-N1850/N1845.pdf. 3, 16.2.16

[37] Charles H. Koelbel, David B. Loveman, Robert S. Schreiber, Guy L. Steele Jr., and Mary E. Zosel. *The High Performance Fortran Handbook*. MIT Press, 1993. 4.1.4