

Parallel Computing  
for  
Science and Engineering

Victor Eijkhout

1st edition 2015

---

**Public draft - open for comments**

---

This book will be open source under CC-BY license.

---

Two of the most common software systems for parallel programming in scientific computing are MPI and OpenMP. They target different types of parallelism, and use very different constructs. Thus, by covering both of them in one book we can offer a treatment of parallelism that spans a large range of possible applications.

# Contents

<b>I MPI</b>	<b>9</b>
1	<b>Getting started with MPI</b> 10
1.1	<i>Distributed memory and message passing</i> 10
1.2	<i>History</i> 10
1.3	<i>Basic model</i> 11
1.4	<i>Making and running an MPI program</i> 12
1.5	<i>Language bindings</i> 12
2	<b>MPI topic 1: Functional parallelism</b> 14
2.1	<i>The SPMD model</i> 14
2.2	<i>Processor identification</i> 16
2.3	<i>Functional parallelism</i> 17
3	<b>MPI topic 2: Global information</b> 18
3.1	<i>Working with global information</i> 18
3.2	<i>Collectives</i> 19
3.3	<i>Rooted collectives: broadcast, reduce</i> 20
3.4	<i>Rooted collectives: gather and scatter</i> 21
3.5	<i>Variable-size-input collectives</i> 22
3.6	<i>Scan operations</i> 22
3.7	<i>User-defined reductions</i> 23
3.8	<i>Reduce-scatter</i> 23
3.9	<i>'All'-type collectives</i> 24
3.10	<i>Non-blocking collectives</i> 24
3.11	<i>Barrier and all-to-all</i> 25
3.12	<i>Performance of collectives</i> 25
3.13	<i>Collectives and synchronization</i> 26
4	<b>MPI topic 3: Distributed data</b> 29
4.1	<i>Distributed computing and distributed data</i> 29
4.2	<i>Local information exchange</i> 31
4.3	<i>Shared-memory-like communication: one-sided communication</i> 43
4.4	<i>Remaining topics in point-to-point communication</i> 50
5	<b>MPI topic 4: Dealing with complicated data</b> 52
5.1	<i>Data types</i> 52
6	<b>MPI topic 5: Sub computations</b> 59

6.1	<i>Subcommunications</i>	59
6.2	<i>Communicators</i>	60
7	<b>MPI topics</b>	65
7.1	<i>Synchronization</i>	65
7.2	<i>Hybrid programming: MPI and threads</i>	65
7.3	<i>Leftover topics</i>	66
7.4	<i>Literature</i>	72
8	<b>MPI Reference</b>	73
8.1	<i>Basics</i>	73
8.2	<i>Data types</i>	76
8.3	<i>Blocking communication</i>	84
8.4	<i>Deadlock-free blocking messages</i>	89
8.5	<i>One-sided communication</i>	96
8.6	<i>Collectives</i>	104
8.7	<i>Communicators</i>	113
8.8	<i>Leftover topics</i>	116
8.9	<i>Error handling</i>	118
8.10	<i>More utility stuff</i>	119
8.11	<i>Multi-threading</i>	120
9	<b>MPI Examples</b>	122
9.1	<i>A</i>	122
9.2	<i>B</i>	124
9.3	<i>C</i>	124
9.4	<i>E</i>	126
9.5	<i>F</i>	127
9.6	<i>G</i>	127
9.7	<i>I</i>	129
9.8	<i>P</i>	130
9.9	<i>R</i>	131
9.10	<i>S</i>	134
9.11	<i>T</i>	137
9.12	<i>W</i>	140
10	<b>MPI Review</b>	144
10.1	<i>Review questions</i>	144

II	<b>OpenMP</b>	145
11	<b>OpenMP tutorial</b>	146
11.1	<i>Basics</i>	146
11.2	<i>Work sharing</i>	152
11.3	<i>Controlling thread data</i>	160
11.4	<i>Reductions</i>	163
11.5	<i>Synchronization</i>	165

11.6	<i>Tasks</i>	171
11.7	<i>Version 4 functionality</i>	173
11.8	<i>Stuff</i>	173
11.9	<i>Performance</i>	179
12	<b>OpenMP Reference</b>	180
12.1	<i>Basics</i>	180
12.2	<i>Parallel regions</i>	181
12.3	<i>Worksharing</i>	182
12.4	<i>Controlling thread data</i>	185
12.5	<i>Synchronization</i>	186
12.6	<i>Internal control variables</i>	188
12.7	<i>Tasks</i>	189
12.8	<i>Stuff</i>	189
13	<b>OpenMP Review</b>	191
13.1	<i>Concepts review</i>	191
13.2	<i>Review questions</i>	192

### III The Rest 199

14	<b>Random number generation</b>	200
15	<b>Hybrid computing</b>	201
15.1	<i>Discussion</i>	201
15.2	<i>Hybrid MPI-plus-threads execution</i>	201
16	<b>Support libraries</b>	203

### IV Tutorials 205

16.1	<i>Debugging</i>	207
16.2	<i>Tracing</i>	216

### V Projects, index 217

17	<b>Class projects</b>	218
17.1	<i>A Style Guide to Project Submissions</i>	218
17.2	<i>Warmup Exercises</i>	220
17.3	<i>Mandelbrot set</i>	224
17.4	<i>Data parallel grids</i>	231
18	<b>Bibliography, index, and list of acronyms</b>	234
18.1	<i>Bibliography</i>	234
18.2	<i>List of acronyms</i>	235
18.3	<i>Index</i>	236



## **PART I**

### **MPI**

# **Chapter 1**

## **Getting started with MPI**

In this chapter you will learn the use of the main tool for distributed memory programming: the Message Passing Interface (MPI) library. The MPI library has about 250 routines, many of which you may never need. Since this is a textbook, not a reference manual, we will focus on the important concepts and give the important routines for each concept. What you learn here should be enough for most common purposes. You are advised to keep a reference document handy, in case there is a specialized routine, or to look up subtleties about the routines you use.

### **1.1 Distributed memory and message passing**

In its simplest form, a distributed memory machine is a collection of single computers hooked up with network cables. In fact, this has a name: a *Beowulf cluster*. As you recognize from that setup, each processor can run an independent program, and has its own memory without direct access to other processors' memory. MPI is the magic that makes multiple instantiations of the same executable run so that they know about each other and can exchange data through the network.

One of the reasons that MPI is so successful as a tool for high performance on clusters is that it is very explicit: the programmer controls many details of the data motion between the processors. Consequently, a capable programmer can write very efficient code with MPI. Unfortunately, that programmer will have to spell things out in considerable detail. For this reason, people sometimes call MPI ‘the assembly language of parallel programming’. If that sounds scary, be assured that things are not that bad. You can get started fairly quickly with MPI, using just the basics, and coming to the more sophisticated tools only when necessary.

Another reason that MPI was a big hit with programmers is that it does not ask you to learn a new language: it is a library that can be interface to C/C++ or Fortran; there are even bindings to Python. A related point is that it is easy to install: there are free implementations that you can download and install on any computer that has a Unix-like operating system, even if that is not a parallel machine.

### **1.2 History**

Before the MPI standard was developed in 1993-4, there were many libraries for distributed memory computing, often proprietary to a vendor platform. MPI standardized the inter-process communication mecha-

nisms. Other features, such as process management in *PVM*, or parallel I/O were omitted. Later versions of the standard have included many of these features.

Since MPI was designed by a large number of academic and commercial participants, it quickly became a standard. A few packages from the pre-MPI era, such as *Charmpp* [5], are still in use since they support mechanisms that do not exist in MPI.

### 1.3 Basic model

Here we sketch the two most common scenarios for using MPI. In the first, the user is working on an interactive machine, which has network access to a number of hosts, typically a network of workstations; see figure 1.1. The user types the command `mpiexec`<sup>1</sup> and supplies

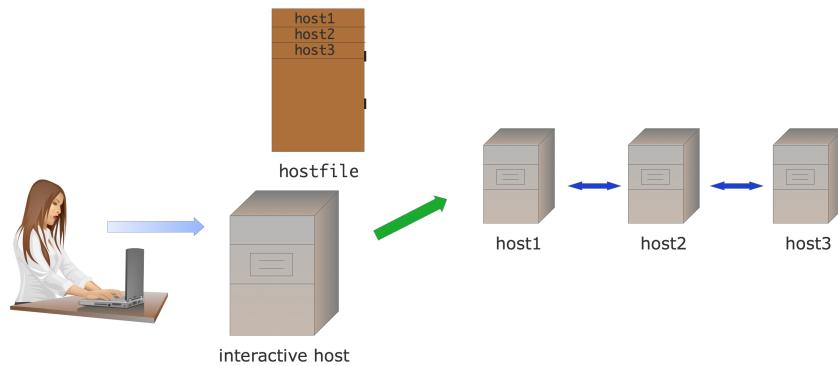


Figure 1.1: Interactive MPI setup

- The number of hosts involved,
- their names, possibly in a hostfile,
- and other parameters, such as whether to include the interactive host; followed by
- the name of the program and its parameters.

The `mpirun` program then makes an `ssh` connection to each of the hosts, giving them sufficient information that they can find each other. All the output of the processors is piped through the `mpirun` program, and appears on the interactive console.

In the second scenario (figure 1.2) the user prepares a *batch job* script with commands, and these will be run when the *batch scheduler* gives a number of hosts to the job. Now the batch script contains the `mpirun` command, or some variant such as `ibrun`, and the hostfile is dynamically generated when the job starts. Since the job now runs at a time when the user may not be logged in, any screen output goes into an output file.

You see that in both scenarios the parallel program is started by the `mpirun` command using an Single Program Multiple Data (SPMD) mode of execution: all hosts execute the same program. It is possible for different hosts to execute different programs, but we will not consider that in this book.

---

1. A command variant is `mpirun`; your local cluster may have a different mechanism.

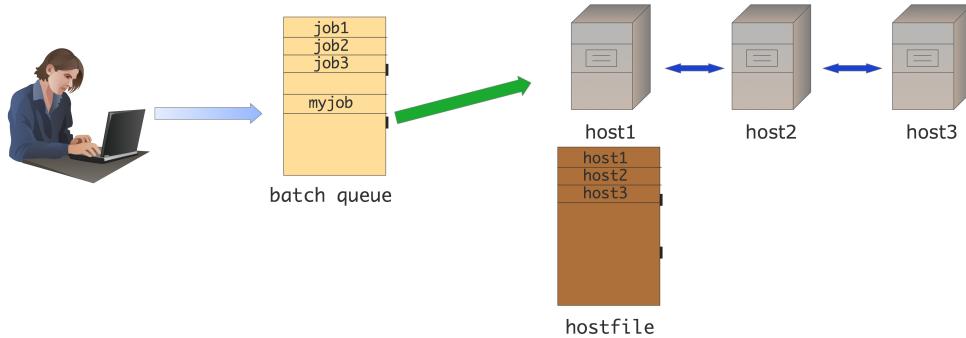


Figure 1.2: Batch MPI setup

## 1.4 Making and running an MPI program

MPI is a library, called from programs in ordinary programming languages such as C/C++ or Fortran. To compile such a program you use your regular compiler:

```
gcc -c my_mpi_prog.c -I/path/to/mpi.h
gcc -o my_mpi_prog my_mpi_prog.o -L/path/to/mpi -lmpich
```

However, MPI libraries may have different names between different architectures, making it hard to have a portable makefile. Therefore, MPI typically has shell scripts around your compiler call:

```
mpicc -c my_mpi_prog.c
mpicc -o my_mpi_prog my_mpi_prog.o
```

MPI programs can be run on many different architectures. Obviously it is your ambition (or at least your dream) to run your code on a cluster with a hundred thousand processors and a fast network. But maybe you only have a small cluster with plain *ethernet*. Or maybe you're sitting in a plane, with just your laptop. An MPI program can be run in all these circumstances – within the limits of your available memory of course.

The way this works is that you do not start your executable directly, but you use a program, typically called `mpirun` or something similar, which makes a connection to all available processors and starts a run of your executable there. So if you have a thousand nodes in your cluster, `mpirun` can start your program once on each, and if you only have your laptop it can start a few instances there. In the latter case you will of course not get great performance, but at least you can test your code for correctness.

## 1.5 Language bindings

### 1.5.1 C/C++

The MPI library is written in C. Thus, its bindings are the most natural for that language.

C++ bindings existed at one point, but they were declared deprecated. The *boost* library has its own version of MPI. A recent effort at idiomatic C++ support is *MPL* <http://numbercrunch.de/blog/2015/08/mp1-a-message-passing-library/>.

### 1.5.2 Fortran

The *Fortran bindings* for MPI look very much like the C ones, except that each routine has a final *error return* parameter.

*Fortran note* Other Fortran-specific differences will be indicated with a note like this.

### 1.5.3 Python

The `mpi4py` package of *python bindings* is not defined by the MPI standards committee. Instead, it is the work of an individual, *Lisandro Dalcin*.

Notable about the Python bindings is that many communication routines exist in two variants:

- a version that can send native Python objects. These routines have lowercase names such as `bcast`; and
- a version that sends *numpy* objects; these routines have names such as `Bcast`. Their syntax can be slightly different.

The first version looks more ‘pythonic’, is easier to write, and can do things like sending python objects, but it is also decidedly less efficient since data is packed and unpacked with `pickle`. As a common sense guideline, use the *numpy* interface in the performance-critical parts of your code, and the native interface only for complicated actions in a setup phase.

Codes with `mpi4py` can be interfaced to other languages through Swig or conversion routines.

Data in *numpy* can be specified as a simple object, or `[data, (count,displ), datatype]`.

## Chapter 2

### MPI topic 1: Functional parallelism

#### 2.1 The SPMD model

MPI programs conform<sup>1</sup> to the Single Program Multiple Data (SPMD) model, where each processor runs the same executable. This running executable we call a *process*.

When MPI was first written, 20 years ago, it was clear what a processor was: it was what was in a computer on someone's desk, or in a rack. If this computer was part of a networked cluster, you called it a *node*. So if you ran an MPI program, each node would have one MPI process; figure 2.1.

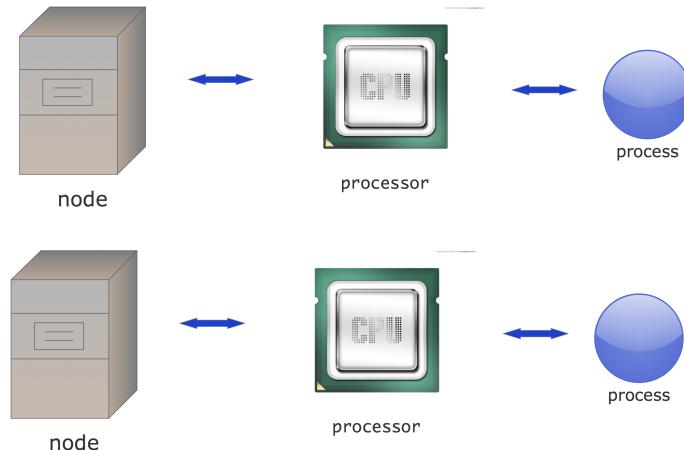


Figure 2.1: Cluster structure as of the mid 1990s

These days the situation is more complicated. You can still talk about a node in a cluster, but now a node can contain more than one processor chip (sometimes called a *socket*), and each processor chip probably has multiple cores. Figure 2.2 shows how you could explore this using a mix of MPI between the nodes, and a shared memory programming system on the nodes.

1. Usually, but not necessarily.

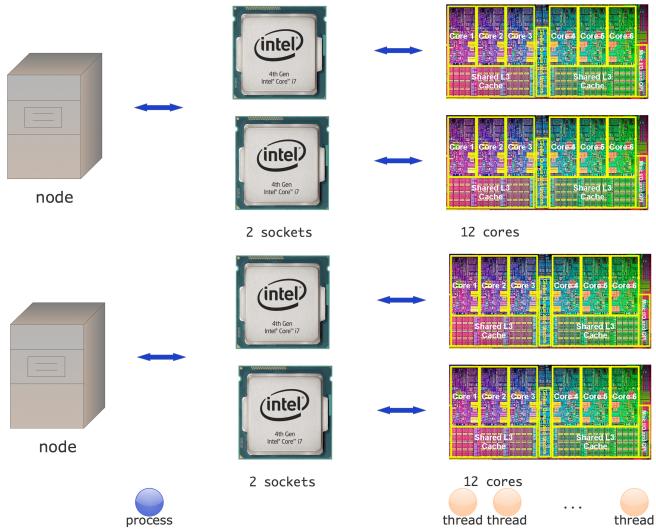


Figure 2.2: Hybrid cluster structure

However, since each core can act like an independent processor, you can also have multiple MPI processes per node. To MPI the cores look like the old completely separate processors. This is the ‘pure MPI’ model of figure 2.3 which we will use in most of this part of the book.

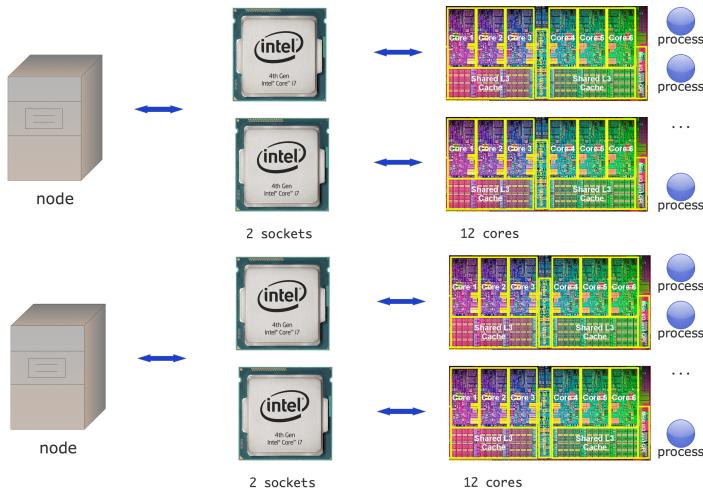


Figure 2.3: MPI-only cluster structure

This is somewhat confusing: the old processors needed MPI programming, because they were physically separated. The cores on a modern processor, on the other hand, share the same memory, and even some caches. In its basic mode MPI ignores all of this: each core receives an MPI process and they communicate as if they are all connected through the same network. In fact, you can't immediately see whether two cores are on the same node or different nodes.

### 2.1.1 Starting and running MPI processes

*The reference for the commands introduced here can be found in section 8.1.1.*

The SPMD model may be initially confusing. Even though there is only a single source, compiled into a single executable, the parallel run comprises a number of independently started MPI processes (see section 1.3 for the mechanism).

The following exercises are designed to give you an intuition for this one-source-many-processes setup. In the first exercise you will see that the mechanism for starting MPI programs starts up independent copies. There is nothing in the source that says ‘and now you become parallel’.

The following exercise shows you that

**Exercise 2.1.** Write a ‘hello world’ program, without any MPI in it, and run it in parallel with `mpiexec` or your local equivalent. Explain the output.

To get a useful MPI program you need at least the calls `MPI_Init` and `MPI_Finalize` surrounding your code. See section 8.1.1 for their syntax.

*Python note* There are no initialize and finalize calls: the `import` statement performs the initialization.

This may look a bit like declaring ‘this is the parallel part of a program’, but that’s not true: again, the whole code is executed multiple times in parallel.

**Exercise 2.2.** Add the commands `MPI_Init` and `MPI_Finalize` to your code. Put three different print statements in your code: one before the init, one between init and finalize, and one after the finalize. Again explain the output.

In the following exercise you will print out the hostname of each MPI process; see section 8.10.1 for the syntax.

**Exercise 2.3.** Now use the command `MPI_Get_processor_name` in between the init and finalize statement, and print out on what processor your process runs. Confirm that you are able to run a program that uses two different nodes.

## 2.2 Processor identification

*The reference for the commands introduced here can be found in section 8.1.4.*

Since all processes in an MPI job are instantiations of the same executable, you’d think that they all execute the exact same instructions, which would not be terribly useful. To distinguish between processors, MPI provides two calls

1. `MPI_Comm_size` reports how many processes there are in all; and
2. `MPI_Comm_rank` states what the number of the process is.

In other words, each process can find out ‘I am process 5 out of a total of 20’.

**Exercise 2.4.** Write a program where each process prints out message reporting its number, and how many processes there are.

Write a second version of this program, where each process opens a unique file and writes to it. *On some clusters this may not be advisable if you have large numbers of processors, since it can overload the file system.*

**Exercise 2.5.** Write a program where only the process with number zero reports on how many processes there are in total.

### 2.3 Functional parallelism

Being able to tell processes apart is already enough for some applications. Based on its rank, a processor can find its section of a search space. For instance, in *Monte Carlo* codes a large number of random samples is generated and some computation performed on each. (This particular example requires each MPI process to run an independent random number generator, which is not entirely trivial.)

**Exercise 2.6.** Is the number  $N = 2,000,000,111$  prime? Let each process test a range of integers, and print out any factor they find. You don't have to test all integers  $< N$ : any factor is at most  $\sqrt{N} \approx 45,200$ .

As another example, in *Boolean satisfiability* problems a number of points in a search space needs to be evaluated. Knowing a process's rank is enough to let it generate its own portion of the search space. The computation of the *Mandelbrot set* can also be considered as a case of functional parallelism. However, the image that is constructed is data that needs to be kept on one processor, which breaks the symmetry of the decomposition.

Of course, at the end of a functionally parallel run you need to summarize the results, for instance printing out some total. The mechanisms for that you will learn next.

## Chapter 3

### MPI topic 2: Global information

#### 3.1 Working with global information

If all processes have individual data, for instance the result of a local computation, you may want to bring that information together, for instance to find the maximal computed value or the sum of all values. Conversely, sometimes one processor has information that needs to be shared with all. For this sort of operation, MPI has *collectives*.

There are various cases, the most common ones are illustrated in figure 3.1.

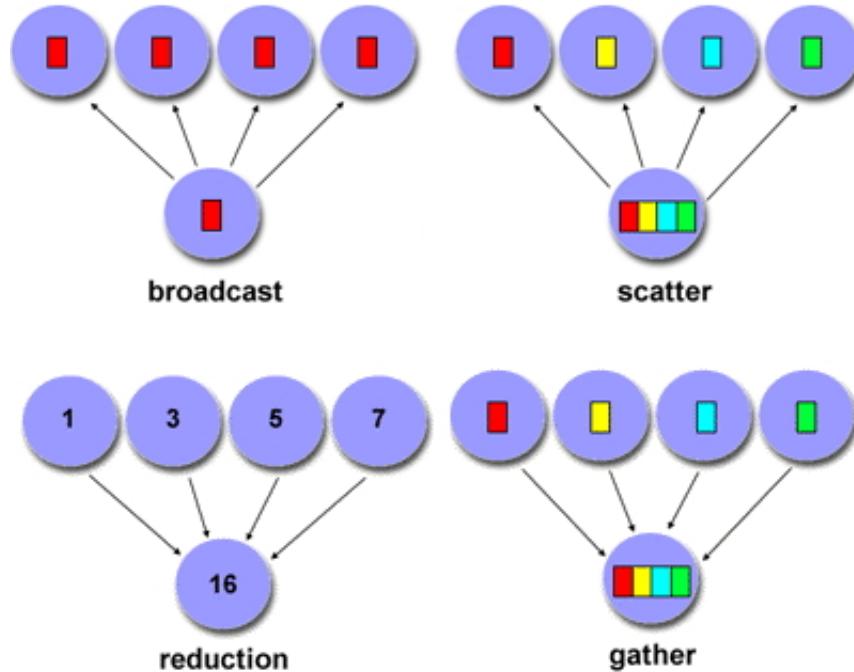


Figure 3.1: The four most common collectives

Above, you saw how each process can perform its own computation with its own result. You may want to summarize these results on one process, known as the *root process*, for instance to print them out. If you

perform an operation on the data from the processors, for instance to compute the maximum value, this is known as a *reduction* (section 3.3). On the other hand, if you need to collect and preserve all computation results, the operation is known as a *gather* (section 3.4).

Conversely, one process can have data that needs to be spread to all others, for instance because it reads it from file. If the same item needs to be sent to all processes, this is known as *broadcast*. If the root process sends individual data to each process, it is called a *scatter*.

**Exercise 3.1.** How would you realize the following scenarios with MPI collectives?

- Let each process compute a random number. You want to print the maximum of these numbers to your screen.
- Each process computes a random number again. Now you want to scale these numbers by their maximum.
- Let each process compute a random number. You want to print on what processor the maximum value is computed.

There are more collectives or variants on the above.

- If you want to gather or scatter information, but the contribution of each processor is of a different size, there are ‘variable’ collectives; they have a *v* in the name (section 3.5).
- Sometimes you want a reduction with partial results, where each processor computes the sum (or other operation) on the values of lower-numbered processors. For this, you use a *scan* collective (section 3.6).
- If every processor needs to broadcast to every other, you use an *all-to-all* operation (section 3.9).
- A barrier is an operation that makes all processes wait until every process has reached the barrier (section 3.11).

Finally, there are some advanced topics in collectives.

- Non-blocking collectives; section 3.10.
- User-defined reduction operators; section 3.7.

## 3.2 Collectives

*The reference for the commands introduced here can be found in section ??.*

Collectives are operations that involve all processes in a communicator. (See section 3.2 for an informal listing.) A collective is a single call, and it blocks on all processors. That does not mean that all processors exit the call at the same time: because of implementational details and network latency they need not be synchronized in their execution. However, semantically we can say that a process can not finish a collective until every other process has at least started the collective.

In addition to these collective operations, there are operations that are said to be ‘collective on their communicator’, but which do not involve data movement. Collective then means that all processors must call this routine; not to do so is an error that will manifest itself in ‘hanging’ code. One such example is MPI\_Win\_fence.

### 3.3 Rooted collectives: broadcast, reduce

*The reference for the commands introduced here can be found in section 8.6.1.*

One simple collective is the broadcast, where one process has some data that needs to be shared with all others. One scenario is that processor zero can parse the commandline arguments of the executable and send the values to all other processors. Another scenario is that you want one processor to read data from file and send it to the other processors: this is likely to be more efficient than having every process open the file.

The broadcast call has the following structure:

```
MPI_Bcast( data..., root , comm);
```

The root is the process that is sending its data. Typically, it will be the root of a broadcast tree. The `comm` argument is a communicator: for now you can use `MPI_COMM_WORLD`. Unlike with send/receive there is no message tag, because collectives are blocking, so you can have only one collective active at a time.

The data in a broadcast (or any other MPI operation for that matter) is specified as

- A buffer. In C this is the address in memory of the data. This means that you broadcast a single scalar as `MPI_Bcast( &value, ... )`, but an array as `MPI_Bcast( array, ... )`.
- The number of items and their datatype. The allowable datatypes are such things as `MPI_INT` and `MPI_FLOAT` for C, and `MPI_INTEGER` and `MPI_REAL` for Fortran, or more complicated types. See section 8.2 for details.

*Python note* In python it is both possible to send objects, and to send more C-like buffers. The two possibilities correspond (see section 1.5.3) to different routine names; the buffers have to be created as numpy objects.

**Exercise 3.2.** If you give a commandline argument to a program, that argument is available as a character string as part of the `argv`, `argc` pair that you typically use as the arguments to your main program. You can use the function `atoi` to convert such a string to integer.

Write a program where process 0 looks for an integer on the commandline, and broadcasts it to the other processes. Initialize the buffer on all processes, and let all processes print out the broadcast number, just to check that you solved the problem correctly.

The reverse of a broadcast is a reduction:

```
MPI_Reduce( senddata, recvdata..., operator,
            root, comm );
```

Now there is a separate buffer for outgoing data, on all processors, and incoming data, only relevant on the root. Also, you have to indicate how the data is to be combined. Popular choices are `MPI_SUM`, `MPI_PROD` and `MPI_MAX`, but complicated operators such as finding the location of the maximum value exist. You can also define your own operators; section 8.6.8.

**Exercise 3.3.** Write a program where each process computes a random number, and process 0 finds and prints the maximum generated value. Let each process print its value, just to check the correctness of your program.

Now let each process scale its value by this maximum.

Collective operations can also take an array argument, instead of just a scalar. In that case, the operation is applied pointwise to each location in the array.

**Exercise 3.4.** Create on each process an array of length 2 integers, and put the values 1, 2 in it on each process. Do a sum reduction on that array. Can you predict what the result should be? Code it. Was your prediction right?

## 3.4 Rooted collectives: gather and scatter

The reference for the commands introduced here can be found in section 8.6.3.

In the `MPI_Scatter` operation, the root spreads information to all other processes. The difference with a broadcast is that it involves individual information from/to every process. Thus, the gather operation typically has an array of items, one coming from each sending process, and scatter has an array, with an

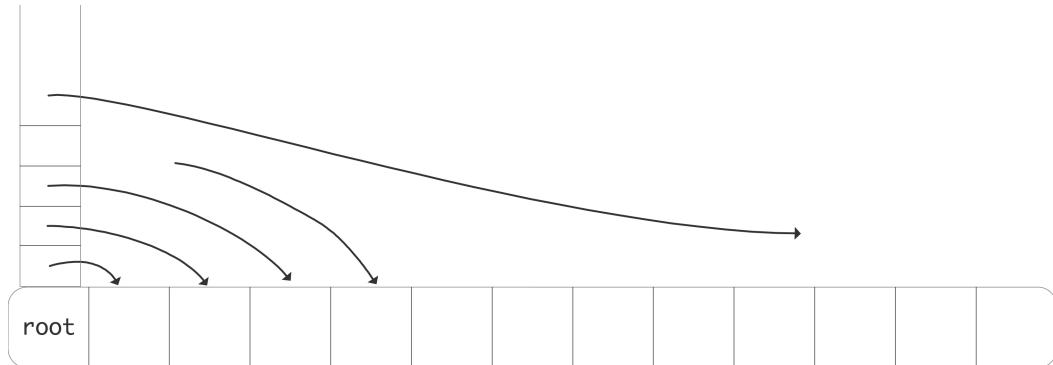


Figure 3.2: A scatter operation

individual item for each receiving process; see figure 3.2.

These gather and scatter collectives have a different parameter list from the broadcast/reduce. The broadcast/reduce involves the same amount of data on each process, so it was enough to have a buffer, datatype, and size. In the gather/scatter calls you have

- a large buffer on the root, with a datatype and size specification, and
- a smaller buffer on each process, with its own type and size specification.

Of course, since we're in SPMD mode, even non-root processes have the argument for the send buffer, but they ignore it. For instance:

```
int MPI_Scatter
    (void* sendbuf, int sendcount, MPI_Datatype sendtype,
     void* recvbuf, int recvcount, MPI_Datatype recvtype,
```

```
int root, MPI_Comm comm)
```

The `sendcount` is not, as you might expect, the total length of the sendbuffer; instead, it is the amount of data sent to each process.

**Exercise 3.5.** Let each process compute a random number. You want to print on what processor the maximum value is computed. What collective do you use? Write a short program.

### 3.5 Variable-size-input collectives

*The reference for the commands introduced here can be found in section 8.6.6.*

In the gather and scatter call above each processor received or sent an identical number of items. In many cases this is appropriate, but sometimes each processor wants or contributes an individual number of items.

Let's take the gather calls as an example. Assume that each processor does a local computation that produces a number of data elements, and this number is different for each processor (or at least not the same for all). In the regular `MPI_Gather` call the root processor had a buffer of size  $nP$ , where  $n$  is the number of elements produced on each processor, and  $P$  the number of processors. The contribution from processor  $p$  would go into locations  $pn, \dots, (p + 1)n - 1$ .

For the variable case, we first need to compute the total required buffer size. This can be done through a simple `MPI_Reduce` with `MPI_SUM` as reduction operator: the buffer size is  $\sum_p n_p$  where  $n_p$  is the number of elements on processor  $p$ . But you can also postpone this calculation for a minute.

The next question is where the contributions of the processor will go into this buffer. For the contribution from processor  $p$  that is  $\sum_{q < p} n_p, \dots, \sum_{q \leq p} n_p - 1$ . To compute this, the root processor needs to have all the  $n_p$  numbers, and it can collect them with an `MPI_Gather` call.

We now have all the ingredients. All the processors specify a send buffer just as with `MPI_Gather`. However, the receive buffer specification on the root is more complicated. It now consists of:

```
outbuffer, array-of-outcounts, array-of-displacements, outtype
```

and you have just seen how to construct that information.

### 3.6 Scan operations

*The reference for the commands introduced here can be found in section 8.6.7.*

The `MPI_Scan` operation also performs a reduction, but it keeps the partial results. That is, if processor  $i$  contains a number  $x_i$ , and  $\oplus$  is an operator, then the scan operation leaves  $x_0 \oplus \dots \oplus x_i$  on processor  $i$ .

```
MPI_Scan( send data, recv data, operator, communicator);
```

This type of operation is often called a *prefix operation*; see HPSC-23.

The MPI\_Scan routine is an *inclusive scan* operation. Often, the more useful variant is the *exclusive scan* MPI\_Exscan

```
MPI_Exscan( send data, recv data, operator, communicator);
```

with the same prototype.

**Exercise 3.6.** The exclusive definition, which computes  $x_0 \oplus x_{i-1}$  on processor  $i$ , can easily be derived from the inclusive operation for operations such as MPI\_PLUS or MPI\_MULT. Are there operators where that is not the case?

The MPI\_Scan operation is often useful with indexing data. Suppose that every processor  $p$  has a local vector where the number of elements  $n_p$  is dynamically determined. In order to translate the local numbering  $0 \dots n_p - 1$  to a global numbering one does a scan with the number of local elements as input. The output is then the global number of the first local variable.

**Exercise 3.7.** Do you use MPI\_Scan or MPI\_Exscan for this operation? How would you describe the result of the other scan operation, given the same input?

It is possible to do a *segmented scan*. Let  $x_i$  be a series of numbers that we want to sum to  $X_i$  as follows. Let  $y_i$  be a series of booleans such that

$$\begin{cases} X_i = x_i & \text{if } y_i = 0 \\ X_i = X_{i-1} + x_i & \text{if } y_i = 1 \end{cases}$$

(This is the basis for the implementation of the *sparse matrix vector product* as prefix operation; see HPSC-23.2.) This means that  $X_i$  sums the segments between locations where  $y_i = 0$  and the first subsequent place where  $y_i = 1$ . To implement this, you need a user-defined operator

$$\begin{pmatrix} X \\ x \\ y \end{pmatrix} = \begin{pmatrix} X_1 \\ x_1 \\ y_1 \end{pmatrix} \oplus \begin{pmatrix} X_2 \\ x_2 \\ y_2 \end{pmatrix} : \begin{cases} X = x_1 + x_2 & \text{if } y_2 == 1 \\ X = x_2 & \text{if } y_2 == 0 \end{cases}$$

This operator is not commutative, and it needs to be declared as such with MPI\_Op\_create; see section 8.6.8

## 3.7 User-defined reductions

*The reference for the commands introduced here can be found in section 8.6.8.*

For use in reductions and scans it is possible to define your own operator.

## 3.8 Reduce-scatter

*The reference for the commands introduced here can be found in section 8.6.4.*

There are several MPI collectives that are functionally equivalent to a combination of others. You have already seen `MPI_Allreduce` which is equivalent to a reduction followed by a broadcast. Often such combinations can be more efficient than using the individual calls; see HPSC-6.1.

Here is another example: `MPI_Reduce_scatter` is equivalent to a reduction on an array of data (meaning a pointwise reduction on each array location) followed by a scatter of this array to the individual processes.

One important example of this command is the *sparse matrix-vector product*; see HPSC-6.5.1 for background information. Each process contains one or more matrix rows, so by looking at indices the process can decide what other processes it needs data from. The problem is for a process to find out what other processes it needs to send data to.

Using `MPI_Reduce_scatter` the process goes as follows:

- Each process creates an array of ones and zeros, describing who it needs data from.
- The reduce part of the reduce-scatter yields an array of requester counts; after the scatter each process knows how many processes request data from it.
- Next, the sender processes need to find out what elements are requested from it. For this, each process sends out arrays of indices.
- The big trick is that each process now knows how many of these requests will be coming in, so it can post precisely that many `MPI_Irecv` calls, with a source of `MPI_ANY_SOURCE`.

## 3.9 ‘All’-type collectives

*The reference for the commands introduced here can be found in section 8.6.5.*

In many applications the result of a collective is needed on all processes. For instance, if  $x, y$  are distributed vector objects, and you want to compute

$$y - (x^t y)x$$

you need the inner product value on all processors. You could do this by writing a reduction followed by a broadcast, but more efficient algorithms exist. Surprisingly, an ‘all-gather’ operation takes as long as a rooted gather (see HPSC-6.1 for details).

Thus, MPI has the following operations:

- `MPI_Allreduce` is equivalent to a `MPI_Reduce` followed by a broadcast.
- `MPI_Allgather` is equivalent to a `MPI_Gather` followed by a broadcast.
- `MPI_Allgatherv` is equivalent to an `MPI_Gatherv` followed by a broadcast.
- `MPI_Alltoall`, `MPI_Alltoallv`.

## 3.10 Non-blocking collectives

*The reference for the commands introduced here can be found in section 8.6.9.*

Above you have seen how the ‘Isend’ and ‘Irecv’ routines can overlap communication with computation. This is not possible with the collectives you have seen so far: they act like blocking sends or receives. However, there are also *non-blocking collectives*. These have roughly the same calling sequence as their blocking counterparts, except that they output an `MPI_Request`. You can then use an `MPI_Wait` call to make sure the collective has completed.

Such operations can be used to increase efficiency. For instance, computing

$$y \leftarrow Ax + (x^t x)y$$

involves a matrix-vector product, which is dominated by computation in the *sparse matrix* case, and an inner product which is typically dominated by the communication cost. You would code this as

```
MPI_Iallreduce( .... x ... , &request);  
// compute the matrix vector product  
MPI_Wait(request);  
// do the addition
```

This can also be used for 3D FFT operations [3]. Occasionally, a non-blocking collective can be used for non-obvious purposes, such as the `MPI_Ibarrier` in [4].

## 3.11 Barrier and all-to-all

The reference for the commands introduced here can be found in section ??.

There are two collectives we have not mentioned yet. A barrier is a call that blocks all processes until they have all reached the barrier call. This call’s simplicity is contrasted with its usefulness, which is very limited. It is almost never necessary to synchronize processes through a barrier: for most purposes it does not matter if processors are out of sync. Conversely, collectives (except the new non-blocking ones) introduce a barrier of sorts themselves.

The all-to-all call is a generalization of a scatter and gather: every process is scattering an array of data, and every process is gathering an array of data. There is also a ‘v’ variant of this routine.

## 3.12 Performance of collectives

It is easy to visualize a broadcast as in figure 3.3: see figure 3.3. the root sends all of its data directly to every other process. While this describes the semantics of the operation, in practice the implementation works quite differently.

The time that a message takes can simply be modeled as

$$\alpha + \beta n,$$

where  $\alpha$  is the *latency*, a one time delay from establishing the communication between two processes, and  $\beta$  is the time-per-byte, or the inverse of the *bandwidth*, and  $n$  the number of bytes sent.

### 3. MPI topic 2: Global information

---

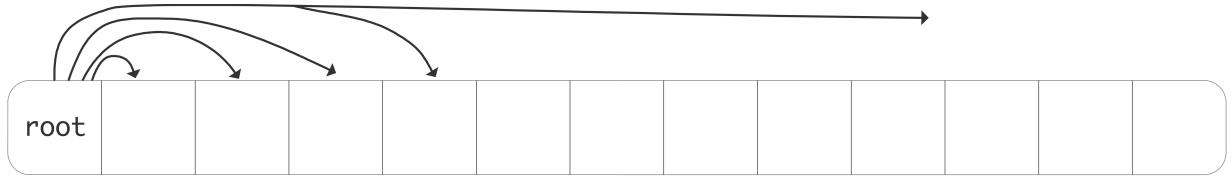


Figure 3.3: A simple broadcast

Under the assumption that a processor can only send one message at a time, the broadcast in figure 3.3 would take a time proportional to the number of processors. One way to ameliorate that is to structure the broadcast in a tree-like fashion. This is depicted in figure 3.4. How does the communication time now

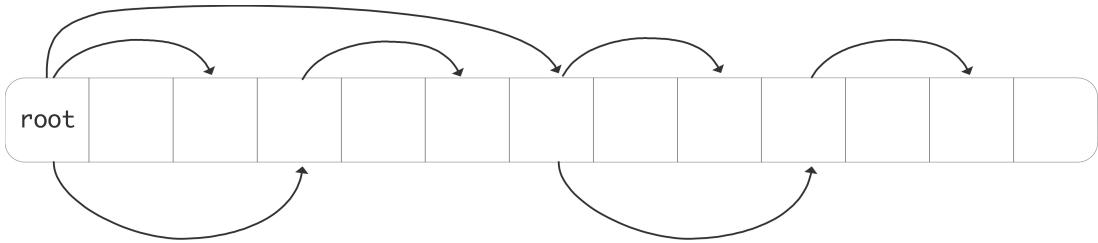


Figure 3.4: A tree-based broadcast

depend on the number of processors? The theory of the complexity of collectives is described in more detail in HPSC-6.1; see also [1].

### 3.13 Collectives and synchronization

Collectives, other than a barrier, have a synchronizing effect between processors. For instance, in

```
MPI_Bcast( ....data... root);  
MPI_Send(....);
```

the send operations on all processors will occur after the root executes the broadcast. Conversely, in a reduce operation the root may have to wait for other processors. This is illustrated in figure 3.5, which gives a TAU trace of a reduction operation on two nodes, with two six-core sockets (processors) each. We see that<sup>1</sup>:

- In each socket, the reduction is a linear accumulation;
- on each node, cores zero and six then combine their result;
- after which the final accumulation is done through the network.

We also see that the two nodes are not perfectly in sync, which is normal for MPI applications. As a result, core 0 on the first node will sit idle until it receives the partial result from core 12, which is on the second node.

---

1. This uses mvapich version 1.6; in version 1.9 the implementation of an on-node reduction has changed to simulate shared memory.

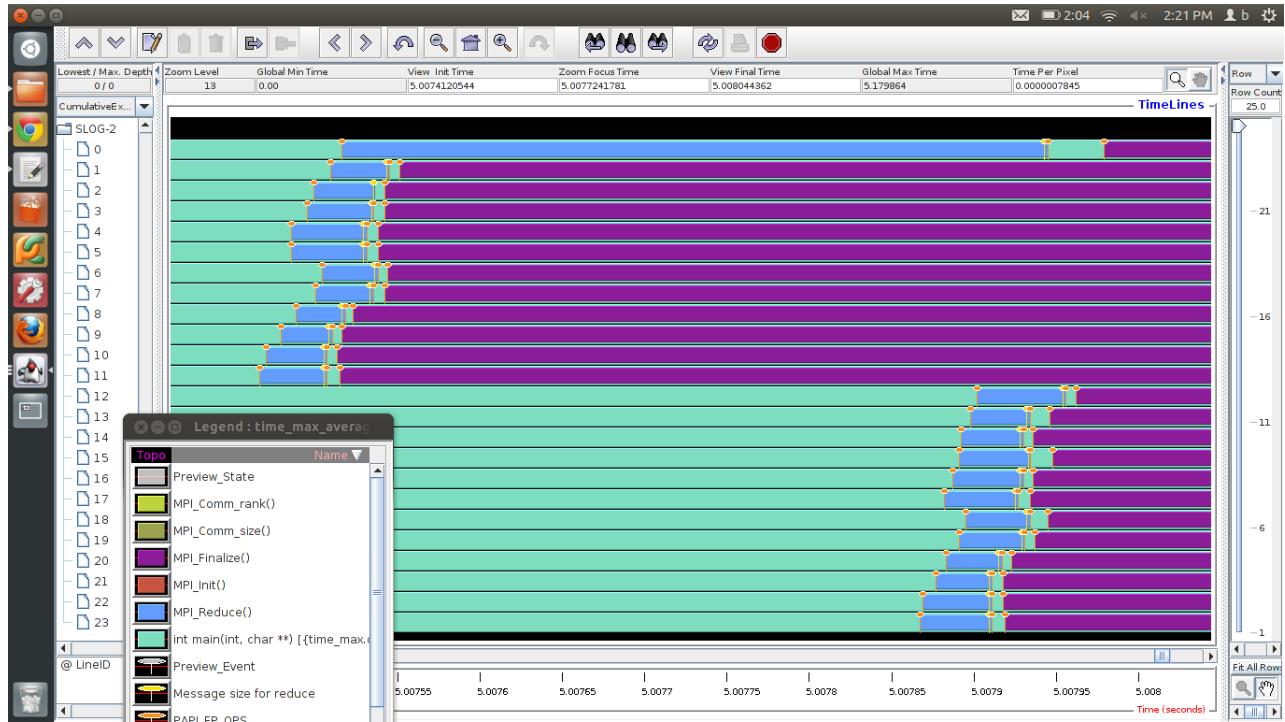


Figure 3.5: Trace of a reduction operation between two dual-socket 12-core nodes

While collectives synchronize in a loose sense, it is not possible to make any statements about events before and after the collectives between processors:

```
...event 1...
MPI_Bcast(....);
...event 2....
```

Consider a specific scenario:

```
switch(rank) {
    case 0:
        MPI_Bcast(buf1, count, type, 0, comm);
        MPI_Send(buf2, count, type, 1, tag, comm);
        break;
    case 1:
        MPI_Recv(buf2, count, type, MPI_ANY_SOURCE, tag, comm, status);
        MPI_Bcast(buf1, count, type, 0, comm);
        MPI_Recv(buf2, count, type, MPI_ANY_SOURCE, tag, comm, status);
        break;
    case 2:
        MPI_Send(buf2, count, type, 1, tag, comm);
        MPI_Bcast(buf1, count, type, 0, comm);
```

### 3. MPI topic 2: Global information

---

```
        break;  
    }
```

Note the MPI\_ANY\_SOURCE parameter in the receive calls on processor 1. One obvious execution of this would be:

1. The send from 2 is caught by processor 1;
2. Everyone executes the broadcast;
3. The send from 0 is caught by processor 1.

However, it is equally possible to have this execution:

1. Processor 0 starts its broadcast, then executes the send;
2. Processor 1's receive catches the data from 0, then it executes its part of the broadcast;
3. Processor 1 catches the data sent by 2, and finally processor 2 does its part of the broadcast.

## Chapter 4

### MPI topic 3: Distributed data

#### 4.1 Distributed computing and distributed data

One reason for using MPI is that sometimes you need to work on more data than can fit in the memory of a single processor. With distributed memory, each processor then gets a part of the whole data structure and only works on that.

So let's say we have a large array, and we want to distribute the data over the processors. That means that, with  $p$  processes and  $n$  elements per processor, we have a total of  $n \cdot p$  elements.

```
int n;  
double data[n];
```

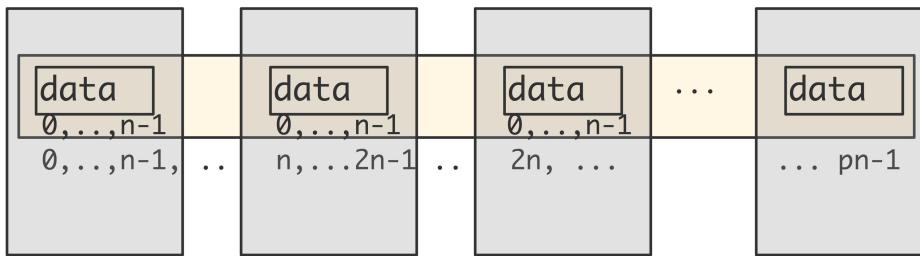


Figure 4.1: Local parts of a distributed array

We sometimes say that *data* is the local part of a *distributed array* with a total size of  $n \cdot p$  elements. However, this array only exists conceptually: each processor has an array with lowest index zero, and you have to translate that yourself to an index in the global array. In other words, you have to write your code in such a way that it acts like you're working with a large array that is distributed over the processors, while actually manipulating only the local arrays on the processors.

Your typical code then looks like

```
int myfirst = .....;  
for (int ilocal=0; ilocal<nlocal; ilocal++) {  
    int iglobal = myfirst+ilocal;  
    array[ilocal] = f(iglobal);
```

}

**Exercise 4.1.** We want to compute  $\sum_{n=1}^N n^2$ , and we do that as follows by filling in an array and summing the elements. (Yes, you can do it without an array, but for purposes of the exercise do it with.)

Read in the global  $N$  parameter, and make sure that it is a multiple of the number  $P$  of processors. Your code should produce an error message and exit immediately if it doesn't.

- Now allocate the local parts: each processor should allocate only  $N/P$  elements.
- On each processor, initialize the local array so that the  $i$ -th location of the distributed array (for  $i = 0, \dots, N - 1$ ) contains  $(i + 1)^2$ .
- Now use a collective operation to compute the sum of the array values. The right value is  $(2N^3 + 3N^2 + N)/6$ . Is that what you get?

To debug your program, first start with  $N = P$ .

(Did you allocate your array as real numbers? Why are integers not a good idea?)

If the array size is not perfectly divisible by the number of processors, we have to come up with a division that is uneven, but not too much. You could for instance, write

```
int Nglobal, // is something large
    Nlocal = Nglobal/ntids,
    excess = Nglobal%ntids;
if (mytid==ntids-1)
    Nlocal += excess;
```

**Exercise 4.2.** Read the section [HPSC-2.10.1](#) about load balancing, and argue that this strategy is not optimal. Can you come up with a better distribution?

One of the more common applications of the reduction operation is the *inner product* computation. Typically, you have two vectors  $x, y$  that have the same distribution, that is, where all processes store equal parts of  $x$  and  $y$ . The computation is then

```
local_inprod = 0;
for (i=0; i<localsize; i++)
    local_inprod += x[i]*y[i];
MPI_Reduce( &local_inprod, &global_inprod, 1,MPI_DOUBLE ... )
```

If all processors need the result, you could then do a broadcast, but it is more efficient to use `MPI_Allreduce`; see section [3.9](#).

**Exercise 4.3.** Implement an inner product routine: let  $x$  be a distributed vector of size  $N$  with elements  $x[i] = i$ , and compute  $x^t x$ . As before, the right value is  $(2N^3 + 3N^2 + N)/6$ .

Use the inner product value to scale to vector so that it has norm 1. Check that your computation is correct.

## 4.2 Local information exchange

Suppose you have an array of numbers  $x_i: i = 0, \dots, N$  and you want to compute  $y_i = (x_{i-1} + x_i + x_{i+1})/3: i = 1, \dots, N - 1$ . As before (see figure 4.1), we give each processor a subset of the  $x_i$ s and  $y_i$ s. Let's define  $i_p$  as the first index of  $y$  that is computed by processor  $p$ . (What is the last index computed by processor  $p$ ? How many indices are computed on that processor?)

We often talk about the *owner computes* model of parallel computing: each processor ‘owns’ certain data items, and it computes their value.

Now let's investigate how processor  $p$  goes about computing  $y_i$  for the  $i$ -values it owns. Let's assume that processor  $p$  also stores the values  $x_i$  for these same indices. Now, it can compute

$$y_{i_p+1} = (x_{i_p} + x_{i_p+1} + x_{i_p+2})/3$$

and likewise  $y_{i_p+2}$  and so on. However, there is a problem with

$$y_{i_p} = (x_{i_p-1} + x_{i_p} + x_{i_p+1})/3$$

since  $x_{i_p}$  is not stored on processor  $p$ : it is stored on  $p - 1$ .

There is a similar story with the last index that  $p$  tries to compute: that involves a value that is only present on  $p + 1$ .

You see that there is a need for processor-to-processor, or technically *point-to-point*, information exchange. MPI realizes this through matched send and receive calls:

- One process does a send to a specific other process;
- the other process does a specific receive from that source.

### 4.2.1 Send example: ping-pong

A simple scenario for information exchange between just two processes is the *ping-pong*: process A sends data to process B, which sends data back to A. This means that process A executes the code

```
MPI_Send( /* to: */ B .... );
MPI_Recv( /* from: */ B ... );
```

while process B executes

```
MPI_Recv( /* from: */ A ... );
MPI_Send( /* to: */ A .... );
```

Since we are programming in SPMD mode, this means our program looks like:

```
if ( /* I am process A */ ) {
    MPI_Send( /* to: */ B .... );
    MPI_Recv( /* from: */ B ... );
} else if ( /* I am process B */ ) {
    MPI_Recv( /* from: */ A ... );
```

```

    MPI_Send( /* to: */ A . . . . );
}

```

Look up the syntax of the send and receive commands in section 8.3, and do the following exercises. You will also need the timer calls of section 7.3.6.

**Exercise 4.4.** Implement the ping-pong program. Add a timer using `MPI_Wtime`. For the status argument of the receive call, use `MPI_STATUS_IGNORE`. Make sure to run each experiment multiple times. (What variance do you see?)

- Run your code with the two communicating processes first on the same node, then on different nodes. Do you see a difference?
- Then modify the program to use longer messages. How does the timing increase with message size?

For bonus points, can you do a regression to determine  $\alpha, \beta$ ?

**Exercise 4.5.** Take your pingpong program and modify it to let half the processors be source and the other half the targets. Does the pingpong time increase?

In the syntax of the `MPI_Recv` command you saw one parameter that the send call lacks: the `MPI_Status` object. This serves the following purpose: the receive call can have a ‘wildcard’ behaviour, for instance specifying that the message can come from any source rather than a specific one. The status object then allows you to find out where the message actually came from.

## 4.2.2 Blocking communication

The reference for the commands introduced here can be found in section 8.3.

The use of `MPI_Send` and `MPI_Recv` is known as *blocking communication*: when your code reaches a send or receive call, it blocks until the call is successfully completed. For a receive call it is clear that the receiving code will wait until the data has actually come in, but for a send call this is more subtle.

You may be tempted to think that the send call puts the data somewhere in the network, and the sending code can progress, as in figure 4.2, left. But this ideal scenario is not realistic: it assumes that somewhere

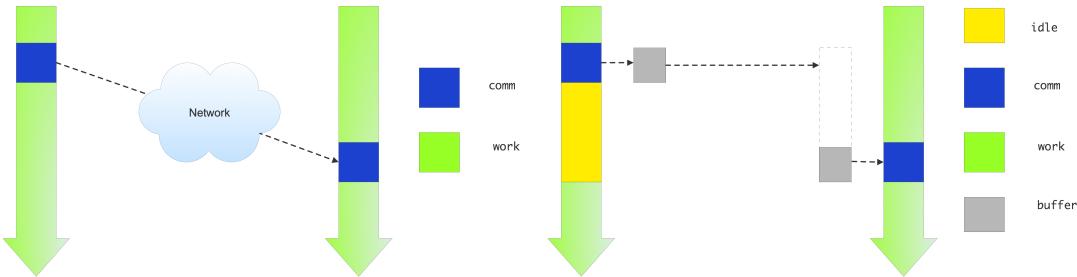


Figure 4.2: Illustration of an ideal (left) and actual (right) send-receive interaction

in the network there is buffer capacity for all messages that are in transit. This is not the case: data resides on the sender, and the sending call blocks, until the receiver has received all of it. (There is an exception for small messages, as explained in the next section.)

### 4.2.3 Problems with blocking communication

Suppose two process need to exchange data, and consider the following pseudo-code, which purports to exchange data between processes 0 and 1:

```
other = 1-mytid; /* if I am 0, other is 1; and vice versa */
receive(source=other);
send(target=other);
```

Imagine that the two processes execute this code. They both issue the send call... and then can't go on, because they are both waiting for the other to issue a receive call. This is known as *deadlock*.

(If you reverse the send and receive call, you should get deadlock, but in practice that code will often work. The reason is that MPI implementations sometimes send small messages regardless of whether the receive has been posted. This relies on the availability of some amount of available buffer space. The size under which this behaviour is used is sometimes referred to as the *eager limit*.)

Formally you can describe deadlock as follows. Draw up a graph where every process is a node, and draw a directed arc from process A to B if A is waiting for B. There is deadlock if this directed graph has a loop.

The solution to the deadlock in the above example is to first do the send from 0 to 1, and then from 1 to 0 (or the other way around). So the code would look like:

```
if ( /* I am processor 0 */ ) {
    send(target=other);
    receive(source=other);
} else {
    receive(source=other);
    send(target=other);
}
```

There is a second, even more subtle problem with blocking communication. Consider the scenario where every processor needs to pass data to its successor, that is, the processor with the next higher rank. The basic idea would be to first send to your successor, then receive from your predecessor. Since the last processor does not have a successor it skips the send, and likewise the first processor skips the receive. The pseudo-code looks like:

```
successor = mytid+1; predecessor = mytid-1;
if ( /* I am not the last processor */ )
    send(target=successor);
if ( /* I am not the first processor */ )
    receive(source=predecessor)
```

This code does not deadlock. All processors but the last one block on the send call, but the last processor executes the receive call. Thus, the processor before the last one can do its send, and subsequently continue to its receive, which enables another send, et cetera.

In one way this code does what you intended to do: it will terminate (instead of hanging forever on a deadlock) and exchange data the right way. However, the execution now suffers from unexpected *serialization*: only one processor is active at any time, so what should have been a parallel operation becomes a sequential

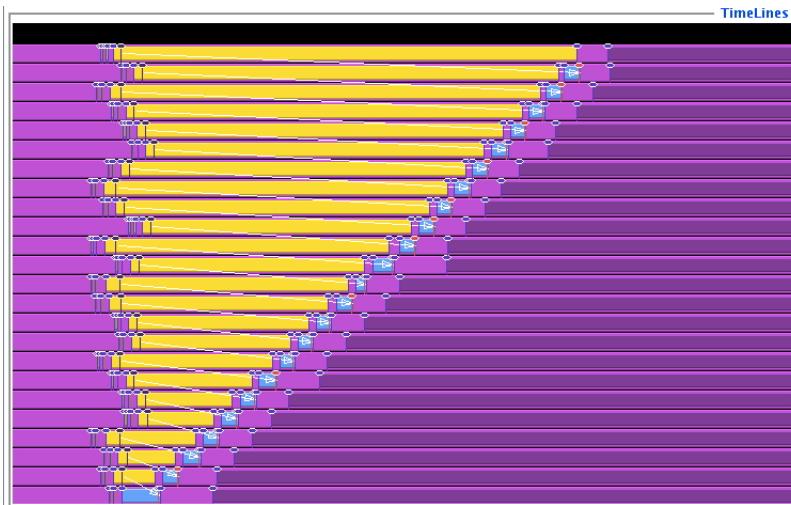


Figure 4.3: Trace of a simple send-recv code

one. This is illustrated in figure 4.3.

**Exercise 4.6.** (Classroom exercise) Each student holds a piece of paper in the right hand – keep your left hand behind your back – and execute the following program:

1. If you are not the rightmost student, turn to the right and give the paper to your right neighbour.
2. If you are not the leftmost student, turn to your left and accept the paper from your left neighbour.

**Exercise 4.7.** Implement the above algorithm using `MPI_Send` and `MPI_Receive` calls.

Run the code, and reproduce the trace output of figure 4.3. See chapter 16.2 on how to use the TAU utility. If you don't have TAU, can you show this serialization behaviour using timings?

It is possible to orchestrate your processes to get an efficient and deadlock-free execution, but doing so is a bit cumbersome.

**Exercise 4.8.** The above solution treated every processor equally. Can you come up with a solution that uses blocking sends and receives, but does not suffer from the serialization behaviour?

There are better solutions which we will explore next.

#### 4.2.4 Pairwise exchange

The reference for the commands introduced here can be found in section 8.4.

Above you saw that with blocking sends the precise ordering of the send and receive calls is crucial. Use the wrong ordering and you get either deadlock, or something that is not efficient at all in parallel. MPI has a way out of this problem that is sufficient for many purposes: the combined send/recv call `MPI_Sendrecv`

```
MPI_Sendrecv( /* send data */ ....
              /* recv data */ .... );
```

The `sendrecv` call works great if every process is paired up. You would then write

```
sendrecv(from=predecessor,to=successor);
```

However, in cases such as the right-shift this is true for all but the first and last. MPI allows for the following variant which makes the code slightly more homogeneous:

```
MPI_Comm_rank( .... &mytid );
if ( /* I am not the first processor */
     predecessor = mytid-1;
else
    predecessor = MPI_PROC_NULL;
if ( /* I am not the last processor */
     successor = mytid+1;
else
    successor = MPI_PROC_NULL;
sendrecv(from=predecessor,to=successor);
```

where the `sendrecv` call is executed by all processors.

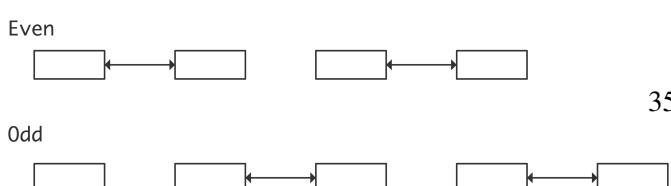
All processors but the last one send to their neighbour; the target value of `MPI_PROC_NULL` for the last processor means a ‘send to the null processor’: no actual send is done. The null processor value is also of use with the `MPI_Sendrecv` call; section 4.2.4

**Exercise 4.9.** Implement the above right-shift scheme using `MPI_Sendrecv`. If you have TAU installed, make a trace. Does it look different from the serialized send/recv code? If you don’t have TAU, run your code with different numbers of processes and show that the runtime is essentially constant.

This call makes it easy to exchange data between two processors: both specify the other as both target and source. However, there need not be any such relation between target and source: it is possible to receive from a predecessor in some ordering, and send to a successor in that ordering; see figure 4.4. Above you saw some examples that had most processors doing both a send and a receive, but some only a send or only a receive. You can still use `MPI_Sendrecv` in this call if you use `MPI_PROC_NULL` for the unused source or target argument.

If the send and receive buffer have the same size, the routine `MPI_Sendrecv_replace` will do an in-place replacement.

The following exercise lets you implement a sorting algorithm with the send-receive call.



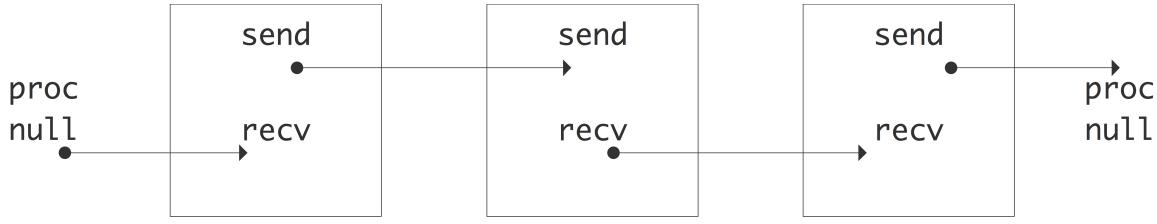


Figure 4.4: An MPI Sendrecv call

**Exercise 4.10.** A very simple sorting algorithm is *exchange sort*: pairs of processors compare data, and if necessary exchange. The elementary step is called a *compare-and-swap*: in a pair of processors each sends their data to the other; one keeps the minimum values, and the other the maximum. For simplicity, in this exercise we give each processor just a single number.

The exchange sort algorithm is split in even and odd stages:

- In the even stage, processors  $2i$  and  $2i + 1$  compare and swap data;
- In the odd stage, processors  $2i + 1$  and  $2i + 2$  compare and swap.

You need to repeat this  $P/2$  times, where  $P$  is the number of processors.

Use

`MPI_PROC_NULL` for the edge cases.

#### 4.2.5 Irregular data exchange

The structure of communication is often a reflection of the structure of the operation. With some regular applications we also get a regular communication pattern. Consider again the above operation:

$$y_i = x_{i-1} + x_i + x_{i+1} : i = 1, \dots, N - 1$$

Doing this in parallel induces communication, as pictured in figure 4.5. We note:

- The data is one-dimensional, and we have a linear ordering of the processors.
- The operation involves neighbouring data points, and we communicate with neighbouring processors.

Above you saw how you can use information exchange between pairs of processors

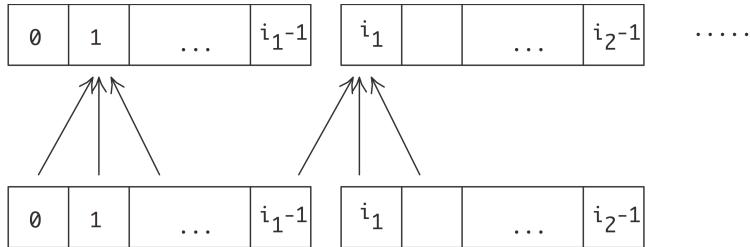


Figure 4.5: Communication in an one-dimensional operation

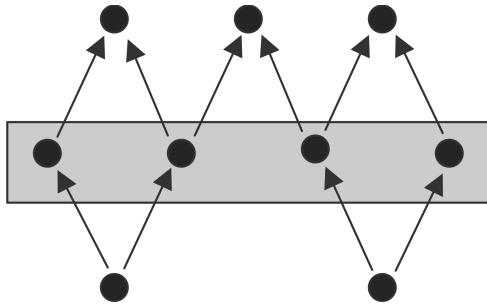


Figure 4.6: Processors with unbalanced send/receive patterns

- using `MPI_Send` and `MPI_Recv`, if you are careful; or
- using `MPI_Sendrecv`, as long as there is indeed some sort of pairing of processors.

However, there are circumstances where it is not possible, not efficient, or simply not convenient, to have such a deterministic setup of the send and receive calls. Figure 4.6 illustrates such a case, where processors are organized in a general graph pattern. Here, the numbers of sends and receive of a processor do not need to match.

In such cases, one wants a possibility to state ‘these are the expected incoming messages’, without having to wait for them in sequence. Likewise, one wants to declare the outgoing messages without having to do them in any particular sequence. Imposing any sequence on the sends and receives is likely to run into the serialization behaviour observed above, or at least be inefficient since processors will be waiting for messages.

#### 4.2.6 Non-blocking communication

*The reference for the commands introduced here can be found in section 8.4.1.*

In the previous section you saw that blocking communication makes programming tricky if you want to avoid deadlock and performance problems. The main advantage of these routines is that you have full control about where the data is: if the send call returns the data has been successfully received, and the send buffer can be used for other purposes or de-allocated.

By contrast, the non-blocking calls `MPI_Irecv` and `MPI_Isend` do not wait for their counterpart: in effect they tell the runtime system ‘here is some data and please send it as follows’ or ‘here is some buffer space, and expect such-and-such data to come’. This is illustrated in figure ??.

While the use of non-blocking routines prevents deadlock, it introduces two new problems:

1. When the send call returns, the actual send may not have been executed, so the send buffer may not be safe to overwrite. When the recv call returns, you do not know for sure that the expected data is in it. Thus, you need a mechanism to make sure that data was actually sent or received.
2. With a blocking send call, you could repeatedly fill the send buffer and send it off.

```
double *buffer;
for ( ... p ... ) {
    buffer = // fill in the data
    MPI_Send( buffer, ... /* to: */ p );
```

To send multiple messages with non-blocking calls you have to allocate multiple buffers.

```
double **buffers;
for ( ... p ... ) {
    buffers[p] = // fill in the data
    MPI_Send( buffers[p], ... /* to: */ p );
```

For the first problem, MPI has two types of routines. The `MPI_Wait...` calls are blocking: when you issue such a call, your execution will wait until the specified requests have been completed. A typical way of using them is:

```
// start non-blocking communication
MPI_Isend( ... ); MPI_Irecv( ... );
// wait for the Isend/Irecv calls to finish in any order
MPI_Wait( ... );
```

**Exercise 4.11.** Now use nonblocking send/receive routines to implement the averaging operation on a distributed array.

There is a second motivation for the `Isend/Irecv` calls: if your hardware supports it, the communication can progress while your program can continue to do useful work:

```
// start non-blocking communication
MPI_Isend( ... ); MPI_Irecv( ... );
// do work that does not depend on incoming data
....
// wait for the Isend/Irecv calls to finish
MPI_Wait( ... );
// now do the work that absolutely needs the incoming data
....
```

This is known as *overlapping computation and communication*, or *latency hiding*.

#### 4.2.6.1 Wait and test calls

The reference for the commands introduced here can be found in section [8.4.1.1](#).

There are several wait calls:

- `MPI_Wait` waits for a single request. If you are indeed waiting for a single nonblocking communication to complete, this is the right routine. If you are waiting for multiple requests you could call this routine in a loop.

```
for (p=0; p<nrequests ; p++)
    MPI_Wait (request [p] , &(status [p])) ;
```

However, this would be inefficient if the first request is fulfilled much later than the others: your waiting process would have lots of idle time. In that case, use one of the following routines.

- `MPI_Waitall` allows you to wait for a number of requests, and it does not matter in what sequence they are satisfied. Using this routine is easier to code than the loop above, and it could be more efficient.
- The ‘waitall’ routine is good if you need all nonblocking communications to be finished before you can proceed with the rest of the program. However, sometimes it is possible to take action as each request is satisfied. In that case you could use `MPI_Waitany` and write:

```
for (p=0; p<nrequests; p++) {
    MPI_Waitany (nrequests, request_array, &index, &status) ;
    // operate on buffer[index]
}
```

Note that this routine takes a single status argument, passed by reference, and not an array of statuses!

- `MPI_Waitsome` is very much like `Waitany`, except that it returns multiple numbers, if multiple requests are satisfied. Now the status argument is an array of `MPI_Status` objects.

Figure 4.7 shows the trace of a non-blocking execution using `MPI_Waitall`.

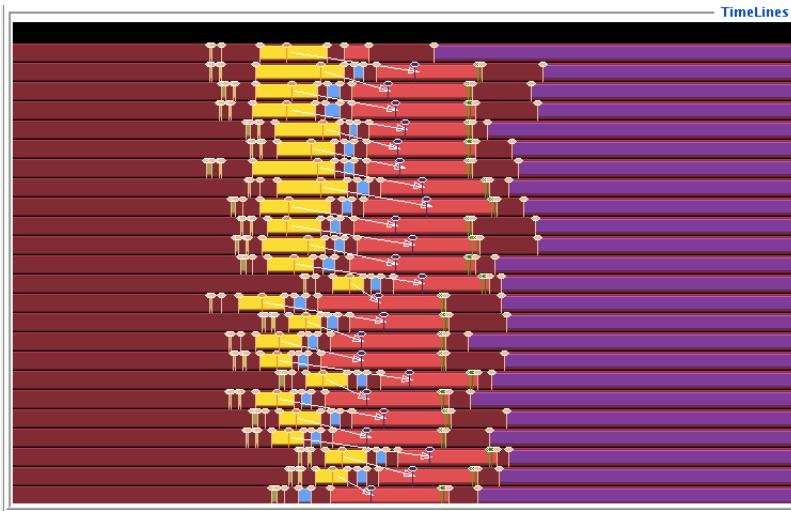


Figure 4.7: A trace of a nonblocking send between neighbouring processors

The `MPI_Wait...` routines are blocking. Thus, they are a good solution if the receiving process can

not do anything until the data (or at least some data) is actually received. The `MPI_Test...` calls are themselves non-blocking: they test for whether one or more requests have been fulfilled, but otherwise immediately return. This can be used in the *master-worker model*: the master process creates tasks, and sends them to whichever worker process has finished its work, but while it waits for the workers it can itself do useful work. Pseudo-code:

```

while ( not done ) {
    // create new inputs for a while
    ....
    // see if anyone has finished
    MPI_Test( .... &index, &flag );
    if ( flag ) {
        // receive processed data and send new
    }
}

```

**Exercise 4.12.** Read section HPSC-6.5 and give pseudo-code for the distributed sparse

matrix-vector product using the above idiom for using `MPI_Test...` calls.

Discuss the advantages and disadvantages of this approach. The answer is not going to be black and white: discuss when you expect which approach to be preferable.

## 4.2.7 More about point-to-point communication

### 4.2.7.1 Message probing

MPI receive calls specify a receive buffer, and its size has to be enough for any data sent. In case you really have no idea how much data is being sent, and you don't want to overallocate the receive buffer, you can use a 'probe' call.

The calls `MPI_Probe`, `MPI_Iprobe`, accept a message, but do not copy the data. Instead, when probing tells you that there is a message, you can use `MPI_Get_count` to determine its size, allocate a large enough receive buffer, and do a regular receive to have the data copied.

### 4.2.7.2 Wildcards in the receive call

*The reference for the commands introduced here can be found in section 8.3.1.*

With some receive calls you know everything about the message in advance: its source, tag, and size. In other cases you want to leave some options open, and inspect the message for them after it was received. To do this, the receive call has a *status* parameter. This status is a property of the actually received message, so `MPI_Irecv` does not have a status parameter, but `MPI_Wait` does.

Here are some of the uses of the status:

**4.2.7.2.1 Source** In some applications it makes sense that a message can come from one of a number of processes. In this case, it is possible to specify `MPI_ANY_SOURCE` as the source. To find out where the message actually came from, you would use the `MPI_SOURCE` field of the status object that is delivered by `MPI_Recv` or the `MPI_Wait...` call after an `MPI_Irecv`.

```
MPI_Recv(recv_buffer+p, 1, MPI_INT, MPI_ANY_SOURCE, 0, comm,
          &status);
sender = status.MPI_SOURCE;
```

There are various scenarios where receiving from ‘any source’ makes sense. One is that of the *master-worker model*. The master task would first send data to the worker tasks, then issues a blocking wait for the data of whichever process finishes first.

If a processor is expecting more than one message from a single other processor, message tags are used to distinguish between them. In that case, a value of `MPI_ANY_TAG` can be used, and the actual tag of a message can be retrieved with

```
int tag = status.MPI_TAG;
```

If the amount of data received is not known a priori, the amount received can be found as

```
MPI_Get_count(&recv_status, MPI_INT, &recv_count);
```

#### 4.2.7.3 Overlap of computation and communication

Non-blocking routines have long held the promise of letting a program *overlap its computation and communication*. The idea was that after posting the non-blocking calls the program could proceed to do non-communication work, while another part of the system would take care of the communication. Unfortunately, a lot of this communication involved activity in user space, so the solution would have been to let it be handled by a separate thread. Until recently, processors were not efficient at doing such multi-threading, so true overlap stayed a promise for the future.

#### 4.2.7.4 More about non-blocking

Above we used `MPI_Irecv`, but we could have used the `MPI_Recv` routine. There is nothing special about a non-blocking or synchronous message once it arrives; the `MPI_Recv` call can match any of the send routines you have seen so far (but not `MPI_Sendrecv`).

### 4.2.8 Synchronous and asynchronous communication

It is easiest to think of blocking as a form of synchronization with the other process, but that is not quite true. Synchronization is a concept in itself, and we talk about *synchronous* communication if there is actual coordination going on with the other process, and *asynchronous* communication if there is not. Blocking then only refers to the program waiting until the user data is safe to reuse; in the synchronous case a blocking call means that the data is indeed transferred, in the asynchronous case it only means that the data has been transferred to some system buffer. The four possible cases are illustrated in figure 4.8.

MPI has a number of routines for synchronous communication, such as `MPI_Ssend`.

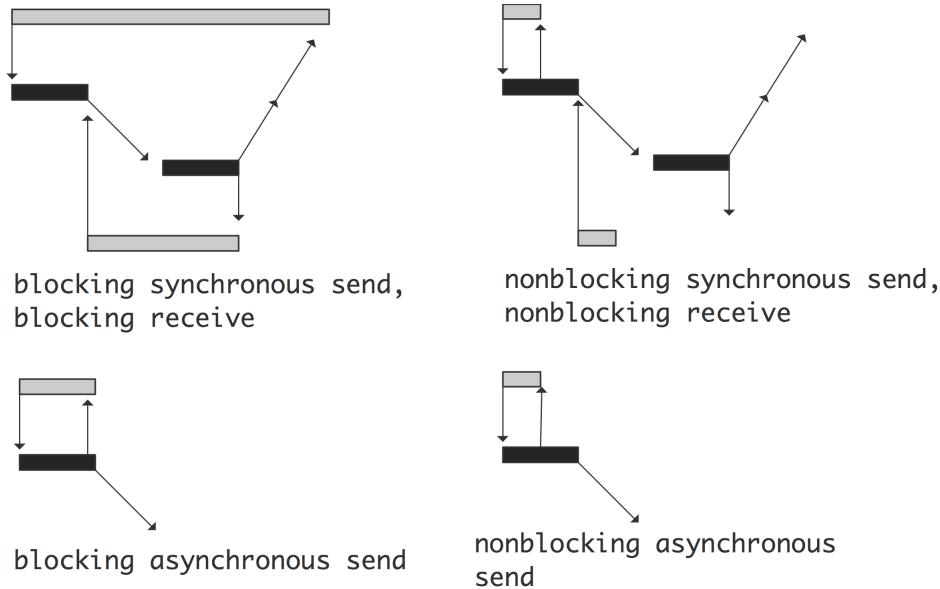


Figure 4.8: Blocking and synchronicity

#### 4.2.9 Buffered communication

*The reference for the commands introduced here can be found in section 8.4.2.*

By now you have probably got the notion that managing buffer space in MPI is important: data has to be somewhere, either in user-allocated arrays or in system buffers. Buffered sends are yet another way of managing buffer space.

1. You allocate your own buffer space, and you attach it to your process;
2. You use the `MPI_Bsend` call for sending;
3. You detach the buffer when you're done with the buffered sends.

There can be only one buffer per process; its size should be enough for all outstanding `MPI_Bsend` calls that are simultaneously outstanding, plus `MPI_BSEND_OVERHEAD`.

#### 4.2.10 Persistent communication

*The reference for the commands introduced here can be found in section 8.4.3.*

An `Irecv` or `Irecv` call as an `MPI_Request` parameter. This is an object that gets created in the send/recv call, and deleted in the wait call. You can imagine that this carries some overhead, and if the same communication is repeated many times you may want to avoid this overhead by reusing the request object.

To do this, MPI has *persistent communication*:

- You describe the communication with `MPI_Send_init`, which has the same calling sequence as `MPI_Isend`, or `MPI_Recv_init`, which has the same calling sequence as `MPI_Irecv`.

- The actual communication is performed by calling `MPI_Start`, for a single request, or `MPI_Startall` for an array or requests.
- Completion of the communication is confirmed with `MPI_Wait` or similar routines as you have seen in the explanation of non-blocking communication.
- The wait call does not release the request object: that is done with `MPI_Request_free`.

## 4.3 Shared-memory-like communication: one-sided communication

The reference for the commands introduced here can be found in section 8.5.

Above, you saw point-to-point operations of the two-sided type: they require the co-operation of a sender and receiver. This co-operation could be loose: you can post a receive with `MPI_ANY_SOURCE` as sender, but there had to be both a send and receive call. In this section, you will see one-sided communication routines where a process can do a ‘put’ or ‘get’ operation, writing data to or reading it from another processor, without that other processor’s involvement.

In one-sided MPI operations, also known as Remote Direct Memory Access (RDMA) or Remote Memory Access (RMA) operations, there are still two processes involved: the *origin*, which is the process that originates the transfer, whether this is a ‘put’ or a ‘get’, and the *target* whose memory is being accessed. Unlike with two-sided operations, the target does not perform an action that is the counterpart of the action on the origin.

That does not mean that the origin can access arbitrary data on the target at arbitrary times. First of all, one-sided communication in MPI is limited to accessing only a specifically declared memory area on the target: the target declares an area of user-space memory that is accessible to other processes. This is known as a *window*. Windows limit how origin processes can access the target’s memory: you can only ‘get’ data from a window or ‘put’ it into a window; all the other memory is not reachable from other processes.

The alternative to having windows is to use *distributed shared memory* or *virtual shared memory*: memory is distributed but acts as if it shared. The so-called Partitioned Global Address Space (PGAS) languages such as Unified Parallel C (UPC) use this model. The MPI RMA model makes it possible to lock a window which makes programming slightly more cumbersome, but the implementation more efficient.

Within one-sided communication, MPI has two modes: active RMA and passive RMA. In *active RMA*, or *active target synchronization*, the target sets boundaries on the time period (the ‘epoch’) during which its window can be accessed. The main advantage of this mode is that the origin program can perform many small transfers, which are aggregated behind the scenes. Active RMA acts much like asynchronous transfer with a concluding `Waitall`.

In *passive RMA*, or *passive target synchronization*, the target process puts no limitation on when its window can be accessed. (PGAS languages such as UPC are based on this model: data is simply read or written at will.) While intuitively it is attractive to be able to write to and read from a target at arbitrary time, there are problems. For instance, it requires a remote agent on the target, which may interfere with execution of the main thread, or conversely it may not be activated at the optimal time. Passive RMA is also very hard to debug and can lead to strange deadlocks.

### 4.3.1 Windows

*The reference for the commands introduced here can be found in section 8.5.1.*

In one-sided communication, each processor can make an area of memory available, called a *window*. This has the following characteristics:

- The window is defined on a communicator, so the create call is collective.
- The window size can be set individually on each process. A zero size is allowed, but since window creation is collective, it is not possible to skip the create call.

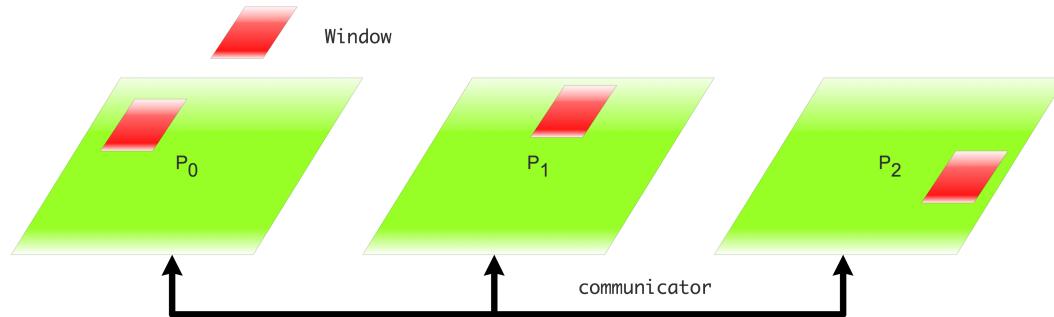


Figure 4.9: Collective definition of a window for one-sided data access

defined with respect to a communicator: each process specifies a memory area. Routine for creating and releasing windows are collective, so each process *has* to call them; see figure 4.9.

```

MPI_Info info;
MPI_Win window;
MPI_Win_create( /* memory area */, info, comm, &window );
MPI_Win_free( &window );

```

(For the `info` parameter you can often use `MPI_INFO_NULL`.) While the creation of a window is collective, each processor can specify its own window size, including zero, and even the type of the elements in it.

The MPI specification allows that the memory of a window can be separate from the regular program memory. The routine `MPI_Alloc_mem` can return a pointer to such privileged memory.

```

MPI_Info info ;
int error ;
error = MPI_Info_create ( & info ) ;
error = MPI_Info_set ( info , "no_locks" , "true" ) ;
/* Use the info object*/
error = MPI_Info_free ( info ) ;

```

### 4.3.2 Active target synchronization: epochs

*The reference for the commands introduced here can be found in section 8.5.3.*

There are two mechanisms for *active target synchronization*, that is, one-sided communications where both sides are involved to the extent that they declare the communication epoch. In this section we look at the first mechanism, which is to use a *fence* operation:

```
MPI_Win_fence (int assert, MPI_Win win)
```

This operation is collective on the communicator of the window. It is comparable to MPI\_Wait calls for non-blocking communication.

The use of fences is somewhat complicated. The interval between two fences is known as an *epoch*. You can give various hints to the system about this epoch versus the ones before and after through the *assert* parameter.

```
MPI_Win_fence( (MPI_MODE_NOPUT | MPI_MODE_NOPRECEDE), win);
MPI_Get( /* operands */, win);
MPI_Win_fence(MPI_MODE_NOSUCCEED, win);
```

In between the two fences the window is exposed, and while it is you should not access it locally. If you absolutely need to access it locally, you can use an RMA operation for that. Also, there can be only one remote process that does a put; multiple accumulate accesses are allowed.

Fences are, together with other window calls, collective operations. That means they imply some amount of synchronization between processes. Consider:

```
MPI_Win_fence( ... win ... ); // start an epoch
if (mytid==0) // do lots of work
else // do almost nothing
MPI_Win_fence( ... win ... ); // end the epoch
```

and assume that all processes execute the first fence more or less at the same time. The zero process does work before it can do the second fence call, but all other processes can call it immediately. However, they can not finish that second fence call until all one-sided communication is finished, which means they wait for the zero process.

```
// putfence.c
MPI_Win the_window;
MPI_Win_create(&window_data, 2*sizeof(int), sizeof(int),
  MPI_INFO_NULL, comm, &the_window);
MPI_Win_fence(0,the_window);
if (mytid==0) {
  MPI_Put( /* data on origin: */  &my_number, 1,MPI_INT,
    /* data on target: */   other,1,      1,MPI_INT,
    the_window);
}
MPI_Win_fence(0,the_window);
MPI_Win_free(&the_window);
```

## 4. MPI topic 3: Distributed data

---

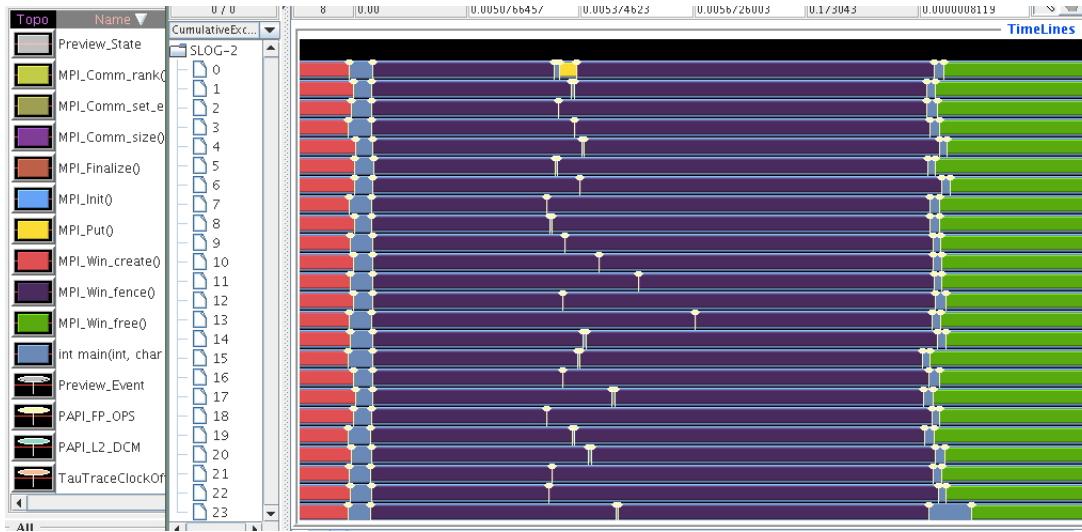


Figure 4.10: A trace of a one-sided communication epoch where process zero only originates a one-sided transfer

As a further restriction, you can not mix Get with Put or Accumulate calls in a single epoch. Hence, we can characterize an epoch as an *access epoch* on the origin, and as an *exposure epoch* on the target.

Assertions are an integer parameter: you can add or logical-or values. The value zero is always correct. For further information, see section 8.5.4.

### 4.3.3 Put, get, accumulate

*The reference for the commands introduced here can be found in section 8.5.2.*

Window areas are accessible to other processes in the communicator by specifying the process rank and an offset from the base of the window.

```
MPI_Put (
    void *origin_addr, int origin_count, MPI_Datatype origin_datatype,
    int target_rank,
    MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype,
    MPI_Win window)
```

**Exercise 4.13.** Write code where process 0 randomly writes in the window on 1 or 2.

```
// randomput_skl.c
MPI_Win_create(&window_data,sizeof(int),sizeof(int),
               MPI_INFO_NULL,comm,&the_window);

for (int c=0; c<10; c++) {
    float randomfraction = (rand() / (double)RAND_MAX);
    if (randomfraction>.5)
```

```

        other = 2;
        else other = 1;
        window_data = 0;
        your_code_goes_here.....
        my_sum += window_data;
    }

    if (mytid>0 && mytid<3)
        printf("Sum on %d: %d\n",mytid,my_sum);
    if (mytid==0) printf("(sum should be 10)\n");
}

```

The MPI\_Get call is very similar; a third one-sided routine is MPI\_Accumulate which does a reduction operation on the results that are being put:

```

MPI_Accumulate (
    void *origin_addr, int origin_count, MPI_Datatype origin_datatype,
    int target_rank,
    MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype,
    MPI_Op op, MPI_Win window)

```

**Exercise 4.14.** Implement an ‘all-gather’ operation using one-sided communication: each processor stores a single number, and you want each processor to build up an array that contains the values from all processors. Note that you do not need a special case for a processor collecting its own value: doing ‘communication’ between a processor and itself is perfectly legal.

Accumulate is a reduction with remote result. As with MPI\_Reduce, the order in which the operands are accumulated is undefined. The same predefined operators are available, but no user-defined ones. There is one extra operator: MPI\_REPLACE, this has the effect that only the last result to arrive is retained.

#### 4.3.4 Put vs Get

```

while (!converged(A)) {
    update(A);
    MPI_Win_fence(MPI_MODE_NOPRECEDE, win);
    for(i=0; i < toneighbors; i++)
        MPI_Put(&frombuf[i], 1, fromtype[i], toneighbor[i],
                todisp[i], 1, totype[i], win);
    MPI_Win_fence((MPI_MODE_NOSTORE | MPI_MODE_NOSUCCEED), win);
}

while (!converged(A)) {
    update_boundary(A);
    MPI_Win_fence((MPI_MODE_NOPUT | MPI_MODE_NOPRECEDE), win);
}

```

```

for(i=0; i < fromneighbors; i++)
    MPI_Get (&tobuf[i], 1, totype[i], fromneighbor[i],
             fromdisp[i], 1, fromtype[i], win);
update_core(A);
MPI_Win_fence (MPI_MODE_NOSUCCEED, win);
}

```

### 4.3.5 More active target synchronization

*The reference for the commands introduced here can be found in section 8.5.5.*

There is a more fine-grained ways of doing *active target synchronization*. While fences corresponded to a global synchronization of one-sided calls, the `MPI_Win_start`, `MPI_Win_complete`, `MPI_Win_post`, `Win_wait` routines are suitable, and possibly more efficient, if only a small number of processor pairs is involved. Which routines you use depends on whether the processor is an *origin* or *target*.

If the current process is going to have the data in its window accessed, you define an *exposure epoch* by:

```

MPI_Win_post( /* group of origin processes */
MPI_Win_wait()

```

This turns the current processor into a target for access operations issued by a different process.

If the current process is going to be issuing one-sided operations, you define an *access epoch* by:

```

MPI_Win_start( /* group of target processes */
// access operations
MPI_Win_complete()

```

This turns the current process into the origin of a number of one-sided access operations.

Both pairs of operations declare a *group of processors*; see section 6.2.3 for how to get such a group from a communicator. On an origin processor you would specify a group that includes the targets you will interact with, on a target processor you specify a group that includes the possible origins.

### 4.3.6 Passive target synchronization

*The reference for the commands introduced here can be found in section 8.5.6.*

In *passive target synchronization* only the origin is actively involved: the target makes no calls whatsoever. This means that the origin process remotely locks the window on the target.

During an access epoch, a process can initiate and finish a one-sided transfer.

```

if (rank == 0) {
    MPI_Win_lock (MPI_LOCK_EXCLUSIVE, 1, 0, win);
    MPI_Put (outbuf, n, MPI_INT, 1, 0, n, MPI_INT, win);
    MPI_Win_unlock (1, win);
}

```

```
}
```

The two lock types are:

- `MPI_LOCK_SHARED` which should be used for Get calls: since multiple processors are allowed to read from a window in the same epoch, the lock can be shared.
- `MPI_LOCK_EXCLUSIVE` which should be used for Put and Accumulate calls: since only one processor is allowed to write to a window during one epoch, the lock should be exclusive.

These routines make MPI behave like a shared memory system; the instructions between locking and unlocking the window effectively become *atomic operations*.

The above mechanism is of limited use. Suppose processor zero has a data structure `work_table` with items that need to be processed. A counter `first_work` keeps track of the lowest numbered item that still needs processing. You can imagine the following *master-worker* scenario:

- Each process connects to the master,
- inspects the `first_work` variable,
- retrieves the corresponding work item, and
- increments the `first_work` variable.

It is important here to avoid a *race condition* (see section HPSC-[2.6.1.5](#)) that would result from a second process reading the `first_work` variable before the first process could have updated it. Therefore, the reading and updating needs to be an *atomic operation*.

Unfortunately, you can not have a put and get call in the same access epoch. For this reason, MPI version 3 has added certain atomic operations, such as `MPI_Fetch_and_op`.

#### 4.3.7 Grouping by shared memory

MPI's one-sided routines take a very symmetric view of processes: each process can access the window of every other process (within a communicator). Of course, in practice there will be a difference in performance depending on whether the origin and target are actually on the same shared memory, or whether they can only communicate through the network. For this reason MPI makes it easy to group processes by shared memory domains using `MPI_Comm_split_type`.

#### 4.3.8 Details

Sometimes an architecture has memory that is shared between processes, or that otherwise is fast for one-sided communication. To put a window in such memory, it can be placed in memory that is especially allocated:

```
MPI_Alloc_mem() and MPI_Free_mem()
```

These calls reduce to `malloc` and `free` if there is no special memory area; SGI is an example where such memory does exist.

### 4.3.9 Implementation

You may wonder how one-sided communication is realized<sup>1</sup>. Can a processor somehow get at another processor's data? Unfortunately, no.

Active target synchronization is implemented in terms of two-sided communication. Imagine that the first fence operation does nothing, unless it concludes prior one-sided operations. The Put and Get calls do nothing involving communication, except for marking with what processors they exchange data. The concluding fence is where everything happens: first a global operation determines which targets need to issue send or receive calls, then the actual sends and receive are executed.

**Exercise 4.15.** Assume that only Get operations are performed during an epoch. Sketch how these are translated to send/receive pairs. The problem here is how the senders find out that they need to send. Show that you can solve this with an `MPI_Scatter_reduce` call.

The previous paragraph noted that a collective operation was necessary to determine the two-sided traffic. Since collective operations induce some amount of synchronization, you may want to limit this.

**Exercise 4.16.** Argue that the mechanism with window post/wait/start/complete operations still needs a collective, but that this is less burdensome.

Passive target synchronization needs another mechanism entirely. Here the target process needs to have a background task (process, thread, daemon,...) running that listens for requests to lock the window. This can potentially be expensive.

## 4.4 Remaining topics in point-to-point communication

### 4.4.1 Subtleties with processor synchronization

Blocking communication involves a complicated dialog between the two processors involved. Processor one says 'I have this much data to send; do you have space for that?', to which processor two replies 'yes, I do; go ahead and send', upon which processor one does the actual send. This back-and-forth (technically known as a *handshake*) takes a certain amount of communication overhead. For this reason, network hardware will sometimes forgo the handshake for small messages, and just send them regardless, knowing that the other process has a small buffer for such occasions.

One strange side-effect of this strategy is that a code that should deadlock according to the MPI specification does not do so. In effect, you may be shielded from your own programming mistake! Of course, if you then run a larger problem, and the small message becomes larger than the threshold, the deadlock will suddenly occur. So you find yourself in the situation that a bug only manifests itself on large problems, which are usually harder to debug. In this case, replacing every `MPI_Send` with a `MPI_Ssend` will force the handshake, even for small messages.

Conversely, you may sometimes wish to avoid the handshake on large messages. MPI has a solution for this: the `MPI_Rsend` ('ready send') routine sends its data immediately, but it needs the receiver to be ready for this. How can you guarantee that the receiving process is ready? You could for instance do the following (this uses non-blocking routines, which are explained below in section 4.2.6):

---

1. For more on this subject, see [6].

```
if ( receiving ) {  
    MPI_Irecv() // post non-blocking receive  
    MPI_Barrier() // synchronize  
else if ( sending ) {  
    MPI_Barrier() // synchronize  
    MPI_Rsend() // send data fast
```

When the barrier is reached, the receive has been posted, so it is safe to do a ready send. However, global barriers are not a good idea. Instead you would just synchronize the two processes involved.

**Exercise 4.17.** Give pseudo-code for a scheme where you synchronize the two processes through the exchange of a blocking zero-size message.

#### 4.4.2 The origin of one-sided communication in ShMem

The *Cray T3E* had a library called *shmemp* which offered a type of shared memory. Rather than having a true global address space it worked by supporting variables that were guaranteed to be identical between processors, and indeed, were guaranteed to occupy the same location in memory. Variables could be declared to be shared a ‘symmetric’ pragma or directive; their values could be retrieved or set by `shmemp_get` and `shmemp_put` calls.

## Chapter 5

### MPI topic 4: Dealing with complicated data

#### 5.1 Data types

In the examples you have seen so far, every time data was sent, it was as a contiguous buffer with elements of a single type. In practice you may want to send heterogeneous data, or non-contiguous data. Figure 5.1 indicates one source of irregular data: with a matrix on *column-major storage*, a column is stored in con-

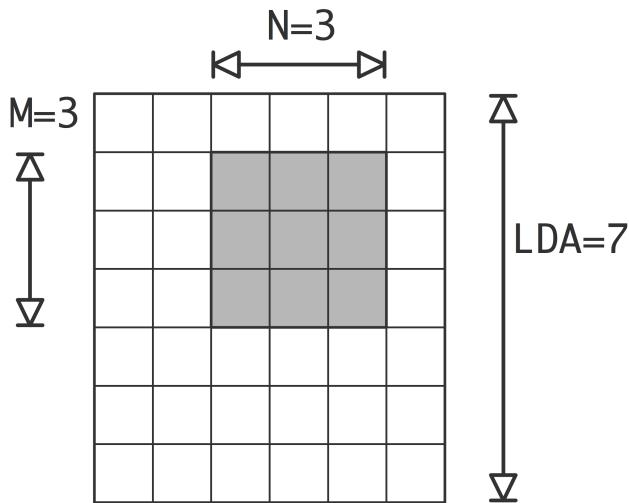


Figure 5.1: Memory layout of a row and column of a matrix in column-major storage

tiguous memory. However, a row of such a matrix is not contiguous; its elements being separated by a *stride* equal to the column length.

**Exercise 5.1.** How would you describe the memory layout of a submatrix, if the whole matrix has size  $M \times N$  and the submatrix  $m \times n$ ?

The datatypes you have dealt with so far are known as *elementary datatypes*; irregular objects are known as *derived datatypes*.

### 5.1.1 Elementary data types

The reference for the commands introduced here can be found in section [8.2.1](#).

MPI has a number of elementary data types, corresponding to the simple data types of programming languages. The names are made to resemble the types of C and Fortran, for instance `MPI_FLOAT` and `MPI_DOUBLE` versus `MPI_REAL` and `MPI_DOUBLE_PRECISION`.

MPI calls accept arrays of elements:

```
double x[20];
MPI_Send( x, 20, MPI_DOUBLE, . . . . . )
```

so for a single element you need to take its address:

```
double x;
MPI_Send( &x, 1, MPI_DOUBLE, . . . . . )
```

### 5.1.2 Derived datatypes

The reference for the commands introduced here can be found in section [8.2.2](#).

MPI allows you to create your own data types, somewhat (but not completely...) analogous to defining structures in a programming language. MPI data types are mostly of use if you want to send multiple items in one message.

There are two problems with using only elementary datatypes as you have seen so far.

- MPI communication routines can only send multiples of a single data type: it is not possible to send items of different types, even if they are contiguous in memory. It would be possible to use the `MPI_BYTE` data type, but this is not advisable.
- It is also ordinarily not possible to send items of one type if they are not contiguous in memory. You could of course send a contiguous memory area that contains the items you want to send, but that is wasteful of bandwidth.

With MPI data types you can solve these problems in several ways.

- You can create a new *contiguous data type* consisting of an array of elements of another data type. There is no essential difference between sending one element of such a type and multiple elements of the component type.
- You can create a *vector data type* consisting of regularly spaced blocks of elements of a component type. This is a first solution to the problem of sending non-contiguous data.
- For not regularly spaced data, there is the *indexed data type*, where you specify an array of index locations for blocks of elements of a component type. The blocks can each be of a different size.
- The *struct data type* can accomodate multiple data types.

And you can combine these mechanisms to get irregularly spaced heterogeneous data, et cetera.

### 5.1.2.1 Datatype signatures

With the primitive types you have seen so far, it pretty much went without saying that if the sender sends an array of doubles, the receiver had to declare the datatype also as doubles. With derived types that is no longer the case: the sender and receiver can declare a different datatype for the send and receive buffer, as long as these have the same *datatype signature*.

The signature of a datatype is the internal representation of that datatype. For instance, if the sender declares a datatype consisting of two doubles, and it sends four elements of that type, the receiver can receive it as two elements of a type consisting of four doubles.

You can also look at the signature as the form ‘under the hood’ in which MPI sends the data.

### 5.1.2.2 Basic calls

*The reference for the commands introduced here can be found in section 8.2.2.1.*

New MPI data types are created by

- MPI\_Type\_contiguous
- MPI\_Type\_create\_subarray
- MPI\_Type\_vector
- MPI\_Type\_struct
- MPI\_Type\_indexed
- MPI\_Type\_hindexed

It is necessary to call `MPI_Type_commit` which makes MPI do the indexing calculations for the data type. When you no longer need the data type, you call `MPI_Type_free`.

### 5.1.2.3 Contiguous type

*The reference for the commands introduced here can be found in section 8.2.2.2.*

The simplest derived type is the ‘contiguous’ type, constructed with `MPI_Type_contiguous`. A contiguous type describes an array of items of an elementary or earlier defined type. There is no difference between sending one item of a contiguous type and multiple items of the constituent type. This is illus-



Figure 5.2: A contiguous datatype is built up out of elements of a constituent type

trated in figure 5.2.

### 5.1.2.4 Vector type

The reference for the commands introduced here can be found in section [8.2.2.3](#).

The simplest non-contiguous datatype is the ‘vector’ type, constructed with `MPI_Type_vector`. A vector type describes a series of blocks, all of equal size, spaced with a constant stride. This is illustrated in

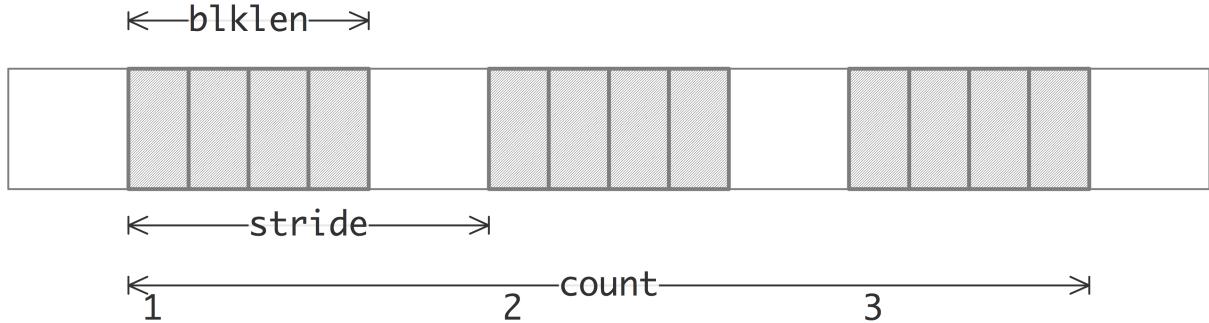


Figure 5.3: A vector datatype is built up out of strided blocks of elements of a constituent type

figure [5.3](#).

As an example of this datatype, consider the example of transposing a matrix, for instance to convert between C and Fortran arrays (see section [HPSC-34.2](#)). Suppose that a processor has a matrix stored in C, row-major, layout, and it needs to send a column to another processor. If the matrix is declared as

```
int M, N; double mat [M] [N]
```

then a column has  $M$  blocks of one element, spaced  $N$  locations apart. In other words:

```
MPI_Datatype MPI_column;
MPI_Type_vector(
    /* count= */ M, /* blocklength= */ 1, /* stride= */ N,
    MPI_DOUBLE, &MPI_column );
```

Sending the first column is easy:

```
MPI_Send( mat, 1, MPI_column, ... );
```

The second column is just a little trickier: you now need to pick out elements with the same stride, but starting at  $A[0][1]$ .

```
MPI_Send( &(mat[0][1]), 1, MPI_column, ... );
```

You can make this marginally more efficient (and harder to read) by replacing the index expression by  $mat+1$ .

**Exercise 5.2.** Suppose you have a matrix of size  $4N \times 4N$ , and you want to send the elements  $A[4*i][4*j]$  with  $i, j = 0, \dots, N - 1$ . How would you send these elements with a single transfer?

**Exercise 5.3.** Allocate a matrix on processor zero, using Fortran column-major storage.  
Using  $P$  sendrecv calls, distribute the rows of this matrix among the processors.

#### 5.1.2.5 Subarray type

The vector datatype can be used for blocks in an array of dimension more than 2 by using it recursively. However, this gets tedious. Instead, there is an explicit subarray type

```
Semantics:
MPI_TYPE_CREATE_SUBARRAY(
    ndims, array_of_sizes, array_of_subsizes,
    array_of_starts, order, oldtype, newtype)
IN ndims: number of array dimensions (positive integer)
IN array_of_sizes: number of elements of type oldtype in each dimension
    of the full array (array of positive integers)
IN array_of_subsizes: number of elements of type oldtype in each
    dimension of the subarray (array of positive integers)
IN array_of_starts: starting coordinates of the subarray in each
    dimension (array of non-negative integers)
IN order: array storage order flag (state)
IN oldtype: array element datatype (handle)
OUT newtype: new datatype (handle)

C:
int MPI_Type_create_subarray(
    int ndims, const int array_of_sizes[],
    const int array_of_subsizes[], const int array_of_starts[],
    int order, MPI_Datatype oldtype, MPI_Datatype *newtype)

Fortran:
MPI_Type_create_subarray(ndims, array_of_sizes, array_of_subsizes,
    array_of_starts, order, oldtype, newtype, ierror)
INTEGER, INTENT(IN) :: ndims, array_of_sizes(ndims),
    array_of_subsizes(ndims), array_of_starts(ndims), order
TYPE(MPI_Datatype), INTENT(IN) :: oldtype
TYPE(MPI_Datatype), INTENT(OUT) :: newtype
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python:
MPI.Datatype.Create_subarray(self, sizes, subsizes, starts, int order=ORDER_C)
```

This describes the dimensionality and extent of the array, and the starting point (the ‘upper left corner’) and extent of the subarray.

#### 5.1.2.6 Indexed type

*The reference for the commands introduced here can be found in section 8.2.2.4.*

The indexed datatype, constructed with `MPI_Type_indexed` can send arbitrarily located elements from an array of a single datatype. You need to supply an array of index locations, plus an array of blocklengths with a separate blocklength for each index. The total number of elements sent is the sum of the blocklengths.

### 5.1.2.7 Struct type

The reference for the commands introduced here can be found in section [8.2.2.5](#).

The structure type, created with `MPI_Type_create_struct`, can contain multiple data types. The

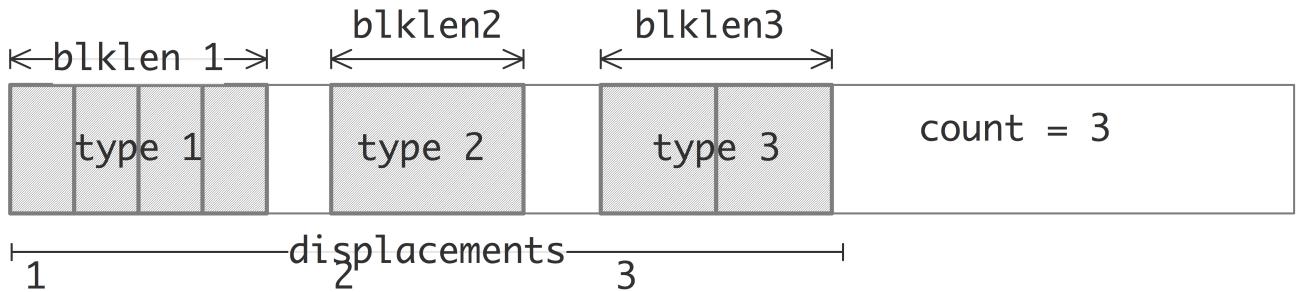


Figure 5.4: The elements of an MPI Struct datatype

specification contains a ‘count’ parameter that specifies how many blocks there are in a single structure. For instance,

```
struct {
    int i;
    float x,y;
} point;
```

has two blocks, one of a single integer, and one of two floats. This is illustrated in figure 5.4.

The structure type is very similar in functionality to `MPI_Type_hindexed`, which uses byte-based indexing. The structure-based type is probably cleaner in use.

### 5.1.3 Big data types

The `size` parameter in MPI send and receive calls is of type integer, meaning that it’s maximally  $2^{31} - 1$ . These day computers are big enough that this is a limitation. Derived types offer some way out: to send  $10^{40}$  elements you would

- create a contiguous type with  $10^{20}$  elements, and
- send  $10^{20}$  elements of that type.

This often works, but it’s not perfect. For instance, the routine `size` returns the total number of basic elements sent (as opposed to `MPI_Get_count` which would return the number of elements of the derived type). Since its output argument is of integer type, it can’t store the right value.

The *MPI 3* standard has addressed this as follows.

- To preserve backwards compatibility, the `size` parameter keeps being of type integer.
- The trick with sending elements of a derived type still works, but
- There are new routines that can return the correct information about the total amount of data; for instance, `MPI_Get_elements_x` returns its result as a `MPI_Count`.

#### 5.1.4 Packing

The reference for the commands introduced here can be found in section [8.2.3](#).

One of the reasons for derived datatypes is dealing with non-contiguous data. In older communication libraries this could only be done by *packing* data from its original containers into a buffer, and likewise unpacking it at the receiver into its destination data structures.

MPI offers this packing facility, partly for compatibility with such libraries, but also for reasons of flexibility. Unlike with derived datatypes, which transfers data atomically, packing routines add data sequentially to the buffer and unpacking takes them sequentially.

This means that one could pack an integer describing how many floating point numbers are in the rest of the packed message. Correspondingly, the unpack routine could then investigate the first integer and based on it unpack the right number of floating point numbers.

MPI offers the following:

- The `MPI_Pack` command adds data to a send buffer;
- the `MPI_Unpack` command retrieves data from a receive buffer;
- the buffer is sent with a datatype of `MPI_PACKED`.

# **Chapter 6**

## **MPI topic 5: Sub computations**

### **6.1 Subcommunications**

In many scenarios you divide a large job over all the available processors. However, your job has two or more parts that can be considered as jobs by themselves. In that case it makes sense to divide your processors into subgroups accordingly.

Supose for instance that you are running a simulation where inputs are generated, a computation is performed on them, and the results of this computation are analyzed or rendered graphically. You could then consider dividing your processors in three groups corresponding to generation, computation, rendering.

As long as you only do sends and receives, this division works fine. However, if one group of processes needs to perform a collective operation, you don't want the other groups involved in this. Thus, you really want the three groups to be really distinct from each other.

In order to make such subsets of processes, MPI has the mechanism of taking a subset of `MPI_COMM_WORLD` and turning that subset into a new communicator.

Now you understand why the MPI collective calls had an argument for the communicator: a collective involves all proceses *of that communicator*. By making a communicator that contains a subset of all available processes, you can do a collective on that subset.

#### **6.1.1 Scenario: climate model**

A climate simulation code has several components, for instance corresponding to land, air, ocean, and ice. You can imagine that each needs a different set of equations and algorithms to simulate. You can then divide your processes, where each subset simulates one component of the climate, occasionally communicating with the other components.

#### **6.1.2 Scenario: quicksort**

The popular quicksort algorithm works by splitting the data into two subsets that each can be sorted individually. If you want to sort in parallel, you could implement this by making two subcommunicators, and sorting the data on these, creating recursively more subcommunicators.

## 6.2 Communicators

*The reference for the commands introduced here can be found in section ??.*

A communicator is an object describing a group of processes. In many applications all processes work together closely coupled, and the only communicator you need is `MPI_COMM_WORLD`. However, there are circumstances where you want one subset of processes to operate independently of another subset. For example:

- If processors are organized in a  $2 \times 2$  grid, you may want to do broadcasts inside a row or column.
- For an application that includes a producer and a consumer part, it makes sense to split the processors accordingly.

In this section we will see mechanisms for defining new communicators and sending messages between communicators.

An important reason for using communicators is the development of software libraries. If the routines in a library use their own communicator (even if it is a duplicate of the ‘outside’ communicator), there will never be a confusion between message tags inside and outside the library.

### 6.2.1 Basics

There are three predefined communicators:

- `MPI_COMM_WORLD` comprises all processes that were started together by `mpirun` (or some related program).
- `MPI_COMM_SELF` is the communicator that contains only the current process.
- `MPI_COMM_NULL` is the invalid communicator. Routines that construct communicators can give this as result if an error occurs.

In some applications you will find yourself regularly creating new communicators, using the mechanisms described below. In that case, you should de-allocate communicators with `MPI_Comm_free` when you’re done with them.

### 6.2.2 Creating new communicators

There are various ways of making new communicators. We discuss three mechanisms, from simple to complicated.

#### 6.2.2.1 Duplicating communicators

*The reference for the commands introduced here can be found in section 8.7.1.*

With `MPI_Comm_dup` you can make an exact duplicate of a communicator. This may seem pointless, but it is actually very useful for the design of software libraries. Imagine that you have a code

```
MPI_Isend(...); MPI_Irecv(...);
// library call
MPI_Waitall(...);
```

and suppose that the library has receive calls. Now it is possible that the receive in the library inadvertently catches the message that was sent in the outer environment.

To prevent this confusion, the library should duplicate the outer communicator, and send all messages with respect to its duplicate. Now messages from the user code can never reach the library software, since they are on different communicators.

### 6.2.2.2 Splitting a communicator

The reference for the commands introduced here can be found in section [8.7.2](#).

Splitting a communicator into multiple disjoint communicators can be done with `MPI_Comm_split`. This uses a ‘colour’:

```
MPI_Comm_split( old_comm, colour, new_comm, ... );
```

and all processes in the old communicator with the same colour wind up in a new communicator together. The old communicator still exists, so processes now have two different contexts in which to communicate.

Here is one example of communicator splitting. Suppose your processors are in a two-dimensional grid:

```
MPI_Comm_rank( MPI_COMM_WORLD, &mytid );
proc_i = mytid % proc_column_length;
proc_j = mytid / proc_column_length;
```

You can now create a communicator per column:

```
MPI_Comm column_comm;
MPI_Comm_split( MPI_COMM_WORLD, proc_j, mytid, &column_comm );
```

and do a broadcast in that column:

```
MPI_Bcast( data, /* tag: */ 0, column_comm );
```

Because of the SPMD nature of the program, you are now doing in parallel a broadcast in every processor column. Such operations often appear in *dense linear algebra*.

**Exercise 6.1.** Organize your processors in a grid, and make subcommunicators for the rows and columns. Do a broadcast from the first row and column through the columns and rows respectively.

If you let the broadcast value be the column/row number, then processor  $(i, j)$  winds up with the numbers  $i$  and  $j$ . Test this.

As an example of communicator splitting, consider the recursive algorithm for *matrix transposition*. Processors are organized in a square grid. The matrix is divided on  $2 \times 2$  block form:

- Swap blocks  $(1, 2)$  and  $(2, 1)$ ; then
- Divide the processors into four subcommunicators, and apply this algorithm recursively on each;
- If the communicator has only one process, transpose the matrix in place.

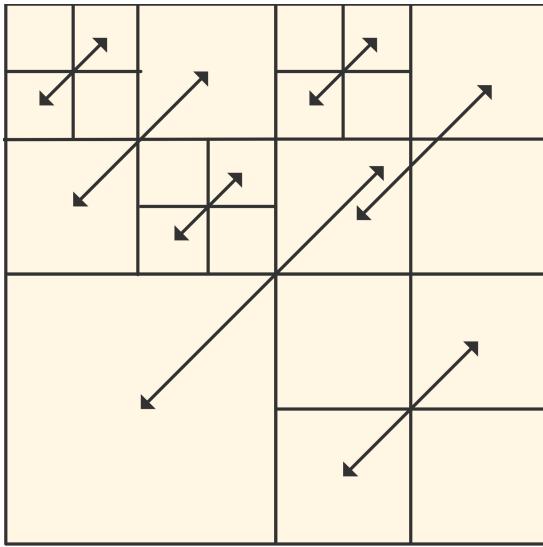


Figure 6.1: Recursive algorithm for matrix transposition

See figure 6.1.

There is an important application of communicator splitting in the context of one-sided communication, grouping processes by whether they access the same shared memory area; see section 4.3.7.

#### 6.2.2.3 Process groups

The most general mechanism is based on groups: you can extract the group from a communicator, combine different groups, and form a new communicator from the resulting group.

The group mechanism is more involved. You get the group from a communicator, or conversely make a communicator from a group with `MPI_Comm_group` and `MPI_Comm_create`:

```
MPI_Comm_group( comm, &group );
MPI_Comm_create( old_comm, group, &new_comm );
```

and groups are manipulated with `MPI_Group_incl`, `MPI_Group_excl`, `MPI_Group_difference` and a few more.

You can name your communicators with `MPI_Comm_set_name`, which could improve the quality of error messages when they arise.

#### 6.2.3 Intra-communicators

We start by exploring the mechanisms for creating a communicator that encompasses a subset of `MPI_COMM_WORLD`.

The most general mechanism for creating communicators is through process groups: you can query the group of processes of a communicator, manipulate groups, and make a new communicator out of a group you have formed.

```
MPI_COMM_GROUP (comm, group, ierr)
MPI_COMM_CREATE (MPI_Comm comm, MPI_Group group, MPI_Comm newcomm, ierr)

MPI_GROUP_UNION(group1, group2, newgroup, ierr)
MPI_GROUP_INTERSECTION(group1, group2, newgroup, ierr)
MPI_GROUP_DIFFERENCE(group1, group2, newgroup, ierr)

MPI_GROUP_INCL(group, n, ranks, newgroup, ierr)
MPI_GROUP_EXCL(group, n, ranks, newgroup, ierr)

MPI_GROUP_SIZE(group, size, ierr)
MPI_GROUP_RANK(group, rank, ierr)
```

#### 6.2.4 Inter-communicators

If two disjoint communicators exist, it may be necessary to communicate between them. This can of course be done by creating a new communicator that overlaps them, but this would be complicated: since the ‘inter’ communication happens in the overlap communicator, you have to translate its ordering into those of the two worker communicators. It would be easier to express messages directly in terms of those communicators, and this can be done with ‘inter-communicators’.

```
MPI_Intercomm_create (local_comm, local_leader, bridge_comm, remote_leader,
```

After this, the intercommunicator can be used in collectives such as

```
MPI_Bcast (buff, count, dtype, root, comm, ierr)
```

- In group A, the root process passes `MPI_ROOT` as ‘root’ value; all others use `MPI_NULL_PROC`.
- In group B, all processes use a ‘root’ value that is the rank of the root process in the root group.

Gather and scatter behave similarly; the allgather is different: all send buffers of group A are concatenated in rank order, and places on all processes of group B.

Inter-communicators can be used if two groups of process work asynchronously with respect to each other; another application is fault tolerance (section 7.3.4).

#### 6.2.5 Process topologies

*The reference for the commands introduced here can be found in section 8.7.3.*

In the communicators you have seen so far, processes are linearly ordered. In some circumstances the problem you are coding has some structure, and expressing the program in terms of that structure would be convenient. For this purpose, MPI can define a virtual *topology*. There are two types:

- regular, Cartesian, grids; and
- general graphs.

### 6.2.5.1 Cartesian grid topology

The reference for the commands introduced here can be found in section [8.7.3.1](#).

A *Cartesian grid* is a structure, typically in 2 or 3 dimensions, of points that have two neighbours in each of the dimensions. Thus, if a Cartesian grid has sizes  $K \times M \times N$ , its points have coordinates  $(k, m, n)$  with  $0 \leq k < K$  et cetera. Most points have six neighbours  $(k \pm 1, m, n)$ ,  $(k, m \pm 1, n)$ ,  $(k, m, n \pm 1)$ ; the exception are the edge points. A grid where edge processors are connected through *wraparound connections* is called a *periodic grid*.

The most common use of Cartesian coordinates is to find the rank of process by referring to it in grid terms. For instance, one could ask ‘what are my neighbours offset by  $(1, 0, 0)$ ,  $(-1, 0, 0)$ ,  $(0, 1, 0)$  et cetera’.

# Chapter 7

## MPI topics

### 7.1 Synchronization

MPI programs conform to the SPMD model, and this means that events in one process can be unrelated in time to events in another process. Any *synchronization* that happens is induced by communication and other MPI mechanisms. By synchronization here we mean any sort of temporal ordering of events in different processes.

You have already seen some mechanisms.

1. In blocking communication, the receive call does not return until the send call has completed.
2. In non-blocking communication, the wait on a receive request will not return until the send has been completed.
3. In one-sided communication, the fence mechanism impose a certain ordering on events.

Another synchronization mechanism is induced by the *barrier* mechanism. However, while an `MPI_Barrier` call guarantees that all processes have reached a certain location in their source, this does not necessarily imply anything about message traffic. Consider this example

Proc 0	Proc 1	Proc 2
Isend to 1	Irecv from any source	
Barrier	Barrier	Barrier
Wait for send request	wait for recv request (another wildcard recv)	Isend to 1 wait for send request

The unexpected behaviour here is that the (first) receive on process 1 can be matched with the send on process 2: the barrier on process 1 only guarantees that the receive instruction was performed, not the actual transfer. For that you need the `MPI_Wait` call, which is after the barrier.

### 7.2 Hybrid programming: MPI and threads

*The reference for the commands introduced here can be found in section 8.11.*

It is not automatic that a program or a library is *thread-safe*. A user can request a certain level of multi-threading with `MPI_Init_thread`, and the system will respond what the highest supported level is.

MPI can be thread-safe on the following levels:

- An MPI implementation can forbid any multi-threading;
- it can allow one thread to make MPI calls;
- it can allow one thread *at a time* to make MPI calls;
- it can allow arbitrary multi-threaded behaviour in MPI calls.

Some points.

- MPI can not distinguish between threads: the communicator rank identifies a process, and is therefore identical for all threads.
- A message sent to a process can be received by any thread that has issued a receive call with the right source/tag specification.
- Multi-threaded calls to an MPI routine have the semantics of an unspecified sequence of calls.
- A blocking MPI call only blocks the thread that makes it.

### 7.3 Leftover topics

#### 7.3.1 Getting message information

In some circumstances the recipient may not know all details of a message.

- If you are expecting multiple incoming messages, it may be most efficient to deal with them in the order in which they arrive. For that, you have to be able to ask ‘who did this message come from, and what is in it’.
- Maybe you know the sender of a message, but the amount of data is unknown. In that case you can overallocate your receive buffer, and after the message is received ask how big it was, or you can ‘probe’ an incoming message and allocate enough data when you find out how much data is being sent.

##### 7.3.1.1 Status object

The receive calls you saw above has a status argument. If you precisely know what is going to be sent, this argument tells you nothing new. Therefore, there is a special value `MPI_STATUS_IGNORE` that you can supply instead of a status object, which tells MPI that the status does not have to be reported. For routines such as `MPI_Waitany` where an array of statuses is needed, you can supply `MPI_STATUSES_IGNORE`.

However, if you expect data from multiple senders, or the amount of data is indeterminate, the status will give you that information.

The `MPI_Status` object is a structure with the following freely accessible members: `MPI_SOURCE`, `MPI_TAG`, and `MPI_ERROR`. There is also opaque information: the amount of data received can be retrieved by a function call to `MPI_Get_count`.

```
int MPI_Get_count(
    MPI_Status *status,
    MPI_Datatype datatype,
    int *count
);
```

This may be necessary since the `count` argument to `MPI_Recv` is the buffer size, not an indication of the actually expected number of data items.

### 7.3.2 Error handling

*The reference for the commands introduced here can be found in section 8.9.*

Errors in normal programs can be tricky to deal with; errors in parallel programs can be even harder. This is because in addition to everything that can go wrong with a single executable (floating point errors, memory violation) you now get errors that come from faulty interaction between multiple executables.

A few examples of what can go wrong:

- MPI errors: an MPI routine can abort for various reasons, such as receiving much more data than its buffer can accomodate. Such errors, as well as the more common type mentioned above, typically cause your whole execution to abort. That is, if one incarnation of your executable aborts, the MPI runtime will kill all others.
- Deadlocks and other hanging executions: there are various scenarios where your processes individually do not abort, but are all waiting for each other. This can happen if two processes are both waiting for a message from each other, and this can be helped by using non-blocking calls. In another scenario, through an error in program logic, one process will be waiting for more messages (including non-blocking ones) than are sent to it.

The MPI library has a general mechanism for dealing with errors that it detects. The default behaviour, where the full run is aborted, is equivalent to your code having the following call<sup>1</sup>:

```
MPI_Comm_set_errhandler(MPI_COMM_WORLD, MPI_ERRORS_ARE_FATAL);
```

Another simple possibility is to specify

```
MPI_Comm_set_errhandler(MPI_COMM_WORLD, MPI_ERRORS_RETURN);
```

which gives you the opportunity to write code that handles the error return value.

In most cases where an MPI error occurs a complete abort is the sensible thing, since there are few ways to recover. The second possibility can for instance be used to print out debugging information:

```
ierr = MPI_Something();
if (ierr!=0) {
    // print out information about what your programming is doing
    MPI_Abort();
```

---

1. The routine `MPI_Errhandler_set` is deprecated.

```
}
```

For instance,

```
Fatal error in MPI_Waitall:  
See the MPI_ERROR field in MPI_Status for the error code
```

You could code this as

```
MPI_Comm_set_errhandler(MPI_COMM_WORLD, MPI_ERRORS_RETURN);  
ierr = MPI_Waitall(2*ntids-2, requests, status);  
if (ierr!=0) {  
    char errtxt[200];  
    for (int i=0; i<2*ntids-2; i++) {  
        int err = status[i].MPI_ERROR; int len=200;  
        MPI_Error_string(err,errtxt,&len);  
        printf("Waitall error: %d %s\n",err,errtxt);  
    }  
    MPI_Abort(MPI_COMM_WORLD, 0);  
}
```

One cases where errors can be handled is that of *MPI file I/O*: if an output file has the wrong permissions, code can possibly progress without writing data, or writing to a temporary file.

### 7.3.3 Fortran issues

The reference for the commands introduced here can be found in section 8.8.2.

MPI is typically written in C, what if you program Fortran?

Assumed shape arrays can be a problem: they need to be copied. That's a problem with Isend.

- Fortran routines have the same prototype as C routines except for the addition of an integer error parameter.
- The call for MPI\_Init in Fortran does not have the commandline arguments; they need to be handled separately.
- The routine MPI\_Sizeof is only available in Fortran, it provides the functionality of the C/C++ operator sizeof.

### 7.3.4 Fault tolerance

Processors are not completely reliable, so it may happen that one ‘breaks’: for software or hardware reasons it becomes unresponsive. For an MPI program this means that it becomes impossible to send data to it, and any collective operation involving it will hang. Can we deal with this case? Yes, but it involves some programming.

First of all, one of the possible MPI error return codes (section ??) is MPI\_ERR\_COMM, which can be returned if a processor in the communicator is unavailable. You may want to catch this error, and add a ‘replacement processor’ to the program. For this, the MPI\_Comm\_spawn can be used:

```
int MPI_Comm_spawn(char *command, char *argv[], int maxprocs, MPI_Info info
                    int root, MPI_Comm comm, MPI_Comm *intercomm,
                    int array_of_errcodes[])
```

But this requires a change of program design: the communicator containing the new process(es) is not part of the old `MPI_COMM_WORLD`, so it is better to set up your code as a collection of inter-communicators to begin with.

### 7.3.5 Context information

*The reference for the commands introduced here can be found in section 8.10.1.*

The `MPI` version is available through two parameters `MPI_VERSION` and `MPI_SUBVERSION` or the function `MPI_Get_version`.

### 7.3.6 Timing

*The reference for the commands introduced here can be found in section 8.10.2.*

Timing of parallel programs is tricky. On each node you can use a timer, typically based on some Operating System (OS) call. MPI supplies its own routine `MPI_Wtime` which gives *wall clock time*. Normally you don't worry about the starting point for this timer: you call it before and after an event and subtract the values.

```
t = MPI_Wtime();
// something happens here
t = MPI_Wtime()-t;
```

If you execute this on a single processor you get fairly reliable timings, except that you would need to subtract the overhead for the timer. This is the usual way to measure timer overhead:

```
t = MPI_Wtime();
// absolutely nothing here
t = MPI_Wtime()-t;
```

#### 7.3.6.1 Global timing

However, if you try to time a parallel application you will most likely get different times for each process, so you would have to take the average or maximum. Another solution is to synchronize the processors by using a *barrier*:

```
MPI_Barrier(comm)
t = MPI_Wtime();
// something happens here
MPI_Barrier(comm)
t = MPI_Wtime()-t;
```

**Exercise 7.1.** This scheme also has some overhead associated with it. How would you measure that?

### 7.3.6.2 Local timing

Now suppose you want to measure the time for a single send. It is not possible to start a clock on the sender and do the second measurement on the receiver, because the two clocks need not be synchronized. Usually a *ping-pong* is done:

```
if ( proc_source ) {
    MPI_Send( /* to target */ );
    MPI_Recv( /* from target */ );
} else if ( proc_target ) {
    MPI_Recv( /* from source */ );
    MPI_Send( /* to source */ );
}
```

**Exercise 7.2.** Why is it generally not a good idea to use processes 0 and 1 for the source and target processor? Can you come up with a better guess?

No matter what sort of timing you are doing, it is good to know the accuracy of your timer. The routine `MPI_Wtick` gives the smallest possible timer increment. If you find that your timing result is too close to this ‘tick’, you need to find a better timer (for CPU measurements there are cycle-accurate timers), or you need to increase your running time, for instance by increasing the amount of data.

### 7.3.7 Profiling

*The reference for the commands introduced here can be found in section ??.*

MPI allows you to write your own profiling interface. To make this possible, every routine `MPI_Something` calls a routine `PMPPI_Something` that does the actual work. You can now write your `MPI_...` routine which calls `PMPPI_...`, and inserting your own profiling calls. As you can see in figure 7.1, normally only the `PMPPI` routines show up in the stack trace.

Does the standard mandate this?

### 7.3.8 Debugging

There are various ways of debugging an MPI program. Typically there are two cases. In the simple case your program can have a serious error in logic which shows up even with small problems and a small number of processors. In the more difficult case your program can only be run on large scale, or the problem only shows up when you run at large scale. For the second case you, unfortunately, need a dedicated debugging tool, and of course the good ones are expensive. In the first case there are some simpler solutions.

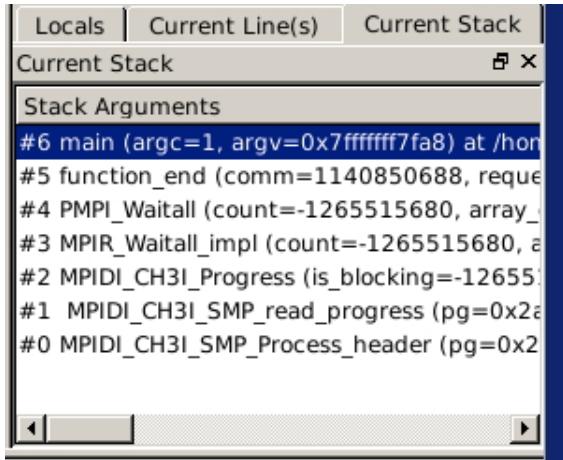


Figure 7.1: A stack trace, showing the MPI calls.

#### 7.3.8.1 Small scale debugging

If your program hangs or crashes even with small numbers of processors, you can try debugging on your local desktop or laptop computer:

```
mpirun -np <n> xterm -e gdb yourprogram
```

This starts up a number of X terminals, each of which runs your program. The magic of `mpirun` makes sure that they all collaborate on a parallel execution of that program. If your program needs commandline arguments, you have to type those in every xterm:

```
run <argument list>
```

See appendix 16.1 for more about debugging with `gdb`.

This approach is not guaranteed to work, since it depends on your ssh setup; see the discussion in <http://www.open-mpi.org/faq/?category=debugging#serial-debuggers>.

#### 7.3.8.2 Large scale debugging

Check out `ddt` or `TotalView`.

#### 7.3.8.3 Memory debugging of MPI programs

The commercial parallel debugging tools typically have a memory debugger. For an open source solution you can use `valgrind`, but that requires some setup during installation. See <http://valgrind.org/docs/manual/mc-manual.html#mc-manual.mpiwrap> for details.

### 7.3.9 Determinism

MPI processes are only synchronized to a certain extent, so you may wonder what guarantees there are that running a code twice will give the same result. You need to consider two cases: first of all, if the two runs are on different numbers of processors there are already numerical problems; see HPSC-3.3.7.

Let us then limit ourselves to two runs on the same set of processors. In that case, MPI is deterministic as long as you do not use wildcards such as `MPI_ANY_SOURCE`. Formally, MPI messages are ‘non-overtaking’: two messages between the same sender-receiver pair will arrive in sequence. Actually, they may not arrive in sequence: they are *matched* in sequence in the user program. If the second message is much smaller than the first, it may actually arrive earlier in the lower transport layer.

### 7.3.10 Progress

Non-blocking communication implies that messages make *progress* while computation is going on. However, communication of this sort can typically not be off-loaded to the network card, so it has to be done by a process. This requires a separate thread of execution, with obvious performance problems. Therefore, in practice overlap may not actually happen, and for the message to make progress it is necessary for the MPI library to become active occasionally. For instance, people have inserted dummy `MPI_Probe` calls.

A similar problem arises with passive target synchronization: it is possible that the origin process may hang until the target process makes an MPI call.

## 7.4 Literature

Online resources:

- MPI 1 Complete reference:  
<http://www.netlib.org/utk/papers/mpi-book/mpi-book.html>
- Official MPI documents:  
<http://www.mpi-forum.org/docs/>
- List of all MPI routines:  
<http://www.mcs.anl.gov/research/projects/mpi/www/www3/>

Tutorial books on MPI:

- Using MPI [2] by some of the original authors.

# Chapter 8

## MPI Reference

This section gives reference information and illustrative examples of the use of MPI. While the code snippets given here should be enough, full programs can be found in the repository for this book <https://bitbucket.org/VictorEijkhout/parallel-computing-book>.

### 8.1 Basics

#### 8.1.1 MPI setup

*This reference section gives the syntax for routines introduced in section 2.1.1.*

If you use MPI commands in a program file, be sure to include the proper header file, *mpi.h* or *mpif.h*.

```
#include "mpi.h" // for C  
#include "mpif.h" ! for Fortran
```

For *Fortran90*, many MPI installations also have an MPI module, so you can write

```
use mpi
```

The internals of these files can be different between MPI installations, so you can not compile one file against one *mpi.h* file and another file, even with the same compiler on the same machine, against a different MPI.

#### 8.1.2 Initialization / finalization

Every MPI program has to start with *MPI initialization*:

```
C:  
int MPI_Init(int *argc, char ***argv)  
  
Fortran:  
MPI_Init(ierror)  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

## 8. MPI Reference

---

where `argc` and `argv` are the arguments of a C language main program:

```
int main(int argc, char **argv) {
    ....
    return 0;
}
```

(It is allowed to pass `NULL` for these arguments.)

The commandline arguments `argc` and `argv` are only guaranteed to be passed to process zero, so the best way to pass commandline information is by a broadcast (section 3.3).

Note that the `MPI_Init` call is one of the few that differs between C and Fortran: the C routine takes the commandline arguments, which Fortran lacks.

If MPI is used in a library, MPI can have already been initialized in a main program. For this reason, one can test where `MPI_Init` has been called with

```
C:
int MPI_Initialized(int *flag)

Fortran:
MPI_Initialized(flag, ierror)
LOGICAL, INTENT(OUT) :: flag
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

The regular way to conclude an MPI program is:

```
C:
int MPI_Finalize(void)

Fortran:
MPI_Finalize(ierror)
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

but an abnormal end to a run can be forced by

```
MPI_Abort(comm,value);
```

This aborts execution on all processes associated with the communicator, but many implementations simply abort all processes. The `value` parameter is returned to the environment.

The corresponding Fortran calls are

```
call MPI_Init(ierr)
// your code
call MPI_Finalize(ierr)
```

You can test whether `MPI_Finalize` has been called with

```
C:  
int MPI_Finalized( int *flag )  
  
Fortran:  
MPI_Finalized(flag, ierror)  
LOGICAL, INTENT(OUT) :: flag  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

### 8.1.2.1 Commandline arguments

The `MPI_Init` routines takes a reference to `argc` and `argv` for the following reason: the `MPI_Init` calls filters out the arguments to `mpirun` or `mpiexec`, thereby lowering the value of `argc` and eliminating some of the `argv` arguments.

On the other hand, the commandline arguments that are meant for `mpiexec` wind up in the `MPI_INFO_ENV` object as a set of key/value pairs.

### 8.1.3 Startup

`MPI_COMM_SPAWN_MULTIPLE`

Once MPI has been initialized, the `MPI_INFO_ENV` object contains:

- command Name of program executed.
- argv Space separated arguments to command.
- maxprocs Maximum number of MPI processes to start.
- soft Allowed values for number of processors.
- host Hostname.
- arch Architecture name.
- wdir Working directory of the MPI process.
- file Value is the name of a file in which additional information is specified.
- thread\_level Requested level of thread support, if requested before the program started execution.

Note that these are the requested values; the running program can for instance have lower thread support.

### 8.1.4 Basic communicator handling

*This reference section gives the syntax for routines introduced in section 2.2.*

There are many calls relating to communicators. The simplest are `MPI_Comm_rank`

```
int MPI_Comm_rank( MPI_Comm comm, int *rank )
```

and `MPI_Comm_size`

```
int MPI_Comm_size( MPI_Comm comm, int *size )
```

Using these calls, the simplest MPI programs does this:

```
// helloworld.c
MPI_Init(&argc,&argv);
MPI_Comm_size(MPI_COMM_WORLD,&ntids);
MPI_Comm_rank(MPI_COMM_WORLD,&mytid);
printf("Hello, this is processor %d out of %d\n",mytid,ntids);
MPI_Finalize();
```

### 8.1.5 Send and receive buffers

The data is specified as a number of elements in a buffer. The same MPI routine can be used with data of different types, so the standard indicates such buffers as *choice*. The specification of this differs per language:

- In C it is an address, so the clean way is to pass it as `(void*) &myvar`.
- Fortran compilers may complain about type mismatches. This can not be helped.

## 8.2 Data types

*This reference section gives the syntax for routines introduced in section ??.*

In this section we discuss the various forms that an `MPI_Datatype` can take: elementary datatypes and derived datatypes. We also discuss packed data, which is not an MPI datatype as such.

### 8.2.1 Elementary types

*This reference section gives the syntax for routines introduced in section 5.1.1.*

C/C++:

<code>MPI_CHAR</code>	only for text data, do not use for small integers
<code>MPI_UNSIGNED_CHAR</code>	
<code>MPI_SIGNED_CHAR</code>	
<code>MPI_SHORT</code>	
<code>MPI_UNSIGNED_SHORT</code>	
<code>MPI_INT</code>	
<code>MPI_UNSIGNED</code>	
<code>MPI_LONG</code>	
<code>MPI_UNSIGNED_LONG</code>	
<code>MPI_FLOAT</code>	
<code>MPI_DOUBLE</code>	
<code>MPI_LONG_DOUBLE</code>	

There is some, but not complete, support for C99 types.

Fortran:

---

MPI_CHARACTER	Character(Len=1)
MPI_LOGICAL	
MPI_INTEGER	
MPI_REAL	
MPI_DOUBLE_PRECISION	
MPI_COMPLEX	
MPI_DOUBLE_COMPLEX	Complex(Kind=Kind(0.d0))

Addresses have type MPI\_Aint or INTEGER (KIND=MPI\_ADDRESS\_KIND) in Fortran. The start of the address range is given in MPI\_BOTTOM.

## 8.2.2 Derived datatypes

*This reference section gives the syntax for routines introduced in section 5.1.2.*

The space taken by a derived type is not immediately obvious from its definition since padding maybe applied. The actual size can be retrieved with MPI\_Type\_extent:

```
int MPI_Type_extent (MPI_Datatype datatype, MPI_Aint *extent)
```

See the example in section 8.2.2.5

### 8.2.2.1 Type create and release calls

*This reference section gives the syntax for routines introduced in section 5.1.2.2.*

A derived type needs to be committed before it can be used:

```
C:
int MPI_Type_commit (MPI_Datatype *datatype)

Fortran:
MPI_Type_commit(datatype, ierror)
TYPE(MPI_Datatype), INTENT(INOUT) :: datatype
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

The commit call is typically used to find an efficient ‘flat’ representation of recursively defined datatypes.

When you no longer need the derived type, its space can be released with MPI\_Type\_free:

```
int MPI_Type_free (MPI_Datatype *datatype)
```

After the type free call

- The definition of the datatype identifier will be changed to MPI\_DATATYPE\_NULL.
- Any communication using this data type, that was already started, will be completed successfully.
- Datatypes that are defined in terms of this data type will still be usable.

### 8.2.2.2 Contiguous type

*This reference section gives the syntax for routines introduced in section 5.1.2.3.*

A contiguous datatype, created with a call to `MPI_Type_contiguous`,

```
Semantics:  
MPI_TYPE_CONTIGUOUS(count, oldtype, newtype)  
IN count: replication count (non-negative integer)  
IN oldtype: old datatype (handle)  
OUT newtype: new datatype (handle)  
  
C:  
int MPI_Type_contiguous(int count, MPI_Datatype oldtype, MPI_Datatype *newtype)  
  
Fortran:  
MPI_Type_contiguous(count, oldtype, newtype, ierror)  
INTEGER, INTENT(IN) :: count  
TYPE(MPI_Datatype), INTENT(IN) :: oldtype  
TYPE(MPI_Datatype), INTENT(OUT) :: newtype  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror  
  
Python:  
Create_contiguous(self, int count)
```

consists of a number of elements of a datatype, contiguous in memory. Sending one element of a contiguous type is fully equivalent to sending a number of elements of the constituent type.

```
// contiguous.c  
MPI_Datatype newvectortype;  
if (mytid==sender) {  
    MPI_Type_contiguous(count,MPI_DOUBLE,&newvectortype);  
    MPI_Type_commit(&newvectortype);  
    MPI_Send(source,1,newvectortype,receiver,0,comm);  
    MPI_Type_free(&newvectortype);  
} else if (mytid==receiver) {  
    MPI_Status recv_status;  
    int recv_count;  
    MPI_Recv(target,count,MPI_DOUBLE, sender,0,comm,  
             &recv_status);  
    MPI_Get_count(&recv_status,MPI_DOUBLE,&recv_count);  
    ASSERT(count==recv_count);  
}
```

### 8.2.2.3 Vector type

*This reference section gives the syntax for routines introduced in section 5.1.2.4.*

The `MPI_Type_vector` type can be used to create a type of regularly spaced blocks of data. All block lengths need to be the same, and the vector type is built out of a single constituent type.

```

Semantics:
MPI_TYPE_VECTOR(count, blocklength, stride, oldtype, newtype)
IN count: number of blocks (non-negative integer)
IN blocklength: number of elements in each block (non-negative integer)
IN stride: number of elements between start of each block (integer)
IN oldtype: old datatype (handle)
OUT newtype: new datatype (handle)

C:
int MPI_Type_vector
    (int count, int blocklength, int stride,
     MPI_Datatype oldtype, MPI_Datatype *newtype)

Fortran:
MPI_Type_vector(count, blocklength, stride, oldtype, newtype, ierror)
INTEGER, INTENT(IN) :: count, blocklength, stride
TYPE(MPI_Datatype), INTENT(IN) :: oldtype
TYPE(MPI_Datatype), INTENT(OUT) :: newtype
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python:
MPI.Datatype.Create_vector(self, int count, int blocklength, int stride)

```

In this example a vector type is created only on the sender, in order to send a strided subset of an array; the receiver receives the data as a contiguous block.

```

// vector.c
source = (double*) malloc(stride*count*sizeof(double));
target = (double*) malloc(count*sizeof(double));
MPI_Datatype newvectortype;
if (mytid==sender) {
    MPI_Type_vector(count,1,stride,MPI_DOUBLE,&newvectortype);
    MPI_Type_commit(&newvectortype);
    MPI_Send(source,1,newvectortype,the_other,0,comm);
    MPI_Type_free(&newvectortype);
} else if (mytid==receiver) {
    MPI_Status recv_status;
    int recv_count;
    MPI_Recv(target,count,MPI_DOUBLE,the_other,0,comm,
             &recv_status);
    MPI_Get_count(&recv_status,MPI_DOUBLE,&recv_count);
    ASSERT(recv_count==count);
}

```

**Exercise 8.1.** Let processor 0 have an array  $x$  of length  $10P$ , where  $P$  is the number of processors. Elements  $0, P, 2P, \dots, 9P$  should go to processor zero,  $1, P+1, 2P+1, \dots$  to processor 1, et cetera. Code this as a sequence of send/recv calls, using a vector datatype for the send, and a contiguous buffer for the receive. For simplicity, skip the send to/from zero. What is the most elegant solution if you want to include that case?  
For testing, define the array as  $x[i] = i$ .

#### 8.2.2.4 Indexed data

*This reference section gives the syntax for routines introduced in section 5.1.2.6.*

The indexed datatype is similar to the vector type, in the sense that it consists of a series of blocks of items, all of the same type. However, where the vector type was described by a single stride and blocklength, with MPI\_Type\_indexed you can specify the location and length of each block.

Semantics:

```
count [in] number of blocks --
      also number of entries in indices and blocklens
blocklens [in] number of elements in each block (array of nonnegative integers)
indices [in] displacement of each block in multiples of old_type
      (array of integers)
old_type [in] old datatype (handle)
newtype [out] new datatype (handle)
```

C:

```
int MPI_Type_indexed(int count, const int array_of_blocklengths[],
                     const int array_of_displacements[], MPI_Datatype oldtype, MPI_Datatype
                     *newtype)
```

Fortran:

```
MPI_Type_indexed(count, array_of_blocklengths, array_of_displacements,
                  oldtype, newtype, ierror)
INTEGER, INTENT(IN) :: count, array_of_blocklengths(count),
array_of_displacements(count)
TYPE(MPI_Datatype), INTENT(IN) :: oldtype
TYPE(MPI_Datatype), INTENT(OUT) :: newtype
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Python:

```
MPI.Datatype.Create_vector(self, blocklengths, displacements )
```

The following example picks items that are on prime number-indexed locations.

```
// indexed.c
displacements = (int*) malloc(count*sizeof(int));
blocklengths = (int*) malloc(count*sizeof(int));
source = (int*) malloc(totalcount*sizeof(int));
target = (int*) malloc(count*sizeof(int));
```

```

MPI_Datatype newvectortype;
if (mytid==sender) {
    MPI_Type_indexed(count,blocklengths,displacements,MPI_INT,&newvectortype);
    MPI_Type_commit(&newvectortype);
    MPI_Send(source,1,newvectortype,the_other,0,comm);
    MPI_Type_free(&newvectortype);
} else if (mytid==receiver) {
    MPI_Status recv_status;
    int recv_count;
    MPI_Recv(target,count,MPI_INT,the_other,0,comm,
             &recv_status);
    MPI_Get_count(&recv_status,MPI_INT,&recv_count);
    ASSERT(recv_count==count);
}

```

You can also `MPI_Type_create_hindexed` which describes blocks of a single old type, but with index locations in bytes, rather than in multiples of the old type.

```

int MPI_Type_create_hindexed
(int count, int blocklens[], MPI_Aint indices[],
 MPI_Datatype old_type,MPI_Datatype *newtype)

```

You can use this to pick all occurrences of a single component out of an array of structures. However, you need to be very careful with the index calculation. Use pointer arithmetic, as in the example in section 8.2.2.5. Another use of this function is in sending an `std<vector>`, that is, a vector object from the *C++ standard library*, if the component type is a pointer. No further explanation here.

### 8.2.2.5 Structure data

*This reference section gives the syntax for routines introduced in section 5.1.2.7.*

The `MPI_Type_create_struct` routine creates a type consisting of blocks of multiple datatypes, much like `MPI_Type_indexed` makes an array of blocks of a single type.

```

int MPI_Type_create_struct(
    int count, int blocklengths[], MPI_Aint displacements[],
    MPI_Datatype types[], MPI_Datatype *newtype);

```

**count** The number of blocks in this datatype. The `blocklengths`, `displacements`, `types` arguments have to be at least of this length.

**blocklengths** array containing the lengths of the blocks of each datatype.

**displacements** array describing the relative location of the blocks of each datatype.

**types** array containing the datatypes; each block in the new type is of a single datatype; there can be multiple blocks consisting of the same type.

## 8. MPI Reference

---

In this example, unlike the previous ones, both sender and receiver create the structure type. With structures it is no longer possible to send as a derived type and receive as a array of a simple type. (It would be possible to send as one structure type and receive as another, as long as they have the same *datatype signature*.)

```
// struct.c
struct object {
    char c;
    double x[2];
    int i;
};

MPI_Datatype newstructuretype;
int structlen = 3;
int blocklengths[structlen]; MPI_Datatype types[structlen];
MPI_Aint displacements[structlen];
// where are the components relative to the structure?
blocklengths[0] = 1; types[0] = MPI_CHAR;
displacements[0] = (size_t)&(myobject.c) - (size_t)&myobject;
blocklengths[1] = 2; types[1] = MPI_DOUBLE;
displacements[1] = (size_t)&(myobject.x[0]) - (size_t)&myobject;
blocklengths[2] = 1; types[2] = MPI_INT;
displacements[2] = (size_t)&(myobject.i) - (size_t)&myobject;
MPI_Type_create_struct(structlen,blocklengths,displacements,types,&newstructuretype);
MPI_Type_commit(&newstructuretype);

{
    MPI_Aint typesize;
    MPI_Type_extent(newstructuretype,&typesize);
    if (mytid==0) printf("Type extent: %d bytes\n",typesize);
}

if (mytid==sender) {
    MPI_Send(&myobject,1,newstructuretype,the_other,0,comm);
} else if (mytid==receiver) {
    MPI_Recv(&myobject,1,newstructuretype,the_other,0,comm,MPI_STATUS_IGNORE)
}
MPI_Type_free(&newstructuretype);
```

Note the displacement calculations in this example, which involve some not so elegant pointer arithmetic. It would have been incorrect to write

```
displacement[0] = 0;
displacement[1] = displacement[0] + sizeof(char);
```

since you do not know the way the *compiler* lays out the structure in memory<sup>1</sup>. The space that MPI takes for a structure type can be queried with `MPI_Type_extent`.

---

1. Homework question: what does the language standard say about this?

---

```
int MPI_Type_extent(
    MPI_Datatype datatype, MPI_Aint *extent);
```

(There is a deprecated function `MPI_Type_struct` with the same functionality.)

### 8.2.3 Packed data

*This reference section gives the syntax for routines introduced in section 5.1.4.*

With `MPI_PACK` data elements can be added to a buffer one at a time. The `position` parameter is updated each time by the packing routine.

```
int MPI_Pack(
    void *inbuf, int incount, MPI_Datatype datatype,
    void *outbuf, int outcount, int *position,
    MPI_Comm comm);
```

Conversely, `MPI_UNPACK` retrieves one element from the buffer at a time. You need to specify the MPI datatype.

```
int MPI_Unpack(
    void *inbuf, int insize, int *position,
    void *outbuf, int outcount, MPI_Datatype datatype,
    MPI_Comm comm);
```

A packed buffer is sent or received with a datatype of `MPI_PACKED`. The sending routine uses the `position` parameter to specify how much data is sent, but the receiving routine does not know this value a priori, so has to specify an upper bound.

```
// pack.c
if (mytid==sender) {
    MPI_Pack(&nstarts, 1, MPI_INT, buffer, buflen, &position, comm);
    for (int i=0; i<nstarts; i++) {
        double value = rand() / (double) RAND_MAX;
        MPI_Pack(&value, 1, MPI_DOUBLE, buffer, buflen, &position, comm);
    }
    MPI_Pack(&nstarts, 1, MPI_INT, buffer, buflen, &position, comm);
    MPI_Send(buffer, position, MPI_PACKED, other, 0, comm);
} else if (mytid==receiver) {
    int irecv_value;
    double xrecv_value;
    MPI_Recv(buffer, buflen, MPI_PACKED, other, 0, comm, MPI_STATUS_IGNORE);
    MPI_Unpack(buffer, buflen, &position, &nstarts, 1, MPI_INT, comm);
    for (int i=0; i<nstarts; i++) {
        MPI_Unpack(buffer, buflen, &position, &xrecv_value, 1, MPI_DOUBLE, comm);
    }
}
```

```
    MPI_Unpack(buffer,buflen,&position,&irecv_value,1,MPI_INT,comm);
    ASSERT(irecv_value==nsends);
}
```

You can precompute the size of the required buffer as follows:

```
int MPI_Pack_size(
    int incount, MPI_Datatype datatype,
    MPI_Comm comm, int *size);
```

Add one time MPI\_BSEND\_OVERHEAD.

### 8.3 Blocking communication

*This reference section gives the syntax for routines introduced in section 4.2.2.*

The blocking send command:

```
C:
int MPI_Send(
    const void* buf, int count, MPI_Datatype datatype,
    int dest, int tag, MPI_Comm comm)

Semantics:
IN buf: initial address of send buffer (choice)
IN count: number of elements in send buffer (non-negative integer)
IN datatype: datatype of each send buffer element (handle)
IN dest: rank of destination (integer)
IN tag: message tag (integer)
IN comm: communicator (handle)

Fortran:
MPI_Send(buf, count, datatype, dest, tag, comm, ierror)
TYPE(*), DIMENSION(..), INTENT(IN) :: buf
INTEGER, INTENT(IN) :: count, dest, tag
TYPE(MPI_Datatype), INTENT(IN) :: datatype
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python native:
Comm.send(self, obj, int dest, int tag=0)
Python numpy:
Comm.Send(self, buf, int dest, int tag=0)
```

This routine may not block for small messages; to force blocking behaviour use MPI\_Ssend with the same argument list. [http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI\\_Ssend.html](http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI_Ssend.html)

The basic blocking receive command:

```
C:  
int MPI_Recv(  
    void* buf, int count, MPI_Datatype datatype,  
    int source, int tag, MPI_Comm comm, MPI_Status *status)  
  
Semantics:  
OUT buf: initial address of receive buffer (choice)  
IN count: number of elements in receive buffer (non-negative integer)  
IN datatype: datatype of each receive buffer element (handle)  
IN source: rank of source or MPI_ANY_SOURCE (integer)  
IN tag: message tag or MPI_ANY_TAG (integer)  
IN comm: communicator (handle)  
OUT status: status object (Status)  
  
Fortran:  
MPI_Recv(buf, count, datatype, source, tag, comm, status, ierror)  
TYPE(*), DIMENSION(..) :: buf  
INTEGER, INTENT(IN) :: count, source, tag  
TYPE(MPI_Datatype), INTENT(IN) :: datatype  
TYPE(MPI_Comm), INTENT(IN) :: comm  
TYPE(MPI_Status) :: status  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror  
  
Python native:  
recvbuf = Comm.recv(self, buf=None, int source=ANY_SOURCE, int tag=ANY_TAG,  
                     Status status=None)  
Python numpy:  
Comm.Recv(self, buf, int source=ANY_SOURCE, int tag=ANY_TAG,  
          Status status=None)
```

The `count` argument indicates the maximum length of a message; the actual length of the received message can be determined from the `status` object. See section ?? for more about the `status` object.

The following code is guaranteed to block, since a `MPI_Recv` always blocks:

```
// recvblock.c  
other = 1-mytid;  
MPI_Recv(&recvbuf,1,MPI_INT,other,0,comm,&status);  
MPI_Send(&sendbuf,1,MPI_INT,other,0,comm);  
printf("This statement will not be reached on %d\n",mytid);
```

On the other hand, if we put the send call before the receive, code may not block for small messages that fall under the *eager limit*.

In this example we send gradually larger messages. From the screen output you can see what the largest message was that fell under the eager limit; after that the code hangs because of a deadlock.

```
// sendblock.c  
other = 1-mytid;
```

```
/* loop over increasingly large messages */
for (int size=1; size<2000000000; size*=10) {
    sendbuf = (int*) malloc(size*sizeof(int));
    recvbuf = (int*) malloc(size*sizeof(int));
    if (!sendbuf || !recvbuf) {
        printf("Out of memory\n"); MPI_Abort(comm,1);
    }
    MPI_Send(sendbuf,size,MPI_INT,other,0,comm);
    MPI_Recv(recvbuf,size,MPI_INT,other,0,comm,&status);
    /* If control reaches this point, the send call
       did not block. If the send call blocks,
       we do not reach this point, and the program will hang.
    */
    if (mytid==0)
        printf("Send did not block for size %d\n",size);
    free(sendbuf); free(recvbuf);
}

// sendblock.F90
other = 1-mytid
size = 1
do
    allocate(sendbuf(size)); allocate(recvbuf(size))
    print *,size
    call MPI_Send(sendbuf,size,MPI_INTEGER,other,0,comm,err)
    call MPI_Recv(recvbuf,size,MPI_INTEGER,other,0,comm,status,err)
    if (mytid==0) then
        print *, "MPI_Send did not block for size",size
    end if
    deallocate(sendbuf); deallocate(recvbuf)
    size = size*10
    if (size>2000000000) goto 20
end do
20   continue

// sendblock.py
size = 1
while size<2000000000:
    sendbuf = np.empty(size, dtype=np.int)
    recvbuf = np.empty(size, dtype=np.int)
    comm.Send(sendbuf, dest=other)
    comm.Recv(recvbuf, source=other)
    if procid<other:
```

```
    print "Send did not block for",size
    size *= 10
```

If you want a code to behave the same for all message sizes, you force the send call to be blocking by using `MPI_Ssend`:

```
// ssendblock.c
other = 1-mytid;
sendbuf = (int*) malloc(sizeof(int));
recvbuf = (int*) malloc(sizeof(int));
size = 1;
MPI_Ssend(sendbuf,size,MPI_INT,other,0,comm);
MPI_Recv(recvbuf,size,MPI_INT,other,0,comm,&status);
printf("This statement is not reached\n");
```

### 8.3.1 Receive status

*This reference section gives the syntax for routines introduced in section 4.2.7.2.*

Any time you receive data, there can be an `MPI_Status` object describing the data that was received.

C:

```
MPI_Status status;
```

Fortran:

```
integer :: status(MPI_STATUS_SIZE)
```

Python:

```
MPI.Status() # returns object
```

*Fortran note* In Fortran there is no `MPI_Status` type, instead an integer array is created by the user.

*Python note* The status object is created as a python object. See also section 8.8.3.

The use of a status parameter can be necessary if you use `MPI_ANY_SOURCE` or `MPI_ANY_TAG` in the description of the receive message. If you are not interested in the status information, you can use the values `MPI_STATUS_IGNORE` for `MPI_Wait` and `MPI_Waitany`, or `MPI_STATUSES_IGNORE` for `MPI_Waitall` and `MPI_Waitsome`.

A receive call has a `count` parameter, but this describes the length of the buffer, not the amount of data expected. That quantity can be retrieved with `MPI_Get_count`.

```
// C:
int MPI_Get_count(MPI_Status *status,MPI_Datatype datatype,
                  int *count)
! Fortran:
```

```
MPI_Get_count (INTEGER status (MPI_STATUS_SIZE), INTEGER datatype,
               INTEGER count, INTEGER ierror)
```

The status object is returned when the message is received. Thus, with `MPI_Recv` it is returned explicitly, but with `MPI_Irecv` it is returned from the `MPI_Wait...` call.

In section ?? we mentioned the master-worker model as one opportunity for inspecting the `MPI_SOURCE` field of the `MPI_Status` object.

```
C:
int status.MPI_SOURCE;

F:

Python:
status.Get_source() # returns int
```

Here is a small example: the tasks perform a variable amount of work (modeled here by a random wait) before sending a message to the master. The master waits for any source, and inspects the status field to report where the message comes from.

```
// anysource.c
if (mytid==ntids-1) {
    int *recv_buffer;
    MPI_Status status;

    recv_buffer = (int*) malloc((ntids-1)*sizeof(int));

    for (int p=0; p<ntids-1; p++) {
        err = MPI_Recv(recv_buffer+p, 1, MPI_INT, MPI_ANY_SOURCE, 0, comm,
                      &status); CHK(err);
        int sender = status.MPI_SOURCE;
        printf("Message from sender=%d: %d\n",
               sender, recv_buffer[p]);
    }
} else {
    float randomfraction = (rand() / (double)RAND_MAX);
    int randomwait = (int) ( ntids * randomfraction );
    printf("process %d waits for %e/%d=%d\n",
           mytid, randomfraction, ntids, randomwait);
    sleep(randomwait);
    err = MPI_Send(&randomwait, 1, MPI_INT, ntids-1, 0, comm); CHK(err);
}
```

## 8.4 Deadlock-free blocking messages

*This reference section gives the syntax for routines introduced in section 4.2.4.*

If messsages are send roughly in pairs, the MPI\_Sendrecv call is an easy way to prevent deadlock. Here you specify both the target of a send and the source of a receive, which can be same in case of a pairwise exchange of data, but they need not be the same. To swap equal-sized buffers you can use MPI\_Sendrecv\_replace.

```
int MPI_Sendrecv(
    void *sendbuf, int sendcount, MPI_Datatype sendtype,
    int dest, int sendtag,
    void *recvbuf, int recvcount, MPI_Datatype recvtype,
    int source, int recvtag,
    MPI_Comm comm, MPI_Status *status)
int MPI_Sendrecv_replace(
    void *buf, int count, MPI_Datatype datatype,
    int dest, int sendtag,
    int source, int recvtag,
    MPI_Comm comm, MPI_Status *status)
```

As an example we set up a ring of three processors: each process sends to its right neighbour, and receives from its left neighbour.

```
// sendrecv.c
right = (mytid+1)%3; left = (mytid+2)%3;
MPI_Sendrecv( &my_data,1,MPI_INTEGER, right,0,
&other_data,1,MPI_INTEGER, left,0,
comm,MPI_STATUS_IGNORE);
```

### 8.4.1 Non-blocking communication

*This reference section gives the syntax for routines introduced in section 4.2.6.*

The non-blocking routines have much the same parameter list as the blocking ones, with the addition of an MPI\_Request parameter. The MPI\_Isend routine does not have an MPI\_Status parameter, which has moved to the ‘wait’ routine.

```
int MPI_Isend(void *buf,
    int count, MPI_Datatype datatype, int dest, int tag,
    MPI_Comm comm, MPI_Request *request)
```

[http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI\\_Isend.html](http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI_Isend.html)

```
int MPI_Irecv(void *buf,
    int count, MPI_Datatype datatype, int source, int tag,
    MPI_Comm comm, MPI_Request *request)
```

## 8. MPI Reference

---

[http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI\\_Irecv.html](http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI_Irecv.html)

*Fortran note* The request parameter is an integer.

There are various ‘wait’ routines. Since you will often do at least one send and one receive, this routine is useful:

```
int MPI_Waitall(int count, MPI_Request array_of_requests[],  
                 MPI_Status array_of_statuses[])
```

[http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI\\_Waitall.html](http://www.mcs.anl.gov/research/projects/mpi/www/www3/MPI_Waitall.html)

Here is a simple code that does a non-blocking exchange between two processors:

```
// irecvnonblock.c  
MPI_Request request[2];  
MPI_Status status[2];  
other = ntids-mytid;  
MPI_Irecv(&recvbuf, 1, MPI_INT, other, 0, comm, &(request[0]));  
MPI_Isend(&sendbuf, 1, MPI_INT, other, 0, comm, &(request[1]));  
MPI_Waitall(2, request, status);
```

It is possible to omit the status array by specifying `MPI_STATUSES_IGNORE`. Other routines are `MPI_Wait` for a single request, and `MPI_Waitsome`, `MPI_Waitany`.

The above fragment is unrealistically simple. In a more general scenario we have to manage send and receive buffers: we need as many buffers as there are simultaneous non-blocking sends and receives.

Instead of waiting for all messages, we can wait for any message to come with `MPI_Waitany`, and process the receive data as it comes in.

```
// irecv_source.c  
if (mytid==ntids-1) {  
    int *recv_buffer;  
    MPI_Request *request; MPI_Status status;  
    recv_buffer = (int*) malloc((ntids-1)*sizeof(int));  
    request = (MPI_Request*) malloc((ntids-1)*sizeof(MPI_Request));  
  
    for (int p=0; p<ntids-1; p++) {  
        ierr = MPI_Irecv(recv_buffer+p, 1, MPI_INT, p, 0, comm,  
                         request+p); CHK(ierr);  
    }  
    for (int p=0; p<ntids-1; p++) {  
        int index, sender;  
        MPI_Waitany(ntids-1, request, &index, &status); //MPI_STATUS_IGNORE);  
        if (index!=status.MPI_SOURCE)  
            printf("Mismatch index %d vs source %d\n", index, status.MPI_SOURCE);  
        printf("Message from %d: %d\n", index, recv_buffer[index]);  
    }  
}
```

```
}
```

Note the `MPI_STATUS_IGNORE` parameter: we know everything about the incoming message, so we do not need to query a status object. Contrast this with the example in section ??.

*Fortran note* The `index` parameter is the index in the array of requests, so it uses *1-based indexing*.

```
// irecv_source.F90
if (mytid==ntids-1) then
    do p=1,ntids-1
        print *, "post"
        call MPI_Irecv(recv_buffer(p), 1, MPI_INTEGER, p-1, 0, comm, &
                      requests(p), err)
    end do
    do p=1,ntids-1
        call MPI_Waitany(ntids-1, requests, index, MPI_STATUS_IGNORE, err)
        write(*,'("Message from",i3,":",i5)') index, recv_buffer(index)
    end do
```

#### 8.4.1.1 Wait and test calls

This reference section gives the syntax for routines introduced in section 4.2.6.1.

Wait calls block until any or all of a set of requests are fulfilled.

Semantics:

```
int MPI_Waitany(
    int count, MPI_Request array_of_requests[], int *index,
    MPI_Status *status)

IN count: list length (non-negative integer)
INOUT array_of_requests: array of requests (array of handles)
OUT index: index of handle for operation that completed (integer)
OUT status: status object (Status)
```

C:

```
MPI_Waitany(count, array_of_requests, index, status, ierror)
```

Fortran:

```
INTEGER, INTENT(IN) :: count
TYPE(MPI_Request), INTENT(INOUT) :: array_of_requests(count)
INTEGER, INTENT(OUT) :: index
TYPE(MPI_Status) :: status
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Python:

```
MPI.Request.Waitany( requests, status=None )
```

```
class method, returns index

Semantics:
MPI_WAITALL( count, array_of_requests, array_of_statuses)
IN countlists length (non-negative integer)
INOUT array_of_requestsarray of requests (array of handles)
OUT array_of_statusesarray of status objects (array of Status)

C:
int MPI_Waitall(int count, MPI_Request array_of_requests[], MPI_Status array_of_st
```

Fortran:

```
MPI_Waitall(count, array_of_requests, array_of_statuses, ierror)
INTEGER, INTENT(IN) :: count
TYPE(MPI_Request), INTENT(INOUT) :: array_of_requests(count)
TYPE(MPI_Status) :: array_of_statuses(*)
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Test calls are non-blocking versions of the corresponding wait calls.

```
C:
int MPI_Testany(
    int count, MPI_Request array_of_requests[],
    int *index, int *flag, MPI_Status *status)

Fortran:

MPI_Testany(count, array_of_requests, index, flag, status, ierror)
INTEGER, INTENT(IN) :: count
TYPE(MPI_Request), INTENT(INOUT) :: array_of_requests(count)
INTEGER, INTENT(OUT) :: index
LOGICAL, INTENT(OUT) :: flag
TYPE(MPI_Status) :: status
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Semantics:
MPI_TESTALL(count, array_of_requests, flag, array_of_statuses)
IN countlists length (non-negative integer)
INOUT array_of_requestsarray of requests (array of handles)
OUT flag(logical)
OUT array_of_statusesarray of status objects (array of Status)

C:
int MPI_Testall(
    int count, MPI_Request array_of_requests[],
    int *flag, MPI_Status array_of_statuses[])

Fortran:
MPI_Testall(count, array_of_requests, flag, array_of_statuses, ierror)
INTEGER, INTENT(IN) :: count
```

```
TYPE(MPI_Request), INTENT(INOUT) :: array_of_requests(count)
LOGICAL, INTENT(OUT) :: flag
TYPE(MPI_Status) :: array_of_statuses(*)
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

### 8.4.2 Buffered communication

*This reference section gives the syntax for routines introduced in section 4.2.9.*

`MPI_Buffer_attach`

```
int MPI_Buffer_attach(
    void *buffer, int size);
```

where the size is indicated in bytes. The possible error codes are

- `MPI_SUCCESS` the routine completed successfully.
- `MPI_ERR_BUFFER` The buffer pointer is invalid; this typically means that you have supplied a null pointer.
- `MPI_ERR_INTERN` An internal error in MPI has been detected.

The buffer is detached with `MPI_Buffer_detach`:

```
int MPI_Buffer_detach(
    void *buffer, int *size);
```

This returns the address and size of the buffer; the call blocks until all buffered messages have been delivered.

You can compute the needed size of the buffer with `MPI_Pack_size`; see section 8.2.3.

`MPI_Bsend`

```
int MPI_Bsend(
    const void *buf, int count, MPI_Datatype datatype,
    int dest, int tag, MPI_Comm comm)
```

The asynchronous version is `MPI_Ibsend`.

You can force delivery by

```
MPI_Buffer_detach( &b, &n );
MPI_Buffer_attach( b, n );
```

### 8.4.3 Persistent communication

*This reference section gives the syntax for routines introduced in section 4.2.10.*

## 8. MPI Reference

---

The calls `MPI_Send_init` and `MPI_Recv_init` for creating a persistent communication have the same syntax as those for non-blocking sends and receives. The difference is that they do not start an actual communication, they only create the request object.

```
C:  
int MPI_Send_init(  
    const void* buf, int count, MPI_Datatype datatype,  
    int dest, int tag, MPI_Comm comm, MPI_Request *request)  
  
Fortran:  
MPI_Send_init(buf, count, datatype, dest, tag, comm, request, ierror)  
TYPE(*), DIMENSION(..), INTENT(IN), ASYNCHRONOUS :: buf  
INTEGER, INTENT(IN) :: count, dest, tag  
TYPE(MPI_Datatype), INTENT(IN) :: datatype  
TYPE(MPI_Comm), INTENT(IN) :: comm  
TYPE(MPI_Request), INTENT(OUT) :: request  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror  
  
Python:  
MPI.Comm.Send_init(self, buf, int dest, int tag=0)  
  
Semantics:  
IN buf: initial address of send buffer (choice)  
IN count: number of elements sent (non-negative integer)  
IN datatype: type of each element (handle)  
IN dest: rank of destination (integer)  
IN tag: message tag (integer)  
IN comm: communicator (handle)  
OUT request: communication request (handle)  
  
C:  
int MPI_Recv_init(  
    void* buf, int count, MPI_Datatype datatype,  
    int source, int tag, MPI_Comm comm, MPI_Request *request)  
  
Fortran:  
MPI_Recv_init(buf, count, datatype, source, tag, comm, request,  
ierror)  
TYPE(*), DIMENSION(..), ASYNCHRONOUS :: buf  
INTEGER, INTENT(IN) :: count, source, tag  
TYPE(MPI_Datatype), INTENT(IN) :: datatype  
TYPE(MPI_Comm), INTENT(IN) :: comm  
TYPE(MPI_Request), INTENT(OUT) :: request  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror  
  
Python:  
MPI.Comm.Recv_init(  
    self, buf, int source=ANY_SOURCE, int tag=ANY_TAG)  
  
Semantics:
```

---

```

OUT buf: initial address of receive buffer (choice)
IN count: number of elements received (non-negative integer)
IN datatype: type of each element (handle)
IN source: rank of source or MPI_ANY_SOURCE (integer)
IN tag: message tag or MPI_ANY_TAG (integer)
IN comm: mcommunicator (handle)
OUT request: communication request (handle)

```

Given these request object, a communication (both send and receive) is then started with `MPI_Start` for a single request or `MPI_Start_all` for multiple requests, given in an array.

```

int MPI_Start(MPI_Request *request)

C:
int MPI_Startall(int count, MPI_Request array_of_requests[])

Fortran:
MPI_Startall(count, array_of_requests, ierror)
INTEGER, INTENT(IN) :: count
TYPE(MPI_Request), INTENT(INOUT) :: array_of_requests(count)
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
MPI_STARTALL(COUNT, ARRAY_OF_REQUESTS, IERROR)
INTEGER COUNT, ARRAY_OF_REQUESTS(*), IERROR

Python:
MPI.Prequest.Startall(type cls, requests)

Semantics:
IN countlist length (non-negative integer)
INOUT array_of_requestsarray of requests (array of handle)

```

These are equivalent to starting an `Irecv` or `Isend`; correspondingly, it is necessary to issue an `MPI_Wait...` call (section 8.4.1) to determine their completion.

After a request object has been used, possibly multiple times, it can be freed; see 8.4.4.

In the following example a ping-pong is implemented with persistent communication.

```

// persist.c
if (mytid==src) {
    MPI_Send_init(send,s,MPI_DOUBLE,tgt,0,comm,requests+0);
    MPI_Recv_init(recv,s,MPI_DOUBLE,tgt,0,comm,requests+1);
    printf("Size %d\n",s);
    t[cnt] = MPI_Wtime();
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Startall(2,requests);
        MPI_Waitall(2,requests,MPI_STATUSES_IGNORE);
    }
    t[cnt] = MPI_Wtime()-t[cnt];
}

```

```
    MPI_Request_free(requests+0); MPI_Request_free(requests+1);
} else if (mytid==tgt) {
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Recv(recv,s,MPI_DOUBLE,src,0,comm,MPI_STATUS_IGNORE);
        MPI_Send(recv,s,MPI_DOUBLE,src,0,comm);
    }
}
```

As with ordinary send commands, there are the variants `MPI_Bsend_init`, `MPI_Ssend_init`, `MPI_Rsend_init`.

### 8.4.4 About `MPI_Request`

An `MPI_Request` object is not actually an object, unlike `MPI_Status`. Instead it is an (opaque) pointer. This means that when you call, for instance, `MPI_Irecv`, MPI will allocate an actual request object, and return its address in the `MPI_Request` variable.

Correspondingly, calls to `MPI_Wait...` or `MPI_Test` free this object. If your application is such that you do not use ‘wait’ call, you can free the request object explicitly with `MPI_Request_free`.

```
int MPI_Request_free(MPI_Request *request)
```

You can inspect the status of a request without freeing the request object with `MPI_Request_get_status`:

```
int MPI_Request_get_status(
    MPI_Request request,
    int *flag,
    MPI_Status *status
);
```

## 8.5 One-sided communication

*This reference section gives the syntax for routines introduced in section 4.3.*

### 8.5.1 Windows and epochs

*This reference section gives the syntax for routines introduced in section 4.3.1.*

C syntax for `MPI_Win_create`

```
C:
int MPI_Win_create
    (void *base, MPI_Aint size, int disp_unit,
     MPI_Info info, MPI_Comm comm, MPI_Win *win)
```

```

Frotran:
MPI_Win_create(base, size, disp_unit, info, comm, win, ierror)
  TYPE(*), DIMENSION(..), ASYNCHRONOUS :: base
  INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) :: size
  INTEGER, INTENT(IN) :: disp_unit
  TYPE(MPI_Info), INTENT(IN) :: info
  TYPE(MPI_Comm), INTENT(IN) :: comm
  TYPE(MPI_Win), INTENT(OUT) :: win
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python:
MPI.Win.Create
  (memory, int disp_unit=1,
   Info info=INFO_NULL, Intracomm comm=COMM_SELF)

```

The data array must not be PARAMETER or static const.

The size parameter is measured in bytes. In C this is easily done with the `sizeof` operator; for doing this calculation in Fortran, see section 8.8.2.3.

The MPI\_Info parameter can be used to pass implementation-dependent information:

```

MPI_Info info;
MPI_Info_create(&info);
MPI_Info_set(info,"no_locks","true");
MPI_Win_create( ... info ... &win);
MPI_Info_free(&info);

```

It is always valid to use MPI\_INFO\_NULL.

```

MPI_Alloc_mem
  int MPI_Alloc_mem(MPI_Aint size, MPI_Info info, void *baseptr)

```

### 8.5.1.1 Window information

```

MPI_Win_get_attr(win, MPI_WIN_BASE, &base, &flag),
MPI_Win_get_attr(win, MPI_WIN_SIZE, &size, &flag),
MPI_Win_get_attr(win, MPI_WIN_DISP_UNIT, &disp_unit, &flag),
MPI_Win_get_attr(win, MPI_WIN_CREATE_FLAVOR, &create_kind, &flag), and
MPI_Win_get_attr(win, MPI_WIN_MODEL, &memory_model, &flag) will return in b

int MPI_Win_get_group(MPI_Win win, MPI_Group *group)
MPI_Win_get_group(win, group, ierror)
  TYPE(MPI_Win), INTENT(IN) :: win
  TYPE(MPI_Group), INTENT(OUT) :: group
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror

```

```
int MPI_Win_set_info(MPI_Win win, MPI_Info info)
MPI_Win_set_info(win, info, ierror)
  TYPE(MPI_Win), INTENT(IN) :: win
  TYPE(MPI_Info), INTENT(IN) :: info
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror

int MPI_Win_get_info(MPI_Win win, MPI_Info *info_used)
MPI_Win_get_info(win, info_used, ierror)
  TYPE(MPI_Win), INTENT(IN) :: win
  TYPE(MPI_Info), INTENT(OUT) :: info_used
  INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

### 8.5.2 Remote memory access

*This reference section gives the syntax for routines introduced in section 4.3.3.*

The MPI\_Put routine is used to put data in the window of a target process

C:

```
int MPI_Put(
  const void *origin_addr, int origin_count, MPI_Datatype origin_datatype,
  int target_rank, MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype,
  MPI_Win win)
```

Semantics:

```
IN origin_addr: initial address of origin buffer (choice)
IN origin_count: number of entries in origin buffer (non-negative integer)
IN origin_datatype: datatype of each entry in origin buffer (handle)
IN target_rank: rank of target (non-negative integer)
IN target_disp: displacement from start of window to target buffer (non-negative integer)
IN target_count: number of entries in target buffer (non-negative integer)
IN target_datatype: datatype of each entry in target buffer (handle)
IN win: window object used for communication (handle)
```

Fortran:

```
MPI_Put(origin_addr, origin_count, origin_datatype,
         target_rank, target_disp, target_count, target_datatype, win, ierror)
TYPE(*), DIMENSION(..), INTENT(IN), ASYNCHRONOUS :: origin_addr
INTEGER, INTENT(IN) :: origin_count, target_rank, target_count
TYPE(MPI_Datatype), INTENT(IN) :: origin_datatype, target_datatype
INTEGER(KIND=MPI_ADDRESS_KIND), INTENT(IN) :: target_disp
TYPE(MPI_Win), INTENT(IN) :: win
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Python:

---

```
win.Put(self, origin, int target_rank, target=None)
```

The data is written in the buffer of the target window, using the window parameters that were specified on the target. Specifically, data is written starting at

```
window_base + target_disp × disp_unit.
```

The MPI\_Get call is very similar.

```
int MPI_Get(void *origin_addr, int origin_count, MPI_Datatype
            origin_datatype, int target_rank, MPI_Aint target_disp,
            int target_count, MPI_Datatype target_datatype, MPI_Win
            win)
```

Here is a single put operation. Note that the window create and window fence calls are collective, so they have to be performed on all processors of the communicator that was used in the create call.

```
// putfence.c
MPI_Win the_window;
MPI_Win_create(&window_data,2*sizeof(int),sizeof(int),
               MPI_INFO_NULL,comm,&the_window);
MPI_Win_fence(0,the_window);
if (mytid==0) {
    MPI_Put( /* data on origin: */   &my_number, 1,MPI_INT,
             /* data on target: */ other,1,      1,MPI_INT,
             the_window);
}
MPI_Win_fence(0,the_window);
MPI_Win_free(&the_window);
```

Very similar, a get operation.

A third one-sided routine is MPI\_Accumulate which does a reduction operation on the results that are being put:

```
MPI_Accumulate (
    void *origin_addr, int origin_count, MPI_Datatype origin_datatype,
    int target_rank,
    MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype,
    MPI_Op op,MPI_Win window)
```

### 8.5.2.1 Request-based operations

Analogous to MPI\_Isend there are request based one-sided operations:

```
C:  
int MPI_Rput(  
    const void *origin_addr, int origin_count, MPI_Datatype origin_datatype,  
    int target_rank, MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype,  
    MPI_Win win, MPI_Request *request)  
  
Semantics:  
IN origin_addr: initial address of origin buffer (choice)  
IN origin_count: number of entries in origin buffer (non-negative integer)  
IN origin_datatype: datatype of each entry in origin buffer (handle)  
IN target_rank: rank of target (non-negative integer)  
IN target_disp: displacement from start of window to target buffer (non-negative integer)  
IN target_count: number of entries in target buffer (non-negative integer)  
IN target_datatype: datatype of each entry in target buffer (handle)  
IN win: window object used for communication (handle)  
OUT request: RMA request (handle)
```

and similarly `MPI_Rget` and `MPI_Raccumulate`.

These only apply to passive target synchronization. Any `MPI_Win_flush...` call also terminates these transfers.

### 8.5.3 Active target synchronization

*This reference section gives the syntax for routines introduced in section 4.3.2.*

```
MPI_Win_fence (int assert, MPI_Win win)
```

### 8.5.4 Assertions

The `MPI_Win_fence` call, as well `MPI_Win_start` and such, take an argument through which assertions can be passed about the activity before, after, and during the epoch. The value zero is always allowed, by you can make your program more efficient by specifying one or more of the following, combined by bitwise OR in C/C++ or `IOR` in Fortran.

**MPI\_WIN\_START** Supports the option:

**MPI\_MODE\_NOCHECK** the matching calls to `MPI_WIN_POST` have already completed on all target processes when the call to `MPI_WIN_START` is made. The nocheck option can be specified in a start call if and only if it is specified in each matching post call. This is similar to the optimization of “ready-send” that may save a handshake when the handshake is implicit in the code. (However, ready-send is matched by a regular receive, whereas both start and post must specify the nocheck option.)

**MPI\_WIN\_POST** supports the following options:

**MPI\_MODE\_NOCHECK** the matching calls to `MPI_WIN_START` have not yet occurred on any origin processes when the call to `MPI_WIN_POST` is made. The nocheck option can be specified by a post call if and only if it is specified by each matching start call.

**MPI\_MODE\_NOSTORE** the local window was not updated by local stores (or local get or receive calls) since last synchronization. This may avoid the need for cache synchronization at the post call.

**MPI\_MODE\_NOPUT** the local window will not be updated by put or accumulate calls after the post call, until the ensuing (wait) synchronization. This may avoid the need for cache synchronization at the wait call.

**MPI\_WIN\_FENCE** supports the following options:

**MPI\_MODE\_NOSTORE** the local window was not updated by local stores (or local get or receive calls) since last synchronization.

**MPI\_MODE\_NOPUT** the local window will not be updated by put or accumulate calls after the fence call, until the ensuing (fence) synchronization.

**MPI\_MODE\_NOPRECEDE** the fence does not complete any sequence of locally issued RMA calls. If this assertion is given by any process in the window group, then it must be given by all processes in the group.

**MPI\_MODE\_NOSUCCEED** the fence does not start any sequence of locally issued RMA calls. If the assertion is given by any process in the window group, then it must be given by all processes in the group.

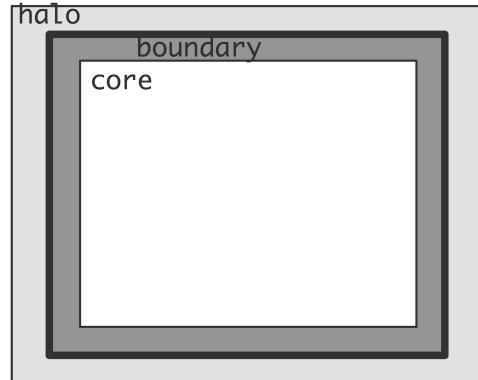
**MPI\_WIN\_LOCK** supports the following option:

**MPI\_MODE\_NOCHECK** no other process holds, or will attempt to acquire a conflicting lock, while the caller holds the window lock. This is useful when mutual exclusion is achieved by other means, but the coherence operations that may be attached to the lock and unlock calls are still required.

As an example, let's look at *halo update*. The array A is updated using the local values and the halo that comes from bordering processors, either through Put or Get operations.

In a first version we separate computation and communication. Each iteration has two fences. Between the two fences in the loop body we do the MPI\_Put operation; between the second and and first one of the next iteration there is only computation, so we add the NOPRECEDE and NOSUCCEED assertions. The NOSTORE assertion states that the local window was not updated: the Put operation only works on remote windows.

```
for ( .... ) {
    update(A);
    MPI_Win_fence(MPI_MODE_NOPRECEDE, win);
    for(i=0; i < toneighbors; i++)
        MPI_Put( ... );
    MPI_Win_fence((MPI_MODE_NOSTORE | MPI_MODE_NOSUCCEED), win);
}
```



Next, we split the update in the core part, which can be done purely from local values, and the boundary, which needs local and halo values. Update of the core can overlap the communication of the halo.

```
for ( .... ) {
    update_boundary (A) ;
    MPI_Win_fence ((MPI_MODE_NOPUT | MPI_MODE_NOPRECEDE), win) ;
    for(i=0; i < fromneighbors; i++)
        MPI_Get( ... );
    update_core (A) ;
    MPI_Win_fence (MPI_MODE_NOSUCCEEDED, win) ;
}
```

The NOPRECEDE and NOSUCCEEDED assertions still hold, but the Get operation implies that instead of NOSTORE in the second fence, we use NOPUT in the first.

### 8.5.5 More active target synchronization

*This reference section gives the syntax for routines introduced in section 4.3.5.*

The ‘fence’ mechanism (section 8.5.3) uses a global synchronization on the communicator of the window, which may lead to performance inefficiencies if processors are not in step with each other. There is a mechanism that is more fine-grained, by using synchronization only on a processor *group*. This takes four different calls, two for starting and two for ending the epoch, separately for target and origin.

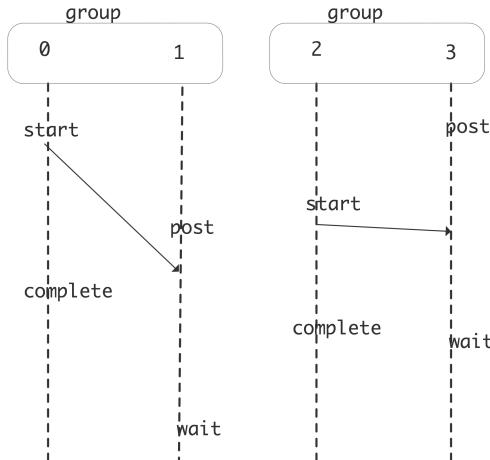


Figure 8.1: Window locking calls in fine-grained active target synchronization

You start and complete an *exposure epoch* with :

```
int MPI_Win_post (MPI_Group group, int assert, MPI_Win win)
int MPI_Win_wait (MPI_Win win)
```

In other words, this turns your window into the *target* for a remote access.

You start and complete an access epoch with :

```
int MPI_Win_start(MPI_Group group, int assert, MPI_Win win)
int MPI_Win_complete(MPI_Win win)
```

In other words, these calls border the access to a remote window, with the current processor being the *origin* of the remote access.

In the following snippet a single processor puts data on one other. Note that they both have their own definition of the group, and that the receiving process only does the post and wait calls.

```
// postwaitwin.c
if (mytid==origin) {
    MPI_Group_incl(all_group,1,&target,&two_group);
    // access
    MPI_Win_start(two_group,0,the_window);
    MPI_Put( /* data on origin: */    &my_number, 1,MPI_INT,
            /* data on target: */   target,0, 1,MPI_INT,
            the_window);
    MPI_Win_complete(the_window);
}

if (mytid==target) {
    MPI_Group_incl(all_group,1,&origin,&two_group);
    // exposure
    MPI_Win_post(two_group,0,the_window);
    MPI_Win_wait(the_window);
}
```

### 8.5.6 Passive target synchronization

*This reference section gives the syntax for routines introduced in section 4.3.6.*

```
MPI_Win_lock (int locktype, int rank, int assert, MPI_Win win)
MPI_Win_unlock (int rank, MPI_Win win)
```

The `MPI_Win_lock` call atomically retrieves an item from the window indicated, and replaces the item on the target by doing an accumulate on it with the data on the origin.

```
int MPI_Fetch_and_op(const void *origin_addr, void *result_addr,
                     MPI_Datatype datatype, int target_rank, MPI_Aint target_disp,
                     MPI_Op op, MPI_Win win)

// passive.cxx
```

```
if (mytid==repository) {
    // Processor zero creates a table of inputs
    // and associates that with the window
}
if (mytid!=repository) {
    float contribution=(float)mytid,table_element;
    int loc=0;
    MPI_Win_lock(MPI_LOCK_EXCLUSIVE,repository,0,the_window);
    // read the table element by getting the result from adding zero
    err = MPI_Fetch_and_op(&contribution,&table_element,MPI_FLOAT,
                           repository,loc,MPI_SUM,the_window); CHK(err);
    MPI_Win_unlock(repository,the_window);
}
```

## 8.6 Collectives

### 8.6.1 Rooted collectives

*This reference section gives the syntax for routines introduced in section 3.3.*

The MPI\_Bcast call has a single data argument. Its value on the root processor is copied to all other processors, where any previous value is overwritten.

```
C:
int MPI_Bcast(
    void* buffer, int count, MPI_Datatype datatype,
    int root, MPI_Comm comm)

Fortran:
MPI_Bcast(buffer, count, datatype, root, comm, ierror)
TYPE(*), DIMENSION(..) :: buffer
INTEGER, INTENT(IN) :: count, root
TYPE(MPI_Datatype), INTENT(IN) :: datatype
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python native:
rbuf = MPI.Comm.bcast(self, obj=None, int root=0)
Python numpy:
MPI.Comm.Bcast(self, buf, int root=0)
```

There is an example in section 8.1.1.

The MPI\_Reduce call combines the values from the individual processors. In order not to overwrite the input value on the root, this call has two data arguments, a send buffer and a receive buffer.

```
C:
int MPI_Reduce(
```

```

const void* sendbuf, void* recvbuf, int count, MPI_Datatype datatype,
MPI_Op op, int root, MPI_Comm comm)

Fortran:
MPI_Reduce(sendbuf, recvbuf, count, datatype, op, root, comm, ierror)
TYPE(*), DIMENSION(..), INTENT(IN) :: sendbuf
TYPE(*), DIMENSION(..) :: recvbuf
INTEGER, INTENT(IN) :: count, root
TYPE(MPI_Datatype), INTENT(IN) :: datatype
TYPE(MPI_Op), INTENT(IN) :: op
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python:
mpi4py.MPI.Comm.Reduce(self, sendbuf, recvbuf, Op op=SUM, int root=0)
mpi4py.MPI.Comm.reduce(self, sendobj=None, recvobj=None, op=SUM, int root=0)

```

On processes that are not the root, the receive buffer is ignored.

```

// reduce.c
float myrandom = (float) rand()/(float)RAND_MAX,
      result;
int target_proc = ntids-1;
// add all the random variables together
MPI_Reduce(&myrandom,&result,1,MPI_FLOAT,MPI_SUM,
           target_proc,comm);
// the result should be approx ntids/2:
if (mytid==target_proc)
    printf("Result %6.3f compared to ntids/2=%5.2f\n",
           result,ntids/2.);

```

On the root, you need two buffers, which could be a significant memory demand in the case of a large array to be reduced. Therefore, you can specify `MPI_IN_PLACE` as the send buffer on the root. The reduction call then uses the value in the receive buffer as the root's contribution to the operation.

```

// reduceinplace.c
float mynumber,result,*sendbuf,*recvbuf;
mynumber = (float) mytid;
int target_proc = ntids-1;
// add all the random variables together
if (mytid==target_proc) {
    sendbuf = (float*)MPI_IN_PLACE; recvbuf = &result;
    result = mynumber;
} else {
    sendbuf = &mynumber;     recvbuf = NULL;
}
MPI_Reduce(sendbuf,recvbuf,1,MPI_FLOAT,MPI_SUM,

```

```
        target_proc,comm);
// the result should be ntids*(ntids-1)/2:
if (mytid==target_proc)
    printf("Result %6.3f compared to n(n-1)/2=%5.2f\n",
           result,ntids*(ntids-1)/2.);
```

In Fortran the code is less elegant because you can not do these address calculations:

```
// reduceinplace.F90
call random_number(mynumber)
target_proc = ntids-1;
! add all the random variables together
if (mytid.eq.target_proc) then
    result = mytid
    call MPI_Reduce(MPI_IN_PLACE,result,1,MPI_REAL,MPI_SUM,&
                    target_proc,comm,err)
else
    mynumber = mytid
    call MPI_Reduce(mynumber,result,1,MPI_REAL,MPI_SUM,&
                    target_proc,comm,err)
end if
! the result should be ntids*(ntids-1)/2:
if (mytid.eq.target_proc) then
    write(*,'("Result ",f5.2," compared to n(n-1)/2=",f5.2)') &
           result,ntids*(ntids-1)/2.
end if
```

## 8.6.2 MPI Operators

The following is the list of predefined MPI\_OP values.

MPI_MAX	maximum	
MPI_MIN	minimum	
MPI_SUM	sum	
MPI_PROD	product	
MPI_LAND	logical and	
MPI_BAND	bitwise and	All except the last two operate on MPI datatypes;
MPI_LOR	logical or	
MPI_BOR	bitwise or	
MPI_LXOR	logical xor	
MPI_BXOR	bitwise xor	
MPI_MAXLOC	max value and location	
MPI_MINLOC	min value and location	
	the last two operate on a value/index pair.	

### 8.6.3 Gather and scatter

*This reference section gives the syntax for routines introduced in section 3.4.*

In the gather and scatter calls, each processor has  $n$  elements of individual data. There is also a root processor that has an array of length  $np$ , where  $p$  is the number of processors. The gather call collects all this data from the processors to the root; the scatter call assumes that the information is initially on the root and it is spread to the individual processors.

The prototype for `MPI_Gather` has two ‘count’ parameters, one for the length of the individual send buffers, and one for the receive buffer. However, confusingly, the second parameter (which is only relevant on the root) does not indicate the total amount of information coming in, but rather the size of *each* contribution. Thus, the two count parameters will usually be the same (at least on the root); they can differ if you use different `MPI_Datatype` values for the sending and receiving processors.

```
int MPI_Gather(
    void *sendbuf, int sendcnt, MPI_Datatype sendtype,
    void *recvbuf, int recvcnt, MPI_Datatype recvtype,
    int root, MPI_Comm comm
);
```

Here is a small example:

```
// gather.c
// we assume that each process has a value "localsize"
// the root process collects these values

if (mytid==root)
    localsizes = (int*) malloc( (ntids+1)*sizeof(int) );

// everyone contributes their info
MPI_Gather(&localsize,1,MPI_INT,
           localsizes,1,MPI_INT,root,comm);
```

This will also be the basis of a more elaborate example in section 8.6.6.

The `MPI_IN_PLACE` option can be used for the send buffer on the root; the data for the root is then assumed to be already in the correct location in the receive buffer.

The `MPI_Scatter` operation is in some sense the inverse of the gather: the root process has an array of length  $np$  where  $p$  is the number of processors and  $n$  the number of elements each processor will receive.

```
int MPI_Scatter
    (void* sendbuf, int sendcount, MPI_Datatype sendtype,
     void* recvbuf, int recvcount, MPI_Datatype recvtype,
     int root, MPI_Comm comm)
```

### 8.6.4 Reduce-scatter

*This reference section gives the syntax for routines introduced in section 3.8.*

The `MPI_Reduce_scatter` command is equivalent to a reduction on an array of data, followed by a scatter of that data to the individual processes.

To be precise, there is an array `recvcounts` where `recvcounts[i]` gives the number of elements that ultimate wind up on process `i`. The result is equivalent to doing a reduction with a length equal to the sum of the `recvcounts[i]` values, followed by a scatter where process `i` receives `recvcounts[i]` values. (Since the amount of data to be scattered depends on the process, this is in fact equivalent to `MPI_Scatterv` rather than a regular scatter.)

```
int MPI_Reduce_scatter
    (void* sendbuf, void* recvbuf, int *recvcounts, MPI_Datatype datatype,
     MPI_Op op, MPI_Comm comm)
```

For instance, if all `recvcounts[i]` values are 1, the sendbuffer has one element for each process, and the receive buffer has length 1.

An important application of this is establishing an irregular communication pattern. Assume that each process knows which other processes it wants to communicate with; the problem is to let the other processes know about this. The solution is to use `MPI_Reduce_scatter` to find out how many processes want to communicate with you, and then wait for precisely that many messages with a source value of `MPI_ANY_SOURCE`.

```
// reducescatter.c
// record what processes you will communicate with
int *recv_requests;
// find how many procs want to communicate with you
MPI_Reduce_scatter
    (recv_requests, &nsend_requests, counts, MPI_INT,
     MPI_SUM, comm);
// send a msg to the selected processes
for (int i=0; i<ntids; i++)
    if (recv_requests[i]>0)
        MPI_Isend(&msg, 1, MPI_INT, /*to:*/ i, 0, comm,
                  mpi_requests+irequest++);
// do as many receives as you know are coming in
for (int i=0; i<nsend_requests; i++)
    MPI_Irecv(&msg, 1, MPI_INT, MPI_ANY_SOURCE, MPI_ANY_TAG, comm,
              mpi_requests+irequest++);
MPI_Waitall(irequest, mpi_requests, MPI_STATUSES_IGNORE);
```

### 8.6.5 ‘All’-type collectives

*This reference section gives the syntax for routines introduced in section 3.9.*

The following collectives construct a result on all processes:

```
C:  
int MPI_Allreduce(const void* sendbuf,  
                  void* recvbuf, int count, MPI_Datatype datatype,  
                  MPI_Op op, MPI_Comm comm)  
  
Semantics:  
IN sendbuf: starting address of send buffer (choice)  
OUT recvbuf: starting address of receive buffer (choice)  
IN count: number of elements in send buffer (non-negative integer)  
IN datatype: data type of elements of send buffer (handle)  
IN op: operation (handle)  
IN comm: communicator (handle)  
  
Fortran:  
MPI_Allreduce(sendbuf, recvbuf, count, datatype, op, comm, ierror)  
TYPE(*), DIMENSION(..), INTENT(IN) :: sendbuf  
TYPE(*), DIMENSION(..) :: recvbuf  
INTEGER, INTENT(IN) :: count  
TYPE(MPI_Datatype), INTENT(IN) :: datatype  
TYPE(MPI_Op), INTENT(IN) :: op  
TYPE(MPI_Comm), INTENT(IN) :: comm  
INTEGER, OPTIONAL, INTENT(OUT) :: ierror  
  
Python native:  
recvobj = MPI.Comm.allreduce(self, sendobj, op=SUM)  
Python numpy:  
MPI.Comm.Allreduce(self, sendbuf, recvbuf, Op op=SUM)
```

*Python note* The receive buffer has to be of the same size as the send buffer.

MPI\_Alltoall

```
int MPI_Alltoall  
(void *sendbuf, int sendcount, MPI_Datatype sendtype,  
 void *recvbuf, int recvcount, MPI_Datatype recvtype,  
 MPI_Comm comm)
```

Each processor has a contribution in their send buffer; the global result is returned in each processor's receive buffer.

If a large amount of data is being communicated, it may be wasteful to have both a (large) send and receive buffer. This problem can be circumvented by using `MPI_IN_PLACE` as the specification of the send buffer. The send data is then assumed to be in the receive buffer. After the reduction it is, of course, overwritten.

```
// allreduceinplace.c  
int nrandoms = 500000;  
float *myrandoms;
```

```
myrandoms = (float*) malloc(nrandoms*sizeof(float));
for (int irand=0; irand<nrandoms; irand++)
    myrandoms[irand] = (float) rand()/(float)RAND_MAX;
// add all the random variables together
MPI_Allreduce(MPI_IN_PLACE,myrandoms,
nrandoms,MPI_FLOAT,MPI_SUM,comm);
// the result should be approx ntids/2:
if (mytid==ntids-1) {
    float sum=0.;
    for (int i=0; i<nrandoms; i++) sum += myrandoms[i];
    sum /= nrandoms*ntids;
    printf("Result %6.9f compared to .5\n",sum);
}
```

### 8.6.6 Variable-size-input collectives

*This reference section gives the syntax for routines introduced in section 3.5.*

There are various calls where processors can have buffers of differing sizes.

- In `MPI_Scatterv` the root process has a different amount of data for each recipient.
- In `MPI_Gatherv`, conversely, each process contributes a different sized send buffer to the received result; `MPI_Allgatherv` does the same, but leaves its result on all processes; `MPI_Alltoallv` does a different variable-sized gather on each process.

```
int MPI_Scatterv
    (void* sendbuf, int *sendcounts, int *displs, MPI_Datatype sendtype,
     void* recvbuf, int recvcount, MPI_Datatype recvtype,
     int root, MPI_Comm comm)

C:
int MPI_Gatherv(
    const void* sendbuf, int sendcount, MPI_Datatype sendtype,
    void* recvbuf, const int recvcounts[], const int displs[],
    MPI_Datatype recvtype, int root, MPI_Comm comm)

Semantics:
IN sendbuf: starting address of send buffer (choice)
IN sendcount: number of elements in send buffer (non-negative integer)
IN sendtype: data type of send buffer elements (handle)
OUT recvbuf: address of receive buffer (choice, significant only at root)
IN recvcounts: non-negative integer array (of length group size) containing the number of elements to receive from each process
IN displs: integer array (of length group size). Entry i specifies the displacement of the start of the data from process i to the receive buffer
IN recvtype: data type of recv buffer elements (significant only at root) (handle)
IN root: rank of receiving process (integer)
IN comm: communicator (handle)
```

```

Fortran:
MPI_Gatherv(sendbuf, sendcount, sendtype, recvbuf, recvcounts, displs, recvtype, r
TYPE(*), DIMENSION(..), INTENT(IN) :: sendbuf
TYPE(*), DIMENSION(..) :: recvbuf
INTEGER, INTENT(IN) :: sendcount, recvcounts(*), displs(*), root
TYPE(MPI_Datatype), INTENT(IN) :: sendtype, recvtype
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python:
Gatherv(self, sendbuf, [recvbuf,counts], int root=0)

int MPI_Allgatherv
    (void *sendbuf, int sendcount, MPI_Datatype sendtype,
     void *recvbuf, int *recvcounts, int *displs,
     MPI_Datatype recvtype, MPI_Comm comm)

MPI_Alltoallv.
int MPI_Alltoallv
    (void *sendbuf, int *sendcnts, int *sdispls, MPI_Datatype sendtype,
     void *recvbuf, int *recvcnts, int *rdispls, MPI_Datatype recvtype,
     MPI_Comm comm)

```

For example, in an `MPI_Gatherv` call each process has an individual number of items to contribute. To gather this, the root process needs to find these individual amounts with an `MPI_Gather` call, and locally construct the offsets array. Note how the offsets array has size `ntids+1`: the final offset value is automatically the total size of all incoming data.

```

// gatherv.c
// we assume that each process has an array "localdata"
// of size "localsize"

// the root process decides how much data will be coming:
// allocate arrays to contain size and offset information
if (mytid==root) {
    localsizes = (int*) malloc( (ntids+1)*sizeof(int) );
    offsets = (int*) malloc( ntids*sizeof(int) );
}
// everyone contributes their info
MPI_Gather(&localsize,1,MPI_INT,
           localsizes,1,MPI_INT,root,comm);
// the root constructs the offsets array
if (mytid==root) {
    offsets[0] = 0;
    for (int i=0; i<ntids; i++)

```

```
    offsets[i+1] = offsets[i]+localsizes[i];
    alldata = (int*) malloc( offsets[ntids]*sizeof(int) );
}
// everyone contributes their data
MPI_Gatherv(localdata,localsize,MPI_INT,
            alldata,localsizes,offsets,MPI_INT,root,comm);
```

### 8.6.7 Scan

*This reference section gives the syntax for routines introduced in section 3.6.*

MPI has two routines for scan, or *prefix*, operations: the inclusive scan

C:

```
int MPI_Scan(const void* sendbuf, void* recvbuf,
              int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
IN sendbuf: starting address of send buffer (choice)
OUT recvbuf: starting address of receive buffer (choice)
IN count: number of elements in input buffer (non-negative integer)
IN datatype: data type of elements of input buffer (handle)
IN op: operation (handle)
IN comm: communicator (handle)
```

Fortran:

```
MPI_Scan(sendbuf, recvbuf, count, datatype, op, comm, ierror)
TYPE(*), DIMENSION(..), INTENT(IN) :: sendbuf
TYPE(*), DIMENSION(..) :: recvbuf
INTEGER, INTENT(IN) :: count
TYPE(MPI_Datatype), INTENT(IN) :: datatype
TYPE(MPI_Op), INTENT(IN) :: op
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Python:

```
res = Intracomm.scan( sendobj=None, recvobj=None, op=MPI.SUM)
res = Intracomm.exscan( sendobj=None, recvobj=None, op=MPI.SUM)
```

and the exclusive scan:

The MPI\_Op operations do not return an error code.

The result of the exclusive scan is undefined on processor 0 (None in python), and on processor 1 it is a copy of the send value of processor 1. In particular, the MPI\_Op need not be called on these two processors.

Scan operations are often useful in index calculations. Suppose that every processor has part of a long array, and it knows only how many element it has. The following bit computes the global index of its first element.

```
// exscan.c
```

```
int my_first=0,localsize;
// localsize = ..... result of local computation ....
// find myfirst location based on the local sizes
err = MPI_Exscan(&localsize,&my_first,
                 1,MPI_INT,MPI_SUM,comm); CHK(err);
```

## 8.6.8 User-defined reductions

*This reference section gives the syntax for routines introduced in section 3.7.*

```
MPI_Op_create( MPI_User_function *func, int commute, MPI_Op *op);
```

## 8.6.9 Non-blocking collectives

*This reference section gives the syntax for routines introduced in section 3.10.*

The same calling sequence as the blocking counterpart, except for the addition of an `MPI_Request` parameter. For instance `MPI_Ibcast`:

```
int MPI_Ibcast(
    void *buffer, int count, MPI_Datatype datatype,
    int root, MPI_Comm comm,
    MPI_Request *request)
```

## 8.7 Communicators

### 8.7.1 Communicator duplication

*This reference section gives the syntax for routines introduced in section 6.2.2.1.*

In section 7.3.9 it was explained that MPI messages are non-overtaking. This may lead to confusing situations, witness the following snippet:

This models a main program that does a simple message exchange, and it makes two calls to library routines. Unbeknown to the user, the library also issues send and receive calls, and they turn out to interfere:

```
// commdup_wrong.cxx
class library {
private:
    MPI_Comm comm;
    int mytid,ntids,other;
    MPI_Request *request;
public:
    library(MPI_Comm incomm) {
        comm = incomm;
```

```

MPI_Comm_rank(comm, &mytid);
other = 1-mytid;
request = new MPI_Request[2];
};

int communication_start();
int communication_end();
};

```

Here

- The main program does a send,
- the library call `function_start` does a send and a receive; because the receive can match either send, it is paired with the first one;
- the main program does a receive, which will be paired with the send of the library call;
- both the main program and the library do a wait call, and in both cases all requests are successfully fulfilled, just not the way you intended.

The solution is to give the library a separate communicator with `MPI_Comm_dup`.

```
int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)
```

Newly created communicators should be released again with `MPI_Comm_free`.

```

// commdup_right.F90
class library {
private:
    MPI_Comm comm;
    int mytid, ntids, other;
    MPI_Request *request;
public:
    library(MPI_Comm incomm) {
        MPI_Comm_dup(incomm, &comm);
        MPI_Comm_rank(comm, &mytid);
        other = 1-mytid;
        request = new MPI_Request[2];
    };
    ~library() {
        MPI_Comm_free(&comm);
    }
    int communication_start();
    int communication_end();
};

```

### 8.7.2 Splitting communicators

*This reference section gives the syntax for routines introduced in section 6.2.2.2.*

The command `MPI_Comm_split` takes a communicator, and divides it into a number of disjoint communicators. It does this by assigning processes to the same subcommunicator if they have the same user-specified ‘colour’ value.

```
int MPI_Comm_split(MPI_Comm comm, int color, int key,
                    MPI_Comm *newcomm)
```

The ranking of processes in the new communicator is determined by a ‘key’ value. Most of the time, there is no reason to use a relative ranking that is different from the global ranking, so the `MPI_Comm_rank` value of the global communicator is a good choice.

```
// mvp2d.cxx
row_number = ntids % ntids_i;
col_number = ntids / ntids_j;
MPI_Comm_split(global_comm, row_number, mytid, &row_comm);
MPI_Comm_split(global_comm, col_number, mytid, &col_comm);
```

There are some predefined colours, named ‘types’, to use in communicator splitting. The routine `MPI_Comm_split_type` looks very much like `MPI_Comm_split`:

```
C:
int MPI_Comm_split_type(
    MPI_Comm comm, int split_type, int key,
    MPI_Info info, MPI_Comm *newcomm)

Fortran:
MPI_Comm_split_type(comm, split_type, key, info, newcomm, ierror)
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, INTENT(IN) :: split_type, key
TYPE(MPI_Info), INTENT(IN) :: info
TYPE(MPI_Comm), INTENT(OUT) :: newcomm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror

Python:
MPI.Comm.Split_type(
    self, int split_type, int key=0, Info info=INFO_NULL)
```

but the `split_type` parameters has to be from the following (short) list:

- `MPI_COMM_TYPE_SHARED`: split the communicator into subcommunicators of processes sharing a memory area.

### 8.7.3 Process topologies

*This reference section gives the syntax for routines introduced in section 6.2.5.*

#### 8.7.3.1 Cartesian grid topology

*This reference section gives the syntax for routines introduced in section 6.2.5.1.*

The cartesian topology is specified by giving `MPI_Cart_create` the sizes of the processor grid along each axis, and whether the grid is periodic along that axis.

```
int MPI_Cart_create(
    MPI_Comm comm_old, int ndims, int *dims, int *periods,
    int reorder, MPI_Comm *comm_cart)
```

Each point in this new communicator has a coordinate and a rank. They can be queried with `MPI_Cart_coord` and `MPI_Cart_rank` respectively.

```
int MPI_Cart_coords(
    MPI_Comm comm, int rank, int maxdims,
    int *coords);
int MPI_Cart_rank(
    MPI_Comm comm, int *coords,
    int *rank);
```

Note that these routines can give the coordinates for any rank, not just for the current process.

```
// cart.c
MPI_Comm comm2d;
ndim = 2; periodic[0] = periodic[1] = 0;
dimensions[0] = idim; dimensions[1] = jdim;
MPI_Cart_create(comm,ndim,dimensions,periodic,1,&comm2d);
MPI_Cart_coords(comm2d,mytid,ndim,coord_2d);
MPI_Cart_rank(comm2d,coord_2d,&rank_2d);
printf("I am %d: (%d,%d); originally %d\n",rank_2d,coord_2d[0],coord_2d[1],
```

The `reorder` parameter to `MPI_Cart_create` indicates whether processes can have a rank in the new communicator that is different from in the old one.

Strangely enough you can only shift in one direction, you can not specify a shift vector.

```
int MPI_Cart_shift(MPI_Comm comm, int direction, int displ, int *source,
                   int *dest)
```

If you specify a processor outside the grid the result is `MPI_PROC_NULL`.

### 8.8 Leftover topics

#### 8.8.1 32-bit size issues

The `size` parameter in MPI routines is defined as an `int`, meaning that it is limited to 32-bit quantities. There are ways around this, such as sending a number of `MPI_Type_contiguous` blocks that add up to more than  $2^{31}$ .

## 8.8.2 Fortran issues

This reference section gives the syntax for routines introduced in section 7.3.3.

### 8.8.2.1 Data types

The equivalent of MPI\_Aint in Fortran is

```
integer(kind=MPI_ADDRESS_KIND) :: winsize
```

### 8.8.2.2 Type issues

Fortran90 is a strongly typed language, so it is not possible to pass argument by reference to their address, as C/C++ do with the void\* type for send and receive buffers. In Fortran this is solved by having separate routines for each datatype, and providing an Interface block in the MPI module. If you manage to request a version that does not exist, the compiler will display a message like

```
There is no matching specific
subroutine for this generic subroutine call [MPI_Send]
```

### 8.8.2.3 Byte calculations

Fortran lacks a sizeof operator to query the sizes of datatypes. Since sometimes exact byte counts are necessary, for instance in one-sided communication, Fortran can use the MPI\_Sizeof routine, for instance for MPI\_Win\_create:

```
call MPI_Sizeof(windowdata,window_element_size,ierr)
window_size = window_element_size*500
call MPI_Win_create( windowdata,window_size,window_element_size,... );
```

## 8.8.3 Python issues

### 8.8.3.1 Byte calculations

The MPI\_Win\_create routine needs a displacement in bytes. Here is a good way for finding the size of numpy datatypes:

```
numpy.dtype('i').itemsize
```

### 8.8.3.2 Arrays of objects

Objects of type MPI\_Status or MPI\_Request often need to be created in an array, for instance when looping through a number of Isend calls. In that case the following idiom may come in handy:

```
requests = [ None ] * nprocs
for p in range(nprocs):
    requests[p] = comm.Irecv( ... )
```

### 8.8.4 Cancelling messages

In section ?? we showed a master-worker example where the master accepts in arbitrary order the messages from the workers. Here we will show a slightly more complicated example, where only the result of the first task to complete is needed. Thus, we issue an `MPI_Recv` with `MPI_ANY_SOURCE` as source. When a result comes, we broadcast its source to all processes. All the other workers then use this information to cancel their message with an `MPI_Cancel` operation.

```
// cancel.c
if (mytid==ntids-1) {
    MPI_Status status;
    ierr = MPI_Recv(dummy,0,MPI_INT, MPI_ANY_SOURCE,0,comm,
                    &status); CHK(ierr);
    first_tid = status.MPI_SOURCE;
    ierr = MPI_Bcast(&first_tid,1,MPI_INT, ntids-1,comm); CHK(ierr);
    printf("first msg came from %d\n",first_tid);
} else {
    float randomfraction = (rand() / (double)RAND_MAX);
    int randomwait = (int) ( ntids * randomfraction );
    MPI_Request request;
    printf("process %d waits for %e/%d=%d\n",
           mytid,randomfraction,ntids,randomwait);
    sleep(randomwait);
    ierr = MPI_Isend(dummy,0,MPI_INT, ntids-1,0,comm,
                     &request); CHK(ierr);
    ierr = MPI_Bcast(&first_tid,1,MPI_INT, ntids-1,comm
                     ); CHK(ierr);
    if (mytid!=first_tid) {
        ierr = MPI_Cancel(&request); CHK(ierr);
    }
}
```

### 8.8.5 Constants

MPI constants such as `MPI_COMM_WORLD` or `MPI_INT` are not necessarily statitally defined, such as by a `#define` statement: the best you can say is that they have a value after `MPI_Init` or `MPI_Init_thread`. That means you can not transfer a compiled MPI file between platforms, or even between compilers on one platform. However, a working MPI source on one MPI implementation will also work on another.

## 8.9 Error handling

*This reference section gives the syntax for routines introduced in section 7.3.2.*

MPI operators (`MPI_Op`) do not return an error code. In case of an error they call `MPI_Abort`; if `MPI_ERRORS_RETURN` is the error handler, errors may be silently ignore.

## 8.10 More utility stuff

### 8.10.1 Context information

*This reference section gives the syntax for routines introduced in section 7.3.5.*

You can query the *hostname* of a processor. This name need not be unique between different processor ranks.

C:

```
int MPI_Get_processor_name(char *name, int *resultlen)
```

Fortran:

```
MPI_Get_processor_name(name, resultlen, ierror)
CHARACTER(LEN=MPI_MAX_PROCESSOR_NAME), INTENT(OUT) :: name
INTEGER, INTENT(OUT) :: resultlen
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

Python:

```
MPI.Get_processor_name()
```

Note that you have to pass in the character storage: the character array must be at least `MPI_MAX_PROCESSOR_NAME` characters long. The actual length of the name is returned in the `resultlen` parameter.

In C and C++,

```
#define MPI_VERSION 2
#define MPI_SUBVERSION 2
```

in Fortran,

```
INTEGER MPI_VERSION, MPI_SUBVERSION
PARAMETER (MPI_VERSION = 2)
PARAMETER (MPI_SUBVERSION = 2)
```

For runtime determination,

```
MPI_GET_VERSION( version, subversion )
OUT version version number (integer)
OUT subversion subversion number (integer)

int MPI_Get_version(int *version, int *subversion)
MPI_GET_VERSION(VERSION, SUBVERSION, IERROR)
INTEGER VERSION, SUBVERSION, IERROR
```

### 8.10.2 Timing

*This reference section gives the syntax for routines introduced in section 7.3.6.*

MPI has a *wall clock* timer: `MPI_Wtime`

```
// C
double MPI_Wtime(void);
! F
DOUBLE PRECISION MPI_WTIME()
```

which gives the number of seconds from a certain point in the past. (Note the absence of the error parameter in the fortran call.)

```
// pingpong.c
int src = 0, tgt = ntids/2;
double t, send=1.1, recv;
if (mytid==src) {
    t = MPI_Wtime();
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Send(&send, 1, MPI_DOUBLE, tgt, 0, comm);
        MPI_Recv(&recv, 1, MPI_DOUBLE, tgt, 0, comm, MPI_STATUS_IGNORE);
    }
    t = MPI_Wtime() - t; t /= NEXPERIMENTS;
    printf("Time for pingpong: %e\n", t);
} else if (mytid==tgt) {
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Recv(&recv, 1, MPI_DOUBLE, src, 0, comm, MPI_STATUS_IGNORE);
        MPI_Send(&recv, 1, MPI_DOUBLE, src, 0, comm);
    }
}
```

The timer has a resolution of `MPI_Wtick`:

```
double MPI_Wtick(void);
```

Timing in parallel is a tricky issue. For instance, most clusters do not have a central clock, so you can not relate start and stop times on one process to those on another. You can test for a global clock as follows :

```
int *v, flag;
MPI_Attr_get( comm, MPI_WTIME_IS_GLOBAL, &v, &flag );
if (mytid==0) printf(``Time synchronized? %d->%d\n'', flag, *v);
```

## 8.11 Multi-threading

*This reference section gives the syntax for routines introduced in section 7.2.*

Hybrid MPI/threaded codes need to replace `MPI_Init` by `MPI_Init_thread`:

```
C:
int MPI_Init_thread(int *argc, char ***argv, int required, int *provided)

Fortran:
MPI_Init_thread(required, provided, ierror)
INTEGER, INTENT(IN) :: required
INTEGER, INTENT(OUT) :: provided
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

With the `required` parameter the user requests a certain level of support, and MPI reports the actual capabilities in the `provided` parameter.

The following constants are defined:

- `MPI_THREAD_SINGLE`: each MPI process can only have a single thread.
- `MPI_THREAD_FUNNELED`: an MPI process can be multithreaded, but all MPI calls need to be done from a single thread.
- `MPI_THREAD_SERIALIZED`: a processes can sustain multiple threads that make MPI calls, but these threads can not be simultaneous: they need to be for instance in an OpenMP *critical section*.
- `MPI_THREAD_MULTIPLE`: processes can be fully generally multi-threaded.

These values are monotonically increasing.

After the initialization call, you can query the support level with `MPI_Query_thread`:

```
C:
int MPI_Query_thread(int *provided)

Fortran:
MPI_Query_thread(provided, ierror)
INTEGER, INTENT(OUT) :: provided
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

In case more than one thread performs communication, the following routine can determine whether a thread is the main thread:

```
C:
int MPI_Is_thread_main(int *flag)

Fortran:
MPI_Is_thread_main(flag, ierror)
LOGICAL, INTENT(OUT) :: flag
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

## Chapter 9

### MPI Examples

#### 9.1 A

`MPI_Allgatherv`

Prior to the actual `gather` call, we need to construct the count and displacement arrays. The easiest way is to use a reduction.

```
// allgatherv.c
MPI_Allgather
( &my_count, 1, MPI_INT,
  recv_counts, 1, MPI_INT, comm );
int accumulate = 0;
for (int i=0; i<ntrids; i++) {
  recv_displs[i] = accumulate; accumulate += recv_counts[i];
int *global_array = (int*) malloc(accumulate*sizeof(int));
MPI_Allgatherv
( my_array, mytid+1, MPI_INT,
  global_array, recv_counts, recv_displs, MPI_INT, comm );
```

In python the receive buffer has to contain the counts and displacements arrays.

```
// allgatherv.py
my_count = np.empty(1, dtype=np.int)
my_count[0] = mycount
comm.Allgather( my_count, recv_counts )

accumulate = 0
for p in range(nprocs):
  recv_displs[p] = accumulate; accumulate += recv_counts[p]
global_array = np.empty(accumulate, dtype=np.float64)
comm.Allgatherv( my_array, [global_array, recv_counts, recv_displs, MPI.DOUBLE]
```

```
MPI_Allreduce
```

We give each process a random number, and sum these numbers together. The result should be approximate 1/2 times the number of processes.

```
// allreduce.c
float myrandom,sumrandom;
myrandom = (float) rand()/(float)RAND_MAX;
// add the random variables together
MPI_Allreduce(&myrandom,&sumrandom,
1,MPI_FLOAT,MPI_SUM,comm);
// the result should be approx ntids/2:
if (mytid==ntids-1)
printf("Result %6.9f compared to .5\n",sumrandom/ntids);
```

Using the MPI\_IN\_PLACE specifier:

```
// allreduceinplace.c
int nrandoms = 500000;
float *myrandoms;
myrandoms = (float*) malloc(nrandoms*sizeof(float));
for (int irand=0; irand<nrandoms; irand++)
    myrandoms[irand] = (float) rand()/(float)RAND_MAX;
// add all the random variables together
MPI_Allreduce(MPI_IN_PLACE,myrandoms,
nrandoms,MPI_FLOAT,MPI_SUM,comm);
// the result should be approx ntids/2:
if (mytid==ntids-1) {
    float sum=0.;
    for (int i=0; i<nrandoms; i++) sum += myrandoms[i];
    sum /= nrandoms*ntids;
    printf("Result %6.9f compared to .5\n",sum);
}
```

For Python we illustrate both the native and the numpy variant. In the numpy variant we create an array for the receive buffer, even though only one element is used.

```
// allreduce.py
random_number = random.randint(1,nprocs*nprocs)
print "[%d] random=%d" % (procid,random_number)

max_random = comm.allreduce(random_number,op=MPI.MAX)
if procid==0:
    print "Python native:\n max=%d" % max_random

myrandom = np.empty(1,dtype=np.int)
```

## 9. MPI Examples

---

```
myrandom[0] = random_number
allrandom = np.empty(nprocs, dtype=np.int)

comm.Allreduce(myrandom, allrandom[:1], op=MPI.MAX)
```

### 9.2 B

#### MPI\_Bcast

In python we illustrate the native and numpy variants. In the native variant the result is given as a function return; in the numpy variant the send buffer is reused.

```
// bcast.py
# first native
if procid==root:
    buffer = [ 5.0 ] * dsize
buffer = comm.bcast(obj=buffer, root=root)
if not reduce( lambda x,y:x and y, [ buffer[i]==5.0 for i in range(len(buffer)) ] ):
    print "Something wrong on proc %d: native buffer <<%s>>" % (procid,str(buffer))

# then with NumPy
buffer = np.arange(dsize, dtype=np.float64)
if procid==root:
    for i in range(dsize):
        buffer[i] = 5.0
comm.Bcast( buffer,root=root )
if not all( buffer==5.0 ):
    print "Something wrong on proc %d: numpy buffer <<%s>>" % (procid,str(buffer))
```

### 9.3 C

#### MPI\_Cancel

Cancelling a send operation:

```
// cancel.c
if (mytid==ntids-1) {
    MPI_Status status;
    ierr = MPI_Recv(dummy,0,MPI_INT, MPI_ANY_SOURCE,0,comm,
                    &status); CHK(ierr);
    first_tid = status.MPI_SOURCE;
    ierr = MPI_Bcast(&first_tid,1,MPI_INT, ntids-1,comm); CHK(ierr);
```

```

        printf("first msg came from %d\n", first_tid);
    } else {
        float randomfraction = (rand() / (double)RAND_MAX);
        int randomwait = (int) ( ntids * randomfraction );
        MPI_Request request;
        printf("process %d waits for %e/%d=%d\n",
               mytid, randomfraction, ntids, randomwait);
        sleep(randomwait);
        ierr = MPI_Isend(dummy, 0, MPI_INT, ntids-1, 0, comm,
                         &request); CHK(ierr);
        ierr = MPI_Bcast(&first_tid, 1, MPI_INT, ntids-1, comm
                         ); CHK(ierr);
        if (mytid!=first_tid) {
            ierr = MPI_Cancel(&request); CHK(ierr);
        }
    }

MPI_Cart...
// cart.c
MPI_Comm comm2d;
ndim = 2; periodic[0] = periodic[1] = 0;
dimensions[0] = idim; dimensions[1] = jdim;
MPI_Cart_create(comm, ndim, dimensions, periodic, 1, &comm2d);
MPI_Cart_coords(comm2d, mytid, ndim, coord_2d);
MPI_Cart_rank(comm2d, coord_2d, &rank_2d);
printf("I am %d: (%d,%d); originally %d\n", rank_2d, coord_2d[0], coord_2d[1],

char mychar = 65+mytid;
MPI_Cart_shift(comm2d, 0, +1, &rank_2d, &rank_right);
MPI_Cart_shift(comm2d, 0, -1, &rank_2d, &rank_left);
MPI_Cart_shift(comm2d, 1, +1, &rank_2d, &rank_up);
MPI_Cart_shift(comm2d, 1, -1, &rank_2d, &rank_down);
int irequest = 0; MPI_Request *requests = malloc(8*sizeof(MPI_Request));
MPI_Isend(&mychar, 1, MPI_CHAR, rank_right, 0, comm, requests+irequest++);
MPI_Isend(&mychar, 1, MPI_CHAR, rank_left, 0, comm, requests+irequest++);
MPI_Isend(&mychar, 1, MPI_CHAR, rank_up, 0, comm, requests+irequest++);
MPI_Isend(&mychar, 1, MPI_CHAR, rank_down, 0, comm, requests+irequest++);
MPI_Irecv( indata+idata++, 1, MPI_CHAR, rank_right, 0, comm, requests+ireques
MPI_Irecv( indata+idata++, 1, MPI_CHAR, rank_left, 0, comm, requests+ireques
MPI_Irecv( indata+idata++, 1, MPI_CHAR, rank_up, 0, comm, requests+ireques
MPI_Irecv( indata+idata++, 1, MPI_CHAR, rank_down, 0, comm, requests+ireques

MPI_Comm_dup

```

## 9. MPI Examples

---

Giving a library its own communicator.

```
// commdup_right.cxx
class library {
private:
    MPI_Comm comm;
    int mytid, ntid, other;
    MPI_Request *request;
public:
    library(MPI_Comm incomm) {
        MPI_Comm_dup(incomm, &comm);
        MPI_Comm_rank(comm, &mytid);
        other = 1-mytid;
        request = new MPI_Request[2];
    };
    ~library() {
        MPI_Comm_free(&comm);
    }
    int communication_start();
    int communication_end();
};

library my_library(comm);
MPI_Isend(&sdata, 1, MPI_INT, other, 1, comm, &(request[0]));
my_library.communication_start();
MPI_Irecv(&rdata, 1, MPI_INT, other, MPI_ANY_TAG,
          comm, &(request[1]));
MPI_Waitall(2, request, status);
my_library.communication_end();
```

## 9.4 E

MPI\_Exscan

Exclusive scan:

```
// exscan.c
int my_first=0, localsize;
// localsize = ..... result of local computation .....
// find myfirst location based on the local sizes
err = MPI_Exscan(&localsize, &my_first,
                  1, MPI_INT, MPI_SUM, comm); CHK(err);
```

```
// exscan.py
localsize = 10+random.randint(1,nprocs)
myfirst = 0
mypartial = comm.exscan(localsize,0)
```

## 9.5 F

### MPI\_Fetch\_and\_op

A root process has a table of data; the other processes do atomic gets and update of that data using *passive target synchronization* through MPI\_Win\_lock.

```
// passive.cxx
if (mytid==repository) {
    // Processor zero creates a table of inputs
    // and associates that with the window
}
if (mytid!=repository) {
    float contribution=(float)mytid,table_element;
    int loc=0;
    MPI_Win_lock(MPI_LOCK_EXCLUSIVE,repository,0,the_window);
    // read the table element by getting the result from adding zero
    err = MPI_Fetch_and_op(&contribution,&table_element,MPI_FLOAT,
                           repository,loc,MPI_SUM,the_window); CHK(err);
    MPI_Win_unlock(repository,the_window);
}
```

## 9.6 G

### MPI\_Gather

Gather data onto a root. Only the root allocates the gather buffer.

```
// gather.c
// we assume that each process has a value "localsize"
// the root process collects these values

if (mytid==root)
    localsizes = (int*) malloc( (ntids+1)*sizeof(int) );

// everyone contributes their info
MPI_Gather(&localsize,1,MPI_INT,
```

## 9. MPI Examples

---

```
    localsizes,1,MPI_INT,root,comm);
```

MPI\_Gatherv

Gather irregularly sized data onto a root. We first need an MPI\_Gather to determine offsets.

```
// gatherv.c
// we assume that each process has an array "localdata"
// of size "localsize"

// the root process decides how much data will be coming:
// allocate arrays to contain size and offset information
if (mytid==root) {
    localsizes = (int*) malloc( (ntids+1)*sizeof(int) );
    offsets = (int*) malloc( ntids*sizeof(int) );
}
// everyone contributes their info
MPI_Gather(&localsize,1,MPI_INT,
           localsizes,1,MPI_INT,root,comm);
// the root constructs the offsets array
if (mytid==root) {
    offsets[0] = 0;
    for (int i=0; i<ntids; i++)
        offsets[i+1] = offsets[i]+localsizes[i];
    alldata = (int*) malloc( offsets[ntids]*sizeof(int) );
}
// everyone contributes their data
MPI_Gatherv(localdata,localsize,MPI_INT,
            alldata,localsizes,offsets,MPI_INT,root,comm);

// gatherv.py
# implicitly using root=0
globalsize = comm.reduce(localsize)
if procid==0:
    print "Global size=%d" % globalsize
collecteddata = np.empty(globalsize,dtype=np.int)
counts = comm.gather(localsize)
comm.Gatherv(localdata, [collecteddata, counts])
```

MPI\_Get

One process does a one-sided get from another. This also illustrates setting size parameters in MPI\_Win\_create. Synchronization is done with MPI\_Win\_fence.

```
// getfence.c
```

---

```

MPI_Win_create(&other_number,2*sizeof(int),sizeof(int),
               MPI_INFO_NULL,comm,&the_window);
MPI_Win_fence(0,the_window);
if (mytid==0) {
    MPI_Get( /* data on origin: */ &my_number, 1,MPI_INT,
            /* data on target: */ other,1,      1,MPI_INT,
            the_window);
}
MPI_Win_fence(0,the_window);

```

We make a null window on processes that do not participate.

```

// getfence.py
if procid==0 or procid==nprocs-1:
    win_mem = np.empty( 1,dtype=np.float64 )
    win = MPI.Win.Create( win_mem,comm=comm )
else:
    win = MPI.Win.Create( None,comm=comm )

# put data on another process
win.Fence()
if procid==0 or procid==nprocs-1:
    putdata = np.empty( 1,dtype=np.float64 )
    putdata[0] = mydata
    print "[%d] putting %e" % (procid,mydata)
    win.Put( putdata,other )
win.Fence()

```

## 9.7 I

`MPI_Init_thread`

The `Init_thread` call takes the requested level of thread support and reports back what the provided level is.

```

// thread.c
MPI_Init_thread(&argc,&argv,MPI_THREAD_MULTIPLE,&threading);
comm = MPI_COMM_WORLD;
MPI_Comm_rank(comm,&mytid);
MPI_Comm_size(comm,&ntids);

if (mytid==0) {
    switch (threading) {

```

```

        case MPI_THREAD_MULTIPLE : printf("Glorious multithreaded MPI\n"); break;
        case MPI_THREAD_SERIALIZED : printf("No simultaneous MPI from threads\n");
        case MPI_THREAD_FUNNELED : printf("MPI from main thread\n"); break;
        case MPI_THREAD_SINGLE : printf("no threading supported\n"); break;
    }
}
MPI_Finalize();

```

## 9.8 P

**MPI\_Pack** Use packing to make a self-documenting message: the first element is an integer describing how many doubles follow. The **MPI\_Unpack** call inspects the integer, then calls unpack on the integers the appropriate number of times.

```

// pack.c
if (mytid==sender) {
    MPI_Pack(&nstarts,1,MPI_INT,buffer,buflen,&position,comm);
    for (int i=0; i<nstarts; i++) {
        double value = rand()/(double)RAND_MAX;
        MPI_Pack(&value,1,MPI_DOUBLE,buffer,buflen,&position,comm);
    }
    MPI_Pack(&nstarts,1,MPI_INT,buffer,buflen,&position,comm);
    MPI_Send(buffer,position,MPI_PACKED,other,0,comm);
} else if (mytid==receiver) {
    int irecv_value;
    double xrecv_value;
    MPI_Recv(buffer,buflen,MPI_PACKED,other,0,comm,MPI_STATUS_IGNORE);
    MPI_Unpack(buffer,buflen,&position,&nstarts,1,MPI_INT,comm);
    for (int i=0; i<nstarts; i++) {
        MPI_Unpack(buffer,buflen,&position,&xrecv_value,1,MPI_DOUBLE,comm);
    }
    MPI_Unpack(buffer,buflen,&position,&irecv_value,1,MPI_INT,comm);
    ASSERT(irecv_value==nstarts);
}

```

### MPI\_Put

A one-sided **MPI\_Put** with active target synchronization through the use of fences. This is more or less the same as the **MPI\_Get** example above.

```

// putfence.c
MPI_Win the_window;
MPI_Win_create(&window_data,2*sizeof(int),sizeof(int),

```

```
    MPI_Info_free(&info);
    MPI_Win_fence(0, the_window);
    if (mytid==0) {
        MPI_Put( /* data on origin: */ &my_number, 1,MPI_INT,
                 /* data on target: */ other,1,      1,MPI_INT,
                 the_window);
    }
    MPI_Win_fence(0, the_window);
    MPI_Win_free(&the_window);

// putfence.py
window_data = np.zeros(2,dtype=np.int)
my_number = np.empty(1,dtype=np.int)
src = 0; tgt = nprocs-1
if procid==src:
    my_number[0] = 37
else:
    my_number[0] = 1

intsize = np.dtype('int').itemsize
win = MPI.Win.Create(window_data,intsize,comm=comm)

win.Fence()
if procid==src:
    # put data in the second element of the window
    win.Put(my_number,tgt,target=1)
win.Fence()
```

## 9.9 R

### MPI\_Recv

Using the `MPI_ANY_SOURCE` specifier. We retrieve the actual source from the `MPI_Status` object through the `MPI_SOURCE` field.

```
// anysource.c
if (mytid==ntids-1) {
    int *recv_buffer;
    MPI_Status status;

    recv_buffer = (int*) malloc((ntids-1)*sizeof(int));

    for (int p=0; p<ntids-1; p++) {
```

## 9. MPI Examples

---

```
    err = MPI_Recv(recv_buffer+p,1,MPI_INT, MPI_ANY_SOURCE, 0,comm,
                   &status); CHK(err);
    int sender = status.MPI_SOURCE;
    printf("Message from sender=%d: %d\n",
           sender,recv_buffer[p]);
}
} else {
    float randomfraction = (rand() / (double)RAND_MAX);
    int randomwait = (int) ( ntids * randomfraction );
    printf("process %d waits for %e/%d=%d\n",
           mytid,randomfraction,ntids,randomwait);
    sleep(randomwait);
    err = MPI_Send(&randomwait,1,MPI_INT, ntids-1,0,comm); CHK(err);
}

// anysource.py
rstatus = MPI.Status()
comm.Recv(rbuf,source=MPI.ANY_SOURCE,status=rstatus)
print "Message came from %d" % rstatus.Get_source()
```

### MPI\_Reduce

A reduction onto a root process.

```
// reduce.c
float myrandom = (float) rand()/(float)RAND_MAX,
      result;
int target_proc = ntids-1;
// add all the random variables together
MPI_Reduce(&myrandom,&result,1,MPI_FLOAT,MPI_SUM,
            target_proc,comm);
// the result should be approx ntids/2:
if (mytid==target_proc)
    printf("Result %6.3f compared to ntids/2=%5.2f\n",
           result,ntids/2.);
```

A reduction with reuse of the receive buffer through MPI\_IN\_PLACE.

```
// reduceinplace.c
float mynumber,result,*sendbuf,*recvbuf;
mynumber = (float) mytid;
int target_proc = ntids-1;
// add all the random variables together
if (mytid==target_proc) {
    sendbuf = (float*)MPI_IN_PLACE; recvbuf = &result;
```

```

        result = mynumber;
    } else {
        sendbuf = &mynumber;      recvbuf = NULL;
    }
MPI_Reduce(sendbuf, recvbuf, 1, MPI_FLOAT, MPI_SUM,
           target_proc, comm);
// the result should be ntids*(ntids-1)/2:
if (mytid==target_proc)
    printf("Result %6.3f compared to n(n-1)/2=%5.2f\n",
           result,ntids*(ntids-1)/2.);

MPI_Reduce_scatter

```

A simple illustration.

```

// reducescatter.c
// record what processes you will communicate with
int *recv_requests;
// find how many procs want to communicate with you
MPI_Reduce_scatter
    (recv_requests,&nsend_requests,counts,MPI_INT,
     MPI_SUM,comm);
// send a msg to the selected processes
for (int i=0; i<ntids; i++)
    if (recv_requests[i]>0)
        MPI_Isend(&msg,1,MPI_INT, /*to:*/ i,0,comm,
                   mpi_requests+irequest++);
// do as many receives as you know are coming in
for (int i=0; i<nsend_requests; i++)
    MPI_Irecv(&msg,1,MPI_INT,MPI_ANY_SOURCE,MPI_ANY_TAG,comm,
               mpi_requests+irequest++);
MPI_Waitall(irequest,mpi_requests,MPI_STATUSES_IGNORE);

```

Use of `MPI_Reduce_scatter` to implement the two-dimensional matrix-vector product. Set up separate row and column communicators with `MPI_Comm_split`, use `MPI_Reduce_scatter` to combine local products.

```

MPI_Allgather(&my_x,1,MPI_DOUBLE,
              local_x,1,MPI_DOUBLE,environ.col_comm);
bli_dgemv( BLIS_NO_TRANSPOSE,
            BLIS_NO_CONJUGATE,
            size_y, size_x,
            &one,
            local_matrix, 1, size_y,
            local_x, 1,

```

```

        &zero,
        local_y, 1 );
// blas_dgemv(CblasColMajor,CblasNoTrans,
// size_y,size_x,1.e0,
// local_matrix,size_y,
// local_x,1,0.e0,local_y,1);
MPI_Reduce_scatter(local_y,&my_y,&ione,MPI_DOUBLE,
MPI_SUM,environ.row_comm);

```

## 9.10 S

MPI\_Scan

In native mode the result is a function return value.

```

// scan.py
mycontrib = 10+random.randint(1,nprocs)
myfirst = 0
mypartial = comm.scan(mycontrib)
sbuf = np.empty(1,dtype=np.int)
rbuf = np.empty(1,dtype=np.int)
sbuf[0] = mycontrib
comm.Scan(sbuf,rbuf)

```

MPI\_Send

A regular ping-pong operation with MPI\_Send and MPI\_Recv. We repeat the experiment multiple times to get a reliable measurement of the time taken.

```

// pingpong.c
int src = 0,tgt = ntids/2;
double t, send=1.1,recv;
if (mytid==src) {
    t = MPI_Wtime();
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Send(&send,1,MPI_DOUBLE,tgt,0,comm);
        MPI_Recv(&recv,1,MPI_DOUBLE,tgt,0,comm,MPI_STATUS_IGNORE);
    }
    t = MPI_Wtime()-t; t /= NEXPERIMENTS;
    printf("Time for pingpong: %e\n",t);
} else if (mytid==tgt) {
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Recv(&recv,1,MPI_DOUBLE,src,0,comm,MPI_STATUS_IGNORE);
        MPI_Send(&recv,1,MPI_DOUBLE,src,0,comm);
    }
}

```

```

        }
    }

// pingpong.py
if mytid==0:
    data = [ 2.*i for i in range(s) ]
    starttime = MPI.Wtime()
    for test in range(ntests):
        comm.send(data,dest=ntids-1)
        rdata = comm.recv(source=ntids-1)
    elapsed = MPI.Wtime()-starttime
    print "Size=%d, elapsed time: %e" % (s,elapsed)
    c = data==rdata
    if not c:
        print "oops",data,rdata
elif mytid==ntids-1:
    for test in range(ntests):
        zdata = comm.recv(source=0)
        comm.send(zdata,dest=0)

// scipingpong.py
if mytid==0:
    data = np.arange(s, dtype=np.float64)
    rdata = np.empty(s, dtype=np.float64)
    for i in range(s):
        data[i] = i+1
    starttime = MPI.Wtime()
    for test in range(ntests):
        comm.Send([data,MPI.DOUBLE],dest=ntids-1)
        comm.Recv([rdata,MPI.DOUBLE],source=ntids-1)
    elapsed = MPI.Wtime()-starttime
    print "Size=%d, elapsed time: %e" % (s,elapsed)
    c = data==rdata #reduce( lambda x,y:x and y, [ data[i]==rdata
    if not c.all():
        print "oops",data,rdata
elif mytid==ntids-1:
    zdata = np.empty(s, dtype=np.float64)
    for test in range(ntests):
        comm.Recv([zdata,MPI.DOUBLE],source=0)
        comm.Send([zdata,MPI.DOUBLE],dest=0)

MPI_Send_init
```

## 9. MPI Examples

---

Persistent communication is setup up on the sending process with `MPI_Send_init` and `MPI_Recv_init`, then performed with `MPI_Startall`. The receiver is using regular sends and receives.

```
// persist.c
if (mytid==src) {
    MPI_Send_init(send,s,MPI_DOUBLE,tgt,0,comm,requests+0);
    MPI_Recv_init(recv,s,MPI_DOUBLE,tgt,0,comm,requests+1);
    printf("Size %d\n",s);
    t[cnt] = MPI_Wtime();
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Startall(2,requests);
        MPI_Waitall(2,requests,MPI_STATUSES_IGNORE);
    }
    t[cnt] = MPI_Wtime() - t[cnt];
    MPI_Request_free(requests+0); MPI_Request_free(requests+1);
} else if (mytid==tgt) {
    for (int n=0; n<NEXPERIMENTS; n++) {
        MPI_Recv(recv,s,MPI_DOUBLE,src,0,comm,MPI_STATUS_IGNORE);
        MPI_Send(recv,s,MPI_DOUBLE,src,0,comm);
    }
}

// persist.py
sendbuf = np.ones(size,dtype=np.int)
recvbuf = np.ones(size,dtype=np.int)
if procid==src:
    print "Size:",size
    times[isize] = MPI.Wtime()
    for n in range(nexperiments):
        requests[0] = comm.Isend(sendbuf[0:size],dest=tgt)
        requests[1] = comm.Irecv(recvbuf[0:size],source=tgt)
        MPI.Request.Waitall(requests)
        sendbuf[0] = sendbuf[0]+1
    times[isize] = MPI.Wtime() - times[isize]
elif procid==tgt:
    for n in range(nexperiments):
        comm.Recv(recvbuf[0:size],source=src)
        comm.Send(recvbuf[0:size],dest=src)

MPI_Sendrecv
```

We set up a ring structure and use `MPI_Sendrecv` to communicate between pairs.

```
// sendrecv.c
right = (mytid+1)%3; left = (mytid+2)%3;
```

```
MPI_Sendrecv( &my_data, 1, MPI_INTEGER, right, 0,
&other_data, 1, MPI_INTEGER, left, 0,
comm, MPI_STATUS_IGNORE);
```

MPI\_Ssend

Using MPI\_Ssend messages that would fall under the *eager limit* do block.

```
// ssendblock.c
other = 1-mytid;
sendbuf = (int*) malloc(sizeof(int));
recvbuf = (int*) malloc(sizeof(int));
size = 1;
MPI_Ssend(sendbuf, size, MPI_INT, other, 0, comm);
MPI_Recv(recvbuf, size, MPI_INT, other, 0, comm, &status);
printf("This statement is not reached\n");
```

## 9.11 T

MPI\_Type\_contiguous

We send a contiguous data type of double and receive it as an array of separate doubles; we use MPI\_Get\_count to ensure that we got the right amount of data.

```
// contiguous.c
MPI_Datatype newvectortype;
if (mytid==sender) {
    MPI_Type_contiguous(count, MPI_DOUBLE, &newvectortype);
    MPI_Type_commit(&newvectortype);
    MPI_Send(source, 1, newvectortype, receiver, 0, comm);
    MPI_Type_free(&newvectortype);
} else if (mytid==receiver) {
    MPI_Status recv_status;
    int recv_count;
    MPI_Recv(target, count, MPI_DOUBLE, sender, 0, comm,
             &recv_status);
    MPI_Get_count(&recv_status, MPI_DOUBLE, &recv_count);
    ASSERT(count==recv_count);
}
```

MPI\_Type\_indexed

We send an indexed data type and receive as separate integers.

## 9. MPI Examples

---

```
// indexed.c
displacements = (int*) malloc(count*sizeof(int));
blocklengths = (int*) malloc(count*sizeof(int));
source = (int*) malloc(totalcount*sizeof(int));
target = (int*) malloc(count*sizeof(int));
MPI_Datatype newvectortype;
if (mytid==sender) {
    MPI_Type_indexed(count,blocklengths,displacements,MPI_INT,&newvectortype);
    MPI_Type_commit(&newvectortype);
    MPI_Send(source,1,newvectortype,the_other,0,comm);
    MPI_Type_free(&newvectortype);
} else if (mytid==receiver) {
    MPI_Status recv_status;
    int recv_count;
    MPI_Recv(target,count,MPI_INT,the_other,0,comm,
             &recv_status);
    MPI_Get_count(&recv_status,MPI_INT,&recv_count);
    ASSERT(recv_count==count);
}

// indexed.py
displacements = np.empty(count,dtype=np.int)
blocklengths = np.empty(count,dtype=np.int)
source = np.empty(totalcount,dtype=np.float64)
target = np.empty(count,dtype=np.float64)
if procid==sender:
    newindextype = MPI.DOUBLE.Create_indexed(blocklengths,displacements)
    newindextype.Commit()
    comm.Send([source,1,newindextype],dest=the_other)
    newindextype.Free()
elif procid==receiver:
    comm.Recv([target,count,MPI.DOUBLE],source=the_other)

MPI_Type_struct
```

A struct data type can consist of different elementary datatypes, so in addition to the displacements and blocklengths we now have an array of MPI datatypes. Also note how the displacement computation is done in bytes.

```
// struct.c
struct object {
    char c;
    double x[2];
    int i;
```

```

};

MPI_Datatype newstructuretype;
int structlen = 3;
int blocklengths[structlen]; MPI_Datatype types[structlen];
MPI_Aint displacements[structlen];
// where are the components relative to the structure?
blocklengths[0] = 1; types[0] = MPI_CHAR;
displacements[0] = (size_t)&(myobject.c) - (size_t)&myobject;
blocklengths[1] = 2; types[1] = MPI_DOUBLE;
displacements[1] = (size_t)&(myobject.x[0]) - (size_t)&myobject;
blocklengths[2] = 1; types[2] = MPI_INT;
displacements[2] = (size_t)&(myobject.i) - (size_t)&myobject;
MPI_Type_create_struct(structlen,blocklengths,displacements,types,&newstructuretype);
MPI_Type_commit(&newstructuretype);

{
    MPI_Aint typesize;
    MPI_Type_extent(newstructuretype,&typesize);
    if (mytid==0) printf("Type extent: %d bytes\n",typesize);
}
if (mytid==sender) {
    MPI_Send(&myobject,1,newstructuretype,the_other,0,comm);
} else if (mytid==receiver) {
    MPI_Recv(&myobject,1,newstructuretype,the_other,0,comm,MPI_STATUS_IGNORE)
}
MPI_Type_free(&newstructuretype);

```

#### MPI\_Type\_vector

Send a strided data object with `Type_vector` and receive it as individual doubles. Use `MPI_Get_count` to inspect the `MPI_Status` object.

```

// vector.c
source = (double*) malloc(stride*count*sizeof(double));
target = (double*) malloc(count*sizeof(double));
MPI_Datatype newvectortype;
if (mytid==sender) {
    MPI_Type_vector(count,1,stride,MPI_DOUBLE,&newvectortype);
    MPI_Type_commit(&newvectortype);
    MPI_Send(source,1,newvectortype,the_other,0,comm);
    MPI_Type_free(&newvectortype);
} else if (mytid==receiver) {
    MPI_Status recv_status;
    int recv_count;
    MPI_Recv(target,count,MPI_DOUBLE,the_other,0,comm,
             &recv_status);
}

```

## 9. MPI Examples

---

```
    MPI_Get_count (&recv_status, MPI_DOUBLE, &recv_count) ;
    ASSERT(recv_count==count) ;
}
```

### 9.12 W

`MPI_Waitall`

Post non-blocking `MPI_Irecv` and `MPI_Isend` to/from all others, then use `MPI_Waitall` on the array of requests.

```
// irecvloop.c
MPI_Request requests =
    (MPI_Request*) malloc( 2*ntids*sizeof(MPI_Request) ) ;
recv_buffers = (int*) malloc( ntids*sizeof(int) ) ;
send_buffers = (int*) malloc( ntids*sizeof(int) ) ;
for (int p=0; p<ntids; p++) {
    int left_p = (p-1) % ntids,
        right_p = (p+1) % ntids;
    send_buffer[p] = ntids-p;
    MPI_Isend(sendbuffer+p,1,MPI_INT, right_p,0, requests+2*p);
    MPI_Irecv(recvbuffer+p,1,MPI_INT, left_p,0, requests+2*p+1);
}
MPI_Waitall(2*ntids,requests,MPI_STATUSES_IGNORE);
```

In python creating the array for the returned requests is somewhat tricky.

```
// irecvloop.py
requests = [ None ] * (2*nprocs)
sendbuffer = np.empty( nprocs, dtype=np.int )
recvbuffer = np.empty( nprocs, dtype=np.int )

for p in range(nprocs):
    left_p = (p-1) % nprocs
    right_p = (p+1) % nprocs
    requests[2*p] = comm.Isend( sendbuffer[p:p+1], dest=left_p)
    requests[2*p+1] = comm.Irecv( recvbuffer[p:p+1], source=right_p)
MPI.Request.Waitall(requests)
```

`MPI_Waitany`

Each process except for the root does a blocking send; the root posts `MPI_Irecv` from all other processors, then loops with `MPI_Waitany` until all requests have come in. Use `MPI_SOURCE` to test the index parameter of the wait call.

```
// irecv_source.c
if (mytid==ntids-1) {
    int *recv_buffer;
    MPI_Request *request; MPI_Status status;
    recv_buffer = (int*) malloc((ntids-1)*sizeof(int));
    request = (MPI_Request*) malloc((ntids-1)*sizeof(MPI_Request));

    for (int p=0; p<ntids-1; p++) {
        ierr = MPI_Irecv(recv_buffer+p, 1, MPI_INT, p, 0, comm,
                         request+p); CHK(ierr);
    }
    for (int p=0; p<ntids-1; p++) {
        int index, sender;
        MPI_Waitany(ntids-1, request, &index, &status); //MPI_STATUS_IGNORE);
        if (index!=status.MPI_SOURCE)
            printf("Mismatch index %d vs source %d\n", index, status.MPI_SOURCE);
        printf("Message from %d: %d\n", index, recv_buffer[index]);
    }
}
```

In python creating the array for the returned requests is somewhat tricky.

```
// irecv_source.py
if procid==nprocs-1:
    receive_buffer = np.empty(nprocs-1, dtype=np.int)
    requests = [ None ] * (nprocs-1)
    for sender in range(nprocs-1):
        requests[sender] = comm.Irecv(receive_buffer[sender:sender+1], source=sender)
    # alternatively: requests = [ comm.Irecv(s) for s in .... ]
    status = MPI.Status()
    for sender in range(nprocs-1):
        ind = MPI.Request.Waitany(requests, status=status)
        if ind!=status.Get_source():
            print "sender mismatch: %d vs %d" % (ind, status.Get_source())
        print "received from", ind
else:
    mywait = random.randint(1, 2*nprocs)
    print "[%d] wait for %d seconds" % (procid, mywait)
    time.sleep(mywait)
    mydata = np.empty(1, dtype=np.int)
    mydata[0] = procid
    comm.Send([mydata, MPI.INT], dest=nprocs-1)
```

MPI\_Win\_lock

## 9. MPI Examples

---

See the `Fetch_and_op` example.

`MPI_Win_start`

A one-sided `MPI_Put` using active target synchronization: use `MPI_Win_start` and `MPI_Win_complete` on the origin, and `MPI_Win_post` and `MPI_Win_wait` on the target.

```
// postwaitwin.c
if (mytid==origin) {
    MPI_Group_incl(all_group,1,&target,&two_group);
    // access
    MPI_Win_start(two_group,0,the_window);
    MPI_Put( /* data on origin: */    &my_number, 1,MPI_INT,
            /* data on target: */   target,0, 1,MPI_INT,
            the_window);
    MPI_Win_complete(the_window);
}

if (mytid==target) {
    MPI_Group_incl(all_group,1,&origin,&two_group);
    // exposure
    MPI_Win_post(two_group,0,the_window);
    MPI_Win_wait(the_window);
}
```

`MPI_Win_create`

See the `MPI_Get` example.

`MPI_Win_fence`

One process does `MPI_Put` operations, randomly on one of two other processes. We use a fence for active target synchronization.

```
// randomput.c
MPI_Win_create(&window_data,sizeof(int),sizeof(int),
               MPI_INFO_NULL,comm,&the_window);

for (int c=0; c<10; c++) {
    float randomfraction = (rand() / (double)RAND_MAX);
    if (randomfraction>.5)
        other = 2;
    else other = 1;
    window_data = 0;
    my_sum += window_data;
}
```

```
if (mytid>0 && mytid<3)
    printf("Sum on %d: %d\n", mytid, my_sum);
if (mytid==0) printf("(sum should be 10)\n");
```

# Chapter 10

## MPI Review

### 10.1 Review questions

If you answer that a statement is false, give a one-line explanation.

1. True or false: `mpicc` is a compiler.
2. True or false: `mpirun` can only be used for interactive parallel runs.
3. What is the function of a hostfile?
4. True or false: in each communicator, processes are numbered consecutively from zero.
5. Describe a deadlock scenario involving three processors.
6. True or false: a message sent with `MPI_Isend` from one processor can be received with an `MPI_Recv` call on another processor.
7. True or false: a message sent with `MPI_Send` from one processor can be received with an `MPI_Irecv` on another processor.
8. Why does the `MPI_Irecv` call not have an `MPI_Status` argument?
9. What is the relation between the concepts of ‘origin’, ‘target’, ‘fence’, and ‘window’ in one-sided communication.
10. What are the three routines for one-sided data transfer?
11. Give an example of a collective call with and without a root processor.
12. Given a distributed array, meaning that every processor has

```
double x[N]; // N can vary per processor
```

give the approximate MPI-based code that computes the maximum value in the array, and leaves the result on every processor.

13. Give two examples of derived datatypes.
14. Give a practical example where the sender uses a different type to send than the receiver uses in the corresponding receive call. Name the types involved.

## **PART II**

### **OPENMP**

# Chapter 11

## OpenMP tutorial

### 11.1 Basics

#### 11.1.1 The OpenMP model

OpenMP is based on two concepts: the use of *threads* and the *fork/join model* of parallelism. For now you can think of a thread as a sort of process: the computer executes a sequence of instructions. The fork/join model says that a thread can split itself ('fork') into a number of threads that are identical copies. At some point these copies go away and the original thread is left ('join'), but while the *team of threads* created by the fork exists, you have parallelism available to you. The part of the execution between fork and join is known as a *parallel region*.

Figure 11.1 gives a simple picture of this: a thread forks into a team of threads, and these threads themselves can fork again.

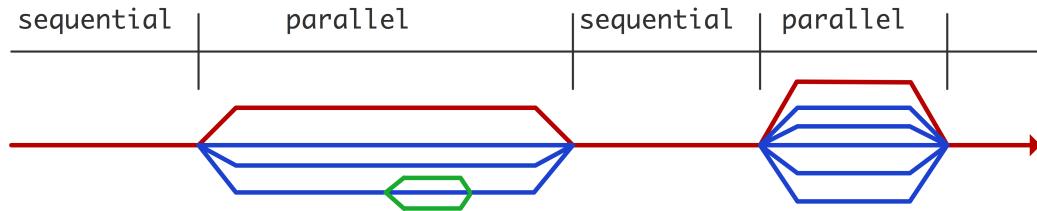


Figure 11.1: Thread creation and deletion during parallel execution

The threads that are forked are all copies of the *master thread*: they have access to all that was computed so far; this is their *shared data*. Of course, if the threads were completely identical the parallelism would be pointless, so they also have private data, and they can identify themselves: they know their thread number. This allows you to do meaningful parallel computations with threads.

This brings us to the third important concept: that of *work sharing* constructs. In a team of threads, initially there will be replicated execution; a work sharing construct divides available parallelism over the threads.

So there you have it: OpenMP uses teams of threads, and inside a parallel region the work is distributed over the threads with a work sharing construct. Threads can access shared data, and they have some private data.

For more on threads, see HPSC-[2.6.1](#).

### 11.1.2 Getting started with OpenMP programming

#### 11.1.2.1 Compiling

Your file or Fortran module needs to contain

```
#include "omp.h"
```

in C, and

```
use omp_lib
```

for Fortran.

OpenMP is handled by extensions to your regular compiler, typically by adding an option to your commandline:

```
# gcc
gcc -o foo foo.c -fopenmp
# Intel compiler
icc -o foo foo.c -openmp
```

If you have separate compile and link stages, you need that option in both.

When you use the openmp compiler option, a *cpp* variable `_OPENMP` will be defined. Thus, you can have conditional compilation by writing

```
#ifdef _OPENMP
...
#else
...
#endif
```

#### 11.1.2.2 Running an OpenMP program

You run an OpenMP program by invoking it the regular way (for instance `./a.out`), but its behaviour is influenced by some *OpenMP environment variables*. The most important one is `OMP_NUM_THREADS`:

```
export OMP_NUM_THREADS=8
```

which sets the number of threads that a program will use. See section [12.6](#) for a list of all environment variables.

### 11.1.3 OpenMP language constructs

The reference for the commands introduced here can be found in section [12.1.2](#).

Of course, OpenMP is not magic, so you have to tell it when something can be done in parallel. This is mostly done through *directives*; additional specifications can be done through library calls.

The first thing we want to do is create a team of threads. This is done with a *parallel region*. Here is a very simple example:

```
// hello.c
#pragma omp parallel
{
    printf("Hello world!\n");
}
```

or in Fortran

```
// hello.F90
 !$omp parallel
    print *, "Hello world!"
 !$omp end parallel
```

This code corresponds to the model we just discussed:

- Immediately preceding the parallel block, one thread will be executing the code. In the main program this is the *initial thread*.
- At the start of the block, a new *team of threads* is created, and the thread that was active before the block becomes the *master thread* of that team.
- After the block only the master thread is active.
- Inside the block there is team of threads: each thread in the team executes the body of the block, and it will have access to all variables of the surrounding environment. How many threads there are can be determined in a number of ways; we will get to that later.

**Exercise 11.1.** Make a full program based on this fragment. Insert different print statements before, inside, and after the parallel region. Run this example. How many times is each print statement executed?

You see that the `parallel` directive

- Is preceded by a special marker: a `#pragma omp` for C/C++, and the `!$OMP sentinel` for Fortran;
- Is followed by a single statement or a block in C/C++, or followed by a block in Fortran which is delimited by an `!$omp end` directive.

Directives look like *cpp directives*, but they are actually handled by the compiler, not the preprocessor.

### 11.1.4 Parallel regions

The reference for the commands introduced here can be found in section [12.2](#).

The simplest way to create parallelism in OpenMP is to use the `parallel` pragma. A block preceded by the `omp parallel` pragma is executed by a newly created team of threads. This is an instance of the *Single Program Multiple Data (SPMD)* model: all threads execute the same segment of code.

```
#pragma omp parallel
{
    // this is executed by a team of threads
}
```

It would be pointless to have the block be executed identically by all threads. One way to get a meaningful parallel code is to use the function `omp_get_thread_num`, to find out which thread you are, and execute work that is individual to that thread. There is also a function `omp_get_num_threads` to find out the total number of threads. Both these functions give a number relative to the current team; recall from figure 11.1 that new teams can be created recursively.

For instance, if you program computes

```
result = f(x) + g(x) + h(x)
```

you could parallelize this as

```
double result = 0;
#pragma omp parallel
{ // this is not the right way to do it....
    int num = omp_get_thread_num();
    if (num==0)      result += f(x);
    else if (num==1) result += g(x);
    else if (num==2) result += h(x);
}
```

This example shows how the three functions are computed in parallel, but other than that **there are many things wrong with this example**. Further in this section we will explain what is wrong here, and what can be done about it.

**Exercise 11.2.** Take the ‘hello world’ program above, and modify it so that you get multiple messages to your screen, saying

```
Hello from thread 0 out of 4!
Hello from thread 1 out of 4!
```

and so on. (The messages may very well appear out of sequence.)

What happens if you set your number of threads larger than the available cores on your computer?

**Exercise 11.3.** What happens if you call `omp_get_thread_num` and `omp_get_num_threads` outside a parallel region?

**Exercise 11.4.** Test nested parallelism. First of all, set `OMP_NESTED` to `TRUE` if that’s not the default on your system. Now write an OpenMP program as follows:

1. Write a subprogram that contains a parallel region.
2. Write a main program with a parallel region; call the subprogram both inside and outside the parallel region.
3. Insert print statements
  - (a) in the main program outside the parallel region,
  - (b) in the parallel region in the main program,
  - (c) in the subprogram outside the parallel region,
  - (d) in the parallel region inside the subprogram.

Run your program and count how many print statements of each type you get.

### 11.1.5 Thread data

In most programming languages, visibility of data is governed by rules on the *scope of variables*: a variable is declared in a block, and it is then visible to any statement in that block and blocks with a *lexical scope* contained in it, but not in surrounding blocks:

```
main () {
    // no variable 'x' define here
    {
        int x = 5;
        if (somecondition) { x = 6; }
        printf("x=%e\n",x); // prints 5 or 6
    }
    printf("x=%e\n",x); // syntax error: 'x' undefined
}
```

Fortran has simpler rules, since it does not have blocks inside blocks.

OpenMP has similar rules concerning data in parallel regions and other OpenMP constructs. First of all, data is visible in enclosed scopes:

```
main() {
    int x;
#pragma omp parallel
    {
        // you can use and set 'x' here
    }
    printf("x=%e\n",x); // value depends on what
                        // happened in the parallel region
}
```

In C, you can redeclare a variable inside a nested scope:

```
{
    int x;
    if (something) {
```

```

        double x; // same name, different entity
    }
    x = ... // this refers to the integer again
}

```

Doing so makes the outer variable inaccessible.

OpenMP has a similar mechanism:

```

{
    int x;
#pragma omp parallel
    {
        double x;
    }
}

```

There is an important difference: each thread in the team gets its own instance of the enclosed variable.

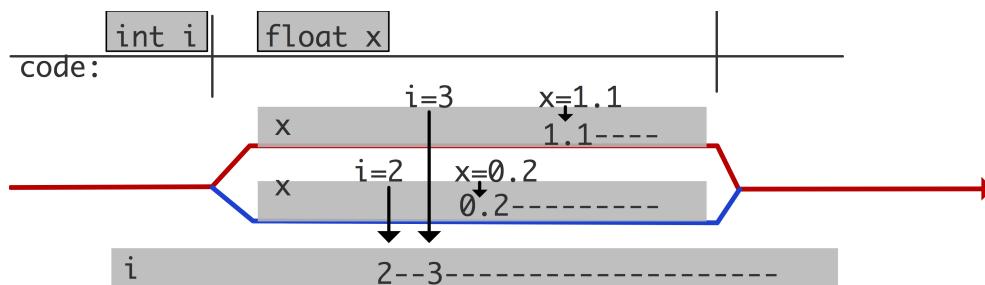


Figure 11.2: Locality of variables in threads

This is illustrated in figure 11.2.

In addition to such scoped variables, which live on a *stack*, there are variables on the *heap*, typically created by a call to `malloc` (in C) or `new` (in C++). Rules for them are more complicated.

Summarizing the above, there are

- *shared variables*, where each thread refers to the same data item, and
- *private variables*, where each thread has its own instance.

In addition to using scoping, OpenMP also uses options on the directives to control whether data is private or shared.

Many of the difficulties of parallel programming with OpenMP stem from the use of shared variables. For instance, if two threads update a shared variable, you not guarantee an order on the updates.

We will discuss all this in detail in section 11.3.

### 11.1.6 Creating parallelism

*The reference for the commands introduced here can be found in section ??.*

The *fork/join model* of OpenMP means that you need some way of indicating where an activity can be forked for independent execution. There are two ways of doing this:

1. You can declare a parallel region and split one thread into a whole team of threads. We will discuss this next in section 11.1.4. The division of the work over the threads is controlled by *work sharing construct* (section 11.2).
2. Alternatively, you can use tasks and indicating one parallel activity at a time. You will see this in section 11.6.

Note that OpenMP only indicates how much parallelism is present; whether independent activities are in fact executed in parallel is a runtime decision. The factors influencing this are discussed in section 11.1.6.

Declaring a parallel region tells OpenMP that a team of threads can be created. The actual size of the team depends on various factors (see section 12.6 for variables and functions mentioned in this section).

- The *environment variable* `OMP_NUM_THREADS` limits the number of threads that can be created.
- If you don't set this variable, you can also set this limit dynamically with the *library routine* `omp_set_num_threads`. This routine takes precedence over the aforementioned environment variable if both are specified.
- A limit on the number of threads can also be set as a clause on a parallel region.

If you specify a greater amount of parallelism than the hardware supports, the runtime system will probably ignore your specification and choose a lower value. To ask how much parallelism is actually used in your parallel region, use `omp_get_num_threads`. To query these hardware limits, use `omp_get_num_procs`.

Another limit on the number of threads is imposed when you use nested parallel regions. This can arise if you have a parallel region in a subprogram which is sometimes called sequentially, sometimes in parallel. The variable `OMP_NESTED` controls whether the inner region will create a team of more than one thread.

## 11.2 Work sharing

*The reference for the commands introduced here can be found in section 12.3.*

The declaration of a *parallel region* establishes a team of threads. This offers the possibility of parallelism, but to actually get meaningful parallel activity you need something more. OpenMP uses the concept of a *work sharing construct*: a way of dividing parallelizable work over a team of threads. The work sharing constructs are:

- `for/do` The threads divide up the loop iterations among themselves; see 11.2.1.
- `sections` The threads divide a fixed number of sections between themselves; see section 11.2.2.
- `single` The section is executed by a single thread; section 11.2.3.
- `task`
- `workshare` Can parallelize Fortran array syntax.

### 11.2.1 Loop parallelism

The reference for the commands introduced here can be found in section [12.3.1](#).

The parallel execution of a loop can be handled a number of different ways. For instance, you can create a parallel region around the loop, and adjust the loop bounds:

```
#pragma omp parallel
{
    int threadnum = omp_get_thread_num(),
        numthreads = omp_get_num_threads();
    int low = N*threadnum/numthreads,
        high = N*(threadnum+1)/numthreads;
    for (i=low; i<high; i++)
        // do something with i
}
```

A more natural option is to use the `parallel for` pragma:

```
#pragma omp parallel
#pragma omp for
for (i=0; i<N; i++) {
    // do something with i
}
```

This has several advantages. For one, you don't have to calculate the loop bounds for the threads yourself, but you can also tell OpenMP to assign the loop iterations according to different schedules (section [11.2.1.1](#)).

Figure [11.3](#) shows the execution on four threads of

```
#pragma omp parallel
{
    code1();
#pragma omp for
    for (i=1; i<=4*N; i++) {
        code2();
    }
    code3();
}
```

The code before and after the loop is executed identically in each thread; the loop iterations are spread over the four threads.

If you parallel region only contains a loop, you can combine the pragmas for the parallel region and distribution of the loop iterations:

```
#pragma omp parallel for
```

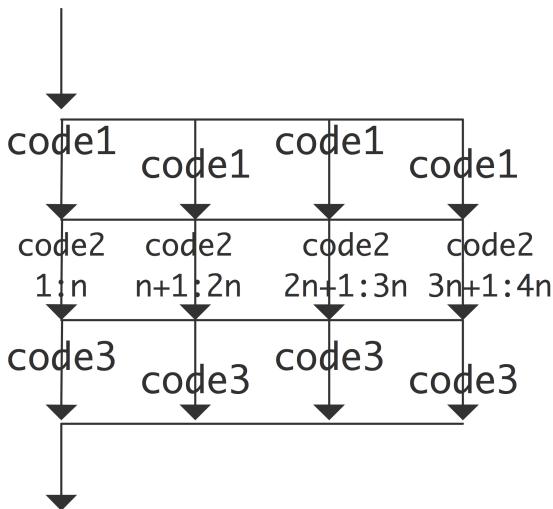


Figure 11.3: Execution of parallel code inside and outside a loop

```
for (i=0; .....
```

The loop has to satisfy a number of basic constraints, such as that you can not jump out of it.

**Exercise 11.5.** Compute  $\pi$  by *numerical integration*. We use the fact that  $\pi$  is the area of the unit circle, and we approximate this by computing the area of a quarter circle using *Riemann sums*.

- Let  $f(x) = \sqrt{1 - x^2}$  be the function that describes the quarter circle for  $x = 0 \dots 1$ ;
- Then we compute

$$\pi/4 \approx \sum_{i=0}^{N-1} \Delta x f(x_i) \quad \text{where } x_i = i\Delta x \text{ and } \Delta x = 1/N$$

Write a program for this, and parallelize it using OpenMP parallel for directives. Use different numbers of cores and compute the speedup you attain over the sequential computation. Is there a performance difference between the OpenMP code with 1 thread and the sequential code?

#### 11.2.1.1 Loop schedules

The reference for the commands introduced here can be found in section [12.3.1.1](#).

Usually you will have many more iterations in a loop than there are threads. Thus, there are several ways you can assign your loop iterations to the threads. OpenMP lets you specify this with the `schedule` clause.

```
#pragma omp for schedule(....)
```

The first distinction we now have to make is between static and dynamic schedules. With static schedules, the iterations are assigned purely based on the number of iterations and the number of threads (and the chunk parameter; see later). In dynamic schedules, on the other hand, iterations are assigned to threads that are unoccupied. Dynamic schedules are a good idea if iterations take an unpredictable amount of time, so that *load balancing* is needed.

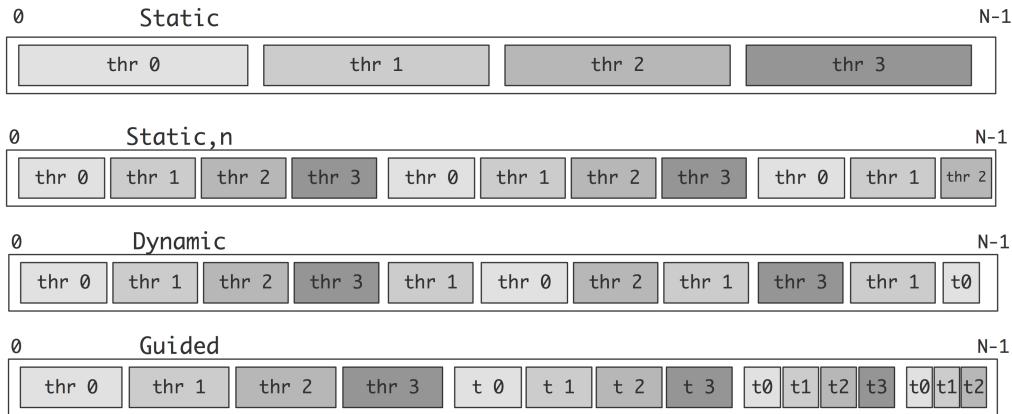


Figure 11.4: Illustration of the loop scheduling strategies

The default static schedule is to assign one consecutive block of iterations to each thread. If you want different sized blocks you can define a chunk size:

```
#pragma omp for schedule(static[,chunk])
```

(where the square brackets indicate an optional argument). With static scheduling, the compiler will split up the loop iterations at compile time, so, provided the iterations take roughly the same amount of time, this is the most efficient at runtime.

The choice of a chunk size is often a balance between the low overhead of having only a few chunks, versus the load balancing effect of having smaller chunks.

Exercise 11.6. Why is a chunk size of 1 typically a bad idea? (Hint: think about cache lines, and read HPSC-1.4.1.2.)

In dynamic scheduling OpenMP will put blocks of iterations (the default chunk size is 1) in a task queue, and the threads take one of these tasks whenever they are finished with the previous.

```
#pragma omp for schedule(static[,chunk])
```

While this schedule may give good load balancing if the iterations take very differing amounts of time to execute, it does carry runtime overhead for managing the queue of iteration tasks.

Finally, there is the guided schedule, which gradually decreases the chunk size. The thinking here is that large chunks carry the least overhead, but smaller chunks are better for load balancing. The various schedules are illustrated in figure 11.4.

If you don't want to decide on a schedule in your code, you can specify the `runtime` schedule. The actual schedule will then at runtime be read from the `OMP_SCHEDULE` environment variable. You can even just leave it to the runtime library by specifying `auto`

**Exercise 11.7.** We continue with exercise 11.5. We add ‘adaptive integration’: where needed, the program refines the step size<sup>1</sup>. This means that the iterations no longer take a predictable amount of time.

```

for (i=0; i<nsteps; i++) {
    double
    x = i*h, x2 = (i+1)*h,
    y = sqrt(1-x*x), y2 = sqrt(1-x2*x2),
    slope = (y-y2)/h;
    if (slope>15) slope = 15;
    int
    samples = 1+(int)slope,
    is;
    for (is=0; is<samples; is++) {
        double
        hs = h/samples,
        xs = x+ is*hs,
        ys = sqrt(1-xs*xs);
        quarterpi += hs*ys;
        nsamples++;
    }
}
pi = 4*quarterpi;

```

1. Use the `omp parallel for` construct to parallelize the loop. As in the previous lab, you may at first see an incorrect result. Use the `reduction` clause to fix this.
2. Your code should now see a decent speedup, using up to 8 cores. However, it is possible to get completely linear speedup. For this you need to adjust the schedule.  
Start by using `schedule(static, $n$)`. Experiment with values for  $n$ . When can you get a better speedup? Explain this.
3. Since this code is somewhat dynamic, try `schedule(dynamic)`. This will actually give a fairly bad result. Why? Use `schedule(dynamic, $n$)` instead, and experiment with values for  $n$ .
4. Finally, use `schedule(guided)`, where OpenMP uses a heuristic. What results does that give?

**Exercise 11.8.** Program the *LU factorization* algorithm without pivoting.

```

for k=1,n:
    A[k,k] = 1./A[k,k]
    for i=k+1,n:
        A[i,k] = A[i,k]/A[k,k]

```

---

1. It doesn't actually do this in a mathematically sophisticated way, so this code is more for the sake of the example.

```

for j=k+1,n:
    A[i,j] = A[i,k]*A[k,j]

```

1. Argue that it is not possible to parallelize the outer loop.
2. Argue that it is possible to parallelize both the  $i$  and  $j$  loops.
3. Parallelize the algorithm by focusing on the  $i$  loop. Why is the algorithm as given here best for a matrix on row-storage? What would you do if the matrix was on column storage?
4. Argue that with the default schedule, if a row is updated by one thread in one iteration, it may very well be updated by another thread in another. Can you find a way to schedule loop iterations so that this does not happen? What practical reason is there for doing so?

### 11.2.1.2 Reductions

So far we have focused on loops with independent iterations. Reductions are a common type of loop with dependencies. There is an extended discussion of reductions in section 11.4

### 11.2.1.3 nowait

The implicit barrier at the end of a work sharing construct can be cancelled with a `nowait` clause. This has the effect that threads that are finished can continue with the next code in the parallel region:

```

#pragma omp parallel
{
#pragma omp for nowait
    for (i=0; i<N; i++) { ... }
    // more parallel code
}

```

In the following example, threads that are finished with the first loop can start on the second. Note that this requires both loops to have the same schedule.

```

#pragma omp parallel
{
    x = local_computation()
#pragma omp for nowait
    for (i=0; i<N; i++) {
        x[i] = ...
    }
#pragma omp for
    for (i=0; i<N; i++) {
        y[i] = ... x[i] ...
    }
}

```

### 11.2.1.4 While loops

OpenMP can only handle ‘for’ loops: *while loops* can not be parallelized. So you have to find a way around that. While loops are for instance used to search through data:

```
while ( a[i]!=0 && i<imax ) {  
    i++; }  
// now i is the first index for which \n{a[i]} is zero.
```

We replace the while loop by a for loop that examines all locations:

```
result = -1;  
#pragma omp parallel for  
for (i=0; i<imax; i++) {  
    if (a[i]!=0 && result<0) result = i;  
}
```

**Exercise 11.9.** Show that this code has a race condition.

You can fix the race condition by making the condition into a critical section; section 11.5.2.1. In this particular example, with a very small amount of work per iteration, that is likely to be inefficient in this case (why?). A more efficient solution uses the `lastprivate` pragma:

```
result = -1;  
#pragma omp parallel for lastprivate(result)  
for (i=0; i<imax; i++) {  
    if (a[i]!=0) result = i;  
}
```

You have now solved a slightly different problem: the `result` variable contains the *last* location where `a[i]` is zero.

### 11.2.2 Sections

*The reference for the commands introduced here can be found in section ??.*

A parallel loop is an example of independent work units that are numbered. If you have a pre-determined number of independent work units, the `sections` is more appropriate. In a `sections` construct can be any number of `section` constructs. These need to be independent, and they can be execute by any available thread in the current team, including having multiple sections done by the same thread.

```
#pragma omp sections  
{  
#pragma omp section  
    // one calculation  
#pragma omp section  
    // another calculation  
}
```

This construct can be used to divide large blocks of independent work. Suppose that in the following line, both  $f(x)$  and  $g(x)$  are big calculations:

```
y = f(x) + g(x)
```

You could then write

```
double y1,y2;
#pragma omp sections
{
#pragma omp section
    y1 = f(x)
#pragma omp section
    y2 = g(x)
}
y = y1+y2;
```

Instead of using two temporaries, you could also use a critical section; see section 11.5.2.1. However, the best solution is have a reduction clause on the sections directive:

```
y = f(x) + g(x)
```

You could then write

```
y = 0;
#pragma omp sections reduction(+:y)
{
#pragma omp section
    y += f(x)
#pragma omp section
    y += g(x)
}
```

### 11.2.3 Single/master

*The reference for the commands introduced here can be found in section ??.*

```
master single
```

The `single` and `master` pragma limit the execution of a block to a single thread. This can for instance be used to print tracing information or doing I/O operations.

```
#pragma omp parallel
{
#pragma omp single
    printf("We are starting this section!\n");
    // parallel stuff
```

```
}
```

Another use of `single` is to perform initializations in a parallel region:

```
int a;
#pragma omp parallel
{
    #pragma omp single
    a = f(); // some computation
    #pragma omp sections
    // various different computations using a
}
```

The point of the `single` directive in this last example is that the computation needs to be done only once, because of the shared memory. Since it's a work sharing construct there is an *implicit barrier* after it, which guarantees that all threads have the correct value in their local memory (see section [11.8.4](#)).

**Exercise 11.10.** What is the difference between this approach and how the same computation would be parallelized in MPI?

The `master` directive, also enforces execution on a single thread, specifically the master thread of the team, but it does not have the synchronization through the implicit barrier.

**Exercise 11.11.** Modify the above code to read:

```
int a;
#pragma omp parallel
{
    #pragma omp master
    a = f(); // some computation
    #pragma omp sections
    // various different computations using a
}
```

This code is no longer correct. Explain.

### 11.2.4 Fortran array syntax parallelization

*The reference for the commands introduced here can be found in section [12.3.4](#).*

The `parallel do` directive is used to parallelize loops, and this applies to both C and Fortran. However, Fortran also has implied loops in its *array syntax*. To parallelize array syntax you can use the `workshare` directive.

## 11.3 Controlling thread data

*The reference for the commands introduced here can be found in section [12.4](#).*

In a parallel region there are two types of data: private and shared. In this sections we will see the various way you can control what category your data falls under; for private data items we also discuss how their values relate to shared data.

### 11.3.1 Shared data

In a parallel region, any data declared outside it will be shared: any thread using a variable `x` will access the same memory location associated with that variable.

Example:

```
int x = 5;
#pragma omp parallel
{
    x = x+1;
    printf("shared: x is %d\n", x);
}
```

All threads increment the same variable, so after the loop it will have a value of five plus the number of threads; or maybe less because of the data races involved. See HPSC-[2.6.1.5](#) for an explanation of the issues involved; see [11.5.2.1](#) for a solution in OpenMP.

Sometimes this global update is what you want; in other cases the variable is intended only for intermediate results in a computation. In that case there are various ways of creating data that is local to a thread, and therefore invisible to other threads.

### 11.3.2 Private data

*The reference for the commands introduced here can be found in section [12.4.2](#).*

In the C/C++ language it is possible to declare variables inside a *lexical scope*; roughly: inside curly braces. This concept extends to OpenMP parallel regions and directives: any variable declared in a block following an OpenMP directive will be local to the executing thread.

Example:

```
int x = 5;
#pragma omp parallel
{
    int x; x = 3;
    printf("local: x is %d\n", x);
}
```

After the parallel region, the outer variable `x` will still have the value 5.

The Fortran language does not have this concept of scope, so you have to use a `private` clause:

```
!$OMP parallel private(x)
```

The `private` directive declares data to have a separate copy in the memory of each thread. Such private variables are initialized as they would be in a main program. Any computed value goes away at the end of the parallel region. (However, see below.) Thus, you should not rely on any initial value, or on the value of the outer variable after the region.

```
int x = 5;
#pragma omp parallel private(x)
{
    x = x+1; // dangerous
    printf("private: x is %d\n", x);
}
printf("after: x is %d\n", x); // also dangerous
```

Private arrays are tricky.

- In C, if an array is statically defined, e.g.,

```
double a[2][5];
```

declaring `private(a)` will indeed put a copy on in each thread. Note that this will lead to an explosion in memory use; in fact, for large arrays you may experience *stack overflow*.

- Dynamically allocated arrays, e.g.,

```
double *a; a = (double*) malloc( some_size );
```

can not be made private with `private(a)`: this only gives each thread a private pointer, but these pointers all point to the same memory location.

### 11.3.3 Default

As remarked, most data in a parallel section is shared. This default behaviour can be controlled by adding a `default` clause:

```
#pragma omp parallel default(shared) private(x)
{ ... }
#pragma omp parallel default(private) shared(matrix)
{ ... }
```

and if you want to play it safe:

```
#pragma omp parallel default(none) private(x) shared(matrix)
{ ... }
```

Setting `default(none)` is useful for debugging. If your code behaves differently in parallel from sequential there is probably a data race. Specifying the status of every variable is a good way to debug this.

### 11.3.4 First and last private

You can force initialization with `firstprivate`.

```
int t=2;
#pragma omp parallel firstprivate(t)
{
    t += f(omp_get_thread_num());
    g(t);
}
```

The variable `t` behaves like a private variable, except that it is initialized to the outside value.

It is possible to preserve a private variable from the last iteration with `lastprivate`:

```
#pragma omp parallel for \
lastprivate(tmp)
for (i=0; i<N; i++) {
    tmp = .....
    x[i] = .... tmp ....
}
..... tmp .....
```

### 11.3.5 Persistent data

*The reference for the commands introduced here can be found in section ??.*

Most data in OpenMP parallel regions is either inherited from the master thread and therefore shared, or temporary within the scope of the region and fully private. There is also a mechanism for *thread-private data*, which is not limited in lifetime to one parallel region. The `threadprivate` pragma is used to declare that each thread is to have a private copy of a variable:

```
#pragma omp threadprivate(var)
```

This variable will then have the lifetime of an ordinary variable, but inside a parallel region the private copy is used.

The typical application for thread-private variables is in *random number generation*. A random number generator needs saved state, since it computes each next value from the current one. To have a parallel generator, each thread will create and initialize a private ‘current value’ variable. This will persist even when the execution is not in a parallel region; it gets updated only in a parallel region.

**Exercise 11.12.** Calculate the area of the *Mandelbrot set* by random sampling. Initialize the random number generator separately for each thread; then use a parallel loop to evaluate the points. Explore performance implications of the different loop scheduling strategies.

## 11.4 Reductions

*The reference for the commands introduced here can be found in section 12.3.1.3.*

Parallel tasks often produce some quantity that needs to be summed or otherwise combined. In section 11.1.4 you saw an example, and it was stated that the solution given there was not very good.

The problem in that example was the *race condition* involving the `result` variable. The simplest solution is to eliminate the race condition by declaring a *critical section*:

```
double result = 0;
#pragma omp parallel
{
    double local_result;
    int num = omp_get_thread_num();
    if (num==0)    local_result = f(x);
    else if (num==1) local_result = g(x);
    else if (num==2) local_result = h(x);
#pragma omp critical
    result += local_result;
}
```

This is a good solution if the amount of serialization in the critical section is small compared to computing the functions  $f, g, h$ . On the other hand, you may not want to do that in a loop:

```
double result = 0;
#pragma omp parallel
{
    double local_result;
#pragma omp for
    for (i=0; i<N; i++) {
        local_result = f(x,i);
#pragma omp critical
        result += local_result;
    } // end of for loop
}
```

**Exercise 11.13.** Can you think of a small modification of this code, that still uses a critical section, that is more efficient? Time both codes.

The easiest way to effect a reduction is of course to use the `reduction` clause. Adding this to an `omp for` or an `omp sections` construct has the following effect:

- OpenMP will make a copy of the reduction variable per thread, initialized to the identity of the reduction operator, for instance 1 for multiplication.
- Each thread will then reduce into its local variable;
- At the end of the loop, the local results are combined, again using the reduction operator, into the global variable.

This is one of those cases where the parallel execution can have a slightly different value from the one that is computed sequentially, because floating point operations are not associative. See HPSC-3.3.7 for more explanation.

If your code can not be easily structure as a reduction, you can realize the above scheme by hand by ‘duplicating’ the global variable and gather the contributions later. This example presumes three threads, and gives each a location of their own to store the result computed on that thread:

```
double result, local_results[3];
#pragma omp parallel
{
    int num = omp_get_thread_num();
    if (num==0)      local_results[num] = f(x)
    else if (num==1) local_results[num] = g(x)
    else if (num==2) local_results[num] = h(x)
}
result = local_results[0]+local_results[1]+local_results[2]
```

While this code is correct, it may be inefficient because of a phenomemon called *false sharing*. Even though the threads write to separate variables, those variables are likely to be on the same *cacheline* (see HPSC-1.4.1.2 for an explanation). This means that the cores will be wasting a lot of time and bandwidth updating each other’s copy of this cacheline.

False sharing can be prevent by giving each thread its own cacheline:

```
double result, local_results[3][8];
#pragma omp parallel
{
    int num = omp_get_thread_num();
    if (num==0)      local_results[num][1] = f(x)
    // et cetera
}
```

A more elegant solution gives each thread a true local variable, and uses a critical section to sum these, at the very end:

```
double result = 0;
#pragma omp parallel
{
    double local_result;
    local_result = .....
#pragma omp critical
    result += local_result;
}
```

## 11.5 Synchronization

In the constructs for declaring parallel regions above, you had little control over in what order threads executed the work they were assigned. This section will discuss *synchronization* constructs: ways of telling

threads to bring a certain order to the sequence in which they do things.

### 11.5.1 Barrier

*The reference for the commands introduced here can be found in section [12.5.1.1](#).*

A barrier defines a point in the code where all active threads will stop until all threads have arrived at that point. With this, you can guarantee that certain calculations are finished. For instance, in this code snippet, computation of `y` can not proceed until another thread has computed its value of `x`.

```
#pragma omp parallel
{
    int mytid = omp_get_thread_num();
    x[mytid] = some_calculation();
    y[mytid] = x[mytid]+x[mytid+1];
}
```

This can be guaranteed with a `barrier` pragma:

```
#pragma omp parallel
{
    int mytid = omp_get_thread_num();
    x[mytid] = some_calculation();
#pragma omp barrier
    y[mytid] = x[mytid]+x[mytid+1];
}
```

### 11.5.2 Mutual exclusion

Sometimes it is necessary to let only one thread execute a piece of code. Such a piece of code is called a *critical section*, and OpenMP has several mechanisms for realizing this.

The most common use of critical sections is to update a variable. Since updating involves reading the old value, and writing back the new, this has the possibility for a *race condition*: another thread reads the current value before the first can update it; the second thread updates to the wrong value.

Critical sections are an easy way to turn an existing code into a correct parallel code. However, there are disadvantages to this, and sometimes a more drastic rewrite is called for.

#### 11.5.2.1 *critical* and *atomic*

*The reference for the commands introduced here can be found in section [12.5.2](#).*

There are two pragmas for critical sections: `critical` and `atomic`. The second one is more limited but has performance advantages.

The typical application of a critical section is to update a variable:

```
#pragma omp parallel
{
    int mytid = omp_get_thread_num();
    double tmp = some_function(mytid);
#pragma omp critical
    sum += tmp;
}
```

**Exercise 11.14.** Consider a loop where each iteration updates a variable.

```
#pragma omp parallel for shared(result)
for ( i ) {
    result += some_function_of(i);
}
```

Discuss qualitatively the difference between:

- turning the update statement into a critical section, versus
- letting the threads accumulate into a private variable `tmp` as above, and summing these after the loop.

Do an Ahmdal-style quantitative analysis of the first case, assuming that you do  $n$  iterations on  $p$  threads, and each iteration has a critical section that takes a fraction  $f$ . Assume the number of iterations  $n$  is a multiple of the number of threads  $p$ . Also assume the default static distribution of loop iterations over the threads.

A `critical` section works by acquiring a lock, which carries a substantial overhead. Furthermore, if your code has multiple critical sections, they are all mutually exclusive: if a thread is in one critical section, the other ones are all blocked.

On the other hand, the syntax for `atomic` sections is limited to the update of a single memory location, but such sections are not exclusive and they can be more efficient, since they assume that there is a hardware mechanism for making them critical.

The problem with `critical` sections being mutually exclusive can be mitigated by naming them:

```
#pragma omp critical (optional_name_in_parens)
```

### 11.5.3 Locks

OpenMP also has the traditional mechanism of a *lock*. A lock is somewhat similar to a critical section: it guarantees that some instructions can only be performed by one process at a time. However, a critical section is indeed about code; a lock is about data. With a lock you make sure that some data elements can only be touched by one process at a time.

One simple example of the use of locks is generation of a *histogram*. A histogram consists of a number of bins, that get updated depending on some data. Here is the basic structure of such a code:

```
int count[100];
float x = some_function();
int ix = (int)x;
if (ix>=100)
    error();
else
    count[ix]++;

```

It would be possible to guard the last line:

```
#pragma omp critical
    count[ix]++;

```

but that is unnecessarily restrictive. If there are enough bins in the histogram, and if the `some_function` takes enough time, there are unlikely to be conflicting writes. The solution then is to create an array of locks, with one lock for each `count` location.

**Exercise 11.15.** In the following code, one process sets array A and then uses it to update B; the other process sets array B and then uses it to update A. Argue that this code can deadlock. How could you fix this?

```
#pragma omp parallel shared(a, b, nthreads, locka, lockb)
    #pragma omp sections nowait
    {
        #pragma omp section
        {
            omp_set_lock(&locka);
            for (i=0; i<N; i++)
                a[i] = ..

            omp_set_lock(&lockb);
            for (i=0; i<N; i++)
                b[i] = .. a[i] ..
            omp_unset_lock(&lockb);
            omp_unset_lock(&locka);
        }

        #pragma omp section
        {
            omp_set_lock(&lockb);
            for (i=0; i<N; i++)
                b[i] = ...

            omp_set_lock(&locka);
            for (i=0; i<N; i++)
                a[i] = .. b[i] ..
        }
    }

```

---

```

        omp_unset_lock(&locka);
        omp_unset_lock(&lockb);
    }
} /* end of sections */
} /* end of parallel region */

```

#### 11.5.4 Example: Fibonacci computation

The *Fibonacci sequence* is recursively defined as

$$F(0) = 1, \quad F(1) = 1, \quad F(n) = F(n-1) + F(n-2) \text{ for } n \geq 2.$$

We start by sketching the basic single-threaded solution. The naive code looks like:

```

int main() {
    value = new int[nmax+1];
    value[0] = 1;
    value[1] = 1;
    fib(10);
}

int fib(int n) {
    int i, j, result;
    if (n>=2) {
        i=fib(n-1); j=fib(n-2);
        value[n] = i+j;
    }
    return value[n];
}

```

However, this is inefficient, since most intermediate values will be computed more than once. We solve this by keeping track of which results are known:

```

...
done = new int[nmax+1];
for (i=0; i<=nmax; i++)
    done[i] = 0;
done[0] = 1;
done[1] = 1;
...
int fib(int n) {
    int i, j;
    if (!done[n]) {
        i = fib(n-1); j = fib(n-2);

```

```

        value[n] = i+j; done[n] = 1;
    }
    return value[n];
}

```

The OpenMP parallel solution calls for two different ideas. First of all, we parallelize the recursion by using tasks (section 11.6):

```

int fib(int n) {
    int i, j;
    if (n>=2) {
#pragma omp task shared(i) firstprivate(n)
        i=fib(n-1);
#pragma omp task shared(j) firstprivate(n)
        j=fib(n-2);
#pragma omp taskwait
        value[n] = i+j;
    }
    return value[n];
}

```

This computes the right solution, but, as in the naive single-threaded solution, it recomputes many of the intermediate values.

A naive addition of the done array leads to data races, and probably an incorrect solution:

```

int fib(int n) {
    int i, j, result;
    if (!done[n]) {
#pragma omp task shared(i) firstprivate(n)
        i=fib(n-1);
#pragma omp task shared(i) firstprivate(n)
        j=fib(n-2);
#pragma omp taskwait
        value[n] = i+j;
        done[n] = 1;
    }
    return value[n];
}

```

For instance, there is no guarantee that the done array is updated later than the value array, so a thread can think that `done[n-1]` is true, but `value[n-1]` does not have the right value yet.

One solution to this problem is to use a lock, and make sure that, for a given index `n`, the values `done[n]` and `value[n]` are never touched by more than one thread at a time:

```

int fib(int n)
{
    int i, j;
    omp_set_lock( &(dolock[n]) );
    if (!done[n]) {
        #pragma omp task shared(i) firstprivate(n)
        i = fib(n-1);
        #pragma omp task shared(j) firstprivate(n)
        j = fib(n-2);
        #pragma omp taskwait
        value[n] = i+j;
        done[n] = 1;
    }
    omp_unset_lock( &(dolock[n]) );
    return value[n];
}

```

This solution is correct, optimally efficient in the sense that it does not recompute anything, and it uses tasks to obtain a parallel execution.

However, the efficiency of this solution is only up to a constant. A lock is still being set, even if a value is already computed and therefore will only be read. This can be solved with a complicated use of critical sections, but we will forego this.

## 11.6 Tasks

*The reference for the commands introduced here can be found in section 12.7.*

Tasks are a mechanism that OpenMP uses under the cover: if you specify something as being parallel, OpenMP will create a ‘block of work’: a section of code plus the data environment in which it occurred. This block is set aside for execution at some later point.

Let’s look at a simple example using the `task` directive.

Code	Execution
<code>x = f();</code>	the variable <code>x</code> gets a value
<code>#pragma omp task { y = g(x); }</code>	a task is created with the current value of <code>x</code>
<code>z = h();</code>	the variable <code>z</code> gets a value

The thread that executes this code segment creates a task, which will later be executed, probably by a different thread. The exact timing of the execution of the task is up to a *task scheduler*, which operates invisible to the user.

Even though the above segment looks like a linear set of statements, it is impossible to say when the code after the `task` directive will be executed. This means that the following code is incorrect:

```

x = f();
#pragma omp task
{ y = g(x); }
z = h(y);

```

Explanation: when the statement computing `z` is executed, the task computing `y` has only been scheduled; it has not necessarily been executed yet.

In order to have a guarantee that a task is finished, you need the `taskwait` directive. The following creates two tasks, which can be executed in parallel, and then waits for the results:

Code	Execution
<code>x = f();</code>	the variable <code>x</code> gets a value
<code>#pragma omp task</code>	
<code>{ y1 = g1(x); }</code>	two tasks are created with the current value of <code>x</code>
<code>#pragma omp task</code>	
<code>{ y2 = g2(x); }</code>	
<code>#pragma omp taskwait</code>	the thread waits until the tasks are finished
<code>z = h(y1)+h(y2);</code>	the variable <code>z</code> is computed using the task results

The `task` pragma is followed by a structured block. Each time the structured block is encountered, a new task is generated. On the other hand `taskwait` is a standalone directive; the code that follows is just code, it is not a structured block belonging to the directive.

Another aspect of the distinction between generating tasks and executing them: usually the tasks are generated by one thread, but executed by many threads. Thus, the typical idiom is:

```

#pragma omp parallel
#pragma omp single
{
    // code that generates tasks
}

```

This makes it possible to execute loops in parallel that do not have the right kind of iteration structure for a `omp parallel for`. As an example, you could traverse and process a linked list:

```

#pragma omp parallel
#pragma omp single
{
    while (!tail(p)) {
        p = p->next();
    #pragma omp task
        process(p)
    }
    #pragma omp taskwait
}

```

One task traverses the linked list creating an independent task for each element in the list. These tasks are then executed in parallel; their assignment to threads is done by the task scheduler.

You can indicate task dependencies in several ways:

1. Using the ‘task wait’ directive you can explicitly indicate the *join* of the *forked* tasks. The instruction after the wait directive will therefore be dependent on the spawned tasks.
2. The `taskgroup` directive, followed by a structured block, ensures completion of all tasks created in the block, even if recursively created.
3. Each OpenMP task can have a `depend` clause, indicating what *data dependency* of the task. By indicating what data is produced or absorbed by the tasks, the scheduler can construct the dependency graph for you.

**Exercise 11.16.** Use the above approach to find the smallest factor of a large number (using 2,000,000,111 as test case): generate a task for each trial factor.

- Make sure you save only the smallest factor.
- Once a factor has been found, you should stop generating tasks. Let tasks that ‘should not have been generated’ print out a message.

### 11.6.1 Tree traversal

## 11.7 Version 4 functionality

### 11.7.1 SIMD

```
simd
```

```
declare simd
```

## 11.8 Stuff

### 11.8.1 Timing

*The reference for the commands introduced here can be found in section 12.8.1.*

OpenMP has a wall clock timer routine `omp_get_wtime` with resolution `omp_get_wtick`.

**Exercise 11.17.** Use the timing routines to demonstrate speedup from using multiple threads.

- Write a code segment that takes a measurable amount of time, that is, it should take a multiple of the tick time.
- Write a parallel loop and measure the speedup. You can for instance do this

```
for (int use_threads=1; use_threads<=nthreads; use_threads++) {  
    #pragma omp parallel for num_threads(use_threads)  
    for (int i=0; i<nthreads; i++) {  
        ....  
    }  
    if (use_threads==1)  
        time1 = tend-tstart;
```

```
else // compute speedup
```

- In order to prevent the compiler from optimizing your loop away, let the body compute a result and use a reduction to preserve these results.

### 11.8.2 Dependency analysis

If two statements refer to the same data item, we say that there is a *data dependency* between the statements. Such dependencies limit the extent to which the execution of the statements can be rearranged. The study of this topic probably started in the 1960s, when processors could execute statements *out of order* to increase throughput. The re-ordering of statements was limited by the fact that the execution had to obey the *program order* semantics: the result had to be as if the statements were executed strictly in the order in which they appear in the program.

These issues of statement ordering, and therefore of data dependencies, arise in OpenMP in two main ways:

1. When a loop is parallelized, the iterations are no longer executed in their program order, so we have to check for dependencies.
2. The introduction of tasks also means that parts of a program can be executed in a different order from in which they appear in a sequential execution.

The easiest case of dependency analysis is that of detecting that loop iterations can be executed independently. Iterations are of course independent if a data item is read in two different iterations, but if the same item is read in one iteration and written in another, or written in two different iterations, we need to do further analysis.

Analysis of *data dependencies* can be performed by a compiler, but compilers take, of necessity, a conservative approach. This means that iterations may be independent, but can not be recognized as such by a compiler. Therefore, OpenMP shifts this responsibility to the programmer.

The three types of dependencies are:

- flow dependencies, or ‘read-after-write’;
- anti dependencies, or ‘write-after-read’; and
- output dependencies, or ‘write-after-write’.

```
for (i) {
    y[i] = t;
    x[i+1] = y[i+1];
    t = x[i];
}
```

#### 11.8.2.1 Flow dependencies

Flow dependencies, or read-afer-write, are not a problem if the read and write occur in the same loop iteration:

```
for (i=0; i<N; i++) {  
    x[i] = .... ;  
    .... = ... x[i] ... ;  
}
```

On the other hand, if the read happens in a later iteration, there is no simple way to parallelize the loop:

```
for (i=0; i<N; i++) {  
    .... = ... x[i] ... ;  
    x[i+1] = .... ;  
}
```

This usually requires rewriting the code.

### 11.8.2.2 Anti dependencies

The simplest case of write-after-read is a reduction:

```
for (i=0; i<N; i++) {  
    t = t + ....  
}
```

This can be dealt with by explicit declaring the loop to be a reduction, or to use any of the other strategies in section 11.4.

If the read and write are on an array the situation is more complicated. The iterations in this fragment

```
for (i=0; i<N; i++) {  
    x[i] = ... x[i+1] ... ;  
}
```

can not be executed in arbitrary order as such. However, conceptually there is no dependency. We can solve this by introducing a temporary array:

```
for (i=0; i<N; i++)  
    xtmp[i] = x[i];  
for (i=0; i<N; i++) {  
    x[i] = ... xtmp[i+1] ... ;  
}
```

This is an example of a transformation that a compiler is unlikely to perform, since it can greatly affect the memory demands of the program. Thus, this is left to the programmer.

### 11.8.2.3 Output dependencies

The case of write-after-write does not occur by itself: if a variable is written twice in sequence without an intervening read, the first write can be removed without changing the meaning of the program. Thus, this

case reduces to a flow dependency.

Other output dependencies can easily be removed. In the following code, `t` can be declared private, thereby removing the dependency.

```
for (i=0; i<N; i++) {
    t = f(i)
    s += t*t;
}
```

If the final value of `t` is wanted, the `lastprivate` can be used.

### 11.8.3 Thread safety

With OpenMP it is relatively easy to take existing code and make it parallel by introducing parallel sections. If you're careful to declare the appropriate variables shared and private, this may work fine. However, your code may include calls to library routines that include a *race condition*; such code is said not to be *thread-safe*.

For example a routine

```
static int isave;
int next_one() {
    int i = isave;
    isave += 1;
    return i;
}

...
for ( .... ) {
    int ivalue = next_one();
}
```

has a clear race condition, as the iterations of the loop may get different `next_one` values, as they are supposed to, or not. This can be solved by using an `critical` pragma for the `next_one` call; another solution is to use an `threadprivate` declaration for `isave`. This is for instance the right solution if the `next_one` routine implements a *random number generator*.

### 11.8.4 Relaxed memory model

*The reference for the commands introduced here can be found in section 12.8.3.*

#### 11.8.4.1 Thread synchronization

Let's do a *producer-consumer* model<sup>2</sup>. This can be implemented with sections, where one section, the producer, sets a flag when data is available, and the other, the consumer, waits until the flag is set.

---

2. This example is from Intel's excellent OMP course by Tim Mattson

```
#pragma omp parallel sections
{
    // the producer
    #pragma omp section
    {
        ... do some producing work ...
        flag = 1;
    }
    // the consumer
    #pragma omp section
    {
        while (flag==0) { }
        ... do some consuming work ...
    }
}
```

One reason this doesn't work, is that the compiler will see that the flag is never used in the producing section, and that is never changed in the consuming section, so it may optimize these statements, to the point of optimizing them away.

The producer then needs to do:

```
... do some producing work ...
#pragma omp flush
#pragma atomic write
    flag = 1;
#pragma omp flush(flag)
```

and the consumer does:

```
#pragma omp flush(flag)
while (flag==0) {
    #pragma omp flush(flag)
}
#pragma omp flush
```

This code strictly speaking has a *race condition* on the `flag` variable. It is better to use an `atomic` pragma here: the producer has

```
#pragma atomic write
    flag = 1;
```

and the consumer:

```
while (1) {
    #pragma omp flush(flag)
```

```
#pragma omp atomic read
    flag_read = flag
    if (flag_read==1) break;
}
```

### 11.8.5 Accelerators

In OpenMP 4.0 there is support for offloading work to an *accelerator* or *co-processor*:

```
#pragma omp target [clauses]
```

with clauses such as

- `data`: place data
- `update`: make data consistent between host and device

### 11.8.6 SIMD

OpenMP 4.0 has a way of indicating that a loop should not be arbitrarily divided over threads, but should be executed over SIMD lanes:

```
#pragma omp simd [clauses]
```

### 11.8.7 Overhead costs

Code parallelization ideally divides the running time of your program by the number of parallel processing entities. In practice, the following factors counteract this.

#### 11.8.7.1 Amdahl effects

Any code will have parts that are not parallelizable. Amdahl's law (see [HPSC-2.2.3](#)) quantizes the effect this has on parallel speedup. In an OpenMP code, the sections that are executed by a single thread will play the role of the sequential part.

#### 11.8.7.2 Thread overhead

At the start of an OpenMP program, a pool of threads is created. This incurs a one-time overhead that will probably be amortized over the total runtime.

Work sharing constructs act as if they create a new team of threads every time. In practice, the program probably keeps a pool of threads around that are dormant in between parallel sections. This means that there is no thread creation overhead associated with the start of a parallel section.

#### 11.8.7.3 Load balance

On the other hand, at the end of a work sharing construct there is a barrier, so an unbalanced load distribution will decrease the parallel efficiency. If loop iterations are not uniform in their running time, it may pay off to use dynamic rather than static scheduling.

On the other other hand, dynamic scheduling has overhead of its own, since it involves the operating system.

#### 11.8.7.4 Synchronization

Various synchronization constructs, such as critical sections, as well as dynamic loop scheduling, are realized through *operating system* functions. These are often quite costly, taking many thousands of cycles. Thus, the cost of a *critical sections* goes far beyond the Amdahl cost of the loss of parallelism. Critical sections should be used only if the parallel work far outweighs it.

## 11.9 Performance

The performance of an OpenMP code can be influenced by the following<sup>3</sup>::

- Amdahl effects
- Communication
- Data affinity
- Load imbalance
- Synchronization

Sequential code must clearly be kept to a minimum.

Cache coherence induces communication. Some of that is unavoidable, but see the next point.

Data is cached, so to minimize communication access to it should be as much as possible on the same core. This is known as affinity.

Load imbalance can be counteracted by using different loop schedules. The loop should be on as high a level as possible.

Barriers are a form of synchronization. They are expensive by themselves, and they expose load imbalance. Implicit barriers happen at the end of worksharing constructs; they can be removed with nowait.

---

3. This section is inspired by a presentation by Alexei Strelchenko.

## Chapter 12

### OpenMP Reference

This section gives reference information and illustrative examples of the use of OpenMP. While the code snippets given here should be enough, full programs can be found in the repository for this book <https://bitbucket.org/VictorEijkhout/parallel-computing-book>.

The definitive information on OpenMP can be found on <http://openmp.org/>; more tutorials can be found at <http://openmp.org/wp/resources/> where the one by Tim Mattson is particularly recommended.

#### 12.1 Basics

##### 12.1.1 OpenMP setup

If you use OMP library routines, you have to make them known to the compiler. In C/C++ you include the header file *omp.h*:

```
#include "omp.h"
```

In Fortran you use a module:

```
use omp_lib
```

##### 12.1.2 Directives

*This reference section gives the syntax for routines introduced in section 11.1.3.*

Directives in C/C++ are case-sensitive. Directives can be broken over multiple lines by escaping the line end.

Directives in Fortran start with a *sentinel*, commonly `!$omp`. If you break a directive over more than one line, all but the last line need to have a continuation character, and each line needs to have the sentinel:

```
!$OMP parallel do &
!%OMP    copyin(x),copyout(y)
```

The directives are case-insensitive. In *Fortran fixed-form* source files, `c$omp` and `*$omp` are allowed too.

### 12.1.3 Code and execution structure

*This reference section gives the syntax for routines introduced in section ??.*

Here are a couple of important concepts:

#### Definition 1

**structured block** An OpenMP directive is followed by an structured block; in C this is a single statement, a compound statement, or a block in braces; In Fortran it is delimited by the directive and its matching ‘end’ directive.

A structured block can not be jumped into, so it can not start with a labeled statement, or contain a jump statement leaving the block.

**construct** An OpenMP construct is the section of code starting with a directive and spanning the following structured block, plus in Fortran the end-directive. This is a lexical concept: it contains the statements directly enclosed, and not any subroutines called from them.

**region of code** A region of code is defined as all statements that are dynamically encountered while executing the code of an OpenMP construct. This is a dynamic concept: unlike a ‘construct’, it does include any subroutines that are called from the code in the structured block.

## 12.2 Parallel regions

*This reference section gives the syntax for routines introduced in section 11.1.4.*

The `parallel` pragma creates a team of threads from the current thread, and makes each thread execute the following block.

To test whether you are in a parallel region, use `omp_in_parallel`.

```
int omp_in_parallel() // C
LOGICAL omp_in_parallel() ! F
```

### 12.2.1 Threading

The number of threads in the team is controlled by `OMP_NUM_THREADS` or the `num_threads` clause on the `parallel` directive.

The number of threads used can differ between parallel regions. This is known as *dynamic mode*. You can set this with `omp_set_dynamic` and query with `omp_get_dynamic`.

If a region should be executed serially under certain conditions, the `if` clause can be used:

```
#pragma omp parallel if (n>1000)
#pragma omp for
for (i=0; i<n; i++) {
    ...
}
```

```
// hellocount.c
#pragma omp parallel
{
    int mythread,nthreads;
    nthreads = omp_get_num_threads();
    mythread = omp_get_thread_num();
    printf("Hello from %d out of %d\n",mythread,nthreads);
}
```

### 12.2.2 Combining with worksharing

If the structured block consists only of a `for/do` or `sections` worksharing construct, you can use the combined directives `parallel for`, `parallel do`, `parallel sections`.

## 12.3 Worksharing

*This reference section gives the syntax for routines introduced in section 11.2.*

The OpenMP *worksharing constructs* serve to distribute work over threads.

### 12.3.1 Loop parallelism

*This reference section gives the syntax for routines introduced in section 11.2.1.*

The `for` (C) and `do` (Fortran) directives tell OpenMP to distribute the iterations of a loop over the threads of a team. These pragmas do not create a team of threads: they take the current team of threads and divide the loop iterations over them. This means that the `omp for` or `omp do` directive needs to be inside a parallel region. It is also possible to have a combined `omp parallel for` or `omp parallel do` directive.

There are some restrictions on the loop: basically, OpenMP needs to be able to determine in advance how many iterations there will be.

- The loop can not contain `break`, `return`, `exit` statements, or `goto` to a label outside the loop.
- The `continue` (C) or `cycle` (F) statement is allowed.
- The index update has to be an increment (or decrement) by a fixed amount.
- The loop index variable is automatically private, and no changes to it inside the loop are allowed.

#### 12.3.1.1 Schedules

*This reference section gives the syntax for routines introduced in section 11.2.1.1.*

The schedule can be declared explicitly, set at runtime through the `OMP_SCHEDULE` environment variable, or left up to the runtime system by specifying `auto`. Especially in the last two cases you may want to enquire what schedule is currently being used with `omp_get_schedule`.

```
int omp_get_schedule(omp_sched_t * kind, int * modifier );
```

Its mirror call is `omp_set_schedule`, which sets the value that is used when `schedule` value `runtime` is used. It is in effect equivalent to setting the environment variable `OMP_SCHEDULE`.

```
void omp_set_schedule (omp_sched_t kind, int modifier);
```

Here are the various schedules you can set with the `schedule` clause:

**affinity** Set by using value `omp_sched_affinity`

**auto** The schedule is left up to the implementation. Set by using value `omp_sched_auto`

**dynamic** value: 2. The modifier parameter is the *chunk* size; default 1. Set by using value `omp_sched_dynamic`

**guided** Value: 3. The modifier parameter is the *chunk* size. Set by using value `omp_sched_guided`

**runtime** Use the value of the `OMP_SCHEDULE` environment variable. Set by using value `omp_sched_runtime`

**static** value: 1. The modifier parameter is the *chunk* size. Set by using value `omp_sched_static`

### 12.3.1.2 Ordered loop iterations

If loop iterations contain code that needs to be executed in the sequence of iterates, the `ordered` clause can be used:

```
#pragma omp for ordered private(t)
for (i=0; i<N; i++) {
    t = F(i) // some expensive calculation
#pragma omp ordered
    a[i+1] = a[i] + t
}
```

Note the separate `ordered` construct inside the loop. Without that, nothing will actually be ordered.

### 12.3.1.3 Reductions

This reference section gives the syntax for routines introduced in section 11.4.

Often, the threads in a team compute partial results which need to be combined, for instance in an addition or multiplication. This combination is known as a *reduction*. You can implement this by using partial results in each thread (see the examples in section 11.4), or by using an atomic update of a shared variable, but the easiest way is to use the `reduction` clause.

The `reduction` clause can be added to various constructs:

- Most commonly it is added to a `for` loop; the results of all iterations are then combined.

- Similarly, it can be added to a `section` construct; the results of the individual `section` blocks are then combined.
- You can even add it directly to the `parallel` directive; this combines the results from the threads in the team that is created.

```
int sum;
#pragma omp parallel for reduction(+:sum)
for (i=0; i<N; i++)
    sum = sum + f(i);
```

**Exercise 12.1.** Write a program to test the fact that the partial results are initialized to the unit of the reduction operator.

Arithmetic reductions: `+`, `*`, `-`, `max`, `min`

Logical operator reductions: `&`, `&&`, `|`, `||`, `^`

Fortran additionally has `max` and `min`.

### 12.3.2 Master and single

The `master` and `single` constructs are quite similar. They both indicate that a structured block is to be executed by only a single thread, and that all other threads that encounter the block skip it. However, the `master` pragma assigns the execution of the block to a specific thread: the master thread. Therefore, it is technically not a worksharing construct, and thus there is no barrier at the end of it.

The `single` pragma does have a barrier, which can be removed with `nowait`.

```
$!omp parallel
..
$!omp single
..
$!omp end single nowait
..
$!omp end parallel
```

### 12.3.3 Sections

```
#pragma omp parallel
{
#pragma omp sections
{
#pragma omp section
    do_thing_A();
#pragma omp section
    do_thing_B();
```

```
#pragma omp section
    do_thing_C();
}
}
```

### 12.3.4 Fortran workshare

*This reference section gives the syntax for routines introduced in section 11.2.4.*

The workshare directive exists only in Fortran. It can be used to parallelize the implied loops in *array syntax*, as well as *forall* loops.

## 12.4 Controlling thread data

*This reference section gives the syntax for routines introduced in section 11.3.*

### 12.4.1 Shared data

Data that existed in the master thread of a team is shared between the team. This is default behaviour. The clause `shared` can be used for completeness, for instance if `default (none)` is declared.

While shared data is readable and writable by every thread, their view of data may not always be consistent. Therefore, reads should be preceded by a `flush` command. Fortunately, in many cases this is done by default; see section 12.8.3.

### 12.4.2 Private data

*This reference section gives the syntax for routines introduced in section 11.3.2.*

Data that is declared private with the `private` directive is put on a separate *stack per thread*. The OpenMP standard does not dictate the size of these stacks, but beware of *stack overflow*. A typical default is a few megabyte; you can control it with the environment variable `OMP_STACKSIZE`

`firstprivate lastprivate copyin`

### 12.4.3 Defaults

You can add clauses to an `omp parallel` pragma to specify explicitly what variables are private and what variables shared. Note the cases with default behaviour:

- Loop variables in an `omp for` are private;
- Local variables in the parallel region are private.

You can alter this default behaviour with the `default` clause:

- The `none` option is good for debugging, because it forces you to specify for each variable in the parallel region whether it's private or shared.
- The `shared` clause means that any private variables need to be declared explicitly.
- The `private` clause means that any shared variables need to be declared explicitly. This value is not available in C.

#### 12.4.4 Threadprivate

*The reference for the commands introduced here can be found in section ??.*

Variables declared outside an OpenMP parallel region can be declared `threadprivate`: each thread will get a private copy as if in a parallel region. However, the variable's parallel lifetime is not limited to a parallel region, but will instead be governed by ordinary scoping rules.

```
int i; double d;  
#pragma omp threadprivate(i,d)
```

*Fortran note* Common blocks can be made thread-private with the syntax

```
$!OMP threadprivate( /blockname/ )
```

### 12.5 Synchronization

- `atomic` Update of a single memory location. Only certain specified syntax patterns are supported. This was added in order to be able to use hardware support for atomic updates.
- `barrier`: section 12.5.1
- `critical`: section 12.5.2
- `ordered`
- `locks`: section 12.5.3
- `flush`: section 12.8.3
- `nowait`

#### 12.5.1 Barriers

A barrier is a location in the code that needs to be reached by all threads before any of them can continue. There is a directive for declaring an explicit barrier, but there are also implicit barriers associated with certain other directives.

##### 12.5.1.1 Explicit barriers

*This reference section gives the syntax for routines introduced in section 11.5.1.*

You saw an example above where explicit barriers are needed. You can sometimes replace an explicit barrier by an implicit one:

```
#pragma omp parallel  
{  
    // do something  
#pragma omp barrier  
    // do something else  
}
```

is equivalent to

```
#pragma omp parallel
{
    // do something
}
#pragma omp parallel
    // do something else
}
```

but the second code has more overhead in creating the team of threads a second time.

#### 12.5.1.2 Implicit barriers

At the end of a parallel region the team of threads is dissolved and only the master thread continues. Therefore, there is an *implicit barrier at the end of a parallel region*.

There is some *barrier behaviour* associated with `omp for` loops and other *worksharing constructs* (see section ??). For instance, there is an *implicit barrier* at the end of the loop. This barrier behaviour can be cancelled with the `nowait` clause.

You will often see the idiom

```
#pragma omp parallel
{
    #pragma omp for nowait
        for (i=0; i<N; i++)
            a[i] = // some expression
    #pragma omp for
        for (i=0; i<N; i++)
            b[i] = ..... a[i] .....
```

Here the `nowait` clause implies that threads can start on the second loop while other threads are still working on the first. Since the two loops use the same schedule here, an iteration that uses `a[i]` can indeed rely on it that that value has been computed.

#### 12.5.2 Critical sections

*This reference section gives the syntax for routines introduced in section 11.5.2.1.*

The pragmas `critical` and `atomic` are two ways to indicate that a section of code can only be executed by one thread at a time.

```
#pragma omp critical [ (name) ] new-line
    structured-block
```

Not required to be in a parallel region?

### 12.5.3 Locks

Create/destroy:

```
void omp_init_lock(omp_lock_t *lock);  
void omp_destroy_lock(omp_lock_t *lock);
```

Set and release:

```
void omp_set_lock(omp_lock_t *lock);  
void omp_unset_lock(omp_lock_t *lock);
```

Since the set call is blocking, there is also

```
omp_test_lock();
```

Unsetting a lock needs to be done by the thread that set it.

Lock operations implicitly have a flush.

## 12.6 Internal control variables

OpenMP has a number of settings that can be set through *environment variables*, and both queried and set through *library routines*. These settings are called *Internal Control Variables (ICVs)*: an OpenMP implementation behaves as if there is an internal variable storing this setting.

First, there are 4 ICVs that behave as if each thread has its own copy of them. The default is implementation-defined unless otherwise noted.

- It may be possible to adjust dynamically the number of threads for a parallel region. Variable: OMP\_DYNAMIC; routines: `omp_set_dynamic`, `omp_get_dynamic`.
- If a code contains *nested parallel regions*, the inner regions may create new teams, or they may be executed by the single thread that encounters them. Variable: OMP\_NESTED; routines `omp_set_nested`, `omp_get_nested`. Allowed values are TRUE and FALSE; the default is false.
- The number of threads used for an encountered parallel region can be controlled. Variable: OMP\_NUM\_THREADS; routines `omp_set_num_threads`, `omp_get_max_threads`.
- The schedule for a parallel loop can be set. Variable: OMP\_SCHEDULE; routines `omp_set_schedule`, `omp_get_schedule`.

Non-obvious syntax:

```
export OMP_SCHEDULE="static,100"
```

Other settings:

- `omp_get_num_threads`: query the number of threads active at the current place in the code; this can be lower than what was set with `omp_set_num_threads`. For a meaningful answer, this should be done in a parallel region.

- `omp_get_thread_num`
- `omp_in_parallel`: test if you are in a parallel region
- `omp_get_num_procs`: query the physical number of cores available.

Other environment variables:

- `OMP_STACKSIZE` controls the amount of space that is allocated as per-thread stack; the space for private variables.
- `OMP_WAIT_POLICY` determines the behaviour of threads that wait, for instance for *critical section*:
  - `ACTIVE` puts the thread in a *spin-lock*, where it actively checks whether it can continue;
  - `PASSIVE` puts the thread to sleep until the Operating System (OS) wakes it up.
- The ‘active’ strategy uses CPU while the thread is waiting; on the other hand, activating it after the wait is instantaneous. With the ‘passive’ strategy, the thread does not use any CPU while waiting, but activating it again is expensive. Thus, the passive strategy only makes sense if threads will be waiting for a (relatively) long time.
- `OMP_PROC_BIND` with values `TRUE` and `FALSE` can bind threads to a processor. On the one hand, doing so can minimize data movement; on the other hand, it may increase load imbalance.

## 12.7 Tasks

*This reference section gives the syntax for routines introduced in section 11.6.*

OpenMP v4 has a `depend` clause.

```
#pragma omp task depend(dependency-type: list)
```

## 12.8 Stuff

### 12.8.1 Timing

*This reference section gives the syntax for routines introduced in section 11.8.1.*

To do OpenMP timing you can use any system utility; however there is a dedicated routine `omp_get_wtime` that express the time since some starting point as a double:

```
double omp_get_wtime(void);
```

The starting point is arbitrary and is different for each program run; however, in one run it is identical for all threads.

To measure a time difference:

```
double tstart,tend,duration;
tstart = omp_get_wtime();
// do stuff
tend = omp_get_wtime();
duration = tend-tstart;
```

The timer resolution is given by:

```
double omp_get_wtick(void);
```

### 12.8.2 Affinity

For performance it can be a good idea to bind threads to specific processors or cores. OpenMP (as of version 3.1) has a mechanism for *thread affinity*: OMP\_PROC\_BIND

```
export OMP_PROC_BIND=true
```

Apart from this, compilers can have proprietary mechanism; e.g., for the intel compiler the variable is

```
export KMP_AFFINITY=compact,0
```

for the sun compiler:

```
export SUNW_MP_PROCBIND=TRUE
```

for gcc (pre-openmp 3.1)

```
export GOMP_CPU_AFFINITY=0-63
```

### 12.8.3 Relaxed memory model

*This reference section gives the syntax for routines introduced in section 11.8.4.*

flush

- There is an implicit flush of all variables at the start and end of a *parallel region*.
- There is a flush at each barrier, whether explicit or implicit, such as at the end of a *work sharing*.
- At entry and exit of a *critical section*
- When a *lock* is set or unset.

# Chapter 13

## OpenMP Review

### 13.1 Concepts review

#### 13.1.1 Basic concepts

- process / thread / thread team
- threads / cores / tasks
- directives / library functions / environment variables

#### 13.1.2 Parallel regions

execution by a team

#### 13.1.3 Work sharing

- loop / sections / single / workshare
- implied barrier
- loop scheduling, reduction
- sections
- single vs master
- (F) workshare

#### 13.1.4 Data scope

- shared vs private, C vs F
- loop variables and reduction variables
- default declaration
- firstprivate, lastprivate

#### 13.1.5 Synchronization

- barriers, implied and explicit
- nowait
- critical sections
- locks, difference with critical

#### 13.1.6 Tasks

- generation vs execution
- dependencies

## 13.2 Review questions

### 13.2.1 Directives

What do the following program output?

```
int main() {
    printf("procs %d\n",
        omp_get_num_procs());
    printf("threads %d\n",
        omp_get_num_threads());
    printf("num %d\n",
        omp_get_thread_num());
    return 0;
}
```

```
int main() {
#pragma omp parallel
{
    printf("procs %d\n",
        omp_get_num_procs());
    printf("threads %d\n",
        omp_get_num_threads());
    printf("num %d\n",
        omp_get_thread_num());
}
return 0;
}
```

**Program** main  
**use** omp.lib  
**print** \*, "Procs:", &  
 omp\_get\_num\_procs()  
**print** \*, "Threads:", &  
 omp\_get\_num\_threads()  
**print** \*, "Num:", &  
 omp\_get\_thread\_num()  
**End Program**

**Program** main  
**use** omp.lib  
*!\$OMP parallel*  
**print** \*, "Procs:", &  
 omp\_get\_num\_procs()  
**print** \*, "Threads:", &  
 omp\_get\_num\_threads()  
**print** \*, "Num:", &  
 omp\_get\_thread\_num()  
*!\$OMP end parallel*  
**End Program**

### 13.2.2 Parallelism

Can the following loops be parallelized? If so, how? (Assume that all arrays are already filled in, and that there are no out-of-bounds errors.)

```
// variant #1
for (i=0; i<N; i++) {
    x[i] = a[i]+b[i+1];
    a[i] = 2*x[i] + c[i+1];
}
```

```
// variant #3
for (i=1; i<N; i++) {
    x[i] = a[i]+b[i+1];
    a[i] = 2*x[i-1] + c[i+1];
}
```

```
// variant #2
for (i=0; i<N; i++) {
    x[i] = a[i]+b[i+1];
    a[i] = 2*x[i+1] + c[i+1];
}
```

```
// variant #4
for (i=1; i<N; i++) {
    x[i] = a[i]+b[i+1];
    a[i+1] = 2*x[i-1] + c[i+1];
}
```

```
! variant #1
do i=1,N
    x(i) = a(i)+b(i+1)
    a(i) = 2*x(i) + c(i+1)
end do
```

```
! variant #3
do i=2,N
    x(i) = a(i)+b(i+1)
    a(i) = 2*x(i-1) + c(i+1)
end do
```

```
! variant #2
do i=1,N
    x(i) = a(i)+b(i+1)
    a(i) = 2*x(i+1) + c(i+1)
end do
```

```
! variant #3
do i=2,N
    x(i) = a(i)+b(i+1)
    a(i+1) = 2*x(i-1) + c(i+1)
end do
```

### 13.2.3 Data and synchronization

#### 13.2.3.1

What is the output of the following fragments? Assume that there are four threads.

```
// variant #1
int nt;
#pragma omp parallel
{
    nt = omp_get_thread_num();
    printf("thread_number: %d\n", nt);
}
```

```
// variant #2
int nt;
#pragma omp parallel private(nt)
{
    nt = omp_get_thread_num();
    printf("thread_number: %d\n", nt);
}
```

```
// variant #3
int nt;
#pragma omp parallel
{
#pragma omp single
{
    nt = omp_get_thread_num();
    printf("thread_number: %d\n", nt);
}}
```

```
// variant #4
int nt;
#pragma omp parallel
{
#pragma omp master
{
    nt = omp_get_thread_num();
    printf("thread_number: %d\n", nt);
}}
```

```
// variant #5
int nt;
#pragma omp parallel
{
#pragma omp critical
{
    nt = omp_get_thread_num();
    printf("thread_number: %d\n", nt);
}}
```

```
! variant #1
integer nt
!$OMP parallel
    nt = omp_get_thread_num()
    print *, "thread_number:", nt
!$OMP end parallel
```

```
! variant #2
integer nt
!$OMP parallel private(nt)
    nt = omp_get_thread_num()
    print *, "thread_number:", nt
!$OMP end parallel
```

```
! variant #3
integer nt
!$OMP parallel
!$OMP single
    nt = omp_get_thread_num()
    print *, "thread_number:", nt
!$OMP end single
!$OMP end parallel
```

```

! variant #4
integer nt
!$OMP parallel
!$OMP master
    nt = omp_get_thread_num()
    print *, "thread_number:", nt
!$OMP end master
!$OMP end parallel

```

```

! variant #5
integer nt
!$OMP parallel
!$OMP critical
    nt = omp_get_thread_num()
    print *, "thread_number:", nt
!$OMP end critical
!$OMP end parallel

```

### 13.2.3.2

The following is an attempt to parallelize a serial code. Assume that all variables and arrays are defined. What errors and potential problems do you see in this code? How would you fix them?

```

#pragma omp parallel
{
    x = f();
    #pragma omp for
    for (i=0; i<N; i++)
        y[i] = g(x, i);
    z = h(y);
}

```

```

!$OMP parallel
    x = f()
    !$OMP do
        do i=1,N
            y(i) = g(x, i)
        end do
    !$OMP end do
    z = h(y)
!$OMP end parallel

```

### 13.2.4 Reductions

#### 13.2.4.1

Is the following code correct? Is it efficient? If not, can you improve it?

```
#pragma omp parallel shared(r)
{
    int x;
    x = f(omp_get_thread_num());
#pragma omp critical
    r += f(x);
}
```

#### 13.2.4.2

Compare two fragments:

<i>// variant 1</i>	<i>// variant 2</i>
<b>#pragma</b> omp parallel reduction(+:s) <b>#pragma</b> omp <b>for</b> <b>for</b> (i=0; i<N; i++) s += f(i);	<b>#pragma</b> omp parallel <b>#pragma</b> omp <b>for</b> reduction(+:s) <b>for</b> (i=0; i<N; i++) s += f(i);

<i>! variant 1</i>	<i>! variant 2</i>
<i>!\$OMP parallel reduction(+:s)</i> <i>!\$OMP do</i> <b>do</b> i=1,N s += f(i); <b>end do</b> <i>!\$OMP end do</i> <i>!\$OMP end parallel</i>	<i>!\$OMP parallel</i> <i>!\$OMP do reduction(+:s)</i> <b>do</b> i=1,N s += f(i); <b>end do</b> <i>!\$OMP end do</i> <i>!\$OMP end parallel</i>

Do they compute the same thing?

### 13.2.5 Data scope

The following program is supposed to initialize as many rows of the array as there are threads.

```
int main() {
    int i , icount , iarray [100][100];
    icount = -1;
#pragma omp parallel private(i)
    {
#pragma omp critical
        { icount++;
            for (i=0; i<100; i++)
                iarray [icount][i] = 1;
        }
    return 0;
}
```

```
Program main
    integer :: i , icount , iarray (100,100)
    icount = 0
!$OMP parallel private(i)
!$OMP critical
    icount = icount + 1
!$OMP end critical
    do i=1,100
        iarray(icount,i) = 1
    end do
!$OMP end parallel
End program
```

Describe the behaviour of the program, with argumentation,

- as given;
- if you add a clause `private(icount)` to the `parallel` directive;
- if you add a clause `firstprivate(icount)`.

What do you think of this solution:

```
#pragma omp parallel private(i) shared(icount)
    {
#pragma omp critical
        { icount++;
            for (i=0; i<100; i++)
                iarray [icount][i] = 1;
        }
    }
return 0;
}
```

```
!$OMP parallel private(i) shared(icount)
!$OMP critical
    icount = icount+1
    do i=1,100
        iarray (icount,i) = 1
    end do
!$OMP critical
!$OMP end parallel
```

### 13.2.6 Tasks

Fix two things in the following example:

```
#pragma omp parallel
#pragma omp single
{
    int x,y,z;
#pragma omp task
```

```
x = f();
#pragma omp task
y = g();
#pragma omp task
z = h();
printf("sum=%d\n",x+y+z);
```

}

```

integer :: x,y,z
!$OMP parallel
!$OMP single

!$OMP task
x = f()
!$OMP end task

!$OMP task
y = g()
!$OMP end task

!$OMP task
z = h()
!$OMP end task

print *, "sum=", x+y+z
!$OMP end single
!$OMP end parallel

```

### 13.2.7 Scheduling

Compare these two fragments. Do they compute the same result? What can you say about their efficiency?

```

#pragma omp parallel
#pragma omp single
{
    for ( i=0; i<N; i++) {
        #pragma omp task
        x[ i ] = f( i )
    }
    #pragma omp taskwait
}

```

```

#pragma omp parallel
#pragma omp for schedule(dynamic)
{
    for ( i=0; i<N; i++) {
        x[ i ] = f( i )
    }
}

```

How would you make the second loop more efficient? Can you do something similar for the first loop?

## **PART III**

### **THE REST**

## Chapter 14

### Random number generation

Here is how you initialize the random number generator uniquely on each process:

C:

```
// Initialize the random number generator
srand((int)(mytid*(double)RAND_MAX/ntids));
// compute a random number
randomfraction = (rand() / (double)RAND_MAX);
```

Fortran:

```
integer :: randsize
integer,allocatable,dimension(:) :: randseed
real :: random_value

call random_seed(size=randsize)
allocate(randseed(randsize))
do i=1,randsize
    randseed(i) = 1023*mytid
end do
call random_seed(put=randseed)
```

## Chapter 15

### Hybrid computing

#### 15.1 Discussion

Hybrid computing decreases the number of messages.

On the other hand it makes the run more synchronous.

New version of Amdahl: sections that MPI-parallel but not OpenMP-parallel.

Allows overdecomposition on the node.

#### 15.2 Hybrid MPI-plus-threads execution

In hybrid execution, the main question is whether all threads are allowed to make MPI calls. To determine this, replace the `MPI_Init` call by `MPI_Init_thread`:

```
int MPI_Init_thread  
  ( int *argc, char ***argv, int required, int *provided )
```

Here the `required` and `provided` parameters can take the following values:

**`MPI_THREAD_SINGLE`** Only a single thread will execute.

**`MPI_THREAD_FUNNELLED`** The program may use multiple threads, but only the main thread will make MPI calls.

**`MPI_THREAD_SERIAL`** The program may use multiple threads, all of which may make MPI calls, but there will never be simultaneous MPI calls in more than one thread.

**`MPI_THREAD_MULTIPLE`** Multiple threads may MPI calls, without restrictions.

The `mpirun` program usually propagates *environment variables*, so the value of `OMP_NUM_THREADS` when you call `mpirun` will be seen by each MPI process.

- It is possible to use blocking sends in threads, and let the threads block. This does away with the need for polling.
- You can not send to a thread number.

Exercise 15.1. Consider the 2D heat equation and explore the mix of MPI/OpenMP parallelism:

- Give each node one MPI process that is fully multi-threaded.
- Give each core an MPI process and don't use multi-threading.

Discuss theoretically why the former can give higher performance. Implement both schemes as special cases of the general hybrid case, and run tests to find the optimal mix.

## **Chapter 16**

### **Support libraries**

ParaMesh

Global Arrays

PETSc

Hdf5 and Silo



## **PART IV**

### **TUTORIALS**

---

here are some tutorials

## 16.1 Debugging

When a program misbehaves, *debugging* is the process of finding out *why*. There are various strategies of finding errors in a program. The crudest one is debugging by print statements. If you have a notion of where in your code the error arises, you can edit your code to insert print statements, recompile, rerun, and see if the output gives you any suggestions. There are several problems with this:

- The edit/compile/run cycle is time consuming, especially since
- often the error will be caused by an earlier section of code, requiring you to edit, compile, and rerun repeatedly. Furthermore,
- the amount of data produced by your program can be too large to display and inspect effectively, and
- if your program is parallel, you probably need to print out data from all processors, making the inspection process very tedious.

For these reasons, the best way to debug is by the use of an interactive *debugger*, a program that allows you to monitor and control the behaviour of a running program. In this section you will familiarize yourself with *gdb*, which is the open source debugger of the *GNU* project. Other debuggers are proprietary, and typically come with a compiler suite. Another distinction is that *gdb* is a commandline debugger; there are graphical debuggers such as *ddd* (a frontend to *gdb*) or *DDT* and *TotalView* (debuggers for parallel codes). We limit ourselves to *gdb*, since it incorporates the basic concepts common to all debuggers.

In this tutorial you will debug a number of simple programs with *gdb* and *valgrind*. The files can be downloaded from <http://tinyurl.com/ISTC-debug-tutorial>.

### 16.1.1 Step 0: compiling for debug

You often need to recompile your code before you can debug it. A first reason for this is that the binary code typically knows nothing about what variable names corresponded to what memory locations, or what lines in the source to what instructions. In order to make the binary executable know this, you have to include the *symbol table* in it, which is done by adding the `-g` option to the compiler line.

Usually, you also need to lower the *compiler optimization level*: a production code will often be compiled with flags such as `-O2` or `-Xhost` that try to make the code as fast as possible, but for debugging you need to replace this by `-O0` ('oh-zero'). The reason is that higher levels will reorganize your code, making it hard to relate the execution to the source<sup>1</sup>.

### 16.1.2 Invoking gdb

There are three ways of using *gdb*: using it to start a program, attaching it to an already running program, or using it to inspect a *core dump*. We will only consider the first possibility.

Here is an example of how to start *gdb* with a program that has no arguments (Fortran users, use `hello.F`):

```
tutorials/gdb/c/hello.c
```

---

1. Typically, actual code motion is done by `-O3`, but at level `-O2` the compiler will inline functions and make other simplifications.

---

```

%% cc -g -o hello hello.c
# regular invocation:
%% ./hello
hello world
# invocation from gdb:
%% gdb hello
GNU gdb 6.3.50-20050815 # ..... version info
Copyright 2004 Free Software Foundation, Inc. .... copyright info ....
(gdb) run
Starting program: /home/eijkhout/tutorials/gdb/hello
Reading symbols for shared libraries +. done
hello world

Program exited normally.
(gdb) quit
%%

```

Important note: the program was compiled with the *debug flag* `-g`. This causes the *symbol table* (that is, the translation from machine address to program variables) and other debug information to be included in the binary. This will make your binary larger than strictly necessary, but it will also make it slower, for instance because the compiler will not perform certain optimizations<sup>2</sup>.

To illustrate the presence of the symbol table do

```

%% cc -g -o hello hello.c
%% gdb hello
GNU gdb 6.3.50-20050815 # ..... version info
(gdb) list

```

and compare it with leaving out the `-g` flag:

```

%% cc -o hello hello.c
%% gdb hello
GNU gdb 6.3.50-20050815 # ..... version info
(gdb) list

```

For a program with commandline input we give the arguments to the `run` command (Fortran users use `say.F`):

`tutorials/gdb/c/say.c`

```

%% cc -o say -g say.c
%% ./say 2

```

---

2. Compiler optimizations are not supposed to change the semantics of a program, but sometimes do. This can lead to the nightmare scenario where a program crashes or gives incorrect results, but magically works correctly with compiled with debug and run in a debugger.

```

hello world
hello world
%% gdb say
.... the usual messages ...
(gdb) run 2
Starting program: /home/eijkhout/tutorials/gdb/c/say 2
Reading symbols for shared libraries +. done
hello world
hello world

Program exited normally.

```

### 16.1.3 Finding errors

Let us now consider some programs with errors.

#### 16.1.3.1 C programs

tutorials/gdb/c/square.c

```

%% cc -g -o square square.c
%% ./square
5000
Segmentation fault

```

The *segmentation fault* (other messages are possible too) indicates that we are accessing memory that we are not allowed to, making the program abort. A debugger will quickly tell us where this happens:

```

%% gdb square
(gdb) run
50000

Program received signal EXC_BAD_ACCESS, Could not access memory.
Reason: KERN_INVALID_ADDRESS at address: 0x000000000000eb4a
0x00007fff824295ca in __svfscanf_l ()

```

Apparently the error occurred in a function `__svfscanf_l`, which is not one of ours, but a system function. Using the backtrace (or `bt`, also `where` or `w`) command we quickly find out how this came to be called:

```

(gdb) backtrace
#0 0x00007fff824295ca in __svfscanf_l ()
#1 0x00007fff8244011b in fscanf ()
#2 0x0000000100000e89 in main (argc=1, argv=0x7fff5fbfc7c0) at square.c:7

```

---

We take a close look at line 7, and see that we need to change nmax to &nmax.

There is still an error in our program:

```
(gdb) run  
50000  
  
Program received signal EXC_BAD_ACCESS, Could not access memory.  
Reason: KERN_PROTECTION_FAILURE at address: 0x000000010000f000  
0x000000010000ebe in main (argc=2, argv=0x7fff5fbfc7a8) at square1.c:9  
9           squares[i] = 1. / (i * i); sum += squares[i];
```

We investigate further:

```
(gdb) print i  
$1 = 11237  
(gdb) print squares[i]  
Cannot access memory at address 0x10000f000
```

and we quickly see that we forgot to allocate squares.

By the way, we were lucky here: this sort of memory errors is not always detected. Starting our programm with a smaller input does not lead to an error:

```
(gdb) run  
50  
Sum: 1.625133e+00  
  
Program exited normally.
```

### 16.1.3.2 Fortran programs

Compile and run the following program:

tutorials/gdb/f/square.F

It should abort with a message such as ‘Illegal instruction’. Running the program in gdb quickly tells you where the problem lies:

```
(gdb) run  
Starting program: tutorials/gdb//fsquare  
Reading symbols for shared libraries +++. done  
  
Program received signal EXC_BAD_INSTRUCTION, Illegal instruction/operand.  
0x000000010000da3 in square () at square.F:7  
7           sum = sum + squares(i)
```

We take a close look at the code and see that we did not allocate squares properly.

### 16.1.4 Memory debugging with Valgrind

Insert the following allocation of `squares` in your program:

```
squares = (float *) malloc( nmax*sizeof(float) );
```

Compile and run your program. The output will likely be correct, although the program is not. Can you see the problem?

To find such subtle memory errors you need a different tool: a memory debugging tool. A popular (because open source) one is *valgrind*; a common commercial tool is *purify*.

**tutorials/gdb/c/square1.c** Compile this program with `cc -o square1 square1.c` and run it with `valgrind square1` (you need to type the input value). You will lots of output, starting with:

```
%% valgrind square1
==53695== Memcheck, a memory error detector
==53695== Copyright (C) 2002-2010, and GNU GPL'd, by Julian Seward et al.
==53695== Using Valgrind-3.6.1 and LibVEX; rerun with -h for copyright info
==53695== Command: a.out
==53695==
10
==53695== Invalid write of size 4
==53695==   at 0x100000EB0: main (square1.c:10)
==53695==   Address 0x10027e148 is 0 bytes after a block of size 40 alloc'd
==53695==   at 0x1000101EF: malloc (vg_replace_malloc.c:236)
==53695==   by 0x100000E77: main (square1.c:8)
==53695==
==53695== Invalid read of size 4
==53695==   at 0x100000EC1: main (square1.c:11)
==53695==   Address 0x10027e148 is 0 bytes after a block of size 40 alloc'd
==53695==   at 0x1000101EF: malloc (vg_replace_malloc.c:236)
==53695==   by 0x100000E77: main (square1.c:8)
```

Valgrind is informative but cryptic, since it works on the bare memory, not on variables. Thus, these error messages take some exegesis. They state that a line 10 writes a 4-byte object immediately after a block of 40 bytes that was allocated. In other words: the code is writing outside the bounds of an allocated array. Do you see what the problem in the code is?

Note that valgrind also reports at the end of the program run how much memory is still in use, meaning not properly freed.

If you fix the array bounds and recompile and rerun the program, valgrind still complains:

```
==53785== Conditional jump or move depends on uninitialised value(s)
==53785==   at 0x10006FC68: __ dtoa (in /usr/lib/libSystem.B.dylib)
==53785==   by 0x10003199F: __ vfprintf (in /usr/lib/libSystem.B.dylib)
==53785==   by 0x1000738AA: vfprintf_l (in /usr/lib/libSystem.B.dylib)
==53785==   by 0x1000A1006: printf (in /usr/lib/libSystem.B.dylib)
==53785==   by 0x100000EF3: main (in ./square2)
```

---

Although no line number is given, the mention of `printf` gives an indication where the problem lies. The reference to an ‘uninitialized value’ is again cryptic: the only value being output is `sum`, and that is not uninitialized: it has been added to several times. Do you see why valgrind calls `is uninitialized` all the same?

### 16.1.5 Stepping through a program

Often the error in a program is sufficiently obscure that you need to investigate the program run in detail. Compile the following program

tutorials/gdb/c/roots.c and run it:

```
%% ./roots
sum: nan
```

Start it in `gdb` as follows:

```
%% gdb roots
GNU gdb 6.3.50-20050815 (Apple version gdb-1469) (Wed May 5 04:36:56 UTC 2005)
Copyright 2004 Free Software Foundation, Inc.
...
(gdb) break main
Breakpoint 1 at 0x100000ea6: file root.c, line 14.
(gdb) run
Starting program: tutorials/gdb/c/roots
Reading symbols for shared libraries +. done

Breakpoint 1, main () at roots.c:14
14          float x=0;
```

Here you have done the following:

- Before calling `run` you set a *breakpoint* at the main program, meaning that the execution will stop when it reaches the main program.
- You then call `run` and the program execution starts;
- The execution stops at the first instruction in `main`.

If execution is stopped at a breakpoint, you can do various things, such as issuing the `step` command:

```
Breakpoint 1, main () at roots.c:14
14          float x=0;
(gdb) step
15          for (i=100; i>-100; i--)
(gdb)
16          x += root(i);
(gdb)
```

(if you just hit return, the previously issued command is repeated). Do a number of steps in a row by hitting return. What do you notice about the function and the loop?

Switch from doing step to doing next. Now what do you notice about the loop and the function?

Set another breakpoint: break 17 and do cont. What happens?

Rerun the program after you set a breakpoint on the line with the sqrt call. When the execution stops there do where and list.

- If you set many breakpoints, you can find out what they are with info breakpoints.
- You can remove breakpoints with delete n where n is the number of the breakpoint.
- If you restart your program with run without leaving gdb, the breakpoints stay in effect.
- If you leave gdb, the breakpoints are cleared but you can save them: save breakpoints <file>. Use source <file> to read them in on the next gdb run.

### 16.1.6 Inspecting values

Run the previous program again in gdb: set a breakpoint at the line that does the sqrt call before you actually call run. When the program gets to line 8 you can do print n. Do cont. Where does the program stop?

If you want to repair a variable, you can do set var=value. Change the variable n and confirm that the square root of the new value is computed. Which commands do you do?

If a problem occurs in a loop, it can be tedious keep typing cont and inspecting the variable with print. Instead you can add a condition to an existing breakpoint: the following:

```
condition 1 if (n<0)
```

or set the condition when you define the breakpoint:

```
break 8 if (n<0)
```

Another possibility is to use ignore 1 50, which will not stop at breakpoint 1 the next 50 times.

Remove the existing breakpoint, redefine it with the condition n<0 and rerun your program. When the program breaks, find for what value of the loop variable it happened. What is the sequence of commands you use?

### 16.1.7 Parallel debugging

Debugging parallel programs is harder than sequential programs, because every sequential bug may show up, plus a number of new types, caused by the interaction of the various processes.

Here are a few possible parallel bugs:

- Processes can deadlock because they are waiting for a message that never comes. This typically happens with blocking send/receive calls due to an error in program logic.
- If an incoming message is unexpectedly larger than anticipated, a memory error can occur.

- 
- A collective call will hang if somehow one of the processes does not call the routine.

There are few low-budget solutions to parallel debugging. The main one is to create an xterm for each process. We will describe this next. There are also commercial packages such as *DDT* and *TotalView*, that offer a GUI. They are very convenient but also expensive. The *Eclipse* project has a parallel package, *Eclipse PTP*, that includes a graphic debugger.

#### 16.1.7.1 MPI debugging with *gdb*

You can not run parallel programs in *gdb*, but you can start multiple *gdb* processes that behave just like MPI processes! The command

```
mpirun -np <NP> xterm -e gdb ./program
```

create a number of xterm windows, each of which execute the commandline *gdb ./program*. And because these xterms have been started with *mpirun*, they actually form a communicator.

#### 16.1.7.2 Full-screen parallel debugging with *DDT*

In this tutorial you will run and diagnose a few incorrect MPI programs using DDT. You can start a session with *ddt yourprogram &*, or use *File > New Session > Run* to specify a program name, and possibly parameters. In both cases you get a dialog where you can specify program parameters. It is also important to check the following:

- You can specify the number of cores here;
- It is usually a good idea to turn on memory checking;
- Make sure you specify the right MPI.

When DDT opens on your main program, it halts at the *MPI\_Init* statement, and need to press the forward arrow, top left of the main window.

##### 16.1.7.2.1 **Problem1**

This program has every process independently generate random numbers, and if the number meets a certain condition, stops execution. There is no problem with this code as such, so let's suppose you simply want to monitor its execution.

- Compile *abort.c*. Don't forget about the *-g -O0* flags; if you use the makefile they are included automatically.
- Run the program with DDT, you'll see that it concludes successfully.
- Set a breakpoint at the *Finalize* statement in the subroutine, by clicking to the left of the line number. Now if you run the program you'll get a message that all processes are stopped at a breakpoint. Pause the execution.
- The 'Stacks' tab will tell you that all processes are the same point in the code, but they are not in fact in the same iteration.
- You can for instance use the 'Input/Output' tabs to see what every process has been doing.
- Alternatively, use the variables pane on the right to examine the *it* variable. You can do that for individual processes, but you can also control click on the *it* variable and choose *View as Array*. Set up the display as a one-dimensional array and check the iteration numbers.

- Activate the barrier statement and rerun the code. Make sure you have no breakpoints. Reason that the code will not complete, but just hang.
- Hit the general Pause button. Now what difference do you see in the ‘Stacks’ tab?

**16.1.7.2.2 Problem2** Compile `problem1.c` and run it in DDT. You’ll get a dialog warning about an abort.

- Pause the program in the dialog. Notice that only the root process is paused. If you want to inspect other processes, press the general pause button. Do this.
- In the bottom panel click on `Stacks`. This gives you the ‘call stack’, which tells you what the processes were doing when you paused them. Where is the root process in the execution? Where are the others?
- From the call stack it is clear what the error was. Fix it and rerun with `File > Restart Session`.

**16.1.7.2.3 Problem2**

### 16.1.8 Further reading

A good tutorial: <http://www.dirac.org/linux/gdb/>.

Reference manual: [http://www.ofb.net-gnu/gdb/gdb\\_toc.html](http://www.ofb.net-gnu/gdb/gdb_toc.html).

---

## 16.2 Tracing

### 16.2.1 TAU profiling and tracing

TAU <http://www.cs.uoregon.edu/Research/tau/home.php> is a utility for profiling and tracing your parallel programs. Profiling is the gathering and displaying of bulk statistics, for instance showing you which routines take the most time, or whether communication takes a large portion of your runtime. When you get concerned about performance, a good profiling tool is indispensable.

Tracing is the construction and displaying of time-dependent information on your program run, for instance showing you if one process lags behind others. For understanding a program's behaviour, and the reasons behind profiling statistics, a tracing tool can be very insightful.

TAU works by adding *instrumentation* to your code: in effect it is a source-to-source translator that takes your code and turns it into one that generates run-time statistics. Doing this instrumentation is fortunately simple: start by having this code fragment in your makefile:

```
ifdef TACC_TAU_DIR
    CC = tau_cc.sh
else
    CC = mpicc
endif

% : %.c
${CC} -o $@ $^
```

To use TAU, do `module load tau`. You have to set the environment variable `TAU_TRACE` to 1; it's advisable to set `TRACEDIR` to some directory for all the TAU output. Likewise set `TAU_PROFILE` to 1 and set `PROFILEDIR`.

**PART V**

**PROJECTS, INDEX**

# Chapter 17

## Class projects

### 17.1 A Style Guide to Project Submissions

Here are some guidelines for how to submit assignments and projects. As a general rule, consider programming as an experimental science, and your writeup as a report on some tests you have done: explain the problem you're addressing, your strategy, your results.

**Structure of your writeup** Most of the exercises in this book test whether you are able to code the solution to a certain problem. That does not mean that turning in the code is sufficient, nor code plus sample output. Turn in a writeup in pdf form that was generated from a text processing program such as Word or (preferably) L<sup>A</sup>T<sub>E</sub>X (for a tutorial, see HPSC-35). Your writeup should have

- The relevant fragments of your code,
- an explanation of your algorithms or solution strategy,
- a discussion of what you observed,
- graphs of runtimes and TAU plots; see 16.2.

*Observe, measure, hypothesize, deduce* In most applications of computing machinery we care about the efficiency with which we find the solution. Thus, make sure that you do measurements. In general, make observations that allow you to judge whether your program behaves the way you would expect it to.

Quite often your program will display unexpected behaviour. It is important to observe this, and hypothesize what the reason might be for your observed behaviour.

*Including code* If you include code samples in your writeup, make sure they look good. For starters, use a mono-spaced font. In L<sup>A</sup>T<sub>E</sub>X, you can use the `verbatim` environment or the `verbbatiminput` command. In that section option the source is included automatically, rather than cut and pasted. This is to be preferred, since your writeup will stay current after you edit the source file.

Including whole source files makes for a long and boring writeup. The code samples in this book were generated as follows. In the source files, the relevant snippet was marked as

```
... boring stuff
#pragma samplex
.. interesting! ..
#pragma end
... more boring stuff
```

The files were then processed with the following command line (actually, included in a makefile, which requires doubling the dollar signs):

```
for f in *.{c,cxx,h} ; do
    cat $x | awk 'BEGIN {f=0}
                    /#pragma end/ {f=0}
                    f==1 {print $0 > file}
                    /pragma/ {f=1; file=$2}
'
done
```

which gives (in this example) a file `samplex`. Other solutions are of course possible.

**Code formatting** Code without proper indentation is very hard to read. Fortunately, most editors have some knowledge of the syntax of the most popular languages. The `emacs` editor will, most of the time, automatically activate the appropriate mode based on the file extension. If this does not happen, you can activate a mode by `ESC x fortran-mode` et cetera, or by putting the string `--*-- fortran --*--` in a comment on the first line of your file.

The `vi` editor also has syntax support: use the commands `:syntax on` to get syntax colouring, and `:set cindent` to get automatic indentation while you're entering text. Some of the more common questions are addressed in <http://stackoverflow.com/questions/97694/auto-indent-spaces-with-c-indent>

**Running your code** A single run doesn't prove anything. For a good report, you need to run your code for more than one input dataset (if available) and in more than one processor configuration. When you choose problem sizes, be aware that an average processor can do a billion operations per second: you need to make your problem large enough for the timings to rise above the level of random variations and startup phenomena.

When you run a code in parallel, beware that on clusters the behaviour of a parallel code will always be different between one node and multiple nodes. On a single node the MPI implementation is likely optimized to use the shared memory. This means that results obtained from a single node run will be unrepresentative. In fact, in timing and scaling tests you will often see a drop in (relative) performance going from one node to two. Therefore you need to run your code in a variety of scenarios, using more than one node.

**Repository organization** If you submit your work through a repository, make sure you organize your submissions in subdirectories, and that you give a clear name to all files.

## 17.2 Warmup Exercises

We start with some simple exercises.

### 17.2.1 Hello world

*The exercises in this section are about the routines introduced in section 2.2; for the reference information see section ??.*

First of all we need to make sure that you have a working setup for parallel jobs. The example program `helloworld.c` does the following:

```
// helloworld.c
MPI_Init(&argc,&argv);
MPI_Comm_size(MPI_COMM_WORLD,&ntids);
MPI_Comm_rank(MPI_COMM_WORLD,&mytid);
printf("Hello, this is processor %d out of %d\n",mytid,ntids);
MPI_Finalize();
```

Compile this program and run it in parallel. Make sure that the processors do *not* all say that they are processor 0 out of 1!

### 17.2.2 Trace output

We want to make trace files of the parallel runs, for which we'll use the TAU utility of the University of Oregon. (For documentation, go to <http://www.cs.uoregon.edu/Research/tau/docs.php>.) Here are the steps:

- Load two modules:

```
module load tau
module load jdk64
```

- Recompile your program with `make yourprog`. You'll notice a lot more output: that is the TAU preprocessor.
- Now run your program, setting environment variables `TAU_TRACE` and `TAU_PROFILE` to 1, and `TRACEDIR` and `PROFILEDIR` to where you want the output to be. Big shortcut: do  
`make submit EXECUTABLE=yourprog`

for a batch job or

```
make idevrun EXECUTABLE=yourprog
```

for an interactive parallel run. These last two set all variables for you. See if you can find where the output went...

- Now you need to postprocess the TAU output. Do `make tau EXECUTABLE=yourprog` and you'll get a file `taulog_yourprog.slog2` which you can view with the `jumpshot` program.

### 17.2.3 Collectives

It is a good idea to be able to collect statistics, so before we do anything interesting, we will look at MPI collectives; section 3.2.

Take a look at `time_max.cxx`. This program sleeps for a random number of seconds:

```
// time_max.cxx
wait = (int) ( 6.*rand() / (double)RAND_MAX );
tstart = MPI_Wtime();
sleep(wait);
tstop = MPI_Wtime();
jitter = tstop-tstart-wait;
```

and measures how long the sleep actually was:

```
if (mytid==0)
    sendbuf = MPI_IN_PLACE;
else sendbuf = (void*)&jitter;
MPI_Reduce(sendbuf, (void*)&jitter, 1, MPI_DOUBLE, MPI_MAX, 0, comm);
```

In the code, this quantity is called ‘jitter’, which is a term for random deviations in a system.

**Exercise 17.1.** Change this program to compute the average jitter by changing the reduction operator.

**Exercise 17.2.** Now compute the standard deviation

$$\sigma = \sqrt{\frac{\sum_i (x_i - m)^2}{n}}$$

where  $m$  is the average value you computed in the previous exercise.

- Solve this exercise twice: once by following the reduce by a broadcast operation and once by using an Allreduce.
- Run your code both on a single cluster node and on multiple nodes, and inspect the TAU trace. Some MPI implementations are optimized for shared memory, so the trace on a single node may not look as expected.
- Can you see from the trace how the allreduce is implemented?

**Exercise 17.3.** Finally, use a gather call to collect all the values on processor zero, and print them out. Is there any process that behaves very differently from the others?

For each exercise, submit code, a TAU trace, and an analysis of what you see in the traces. Submit your work by leaving a code, graphics, and a writeup in your repository.

### 17.2.4 Linear arrays of processors

In this section you are going to write a number of variations on a very simple operation: all processors pass a data item to the processor with the next higher number.

- In the file `linear-serial.c` you will find an implementation using blocking send and receive calls.
- You will change this code to use non-blocking sends and receives; they require an `MPI_Wait` call to finalize them.
- Next, you will use `MPI_Sendrecv` to arrive at a synchronous, but deadlock-free implementation.
- Finally, you will use two different one-sided scenarios.

In the reference code `linear-serial.c`, each process defines two buffers:

```
// linear-serial.c
int my_number = mytid, other_number=-1.;
```

where `other_number` is the location where the data from the left neighbour is going to be stored.

To check the correctness of the program, there is a gather operation on processor zero:

```
int *gather_buffer=NULL;
if (mytid==0) {
    gather_buffer = (int*) malloc(ntids*sizeof(int));
    if (!gather_buffer) MPI_Abort(comm,1);
}
MPI_Gather(&other_number,1,MPI_INT,
           gather_buffer,1,MPI_INT, 0,comm);
if (mytid==0) {
    int i,error=0;
    for (i=0; i<ntids; i++)
        if (gather_buffer[i]!=i-1) {
            printf("Processor %d was incorrect: %d should be %d\n",
                   i,gather_buffer[i],i-1);
            error =1;
        }
    if (!error) printf("Success!\n");
    free(gather_buffer);
}
```

#### 17.2.4.1 Coding with blocking calls

Passing data to a neighbouring processor should be a very parallel operation. However, if we code this naively, with `MPI_Send` and `MPI_Recv`, we get an unexpected serial behaviour, as was explained in section 4.2.2.

```
if (mytid<ntids-1)
    MPI_Ssend( /* data: */ &my_number,1,MPI_INT,
               /* to: */ mytid+1, /* tag: */ 0, comm);
if (mytid>0)
```

```
MPI_Recv( /* data: */ &other_number, 1, MPI_INT,  
          /* from: */ mytid-1, 0, comm, &status);
```

(Note that this uses an `Ssend`; see section 4.4.1 for the explanation why.)

**Exercise 17.4.** Compile and run this code, and generate a TAU trace file. Confirm that the execution is serial. Does replacing the `Ssend` by `Send` change this?

Let's clean up the code a little.

**Exercise 17.5.** First write this code more elegantly by using `MPI_PROC_NULL`.

#### 17.2.4.2 A better blocking solution

The easiest way to prevent the serialization problem of the previous exercises is to use the `MPI_Sendrecv` call. This routine acknowledges that often a processor will have a receive call whenever there is a send. For border cases where a send or receive is unmatched you can use `MPI_PROC_NULL`.

**Exercise 17.6.** Rewrite the code using `MPI_Sendrecv`. Confirm with a TAU trace that execution is no longer serial.

Note that the `Sendrecv` call itself is still blocking, but at least the ordering of its constituent send and recv are no longer ordered in time.

#### 17.2.4.3 Non-blocking calls

The other way around the blocking behaviour is to use `Irecv` and `Isend` calls, which do not block. Of course, now you need a guarantee that these send and receive actions are concluded; in this case, use `MPI_Waitall`.

**Exercise 17.7.** Implement a fully parallel version by using `MPI_Isend` and `MPI_Irecv`.

#### 17.2.4.4 One-sided communication

Another way to have non-blocking behaviour is to use one-sided communication. During a `Put` or `Get` operation, execution will only block while the data is being transferred out of or into the origin process, but it is not blocked by the target. Again, you need a guarantee that the transfer is concluded; here use `MPI_Win_fence`.

**Exercise 17.8.** Write two versions of the code: one using `MPI_Put` and one with `MPI_Get`. Make TAU traces.

Investigate blocking behaviour through TAU visualizations.

**Exercise 17.9.** If you transfer a large amount of data, and the target processor is occupied, can you see any effect on the origin? Are the fences synchronized?

### 17.3 Mandelbrot set

If you've never heard the name *Mandelbrot set*, you probably recognize the picture; figure 17.1 Its formal

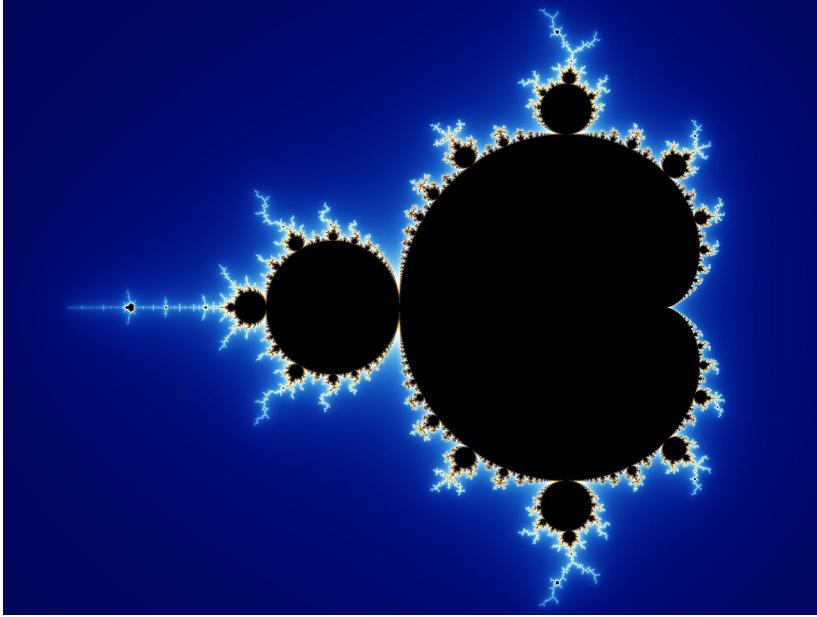


Figure 17.1: The Mandelbrot set

definition is as follows:

A point  $c$  in the complex plane is part of the Mandelbrot set if the series  $x_n$  defined by

$$\begin{cases} x_0 = 0 \\ x_{n+1} = x_n^2 + c \end{cases}$$

satisfies

$$\forall n : |x_n| \leq 2.$$

It is easy to see that only points  $c$  in the bounding circle  $|c| < 2$  qualify, but apart from that it's hard to say much without a lot more thinking. Or computing; and that's what we're going to do.

In this set of exercises you are going to take an example program `mandel_main.cxx` and extend it to use a variety of MPI programming constructs. This program has been set up as a *master-worker* model: there is one master processor (for a change this is the last processor, rather than zero) which gives out work to, and accepts results from, the worker processors. It then takes the results and constructs an image file from them.

#### 17.3.1 Invocation

The `mandel_main` program is called as

```
mpirun -np 123 mandel_main steps 456 iters 789
```

where the `steps` parameter indicates how many steps in  $x, y$  direction there are in the image, and `iters` gives the maximum number of iterations in the `belong` test.

If you forget the parameter, you can call the program with

```
mandel_serial -h
```

and it will print out the usage information.

### 17.3.2 Tools

The driver part of the Mandelbrot program is simple. There is a circle object that can generate coordinates

```
// mandel.h
class circle {
public :
    circle(int pxls,int bound,int bs);
    void next_coordinate(struct coordinate& xy);
    int is_valid_coordinate(struct coordinate xy);
    void invalid_coordinate(struct coordinate& xy);
```

and a global routine that tests whether a coordinate is in the set, at least up to an iteration bound. It returns zero if the series from the given starting point has not diverged, or the iteration number in which it diverged if it did so.

```
int belongs(struct coordinate xy,int itbound) {
    double x=xy.x, y=xy.y; int it;
    for (it=0; it<itbound; it++) {
        double xx,yy;
        xx = x*x - y*y + xy.x;
        yy = 2*x*y + xy.y;
        x = xx; y = yy;
        if (xx+yy>4.) {
            return it;
        }
    }
    return 0;
}
```

In the former case, the point could be in the Mandelbrot set, and we colour it black, in the latter case we give it a colour depending on the iteration number.

```
if (iteration==0)
    memset(colour,0,3*sizeof(float));
```

```

else {
    float rfloat = ((float) iteration) / workcircle->infty;
    colour[0] = rfloat;
    colour[1] = MAX((float)0, (float)(1-2*rfloat));
    colour[2] = MAX((float)0, (float)(2*(rfloat-.5)));
}

```

We use a fairly simple code for the worker processes: they execute a loop in which they wait for input, process it, return the result.

```

void queue::wait_for_work(MPI_Comm comm, circle *workcircle) {
    MPI_Status status; int ntids;
    MPI_Comm_size(comm, &ntids);
    int stop = 0;

    while (!stop) {
        struct coordinate xy;
        int res;

        MPI_Recv(&xy, 1, coordinate_type, ntids-1, 0, comm, &status);
        stop = !workcircle->is_valid_coordinate(xy);
        if (stop) break; //res = 0;
        else {
            res = belongs(xy, workcircle->infty);
        }
        MPI_Send(&res, 1, MPI_INT, ntids-1, 0, comm);
    }
    return;
}

```

A very simple solution using blocking sends on the master is given:

```

// mandel_serial.cxx
class serialqueue : public queue {
private :
    int free_processor;
public :
    serialqueue(MPI_Comm queue_comm, circle *workcircle)
        : queue(queue_comm, workcircle) {
        free_processor=0;
    };
    /**
     * The 'addtask' routine adds a task to the queue. In this
     * simple case it immediately sends the task to a worker
     * and waits for the result, which is added to the image.
     */
}

```

```

This routine is only called with valid coordinates;
the calling environment will stop the process once
an invalid coordinate is encountered.

*/
int addtask(struct coordinate xy) {
    MPI_Status status; int contribution, err;

    err = MPI_Send(&xy, 1, coordinate_type,
        free_processor, 0, comm); CHK(err);
    err = MPI_Recv(&contribution, 1, MPI_INT,
        free_processor, 0, comm, &status); CHK(err);

    coordinate_to_image(xy, contribution);
    total_tasks++;
    free_processor = (free_processor+1)% (ntids-1);

    return 0;
}

```

**Exercise 17.10.** Explain why this solution is very inefficient. Make a trace of its execution that bears this out.

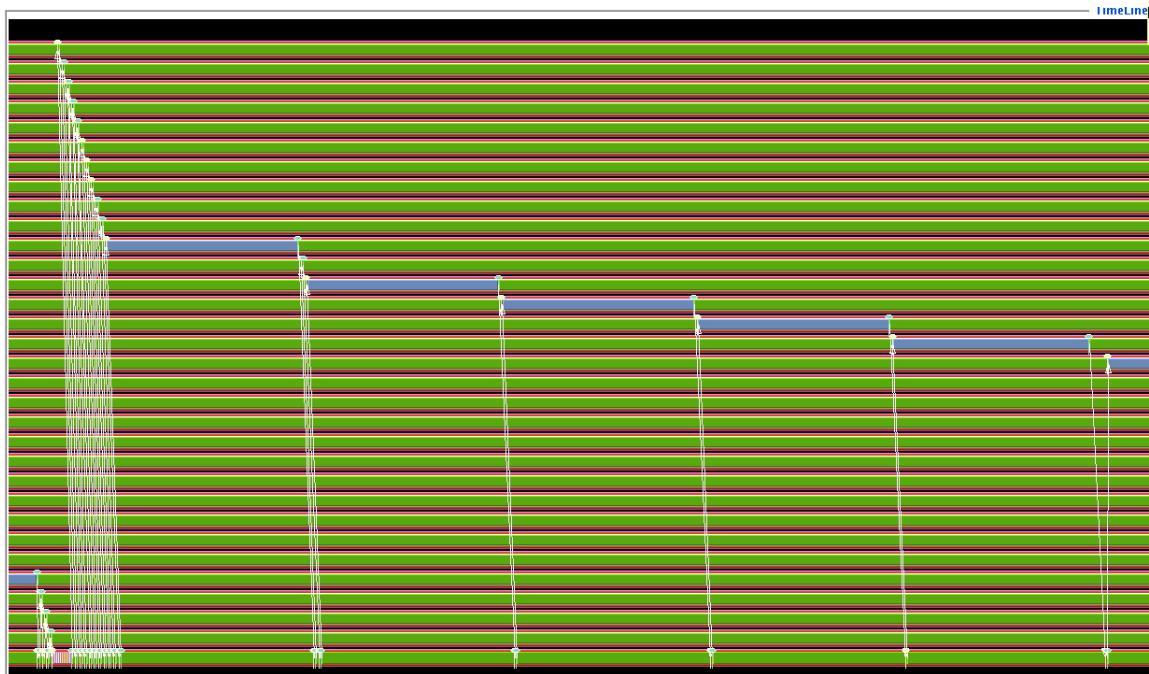


Figure 17.2: Trace of a serial Mandelbrot calculation

### 17.3.3 Bulk task scheduling

The previous section showed a very inefficient solution, but that was mostly intended to set up the code base. If all tasks take about the same amount of time, you can give each process a task, and then wait on them all to finish. A first way to do this is with non-blocking sends.

**Exercise 17.11.** Code a solution where you give a task to all worker processes using non-blocking sends and receives, and then wait for these tasks with MPI\_Waitall to finish before you give a new round of data to all workers. Make a trace of the execution of this and report on the total time.  
 You can do this by writing a new class that inherits from queue, and that provides its own addtask method:

```
// mandel_bulk.cxx
class bulkqueue : public queue {
public :
    bulkqueue(MPI_Comm queue_comm,circle *workcircle)
        : queue(queue_comm,workcircle) {
```

You will also have to override the complete method: when the circle object indicates that all coordinates have been generated, not all workers will be busy, so you need to supply the proper MPI\_Waitall call.

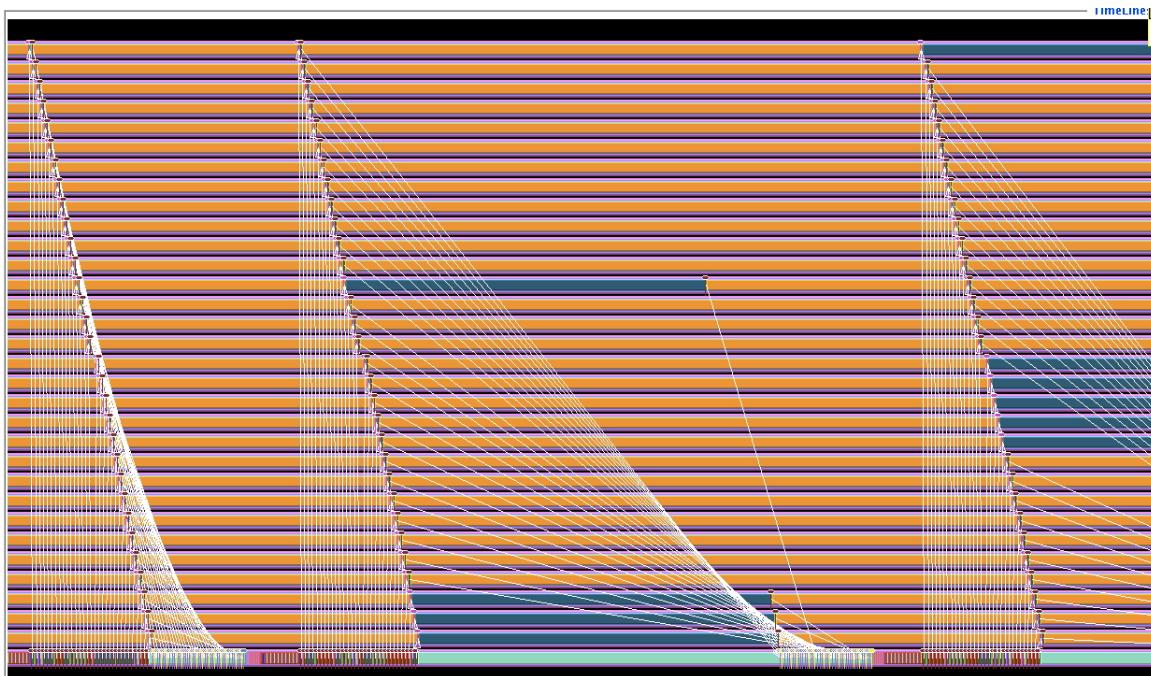


Figure 17.3: Trace of a bulk Mandelbrot calculation

### 17.3.4 Collective task scheduling

Another implementation of the bulk scheduling of the previous section would be through using collectives.

**Exercise 17.12.** Code a solution which uses scatter to distribute data to the worker tasks, and gather to collect the results. Is this solution more or less efficient than the previous?

### 17.3.5 Asynchronous task scheduling

At the start of section 17.3.3 we said that bulk scheduling mostly makes sense if all tasks take similar time to complete. In the Mandelbrot case this is clearly not the case.

**Exercise 17.13.** Code a fully dynamic solution that uses MPI\_Probe or MPI\_Waitany.  
Make an execution trace and report on the total running time.

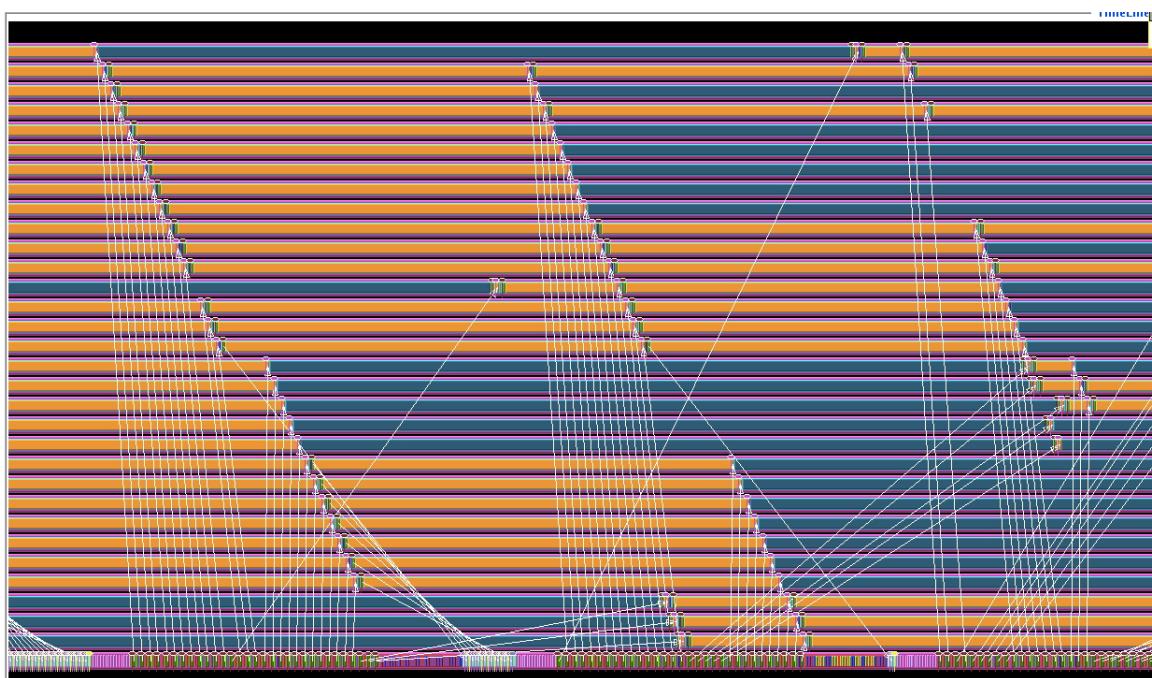


Figure 17.4: Trace of an asynchronous Mandelbrot calculation

### 17.3.6 One-sided solution

Let us reason about whether it is possible (or advisable) to code a one-sided solution to computing the Mandelbrot set. With active target synchronization you could have an exposure window on the host to which the worker tasks would write. To prevent conflicts you would allocate an array and have each worker write to a separate location in it. The problem here is that the workers may not be sufficiently synchronized because of the differing time for computation.

Consider then passive target synchronization. Now the worker tasks could write to the window on the master whenever they have something to report; by locking the window they prevent other tasks from interfering. After a worker writes a result, it can get new data from an array of all coordinates on the master.

It is hard to get results into the image as they become available. For this, the master would continuously have to scan the results array. Therefore, constructing the image is easiest done when all tasks are concluded.

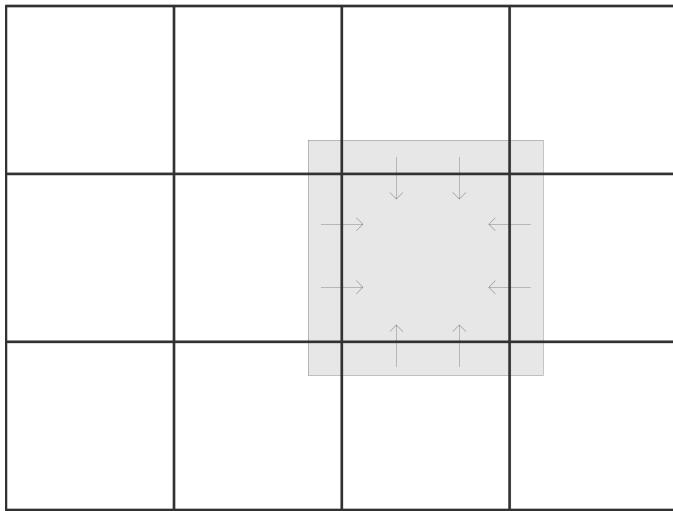


Figure 17.5: A grid divided over processors, with the ‘ghost’ region indicated

## 17.4 Data parallel grids

In this section we will gradually build a semi-realistic example program. To get you started some pieces have already been written: as a starting point look at `code/mpi/c/grid.cxx`.

### 17.4.1 Description of the problem

With this example you will investigate several strategies for implementing a simple iterative method. Let’s say you have a two-dimensional grid of datapoints  $G = \{g_{ij} : 0 \leq i < n_i, 0 \leq j < n_j\}$  and you want to compute  $G'$  where

$$g'_{ij} = 1/4 \cdot (g_{i+1,j} + g_{i-1,j} + g_{i,j+1} + g_{i,j-1}). \quad (17.1)$$

This is easy enough to implement sequentially, but in parallel this requires some care.

Let’s divide the grid  $G$  and divide it over a two-dimension grid of  $p_i \times p_j$  processors. (Other strategies exist, but this one scales best; see section [HPSC-6.5](#).) Formally, we define two sequences of points

$$0 = i_0 < \dots < i_{p_i} < i_{p_i+1} = n_i, \quad 0 < j_0 < \dots < j_{p_j} < j_{p_j+1} = n_j$$

and we say that processor  $(p, q)$  computes  $g_{ij}$  for

$$i_p \leq i < i_{p+1}, \quad j_q \leq j < j_{q+1}.$$

From formula (17.1) you see that the processor then needs one row of points on each side surrounding its part of the grid. A picture makes this clear; see figure 17.5. These elements surrounding the processor’s own part are called the *halo* or *ghost region* of that processor.

The problem is now that the elements in the halo are stored on a different processor, so communication is needed to gather them. In the upcoming exercises you will have to use different strategies for doing so.

### 17.4.2 Code basics

The program needs to read the values of the grid size and the processor grid size from the commandline, as well as the number of iterations. This routine does some error checking: if the number of processors does not add up to the size of MPI\_COMM\_WORLD, a nonzero error code is returned.

```
ierr = parameters_from_commandline
      (argc, argv, comm, &ni, &nj, &pi, &pj, &nit);
if (ierr) return MPI_Abort(comm, 1);
```

From the processor parameters we make a processor grid object:

```
processor_grid *pgrid = new processor_grid(comm, pi, pj);
```

and from the numerical parameters we make a number grid:

```
number_grid *grid = new number_grid(pgrid, ni, nj);
```

Number grids have a number of methods defined. To set the value of all the elements belonging to a processor to that processor's number:

```
grid->set_test_values();
```

To set random values:

```
grid->set_random_values();
```

If you want to visualize the whole grid, the following call gathers all values on processor zero and prints them:

```
grid->gather_and_print();
```

Next we need to look at some data structure details.

The definition of the `number_grid` object starts as follows:

```
class number_grid {
public:
    processor_grid *pgrid;
    double *values, *shadow;
```

where `values` contains the elements owned by the processor, and `shadow` is intended to contain the values plus the ghost region. So how does `shadow` receive those values? Well, the call looks like

```
grid->build_shadow();
```

and you will need to supply the implementation of that. Once you've done so, there is a routine that prints out the shadow array of each processor

```
grid->print_shadow();
```

This routine does the sequenced printing that you implemented in exercise ??.

In the file `code/mpi/c/grid_impl.cxx` you can see several uses of the macro `INDEX`. This translates from a two-dimensional coordinate system to one-dimensional. Its main use is letting you use  $(i, j)$  coordinates for indexing the processor grid and the number grid: for processors you need the translation to the linear rank, and for the grid you need the translation to the linear array that holds the values.

A good example of the use of `INDEX` is in the `number_grid::relax` routine: this takes points from the shadow array and averages them into a point of the `values` array. (To understand the reason for this particular averaging, see HPSC-4.2.3 and HPSC-5.5.3.) Note how the `INDEX` macro is used to index in a  $\text{ilength} \times \text{jlength}$  target array `values`, while reading from a  $(\text{ilength} + 2) \times (\text{jlength} + 2)$  source array `shadow`.

```
for (i=0; i<ilength; i++) {
    for (j=0; j<jlength; j++) {
        int c=0;
        double new_value=0.;
        for (c=0; c<5; c++) {
            int ioff=i+1+ioffsets[c], joff=j+1+joffsets[c];
            new_value += coefficients[c] *
                shadow[ INDEX(ioff, joff, ilength+2, jlength+2) ];
        }
        values[ INDEX(i, j, ilength, jlength) ] = new_value/8.;
    }
}
```

## Chapter 18

### Bibliography, index, and list of acronyms

#### 18.1 Bibliography

- [1] Ernie Chan, Marcel Heimlich, Avi Purkayastha, and Robert van de Geijn. Collective communication: theory, practice, and experience. *Concurrency and Computation: Practice and Experience*, 19:1749–1783, 2007.
- [2] W. Gropp, E. Lusk, and A. Skjellum. *Using MPI*. The MIT Press, 1994.
- [3] Torsten Hoefler, Prabhanjan Kambadur, Richard L. Graham, Galen Shipman, and Andrew Lumsdaine. A case for standard non-blocking collective operations. In *Proceedings, Euro PVM/MPI*, Paris, France, October 2007.
- [4] Torsten Hoefler, Christian Siebert, and Andrew Lumsdaine. Scalable communication protocols for dynamic sparse data exchange. *SIGPLAN Not.*, 45(5):159–168, January 2010.
- [5] L. V. Kale and S. Krishnan. Charm++: Parallel programming with message-driven objects. In *Parallel Programming using C++*, G. V. Wilson and P. Lu, editors, pages 175–213. MIT Press, 1996.
- [6] R. Thakur, W. Gropp, and B. Toonen. Optimizing the synchronization operations in MPI one-sided communication. *Int'l Journal of High Performance Computing Applications*, 19:119–128, 2005.

## 18.2 List of acronyms

**AVX** Advanced Vector Extensions  
**BSP** Bulk Synchronous Parallel  
**CAF** Co-array Fortran  
**CUDA** Compute-Unified Device Architecture  
**DAG** Directed Acyclic Graph  
**DSP** Digital Signal Processing  
**FPU** Floating Point Unit  
**FFT** Fast Fourier Transform  
**FSA** Finite State Automaton  
**GPU** Graphics Processing Unit  
**HPC** High-Performance Computing  
**HPF** High Performance Fortran  
**ICV** Internal Control Variable  
**MIC** Many Integrated Cores  
**MIMD** Multiple Instruction Multiple Data  
**MPI** Message Passing Interface  
**MTA** Multi-Threaded Architecture  
**NUMA** Non-Uniform Memory Access  
**OS** Operating System  
**PGAS** Partitioned Global Address Space

**PDE** Partial Differential Equation  
**PRAM** Parallel Random Access Machine  
**RDMA** Remote Direct Memory Access  
**RMA** Remote Memory Access  
**SAN** Storage Area Network  
**SaaS** Software as-a Service  
**SFC** Space-Filling Curve  
**SIMD** Single Instruction Multiple Data  
**SIMT** Single Instruction Multiple Thread  
**SM** Streaming Multiprocessor  
**SMP** Symmetric Multi Processing  
**SOR** Successive Over-Relaxation  
**SP** Streaming Processor  
**SPMD** Single Program Multiple Data  
**SPD** symmetric positive definite  
**SSE** SIMD Streaming Extensions  
**TLB** Translation Look-aside Buffer  
**UMA** Uniform Memory Access  
**UPC** Unified Parallel C  
**WAN** Wide Area Network

## 18.3 Index

\_OPENMP, 147

active target synchronization, 43, 45, 48  
all-to-all, 19  
anti dependency, see data dependencies  
argc, 74, 75  
argv, 74, 75  
atomic operation, 49  
atomic operations, 49

bandwidth, 25  
barrier, 69  
    in MPI, 65  
batch  
    job, 11  
    scheduler, 11  
Beowulf cluster, 10  
Boolean satisfiability, 17  
boost, 13  
breakpoint, 212  
broadcast, 19

C  
    bindings, 12–13  
C++  
    bindings, 12–13  
    standard library, 81  
        vector, 81  
C99, 76  
cacheline, 165  
Charmpp, 11  
choice, 76  
chunk, 183  
    chunk, 183  
collectives, 18–28  
    non-blocking, 25  
column-major storage, 52  
communication  
    asynchronous, 41  
    blocking, 32–34  
    non-blocking, 37–40  
    one-sided, 43–50

one-sided, implementation of, 50  
overlap with computation, 41  
persistent, 42–43, 94–96  
synchronous, 41  
two-sided, 43

communicator, 60–64  
compare-and-swap, 36  
compiler, 82  
    optimization level, 207  
construct, 181  
contiguous  
    data type, 53  
core, 14  
core dump, 207  
cpp, 147  
Cray  
    T3E, 51  
critical section  
    flush at, 190  
critical section, 121, 164, 166, 189  
critical sections  
    cost of a, 179

Dalcin  
    Lisandro, 13  
data dependencies, 174–176  
data dependency, 173  
datatype, 52–58  
    derived, 53–57, 77–83  
    elementary, 53, 76–77  
    signature, 54, 82  
datatypes  
    derived, 52  
    elementary, 52  
ddd, 207  
DDT, 207, 214–215  
ddt, 71  
deadlock, 33, 50, 213  
debug flag, 208  
debugger, 207  
debugging, 207–215  
    parallel, 213

dense linear algebra, 61  
directives, 148, 180  
  cpp, 148  
distributed array, 29  
distributed shared memory, 43  
dynamic mode, 181  
  
eager limit, 33, 85, 137  
Eclipse, 214  
  PTP, 214  
emacs, 219  
environment variables, 201  
epoch, 45  
  access, 46, 48, 103  
  exposure, 46, 48, 102  
error return, 13  
ethernet, 12  
  
false sharing, 165  
fence, 45  
Fibonacci sequence, 169–171  
flow dependency, see data dependencies  
fork/join model, 146, 152, 173  
Fortran  
  1-based indexing, 91  
  array syntax, 160, 185  
  bindings, 13  
  fixed-form source, 180  
  forall loops, 185  
  Fortran90, 73  
  
gather, 19  
gdb, 207–214  
ghost region, 231  
GNU, 207  
  gdb, see gdb  
grid  
  Cartesian, 64  
  periodic, 64  
group, 102  
group of  
  processors, 48  
  
halo, 231  
  update, 101  
  
handshake, 50  
heap, 151  
histogram, 167  
hostname, 119  
  
I/O  
  in OpenMP, 159  
ibrwn, 11  
implicit barrier, 187  
  after single directive, 160  
indexed  
  data type, 53  
inner product, 30  
instrumentation, 216  
Internal Control Variable (ICV), 188–189  
  
latency, 25  
  hiding, 38  
lexical scope, 161  
load balancing, 155  
lock, 167  
  flush at, 190  
LU factorization, 156  
  
malloc, 151  
Mandelbrot set, 17, 163, 224  
master-worker, 49, 224  
master-worker model, 40, 41  
matching, 72  
matrix  
  sparse, 25  
  transposition, 61  
matrix-vector product  
  sparse, 24  
Monte Carlo codes, 17  
MPI  
  3, 57  
  C/C++ bindings, see C/C++, bindings  
  Fortran issues, 68, 117  
  I/O, 68  
  initialization, 73  
  Python issues, 117  
  semantics, 72  
  version, 69

mpi.h, 73  
mpi4py, 13  
MPI\_Abort, 74  
MPI\_Accumulate, 47, 99  
MPI\_Aint  
    in Fortran, 117  
MPI\_Aint, 77  
MPI\_Allgather, 24, 109  
MPI\_Allgatherv, 24, 110, 122  
MPI\_Alloc\_mem, 44, 97  
MPI\_Allreduce, 24, 30, 109, 123  
MPI\_Alltoall, 24, 109  
MPI\_Alltoallv, 24, 110, 111  
MPI\_ANY\_SOURCE, 24, 40, 72, 87, 108, 118, 131  
MPI\_ANY\_TAG, 41, 87  
MPI\_Barrier, 65, 69  
MPI\_Bcast, 104, 104, 124  
MPI\_BOTTOM, 77  
MPI\_Bsend, 42, 93  
MPI\_Bsend\_init, 96  
MPI\_BSEND\_OVERHEAD, 42, 84  
MPI\_Buffer\_attach, 93  
MPI\_Buffer\_detach, 93  
MPI\_Cancel, 118, 124  
MPI\_Cart..., 125  
MPI\_Cart\_coord, 116  
MPI\_Cart\_create, 116  
MPI\_Cart\_rank, 116  
MPI\_Comm\_create, 62  
MPI\_Comm\_dup, 60, 114, 125  
MPI\_Comm\_free, 60, 114  
MPI\_Comm\_group, 62  
MPI\_COMM\_NULL, 60  
MPI\_Comm\_rank, 16, 75  
MPI\_COMM\_SELF, 60  
MPI\_Comm\_set\_errhandler, 67  
MPI\_Comm\_set\_name, 62  
MPI\_Comm\_size, 16, 75  
MPI\_Comm\_split, 61, 115, 133  
MPI\_Comm\_split\_type, 49, 115, 115  
MPI\_COMM\_TYPE\_SHARED, 115  
MPI\_COMM\_WORLD, 59, 60  
MPI\_Datatype, 76, 107  
MPI\_DATATYPE\_NULL, 77  
MPI\_ERR\_BUFFER, 93  
MPI\_ERR\_INTERN, 93  
MPI\_ERROR, 68  
MPI\_Error\_string, 68  
MPI\_ERRORS\_ARE\_FATAL, 67  
MPI\_ERRORS\_RETURN, 67, 119  
MPI\_Exscan, 23, 112, 112, 126  
MPI\_Fetch\_and\_op, 49, 103, 127  
MPI\_Finalize, 16, 74  
MPI\_Finalized, 74  
MPI\_Gather, 22, 24, 107, 127  
MPI\_Gatherv, 24, 110, 110, 111, 128  
MPI\_Get, 47, 99, 128  
MPI\_Get\_count, 40, 57, 66, 87, 137, 139  
MPI\_Get\_elements, 57  
MPI\_Get\_elements\_x, 58  
MPI\_Get\_processor\_name, 16, 119  
MPI\_Get\_version, 69  
MPI\_Group\_difference, 62  
MPI\_Group\_excl, 62  
MPI\_Group\_incl, 62  
MPI\_Ibarrier, 25  
MPI\_Ibcast, 113  
MPI\_Ibsend, 93  
MPI\_IN\_PLACE, 105, 109, 123, 132  
MPI\_Info, 97  
MPI\_INFO\_ENV, 75  
MPI\_INFO\_NULL, 44, 97  
MPI\_Init  
    in Fortran, 68  
MPI\_Init, 16, 73, 75, 118, 121  
MPI\_Init\_thread, 66, 118, 121, 129, 201  
MPI\_Initialized, 74  
MPI\_Iprobe, 40  
MPI\_Irecv, 37, 89, 140  
MPI\_Is\_thread\_main, 121  
MPI\_Isend, 37, 89, 140  
MPI\_LOCK\_EXCLUSIVE, 49  
MPI\_LOCK\_SHARED, 49  
MPI\_MAX, 20  
MPI\_MAX\_PROCESSOR\_NAME, 119  
MPI\_OP, 106

MPI\_Op, 119  
MPI\_Op\_create, 23  
MPI\_PACK, 83  
MPI\_Pack, 58, 130  
MPI\_Pack\_size, 93  
MPI\_PACKED, 58, 83  
MPI\_Probe, 40  
MPI\_PROC\_NULL, 35, 35, 116, 223  
MPI\_PROD, 20  
MPI\_Put, 46, 98, 130, 142  
MPI\_Query\_thread, 121  
MPI\_Raccumulate, 100  
MPI\_Recv, 32, 37, 84, 131, 134  
MPI\_Recv\_init, 42, 94, 94, 136  
MPI\_Reduce, 22, 24, 104, 104, 132  
MPI\_Reduce\_scatter, 24, 108, 133  
MPI\_REPLACE, 47  
MPI\_Request, 25, 89, 96, 113  
MPI\_Request\_free, 43, 96  
MPI\_Request\_get\_status, 96  
MPI\_Rget, 100  
MPI\_Rput, 99  
MPI\_Rsend, 50  
MPI\_Rsend\_init, 96  
MPI\_Scan, 22, 112, 134  
MPI\_Scatter, 21, 107  
MPI\_Scatter\_reduce, 50  
MPI\_Scatterv, 108, 110  
MPI\_Send, 32, 37, 84, 134  
MPI\_Send\_init, 42, 94, 94, 135, 136  
MPI\_Sendrecv, 35, 37, 89, 136, 223  
MPI\_Sendrecv\_replace, 35, 89  
MPI\_Sizeof, 68, 117  
MPI\_SOURCE, 40, 88, 88, 131, 140  
MPI\_Ssend, 41, 50, 137  
MPI\_Ssend\_init, 96  
MPI\_Start, 43, 95  
MPI\_Start\_all, 95  
MPI\_Startall, 43, 95, 136  
MPI\_Status, 32, 66, 87, 87–89, 131  
MPI\_STATUS\_IGNORE, 32, 66, 87, 91  
MPI\_STATUSES\_IGNORE, 66, 87, 90  
MPI\_SUBVERSION, 69  
MPI\_SUM, 20, 22  
MPI\_Test, 96  
MPI\_Test..., 40  
MPI\_Testall, 92  
MPI\_Testany, 92  
MPI\_THREAD\_FUNNELED, 121  
MPI\_THREAD\_FUNNELLED, 201  
MPI\_THREAD\_MULTIPLE, 121, 201  
MPI\_THREAD\_SERIAL, 201  
MPI\_THREAD\_SERIALIZED, 121  
MPI\_THREAD\_SINGLE, 121, 201  
MPI\_Type\_commit, 54, 77  
MPI\_Type\_contiguous, 54, 78, 78, 116, 137  
MPI\_Type\_create\_hindexed, 81  
MPI\_Type\_create\_struct, 57, 81  
MPI\_Type\_create\_subarray, 54, 56  
MPI\_Type\_extent, 77, 82  
MPI\_Type\_free, 54, 77  
MPI\_Type\_hindexed, 54, 57  
MPI\_Type\_indexed, 54, 57, 80, 80, 137  
MPI\_Type\_struct, 54, 83, 138  
MPI\_Type\_vector, 54, 55, 79, 79, 139  
MPI\_UNPACK, 83  
MPI\_Unpack, 58, 130  
MPI\_VERSION, 69  
MPI\_Wait, 25, 39, 43, 87  
MPI\_Wait..., 38, 96  
MPI\_Waitall, 39, 87, 92, 140  
MPI\_Waitany, 39, 66, 87, 90, 91, 140  
MPI\_Waitsome, 39, 87  
MPI\_Win\_complete, 103, 142  
MPI\_Win\_create, 96, 96, 117, 128, 142  
MPI\_Win\_fence, 45, 100, 128, 142, 223  
MPI\_Win\_flush..., 100  
MPI\_Win\_lock, 103, 127, 141  
MPI\_Win\_post, 102, 142  
MPI\_Win\_start, 103, 142  
MPI\_Win\_unlock, 103  
MPI\_Win\_wait, 102, 142  
MPI\_Wtick, 70, 120  
MPI\_Wtime, 69, 120  
MPI\_WTIME\_IS\_GLOBAL, 120  
mpiexec, 11, 75

mpiexec, 16  
mpif.h, 73  
mpirun, 11, 12, 60, 75  
    and environment variables, 201  
MPL, 13  
  
new, 151  
node, 14  
num\_threads, 181  
numerical integration, 154  
numpy, 13, 117  
  
omp  
    atomic, 166, 177, 187  
    barrier, 166  
    copyin, 185  
    critical, 166, 176, 187  
    declare simd, 173  
    default, 162  
    do, 182  
    firstprivate, 185  
    flush, 185, 190  
    for, 182  
        barrier behaviour, 187  
    if, 181  
    lastprivate, 158, 176, 185  
    master, 159, 184  
    ordered, 183  
    parallel, 149, 181  
    parallel do, 182  
    parallel for, 153, 182  
    parallel sections, 182  
    private, 162, 185  
    section, 158  
    sections, 158  
    simd, 173  
    single, 159, 184  
    task, 171, 172  
    taskgroup, 173  
    taskwait, 172  
    threadprivate, 163, 176, 186  
    workshare, 160, 185  
  
omp clause  
    default, 185  
none, 185  
private, 185  
shared, 185  
depend, 173, 189  
firstprivate, 162  
lastprivate, 163  
nowait, 157, 187  
private, 161  
reduction, 183  
schedule  
    auto, 156  
    chunk, 155  
    guided, 155  
    runtime, 156  
    shared, 185  
omp.h, 180  
OMP\_DYNAMIC, 188  
omp\_get\_dynamic, 181, 188  
omp\_get\_max\_threads, 188  
omp\_get\_nested, 188  
omp\_get\_num\_procs, 152, 189  
omp\_get\_num\_threads, 149, 152, 188  
omp\_get\_schedule, 183, 188  
omp\_get\_thread\_num, 149, 189  
omp\_get\_wtick, 173  
omp\_get\_wtime, 173, 189  
omp\_in\_parallel, 181, 189  
OMP\_NESTED, 152, 188  
OMP\_NUM\_THREADS, 147, 152, 188, 201  
OMP\_PROC\_BIND, 189, 190  
omp\_sched\_affinity, 183  
omp\_sched\_auto, 183  
omp\_sched\_dynamic, 183  
omp\_sched\_guided, 183  
omp\_sched\_runtime, 183  
omp\_sched\_static, 183  
OMP\_SCHEDULE, 156, 183, 188  
omp\_set\_dynamic, 181, 188  
omp\_set\_nested, 188  
omp\_set\_num\_threads, 152, 188  
omp\_set\_schedule, 183, 188  
OMP\_STACKSIZE, 185, 189  
OMP\_WAIT\_POLICY, 189

OpenMP  
    accelerator support in, 178  
    co-processor support in, 178  
    compiling, 147  
    environment variables, 147, 152, 188–189  
    library routines, 152  
    library routines, 188–189  
    running, 147  
    version 3.1  
        thread affinity, 190  
operating system, 179  
origin, 43, 48, 103  
out-of-order execution, 174  
output dependency, see data dependencies  
overlapping computation and communication, see  
    latency, hiding  
owner computes, 31

packing, 58  
parallel region  
    barrier at the end of, 187  
parallel region, 146, 148, 152  
    flush at, 190  
parallel regions  
    nested, 188  
passive target synchronization, 43, 48, 127  
persistent communication, see communication, per-  
    sistent  
ping-pong, 31, 70  
MPI\_..., 70  
point-to-point, 31  
prefix, 112  
prefix operation, 23  
private variables, 151  
process, 14  
producer-consumer, 176  
program order, 174  
progress, 72  
purify, 211  
PVM, 11  
Python  
    bindings, 13  
race condition, 49, 164, 166, 176, 177

random number generation, 163  
random number generator, 176  
reduction, 19, 183  
reduction, 164  
region of code, 181  
Riemann sums, 154  
RMA  
    active, 43  
    passive, 43  
root process, 18

scan, 19  
    exclusive, 23  
    inclusive, 23  
scatter, 19  
schedule  
    clause, 183  
schedule, 154  
scope  
    lexical, 150  
    of variables, 150  
segmentation fault, 209  
segmented scan, 23  
sentinel, 148, 180  
serialization, 34  
shared data, 146  
shared variables, 151  
shmem, 51  
Single Program Multiple Data (SPMD), 149  
sizeof, 68, 97  
socket, 14  
sort  
    exchange, 36  
sparse matrix vector product, 23  
spin-lock, 189  
ssh, 11  
stack, 151, 189  
    overflow, 162, 185  
    per thread, 185  
status  
    of receive call, 40  
stride, 52  
struct  
    data type, 53

- structured block, 181
- symbol table, 207, 208
- synchronization
  - in MPI, 65
  - in OpenMP, 165–171
- target, 43, 48, 103
  - active synchronization, see active target synchronization
  - passive synchronization, see passive target synchronization
- task
  - scheduler, 171
- TAU, 216
- thread
  - affinity, 190
  - private data, 163
- thread-safe, 66, 176
- threads, 146
  - initial, 148
  - master, 146, 148
  - team of, 146, 148
- timing
  - MPI, 69–70
  - OpenMP, 189
- topology, 63
- TotalView, 71, 207, 214
- valgrind, 71, 211–212
- vector
  - data type, 53
- vi, 219
- virtual shared memory, 43
- wall clock, 120
- wall clock time, 69
- while loops, 158
- window, 43–46
- work sharing, 146
- work sharing construct, 152
- workshare
  - flush after, 190
- worksharing constructs, 182–185
  - implied barriers at, 187
- wraparound connections, 64