

D R A F T

Document for a Standard Message-Passing Interface

Message Passing Interface Forum

December 26, 2010

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

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

Chapter 14

Profiling Interface

14.1 Requirements

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

1. provide a mechanism through which all of the MPI defined [functions]functions, except those allowed as macros (See Section 2.6.5[)], **may be accessed with a name shift**. This requires, in C and Fortran, an alternate entry point name, with the prefix `PMPI_` for each MPI function. The profiling interface in C++ is described in Section 16.1.10. For routines implemented as macros, it is still required that the `PMPI_` version be supplied and work as expected, but it is not possible to replace at link time the `MPI_` version with a user-defined version.
2. ensure that those MPI functions that are not replaced may still be linked into an executable image without causing name clashes.
3. document the implementation of different language bindings of the MPI interface if they are layered on top of each other, so that the profiler developer knows whether she must implement the profile interface for each binding, or can [economise]economize by implementing it only for the lowest level routines.
4. where the implementation of different language bindings is done through a layered approach ([e.g.]e.g., the Fortran binding is a set of “wrapper” functions that call the C implementation), ensure that these wrapper functions are separable from the rest of the library.

This separability is necessary to allow a separate profiling library to be correctly implemented, since (at least with Unix linker semantics) the profiling library must contain these wrapper functions if it is to perform as expected. This requirement allows the person who builds the profiling library to extract these functions from the original MPI library and add them into the profiling library without bringing along any other unnecessary code.

5. provide a no-op routine `MPI_PCONTROL` in the MPI library.

14.2 Discussion

The objective of the MPI profiling interface is to ensure that it is relatively easy for authors of profiling (and other similar) tools to interface their codes to MPI implementations on different machines.

Since MPI is a machine independent standard with many different implementations, it is unreasonable to expect that the authors of profiling tools for MPI will have access to the source code that implements MPI on any particular machine. It is therefore necessary to provide a mechanism by which the implementors of such tools can collect whatever performance information they wish *without* access to the underlying implementation.

We believe that having such an interface is important if MPI is to be attractive to end users, since the availability of many different tools will be a significant factor in attracting users to the MPI standard.

The profiling interface is just that, an interface. It says *nothing* about the way in which it is used. There is therefore no attempt to lay down what information is collected through the interface, or how the collected information is saved, filtered, or displayed.

While the initial impetus for the development of this interface arose from the desire to permit the implementation of profiling tools, it is clear that an interface like that specified may also prove useful for other purposes, such as “internetworking” multiple MPI implementations. Since all that is defined is an interface, there is no objection to its being used wherever it is useful.

As the issues being addressed here are intimately tied up with the way in which executable images are built, which may differ greatly on different machines, the examples given below should be treated solely as one way of implementing the objective of the MPI profiling interface. The actual requirements made of an implementation are those detailed in the Requirements section above, the whole of the rest of this chapter is only present as justification and discussion of the logic for those requirements.

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

14.3 Logic of the Design

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

14.3.1 Miscellaneous Control of Profiling

There is a clear requirement for the user code to be able to control the profiler dynamically at run time. This is normally used for (at least) the purposes of

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

These requirements are met by use of the MPI_PCONTROL.

```
MPI_PCONTROL(level, ...)
```

```
IN      level          Profiling level
```

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

```
MPI_PCONTROL(LEVEL)
    INTEGER LEVEL
```

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

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

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

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

- `level==0` Profiling is disabled.
- `level==1` Profiling is enabled at a normal default level of detail.
- `level==2` Profile buffers are `[flushed. (This may be a no-op in some profilers).]flushed, which may be a no-op in some profilers.`
- All other values of `level` have profile library defined effects and additional arguments.

We also request that the default state after MPI_INIT has been called is for profiling to be enabled at the normal default level. (i.e. as if MPI_PCONTROL had just been called with the argument 1). This allows users to link with a profiling library and obtain profile output without having to modify their source code at all.

The provision of MPI_PCONTROL as a no-op in the standard MPI library `[allows them to modify their source code to obtain]supports the collection of` more detailed profiling information`[, but still be able to link exactly the]with source [same code]code that can still link` against the standard MPI library.

14.4 Examples

14.4.1 Profiler Implementation

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

```
static int totalBytes = 0;
static double totalTime = 0.0;
```

```

1  int MPI_Send(void* buffer, int count, MPI_Datatype datatype,
2              int dest, int tag, MPI_Comm comm)
3  {
4      double tstart = MPI_Wtime();      /* Pass on all the arguments */
5      int extent;
6      int result    = PMPI_Send(buffer, count, datatype, dest, tag, comm);
7
8      MPI_Type_size(datatype, &extent); /* Compute size */
9      totalBytes += count*extent;
10
11     totalTime += MPI_Wtime() - tstart;      /* and time      */
12
13     return result;
14 }

```

14.4.2 MPI Library Implementation

ticket0. 17 [On a Unix system, in which the MPI library is implemented in C, then] If the MPI library
 ticket0. 18 is implemented in C on a Unix system, then there [there are various possible options, of
 ticket0. 19 which two of the most obvious] are various options, including the two presented here, for
 ticket0. 20 supporting [are presented here. Which is better depends on whether the linker and] the
 ticket0. 21 name-shift requirement. The choice between these two options [compiler support weak
 ticket0. 22 symbols.] depends partly on whether the linker and compiler support weak symbols.
 ticket0. 23

24 Systems with Weak Symbols

ticket0. 26 If the compiler and linker support weak external symbols ([e.g.]e.g., Solaris 2.x, other system
 ticket0. 27 V.4 machines), then only a single library is required through the use of `#pragma weak` thus
 ticket0. 28

```

29 #pragma weak MPI_Example = PMPI_Example
30
31 int PMPI_Example(/* appropriate args */)
32 {
33     /* Useful content */
34 }
35

```

36 The effect of this `#pragma` is to define the external symbol `MPI_Example` as a weak
 37 definition. This means that the linker will not complain if there is another definition of the
 38 symbol (for instance in the profiling library), however if no other definition exists, then the
 39 linker will use the weak definition.
 40

41 Systems Without Weak Symbols

42 In the absence of weak symbols then one possible solution would be to use the `C` macro
 43 pre-processor thus
 44

```

45 #ifndef PROFILELIB
46 #   ifdef __STDC__
47 #       define FUNCTION(name) P##name
48

```

```

#   else
#       define FUNCTION(name) P/**/name
#   endif
#else
#   define FUNCTION(name) name
#endif

```

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

```

int FUNCTION(MPI_Example)(/* appropriate args */)
{
    /* Useful content */
}

```

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

It is required that the standard MPI library be built in such a way that the inclusion of MPI functions can be achieved one at a time. This is a somewhat unpleasant requirement, since it may mean that each external function has to be compiled from a separate file. However this is necessary so that the author of the profiling library need only define those MPI functions that she wishes to intercept, references to any others being fulfilled by the normal MPI library. Therefore the link step can look something like this

```
% cc ... -lmyprof -lpmpi -lmpi
```

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

14.4.3 Complications

Multiple Counting

Since parts of the MPI library may themselves be implemented using more basic MPI functions ([e.g.]e.g., a portable implementation of the collective operations implemented using point to point communications), there is potential for profiling functions to be called from within an MPI function that was called from a profiling function. This could lead to “double counting” of the time spent in the inner routine. Since this effect could actually be useful under some circumstances ([e.g.]e.g., it might allow one to answer the question “How much time is spent in the point to point routines when they’re called from collective functions?”), we have decided not to enforce any restrictions on the author of the MPI library that would overcome this. Therefore the author of the profiling library should be aware of this problem, and guard against it herself. In a single threaded world this is easily achieved through use of a static variable in the profiling code that remembers if you are already inside a profiling routine. It becomes more complex in a multi-threaded environment (as does the meaning of the times recorded[!]).

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 ticket0.
35
36
37
38 ticket0.
39
40
41
42
43
44
45 ticket0.
46 ticket0.
47
48

1 Linker Oddities

2 The Unix linker traditionally operates in one `[pass :]pass`: the effect of this is that functions
 3 from libraries are only included in the image if they are needed at the time the library is
 4 scanned. When combined with weak symbols, or multiple definitions of the same function,
 5 this can cause odd (and unexpected) effects.

6 Consider, for instance, an implementation of MPI in which the Fortran binding is
 7 achieved by using wrapper functions on top of the C implementation. The author of the
 8 profile library then assumes that it is reasonable only to provide profile functions for the C
 9 binding, since Fortran will eventually call these, and the cost of the wrappers is assumed
 10 to be small. However, if the wrapper functions are not in the profiling library, then none
 11 of the profiled entry points will be undefined when the profiling library is called. Therefore
 12 none of the profiling code will be included in the image. When the standard MPI library
 13 is scanned, the Fortran wrappers will be resolved, and will also pull in the base versions of
 14 the MPI functions. The overall effect is that the code will link successfully, but will not be
 15 profiled.

16 To overcome this we must ensure that the Fortran wrapper functions are included in
 17 the profiling version of the library. We ensure that this is possible by requiring that these
 18 be separable from the rest of the base MPI library. This allows them to be aared out of the
 19 base library and into the profiling one.

22 14.5 Multiple Levels of Interception

23 The scheme given here does not directly support the nesting of profiling functions, since it
 24 provides only a single alternative name for each MPI function. Consideration was given to
 25 an implementation that would allow multiple levels of call interception, however we were
 26 unable to construct an implementation of this that did not have the following disadvantages

- 28 • assuming a particular implementation language[.],
- 29
- 30 • imposing a run time cost even when no profiling was taking place.

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

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

37 [Unfortunately such an implementation may require more cooperation between the
 38 different profiling libraries than is required for the single level implementation detailed
 39 above.]Note, however, that it is possible to use the scheme above to implement a multi-level
 40 system, since the function called by the user may call many different profiling functions
 41 before calling the underlying MPI function. This capability has been demonstrated in the
 42 P^NMPI tool infrastructure [1].

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

Bibliography

- [1] Martin Schulz and Bronis R. de Supinski. P^NMPI Tools: A Whole Lot Greater Than the Sum of Their Parts. In *ACM/IEEE Supercomputing Conference (SC)*, pages 1–10. ACM, 2007. [14.5](#)

Index

EXAMPLES:Profiling interface, [3](#)

MPI_INIT, [3](#)

MPI_PCONTROL, [1](#), [3](#)

MPI_PCONTROL(level, ...), [3](#)

MPI_SEND, [3](#)

MPI_TYPE_SIZE, [3](#)

MPI_WTIME, [3](#)

PMPI_, [1](#)