

MPI: A Message-Passing Interface Standard

Version 3.0

⊤ (Fin2)

⊥ (Fin2)

Message Passing Interface Forum

Draft October 8, 2010

Contents

1	Tool Interfaces for MPI	1
1.1	Introduction	1
1.2	MPIT Performance Interface	1
1.2.1	Initialization and Finalization	2
1.2.2	Type System	3
1.2.3	Verbosity Levels	5
1.2.4	Control Variables	6
	Control Variable Query Functions	6
	Control Variable Access Functions	8
1.2.5	Performance Variables	9
	Performance Variable Classes	9
	Performance Variable Query Functions	10
	Performance Experiment Sessions	12
	Performance Variable Activation	13
	Starting and Stopping of Performance Variables	13
	Performance Variable Access Functions	14
1.2.6	Performance and Control Variable Taxonomic Information	16
1.2.7	Return and Error Codes	18
	Return Codes for Type Functions	19
	Return Codes for Control Variable Access Functions	19
	Return Codes for Performance Variable Access and Control	19
	Return Codes for Taxonomy Functions	19
1.2.8	Profiling Interface	19
	Bibliography	21
	Examples Index	22
	MPI Constant and Predefined Handle Index	22
	MPI Declarations Index	23
	MPI Callback Function Prototype Index	23
	MPI Function Index	23

List of Figures

List of Tables

1.1	MPIT datatypes and their MPI equivalences.	3
1.2	MPIT type classes.	5
1.3	MPIT verbosity levels.	5
1.4	Scopes for MPIT control variables.	8
1.5	Return codes used by any MPIT function.	19
1.6	Return codes used by MPIT type functions.	19
1.7	Return codes used by MPIT control variable access functions.	19
1.8	Return codes used by MPIT performance variable access, start, stop, or activation functions.	20
1.9	Return codes used MPIT taxonomy functions.	20

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. This includes a profiling interface (Section ??), PMPI, to transparently intercept and inspect any MPI call; and an information interface (Section 1.2), MPIT, 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. Additionally, several other tool interfaces exist that define interfaces that are primarily intended to be used from external processes. An example for the latter is the MPIR process acquisition interface, which is used by debuggers and performance analysis tools to detect and locate all MPI processes belonging to a given job. Currently, these interfaces are not included in MPI standard, but rather described in MPI forum white papers, which are published on the MPI forum's website.

1.2 MPIT Performance Interface

Open questions / Todos:

- Versioning - should this be part of MPI or MPIT
- Change the get info calls to use structs
- String returns in Taxonomy section
- Iterators in Taxonomy section

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

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

1 from the MPI functionality. This is particularly true for the initialization and finalization of
 2 MPIT, which is provided through a separate set of routines. However, all conventions and
 3 principles governing the MPI API also apply to the MPIT interface and the MPIT interface
 4 shall be defined in the same header or API definition file(s) as the regular MPI routines
 5 (e.g., *mpi.h* where appropriate).

6 The interface is split into two parts: the first part provides information about control
 7 variables used by the MPI library to fine tune its performance. The second part provides ac-
 8 cess to performance variables that can provide insight into internal performance information
 9 of the underlying MPI implementation.

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

14 On success all MPIT routines return MPIT_SUCCESS, otherwise they return an appro-
 15 priate error code. Details on error codes can be found in Section 1.2.7. However, errors
 16 returned by the MPIT interface shall not be fatal nor have any impact on the execution of
 17 MPI routines.

18
 19 *Advice to users.* The number and type of control variables and performance variables
 20 can vary between MPI libraries, platforms, and even different builds of the same library
 21 on the same platform. Hence, any application relying on a particular variable will no
 22 longer be portable.

23 This interface is primarily intended for performance monitoring tools, as well as sup-
 24 port tools and libraries controlling the application's environment. Application pro-
 25 grammers should either avoid using it and avoid being dependent on the existence of
 26 a particular control or performance variable. (*End of advice to users.*)
 27

28 1.2.1 Initialization and Finalization

29
 30 Since the MPIT interface is implemented in a separate name space and hence is independent
 31 of the core MPI functions, it requires a separate set of initialization and finalization routines.

32
 33 MPIT_INIT()

34
 35 `int MPIT_Init()`
 36

37 All programs or tools that use the MPIT interface must initialize the MPIT interface
 38 before calling any MPIT routine. The only exception to this rule is that the function
 39 MPIT_INITCOUNT can be called at any time.

40 A user can initialize the MPIT interface by calling MPIT_INIT, which can be called
 41 multiple times.
 42

43
 44 MPIT_FINALIZE()

45
 46 `int MPIT_Finalize()`
 47

48 This routine finalizes the use of the MPIT interface and may be called as often as the
 corresponding MPIT_INIT 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 and invalidate all open sessions (for the concept of Sessions see Section 1.2.5).

`MPIT_INITCOUNT(num)`

OUT num returns the number of times `MPIT_INIT` has been called minus the times `MPIT_Finalize` has been called up to the current point of execution.

`int MPIT_Initcount(int *num)`

Note, that the MPIT functions are independent of the MPI functions. This means that MPIT functions can be called before `MPI_INIT` and after `MPI_FINALIZE`.

1.2.2 Type System

The MPIT interface provides its own type system. All types are represented by a variable or constant of type `MPIT_Datatype`. The Table 1.1 lists all available constants that can be used to identify a type for MPIT calls.

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

Table 1.1: MPIT datatypes and their MPI equivalences.

Conforming implementations of MPIT have to ensure that the MPIT types are equivalent to the listed MPI datatypes for any section of the code for which both MPI and MPIT have been initialized. In particular, this requires that the size of variables of these types are equal and that it is possible to send and receive data of a particular MPIT type with regular MPI operations using the equivalent MPI type.

In addition to the predefined datatypes listed in the table, an MPI implementation may provide an additional set of enumeration datatypes to describe variables with a fixed set of discrete values. These types are represented through integer variables and have `MPI_INT` as

1 their equivalent MPI type. Their values range from 0 to $N - 1$, with a fixed N that can be
 2 queried using `MPIT_TYPE_ENUMQUERY`.

3
 4
 5 `MPIT_TYPE_ENUMQUERY(datatype,size,name,name_len)`

6	IN	datatype	MPIT datatype to be queried
7	OUT	size	number of elements representable with this enumera-
8			tion datatype
9	OUT	name	buffer to return the name of the type
10	INOUT	name_len	length of the string and/or buffer for name

11
 12
 13 `int MPIT_Type_Enumquery(MPIT_Datatype datatype, int *size, char *name, int`
 14 `*name_len)`

15
 16 This routine returns, if `datatype` represents a valid enumeration type, the size of the
 17 enumeration as well as a name for it.

18 The argument `name` provides a buffer to return the string describing the name of the
 19 type. The user has to pass the size of the buffer as the `name_len` argument. On return, the
 20 function deposits at most `name_len-1` characters of the requested string into the buffer `name`
 21 followed by a terminating zero character. Additionally, the function writes the length of the
 22 returned string (including the terminating zero character) into `name_len`. If the returned
 23 value is smaller than the argument supplied to the function, the string has been truncated
 24 due to insufficient buffer resources. If the user passes `NULL` as the buffer argument or
 25 passes `-1` as `name_len`, the function does not return the string and only returns the length
 26 of the string in `name_len`.

27 Names for the individual items in each enumeration can be queried using
 28 `MPIT_TYPE_ENUMITEM`.

29
 30
 31 `MPIT_TYPE_ENUMITEM(datatype,item,name,name_len)`

32	IN	datatype	MPIT datatype to be queried
33	IN	item	item number in the MPIT datatype to be queried
34	OUT	name	buffer to return the name of the enumeration item
35	INOUT	name_len	length of the string and/or buffer for name

36
 37
 38 `int MPIT_Type_Enumitem(MPIT_Datatype datatype, int item, char *name, int`
 39 `*name_len)`

40
 41 The argument `name` provides a buffer to return the string describing the name of the
 42 enumeration item. The user has to pass the size of the buffer as the `name_len` argument.
 43 On return, the function deposits at most `name_len-1` characters of the requested string into
 44 the buffer `name` followed by a terminating zero character. Additionally, the function writes
 45 the length of the returned string (including the terminating zero character) into `name_len`.
 46 If the returned value is smaller than the argument supplied to the function, the string has
 47 been truncated due to insufficient buffer resources. If the user passes `NULL` as the buffer
 48

argument or passes -1 as `name_len`, the function does not return the string and only returns the length of the string in `name_len`.

`MPIT_TYPE_GETCLASS(datatype, typeclass)`

IN	<code>datatype</code>	MPIT datatype to be queried
OUT	<code>typeclass</code>	Class of the type passed in

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

This routine returns the class of the type for the datatype provided. This allows users of MPIT to distinguish whether a datatype used is an enumeration type or is one of the predefined types listed above. On return, the `typeclass` argument is set to one of the following constants:

<code>MPIT_TYPECLASS_PREDEFINED</code>	the datatype is a predefined datatype
<code>MPIT_TYPECLASS_ENUMERATION</code>	the datatype is an enumeration datatype

Table 1.2: MPIT type classes.

1.2.3 Verbosity Levels

The MPIT interface provides users access to internal performance data through a set of control and performance variables, which are defined by the MPI implementation. Since the number of variables can be large for particular implementations, every variable exported by the MPIT interface has to be associated with one of the following verbosity levels.

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

Table 1.3: MPIT verbosity levels.

The classification into several verbosity classes is optional for MPI implementations. Alternatively, all variables can be assigned to a single verbosity level. In this case it is recommended to assign all variables to the level `MPIT_VERBOSITY_USER_BASIC`.

However, MPI implementations using verbosity levels should first classify all variables according to the intended target audience (end user, performance optimization, or MPI developer) and then distinguish three level of verbosity (basic, detailed, and verbose) within each class.

1.2.4 Control Variables

The first set of routines in the MPIT interface focuses on the ability to list, query, and possibly set all 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 library. On UNIX systems, such variables can often be set using environment variables, although many other configurations mechanisms might be used (e.g., configuration files, 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 MPI implementation does not export any control variables, otherwise the provided control variables are numbered from 1 to N . An MPI 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, MPI implementations are not allowed to change the number of a control variable or delete it once it has been added to the set.

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

```
MPIT_CTRLVAR_GETNUM(num)
```

```
OUT    num                returns number of control variables
```

```
int MPIT_CTRLVAR_Getum(int *num)
```

The name of individual variables (with numbers between 1 and N acquired by calling MPIT_CTRLVAR_GETNUM) can then be queried with the following function along with any associated information.

```
MPIT_CTRLVAR_GETINFO(num, name, name_len, verbosity, datatype, count, desc, desc_len,
scope, comm)
```

IN	num	number of the control variable to be queried	1
OUT	len	buffer to return the name of the control variable	2
INOUT	len_len	length of the string and/or buffer for len	3
OUT	verbosity	verbosity level of this variable	4
OUT	datatype	MPIT type of the information stored in the control variable	5
OUT	count	number of elements returned	6
OUT	desc	buffer to return a description of the control variable	7
INOUT	desc_len	length of the string and/or buffer for desc	8
OUT	scope	scope of when changes to this variable are possible	9
OUT	comm	communicator that collective write operations to this variable have to be executed on	10

```
int MPIT_Ctrlvar_Getinfo(int num, char *name, int *name_len, int
    *verbosity, MPIT_Datatype *datatype, int *count, char *desc,
    int *desc_len, int *scope, MPI_Comm *comm)
```

The argument `name` provides a buffer to return the string describing the name of the control variable. The user has to pass the size of the buffer as the `name_len` argument. On return, the function deposits at most `name_len-1` characters of the requested string into the buffer `name` followed by a terminating zero character. Additionally, the function writes the length of the returned string (including the terminating zero character) into `name_len`. If the returned value is smaller than the argument supplied to the function, the string has been truncated due to insufficient buffer resources. If the user passes `NULL` as the buffer argument or passes `-1` as `name_len`, the function does not return the string and only returns the length of the string in `name_len`.

The argument `verbosity` returns the verbosity level (see Section 1.2.3) assigned by the MPI implementation to the variable.

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

The argument `desc` provides a buffer to return the string describing a description of the control variable. The user has to pass the size of the buffer as the `desc_len` argument. On return, the function deposits at most `desc_len-1` characters of the requested string into the buffer `desc` followed by a terminating zero character. Additionally, the function writes the length of the returned string (including the terminating zero character) into `desc_len`. If the returned value is smaller than the argument supplied to the function, the string has been truncated due to insufficient buffer resources. If the user passes `NULL` as the buffer argument or passes `-1` as `desc_len`, the function does not return the string and only returns the length of the string in `desc_len`.

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

The 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 return from `MPIT_CTRLVAR_GETINFO` it will be set to one of the constants listed in Table 1.4. If setting this variable requires a collective operation, the communicator on which this collective operation has to be executed, is returned as `comm`. If such an operation is not collective, the implementation should return `MPI_COMM_SELF`.

Scope Constant	Description
<code>MPIT_SCOPE_READONLY</code>	only read-only, cannot be written
<code>MPIT_SCOPE_LOCAL</code>	may be writeable, writing is not a collective operation
<code>MPIT_SCOPE_GLOBAL</code>	may be writeable, writing is a collective operation

Table 1.4: Scopes for MPIT control variables.

Note that the scope of a variable only indicates when a variable might be changeable; it is not a guarantee that can be changed at any time. If it can not be changed at a time the user tries to set it, the MPIT implementation is allowed to return an error code as the result of the write operation.

Control Variable Access Functions

`MPIT_CTRLVAR_READ(num, buf)`

IN	<code>num</code>	number of control variable to be read
OUT	<code>buf</code>	initial address of storage location for variable value

```
int MPIT_Ctrlvar_Read(int num, void* buf)
```

The `MPIT_CTRLVAR_READ` queries the value of the control variable with the number `num` 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 control variable (based on the returned type and count during the `MPIT_CTRLVAR_GETINFO` call).

`MPIT_CTRLVAR_WRITE(num, buf, comm)`

IN	<code>num</code>	number of control variable to be read
IN	<code>buf</code>	initial address of storage location for variable value
IN	<code>comm</code>	communicator for which this operation is collective on

```
int MPIT_Ctrlvar_Write(int num, void* buf, MPI_Comm comm)
```

The `MPIT_CTRLVAR_WRITE` sets the value of the control variable with the number `num` to the data stored 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 control variable (based on the returned type and count during the query `MPIT_CTRLVAR_GETINFO` call).

The operation is collective with respect to the communicator `comm`. The user is responsible that the right communicator, i.e., the one returned by `MPIT_CTRLVAR_GETINFO`, is

passed as the `comm` argument and that this operation is called as a collective operation on all processes in the communicator. The same ordering constraints as for MPI collectives apply. If this operation is local and not collective, the user is required to pass `MPI_COMM_SELF`.

If it is not possible to change the variable at the time the call is made, the functions returns either `MPIT_ERR_SETNOW`, if there could 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 time.

1.2.5 Performance Variables

The second set of functions included in the MPIT interface 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 like 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, which describes the basic semantics of the variable. 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 library or a component of the MPI library is in. The value of this kind of variable can change at any time to any value within the type definition. Variables of this class are expected to be represented by an enumeration type. Variables of this class don't have a default starting value, since the variable reflects a current state of the library.

- `MPIT_PERFVAR_CLASS_UTILIZATION`

The value of a performance variable in this class represent the percentage utilization of a finite resource in the MPI library. The value of this kind of variable can change at any time and should be returned as a `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). Variables of this class don't have a default starting value, since the variable reflects a current state of the library.

- `MPIT_PERFVAR_CLASS_RESOURCE`

A performance variable in this class represents a value that describes the absolute utilization level of a resource within the MPI library. The value of this kind of variable can change at any time and values returned from variables in this class must be non-negative and are represented by one of the following types: `MPIT_BYTE`, `MPIT_SHORT`, `MPIT_INT`, `MPIT_LONG`, `MPIT_LONG_LONG`, `MPIT_FLOAT` or `MPIT_DOUBLE`. Variables of this class don't have a default starting value, since the variable reflects a current state of the library.

- `MPIT_PERFVAR_CLASS_HIGHWATERMARK`

A performance variable in this class represents a value that describes the high watermark absolute utilization of a resource within the MPI library. The value of this kind of variable is monotonically growing (from the initialization or reset of the variable). It

1 must be non-negative and represented by one of the following types: MPIT_BYTE,
 2 MPIT_SHORT, MPIT_INT, MPIT_LONG, MPIT_LONG_LONG, MPIT_FLOAT
 3 or MPIT_DOUBLE. The default starting value for variables of this class is the cur-
 4 rent absolute utilization of the resource.

5
 6 • MPIT_PERFVAR_CLASS_LOWWATERMARK

7 A performance variable in this class represents a value that describes the low water-
 8 mark absolute utilization of a resource within the MPI library. The value of this kind
 9 of variable is monotonically shrinking (from the initialization or reset of the variable).
 10 It must be non-negative and represented by one of the following types: MPIT_BYTE,
 11 MPIT_SHORT, MPIT_INT, MPIT_LONG, MPIT_LONG_LONG, MPIT_FLOAT
 12 or MPIT_DOUBLE. The default starting value for variables of this class is the cur-
 13 rent absolute utilization of the resource.

14
 15 • MPIT_PERFVAR_CLASS_COUNTER

16 A performance variable in this class counts the number of occurrences of a specific
 17 event during the execution time of an application. The value of this kind of variable is
 18 monotonically increasing (from the initialization or reset of the performance variable).
 19 It must be non-negative and represented by one of the following types: MPIT_SHORT,
 20 MPIT_INT, MPIT_LONG, MPIT_LONG_LONG. The default starting value for vari-
 21 ables of this class is 0.

22
 23 • MPIT_PERFVAR_CLASS_AGGREGATE

24 The value of a performance variable in this class is an an aggregated value of over time.
 25 This class is similar to the counter class, but instead of counting individual events, the
 26 value can be incremented by arbitrary amounts. The value of this kind of variable is
 27 monotonically increasing (from the initialization or reset of the performance variable).
 28 It must be non-negative and represented by one of the following types: MPIT_SHORT,
 29 MPIT_INT, MPIT_LONG, MPIT_LONG_LONG, MPIT_FLOAT, MPI_DOUBLE.
 The default starting value for variables of this class is 0.

30
 31 • MPIT_PERFVAR_CLASS_TIMER

32 The value of a performance variable in this class represents the aggregated time that
 33 the MPI library spends executing a particular event. The value of this kind of vari-
 34 able is monotonically increasing (from the initialization or reset of the performance
 35 variable). It must be non-negative and represented by one of the following types:
 36 MPIT_INT, MPIT_LONG, MPIT_LONG_LONG, MPIT_FLOAT, MPIT_DOUBLE.
 37 The default starting value for variables if this class is 0.

38
 39 Performance Variable Query Functions

40 Each MPI implementation exports a set of N performance variables through MPIT. If N is
 41 zero, then the MPI implementation does not export any performance variables, otherwise
 42 the provided performance variables are numbered from 1 to N . An MPI implementation
 43 is allowed to increase the number of performance variables during the execution of an MPI
 44 application, e.g., when new variables become available through dynamic loading. However,
 45 MPI implementations are not allowed to change the number of a performance variable or
 46 delete it once it has been added to the set.

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

MPIT_PERFVAR_GETNUM(num) 1

OUT num returns number of performance variables 2

int MPIT_PERFVAR_Getum(int *num) 3

The name of individual variables (with numbers between 1 and N acquired by calling MPIT_PERFVAR_GETNUM) can then be queried with the following function along with any associated information. 4

MPIT_PERFVAR_GETINFO(num, name, name_len, verbosity, varclass, datatype, count, desc, desc_len, readonly, continuous) 5

IN num number of the performance variable to be queried 6

OUT len buffer to return the name of the performance variable 7

INOUT len_len length of the string and/or buffer for len 8

OUT verbosity verbosity level of this variable 9

OUT varclass class of performance variable 10

OUT datatype MPIT type of the information stored in the performance variable 11

OUT count number of elements returned 12

OUT desc buffer to return a description of the control variable 13

INOUT desc_len length of the string and/or buffer for desc 14

OUT readonly flags indicating whether variable can be written/reset 15

OUT continuous flags indicating whether variable can be started/stopped or is continuously activated 16

```
int MPIT_Perfvar_Getinfo(int num, char *name, int *name_len, int
    *verbosity, int *varclass, MPIT_Datatype *datatype, int
    *count, char *desc, int *desc_len, int *readonly, int
    *continuous) 17
```

The argument name provides a buffer to return the string describing the name of the control variable. The user has to pass the size of the buffer as the name_len argument. On return, the function deposits at most name_len-1 characters of the requested string into the buffer name followed by a terminating zero character. Additionally, the function writes the length of the returned string (including the terminating zero character) into name_len. If the returned value is smaller than the argument supplied to the function, the string has been truncated due to insufficient buffer resources. If the user passes NULL as the buffer argument or passes -1 as name_len, the function does not return the string and only returns the length of the string in name_len. 18

The argument verbosity returns the verbosity level (see Section 1.2.3) assigned by the MPI implementation to the variable. 19

1 The class of the performance variable is returned in the parameter `varclass` and can be
 2 one of the constants defined in Section 1.2.5.

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

5 The argument `desc` provides a buffer to return the string describing a description of
 6 the control variable. The user has to pass the size of the buffer as the `desc_len` argument.
 7 On return, the function deposits at most `desc_len-1` characters of the requested string into
 8 the buffer `desc` followed by a terminating zero character. Additionally, the function writes
 9 the length of the returned string (including the terminating zero character) into `desc_len`.
 10 If the returned value is smaller than the argument supplied to the function, the string has
 11 been truncated due to insufficient buffer resources. If the user passes `NULL` as the buffer
 12 argument or passes `-1` as `desc_len`, the function does not return the string and only returns
 13 the length of the string in `desc_len`.

14 Returning a description is optional. If an MPI library decides not to return a descrip-
 15 tion, the first character for `desc` must be set to the null character and `desc_len` must be set
 16 to one at the return from this function.

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

20 Upon return, the argument `continuous` will be set to one if the variable can be started
 21 and stopped by the user, or zero if the variable is automatically activated during `MPIT_INIT`
 22 and can not be stopped by the user.

24 Performance Experiment Sessions

25 Within a single program, multiple components can use the MPIT interface. To avoid col-
 26 lisions with respect to accesses to performance variables, users of the MPIT interface must
 27 first create a session. All subsequent calls accessing performance variables are then within
 28 the context of this session. Any call executed in a session shall not influence the results in
 29 any other session.

32 `MPIT_PERFVAR_SESSIONCREATE(session)`

33 IN `session` identifier of performance experiment session

36 `int MPIT_Perfvar_Sessioncreate(int *session)`

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

41 `MPIT_PERFVAR_SESSIONFREE(session)`

42 IN `session` identifier of performance experiment session

45 `int MPIT_Perfvar_Sessionfree(int session)`

46 This call frees an existing session, i.e., calls to MPIT can no longer be made within the
 47 freed session. After the call, all active performance variables in this context are deactivated.

48


```
int MPIT_Perfvar_Sessionfree(int session)
```

Performance Variable Activation

Before a performance variable can be used, i.e., started, stopped, read, written, or reset, it must first be activated. Only activated performance variables can be passed to start, stop or access functions discussed in the next sections.

```
MPIT_PERFVAR_ACTIVATE(session,num)
```

IN	session	identifier of performance experiment session
----	---------	--

IN	num	number of the performance variable
----	-----	------------------------------------

```
int MPIT_Perfvar_Activate(int session, int num)
```

This routine activates the performance variable `num` with respect to session `session`. If this variable is not yet activated, the variable will be reset to its default value. Calling this function on already activated variables (within the same session) has no affect.

```
MPIT_PERFVAR_DEACTIVATE(session,num)
```

IN	session	identifier of performance experiment session
----	---------	--

IN	num	number of the performance variable
----	-----	------------------------------------

```
int MPIT_Perfvar_Deactivate(int session, int num)
```

This routine deactivates the performance variable `num` with respect to session `session`.

Advice to implementors. The extra step of activating performance variables allows MPIT implementations to selectively enable counters and only monitor activated events. This can be used to minimize the overhead of any performance monitor when not used. (*End of advice to implementors.*)

Starting and Stopping of Performance Variables

Performance variables that have the `continuous` flag set during the query operation are continuously operating after a call to `MPIT_PERFVAR_ACTIVATE` and can not be stopped or paused by the user. All other variables are in a stopped state after their first activation within a session, i.e., they are not updated as the program executes, and have to be started by the user.

```
MPIT_PERFVAR_START(session,num)
```

IN	session	Identifier of performance experiment session
----	---------	--

IN	num	number of the performance variable
----	-----	------------------------------------

```
int MPIT_Perfvar_Start(int session, int num)
```

1 This function starts the performance variable with the number `num` in the session
 2 `session`. The variable has to be activated before making this call using the function
 3 `MPIT_PERFVAR_ACTIVATE`.

4 If the constant `MPIT_PERFVAR_ALL` is passed in `num`, the MPI library attempts to
 5 start all activated variables within the session identified by `session`. In this case, the routine
 6 returns `MPI_SUCCESS` if all variables are started successfully; continuous variables, variables
 7 that are already started, and not activated variables are ignored when used with
 8 `MPIT_PERFVAR_ALL`.

9
 10
 11 `MPIT_PERFVAR_STOP(session, num)`

12	IN	<code>session</code>	Identifier of performance experiment session
13	IN	<code>num</code>	number of the performance variable

14
 15
 16 `int MPIT_Perfvar_Stop(int session, int num)`

17 This function stops the performance variable with the number `num` in the session
 18 `session`. The variable has to be activated before making this call using the function
 19 `MPIT_PERFVAR_ACTIVATE`.

20 If the constant `MPIT_PERFVAR_ALL` is passed in `num`, the MPI library attempts
 21 to stop all activated variables within the session identified by `session`. In this case, the
 22 routine returns `MPI_SUCCESS` if all variables are stopped successfully; continuous variables,
 23 variables that are already stopped, and not activated variables are ignored when used with
 24 `MPIT_PERFVAR_ALL`.

25
 26 *Advice to implementors.* Although MPI places no requirements on the interaction
 27 with external mechanisms such as signal handlers, it is strongly recommended that the
 28 routines in this section to start and stop performance variables should be safe to call
 29 in asynchronous contexts. Examples of asynchronous contexts include signal handlers
 30 and interrupt handlers. Such safety permits the development of sampling-based tools.
 31 High quality implementations should strive to make the results of any such interactions
 32 intuitive to users, and attempt to document restrictions where deemed necessary. (*End*
 33 *of advice to implementors.*)

34 35 Performance Variable Access Functions

36
 37
 38 `MPIT_PERFVAR_READ(session, num, buf)`

39	IN	<code>session</code>	Identifier of performance experiment session
40	IN	<code>num</code>	number of the performance variable
41	OUT	<code>buf</code>	initial address of storage location for variable value

42
 43
 44
 45 `int MPIT_Perfvar_Read(int session, int num, void* buf)`

46 The `MPIT_PERFVAR_READ` call queries the value of the performance variable with
 47 the number `num` in the session `session` and stores the result in the buffer `buf`. The user is
 48 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_GETINFO call). The variable has to be activated before making this call using the function MPIT_PERFVAR_ACTIVATE.

MPIT_PERFVAR_WRITE(session, num, buf)

IN	session	Identifier of performance experiment session
IN	num	number of the performance variable
IN	buf	initial address of storage location for variable value

int MPIT_Perfvar_write(int session, int num, void* buf)

The MPIT_PERFVAR_WRITE call attempts to write the value of the performance variable with the number `num` 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_GETINFO call). The variable has to be activated before making this call using the function MPIT_PERFVAR_ACTIVATE.

If it is not possible to change the variable the function returns MPIT_ERR_PERFVAR_WRITE.

MPIT_PERFVAR_RESET(session,num)

IN	session	Identifier of performance experiment session
IN	num	number of the performance variable

int MPIT_Perfvar_reset(int session, int num)

The MPIT_PERFVAR_RESET call sets the value of the performance variable to its default starting value. If it is not possible to change the variable the function returns MPIT_ERR_PERFVAR_WRITE.

If the constant MPIT_PERFVAR_ALL is passed in `num`, the MPI library attempts to reset all activated variables within the session identified by `session`. In this case, the routine returns MPIT_SUCCESS if all variables are reset successfully; readonly variables, and not activated variables are ignored when used with MPIT_PERFVAR_ALL .

MPIT_PERFVAR_READRESET(session,num, buf)

IN	session	Identifier of performance experiment session
IN	num	number of the performance variable
OUT	buf	initial address of storage location for variable value

int MPIT_Perfvar_Readreset(int session, int num, void* buf)

The MPIT_PERFVAR_READRESET call atomically queries the value of the performance variable, stores the result in the buffer `buf`, and then sets the value of the performance variable to its default starting value. 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 query call). If it is not possible to change the variable the function returns `MPIT_ERR_PERFVAR_WRITE`. In this case, the value returned in `buf` is the same as if the variable would have been read by the `MPIT_PERFVAR_READ` call.

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

1.2.6 Performance and Control Variable Taxonomic Information

MPI implementations can optionally provide information that describes the relationship of performance and control variables to each other. For this, an MPI implementation can define names that represent sets of variables and then associate each performance/control variable with zero or more sets. Sets may contain zero or more performance/control variables and zero or more other sets. Sets may not contain themselves either directly or indirectly. More formally, these sets and performance/control variables form a directed acyclic graph (DAG). This information is accessible via several interrogative routines.

`MPIT_TAXON_QUERY_SET_SETS(iterator, setname, name, namelen, type)`

INOUT	iterator	iterator variable passed in by user (iterator)
IN	setname	name of the set to be queried (string)
OUT	name	name of the set returned on this iteration (string)
OUT	namelen	length of the name of the set returned on this iteration (string)
OUT	type	type of the set returned this iteration (integer)

```
int MPIT_Taxon_query_set_sets(MPIT_Taxonquery_iterator *iterator, char
                             *setname, char *name, int *namelen, int *type);
```

Iterate over all sets contained in the set identified by `setname`. A unique identifying name for the contained set is returned in `name` and `namelen` is set to the number of characters written. The value of `namelen` cannot be larger than `MPIT_MAX_SET_NAME-1`. The set of all root sets (sets that no other set contains) available in the implementation can be queried by using a `setname` of `MPIT_ROOT_SETS`. The `type` parameter is set to `MPIT_TYPE_CTRLVAR` if the variable is a control variable and to `MPIT_TYPE_PERFVAR` if it is a performance variable.

On the first call to `MPIT_TAXON_QUERY_SET_SETS`, the caller must initialize a variable to `MPIT_TAXON_QUERY_START` and pass this variable as the `iter` parameter. Subsequent calls require the user to pass the returned value `iter` to query further taxonomic information. Once all taxonomic information is returned, the call to

MPIT_TAXON_QUERY_SET_SETS returns MPIT_END and sets iter to MPIT_TAXON_QUERY_END.

MPIT_TAXON_QUERY_VARIABLE_SETS(iterator, varname, name, namelen, type)

INOUT	iterator	iterator variable passed in by user (iterator)
IN	varname	name of the variable to be queried (string)
OUT	name	name of the set returned this iteration (string)
OUT	namelen	length of the name of the set returned this iteration (integer)
OUT	type	type of the set returned this iteration (integer)

```
int MPIT_Taxon_query_variable_sets(MPIT_Taxonquery_iterator *iterator, char
    *varname, char *name, int *namelen, int *type);
```

Iterate over all sets that contain the performance/control variable identified by `varname`. A unique identifying name for the set is returned in `name` and `namelen` is set to the number of characters written. The value of `namelen` cannot be larger than `MPIT_MAX_SET_NAME-1`. The `type` parameter is set to `MPIT_TYPE_CTRLVAR` if the variable is a control variable and to `MPIT_TYPE_PERFVAR` if it is a performance variable.

On the first call to `MPIT_TAXON_QUERY_VARIABLE_SETS`, the caller must initialize a variable to `MPIT_TAXON_QUERY_START` and pass this variable as the `iter` parameter. Subsequent calls require the user to pass the returned value `iter` to query further taxonomic information. Once all taxonomic information is returned, the call to `MPIT_TAXON_QUERY_VARIABLE_SETS` returns `MPIT_END` and sets `iter` to `MPIT_TAXON_QUERY_END`.

MPIT_TAXON_QUERY_SET_VARIABLES(iterator, setname, name, namelen, type)

INOUT	iterator	iterator variable passed in by user (iterator)
IN	setname	name of the set to be queried (string)
OUT	name	name of the variable returned this iteration (string)
OUT	namelen	length of the name of the variable returned this iteration (integer)
OUT	type	type of the variable returned this iteration (integer)

```
int MPIT_Taxon_query_set_variables(MPIT_Taxonquery_iterator *iterator, char
    *setname, char *name, int *namelen, int *type);
```

Iterate over all variables directly contained in the set identified by `setname`. That is, variables contained indirectly by a contained set will not be returned by this call. A unique identifying name for the variable is returned in `name` and `namelen` is set to the number of characters written. The value of `namelen` cannot be larger than `MPIT_MAX_SET_NAME-1`. The `type` parameter is set to `MPIT_TYPE_CTRLVAR` if the variable is a control variable and to `MPIT_TYPE_PERFVAR` if it is a performance variable.

1 On the first call to `MPIT_TAXON_QUERY_SET_VARIABLES`, the caller must initialize
 2 a variable to `MPIT_TAXON_QUERY_START` and pass this variable as the
 3 `iter` parameter. Subsequent calls require the user to pass the returned value `iter` to query
 4 further taxonomic information. Once all taxonomic information is returned, the call to
 5 `MPIT_TAXON_QUERY_SET_VARIABLES` returns `MPIT_END` and sets `iter` to
 6 `MPIT_TAXONQUERY_END`.

7
 8
 9 `MPIT_TAXON_CHANGED(flag)`

10 OUT flag true if the taxonomic information has changed since
 11 the last call to a query function (boolean)

12
 13 `int MPIT_Taxon_changed(int *flag);`
 14

15 This routine returns true in the `flag` argument if the list of available performance/control
 16 variables or sets has changed since the last time the user has called any of the
 17 `MPIT_TAXON_QUERY_*` routines with the argument `MPIT_TAXON_QUERY_START` as the
 18 first argument. If the user has not yet called any such routines, the argument `flag` will
 19 contain the value true.

20
 21 `MPIT_TAXON_DESCRIBE_SET(name, desc, desclen)`
 22

23 IN name name of the set to describe (string)
 24 OUT desc description of the set (string)
 25 OUT desclen length of the string returned in `desc` (int)

26
 27
 28 `int MPIT_Taxon_describe_set(char *name, char *desc, int desclen);`

29 Retrieve the description for the set identified by `name` and store it in `desc`. The `desclen`
 30 parameter is set to the number of characters written. The value of `desclen` cannot be larger
 31 than `MPIT_MAX_SET_DESC-1`.
 32

33 **Set and Variable Names** MPI does not specify the character encoding of strings in the
 34 MPIT interface. The only requirement is that strings are terminated with a null character.

35 MPI reserves all set and variable names with the prefixes “MPI_” and “MPIT_” for
 36 its own use.
 37

38 1.2.7 Return and Error Codes

39 All MPIT functions return a return or error code. The following constants are available
 40 for this for the specific calls. None of the error codes returned by a MPIT routine shall be
 41 considered fatal to the overall MPI implementation or shall invoke an MPI error handler.
 42 In any case, the execution of the MPI program shall continue as if the call would have
 43 succeeded. However, the MPIT implementation is not required to check all user provided
 44 parameters; if a user passes illegal parameter values to any MPIT routine that are not caught
 45 by the implementation, the behavior of the library is undefined.
 46
 47
 48

Return Code	Description
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

Table 1.5: Return codes used by any MPIT function.

Return Code	Description
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 number queried is out of range (for MPIT_TYPE_ENUMITEM only)

Table 1.6: Return codes used by MPIT type functions.

Return Code	Description
MPIT_ERR_SETNOTNOW	Variable cannot be set at this moment
MPIT_ERR_SETNEVER	Variable cannot be set until end of execution
MPIT_ERR_INVALIDVAR	Control variable does not exist

Table 1.7: Return codes used by MPIT control variable access functions.

Return Codes for all MPIT Functions

The return codes in Table 1.5 apply to all MPIT functions.

Return Codes for Type Functions

The return codes in Table 1.6 apply to MPIT_TYPE_ENUMQUERY, MPIT_TYPE_GETCLASS and MPIT_TYPE_ENUMITEM.

Return Codes for Control Variable Access Functions

The return codes in Table 1.7 apply to MPIT_CONFIG_READ and MPIT_CONFIG_WRITE.

Return Codes for Performance Variable Access and Control

The return codes in Table 1.8 apply to MPIT_PERFVAR_START , MPIT_ERR_PERFVAR_STOP , MPIT_PERFVAR_READ , MPIT_ERR_PERFVAR_WRITE , MPIT_PERFVAR_RESET , and MPIT_PERFVAR_READRESET .

Return Codes for Taxonomy Functions

The return codes in Table 1.9 apply to MPIT_TAXON_ routines.

1.2.8 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 has to pro-

Return Code	Description
MPIT_ERR_INVALIDVAR	Performance variable does not exist
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

Table 1.8: Return codes used by MPIT performance variable access, start, stop, or activation functions.

Return Code	Description
MPIT_ERR_NOSET	The set does not exist
MPIT_ERR_NODATA	No description for this set available

Table 1.9: Return codes used MPIT taxonomy functions.

vide matching PMPIT calls for every MPIT call. All rules, guidelines, and recommendations from Section ?? apply equally to PMPIT calls.

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.

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

MPI Constant and Predefined Handle Index

This index lists predefined MPI constants and handles.

MPI_BYTE, 3
MPI_CHAR, 3
MPI_COMM_SELF, 8
MPI_DOUBLE, 3
MPI_FLOAT, 3
MPI_INT, 3
MPI_LOGICAL, 3
MPI_LONG, 3
MPI_LONG_LONG, 3
MPI_SHORT, 3
MPI_SUCCESS, 14
MPI_VERBOSITY_USER_BASIC, 5
MPIT_BYTE, 3
MPIT_CHAR, 3
MPIT_DOUBLE, 3
MPIT_END, 17, 18
MPIT_ERR_CANTINIT, 19
MPIT_ERR_INVALIDITEM, 19
MPIT_ERR_INVALIDSESSION, 20
MPIT_ERR_INVALIDTYPE, 19
MPIT_ERR_INVALIDVAR, 19, 20
MPIT_ERR_MEMORY, 19
MPIT_ERR_NODATA, 20
MPIT_ERR_NOSET, 20
MPIT_ERR_NOSTARTSTOP, 20
MPIT_ERR_NOTINITIALIZED, 19
MPIT_ERR_NOWRITE, 20
MPIT_ERR_PERFVAR_WRITE, 15, 16
MPIT_ERR_PREDEFINED, 19
MPIT_ERR_SETNEVER, 9, 19
MPIT_ERR_SETNOTNOW, 9, 19
MPIT_FLOAT, 3
MPIT_INT, 3
MPIT_LOGICAL, 3
MPIT_LONG, 3
MPIT_LONG_LONG, 3
MPIT_MAX_SET_DESC, 18
MPIT_MAX_SET_NAME, 16, 17
MPIT_PERFVAR_ALL, 14, 15
MPIT_ROOT_SETS, 16
MPIT_SCOPE_GLOBAL, 8
MPIT_SCOPE_LOCAL, 8
MPIT_SCOPE_READONLY, 8
MPIT_SHORT, 3
MPIT_SUCCESS, 2, 15, 19
MPIT_TAXON_QUERY_END, 17
MPIT_TAXON_QUERY_START, 16–18
MPIT_TAXONQUERY_END, 18
MPIT_TYPE_CTRLVAR, 16, 17
MPIT_TYPE_PERFVAR, 16, 17
MPIT_TYPECLASS_ENUMERATION, 5
MPIT_TYPECLASS_PREDEFINED, 5
MPIT_VERBOSITY_MPIDEV_BASIC, 5
MPIT_VERBOSITY_MPIDEV_DETAILED, 5
MPIT_VERBOSITY_MPIDEV_VERBOSE, 5
MPIT_VERBOSITY_TUNER_BASIC, 5
MPIT_VERBOSITY_TUNER_DETAILED, 5
MPIT_VERBOSITY_TUNER_VERBOSE, 5
MPIT_VERBOSITY_USER_BASIC, 5
MPIT_VERBOSITY_USER_DETAILED, 5
MPIT_VERBOSITY_USER_VERBOSE, 5

MPI Function Index

The underlined page numbers refer to the function definitions.

MPI_COMM_SELF, 9

MPI_FINALIZE, 3

MPI_INIT, 3