

# MPI: A Message-Passing Interface Standard

Version 3.0

⊤ (Fin2)

⊥ (Fin2)

Message Passing Interface Forum

Draft January 24th, 2011

# Contents

<b>1</b>	<b>Tool Interfaces for MPI</b>	<b>1</b>
1.1	Introduction . . . . .	1
1.2	MPIT Performance Interface . . . . .	1
1.2.1	Verbosity Levels . . . . .	2
1.2.2	Associations between MPIT Variables and MPI Resources . . . . .	3
1.2.3	String Arguments . . . . .	5
1.2.4	Initialization and Finalization . . . . .	6
1.2.5	Type System . . . . .	6
1.2.6	Control Variables . . . . .	9
	Control Variable Query Functions . . . . .	9
	Handle Allocation and Deallocation . . . . .	11
	Control Variable Access Functions . . . . .	12
1.2.7	Performance Variables . . . . .	13
	Performance Variable Classes . . . . .	13
	Performance Variable Query Functions . . . . .	14
	Performance Experiment Sessions . . . . .	16
	Handle Allocation and Deallocation . . . . .	17
	Starting and Stopping of Performance Variables . . . . .	18
	Performance Variable Access Functions . . . . .	19
1.2.8	Variable Categorization . . . . .	21
1.2.9	Return and Error Codes . . . . .	23
1.2.10	Profiling Interface . . . . .	24
	<b>Bibliography</b>	<b>26</b>
	<b>Examples Index</b>	<b>27</b>
	<b>MPIT Constant and Predefined Handle Index</b>	<b>27</b>
	<b>MPIT Declarations Index</b>	<b>29</b>
	<b>MPIT Callback Function Prototype Index</b>	<b>29</b>
	<b>MPIT Function Index</b>	<b>29</b>

# List of Figures

# List of Tables

1.1	MPIT verbosity levels. . . . .	2
1.2	Constant to identify associations of MPIT control variables. . . . .	3
1.3	Predefined MPIT datatypes and their MPI equivalents. . . . .	7
1.4	MPIT type classes. . . . .	7
1.5	Scopes for MPIT control variables. . . . .	11
1.6	Return and error codes used MPIT functions. . . . .	25

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

# Chapter 1

## Tool Interfaces for MPI

### 1.1 Introduction

This chapter discusses a set of interfaces that allows tools such as debuggers, performance analyzers, and others to extract information about the operation of MPI processes. Specifically, this chapter defines the PMPI profiling interface (Section ??) to transparently intercept and inspect any MPI call; and the MPIT tool information interface (Section 1.2) to query 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 as any other MPI function.

### 1.2 MPIT Performance Interface

To optimize MPI applications or their runtime behavior, it is often advantageous to understand the performance switches an MPI implementation offers to the user as well as to monitor properties and timing information from within the MPI implementation. The MPIT interface described in this section provides access to this information.

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

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

To avoid restrictions on the MPI implementation, the MPIT interface allows the implementation to specify which control and performance variables exist. For both types of variables, the interface provides the ability to query the variables offered by the particular MPI implementation, along with additional semantics and descriptions.

To avoid conflicts between the standard MPI functionality and the tools-oriented functionality introduced with MPIT, the MPIT interface is contained in its own name space. All identifiers covered by this interface carry the prefix MPIT and can be used independently from the MPI functionality. This includes initialization and finalization of MPIT, which is

provided through a separate set of routines. Consequently, MPIT routines can be called before `MPI_INIT` and after `MPI_FINALIZE`.

On success all MPIT routines return `MPIT_SUCCESS`, otherwise they return an appropriate error code. Details on error codes can be found in Section 1.2.9. However, errors returned by the MPIT interface are not fatal nor have any impact on the execution of MPI routines.

*Advice to users.* The number and type of control variables and performance variables can vary between MPI implementations, platforms, and even different builds of the same implementation on the same platform. Hence, any application relying on a particular variable will no longer be portable.

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

Since the MPIT interface mostly focuses on tools and support libraries, the MPIT implementations are only required to provide C Bindings. Except where otherwise notes, all conventions and principles governing the C Bindings of the MPI API also apply to the MPIT interface and the MPIT interface must be defined in the same header or API definition file(s) as the regular MPI routines.

### 1.2.1 Verbosity Levels

The MPIT interface provides users access to internal configuration and performance information through a set of control and performance variables, which are defined by the MPIT implementation. Since the number of variables can be large for particular implementations, every variable exported by the MPIT interface must be associated with one of the following verbosity levels to indicate which group of users of MPIT or tools built on top of MPIT are targeted with the information provided by the respective variable.

MPIT_VERBOSITY_USER_BASIC	Basic information of interest for end users
MPIT_VERBOSITY_USER_DETAILED	Detailed information of interest for end users
MPIT_VERBOSITY_USER_VERBOSE	All information of interest for end users
MPIT_VERBOSITY_TUNER_BASIC	Basic information required for tuning
MPIT_VERBOSITY_TUNER_DETAILED	Detailed information required for tuning
MPIT_VERBOSITY_TUNER_VERBOSE	All information required for tuning
MPIT_VERBOSITY_MPIDEV_BASIC	Basic low-level information for MPI developers
MPIT_VERBOSITY_MPIDEV_DETAILED	Detailed low-level information for MPI developers
MPIT_VERBOSITY_MPIDEV_VERBOSE	All low-level information for MPI developers

Table 1.1: MPIT verbosity levels.

Implementations must assign each variable to one of the verbosity levels. MPI implementations should sort all variables according to the intended target audience (end user, performance *timizers*, or MPI developer) and then distinguish three levels of verbosity (basic, detailed, and verbose) within each audience.

*Advice to implementors.* If an MPIT implementation only uses a single verbosity level for all variables, it is recommended to assign all variables to the level `MPI_VERBOSITY_USER_BASIC`. If an MPIT implementation only uses a single verbosity level for all variables intended for each target audience, it is recommended to assign all variables to corresponding basic level. (*End of advice to implementors.*)

### 1.2.2 Associations between MPIT Variables and MPI Resources

Each variable provides access to a particular control setting or performance property provided by the MPI implementation. The meaning of these variables can refer to the complete MPI library as a global variable or can be associated with a particular MPI resource, such as a communicator, datatype, or one-sided communication window. In the latter case, the variable is associated with exactly one MPI resource type. Before it can be used, it must be bound to an instance of an MPI resource of that type. Table 1.2 lists all types of MPI resources supported by MPIT along with a corresponding constant used by the MPIT interface to identify that resource type.

Constant	Associated MPI resource
<code>MPIT_MPI_RESOURCE_TYPE_GLOBAL</code>	N/A — global meaning
<code>MPIT_MPI_RESOURCE_TYPE_COMMUNICATOR</code>	MPI communicators
<code>MPIT_MPI_RESOURCE_TYPE_DATATYPE</code>	MPI datatypes
<code>MPIT_MPI_RESOURCE_TYPE_ERRORHANDLER</code>	MPI error handler
<code>MPIT_MPI_RESOURCE_TYPE_FILE</code>	MPI file handles
<code>MPIT_MPI_RESOURCE_TYPE_GROUP</code>	MPI groups
<code>MPIT_MPI_RESOURCE_TYPE_OPERATOR</code>	MPI reduction operators
<code>MPIT_MPI_RESOURCE_TYPE_REQUEST</code>	MPI requests
<code>MPIT_MPI_RESOURCE_TYPE_WINDOW</code>	MPI windows for one-sided communication

Table 1.2: Constant to identify associations of MPIT control variables.

*Rationale.* Certain variables have meanings that are limited to a particular MPI resource. Examples are the number of send or receive operations using a particular datatype, the number of times an error handler has been called, or the communication protocol and eager limit used for a particular communicator. Creating a separate variable for each MPI resource, e.g., for each communicator, would cause the number of variables to grow unboundedly since they cannot be reused to avoid naming conflicts. By associating variables with MPI resource types, only a single variable must be specified and maintained by the MPI implementation, which can then be reused on as many instances of this MPI resource type as created during the program's execution. (*End of rationale.*)

In order to instantiate a variable with a particular MPI resource instance, the user must be able to convert a reference to a resource of each supported type to a generic reference, which can then be passed to the MPIT routine responsible for instantiating the MPIT variable. For this purpose, the interface offers the following conversion routines, which each take a reference to an MPI resource and return a reference to a generic MPI resource of type `MPIT_MPI_Resource`.

```
1 MPIT_MPI_RESOURCE_COMMUNICATOR(communicator, resource)
2   IN      communicator      Reference to an MPI communicator
3   OUT      resource         Reference to a generic MPI resource
4
5
6 int MPIT_MPI_Resource_communicator(MPI_Comm communicator, MPIT_MPI_Resource
7     *resource)
8
9
```

```
10 MPIT_MPI_RESOURCE_DATATYPE(datatype, resource)
11   IN      datatype         Reference to an MPI datatype
12   OUT      resource         Reference to a generic MPI resource
13
14
15 int MPIT_MPI_Resource_datatype(MPI_Datatype datatype, MPIT_MPI_Resource
16     *resource)
17
18
```

```
19 MPIT_MPI_RESOURCE_ERRORHANDLER(errorhandler, resource)
20   IN      errorhandler     Reference to an MPI error handler
21   OUT      resource         Reference to a generic MPI resource
22
23
24 int MPIT_MPI_Resource_errorhandler(MPI_Errorhandler errorhandler,
25     MPIT_MPI_Resource *resource)
26
27
```

```
28 MPIT_MPI_RESOURCE_GROUP(group, resource)
29   IN      group            Reference to an MPI group
30   OUT      resource         Reference to a generic MPI resource
31
32
33 int MPIT_MPI_Resource_group(MPI_Group group, MPIT_MPI_Resource *resource)
34
35
```

```
36 MPIT_MPI_RESOURCE_FILE(file, resource)
37   IN      file             Reference to an MPI files
38   OUT      resource         Reference to a generic MPI resource
39
40
41 int MPIT_MPI_Resource_file(MPI_File file, MPIT_MPI_Resource *resource)
42
43
44
45
46
47
48
```



MPIT\_MPI\_RESOURCE\_OPERATIONS(operation, resource)

IN	operation	Reference to an MPI reduction operation
OUT	resource	Reference to a generic MPI resource

```
int MPIT_MPI_Resource_operation(MPI_Op operation, MPIT_MPI_Resource
                               *resource)
```

MPIT\_MPI\_RESOURCE\_REQUEST(request, resource)

IN	request	Reference to an MPI asynchronous communication request
OUT	resource	Reference to a generic MPI resource

```
int MPIT_MPI_Resource_request(MPI_Request request, MPIT_MPI_Resource
                              *resource)
```

MPIT\_MPI\_RESOURCE\_WINDOW(window, resource)

IN	window	Reference to an MPI one-sided communication window
OUT	resource	Reference to a generic MPI resource

```
int MPIT_MPI_Resource_window(MPI_Win window, MPIT_MPI_Resource *resource)
```

Additionally, the MPIT interface provides the constant `MPIT_MPI_RESOURCE_GLOBAL` of type `MPIT_MPI_Resource*` that can be used in routines that expect a reference to an MPI resource if the resource type is `MPIT_MPI_RESOURCE_TYPE_GLOBAL`.

*Rationale.* The MPIT interface provides a separate routine for each MPI resource type to provide type safety. The alternative, a single conversion routine that takes a parameter of type `void*` for arguments of any MPI resource type, would not provide this kind of safety. (*End of rationale.*)

### 1.2.3 String Arguments

Several MPIT function return one or more strings. These functions have two arguments for each string to be returned, one identifying a pointer to the buffer in which the string will be returned, and one to pass the length of the buffer. The latter is used as an IN/OUT argument. The user is responsible for the memory allocation of the buffer and must pass the size of the buffer 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. The buffer is always null-terminated. 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.

MPIT does not specify the character encoding of strings in the interface. The only requirement is that strings are terminated with a null character. MPI reserves all type, enumeration type item, variables, and category names with the prefixes “MPI\_” and “MPIT\_” for its own use.

#### 1.2.4 Initialization and Finalization

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

MPIT\_INIT()

```
int MPIT_Init(void)
```

All programs or tools that use the MPIT interface **must** initialize the MPIT interface before calling any other MPIT routine. A user can initialize the MPIT interface by calling MPIT\_INIT, which can be called multiple times.

MPIT\_FINALIZE( )

```
int MPIT_Finalize(void)
```

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

Once MPIT\_FINALIZE is called the same number of times as the routine MPIT\_INIT up to the current point of execution, the MPIT interface is no longer initialized. Further, the call to MPIT\_FINALIZE that ends the initialization of MPIT may clean up all MPIT state, invalidate all open sessions (for the concept of Sessions see Section 1.2.7), and all handles that have been allocated by MPIT. MPIT can be reinitialized by subsequent calls to MPIT\_INIT.

At the end of the program execution, **unless MPI\_ABORT is called**, an application **must** have called MPIT\_INIT and MPIT\_FINALIZE an equal number of times.

#### 1.2.5 Type System

Since the initialization of MPIT is separate from the initialization of MPI, it can not be guaranteed that MPI data types are available at any time during the usage of MPIT. Therefore, the MPIT interface provides a separate type system. All types are represented by a variable or constant of type `MPIT_Datatype` and are classified into two type classes: predefined and enumeration types. The Table 1.3 lists all available constants that can be used to identify or describe a predefined **type** for MPIT calls.

`MPIT_TYPE_GET_CLASS(datatype, typeclass)`

IN        `datatype`                    MPIT datatype to be queried  
 OUT      `typeclass`                    class of the type passed in

`int MPIT_Type_get_class(MPIT_Datatype datatype, int *typeclass)`

This routine returns the type class for the datatype provided by the argument `datatype`. This allows users of MPIT to distinguish whether a datatype is an enumeration type, e.g., to represent the state of a resource, or is one of the predefined types listed in Table 1.3. On return, the `typeclass` argument is set to one of the constants listed in Table 1.4, if `datatype` represents a valid type.

MPIT Datatype	Equivalent MPI Datatype
MPIT_LOGICAL	MPI_LOGICAL
MPIT_BYTE	MPI_BYTE
MPIT_SHORT	MPI_SHORT
MPIT_INT	MPI_INT
MPIT_LONG	MPI_LONG
MPIT_LONG_LONG	MPI_LONG_LONG
MPIT_CHAR	MPI_CHAR
MPIT_FLOAT	MPI_FLOAT
MPIT_DOUBLE	MPI_DOUBLE

Table 1.3: Predefined MPIT datatypes and their MPI equivalents.

MPIT_TYPECLASS_PREDEFINED	the datatype is a predefined datatype
MPIT_TYPECLASS_ENUMERATION	the datatype is an enumeration datatype

Table 1.4: MPIT type classes.

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

*Rationale.* The concept of equivalent MPIT and MPI datatypes allows to safely communicate values of MPIT datatypes using regular MPI messages. (*End of rationale.*)

The function `MPIT_TYPE_GET_SIZE` can be used to query the storage size for each MPIT datatype.

1 `MPIT_TYPE_GET_SIZE(datatype, size)`

2     IN       datatype           MPIT datatype to be queried  
3  
4     OUT       size               Number of bytes required to store a value of datatype  
5                               size  
6

7 `int MPIT_Type_get_size(MPIT_Datatype datatype, int *size)`

8  
9     The second type class, enumeration types, describe variables with a fixed set of discrete  
10 values. These types are represented through integer variables and have `MPI_INT` as their  
11 equivalent MPI type. Their values range from 0 to  $N - 1$ , with a fixed  $N$  that can be queried  
12 using `MPIT_TYPE_ENUM_QUERY`.

13  
14 `MPIT_TYPE_ENUM_GET_INFO(datatype, num, name, name_len)`

15     IN       datatype           MPIT datatype to be queried  
16  
17     OUT       num                number of discrete values represented by this enumer-  
18                               ation datatype  
19     OUT       name               buffer to return the name of the enumeration type  
20  
21     INOUT     name\_len           length of the string and/or buffer for name

22  
23 `int MPIT_Type_enum_get_info (MPIT_Datatype datatype, int *num, char *name,`  
24 `int *name_len)`

25     This routine returns, if `datatype` represents a valid enumeration type, the size of the  
26 enumeration as well as a name for it.

27     The arguments `name` and `name_len` are used to return the name of the type as described  
28 in Section 1.2.3. .

29     If completed successfully, the routine is required to return a name of at least length  
30 one, which is unique with respect to all other names for MPIT datatypes used by the MPI  
31 implementation.

32     Names for the individual items in each enumeration type can be queried using  
33 `MPIT_TYPE_ENUM_GET_ITEM`.

34  
35  
36 `MPIT_TYPE_ENUM_GET_ITEM(datatype, item, name, name_len)`

37     IN       datatype           MPIT datatype to be queried  
38  
39     IN       item                item number in the MPIT datatype to be queried  
40     OUT       name               buffer to return the name of the enumeration item  
41  
42     INOUT     name\_len           length of the string and/or buffer for name

43  
44 `int MPIT_Type_enum_get_item (MPIT_Datatype datatype, int item, char *name,`  
45 `int *name_len)`

46     The arguments `name` and `name_len` are used to return the name of the enumeration  
47 item as described in Section 1.2.3.  
48

If completed successfully, the routine is required to return a name of at least length one, which is unique with respect to all other names of items for the same MPIT enumeration type.

### 1.2.6 Control Variables

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

#### Control Variable Query Functions

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

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

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

`MPIT_CONTROLVAR_GET_NUM(num)`

OUT      num                                      returns number of control variables

```
int MPIT_Controlvar_get_num (int *num)
```

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

```

1  MPIT_CONTROLVAR_GET_INFO(index, name, name_len, verbosity, datatype, count, desc,
2  desc_len, assoc, attributes)

```

3	IN	index	index of the control variable to be queried
4			
5	OUT	name	buffer to return the name of the control variable
6	INOUT	name_len	length of the string and/or buffer for name
7			
8	OUT	verbosity	verbosity level of this variable
9	OUT	datatype	MPIT type of the information stored in the control variable
10			
11	OUT	count	number of elements returned
12	OUT	desc	buffer to return a description of the control variable
13			
14			
15	INOUT	desc_len	length of the string and/or buffer for desc
16	OUT	assoc	type of MPI resource this variable is associated with
17	OUT	attributes	additional attributes defining this variable
18			

```

19
20 int MPIT_Controlvar_get_info(int index, char *name, int *name_len, int
21     *verbosity, MPIT_Datatype *datatype, int *count, char *desc,
22     int *desc_len, int *assoc, MPIT_Controlvar_attributes
23     *attributes)

```

24 After a successful call to `MPIT_CONTROLVAR_GET_INFO` for a particular variable, subsequent calls to this routine querying information about the same variable **must** return the same information. An MPIT implementation is not allowed to alter it at runtime.

27 The arguments `name` and `name_len` are used to return the name of the control variable as described in Section 1.2.3.

29 **If completed successfully, the routine is required to return a name of at least length one, which is unique with respect to all other names for MPIT control variables used by the MPI implementation.**

32 The argument `verbosity` returns the verbosity level (see Section 1.2.1) assigned by the MPI implementation to the variable.

34 The argument `datatype` returns the MPIT datatype in which the value for this control variable will be returned. The value consists of `count` elements of this type.

36 The arguments `desc` and `desc_len` are used to return a description of the control variable as described in Section 1.2.3.

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

41 **The parameter `assoc` returns the type of MPI resource the variable is associated with (see Section 1.2.2).**

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

46

47

48

MPIT\_CONTROLVAR\_ATTR\_GET\_SCOPE(attributes, scope)

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

```
int MPIT_Controlvar_attr_get_scope(MPIT_Controlvar_attributes *attributes,
                                   int *scope)
```

The scope of a variable determines whether it might be changeable through the MPIT interface and whether changing this variable is a local or a collective operation. On successful return from MPIT\_CONTROLVAR\_ATTR\_GET\_SCOPE, the argument `scope` will be set to one of the constants listed in Table 1.5.

Scope Constant	Description
MPIT_SCOPE_READONLY	read-only, cannot be written
MPIT_SCOPE_LOCAL	may be writeable, writing is a local operation
MPIT_SCOPE_GLOBAL	may be writeable, writing is a global operation

Table 1.5: Scopes for MPIT control variables.

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

### Handle Allocation and Deallocation

Before reading or writing the value of a variable, a user must first allocate a handle for it by binding it to an instance of an MPI resource (see also Section 1.2.2). The type of the resource is returned by a previous call to MPIT\_CONTROLVAR\_GET\_INFO.

MPIT\_CONTROLVAR\_HANDLE\_ALLOCATE(index, resource, handle)

IN	index	index of control variable for which handle is to be allocated
IN	resource	reference to an MPI resource
OUT	handle	allocated handle

```
int MPIT_Controlvar_handle_allocate(int index, MPIT_MPI_Resource resource,
                                    MPIT_Controlvar_handle *handle)
```

The reference to the resource instance passed through the argument `resource` can be generated by converting an MPI resource reference to a generic MPIT resource reference of type `MPIT_MPI_Resource` using the conversions functions described in Section 1.2.2).

A call to this routine, if successfully completed, allocates a handle for the control variable specified by the argument `index` and binds this variable to the instance of an MPI resource passed in the argument `resource`. The type of resource passed into this routine

1 **must match** the type of resources for this variable as returned by a prior call to  
 2 MPIT\_CONTROLVAR\_GET\_INFO.

3  
 4  
 5 MPIT\_CONTROLVAR\_HANDLE\_FREE(handle)

6     INOUT    handle                            handle to be freed

7  
 8  
 9 `int MPIT_Controlvar_handle_free(MPIT_Controlvar_handle *handle)`

10     If a handle is **no** longer needed, a user of MPIT should call  
 11 MPIT\_CONTROLVAR\_HANDLE\_FREE to free the handle and the associated resources. **On**  
 12 **a successful return, MPIT sets the handle to MPIT\_CONTROLVAR\_HANDLE\_NULL.**

13  
 14 Control Variable Access Functions

15  
 16  
 17 MPIT\_CONTROLVAR\_READ(handle, buf)

18     IN        handle                            handle to the control variable to be read

19     OUT       buf                                initial address of storage location for variable value

20  
 21  
 22 `int MPIT_Controlvar_read(MPIT_Controlvar_handle handle, void* buf)`

23  
 24     The MPIT\_CONTROLVAR\_READ queries the value of the control variable identified  
 25 by the argument `handle` and stores the result in the buffer `buf`. The user is responsible  
 26 to ensure that the buffer is of the appropriate size and fits the entire value of the control  
 27 variable (based on the returned type and count from a prior corresponding call to  
 28 **MPIT\_CONTROLVAR\_GET\_INFO**).

29  
 30 MPIT\_CONTROLVAR\_WRITE(handle, buf)

31     IN        handle                            handle to the control variable to be written

32     IN        buf                                initial address of storage location for variable value

33  
 34  
 35 `int MPIT_Controlvar_write(MPIT_Controlvar_handle handle, void* buf)`

36  
 37     The MPIT\_CONTROLVAR\_WRITE sets the value of the control variable identified by  
 38 the argument `handle` to the data stored in the buffer `buf`. The user is responsible to ensure  
 39 that the buffer is of the appropriate size and fits the entire value of the control variable  
 40 (based on the returned type and count from a prior corresponding call to  
 41 **MPIT\_CONTROLVAR\_GET\_INFO**.)

42     If the variable has a global scope (as returned by a prior corresponding  
 43 **MPIT\_CONTROLVAR\_ATTR\_GET\_SCOPE** call), any write call to this variable **must** be  
 44 issued on all **connected MPI processes**. The user is responsible to ensure that the writes in  
 45 all processes are consistent.

46     If it is not possible to change the variable at the time the call is made, the functions  
 47 returns either **MPIT\_ERR\_SETNOTNOW**, if there may be a later time at which the variable



could be set, or `MPIT_ERR_SETNEVER`, if the variable cannot be set for the remainder of the application's execution.

### 1.2.7 Performance Variables

The following section focuses on the ability to list and query performance variables provided by the MPI implementation. Performance variables provide insight into MPI implementation specific internals and can represent information such as the state a component is in, aggregated timing data for submodules, or queue sizes and lengths.

#### Performance Variable Classes

Each reported performance variable is associated with a class of performance variables describing its the basic semantics. The class of a variable also defines its basic behavior, when and how an MPI implementation can change its value and what the initial or starting value of this variable is when it is either used for the first time or reset. Further, it also defines which types can be used to represent it. These classes are defined by the following constants:

- `MPIT_PERFVAR_CLASS_STATE`

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

- `MPIT_PERFVAR_CLASS_RESOURCE_LEVEL`

A performance variable in this class represents a value that describes the utilization level of a resource within the MPI implementation. The value of a variable of this class can change at any time to match the current utilization level of the resource. Values returned from variables in this class are represented by one of the following types: `MPIT_BYTE`, `MPIT_SHORT`, `MPIT_INT`, `MPIT_LONG`, `MPIT_LONG_LONG`, `MPIT_FLOAT` or `MPIT_DOUBLE`. The default starting value is the current utilization level of the resource.

- `MPIT_PERFVAR_CLASS_RESOURCE_PERCENTAGE`

The value of a performance variable in this class represent the percentage utilization of a finite resource in the MPI implementation. The value of a variable of this class can change at any time to match the current utilization level of the resource. It should be returned as an `MPIT_FLOAT` or `MPIT_DOUBLE` type. The value must always be between 0.0 (resource not used at all) and 1.0 (resource completely used). The default starting value is the current percentage utilization level of the resource.

- `MPIT_PERFVAR_CLASS_RESOURCE_HIGHWATERMARK`

A performance variable in this class represents a value that describes the high watermark utilization of a resource within the MPI implementation. The value of a variable of this class is monotonically growing (from the initialization or reset of the variable). It can be represented by one of the following types: `MPIT_BYTE`, `MPIT_SHORT`, `MPIT_INT`, `MPIT_LONG`, `MPIT_LONG_LONG`, `MPIT_FLOAT`

or MPIT\_DOUBLE. The default starting value is the current utilization level of the resource.

- **MPIT\_PERFVAR\_CLASS\_RESOURCE\_LOWWATERMARK**

A performance variable in this class represents a value that describes the low watermark utilization of a resource within the MPI implementation. The value of a variable of this class is monotonically decreasing (from the initialization or reset of the variable). It can be represented by one of the following types: MPIT\_BYTE, MPIT\_SHORT, MPIT\_INT, MPIT\_LONG, MPIT\_LONG\_LONG, MPIT\_FLOAT or MPIT\_DOUBLE. The default starting value is the current utilization level of the resource.

- **MPIT\_PERFVAR\_CLASS\_EVENT\_COUNTER**

A performance variable in this class counts the number of occurrences of a specific event during the execution time of an application (e.g., the number of memory allocations within an MPI library). The value of a variable of this class is monotonically increasing (from the initialization or reset of the performance variable) by one for each specific event that is observed. Values must be non-negative and represented by one of the following types: MPIT\_SHORT, MPIT\_INT, MPIT\_LONG, MPIT\_LONG\_LONG. The default starting value for variables of this class is 0.

- **MPIT\_PERFVAR\_CLASS\_EVENT\_AGGREGATE**

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

- **MPIT\_PERFVAR\_CLASS\_EVENT\_TIMER**

The value of a performance variable in this class represents the aggregated time that the MPI implementation spends executing a particular event. This class has the same basic semantics as MPIT\_PERFVAR\_CLASS\_EVENT\_AGGREGATE, but explicitly records a timing value. The value of a variable of this class is monotonically increasing (from the initialization or reset of the performance variable). It must be non-negative and represented by one of the following types: MPIT\_INT, MPIT\_LONG, MPIT\_LONG\_LONG, MPIT\_FLOAT, MPIT\_DOUBLE. The default starting value for variables of this class is 0.

## Performance Variable Query Functions

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

An MPIT implementation is allowed to increase the number of performance variables during the execution of an MPI application, e.g., when new variables become available

through dynamic loading. However, MPIT implementations are not allowed to change the `index` of a performance variable or delete a variable once it has been added to the set.

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

`MPIT_PERFVAR_GET_NUM(num)`

OUT      `num`                      returns number of performance variables

```
int MPIT_Perfvar_get_num(int *num)
```

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

`MPIT_PERFVAR_GET_INFO(index, name, name_len, verbosity, varclass, datatype, count, desc, desc_len, assoc, attributes)`

IN	<code>index</code>	index of the performance variable to be queried
OUT	<code>name</code>	buffer to return the name of the performance variable
INOUT	<code>name_len</code>	length of the string and/or buffer for name
OUT	<code>verbosity</code>	verbosity level of this variable
OUT	<code>varclass</code>	class of performance variable
OUT	<code>datatype</code>	MPIT type of the information stored in the performance variable
OUT	<code>count</code>	number of elements returned
OUT	<code>desc</code>	buffer to return a description of the performance variable
INOUT	<code>desc_len</code>	length of the string and/or buffer for desc
OUT	<code>assoc</code>	type of MPI resource this variable is associated with
OUT	<code>attributes</code>	additional attributes defining this variable

```
int MPIT_Perfvar_get_info(int num, char *name, int *name_len, int
    *verbosity, int *varclass, MPIT_Datatype *datatype, int
    *count, char *desc, int *desc_len, int *assoc,
    MPIT_Perfvar_attributes *attributes)
```

After a successful call to `MPIT_PERFVAR_GET_INFO` for a particular variable, subsequent calls to this routine querying information about the same variable must return the same information. An MPIT implementation is not allowed to alter it at runtime.

The arguments `name` and `name_len` are used to return the name of the performance variable as described in Section 1.2.3.

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

The argument `verbosity` returns the verbosity level (see Section 1.2.1) assigned by the MPI implementation to the variable.

The class of the performance variable is returned in the parameter `varclass` and can be one of the constants defined in Section 1.2.7.

The argument `datatype` returns the `MPIT` datatype in which the value for this performance variable will be returned. The value consists of `count` elements of this type.

The arguments `desc` and `desc_len` are used to return a description of the control variable as described in Section 1.2.3.

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

The parameter `assoc` returns the type of MPI resource the variable is associated with (see Section 1.2.2).

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

```
MPIT_PERFVAR_ATTR_GET_READONLY(attributes, readonly)
```

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

```
int MPIT_Perfvar_attr_get_readonly(MPIT_Perfvar_attributes *attributes, int
    *readonly)
```

Upon return, the argument `readonly` will be set to null if the variable can be written or reset by the user, or one if the variable is only initialized at `MPIT_INIT` and can only be read after that.

```
MPIT_PERFVAR_ATTR_GET_CONTINUOUS(attributes, continuous)
```

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

```
int MPIT_Perfvar_attr_get_continuous(MPIT_Perfvar_attributes *attributes,
    int *continuous)
```

Upon return, the argument `continuous` will be set to null if the variable can be started and stopped by the user, or one if the variable is automatically active and can not be stopped by the user.

## Performance Experiment Sessions

Within a single program, multiple components can use the `MPIT` interface. To avoid collisions with respect to accesses to performance variables, users of the `MPIT` interface **must first create a session**. All subsequent calls accessing performance variables are then within

the context of this session. Any call executed in a session **must** not influence the results in any other session.

#### MPIT\_PERFVAR\_SESSION\_CREATE(session)

OUT      session                                      identifier of performance experiment session

```
int MPIT_Perfvar_session_create(MPIT_Perfvar_session *session)
```

This call creates a new session for accessing performance variables. An identifier of the current section is returned in `session` using the type `MPIT_Perfvar_session`.

#### MPIT\_PERFVAR\_SESSION\_FREE(session)

INOUT    session                                      identifier of performance experiment session

```
int MPIT_Perfvar_session_free(MPIT_Perfvar_session *session)
```

This call frees an existing session, i.e., calls to MPIT can no longer be made within the context of the freed session. This call also frees all handles that have been allocated within the specified session — see below for handle allocation and freeing. On a successful return, MPIT sets the handle to `MPIT_PERFVAR_SESSION_NULL`.

#### Handle Allocation and Deallocation

Before using a performance variable, a user must first allocate a handle for it by binding it to an instance of an MPI resource (see also Section 1.2.2). The type of the resource is returned by a previous call to `MPIT_PERFVAR_GET_INFO`.

#### MPIT\_PERFVAR\_HANDLE\_ALLOCATE(session, index, resource, handle)

IN        session                                      identifier of performance experiment session  
 IN        index                                        index of performance variable for which handle is to be allocated  
 IN        resource                                    reference to an MPI resource  
 OUT      handle                                        allocated handle

```
int MPIT_Perfvar_handle_allocate(MPIT_Perfvar_session session, int index,
                                  MPIT_MPI_Resource resource, MPIT_Perfvar_handle *handle)
```

A call to this routine, if successfully completed, allocates a handle for the performance variable specified by the argument `index`, binds this variable to the instance of an MPI resource passed in the argument `resource`, and resets the value of the variable to its default value (as specified in Section 1.2.7). The type of resource passed into this routine must match the type of resources for this variable as returned by a prior call to `MPIT_PERFVAR_GET_INFO`.

The reference to the resource instance passed through the argument `resource` can be generated by converting an MPI resource reference to a generic MPIT resource reference of type `MPIT_MPI_Resource` using the conversions functions described in Section 1.2.2).

```
MPIT_PERFVAR_HANDLE_FREE(session,handle)
```

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

```
int MPIT_Perfvar_handle_free(MPIT_Perfvar_session session,
                             MPIT_Perfvar_handle *handle)
```

If a handle is no longer needed, a user of MPIT should call `MPIT_PERFVAR_HANDLE_FREE` to free the handle and the associated resources. On a successful return, MPIT sets the handle to `MPIT_PERFVAR_HANDLE_NULL`.

### Starting and Stopping of Performance Variables

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

```
MPIT_PERFVAR_START(session, handle)
```

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

```
int MPIT_Perfvar_start(MPIT_Perfvar_session session, MPIT_Perfvar_handle
                      handle)
```

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

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

```
MPIT_PERFVAR_STOP(session, handle)
```

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

```
int MPIT_Perfvar_stop(MPIT_Perfvar_session session, MPIT_Perfvar_handle
                    handle)
```

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

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

### Performance Variable Access Functions

```
MPIT_PERFVAR_READ(session, handle, buf)
```

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

```
int MPIT_Perfvar_read(MPIT_Perfvar_session session, MPIT_Perfvar_handle
                    handle, void* buf)
```

The `MPIT_PERFVAR_READ` call queries the value of the performance variable with the handle `handle` in the session `session` and stores the result in the buffer `buf`. The user is responsible to ensure that the buffer is of the appropriate size and fits the entire value of the performance variable (based on the returned type and count during the `MPIT_PERFVAR_GET_INFO` call).

Note that the constant `MPIT_PERFVAR_ALL_HANDLES` can not be used as an argument for the MPIT function `MPIT_PERFVAR_READ`, since this would require the function to return a set of variable values instead of just one.

```
MPIT_PERFVAR_WRITE(session,handle, buf)
```

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

```
int MPIT_Perfvar_write(MPIT_Perfvar_session session, MPIT_Perfvar_handle
                    handle, void* buf)
```

The `MPIT_PERFVAR_WRITE` call attempts to write the value of the performance variable with the handle `handle` in the session `session`. The value to be written is passed in the buffer `buf`. The user is responsible to ensure that the buffer is of the appropriate size and fits the entire value of the performance variable (based on the returned type and count during the `MPIT_PERFVAR_GET_INFO` call).

1 If it is not possible to change the variable the function returns  
 2 MPIT\_ERR\_PERFVAR\_WRITE.

3 Note that the constant MPIT\_PERFVAR\_ALL\_HANDLES can not be used as an argument  
 4 for the MPIT function MPIT\_PERFVAR\_WRITE, since this would require the function to  
 5 accept a set of variable values instead of just one.

6  
 7  
 8 **MPIT\_PERFVAR\_RESET(session, handle)**

9     **IN**        **session**                    identifier of performance experiment session  
 10     **IN**        **handle**                    handle of a performance variable

11  
 12  
 13 **int MPIT\_Perfvar\_reset(MPIT\_Perfvar\_session session, MPIT\_Perfvar\_handle**  
 14                                    **handle)**

15 The MPIT\_PERFVAR\_RESET call sets of the performance variable with the handle  
 16 handle to its default starting value (as specified in Section 1.2.7). If it is not possible to  
 17 change the variable the function returns MPIT\_ERR\_PERFVAR\_WRITE.

18 If the constant MPIT\_PERFVAR\_ALL\_HANDLES is passed in handle, the MPI implementa-  
 19 tion attempts to reset all variables within the session identified by session for which handles  
 20 have been allocated. In this case, the routine returns MPIT\_SUCCESS if all variables are reset  
 21 successfully, otherwise MPIT\_ERR\_NOWRITE is returned. Readonly variables are ignored  
 22 when used with MPIT\_PERFVAR\_ALL\_HANDLES .

23  
 24  
 25 **MPIT\_PERFVAR\_READRESET(session, handle, buf)**

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

29  
 30  
 31 **int MPIT\_Perfvar\_readreset(MPIT\_Perfvar\_session session,**  
 32                                    **MPIT\_Perfvar\_handle handle, void\* buf)**

33 The MPIT\_PERFVAR\_READRESET call atomically queries the value of the performance  
 34 variable, stores the result in the buffer buf, and then sets the value of the performance  
 35 variable to its default starting value (as specified in Section 1.2.7). The user is responsible  
 36 to ensure that the buffer is of the appropriate size and fits the entire value of the performance  
 37 variable (based on the returned type and count during the query call). If it is not possible  
 38 to change the variable the function returns MPIT\_ERR\_PERFVAR\_WRITE. In this case, the  
 39 value returned in buf is the same as if the variable would have been read by the  
 40 MPIT\_PERFVAR\_READ call.

41 Note that the constant MPIT\_PERFVAR\_ALL\_HANDLES can not be used as an argument  
 42 for the MPIT function MPIT\_PERFVAR\_READRESET, since this would require the function  
 43 to return a set of variable values instead of just one.

44  
 45 *Advice to implementors.* Although MPI places no requirements on the interaction  
 46 with external mechanisms such as signal handlers, it is strongly recommended that all  
 47 routines to start, stop, read, write, and reset performance variables should be safe to  
 48



call in asynchronous contexts. Examples of asynchronous contexts include signal handlers and interrupt handlers. Such safety permits the development of sampling-based tools. High quality implementations should strive to make the results of any such interactions intuitive to users, and attempt to document restrictions where deemed necessary. (*End of advice to implementors.*)

### 1.2.8 Variable Categorization

MPI implementations can optionally group performance and control variables into categories to express logical relationships between various variables. **Categories can also contain other categories to form a hierarchical grouping. Categories can never include themselves either directly or transitively within other included categories.**

*Rationale.* **The ability to include categories in other categories enables the creation of a hierarchical grouping of variables. The restriction that categories can not include themselves directly or transitively guarantees that this structure is strictly hierarchical and does not contain any loops.** (*End of rationale.*)

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

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 MPIT interface. If  $N = 0$ , then the MPI implementation does not export any categories. **This index number is used by MPIT by subsequent calls to identify the individual categories.**

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

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

**MPIT\_CATEGORY\_GET\_NUM(num)**

**OUT     num                             current number of categories**

**int MPIT\_Category\_get\_num(int \*num)**

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

```

1 MPIT_CATEGORY_GET_INFO(index, name, name_len, desc, desc_len, num_controlvars, num_perfvars,
2 num_categories)

```

3	IN	index	index of the category to be queried, in the range [0, $N-1$ ]
4			
5			
6	OUT	name	buffer to return the name of the category
7	INOUT	name_len	length of the string and/or buffer for name
8	OUT	desc	buffer to return the description of the category
9			
10	INOUT	desc_len	length of the string and/or buffer for desc
11	OUT	num_controlvars	number of control variables in the category
12	OUT	num_perfvars	number of performance variables in the category
13			
14	OUT	num_categories	number of MPIT categories contained in the category

```

15
16 int MPIT_Category_get_info(int index, char *name, int *name_len, char
17 *desc, int *desc_len, int *num_controlvars, int
18 *num_perfvars, int *num_categories)

```

19 The arguments `name` and `name_len` are used to return the name of the category as  
20 described in Section 1.2.3.

21 If completed successfully, the routine is required to return a name of at least length  
22 one, which is unique with respect to all other names for MPIT categories used by the MPIT  
23 implementation.

24 The arguments `desc` and `desc_len` are used to return the description of the category as  
25 described in Section 1.2.3.

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

29 On successful completion, the function returns the number of control variables (  
30 `num_controlvars`), performance variables (`num_perfvars`) and other categories (  
31 `num_categories`) contained in the queried category.

```

32
33
34 MPIT_CATEGORY_GET_CONTROLVARS(cat_index, len, indices)

```

35	IN	cat_index	index of the category to be queried, in the range [0, $N-1$ ]
36			
37			
38	IN	len	the length of the kinds and indices arrays
39	OUT	indices	an integer array of size <code>len</code> , indicating variable indices

```

40
41 int MPIT_Category_get_controlvars(int cat_index, int len, int indices[])

```

42 `MPIT_CATEGORY_GET_CONTROLVARS` can be used to query which control variables  
43 contained in a particular category. A category may contain zero or more control variables.  
44

45  
46  
47  
48

`MPIT_CATEGORY_GET_PERFVARS(cat_index,len,indices)`

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

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

`MPIT_CATEGORY_GET_PERFVARS` can be used to query which performance variables contained in a particular category. A category may contain zero or more performance variables.

`MPIT_CATEGORY_GET_CATEGORIES(cat_index,len,indices)`

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

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

`MPIT_CATEGORY_GET_CATEGORIES` can be used to query which other categories contained in a particular category. A category may contain zero or more other categories.

The index values returned in `indices` by `MPIT_CATEGORY_GET_CONTROLVARS`, `MPIT_CATEGORY_GET_PERFVARS` or `MPIT_CATEGORY_GET_CATEGORIES` can be used as input to `MPIT_CONTROLVAR_GET_INFO`, `MPIT_PERFVAR_GET_INFO` or `MPIT_CATEGORY_GET_INFO` respectively.

The user is responsible for allocating the arrays passed into the `functions` `MPIT_CATEGORY_GET_CONTROLVARS`, `MPIT_CATEGORY_GET_PERFVARS` and `MPIT_CATEGORY_GET_CATEGORIES`.

The `functions` will only write up to `len` elements into the `respective array`. If the category contains more than `len` variables or other categories `respectively` the function returns an arbitrary subset; if it contains less than `len` variables or other categories `respectively`, all will be returned and the remaining array entries will not be modified.

### 1.2.9 Return and Error Codes

All MPIT functions return a return or error code. The `constants in Table 1.6 are defined for this purpose`. None of the error codes returned by an MPIT routine is fatal to the overall MPI implementation or `invokes` an MPI error handler. In any case, the execution of the MPI program `continues` as if the call would have succeeded. However, the MPIT implementation is not required to check all user provided parameters; if a user passes illegal parameter values to any MPIT routine that are not caught by the implementation, the behavior of the implementation is undefined.

### 1.2.10 Profiling Interface

All requirements for the profiling interfaces, as described in Section ??, also apply to the MPIT interface. In particular, this means that a complying MPI implementation **must** provide matching PMPIT calls for every MPIT call. All rules, guidelines, and recommendations from Section ?? apply equally to PMPIT calls.

Return Code	Description
Return Codes for all MPIT Functions	
MPIT_SUCCESS	No error, call completed
MPIT_ERR_MEMORY	Out of memory
MPIT_ERR_NOTINITIALIZED	MPIT not initialized
MPIT_ERR_CANTINIT	MPIT not in the state to be initialized
Return Codes for Type Functions: MPIT_TYPE_*	
MPIT_ERR_PREDEFINED	Datatype is a predefined type and not an enumeration
MPIT_ERR_INVALIDTYPE	Datatype is not a valid datatype
MPIT_ERR_INVALIDITEM	The item index queried is out of range (for MPIT_TYPE_ENUMITEM only)
Return Codes for variable and category query functions: MPIT_*_GET_INFO	
MPIT_ERR_INVALIDINDEX	The variable or category index is invalid
Return Codes for Handle Functions: MPIT_*_ALLOCATE,FREE	
MPIT_ERR_INVALIDINDEX	The variable index is invalid
MPIT_ERR_INVALIDHANDLE	The handle is invalid
MPIT_ERR_OUTOFHANDLES	No more handles available
Return Codes for Session Functions: MPIT_PERFVAR_SESSION_*	
MPIT_ERR_OUTOFSESSIONS	No more sessions available
MPIT_ERR_INVALIDSESSION	Session argument is not a valid session
Return Codes for Control Variable Access Functions: MPIT_CONTROLVAR_READ,WRITE	
MPIT_ERR_SETNOTNOW	Variable cannot be set at this moment
MPIT_ERR_SETNEVER	Variable cannot be set until end of execution
MPIT_ERR_INVALIDDVAR	Control variable does not exist
MPIT_ERR_INVALIDHANDLE	The handle is invalid
Return Codes for Performance Variable Access and Control: MPIT_PERFVAR_START,STOP,READ,WRITE,RESET,READRESET	
MPIT_ERR_INVALIDHANDLE	The handle is invalid
MPIT_ERR_INVALIDSESSION	Session argument is not a valid session
MPIT_ERR_NOSTARTSTOP	Variable can not be started or stopped for MPIT_PERFVAR_START and MPIT_PERFVAR_STOP
MPIT_ERR_NOWRITE	Variable can not be written or reset for MPIT_PERFVAR_WRITE and MPIT_PERFVAR_RESET
Return Codes for Category Functions: MPIT_CATEGORY_*	
MPIT_ERR_INVALIDCATEGORY	The specified category index does not exist

Table 1.6: Return and error codes used MPIT functions.

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

# Bibliography

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

# MPIT Constant and Predefined Handle Index

This index lists predefined MPIT constants and handles.

MPIT\_BYTE, 7  
 MPIT\_CHAR, 7  
 MPIT\_CONTROLVAR\_HANDLE\_NULL, 12  
 MPIT\_DOUBLE, 7  
 MPIT\_ERR\_CANTINIT, 25  
 MPIT\_ERR\_INVALIDCATEGORY, 25  
 MPIT\_ERR\_INVALIDHANDLE, 25  
 MPIT\_ERR\_INVALIDINDEX, 25  
 MPIT\_ERR\_INVALIDITEM, 25  
 MPIT\_ERR\_INVALIDSESSION, 25  
 MPIT\_ERR\_INVALIDTYPE, 25  
 MPIT\_ERR\_INVALIDVAR, 25  
 MPIT\_ERR\_MEMORY, 25  
 MPIT\_ERR\_NOSTARTSTOP, 18, 19, 25  
 MPIT\_ERR\_NOTINITIALIZED, 25  
 MPIT\_ERR\_NOWRITE, 20, 25  
 MPIT\_ERR\_OUTOFHANDLES, 25  
 MPIT\_ERR\_OUTOFSESSIONS, 25  
 MPIT\_ERR\_PERFVAR\_WRITE, 20  
 MPIT\_ERR\_PREDEFINED, 25  
 MPIT\_ERR\_SETNEVER, 13, 25  
 MPIT\_ERR\_SETNOTNOW, 12, 25  
 MPIT\_FLOAT, 7  
 MPIT\_INT, 7  
 MPIT\_LOGICAL, 7  
 MPIT\_LONG, 7  
 MPIT\_LONG\_LONG, 7  
 MPIT\_MPI\_RESOURCE\_GLOBAL, 5  
 MPIT\_MPI\_RESOURCE\_TYPE\_COMMUNICATION, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_DATATYPE, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_ERRORHANDLING, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_FILE, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_GLOBAL, 3, 5  
 MPIT\_MPI\_RESOURCE\_TYPE\_GROUP, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_OPERATOR, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_REQUEST, 3  
 MPIT\_MPI\_RESOURCE\_TYPE\_WINDOW, 3  
 MPIT\_PERFVAR\_ALL\_HANDLES, 18–20  
 MPIT\_PERFVAR\_CLASS\_EVENT\_AGGREGATE, 14  
 MPIT\_PERFVAR\_CLASS\_EVENT\_COUNTER, 14  
 MPIT\_PERFVAR\_CLASS\_EVENT\_TIMER, 14  
 MPIT\_PERFVAR\_CLASS\_RESOURCE\_HIGHWATERMARK, 13  
 MPIT\_PERFVAR\_CLASS\_RESOURCE\_LEVEL, 13  
 MPIT\_PERFVAR\_CLASS\_RESOURCE\_LOWWATERMARK, 14  
 MPIT\_PERFVAR\_CLASS\_RESOURCE\_PERCENTAGE, 13  
 MPIT\_PERFVAR\_CLASS\_STATE, 13  
 MPIT\_PERFVAR\_HANDLE\_NULL, 18  
 MPIT\_PERFVAR\_SESSION\_NULL, 17  
 MPIT\_SCOPE\_GLOBAL, 11  
 MPIT\_SCOPE\_LOCAL, 11  
 MPIT\_SCOPE\_READONLY, 11  
 MPIT\_SHORT, 7  
 MPIT\_SUCCESS, 2, 20, 25

1	MPIT_TYPECLASS_ENUMERATION, 7
2	MPIT_TYPECLASS_PREDEFINED, 7
3	MPIT_VERBOSITY_MPIDEV_BASIC, 2
4	MPIT_VERBOSITY_MPIDEV_DETAILED,
5	2
6	MPIT_VERBOSITY_MPIDEV_VERBOSE,
7	2
8	MPIT_VERBOSITY_TUNER_BASIC, 2
9	MPIT_VERBOSITY_TUNER_DETAILED,
10	2
11	MPIT_VERBOSITY_TUNER_VERBOSE,
12	2
13	MPIT_VERBOSITY_USER_BASIC, 2
14	MPIT_VERBOSITY_USER_DETAILED,
15	2
16	MPIT_VERBOSITY_USER_VERBOSE,
17	2
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	



# MPIT Function Index

The underlined page numbers refer to the function definitions.

MPIT\_CATEGORY\_GET\_CATEGORIES, [23](#), [23](#)  
MPIT\_CATEGORY\_GET\_CONTROLVARS, [22](#), [22](#), [23](#)  
MPIT\_CATEGORY\_GET\_INFO, [22](#), [23](#)  
MPIT\_CATEGORY\_GET\_NUM, [21](#)  
MPIT\_CATEGORY\_GET\_PERFVAR, [23](#), [23](#)  
MPIT\_CONTROLVAR\_ATTR\_GET\_SCOPE, [11](#), [11](#), [12](#)  
MPIT\_CONTROLVAR\_GET\_INFO, [9](#), [10](#), [10](#), [11](#), [12](#), [23](#)  
MPIT\_CONTROLVAR\_GET\_NUM, [9](#)  
MPIT\_CONTROLVAR\_HANDLE\_ALLOCATE, [11](#)  
MPIT\_CONTROLVAR\_HANDLE\_FREE, [12](#), [12](#)  
MPIT\_CONTROLVAR\_READ, [12](#), [12](#)  
MPIT\_CONTROLVAR\_WRITE, [12](#), [12](#)  
MPIT\_FINALIZE, [6](#), [6](#)  
MPIT\_INIT, [6](#), [6](#), [16](#)  
MPIT\_MPI\_RESOURCE\_COMMUNICATOR, [4](#)  
MPIT\_MPI\_RESOURCE\_DATATYPE, [4](#)  
MPIT\_MPI\_RESOURCE\_ERRORHANDLER, [4](#)  
MPIT\_MPI\_RESOURCE\_FILE, [4](#)  
MPIT\_MPI\_RESOURCE\_GROUP, [4](#)  
MPIT\_MPI\_RESOURCE\_OPERATIONS, [5](#)  
MPIT\_MPI\_RESOURCE\_REQUEST, [5](#)  
MPIT\_MPI\_RESOURCE\_WINDOW, [5](#)  
MPIT\_PERFVAR\_ATTR\_GET\_CONTINUOUS, [16](#)  
MPIT\_PERFVAR\_ATTR\_GET\_READONLY, [16](#)  
MPIT\_PERFVAR\_GET\_INFO, [15](#), [15](#), [17](#), [19](#), [23](#)  
MPIT\_PERFVAR\_GET\_NUM, [15](#)  
MPIT\_PERFVAR\_HANDLE\_ALLOCATE, [17](#)  
MPIT\_PERFVAR\_HANDLE\_FREE, [18](#), [18](#)  
MPIT\_PERFVAR\_READ, [19](#), [19](#), [20](#)  
MPIT\_PERFVAR\_READRESET, [20](#), [20](#)  
MPIT\_PERFVAR\_RESET, [20](#), [20](#), [25](#)  
MPIT\_PERFVAR\_SESSION\_CREATE, [17](#)  
MPIT\_PERFVAR\_SESSION\_FREE, [17](#)  
MPIT\_PERFVAR\_START, [18](#), [25](#)  
MPIT\_PERFVAR\_STOP, [18](#), [25](#)  
MPIT\_PERFVAR\_WRITE, [19](#), [19](#), [20](#), [25](#)  
MPIT\_TYPE\_ENUM\_GET\_INFO, [8](#)  
MPIT\_TYPE\_ENUM\_GET\_ITEM, [8](#), [8](#)  
MPIT\_TYPE\_ENUM\_QUERY, [8](#)  
MPIT\_TYPE\_ENUMITEM, [25](#)  
MPIT\_TYPE\_GET\_CLASS, [7](#)  
MPIT\_TYPE\_GET\_SIZE, [7](#), [8](#)  
MPIT\_FINALIZE, [6](#)

