

*D R A F T*

# Document for a Standard Message-Passing Interface

Message Passing Interface Forum

July 30, 2009

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.

# Chapter 4

## Datatypes

Basic datatypes were introduced in Section 3.2.2 Message Data on page 29 and in Section 3.3 Data Type Matching and Data Conversion on page 36. In this chapter, this model is extended to describe any data layout. We consider general datatypes that allow one to transfer efficiently heterogeneous and noncontiguous data. We conclude with the description of calls for explicit packing and unpacking of messages.

### 4.1 Derived Datatypes

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

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

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

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

- A sequence of basic datatypes

- A sequence of integer (byte) displacements

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

Let

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

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

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

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

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

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

General datatypes can be used in all send and receive operations. We discuss, in Section 4.1.11, the case where the second argument `count` has value  $> 1$ .

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

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

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

then

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

If  $type_i$  requires alignment to a byte address that is a multiple of  $k_i$ , then  $\epsilon$  is the least [nonnegative]non-negative increment needed to round  $extent(Typemap)$  to the next multiple of  $\max_i k_i$ . The complete definition of **extent** is given on page 20.

**Example 4.1** Assume that  $Type = \{(double, 0), (char, 8)\}$  (a `double` at displacement zero, followed by a `char` at displacement eight). Assume, furthermore, that doubles have to be strictly aligned at addresses that are multiples of eight. Then, the extent of this datatype is 16 (9 rounded to the next multiple of 8). A datatype that consists of a character immediately followed by a double will also have an extent of 16.

*Rationale.* The definition of extent is motivated by the assumption that the amount of padding added at the end of each structure in an array of structures is the least needed to fulfill alignment constraints. More explicit control of the extent is provided in Section 4.1.6. Such explicit control is needed in cases where the assumption does not hold, for example, where union types are used. (*End of rationale.*)

#### 4.1.1 Type Constructors with Explicit Addresses

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

#### 4.1.2 Datatype Constructors

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

`MPI_TYPE_CONTIGUOUS(count, oldtype, newtype)`

IN	count	replication count ([nonnegative]non-negative integer)	ticket74.
IN	oldtype	old datatype (handle)	
OUT	newtype	new datatype (handle)	

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

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

```
{MPI::Datatype MPI::Datatype::Create_contiguous(int count) const (binding
                        deprecated, see Section ??) }
```

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

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

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

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

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

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

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

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

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

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

ticket74.	IN	count	number of blocks ([nonnegative]non-negative integer)
ticket74.	IN	blocklength	number of elements in each block ([nonnegative]non-negative integer)
	IN	stride	number of elements between start of each block (integer)
	IN	oldtype	old datatype (handle)
	OUT	newtype	new datatype (handle)

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

IN	count	number of blocks ([nonnegative]non-negative integer)	ticket74.
IN	blocklength	number of elements in each block ([nonnegative]non-negative integer)	ticket74.
IN	stride	number of bytes between start of each block (integer)	
IN	oldtype	old datatype (handle)	
OUT	newtype	new datatype (handle)	

```

1  int MPI_Type_create_hvector(int count, int blocklength, MPI_Aint stride,
2      MPI_Datatype oldtype, MPI_Datatype *newtype)
3
4  MPI_TYPE_CREATE_HVECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE,
5      IERROR)
6      INTEGER COUNT, BLOCKLENGTH, OLDTYPE, NEWTYPE, IERROR
7      INTEGER(KIND=MPI_ADDRESS_KIND) STRIDE

```

```

8  {MPI::Datatype MPI::Datatype::Create_hvector(int count, int blocklength,
9      MPI::Aint stride) const (binding deprecated, see Section ??) }

```

This function replaces MPI\_TYPE\_HVECTOR, whose use is deprecated. See also Chapter 15.

Assume that `oldtype` has type map,

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

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

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

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



MPI\_TYPE\_INDEXED( count, array\_of\_blocklengths, array\_of\_displacements, oldtype, newtype)

ticket74.	IN	count	number of blocks – also number of entries in array_of_displacements and array_of_blocklengths ([nonnegative]non-negative integer)	1
				2
ticket74.	IN	array_of_blocklengths	number of elements per block (array of [nonnegative]non-negative integers)	3
				4
	IN	array_of_displacements	displacement for each block, in multiples of oldtype extent (array of integer)	5
				6
	IN	oldtype	old datatype (handle)	7
				8
	OUT	newtype	new datatype (handle)	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

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

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

```
{MPI::Datatype MPI::Datatype::Create_indexed(int count,
const int array_of_blocklengths[],
const int array_of_displacements[]) const (binding deprecated, see
Section ??) }
```

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

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

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

In general, assume that oldtype has type map,

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

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

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

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

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

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

and

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

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

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

<b>IN</b>	<b>count</b>	number of blocks — also number of entries in <code>array_of_displacements</code> and <code>array_of_blocklengths</code> ( <b>[nonnegative]non-negative integer</b> )
<b>IN</b>	<b>array_of_blocklengths</b>	number of elements in each block (array of <b>[nonnegative]non-negative integers</b> )
<b>IN</b>	<b>array_of_displacements</b>	byte displacement of each block (array of integer)
<b>IN</b>	<b>oldtype</b>	old datatype (handle)
<b>OUT</b>	<b>newtype</b>	new datatype (handle)

```

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

```

```

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

```

```

{MPI::Datatype MPI::Datatype::Create_hindexed(int count,
        const int array_of_blocklengths[],
        const MPI::Aint array_of_displacements[]) const (binding deprecated, see Section ??) }

```

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

Assume that `oldtype` has type map,

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

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

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

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

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

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

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

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

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

```
{MPI::Datatype MPI::Datatype::Create_indexed_block(int count,
    int blocklength, const int array_of_displacements[]) const
    (binding deprecated, see Section ??) }
```

ticket150.

ticket150.

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

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

ticket74.	8	IN	count	number of blocks ([nonnegative]non-negative integer) — also number of entries in arrays <code>array_of_types</code> , <code>array_of_displacements</code> and <code>array_of_blocklengths</code>
ticket74.	11	IN	array_of_blocklength	number of elements in each block (array of [nonnegative]non-negative integer)
	13	IN	array_of_displacements	byte displacement of each block (array of integer)
	15	IN	array_of_types	type of elements in each block (array of handles to datatype objects)
	17	OUT	newtype	new datatype (handle)

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

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

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

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

**Example 4.6** Let `type1` have type map,

`{(double, 0), (char, 8)}`,

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

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

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

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

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

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

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

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

#### 4.1.3 Subarray Datatype Constructor

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

IN	<code>ndims</code>	number of array dimensions (positive integer)
IN	<code>array_of_sizes</code>	number of elements of type <code>oldtype</code> in each dimension of the full array (array of positive integers)
IN	<code>array_of_subsizes</code>	number of elements of type <code>oldtype</code> in each dimension of the subarray (array of positive integers)
IN	<code>array_of_starts</code>	starting coordinates of the subarray in each dimension (array of <span style="color: blue;">nonnegative</span> <span style="color: red;">non-negative</span> integers)
IN	<code>order</code>	array storage order flag (state)
IN	<code>oldtype</code>	array element datatype (handle)
OUT	<code>newtype</code>	new datatype (handle)

ticket74.

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

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

ticket150.

```

1  {MPI::Datatype MPI::Datatype::Create_subarray(int ndims,
2      const int array_of_sizes[], const int array_of_subsizes[],
3      const int array_of_starts[], int order) const (binding
4      deprecated, see Section ??) }
5

```

The subarray type constructor creates an MPI datatype describing an  $n$ -dimensional subarray of an  $n$ -dimensional array. The subarray may be situated anywhere within the full array, and may be of any nonzero size up to the size of the larger array as long as it is confined within this array. This type constructor facilitates creating filetypes to access arrays distributed in blocks among processes to a single file that contains the global array, see MPI I/O, especially Section 13.1.1 on page 383.

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

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

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

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

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

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

**MPI\_ORDER\_C** The ordering used by C arrays, (i.e., row-major order)

**MPI\_ORDER\_FORTRAN** The ordering used by Fortran arrays, (i.e., column-major order)

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

```

newtype = Subarray(ndims, {size0, size1, ..., sizendims-1},
                    {subsize0, subsize1, ..., subsizendims-1},
                    {start0, start1, ..., startndims-1}, oldtype)

```

Let the typemap of `oldtype` have the form:

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

where  $type_i$  is a predefined MPI datatype, and let  $ex$  be the extent of `oldtype`. Then we define the `Subarray()` function recursively using the following three equations. Equation 4.2 defines the base step. Equation 4.3 defines the recursion step when `order = MPI_ORDER_FORTRAN`, and Equation 4.4 defines the recursion step when `order = MPI_ORDER_C`.

$$\begin{aligned}
& \text{Subarray}(1, \{size_0\}, \{subsize_0\}, \{start_0\}, \\
& \quad \{(type_0, disp_0), (type_1, disp_1), \dots, (type_{n-1}, disp_{n-1})\}) \\
& = \{(\text{MPI\_LB}, 0), \\
& \quad (type_0, disp_0 + start_0 \times ex), \dots, (type_{n-1}, disp_{n-1} + start_0 \times ex), \\
& \quad (type_0, disp_0 + (start_0 + 1) \times ex), \dots, (type_{n-1}, \\
& \quad \quad disp_{n-1} + (start_0 + 1) \times ex), \dots \\
& \quad (type_0, disp_0 + (start_0 + subsize_0 - 1) \times ex), \dots, \\
& \quad \quad (type_{n-1}, disp_{n-1} + (start_0 + subsize_0 - 1) \times ex), \\
& \quad (\text{MPI\_UB}, size_0 \times ex)\}
\end{aligned} \tag{4.2}$$

$$\begin{aligned}
& \text{Subarray}(ndims, \{size_0, size_1, \dots, size_{ndims-1}\}, \\
& \quad \{subsize_0, subsize_1, \dots, subsize_{ndims-1}\}, \\
& \quad \{start_0, start_1, \dots, start_{ndims-1}\}, \text{oldtype}) \\
& = \text{Subarray}(ndims - 1, \{size_1, size_2, \dots, size_{ndims-1}\}, \\
& \quad \{subsize_1, subsize_2, \dots, subsize_{ndims-1}\}, \\
& \quad \{start_1, start_2, \dots, start_{ndims-1}\}, \\
& \quad \text{Subarray}(1, \{size_0\}, \{subsize_0\}, \{start_0\}, \text{oldtype}))
\end{aligned} \tag{4.3}$$

$$\begin{aligned}
& \text{Subarray}(ndims, \{size_0, size_1, \dots, size_{ndims-1}\}, \\
& \quad \{subsize_0, subsize_1, \dots, subsize_{ndims-1}\}, \\
& \quad \{start_0, start_1, \dots, start_{ndims-1}\}, \text{oldtype}) \\
& = \text{Subarray}(ndims - 1, \{size_0, size_1, \dots, size_{ndims-2}\}, \\
& \quad \{subsize_0, subsize_1, \dots, subsize_{ndims-2}\}, \\
& \quad \{start_0, start_1, \dots, start_{ndims-2}\}, \\
& \quad \text{Subarray}(1, \{size_{ndims-1}\}, \{subsize_{ndims-1}\}, \{start_{ndims-1}\}, \text{oldtype}))
\end{aligned} \tag{4.4}$$

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

#### 4.1.4 Distributed Array Datatype Constructor

The distributed array type constructor supports HPF-like [1] data distributions. However, unlike in HPF, the storage order may be specified for C arrays as well as for Fortran arrays.

*Advice to users.* One can create an HPF-like file view using this type constructor as follows. Complementary filetypes are created by having every process of a group call this constructor with identical arguments (with the exception of `rank` which should be set appropriately). These filetypes (along with identical `disp` and `etype`) are then used to define the view (via `MPI_FILE_SET_VIEW`), see MPI I/O, especially Section 13.1.1 on page 383 and Section 13.3 on page 395. Using this view, a collective data access operation (with identical offsets) will yield an HPF-like distribution pattern. (*End of advice to users.*)

```

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

```

```

26 int MPI_Type_create_darray(int size, int rank, int ndims,
27     int array_of_gsizes[], int array_of_distrib[], int
28     array_of_dargs[], int array_of_psize[], int order,
29     MPI_Datatype oldtype, MPI_Datatype *newtype)
30
31 MPI_TYPE_CREATE_DARRAY(SIZE, RANK, NDIMS, ARRAY_OF_GSIZES,
32     ARRAY_OF_DISTRIBS, ARRAY_OF_DARGS, ARRAY_OF_PSIZE, ORDER,
33     OLDTYPE, NEWTYPE, IERROR)
34
35 INTEGER SIZE, RANK, NDIMS, ARRAY_OF_GSIZES(*), ARRAY_OF_DISTRIBS(*),
36 ARRAY_OF_DARGS(*), ARRAY_OF_PSIZE(*), ORDER, OLDTYPE, NEWTYPE, IERROR

```

```

37 {MPI::Datatype MPI::Datatype::Create_darray(int size, int rank, int ndims,
38     const int array_of_gsizes[], const int array_of_distrib[],
39     const int array_of_dargs[], const int array_of_psize[],
40     int order) const (binding deprecated, see Section ??) }

```

MPI\_TYPE\_CREATE\_DARRAY can be used to generate the datatypes corresponding to the distribution of an  $ndims$ -dimensional array of `oldtype` elements onto an  $ndims$ -dimensional grid of logical processes. Unused dimensions of `array_of_psize` should be set to 1. (See Example 4.7, page 17.) For a call to MPI\_TYPE\_CREATE\_DARRAY to be correct, the equation  $\prod_{i=0}^{ndims-1} array\_of\_psize[i] = size$  must be satisfied. The ordering of processes in the process grid is assumed to be row-major, as in the case of virtual Cartesian process topologies.

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



corresponding process topology functions, see Chapter 7 on page 251. (*End of advice to users.*)

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

- MPI\_DISTRIBUTE\_BLOCK - Block distribution
- MPI\_DISTRIBUTE\_CYCLIC - Cyclic distribution
- MPI\_DISTRIBUTE\_NONE - Dimension not distributed.

The constant MPI\_DISTRIBUTE\_DFLT\_DARG specifies a default distribution argument. The distribution argument for a dimension that is not distributed is ignored. For any dimension  $i$  in which the distribution is MPI\_DISTRIBUTE\_BLOCK, it is erroneous to specify  $\text{array\_of\_dargs}[i] * \text{array\_of\_psizes}[i] < \text{array\_of\_gsizes}[i]$ .

For example, the HPF layout `ARRAY(CYCLIC(15))` corresponds to MPI\_DISTRIBUTE\_CYCLIC with a distribution argument of 15, and the HPF layout `ARRAY(BLOCK)` corresponds to MPI\_DISTRIBUTE\_BLOCK with a distribution argument of MPI\_DISTRIBUTE\_DFLT\_DARG.

The `order` argument is used as in MPI\_TYPE\_CREATE\_SUBARRAY to specify the storage order. Therefore, arrays described by this type constructor may be stored in Fortran (column-major) or C (row-major) order. Valid values for `order` are MPI\_ORDER\_FORTRAN and MPI\_ORDER\_C.

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

Without loss of generality, it suffices to define the typemap for the MPI\_DISTRIBUTE\_CYCLIC case where MPI\_DISTRIBUTE\_DFLT\_DARG is not used.

MPI\_DISTRIBUTE\_BLOCK and MPI\_DISTRIBUTE\_NONE can be reduced to the MPI\_DISTRIBUTE\_CYCLIC case for dimension  $i$  as follows.

MPI\_DISTRIBUTE\_BLOCK with  $\text{array\_of\_dargs}[i]$  equal to MPI\_DISTRIBUTE\_DFLT\_DARG is equivalent to MPI\_DISTRIBUTE\_CYCLIC with  $\text{array\_of\_dargs}[i]$  set to

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

If  $\text{array\_of\_dargs}[i]$  is not MPI\_DISTRIBUTE\_DFLT\_DARG, then MPI\_DISTRIBUTE\_BLOCK and MPI\_DISTRIBUTE\_CYCLIC are equivalent.

MPI\_DISTRIBUTE\_NONE is equivalent to MPI\_DISTRIBUTE\_CYCLIC with  $\text{array\_of\_dargs}[i]$  set to  $\text{array\_of\_gsizes}[i]$ .

Finally, MPI\_DISTRIBUTE\_CYCLIC with  $\text{array\_of\_dargs}[i]$  equal to MPI\_DISTRIBUTE\_DFLT\_DARG is equivalent to MPI\_DISTRIBUTE\_CYCLIC with  $\text{array\_of\_dargs}[i]$  set to 1.

For MPI\_ORDER\_FORTRAN, an  $\text{ndims}$ -dimensional distributed array (`newtype`) is defined by the following code fragment:

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

```

1         oldtype[i]);
2     }
3     newtype = oldtype[ndims];
4
5     For MPI_ORDER_C, the code is:
6
7     oldtype[0] = oldtype;
8     for ( i = 0; i < ndims; i++ ) {
9         oldtype[i + 1] = cyclic(array_of_dargs[ndims - i - 1],
10                                array_of_gsizes[ndims - i - 1],
11                                r[ndims - i - 1],
12                                array_of_psize[ndims - i - 1],
13                                oldtype[i]);
14     }
15     newtype = oldtype[ndims];
16

```

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

```

20     t_rank = rank;
21     t_size = 1;
22     for (i = 0; i < ndims; i++)
23         t_size *= array_of_psize[i];
24     for (i = 0; i < ndims; i++) {
25         t_size = t_size / array_of_psize[i];
26         r[i] = t_rank / t_size;
27         t_rank = t_rank % t_size;
28     }
29

```

Let the typemap of oldtype have the form:

```

31     {(type0, disp0), (type1, disp1), ..., (typen-1, dispn-1)}
```

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

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

```

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

```

```

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

```

where *count* is defined by this code fragment:

```

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

```

Here, *nblocks* is the number of blocks that must be distributed among the processors. Finally, *darg<sub>last</sub>* is defined by this code fragment:

```

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

```

**Example 4.7** Consider generating the filetypes corresponding to the HPF distribution:

```

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

```

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

```

1      ndims = 3
2      array_of_gsizes(1) = 100
3      array_of_distribs(1) = MPI_DISTRIBUTE_CYCLIC
4      array_of_dargs(1) = 10
5      array_of_gsizes(2) = 200
6      array_of_distribs(2) = MPI_DISTRIBUTE_NONE
7      array_of_dargs(2) = 0
8      array_of_gsizes(3) = 300
9      array_of_distribs(3) = MPI_DISTRIBUTE_BLOCK
10     array_of_dargs(3) = [ticket116.] [MPI_DISTRIBUTE_DFLT_ARG] MPI_DISTRIBUTE_DFLT_DARG
11     array_of_psize(1) = 2
12     array_of_psize(2) = 1
13     array_of_psize(3) = 3
14     call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
15     call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
16     call MPI_TYPE_CREATE_DARRAY(size, rank, ndims, array_of_gsizes, &
17                                array_of_distribs, array_of_dargs, array_of_psize, &
18                                MPI_ORDER_FORTRAN, oldtype, newtype, ierr)
19

```

#### 4.1.5 Address and Size Functions

The displacements in a general datatype are relative to some initial buffer address. **Absolute addresses** can be substituted for these displacements: we treat them as displacements relative to “address zero,” the start of the address space. This initial address zero is indicated by the constant `MPI_BOTTOM`. Thus, a datatype can specify the absolute address of the entries in the communication buffer, in which case the `buf` argument is passed the value `MPI_BOTTOM`.

The address of a location in memory can be found by invoking the function `MPI_GET_ADDRESS`.

```

31 MPI_GET_ADDRESS(location, address)

```

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

```

36 int MPI_Get_address(void *location, MPI_Aint *address)

```

```

38 MPI_GET_ADDRESS(LOCATION, ADDRESS, IERROR)

```

```

39 <type> LOCATION(*)

```

```

40 INTEGER IERROR

```

```

41 INTEGER(KIND=MPI_ADDRESS_KIND) ADDRESS

```

```

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

```

This function replaces `MPI_ADDRESS`, whose use is deprecated. See also Chapter 15. Returns the (byte) address of `location`.

*Advice to users.* Current Fortran MPI codes will run unmodified, and will port to any system. However, they may fail if addresses larger than  $2^{32} - 1$  are used

in the program. New codes should be written so that they use the new functions. This provides compatibility with C/C++ and avoids errors on 64 bit architectures. However, such newly written codes may need to be (slightly) rewritten to port to old Fortran 77 environments that do not support KIND declarations. (*End of advice to users.*)

**Example 4.8** Using MPI\_GET\_ADDRESS for an array.

```
REAL A(100,100)
INTEGER(KIND=MPI_ADDRESS_KIND) I1, I2, DIFF
CALL MPI_GET_ADDRESS(A(1,1), I1, IERROR)
CALL MPI_GET_ADDRESS(A(10,10), I2, IERROR)
DIFF = I2 - I1
! The value of DIFF is 909*sizeofreal; the values of I1 and I2 are
! implementation dependent.
```

*Advice to users.* C users may be tempted to avoid the usage of MPI\_GET\_ADDRESS and rely on the availability of the address operator &. Note, however, that & *cast-expression* is a pointer, not an address. ISO C does not require that the value of a pointer (or the pointer cast to int) be the absolute address of the object pointed at — although this is commonly the case. Furthermore, referencing may not have a unique definition on machines with a segmented address space. The use of MPI\_GET\_ADDRESS to “reference” C variables guarantees portability to such machines as well. (*End of advice to users.*)

*Advice to users.* To prevent problems with the argument copying and register optimization done by Fortran compilers, please note the hints in subsections “Problems Due to Data Copying and Sequence Association,” and “A Problem with Register Optimization” in Section 16.2.2 on pages 473 and 476. (*End of advice to users.*)

The following auxiliary function provides useful information on derived datatypes.

MPI\_TYPE\_SIZE(datatype, size)

IN	datatype	datatype (handle)
OUT	size	datatype size (integer)

```
int MPI_Type_size(MPI_Datatype datatype, int *size)
```

```
MPI_TYPE_SIZE(DATATYPE, SIZE, IERROR)
```

```
INTEGER DATATYPE, SIZE, IERROR
```

```
{int MPI::Datatype::Get_size() const (binding deprecated, see Section ??) }
```

MPI\_TYPE\_SIZE returns the total size, in bytes, of the entries in the type signature associated with datatype; i.e., the total size of the data in a message that would be created with this datatype. Entries that occur multiple times in the datatype are counted with their multiplicity.

ticket150.

ticket150.

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

It is often convenient to define explicitly the lower bound and upper bound of a type map, and override the definition given on page 20. This allows one to define a datatype that has “holes” at its beginning or its end, or a datatype with entries that extend above the upper bound or below the lower bound. Examples of such usage are provided in Section 4.1.14. Also, the user may want to override the alignment rules that are used to compute upper bounds and extents. E.g., a C compiler may allow the user to override default alignment rules for some of the structures within a program. The user has to specify explicitly the bounds of the datatypes that match these structures.

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

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

In general, if

$$\text{Typemap} = \{(\text{type}_0, \text{disp}_0), \dots, (\text{type}_{n-1}, \text{disp}_{n-1})\},$$

then the **lower bound** of  $\text{Typemap}$  is defined to be

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

Similarly, the **upper bound** of  $\text{Typemap}$  is defined to be

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

Then

$$\text{extent}(\text{Typemap}) = \text{ub}(\text{Typemap}) - \text{lb}(\text{Typemap})$$

If  $\text{type}_i$  requires alignment to a byte address that is a multiple of  $k_i$ , then  $\epsilon$  is the least [nonnegative]non-negative increment needed to round  $\text{extent}(\text{Typemap})$  to the next multiple of  $\max_i k_i$ .

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

## 4.1.7 Extent and Bounds of Datatypes

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

`MPI_TYPE_GET_EXTENT(datatype, lb, extent)`

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

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

```
MPI_TYPE_GET_EXTENT(DATATYPE, LB, EXTENT, IERROR)
```

```
    INTEGER DATATYPE, IERROR
```

```
    INTEGER(KIND = MPI_ADDRESS_KIND) LB, EXTENT
```

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

Returns the lower bound and the extent of `datatype` (as defined in Section 4.1.6 on page 20).

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

`MPI_TYPE_CREATE_RESIZED(oldtype, lb, extent, newtype)`

IN	oldtype	input datatype (handle)
IN	lb	new lower bound of datatype (integer)
IN	extent	new extent of datatype (integer)
OUT	newtype	output datatype (handle)

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

```
MPI_TYPE_CREATE_RESIZED(OLDTYPE, LB, EXTENT, NEWTYPE, IERROR)
```

```
    INTEGER OLDTYPE, NEWTYPE, IERROR
```

```
    INTEGER(KIND=MPI_ADDRESS_KIND) LB, EXTENT
```

```

1      {MPI::Datatype MPI::Datatype::Create_resized(const MPI::Aint lb,
ticket150. 2          const MPI::Aint extent) const (binding deprecated, see Section ??) }
3

```

Returns in `newtype` a handle to a new datatype that is identical to `oldtype`, except that the lower bound of this new datatype is set to be `lb`, and its upper bound is set to be `lb + extent`. Any previous `lb` and `ub` markers are erased, and a new pair of lower bound and upper bound markers are put in the positions indicated by the `lb` and `extent` arguments. This affects the behavior of the datatype when used in communication operations, with `count > 1`, and when used in the construction of new derived datatypes.

*Advice to users.* It is strongly recommended that users use these two new functions, rather than the old MPI-1 functions to set and access lower bound, upper bound and extent of datatypes. *(End of advice to users.)*

#### 4.1.8 True Extent of Datatypes

Suppose we implement gather (see also Section 5.5 on page 139) as a spanning tree implemented on top of point-to-point routines. Since the receive buffer is only valid on the root process, one will need to allocate some temporary space for receiving data on intermediate nodes. However, the datatype extent cannot be used as an estimate of the amount of space that needs to be allocated, if the user has modified the extent using the `MPI_UB` and `MPI_LB` values. A function is provided which returns the true extent of the datatype.

```

23
24 MPI_TYPE_GET_TRUE_EXTENT(datatype, true_lb, true_extent)

```

IN	<code>datatype</code>	datatype to get information on (handle)
OUT	<code>true_lb</code>	true lower bound of datatype (integer)
OUT	<code>true_extent</code>	true size of datatype (integer)

```

29
30 int MPI_Type_get_true_extent(MPI_Datatype datatype, MPI_Aint *true_lb,
31     MPI_Aint *true_extent)

```

```

32 MPI_TYPE_GET_TRUE_EXTENT(DATATYPE, TRUE_LB, TRUE_EXTENT, IERROR)
33     INTEGER DATATYPE, IERROR
34     INTEGER(KIND = MPI_ADDRESS_KIND) TRUE_LB, TRUE_EXTENT

```

```

ticket150. 35
36 {void MPI::Datatype::Get_true_extent(MPI::Aint& true_lb,
ticket150. 37     MPI::Aint& true_extent) const (binding deprecated, see Section ??) }
38

```

`true_lb` returns the offset of the lowest unit of store which is addressed by the datatype, i.e., the lower bound of the corresponding typemap, ignoring `MPI_LB` markers. `true_extent` returns the true size of the datatype, i.e., the extent of the corresponding typemap, ignoring `MPI_LB` and `MPI_UB` markers, and performing no rounding for alignment. If the typemap associated with `datatype` is

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

Then

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



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

and

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

(Readers should compare this with the definitions in Section 4.1.6 on page 20 and Section 4.1.7 on page 21, which describe the function `MPI_TYPE_GET_EXTENT`.)

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

#### 4.1.9 Commit and Free

A datatype object has to be **committed** before it can be used in a communication. As an argument in datatype constructors, uncommitted and also committed datatypes can be used. There is no need to commit basic datatypes. They are “pre-committed.”

`MPI_TYPE_COMMIT(datatype)`

INOUT     datatype                             datatype that is committed (handle)

`int MPI_Type_commit(MPI_Datatype *datatype)`

`MPI_TYPE_COMMIT(DATATYPE, IERROR)`

INTEGER DATATYPE, IERROR

`{void MPI::Datatype::Commit() (binding deprecated, see Section ??) }`

The commit operation commits the datatype, that is, the formal description of a communication buffer, not the content of that buffer. Thus, after a datatype has been committed, it can be repeatedly reused to communicate the changing content of a buffer or, indeed, the content of different buffers, with different starting addresses.

*Advice to implementors.* The system may “compile” at commit time an internal representation for the datatype that facilitates communication, e.g. change from a compacted representation to a flat representation of the datatype, and select the most convenient transfer mechanism. (*End of advice to implementors.*)

`MPI_TYPE_COMMIT` will accept a committed datatype; in this case, it is equivalent to a no-op.

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

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

```

1  CALL MPI_TYPE_VECTOR(3, 5, 4, MPI_REAL, type1, ierr)
2      ! new uncommitted type object created
3  CALL MPI_TYPE_COMMIT(type1, ierr)
4      ! now type1 can be used anew for communication

```

```

7  MPI_TYPE_FREE(datatype)

```

```

9      INOUT    datatype                datatype that is freed (handle)

```

```

11 int MPI_Type_free(MPI_Datatype *datatype)

```

```

12 MPI_TYPE_FREE(DATATYPE, IERROR)

```

```

13     INTEGER DATATYPE, IERROR

```

```

15 {void MPI::Datatype::Free() (binding deprecated, see Section ??) }

```

Marks the datatype object associated with `datatype` for deallocation and sets `datatype` to `MPI_DATATYPE_NULL`. Any communication that is currently using this datatype will complete normally. Freeing a datatype does not affect any other datatype that was built from the freed datatype. The system behaves as if input datatype arguments to derived datatype constructors are passed by value.

*Advice to implementors.* The implementation may keep a reference count of active communications that use the datatype, in order to decide when to free it. Also, one may implement constructors of derived datatypes so that they keep pointers to their datatype arguments, rather than copying them. In this case, one needs to keep track of active datatype definition references in order to know when a datatype object can be freed. (*End of advice to implementors.*)

#### 4.1.10 Duplicating a Datatype

```

33 MPI_TYPE_DUP(type, newtype)

```

```

34     IN        type                datatype (handle)

```

```

35     OUT       newtype            copy of type (handle)

```

```

37 int MPI_Type_dup(MPI_Datatype type, MPI_Datatype *newtype)

```

```

39 MPI_TYPE_DUP(TYPE, NEWTYPE, IERROR)

```

```

40     INTEGER TYPE, NEWTYPE, IERROR

```

```

42 {MPI::Datatype MPI::Datatype::Dup() const (binding deprecated, see Section ??) }

```

`MPI_TYPE_DUP` is a type constructor which duplicates the existing `type` with associated key values. For each key value, the respective copy callback function determines the attribute value associated with this key in the new communicator; one particular action that a copy callback may take is to delete the attribute from the new datatype. Returns in `newtype` a new datatype with exactly the same properties as `type`

and any copied cached information, see Section 6.7.4 on page 239. The new datatype has identical upper bound and lower bound and yields the same net result when fully decoded with the functions in Section 4.1.13. The `newtype` has the same committed state as the old type.

#### 4.1.11 Use of General Datatypes in Communication

Handles to derived datatypes can be passed to a communication call wherever a datatype argument is required. A call of the form `MPI_SEND(buf, count, datatype, ...)`, where `count > 1`, is interpreted as if the call was passed a new datatype which is the concatenation of `count` copies of `datatype`. Thus, `MPI_SEND(buf, count, datatype, dest, tag, comm)` is equivalent to,

```
MPI_TYPE_CONTIGUOUS(count, datatype, newtype)
MPI_TYPE_COMMIT(newtype)
MPI_SEND(buf, 1, newtype, dest, tag, comm).
```

Similar statements apply to all other communication functions that have a `count` and `datatype` argument.

Suppose that a send operation `MPI_SEND(buf, count, datatype, dest, tag, comm)` is executed, where `datatype` has type map,

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

and extent *extent*. (Empty entries of “pseudo-type” `MPI_UB` and `MPI_LB` are not listed in the type map, but they affect the value of *extent*.) The send operation sends  $n \cdot \text{count}$  entries, where entry  $i \cdot n + j$  is at location  $addr_{i,j} = \text{buf} + \text{extent} \cdot i + disp_j$  and has type  $type_j$ , for  $i = 0, \dots, \text{count} - 1$  and  $j = 0, \dots, n - 1$ . These entries need not be contiguous, nor distinct; their order can be arbitrary.

The variable stored at address  $addr_{i,j}$  in the calling program should be of a type that matches  $type_j$ , where type matching is defined as in Section 3.3.1. The message sent contains  $n \cdot \text{count}$  entries, where entry  $i \cdot n + j$  has type  $type_j$ .

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

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

with extent *extent*. (Again, empty entries of “pseudo-type” `MPI_UB` and `MPI_LB` are not listed in the type map, but they affect the value of *extent*.) This receive operation receives  $n \cdot \text{count}$  entries, where entry  $i \cdot n + j$  is at location  $\text{buf} + \text{extent} \cdot i + disp_j$  and has type  $type_j$ . If the incoming message consists of  $k$  elements, then we must have  $k \leq n \cdot \text{count}$ ; the  $i \cdot n + j$ -th element of the message should have a type that matches  $type_j$ .

Type matching is defined according to the type signature of the corresponding datatypes, that is, the sequence of basic type components. Type matching does not depend on some aspects of the datatype definition, such as the displacements (layout in memory) or the intermediate types used.

**Example 4.11** This example shows that type matching is defined in terms of the basic types that a derived type consists of.

```

1  ...
2  CALL MPI_TYPE_CONTIGUOUS( 2, MPI_REAL, type2, ...)
3  CALL MPI_TYPE_CONTIGUOUS( 4, MPI_REAL, type4, ...)
4  CALL MPI_TYPE_CONTIGUOUS( 2, type2, type22, ...)
5  ...
6  CALL MPI_SEND( a, 4, MPI_REAL, ...)
7  CALL MPI_SEND( a, 2, type2, ...)
8  CALL MPI_SEND( a, 1, type22, ...)
9  CALL MPI_SEND( a, 1, type4, ...)
10 ...
11 CALL MPI_RECV( a, 4, MPI_REAL, ...)
12 CALL MPI_RECV( a, 2, type2, ...)
13 CALL MPI_RECV( a, 1, type22, ...)
14 CALL MPI_RECV( a, 1, type4, ...)

```

Each of the sends matches any of the receives.

A datatype may specify overlapping entries. The use of such a datatype in a receive operation is erroneous. (This is erroneous even if the actual message received is short enough not to write any entry more than once.)

Suppose that `MPI_RECV(buf, count, datatype, dest, tag, comm, status)` is executed, where `datatype` has type map,

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

The received message need not fill all the receive buffer, nor does it need to fill a number of locations which is a multiple of  $n$ . Any number,  $k$ , of basic elements can be received, where  $0 \leq k \leq \text{count} \cdot n$ . The number of basic elements received can be retrieved from `status` using the query function `MPI_GET_ELEMENTS`.

`MPI_GET_ELEMENTS( status, datatype, count)`

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

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

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

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

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

The previously defined function, `MPI_GET_COUNT` (Section 3.2.5), has a different behavior. It returns the number of “top-level entries” received, i.e. the number of “copies” of type `datatype`. In the previous example, `MPI_GET_COUNT` may return any integer value  $k$ , where  $0 \leq k \leq \text{count}$ . If `MPI_GET_COUNT` returns  $k$ , then the number of basic elements received (and the value returned by `MPI_GET_ELEMENTS`) is  $n \cdot k$ . If the number of basic elements received is not a multiple of  $n$ , that is, if the receive operation has not

ticket150.

ticket150.

received an integral number of `datatype` “copies,” then `MPI_GET_COUNT` returns the value `MPI_UNDEFINED`. The `datatype` argument should match the argument provided by the receive call that set the `status` variable.

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

```

...
CALL MPI_TYPE_CONTIGUOUS(2, MPI_REAL, Type2, ierr)
CALL MPI_TYPE_COMMIT(Type2, ierr)
...
CALL MPI_COMM_RANK(comm, rank, ierr)
IF (rank.EQ.0) THEN
    CALL MPI_SEND(a, 2, MPI_REAL, 1, 0, comm, ierr)
    CALL MPI_SEND(a, 3, MPI_REAL, 1, 0, comm, ierr)
ELSE IF (rank.EQ.1) THEN
    CALL MPI_RECV(a, 2, Type2, 0, 0, comm, stat, ierr)
    CALL MPI_GET_COUNT(stat, Type2, i, ierr)      ! returns i=1
    CALL MPI_GET_ELEMENTS(stat, Type2, i, ierr)   ! returns i=2
    CALL MPI_RECV(a, 2, Type2, 0, 0, comm, stat, ierr)
    CALL MPI_GET_COUNT(stat, Type2, i, ierr)      ! returns i=MPI_UNDEFINED
    CALL MPI_GET_ELEMENTS(stat, Type2, i, ierr)   ! returns i=3
END IF

```

The function `MPI_GET_ELEMENTS` can also be used after a probe to find the number of elements in the probed message. Note that the two functions `MPI_GET_COUNT` and `MPI_GET_ELEMENTS` return the same values when they are used with basic datatypes.

*Rationale.* The extension given to the definition of `MPI_GET_COUNT` seems natural: one would expect this function to return the value of the `count` argument, when the receive buffer is filled. Sometimes `datatype` represents a basic unit of data one wants to transfer, for example, a record in an array of records (structures). One should be able to find out how many components were received without bothering to divide by the number of elements in each component. However, on other occasions, `datatype` is used to define a complex layout of data in the receiver memory, and does not represent a basic unit of data for transfers. In such cases, one needs to use the function `MPI_GET_ELEMENTS`. (*End of rationale.*)

*Advice to implementors.* The definition implies that a receive cannot change the value of storage outside the entries defined to compose the communication buffer. In particular, the definition implies that padding space in a structure should not be modified when such a structure is copied from one process to another. This would prevent the obvious optimization of copying the structure, together with the padding, as one contiguous block. The implementation is free to do this optimization when it does not impact the outcome of the computation. The user can “force” this optimization by explicitly including padding as part of the message. (*End of advice to implementors.*)

#### 4.1.12 Correct Use of Addresses

Successively declared variables in C or Fortran are not necessarily stored at contiguous locations. Thus, care must be exercised that displacements do not cross from one variable

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

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

1. The function `MPI_GET_ADDRESS` returns a valid address, when passed as argument a variable of the calling program.
2. The `buf` argument of a communication function evaluates to a valid address, when passed as argument a variable of the calling program.
3. If  $v$  is a valid address, and  $i$  is an integer, then  $v+i$  is a valid address, provided  $v$  and  $v+i$  are in the same sequential storage.
4. If  $v$  is a valid address then `MPI_BOTTOM + v` is a valid address.

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

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

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

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

#### 4.1.13 Decoding a Datatype

MPI datatype objects allow users to specify an arbitrary layout of data in memory. There are several cases where accessing the layout information in opaque datatype objects would be useful. The opaque datatype object has found a number of uses outside MPI. Furthermore, a number of tools wish to display internal information about a datatype. To achieve this, datatype decoding functions are provided. The two functions in this section are used

together to decode datatypes to recreate the calling sequence used in their initial definition. These can be used to allow a user to determine the type map and type signature of a datatype.

`MPI_TYPE_GET_ENVELOPE(datatype, num_integers, num_addresses, num_datatypes, combiner)`

IN	datatype	datatype to access (handle)	
OUT	num_integers	number of input integers used in the call constructing combiner ([nonnegative]non-negative integer)	ticket74.
OUT	num_addresses	number of input addresses used in the call constructing combiner ([nonnegative]non-negative integer)	ticket74.
OUT	num_datatypes	number of input datatypes used in the call constructing combiner ([nonnegative]non-negative integer)	ticket74.
OUT	combiner	combiner (state)	

```
int MPI_Type_get_envelope(MPI_Datatype datatype, int *num_integers,
                          int *num_addresses, int *num_datatypes, int *combiner)
```

```
MPI_TYPE_GET_ENVELOPE(DATATYPE, NUM_INTEGERS, NUM_ADDRESSES, NUM_DATATYPES,
                       COMBINER, IERROR)
```

```
INTEGER DATATYPE, NUM_INTEGERS, NUM_ADDRESSES, NUM_DATATYPES, COMBINER,
IERROR
```

```
{void MPI::Datatype::Get_envelope(int& num_integers, int& num_addresses,
                                  int& num_datatypes, int& combiner) const (binding deprecated, see
                                  Section ??) }
```

For the given datatype, `MPI_TYPE_GET_ENVELOPE` returns information on the number and type of input arguments used in the call that created the datatype. The number-of-arguments values returned can be used to provide sufficiently large arrays in the decoding routine `MPI_TYPE_GET_CONTENTS`. This call and the meaning of the returned values is described below. The combiner reflects the MPI datatype constructor call that was used in creating datatype.

*Rationale.* By requiring that the combiner reflect the constructor used in the creation of the datatype, the decoded information can be used to effectively recreate the calling sequence used in the original creation. One call is effectively the same as another when the information obtained from `MPI_TYPE_GET_CONTENTS` may be used with either to produce the same outcome. C calls `MPI_Type_hindexed` and `MPI_Type_create_hindexed` are always effectively the same while the Fortran call `MPI_TYPE_HINDEXED` will be different than either of these in some MPI implementations. This is the most useful information and was felt to be reasonable even though it constrains implementations to remember the original constructor sequence even if the internal representation is different.

The decoded information keeps track of datatype duplications. This is important as one needs to distinguish between a predefined datatype and a dup of a predefined

datatype. The former is a constant object that cannot be freed, while the latter is a derived datatype that can be freed. (*End of rationale.*)

The list below has the values that can be returned in `combiner` on the left and the call associated with them on the right.

<code>MPI_COMBINER_NAMED</code>	a named predefined datatype
<code>MPI_COMBINER_DUP</code>	<code>MPI_TYPE_DUP</code>
<code>MPI_COMBINER_CONTIGUOUS</code>	<code>MPI_TYPE_CONTIGUOUS</code>
<code>MPI_COMBINER_VECTOR</code>	<code>MPI_TYPE_VECTOR</code>
<code>MPI_COMBINER_HVECTOR_INTEGER</code>	<code>MPI_TYPE_HVECTOR</code> from Fortran
<code>MPI_COMBINER_HVECTOR</code>	<code>MPI_TYPE_HVECTOR</code> from C or C++ and in some case Fortran or <code>MPI_TYPE_CREATE_HVECTOR</code>
<code>MPI_COMBINER_INDEXED</code>	<code>MPI_TYPE_INDEXED</code>
<code>MPI_COMBINER_HINDEXED_INTEGER</code>	<code>MPI_TYPE_HINDEXED</code> from Fortran
<code>MPI_COMBINER_HINDEXED</code>	<code>MPI_TYPE_HINDEXED</code> from C or C++ and in some case Fortran or <code>MPI_TYPE_CREATE_HINDEXED</code>
<code>MPI_COMBINER_INDEXED_BLOCK</code>	<code>MPI_TYPE_CREATE_INDEXED_BLOCK</code>
<code>MPI_COMBINER_STRUCT_INTEGER</code>	<code>MPI_TYPE_STRUCT</code> from Fortran
<code>MPI_COMBINER_STRUCT</code>	<code>MPI_TYPE_STRUCT</code> from C or C++ and in some case Fortran or <code>MPI_TYPE_CREATE_STRUCT</code>
<code>MPI_COMBINER_SUBARRAY</code>	<code>MPI_TYPE_CREATE_SUBARRAY</code>
<code>MPI_COMBINER_DARRAY</code>	<code>MPI_TYPE_CREATE_DARRAY</code>
<code>MPI_COMBINER_F90_REAL</code>	<code>MPI_TYPE_CREATE_F90_REAL</code>
<code>MPI_COMBINER_F90_COMPLEX</code>	<code>MPI_TYPE_CREATE_F90_COMPLEX</code>
<code>MPI_COMBINER_F90_INTEGER</code>	<code>MPI_TYPE_CREATE_F90_INTEGER</code>
<code>MPI_COMBINER_RESIZED</code>	<code>MPI_TYPE_CREATE_RESIZED</code>

Table 4.1: `combiner` values returned from `MPI_TYPE_GET_ENVELOPE`

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

For deprecated calls with address arguments, we sometimes need to differentiate whether the call used an integer or an address size argument. For example, there are two combin-ers for hvector: `MPI_COMBINER_HVECTOR_INTEGER` and `MPI_COMBINER_HVECTOR`. The former is used if it was the MPI-1 call from Fortran, and the latter is used if it was the MPI-1 call from C or C++. However, on systems where `MPI_ADDRESS_KIND = MPI_INTEGER_KIND` (i.e., where integer arguments and address size arguments are the same), the combiner `MPI_COMBINER_HVECTOR` may be returned for a datatype constructed by a call to `MPI_TYPE_HVECTOR` from Fortran. Similarly, `MPI_COMBINER_HINDEXED` may be returned for a datatype constructed by a call to `MPI_TYPE_HINDEXED` from Fortran, and `MPI_COMBINER_STRUCT` may be returned for a datatype constructed by a call to `MPI_TYPE_STRUCT` from Fortran. On such systems, one need not differentiate construc-tors that take address size arguments from constructors that take integer arguments, since these are the same. The preferred calls all use address sized arguments so two combin-ers



are not required for them.

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

The actual arguments used in the creation call for a `datatype` can be obtained from the call:

`MPI_TYPE_GET_CONTENTS(datatype, max_integers, max_addresses, max_datatypes, array_of_integers, array_of_addresses, array_of_datatypes)`

IN	<code>datatype</code>	datatype to access (handle)	
IN	<code>max_integers</code>	number of elements in <code>array_of_integers</code> ([nonnegative]non-negative integer)	ticket74.
IN	<code>max_addresses</code>	number of elements in <code>array_of_addresses</code> ([nonnegative]non-negative integer)	ticket74.
IN	<code>max_datatypes</code>	number of elements in <code>array_of_datatypes</code> ([nonnegative]non-negative integer)	ticket74.
OUT	<code>array_of_integers</code>	contains integer arguments used in constructing datatype (array of integers)	
OUT	<code>array_of_addresses</code>	contains address arguments used in constructing datatype (array of integers)	
OUT	<code>array_of_datatypes</code>	contains datatype arguments used in constructing datatype (array of handles)	

```
int MPI_Type_get_contents(MPI_Datatype datatype, int max_integers,
    int max_addresses, int max_datatypes, int array_of_integers[],
    MPI_Aint array_of_addresses[],
    MPI_Datatype array_of_datatypes[])
```

```
MPI_TYPE_GET_CONTENTS(DATATYPE, MAX_INTEGERS, MAX_ADDRESSES, MAX_DATATYPES,
    ARRAY_OF_INTEGERS, ARRAY_OF_ADDRESSES, ARRAY_OF_DATATYPES,
    IERROR)
```

```
INTEGER DATATYPE, MAX_INTEGERS, MAX_ADDRESSES, MAX_DATATYPES,
    ARRAY_OF_INTEGERS(*), ARRAY_OF_DATATYPES(*), IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) ARRAY_OF_ADDRESSES(*)
```

```
{void MPI::Datatype::Get_contents(int max_integers, int max_addresses,
    int max_datatypes, int array_of_integers[],
    MPI::Aint array_of_addresses[],
    MPI::Datatype array_of_datatypes[]) const (binding deprecated, see
    Section ??) }
```

`datatype` must be a predefined unnamed or a derived datatype; the call is erroneous if `datatype` is a predefined named datatype.

The values given for `max_integers`, `max_addresses`, and `max_datatypes` must be at least as large as the value returned in `num_integers`, `num_addresses`, and `num_datatypes`, respectively, in the call `MPI_TYPE_GET_ENVELOPE` for the same `datatype` argument.

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

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

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

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

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

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

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

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

*Rationale.* By having all address arguments returned in the `array_of_addresses` argument, the result from a C and Fortran decoding of a `datatype` gives the result in the same argument. It is assumed that an integer of type

INTEGER(KIND=MPI\_ADDRESS\_KIND) will be at least as large as the INTEGER argument used in datatype construction with the old MPI-1 calls so no loss of information will occur. (*End of rationale.*)

The following defines what values are placed in each entry of the returned arrays depending on the datatype constructor used for datatype. It also specifies the size of the arrays needed which is the values returned by MPI\_TYPE\_GET\_ENVELOPE. In Fortran, the following calls were made:

```

PARAMETER (LARGE = 1000)
INTEGER TYPE, NI, NA, ND, COMBINER, I(LARGE), D(LARGE), IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) A(LARGE)
! CONSTRUCT DATATYPE TYPE (NOT SHOWN)
CALL MPI_TYPE_GET_ENVELOPE(TYPE, NI, NA, ND, COMBINER, IERROR)
IF ((NI .GT. LARGE) .OR. (NA .GT. LARGE) .OR. (ND .GT. LARGE)) THEN
    WRITE (*, *) "NI, NA, OR ND = ", NI, NA, ND, &
    " RETURNED BY MPI_TYPE_GET_ENVELOPE IS LARGER THAN LARGE = ", LARGE
    CALL MPI_ABORT(MPI_COMM_WORLD, 99[ticket116.], IERROR)
ENDIF
CALL MPI_TYPE_GET_CONTENTS(TYPE, NI, NA, ND, I, A, D, IERROR)

```

or in C the analogous calls of:

```

#define LARGE 1000
int ni, na, nd, combiner, i[LARGE];
MPI_Aint a[LARGE];
MPI_Datatype type, d[LARGE];
/* construct datatype type (not shown) */
MPI_Type_get_envelope(type, &ni, &na, &nd, &combiner);
if ((ni > LARGE) || (na > LARGE) || (nd > LARGE)) {
    fprintf(stderr, "ni, na, or nd = %d %d %d returned by ", ni, na, nd);
    fprintf(stderr, "MPI_Type_get_envelope is larger than LARGE = %d\n",
        LARGE);
    MPI_Abort(MPI_COMM_WORLD, 99);
};
MPI_Type_get_contents(type, ni, na, nd, i, a, d);

```

The C++ code is in analogy to the C code above with the same values returned.

In the descriptions that follow, the lower case name of arguments is used.

If combiner is MPI\_COMBINER\_NAMED then it is erroneous to call MPI\_TYPE\_GET\_CONTENTS.

If combiner is MPI\_COMBINER\_DUP then

Constructor argument	C & C++ location	Fortran location
oldtype	d[0]	D(1)

and ni = 0, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_CONTIGUOUS then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
oldtype	d[0]	D(1)

and ni = 1, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_VECTOR then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
blocklength	i[1]	I(2)
stride	i[2]	I(3)
oldtype	d[0]	D(1)

and ni = 3, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_HVECTOR\_INTEGER or MPI\_COMBINER\_HVECTOR then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
blocklength	i[1]	I(2)
stride	a[0]	A(1)
oldtype	d[0]	D(1)

and ni = 2, na = 1, nd = 1.

If combiner is MPI\_COMBINER\_INDEXED then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
array_of_blocklengths	i[1] to i[i[0]]	I(2) to I(I(1)+1)
array_of_displacements	i[i[0]+1] to i[2*i[0]]	I(I(1)+2) to I(2*I(1)+1)
oldtype	d[0]	D(1)

and ni = 2\*count+1, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_HINDEXED\_INTEGER or MPI\_COMBINER\_HINDEXED then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
array_of_blocklengths	i[1] to i[i[0]]	I(2) to I(I(1)+1)
array_of_displacements	a[0] to a[i[0]-1]	A(1) to A(I(1))
oldtype	d[0]	D(1)

and ni = count+1, na = count, nd = 1.

If combiner is MPI\_COMBINER\_INDEXED\_BLOCK then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
blocklength	i[1]	I(2)
array_of_displacements	i[2] to i[i[0]+1]	I(3) to I(I(1)+2)
oldtype	d[0]	D(1)

and ni = count+2, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_STRUCT\_INTEGER or MPI\_COMBINER\_STRUCT then

Constructor argument	C & C++ location	Fortran location
count	i[0]	I(1)
array_of_blocklengths	i[1] to i[i[0]]	I(2) to I(I(1)+1)
array_of_displacements	a[0] to a[i[0]-1]	A(1) to A(I(1))
array_of_types	d[0] to d[i[0]-1]	D(1) to D(I(1))

and ni = count+1, na = count, nd = count.

If combiner is MPI\_COMBINER\_SUBARRAY then

Constructor argument	C & C++ location	Fortran location
ndims	i[0]	I(1)
array_of_sizes	i[1] to i[i[0]]	I(2) to I(I(1)+1)
array_of_subsizes	i[i[0]+1] to i[2*i[0]]	I(I(1)+2) to I(2*I(1)+1)
array_of_starts	i[2*i[0]+1] to i[3*i[0]]	I(2*I(1)+2) to I(3*I(1)+1)
order	i[3*i[0]+1]	I(3*I(1)+2)
oldtype	d[0]	D(1)

and ni = 3\*ndims+2, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_DARRAY then

Constructor argument	C & C++ location	Fortran location
size	i[0]	I(1)
rank	i[1]	I(2)
ndims	i[2]	I(3)
array_of_gsizes	i[3] to i[i[2]+2]	I(4) to I(I(3)+3)
array_of_distribs	i[i[2]+3] to i[2*i[2]+2]	I(I(3)+4) to I(2*I(3)+3)
array_of_dargs	i[2*i[2]+3] to i[3*i[2]+2]	I(2*I(3)+4) to I(3*I(3)+3)
array_of_psize	i[3*i[2]+3] to i[4*i[2]+2]	I(3*I(3)+4) to I(4*I(3)+3)
order	i[4*i[2]+3]	I(4*I(3)+4)
oldtype	d[0]	D(1)

and ni = 4\*ndims+4, na = 0, nd = 1.

If combiner is MPI\_COMBINER\_F90\_REAL then

Constructor argument	C & C++ location	Fortran location
p	i[0]	I(1)
r	i[1]	I(2)

and ni = 2, na = 0, nd = 0.

If combiner is MPI\_COMBINER\_F90\_COMPLEX then

Constructor argument	C & C++ location	Fortran location
p	i[0]	I(1)
r	i[1]	I(2)

and ni = 2, na = 0, nd = 0.

If combiner is MPI\_COMBINER\_F90\_INTEGER then

Constructor argument	C & C++ location	Fortran location
r	i[0]	I(1)

and ni = 1, na = 0, nd = 0.

If combiner is MPI\_COMBINER\_RESIZED then

Constructor argument	C & C++ location	Fortran location
lb	a[0]	A(1)
extent	a[1]	A(2)
oldtype	d[0]	D(1)

and ni = 0, na = 2, nd = 1.

#### 4.1.14 Examples

The following examples illustrate the use of derived datatypes.

**Example 4.13** Send and receive a section of a 3D array.

```

1  REAL a(100,100,100), e(9,9,9)
2
3  INTEGER oneslice, twoslice, threeslice, sizeofreal, myrank, ierr
4  INTEGER status(MPI_STATUS_SIZE)
5
6
7  C    extract the section a(1:17:2, 3:11, 2:10)
8  C    and store it in e(:, :, :).
9
10 CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)
11
12 CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)
13
14 C    create datatype for a 1D section
15 CALL MPI_TYPE_VECTOR( 9, 1, 2, MPI_REAL, oneslice, ierr)
16
17 C    create datatype for a 2D section
18 CALL MPI_TYPE_HVECTOR(9, 1, 100*sizeofreal, oneslice, twoslice, ierr)
19
20 C    create datatype for the entire section
21 CALL MPI_TYPE_HVECTOR( 9, 1, 100*100*sizeofreal, twoslice,
22                        threeslice, ierr)
23
24 CALL MPI_TYPE_COMMIT( threeslice, ierr)
25 CALL MPI_SENDRECV(a(1,3,2), 1, threeslice, myrank, 0, e, 9*9*9,
26                  MPI_REAL, myrank, 0, MPI_COMM_WORLD, status, ierr)
27

```

**Example 4.14** Copy the (strictly) lower triangular part of a matrix.

```

29  REAL a(100,100), b(100,100)
30  INTEGER disp(100), blocklen(100), ltype, myrank, ierr
31  INTEGER status(MPI_STATUS_SIZE)
32
33 C    copy lower triangular part of array a
34 C    onto lower triangular part of array b
35
36 CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)
37
38 C    compute start and size of each column
39 DO i=1, 100
40     disp(i) = 100*(i-1) + i
41     block[ticket116.]len(i) = 100-i
42 END DO
43
44 C    create datatype for lower triangular part
45 CALL MPI_TYPE_INDEXED( 100, block[ticket116.]len, disp, MPI_REAL, ltype, ierr)
46
47 CALL MPI_TYPE_COMMIT(ltype, ierr)
48

```

```

CALL MPI_SENDRECV( a, 1, ltype, myrank, 0, b, 1,
                  ltype, myrank, 0, MPI_COMM_WORLD, status, ierr)

```

**Example 4.15** Transpose a matrix.

```

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

C      transpose matrix a onto b

CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)

CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)

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

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

CALL MPI_TYPE_COMMIT( xpose, ierr)

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

```

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

```

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

CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierr)

C      transpose matrix a onto b

CALL MPI_TYPE_EXTENT( MPI_REAL, sizeofreal, ierr)

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

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

```

```

1      blocklen(2) = 1
2      CALL MPI_TYPE_STRUCT( 2, blocklen, disp, type, row1, ierr)
3
4      CALL MPI_TYPE_COMMIT( row1, ierr)
5
6  C      send 100 rows and receive in column major order
7      CALL MPI_SENDRECV( a, 100, row1, myrank, 0, b, 100*100,
8                          MPI_REAL, myrank, 0, MPI_COMM_WORLD, status, ierr)
9

```

**Example 4.17** We manipulate an array of structures.

```

11  struct Partstruct
12  {
13      int      kind; /* particle class */
14      double d[6]; /* particle coordinates */
15      char    b[7]; /* some additional information */
16  };
17
18  struct Partstruct    particle[1000];
19
20  int      i, dest, rank[ticket116.], tag;
21  MPI_Comm comm;
22
23
24
25  /* build datatype describing structure */
26
27  MPI_Datatype Particletype;
28  MPI_Datatype type[3] = {MPI_INT, MPI_DOUBLE, MPI_CHAR};
29  int      blocklen[3] = {1, 6, 7};
30  MPI_Aint disp[3];
31  MPI_Aint base;
32
33
34  /* compute displacements of structure components */
35
36  MPI_Address( particle, disp);
37  MPI_Address( particle[0].d, disp+1);
38  MPI_Address( particle[0].b, disp+2);
39  base = disp[0];
40  for (i=0; i < 3; i++) disp[i] -= base;
41
42  MPI_Type_struct( 3, blocklen, disp, type, &Particletype);
43
44  /* If compiler does padding in mysterious ways,
45     the following may be safer */
46
47  MPI_Datatype type1[4] = {MPI_INT, MPI_DOUBLE, MPI_CHAR, MPI_UB};
48  int      blocklen1[4] = {1, 6, 7, 1};

```



```

MPI_Aint      disp1[4];
1
2
/* compute displacements of structure components */
3
4
MPI_Address( particle, disp1);
5
MPI_Address( particle[0].d, disp1+1);
6
MPI_Address( particle[0].b, disp1+2);
7
MPI_Address( particle+1, disp1+3);
8
base = disp1[0];
9
for (i=0; i < 4; i++) disp1[i] -= base;
10
11
/* build datatype describing structure */
12
13
MPI_Type_struct( 4, blocklen1, disp1, type1, &Particletype);
14
15
16
/* 4.1:
17
send the entire array */
18
19
MPI_Type_commit( &Particletype);
20
MPI_Send( particle, 1000, Particletype, dest, tag, comm);
21
22
23
/* 4.2:
24
send only the entries of class zero particles,
25
preceded by the number of such entries */
26
27
MPI_Datatype Zparticles; /* datatype describing all particles
28
                           with class zero (needs to be recomputed
29
                           if classes change) */
30
MPI_Datatype Ztype;
31
32
MPI_Aint      zdisp[1000];
33
int           zblock[1000], j, k;
34
[ticket116.]int      zzblock[2] = {1,1};
35
MPI_Aint      zzdisp[2];
36
MPI_Datatype  zztype[2];
37
38
/* compute displacements of class zero particles */
39
j = 0;
40
for(i=0; i < 1000; i++)
41
    if (particle[i].kind == 0)
42
    {
43
        zdisp[j] = i;
44
        zblock[j] = 1;
45
        j++;
46
    }
47
48

```

```

1  /* create datatype for class zero particles */
2  MPI_Type_indexed( j, zblock, zdisp, Particletype, &Zparticles);
3
4  /* prepend particle count */
5  MPI_Address(&j, zzdisp);
6  MPI_Address(particle, zzdisp+1);
7  zztype[0] = MPI_INT;
8  zztype[1] = Zparticles;
9  MPI_Type_struct(2, zzbblock, zzdisp, zztype, &Ztype);
10
11 MPI_Type_commit( &Ztype);
12 MPI_Send( MPI_BOTTOM, 1, Ztype, dest, tag, comm);
13
14
15      /* A probably more efficient way of defining Zparticles */
16
17 /* consecutive particles with [ticket116.][index]kind zero are handled as one block */
18 j=0;
19 for (i=0; i < 1000; i++)
20     if (particle[i].[ticket116.][index]kind == 0)
21     {
22         for (k=i+1; (k < 1000)&&(particle[k].[ticket116.][index]kind == 0) ; k++);
23         zdisp[j] = i;
24         zblock[j] = k-i;
25         j++;
26         i = k;
27     }
28 MPI_Type_indexed( j, zblock, zdisp, Particletype, &Zparticles);
29
30
31      /* 4.3:
32      send the first two coordinates of all entries */
33
34 MPI_Datatype Allpairs;      /* datatype for all pairs of coordinates */
35
36 MPI_Aint sizeofentry;
37
38 MPI_Type_extent( Particletype, &sizeofentry);
39
40      /* sizeofentry can also be computed by subtracting the address
41      of particle[0] from the address of particle[1] */
42
43 MPI_Type_hvector( 1000, 2, sizeofentry, MPI_DOUBLE, &Allpairs);
44 MPI_Type_commit( &Allpairs);
45 MPI_Send( particle[0].d, 1, Allpairs, dest, tag, comm);
46
47      /* an alternative solution to 4.3 */
48

```

```

MPI_Datatype Onepair;    /* datatype for one pair of coordinates, with
                           the extent of one particle entry */
MPI_Aint disp2[3];
MPI_Datatype type2[3] = {MPI_LB, MPI_DOUBLE, MPI_UB};
int blocklen2[3] = {1, 2, 1};

MPI_Address( particle, disp2);
MPI_Address( particle[0].d, disp2+1);
MPI_Address( particle+1, disp2+2);
base = disp2[0];
for (i=0; i<2; i++) disp2[i] -= base;

MPI_Type_struct( 3, blocklen2, disp2, type2, &Onepair);
MPI_Type_commit( &Onepair);
MPI_Send( particle[0].d, 1000, Onepair, dest, tag, comm);

```

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

```

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

struct Partstruct particle[1000];

/* build datatype describing first array entry */

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

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

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

/* 5.1:
   send the entire array */

MPI_Type_commit( &Particletype);

```

```

1  MPI_Send( MPI_BOTTOM, 1000, Particletype, dest, tag[ticket116.], MPI_COMM_WORLD);
2
3
4      /* 5.2:
5      send the entries of class zero,
6      preceded by the number of such entries */
7
8  MPI_Datatype Zparticles, Ztype;
9
10 MPI_Aint      zdisp[1000];
11 int          zblock[1000], i, j, k;
12 int          zzblock[2] = {1,1};
13 MPI_Datatype zztype[2];
14 MPI_Aint      zzdisp[2];
15
16 j=0;
17 for (i=0; i < 1000; i++)
18     if (particle[i].[ticket116.][index]kind == 0)
19     {
20         for (k=i+1; (k < 1000)&&(particle[k].[ticket116.][index]kind [ticket16.]== 0) ; k++)
21             zdisp[j] = i;
22             zblock[j] = k-i;
23             j++;
24             i = k;
25     }
26 MPI_Type_indexed( j, zblock, zdisp, Particletype, &Zparticles);
27 /* Zparticles describe particles with class zero, using
28    their absolute addresses*/
29
30 /* prepend particle count */
31 MPI_Address(&j, zzdisp);
32 zzdisp[1] = MPI_BOTTOM;
33 zztype[0] = MPI_INT;
34 zztype[1] = Zparticles;
35 MPI_Type_struct(2, zzblock, zzdisp, zztype, &Ztype);
36
37 MPI_Type_commit( &Ztype);
38 MPI_Send( MPI_BOTTOM, 1, Ztype, dest, tag, [ticket116.][comm]MPI_COMM_WORLD);
39
40

```

#### Example 4.19 Handling of unions.

```

42 union {
43     int      ival;
44     float    fval;
45 } u[1000];
46
47
48 int      utype;

```

```

1
2  /* All entries of u have identical type; variable
3     utype keeps track of their current type */
4
5  MPI_Datatype  type[2];
6  int          blocklen[2] = {1,1}[ticket116.], dest, tag;
7  MPI_Aint     disp[2];
8  MPI_Datatype  mpi_utype[2];
9  MPI_Aint     i,j;
10
11 /* compute an MPI datatype for each possible union type;
12    assume values are left-aligned in union storage. */
13
14 MPI_Address( u, &i);
15 MPI_Address( u+1, &j);
16 disp[0] = 0; disp[1] = j-i;
17 type[1] = MPI_UB;
18
19 type[0] = MPI_INT;
20 MPI_Type_struct(2, blocklen, disp, type, &mpi_utype[0]);
21
22 type[0] = MPI_FLOAT;
23 MPI_Type_struct(2, blocklen, disp, type, &mpi_utype[1]);
24
25 for(i=0; i<2; i++) MPI_Type_commit(&mpi_utype[i]);
26
27 /* actual communication */
28
29 MPI_Send(u, 1000, mpi_utype[utype], dest, tag, [ticket116.][comm]MPI_COMM_WORLD);
30

```

**Example 4.20** This example shows how a datatype can be decoded. The routine `printdatatype` prints out the elements of the datatype. Note the use of `MPI_Type_free` for datatypes that are not predefined.

```

31
32
33
34 /*
35    Example of decoding a datatype.
36
37    Returns 0 if the datatype is predefined, 1 otherwise
38 */
39 #include <stdio.h>
40 #include <stdlib.h>
41 #include "mpi.h"
42 int printdatatype( MPI_Datatype datatype )
43 {
44     int *array_of_ints;
45     MPI_Aint *array_of_adds;
46     MPI_Datatype *array_of_dtypes;
47     int num_ints, num_adds, num_dtypes, combiner;
48

```

```

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

```

## 4.2 Pack and Unpack

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

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

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

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

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

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

Packs the message in the send buffer specified by `inbuf`, `incount`, `datatype` into the buffer space specified by `outbuf` and `outsize`. The input buffer can be any communication buffer allowed in `MPI_SEND`. The output buffer is a contiguous storage area containing `outsize` bytes, starting at the address `outbuf` (length is counted in bytes, not elements, as if it were a communication buffer for a message of type `MPI_PACKED`).

The input value of `position` is the first location in the output buffer to be used for packing. `position` is incremented by the size of the packed message, and the output value of `position` is the first location in the output buffer following the locations occupied by the packed message. The `comm` argument is the communicator that will be subsequently used for sending the packed message.

```

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

```

```

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

```

```

15 MPI_UNPACK(INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT, DATATYPE, COMM,
16            IERROR)

```

```

17     <type> INBUF(*), OUTBUF(*)

```

```

18     INTEGER INSIZE, POSITION, OUTCOUNT, DATATYPE, COMM, IERROR

```

```

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

```

Unpacks a message into the receive buffer specified by `outbuf`, `outcount`, `datatype` from the buffer space specified by `inbuf` and `insize`. The output buffer can be any communication buffer allowed in `MPI_RECV`. The input buffer is a contiguous storage area containing `insize` bytes, starting at address `inbuf`. The input value of `position` is the first location in the input buffer occupied by the packed message. `position` is incremented by the size of the packed message, so that the output value of `position` is the first location in the input buffer after the locations occupied by the message that was unpacked. `comm` is the communicator used to receive the packed message.

*Advice to users.* Note the difference between `MPI_RECV` and `MPI_UNPACK`: in `MPI_RECV`, the `count` argument specifies the maximum number of items that can be received. The actual number of items received is determined by the length of the incoming message. In `MPI_UNPACK`, the `count` argument specifies the actual number of items that are unpacked; the “size” of the corresponding message is the increment in `position`. The reason for this change is that the “incoming message size” is not predetermined since the user decides how much to unpack; nor is it easy to determine the “message size” from the number of items to be unpacked. In fact, in a heterogeneous system, this number may not be determined *a priori*. (*End of advice to users.*)

To understand the behavior of pack and unpack, it is convenient to think of the data part of a message as being the sequence obtained by concatenating the successive values sent in that message. The pack operation stores this sequence in the buffer space, as if sending the message to that buffer. The unpack operation retrieves this sequence from buffer space, as if receiving a message from that buffer. (It is helpful to think of internal Fortran files or `sscanf` in C, for a similar function.)



Several messages can be successively packed into one **packing unit**. This is effected by several successive **related** calls to `MPI_PACK`, where the first call provides `position = 0`, and each successive call inputs the value of `position` that was output by the previous call, and the same values for `outbuf`, `outcount` and `comm`. This packing unit now contains the equivalent information that would have been stored in a message by one send call with a send buffer that is the “concatenation” of the individual send buffers.

A packing unit can be sent using type `MPI_PACKED`. Any point to point or collective communication function can be used to move the sequence of bytes that forms the packing unit from one process to another. This packing unit can now be received using any receive operation, with any datatype: the type matching rules are relaxed for messages sent with type `MPI_PACKED`.

A message sent with any type (including `MPI_PACKED`) can be received using the type `MPI_PACKED`. Such a message can then be unpacked by calls to `MPI_UNPACK`.

A packing unit (or a message created by a regular, “typed” send) can be unpacked into several successive messages. This is effected by several successive related calls to `MPI_UNPACK`, where the first call provides `position = 0`, and each successive call inputs the value of `position` that was output by the previous call, and the same values for `inbuf`, `insize` and `comm`.

The concatenation of two packing units is not necessarily a packing unit; nor is a substring of a packing unit necessarily a packing unit. Thus, one cannot concatenate two packing units and then unpack the result as one packing unit; nor can one unpack a substring of a packing unit as a separate packing unit. Each packing unit, that was created by a related sequence of pack calls, or by a regular send, must be unpacked as a unit, by a sequence of related unpack calls.

*Rationale.* The restriction on “atomic” packing and unpacking of packing units allows the implementation to add at the head of packing units additional information, such as a description of the sender architecture (to be used for type conversion, in a heterogeneous environment) (*End of rationale.*)

The following call allows the user to find out how much space is needed to pack a message and, thus, manage space allocation for buffers.

`MPI_PACK_SIZE(incount, datatype, comm, size)`

IN	<code>incount</code>	count argument to packing call (non-negative integer)
IN	<code>datatype</code>	datatype argument to packing call (handle)
IN	<code>comm</code>	communicator argument to packing call (handle)
OUT	<code>size</code>	upper bound on size of packed message, in bytes (non-negative integer)

```
int MPI_Pack_size(int incount, MPI_Datatype datatype, MPI_Comm comm,
                  int *size)
```

```
MPI_PACK_SIZE(INCOUNT, DATATYPE, COMM, SIZE, IERROR)
INTEGER INCOUNT, DATATYPE, COMM, SIZE, IERROR
```

```

1  {int MPI::Datatype::Pack_size(int incount, const MPI::Comm& comm) const
2      (binding deprecated, see Section ??) }

```

ticket150.

A call to `MPI_PACK_SIZE(incount, datatype, comm, size)` returns in `size` an upper bound on the increment in position that is effected by a call to `MPI_PACK(inbuf, incount, datatype, outbuf, outcount, position, comm)`.

*Rationale.* The call returns an upper bound, rather than an exact bound, since the exact amount of space needed to pack the message may depend on the context (e.g., first message packed in a packing unit may take more space). (*End of rationale.*)

**Example 4.21** An example using `MPI_PACK`.

```

12  int      position, i, j, a[2];
13  char     buff[1000];
14  [ticket116.]MPI_Status status;
15
16  MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
17  if (myrank == 0)
18  {
19      /*[ticket116.] [ ]* SENDER CODE */
20
21      position = 0;
22      MPI_Pack(&i, 1, MPI_INT, buff, 1000, &position, MPI_COMM_WORLD);
23      MPI_Pack(&j, 1, MPI_INT, buff, 1000, &position, MPI_COMM_WORLD);
24      MPI_Send( buff, position, MPI_PACKED, 1, 0, MPI_COMM_WORLD);
25  }
26  else /* RECEIVER CODE */
27      MPI_Recv( a, 2, MPI_INT, 0, 0, MPI_COMM_WORLD[ticket116.], &status);
28

```

**Example 4.22** An elaborate example.

```

31  int      position, i;
32  float    a[1000];
33  char     buff[1000];
34  [ticket116.] [...]
35
36  MPI_Comm_rank([ticket116.] [MPI_Comm_world]MPI_COMM_WORLD, &myrank);
37  if (myrank == 0)
38  {
39      /*[ticket116.] [ ]* SENDER CODE */
40
41      int len[2];
42      MPI_Aint disp[2];
43      MPI_Datatype type[2], newtype;
44
45      /* build datatype for i followed by a[0]...a[i-1] */
46
47      len[0] = 1;
48      len[1] = i;

```

```

MPI_Address( &i, disp);
MPI_Address( a, disp+1);
type[0] = MPI_INT;
type[1] = MPI_FLOAT;
MPI_Type_struct( 2, len, disp, type, &newtype);
MPI_Type_commit( &newtype);

/* Pack i followed by a[0]...a[i-1]*/

position = 0;
MPI_Pack( MPI_BOTTOM, 1, newtype, buff, 1000, &position, MPI_COMM_WORLD);

/* Send */

MPI_Send( buff, position, MPI_PACKED, 1, 0,
          MPI_COMM_WORLD)[ticket116.];

/* *****
   One can replace the last three lines with
   MPI_Send( MPI_BOTTOM, 1, newtype, 1, 0, MPI_COMM_WORLD);
   ***** */
}
else if (myrank == 1)
{
    /* RECEIVER CODE */

    MPI_Status status;

    /* Receive */

    MPI_Recv( buff, 1000, MPI_PACKED, 0, 0, [ticket116.]MPI_COMM_WORLD, &status);

    /* Unpack i */

    position = 0;
    MPI_Unpack(buff, 1000, &position, &i, 1, MPI_INT, MPI_COMM_WORLD);

    /* Unpack a[0]...a[i-1] */
    MPI_Unpack(buff, 1000, &position, a, i, MPI_FLOAT, MPI_COMM_WORLD);
}

```

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

```

int  count, gsize, counts[64], totalcount, k1, k2, k,
     displs[64], position, concat_pos;
char chr[100], *lbuf, *rbuf, *cbuf;
[ticket116.][....]

```

```

1  MPI_Comm_size(comm, &gsize);
2  MPI_Comm_rank(comm, &myrank);
3
4      /* allocate local pack buffer */
5  MPI_Pack_size(1, MPI_INT, comm, &k1);
6  MPI_Pack_size(count, MPI_CHAR, comm, &k2);
7  k = k1+k2;
8  lbuf = (char *)malloc(k);
9
10     /* pack count, followed by count characters */
11 position = 0;
12 MPI_Pack(&count, 1, MPI_INT, lbuf, k, &position, comm);
13 MPI_Pack(chr, count, MPI_CHAR, lbuf, k, &position, comm);
14
15 if (myrank != root) {
16     /* gather at root sizes of all packed messages */
17     MPI_Gather( &position, 1, MPI_INT, NULL, [ticket116.] [NULL]0,
18               [ticket116.] [NULL]MPI_DATATYPE_NULL, root, comm);
19
20     /* gather at root packed messages */
21     MPI_Gatherv( lbuf, position, MPI_PACKED, NULL,
22                 NULL, NULL, NULL, root, comm);
23
24 } else { /* root code */
25     /* gather sizes of all packed messages */
26     MPI_Gather( &position, 1, MPI_INT, counts, 1,
27               MPI_INT, root, comm);
28
29     /* gather all packed messages */
30     displs[0] = 0;
31     for (i=1; i < gsize; i++)
32         displs[i] = displs[i-1] + counts[i-1];
33     totalcount = di[ticket116.]spls[gsize-1] + counts[gsize-1];
34     rbuf = (char *)malloc(totalcount);
35     cbuf = (char *)malloc(totalcount);
36     MPI_Gatherv( lbuf, position, MPI_PACKED, rbuf,
37                 counts, displs, MPI_PACKED, root, comm);
38
39     /* unpack all messages and concatenate strings */
40     concat_pos = 0;
41     for (i=0; i < gsize; i++) {
42         position = 0;
43         MPI_Unpack( rbuf+displs[i], totalcount-displs[i],
44                   &position, &count, 1, MPI_INT, comm);
45         MPI_Unpack( rbuf+displs[i], totalcount-displs[i],
46                   &position, cbuf+concat_pos, count, MPI_CHAR, comm);
47         concat_pos += count;
48     }

```

```

    cbuf[concat_pos] = [ticket116.][‘]’\0’;
}

```

### 4.3 Canonical MPI\_PACK and MPI\_UNPACK

These functions read/write data to/from the buffer in the “external32” data format specified in Section 13.5.2, and calculate the size needed for packing. Their first arguments specify the data format, for future extensibility, but currently the only valid value of the `datarep` argument is “external32.”

*Advice to users.* These functions could be used, for example, to send typed data in a portable format from one MPI implementation to another. (*End of advice to users.*)

The buffer will contain exactly the packed data, without headers. `MPI_BYTE` should be used to send and receive data that is packed using `MPI_PACK_EXTERNAL`.

*Rationale.* `MPI_PACK_EXTERNAL` specifies that there is no header on the message and further specifies the exact format of the data. Since `MPI_PACK` may (and is allowed to) use a header, the datatype `MPI_PACKED` cannot be used for data packed with `MPI_PACK_EXTERNAL`. (*End of rationale.*)

`MPI_PACK_EXTERNAL(datarep, inbuf, incount, datatype, outbuf, outsize, position )`

IN	<code>datarep</code>	data representation (string)
IN	<code>inbuf</code>	input buffer start (choice)
IN	<code>incount</code>	number of input data items (integer)
IN	<code>datatype</code>	datatype of each input data item (handle)
OUT	<code>outbuf</code>	output buffer start (choice)
IN	<code>outsize</code>	output buffer size, in bytes (integer)
INOUT	<code>position</code>	current position in buffer, in bytes (integer)

```

int MPI_Pack_external(char *datarep, void *inbuf, int incount,
    MPI_Datatype datatype, void *outbuf, MPI_Aint outsize,
    MPI_Aint *position)

```

```

MPI_PACK_EXTERNAL(DATAREP, INBUF, INCOUNT, DATATYPE, OUTBUF, OUTSIZE,
    POSITION, IERROR)

```

```

INTEGER INCOUNT, DATATYPE, IERROR
INTEGER(KIND=MPI_ADDRESS_KIND) OUTSIZE, POSITION
CHARACTER*(*) DATAREP
<type> INBUF(*), OUTBUF(*)

```

```

{void MPI::Datatype::Pack_external(const char* datarep, const void* inbuf,
    int incount, void* outbuf, MPI::Aint outsize,
    MPI::Aint& position) const (binding deprecated, see Section ??) }

```

ticket150.

ticket150.

```

1 MPI_UNPACK_EXTERNAL(datarep, inbuf, insize, position, outbuf, outsize, position )
2     IN      datarep      data representation (string)
3     IN      inbuf        input buffer start (choice)
4     IN      insize        input buffer size, in bytes (integer)
5     INOUT   position      current position in buffer, in bytes (integer)
6     OUT     outbuf        output buffer start (choice)
7     IN      outcount      number of output data items (integer)
8     IN      datatype      datatype of output data item (handle)

```

```

12 int MPI_Unpack_external(char *datarep, void *inbuf, MPI_Aint insize,
13     MPI_Aint *position, void *outbuf, int outcount,
14     MPI_Datatype datatype)
15
16 MPI_UNPACK_EXTERNAL(DATAREP, INBUF, INSIZE, POSITION, OUTBUF, OUTCOUNT,
17     DATATYPE, IERROR)
18     INTEGER OUTCOUNT, DATATYPE, IERROR
19     INTEGER(KIND=MPI_ADDRESS_KIND) INSIZE, POSITION
20     CHARACTER*(*) DATAREP
21     <type> INBUF(*), OUTBUF(*)

```

```

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

```

```

28 MPI_PACK_EXTERNAL_SIZE( datarep, incount, datatype, size )
29
30     IN      datarep      data representation (string)
31     IN      incount      number of input data items (integer)
32     IN      datatype      datatype of each input data item (handle)
33     OUT     size          output buffer size, in bytes (integer)

```

```

36 int MPI_Pack_external_size(char *datarep, int incount,
37     MPI_Datatype datatype, MPI_Aint *size)
38
39 MPI_PACK_EXTERNAL_SIZE(DATAREP, INCOUNT, DATATYPE, SIZE, IERROR)
40     INTEGER INCOUNT, DATATYPE, IERROR
41     INTEGER(KIND=MPI_ADDRESS_KIND) SIZE
42     CHARACTER*(*) DATAREP

```

```

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

```

# Bibliography

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

# Index

CONST:&, 19  
 CONST:int, 19  
 CONST:MPI::Aint, 3, 3, 5, 8, 10, 18, 21, 22, 31, 51, 52  
 CONST:MPI::Datatype, 3  
 CONST:MPI::Status, 26  
 CONST:MPI\_ADDRESS\_KIND, 30  
 CONST:MPI\_Aint, 3, 3, 5, 8, 10, 18, 21, 22, 31, 51, 52  
 CONST:MPI\_BOTTOM, 18, 28  
 CONST:MPI\_BYTE, 51  
 CONST:MPI\_CHAR, 10  
 CONST:MPI\_COMBINER\_CONTIGUOUS, 30, 33  
 CONST:MPI\_COMBINER\_DARRAY, 30, 35  
 CONST:MPI\_COMBINER\_DUP, 30, 33  
 CONST:MPI\_COMBINER\_F90\_COMPLEX, 30, 35  
 CONST:MPI\_COMBINER\_F90\_INTEGER, 30, 35  
 CONST:MPI\_COMBINER\_F90\_REAL, 30, 35  
 CONST:MPI\_COMBINER\_HINDEXED, 30, 34  
 CONST:MPI\_COMBINER\_HINDEXED\_INTEGER, 30, 34  
 CONST:MPI\_COMBINER\_HVECTOR, 30, 34  
 CONST:MPI\_COMBINER\_HVECTOR\_INTEGER, 30, 34  
 CONST:MPI\_COMBINER\_INDEXED, 30, 34  
 CONST:MPI\_COMBINER\_INDEXED\_BLOCK, 30, 34  
 CONST:MPI\_COMBINER\_NAMED, 30, 33  
 CONST:MPI\_COMBINER\_RESIZED, 30, 35  
 CONST:MPI\_COMBINER\_STRUCT, 30, 34  
 CONST:MPI\_COMBINER\_STRUCT\_INTEGER, 30, 34  
 CONST:MPI\_COMBINER\_SUBARRAY, 30, 34  
 CONST:MPI\_COMBINER\_VECTOR, 30, 34  
 CONST:MPI\_Datatype, 3  
 CONST:MPI\_DATATYPE\_NULL, 24  
 CONST:MPI\_DISTRIBUTE\_BLOCK, 15  
 CONST:MPI\_DISTRIBUTE\_CYCLIC, 15  
 CONST:MPI\_DISTRIBUTE\_DFLT\_DARG, 15  
 CONST:MPI\_DISTRIBUTE\_NONE, 15  
 CONST:MPI\_FLOAT, 10  
 CONST:MPI\_INT, 2  
 CONST:MPI\_INTEGER\_KIND, 30  
 CONST:MPI\_LB, 13, 16, 20–22, 25  
 CONST:MPI\_ORDER\_C, 12, 15, 16  
 CONST:MPI\_ORDER\_FORTRAN, 12, 15  
 CONST:MPI\_PACKED, 45, 47, 51  
 CONST:MPI\_Status, 26  
 CONST:MPI\_UB, 13, 17, 20–22, 25  
 CONST:MPI\_UNDEFINED, 27  
 CONST:v, 28  
 EXAMPLES:Datatype  
     3D array, 36  
     absolute addresses, 41  
     array of structures, 38  
     elaborate example, 48, 49  
     matching type, 25  
     matrix transpose, 37  
     union, 42  
 EXAMPLES:MPI\_ADDRESS, 19  
 EXAMPLES:MPI\_Address, 38, 41, 42, 48  
 EXAMPLES:MPI\_Aint, 38  
 EXAMPLES:MPI\_Gather, 49  
 EXAMPLES:MPI\_Gatherv, 49  
 EXAMPLES:MPI\_GET\_ADDRESS, 19  
 EXAMPLES:MPI\_Get\_address, 38, 41, 42, 48  
 EXAMPLES:MPI\_GET\_COUNT, 27  
 EXAMPLES:MPI\_GET\_ELEMENTS, 27  
 EXAMPLES:MPI\_Pack, 48, 49  
 EXAMPLES:MPI\_Pack\_size, 49



- EXAMPLES:MPI\_RECV, [25](#)  
 EXAMPLES:MPI\_SEND, [25](#)  
 EXAMPLES:MPI\_Send, [38](#), [41](#), [42](#), [48](#)  
 EXAMPLES:MPI\_SENDR, [36](#), [37](#)  
 EXAMPLES:MPI\_TYPE\_COMMIT, [23](#), [36](#),  
     [37](#)  
 EXAMPLES:MPI\_Type\_commit, [38](#), [41](#), [42](#),  
     [48](#)  
 EXAMPLES:MPI\_TYPE\_CONTIGUOUS, [3](#),  
     [20](#), [25](#), [27](#)  
 EXAMPLES:MPI\_TYPE\_CREATE\_DARRAY, [46](#)  
     [17](#)  
 EXAMPLES:MPI\_TYPE\_CREATE\_HVECTOR, [36](#), [37](#)  
 EXAMPLES:MPI\_Type\_create\_hvector, [38](#),  
     [41](#)  
 EXAMPLES:MPI\_TYPE\_CREATE\_STRUCT, [10](#), [20](#), [37](#)  
 EXAMPLES:MPI\_Type\_create\_struct, [38](#), [41](#),  
     [42](#), [48](#)  
 EXAMPLES:MPI\_TYPE\_EXTENT, [36](#), [37](#)  
 EXAMPLES:MPI\_Type\_extent, [38](#)  
 EXAMPLES:MPI\_Type\_get\_contents, [43](#)  
 EXAMPLES:MPI\_Type\_get\_envelope, [43](#)  
 EXAMPLES:MPI\_TYPE\_HVECTOR, [36](#), [37](#)  
 EXAMPLES:MPI\_Type\_hvector, [38](#), [41](#)  
 EXAMPLES:MPI\_TYPE\_INDEXED, [7](#), [36](#)  
 EXAMPLES:MPI\_Type\_indexed, [38](#), [41](#)  
 EXAMPLES:MPI\_TYPE\_STRUCT, [10](#), [20](#),  
     [37](#)  
 EXAMPLES:MPI\_Type\_struct, [38](#), [41](#), [42](#),  
     [48](#)  
 EXAMPLES:MPI\_TYPE\_VECTOR, [4](#), [5](#), [36](#),  
     [37](#)  
 EXAMPLES:MPI\_Unpack, [48](#), [49](#)  
 EXAMPLES:Typemap, [3-5](#), [7](#), [10](#), [17](#)  
  
 MPI\_ADDRESS, [18](#)  
 MPI\_FILE\_SET\_VIEW, [13](#)  
 MPI\_GET\_ADDRESS, [3](#), [18](#), [19](#), [28](#)  
 MPI\_GET\_ADDRESS(location, address), [18](#)  
 MPI\_GET\_COUNT, [26](#), [27](#)  
 MPI\_GET\_ELEMENTS, [26](#), [27](#)  
 MPI\_GET\_ELEMENTS( status, datatype,  
     count), [26](#)  
 MPI\_PACK, [48](#), [51](#)  
 MPI\_PACK(inbuf, incount, datatype, out-  
     buf, outcount, position, comm), [48](#)  
 MPI\_PACK(inbuf, incount, datatype, out-  
     buf, outsize, position, comm), [45](#)  
 MPI\_PACK\_EXTERNAL, [51](#)  
 MPI\_PACK\_EXTERNAL(datarep, inbuf, in-  
     count, datatype, outbuf, outsize, po-  
     sition ), [51](#)  
 MPI\_PACK\_EXTERNAL\_SIZE( datarep, in-  
     count, datatype, size ), [52](#)  
 MPI\_PACK\_SIZE(incount, datatype, comm,  
     size), [47](#), [48](#)  
 MPI\_RECV, [46](#)  
 MPI\_RECV(buf, 1, datatype,...), [2](#)  
 MPI\_RECV(buf, count, datatype, dest, tag,  
     comm, status), [26](#)  
 MPI\_RECV(buf, count, datatype, source, tag,  
     comm, status), [25](#)  
 MPI\_SEND, [45](#)  
 MPI\_SEND(buf, 1, datatype,...), [2](#)  
 MPI\_SEND(buf, count, datatype , ...), [25](#)  
 MPI\_SEND(buf, count, datatype, dest, tag,  
     comm), [25](#)  
 MPI\_TYPE\_COMMIT, [23](#)  
 MPI\_TYPE\_COMMIT(datatype), [23](#)  
 MPI\_TYPE\_CONTIGUOUS, [3](#), [30](#)  
 MPI\_TYPE\_CONTIGUOUS(2, type1, type2),  
     [20](#)  
 MPI\_TYPE\_CONTIGUOUS(count, oldtype,  
     newtype), [3](#), [5](#)  
 MPI\_TYPE\_CREATE\_DARRAY, [14](#), [30](#)  
 MPI\_TYPE\_CREATE\_DARRAY(size, rank,  
     ndims, array\_of\_gsizes, array\_of\_distrib-  
     array\_of\_dargs, array\_of\_psize, or-  
     der, oldtype, newtype), [14](#)  
 MPI\_TYPE\_CREATE\_F90\_COMPLEX, [30](#),  
     [32](#)  
 MPI\_TYPE\_CREATE\_F90\_INTEGER, [30](#),  
     [32](#)  
 MPI\_TYPE\_CREATE\_F90\_REAL, [30](#), [32](#)  
 MPI\_TYPE\_CREATE\_HINDEXED, [3](#), [8](#), [10](#),  
     [30](#)  
 MPI\_TYPE\_CREATE\_HINDEXED( count,  
     array\_of\_blocklengths, array\_of\_displacements,  
     oldtype, newtype), [8](#)  
 MPI\_TYPE\_CREATE\_HINDEXED(count, B,  
     D, oldtype, newtype), [11](#)  
 MPI\_TYPE\_CREATE\_HVECTOR, [3](#), [5](#), [30](#)  
 MPI\_TYPE\_CREATE\_HVECTOR( count, block-  
     length, stride, oldtype, newtype), [5](#)

MPI\_TYPE\_CREATE\_INDEXED\_BLOCK, MPI\_TYPE\_INDEXED(count, B, D, oldtype,  
 30 newtype), 8  
 MPI\_TYPE\_CREATE\_INDEXED\_BLOCK(count, MPI\_TYPE\_LB, 21  
 blocklength, array\_of\_displacements, MPI\_TYPE\_SIZE, 19  
 oldtype, newtype), 9 MPI\_TYPE\_SIZE(datatype, size), 19  
 MPI\_TYPE\_CREATE\_RESIZED, 30 MPI\_TYPE\_STRUCT, 10, 30  
 MPI\_TYPE\_CREATE\_RESIZED(oldtype, lb, MPI\_TYPE\_STRUCT(3, B, D, T, newtype),  
 extent, newtype), 21 10  
 MPI\_TYPE\_CREATE\_STRUCT, 3, 30 MPI\_TYPE\_STRUCT(3, B, D, T, type1),  
 MPI\_TYPE\_CREATE\_STRUCT(count, ar- 20  
 ray\_of\_blocklengths, array\_of\_displacements, MPI\_TYPE\_UB, 21  
 array\_of\_types, newtype), 10 MPI\_TYPE\_VECTOR, 4, 5, 30  
 MPI\_TYPE\_CREATE\_STRUCT(count, B, MPI\_TYPE\_VECTOR( 2, 3, 4, oldtype, new-  
 D, T, newtype), 11 type), 4  
 MPI\_TYPE\_CREATE\_SUBARRAY, 13, 15, MPI\_TYPE\_VECTOR( count, blocklength,  
 30 stride, oldtype, newtype), 4  
 MPI\_TYPE\_CREATE\_SUBARRAY(ndims, MPI\_TYPE\_VECTOR(1, count, n, oldtype,  
 array\_of\_sizes, array\_of\_subsizes, ar- newtype), 5  
 ray\_of\_starts, order, oldtype, new- MPI\_TYPE\_VECTOR(3, 1, -2, oldtype, new-  
 type), 11 type), 5  
 MPI\_TYPE\_DUP, 24, 30 MPI\_TYPE\_VECTOR(count, 1, 1, oldtype,  
 MPI\_TYPE\_DUP(type, newtype), 24 newtype), 5  
 MPI\_TYPE\_EXTENT, 21 MPI\_TYPE\_VECTOR(count, blocklength, stride,  
 MPI\_TYPE\_FREE, 32 oldtype, newtype), 8  
 MPI\_TYPE\_FREE(datatype), 24 MPI\_UNPACK, 46, 47, 51  
 MPI\_TYPE\_GET\_CONTENTS, 29, 32, 33 MPI\_UNPACK(inbuf, insize, position, out-  
 MPI\_TYPE\_GET\_CONTENTS(datatype, max\_integers, inbuf, outcount, datatype, comm), 46  
 max\_addresses, max\_datatypes, ar- MPI\_UNPACK\_EXTERNAL(datarep, inbuf,  
 ray\_of\_integers, array\_of\_addresses, insize, position, outbuf, outsize, po-  
 array\_of\_datatypes), 31 sition ), 52  
 MPI\_TYPE\_GET\_ENVELOPE, 29, 32, 33  
 MPI\_TYPE\_GET\_ENVELOPE(datatype, num\_integers,  
 num\_addresses, num\_datatypes, com-  
 biner), 29  
 MPI\_TYPE\_GET\_EXTENT, 23  
 MPI\_TYPE\_GET\_EXTENT(datatype, lb, ex-  
 tent), 21  
 MPI\_TYPE\_GET\_TRUE\_EXTENT(datatype,  
 true\_lb, true\_extent), 22  
 MPI\_TYPE\_HINDEXED, 8, 30  
 MPI\_TYPE\_HVECTOR, 6, 30  
 MPI\_TYPE\_INDEXED, 6, 8, 9, 30  
 MPI\_TYPE\_INDEXED( count, array\_of\_blocklengths,  
 array\_of\_displacements, oldtype, new-  
 type), 7  
 MPI\_TYPE\_INDEXED(2, B, D, oldtype, new-  
 type), 7