

# MPI: A Message-Passing Interface Standard

Version 3.0

⊤ (Fin2)

⊥ (Fin2)

Message Passing Interface Forum

Draft March 14th, 2011

# Contents

<b>1</b>	<b>Tool Interfaces</b>	<b>1</b>
1.1	Introduction . . . . .	1
1.2	Profiling Interface . . . . .	1
1.2.1	Requirements . . . . .	1
1.2.2	Discussion . . . . .	2
1.2.3	Logic of the Design . . . . .	2
	Miscellaneous Control of Profiling . . . . .	3
1.2.4	Profiler Implementation Example . . . . .	4
1.2.5	MPI Library Implementation Example . . . . .	4
1.2.6	Complications . . . . .	5
	Multiple Counting . . . . .	5
	Linker Oddities . . . . .	6
1.2.7	Multiple Levels of Interception . . . . .	6
1.3	MPI_T Tool Information Interface . . . . .	6
1.3.1	Verbosity Levels . . . . .	7
1.3.2	Binding of MPI_T Variables to MPI Objects . . . . .	8
1.3.3	String Arguments . . . . .	9
1.3.4	Initialization and Finalization . . . . .	9
1.3.5	Datatype System . . . . .	10
1.3.6	Control Variables . . . . .	13
	Control Variable Query Functions . . . . .	13
	Handle Allocation and Deallocation . . . . .	15
	Control Variable Access Functions . . . . .	16
1.3.7	Performance Variables . . . . .	17
	Performance Variable Classes . . . . .	17
	Performance Variable Query Functions . . . . .	19
	Performance Experiment Sessions . . . . .	21
	Handle Allocation and Deallocation . . . . .	22
	Starting and Stopping of Performance Variables . . . . .	23
	Performance Variable Access Functions . . . . .	24
1.3.8	Variable Categorization . . . . .	26
1.3.9	Return and Error Codes . . . . .	29
1.3.10	Profiling Interface . . . . .	29
	<b>Bibliography</b>	<b>31</b>
	<b>Examples Index</b>	<b>32</b>

# List of Figures

# List of Tables

1.1	MPI_T verbosity levels. . . . .	8
1.2	Constants to identify associations of MPI_T control variables. . . . .	8
1.3	Predefined MPI_T datatypes and their MPI equivalents. . . . .	11
1.4	MPI_T datatype classes. . . . .	11
1.5	Scopes for MPI_T control variables. . . . .	15
1.6	Return and error codes used MPI_T functions. . . . .	30

# Chapter 1

## Tool Interfaces

### 1.1 Introduction

This chapter discusses a set of interfaces that allows debuggers, performance analyzers, and other tools to extract information about the operation of MPI processes. Specifically, this chapter defines both the PMPI profiling interface (Section 1.2) for transparently intercepting and inspecting any MPI call, and the MPI\_T tool information interface (Section 1.3) for querying MPI control and performance variables. The interfaces described in this chapter are all defined in the context of an MPI process, i.e., are callable from the same code that invokes other MPI functions.

### 1.2 Profiling Interface

#### 1.2.1 Requirements

To meet [the]the requirements for the MPI profiling interface, an implementation of the MPI\_T functions *must*

1. provide a mechanism through which all of the MPI defined [functions]functions, except those allowed as macros (See Section ??[?]), may be accessed with a name shift. This requires, in C and Fortran, an alternate entry point name, with the prefix PMPI\_ for each MPI function. The profiling interface in C++ is described in Section ??. For routines implemented as macros, it is still required that the PMPI\_ version be supplied and work as expected, but it is not possible to replace at link time the MPI\_ version with a user-defined version.
2. ensure that those MPI functions that are not replaced may still be linked into an executable image without causing name clashes.
3. document the implementation of different language bindings of the MPI interface if they are layered on top of each other, so that the profiler developer knows whether she must implement the profile interface for each binding, or can [economise]economize by implementing it only for the lowest level routines.
4. where the implementation of different language bindings is done through a layered approach ([e.g.]e.g., the Fortran binding is a set of “wrapper” functions that call the

C implementation), ensure that these wrapper functions are separable from the rest of the library.

This separability is necessary to allow a separate profiling library to be correctly implemented, since (at least with Unix linker semantics) the profiling library must contain these wrapper functions if it is to perform as expected. This requirement allows the person who builds the profiling library to extract these functions from the original MPI library and add them into the profiling library without bringing along any other unnecessary code.

5. provide a no-op routine `MPI_PCONTROL` in the MPI library.

### 1.2.2 Discussion

The objective of the MPI profiling interface is to ensure that it is relatively easy for authors of profiling (and other similar) tools to interface their codes to MPI implementations on different machines.

Since MPI is a machine independent standard with many different implementations, it is unreasonable to expect that the authors of profiling tools for MPI will have access to the source code that implements MPI on any particular machine. It is therefore necessary to provide a mechanism by which the implementors of such tools can collect whatever performance information they wish *without* access to the underlying implementation.

We believe that having such an interface is important if MPI is to be attractive to end users, since the availability of many different tools will be a significant factor in attracting users to the MPI standard.

The profiling interface is just that, an interface. It says *nothing* about the way in which it is used. There is therefore no attempt to lay down what information is collected through the interface, or how the collected information is saved, filtered, or displayed.

While the initial impetus for the development of this interface arose from the desire to permit the implementation of profiling tools, it is clear that an interface like that specified may also prove useful for other purposes, such as “internetworking” multiple MPI implementations. Since all that is defined is an interface, there is no objection to its being used wherever it is useful.

As the issues being addressed here are intimately tied up with the way in which executable images are built, which may differ greatly on different machines, the examples given below should be treated solely as one way of implementing the objective of the MPI profiling interface. The actual requirements made of an implementation are those detailed in the Requirements section above, the whole of the rest of this chapter is only present as justification and discussion of the logic for those requirements.

The examples below show one way in which an implementation could be constructed to meet the requirements on a Unix system (there are doubtless others that would be equally valid).

### 1.2.3 Logic of the Design

Provided that an MPI implementation meets the requirements above, it is possible for the implementor of the profiling system to intercept all of the MPI calls that are made by the user program. She can then collect whatever information she requires before calling the underlying MPI implementation (through its name shifted entry points) to achieve the desired effects.

## Miscellaneous Control of Profiling

There is a clear requirement for the user code to be able to control the profiler dynamically at run time. This is normally used for (at least) the purposes of

- Enabling and disabling profiling depending on the state of the calculation.
- Flushing trace buffers at non-critical points in the `[calculation]` calculation.
- Adding user events to a trace file.

These requirements are met by use of the `MPI_PCONTROL`.

`MPI_PCONTROL(level, ...)`

IN            level                            Profiling level

`int MPI_Pcontrol(const int level, ...)`

`MPI_PCONTROL(LEVEL)`

INTEGER LEVEL

`{void MPI::Pcontrol(const int level, ...) (binding deprecated, see Section ??) }`

MPI libraries themselves make no use of this routine, and simply return immediately to the user code. However the presence of calls to this routine allows a profiling package to be explicitly called by the user.

Since MPI has no control of the implementation of the profiling code, we are unable to specify precisely the semantics that will be provided by calls to `MPI_PCONTROL`. This vagueness extends to the number of arguments to the function, and their datatypes.

However to provide some level of portability of user codes to different profiling libraries, we request the following meanings for certain values of `level`.

- `level==0` Profiling is disabled.
- `level==1` Profiling is enabled at a normal default level of detail.
- `level==2` Profile buffers are `[flushed. (This may be a no-op in some profilers).]` flushed, which may be a no-op in some profilers.
- All other values of `level` have profile library defined effects and additional arguments.

We also request that the default state after `MPI_INIT` has been called is for profiling to be enabled at the normal default level. (i.e. as if `MPI_PCONTROL` had just been called with the argument 1). This allows users to link with a profiling library and obtain profile output without having to modify their source code at all.

The provision of `MPI_PCONTROL` as a no-op in the standard MPI library `[allows them to modify their source code to obtain]` supports the collection of more detailed profiling information, `[but still be able to link exactly the]` with source `[same code]` code that can still link against the standard MPI library.

### 1.2.4 Profiler Implementation Example

[Suppose that the profiler wishes to]A profiler can accumulate the total amount of data sent by the [MPI\_SEND]MPI\_SEND function, along with the total elapsed time spent in the [function. This could trivially be achieved thus]function, as follows:

```
static int totalBytes = 0;
static double totalTime = 0.0;

int MPI_Send(void* buffer, int count, MPI_Datatype datatype,
             int dest, int tag, MPI_Comm comm)
{
    double tstart = MPI_Wtime();      /* Pass on all the arguments */
    int extent;
    int result = PMPI_Send(buffer, count, datatype, dest, tag, comm);

    MPI_Type_size(datatype, &extent); /* Compute size */
    totalBytes += count*extent;

    totalTime += MPI_Wtime() - tstart; /* and time */

    return result;
}
```

### 1.2.5 MPI Library Implementation Example

[On a Unix system, in which the MPI library is implemented in C, then]If the MPI library is implemented in C on a Unix system, then there [there are various possible options, of which two of the most obvious]are various options, including the two presented here, for supporting [are presented here. Which is better depends on whether the linker and]the name-shift requirement. The choice between these two options [compiler support weak symbols.]depends partly on whether the linker and compiler support weak symbols.

**Systems with Weak Symbols** If the compiler and linker support weak external symbols ([e.g.]e.g., Solaris 2.x, other system V.4 machines), then only a single library is required through the use of `#pragma weak` thus

```
#pragma weak MPI_Example = PMPI_Example

int PMPI_Example(/* appropriate args */)
{
    /* Useful content */
}
```

The effect of this `#pragma` is to define the external symbol `MPI_Example` as a weak definition. This means that the linker will not complain if there is another definition of the symbol (for instance in the profiling library), however if no other definition exists, then the linker will use the weak definition.



Systems Without Weak Symbols In the absence of weak symbols then one possible solution would be to use the C macro pre-processor thus

```

#ifdef PROFILELIB
#   ifdef __STDC__
#       define FUNCTION(name) P##name
#   else
#       define FUNCTION(name) P/**/name
#   endif
#else
#   define FUNCTION(name) name
#endif

```

Each of the user visible functions in the library would then be declared thus

```

int FUNCTION(MPI_Example)(/* appropriate args */)
{
    /* Useful content */
}

```

The same source file can then be compiled to produce both versions of the library, depending on the state of the PROFILELIB macro symbol.

It is required that the standard MPI library be built in such a way that the inclusion of MPI functions can be achieved one at a time. This is a somewhat unpleasant requirement, since it may mean that each external function has to be compiled from a separate file. However this is necessary so that the author of the profiling library need only define those MPI functions that she wishes to intercept, references to any others being fulfilled by the normal MPI library. Therefore the link step can look something like this

```
% cc ... -lmyprof -lpmi -lmpi
```

Here libmyprof.a contains the profiler functions that intercept some of the MPI functions<sup>[1]</sup>, libpmi.a contains the “name shifted” MPI functions, and libmpi.a contains the normal definitions of the MPI functions.

### 1.2.6 Complications

#### Multiple Counting

Since parts of the MPI library may themselves be implemented using more basic MPI functions ([e.g.]e.g., a portable implementation of the collective operations implemented using point to point communications), there is potential for profiling functions to be called from within an MPI function that was called from a profiling function. This could lead to “double counting” of the time spent in the inner routine. Since this effect could actually be useful under some circumstances ([e.g.]e.g., it might allow one to answer the question “How much time is spent in the point to point routines when they’re called from collective functions?”), we have decided not to enforce any restrictions on the author of the MPI library that would overcome this. Therefore the author of the profiling library should be aware of this problem, and guard against it herself. In a single threaded world this is easily achieved through use of a static variable in the profiling code that remembers if you are already inside a profiling routine. It becomes more complex in a multi-threaded environment (as does the meaning of the times recorded<sup>[!]</sup>).

⊤ (Fin2)<sup>48</sup>  
⊥ (Fin2)  
⊤ (Fin2)  
⊥ (Fin2)

## Linker Oddities

The Unix linker traditionally operates in one `[pass:]pass`: the effect of this is that functions from libraries are only included in the image if they are needed at the time the library is scanned. When combined with weak symbols, or multiple definitions of the same function, this can cause odd (and unexpected) effects.

Consider, for instance, an implementation of MPI in which the Fortran binding is achieved by using wrapper functions on top of the C implementation. The author of the profile library then assumes that it is reasonable only to provide profile functions for the C binding, since Fortran will eventually call these, and the cost of the wrappers is assumed to be small. However, if the wrapper functions are not in the profiling library, then none of the profiled entry points will be undefined when the profiling library is called. Therefore none of the profiling code will be included in the image. When the standard MPI library is scanned, the Fortran wrappers will be resolved, and will also pull in the base versions of the MPI functions. The overall effect is that the code will link successfully, but will not be profiled.

To overcome this we must ensure that the Fortran wrapper functions are included in the profiling version of the library. We ensure that this is possible by requiring that these be separable from the rest of the base MPI library. This allows them to be aared out of the base library and into the profiling one.

### 1.2.7 Multiple Levels of Interception

The scheme given here does not directly support the nesting of profiling functions, since it provides only a single alternative name for each MPI function. Consideration was given to an implementation that would allow multiple levels of call interception, however we were unable to construct an implementation of this that did not have the following disadvantages

- assuming a particular implementation language`[.]`,
- imposing a run time cost even when no profiling was taking place.

Since one of the objectives of MPI is to permit efficient, low latency implementations, and it is not the business of a standard to require a particular implementation language, we decided to accept the scheme outlined above.

[Note, however, that it is possible to use the scheme above to implement a multi-level system, since the function called by the user may call many different profiling functions before calling the underlying MPI function.]

[Unfortunately such an implementation may require more cooperation between the different profiling libraries than is required for the single level implementation detailed above.]Note, however, that it is possible to use the scheme above to implement a multi-level system, since the function called by the user may call many different profiling functions before calling the underlying MPI function. This capability has been demonstrated in the  $P^N$ MPI tool infrastructure `[?]`.

## 1.3 MPI\_T Tool Information Interface

To optimize MPI applications or their runtime behavior, it is often advantageous to understand the performance switches an MPI implementation offers to the user as well as to monitor properties and timing information from within the MPI implementation.

The `MPI_T` interface described in this section provides a mechanism for the MPI implementation to expose a set of variables, each of which represent a particular property, setting, or performance measurement from within the MPI implementation. The `MPI_T` interface provides the necessary routines to find all variables that exist in the particular MPI implementation, query their properties, retrieve descriptions about their meaning and access and, if appropriate, alter their values.

The interface is split into two parts: the first part provides information about control variables used by the MPI implementation to fine tune its configuration. The second part provides access to performance variables that can provide insight into internal performance information of the underlying MPI implementation.

To avoid restrictions on the MPI implementation, the `MPI_T` interface allows the implementation to specify which control and performance variables exist. Additionally, the `MPI_T` interface can obtain metadata about each available variable, such as its datatype and size, a textual description, etc.

To avoid conflicts between the standard MPI functionality and the tools-oriented functionality introduced with `MPI_T`, the `MPI_T` interface is contained in its own name space. All identifiers covered by this interface carry the prefix `MPI_T` and can be used independently from the MPI functionality. This includes initialization and finalization of `MPI_T`, which is provided through a separate set of routines. Consequently, `MPI_T` routines can be called before `MPI_INIT` and after `MPI_FINALIZE`.

On success, all `MPI_T` routines return `MPI_T_SUCCESS`, otherwise they return an appropriate error code. Details on error codes can be found in Section 1.3.9. However, errors returned by the `MPI_T` interface are not fatal and do not have any impact on the execution of MPI routines.

*Advice to users.* The number and type of control variables and performance variables can vary between MPI implementations, platforms, and even different builds of the same implementation on the same platform. Hence, any application relying on a particular variable will not be portable.

This interface is primarily intended for performance monitoring tools, support tools, and libraries controlling the application's environment. Application programmers should either avoid using the `MPI_T` interface or avoid being dependent on the existence of a particular control or performance variable. (*End of advice to users.*)

Since the `MPI_T` interface mostly focuses on tools and support libraries, `MPI_T` implementations are only required to provide C bindings. Except where otherwise noted, all conventions and principles governing the C bindings of the MPI API also apply to the `MPI_T` interface. The `MPI_T` interface is available by including the `mpi.h` header file.

### 1.3.1 Verbosity Levels

The `MPI_T` interface provides users access to internal configuration and performance information through a set of control and performance variables defined by the `MPI_T` implementation. Since some implementations may export a large number of variables, variables are classified by a verbosity level that categorizes both their intended audience (end users, performance tuners or MPI implementors) and a relative measure of [complexity]level of detail (basic, detailed or all). See Table 1.3.1

*Advice to implementors.* If an `MPI_T` implementation chooses to use only a single verbosity level for all variables, it is recommended that

MPI_T_VERBOSE_USER_BASIC	Basic information of interest for end users
MPI_T_VERBOSE_USER_DETAIL	Detailed information of interest for end users
MPI_T_VERBOSE_USER_ALL	All information of interest for end users
MPI_T_VERBOSE_TUNER_BASIC	Basic information required for tuning
MPI_T_VERBOSE_TUNER_DETAIL	Detailed information required for tuning
MPI_T_VERBOSE_TUNER_ALL	All information required for tuning
MPI_T_VERBOSE_MPIDEV_BASIC	Basic low-level information for MPI implementors
MPI_T_VERBOSE_MPIDEV_DETAIL	Detailed low-level information for MPI implementors
MPI_T_VERBOSE_MPIDEV_ALL	All low-level information for MPI implementors

Table 1.1: MPI\_T verbosity levels.

MPI\_T\_VERBOSE\_USER\_BASIC be used. If an MPI\_T implementation only uses a single [complexity]level of detail value for all variables in each target audience, it is recommended that all variables be assigned to corresponding BASIC level. (*End of advice to implementors.*)

### 1.3.2 Binding of MPI\_T Variables to MPI Objects

Each MPI\_T variable provides access to a particular control setting or performance property provided by the MPI implementation. [ These variables can apply globally to no specific object or can refer to a particular MPI object such as a communicator, datatype, or one-sided communication window. In the latter case, ] A variable may refer to a particular MPI object such as a communicator, datatype, or one-sided communication window, or the variable may refer more generally to the MPI environment of the process. In the first case, the variable must be bound to exactly one MPI object before it can be used. Table 1.2 lists all MPI object types to which an MPI\_T variable can be bound, together with matching constant that are used by MPI\_T routines to identify the object type.

Constant	MPI object
MPI_T_BIND_NO_OBJECT	N/A; applies globally to entire MPI process
MPI_T_BIND_MPI_COMMUNICATOR	MPI communicators
MPI_T_BIND_MPI_DATATYPE	MPI datatypes
MPI_T_BIND_MPI_ERRORHANDLER	MPI error handlers
MPI_T_BIND_MPI_FILE	MPI file handles
MPI_T_BIND_MPI_GROUP	MPI groups
MPI_T_BIND_MPI_OPERATOR	MPI reduction operators
MPI_T_BIND_MPI_REQUEST	MPI requests
MPI_T_BIND_MPI_WINDOW	MPI windows for one-sided communication

Table 1.2: Constants to identify associations of MPI\_T control variables.

*Rationale.* Some variables have meanings tied to a specific MPI object. Examples include the number of send or receive operations using a particular datatype, the number of times an error handler has been called, or or the communication protocol

and “eager limit” used for a particular communicator. Creating a new `MPI_T` variable for each MPI object could cause the number of variables to grow without bound since they cannot be reused to avoid naming conflicts. By associating `MPI_T` variables with a specific MPI object, only a single variable must be specified and maintained by the MPI implementation, which can then be reused on as many MPI objects of the respective type as created during the program’s execution. (*End of rationale.*)

### 1.3.3 String Arguments

Several `MPI_T` functions return one or more strings. These functions have two arguments for each string to be returned: an `OUT` parameter that identifies a pointer to the buffer in which the string will be returned, and an `IN/OUT` parameter to pass the length of the buffer. The user is responsible for the memory allocation of the buffer and must pass `[T the size of the buffer as the length argument (n). ] the size of the buffer (n) as the length argument.` Let  $n$  be the length value specified to the function. On return, the function writes at most  $n - 1$  of the string’s characters into the buffer, followed by a null terminator. If the returned string’s length is greater than or equal to  $n$ , the string will be truncated to  $n - 1$  characters. In this case, the length of the string plus one (for the terminating null character) is returned in the length argument. If the user passes the null pointer as the buffer argument or passes 0 as the length argument, the function does not return the string and only returns the length of the string plus one in the length argument. `If the [users]user T passes the null pointer as the length argument, the buffer argument is ignored and nothing is returned.`

`MPI_T` does not specify the character encoding of strings in the interface. The only requirement is that strings are terminated with a null character. `MPI_T` reserves all datatype, enumeration datatype items, variables and category names with the prefix `MPI_T` for its own use.

### 1.3.4 Initialization and Finalization

Since the `MPI_T` interface is implemented in a separate name space and is independent of the core MPI functions, it requires a separate set of initialization and finalization routines.

`MPI_T_INIT_THREAD(required, provided)`

IN	required	desired level of thread support (integer)
OUT	provided	provided level of thread support (integer)

`int MPI_T_Init_thread(int required, int *provided)`

All programs or tools that use the `MPI_T` interface must initialize the `MPI_T` interface before calling any other `MPI_T` routine. A user can initialize the `MPI_T` interface by calling `MPI_T_INIT_THREAD`, which can be called multiple times. In addition, this routine initializes the thread environment. The argument `required` is used to specify the desired level of thread support. The possible values and their semantics are identical to the ones that can be used with `MPI_INIT_THREAD` listed in Section ???. The call returns in `provided` information about the actual level of thread support that will be provided by `MPI_T`. It can be one of the four values listed [above] in Section ??.

⌈ (Fin2)  
⌊ (Fin2)

$\top$  (Fin2) 1 *Advice to users.* [ The MPI specification does not require all MPI ranks to be  
 2 executed before the call to MPI\_INIT. If users use MPI\_T before MPI\_INIT, they need  
 3 to call MPI\_T\_INIT\_THREAD on every process that is active as this time. ] The MPI  
 4 specification does not require all MPI processes to exist before the call to MPI\_INIT.  
 5 If MPI\_T is used before MPI\_INIT has been called, MPI\_T\_INIT\_THREAD must be  
 $\perp$  (Fin2) 6 called on each process that exists. Processes created by the MPI implementation  
 7 during MPI\_INIT inherit the status of MPI\_T (whether it is initialized or not as well  
 8 as all active handles) from the process they are created from. (End of advice to  
 9 users.)  
 10  
 11

12 *Advice to implementors.* If MPI\_T\_INIT\_THREAD is called before  
 13 MPI\_INIT\_THREAD, it is possible that the requested and granted thread level for  
 14 MPI\_T\_INIT\_THREAD influences the behavior and return value of  
 15 MPI\_INIT\_THREAD. The same is true for the reverse order. (End of advice to imple-  
 16 mentors.)  
 17  
 18

## 19 MPI\_T\_FINALIZE( )

20  
 21 int MPI\_T\_Finalize(void)  
 22

23 This routine finalizes the use of the MPI\_T interface and may be called as often  
 24 as the corresponding MPI\_T\_INIT\_THREAD routine up to the current point of execution.  
 25 Calling it more times is erroneous. As long as the number of calls to MPI\_T\_FINALIZE  
 26 is smaller than the number of calls to MPI\_T\_INIT\_THREAD up to the current point of  
 27 execution, the MPI\_T interface remains initialized and calls to all MPI\_T routines are  
 28 permissible. Further, additional calls to MPI\_T\_INIT\_THREAD after one or more calls to  
 29 MPI\_T\_FINALIZE are permissible.

30 Once MPI\_T\_FINALIZE is called the same number of times as the routine  
 31 MPI\_T\_INIT\_THREAD up to the current point of execution, the MPI\_T interface is no  
 32 longer initialized. Further, the call to MPI\_T\_FINALIZE that ends the initialization of  
 $\top$  (Fin2) 33 MPI\_T may clean up all MPI\_T state, invalidate all open sessions ([for the concept of  
 $\perp$  (Fin2) 34 Sessions] see Section 1.3.7), and all handles that have been allocated by MPI\_T. MPI\_T  
 35 can be reinitialized by subsequent calls to MPI\_T\_INIT\_THREAD.  
 36

37 At the end of the program execution, unless MPI\_ABORT is called, an application must  
 38 have called MPI\_T\_INIT\_THREAD and MPI\_T\_FINALIZE an equal number of times.  
 39

## 40 1.3.5 Datatype System

41 Since the initialization of MPI\_T is separate from the initialization of MPI, it can not  
 42 be guaranteed that MPI datatypes are available at any time during the usage of MPI\_T.  
 43 Therefore, the MPI\_T interface provides a separate datatype system.

44 The MPI\_T interface requires a significantly simpler type system than MPI itself. All  
 45 datatypes are represented by a value of type MPI\_T\_Datatype and are classified into two  
 46 datatype classes: predefined and enumeration datatypes.  
 47  
 48

```
MPI_T_DATATYPE_GET_CLASS(datatype, datatypeclass)
```

```
IN      datatype      MPI_T datatype to be queried
OUT     datatypeclass  class of the datatype passed in
```

```
int MPI_T_Datatype_get_class(MPI_T_Datatype datatype, int *datatypeclass)
```

This routine returns the datatype class for the datatype provided by the argument `datatype`. This allows users of `MPI_T` to distinguish whether a datatype is an enumeration datatype, e.g., to represent the state of a resource, or is one of the predefined datatypes listed in Table 1.3.5. On return, the `datatypeclass` argument is set to one of the constants listed in Table 1.3.5, if `datatype` represents a valid datatype.

MPI_T Datatype	Equivalent MPI Datatype
MPI_T_INT	MPI_INT
MPI_T_LONG_LONG	MPI_LONG_LONG
MPI_T_CHAR	MPI_CHAR
MPI_T_DOUBLE	MPI_DOUBLE

Table 1.3: Predefined MPI\_T datatypes and their MPI equivalents.

MPI_T_DATATYPE_PREDEFINED	the datatype is a predefined datatype
MPI_T_DATATYPE_ENUMERATION	the datatype is an enumeration datatype

Table 1.4: MPI\_T datatype classes.

Conforming implementations of `MPI_T` must ensure that the `MPI_T` datatypes are equivalent to the listed MPI datatypes for any section of the code in which both MPI and `MPI_T` can be used. In particular, this requires that the size of a value represented by an `MPI_T` datatype and its equivalent MPI datatype are equal and that it is possible to communicate the value of a particular `MPI_T` datatype using the equivalent MPI datatype through regular MPI operations.

*Rationale.* The concept of equivalent `MPI_T` and MPI datatypes allows tools to safely communicate values of `MPI_T` datatypes using MPI message passing functionality. (End of rationale.)

The function `MPI_T_DATATYPE_GET_SIZE` can be used to query the storage size of a variable for each `MPI_T` datatype.

```
MPI_T_DATATYPE_GET_SIZE(datatype, size)
```

```
IN      datatype      MPI_T datatype to be queried
OUT     size          Number of bytes required to store a value of datatype
                        datatype
```

```
int MPI_T_Datatype_get_size(MPI_T_Datatype datatype, int *size)
```



The second datatype class, enumeration datatypes, describes variables with a fixed set of discrete values. These datatypes are represented by integer variables and have MPI\_INT as their equivalent MPI datatype. Their values range from 0 to  $N - 1$ , with a fixed  $N$  that can be queried using MPI\_T\_DATATYPE\_ENUM\_GET\_INFO.

MPI\_T\_DATATYPE\_ENUM\_GET\_INFO(datatype, num, name, name\_len)

IN	datatype	MPI_T datatype to be queried
OUT	num	number of discrete values represented by this enumeration datatype
OUT	name	buffer to return the string containing the name of the enumeration datatype
INOUT	name_len	length of the string and/or buffer for name

```
int MPI_T_Datatype_enum_get_info(MPI_T_Datatype datatype, int *num, char
                                *name, int *name_len)
```

[ This routine returns, if datatype represents a valid enumeration datatype,  $N$  representing the range of the enumeration 0 to  $N - 1$  as well as a string with a name for it. ] If datatype is a valid enumeration datatype, this routine returns the enumeration range and the name of the datatype. If the range runs from 0 to  $N - 1$ , the value  $N$  is returned in num.  $N$  has to be at least one, i.e., has to represent at least one item. The integer values in this range are used to represent the  $N$  items that can be represented by this enumeration type.

The arguments name and name\_len are used to return the name of the datatype as described in Section 1.3.3.

The routine is required to return a name of at least length one. This name must be unique with respect to all other names for MPI\_T datatypes used by the MPI implementation.

Names for the individual items in each enumeration datatype can be queried using MPI\_T\_DATATYPE\_ENUM\_GET\_ITEM.

MPI\_T\_DATATYPE\_ENUM\_GET\_ITEM(datatype, item, name, name\_len)

IN	datatype	MPI_T datatype to be queried
IN	item	item number in the MPI_T datatype to be queried
OUT	name	buffer to return the string containing the name of the enumeration item
INOUT	name_len	length of the string and/or buffer for name

```
int MPI_T_Datatype_enum_get_item(MPI_T_Datatype datatype, int item, char
                                *name, int *name_len)
```

The arguments name and name\_len are used to return the name of the enumeration item as described in Section 1.3.3.



If completed successfully, the routine is required to return a name of at least length one. This name must be unique with respect to all other names of items for the same MPI\_T enumeration datatype.

### 1.3.6 Control Variables

The routines described in this section of the MPI\_T interface specification focus on the ability to list, query, and possibly set control variables exposed by the MPI implementation. These variables can typically be used by the user to fine tune properties and configuration settings of the MPI implementation. On many systems, such variables can be set using environment variables, although other configuration mechanisms may be available, such as configuration files or central configuration registries. A typical example that is available in several existing MPI implementations is the ability to specify an “eager limit”, i.e., an upper bound on the size of messages sent or received using an eager protocol.

#### Control Variable Query Functions

An MPI implementation exports a set of  $N$  control variables through MPI\_T. If  $N$  is zero, then the MPI\_T implementation does not export any control variables, otherwise the provided control variables are indexed from 0 to  $N - 1$ . This index number is used in subsequent MPI\_T calls to identify the individual variables.

An MPI\_T implementation is allowed to increase the number of control variables during the execution of an MPI application when new variables become available through dynamic loading. However, MPI\_T implementations are not allowed to change the index of a control variable or delete a variable once it has been added to the set.

The following function can be used to query the number of control variables,  $N$ :

**MPI\_T\_CVAR\_GET\_NUM(num)**

OUT      num      returns number of control variables

**int MPI\_T\_Cvar\_get\_num(int \*num)**

The function MPI\_T\_CVAR\_GET\_INFO provides access to additional information for each variable.

```

1 MPI_T_CVAR_GET_INFO(index, name, name_len, verbosity, datatype, count, desc, desc_len,
2 bind, attributes)
3
4     IN      index      index of the control variable to be queried
5
6     OUT     name       buffer to return the string containing the name of the
7                        control variable
8
9     INOUT   name_len   length of the string and/or buffer for name
10
11     OUT     verbosity  verbosity level of this variable
12
13     OUT     datatype   MPI_T datatype of the information stored in the con-
14                        trol variable
15
16     OUT     count      number of elements returned
17
18     OUT     desc       buffer to return the string containing a description of
19                        the control variable
20
21     INOUT   desc_len   length of the string and/or buffer for desc
22
23     OUT     bind       type of MPI object to which this variable must be
24                        bound
25
26     OUT     attributes additional attributes defining this variable
27
28
29 int MPI_T_Cvar_get_info(int index, char *name, int *name_len, int
30                        *verbosity, MPI_T_Datatype *datatype, int *count, char *desc,
31                        int *desc_len, int *bind, MPI_T_Cvar_attributes *attributes)
32

```

After a successful call to `MPI_T_CVAR_GET_INFO` for a particular variable, subsequent calls to this routine querying information about the same variable must return the same information. An `MPI_T` implementation is not allowed to alter any of the returned values.

The arguments `name` and `name_len` are used to return the name of the control variable as described in Section 1.3.3.

If completed successfully, the routine is required to return a name of at least length one. The name must be unique with respect to all other names for `MPI_T` control variables used by the MPI implementation.

The argument `verbosity` returns the verbosity level [ (see Section 1.3.1) assigned by the MPI implementation to the variable. ] of the variable (see Section 1.3.1).

The argument `datatype` returns the `MPI_T` datatype [ in which the value for this control variable is returned. ] that is used to represent the control variable.

The arguments `desc` and `desc_len` are used to return a description of the control variable as described in Section 1.3.3.

Returning a description is optional. If an MPI implementation decides not to return a description, the first character for `desc` must be set to the null character and `desc_len` must be set to one at the return of this call.

The parameter `bind` returns the type of the MPI object to which the variable must be bound or the value `MPI_T_BIND_NO_OBJECT` (see Section 1.3.1).

Additional information about the variable is returned through the `attributes` argument using an opaque structure of type `MPI_T_Cvar_attributes` and can be queried using the following accessor function.

[

*Rationale.* The use of opaque attributes enables extensions of the MPI\_T specification in subsequent versions of the MPI standard without having to redefine or alter the query function. Instead new information can be added by adding new accessor functions. (*End of rationale.*)

[ONLY MOVED]

MPI\_T\_CVAR\_ATTR\_GET\_SCOPE(attributes, scope)

IN attributes attributes returned by a previous query call  
OUT scope scope of when changes to this variable are possible

int MPI\_T\_Cvar\_attr\_get\_scope(MPI\_T\_Cvar\_attributes attributes, int \*scope)

The scope of a variable determines whether an operation either local to the [processor] process or collective across multiple processes can change a variable through the MPI\_T interface. On successful return from MPI\_T\_CVAR\_ATTR\_GET\_SCOPE, the argument scope will be set to one of the constants listed in Table 1.3.6.

Scope Constant	Description
MPI_T_SCOPE_READONLY	read-only, cannot be written
MPI_T_SCOPE_LOCAL	may be writeable, writing is a local operation
MPI_T_SCOPE_GLOBAL	may be writeable, writing is a global operation

Table 1.5: Scopes for MPI\_T control variables.

*Advice to users.* The scope of a variable only indicates if a variable might be changeable; it is not a guarantee that it can be changed at any time. If it cannot be changed at a time the user tries to write to it, the MPI\_T implementation is allowed to return an error code as the result of the write operation. (*End of advice to users.*)

[ONLY MOVED]

*Rationale.* The use of opaque attributes enables extensions of the MPI\_T specification in subsequent versions of the MPI standard without having to redefine or alter the query function. Instead new information can be added by adding new accessor functions. (*End of rationale.*)

## Handle Allocation and Deallocation

Before reading or writing the value of a variable, a user must first allocate a handle for it by binding it to an MPI object (see also Section 1.3.1).

*Rationale.* MPI\_T handles are distinct from MPI handles because they must be usable before [MPI\_Init]MPI\_INIT and after [MPI\_Finalize]MPI\_FINALIZE. Further, accessing handles, in particular for performance variables, can be time critical and having a separate handle space enables optimizations. (*End of rationale.*)

1 MPI\_T\_CVAR\_HANDLE\_ALLOC(index, object, handle)

2 IN index index of control variable for which handle is to be al-  
3 located

4 IN obj\_handle reference to a handle of the MPI object to which this  
5 variable is supposed to be bound

6 OUT handle allocated handle

7  
8  
9 int MPI\_T\_Cvar\_handle\_alloc(int index, void \*obj\_handle, MPI\_T\_Cvar\_handle  
10 \*handle)

11  
12 [ This routine allocates a handle for the control variable specified by the argument index  
13 and binds this variable to the MPI object referenced by the pointer to its handle passed  
14 in the argument obj\_handle . The type of the MPI handle passed into this routine must  
15 match the type returned by the bind parameter in a prior call to MPI\_T\_CVAR\_GET\_INFO.  
16 If the type of the object is identified as MPI\_T\_BIND\_NO\_OBJECT, i.e., the variable refers  
17 to the entire MPI process, the argument object is ignored. ] This routine binds the control  
18 variable specified by the argument index to the MPI object referenced by the handle passed  
19 in argument obj\_handle and returns an allocated variable handle in the argument handle.  
20 The value of index should be in the range 0 to  $N - 1$ , where  $N$  is the number of available  
21 control variables as determined from a prior call to MPI\_T\_CVAR\_GET\_NUM. The value of  
22 obj\_handle must be the memory address of the object's MPI handle, and the type of MPI  
23 object it references must be consistent with the type returned in the bind argument in a  
24 prior call to MPI\_T\_CVAR\_GET\_INFO.

25  
26 MPI\_T\_CVAR\_HANDLE\_FREE(handle)

27 INOUT handle handle to be freed

28  
29  
30 int MPI\_T\_Cvar\_handle\_free(MPI\_T\_Cvar\_handle \*handle)

31  
32 When a handle is no longer needed, a user of MPI\_T should call  
33 MPI\_T\_CVAR\_HANDLE\_FREE to free the handle and the associated resources in the MPI\_T  
34 implementation. On a successful return, MPI\_T sets the handle to  
35 MPI\_T\_CVAR\_HANDLE\_NULL.

### 36 Control Variable Access Functions

37  
38  
39 MPI\_T\_CVAR\_READ(handle, buf)

40 IN handle handle to the control variable to be read

41 OUT buf initial address of storage location for variable value

42  
43  
44  
45 int MPI\_T\_Cvar\_read(MPI\_T\_Cvar\_handle handle, void\* buf)

46  
47 The MPI\_T\_CVAR\_READ queries the value of the control variable identified by the  
48 argument handle and stores the result in the buffer buf. The user is responsible to ensure

that the buffer is of the appropriate size [and fits]to hold the entire value of the control variable (based on the returned datatype and count from a prior corresponding call to MPI\_T\_CVAR\_GET\_INFO).

MPI\_T\_CVAR\_WRITE(handle, buf)

IN handle handle to the control variable to be written  
IN buf initial address of storage location for variable value

int MPI\_T\_Cvar\_write(MPI\_T\_Cvar\_handle handle, void\* buf)

The MPI\_T\_CVAR\_WRITE sets the value of the control variable identified by the argument handle to the data stored in the buffer buf. The user is responsible to ensure that the buffer is of the appropriate size [and fits]to hold the entire value of the control variable (based on the returned datatype and count from a prior corresponding call to MPI\_T\_CVAR\_GET\_INFO).

If the variable has a global scope (as returned by a prior corresponding MPI\_T\_CVAR\_ATTR\_GET\_SCOPE call), any write call to this variable must be issued consistently in all connected (as defined in Section ??) MPI processes. The user is responsible to ensure that the writes in all processes are consistent.

If it is not possible to change the variable at the time the call is made, the function returns either MPI\_T\_ERR\_SETNOW, if there may be a later time at which the variable could be set, or MPI\_T\_ERR\_SETNEVER, if the variable cannot be set for the remainder of the application's execution.

### 1.3.7 Performance Variables

The following section focuses on the ability to list and query performance variables provided by the MPI implementation. Performance variables provide insight into MPI implementation specific internals and can represent information such as the state of the MPI implementation (e.g., waiting blocked, receiving, not active), aggregated timing data for submodules, or queue sizes and lengths. Performance variables are always local to [a single]an MPI process.

*Rationale.* The interface for performance variables is separate from the interface for control variables, since performance variables have different requirements and parameters. By keeping them separate, the interface [has a]provides cleaner semantics and allows for more performance optimization opportunities. (End of rationale.)

#### Performance Variable Classes

[ Each performance variable is associated with a class describing its [the] basic semantics. The class of a variable also defines its basic behavior, when and how an MPI implementation can change its value, and what the initial value of this variable is at the time it is either used for the first time or reset. Each performance variable is associated with a class that describes its basic semantics, basic behavior, its starting value, and when and how an MPI implementation can changes its value. The starting value is the value the variable assumes when it is used for the first time or whenever it is reset. In the following this referred to as the starting value. Further, [it]the class of the variable also defines which datatypes can

⊥ (Fin2)	1	be used to represent it.]	These classes are defined by the following constants:
	2		
	3		
	4		
⊤ (Fin2)	5		
⊥ (Fin2)	6		
⊤ (Fin2)	7		
⊥ (Fin2)	8		
⊤ (Fin2)	9		
⊥ (Fin2)	10		
	11		
	12		
	13		
	14		
⊤ (Fin2)	15		
⊥ (Fin2)	16		
⊤ (Fin2)	17		
⊥ (Fin2)	18		
	19		
	20		
	21		
	22		
⊤ (Fin2)	23		
⊥ (Fin2)	24		
⊤ (Fin2)	25		
⊥ (Fin2)	26		
	27		
	28		
	29		
	30		
	31		
⊤ (Fin2)	32		
⊥ (Fin2)	33		
⊤ (Fin2)	34		
⊥ (Fin2)	35		
	36		
	37		
	38		
	39		
⊤ (Fin2)	40		
⊥ (Fin2)	41		
⊤ (Fin2)	42		
⊥ (Fin2)	43		
⊤ (Fin2)	44		
⊥ (Fin2)	45		
	46		
	47		
⊤ (Fin2)	48		
⊥ (Fin2)			

- **MPI\_T\_PVAR\_CLASS\_STATE**

A performance variable in this class represents a set of discrete states. Variables of this class [are expected to be] are represented by an enumeration datatype and can be set by the MPI implementation at any time. The [default] starting value is the current state [ (at the time the starting value is set) of the implementation. ] of the implementation at the time the starting value is set .

- **MPI\_T\_PVAR\_CLASS\_LEVEL**

A performance variable in this class represents a value that describes the utilization level of a resource. The value of a variable of this class can change at any time to match the current utilization level of the resource. Values returned from variables in this class are represented by one of the following datatypes: **MPI\_T\_INT**, **MPI\_T\_LONG\_LONG**, **MPI\_T\_DOUBLE**. The [default] starting value is the current utilization level [ (at the time the starting value is set) of the resource. ] of the resource at the time the starting value is set.

- **MPI\_T\_PVAR\_CLASS\_PERCENTAGE**

The value of a performance variable in this class represents the percentage utilization of a finite resource. The value of a variable of this class can change at any time to match the current utilization level of the resource. It will be returned as an **MPI\_T\_DOUBLE** datatype. The value must always be between 0.0 (resource not used at all) and 1.0 (resource completely used). The [default] starting value is the current percentage utilization level [ (at the time the starting value is set) of the resource. ] of the resource at the time the starting value is set.

- **MPI\_T\_PVAR\_CLASS\_HIGHWATERMARK**

A performance variable in this class represents a value that describes the high watermark utilization of a resource. The value of a variable of this class grows monotonically from the initialization or reset of the variable. It can be represented by one of the following datatypes: **MPI\_T\_INT**, **MPI\_T\_LONG\_LONG**, **MPI\_T\_DOUBLE**. The [default] starting value is the current utilization level [ (at the time the starting value is set) of the resource. ] of the resource at the time the starting value is set.

- **MPI\_T\_PVAR\_CLASS\_LOWWATERMARK**

A performance variable in this class represents a value that describes the low watermark utilization of a resource. The value of a variable of this class decreases monotonically from the initialization or reset of the variable. It can be represented by one of the following datatypes: **MPI\_T\_INT**, **MPI\_T\_LONG\_LONG**, **MPI\_T\_DOUBLE**. The [default] starting value is the current utilization level [ (at the time the starting value is set) of the resource. ] of the resource at the time the starting value is set.

- **MPI\_T\_PVAR\_CLASS\_COUNTER**

A performance variable in this class counts the number of occurrences of a specific event [during the execution time of an application] (e.g., the number of memory allocations within an MPI library). The value of a variable of this class increases monotonically from the initialization or reset of the performance variable by one for each specific event that is observed. Values must be non-negative and represented by one of the following datatypes: **MPI\_T\_INT**, **MPI\_T\_LONG\_LONG**. The [default]

starting value for variables of this class is 0.

- **MPI\_T\_PVAR\_CLASS\_AGGREGATE**

The value of a performance variable in this class is an aggregated value that represents a sum of arguments processed during a specific event (e.g., the amount of memory allocated by all memory allocations). This class is similar to the counter class, but instead of counting individual events, the value can be incremented by arbitrary amounts. The value of a variable of this class **increases** monotonically from the initialization or reset of the performance variable. It must be non-negative and represented by one of the following datatypes: **MPI\_T\_INT**, **MPI\_T\_LONG\_LONG**, **MPI\_T\_DOUBLE**. The **[default]** starting value for variables of this class is 0.

- **MPI\_T\_PVAR\_CLASS\_TIMER**

The value of a performance variable in this class represents the aggregated time that the MPI implementation spends executing a particular event **or type of event**. This class has the same basic semantics as **MPI\_T\_PVAR\_CLASS\_AGGREGATE**, but explicitly records a timing value. The value of a variable of this class **increases** monotonically from the initialization or reset of the performance variable. It must be non-negative and represented by one of the following datatypes: **MPI\_T\_INT**, **MPI\_T\_LONG\_LONG**, **MPI\_T\_DOUBLE**. The **[default]** starting value for variables **[if]** of this class is 0. **If the type MPI\_T\_DOUBLE is used, the units representing time in this datatype must match the units used by MPI\_WTIME.**

### Performance Variable Query Functions

An MPI implementation exports a set of  $N$  performance variables through **MPI\_T**. If  $N$  is zero, then the **MPI\_T** implementation does not export any performance variables, otherwise the provided performance variables are indexed from 0 to  $N - 1$ . This index number is used in subsequent **MPI\_T** calls to identify the individual variables.

An **MPI\_T** implementation is allowed to increase the number of performance variables during the execution of an MPI application when new variables become available through dynamic loading. However, **MPI\_T** implementations are not allowed to change the index of a performance variable or delete a variable once it has been added to the set.

The following function can be used to query the number of performance variables,  $N$ :

**MPI\_T\_PVAR\_GET\_NUM(num)**

OUT      num      returns number of performance variables

**int MPI\_T\_Pvar\_get\_num(int \*num)**

The function **MPI\_T\_PVAR\_GET\_INFO** provides access to additional information for each variable.

```
1 MPI_T_PVAR_GET_INFO(index, name, name_len, verbosity, varclass, datatype, count, desc,
2 desc_len, bind, attributes)
```

3	IN	index	index of the performance variable to be queried
4	OUT	name	buffer to return the string containing the name of the performance variable
5			
6			
7	INOUT	name_len	length of the string and/or buffer for name
8	OUT	verbosity	verbosity level of this variable
9			
10	OUT	var_class	class of performance variable
11	OUT	datatype	MPI_T datatype of the information stored in the performance variable
12			
13	OUT	count	number of elements returned
14	OUT	desc	buffer to return the string containing a description of the performance variable
15			
16			
17	INOUT	desc_len	length of the string and/or buffer for desc
18	OUT	bind	type of MPI object to which this variable must be bound
19			
20			
21	OUT	attributes	additional attributes defining this variable

```
22
23 int MPI_T_Pvar_get_info(int num, char *name, int *name_len, int *verbosity,
24 int *var_class, MPI_T_Datatype *datatype, int *count, char
25 *desc, int *desc_len, int *bind, MPI_T_Pvar_attributes
26 *attributes)
```

After a successful call to `MPI_T_PVAR_GET_INFO` for a particular variable, subsequent calls to this routine querying information about the same variable must return the same information. An `MPI_T` implementation is not allowed to alter any of the returned values.

The arguments `name` and `name_len` are used to return the name of the performance variable as described in Section 1.3.3.

If completed successfully, the routine is required to return a name of at least length one. This name must be unique with respect to all other names for `MPI_T` performance variables used by the MPI implementation.

The argument `verbosity` returns the verbosity level [ (see Section 1.3.1) assigned by the MPI implementation to the variable. ] of the variable (see Section 1.3.1).

The class of the performance variable is returned in the parameter `var_class` [ and can ]. The class must be one of the constants defined in Section 1.3.7.

The argument `datatype` returns the `MPI_T` datatype [ in which the value for this performance variable is returned. ] that is used to represent the performance variable. The value consists of `count` elements of this datatype.

The arguments `desc` and `desc_len` are used to return a description of the [control]performance variable as described in Section 1.3.3.

Returning a description is optional. If an MPI implementation decides not to return a description, the first character for `desc` must be set to the null character and `desc_len` must be set to one at the return from this function.



The parameter `bind` returns the type of the MPI object to which the variable must be bound or the value `MPI_T_BIND_NO_OBJECT` (see Section 1.3.1).

Additional information about the variable is returned through the `attributes` argument using an opaque structure of type `MPI_T_Pvar_attributes` and can be queried using the following accessor functions.

#### `MPI_T_PVAR_ATTR_GET_READONLY(attributes, readonly)`

IN	<code>attributes</code>	attributes returned by a previous query call
OUT	<code>readonly</code>	flag indicating whether a variable can be written/reset

```
int MPI_T_Pvar_attr_get_readonly(MPI_T_Pvar_attributes attributes, int
                                *readonly)
```

Upon return, the argument `readonly` [will be] is set to zero if the variable can be written or reset by the user[, or]. It is set to one if the variable can only be read.

#### `MPI_T_PVAR_ATTR_GET_CONTINUOUS(attributes, continuous)`

IN	<code>attributes</code>	attributes returned by a previous query call
OUT	<code>continuous</code>	flag indicating whether a variable can be started and stopped or is continuously active

```
int MPI_T_Pvar_attr_get_continuous(MPI_T_Pvar_attributes attributes, int
                                   *continuous)
```

Upon return, the argument `continuous` [will be] is set to zero if the variable can be started and stopped by the user, i.e, it is possible for the user to control if and when the value of a variable is updated [, or]. It is set to one if the variable is [automatically] always active and [can not by] cannot be controlled by the user.

#### Performance Experiment Sessions

Within a single program, multiple components can use the `MPI_T` interface. To avoid collisions with respect to accesses to performance variables, users of the `MPI_T` interface must first create a session. All subsequent calls accessing performance variables are then within the context of this session. Any call executed in a session must not influence the results in any other session.

#### `MPI_T_PVAR_SESSION_CREATE(session)`

OUT	<code>session</code>	identifier of performance session
-----	----------------------	-----------------------------------

```
int MPI_T_Pvar_session_create(MPI_T_Pvar_session *session)
```

This call creates a new session for accessing performance variables and returns an identifier for this session in the argument `session`.

1 `MPI_T_PVAR_SESSION_FREE(session)`

2     INOUT     session                             identifier of performance experiment session

4  
5 `int MPI_T_Pvar_session_free(MPI_T_Pvar_session *session)`

6     This call frees an existing session[, i.e., calls]Calls to `MPI_T` can no longer be made  
7 within the [ context of the MPI freed session. ]context of a session after it is freed. This call  
8 also frees all handles that have been allocated within the specified session ( see below for  
9 handle allocation and freeing). On a successful return, `MPI_T` sets the session identifier to  
10 `MPI_T_PVAR_SESSION_NULL`.

## 12 Handle Allocation and Deallocation

13  
14 Before using a performance variable, a user must first allocate a handle for it by binding  
15 it to an MPI object (see also Section 1.3.1). [ The type of the MPI object is returned by a  
16 previous call to `MPI_T_PVAR_GET_INFO` in the bind argument. ]

18  
19 `MPI_T_PVAR_HANDLE_ALLOC(session, index, objhandle, handle)`

20     IN         session                             identifier of performance experiment session

21     IN         index                               index of performance variable for which handle is to  
22   be allocated

23     IN         obj\_handle                         reference to a handle of the MPI object to which this  
24   variable is supposed to be bound

25  
26     OUT        handle                             allocated handle

27  
28 `int MPI_T_Pvar_handle_alloc(MPI_T_Pvar_session session, int index, void`  
29 `*obj_handle, MPI_T_Pvar_handle *handle)`

30  
31     [ A call to this routine allocates a handle for the performance variable specified by  
32 the argument `index` and binds this variable to the MPI object referenced by the pointer to  
33 its handle passed in the argument `obj_handle` . The type of the MPI object passed into  
34 this routine must match the type of the MPI object for this variable as returned by a prior  
35 call to `MPI_T_PVAR_GET_INFO`. If the type of the object is identified as  
36 `MPI_T_BIND_NO_OBJECT`, i.e., the variable refers to the entire MPI implementation the  
37 argument object is ignored. ] This routine binds the performance variable specified  
38 by the argument `index` to the MPI object referenced by the handle passed in argument  
39 `obj_handle` and returns an allocated variable handle in the argument `handle`. The value of  
40 `index` should be in the range 0 to  $N - 1$ , where  $N$  is the number of available control variables  
41 as determined from a prior call to `MPI_T_PVAR_GET_NUM`. The value of `obj_handle` must  
42 be the memory address of the object's MPI handle, and the type of MPI object it references  
43 must be consistent with the type returned in the bind argument in a prior call to  
44 `MPI_T_PVAR_GET_INFO`.

MPI\_T\_PVAR\_HANDLE\_FREE(session, handle)

IN session identifier of performance experiment session  
 INOUT handle handle to be freed

```
int MPI_T_Pvar_handle_free(MPI_T_Pvar_session session, MPI_T_Pvar_handle
                          *handle)
```

When a handle is no longer needed, a user of MPI\_T should call MPI\_T\_PVAR\_HANDLE\_FREE to free the handle and the associated resources in the MPI\_T implementation. On a successful return, MPI\_T sets the handle to MPI\_T\_PVAR\_HANDLE\_NULL.

### Starting and Stopping of Performance Variables

[ Performance variables that have the continuous flag set during the query operation are continuously operating once a handle has been allocated and can be queried any time. They cannot be stopped or paused by the user. All other variables are in a stopped state after their handle has been allocated; their values are not updated as the program executes, and must be started by the user. ] Performance variables that have the continuous flag set during the query operation are continuously operating once a handle has been allocated. Such variables may be queried at any time, but they cannot be stopped or paused by the user. All other variables are in a stopped state after their handle has been allocated; their values are not updated until they have been started by the user.

MPI\_T\_PVAR\_START(session, handle)

IN session identifier of performance experiment session  
 IN handle handle of a performance variable

```
int MPI_T_Pvar_start(MPI_T_Pvar_session session, MPI_T_Pvar_handle handle)
```

This functions starts the performance variable with the handle handle in the session session.

If the constant MPI\_T\_PVAR\_ALL\_HANDLES is passed in handle, the MPI implementation attempts to start all variables within the session identified by the parameter session for which handles have been allocated. In this case, the routine returns MPI\_T\_SUCCESS if all variables are started successfully, otherwise MPI\_T\_ERR\_NOSTARTSTOP is returned. Continuous variables and variables that are already started are ignored when [ used with MPI\_T\_PVAR\_ALL\_HANDLES. ] MPI\_T\_PVAR\_ALL\_HANDLES is specified.

MPI\_T\_PVAR\_STOP(session, handle)

IN session identifier of performance experiment session  
 IN handle handle of a performance variable

```
int MPI_T_Pvar_stop(MPI_T_Pvar_session session, MPI_T_Pvar_handle handle)
```

This function stops the performance variable with the handle `handle` in the session `session`.

If the constant `MPI_T_PVAR_ALL_HANDLES` is passed in `handle`, the MPI implementation attempts to stop all variables within the session identified by the parameter `session` for which handles have been allocated. In this case, the routine returns `MPI_SUCCESS` if all variables are stopped successfully, otherwise `MPI_T_ERR_NOSTARTSTOP` is returned. Continuous variables and variables that are already stopped are ignored when used with `MPI_T_PVAR_ALL_HANDLES`. `MPI_T_PVAR_ALL_HANDLES` is specified.

## Performance Variable Access Functions

### `MPI_T_PVAR_READ(session, handle, buf)`

IN	<code>session</code>	identifier of performance experiment session
IN	<code>handle</code>	handle of a performance variable
OUT	<code>buf</code>	initial address of storage location for variable value

```
int MPI_T_Pvar_read(MPI_T_Pvar_session session, MPI_T_Pvar_handle handle,
                    void* buf)
```

The `MPI_T_PVAR_READ` call queries the value of the performance variable with the handle `handle` in the session identified by the parameter `session` and stores the result in the buffer `buf`. The user is responsible to ensure that the buffer is of the appropriate size [and fits] to hold the entire value of the performance variable (based on the returned datatype and count during the `MPI_T_PVAR_GET_INFO` call).

[Note that the constant `MPI_T_PVAR_ALL_HANDLES` cannot be used as an argument for the MPI\_T function `MPI_T_PVAR_READ`, since this would require the function to return a set of variable values instead of just one. ] The constant `MPI_T_PVAR_ALL_HANDLES` cannot be used as an argument for the MPI\_T function `MPI_T_PVAR_READ`.

### `MPI_T_PVAR_WRITE(session, handle, buf)`

IN	<code>session</code>	identifier of performance experiment session
IN	<code>handle</code>	handle of a performance variable
IN	<code>buf</code>	initial address of storage location for variable value

```
int MPI_T_Pvar_write(MPI_T_Pvar_session session, MPI_T_Pvar_handle handle,
                     void* buf)
```

The `MPI_T_PVAR_WRITE` call attempts to write the value of the performance variable with the handle `handle` in the session identified by the parameter `session`. The value to be written is passed in the buffer `buf`. The user is responsible to ensure that the buffer is of the appropriate size [and fits] to hold the entire value of the performance variable (based on the returned datatype and count during the `MPI_T_PVAR_GET_INFO` call).

If it is not possible to change the variable, the function returns `MPI_T_ERR_PVAR_WRITE`.

[ Note that the constant `MPI_T_PVAR_ALL_HANDLES` cannot be used as an argument for the `MPI_T` function `MPI_T_PVAR_WRITE`, since this would require the function to accept a set of variable values instead of just one. ] The constant `MPI_T_PVAR_ALL_HANDLES` cannot be used as an argument for the `MPI_T` function `MPI_T_PVAR_WRITE`.

`MPI_T_PVAR_RESET(session, handle)`

IN	session	identifier of performance experiment session
IN	handle	handle of a performance variable

`int MPI_T_Pvar_reset(MPI_T_Pvar_session session, MPI_T_Pvar_handle handle)`

The `MPI_T_PVAR_RESET` call sets [of] the performance variable with the handle `handle` to its starting value specified in Section 1.3.7. If it is not possible to change the variable, the function returns `MPI_T_ERR_PVAR_WRITE`.

If the constant `MPI_T_PVAR_ALL_HANDLES` is passed in `handle`, the MPI implementation attempts to reset all variables within the session identified by the parameter `session` for which handles have been allocated. In this case, the routine returns `MPI_T_SUCCESS` if all variables are reset successfully, otherwise `MPI_T_ERR_NOWRITE` is returned. Readonly variables are ignored when [ used with `MPI_T_PVAR_ALL_HANDLES`. ] `MPI_T_PVAR_ALL_HANDLES` is specified.

`MPI_T_PVAR_READRESET(session, handle, buf)`

IN	session	identifier of performance experiment session
IN	handle	handle of a performance variable
OUT	buf	initial address of storage location for variable value

`int MPI_T_Pvar_readreset(MPI_T_Pvar_session session, MPI_T_Pvar_handle handle, void* buf)`

This call combines the functionality of `MPI_T_PVAR_READ` and `MPI_T_PVAR_RESET` with the same semantics as if these two calls were called separately.

[ Note that the constant `MPI_T_PVAR_ALL_HANDLES` can not be used as an argument for the `MPI_T` function `MPI_T_PVAR_READRESET`, since this would require the function to return a set of variable values instead of just one. ] The constant `MPI_T_PVAR_ALL_HANDLES` can not be used as an argument for the `MPI_T` function `MPI_T_PVAR_READRESET`.

*Advice to implementors.* Although MPI places no requirements on the interaction with external mechanisms such as signal handlers, it is strongly recommended that all routines to start, stop, read, write, and reset performance variables should be safe to call in asynchronous contexts. Examples of asynchronous contexts include signal handlers and interrupt handlers. Such safety permits the development of sampling-based tools. High quality implementations should strive to make the results of any such interactions intuitive to users, and document known restrictions. (*End of advice to implementors.*)

### 1.3.8 Variable Categorization

MPI implementations can optionally group performance and control variables into categories to express logical relationships between various variables. For example, an `MPI_T` implementation could group all control and performance variables that refer to message transfers in the MPI implementation and thereby [distinguishing] distinguish them from variables that refer to local resources such as memory allocations or other interactions with the operating system.

Categories can also contain other categories to form a hierarchical grouping. Categories can never include themselves, either directly or transitively within other included categories.

Expanding on the example above, this allows `MPI_T` to refine the grouping of variables referring to message transfers into variables to control and monitor message queues, message matching activities and communication protocols. Each of these groups of variables would be represented by a separate category and these categories would then be listed in a single category representing variables for message transfers.

The category information may be queried in a fashion similar to the mechanism for querying variable information. The MPI implementation exports a set of  $N$  categories via the `MPI_T` interface. If  $N = 0$ , then the MPI implementation does not export any categories, otherwise the provided performance variables are indexed from 0 to  $N - 1$ . This index number is used in subsequent `MPI_T` calls to identify the individual categories.

An MPI implementation is permitted to increase the number of categories during the execution of an MPI program when new categories become available through dynamic loading. However, MPI implementations are not allowed to change the index of a category or delete it once it has been added to the set.

The following function can be used to query the number of control variables,  $N$ .

```
MPI_T_CATEGORY_GET_NUM(num)
```

```
OUT      num                current number of categories
```

```
int MPI_T_Category_get_num(int *num)
```

Individual category information can then be queried by calling the following function:

`MPI_T_CATEGORY_GET_INFO(index, name, name_len, desc, desc_len, num_controlvars, num_perfvars, num_categories)`

IN	index	index of the category to be queried
OUT	name	buffer to return the string containing the name of the category
INOUT	name_len	length of the string and/or buffer for name
OUT	desc	buffer to return the string containing the description of the category
INOUT	desc_len	length of the string and/or buffer for desc
OUT	num_controlvars	number of control variables in the category
OUT	num_perfvars	number of performance variables in the category
OUT	num_categories	number of MPI_T categories contained in the category

```
int MPI_T_Category_get_info(int index, char *name, int *name_len, char
                           *desc, int *desc_len, int *num_controlvars, int
                           *num_perfvars, int *num_categories)
```

The arguments `name` and `name_len` are used to return the name of the category as described in Section 1.3.3.

The routine is required to return a name of at least length one. This name must be unique with respect to all other names for MPI\_T categories used by the MPI\_T implementation.

The arguments `desc` and `desc_len` are used to return the description of the category as described in Section 1.3.3.

Returning a description is optional. If an MPI implementation decides not to return a description, the first character for `desc` must be set to the null character and `desc_len` must be set to one at the return of this call.

The function returns the number of control variables, performance variables and other categories contained in the queried category in the arguments `num_controlvars`, `num_perfvars`, and `num_categories` respectively.

*Advice to implementors.* To avoid confusion and to simplify the interpretation of the categories provided by a particular implementation, it is recommended that categories should either only contain other categories or only control and performance variables. Mixing categories and control and performance variables within a single category is not recommended. *(End of advice to implementors.)*

```
1 MPI_T_CATEGORY_GET_CVARS(cat_index, len, indices)
```

```
2     IN      cat_index      index of the category to be queried, in the range [0, N-
3                               1]
4
5     IN      len            the length of the indices array
6     OUT     indices        an integer array of size len, indicating control variable
7                               indices
8
```

```
9 int MPI_T_Category_get_cvars(int cat_index, int len, int indices[])
10
```

11 `MPI_T_CATEGORY_GET_CVARS` can be used to query which control variables are  
 12 contained in a particular category. A category contains zero or more control variables.

```
14 MPI_T_CATEGORY_GET_PVARS(cat_index, len, indices)
```

```
16     IN      cat_index      index of the category to be queried, in the range [0, N-
17                               1]
18
19     IN      len            the length of the indices array
20     OUT     indices        an integer array of size len, indicating performance
21                               variable indices
22
```

```
23 int MPI_T_Category_get_pvars(int cat_index, int len, int indices[])
24
```

25 `MPI_T_CATEGORY_GET_PVARS` can be used to query which performance variables  
 26 are contained in a particular category. A category contains zero or more performance  
 27 variables.

```
29 MPI_T_CATEGORY_GET_CATEGORIES(cat_index, len, indices)
```

```
30     IN      cat_index      index of the category to be queried, in the range [0, N-
31                               1]
32
33     IN      len            the length of the indices array
34     OUT     indices        an integer array of size len, indicating category indices
35
```

```
36 int MPI_T_Category_get_categories(int cat_index, int len, int indices[])
37
```

38 `MPI_T_CATEGORY_GET_CATEGORIES` can be used to query which other categories  
 39 are contained in a particular category. A category contains zero or more other categories.

40 The index values returned in `indices` by `MPI_T_CATEGORY_GET_CVARS`,  
 41 `MPI_T_CATEGORY_GET_PVARS` and `MPI_T_CATEGORY_GET_CATEGORIES` can be used  
 42 as input to `MPI_T_CVAR_GET_INFO`, `MPI_T_PVAR_GET_INFO` and  
 43 `MPI_T_CATEGORY_GET_INFO` respectively.

44 The user is responsible for allocating the arrays passed into the functions  
 45 `MPI_T_CATEGORY_GET_CVARS`, `MPI_T_CATEGORY_GET_PVARS` and  
 46 `MPI_T_CATEGORY_GET_CATEGORIES`. [ The functions will only write up to `len` elements  
 47 into the respective array. If the category contains more than `len` variables or other categories  
 48 respectively, the function returns an arbitrary subset; if it contains less than `len` variables



or other categories respectively, all will be returned and the remaining array entries will not be modified. ] Starting from array index 0, each function writes up to len elements into the array. If the category contains more than len elements, the function returns an arbitrary subset of size len. Otherwise, the entire set of elements is returned in the beginning entries of the array, and any remaining array entries are not modified.”

⊥ (Fin2)

### 1.3.9 Return and Error Codes

All MPI\_T functions return an integer error code (See Table 1.3.9). [ None of the error codes returned by an MPI\_T routine are fatal to the MPI process or invoke an MPI error handler. ] errors returned by MPI\_T routines are not fatal and do not invoke MPI error handlers. The execution of the MPI process continues as if the MPI\_T call would have succeeded. However, the MPI\_T implementation is not required to check all user provided parameters; if a user passes invalid parameter values to any MPI\_T routine the behavior of the implementation is undefined.

### 1.3.10 Profiling Interface

All requirements for the profiling interfaces, as described in Section 1.2, also apply to the MPI\_T interface. In particular, this means that [a complying]compliant MPI\_T implementation[s ] must provide matching PMPI\_T calls for every MPI\_T call. All rules, guidelines, and recommendations from Section 1.2 apply equally to PMPI\_T calls.

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

⊥ (Fin2)

Return Code	Description
Return Codes for all <code>MPI_T</code> Functions	
<code>MPI_T_SUCCESS</code>	No error, call completed
<code>MPI_T_ERR_MEMORY</code>	Out of memory
<code>MPI_T_ERR_NOTINITIALIZED</code>	<code>MPI_T</code> not initialized
<code>MPI_T_ERR_CANTINIT</code>	<code>MPI_T</code> not in the state to be initialized
Return Codes for Datatype Functions: <code>MPI_T_DATATYPE_*</code>	
<code>MPI_T_ERR_PREDEFINED</code>	Datatype is a predefined datatype and not an enumeration
<code>MPI_T_ERR_INVALIDDATATYPE</code>	Datatype is not a valid datatype
<code>MPI_T_ERR_INVALIDITEM</code>	The item index queried is out of range (for <code>MPI_T_DATATYPE_ENUMITEM</code> only)
Return Codes for variable and category query functions: <code>MPI_T_*.GET_INFO</code>	
<code>MPI_T_ERR_INVALIDINDEX</code>	The variable or category index is invalid
Return Codes for Handle Functions: <code>MPI_T_*.ALLOCATE,FREE</code>	
<code>MPI_T_ERR_INVALIDINDEX</code>	The variable index is invalid
<code>MPI_T_ERR_INVALIDHANDLE</code>	The handle is invalid
<code>MPI_T_ERR_OUTOFHANDLES</code>	No more handles available
Return Codes for Session Functions: <code>MPI_T_PVAR_SESSION_*</code>	
<code>MPI_T_ERR_OUTOFSESSIONS</code>	No more sessions available
<code>MPI_T_ERR_INVALIDSESSION</code>	Session argument is not a valid session
Return Codes for Control Variable Access Functions: <code>MPI_T_CVAR_READ, WRITE</code>	
<code>MPI_T_ERR_SETNOTNOW</code>	Variable cannot be set at this moment
<code>MPI_T_ERR_SETNEVER</code>	Variable cannot be set until end of execution
<code>MPI_T_ERR_INVALIDVAR</code>	Control variable does not exist
<code>MPI_T_ERR_INVALIDHANDLE</code>	The handle is invalid
Return Codes for Performance Variable Access and Control: <code>MPI_T_PVAR_START, STOP, READ, WRITE, RESET, READRESET</code>	
<code>MPI_T_ERR_INVALIDHANDLE</code>	The handle is invalid
<code>MPI_T_ERR_INVALIDSESSION</code>	Session argument is not a valid session
<code>MPI_T_ERR_NOSTARTSTOP</code>	Variable can not be started or stopped for <code>MPI_T_PVAR_START</code> and <code>MPI_T_PVAR_STOP</code>
<code>MPI_T_ERR_NOWRITE</code>	Variable can not be written or reset for <code>MPI_T_PVAR_WRITE</code> and <code>MPI_T_PVAR_RESET</code>
Return Codes for Category Functions: <code>MPI_T_CATEGORY_*</code>	
<code>MPI_T_ERR_INVALIDCATEGORY</code>	The specified category index does not exist

Table 1.6: Return and error codes used `MPI_T` functions.

# Bibliography

- [1] mpi-debug: Finding Processes. <http://www-unix.mcs.anl.gov/mpi/mpi-debug/>.
- [2] James Cownie and William Gropp. A Standard Interface for Debugger Access to Message Queue Information in MPI. In *Proceedings of the 6th European PVM/MPI Users' Group Meeting on Recent Advances in Parallel Virtual Machine and Message Passing Interface*, pages 51–58, Barcelona, Spain, September 1999.

# Examples Index

This index lists code examples throughout the text. Some examples are referred to by content; others are listed by the major [MPI\\_T](#) function that they are demonstrating. [MPI\\_T](#) functions listed in all capital letter are Fortran examples; [MPI\\_T](#) functions listed in mixed case are C/C++ examples.

Profiling interface, [4](#)