

D R A F T

Document for a Standard Message-Passing Interface

Message Passing Interface Forum

July 4, 2012

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

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

1 ticket266.
2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

Chapter 14

Tool Support

14.1 Introduction

This chapter discusses [a set of interfaces]interfaces that [allows]allow debuggers, performance analyzers, and other tools to extract information about the operation of MPI processes. Specifically, this chapter defines both the MPI profiling interface (Section 14.2), which supports the transparent interception and inspection of MPI calls, and the MPI tool information interface (Section 14.3), which supports the inspection and manipulation of 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.

14.2 Profiling Interface

[WAS: Chapter]

14.2.1 Requirements

[WAS: Section]

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

1. provide a mechanism through which all of the MPI defined [functions]functions, except those allowed as macros (See Section 2.6.5[)], 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 in each provided language binding and language support method. The profiling interface in C++ is described in Section 16.1.10. 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.

For Fortran, the different support methods cause several linker names. Therefore, several profiling routines (with these linker names) are needed for each Fortran MPI routine, as described in Section 16.2.5 on page 647.

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.

14.2.2 Discussion

[\[WAS: Section \]](#)

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 [\[chapter\]section](#) is only present as justification and discussion of the logic for those requirements.

1 The examples below show one way in which an implementation could be constructed to
 2 meet the requirements on a Unix system (there are doubtless others that would be equally
 ticket266. 3 valid).

4 14.2.3 Logic of the Design

5 [WAS: Section]

6 Provided that an MPI implementation meets the requirements above, it is possible for
 ticket0-new. 7 the implementor of the profiling system to intercept [all of the]the MPI calls that are made
 8 by the user program. She can then collect whatever information she requires before calling
 9 the underlying MPI implementation (through its name shifted entry points) to achieve the
 ticket266. 10 desired effects.

11 14.2.4 Miscellaneous Control of Profiling

12 [WAS: Subsection, Now still a subsection to remove single subsection]

13 There is a clear requirement for the user code to be able to control the profiler dy-
 ticket0-new. 14 namically at run time. [This]This capability is normally used for (at least) the purposes
 15 of

- 16 • Enabling and disabling profiling depending on the state of the calculation.
- 17 ticket0. 18 • Flushing trace buffers at non-critical points in the [calculation]calculation.
- 19 • Adding user events to a trace file.

20 These requirements are met by use of [the MPI_PCONTROL]MPI_PCONTROL.

21 MPI_PCONTROL(level, ...)

22 IN level Profiling level

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

24 ticket-248T. 25 MPI_Pcontrol(level) BIND(C)
 26 INTEGER, INTENT(IN) :: level

27 MPI_PCONTROL(LEVEL)
 28 INTEGER LEVEL

29 {void MPI::Pcontrol(const int level, ...) (binding deprecated, see Section 15.2) }

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

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

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

- 38 • 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.

[WAS: Subsection Examples]

14.2.5 Profiler Implementation [Example]

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

Example 14.1

```
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 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;
}
```

14.2.6 MPI Library Implementation [Examples]

[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

1 name-shift requirement. The choice between these two options [compiler support weak
2 symbols.] depends partly on whether the linker and compiler support weak symbols. ticket0.

3 4 Systems with Weak Symbols

5 If the compiler and linker support weak external symbols ([e.g.]e.g., Solaris 2.x, other system
6 V.4 machines), then only a single library is required [through the use of #pragma weak
7 thus]as the following example shows: ticket0-new.

8 9 **Example 14.2**

```
10 #pragma weak MPI_Example = PMPI_Example
11
12
13 int PMPI_Example(/* appropriate args */)
14 {
15     /* Useful content */
16 }
```

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

21 22 23 Systems Without Weak Symbols

24 In the absence of weak symbols then one possible solution would be to use the C macro
25 pre-processor [thus]as the following example shows: ticket0-new.

26 27 **Example 14.3**

```
28 #ifdef PROFILELIB
29 #   ifdef __STDC__
30 #       define FUNCTION(name) P##name
31 #   else
32 #       define FUNCTION(name) P/**/name
33 #   endif
34 #else
35 #   define FUNCTION(name) name
36 #endif
37
```

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

```
39
40 int FUNCTION(MPI_Example)(/* appropriate args */)
41 {
42     /* Useful content */
43 }
```

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

46 It is required that the standard MPI library be built in such a way that the inclusion of
47 MPI functions can be achieved one at a time. This is a somewhat unpleasant requirement,
48

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 -lmpi -lmpi
```

Here `libmyprof.a` contains the profiler functions that intercept some of the MPI functions`[.]`, `libmpmi.a` contains the “name shifted” MPI functions, and `libmpi.a` contains the normal definitions of the MPI functions.

14.2.7 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.]it. 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 [!]).

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 Fortran Support Methods

2 The different Fortran support methods and possible options for the support of subarrays
 3 (depending on whether the compiler can support `TYPE(*)`, `DIMENSION(...)` choice buffers)
 4 imply different linker names for the same Fortran MPI routine. The rules and implications
 5 for the profiling interface are described in Section 16.2.5 on page 647.

7 14.2.8 Multiple Levels of Interception

8 [WAS: Section] The scheme given here does not directly support the nesting of profiling
 9 functions, since it provides only a single alternative name for each MPI function. Consider-
 10 ation was given to an implementation that would allow multiple levels of call interception,
 11 however we were unable to construct an implementation of this that did not have the fol-
 12 lowing disadvantages

- 13
- 14
- 15 ticket0. • assuming a particular implementation language[.],
- 16 • imposing a run time cost even when no profiling was taking place.
- 17

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

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

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

31 14.3 The MPI Tool Information Interface

32 MPI implementations often use internal variables to control their operation and perfor-
 33 mance. Understanding and manipulating these variables can provide a more efficient exe-
 34 cution environment or improve performance for many applications. This section describes
 35 the MPI tool information interface, which provides a mechanism for MPI implementors to
 36 expose [a set of variables]variables, each of which represents a particular property, setting,
 37 or performance measurement from within the MPI implementation. The interface is split
 38 into two parts: the first part provides information about and supports the setting of control
 39 variables through which the MPI implementation tunes its configuration. The second part
 40 provides access to performance variables that can provide insight into internal performance
 41 information of the MPI implementation.

42 To avoid restrictions on the MPI implementation, the MPI tool information interface
 43 allows the implementation to specify which control and performance variables exist. Ad-
 44 ditionally, the user of the MPI tool information interface can obtain metadata about each
 45 available variable, such as its datatype, and a textual description. The MPI tool information
 46 interface provides the necessary routines to find all variables that exist in a particular MPI
 47
 48

implementation, to query their properties, to retrieve descriptions about their meaning, and to access and, if appropriate, to alter their values.

The MPI tool information interface can be used independently from the MPI communication functionality. In particular, the routines of this interface can be called before MPI_INIT (or equivalent) and after MPI_FINALIZE. In order to support this behavior cleanly, the MPI tool information interface uses separate initialization and finalization routines. All identifiers used in the MPI tool information interface have the prefix MPI_T_.

On success, all MPI tool information interface routines return MPI_SUCCESS, otherwise they return an appropriate and unique return code indicating the reason why the call was not successfully completed. Details on return codes can be found in Section 14.3.9. However, unsuccessful calls to the MPI tool information interface are not fatal and do not impact the execution of subsequent MPI routines.

Since the MPI tool information interface primarily focuses on tools and support libraries, MPI implementations are only required to provide C bindings for functions introduced in this [Section 14.3]section. Except where otherwise noted, all conventions and principles governing the C bindings of the MPI API also apply to the MPI tool information interface, which is available by including the mpi.h header file. All routines in this interface have local semantics.

Advice to users. The number and type of control variables and performance variables can vary between MPI implementations, platforms and different builds of the same implementation on the same platform as well as between runs. Hence, any application relying on a particular variable will not be portable. Further, there is no guarantee that number of variables, variable indices, and variable names are the same across processes.

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

14.3.1 Verbosity Levels

The MPI tool information interface provides access to internal configuration and performance information through a set of control and performance variables defined by the MPI 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 level of detail (basic, detailed or all). These verbosity levels are described by a single integer. Table 14.1 lists the constants for all possible verbosity levels. The values of the constants are monotonic in the order listed in the table; i.e., MPI_T_VERBOSITY_USER_BASIC < MPI_T_VERBOSITY_USER_DETAIL < ... < MPI_T_VERBOSITY_MPIDEV_ALL.

14.3.2 Binding MPI Tool Information Interface Variables to MPI Objects

Each MPI tool information interface variable provides access to a particular control setting or performance property of the MPI implementation. A variable may refer to a specific MPI object such as a communicator, datatype, or one-sided communication window, or the

MPI_T_VERBOSITY_USER_BASIC	Basic information of interest to users
MPI_T_VERBOSITY_USER_DETAIL	Detailed information of interest to users
MPI_T_VERBOSITY_USER_ALL	All information of interest to users
MPI_T_VERBOSITY_TUNER_BASIC	Basic information required for tuning
MPI_T_VERBOSITY_TUNER_DETAIL	Detailed information required for tuning
MPI_T_VERBOSITY_TUNER_ALL	All information required for tuning
MPI_T_VERBOSITY_MPIDEV_BASIC	Basic information for MPI implementors
MPI_T_VERBOSITY_MPIDEV_DETAIL	Detailed information for MPI implementors
MPI_T_VERBOSITY_MPIDEV_ALL	All information for MPI implementors

Table 14.1: MPI tool information interface verbosity levels.

variable may refer more generally to the MPI environment of the process. Except for the last case, the variable must be bound to exactly one MPI object before it can be used. Table 14.2 lists all MPI object types to which an MPI tool information interface variable can be bound, together with the matching constant that MPI tool information interface routines return 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_COMM	MPI communicators
MPI_T_BIND_MPI_DATATYPE	MPI datatypes
MPI_T_BIND_MPI_ERRHANDLER	MPI error handlers
MPI_T_BIND_MPI_FILE	MPI file handles
MPI_T_BIND_MPI_GROUP	MPI groups
MPI_T_BIND_MPI_OP	MPI reduction operators
MPI_T_BIND_MPI_REQUEST	MPI requests
MPI_T_BIND_MPI_WIN	MPI windows for one-sided communication
MPI_T_BIND_MPI_MESSAGE	MPI message object
MPI_T_BIND_MPI_INFO	MPI info object

Table 14.2: Constants to identify associations of variables.

Rationale. Some variables have meanings tied to a specific MPI object. Examples include the number of send or receive operations [using] that use a particular datatype, the number of times a particular error handler has been called, or the communication protocol and “eager limit” used for a particular communicator. Creating a new MPI tool information interface variable for each MPI object would cause the number of variables to grow without [bounds] bound, since they cannot be reused to avoid naming conflicts. By associating MPI tool information interface variables with a specific MPI object, the MPI implementation only must specify and maintain a single variable, which can then be applied to as many MPI objects of the respective type as created during the program’s execution. (*End of rationale.*)

14.3.3 Convention for Returning Strings

Several MPI tool information interface 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 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 user passes the null pointer as the length argument, the buffer argument is ignored and nothing is returned.

14.3.4 Initialization and Finalization

The MPI tool information interface 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 tool information interface must initialize the MPI tool information interface in the processes that will use the interface before calling any other of its routines. A user can initialize the MPI tool information interface by calling `MPI_T_INIT_THREAD`, which can be called multiple times. In addition, this routine initializes the thread environment for all routines in the MPI tool information interface. Calling this routine when the MPI tool information interface is already initialized has no effect beyond increasing the reference count of how often the interface has been initialized. 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 12.4. The call returns in `provided` information about the actual level of thread support that will be provided by the MPI implementation for calls to MPI tool information interface routines. It can be one of the four values listed in Section 12.4.

The MPI specification does not require all MPI processes to exist before the call to `MPI_INIT`. If the MPI tool information interface is used before `MPI_INIT` has been called, `MPI_T_INIT_THREAD` must be called on each process that will use the MPI tool information interface. Processes created by the MPI implementation during `MPI_INIT` inherit the status of the MPI tool information interface (whether it is initialized or not as well as all active sessions and handles) from the process from which they are created.

Processes created at runtime as a result of calls [MPI's] to MPI's dynamic process man-

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

ticket0-new.

agement require their own initialization before they can use the MPI tool information interface.

Advice to users. If `MPI_T_INIT_THREAD` is called before `MPI_INIT_THREAD`, the requested and granted thread level for `MPI_T_INIT_THREAD` may influence the behavior and return value of `MPI_INIT_THREAD`. The same is true for the reverse order. (*End of advice to users.*)

Advice to implementors. MPI implementations should strive to make as many control or performance variables available before `MPI_INIT` (instead of adding them within `MPI_INIT`) to allow tools the most flexibility. In particular, control variables should be available before `MPI_INIT` if their value cannot be changed after `MPI_INIT`. (*End of advice to implementors.*)

MPI_T_FINALIZE()

```
int MPI_T_finalize(void)
```

This routine finalizes the use of the MPI tool information interface and may be called as often as the corresponding `MPI_T_INIT_THREAD` routine up to the current point of execution. Calling it more times returns a corresponding error code. As long as the number of calls to `MPI_T_FINALIZE` is smaller than the number of calls to `MPI_T_INIT_THREAD` up to the current point of execution, the MPI tool information interface remains initialized and calls to its routines are permissible. Further, additional calls to `MPI_T_INIT_THREAD` after one or more calls to `MPI_T_FINALIZE` are permissible.

Once `MPI_T_FINALIZE` is called the same number of times as the routine `MPI_T_INIT_THREAD` up to the current point of execution, the MPI tool information interface is no longer initialized. The interface can be reinitialized by subsequent calls to `MPI_T_INIT_THREAD`.

At the end of the program execution, unless `MPI_ABORT` is called, an application must have called `MPI_T_INIT_THREAD` and `MPI_T_FINALIZE` an equal number of times.

14.3.5 Datatype System

All variables managed through the MPI tool information interface represent their values through typed buffers of a given length and type using an MPI datatype (similar to regular send/receive buffers). Since the initialization of the MPI tool information interface is separate from the initialization of MPI, MPI tool information interface routines can be called before `MPI_INIT`. Consequently, these routines can also use MPI datatypes before `MPI_INIT`. Therefore, within the context of the MPI tool information interface, it is permissible to use a subset of MPI datatypes as specified below before a call to `MPI_INIT` (or equivalent).

Rationale. The MPI tool information interface relies mainly on unsigned datatypes for integer values since most variables are expected to represent counters or resource sizes. `MPI_INT` is provided for additional flexibility and is expected to be used mainly for control variables and enumeration types (see below).

Providing all basic datatypes, in particular providing all signed and unsigned variants of integer types, would lead to a larger number of types, which tools need to interpret.

MPI_INT	1
MPI_UNSIGNED	2
MPI_UNSIGNED_LONG	3
MPI_UNSIGNED_LONG_LONG	4
MPI_COUNT	5
MPI_CHAR	6
MPI_DOUBLE	7

Table 14.3: MPI datatypes that can be used by the MPI tool information interface.

This would cause unnecessary complexity in the implementation of tools based on the MPI tool information interface. (*End of rationale.*)

The MPI tool information interface only relies on a subset of the basic MPI datatypes and does not use any derived MPI datatypes. Table 14.3 lists all MPI datatypes that can be returned by the MPI tool information interface to represent its variables.

Rationale. The MPI tool information interface requires a significantly simpler type system than MPI itself. Therefore, only its required subset must be present before MPI_INIT (or equivalent) and MPI implementations do not need to initialize the complete MPI datatype system. (*End of rationale.*)

For variables of type MPI_INT, an MPI implementation can provide additional information by associating names with a fixed number of values. We refer to this information in the following as an enumeration. In this case, the respective calls that provide additional metadata for each control or performance variable, i.e., MPI_T_CVAR_GET_INFO (Section 14.3.6) and MPI_T_PVAR_GET_INFO (Section 14.3.7), return a handle of type MPI_T_enum that can be passed to the following functions to extract additional information. Thus, the MPI implementation can describe variables with a fixed set of values that each represents a particular state. Each enumeration type can have N different values, with a fixed N that can be queried using MPI_T_ENUM_GET_INFO.

MPI_T_ENUM_GET_INFO(enumtype, num, name, name_len)

IN	enumtype	enumeration to be queried (handle)
OUT	num	number of discrete values represented by this enumeration (integer)
OUT	name	buffer to return the string containing the name of the enumeration (string)
INOUT	name_len	length of the string and/or buffer for name (integer)

```
int MPI_T_enum_get_info(MPI_T_enum enumtype, int *num, char *name, int
                        *name_len)
```

If enumtype is a valid enumeration, this routine returns the number of items represented by this enumeration type. range and the name of the enumeration. N must be greater than 0, i.e., the enumeration must represent at least one value.

1 The arguments `name` and `name_len` are used to return the name of the enumerations
2 as described in Section 14.3.3.

3 The routine is required to return a name of at least length one. This name must be
4 unique with respect to all other names for enumerations that the MPI implementation uses.

5 Names associated with individual values in each enumeration `enumtype` can be queried
6 using `MPI_T_ENUM_GET_ITEM`.

7
8
9 `MPI_T_ENUM_GET_ITEM(enumtype, index, value, name, name_len)`

10	IN	<code>enumtype</code>	enumeration to be queried (handle)
11	IN	<code>index</code>	number of the value to be queried in this enumeration 12 (integer)
13	OUT	<code>value</code>	variable value (integer)
14	OUT	<code>name</code>	buffer to return the string containing the name of the 15 enumeration item (string)
16			
17	INOUT	<code>name_len</code>	length of the string and/or buffer for name (integer)
18			

19 `int MPI_T_enum_get_item(MPI_T_enum enumtype, int [intex]index, int value,`
20 `char *name, int *name_len)`

21
22 The arguments `name` and `name_len` are used to return the name of the enumeration
23 item as described in Section 14.3.3.

24 If completed successfully, the routine returns the name/value pair [describing] that de-
25 scribes the enumeration at the specified index. The call is further required to return a name
26 of at least length one. This name must be unique with respect to all other names of items
27 for the same enumeration.

28 29 14.3.6 Control Variables

30
31 The routines described in this section of the MPI tool information interface specification
32 focus on the ability to list, query, and possibly set control variables exposed by the MPI
33 implementation. These variables can typically be used by the user to fine tune properties
34 and configuration settings of the MPI implementation. On many systems, such variables
35 can be set using environment variables, although other configuration mechanisms may be
36 available, such as configuration files or central configuration registries. A typical example
37 that is available in several existing MPI implementations is the ability to specify an “eager
38 limit”, i.e., an upper bound on the size of messages sent or received using an eager protocol.

39 40 Control Variable Query Functions

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

45 An MPI implementation is allowed to increase the number of control variables during
46 the execution of an MPI application when new variables become available through dynamic
47 loading. However, MPI implementations are not allowed to change the index of a control
48

variable or [\[delete\]](#) to delete a variable once it has been added to the set. When variables become inactive, e.g., through dynamic unloading, accessing its value should return a corresponding error code.

Advice to users. While the MPI tool information interface guarantees that indices or variable properties do not change during a particular run of an MPI program, it does not provide a similar guarantee between runs. (*End of advice to users.*)

The following function can be used to query the number of control variables, *num_cvar*:

```
MPI_T_CVAR_GET_NUM(num_cvar)
```

OUT	num_cvar	returns number of control variables (integer)
-----	----------	---

```
int MPI_T_cvar_get_num(int *num_cvar)
```

The function `MPI_T_CVAR_GET_INFO` provides access to additional information for each variable.

```
MPI_T_CVAR_GET_INFO(cvar_index, name, name_len, verbosity, datatype, enumtype, desc,
                    desc_len, bind, scope)
```

IN	cvar_index	index of the control variable to be queried, value between 0 and <i>num_cvar</i> - 1 (integer)
OUT	name	buffer to return the string containing the name of the control variable (string)
INOUT	name_len	length of the string and/or buffer for <i>name</i> (integer)
OUT	verbosity	verbosity level of this variable (integer)
OUT	datatype	MPI datatype of the information stored in the control variable (handle)
OUT	enumtype	optional descriptor for enumeration information (handle)
OUT	desc	buffer to return the string containing a description of the control variable (string)
INOUT	desc_len	length of the string and/or buffer for <i>desc</i> (integer)
OUT	bind	type of MPI object to which this variable must be bound (integer)
OUT	scope	scope of when changes to this variable are possible (integer)

```
int MPI_T_cvar_get_info(int cvar_index, char *name, int *name_len, int
                       *verbosity, MPI_Datatype *datatype, MPI_T_enum *enumtype, char
                       *desc, int *desc_len, int *bind, int *scope)
```

After a successful call to `MPI_T_CVAR_GET_INFO` for a particular variable, subsequent calls to this routine [\[querying\]](#) that query information about the same variable must return

the same information. An MPI 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 14.3.3.

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 control variables used by the MPI implementation.

The argument `verbosity` returns the verbosity level of the variable (see Section 14.3.1).

The argument `datatype` returns the MPI datatype that is used to represent the control variable.

If the variable is of type `MPI_INT`, MPI can optionally specify an enumeration for the values represented by this variable and return it in `enumtype`. In this case, MPI returns an enumeration identifier, which can then be used to gather more information as described in Section 14.3.5. If the datatype is not `MPI_INT` or the argument `enumtype` is the constant `MPI_T_ENUM_NULL`, no enumeration type is returned.

The arguments `desc` and `desc_len` are used to return a description of the control variable as described in Section 14.3.3.

Returning a description is optional. If an MPI implementation [decides] does 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 14.3.2).

The scope of a variable determines whether changing a variable's value is either local to the process or must be done by the user across multiple processes. The latter is further split into variables that require changes in a group of processes and those that require collective changes among all connected processes. Both cases can require all processes [to either] either to be set to consistent (but potentially different) values or to equal values on every participating process. The description provided with the variable must contain an explanation about the requirements and/or restrictions for setting the particular variable.

On successful return from `MPI_T_CVAR_GET_INFO`, the argument `scope` will be set to one of the constants listed in Table 14.4.

Scope Constant	Description
<code>MPI_T_SCOPE_READONLY</code>	read-only, cannot be written
<code>MPI_T_SCOPE_LOCAL</code>	may be writeable, writing is a local operation
<code>MPI_T_SCOPE_GROUP</code>	may be writeable, must be done to a group of processes, all processes in a group must be set to consistent values
<code>MPI_T_SCOPE_GROUP_EQ</code>	may be writeable, must be done to a group of processes, all processes in a group must be set to the same value
<code>MPI_T_SCOPE_ALL</code>	may be writeable, must be done to all processes, all connected processes must be set to consistent values
<code>MPI_T_SCOPE_ALL_EQ</code>	may be writeable, must be done to all processes, all connected processes must be set to the same value

Table 14.4: Scopes for 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. (*End of advice to users.*)

Example: Printing All Control Variables

Example 14.4

The following example shows how the MPI tool information interface can be used to query and `[print]` to print the names of all available control variables.

```
#include <stdio.h>
#include <stdlib.h>
#include <mpi.h>

int main(int argc, char **argv) {
    int i, err, num, namelen, bind, verbose, scope;
    int threadsupport;
    char name[100];
    MPI_Datatype datatype;

    err=MPI_T_init_thread(MPI_THREAD_SINGLE,&threadsupport);
    if (err!=MPI_SUCCESS)
        return err;

    err=MPI_T_cvar_get_num(&num);
    if (err!=MPI_SUCCESS)
        return err;

    for (i=0; i<num; i++) {
        namelen=100;
        err=MPI_T_cvar_get_info(i, name, &namelen,
                               &verbose, &datatype, MPI_T_ENUM_NULL,
                               NULL, NULL, /*no description */
                               &bind, &scope);
        if (err!=MPI_SUCCESS) return err;
        printf("Var %i: %s\n", i, name);
    }

    err=MPI_T_finalize();
    if (err!=MPI_SUCCESS)
        return 1;
    else
        return 0;
}
```

Handle Allocation and Deallocation

Before reading or writing the value of a variable, a user must first allocate a handle of type `MPI_T_cvar_handle` for the variable by binding it to an MPI object (see also Section 14.3.2).

Rationale. Handles used in the MPI tool information interface are distinct from handles used in the remaining parts of the MPI standard because they must be usable before `MPI_INIT` and after `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.*)

`MPI_T_CVAR_HANDLE_ALLOC(cvar_index, object, handle, count)`

IN	<code>cvar_index</code>	index of control variable for which handle is to be allocated (index)
IN	<code>obj_handle</code>	reference to a handle of the MPI object to which this variable is supposed to be bound (pointer)
OUT	<code>handle</code>	allocated handle (handle)
OUT	<code>count</code>	number of elements used to represent this variable (integer)

```
int MPI_T_cvar_handle_alloc(int cvar_index, void *obj_handle,
                           MPI_T_cvar_handle *handle, int *count)
```

This routine binds the control variable specified by the argument `index` to an MPI object. The object is passed in the argument `obj_handle` as an address to a local variable that stores the object's handle. The handle allocated to reference the variable is returned in the argument `handle`. Upon successful return, `count` contains the number of elements (of the datatype returned by a previous `MPI_T_CVAR_GET_INFO` call) used to represent this variable.

Advice to users. The `count` can be different based on the MPI object to which [it]the control variable was bound. For example, variables bound to communicators could have a count that matches the size of the communicator.

It is not portable to pass references to predefined MPI object handles, such as `MPI_COMM_WORLD` to this routine, since their implementation depends on the MPI library. Instead, such object handles should be stored in a local variable and the address of this local variables should be passed into `MPI_T_CVAR_HANDLE_ALLOC`. (*End of advice to users.*)

The value of `cvar_index` should be in the range 0 to `num_cvar - 1`, where `num_cvar` is the number of available control variables as determined from a prior call to `MPI_T_CVAR_GET_NUM`. The type of the MPI object it references must be consistent with the type returned in the `bind` argument in a prior call to `MPI_T_CVAR_GET_INFO`.

In the case [the]that the `bind` argument returned by `MPI_T_CVAR_GET_INFO` equals `MPI_T_BIND_NO_OBJECT`, the argument `obj_handle` is ignored.

`MPI_T_CVAR_HANDLE_FREE(handle)`

INOUT	<code>handle</code>	handle to be freed (handle)
-------	---------------------	-----------------------------

```
int MPI_T_cvar_handle_free(MPI_T_cvar_handle *handle)
```

When a handle is no longer needed, a user of the MPI tool information interface should call `MPI_T_CVAR_HANDLE_FREE` to free the handle and the associated resources in the MPI implementation. On a successful return, MPI sets the handle to `MPI_T_CVAR_HANDLE_NULL`.

Control Variable Access Functions

`MPI_T_CVAR_READ(handle, buf)`

IN	<code>handle</code>	handle to the control variable to be read (handle)
OUT	<code>buf</code>	initial address of storage location for variable value (choice)

```
int MPI_T_cvar_read(MPI_T_cvar_handle handle, void* buf)
```

This routine queries the value of the control variable identified by the argument `handle` and stores the result in the buffer identified by the parameter `buf`. The user must ensure that the buffer is of the appropriate size to hold the entire value of the control variable (based on the returned datatype and count from prior corresponding calls to `MPI_T_CVAR_GET_INFO` and `MPI_T_CVAR_HANDLE_ALLOC`, respectively).

`MPI_T_CVAR_WRITE(handle, buf)`

IN	<code>handle</code>	handle to the control variable to be written (handle)
IN	<code>buf</code>	initial address of storage location for variable value (choice)

```
int MPI_T_cvar_write(MPI_T_cvar_handle handle, const void* buf)
```

This routine sets the value of the control variable identified by the argument `handle` to the data stored in the buffer identified by the parameter `buf`. The user must ensure that the buffer is of the appropriate size to hold the entire value of the control variable (based on the returned datatype and count from prior corresponding calls to `MPI_T_CVAR_GET_INFO` and `MPI_T_CVAR_HANDLE_ALLOC`, respectively).

If the variable has a global scope (as returned by a prior corresponding `MPI_T_CVAR_GET_INFO` call) any write call to this variable must be issued by the user in all connected (as defined in Section 10.5.4) MPI processes. If the variable has [a group scope]group scope, any write call to this variable must be issued by the user in all MPI processes in the group, which must be described by the MPI implementation in the description by the `MPI_T_CVAR_GET_INFO`.

In both cases, the user must ensure that the writes in all processes are consistent. If the scope is either `MPI_T_SCOPE_GLOBAL_EQ` or `MPI_T_SCOPE_GROUP_EQ` this means that the variable in all processes must be set to the same value.

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

1 Example: Reading the Value of a Control Variable

2
3 **Example 14.5**

4 The following example shows a routine that can be used to query the value with a
5 control variable with a given index. The example assumes that the variable is intended to
6 be bound to an MPI communicator.

```
7
8 int getValue_int_comm(int index, MPI_Comm comm, int *val) {
9     int err,count;
10    MPI_T_cvar_handle handle;
11
12    /* This example assumes that the variable index */
13    /* can be bound to a communicator */
14
15    err=MPI_T_cvar_handle_alloc(index,&comm,&handle,&count);
16    if (err!=MPI_SUCCESS) return err;
17
18    /* The following assumes that the variable is */
19    /* represented by a single integer */
20
21    err=MPI_T_cvar_read(handle,val);
22    if (err!=MPI_SUCCESS) return err;
23
24    err=MPI_T_cvar_handle_free(&handle);
25    return err;
26 }
27
```

28 **14.3.7 Performance Variables**

29
30 ticket0-new. The following section focuses on the ability to list and [\[query\]](#) to query performance vari-
31 ables provided by the MPI implementation. Performance variables provide insight into MPI
32 implementation specific internals and can represent information such as the state of the
33 MPI implementation (e.g., waiting blocked, receiving, not active), aggregated timing data
34 for submodules, or queue sizes and lengths.

35
36 *Rationale.* The interface for performance variables is separate from the interface for
37 control variables, since performance variables have different requirements and param-
38 eters. By keeping them separate, the interface provides cleaner semantics and allows
39 for more performance optimization opportunities. (*End of rationale.*)

40
41 **Performance Variable Classes**

42 Each performance variable is associated with a class that describes its basic semantics,
43 possible datatypes, basic behavior, its starting value, whether it can overflow, and when
44 and how an MPI implementation can change the variable's value. The starting value is
45 ticket0-new. [\[the variable assumes\]](#) that the variable has when it is used for the first time or
46 whenever it is reset.

Advice to users. If a performance variable belongs to a class that can overflow, it is up to the user to [appropriately protect against this]protect against this overflow, e.g., by frequently reading and resetting the variable value. (*End of advice to users.*)

Advice to implementors. MPI implementations should use large enough datatypes for each performance variable to avoid overflows under normal circumstances. (*End of advice to implementors.*)

The classes are defined by the following constants:

- **MPI_T_PVAR_CLASS_STATE**

A performance variable in this class represents a set of discrete states. Variables of this class are represented by `MPI_INT` and can be set by the MPI implementation at any time. Variables of this type should be described further using an enumeration, as discussed in Section 14.3.5. The starting value is the current state of the implementation at the time [the]that the starting value is set. MPI implementations must ensure that variables of this class cannot overflow.

- **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 non-negative and represented by one of the following datatypes: `MPI_UNSIGNED`, `MPI_UNSIGNED_LONG`, `MPI_UNSIGNED_LONG_LONG`, `MPI_DOUBLE`. The starting value is the current utilization level of the resource at the time [the]that the starting value is set. MPI implementations must ensure that variables of this class cannot overflow.

- **MPI_T_PVAR_CLASS_SIZE**

A performance variable in this class represents a value that is the fixed size of a resource. Values returned from variables in this class are non-negative and represented by one of the following datatypes: `MPI_UNSIGNED`, `MPI_UNSIGNED_LONG`, `MPI_UNSIGNED_LONG_LONG`, `MPI_DOUBLE`. The starting value is the current utilization level of the resource at the time [the]that the starting value is set. MPI implementations must ensure that variables of this class cannot overflow.

- **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_DOUBLE` datatype. The value must always be between 0.0 (resource not used at all) and 1.0 (resource completely used). The starting value is the current percentage utilization level of the resource at the time [the]that the starting value is set. MPI implementations must ensure that variables of this class cannot overflow.

- **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 is non-negative and grows monotonically from the initialization or reset of the variable. It can be represented by one of the following datatypes: `MPI_UNSIGNED`,

1 MPI_UNSIGNED_LONG, MPI_UNSIGNED_LONG_LONG, MPI_DOUBLE. The
 ticket0-new. 2 starting value is the current utilization level of the resource at the time [the]that
 3 the starting value is set. MPI implementations must ensure that variables of this class
 4 cannot overflow.

5
 6 • MPI_T_PVAR_CLASS_LOWWATERMARK

7 A performance variable in this class represents a value that describes the low wa-
 8 termark utilization of a resource. The value of a variable of this class is non-
 9 negative and decreases monotonically from the initialization or reset of the vari-
 10 able. It can be represented by one of the following datatypes: MPI_UNSIGNED,
 ticket0-new. 11 MPI_UNSIGNED_LONG, MPI_UNSIGNED_LONG_LONG, MPI_DOUBLE. The
 12 starting value is the current utilization level of the resource at the time [the]that
 13 the starting value is set. MPI implementations must ensure that variables of this class
 14 cannot overflow.

15
 16 • MPI_T_PVAR_CLASS_COUNTER

17 A performance variable in this class counts the number of occurrences of a specific
 18 event (e.g., the number of memory allocations within an MPI library). The value of
 19 a variable of this class increases monotonically from the initialization or reset of the
 20 performance variable by one for each specific event that is observed. Values must be
 21 non-negative and represented by one of the following datatypes: MPI_UNSIGNED,
 22 MPI_UNSIGNED_LONG, MPI_UNSIGNED_LONG_LONG. The starting value for
 23 variables of this class is 0. Variables of this class can overflow.

24
 25 • MPI_T_PVAR_CLASS_AGGREGATE

26 The value of a performance variable in this class is an aggregated value that repre-
 27 sents a sum of arguments processed during a specific event (e.g., the amount of mem-
 28 ory allocated by all memory allocations). This class is similar to the counter class,
 29 but instead of counting individual events, the value can be incremented by arbitrary
 30 amounts. The value of a variable of this class increases monotonically from the initial-
 31 ization or reset of the performance variable. It must be non-negative and represented
 32 by one of the following datatypes: MPI_UNSIGNED, MPI_UNSIGNED_LONG,
 33 MPI_UNSIGNED_LONG_LONG, MPI_DOUBLE. The starting value for variables
 34 of this class is 0. Variables of this class can overflow.

35
 36 • MPI_T_PVAR_CLASS_TIMER

37 The value of a performance variable in this class represents the aggregated time
 38 that the MPI implementation spends executing a particular event, type of event,
 39 or section of the MPI library. This class has the same basic semantics as
 40 MPI_T_PVAR_CLASS_AGGREGATE, but explicitly records a timing value. The
 41 value of a variable of this class increases monotonically from the initialization
 42 or reset of the performance variable. It must be non-negative and represented
 ticket0-new. 43 by one of the following datatypes: MPI_UNSIGNED, MPI_UNSIGNED_LONG,
 44 MPI_UNSIGNED_LONG_LONG, MPI_DOUBLE. The starting value for variables
 45 of this class is 0. If the type MPI_DOUBLE is used, the units [representing]that
 46 represent time in this datatype must match the units used by MPI_WTIME. Oth-
 47 erwise, the time units should be documented, e.g., in the description returned by
 48 MPI_T_PVAR_GET_INFO. Variables of this class can overflow.

- `MPI_T_PVAR_CLASS_GENERIC`

This class can be used to describe a variable that does not fit into any of the other classes. For variables in this class, the starting value is variable specific and implementation defined.

Performance Variable Query Functions

An MPI implementation exports a set of N performance variables through the MPI tool information interface. If N is zero, then the MPI 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 calls to identify the individual variables.

An MPI 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 implementations are not allowed to change the index of a performance variable or `[delete]` to delete a variable once it has been added to the set. When variables become inactive, e.g., through dynamic unloading, accessing its value should return a corresponding error code.

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

`MPI_T_PVAR_GET_NUM(num_pvar)`

OUT `num_pvar` returns number of performance variables (integer)

```
int MPI_T_pvar_get_num(int *num_pvar)
```

The function `MPI_T_PVAR_GET_INFO` provides access to additional information for each variable.

```

1 MPI_T_PVAR_GET_INFO(pvar_index, name, name_len, verbosity, varclass, datatype, enum-
2     type, desc, desc_len, bind, readonly, continuous, atomic)
3     IN      pvar_index      index of the performance variable to be queried be-
4     between 0 and num_pvar - 1 (integer)
5     OUT     name            buffer to return the string containing the name of the
6     performance variable (string)
7
8     INOUT   name_len       length of the string and/or buffer for name (integer)
9     OUT     verbosity      verbosity level of this variable (integer)
10    OUT     var_class      class of performance variable (integer)
11    OUT     datatype       MPI datatype of the information stored in the perfor-
12    mance variable (handle)
13
14    OUT     enumtype       optional descriptor for enumeration information (han-
15    dle)
16
17    OUT     desc           buffer to return the string containing a description of
18    the performance variable (string)
19
20    INOUT   desc_len      length of the string and/or buffer for desc (integer)
21    OUT     bind           type of MPI object to which this variable must be
22    bound (integer)
23
24    OUT     readonly      flag indicating whether the variable can be written/reset
25    (integer)
26
27    OUT     continuous    flag indicating whether the variable can be started and
28    stopped or is continuously active (integer)
29
30    OUT     atomic        flag indicating whether the variable can be atomically
31    read and reset (integer)
32
33    int MPI_T_pvar_get_info(int pvar_index, char *name, int *name_len, int
34    *verbosity, int *var_class, MPI_Datatype *datatype, MPI_T_enum
35    *enumtype, char *desc, int *desc_len, int *bind, int
36    *readonly, int *continuous, int *atomic)

```

After a successful call to `MPI_T_PVAR_GET_INFO` for a particular variable, subsequent calls to this routine [querying] that query information about the same variable must return the same information. An MPI 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 14.3.3. If completed successfully, the routine is required to return a name of at least length one.

The argument `verbosity` returns the verbosity level of the variable (see Section 14.3.1).

The class of the performance variable is returned in the parameter `var_class`. The class must be one of the constants defined in Section 14.3.7.

The combination of the name and the class of the performance variable must be unique with respect to all other names for performance variables used by the MPI implementation.

1 `MPI_T_PVAR_SESSION_FREE(session)`

2 INOUT `session` identifier of performance experiment session (handle)

4 `int MPI_T_pvar_session_free(MPI_T_pvar_session *session)`

6 This call frees an existing session. Calls to the MPI tool information interface can no longer be made within the context of a session after it is freed. On a successful return, MPI sets the session identifier to `MPI_T_PVAR_SESSION_NULL`.

10 Handle Allocation and Deallocation

12 Before using a performance variable, a user must first allocate a handle of type `MPI_T_pvar_handle` for the variable by binding it to an MPI object (see also Section 14.3.2).

15 `MPI_T_PVAR_HANDLE_ALLOC(session, pvar_index, obj_handle, handle, count)`

17 IN `session` identifier of performance experiment session (handle)

18 IN `pvar_index` index of performance variable for which handle is to be allocated (integer)

20 IN `obj_handle` reference to a handle of the MPI object to which this variable is supposed to be bound (pointer)

23 OUT `handle` allocated handle (handle)

24 OUT `count` number of elements used to represent this variable (integer)

27 `int MPI_T_pvar_handle_alloc(MPI_T_pvar_session session, int pvar_index, void *obj_handle, MPI_T_pvar_handle *handle, int *count)`

30 This routine binds the performance variable specified by the argument `index` to an MPI object in the session identified by the parameter `session`. The object is passed in the argument `obj_handle` as an address to a local variable that stores the object's handle. The handle allocated to reference the variable is returned in the argument `handle`. Upon successful return, `count` contains the number of elements (of the datatype returned by a previous `MPI_T_PVAR_GET_INFO` call) used to represent this variable.

36 *Advice to users.* The `count` can be different based on the MPI object, to which [is] the performance variable was bound. For example, variables bound to communicators could have a count that matches the size of the communicator.

40 It is not portable to pass references to predefined MPI object handles, such as `MPI_COMM_WORLD`, to this routine, since their implementation depends on the MPI library. Instead, such object handles should be stored in a local variable and the address of this local [variables] variable should be passed into `MPI_T_PVAR_HANDLE_ALLOC`.
 44 *(End of advice to users.)*

46 The value of `index` should be in the range 0 to `num_pvar - 1`, where `num_pvar` is the number of available [control] performance variables as determined from a prior call to

MPI_T_PVAR_GET_NUM. The type of the MPI object it references must be consistent with the type returned in the `bind` argument in a prior call to `MPI_T_PVAR_GET_INFO`.

In the case the `bind` argument equals `MPI_T_BIND_NO_OBJECT`, the argument `obj_handle` is ignored.

`MPI_T_PVAR_HANDLE_FREE(session, handle)`

IN	<code>session</code>	identifier of performance experiment session (handle)
INOUT	<code>handle</code>	handle to be freed (handle)

```
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 the MPI tool information interface should call `MPI_T_PVAR_HANDLE_FREE` to free the handle in the session identified by the parameter `session` and the associated resources in the MPI implementation. On a successful return, MPI 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. Such variables may be queried at any time, but they cannot be started or stopped 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	<code>session</code>	identifier of performance experiment session (handle)
IN	<code>handle</code>	handle of a performance variable (handle)

```
int MPI_T_pvar_start(MPI_T_pvar_session session, MPI_T_pvar_handle handle)
```

This functions starts the performance variable with the handle identified by the parameter `handle` in the session identified by the parameter `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_SUCCESS` if all variables are started successfully, otherwise `MPI_T_ERR_PVAR_NOSTARTSTOP` is returned. Continuous variables and variables that are already started are ignored when `MPI_T_PVAR_ALL_HANDLES` is specified.

1 MPI_T_PVAR_STOP(session, handle)

2 IN session identifier of performance experiment session (handle)
 3
 4 IN handle handle of a performance variable (handle)

5
 6 int MPI_T_pvar_stop(MPI_T_pvar_session session, MPI_T_pvar_handle handle)

7 This function stops the performance variable with the handle identified by the parameter `handle` in the session identified by the parameter `session`.

8
 9 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_PVAR_NOSTARTSTOP` is returned. Continuous variables and variables that are already stopped are ignored when `MPI_T_PVAR_ALL_HANDLES` is specified.

16 Performance Variable Access Functions

19 MPI_T_PVAR_READ(session, handle, buf)

21 IN session identifier of performance experiment session (handle)
 22 IN handle handle of a performance variable (handle)
 23
 24 OUT buf initial address of storage location for variable value
 25 (choice)

26
 27 int MPI_T_pvar_read(MPI_T_pvar_session session, MPI_T_pvar_handle handle,
 28 void* buf)

29 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 identified by the parameter `buf`. The user is responsible to ensure that the buffer is of the appropriate size to hold the entire value of the performance variable (based on the datatype and count returned by the corresponding previous calls to `MPI_T_PVAR_GET_INFO` and `MPI_T_PVAR_HANDLE_ALLOC`, respectively).

35 The constant `MPI_T_PVAR_ALL_HANDLES` cannot be used as an argument for the function `MPI_T_PVAR_READ`.

39 MPI_T_PVAR_WRITE(session, handle, buf)

40 IN session identifier of performance experiment session (handle)
 41 IN handle handle of a performance variable (handle)
 42
 43 IN buf initial address of storage location for variable value
 44 (choice)

45
 46 int MPI_T_pvar_write(MPI_T_pvar_session session, MPI_T_pvar_handle handle,
 47 const void* buf)

The `MPI_T_PVAR_WRITE` call attempts to write the value of the performance variable with the handle identified by the parameter `handle` in the session identified by the parameter `session`. The value to be written is passed in the buffer identified by the parameter `buf`. The user must ensure that the buffer is of the appropriate size to hold the entire value of the performance variable (based on the datatype and count returned by the corresponding previous calls to `MPI_T_PVAR_GET_INFO` and `MPI_T_PVAR_HANDLE_ALLOC`, respectively).

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

The constant `MPI_T_PVAR_ALL_HANDLES` cannot be used as an argument for the function `MPI_T_PVAR_WRITE`.

`MPI_T_PVAR_RESET(session, handle)`

IN	<code>session</code>	identifier of performance experiment session (handle)
IN	<code>handle</code>	handle of a performance variable (handle)

`int MPI_T_pvar_reset(MPI_T_pvar_session session, MPI_T_pvar_handle handle)`

The `MPI_T_PVAR_RESET` call sets the performance variable with the handle identified by the parameter `handle` to its starting value specified in Section 14.3.7. If it is not possible to change the variable, the function returns `MPI_T_ERR_PVAR_NOWRITE`.

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_SUCCESS` if all variables are reset successfully, otherwise `MPI_T_ERR_PVAR_NOWRITE` is returned. Read-only variables are ignored when `MPI_T_PVAR_ALL_HANDLES` is specified.

`MPI_T_PVAR_READRESET(session, handle, buf)`

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

`int MPI_T_pvar_readreset(MPI_T_pvar_session session, MPI_T_pvar_handle handle, void* buf)`

This call atomically 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. If atomic operations on this variable are not supported, this routine returns `MPI_ERR_NOATOMIC`.

The constant `MPI_T_PVAR_ALL_HANDLES` [can not] cannot be used as an argument for the function `MPI_T_PVAR_READRESET`.

Advice to implementors. Sampling based tools rely on the ability to call the MPI tool information interface, in particular routines to start, stop, read, write and reset performance variables, from any program context, including asynchronous contexts

such as signal handlers. MPI implementations should strive, if possible in their particular environment, to enable these usage scenarios for all or a subset of the routines mentioned above. If implementing only a subset, the read, write, and reset routines are typically the most critical for sampling based tools. An MPI implementation should clearly document any restrictions on the program contexts in which the MPI tool information interface can be used. Restrictions might include guaranteeing usage outside of all signals or outside a specific set of signals. Any restrictions could be documented, for example, through the description returned by `MPI_T_PVAR_GET_INFO`. (*End of advice to implementors.*)

Rationale. All routines to read, [write or reset] to write or to reset performance variables require the session argument. [This] This requirement keeps the interface consistent and allows the [use] use of `MPI_T_PVAR_ALL_HANDLES` where appropriate. Further, this opens up additional performance optimizations for the implementation of handles. (*End of rationale.*)

Example: Tool to Detect Receives with Long Unexpected Message Queues

Example 14.6

The following example shows a sample tool to identify receive operations that occur during times with long message queues. This examples assumes that the MPI implementation exports a variable with the name "MPI_T_UMQ_LENGTH" to represent the current length of the unexpected message queue. The tool is implemented as a PMPI tool using the MPI profiling interface.

The tool consists of three parts: (1) the initialization (by intercepting the call to `MPI_INIT`), (2) the test for long unexpected message queues (by intercepting calls to `MPI_RECV`), and (3) the clean up phase (by intercepting the call to `[MPI_FINALIZE.] MPI_FINALIZE`). To capture all receives, the example would have to be extended to have similar wrappers for all receive operations.

Part 1— Initialization: During initialization, the tool searches for the variable and, once the right index is found, allocates a session and a handle for the variable with the found index, and starts the performance variable.

```
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <mpi.h>

/* Global variables for the tool */
static MPI_T_pvar_session session;
static MPI_T_pvar_handle handle;

int MPI_Init(int *argc, char ***argv) {
    int err, num, i, index, namelen, verbosity;
    int var_class, bind, threadsup;
    int readonly, continuous, atomic, count;
    char name[17];
```

```

MPI_Comm comm;
MPI_Datatype datatype;
MPI_T_enum enumtype;

err=PMPI_Init(argc,argv);
if (err!=MPI_SUCCESS) return err;

err=PMPI_T_init_thread(MPI_THREAD_SINGLE,&threadsup);
if (err!=MPI_SUCCESS) return err;

err=PMPI_T_pvar_get_num(&num);
if (err!=MPI_SUCCESS) return err;
index=-1;
i=0;
while ((i<num) && (index<0)) {
    namelen=17;
    err=PMPI_T_pvar_get_info(i, name, namelen, &verbosity,
        &var_class, &datatype, &enumtype, &bind,
        &readonly, &continuous, &atomic);
    if (strcmp(name,"MPI_T_UMQ_LENGTH")==0) index=i;
    i++; }

/* this could be handled in a more flexible way for a generic tool */
assert(index>=0);
assert(var_class==MPI_T_PVAR_CLASS_LEVEL);
assert(datatype==MPI_INT);
assert(bind==MPI_T_BIND_MPI_COMM);

/* Create a session */
err=PMPI_T_pvar_session_create(&session);
if (err!=MPI_SUCCESS) return err;

/* Get a handle and bind to MPI_COMM_WORLD */
comm=MPI_COMM_WORLD;
err=PMPI_T_pvar_handle_alloc(session, index, &comm, &handle, &count);
if (err!=MPI_SUCCESS) return err;

/* this could be handled in a more flexible way for a generic tool */
assert(count==1);

/* Start variable */
err=PMPI_T_pvar_start(session, handle);
if (err!=MPI_SUCCESS) return err;

return MPI_SUCCESS;
}

```

1 Part 2 — Testing the Queue Lengths During Receives: During every receive operation, the
 2 tool reads the unexpected queue length through the matching performance variable and
 3 compares it against a predefined threshold.

```

4 #define THRESHOLD 5
5
6
7 int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int source, int tag,
8             MPI_Comm comm, MPI_Status *status)
9 {
10     int value, err;
11
12     if (comm==MPI_COMM_WORLD) {
13         err=PMPI_T_pvar_read(session, handle, &value);
14         if ((err==MPI_SUCCESS) && (value>THRESHOLD))
15             {
16                 /* tool identified receive called with long UMQ */
17                 /* execute tool functionality, */
18                 /* e.g., gather and print call stack */
19             }
20     }
21
22     return PMPI_Recv(buf, count, datatype, source, tag, comm, status);
23 }

```

25 Part 3 — Termination: In the wrapper for MPI_FINALIZE, the MPI tool information inter-
 26 face is finalized.

```

27
28 int MPI_Finalize()
29 {
30     int err;
31     err=PMPI_T_handle_free(&session, &handle);
32     err=PMPI_T_session_free(&session);
33     err=PMPI_T_finalize();
34     return PMPI_Finalize();
35 }

```

37 14.3.8 Variable Categorization

38 MPI implementations can optionally group performance and control variables into categories
 39 to express logical relationships between various variables. For example, an MPI implemen-
 40 tation could group all control and performance variables that refer to message transfers in
 41 the MPI implementation and thereby distinguish them from variables that refer to local
 42 resources such as memory allocations or other interactions with the operating system.

43 Categories can also contain other categories to form a hierarchical grouping. Categories
 44 can never include themselves, either directly or transitively within other included categories.
 45 Expanding on the example above, this allows MPI to refine the grouping of variables referring
 46 to message transfers into variables to control and [monitor]to monitor message queues,
 47 message matching activities and communication protocols. Each of these groups of variables
 48

ticket0-new.

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 tool information interface. If $N = 0$, then the MPI implementation does not export any categories, otherwise the provided categories are indexed from 0 to $N - 1$. This index number is used in subsequent calls to functions of the MPI tool information interface 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.

Similarly, MPI implementations are allowed to add variables to categories, but they are not allowed to remove variables from categories or change the order in which they are returned.

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

```
MPI_T_CATEGORY_GET_NUM(num_cat)
```

```
OUT    num_cat          current number of categories (integer)
```

```
int MPI_T_category_get_num(int *num_cat)
```

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

```
MPI_T_CATEGORY_GET_INFO(cat_index, name, name_len, desc, desc_len, num_cvars, num_pvars,
                        num_categories)
```

```
IN     cat_index       index of the category to be queried (integer)
OUT    name            buffer to return the string containing the name of the
                        category (string)
INOUT  name_len       length of the string and/or buffer for name (integer)
OUT    desc            buffer to return the string containing the description
                        of the category (string)
INOUT  desc_len       length of the string and/or buffer for desc (integer)
OUT    num_cvars       number of control variables in the category (integer)
OUT    num_pvars       number of performance variables in the category (in-
                        teger)
OUT    num_categories  number of categories contained in the category (inte-
                        ger)
```

```
int MPI_T_category_get_info(int cat_index, char *name, int *name_len, char
                            *desc, int *desc_len, int *num_cvars, int *num_pvars, int
                            *num_categories)
```

1 The arguments `name` and `name_len` are used to return the name of the category as
 2 described in Section 14.3.3.

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

5 The arguments `desc` and `desc_len` are used to return the description of the category as
 6 described in Section 14.3.3.

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

10 The function returns the number of control variables, performance variables and other
 11 categories contained in the queried category in the arguments `num_cvars`, `num_pvars`, and
 12 `num_categories`, respectively.

13
 14 `MPI_T_CATEGORY_GET_CVARS(cat_index, len, indices)`

16	IN	<code>cat_index</code>	index of the category to be queried, in the range $[0, N-1]$ (integer)
17			
18	IN	<code>len</code>	the length of the indices array (integer)
19			
20	OUT	<code>indices</code>	an integer array of size <code>len</code> , indicating control variable 21 indices (array of integers)

22
 23 `int MPI_T_category_get_cvars(int cat_index, int len, int indices[])`

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

26
 27
 28 `MPI_T_CATEGORY_GET_PVARS(cat_index,len,indices)`

29	IN	<code>cat_index</code>	index of the category to be queried, in the range $[0, N-1]$ (integer)
30			
31			
32	IN	<code>len</code>	the length of the indices array (integer)
33			
34	OUT	<code>indices</code>	an integer array of size <code>len</code> , indicating performance 35 variable indices (array of integers)

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

38 `MPI_T_CATEGORY_GET_PVARS` can be used to query which performance variables
 39 are contained in a particular category. A category contains zero or more performance
 40 variables.

`MPI_T_CATEGORY_GET_CATEGORIES(cat_index,len,indices)`

IN	<code>cat_index</code>	index of the category to be queried, in the range $[0, N-1]$ (integer)
IN	<code>len</code>	the length of the indices array (integer)
OUT	<code>indices</code>	an integer array of size <code>len</code> , indicating category indices (array of integers)

```
int MPI_T_category_get_categories(int cat_index, int len, int indices[])
```

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

As mentioned above, MPI implementations can grow the number of categories as well as the number of variables or other categories within a category. In order to allow users of the MPI tool information interface [\[to quickly check\]](#) to check quickly whether new categories have been added or new variables or categories have been added to a category, MPI maintains a virtual timestamp. This timestamp is monotonically increasing during the execution and is returned by the following function:

`MPI_T_CATEGORY_CHANGED(stamp)`

OUT	<code>stamp</code>	a virtual time stamp to indicate the last change to the categories (integer)
-----	--------------------	--

```
int MPI_T_category_changed(int *stamp)
```

If two subsequent calls to this routine return the same timestamp, it is guaranteed that the category information has not changed between the two calls. If the timestamp retrieved from the second call is higher, then some categories have been added or expanded.

Advice to users. The timestamp value is purely virtual and only intended to check for changes in the category information. It should not be used for any other purpose. *(End of advice to users.)*

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

The user is responsible for allocating the arrays passed into the functions `MPI_T_CATEGORY_GET_CVARS`, `MPI_T_CATEGORY_GET_PVARS` and `MPI_T_CATEGORY_GET_CATEGORIES`. 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.

14.3.9 Return Codes for the MPI tool information interface

All functions defined as part of the MPI tool information interface return an integer return code (see Table 14.5) to indicate whether the function [\[has\]](#) was completed successfully or

ticket0-new. 1 [aborted its execution] was aborted. In the latter case the return code indicates the reason
2 for not completing the routine. None of the return codes returned by a routine impact the
3 execution of the MPI process and do not invoke MPI error handlers. The execution of the
4 MPI process continues as if the call would have completed. However, the MPI implementa-
5 tion is not required to check all user provided parameters; if a user passes invalid parameter
6 values to any routine the behavior of the implementation is undefined.

7 All return codes with the prefix MPI_T_ must be unique values and cannot overlap
8 with any other return values returned by the MPI implementation.

9 10 14.3.10 Profiling Interface

11 All requirements for the profiling interfaces, as described in Section 14.2, also apply to
12 the MPI tool information interface. All rules, guidelines, and recommendations from Sec-
13 tion 14.2 apply equally to calls defined as part of the MPI tool information interface.
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

Return Code	Description
Return Codes for all Functions in the MPI tool information interface	
MPI_SUCCESS	Call completed successfully
MPI_T_ERR_MEMORY	Out of memory
MPI_T_ERR_NOTINITIALIZED	Interface not initialized
MPI_T_ERR_CANTINIT	Interface not in the state to be initialized
Return Codes for Datatype Functions: MPI_T_ENUM_*	
MPI_T_ERR_INVALIDINDEX	The enumeration index is invalid or has been deleted.
MPI_T_ERR_INVALIDITEM	The item index queried is out of range (for MPI_T_ENUMITEM only)
Return Codes for variable and category query functions: MPI_T_*_GET_INFO	
MPI_T_ERR_INVALIDINDEX	The variable or category index is invalid
Return Codes for Handle Functions: MPI_T_*_{ALLOCATE FREE}	
MPI_T_ERR_INVALIDINDEX	The variable index is invalid or has been deleted
MPI_T_ERR_INVALIDHANDLE	The handle is invalid
MPI_T_ERR_OUTOFHANDLES	No more handles available
Return Codes for Session Functions: MPI_T_PVAR_SESSION_*	
MPI_T_ERR_OUTOFSESSIONS	No more sessions available
MPI_T_ERR_INVALIDSESSION	Session argument is not a valid session
Return Codes for Control Variable Access Functions: MPI_T_CVAR_READ, WRITE	
MPI_T_ERR_CVAR_SETNOTNOW	Variable cannot be set at this moment
MPI_T_ERR_CVAR_SETNEVER	Variable cannot be set until end of execution
MPI_T_ERR_INVALIDHANDLE	The handle is invalid
Return Codes for Performance Variable Access and Control: MPI_T_PVAR_{START STOP READ WRITE RESET READREST}	
MPI_T_ERR_INVALIDHANDLE	The handle is invalid
MPI_T_ERR_INVALIDSESSION	Session argument is not a valid session
MPI_T_ERR_PVAR_NOSTARTSTOP	Variable cannot be started or stopped (for MPI_T_PVAR_START and MPI_T_PVAR_STOP)
MPI_T_ERR_PVAR_NOWRITE	Variable cannot be written or reset (for MPI_T_PVAR_WRITE and MPI_T_PVAR_RESET)
MPI_T_NOATOMIC	Variable cannot be read and written atomically (for MPI_T_PVAR_READRESET)
Return Codes for Category Functions: MPI_T_CATEGORY_*	
MPI_T_ERR_INVALIDINDEX	The category index is invalid

Table 14.5: Return codes used in functions of the MPI tool information interface.

Index

CONST:MPI_CHAR, 13
CONST:MPI_COMM_WORLD, 18, 26
CONST:MPI_COUNT, 13
CONST:MPI_DOUBLE, 13, 22
CONST:MPI_ERR_NOATOMIC, 29
CONST:MPI_INT, 12, 13, 16, 21, 25
CONST:MPI_SUCCESS, 9, 27–29, 37
CONST:MPI_T_BIND_MPI_COMM, 10
CONST:MPI_T_BIND_MPI_DATATYPE, 10
CONST:MPI_T_BIND_MPI_ERRHANDLER, 10
CONST:MPI_T_BIND_MPI_FILE, 10
CONST:MPI_T_BIND_MPI_GROUP, 10
CONST:MPI_T_BIND_MPI_INFO, 10
CONST:MPI_T_BIND_MPI_MESSAGE, 10
CONST:MPI_T_BIND_MPI_OP, 10
CONST:MPI_T_BIND_MPI_REQUEST, 10
CONST:MPI_T_BIND_MPI_WIN, 10
CONST:MPI_T_BIND_NO_OBJECT, 10, 16, 18, 25, 27
CONST:MPI_T_cvar_handle, 17
CONST:MPI_T_CVAR_HANDLE_NULL, 19
CONST:MPI_T_enum, 13
CONST:MPI_T_ENUM_NULL, 16, 25
CONST:MPI_T_ERR_CANTINIT, 37
CONST:MPI_T_ERR_CVAR_SETNEVER, 19, 37
CONST:MPI_T_ERR_CVAR_SETNOTNOW, 19, 37
CONST:MPI_T_ERR_INVALIDHANDLE, 37
CONST:MPI_T_ERR_INVALIDINDEX, 37
CONST:MPI_T_ERR_INVALIDITEM, 37
CONST:MPI_T_ERR_INVALIDSESSION, 37
CONST:MPI_T_ERR_MEMORY, 37
CONST:MPI_T_ERR_NOTINITIALIZED, 37
CONST:MPI_T_ERR_OUTOFHANDLES, 37
CONST:MPI_T_ERR_OUTOFSESSIONS, 37
CONST:MPI_T_ERR_PVAR_NOSTARTSTOP, 27, 28, 37
CONST:MPI_T_ERR_PVAR_NOWRITE, 29, 37
CONST:MPI_T_NOATOMIC, 37
CONST:MPI_T_PVAR_ALL_HANDLES, 27–30
CONST:MPI_T_PVAR_CLASS_AGGREGATE, 22
CONST:MPI_T_PVAR_CLASS_COUNTER, 22
CONST:MPI_T_PVAR_CLASS_GENERIC, 23
CONST:MPI_T_PVAR_CLASS_HIGHWATERMARK, 21
CONST:MPI_T_PVAR_CLASS_LEVEL, 21
CONST:MPI_T_PVAR_CLASS_LOWWATERMARK, 22
CONST:MPI_T_PVAR_CLASS_PERCENTAGE, 21
CONST:MPI_T_PVAR_CLASS_SIZE, 21
CONST:MPI_T_PVAR_CLASS_STATE, 21
CONST:MPI_T_PVAR_CLASS_TIMER, 22
CONST:MPI_T_pvar_handle, 26
CONST:MPI_T_PVAR_HANDLE_NULL, 27
CONST:MPI_T_pvar_session, 25
CONST:MPI_T_PVAR_SESSION_NULL, 26
CONST:MPI_T_SCOPE_ALL, 16
CONST:MPI_T_SCOPE_ALL_EQ, 16
CONST:MPI_T_SCOPE_GLOBAL_EQ, 19
CONST:MPI_T_SCOPE_GROUP, 16
CONST:MPI_T_SCOPE_GROUP_EQ, 16, 19
CONST:MPI_T_SCOPE_LOCAL, 16
CONST:MPI_T_SCOPE_READONLY, 16
CONST:MPI_T_VERBOSITY_MPIDEV_ALL, 9, 10
CONST:MPI_T_VERBOSITY_MPIDEV_BASIC, 10
CONST:MPI_T_VERBOSITY_MPIDEV_DETAIL, 10
CONST:MPI_T_VERBOSITY_TUNER_ALL, 10
CONST:MPI_T_VERBOSITY_TUNER_BASIC,

		MPI_T_CATEGORY_GET_PVARS(cat_index, len, indices),	
10			
CONST:MPI_T_VERBOSITY_TUNER_DETAIL,	34		2
10		MPI_T_CVAR_GET_INFO, 13 , 15 , 16 , 18 ,	3
CONST:MPI_T_VERBOSITY_USER_ALL,	19 , 35		4
10		MPI_T_CVAR_GET_INFO(cvar_index, name,	5
CONST:MPI_T_VERBOSITY_USER_BASIC,		name_len, verbosity, datatype, enum-	6
9 , 10		type, desc, desc_len, bind, scope),	7
CONST:MPI_T_VERBOSITY_USER_DETAIL,	15		8
9 , 10		MPI_T_CVAR_GET_NUM, 18	9
CONST:MPI_UNSIGNED, 13		MPI_T_CVAR_GET_NUM(num_cvar), 15	10
CONST:MPI_UNSIGNED_LONG, 13		MPI_T_CVAR_HANDLE_ALLOC, 18 , 19	11
CONST:MPI_UNSIGNED_LONG_LONG, 13		MPI_T_CVAR_HANDLE_ALLOC(cvar_index,	12
		object, handle, count), 18	13
EXAMPLES:Basic tool using performance vari-		MPI_T_CVAR_HANDLE_FREE, 19	14
ables in the MPI tool information inter-		MPI_T_CVAR_HANDLE_FREE(handle), 18	15
face, 30		MPI_T_CVAR_READ(handle, buf), 19	16
EXAMPLES:Profiling interface, 5		MPI_T_CVAR_WRITE(handle, buf), 19	17
EXAMPLES:Reading the value of a control		MPI_T_ENUM_GET_INFO, 13	18
variable in the MPI tool information		MPI_T_ENUM_GET_INFO(enumtype, num,	19
interface, 20		name, name_len), 13	20
EXAMPLES:Using MPI_T_CVAR_GET_INFO		MPI_T_ENUM_GET_ITEM, 14	21
to list all names of control variables.,		MPI_T_ENUM_GET_ITEM(enumtype, in-	22
17		dex, value, name, name_len), 14	23
		MPI_T_ENUMITEM, 37	24
MPI_ABORT, 12		MPI_T_FINALIZE, 12	25
MPI_FINALIZE, 9 , 18 , 30 , 32		MPI_T_FINALIZE(), 12	26
MPI_INIT, 5 , 9 , 11–13 , 18 , 30		MPI_T_INIT_THREAD, 11 , 12	27
MPI_INIT_THREAD, 11 , 12		MPI_T_INIT_THREAD(required, provided),	28
MPI_PCONTROL, 3–5		11	29
MPI_PCONTROL(level, ...), 4		MPI_T_PVAR_GET_INFO, 13 , 22–24 , 26–	30
MPI_RECV, 30		30 , 35	31
MPI_SEND, 5		MPI_T_PVAR_GET_INFO(pvar_index, name,	32
MPI_T_CATEGORY_CHANGED(stamp), 35		name_len, verbosity, varclass, datatype,	33
MPI_T_CATEGORY_GET_CATEGORIES,		enumtype, desc, desc_len, bind, read-	34
35		only, continuous, atomic), 24	35
MPI_T_CATEGORY_GET_CATEGORIES(cat_index, len, indices),	35	MPI_T_PVAR_GET_NUM, 27	36
35		MPI_T_PVAR_GET_NUM(num_pvar), 23	37
MPI_T_CATEGORY_GET_CVARS, 34 , 35		MPI_T_PVAR_HANDLE_ALLOC, 26 , 28 ,	38
MPI_T_CATEGORY_GET_CVARS(cat_index,		29	39
len, indices), 34		MPI_T_PVAR_HANDLE_ALLOC(session, pvar	40
MPI_T_CATEGORY_GET_INFO, 35		index,	
MPI_T_CATEGORY_GET_INFO(cat_index,		obj_handle, handle, count), 26	41
name, name_len, desc, desc_len, num_		MPI_T_PVAR_HANDLE_FREE, 27	42
cvvars,		MPI_T_PVAR_HANDLE_FREE(session, han-	43
num_pvars, num_categories), 33		dle), 27	44
MPI_T_CATEGORY_GET_NUM(num_cat),		MPI_T_PVAR_READ, 28 , 29	45
33		MPI_T_PVAR_READ(session, handle, buf),	46
MPI_T_CATEGORY_GET_PVARS, 34 , 35		28	47
		MPI_T_PVAR_READRESET, 25 , 29 , 37	48

1 MPI_T_PVAR_READRESET(session, han-
2 dle, buf), [29](#)
3 MPI_T_PVAR_RESET, [29](#), [37](#)
4 MPI_T_PVAR_RESET(session, handle), [29](#)
5 MPI_T_PVAR_SESSION_CREATE(session),
6 [25](#)
7 MPI_T_PVAR_SESSION_FREE(session), [26](#)
8 MPI_T_PVAR_START, [37](#)
9 MPI_T_PVAR_START(session, handle), [27](#)
10 MPI_T_PVAR_STOP, [37](#)
11 MPI_T_PVAR_STOP(session, handle), [28](#)
12 MPI_T_PVAR_WRITE, [29](#), [37](#)
13 MPI_T_PVAR_WRITE(session,handle, buf),
14 [28](#)
15 MPI_TYPE_SIZE, [5](#)
16 MPI_WTIME, [5](#), [22](#)
17
18 PMPI_, [2](#)
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