

# MPI3: Supporting Multiple Endpoints per Process

Hybrid Programming Working Group

September 4, 2011

## 0.1 Introduction

### 0.1.1 Rationale

Many applications run best when shared memory is used for communication inside a node and MPI is used between nodes, *with more than one MPI process per node*. We expect this model to become increasingly common as the number of cores per node increases: A larger number of concurrent threads per node will often lead to a higher message rate that is hard to achieve within one MPI process, for a variety of reasons. For example, message matching time tends to increase linearly with the number of pending requests; the problem is inherent to the matching semantics of MPI. The use of multiple MPI processes can also reduce intra-node communication, especially on NUMA nodes.

This model can be supported in two ways:

1. Run multiple OS processes at each node, each with one MPI rank, and use (Unix V like) shared memory segments to communicate with shared memory between the processes.
2. Run one OS process per node, but provide support for multiple “MPI processes” within one address space.

The use of Unix V like shared memory segments introduces various overheads and complexities, and is not supported by higher level shared memory programming models, such as OpenMP, and by debuggers and performance tools. This document describes an MPI extension for supporting the second approach.

Most MPI libraries, but not all, equate an MPI process with an OS process. This equivalence is not an MPI requirement: The MPI standard says [3, §2.7]:

An MPI program consists of autonomous processes, executing their own code, in a MIMD style. The codes executed by each process need not be identical. The processes communicate via calls to MPI communication primitives. Typically, each process executes in its own address space, although shared-memory implementations of MPI are possible.

The standard does not define “process” and does not equate it with an OS process. In fact, there are MPI implementations where an MPI agent is an OS thread [1, 5], or even a task that can be dynamically scheduled by a run-time on different threads, for load balancing [2].

### 0.1.2 Proposed Approach

We use the following terminology

**An OS process** consists of an address space, a program, one or more threads and, possibly, additional system resources, such as MPI endpoints.

**A thread** consists of a stack, program counter and other state associated with an execution.

**An MPI endpoint** consists of a set of resources that enable MPI calls. Such an endpoint is identified in MPI by a rank in MPI\_COMM\_WORLD or a rank in derived communicators.

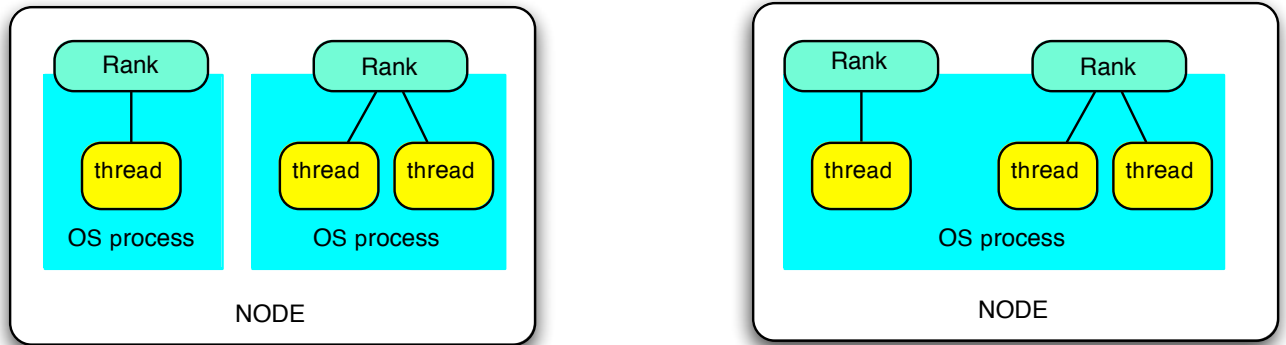


Figure 1: Current and new design for MPI

An **MPI process** consists of an MPI endpoint and a set of threads that can perform MPI calls using that endpoint. All references to a “process” in other parts of the MPI standard should be understood as referring to an “MPI process”.

Most MPI implementations can support multiple MPI processes within one OS image, each running in a distinct address space, as shown in Figure 1, left. The same mechanisms can be used to create multiple endpoints within one address space, as shown in Figure 1, right.

Changing the nature of an MPI process does not change in any way the semantics of MPI. It also requires very few changes in the MPI software stack – the changes will mostly be in the initialization code that allocates communication resources to an MPI process.

Currently, an MPI library creates an MPI endpoint for each OS process at initialization time; all threads running in this process are implicitly associated with this endpoint.

In the new design, multiple MPI endpoints may be bound to the same OS process. The binding of MPI endpoints to processes will occur when MPI is initialized. On the other hand, the binding of threads to endpoints can change during execution. At any point in time, a thread will be attached to at most one endpoint; an MPI send executed by such a thread will use the (sender) rank of that endpoint.

### 0.1.3 Outline

The remainder of the document is organized as follows:

Proposed extensions to MPI to enable the association of multiple MPI endpoints with distinct threads within one process are described in Section 0.2. We describe this proposal assuming the existence of threads, but without making assumptions on the properties of thread, other than that they run within a shared address space.

Section 0.3 describes how the MPI endpoints constructs interoperate with shared memory programming models such as Posix thread libraries, OpenMP and PGAS languages.

## 0.2 MPI Support of Multiple Endpoints per Process

### 0.2.1 MPI Endpoints

An *MPI endpoint* is a (handle to a) set of resources that supports the independent execution of MPI communications. These can be physical resources (e.g., registers mapped into the address space of the process), or logical resources. An endpoint corresponds to a rank in an MPI communicator. A thread can be *attached* to an endpoint, at which point it can make MPI calls using the resources of the endpoint. Each process owns one or more endpoints; the association is static. Each thread running within this process can be attached to at most one endpoint; the association is dynamic.

In MPI with no endpoints, there is a fixed one-to-one correspondence between MPI processes and ranks in `MPI_COMM_WORLD`; when a thread executes an MPI call with argument `MPI_COMM_WORLD` then the rank of the caller is the rank of the MPI process containing this thread.

Similarly, in our proposal, there is a one-to-one correspondence between MPI endpoints and ranks in `MPI_COMM_ENDPOINTS`; when a thread that is attached to an MPI endpoint executes an MPI call with argument `MPI_COMM_ENDPOINTS` then the caller's rank is the rank of the attached endpoint. Thus, if a thread is attached to endpoint 5, then a call by that thread to `MPI_SEND(..., MPI_COMM_ENDPOINTS)` will appear as a send by the MPI process with rank 5 in `MPI_COMM_ENDPOINTS`. Similarly, if a thread executes

```
MPI_Comm_Dup(MPI_COMM_ENDPOINTS, newcomm);  
MPI_Send(..., newcomm);
```

Then the send appears to be executed by the MPI process with rank 5 in `newcomm`.

### 0.2.2 Initialization

The number of available endpoints per OS process is determined by the MPI job startup command, e.g., `mpirun`; see Section 0.2.4.

A call to an MPI initialization routine `MPI_INIT` or `MPI_THREAD_INIT` generates four communicators:

<b>MPI_COMM_ENDPOINTS</b>	communicator that includes all endpoints
<b>MPI_COMM_WORLD</b>	Communicator that includes the first endpoint of each OS process
<b>MPI_COMM_PROCESS</b>	Communicator that includes all endpoints within the local OS process
<b>MPI_COMM_SELF</b>	Communicator that contains exactly one endpoint

Endpoints within each process are ordered by endpoint number and have consecutive ranks in `MPI_COMM_ENDPOINTS` and `MPI_COMM_PROCESS`. This is illustrated in Figure 2.

All MPI programs must contain exactly one call per OS process to an MPI initialization routine.

*Advice to users.* `MPI_COMM_WORLD` and derived communicators should be used by code that is unaware of multiple endpoints per process; `MPI_COMM_ENDPOINTS` and derived communicators should be used by code that is aware of multiple endpoints per process. (*End of advice to users.*)



*Advice to users.* Note that a thread attaches to an endpoint, i.e., migrates to an MPI process. The endpoint can have different ranks in different communicators. Once attached to endpoint, the thread can communicate using any of these communicators. Thus, in the example of Figure 2, a thread in the second OS process can execute the following code:

```

1 ...
2 MPI_Thread_attach(2, MPLCOMMPROCESS);
3 MPI_Send(..., MPLCOMMPROCESS); // sender rank is 2
4 MPI_Send(..., MPLCOMMENDPOINTS); // sender rank is 4
5 ...

```

*(End of advice to users.)*

## 0.2.4 Communication With Endpoints

Unless said otherwise, whenever “process” is mentioned in the MPI standard, it should be understood to mean “MPI process”, i.e., an endpoint and all the threads attached to that endpoint. The rules listed below follow from this interpretation.

A thread must be attached to an endpoint (i.e. must be part of an MPI process) in order to make MPI calls other than `MPI_INIT`, `MPI_INIT_THREAD`, `MPI_GET_VERSION`, `MPI_INITIALIZED`, `MPI_FINALIZED` and `MPI_THREAD_ATTACH`. A nonblocking call started on an endpoint must be completed on the same endpoint. The rules and restrictions specified by the MPI standard [3, §12.4] for threads continue to apply. In particular, a blocking MPI call will block the thread that executes the call, but will not affect other threads. The blocked thread will continue execution when the call completes. Since each thread can be attached to only one endpoint, deadlock situations do not arise. Two distinct threads should not block on the same request.

*Advice to implementors.* The support of multiple MPI processes within one address space requires that the MPI library maintain, for each thread, the identity of the endpoint the thread is attached to, if any. No other changes are needed. Communication using the `MPI_THREAD_FUNNELED` model, with  $k$  endpoints in one OS process within an OS image, should be performing as well or better than communication with  $k$  single-threaded OS processes, each with one endpoint, in that OS image. *(End of advice to implementors.)*

## Progress

The MPI standard specifies situations where progress on an MPI call at an agent might depend on the execution of matching MPI calls at other agents. Thus, a blocking send operation might not return until a matching receive is executed; a blocking call to a collective operation call might not return until the call is invoked by all other processes in the communicator; and so on. On the other hand, a non-blocking send or non-blocking collective will return irrespective of activities at other MPI processes.

These rules are extended to the situation where an OS process may have multiple endpoints: a blocking send on an endpoint might not return until a matching receive has occurred at the destination endpoint; and a blocking call to a collective operation might not return until the operation is invoked at all endpoints of the communicator (including

endpoints at the same OS process). On the other hand, a non-blocking send or a non-blocking call to a collective operation will return, irrespective of the activities of threads attached to other endpoints (including threads in the same address space).

The same rules dictate progress when an endpoint has no attached thread. An endpoint with no attached thread might prevent progress of an MPI call if the progress of that call depends on the execution of a matching MPI call at that endpoint, but will not prevent progress of other MPI calls. Thus, a blocking send might not return if no thread is attached to the destination endpoint; a blocking call to a collective operation call might not return if one of the endpoints in the communicator has no attached thread. However, a non-blocking send will return even if there is no thread attached to the destination endpoint; and a non-blocking call to a collective operation will return even if one of the endpoints in the communicator has no attached thread.

## Thread Support

If the level of thread support is `MPI_THREAD_FUNNELED` then only one thread can execute MPI calls on each endpoint; different threads within the same OS process can execute concurrently MPI calls on different endpoints. If the level of thread support is `MPI_THREAD_SERIALIZED` then no two threads can make concurrent MPI calls on the same endpoint.

## Finalize

`MPI_FINALIZE` must be invoked once at each OS process. The call should be invoked only after all pending MPI calls at that process have completed. (This is an exception to the rule – we do not require a separate call for each MPI process.)

## Memory Allocation

Memory allocated by `MPI_ALLOC_MEM` [3, §6.2] can be used only for communication with the endpoint the calling thread is attached to.

*Rationale.* Endpoints may be supported by distinct adapters, each requiring different memory areas for efficient communication. (*End of rationale.*)

## Error Handling

Error handlers are attached to endpoints. A communicator may have different error handlers attached to the different endpoints of that communicator within the same address space.

The same rule applies to error handlers attached to windows or files.

## Process Manager Interface

The function `MPI_COMM_SPAWN` function can be used to spawn OS processes with multiple endpoints, with the same number of endpoints at each OS process. The number of spawned OS processes is specified by the argument `maxprocs`. The number of endpoints per OS process is specified by the value of the reserved key `num_endpoints` in the `info` argument. If the key `num_endpoints` is not provided in the `info` argument, then the new communicator has one endpoint per OS process. The call returns an error of class `MPI_ERR_SPAWN` if it cannot generate the required number of endpoints per process. If it succeeds, then it returns

an intercommunicator that contains the parent endpoints in the local group and the first endpoint of each spawned process in the remote group. This intercommunicator has, as local group, the endpoints of the `comm` argument of `MPI_COMM_SPAWN`, and, as a remote group, the endpoints of newly created `MPI_COMM_WORLD`. This intercommunicator will be returned at a newly spawned process by the call to `MPI_GET_PARENT` (for threads attached to the first endpoint of the OS process). One error code per OS process is returned in `array_of_errcodes` argument. The MPI initialization call at the newly spawned OS processes return four communicators, as specified in Section 0.2.2.

The function `MPI_COMM_SPAWN_MULTIPLE` is extended in a similar manner. The function is passed multiple array arguments. The values associated with the key `num_endpoints` in the  $i$ -th entry of the `info` array argument specifies the number of endpoints to generate in each of the processes that execute the  $i$ -th command.

### Portable Process Startup

*Advice to implementors.* If the implementation provides an `mpiexec` function then it is advised that the number of endpoints per OS process be indicated by the parameter `-endpoints i i`. (*End of advice to implementors.*)

### Windows

An invocation to `MPI_WIN_CREATE(base, size, disp_unit, info, comm, win)` may return a different window for different endpoints in the same address space. Windows associated with different endpoints in the same address space may overlap. However, the outcome of a code where conflicting accesses occur to a location that appears in two windows is undefined.

### I/O

The invocation to `MPI_FILE_OPEN` returns a distinct file handle at each endpoint. Note that the function is collective and must be invoked at all endpoints of the communicator with the same file name and access mode arguments.

An invocation to `MPI_FILE_SET_VIEW` can set a different view of the file for each file handle argument (passing different `disp`, `filetype` or `info` arguments) – hence a different view at each endpoint within the same address space.

Data access calls that use individual file pointers (such as `MPI_FILE_READ`) maintain a distinct file pointer for each file handle; hence different endpoints within the same address space are associated with distinct individual file pointers.

## 0.3 Interoperability

This section discusses how the endpoint constructs can interact with other programming models, such as Posix threads, OpenMP and PGAS languages.

### 0.3.1 Posix Binding

An MPI library is compatible with a POSIX thread library if the behavior described in the previous section holds for threads spawned by the POSIX thread library.



### 0.3.2 example

We present below a simple “hello world” program. The program creates multiple endpoints at each process, and attaches one thread to each endpoint. The thread attached to endpoint zero gathers the thread ids of all the spawned threads and prints them.

Listing 1: Simple MPI hybrid program

```
1 #include <mpi.h>
2 #include <posix.h>
3 #include <stdio.h>
4 #define max_endpoints 32
5
6 pthread_t thread[max_endpoints];
7
8 /* code executed by each thread */
9 void *foo(void *id)
10 {
11     int i = (int)id;
12     long tid = (long)thread[i];
13     int rank, size;
14     long *recvbuf;
15
16     MPI_Thread_attach(i, MPLCOMMPROCESS);
17     MPI_Comm_size(MPLCOMMENDPOINTS, size);
18     MPI_Comm_rank(MPLCOMMENDPOINTS, rank);
19     if (rank == 0)
20         recvbuf = (long *)malloc(size*sizeof(long));
21     MPI_Gather(&tid, 1, MPLLONG, &recvbuf, 1, MPLLONG, 0,
22               MPIENDPOINTS);
23     if (rank==0)
24     {
25         printf("number_of_endpoints_is_%d;_their_thread_ids_are:",
26               size);
27         for (int j=0; j<size; j++)
28             printf("_%d_", *(recvbuf+j));
29     }
30     pthread_exit(NULL);
31 }
32
33 int main()
34 {
35     MPI_Init_endpoint(NULL, NULL, max_endpoints, max_thread_level,
36                       size, rank);
37
38     /* create a thread for each endpoint */
39     for (int i=0; i < max_endpoints; i++)
40         pthread_create(&thread[i], NULL, foo, (void *)i);
41     pthread_exit(NULL);
```

### 0.3.3 OpenMP

An MPI library is compatible with OpenMP if the behavior described in the previous section holds for threads as defined in the OpenMP standard.

### 0.3.4 PGAS languages

A UPC or Fortran implementation that is compatible with MPI will provide a mechanism for identifying which threads (in UPC) or images (in Fortran) are within the same address space. The MPI library should support the creation of at least one MPI endpoint per thread (in UPC) or image (in Fortran).

### 0.3.5 Porting Codes to Hybrid MPI

Codes written with the current MPI interfaces port without change, and use one endpoint per OS process. The transition from current MPI codes to codes using multiple endpoints per OS process can entail one or more of the following scenarios:

1. Code is written to utilize multiple endpoints per OS process and leverage shared memory communication within OS processes.
2. Libraries are written so that they can be invoked in parallel by one thread per endpoint; they compute correctly irrespective of the number of endpoints per OS process. Such portability is important for MPI libraries invoked from UPC or Fortran: The library can have one MPI process per UPC thread (or Fortran image), and its behavior will not depend on the number of UPC threads per OS process.
3. Code written to use a single endpoint per OS process invokes a library written to use multiple endpoints per OS process. This will enable recoding only critical kernels, with no changes to the overall program logic.

MPI code can be written so that it has the same outcome, whether each MPI process is a distinct OS process, or multiple MPI processes run within the same address space. The only part of the code that has to be aware of the system configuration (i.e., the number of MPI processes per OS process) is the initialization code that creates the endpoints and attaches threads to endpoints. To achieve this portability one needs to ensure that threads belonging to different MPI processes do not interfere with each other (i.e., do not communicate inadvertently). Consider the simple case of a code with single-threaded MPI processes written using C, C++ or Fortran. Problems arise if the code uses mutable static variables: When the threads are in distinct address spaces then the threads have distinct instances of these variables; but when they run in the same address space, they would share the same instance.

In C and C++, static heap variables can be made thread-private by declaring them with the storage class keyword `__thread`. This storage specifier ensures that there will be one separate instance of the declared variable for each thread. While not standard, this extension is widely supported [4, §5.54]. This extension is not currently supported in Fortran – we hope this will change. Alternatively, a preprocessor can be used for replacing

each static heap variable with a dynamic one that is allocated separately by each thread – see, e.g., [2]. Additional care must be taken to ensure that the library code invokes only thread-safe libraries.

# Bibliography

- [1] E.D. Demaine. A threads-only MPI implementation for the development of parallel programs. In *Proceedings of the 11th International Symposium on High Performance Computing Systems*, pages 153–163. Citeseer, 1997. [0.1.1](#)
- [2] C. Huang, O. Lawlor, and L.V. Kale. Adaptive mpi. *Lecture notes in computer science*, pages 306–322, 2003. [0.1.1](#), [0.3.5](#)
- [3] MPI Forum. MPI: A Message-Passing Interface Standard V2.2, 2009. [0.1.1](#), [0.2.4](#), [0.2.4](#)
- [4] Richard M. Stallman and the GCC Developer Community. Using the GNU Compiler Collection (for GCC version 4.4.2). [0.3.5](#)
- [5] H. Tang, K. Shen, and T. Yang. Program transformation and runtime support for threaded MPI execution on shared-memory machines. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 22(4):673–700, 2000. [0.1.1](#)