

D R A F T

Document for a Standard Message-Passing Interface

Message Passing Interface Forum

August 23, 2011

This work was supported in part by NSF and ARPA under NSF contract CDA-9115428 and Esprit under project HPC Standards (21111).

This is the result of a LaTeX run of a draft of a single chapter of the MPIF Final Report document.

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

Chapter 6

Groups, Contexts, Communicators, and Caching

6.1 Introduction

This chapter introduces MPI features that support the development of parallel libraries. Parallel libraries are needed to encapsulate the distracting complications inherent in parallel implementations of key algorithms. They help to ensure consistent correctness of such procedures, and provide a “higher level” of portability than MPI itself can provide. As such, libraries prevent each programmer from repeating the work of defining consistent data structures, data layouts, and methods that implement key algorithms (such as matrix operations). Since the best libraries come with several variations on parallel systems (different data layouts, different strategies depending on the size of the system or problem, or type of floating point), this too needs to be hidden from the user.

We refer the reader to [4] and [1] for further information on writing libraries in MPI, using the features described in this chapter.

6.1.1 Features Needed to Support Libraries

The key features needed to support the creation of robust parallel libraries are as follows:

- Safe communication space, that guarantees that libraries can communicate as they need to, without conflicting with communication extraneous to the library,
- Group scope for collective operations, that allow libraries to avoid unnecessarily synchronizing uninvolved processes (potentially running unrelated code),
- Abstract process naming to allow libraries to describe their communication in terms suitable to their own data structures and algorithms,
- The ability to “adorn” a set of communicating processes with additional user-defined attributes, such as extra collective operations. This mechanism should provide a means for the user or library writer effectively to extend a message-passing notation.

In addition, a unified mechanism or object is needed for conveniently denoting communication context, the group of communicating processes, to house abstract process naming, and to store adornments.

6.1.2 MPI's Support for Libraries

The corresponding concepts that MPI provides, specifically to support robust libraries, are as follows:

- **Contexts** of communication,
- **Groups** of processes,
- **Virtual topologies**,
- **Attribute caching**,
- **Communicators**.

Communicators (see [2, 3, 5]) encapsulate all of these ideas in order to provide the appropriate scope for all communication operations in MPI. Communicators are divided into two kinds: intra-communicators for operations within a single group of processes and inter-communicators for operations between two groups of processes.

Caching. Communicators (see below) provide a “caching” mechanism that allows one to associate new attributes with communicators, on a par with MPI built-in features. This can be used by advanced users to adorn communicators further, and by MPI to implement some communicator functions. For example, the virtual-topology functions described in Chapter 7 are likely to be supported this way.

Groups. Groups define an ordered collection of processes, each with a rank, and it is this group that defines the low-level names for inter-process communication (ranks are used for sending and receiving). Thus, groups define a scope for process names in point-to-point communication. In addition, groups define the scope of collective operations. Groups may be manipulated separately from communicators in MPI, but only communicators can be used in communication operations.

Intra-communicators. The most commonly used means for message passing in MPI is via intra-communicators. Intra-communicators contain an instance of a group, contexts of communication for both point-to-point and collective communication, and the ability to include virtual topology and other attributes. These features work as follows:

- **Contexts** provide the ability to have separate safe “universes” of message-passing in MPI. A context is akin to an additional tag that differentiates messages. The system manages this differentiation process. The use of separate communication contexts by distinct libraries (or distinct library invocations) insulates communication internal to the library execution from external communication. This allows the invocation of the library even if there are pending communications on “other” communicators, and avoids the need to synchronize entry or exit into library code. Pending point-to-point communications are also guaranteed not to interfere with collective communications within a single communicator.
- **Groups** define the participants in the communication (see above) of a communicator.

- A **virtual topology** defines a special mapping of the ranks in a group to and from a topology. Special constructors for communicators are defined in Chapter 7 to provide this feature. Intra-communicators as described in this chapter do not have topologies.
- **Attributes** define the local information that the user or library has added to a communicator for later reference.

Advice to users. The practice in many communication libraries is that there is a unique, predefined communication universe that includes all processes available when the parallel program is initiated; the processes are assigned consecutive ranks. Participants in a point-to-point communication are identified by their rank; a collective communication (such as broadcast) always involves all processes. This practice can be followed in MPI by using the predefined communicator `MPI_COMM_WORLD`. *Users who are satisfied with this practice can plug in `MPI_COMM_WORLD` wherever a communicator argument is required, and can consequently disregard the rest of this chapter. (End of advice to users.)*

Inter-communicators. The discussion has dealt so far with **intra-communication**: communication within a group. MPI also supports **inter-communication**: communication between two non-overlapping groups. When an application is built by composing several parallel modules, it is convenient to allow one module to communicate with another using local ranks for addressing within the second module. This is especially convenient in a client-server computing paradigm, where either client or server are parallel. The support of inter-communication also provides a mechanism for the extension of MPI to a dynamic model where not all processes are preallocated at initialization time. In such a situation, it becomes necessary to support communication across “universes.” Inter-communication is supported by objects called **inter-communicators**. These objects bind two groups together with communication contexts shared by both groups. For inter-communicators, these features work as follows:

- **Contexts** provide the ability to have a separate safe “universe” of message-passing between the two groups. A send in the local group is always a receive in the remote group, and vice versa. The system manages this differentiation process. The use of separate communication contexts by distinct libraries (or distinct library invocations) insulates communication internal to the library execution from external communication. This allows the invocation of the library even if there are pending communications on “other” communicators, and avoids the need to synchronize entry or exit into library code.
- A local and remote group specify the recipients and destinations for an inter-communicator.
- Virtual topology is undefined for an inter-communicator.
- As before, attributes cache defines the local information that the user or library has added to a communicator for later reference.

MPI provides mechanisms for creating and manipulating inter-communicators. They are used for point-to-point and collective communication in an related manner to intra-communicators. Users who do not need inter-communication in their applications can safely

1 ignore this extension. Users who require inter-communication between overlapping groups
 2 must layer this capability on top of MPI.

3 4 6.2 Basic Concepts

5
6 In this section, we turn to a more formal definition of the concepts introduced above.

7 8 6.2.1 Groups

9
10 A **group** is an ordered set of process identifiers (henceforth processes); processes are
 11 implementation-dependent objects. Each process in a group is associated with an inte-
 12 ger **rank**. Ranks are contiguous and start from zero. Groups are represented by opaque
 13 **group objects**, and hence cannot be directly transferred from one process to another. A
 14 group is used within a communicator to describe the participants in a communication “uni-
 15 verse” and to rank such participants (thus giving them unique names within that “universe”
 16 of communication).

17 There is a special pre-defined group: `MPI_GROUP_EMPTY`, which is a group with no
 18 members. The predefined constant `MPI_GROUP_NULL` is the value used for invalid group
 19 handles.

20
21 *Advice to users.* `MPI_GROUP_EMPTY`, which is a valid handle to an empty group,
 22 should not be confused with `MPI_GROUP_NULL`, which in turn is an invalid handle.
 23 The former may be used as an argument to group operations; the latter, which is
 24 returned when a group is freed, is not a valid argument. (*End of advice to users.*)

25
26 *Advice to implementors.* A group may be represented by a virtual-to-real process-
 27 address-translation table. Each communicator object (see below) would have a pointer
 28 to such a table.

29 Simple implementations of MPI will enumerate groups, such as in a table. However,
 30 more advanced data structures make sense in order to improve scalability and memory
 31 usage with large numbers of processes. Such implementations are possible with MPI.
 32 (*End of advice to implementors.*)

33 34 6.2.2 Contexts

35
36 A **context** is a property of communicators (defined next) that allows partitioning of the
 37 communication space. A message sent in one context cannot be received in another context.
 38 Furthermore, where permitted, collective operations are independent of pending point-to-
 39 point operations. Contexts are not explicit MPI objects; they appear only as part of the
 40 realization of communicators (below).

41
42 *Advice to implementors.* Distinct communicators in the same process have distinct
 43 contexts. A context is essentially a system-managed tag (or tags) needed to make
 44 a communicator safe for point-to-point and MPI-defined collective communication.
 45 Safety means that collective and point-to-point communication within one commu-
 46 nicator do not interfere, and that communication over distinct communicators don’t
 47 interfere.

A possible implementation for a context is as a supplemental tag attached to messages on send and matched on receive. Each intra-communicator stores the value of its two tags (one for point-to-point and one for collective communication). Communicator-generating functions use a collective communication to agree on a new group-wide unique context.

Analogously, in inter-communication, two context tags are stored per communicator, one used by group A to send and group B to receive, and a second used by group B to send and for group A to receive.

Since contexts are not explicit objects, other implementations are also possible. (*End of advice to implementors.*)

6.2.3 Intra-Communicators

Intra-communicators bring together the concepts of group and context. To support implementation-specific optimizations, and application topologies (defined in the next chapter, Chapter 7), communicators may also “cache” additional information (see Section 6.7). MPI communication operations reference communicators to determine the scope and the “communication universe” in which a point-to-point or collective operation is to operate.

Each communicator contains a group of valid participants; this group always includes the local process. The source and destination of a message is identified by process rank within that group.

For collective communication, the intra-communicator specifies the set of processes that participate in the collective operation (and their order, when significant). Thus, the communicator restricts the “spatial” scope of communication, and provides machine-independent process addressing through ranks.

Intra-communicators are represented by opaque **intra-communicator objects**, and hence cannot be directly transferred from one process to another.

6.2.4 Predefined Intra-Communicators

An initial intra-communicator `MPI_COMM_WORLD` of all processes the local process can communicate with after initialization (itself included) is defined once `MPI_INIT` or `MPI_INIT_THREAD` has been called. In addition, the communicator `MPI_COMM_SELF` is provided, which includes only the process itself.

The predefined constant `MPI_COMM_NULL` is the value used for invalid communicator handles.

In a static-process-model implementation of MPI, all processes that participate in the computation are available after MPI is initialized. For this case, `MPI_COMM_WORLD` is a communicator of all processes available for the computation; this communicator has the same value in all processes. In an implementation of MPI where processes can dynamically join an MPI execution, it may be the case that a process starts an MPI computation without having access to all other processes. In such situations, `MPI_COMM_WORLD` is a communicator incorporating all processes with which the joining process can immediately communicate. Therefore, `MPI_COMM_WORLD` may simultaneously represent disjoint groups in different processes.

All MPI implementations are required to provide the `MPI_COMM_WORLD` communicator. It cannot be deallocated during the life of a process. The group corresponding to this communicator does not appear as a pre-defined constant, but it may be accessed using

1 MPI_COMM_GROUP (see below). MPI does not specify the correspondence between the
 2 process rank in MPI_COMM_WORLD and its (machine-dependent) absolute address. Neither
 3 does MPI specify the function of the host process, if any. Other implementation-dependent,
 4 predefined communicators may also be provided.

6 6.3 Group Management

8 This section describes the manipulation of process groups in MPI. These operations are
 9 local and their execution does not require interprocess communication.

11 6.3.1 Group Accessors

14 MPI_GROUP_SIZE(group, size)

16 IN group group (handle)
 17 OUT size number of processes in the group (integer)

19 int MPI_Group_size(MPI_Group group, int *size)

21 MPI_GROUP_SIZE(GROUP, SIZE, IERROR)

22 INTEGER GROUP, SIZE, IERROR

23 {int MPI::Group::Get_size() const(*binding deprecated, see Section 15.2*) }

26 MPI_GROUP_RANK(group, rank)

28 IN group group (handle)
 29 OUT rank rank of the calling process in group, or
 30 MPI_UNDEFINED if the process is not a member (in-
 31 teger)

33 int MPI_Group_rank(MPI_Group group, int *rank)

35 MPI_GROUP_RANK(GROUP, RANK, IERROR)

36 INTEGER GROUP, RANK, IERROR

37 {int MPI::Group::Get_rank() const(*binding deprecated, see Section 15.2*) }

MPI_GROUP_TRANSLATE_RANKS (group1, n, ranks1, group2, ranks2)			1
IN	group1	group1 (handle)	2
IN	n	number of ranks in ranks1 and ranks2 arrays (integer)	3
IN	ranks1	array of zero or more valid ranks in group1	4
IN	group2	group2 (handle)	5
OUT	ranks2	array of corresponding ranks in group2, MPI_UNDEFINED when no correspondence exists.	6
			7
			8
			9

```
int MPI_Group_translate_ranks (MPI_Group group1, int n, int *ranks1,
                             MPI_Group group2, int *ranks2)
```

```
MPI_GROUP_TRANSLATE_RANKS(GROUP1, N, RANKS1, GROUP2, RANKS2, IERROR)
    INTEGER GROUP1, N, RANKS1(*), GROUP2, RANKS2(*), IERROR
```

```
{static void MPI::Group::Translate_ranks (const MPI::Group& group1, int n,
    const int ranks1[], const MPI::Group& group2,
    int ranks2[]) (binding deprecated, see Section 15.2) }
```

This function is important for determining the relative numbering of the same processes in two different groups. For instance, if one knows the ranks of certain processes in the group of MPI_COMM_WORLD, one might want to know their ranks in a subset of that group.

MPI_PROC_NULL is a valid rank for input to MPI_GROUP_TRANSLATE_RANKS, which returns MPI_PROC_NULL as the translated rank.

MPI_GROUP_COMPARE(group1, group2, result)			26
IN	group1	first group (handle)	27
IN	group2	second group (handle)	28
OUT	result	result (integer)	29
			30
			31

```
int MPI_Group_compare(MPI_Group group1, MPI_Group group2, int *result)
```

```
MPI_GROUP_COMPARE(GROUP1, GROUP2, RESULT, IERROR)
    INTEGER GROUP1, GROUP2, RESULT, IERROR
```

```
{static int MPI::Group::Compare(const MPI::Group& group1,
    const MPI::Group& group2) (binding deprecated, see Section 15.2) }
```

MPI_IDENT results if the group members and group order is exactly the same in both groups. This happens for instance if group1 and group2 are the same handle. MPI_SIMILAR results if the group members are the same but the order is different. MPI_UNEQUAL results otherwise.

6.3.2 Group Constructors

Group constructors are used to subset and superset existing groups. These constructors construct new groups from existing groups. These are local operations, and distinct groups may be defined on different processes; a process may also define a group that does not include itself. Consistent definitions are required when groups are used as arguments in

communicator-building functions. MPI does not provide a mechanism to build a group from scratch, but only from other, previously defined groups. The base group, upon which all other groups are defined, is the group associated with the initial communicator MPI_COMM_WORLD (accessible through the function MPI_COMM_GROUP).

Rationale. In what follows, there is no group duplication function analogous to MPI_COMM_DUP, defined later in this chapter. There is no need for a group duplicator. A group, once created, can have several references to it by making copies of the handle. The following constructors address the need for subsets and supersets of existing groups. (*End of rationale.*)

Advice to implementors. Each group constructor behaves as if it returned a new group object. When this new group is a copy of an existing group, then one can avoid creating such new objects, using a reference-count mechanism. (*End of advice to implementors.*)

```
MPI_COMM_GROUP(comm, group)
```

```
IN      comm      communicator (handle)
```

```
OUT     group     group corresponding to comm (handle)
```

```
int MPI_Comm_group(MPI_Comm comm, MPI_Group *group)
```

```
MPI_COMM_GROUP(COMM, GROUP, IERROR)
```

```
INTEGER COMM, GROUP, IERROR
```

```
{MPI::Group MPI::Comm::Get_group() const(binding deprecated, see Section 15.2) }
```

MPI_COMM_GROUP returns in group a handle to the group of comm.

```
MPI_GROUP_UNION(group1, group2, newgroup)
```

```
IN      group1    first group (handle)
```

```
IN      group2    second group (handle)
```

```
OUT     newgroup  union group (handle)
```

```
int MPI_Group_union(MPI_Group group1, MPI_Group group2,
```

```
                MPI_Group *newgroup)
```

```
MPI_GROUP_UNION(GROUP1, GROUP2, NEWGROUP, IERROR)
```

```
INTEGER GROUP1, GROUP2, NEWGROUP, IERROR
```

```
{static MPI::Group MPI::Group::Union(const MPI::Group& group1,
```

```
                const MPI::Group& group2)(binding deprecated, see Section 15.2) }
```

```

MPI_GROUP_INTERSECTION(group1, group2, newgroup) 1
    IN      group1          first group (handle) 2
    IN      group2          second group (handle) 3
    OUT     newgroup        intersection group (handle) 4
                                                    5
                                                    6
int MPI_Group_intersection(MPI_Group group1, MPI_Group group2, 7
                          MPI_Group *newgroup) 8
MPI_GROUP_INTERSECTION(GROUP1, GROUP2, NEWGROUP, IERROR) 9
    INTEGER GROUP1, GROUP2, NEWGROUP, IERROR 10
{static MPI::Group MPI::Group::Intersect(const MPI::Group& group1, 11
                                          const MPI::Group& group2) (binding deprecated, see Section 15.2) } 12
                                                    13
                                                    14
                                                    15
MPI_GROUP_DIFFERENCE(group1, group2, newgroup) 16
    IN      group1          first group (handle) 17
    IN      group2          second group (handle) 18
    OUT     newgroup        difference group (handle) 19
                                                    20
                                                    21
int MPI_Group_difference(MPI_Group group1, MPI_Group group2, 22
                          MPI_Group *newgroup) 23
MPI_GROUP_DIFFERENCE(GROUP1, GROUP2, NEWGROUP, IERROR) 24
    INTEGER GROUP1, GROUP2, NEWGROUP, IERROR 25
{static MPI::Group MPI::Group::Difference(const MPI::Group& group1, 26
                                           const MPI::Group& group2) (binding deprecated, see Section 15.2) } 27
                                                    28
                                                    29
The set-like operations are defined as follows: 30
                                                    31
union All elements of the first group (group1), followed by all elements of second group 32
        (group2) not in first. 33
intersect all elements of the first group that are also in the second group, ordered as in 34
        first group. 35
difference all elements of the first group that are not in the second group, ordered as in 36
        the first group. 37
Note that for these operations the order of processes in the output group is determined 38
primarily by order in the first group (if possible) and then, if necessary, by order in the 39
second group. Neither union nor intersection are commutative, but both are associative. 40
The new group can be empty, that is, equal to MPI_GROUP_EMPTY. 41
                                                    42
                                                    43
                                                    44
                                                    45
                                                    46
                                                    47
                                                    48

```

1 MPI_GROUP_INCL(group, n, ranks, newgroup)

2	IN	group	group (handle)
3			
4	IN	n	number of elements in array ranks (and size of
5			newgroup) (integer)
6	IN	ranks	ranks of processes in group to appear in
7			newgroup (array of integers)
8	OUT	newgroup	new group derived from above, in the order defined by
9			ranks (handle)

10
11 int MPI_Group_incl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup)

12
13 MPI_GROUP_INCL(GROUP, N, RANKS, NEWGROUP, IERROR)
14 INTEGER GROUP, N, RANKS(*), NEWGROUP, IERROR

15 {MPI::Group MPI::Group::Incl(int n, const int ranks[]) const (*binding*
16 *deprecated, see Section 15.2*) }

17
18 The function MPI_GROUP_INCL creates a group `newgroup` that consists of the
19 `n` processes in `group` with ranks `rank[0], ..., rank[n-1]`; the process with rank `i` in `newgroup`
20 is the process with rank `ranks[i]` in `group`. Each of the `n` elements of `ranks` must be a valid
21 rank in `group` and all elements must be distinct, or else the program is erroneous. If `n = 0`,
22 then `newgroup` is `MPI_GROUP_EMPTY`. This function can, for instance, be used to reorder
23 the elements of a group. See also `MPI_GROUP_COMPARE`.

24
25
26 MPI_GROUP_EXCL(group, n, ranks, newgroup)

27	IN	group	group (handle)
28			
29	IN	n	number of elements in array ranks (integer)
30	IN	ranks	array of integer ranks in group not to appear in
31			newgroup
32	OUT	newgroup	new group derived from above, preserving the order
33			defined by group (handle)

34
35 int MPI_Group_excl(MPI_Group group, int n, int *ranks, MPI_Group *newgroup)

36
37 MPI_GROUP_EXCL(GROUP, N, RANKS, NEWGROUP, IERROR)
38 INTEGER GROUP, N, RANKS(*), NEWGROUP, IERROR

39 {MPI::Group MPI::Group::Excl(int n, const int ranks[]) const (*binding*
40 *deprecated, see Section 15.2*) }

41
42 The function MPI_GROUP_EXCL creates a group of processes `newgroup` that is obtained
43 by deleting from `group` those processes with ranks `ranks[0] ... ranks[n-1]`. The ordering of
44 processes in `newgroup` is identical to the ordering in `group`. Each of the `n` elements of `ranks`
45 must be a valid rank in `group` and all elements must be distinct; otherwise, the program is
46 erroneous. If `n = 0`, then `newgroup` is identical to `group`.

MPI_GROUP_RANGE_INCL(group, n, ranges, newgroup)	1
IN group	2
IN n	3
IN ranges	4
	5
	6
	7
OUT newgroup	8
	9
	10
int MPI_Group_range_incl(MPI_Group group, int n, int ranges[][3],	11
MPI_Group *newgroup)	12
	13
MPI_GROUP_RANGE_INCL(GROUP, N, RANGES, NEWGROUP, IERROR)	14
INTEGER GROUP, N, RANGES(3,*), NEWGROUP, IERROR	15
{MPI::Group MPI::Group::Range_incl(int n, const int ranges[][3])	16
const(binding deprecated, see Section 15.2) }	17
	18
If ranges consist of the triplets	19
(<i>first</i> ₁ , <i>last</i> ₁ , <i>stride</i> ₁), ..., (<i>first</i> _{<i>n</i>} , <i>last</i> _{<i>n</i>} , <i>stride</i> _{<i>n</i>})	20
	21
then newgroup consists of the sequence of processes in group with ranks	22
<i>first</i> ₁ , <i>first</i> ₁ + <i>stride</i> ₁ , ..., <i>first</i> ₁ + $\left\lfloor \frac{\textit{last}_1 - \textit{first}_1}{\textit{stride}_1} \right\rfloor \textit{stride}_1, \dots$	23
	24
<i>first</i> _{<i>n</i>} , <i>first</i> _{<i>n</i>} + <i>stride</i> _{<i>n</i>} , ..., <i>first</i> _{<i>n</i>} + $\left\lfloor \frac{\textit{last}_n - \textit{first}_n}{\textit{stride}_n} \right\rfloor \textit{stride}_n$.	25
	26
	27
Each computed rank must be a valid rank in group and all computed ranks must be	28
distinct, or else the program is erroneous. Note that we may have <i>first</i> _{<i>i</i>} > <i>last</i> _{<i>i</i>} , and <i>stride</i> _{<i>i</i>}	29
may be negative, but cannot be zero.	30
The functionality of this routine is specified to be equivalent to expanding the array	31
of ranges to an array of the included ranks and passing the resulting array of ranks and	32
other arguments to MPI_GROUP_INCL. A call to MPI_GROUP_INCL is equivalent to a call	33
to MPI_GROUP_RANGE_INCL with each rank <i>i</i> in ranks replaced by the triplet (<i>i</i> , <i>i</i> , 1) in	34
the argument ranges.	35
	36
MPI_GROUP_RANGE_EXCL(group, n, ranges, newgroup)	37
IN group	38
IN n	39
IN ranges	40
	41
	42
	43
	44
OUT newgroup	45
	46
	47
	48

```

1 int MPI_Group_range_excl(MPI_Group group, int n, int ranges[][3],
2     MPI_Group *newgroup)
3 MPI_GROUP_RANGE_EXCL(GROUP, N, RANGES, NEWGROUP, IERROR)
4     INTEGER GROUP, N, RANGES(3,*), NEWGROUP, IERROR
5
6 {MPI::Group MPI::Group::Range_excl(int n, const int ranges[][3])
7     const(binding deprecated, see Section 15.2) }
8

```

Each computed rank must be a valid rank in `group` and all computed ranks must be distinct, or else the program is erroneous.

The functionality of this routine is specified to be equivalent to expanding the array of ranges to an array of the excluded ranks and passing the resulting array of ranks and other arguments to `MPI_GROUP_EXCL`. A call to `MPI_GROUP_EXCL` is equivalent to a call to `MPI_GROUP_RANGE_EXCL` with each rank `i` in `ranks` replaced by the triplet `(i,i,1)` in the argument `ranges`.

Advice to users. The range operations do not explicitly enumerate ranks, and therefore are more scalable if implemented efficiently. Hence, we recommend MPI programmers to use them whenever possible, as high-quality implementations will take advantage of this fact. (*End of advice to users.*)

Advice to implementors. The range operations should be implemented, if possible, without enumerating the group members, in order to obtain better scalability (time and space). (*End of advice to implementors.*)

6.3.3 Group Destructors

```

28 MPI_GROUP_FREE(group)
29     INOUT    group                group (handle)
30
31
32 int MPI_Group_free(MPI_Group *group)
33 MPI_GROUP_FREE(GROUP, IERROR)
34     INTEGER GROUP, IERROR
35
36 {void MPI::Group::Free() (binding deprecated, see Section 15.2) }
37

```

This operation marks a group object for deallocation. The handle `group` is set to `MPI_GROUP_NULL` by the call. Any on-going operation using this group will complete normally.

Advice to implementors. One can keep a reference count that is incremented for each call to `MPI_COMM_GROUP`, `MPI_COMM_CREATE` and `MPI_COMM_DUP`, and decremented for each call to `MPI_GROUP_FREE` or `MPI_COMM_FREE`; the group object is ultimately deallocated when the reference count drops to zero. (*End of advice to implementors.*)

6.4 Communicator Management

This section describes the manipulation of communicators in MPI. Operations that access communicators are local and their execution does not require interprocess communication. Operations that create communicators are collective and may require interprocess communication.

Advice to implementors. High-quality implementations should amortize the overheads associated with the creation of communicators (for the same group, or subsets thereof) over several calls, by allocating multiple contexts with one collective communication. (*End of advice to implementors.*)

6.4.1 Communicator Accessors

The following are all local operations.

`MPI_COMM_SIZE(comm, size)`

IN	comm	communicator (handle)
OUT	size	number of processes in the group of comm (integer)

`int MPI_Comm_size(MPI_Comm comm, int *size)`

`MPI_COMM_SIZE(COMM, SIZE, IERROR)`
`INTEGER COMM, SIZE, IERROR`

`{int MPI::Comm::Get_size() const(binding deprecated, see Section 15.2) }`

Rationale. This function is equivalent to accessing the communicator's group with `MPI_COMM_GROUP` (see above), computing the size using `MPI_GROUP_SIZE`, and then freeing the temporary group via `MPI_GROUP_FREE`. However, this function is so commonly used, that this shortcut was introduced. (*End of rationale.*)

Advice to users. This function indicates the number of processes involved in a communicator. For `MPI_COMM_WORLD`, it indicates the total number of processes available (for this version of MPI, there is no standard way to change the number of processes once initialization has taken place).

This call is often used with the next call to determine the amount of concurrency available for a specific library or program. The following call, `MPI_COMM_RANK` indicates the rank of the process that calls it in the range from `0 . . . size-1`, where `size` is the return value of `MPI_COMM_SIZE`. (*End of advice to users.*)

`MPI_COMM_RANK(comm, rank)`

IN	comm	communicator (handle)
OUT	rank	rank of the calling process in group of comm (integer)

```

1  int MPI_Comm_rank(MPI_Comm comm, int *rank)
2
3  MPI_COMM_RANK(COMM, RANK, IERROR)
4      INTEGER COMM, RANK, IERROR
5
6  {int MPI::Comm::Get_rank() const(binding deprecated, see Section 15.2) }

```

Rationale. This function is equivalent to accessing the communicator’s group with MPI_COMM_GROUP (see above), computing the rank using MPI_GROUP_RANK, and then freeing the temporary group via MPI_GROUP_FREE. However, this function is so commonly used, that this shortcut was introduced. (*End of rationale.*)

Advice to users. This function gives the rank of the process in the particular communicator’s group. It is useful, as noted above, in conjunction with MPI_COMM_SIZE.

Many programs will be written with the master-slave model, where one process (such as the rank-zero process) will play a supervisory role, and the other processes will serve as compute nodes. In this framework, the two preceding calls are useful for determining the roles of the various processes of a communicator. (*End of advice to users.*)

```

21 MPI_COMM_COMPARE(comm1, comm2, result)
22     IN      comm1          first communicator (handle)
23     IN      comm2          second communicator (handle)
24     OUT     result         result (integer)
25
26
27 int MPI_Comm_compare(MPI_Comm comm1, MPI_Comm comm2, int *result)
28
29 MPI_COMM_COMPARE(COMM1, COMM2, RESULT, IERROR)
30     INTEGER COMM1, COMM2, RESULT, IERROR
31
32 {static int MPI::Comm::Compare(const MPI::Comm& comm1,
33     const MPI::Comm& comm2)(binding deprecated, see Section 15.2) }

```

MPI_IDENT results if and only if comm1 and comm2 are handles for the same object (identical groups and same contexts). MPI_CONGRUENT results if the underlying groups are identical in constituents and rank order; these communicators differ only by context. MPI_SIMILAR results if the group members of both communicators are the same but the rank order differs. MPI_UNEQUAL results otherwise.

6.4.2 Communicator Constructors

The following are collective functions that are invoked by all processes in the group or groups associated with comm.

Rationale. Note that there is a chicken-and-egg aspect to MPI in that a communicator is needed to create a new communicator. The base communicator for all MPI communicators is predefined outside of MPI, and is MPI_COMM_WORLD. This model was arrived at after considerable debate, and was chosen to increase “safety” of programs written in MPI. (*End of rationale.*)

The MPI interface provides four communicator construction routines that apply to both intracommunicators and intercommunicators. The construction routine `MPI_INTERCOMM_CREATE` (discussed later) applies only to intercommunicators.

An intracommunicator involves a single group while an intercommunicator involves two groups. Where the following discussions address intercommunicator semantics, the two groups in an intercommunicator are called the *left* and *right* groups. A process in an intercommunicator is a member of either the left or the right group. From the point of view of that process, the group that the process is a member of is called the *local* group; the other group (relative to that process) is the *remote* group. The left and right group labels give us a way to describe the two groups in an intercommunicator that is not relative to any particular process (as the local and remote groups are).

`MPI_COMM_DUP(comm, newcomm)`

IN	<code>comm</code>	communicator (handle)
OUT	<code>newcomm</code>	copy of <code>comm</code> (handle)

`int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)`

`MPI_COMM_DUP(COMM, NEWCOMM, IERROR)`

INTEGER COMM, NEWCOMM, IERROR

{MPI::Intracomm MPI::Intracomm::Dup() const(*binding deprecated, see Section 15.2*) }

{MPI::Intercomm MPI::Intercomm::Dup() const(*binding deprecated, see Section 15.2*) }

{MPI::Cartcomm MPI::Cartcomm::Dup() const(*binding deprecated, see Section 15.2*) }

{MPI::Graphcomm MPI::Graphcomm::Dup() const(*binding deprecated, see Section 15.2*) }

{MPI::Distgraphcomm MPI::Distgraphcomm::Dup() const(*binding deprecated, see Section 15.2*) }

{MPI::Comm& MPI::Comm::Clone() const = 0(*binding deprecated, see Section 15.2*) }

{MPI::Intracomm& MPI::Intracomm::Clone() const(*binding deprecated, see Section 15.2*) }

{MPI::Intercomm& MPI::Intercomm::Clone() const(*binding deprecated, see Section 15.2*) }

{MPI::Cartcomm& MPI::Cartcomm::Clone() const(*binding deprecated, see Section 15.2*) }

{MPI::Graphcomm& MPI::Graphcomm::Clone() const(*binding deprecated, see Section 15.2*) }

{MPI::Distgraphcomm& MPI::Distgraphcomm::Clone() const(*binding deprecated, see Section 15.2*) }

1 MPI_COMM_DUP Duplicates the existing communicator `comm` with associated key val-
 2 ues. For each key value, the respective copy callback function determines the attribute value
 3 associated with this key in the new communicator; one particular action that a copy call-
 4 back may take is to delete the attribute from the new communicator. Returns in `newcomm`
 5 a new communicator with the same group or groups, any copied cached information, but a
 6 new context (see Section 6.7.1). Please see Section 16.1.7 on page 496 for further discussion
 7 about the C++ bindings for `Dup()` and `Clone()`.

8
 9 *Advice to users.* This operation is used to provide a parallel library call with a dupli-
 10 cate communication space that has the same properties as the original communicator.
 11 This includes any attributes (see below), and topologies (see Chapter 7). This call is
 12 valid even if there are pending point-to-point communications involving the commu-
 13 nicator `comm`. A typical call might involve a `MPI_COMM_DUP` at the beginning of
 14 the parallel call, and an `MPI_COMM_FREE` of that duplicated communicator at the
 15 end of the call. Other models of communicator management are also possible.

16 This call applies to both intra- and inter-communicators. (*End of advice to users.*)

17
 18 *Advice to implementors.* One need not actually copy the group information, but only
 19 add a new reference and increment the reference count. Copy on write can be used
 20 for the cached information. (*End of advice to implementors.*)

21
 22
 23
 24 `MPI_COMM_CREATE(comm, group, newcomm)`

25 IN comm communicator (handle)

26 IN group Group, which is a subset of the group of
 27 comm (handle)

28 OUT newcomm new communicator (handle)

29
 30
 31 `int MPI_Comm_create(MPI_Comm comm, MPI_Group group, MPI_Comm *newcomm)`

32 `MPI_COMM_CREATE(COMM, GROUP, NEWCOMM, IERROR)`

33 `INTEGER COMM, GROUP, NEWCOMM, IERROR`

34
 35 `{MPI::Intercomm MPI::Intercomm::Create(const MPI::Group& group)`
 36 `const(binding deprecated, see Section 15.2) }`

37
 38 `{MPI::Intracomm MPI::Intracomm::Create(const MPI::Group& group)`
 39 `const(binding deprecated, see Section 15.2) }`

40 If `comm` is an intracommunicator, this function returns a new communicator `newcomm` with
 41 communication group defined by the `group` argument. No cached information propagates
 42 from `comm` to `newcomm`. Each process must call with a `group` argument that is a subgroup
 43 of the `group` associated with `comm`; this could be `MPI_GROUP_EMPTY`. The processes may
 44 specify different values for the `group` argument. If a process calls with a non-empty `group`
 45 then all processes in that `group` must call the function with the same `group` as argument,
 46 that is the same processes in the same order. Otherwise the call is erroneous. This implies
 47 that the set of groups specified across the processes must be disjoint. If the calling process
 48 is a member of the group given as `group` argument, then `newcomm` is a communicator with

`group` as its associated group. In the case that a process calls with a `group` to which it does not belong, e.g., `MPI_GROUP_EMPTY`, then `MPI_COMM_NULL` is returned as `newcomm`. The function is collective and must be called by all processes in the group of `comm`.

Rationale. The interface supports the original mechanism from MPI-1.1, which required the same `group` in all processes of `comm`. It was extended in MPI-2.2 to allow the use of disjoint subgroups in order to allow implementations to eliminate unnecessary communication that `MPI_COMM_SPLIT` would incur when the user already knows the membership of the disjoint subgroups. (*End of rationale.*)

Rationale. The requirement that the entire group of `comm` participate in the call stems from the following considerations:

- It allows the implementation to layer `MPI_COMM_CREATE` on top of regular collective communications.
- It provides additional safety, in particular in the case where partially overlapping groups are used to create new communicators.
- It permits implementations sometimes to avoid communication related to context creation.

(*End of rationale.*)

Advice to users. `MPI_COMM_CREATE` provides a means to subset a group of processes for the purpose of separate MIMD computation, with separate communication space. `newcomm`, which emerges from `MPI_COMM_CREATE` can be used in subsequent calls to `MPI_COMM_CREATE` (or other communicator constructors) further to subdivide a computation into parallel sub-computations. A more general service is provided by `MPI_COMM_SPLIT`, below. (*End of advice to users.*)

Advice to implementors. When calling `MPI_COMM_DUP`, all processes call with the same `group` (the `group` associated with the communicator). When calling `MPI_COMM_CREATE`, the processes provide the same `group` or disjoint subgroups. For both calls, it is theoretically possible to agree on a group-wide unique context with no communication. However, local execution of these functions requires use of a larger context name space and reduces error checking. Implementations may strike various compromises between these conflicting goals, such as bulk allocation of multiple contexts in one collective operation.

Important: If new communicators are created without synchronizing the processes involved then the communication system should be able to cope with messages arriving in a context that has not yet been allocated at the receiving process. (*End of advice to implementors.*)

If `comm` is an intercommunicator, then the output communicator is also an intercommunicator where the local group consists only of those processes contained in `group` (see Figure 6.1). The `group` argument should only contain those processes in the local group of the input intercommunicator that are to be a part of `newcomm`. All processes in the same local group of `comm` must specify the same value for `group`, i.e., the same members in the same order. If either `group` does not specify at least one process in the local group of the intercommunicator, or if the calling process is not included in the `group`, `MPI_COMM_NULL` is returned.

Rationale. In the case where either the left or right group is empty, a null communicator is returned instead of an intercommunicator with `MPI_GROUP_EMPTY` because the side with the empty group must return `MPI_COMM_NULL`. (*End of rationale.*)

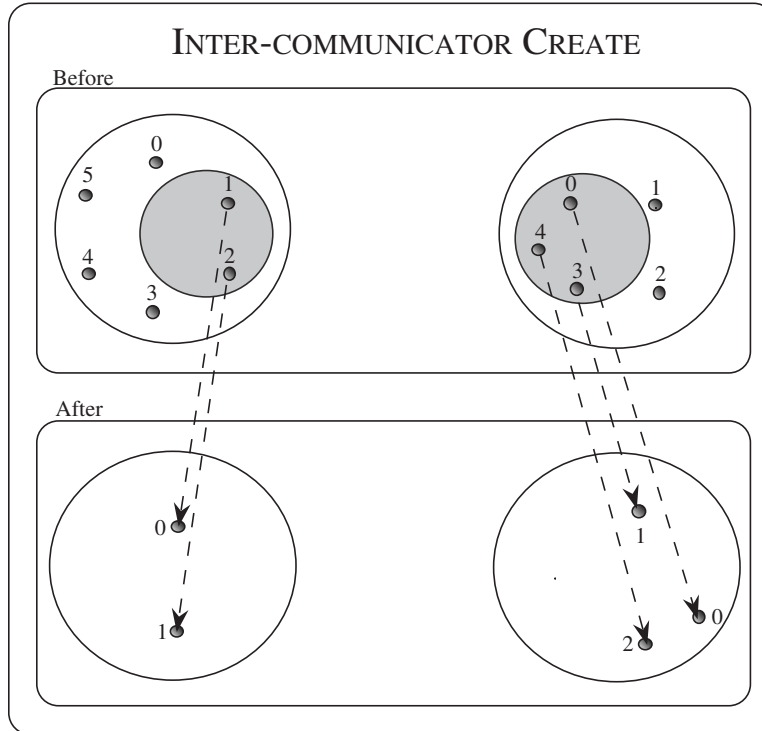


Figure 6.1: Intercommunicator create using `MPI_COMM_CREATE` extended to intercommunicators. The input groups are those in the grey circle.

Example 6.1 The following example illustrates how the first node in the left side of an intercommunicator could be joined with all members on the right side of an intercommunicator to form a new intercommunicator.

```

34     MPI_Comm  inter_comm, new_inter_comm;
35     MPI_Group local_group, group;
36     int      rank = 0; /* rank on left side to include in
37                       new inter-comm */
38
39     /* Construct the original intercommunicator: "inter_comm" */
40     ...
41
42     /* Construct the group of processes to be in new
43        intercommunicator */
44     if (/* I'm on the left side of the intercommunicator */) {
45         MPI_Comm_group ( inter_comm, &local_group );
46         MPI_Group_incl ( local_group, 1, &rank, &group );
47         MPI_Group_free ( &local_group );
48     }

```

```

else
    MPI_Comm_group ( inter_comm, &group );
    MPI_Comm_create ( inter_comm, group, &new_inter_comm );
    MPI_Group_free( &group );
MPI_COMM_SPLIT(comm, color, key, newcomm)
IN      comm          communicator (handle)
IN      color         control of subset assignment (integer)
IN      key           control of rank assignment (integer)
OUT     newcomm       new communicator (handle)
int MPI_Comm_split(MPI_Comm comm, int color, int key, MPI_Comm *newcomm)
MPI_COMM_SPLIT(COMM, COLOR, KEY, NEWCOMM, IERROR)
    INTEGER COMM, COLOR, KEY, NEWCOMM, IERROR
{MPI::Intercomm MPI::Intercomm::Split(int color, int key) const(binding
    deprecated, see Section 15.2) }
{MPI::Intracomm MPI::Intracomm::Split(int color, int key) const(binding
    deprecated, see Section 15.2) }

```

This function partitions the group associated with `comm` into disjoint subgroups, one for each value of `color`. Each subgroup contains all processes of the same color. Within each subgroup, the processes are ranked in the order defined by the value of the argument `key`, with ties broken according to their rank in the old group. A new communicator is created for each subgroup and returned in `newcomm`. A process may supply the color value `MPI_UNDEFINED`, in which case `newcomm` returns `MPI_COMM_NULL`. This is a collective call, but each process is permitted to provide different values for `color` and `key`.

With an intracommunicator `comm`, a call to `MPI_COMM_CREATE(comm, group, newcomm)` is equivalent to a call to `MPI_COMM_SPLIT(comm, color, key, newcomm)`, where processes that are members of their `group` argument provide `color = number of the group` (based on a unique numbering of all disjoint groups) and `key = rank in group`, and all processes that are not members of their `group` argument provide `color = MPI_UNDEFINED`.

The value of `color` must be non-negative.

Advice to users. This is an extremely powerful mechanism for dividing a single communicating group of processes into k subgroups, with k chosen implicitly by the user (by the number of colors asserted over all the processes). Each resulting communicator will be non-overlapping. Such a division could be useful for defining a hierarchy of computations, such as for multigrid, or linear algebra. For intracommunicators, `MPI_COMM_SPLIT` provides similar capability as `MPI_COMM_CREATE` to split a communicating group into disjoint subgroups. `MPI_COMM_SPLIT` is useful when some processes do not have complete information of the other members in their group, but all processes know (the color of) the group to which they belong. In this

1 case, the MPI implementation discovers the other group members via communica-
 2 tion. `MPI_COMM_CREATE` is useful when all processes have complete information
 3 of the members of their group. In this case, MPI can avoid the extra communication
 4 required to discover group membership.

5 Multiple calls to `MPI_COMM_SPLIT` can be used to overcome the requirement that
 6 any call have no overlap of the resulting communicators (each process is of only one
 7 color per call). In this way, multiple overlapping communication structures can be
 8 created. Creative use of the `color` and `key` in such splitting operations is encouraged.

9 Note that, for a fixed color, the keys need not be unique. It is `MPI_COMM_SPLIT`'s
 10 responsibility to sort processes in ascending order according to this key, and to break
 11 ties in a consistent way. If all the keys are specified in the same way, then all the
 12 processes in a given color will have the relative rank order as they did in their parent
 13 group.

14 Essentially, making the key value zero for all processes of a given color means that one
 15 doesn't really care about the rank-order of the processes in the new communicator.
 16 (*End of advice to users.*)

17
 18 *Rationale.* `color` is restricted to be non-negative, so as not to conflict with the value
 19 assigned to `MPI_UNDEFINED`. (*End of rationale.*)

20
 21 The result of `MPI_COMM_SPLIT` on an intercommunicator is that those processes on the
 22 left with the same `color` as those processes on the right combine to create a new intercom-
 23 municator. The `key` argument describes the relative rank of processes on each side of the
 24 intercommunicator (see Figure 6.2). For those colors that are specified only on one side of
 25 the intercommunicator, `MPI_COMM_NULL` is returned. `MPI_COMM_NULL` is also returned
 26 to those processes that specify `MPI_UNDEFINED` as the color.

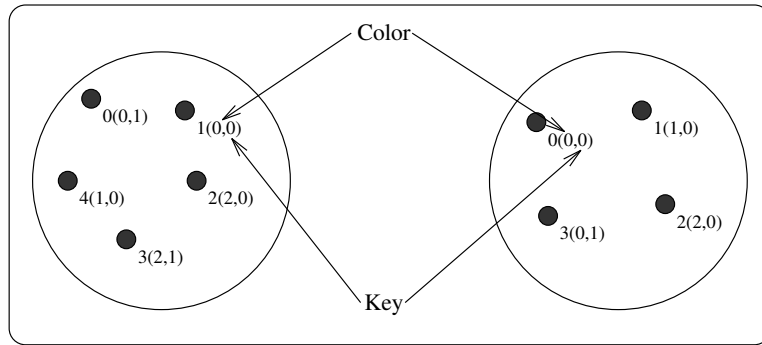
27
 28 *Advice to users.* For intercommunicators, `MPI_COMM_SPLIT` is more general than
 29 `MPI_COMM_CREATE`. A single call to `MPI_COMM_SPLIT` can create a set of disjoint
 30 intercommunicators, while a call to `MPI_COMM_CREATE` creates only one. (*End of*
 31 *advice to users.*)

32
 33 **Example 6.2** (Parallel client-server model). The following client code illustrates how clients
 34 on the left side of an intercommunicator could be assigned to a single server from a pool of
 35 servers on the right side of an intercommunicator.

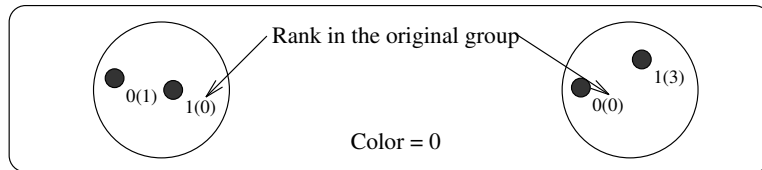
```

36
37     /* Client code */
38     MPI_Comm  multiple_server_comm;
39     MPI_Comm  single_server_comm;
40     int       color, rank, num_servers;
41
42     /* Create intercommunicator with clients and servers:
43        multiple_server_comm */
44     ...
45
46     /* Find out the number of servers available */
47     MPI_Comm_remote_size ( multiple_server_comm, &num_servers );
48

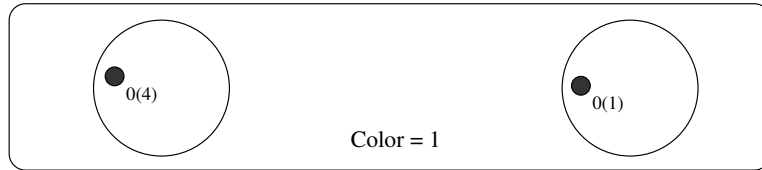
```



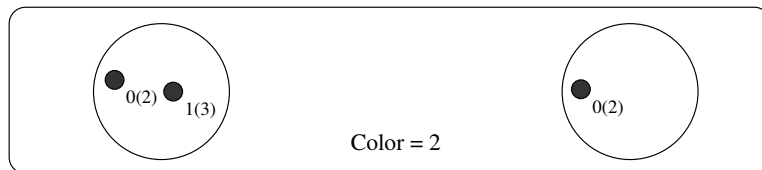
Input Intercommunicator (comm)



Color = 0



Color = 1



Color = 2

Disjoint output communicators (newcomm)
(one per color)

Figure 6.2: Intercommunicator construction achieved by splitting an existing intercommunicator with MPI_COMM_SPLIT extended to intercommunicators.

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      /* Determine my color */
2      MPI_Comm_rank ( multiple_server_comm, &rank );
3      color = rank % num_servers;
4
5      /* Split the intercommunicator */
6      MPI_Comm_split ( multiple_server_comm, color, rank,
7                      &single_server_comm );

```

The following is the corresponding server code:

```

9
10     /* Server code */
11     MPI_Comm  multiple_client_comm;
12     MPI_Comm  single_server_comm;
13     int       rank;
14
15     /* Create intercommunicator with clients and servers:
16        multiple_client_comm */
17     ...
18
19     /* Split the intercommunicator for a single server per group
20        of clients */
21     MPI_Comm_rank ( multiple_client_comm, &rank );
22     MPI_Comm_split ( multiple_client_comm, rank, 0,
23                     &single_server_comm );

```

MPI_COMM_SPLIT_TYPE(comm, type, key, newcomm)

27	IN	comm	communicator (handle)
28	IN	type	subset communicator type (integer)
29	IN	key	control of rank assignment (integer)
30	IN	key	control of rank assignment (integer)
31	OUT	newcomm	new communicator (handle)

32

```

33 int MPI_Comm_split_type(MPI_Comm comm, int type, int key,
34                        MPI_Comm *newcomm)

```

35

```

36 MPI_COMM_SPLIT_TYPE(COMM, TYPE, KEY, NEWCOMM, IERROR)
37     INTEGER COMM, TYPE, KEY, NEWCOMM, IERROR

```

38

This function partitions the group associated with `comm` into disjoint subgroups, based on the type specified by `type`. Each subgroup contains all processes of the same type. Within each subgroup, the processes are ranked in the order defined by the value of the argument `key`, with ties broken according to their rank in the old group. A new communicator is created for each subgroup and returned in `newcomm`. A process may supply the type value `MPI_UNDEFINED`, in which case `newcomm` returns `MPI_COMM_NULL`. This is a collective call, but each process is permitted to provide different values for `type` and `key`.

45

Two types are predefined by MPI:

46

`MPI_COMM_TYPE_SHM` — this type splits the communicator into subcommunicators, each of which can create a shared memory region using `MPI_WIN_ALLOCATE_SHARED`.

48

MPI_COMM_TYPE_PROCESS – this type splits the communicator into subcommunicators, each of which contains all the endpoints in the parent communicator that are created by the same OS process.

6.4.3 Communicator Destructors

`MPI_COMM_FREE(comm)`

INOUT comm communicator to be destroyed (handle)

`int MPI_Comm_free(MPI_Comm *comm)`

`MPI_COMM_FREE(COMM, IERROR)`

INTEGER COMM, IERROR

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

This collective operation marks the communication object for deallocation. The handle is set to `MPI_COMM_NULL`. Any pending operations that use this communicator will complete normally; the object is actually deallocated only if there are no other active references to it. This call applies to intra- and inter-communicators. The delete callback functions for all cached attributes (see Section 6.7) are called in arbitrary order.

Advice to implementors. A reference-count mechanism may be used: the reference count is incremented by each call to `MPI_COMM_DUP`, and decremented by each call to `MPI_COMM_FREE`. The object is ultimately deallocated when the count reaches zero.

Though collective, it is anticipated that this operation will normally be implemented to be local, though a debugging version of an MPI library might choose to synchronize. (*End of advice to implementors.*)

6.5 Motivating Examples

6.5.1 Current Practice #1

Example #1a:

```
int main(int argc, char **argv)
{
    int me, size;
    ...
    MPI_Init ( &argc, &argv );
    MPI_Comm_rank (MPI_COMM_WORLD, &me);
    MPI_Comm_size (MPI_COMM_WORLD, &size);

    (void)printf ("Process %d size %d\n", me, size);
    ...
    MPI_Finalize();
}
```

1 Example #1a is a do-nothing program that initializes itself legally, and refers to the “all”
 2 communicator, and prints a message. It terminates itself legally too. This example does
 3 not imply that MPI supports `printf`-like communication itself.

4 Example #1b (supposing that `size` is even):

```

5        int main(int argc, char **argv)
6        {
7            int me, size;
8            int SOME_TAG = 0;
9            ...
10          MPI_Init(&argc, &argv);
11
12          MPI_Comm_rank(MPI_COMM_WORLD, &me);    /* local */
13          MPI_Comm_size(MPI_COMM_WORLD, &size); /* local */
14
15          if((me % 2) == 0)
16          {
17            /* send unless highest-numbered process */
18            if((me + 1) < size)
19                MPI_Send(..., me + 1, SOME_TAG, MPI_COMM_WORLD);
20          }
21          else
22            MPI_Recv(..., me - 1, SOME_TAG, MPI_COMM_WORLD, &status);
23
24          ...
25          MPI_Finalize();
26        }
27

```

28 Example #1b schematically illustrates message exchanges between “even” and “odd” pro-
 29 cesses in the “all” communicator.

31 6.5.2 Current Practice #2

```

32        int main(int argc, char **argv)
33        {
34            int me, count;
35            void *data;
36            ...
37
38          MPI_Init(&argc, &argv);
39          MPI_Comm_rank(MPI_COMM_WORLD, &me);
40
41          if(me == 0)
42          {
43            /* get input, create buffer ‘‘data’’ */
44            ...
45          }
46
47          MPI_Bcast(data, count, MPI_BYTE, 0, MPI_COMM_WORLD);
48

```

```

...
MPI_Finalize();
}

```

This example illustrates the use of a collective communication.

6.5.3 (Approximate) Current Practice #3

```

int main(int argc, char **argv)
{
    int me, count, count2;
    void *send_buf, *recv_buf, *send_buf2, *recv_buf2;
    MPI_Group MPI_GROUP_WORLD, grprem;
    MPI_Comm commslave;
    static int ranks[] = {0};
    ...
    MPI_Init(&argc, &argv);
    MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);
    MPI_Comm_rank(MPI_COMM_WORLD, &me); /* local */

    MPI_Group_excl(MPI_GROUP_WORLD, 1, ranks, &grprem); /* local */
    MPI_Comm_create(MPI_COMM_WORLD, grprem, &commslave);

    if(me != 0)
    {
        /* compute on slave */
        ...
        MPI_Reduce(send_buf,recv_buff,count, MPI_INT, MPI_SUM, 1, commslave);
        ...
        MPI_Comm_free(&commslave);
    }
    /* zero falls through immediately to this reduce, others do later... */
    MPI_Reduce(send_buf2, recv_buff2, count2,
               MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);

    MPI_Group_free(&MPI_GROUP_WORLD);
    MPI_Group_free(&grprem);
    MPI_Finalize();
}

```

This example illustrates how a group consisting of all but the zeroth process of the “all” group is created, and then how a communicator is formed (`commslave`) for that new group. The new communicator is used in a collective call, and all processes execute a collective call in the `MPI_COMM_WORLD` context. This example illustrates how the two communicators (that inherently possess distinct contexts) protect communication. That is, communication in `MPI_COMM_WORLD` is insulated from communication in `commslave`, and vice versa.

1 In summary, “group safety” is achieved via communicators because distinct contexts
 2 within communicators are enforced to be unique on any process.

3 4 6.5.4 Example #4

5 The following example is meant to illustrate “safety” between point-to-point and collective
 6 communication. MPI guarantees that a single communicator can do safe point-to-point and
 7 collective communication.
 8

```

9  #define TAG_ARBITRARY 12345
10 #define SOME_COUNT      50
11
12 int main(int argc, char **argv)
13 {
14     int me;
15     MPI_Request request[2];
16     MPI_Status status[2];
17     MPI_Group MPI_GROUP_WORLD, subgroup;
18     int ranks[] = {2, 4, 6, 8};
19     MPI_Comm the_comm;
20     ...
21     MPI_Init(&argc, &argv);
22     MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);
23
24     MPI_Group_incl(MPI_GROUP_WORLD, 4, ranks, &subgroup); /* local */
25     MPI_Group_rank(subgroup, &me); /* local */
26
27     MPI_Comm_create(MPI_COMM_WORLD, subgroup, &the_comm);
28
29     if(me != MPI_UNDEFINED)
30     {
31         MPI_Irecv(buff1, count, MPI_DOUBLE, MPI_ANY_SOURCE, TAG_ARBITRARY,
32                 the_comm, request);
33         MPI_Isend(buff2, count, MPI_DOUBLE, (me+1)%4, TAG_ARBITRARY,
34                 the_comm, request+1);
35         for(i = 0; i < SOME_COUNT, i++)
36             MPI_Reduce(..., the_comm);
37         MPI_Waitall(2, request, status);
38
39         MPI_Comm_free(&the_comm);
40     }
41
42     MPI_Group_free(&MPI_GROUP_WORLD);
43     MPI_Group_free(&subgroup);
44     MPI_Finalize();
45 }
46
47
48
```

6.5.5 Library Example #1

The main program:

```

int main(int argc, char **argv)
{
    int done = 0;
    user_lib_t *libh_a, *libh_b;
    void *dataset1, *dataset2;
    ...
    MPI_Init(&argc, &argv);
    ...
    init_user_lib(MPI_COMM_WORLD, &libh_a);
    init_user_lib(MPI_COMM_WORLD, &libh_b);
    ...
    user_start_op(libh_a, dataset1);
    user_start_op(libh_b, dataset2);
    ...
    while(!done)
    {
        /* work */
        ...
        MPI_Reduce(..., MPI_COMM_WORLD);
        ...
        /* see if done */
        ...
    }
    user_end_op(libh_a);
    user_end_op(libh_b);

    ...

    uninit_user_lib(libh_a);
    uninit_user_lib(libh_b);
    MPI_Finalize();
}

```

The user library initialization code:

```

void init_user_lib(MPI_Comm comm, user_lib_t **handle)
{
    user_lib_t *save;

    user_lib_initsave(&save); /* local */
    MPI_Comm_dup(comm, &(save -> comm));

    /* other inits */
    ...

    *handle = save;
}

```

1 User start-up code:

```

2        void user_start_op(user_lib_t *handle, void *data)
3        {
4            MPI_Irecv( ..., handle->comm, &(handle -> irecv_handle) );
5            MPI_Isend( ..., handle->comm, &(handle -> isend_handle) );
6        }
7

```

8 User communication clean-up code:

```

9
10       void user_end_op(user_lib_t *handle)
11       {
12           MPI_Status status;
13           MPI_Wait(handle -> isend_handle, &status);
14           MPI_Wait(handle -> irecv_handle, &status);
15       }
16

```

17 User object clean-up code:

```

18       void uninit_user_lib(user_lib_t *handle)
19       {
20           MPI_Comm_free(&(handle -> comm));
21           free(handle);
22       }
23

```

24 6.5.6 Library Example #2

25 The main program:

```

26       int main(int argc, char **argv)
27       {
28           int ma, mb;
29           MPI_Group MPI_GROUP_WORLD, group_a, group_b;
30           MPI_Comm comm_a, comm_b;
31
32           static int list_a[] = {0, 1};
33           #if defined(EXAMPLE_2B) | defined(EXAMPLE_2C)
34           static int list_b[] = {0, 2 ,3};
35           #else/* EXAMPLE_2A */
36           static int list_b[] = {0, 2};
37           #endif
38           int size_list_a = sizeof(list_a)/sizeof(int);
39           int size_list_b = sizeof(list_b)/sizeof(int);
40
41           ...
42           MPI_Init(&argc, &argv);
43           MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);
44
45           MPI_Group_incl(MPI_GROUP_WORLD, size_list_a, list_a, &group_a);
46           MPI_Group_incl(MPI_GROUP_WORLD, size_list_b, list_b, &group_b);
47
48

```

```

1
2 MPI_Comm_create(MPI_COMM_WORLD, group_a, &comm_a);
3 MPI_Comm_create(MPI_COMM_WORLD, group_b, &comm_b);
4
5 if(comm_a != MPI_COMM_NULL)
6     MPI_Comm_rank(comm_a, &ma);
7 if(comm_b != MPI_COMM_NULL)
8     MPI_Comm_rank(comm_b, &mb);
9
10 if(comm_a != MPI_COMM_NULL)
11     lib_call(comm_a);
12
13 if(comm_b != MPI_COMM_NULL)
14 {
15     lib_call(comm_b);
16     lib_call(comm_b);
17 }
18
19 if(comm_a != MPI_COMM_NULL)
20     MPI_Comm_free(&comm_a);
21 if(comm_b != MPI_COMM_NULL)
22     MPI_Comm_free(&comm_b);
23 MPI_Group_free(&group_a);
24 MPI_Group_free(&group_b);
25 MPI_Group_free(&MPI_GROUP_WORLD);
26 MPI_Finalize();
27 }

```

The library:

```

28
29
30 void lib_call(MPI_Comm comm)
31 {
32     int me, done = 0;
33     MPI_Status status;
34     MPI_Comm_rank(comm, &me);
35     if(me == 0)
36         while(!done)
37         {
38             MPI_Recv(..., MPI_ANY_SOURCE, MPI_ANY_TAG, comm, &status);
39             ...
40         }
41     else
42     {
43         /* work */
44         MPI_Send(..., 0, ARBITRARY_TAG, comm);
45         ....
46     }
47 #ifdef EXAMPLE_2C
48     /* include (resp, exclude) for safety (resp, no safety): */

```

```

1     MPI_Barrier(comm);
2 #endif
3     }

```

The above example is really three examples, depending on whether or not one includes rank 3 in `list_b`, and whether or not a synchronize is included in `lib_call`. This example illustrates that, despite contexts, subsequent calls to `lib_call` with the same context need not be safe from one another (colloquially, “back-masking”). Safety is realized if the `MPI_Barrier` is added. What this demonstrates is that libraries have to be written carefully, even with contexts. When rank 3 is excluded, then the synchronize is not needed to get safety from back masking.

Algorithms like “reduce” and “allreduce” have strong enough source selectivity properties so that they are inherently okay (no backmasking), provided that MPI provides basic guarantees. So are multiple calls to a typical tree-broadcast algorithm with the same root or different roots (see [5]). Here we rely on two guarantees of MPI: pairwise ordering of messages between processes in the same context, and source selectivity — deleting either feature removes the guarantee that backmasking cannot be required.

Algorithms that try to do non-deterministic broadcasts or other calls that include wildcard operations will not generally have the good properties of the deterministic implementations of “reduce,” “allreduce,” and “broadcast.” Such algorithms would have to utilize the monotonically increasing tags (within a communicator scope) to keep things straight.

All of the foregoing is a supposition of “collective calls” implemented with point-to-point operations. MPI implementations may or may not implement collective calls using point-to-point operations. These algorithms are used to illustrate the issues of correctness and safety, independent of how MPI implements its collective calls. See also Section 6.9.

6.6 Inter-Communication

This section introduces the concept of inter-communication and describes the portions of MPI that support it. It describes support for writing programs that contain user-level servers.

All communication described thus far has involved communication between processes that are members of the same group. This type of communication is called “intra-communication” and the communicator used is called an “intra-communicator,” as we have noted earlier in the chapter.

In modular and multi-disciplinary applications, different process groups execute distinct modules and processes within different modules communicate with one another in a pipeline or a more general module graph. In these applications, the most natural way for a process to specify a target process is by the rank of the target process within the target group. In applications that contain internal user-level servers, each server may be a process group that provides services to one or more clients, and each client may be a process group that uses the services of one or more servers. It is again most natural to specify the target process by rank within the target group in these applications. This type of communication is called “inter-communication” and the communicator used is called an “inter-communicator,” as introduced earlier.

An inter-communication is a point-to-point communication between processes in different groups. The group containing a process that initiates an inter-communication operation is called the “local group,” that is, the sender in a send and the receiver in a receive. The

group containing the target process is called the “remote group,” that is, the receiver in a send and the sender in a receive. As in intra-communication, the target process is specified using a (**communicator**, **rank**) pair. Unlike intra-communication, the rank is relative to a second, remote group.

All inter-communicator constructors are blocking and require that the local and remote groups be disjoint.

Advice to users. The groups must be disjoint for several reasons. Primarily, this is the intent of the intercommunicators — to provide a communicator for communication between disjoint groups. This is reflected in the definition of `MPI_INTERCOMM_MERGE`, which allows the user to control the ranking of the processes in the created intracommunicator; this ranking makes little sense if the groups are not disjoint. In addition, the natural extension of collective operations to intercommunicators makes the most sense when the groups are disjoint. (*End of advice to users.*)

Here is a summary of the properties of inter-communication and inter-communicators:

- The syntax of point-to-point and collective communication is the same for both inter- and intra-communication. The same communicator can be used both for send and for receive operations.
- A target process is addressed by its rank in the remote group, both for sends and for receives.
- Communications using an inter-communicator are guaranteed not to conflict with any communications that use a different communicator.
- A communicator will provide either intra- or inter-communication, never both.

The routine `MPI_COMM_TEST_INTER` may be used to determine if a communicator is an inter- or intra-communicator. Inter-communicators can be used as arguments to some of the other communicator access routines. Inter-communicators cannot be used as input to some of the constructor routines for intra-communicators (for instance, `MPI_CART_CREATE`).

Advice to implementors. For the purpose of point-to-point communication, communicators can be represented in each process by a tuple consisting of:

group
send_context
receive_context
source

For inter-communicators, **group** describes the remote group, and **source** is the rank of the process in the local group. For intra-communicators, **group** is the communicator group (remote=local), **source** is the rank of the process in this group, and **send context** and **receive context** are identical. A group can be represented by a rank-to-absolute-address translation table.

The inter-communicator cannot be discussed sensibly without considering processes in both the local and remote groups. Imagine a process **P** in group \mathcal{P} , which has an inter-communicator $\mathbf{C}_{\mathcal{P}}$, and a process **Q** in group \mathcal{Q} , which has an inter-communicator $\mathbf{C}_{\mathcal{Q}}$. Then

- 1 • **C_P.group** describes the group **Q** and **C_Q.group** describes the group **P**.
- 2 • **C_P.send_context** = **C_Q.receive_context** and the context is unique in **Q**;
- 3 **C_P.receive_context** = **C_Q.send_context** and this context is unique in **P**.
- 4 • **C_P.source** is rank of **P** in **P** and **C_Q.source** is rank of **Q** in **Q**.

6 Assume that **P** sends a message to **Q** using the inter-communicator. Then **P** uses
 7 the **group** table to find the absolute address of **Q**; **source** and **send_context** are
 8 appended to the message.

9 Assume that **Q** posts a receive with an explicit source argument using the inter-
 10 communicator. Then **Q** matches **receive_context** to the message context and source
 11 argument to the message source.

12 The same algorithm is appropriate for intra-communicators as well.

13 In order to support inter-communicator accessors and constructors, it is necessary to
 14 supplement this model with additional structures, that store information about the
 15 local communication group, and additional safe contexts. (*End of advice to imple-*
 16 *mentors.*)

19 6.6.1 Inter-communicator Accessors

22 MPI_COMM_TEST_INTER(comm, flag)

24 IN comm communicator (handle)
 25 OUT flag (logical)

27 int MPI_Comm_test_inter(MPI_Comm comm, int *flag)

29 MPI_COMM_TEST_INTER(COMM, FLAG, IERROR)

30 INTEGER COMM, IERROR

31 LOGICAL FLAG

32 {bool MPI::Comm::Is_inter() const(*binding deprecated, see Section 15.2*) }

34 This local routine allows the calling process to determine if a communicator is an inter-
 35 communicator or an intra-communicator. It returns true if it is an inter-communicator,
 36 otherwise false.

37 When an inter-communicator is used as an input argument to the communicator ac-
 38 cessors described above under intra-communication, the following table describes behavior.

MPI_COMM_SIZE	returns the size of the local group.
MPI_COMM_GROUP	returns the local group.
MPI_COMM_RANK	returns the rank in the local group

45 Table 6.1: MPI_COMM_* Function Behavior (in Inter-Communication Mode)

46 Furthermore, the operation MPI_COMM_COMPARE is valid for inter-communicators. Both
 47 communicators must be either intra- or inter-communicators, or else MPI_UNEQUAL results.
 48

Both corresponding local and remote groups must compare correctly to get the results MPI_CONGRUENT and MPI_SIMILAR. In particular, it is possible for MPI_SIMILAR to result because either the local or remote groups were similar but not identical.

The following accessors provide consistent access to the remote group of an inter-communicator:

The following are all local operations.

MPI_COMM_REMOTE_SIZE(comm, size)

IN	comm	inter-communicator (handle)
OUT	size	number of processes in the remote group of comm (integer)

```
int MPI_Comm_remote_size(MPI_Comm comm, int *size)
```

```
MPI_COMM_REMOTE_SIZE(COMM, SIZE, IERROR)
```

```
INTEGER COMM, SIZE, IERROR
```

```
{int MPI::Intercomm::Get_remote_size() const(binding deprecated, see Section 15.2)
    }
```

MPI_COMM_REMOTE_GROUP(comm, group)

IN	comm	inter-communicator (handle)
OUT	group	remote group corresponding to comm (handle)

```
int MPI_Comm_remote_group(MPI_Comm comm, MPI_Group *group)
```

```
MPI_COMM_REMOTE_GROUP(COMM, GROUP, IERROR)
```

```
INTEGER COMM, GROUP, IERROR
```

```
{MPI::Group MPI::Intercomm::Get_remote_group() const(binding deprecated, see Section 15.2) }
```

Rationale. Symmetric access to both the local and remote groups of an inter-communicator is important, so this function, as well as MPI_COMM_REMOTE_SIZE have been provided. (*End of rationale.*)

6.6.2 Inter-communicator Operations

This section introduces four blocking inter-communicator operations.

MPI_INTERCOMM_CREATE is used to bind two intra-communicators into an inter-communicator; the function MPI_INTERCOMM_MERGE creates an intra-communicator by merging the local and remote groups of an inter-communicator. The functions MPI_COMM_DUP and MPI_COMM_FREE, introduced previously, duplicate and free an inter-communicator, respectively.

Overlap of local and remote groups that are bound into an inter-communicator is prohibited. If there is overlap, then the program is erroneous and is likely to deadlock. (If

1 a process is multithreaded, and MPI calls block only a thread, rather than a process, then
 2 “dual membership” can be supported. It is then the user’s responsibility to make sure that
 3 calls on behalf of the two “roles” of a process are executed by two independent threads.)

4 The function `MPI_INTERCOMM_CREATE` can be used to create an inter-communicator
 5 from two existing intra-communicators, in the following situation: At least one selected
 6 member from each group (the “group leader”) has the ability to communicate with the
 7 selected member from the other group; that is, a “peer” communicator exists to which both
 8 leaders belong, and each leader knows the rank of the other leader in this peer communicator.
 9 Furthermore, members of each group know the rank of their leader.

10 Construction of an inter-communicator from two intra-communicators requires separate
 11 collective operations in the local group and in the remote group, as well as a point-to-point
 12 communication between a process in the local group and a process in the remote group.

13 In standard MPI implementations (with static process allocation at initialization), the
 14 `MPI_COMM_WORLD` communicator (or preferably a dedicated duplicate thereof) can be this
 15 peer communicator. For applications that have used `spawn` or `join`, it may be necessary to
 16 first create an intracommunicator to be used as peer.

17 The application topology functions described in Chapter 7 do not apply to inter-
 18 communicators. Users that require this capability should utilize
 19 `MPI_INTERCOMM_MERGE` to build an intra-communicator, then apply the graph or carte-
 20 sian topology capabilities to that intra-communicator, creating an appropriate topology-
 21 oriented intra-communicator. Alternatively, it may be reasonable to devise one’s own ap-
 22 plication topology mechanisms for this case, without loss of generality.

23
 24
 25 `MPI_INTERCOMM_CREATE(local_comm, local_leader, peer_comm, remote_leader, tag,`
 26 `newintercomm)`

27	IN	<code>local_comm</code>	local intra-communicator (handle)
28	IN	<code>local_leader</code>	rank of local group leader in <code>local_comm</code> (integer)
29	IN	<code>peer_comm</code>	“peer” communicator; significant only at the
30			<code>local_leader</code> (handle)
31			
32	IN	<code>remote_leader</code>	rank of remote group leader in <code>peer_comm</code> ; significant
33			only at the <code>local_leader</code> (integer)
34	IN	<code>tag</code>	“safe” tag (integer)
35	OUT	<code>newintercomm</code>	new inter-communicator (handle)
36			

37
 38 `int MPI_Intercomm_create(MPI_Comm local_comm, int local_leader,`
 39 `MPI_Comm peer_comm, int remote_leader, int tag,`
 40 `MPI_Comm *newintercomm)`

41 `MPI_INTERCOMM_CREATE(LOCAL_COMM, LOCAL_LEADER, PEER_COMM, REMOTE_LEADER,`
 42 `TAG, NEWINTERCOMM, IERROR)`

43 `INTEGER LOCAL_COMM, LOCAL_LEADER, PEER_COMM, REMOTE_LEADER, TAG,`
 44 `NEWINTERCOMM, IERROR`

45
 46 `{MPI::Intercomm MPI::Intracomm::Create_intercomm(int local_leader, const`
 47 `MPI::Comm& peer_comm, int remote_leader, int tag) const(binding`
 48 `deprecated, see Section 15.2) }`

This call creates an inter-communicator. It is collective over the union of the local and remote groups. Processes should provide identical `local_comm` and `local_leader` arguments within each group. Wildcards are not permitted for `remote_leader`, `local_leader`, and `tag`.

This call uses point-to-point communication with communicator `peer_comm`, and with tag `tag` between the leaders. Thus, care must be taken that there be no pending communication on `peer_comm` that could interfere with this communication.

Advice to users. We recommend using a dedicated peer communicator, such as a duplicate of `MPI_COMM_WORLD`, to avoid trouble with peer communicators. (*End of advice to users.*)

`MPI_INTERCOMM_MERGE(intercomm, high, newintracomm)`

IN	<code>intercomm</code>	Inter-Communicator (handle)
IN	<code>high</code>	(logical)
OUT	<code>newintracomm</code>	new intra-communicator (handle)

```
int MPI_Intercomm_merge(MPI_Comm intercomm, int high,
                        MPI_Comm *newintracomm)
```

```
MPI_INTERCOMM_MERGE(INTERCOMM, HIGH, INTRACOMM, IERROR)
INTEGER INTERCOMM, INTRACOMM, IERROR
LOGICAL HIGH
```

```
{MPI::Intracomm MPI::Intercomm::Merge(bool high) const(binding deprecated, see
Section 15.2) }
```

This function creates an intra-communicator from the union of the two groups that are associated with `intercomm`. All processes should provide the same `high` value within each of the two groups. If processes in one group provided the value `high = false` and processes in the other group provided the value `high = true` then the union orders the “low” group before the “high” group. If all processes provided the same `high` argument then the order of the union is arbitrary. This call is blocking and collective within the union of the two groups.

The error handler on the new intercommunicator in each process is inherited from the communicator that contributes the local group. Note that this can result in different processes in the same communicator having different error handlers.

Advice to implementors. The implementation of `MPI_INTERCOMM_MERGE`, `MPI_COMM_FREE` and `MPI_COMM_DUP` are similar to the implementation of `MPI_INTERCOMM_CREATE`, except that contexts private to the input inter-communicator are used for communication between group leaders rather than contexts inside a bridge communicator. (*End of advice to implementors.*)

6.6.3 Inter-Communication Examples

Example 1: Three-Group “Pipeline”

Groups 0 and 1 communicate. Groups 1 and 2 communicate. Therefore, group 0 requires one inter-communicator, group 1 requires two inter-communicators, and group 2 requires 1

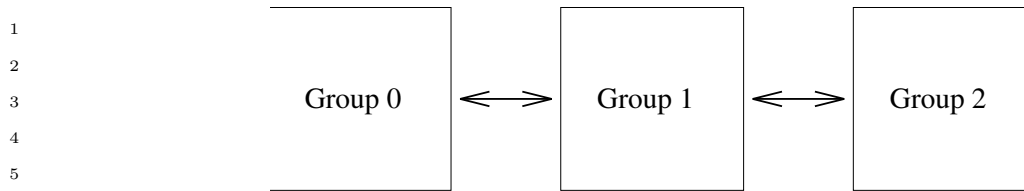


Figure 6.3: Three-group pipeline[ticket0.][.]

inter-communicator.

```

11 int main(int argc, char **argv)
12 {
13     MPI_Comm myComm; /* intra-communicator of local sub-group */
14     MPI_Comm myFirstComm; /* inter-communicator */
15     MPI_Comm mySecondComm; /* second inter-communicator (group 1 only) */
16     int membershipKey;
17     int rank;
18
19     MPI_Init(&argc, &argv);
20     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
21
22     /* User code must generate membershipKey in the range [0, 1, 2] */
23     membershipKey = rank % 3;
24
25     /* Build intra-communicator for local sub-group */
26     MPI_Comm_split(MPI_COMM_WORLD, membershipKey, rank, &myComm);
27
28     /* Build inter-communicators. Tags are hard-coded. */
29     if (membershipKey == 0)
30     {
31         /* Group 0 communicates with group 1. */
32         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,
33                               1, &myFirstComm);
34     }
35     else if (membershipKey == 1)
36     {
37         /* Group 1 communicates with groups 0 and 2. */
38         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 0,
39                               1, &myFirstComm);
40         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 2,
41                               12, &mySecondComm);
42     }
43     else if (membershipKey == 2)
44     {
45         /* Group 2 communicates with group 1. */
46         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,
47                               12, &myFirstComm);
48     }
49
50     /* Do work ... */

```

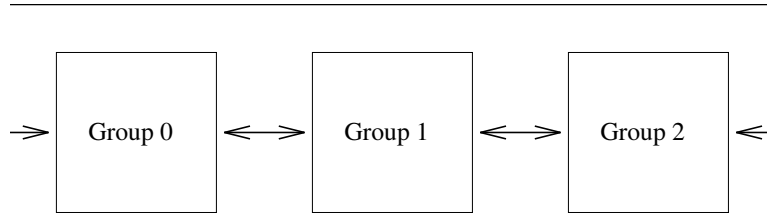


Figure 6.4: Three-group ring[ticket0.][.]

```

switch(membershipKey) /* free communicators appropriately */
{
case 1:
    MPI_Comm_free(&mySecondComm);
case 0:
case 2:
    MPI_Comm_free(&myFirstComm);
    break;
}

MPI_Finalize();
}

```

Example 2: Three-Group “Ring”

Groups 0 and 1 communicate. Groups 1 and 2 communicate. Groups 0 and 2 communicate. Therefore, each requires two inter-communicators.

```

int main(int argc, char **argv)
{
    MPI_Comm    myComm;        /* intra-communicator of local sub-group */
    MPI_Comm    myFirstComm; /* inter-communicators */
    MPI_Comm    mySecondComm;
    MPI_Status  status;
    int  membershipKey;
    int  rank;

    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    ...

    /* User code must generate membershipKey in the range [0, 1, 2] */
    membershipKey = rank % 3;

    /* Build intra-communicator for local sub-group */
    MPI_Comm_split(MPI_COMM_WORLD, membershipKey, rank, &myComm);

    /* Build inter-communicators. Tags are hard-coded. */
    if (membershipKey == 0)

```

```

1      {          /* Group 0 communicates with groups 1 and 2. */
2          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,
3                               1, &myFirstComm);
4          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 2,
5                               2, &mySecondComm);
6      }
7      else if (membershipKey == 1)
8      {          /* Group 1 communicates with groups 0 and 2. */
9          MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 0,
10                              1, &myFirstComm);
11         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 2,
12                              12, &mySecondComm);
13     }
14     else if (membershipKey == 2)
15     {          /* Group 2 communicates with groups 0 and 1. */
16         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 0,
17                              2, &myFirstComm);
18         MPI_Intercomm_create( myComm, 0, MPI_COMM_WORLD, 1,
19                              12, &mySecondComm);
20     }
21
22     /* Do some work ... */
23
24     /* Then free communicators before terminating... */
25     MPI_Comm_free(&myFirstComm);
26     MPI_Comm_free(&mySecondComm);
27     MPI_Comm_free(&myComm);
28     MPI_Finalize();
29 }

```

6.7 Caching

MPI provides a “caching” facility that allows an application to attach arbitrary pieces of information, called **attributes**, to three kinds of MPI objects, communicators, windows and datatypes. More precisely, the caching facility allows a portable library to do the following:

- pass information between calls by associating it with an MPI intra- or inter-communicator, window or datatype,
- quickly retrieve that information, and
- be guaranteed that out-of-date information is never retrieved, even if the object is freed and its handle subsequently reused by MPI.

The caching capabilities, in some form, are required by built-in MPI routines such as collective communication and application topology. Defining an interface to these capabilities as part of the MPI standard is valuable because it permits routines like collective communication and application topologies to be implemented as portable code, and also because it makes MPI more extensible by allowing user-written routines to use standard MPI calling sequences.

Advice to users. The communicator `MPI_COMM_SELF` is a suitable choice for posting process-local attributes, via this attributing-caching mechanism. (*End of advice to users.*)

Rationale. In one extreme one can allow caching on all opaque handles. The other extreme is to only allow it on communicators. Caching has a cost associated with it and should only be allowed when it is clearly needed and the increased cost is modest. This is the reason that windows and datatypes were added but not other handles. (*End of rationale.*)

One difficulty is the potential for size differences between Fortran integers and C pointers. To overcome this problem with attribute caching on communicators, functions are also given for this case. The functions to cache on datatypes and windows also address this issue. For a general discussion of the address size problem, see Section 16.3.6.

Advice to implementors. High-quality implementations should raise an error when a keyval that was created by a call to `MPI_XXX_CREATE_KEYVAL` is used with an object of the wrong type with a call to `MPI_YYY_GET_ATTR`, `MPI_YYY_SET_ATTR`, `MPI_YYY_DELETE_ATTR`, or `MPI_YYY_FREE_KEYVAL`. To do so, it is necessary to maintain, with each keyval, information on the type of the associated user function. (*End of advice to implementors.*)

6.7.1 Functionality

Attributes can be attached to communicators, windows, and datatypes. Attributes are local to the process and specific to the communicator to which they are attached. Attributes are not propagated by MPI from one communicator to another except when the communicator is duplicated using `MPI_COMM_DUP` (and even then the application must give specific permission through callback functions for the attribute to be copied).

Advice to users. Attributes in C are of type `void *`. Typically, such an attribute will be a pointer to a structure that contains further information, or a handle to an MPI object. In Fortran, attributes are of type `INTEGER`. Such attribute can be a handle to an MPI object, or just an integer-valued attribute. (*End of advice to users.*)

Advice to implementors. Attributes are scalar values, equal in size to, or larger than a C-language pointer. Attributes can always hold an MPI handle. (*End of advice to implementors.*)

The caching interface defined here requires that attributes be stored by MPI opaquely within a communicator, window, and datatype. Accessor functions include the following:

- 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
```

```
{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.

The Fortran callback functions are:

```
SUBROUTINE COMM_COPY_ATTR_FN(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_FN(COMM, COMM_KEYVAL, ATTRIBUTE_VAL,
    EXTRA_STATE, IERROR)
    INTEGER COMM, 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.

1 *Advice to users.* Even though both formal arguments `attribute_val_in` and
2 `attribute_val_out` are of type `void *`, their usage differs. The C copy function is passed
3 by MPI in `attribute_val_in` the *value* of the attribute, and in `attribute_val_out` the
4 *address* of the attribute, so as to allow the function to return the (new) attribute
5 value. The use of type `void *` for both is to avoid messy type casts.

6 A valid copy function is one that completely duplicates the information by making
7 a full duplicate copy of the data structures implied by an attribute; another might
8 just make another reference to that data structure, while using a reference-count
9 mechanism. Other types of attributes might not copy at all (they might be specific
10 to `oldcomm` only). (*End of advice to users.*)

11
12 *Advice to implementors.* A C interface should be assumed for copy and delete
13 functions associated with key values created in C; a Fortran calling interface should
14 be assumed for key values created in Fortran. (*End of advice to implementors.*)

15
16 Analogous to `comm_copy_attr_fn` is a callback deletion function, defined as follows.
17 The `comm_delete_attr_fn` function is invoked when a communicator is deleted by
18 `MPI_COMM_FREE` or when a call is made explicitly to `MPI_COMM_DELETE_ATTR`.
19 `comm_delete_attr_fn` should be of type `MPI_Comm_delete_attr_function`.

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

24 The argument `comm_delete_attr_fn` may be specified as `MPI_COMM_NULL_DELETE_FN`
25 from either C, C++, or Fortran. `MPI_COMM_NULL_DELETE_FN` is a function that
26 does nothing, other than returning `MPI_SUCCESS`. `MPI_COMM_NULL_DELETE_FN` re-
27 places `MPI_NULL_DELETE_FN`, whose use is deprecated.

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

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

33
34 *Advice to implementors.* To be able to use the predefined C functions
35 `MPI_COMM_NULL_COPY_FN` or `MPI_COMM_DUP_FN` as `comm_copy_attr_fn` argu-
36 ment and/or `MPI_COMM_NULL_DELETE_FN` as the `comm_delete_attr_fn` argument
37 in a call to the C++ routine `MPI::Comm::Create_keyval`, this routine may be over-
38 loaded with 3 additional routines that accept the C functions as the first, the second,
39 or both input arguments (instead of an argument that matches the C++ prototype).
40 (*End of advice to implementors.*)

41
42 *Advice to users.* If a user wants to write a “wrapper” routine that internally calls
43 `MPI::Comm::Create_keyval` and `comm_copy_attr_fn` and/or `comm_delete_attr_fn` are
44 arguments of this wrapper routine, and if this wrapper routine should be callable with
45 both user-defined C++ copy and delete functions and with the predefined C functions,
46 then the same overloading as described above in the advice to implementors may be
47 necessary. (*End of advice to users.*)

48

```

MPI_COMM_FREE_KEYVAL(comm_keyval) 1
    INOUT    comm_keyval            key value (integer) 2
                                                3
int MPI_Comm_free_keyval(int *comm_keyval) 4
                                                5
MPI_COMM_FREE_KEYVAL(COMM_KEYVAL, IERROR) 6
    INTEGER COMM_KEYVAL, IERROR 7
{static void MPI::Comm::Free_keyval(int& comm_keyval) (binding deprecated, see 8
    Section 15.2) } 9
                                                10

```

Frees an extant attribute key. This function sets the value of `keyval` to `MPI_KEYVAL_INVALID`. Note that it is not erroneous to free an attribute key that is in use, because the actual free does not transpire until after all references (in other communicators on the process) to the key have been freed. These references need to be explicitly freed by the program, either via calls to `MPI_COMM_DELETE_ATTR` that free one attribute instance, or by calls to `MPI_COMM_FREE` that free all attribute instances associated with the freed communicator.

This call is identical to the MPI-1 call `MPI_KEYVAL_FREE` but is needed to match the new communicator-specific creation function. The use of `MPI_KEYVAL_FREE` is deprecated.

```

MPI_COMM_SET_ATTR(comm, comm_keyval, attribute_val) 21
    INOUT    comm                    communicator from which attribute will be attached 23
                                                (handle) 24
    IN      comm_keyval              key value (integer) 25
    IN      attribute_val            attribute value 27
                                                28
int MPI_Comm_set_attr(MPI_Comm comm, int comm_keyval, void *attribute_val) 29
MPI_COMM_SET_ATTR(COMM, COMM_KEYVAL, ATTRIBUTE_VAL, IERROR) 30
    INTEGER COMM, COMM_KEYVAL, IERROR 32
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 33
{void MPI::Comm::Set_attr(int comm_keyval, const void* attribute_val) 34
    const (binding deprecated, see Section 15.2) } 35
                                                36

```

This function stores the stipulated attribute value `attribute_val` for subsequent retrieval by `MPI_COMM_GET_ATTR`. If the value is already present, then the outcome is as if `MPI_COMM_DELETE_ATTR` was first called to delete the previous value (and the callback function `comm_delete_attr_fn` was executed), and a new value was next stored. The call is erroneous if there is no key with value `keyval`; in particular `MPI_KEYVAL_INVALID` is an erroneous key value. The call will fail if the `comm_delete_attr_fn` function returned an error code other than `MPI_SUCCESS`.

This function replaces `MPI_ATTR_PUT`, whose use is deprecated. The C binding is identical. The Fortran binding differs in that `attribute_val` is an address-sized integer.

```

1 MPI_COMM_GET_ATTR(comm, comm_keyval, attribute_val, flag)
2     IN      comm      communicator to which the attribute is attached (han-
3                dle)
4     IN      comm_keyval  key value (integer)
5     OUT     attribute_val  attribute value, unless flag = false
6     OUT     flag         false if no attribute is associated with the key (logical)
7
8

```

```

9
10 int MPI_Comm_get_attr(MPI_Comm comm, int comm_keyval, void *attribute_val,
11                      int *flag)
12

```

```

13 MPI_COMM_GET_ATTR(COMM, COMM_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)
14     INTEGER COMM, COMM_KEYVAL, IERROR
15     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
16     LOGICAL FLAG
17

```

```

18 {bool MPI::Comm::Get_attr(int comm_keyval, void* attribute_val)
19     const(binding deprecated, see Section 15.2) }
20

```

Retrieves attribute value by key. The call is erroneous if there is no key with value `keyval`. On the other hand, the call is correct if the key value exists, but no attribute is attached on `comm` for that key; in such case, the call returns `flag = false`. In particular `MPI_KEYVAL_INVALID` is an erroneous key value.

Advice to users. The call to `MPI_Comm_set_attr` passes in `attribute_val` the *value* of the attribute; the call to `MPI_Comm_get_attr` passes in `attribute_val` the *address* of the location where the attribute value is to be returned. Thus, if the attribute value itself is a pointer of type `void*`, then the actual `attribute_val` parameter to `MPI_Comm_set_attr` will be of type `void*` and the actual `attribute_val` parameter to `MPI_Comm_get_attr` will be of type `void**`. (*End of advice to users.*)

Rationale. The use of a formal parameter `attribute_val` or type `void*` (rather than `void**`) avoids the messy type casting that would be needed if the attribute value is declared with a type other than `void*`. (*End of rationale.*)

This function replaces `MPI_ATTR_GET`, whose use is deprecated. The C binding is identical. The Fortran binding differs in that `attribute_val` is an address-sized integer.

```

37
38 MPI_COMM_DELETE_ATTR(comm, comm_keyval)
39
40     INOUT   comm      communicator from which the attribute is deleted (han-
41                dle)
42     IN      comm_keyval  key value (integer)
43

```

```

44 int MPI_Comm_delete_attr(MPI_Comm comm, int comm_keyval)
45

```

```

46 MPI_COMM_DELETE_ATTR(COMM, COMM_KEYVAL, IERROR)
47     INTEGER COMM, COMM_KEYVAL, IERROR
48

```

```
{void MPI::Comm::Delete_attr(int comm_keyval) (binding deprecated, see
    Section 15.2) }
```

Delete attribute from cache by key. This function invokes the attribute delete function `comm_delete_attr_fn` specified when the keyval was created. The call will fail if the `comm_delete_attr_fn` function returns an error code other than `MPI_SUCCESS`.

Whenever a communicator is replicated using the function `MPI_COMM_DUP`, all callback copy functions for attributes that are currently set are invoked (in arbitrary order). Whenever a communicator is deleted using the function `MPI_COMM_FREE` all callback delete functions for attributes that are currently set are invoked.

This function is the same as `MPI_ATTR_DELETE` but is needed to match the new communicator specific functions. The use of `MPI_ATTR_DELETE` is deprecated.

6.7.3 Windows

The new functions for caching on windows are:

```
MPI_WIN_CREATE_KEYVAL(win_copy_attr_fn, win_delete_attr_fn, win_keyval, extra_state)
```

IN	<code>win_copy_attr_fn</code>	copy callback function for <code>win_keyval</code> (function)
IN	<code>win_delete_attr_fn</code>	delete callback function for <code>win_keyval</code> (function)
OUT	<code>win_keyval</code>	key value for future access (integer)
IN	<code>extra_state</code>	extra state for callback functions

```
int MPI_Win_create_keyval(MPI_Win_copy_attr_function *win_copy_attr_fn,
    MPI_Win_delete_attr_function *win_delete_attr_fn,
    int *win_keyval, void *extra_state)
```

```
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
```

```
{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:

```

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

```

```

4      and

```

```

5  typedef int MPI_Win_delete_attr_function(MPI_Win win, int win_keyval,
6      void *attribute_val, void *extra_state);
7

```

8 The Fortran callback functions are:

```

9  SUBROUTINE WIN_COPY_ATTR_FN(OLDWIN, WIN_KEYVAL, EXTRA_STATE,
10     ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
11     INTEGER OLDWIN, WIN_KEYVAL, IERROR
12     INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE, ATTRIBUTE_VAL_IN,
13     ATTRIBUTE_VAL_OUT
14     LOGICAL FLAG

```

```

15     and

```

```

16  SUBROUTINE WIN_DELETE_ATTR_FN(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, EXTRA_STATE,
17     IERROR)
18     INTEGER WIN, WIN_KEYVAL, IERROR
19     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE
20

```

21 The C++ callbacks are:

```

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

```

```

26     and

```

```

27  {typedef int MPI::Win::Delete_attr_function(MPI::Win& win, int win_keyval,
28     void* attribute_val, void* extra_state); (binding deprecated, see
29     Section 15.2)}
30

```

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

```

34
35
36 MPI_WIN_FREE_KEYVAL(win_keyval)

```

```

37     INOUT    win_keyval                key value (integer)
38

```

```

39 int MPI_Win_free_keyval(int *win_keyval)

```

```

40
41 MPI_WIN_FREE_KEYVAL(WIN_KEYVAL, IERROR)
42     INTEGER WIN_KEYVAL, IERROR

```

```

43 {static void MPI::Win::Free_keyval(int& win_keyval) (binding deprecated, see
44     Section 15.2 ) }
45

```

```

46
47
48

```



```

MPI_WIN_SET_ATTR(win, win_keyval, attribute_val) 1
    INOUT win window to which attribute will be attached (handle) 2
    IN win_keyval key value (integer) 3
    IN attribute_val attribute value 4
    5
    6
int MPI_Win_set_attr(MPI_Win win, int win_keyval, void *attribute_val) 7
MPI_WIN_SET_ATTR(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, IERROR) 8
    INTEGER WIN, WIN_KEYVAL, IERROR 9
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 10
{void MPI::Win::Set_attr(int win_keyval, const void* attribute_val) (binding 12
    deprecated, see Section 15.2) } 13
    14
    15
MPI_WIN_GET_ATTR(win, win_keyval, attribute_val, flag) 16
    IN win window to which the attribute is attached (handle) 17
    IN win_keyval key value (integer) 18
    OUT attribute_val attribute value, unless flag = false 19
    OUT flag false if no attribute is associated with the key (logical) 20
    21
    22
int MPI_Win_get_attr(MPI_Win win, int win_keyval, void *attribute_val, 23
    int *flag) 24
MPI_WIN_GET_ATTR(WIN, WIN_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR) 25
    INTEGER WIN, WIN_KEYVAL, IERROR 26
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 27
    LOGICAL FLAG 28
    29
    30
{bool MPI::Win::Get_attr(int win_keyval, void* attribute_val) const (binding 31
    deprecated, see Section 15.2) } 32
    33
    34
MPI_WIN_DELETE_ATTR(win, win_keyval) 35
    INOUT win window from which the attribute is deleted (handle) 36
    IN win_keyval key value (integer) 37
    38
    39
int MPI_Win_delete_attr(MPI_Win win, int win_keyval) 40
MPI_WIN_DELETE_ATTR(WIN, WIN_KEYVAL, IERROR) 41
    INTEGER WIN, WIN_KEYVAL, IERROR 42
    43
{void MPI::Win::Delete_attr(int win_keyval) (binding deprecated, see Section 15.2) 44
    } 45
    46
    47
    48

```

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
```

```
{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 type,
    int type_keyval, void *attribute_val, void *extra_state);
```

The Fortran callback functions are:

```
SUBROUTINE TYPE_COPY_ATTR_FN(OLDTYPE, TYPE_KEYVAL, EXTRA_STATE,
    ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT, FLAG, IERROR)
```

```
INTEGER OLDTYPE, TYPE_KEYVAL, IERROR
```

```
INTEGER(KIND=MPI_ADDRESS_KIND) EXTRA_STATE,
```

```

    ATTRIBUTE_VAL_IN, ATTRIBUTE_VAL_OUT 1
    LOGICAL FLAG 2
    and 3
    SUBROUTINE TYPE_DELETE_ATTR_FN(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, 4
        EXTRA_STATE, IERROR) 5
    INTEGER TYPE, TYPE_KEYVAL, IERROR 6
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL, EXTRA_STATE 7
    8
    The C++ callbacks are: 9
    {typedef int 10
        MPI::Datatype::Copy_attr_function(const MPI::Datatype& oldtype, 11
            int type_keyval, void* extra_state, 12
            const void* attribute_val_in, void* attribute_val_out, 13
            bool& flag); (binding deprecated, see Section 15.2)} 14
    and 15
    {typedef int MPI::Datatype::Delete_attr_function(MPI::Datatype& type, 16
        int type_keyval, void* attribute_val, void* extra_state); 17
        (binding deprecated, see Section 15.2)} 18
    19
    If an attribute copy function or attribute delete function returns other than 20
    MPI_SUCCESS, then the call that caused it to be invoked (for example, MPI_TYPE_FREE), 21
    is erroneous. 22
    23
    MPI_TYPE_FREE_KEYVAL(type_keyval) 24
    INOUT type_keyval key value (integer) 25
    26
    int MPI_Type_free_keyval(int *type_keyval) 27
    28
    MPI_TYPE_FREE_KEYVAL(TYPE_KEYVAL, IERROR) 29
    INTEGER TYPE_KEYVAL, IERROR 30
    31
    {static void MPI::Datatype::Free_keyval(int& type_keyval) (binding deprecated, 32
        see Section 15.2) } 33
    34
    35
    MPI_TYPE_SET_ATTR(type, type_keyval, attribute_val) 36
    INOUT type datatype to which attribute will be attached (handle) 37
    IN type_keyval key value (integer) 38
    IN attribute_val attribute value 39
    40
    int MPI_Type_set_attr(MPI_Datatype type, int type_keyval, 41
        void *attribute_val) 42
    43
    MPI_TYPE_SET_ATTR(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, IERROR) 44
    INTEGER TYPE, TYPE_KEYVAL, IERROR 45
    INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL 46
    47
    48

```

```

1  {void MPI::Datatype::Set_attr(int type_keyval, const void*
2      attribute_val) (binding deprecated, see Section 15.2) }
3
4
5  MPI_TYPE_GET_ATTR(type, type_keyval, attribute_val, flag)
6
7      IN      type                datatype to which the attribute is attached (handle)
8
9      IN      type_keyval         key value (integer)
10
11     OUT     attribute_val       attribute value, unless flag = false
12
13     OUT     flag                false if no attribute is associated with the key (logical)
14
15 int MPI_Type_get_attr(MPI_Datatype type, int type_keyval, void
16     *attribute_val, int *flag)
17
18 MPI_TYPE_GET_ATTR(TYPE, TYPE_KEYVAL, ATTRIBUTE_VAL, FLAG, IERROR)
19     INTEGER TYPE, TYPE_KEYVAL, IERROR
20     INTEGER(KIND=MPI_ADDRESS_KIND) ATTRIBUTE_VAL
21     LOGICAL FLAG
22
23 {bool MPI::Datatype::Get_attr(int type_keyval, void* attribute_val)
24     const(binding deprecated, see Section 15.2) }
25
26
27 MPI_TYPE_DELETE_ATTR(type, type_keyval)
28
29     INOUT   type                datatype from which the attribute is deleted (handle)
30
31     IN      type_keyval         key value (integer)
32
33 int MPI_Type_delete_attr(MPI_Datatype type, int type_keyval)
34
35 MPI_TYPE_DELETE_ATTR(TYPE, TYPE_KEYVAL, IERROR)
36     INTEGER TYPE, TYPE_KEYVAL, IERROR
37
38 {void MPI::Datatype::Delete_attr(int type_keyval) (binding deprecated, see
39     Section 15.2) }
40
41
42
43
44
45
46
47
48

```

6.7.5 Error Class for Invalid Keyval

Key values for attributes are system-allocated, by `MPI_{TYPE,COMM,WIN}_CREATE_KEYVAL`. Only such values can be passed to the functions that use key values as input arguments. In order to signal that an erroneous key value has been passed to one of these functions, there is a new MPI error class: `MPI_ERR_KEYVAL`. It can be returned by `MPI_ATTR_PUT`, `MPI_ATTR_GET`, `MPI_ATTR_DELETE`, `MPI_KEYVAL_FREE`, `MPI_{TYPE,COMM,WIN}_DELETE_ATTR`, `MPI_{TYPE,COMM,WIN}_SET_ATTR`, `MPI_{TYPE,COMM,WIN}_GET_ATTR`, `MPI_{TYPE,COMM,WIN}_FREE_KEYVAL`, `MPI_COMM_DUP`, `MPI_COMM_DISCONNECT`, and `MPI_COMM_FREE`. The last three are included because `keyval` is an argument to the copy and delete functions for attributes.

6.7.6 Attributes Example

Advice to users. This example shows how to write a collective communication operation that uses caching to be more efficient after the first call. The coding style assumes that MPI function results return only error statuses. (*End of advice to users.*)

```

/* key for this module's stuff: */
static int gop_key = MPI_KEYVAL_INVALID;

typedef struct
{
    int ref_count;          /* reference count */
    /* other stuff, whatever else we want */
} gop_stuff_type;

Efficient_Collective_Op (comm, ...)
MPI_Comm comm;
{
    gop_stuff_type *gop_stuff;
    MPI_Group      group;
    int            foundflag;

    MPI_Comm_group(comm, &group);

    if (gop_key == MPI_KEYVAL_INVALID) /* get a key on first call ever */
    {
        if ( ! MPI_Comm_create_keyval( gop_stuff_copier,
                                       gop_stuff_destructor,
                                       &gop_key, (void *)0));
        /* get the key while assigning its copy and delete callback
           behavior. */

        MPI_Abort (comm, 99);
    }

    MPI_Comm_get_attr (comm, gop_key, &gop_stuff, &foundflag);
    if (foundflag)
    { /* This module has executed in this group before.
       We will use the cached information */
    }
    else
    { /* This is a group that we have not yet cached anything in.
       We will now do so.
       */

        /* First, allocate storage for the stuff we want,
           and initialize the reference count */

        gop_stuff = (gop_stuff_type *) malloc (sizeof(gop_stuff_type));

```

```

1      if (gop_stuff == NULL) { /* abort on out-of-memory error */ }
2
3      gop_stuff -> ref_count = 1;
4
5      /* Second, fill in *gop_stuff with whatever we want.
6         This part isn't shown here */
7
8      /* Third, store gop_stuff as the attribute value */
9      MPI_Comm_set_attr ( comm, gop_key, gop_stuff);
10     }
11     /* Then, in any case, use contents of *gop_stuff
12        to do the global op ... */
13 }
14
15 /* The following routine is called by MPI when a group is freed */
16
17 gop_stuff_destructor (comm, keyval, gop_stuff, extra)
18 MPI_Comm comm;
19 int keyval;
20 gop_stuff_type *gop_stuff;
21 void *extra;
22 {
23     if (keyval != gop_key) { /* abort -- programming error */ }
24
25     /* The group's being freed removes one reference to gop_stuff */
26     gop_stuff -> ref_count -= 1;
27
28     /* If no references remain, then free the storage */
29     if (gop_stuff -> ref_count == 0) {
30         free((void *)gop_stuff);
31     }
32 }
33
34 /* The following routine is called by MPI when a group is copied */
35 gop_stuff_copier (comm, keyval, extra, gop_stuff_in, gop_stuff_out, flag)
36 MPI_Comm comm;
37 int keyval;
38 gop_stuff_type *gop_stuff_in, *gop_stuff_out;
39 void *extra;
40 {
41     if (keyval != gop_key) { /* abort -- programming error */ }
42
43     /* The new group adds one reference to this gop_stuff */
44     gop_stuff -> ref_count += 1;
45     gop_stuff_out = gop_stuff_in;
46 }
47
48

```

6.8 Naming Objects

There are many occasions on which it would be useful to allow a user to associate a printable identifier with an MPI communicator, window, or datatype, for instance error reporting, debugging, and profiling. The names attached to opaque objects do not propagate when the object is duplicated or copied by MPI routines. For communicators this can be achieved using the following two functions.

`MPI_COMM_SET_NAME` (`comm`, `comm_name`)

INOUT	<code>comm</code>	communicator whose identifier is to be set (handle)
IN	<code>comm_name</code>	the character string which is remembered as the name (string)

`int MPI_Comm_set_name(MPI_Comm comm, char *comm_name)`

`MPI_COMM_SET_NAME(COMM, COMM_NAME, IERROR)`

INTEGER `COMM`, `IERROR`
 CHARACTER*(*) `COMM_NAME`

{void MPI::Comm::Set_name(const char* comm_name) (*binding deprecated, see Section 15.2*) }

`MPI_COMM_SET_NAME` allows a user to associate a name string with a communicator. The character string which is passed to `MPI_COMM_SET_NAME` will be saved inside the MPI library (so it can be freed by the caller immediately after the call, or allocated on the stack). Leading spaces in `name` are significant but trailing ones are not.

`MPI_COMM_SET_NAME` is a local (non-collective) operation, which only affects the name of the communicator as seen in the process which made the `MPI_COMM_SET_NAME` call. There is no requirement that the same (or any) name be assigned to a communicator in every process where it exists.

Advice to users. Since `MPI_COMM_SET_NAME` is provided to help debug code, it is sensible to give the same name to a communicator in all of the processes where it exists, to avoid confusion. (*End of advice to users.*)

The length of the name which can be stored is limited to the value of `MPI_MAX_OBJECT_NAME` in Fortran and `MPI_MAX_OBJECT_NAME-1` in C and C++ to allow for the null terminator. Attempts to put names longer than this will result in truncation of the name. `MPI_MAX_OBJECT_NAME` must have a value of at least 64.

Advice to users. Under circumstances of store exhaustion an attempt to put a name of any length could fail, therefore the value of `MPI_MAX_OBJECT_NAME` should be viewed only as a strict upper bound on the name length, not a guarantee that setting names of less than this length will always succeed. (*End of advice to users.*)

Advice to implementors. Implementations which pre-allocate a fixed size space for a name should use the length of that allocation as the value of `MPI_MAX_OBJECT_NAME`. Implementations which allocate space for the name from the heap should still define

1 MPI_MAX_OBJECT_NAME to be a relatively small value, since the user has to allocate
 2 space for a string of up to this size when calling MPI_COMM_GET_NAME. (*End of*
 3 *advice to implementors.*)
 4

5
 6 MPI_COMM_GET_NAME (comm, comm_name, resultlen)

7
 8 IN comm communicator whose name is to be returned (handle)
 9 OUT comm_name the name previously stored on the communicator, or
 10 an empty string if no such name exists (string)
 11 OUT resultlen length of returned name (integer)
 12

13 int MPI_Comm_get_name(MPI_Comm comm, char *comm_name, int *resultlen)

14 MPI_COMM_GET_NAME(COMM, COMM_NAME, RESULTLEN, IERROR)

15 INTEGER COMM, RESULTLEN, IERROR

16 CHARACTER*(*) COMM_NAME

17
 18 {void MPI::Comm::Get_name(char* comm_name, int& resultlen) const(*binding*
 19 *deprecated, see Section 15.2*) }
 20

21 MPI_COMM_GET_NAME returns the last name which has previously been associated
 22 with the given communicator. The name may be set and got from any language. The same
 23 name will be returned independent of the language used. name should be allocated so that
 24 it can hold a resulting string of length MPI_MAX_OBJECT_NAME characters.

25 MPI_COMM_GET_NAME returns a copy of the set name in name.

26 In C, a null character is additionally stored at name[resultlen]. resultlen cannot be
 27 larger than MPI_MAX_OBJECT_NAME-1. In Fortran, name is padded on the right with
 28 blank characters. resultlen cannot be larger than MPI_MAX_OBJECT_NAME.

29 If the user has not associated a name with a communicator, or an error occurs,
 30 MPI_COMM_GET_NAME will return an empty string (all spaces in Fortran, "" in C and
 31 C++). The three predefined communicators will have predefined names associated with
 32 them. Thus, the names of MPI_COMM_WORLD, MPI_COMM_SELF, and the communicator
 33 returned by MPI_COMM_GET_PARENT (if not MPI_COMM_NULL) will have the default of
 34 MPI_COMM_WORLD, MPI_COMM_SELF, and MPI_COMM_PARENT. The fact that the system
 35 may have chosen to give a default name to a communicator does not prevent the user from
 36 setting a name on the same communicator; doing this removes the old name and assigns
 37 the new one.
 38

39 *Rationale.* We provide separate functions for setting and getting the name of a com-
 40 municator, rather than simply providing a predefined attribute key for the following
 41 reasons:

- 42 • It is not, in general, possible to store a string as an attribute from Fortran.
- 43 • It is not easy to set up the delete function for a string attribute unless it is known
 44 to have been allocated from the heap.
- 45 • To make the attribute key useful additional code to call `strdup` is necessary. If
 46 this is not standardized then users have to write it. This is extra unneeded work
 47 which we can easily eliminate.
 48

- The Fortran binding is not trivial to write (it will depend on details of the Fortran compilation system), and will not be portable. Therefore it should be in the library rather than in user code.

(End of rationale.)

Advice to users. The above definition means that it is safe simply to print the string returned by `MPI_COMM_GET_NAME`, as it is always a valid string even if there was no name.

Note that associating a name with a communicator has no effect on the semantics of an MPI program, and will (necessarily) increase the store requirement of the program, since the names must be saved. Therefore there is no requirement that users use these functions to associate names with communicators. However debugging and profiling MPI applications may be made easier if names are associated with communicators, since the debugger or profiler should then be able to present information in a less cryptic manner. *(End of advice to users.)*

The following functions are used for setting and getting names of datatypes.

`MPI_TYPE_SET_NAME` (type, type_name)

INOUT	type	datatype whose identifier is to be set (handle)
IN	type_name	the character string which is remembered as the name (string)

```
int MPI_Type_set_name(MPI_Datatype type, char *type_name)
```

```
MPI_TYPE_SET_NAME(TYPE, TYPE_NAME, IERROR)
```

```
INTEGER TYPE, IERROR
CHARACTER*(*) TYPE_NAME
```

```
{void MPI::Datatype::Set_name(const char* type_name) (binding deprecated, see Section 15.2) }
```

`MPI_TYPE_GET_NAME` (type, type_name, resultlen)

IN	type	datatype whose name is to be returned (handle)
OUT	type_name	the name previously stored on the datatype, or a empty string if no such name exists (string)
OUT	resultlen	length of returned name (integer)

```
int MPI_Type_get_name(MPI_Datatype type, char *type_name, int *resultlen)
```

```
MPI_TYPE_GET_NAME(TYPE, TYPE_NAME, RESULTLEN, IERROR)
```

```
INTEGER TYPE, RESULTLEN, IERROR
CHARACTER*(*) TYPE_NAME
```

```
{void MPI::Datatype::Get_name(char* type_name, int& resultlen) const (binding deprecated, see Section 15.2) }
```

1 Named predefined datatypes have the default names of the datatype name. For exam-
 2 ple, MPI_WCHAR has the default name of MPI_WCHAR.

3 The following functions are used for setting and getting names of windows.

4
 5 MPI_WIN_SET_NAME (win, win_name)

6
 7 INOUT win window whose identifier is to be set (handle)
 8 IN win_name the character string which is remembered as the name
 9 (string)

10
 11 int MPI_Win_set_name(MPI_Win win, char *win_name)

12 MPI_WIN_SET_NAME(WIN, WIN_NAME, IERROR)

13 INTEGER WIN, IERROR
 14 CHARACTER*(*) WIN_NAME

15
 16 {void MPI::Win::Set_name(const char* win_name) (*binding deprecated, see*
 17 *Section 15.2*) }

18
 19
 20 MPI_WIN_GET_NAME (win, win_name, resultlen)

21
 22 IN win window whose name is to be returned (handle)
 23 OUT win_name the name previously stored on the window, or a empty
 24 string if no such name exists (string)
 25 OUT resultlen length of returned name (integer)

26
 27 int MPI_Win_get_name(MPI_Win win, char *win_name, int *resultlen)

28 MPI_WIN_GET_NAME(WIN, WIN_NAME, RESULTLEN, IERROR)

29 INTEGER WIN, RESULTLEN, IERROR
 30 CHARACTER*(*) WIN_NAME

31
 32 {void MPI::Win::Get_name(char* win_name, int& resultlen) const (*binding*
 33 *deprecated, see Section 15.2*) }

34 35 36 37 6.9 Formalizing the Loosely Synchronous Model

38
 39 In this section, we make further statements about the loosely synchronous model, with
 40 particular attention to intra-communication.

41 42 6.9.1 Basic Statements

43
 44 When a caller passes a communicator (that contains a context and group) to a callee, that
 45 communicator must be free of side effects throughout execution of the subprogram: there
 46 should be no active operations on that communicator that might involve the process. This
 47 provides one model in which libraries can be written, and work “safely.” For libraries
 48 so designated, the callee has permission to do whatever communication it likes with the

communicator, and under the above guarantee knows that no other communications will interfere. Since we permit good implementations to create new communicators without synchronization (such as by preallocated contexts on communicators), this does not impose a significant overhead.

This form of safety is analogous to other common computer-science usages, such as passing a descriptor of an array to a library routine. The library routine has every right to expect such a descriptor to be valid and modifiable.

6.9.2 Models of Execution

In the loosely synchronous model, transfer of control to a **parallel procedure** is effected by having each executing process invoke the procedure. The invocation is a collective operation: it is executed by all processes in the execution group, and invocations are similarly ordered at all processes. However, the invocation need not be synchronized.

We say that a parallel procedure is *active* in a process if the process belongs to a group that may collectively execute the procedure, and some member of that group is currently executing the procedure code. If a parallel procedure is active in a process, then this process may be receiving messages pertaining to this procedure, even if it does not currently execute the code of this procedure.

Static communicator allocation

This covers the case where, at any point in time, at most one invocation of a parallel procedure can be active at any process, and the group of executing processes is fixed. For example, all invocations of parallel procedures involve all processes, processes are single-threaded, and there are no recursive invocations.

In such a case, a communicator can be statically allocated to each procedure. The static allocation can be done in a preamble, as part of initialization code. If the parallel procedures can be organized into libraries, so that only one procedure of each library can be concurrently active in each processor, then it is sufficient to allocate one communicator per library.

Dynamic communicator allocation

Calls of parallel procedures are well-nested if a new parallel procedure is always invoked in a subset of a group executing the same parallel procedure. Thus, processes that execute the same parallel procedure have the same execution stack.

In such a case, a new communicator needs to be dynamically allocated for each new invocation of a parallel procedure. The allocation is done by the caller. A new communicator can be generated by a call to `MPI_COMM_DUP`, if the callee execution group is identical to the caller execution group, or by a call to `MPI_COMM_SPLIT` if the caller execution group is split into several subgroups executing distinct parallel routines. The new communicator is passed as an argument to the invoked routine.

The need for generating a new communicator at each invocation can be alleviated or avoided altogether in some cases: If the execution group is not split, then one can allocate a stack of communicators in a preamble, and next manage the stack in a way that mimics the stack of recursive calls.

1 One can also take advantage of the well-ordering property of communication to avoid
2 confusing caller and callee communication, even if both use the same communicator. To do
3 so, one needs to abide by the following two rules:

- 4 • messages sent before a procedure call (or before a return from the procedure) are also
5 received before the matching call (or return) at the receiving end;
- 6 • messages are always selected by source (no use is made of `MPI_ANY_SOURCE`).

ticket0. 9 The General [c]Case

10
11 In the general case, there may be multiple concurrently active invocations of the same
12 parallel procedure within the same group; invocations may not be well-nested. A new
13 communicator needs to be created for each invocation. It is the user's responsibility to make
14 sure that, should two distinct parallel procedures be invoked concurrently on overlapping
15 sets of processes, then communicator creation be properly coordinated.

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

Bibliography

- [1] Purushotham V. Bangalore, Nathan E. Doss, and Anthony Skjellum. MPI++: Issues and Features. In *OON-SKI '94*, page in press, 1994. [6.1](#)
- [2] D. Feitelson. Communicators: Object-based multiparty interactions for parallel programming. Technical Report 91-12, Dept. Computer Science, The Hebrew University of Jerusalem, November 1991. [6.1.2](#)
- [3] A. Skjellum and A. Leung. Zipcode: a portable multicomputer communication library atop the reactive kernel. In D. W. Walker and Q. F. Stout, editors, *Proceedings of the Fifth Distributed Memory Concurrent Computing Conference*, pages 767–776. IEEE Press, 1990. [6.1.2](#)
- [4] Anthony Skjellum, Nathan E. Doss, and Purushotham V. Bangalore. Writing Libraries in MPI. In Anthony Skjellum and Donna S. Reese, editors, *Proceedings of the Scalable Parallel Libraries Conference*, pages 166–173. IEEE Computer Society Press, October 1993. [6.1](#)
- [5] Anthony Skjellum, Steven G. Smith, Nathan E. Doss, Alvin P. Leung, and Manfred Morari. The Design and Evolution of Zipcode. *Parallel Computing*, 20(4):565–596, April 1994. [6.1.2](#), [6.5.6](#)

Index

- comm_copy_attr_fn, [41](#), [42](#)
- comm_delete_attr_fn, [42](#), [43](#)
- CONST:false, [32](#), [44](#), [47](#), [50](#)
- CONST:flag = 0, [41](#)
- CONST:flag = 1, [41](#)
- CONST:INTEGER, [39](#)
- CONST:MPI::Comm, [8](#), [13–16](#), [19](#), [22](#), [23](#), [32–35](#), [41](#), [43](#), [44](#)
- CONST:MPI::Group, [6](#), [6](#), [7–12](#), [16](#), [33](#)
- CONST:MPI::Win, [45–47](#), [56](#)
- CONST:MPI_ANY_SOURCE, [58](#)
- CONST:MPI_Comm, [8](#), [13–16](#), [19](#), [22](#), [23](#), [32–35](#), [41](#), [43](#), [44](#)
- CONST:MPI_COMM_NULL, [5](#), [16](#), [17](#), [19](#), [20](#), [22](#), [23](#), [54](#)
- CONST:MPI_COMM_PARENT, [54](#)
- CONST:MPI_COMM_SELF, [5](#), [39](#), [54](#)
- CONST:MPI_COMM_TYPE_PROCESS, [22](#)
- CONST:MPI_COMM_TYPE_SHM, [22](#)
- CONST:MPI_COMM_WORLD, [5–7](#), [13](#), [14](#), [25](#), [34](#), [35](#), [54](#)
- CONST:MPI_CONGRUENT, [14](#), [32](#)
- CONST:MPI_ERR_KEYVAL, [50](#)
- CONST:MPI_Group, [6](#), [6](#), [7–12](#), [16](#), [33](#)
- CONST:MPI_GROUP_EMPTY, [4](#), [9](#), [10](#), [16](#), [17](#)
- CONST:MPI_GROUP_NULL, [4](#), [12](#)
- CONST:MPI_IDENT, [7](#), [14](#)
- CONST:MPI_KEYVAL_INVALID, [42–44](#)
- CONST:MPI_MAX_OBJECT_NAME, [53](#), [54](#)
- CONST:MPI_PROC_NULL, [7](#)
- CONST:MPI_SIMILAR, [7](#), [14](#), [32](#)
- CONST:MPI_SUCCESS, [41–43](#), [45](#), [46](#), [48](#), [49](#)
- CONST:MPI_UNDEFINED, [6](#), [7](#), [20](#)
- CONST:MPI_UNEQUAL, [7](#), [14](#), [32](#)
- CONST:MPI_WCHAR, [56](#)
- CONST:MPI_Win, [45–47](#), [56](#)
- CONST:true, [32](#)
- CONST:void *, [39](#), [42](#)
- CONST:void*, [44](#)
- CONST:void**, [44](#)
- EXAMPLES:Intercommunicator, [18](#), [20](#)
- EXAMPLES:MPI_Comm_create, [18](#)
- EXAMPLES:MPI_Comm_group, [18](#)
- EXAMPLES:MPI_Comm_remote_size, [20](#)
- EXAMPLES:MPI_Comm_split, [20](#)
- EXAMPLES:MPI_Group_free, [18](#)
- EXAMPLES:MPI_Group_incl, [18](#)
- MPI_ATTR_DELETE, [45](#), [50](#)
- MPI_ATTR_GET, [44](#), [50](#)
- MPI_ATTR_PUT, [43](#), [50](#)
- MPI_CART_CREATE, [31](#)
- MPI_COMM_COMPARE, [32](#)
- MPI_COMM_COMPARE(comm1, comm2, result), [14](#)
- MPI_COMM_CREATE, [12](#), [17–20](#)
- MPI_COMM_CREATE(comm, group, newcomm), [16](#), [19](#)
- MPI_COMM_CREATE_KEYVAL, [39](#), [42](#), [50](#)
- MPI_COMM_CREATE_KEYVAL(comm_copy_attr_fn, comm_delete_attr_fn, comm_keyval, extra_state), [40](#)
- MPI_COMM_DELETE_ATTR, [39](#), [42](#), [43](#), [50](#)
- MPI_COMM_DELETE_ATTR(comm, comm_keyval), [44](#)
- MPI_COMM_DISCONNECT, [50](#)
- MPI_COMM_DUP, [8](#), [12](#), [15–17](#), [23](#), [33](#), [35](#), [39](#), [41](#), [45](#), [50](#), [57](#)
- MPI_COMM_DUP(comm, newcomm), [15](#)
- MPI_COMM_DUP_FN, [41](#), [41](#), [42](#)
- MPI_COMM_FREE, [12](#), [16](#), [23](#), [33](#), [35](#), [42](#), [43](#), [45](#), [50](#)
- MPI_COMM_FREE(comm), [23](#)
- MPI_COMM_FREE_KEYVAL, [39](#), [50](#)
- MPI_COMM_FREE_KEYVAL(comm_keyval), [43](#)

MPI_COMM_GET_ATTR, 39 , 43 , 50	MPI_GROUP_INTERSECTION(group1, group2, newgroup), 9	2
MPI_COMM_GET_ATTR(comm, comm_keyval, attribute_val, flag), 44	MPI_GROUP_RANGE_EXCL, 12	3
MPI_COMM_GET_NAME, 54 , 55	MPI_GROUP_RANGE_EXCL(group, n, ranges, newgroup), 11	5
MPI_COMM_GET_NAME(comm, comm_name, resultlen), 54	MPI_GROUP_RANGE_INCL, 11	6
MPI_COMM_GET_PARENT, 54	MPI_GROUP_RANGE_INCL(group, n, ranges, newgroup), 11	8
MPI_COMM_GROUP, 6 , 8 , 12–14 , 32	MPI_GROUP_RANK, 14	9
MPI_COMM_GROUP(comm, group), 8	MPI_GROUP_RANK(group, rank), 6	10
MPI_COMM_NULL_COPY_FN, 41 , 41 , 42	MPI_GROUP_SIZE, 13	11
MPI_COMM_NULL_DELETE_FN, 42 , 42	MPI_GROUP_SIZE(group, size), 6	12
MPI_COMM_RANK, 13 , 32	MPI_GROUP_TRANSLATE_RANKS, 7	13
MPI_COMM_RANK(comm, rank), 13	MPI_GROUP_TRANSLATE_RANKS(group1, n, ranks1, group2, ranks2), 7	15
MPI_COMM_REMOTE_GROUP(comm, group), 33	MPI_GROUP_UNION(group1, group2, newgroup), 8	17
MPI_COMM_REMOTE_SIZE, 33	MPI_INIT, 5	18
MPI_COMM_REMOTE_SIZE(comm, size), 33	MPI_INIT_THREAD, 5	19
MPI_COMM_SET_ATTR, 39 , 42 , 50	MPI_INTERCOMM_CREATE, 15 , 33 , 35	20
MPI_COMM_SET_ATTR(comm, comm_keyval, attribute_val), 43	MPI_INTERCOMM_CREATE(local_comm, local_leader, peer_comm, remote_leader, tag, newintercomm), 34	23
MPI_COMM_SET_NAME, 53	MPI_INTERCOMM_MERGE, 31 , 33–35	24
MPI_COMM_SET_NAME(comm, comm_name), 53	MPI_INTERCOMM_MERGE(intercomm, high, newintracomm), 35	26
MPI_COMM_SIZE, 13 , 14 , 32	MPI_KEYVAL_CREATE, 40 , 42	27
MPI_COMM_SIZE(comm, size), 13	MPI_KEYVAL_FREE, 43 , 50	28
MPI_COMM_SPLIT, 17 , 19–21 , 57	MPI_NULL_COPY_FN, 41	29
MPI_COMM_SPLIT(comm, color, key, newcomm), 19 , 19	MPI_NULL_DELETE_FN, 42	30
MPI_COMM_SPLIT_TYPE(comm, type, key, newcomm), 22	MPI_TYPE_CREATE_KEYVAL, 39 , 50	31
MPI_COMM_TEST_INTER, 31	MPI_TYPE_CREATE_KEYVAL(type_copy_attr_fn, type_delete_attr_fn, type_keyval, extra_state), 48	34
MPI_COMM_TEST_INTER(comm, flag), 32	MPI_TYPE_DELETE_ATTR, 39 , 50	35
MPI_DUP_FN, 41	MPI_TYPE_DELETE_ATTR(type, type_keyval), 50	37
MPI_GROUP_COMPARE, 10	MPI_TYPE_DUP_FN, 48 , 48	38
MPI_GROUP_COMPARE(group1, group2, result), 7	MPI_TYPE_FREE, 49	39
MPI_GROUP_DIFFERENCE(group1, group2, newgroup), 9	MPI_TYPE_FREE_KEYVAL, 39 , 50	40
MPI_GROUP_EXCL, 10 , 12	MPI_TYPE_FREE_KEYVAL(type_keyval), 49	42
MPI_GROUP_EXCL(group, n, ranks, newgroup), 10	MPI_TYPE_GET_ATTR, 39 , 50	43
MPI_GROUP_FREE, 12–14	MPI_TYPE_GET_ATTR(type, type_keyval, attribute_val, flag), 50	45
MPI_GROUP_FREE(group), 12	MPI_TYPE_GET_NAME(type, type_name, resultlen), 55	47
MPI_GROUP_INCL, 10 , 11	MPI_TYPE_NULL_COPY_FN, 48 , 48	48
MPI_GROUP_INCL(group, n, ranks, newgroup), 10		

```

1  MPI_TYPE_NULL_DELETE_FN, 48      TYPEDEF:MPI_Win_delete_attr_function(MPI_Win win,
2  MPI_TYPE_SET_ATTR, 39, 50          int win_keyval, void *attribute_val,
3  MPI_TYPE_SET_ATTR(type, type_keyval,  void *extra_state), 46
4      attribute_val), 49
5  MPI_TYPE_SET_NAME (type, type_name),
6      55
7  MPI_WIN_ALLOCATE_SHARED, 22
8  MPI_WIN_CREATE_KEYVAL, 39, 50
9  MPI_WIN_CREATE_KEYVAL(win_copy_attr_fn,
10     win_delete_attr_fn, win_keyval, ex-
11     tra_state), 45
12  MPI_WIN_DELETE_ATTR, 39, 50
13  MPI_WIN_DELETE_ATTR(win, win_keyval),
14     47
15  MPI_WIN_DUP_FN, 45, 45
16  MPI_WIN_FREE, 46
17  MPI_WIN_FREE_KEYVAL, 39, 50
18  MPI_WIN_FREE_KEYVAL(win_keyval), 46
19  MPI_WIN_GET_ATTR, 39, 50
20  MPI_WIN_GET_ATTR(win, win_keyval, at-
21     tribute_val, flag), 47
22  MPI_WIN_GET_NAME (win, win_name, re-
23     sultlen), 56
24  MPI_WIN_NULL_COPY_FN, 45, 45
25  MPI_WIN_NULL_DELETE_FN, 45
26  MPI_WIN_SET_ATTR, 39, 50
27  MPI_WIN_SET_ATTR(win, win_keyval, at-
28     tribute_val), 47
29  MPI_WIN_SET_NAME (win, win_name), 56
30
31  TYPEDEF:MPI_Comm_copy_attr_function(MPI_Comm old-
32     comm, int comm_keyval, void *ex-
33     tra_state, void *attribute_val_in, void *at-
34     tribute_val_out, int *flag), 41
35  TYPEDEF:MPI_Comm_delete_attr_function(MPI_Comm
36     comm, int comm_keyval, void *at-
37     tribute_val, void *extra_state), 41
38  TYPEDEF:MPI_Type_copy_attr_function(MPI_Datatype old-
39     type, int type_keyval, void *extra_state,
40     void *attribute_val_in, void *attribute_val_out,
41     int *flag), 48
42  TYPEDEF:MPI_Type_delete_attr_function(MPI_Datatype type,
43     int type_keyval, void *attribute_val,
44     void *extra_state), 48
45  TYPEDEF:MPI_Win_copy_attr_function(MPI_Win old-
46     win, int win_keyval, void *extra_state,
47     void *attribute_val_in, void *attribute_val_out,
48     int *flag), 46

```