

```

1 MPI_WTIME()
2
3 double MPI_Wtime(void)
4
5 DOUBLE PRECISION MPI_WTIME()
6 {double MPI::Wtime() (binding deprecated, see Section 15.2) }

```

MPI_WTIME returns a floating-point number of seconds, representing elapsed wall-clock time since some time in the past.

The “time in the past” is guaranteed not to change during the life of the process. The user is responsible for converting large numbers of seconds to other units if they are preferred.

This function is portable (it returns seconds, not “ticks”), it allows high-resolution, and carries no unnecessary baggage. One would use it like this:

```

15 {
16     double starttime, endtime;
17     starttime = MPI_Wtime();
18     .... stuff to be timed ...
19     endtime   = MPI_Wtime();
20     printf("That took %f seconds\n",endtime-starttime);
21 }

```

The times returned are local to the node that called them. There is no requirement that different nodes return “the same time.” (But see also the discussion of MPI_WTIME_IS_GLOBAL).

```

28 MPI_WTICK()
29
30 double MPI_Wtick(void)
31
32 DOUBLE PRECISION MPI_WTICK()
33 {double MPI::Wtick() (binding deprecated, see Section 15.2) }

```

MPI_WTICK returns the resolution of MPI_WTIME in seconds. That is, it returns, as a double precision value, the number of seconds between successive clock ticks. For example, if the clock is implemented by the hardware as a counter that is incremented every millisecond, the value returned by MPI_WTICK should be 10^{-3} .

8.7 Startup

One goal of MPI is to achieve *source code portability*. By this we mean that a program written using MPI and complying with the relevant language standards is portable as written, and must not require any source code changes when moved from one system to another. This explicitly does *not* say anything about how an MPI program is started or launched from the command line, nor what the user must do to set up the environment in which an MPI program will run. However, an implementation may require some setup to be performed

before other MPI routines may be called. To provide for this, MPI includes an initialization routine MPI_INIT.

MPI_INIT()

```
int MPI_Init(int *argc, char ***argv)
```

```
MPI_INIT(IERROR)
```

```
INTEGER IERROR
```

```
{void MPI::Init(int& argc, char**& argv) (binding deprecated, see Section 15.2) }
```

```
{void MPI::Init() (binding deprecated, see Section 15.2) }
```

[All MPI programs must contain exactly one call to an MPI initialization routine: MPI_INIT or MPI_INIT_THREAD. Subsequent calls to any initialization routines are erroneous. The only MPI functions that may be invoked before the MPI initialization routines are called are MPI_GET_VERSION, MPI_INITIALIZED, and MPI_FINALIZED.] Each MPI process must call an MPI initialization routine, MPI_INIT or MPI_INIT_THREAD, exactly once. Subsequent calls by the process to any initialization routine are erroneous. The only MPI functions that may be invoked by a process before the MPI initialization routine completed are MPI_GET_VERSION, MPI_INITIALIZED, and MPI_FINALIZED.

The version for ISO C accepts the argc and argv that are provided by the arguments to main or NULL:

```
int main(int argc, char **argv)
```

```
{
```

```
    MPI_Init(&argc, &argv);
```

```
    /* parse arguments */
```

```
    /* main program */
```

```
    MPI_Finalize();    /* see below */
```

```
}
```

The Fortran version takes only IERROR.

Conforming implementations of MPI are required to allow applications to pass NULL for both the argc e argv arguments of main in C. [and C++. In C++, there is an alternative binding for MPI::Init that does not have these arguments at all.]

Rationale. In some applications, libraries may be making the call to MPI_Init, and may not have access to argc and argv from main. It is anticipated that applications requiring special information about the environment or information supplied by mpiexec can get that information from environment variables. (*End of rationale.*)

MPI_FINALIZE()

```
int MPI_Finalize(void)
```

```

1 MPI_FINALIZE(IERROR)
2     INTEGER IERROR
3
4 {void MPI::Finalize() (binding deprecated, see Section 15.2) }
```

ticket313. This routine cleans up all MPI state. [Each process must call `MPI_FINALIZE` before it exits. Unless there has been a call to `MPI_ABORT`, before each process exits process must ensure that all pending nonblocking communications are (locally) complete before calling `MPI_FINALIZE`. Further, at the instant at which the last process calls `MPI_FINALIZE`, all pending sends must be matched by a receive, and all pending receives must be matched by a send.

For example, the following program is correct] If an MPI program terminates normally (i.e., not due to a call to `MPI_ABORT` or an unrecovered error) then the following must hold:

- Each process must call `MPI_FINALIZE` before the process exits.
- When the last process calls `MPI_FINALIZE`, all non-local MPI calls at each process have been matched by MPI calls at the other processes that are needed to complete the relevant operation: For each send, the matching receive has occurred, each collective operation has been invoked at all involved processes, etc.

The following examples illustrates these rules

Example 8.3 The following code is correct

Process 0	Process 1
-----	-----
MPI_Init();	MPI_Init();
MPI_Send(dest=1);	MPI_Recv(src=0);
MPI_Finalize();	MPI_Finalize();

Example 8.4 Without a matching receive, the program is erroneous

Process 0	Process 1
-----	-----
MPI_Init();	MPI_Init();
MPI_Send (dest=1);	
MPI_Finalize();	MPI_Finalize();

ticket313. [deleted in April Since `MPI_FINALIZE` is a collective call, a correct MPI program will naturally ensure that all participants in pending collective operations have made the call before calling `MPI_FINALIZE`.

A successful return from a blocking communication operation or from `MPI_WAIT` or `MPI_TEST` tells the user that the buffer can be reused and means that the communication is completed by the user, but does not guarantee that the local process has no more work to do. A successful return from `MPI_REQUEST_FREE` with a request handle generated by an `MPI_ISEND` nullifies the handle but provides no assurance of operation completion. The `MPI_ISEND` is complete only when it is known by some means that a matching receive has

completed. `MPI_FINALIZE` guarantees that all local actions required by communications the user has completed will, in fact, occur before it returns.

`MPI_FINALIZE` guarantees nothing about pending communications that have not been completed (completion is assured only by `MPI_WAIT`, `MPI_TEST`, or `MPI_REQUEST_FREE` combined with some other verification of completion).]

Example 8.5 This program is correct HEADER SKIP ENDHEADER

```
rank 0                                rank 1
=====
...
MPI_Isend();                          MPI_Recv();
MPI_Request_free();                   MPI_Barrier();
MPI_Barrier();                        MPI_Finalize();
MPI_Finalize();                       exit();
exit();
```

Example 8.6 This program is erroneous and its behavior is undefined: HEADER SKIP ENDHEADER

```
rank 0                                rank 1
=====
...
MPI_Isend();                          MPI_Recv();
MPI_Request_free();                   MPI_Finalize();
MPI_Finalize();                       exit();
exit();
```

Example 8.7 This program is erroneous: The `MPI_Isend` call is not guaranteed to be locally complete before process 0 calls `MPI_Finalize`

```
Process 0                               Process 1
-----
MPI_Isend();                            MPI_Recv();
MPI_Request_free();                     MPI_Barrier();
MPI_Barrier();                          MPI_Finalize();
MPI_Finalize();
```

[If no `MPI_BUFFER_DETACH` occurs between an `MPI_BSEND` (or other buffered send) and `MPI_FINALIZE`, the `MPI_FINALIZE` implicitly supplies the `MPI_BUFFER_DETACH`.

Example 8.8 This program is correct, and after the `MPI_Finalize`, it is as if the buffer had been detached. HEADER SKIP ENDHEADER

```

1      rank 0                                rank 1
2      =====
3      ...                                    ...
4      buffer = malloc(1000000);             MPI_Recv();
5      MPI_Buffer_attach();                  MPI_Finalize();
6      MPI_Bsend();                          exit();
7      MPI_Finalize();
8      free(buffer);
9      exit();

```

ticket313.] While the user must ensure that communications can complete before MPI is finalized, it needs not free resources allocated by MPI (buffers, windows, requests, communicators, etc.); the MPI_FINALIZE function will do so.

Example 8.9 This program is correct, and after the MPI_Finalize, it is as if the buffer had been detached.

```

17      Process 0                                Process 1
18      -----                                -----
19      buffer = malloc(1000000);             MPI_Recv();
20      MPI_Buffer_attach();                  MPI_Finalize();
21      MPI_Bsend();                          exit();
22      MPI_Finalize();
23      free(buffer);
24      exit();

```

ticket313.

Example 8.10 In this example, MPI_lprobe() must return a FALSE flag. MPI_Test_cancelled() must return a TRUE flag, independent of the relative order of execution of MPI_Cancel() in process 0 and MPI_Finalize() in process 1.

The MPI_lprobe() call is there to make sure the implementation knows that the “tag1” message exists at the destination, without being able to claim that the user knows about it.

HEADER SKIP ENDHEADER

```

36      rank 0                                rank 1
37      =====
38      MPI_Init();                          MPI_Init();
39      MPI_Isend(tag1);                      MPI_Barrier();
40      MPI_Barrier();                       MPI_Iprobe(tag2);
41                                           MPI_Barrier();
42      MPI_Barrier();                       MPI_Finalize();
43                                           exit();
44
45      MPI_Cancel();
46      MPI_Wait();
47      MPI_Test_cancelled();
48      MPI_Finalize();

```

```
exit();
```

```
]
```

Example 8.11 This program is correct. The cancel operation must succeed, since the send cannot complete normally.

```

Process 0                                Process 1
-----                                -----
MPI_Isend(tag1);                          MPI_Finalize();
MPI_Cancel();
MPI_Wait();
MPI_Finalize();

```

```
[
```

Advice to implementors. An implementation may need to delay the return from MPI_FINALIZE until all potential future message cancellations have been processed. One possible solution is to place a barrier inside MPI_FINALIZE (*End of advice to implementors.*)

```
]
```

Advice to implementors. An implementation may need to delay the return from MPI_FINALIZE on a process even if all communications related to MPI calls by that process have completed; the process may still receive cancel requests for messages it has completed receiving. One possible solution is to place a barrier inside MPI_FINALIZE (*End of advice to implementors.*)

Once MPI_FINALIZE returns, no MPI routine (not even MPI_INIT) may be called, except for MPI_GET_VERSION, MPI_INITIALIZED, and MPI_FINALIZED. Each process must complete any pending communication it initiated before it calls MPI_FINALIZE. If the call returns, each process may continue local computations, or exit, without participating in further MPI communication with other processes. MPI_FINALIZE is collective over all connected processes. If no processes were spawned, accepted or connected then this means over MPI_COMM_WORLD; otherwise it is collective over the union of all processes that have been and continue to be connected, as explained in Section 10.5.4 on page 358.

Advice to implementors. Even though a process has completed all the communication it initiated, such communication may not yet be completed from the viewpoint of the underlying MPI system. E.g., a blocking send may have completed, even though the data is still buffered at the sender. The MPI implementation must ensure that a process has completed any involvement in MPI communication before MPI_FINALIZE returns. Thus, if a process exits after the call to MPI_FINALIZE, this will not cause an ongoing communication to fail. (*End of advice to implementors.*)

Although it is not required that all processes return from MPI_FINALIZE, it is required that at least process 0 in MPI_COMM_WORLD return, so that users can know that the MPI

1 portion of the computation is over. In addition, in a POSIX environment, they may desire
 2 to supply an exit code for each process that returns from MPI_FINALIZE.

3
 4 **Example 8.12** The following illustrates the use of requiring that at least one process
 5 return and that it be known that process 0 is one of the processes that return. One wants
 6 code like the following to work no matter how many processes return.

```
7
8     ...
9     MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
10    ...
11    MPI_Finalize();
12    if (myrank == 0) {
13        resultfile = fopen("outfile","w");
14        dump_results(resultfile);
15        fclose(resultfile);
16    }
17    exit(0);
```

18
 19
 20 MPI_INITIALIZED(flag)

21 OUT flag Flag is true if MPI_INIT has been called and false
 22 otherwise.

23
 24 int MPI_Initialized(int *flag)

25
 26 MPI_INITIALIZED(FLAG, IERROR)

27 LOGICAL FLAG
 28 INTEGER IERROR

29 {bool MPI::Is_initialized() (*binding deprecated, see Section 15.2*) }

30
 31 This routine may be used to determine whether MPI_INIT has been called.
 32 MPI_INITIALIZED returns true if the calling process has called MPI_INIT. Whether
 33 MPI_FINALIZE has been called does not affect the behavior of MPI_INITIALIZED. It is one
 34 of the few routines that may be called before MPI_INIT is called.

35
 36
 37 MPI_ABORT(comm, errorcode)

38 IN comm communicator of tasks to abort
 39 IN errorcode error code to return to invoking environment

40
 41 int MPI_Abort(MPI_Comm comm, int errorcode)

42
 43 MPI_ABORT(COMM, ERRORCODE, IERROR)

44 INTEGER COMM, ERRORCODE, IERROR

45
 46 {void MPI::Comm::Abort(int errorcode) (*binding deprecated, see Section 15.2*) }

47 This routine makes a “best attempt” to abort all tasks in the group of comm. This
 48 function does not require that the invoking environment take any action with the error

code. However, a Unix or POSIX environment should handle this as a `return errorcode` from the main program.

It may not be possible for an MPI implementation to abort only the processes represented by `comm` if this is a subset of the processes. In this case, the MPI implementation should attempt to abort all the connected processes but should not abort any unconnected processes. If no processes were spawned, accepted or connected then this has the effect of aborting all the processes associated with `MPI_COMM_WORLD`.

Rationale. The communicator argument is provided to allow for future extensions of MPI to environments with, for example, dynamic process management. In particular, it allows but does not require an MPI implementation to abort a subset of `MPI_COMM_WORLD`. (*End of rationale.*)

Advice to users. Whether the errorcode is returned from the executable or from the MPI process startup mechanism (e.g., `mpiexec`), is an aspect of quality of the MPI library but not mandatory. (*End of advice to users.*)

Advice to implementors. Where possible, a high-quality implementation will try to return the errorcode from the MPI process startup mechanism (e.g. `mpiexec` or singleton init). (*End of advice to implementors.*)

8.7.1 Allowing User Functions at Process Termination

There are times in which it would be convenient to have actions happen when an MPI process finishes. For example, a routine may do initializations that are useful until the MPI job (or that part of the job that being terminated in the case of dynamically created processes) is finished. This can be accomplished in MPI by attaching an attribute to `MPI_COMM_SELF` with a callback function. When `MPI_FINALIZE` is called, it will first execute the equivalent of an `MPI_COMM_FREE` on `MPI_COMM_SELF`. This will cause the delete callback function to be executed on all keys associated with `MPI_COMM_SELF`, in the reverse order that they were set on `MPI_COMM_SELF`. If no key has been attached to `MPI_COMM_SELF`, then no callback is invoked. The “freeing” of `MPI_COMM_SELF` occurs before any other parts of MPI are affected. Thus, for example, calling `MPI_FINALIZED` will return `false` in any of these callback functions. Once done with `MPI_COMM_SELF`, the order and rest of the actions taken by `MPI_FINALIZE` is not specified.

Advice to implementors. Since attributes can be added from any supported language, the MPI implementation needs to remember the creating language so the correct callback is made. Implementations that use the attribute delete callback on `MPI_COMM_SELF` internally should register their internal callbacks before returning from `MPI_INIT / MPI_INIT_THREAD`, so that libraries or applications will not have portions of the MPI implementation shut down before the application-level callbacks are made. (*End of advice to implementors.*)

8.7.2 Determining Whether MPI Has Finished

One of the goals of MPI was to allow for layered libraries. In order for a library to do this cleanly, it needs to know if MPI is active. In MPI the function `MPI_INITIALIZED` was

provided to tell if MPI had been initialized. The problem arises in knowing if MPI has been finalized. Once MPI has been finalized it is no longer active and cannot be restarted. A library needs to be able to determine this to act accordingly. To achieve this the following function is needed:

MPI_FINALIZED(flag)

OUT flag true if MPI was finalized (logical)

int MPI_Finalized(int *flag)

MPI_FINALIZED(FLAG, IERROR)

LOGICAL FLAG

INTEGER IERROR

{bool MPI::Is_finalized() (*binding deprecated, see Section 15.2*) }

This routine returns true if MPI_FINALIZE has completed. It is legal to call MPI_FINALIZED before MPI_INIT and after MPI_FINALIZE.

Advice to users. MPI is “active” and it is thus safe to call MPI functions if MPI_INIT has completed and MPI_FINALIZE has not completed. If a library has no other way of knowing whether MPI is active or not, then it can use MPI_INITIALIZED and MPI_FINALIZED to determine this. For example, MPI is “active” in callback functions that are invoked during MPI_FINALIZE. (*End of advice to users.*)

8.8 Portable MPI Process Startup

A number of implementations of MPI provide a startup command for MPI programs that is of the form

```
mpirun <mpirun arguments> <program> <program arguments>
```

Separating the command to start the program from the program itself provides flexibility, particularly for network and heterogeneous implementations. For example, the startup script need not run on one of the machines that will be executing the MPI program itself.

Having a standard startup mechanism also extends the portability of MPI programs one step further, to the command lines and scripts that manage them. For example, a validation suite script that runs hundreds of programs can be a portable script if it is written using such a standard startup mechanism. In order that the “standard” command not be confused with existing practice, which is not standard and not portable among implementations, instead of `mpirun` MPI specifies `mpiexec`.

While a standardized startup mechanism improves the usability of MPI, the range of environments is so diverse (e.g., there may not even be a command line interface) that MPI cannot mandate such a mechanism. Instead, MPI specifies an `mpiexec` startup command and recommends but does not require it, as advice to implementors. However, if an implementation does provide a command called `mpiexec`, it must be of the form described below.

It is suggested that

```
mpiexec -n <numprocs> <program>
```

be at least one way to start <program> with an initial MPI_COMM_WORLD whose group contains <numprocs> processes. Other arguments to mpiexec may be implementation-dependent.

Advice to implementors. Implementors, if they do provide a special startup command for MPI programs, are advised to give it the following form. The syntax is chosen in order that mpiexec be able to be viewed as a command-line version of MPI_COMM_SPAWN (See Section 10.3.4).

Analogous to MPI_COMM_SPAWN, we have

```
[
    mpiexec -n    <maxprocs>
            -soft <      >
            -host <      >
            -arch <      >
            -wdir <      >
            -path <      >
            -file <      >
            ...
            <command line>
]
```

```
mpiexec -n    <maxprocs>
            -soft <      >
            -host <      >
            -arch <      >
            -wdir <      >
            -path <      >
            -file <      >
            -asp  <      >
            ...
            <command line>
```

for the case where a single command line for the application program and its arguments will suffice. See Section 10.3.4 for the meanings of these arguments. For the case corresponding to MPI_COMM_SPAWN_MULTIPLE there are two possible formats:

Form A:

```
mpiexec { <above arguments> } : { ... } : { ... } : ... : { ... }
```

As with MPI_COMM_SPAWN, all the arguments are optional. (Even the -n x argument is optional; the default is implementation dependent. It might be 1, it might be taken from an environment variable, or it might be specified at compile time.) The names and meanings of the arguments are taken from the keys in the info argument

1 to MPI_COMM_SPAWN. There may be other, implementation-dependent arguments
2 as well.

3 Note that Form A, though convenient to type, prevents colons from being program
4 arguments. Therefore an alternate, file-based form is allowed:

5 Form B:

```
6
7     mpiexec -configfile <filename>
```

8
9 where the lines of <filename> are of the form separated by the colons in Form A.
10 Lines beginning with '#' are comments, and lines may be continued by terminating
11 the partial line with '\'.
12

13 **Example 8.13** Start 16 instances of myprog on the current or default machine:

```
14
15     mpiexec -n 16 myprog
```

16
17 **Example 8.14** Start 10 processes on the machine called ferrari:

```
18
19     mpiexec -n 10 -host ferrari myprog
```

20
21 **Example 8.15** Start three copies of the same program with different command-line
22 arguments:
23

```
24     mpiexec myprog infile1 : myprog infile2 : myprog infile3
```

25
26 **Example 8.16** Start the ocean program on five Suns and the atmos program on 10
27 RS/6000's:
28

```
29     mpiexec -n 5 -arch sun ocean : -n 10 -arch rs6000 atmos
```

30
31 It is assumed that the implementation in this case has a method for choosing hosts of
32 the appropriate type. Their ranks are in the order specified.
33

34 **Example 8.17** Start the ocean program on five Suns and the atmos program on 10
35 RS/6000's (Form B):

```
36     mpiexec -configfile myfile
```

37
38 where myfile contains

```
39     -n 5 -arch sun    ocean
40     -n 10 -arch rs6000 atmos
```

41
42
43 **Example 8.18** Start 12 MPI processes of the foo program, with 4 MPI processes in
44 each address space:

```
45
46     mpiexec -asp 4 -n 12 foo
```

47 (*End of advice to implementors.*)
48

ticket310.