
Programming With the Message Passing Paradigm

Based on several resources including but not limited to
``Introduction to Parallel Computing'', Addison Wesley (2003).

MPI Overview

1) MPI in a nutshell

- ❑ Introduction
- ❑ Programmer view
- ❑ First MPI example

2) Message passing operations

- ❑ Non-buffered
- ❑ Buffered
- ❑ Blocking
- ❑ Non-blocking

3) MPI basics

- ❑ Six fundamental functions
- ❑ Terminologies

4) Advanced MPI programming

- ❑ Avoid deadlock
- ❑ Topologies
- ❑ Collective communications
- ❑ Groups and communicators

What is MPI ?

- MPI is a parallel programming paradigm which derives parallelism out of a set of concurrent sequential processes (CSP), that communicate and synchronize their work in order to complete a common task as fast as possible.
- As a library, MPI can be linked and used with any existing programming language such as C, C++, Fortran, ...

Fundamentals of MPI

- MPI is a standard specification for a library of message passing functions
- MPI is a library, not a language. It specifies:
 - The names, calling sequences and results of subroutines to be called from Fortran programs.
 - The functions to be called from C programs
 - The classes and methods that make up the MPI C++ library
 - ..
- The programs that users write are compiled with ordinary compilers and linked with the MPI library.

Fundamentals of MPI (cont.)

- MPI is a specification, not a particular implementation. A correct MPI program should be able to run on all MPI implementations without change

The MPI Forum

- At Supercomping'92, a committee was formed to define a message-passing standard with the following goals
 - Define a portable standard for message passing
 - Operate in a completely open way, allowing anyone to join the discussion
 - Finish it in one year
- The MPI standard was completed on May 1994.
- The forum reconvened during 1995-97 to extend MPI to include remote memory operations, parallel I/O and dynamic process management.

MPI History

- MPI-1 finished in May 1994 (1.0)
- Clarifications (1.1) June 1995
- MPI-1 achieves critical mass in 1996
- MPI-2 started in 1995, finished July 1997
 - MPI-1.2 issued at same time : small corrections/clarifications
 - MPI-2 is additions to MPI-1, not a change
- MPI-3 is here, lots of PGAS concepts

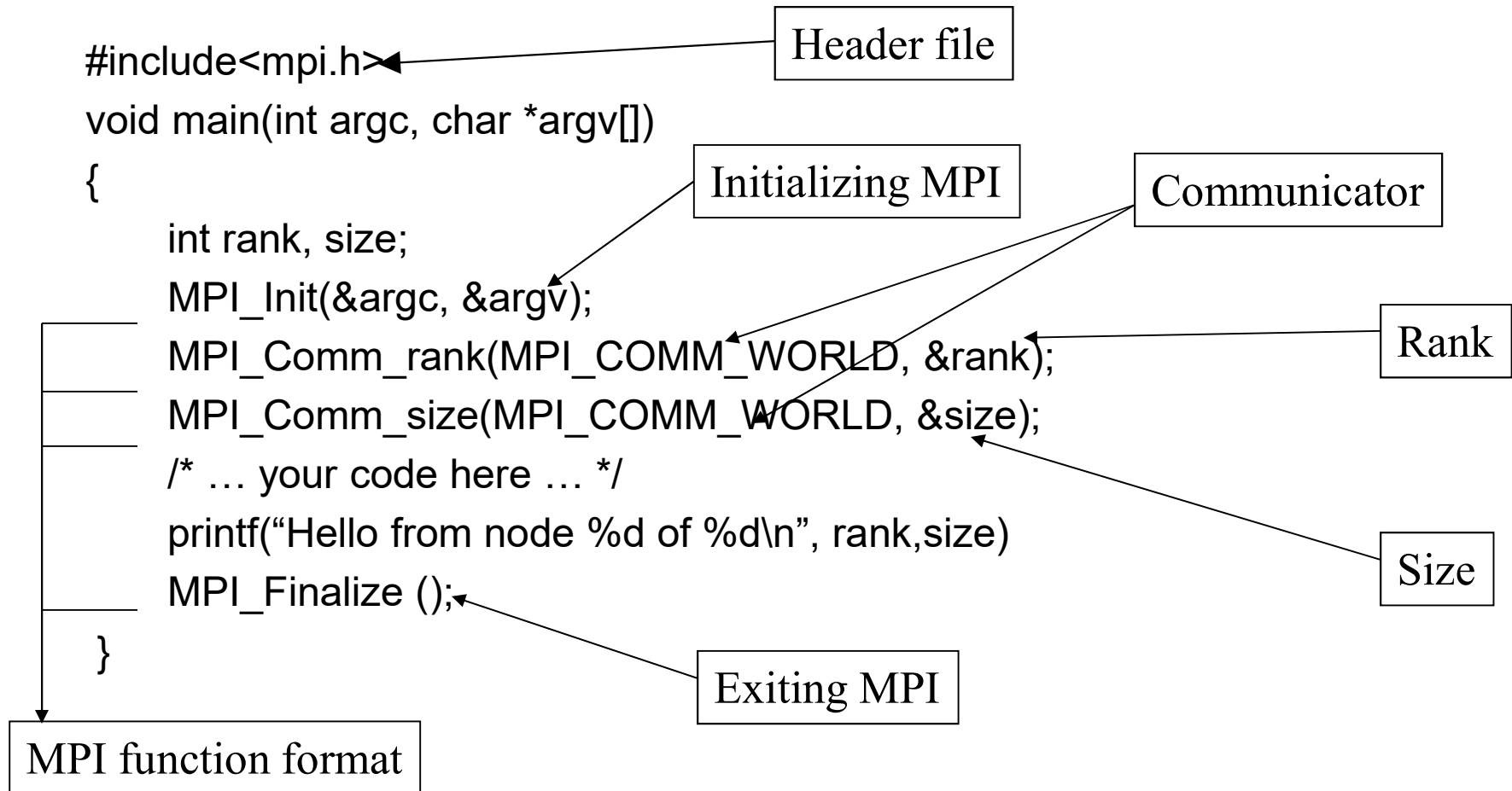
Is MPI Large or Small?

- MPI is Large (more than 200 functions)
 - Many feature requires extensive API
 - Complexity of use not related to number of functions
- MPI is small (6 functions)
 - All that's needed to get started are only 6 functions

Programmer View of MPI

- Group of processes that are allowed to communicate with each other, through send and receive calls
- Single Program Multiple Data (SPMD) operation.
- Synchronization mechanisms
- A communicator argument, most often `MPI_COMM_WORLD`, defines your group of processes
- A process can use a call to find out its rank (id)
- A process can use a call to find out the size of its group
- The program (process before execution) is made of a typical C (FORTRAN) program along with these calls

Anatomy of An MPI Programs



“Hello World” output

- To Compile:
mpicc hello.c -o hello
- To run with 4 processes:
mpirun -np 4 hello
- Output:
Hello from node 2 of 4
Hello from node 1 of 4
Hello from node 3 of 4
Hello from node 0 of 4

* Note - Order of output is not specified by MPI

MPI Basic Functions

<code>int MPI_Init(int argc *, char **argv)</code>	Initialize the MPI environment; must be called before calling any other MPI function
<code>int MPI_Send(void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)</code>	Send contents of <code>buf</code> , containing <code>count</code> instances of <code>datatype</code> to process specified by <code>envelope</code> information
<code>int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status *status)</code>	Receive <code>count</code> instances of <code>datatype</code> into <code>buf</code> from process given by <code>envelope</code> , and return information on <code>status</code>
<code>int MPI_Comm_rank(MPI_Comm comm, int *rank)</code>	Return in <code>rank</code> the process number of the calling process in group of <code>comm</code>
<code>int MPI_Comm_size(MPI_Comm comm, int *size)</code>	Return in <code>size</code> the number of processes in group associated with <code>comm</code>
<code>int MPI_Finalize(void)</code>	Conclude operation and clean up MPI; no MPI functions can be called after.

Features of MPI

- MPI features include:
 - A set of functions to achieve:
 - Point-to-Point Communication (Blocking and Non Blocking)
 - Collective communication (One-to-All, All-to-All,...)
 - Strong Typing
 - Tools for definition of Virtual Topologies
 - A set of tools for performance monitoring

MPI Overview

1) MPI in a nutshell

- ❑ Introduction
- ❑ Programmer view
- ❑ First MPI example

2) Message passing operations

- ❑ Non-buffered
- ❑ Buffered
- ❑ Blocking
- ❑ Non-blocking

3) MPI basics

- ❑ Six fundamental functions
- ❑ Terminologies

4) Advanced MPI programming

- ❑ Avoid deadlock
- ❑ Topologies
- ❑ Collective communications
- ❑ Groups and communicators

Principles of Message-Passing Programming

- The logical view of a machine supporting the message-passing paradigm consists of p processes, each with its own exclusive address space.
- Each data element must belong to one of the partitions of the space; hence, data must be explicitly partitioned and placed.
- All interactions (read-only or read/write) require cooperation of two processes (two-sided) the process that has the data and the process that wants to access the data.
- These two constraints, while onerous, make underlying costs very explicit to the programmer.

Principles of Message-Passing Programming

- Message-passing programs are often written using the *asynchronous* or *loosely synchronous* paradigms.
- In the asynchronous paradigm, all concurrent tasks execute asynchronously.
- In the loosely synchronous model, tasks or subsets of tasks synchronize to perform interactions. Between these interactions, tasks execute completely asynchronously.
- Most message-passing programs are written using the *single program multiple data* (SPMD) model.

The Building Blocks: Send and Receive Operations

- The prototypes of these operations are as follows:

```
send(void *sendbuf, int nelems, int dest)
receive(void *recvbuf, int nelems, int source)
```

- Consider the following code segments:

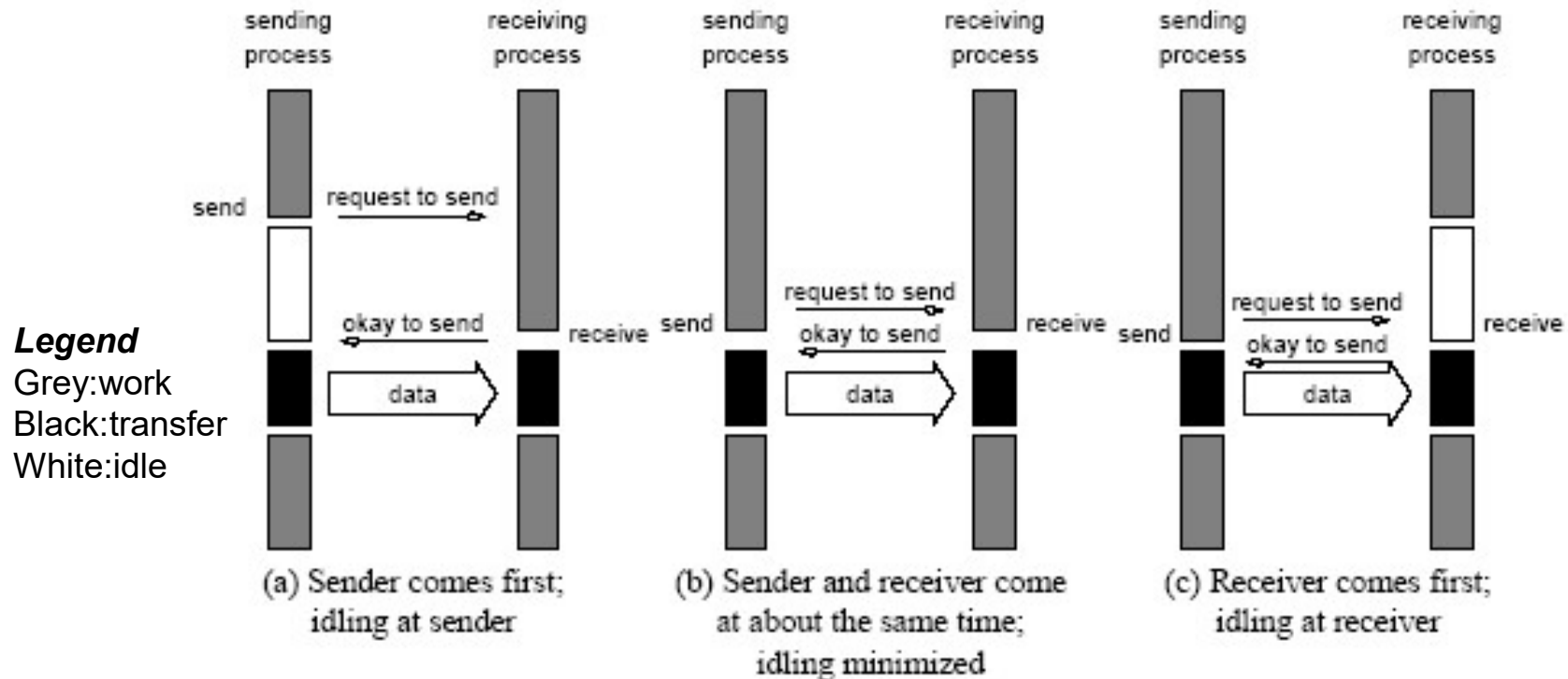
P0	P1
a = 100;	receive(&a, 1, 0)
send(&a, 1, 1);	printf("%d\n", a);
a = 0;	

- The semantics of the send operation require that the value received by process P1 must be 100 as opposed to 0.
- This motivates the design of the send and receive protocols.

Non-Buffered Blocking Message Passing Operations

- A simple method for forcing send/receive semantics is for the send operation to return only when it is safe to do so.
- In the non-buffered blocking send, the operation does not return until the matching receive has been encountered at the receiving process.
- Idling and deadlocks are major issues with non-buffered blocking sends.
- In buffered blocking sends, the sender simply copies the data into the designated buffer and returns after the copy operation has been completed. The data is copied at a buffer at the receiving end as well.
- Buffering alleviates idling at the expense of copying overheads.

Non-Buffered Blocking Message Passing Operations



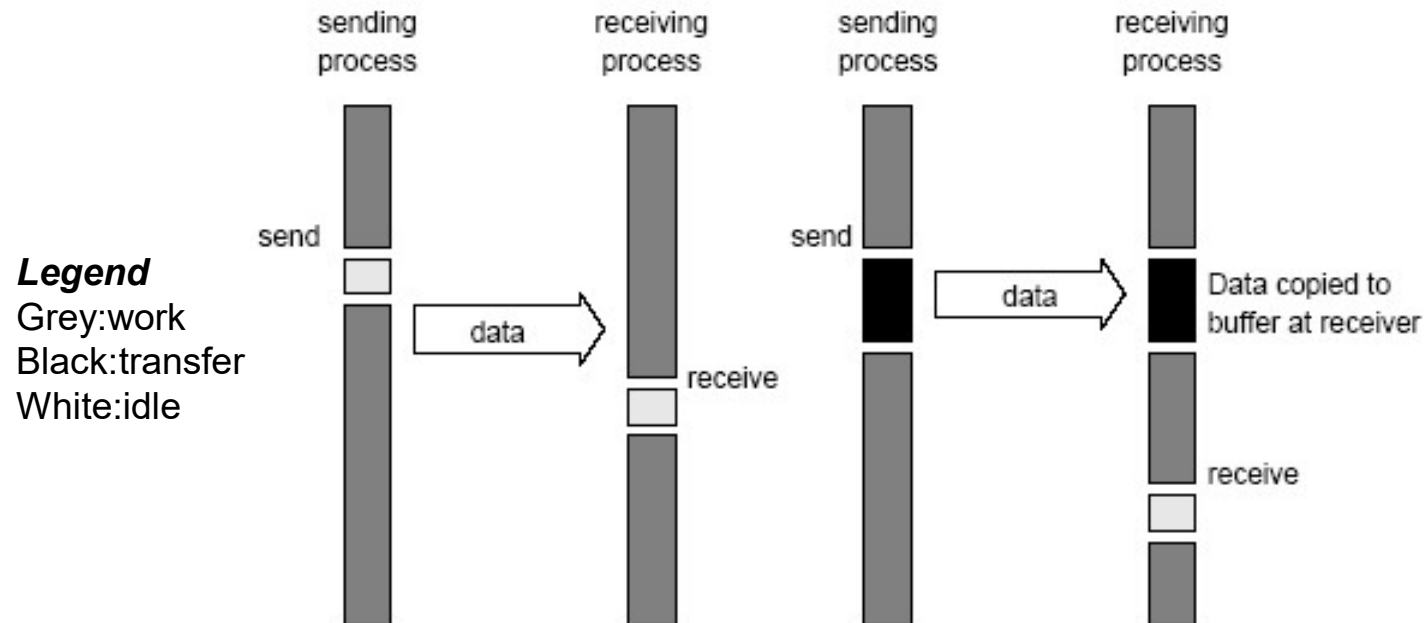
Handshake for a blocking non-buffered send/receive operation.

It is easy to see that in cases where sender and receiver do not reach communication point at similar times, there can be considerable idling overheads.

Buffered Blocking Message Passing Operations

- A simple solution to the idling and deadlocking problem outlined above is to rely on buffers at the sending and receiving ends.
- The sender simply copies the data into the designated buffer and returns after the copy operation has been completed.
- The data must be buffered at the receiving end as well.
- Buffering trades off idling overhead for buffer copying overhead.

Buffered Blocking Message Passing Operations



Blocking buffered transfer protocols: (a) in the presence of communication hardware with buffers at send and receive ends; and (b) in the absence of communication hardware, sender interrupts receiver and deposits data in buffer at receiver end.

Buffered Blocking Message Passing Operations

Bounded buffer sizes can have significant impact on performance.

P0	P1
for (i = 0; i < 1000; i++){	for (i = 0; i < 1000; i++){
produce_data(&a);	receive(&a, 1, 0);
send(&a, 1, 1);	consume_data(&a);
}	}

What if consumer was much slower than producer?

Buffered Blocking Message Passing Operations

Deadlocks are still possible with buffering since receive operations block.

P0

```
receive(&a, 1, 1);  
send(&b, 1, 1);
```

P1

```
receive(&a, 1, 0);  
send(&b, 1, 0);
```

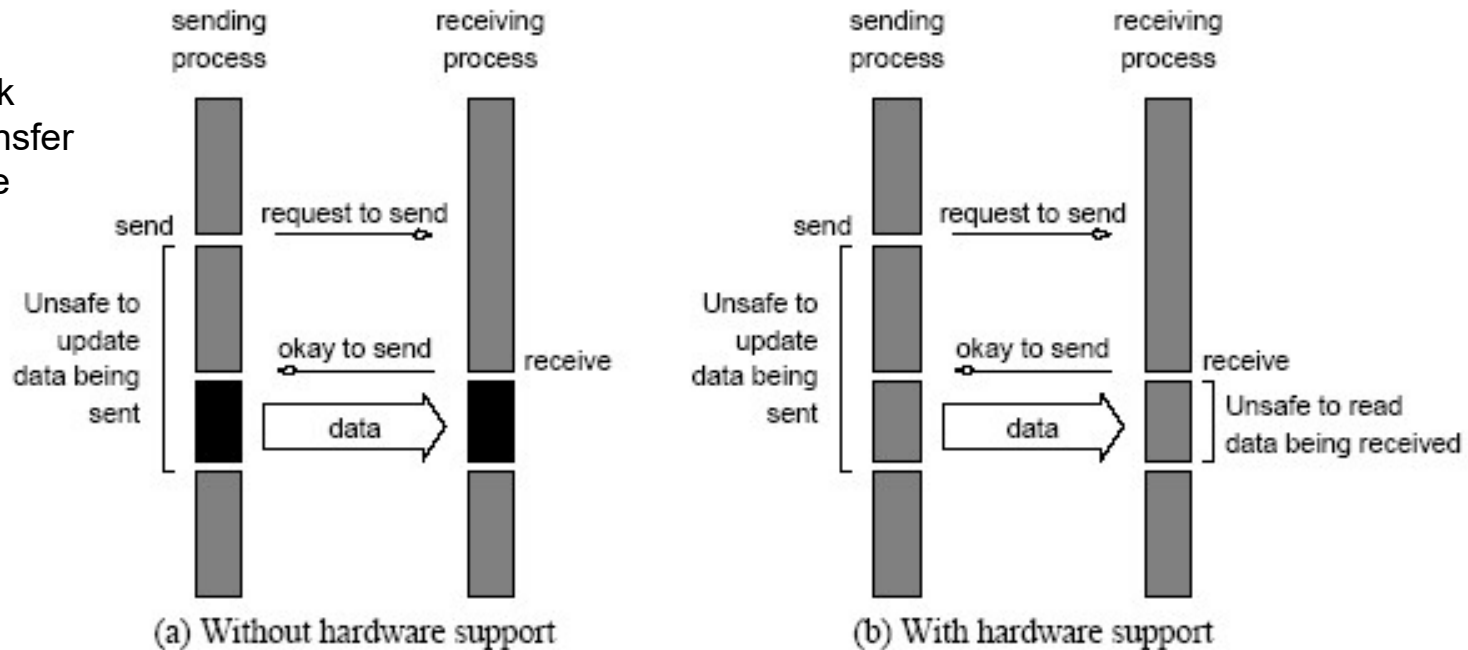
Non-Blocking Message Passing Operations

- The programmer must ensure semantics of the send and receive.
- This class of non-blocking protocols returns from the send or receive operation before it is semantically safe to do so.
- Non-blocking operations are generally accompanied by a check-status operation.
- When used correctly, these primitives are capable of overlapping communication overheads with useful computations.
- Message passing libraries typically provide both blocking and non-blocking primitives.

Non-Blocking Message Passing Operations

Legend

Grey:work
Black:transfer
White:idle



Non-blocking non-buffered send and receive operations (a) in absence of communication hardware; (b) in presence of communication hardware.

Send and Receive Protocols

	Blocking Operations	Non-Blocking Operations
Buffered	<p>Sending process returns after data has been copied into communication buffer</p>	<p>Sending process returns after initiating DMA transfer to buffer. This operation may not be completed on return</p>
Non-Buffered	<p>Sending process blocks until matching receive operation has been encountered</p> <p>Send and Receive semantics assured by corresponding operation</p>	<p>Programmer must explicitly ensure semantics by polling to verify completion</p>

Space of possible protocols for send and receive operations.

MPI Overview

1) MPI in a nutshell

- ❑ Introduction
- ❑ Programmer view
- ❑ First MPI example

2) Message passing operations

- ❑ Non-buffered
- ❑ Buffered
- ❑ Blocking
- ❑ Non-blocking

3) MPI basics

- ❑ Six fundamental functions
- ❑ Terminologies

4) Advanced MPI programming

- ❑ Avoid deadlock
- ❑ Topologies
- ❑ Collective communications
- ❑ Groups and communicators

MPI: the Message Passing Interface

- MPI defines a standard library for message-passing that can be used to develop portable message-passing programs using either C or Fortran.
- The MPI standard defines both the syntax as well as the semantics of a core set of library routines.
- Vendor implementations of MPI are available on almost all commercial parallel computers.
- It is possible to write fully-functional message-passing programs by using only the six routines.

MPI: the Message Passing Interface

The minimal set of MPI routines.

<code>MPI_Init</code>	Initializes MPI.
<code>MPI_Finalize</code>	Terminates MPI.
<code>MPI_Comm_size</code>	Determines the number of processes.
<code>MPI_Comm_rank</code>	Determines the label of calling process.
<code>MPI_Send</code>	Sends a message.
<code>MPI_Recv</code>	Receives a message.

Starting and Terminating the MPI Library

- `MPI_Init` is called prior to any calls to other MPI routines. Its purpose is to initialize the MPI environment.
- `MPI_Finalize` is called at the end of the computation, and it performs various clean-up tasks to terminate the MPI environment.
- The prototypes of these two functions are:

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

- `MPI_Init` also strips off any MPI related command-line arguments.
- All MPI routines, data-types, and constants are prefixed by “MPI_”. The return code for successful completion is `MPI_SUCCESS`.

Communicators

- A communicator defines a *communication domain* - a set of processes that are allowed to communicate with each other.
- Information about communication domains is stored in variables of type `MPI_Comm`.
- Communicators are used as arguments to all message transfer MPI routines.
- A process can belong to many different (possibly overlapping) communication domains.
- MPI defines a default communicator called `MPI_COMM_WORLD` which includes all the processes.

Querying Information

- The `MPI_Comm_size` and `MPI_Comm_rank` functions are used to determine the number of processes and the label of the calling process, respectively.

- The calling sequences of these routines are as follows:

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

- The rank of a process is an integer that ranges from zero up to the size of the communicator minus one.

Sending and Receiving Messages

- The basic functions for sending and receiving messages in MPI are the `MPI_Send` and `MPI_Recv`, respectively.
- The calling sequences of these routines are as follows:

```
int MPI_Send(void *buf, int count, MPI_Datatype
             datatype, int dest, int tag, MPI_Comm comm)
int MPI_Recv(void *buf, int count, MPI_Datatype
             datatype, int source, int tag,
             MPI_Comm comm, MPI_Status *status)
```

- MPI provides equivalent datatypes for all C datatypes. This is done for portability reasons.
- The datatype `MPI_BYTE` corresponds to a byte (8 bits) and `MPI_PACKED` corresponds to a collection of data items that has been created by packing non-contiguous data.
- The message-tag can take values ranging from zero up to the MPI defined constant `MPI_TAG_UB`.

MPI Datatypes

MPI Datatype	C Datatype
MPI_CHAR	signed char
MPI_SHORT	signed short int
MPI_INT	signed int
MPI_LONG	signed long int
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED_SHORT	unsigned short int
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long int

MPI Datatypes (cont.)

MPI Datatype	C Datatype
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double
MPI_BYTE	
MPI_PACKED	

Sending and Receiving Messages

- MPI allows specification of wildcard arguments for both source and tag.
- If source is set to `MPI_ANY_SOURCE`, then any process of the communication domain can be the source of the message.
- If tag is set to `MPI_ANY_TAG`, then messages with any tag are accepted.
- On the receive side, the message must be of length equal to or less than the length field specified.

Sending and Receiving Messages

- On the receiving end, the status variable can be used to get information about the `MPI_Recv` operation.
- The corresponding data structure contains:

```
typedef struct MPI_Status {  
    int MPI_SOURCE;  
    int MPI_TAG;  
    int MPI_ERROR; };
```

- The `MPI_Get_count` function returns the precise count of data items received.

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

MPI Overview

1) MPI in a nutshell

- ❑ Introduction
- ❑ Programmer view
- ❑ First MPI example

2) Message passing operations

- ❑ Non-buffered
- ❑ Buffered
- ❑ Blocking
- ❑ Non-blocking

3) MPI basics

- ❑ Six fundamental functions
- ❑ Terminologies

4) Advanced MPI programming

- ❑ Avoid deadlock
- ❑ Topologies
- ❑ Collective communications
- ❑ Groups and communicators

Avoiding Deadlocks

Consider:

```
int a[10], b[10], myrank;
MPI_Status status;
...
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank == 0) {
    MPI_Send(a, 10, MPI_INT, 1, 1, MPI_COMM_WORLD);
    MPI_Send(b, 10, MPI_INT, 1, 2, MPI_COMM_WORLD);
}
else if (myrank == 1) {
    MPI_Recv(b, 10, MPI_INT, 0, 2, MPI_COMM_WORLD);
    MPI_Recv(a, 10, MPI_INT, 0, 1, MPI_COMM_WORLD);
}
...
```

If MPI_Send is blocking, there is a deadlock.

Avoiding Deadlocks

Consider the following piece of code, in which process i sends a message to process $i + 1$ (modulo the number of processes) and receives a message from process $i - 1$ (modulo the number of processes).

```
int a[10], b[10], npes, myrank;
MPI_Status status;
...
MPI_Comm_size(MPI_COMM_WORLD, &npes);
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
         MPI_COMM_WORLD);
MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
         MPI_COMM_WORLD);
...
```

Once again, we have a deadlock if MPI_SEND is blocking.

Avoiding Deadlocks

We can break the circular wait to avoid deadlocks as follows:

```
int a[10], b[10], npes, myrank;
MPI_Status status;
...
MPI_Comm_size(MPI_COMM_WORLD, &npes);
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank%2 == 1) {
    MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
             MPI_COMM_WORLD);
    MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
             MPI_COMM_WORLD);
}
else {
    MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
             MPI_COMM_WORLD);
    MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
             MPI_COMM_WORLD);
}
...

```

Sending and Receiving Messages Simultaneously

To exchange messages, MPI provides the following function:

```
int MPI_Sendrecv(void *sendbuf, int sendcount,
    MPI_Datatype senddatatype, int dest, int
    sendtag, void *recvbuf, int recvcount,
    MPI_Datatype recvdatatype, int source, int recvtag,
    MPI_Comm comm, MPI_Status *status)
```

The arguments include arguments to the send and receive functions. If we wish to use the same buffer for both send and receive, we can use:

```
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)
```

Topologies and Embeddings

- MPI allows a programmer to organize processors into logical k -d meshes.
- The processor ids in `MPI_COMM_WORLD` can be mapped to other communicators (corresponding to higher-dimensional meshes) in many ways.
- The goodness of any such mapping is determined by the interaction pattern of the underlying program and the topology of the machine.
- MPI does not provide the programmer any control over these mappings.

Topologies and Embeddings

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

(a) Row-major mapping

0	4	8	12
1	5	9	13
2	6	10	14
3	7	11	15

(b) Column-major mapping

0	3	4	5
1	2	7	6
14	13	8	9
15	12	11	10

(c) Space-filling curve mapping

0	1	3	2
4	5	7	6
12	13	15	14
8	9	11	10

(d) Hypercube mapping

Different ways to map a set of processes to a two-dimensional grid. **(a)** and **(b)** show a row- and column-wise mapping of these processes, **(c)** shows a mapping that follows a space-filling curve (dotted line), and **(d)** shows a mapping in which neighboring processes are directly connected in a hypercube.

Creating and Using Cartesian Topologies

- We can create cartesian topologies using the function:

```
int MPI_Cart_create(MPI_Comm comm_old, int ndims,  
                   int *dims, int *periods, int reorder,  
                   MPI_Comm *comm_cart)
```

This function takes the processes in the old communicator and creates a new communicator with `dims` dimensions.

- Each processor can now be identified in this new cartesian topology by a vector of dimension `dims`.

Creating and Using Cartesian Topologies

- Since sending and receiving messages still require (one-dimensional) ranks, MPI provides routines to convert ranks to cartesian coordinates and vice-versa.

```
int MPI_Cart_coord(MPI_Comm comm_cart, int rank, int maxdims,  
                  int *coords)
```

```
int MPI_Cart_rank(MPI_Comm comm_cart, int *coords, int *rank)
```

- The most common operation on cartesian topologies is a shift. To determine the rank of source and destination of such shifts, MPI provides the following function:

```
int MPI_Cart_shift(MPI_Comm comm_cart, int dir, int s_step,  
                  int *rank_source, int *rank_dest)
```

Overlapping Communication with Computation

- In order to overlap communication with computation, MPI provides a pair of functions for performing non-blocking send and receive operations.

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

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

- These operations return before the operations have been completed. Function `MPI_Test` tests whether or not the non-blocking send or receive operation identified by its request has finished.

```
int MPI_Test(MPI_Request *request, int *flag,
              MPI_Status *status)
```

- `MPI_Wait` waits for the operation to complete.

```
int MPI_Wait(MPI_Request *request, MPI_Status *status)
```

Avoiding Deadlocks

Using non-blocking operations remove most deadlocks. Consider:

```
int a[10], b[10], myrank;
MPI_Status status;
...
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
if (myrank == 0) {
    MPI_Send(a, 10, MPI_INT, 1, 1, MPI_COMM_WORLD);
    MPI_Recv(a, 10, MPI_INT, 1, 2, &status, MPI_COMM_WORLD);
}
else if (myrank == 1) {
    MPI_Send(b, 10, MPI_INT, 0, 2, MPI_COMM_WORLD);
    MPI_Recv(b, 10, MPI_INT, 0, 1, &status, MPI_COMM_WORLD);
}
...
```

Replacing either the send or the receive operations with non-blocking counterparts fixes this deadlock.

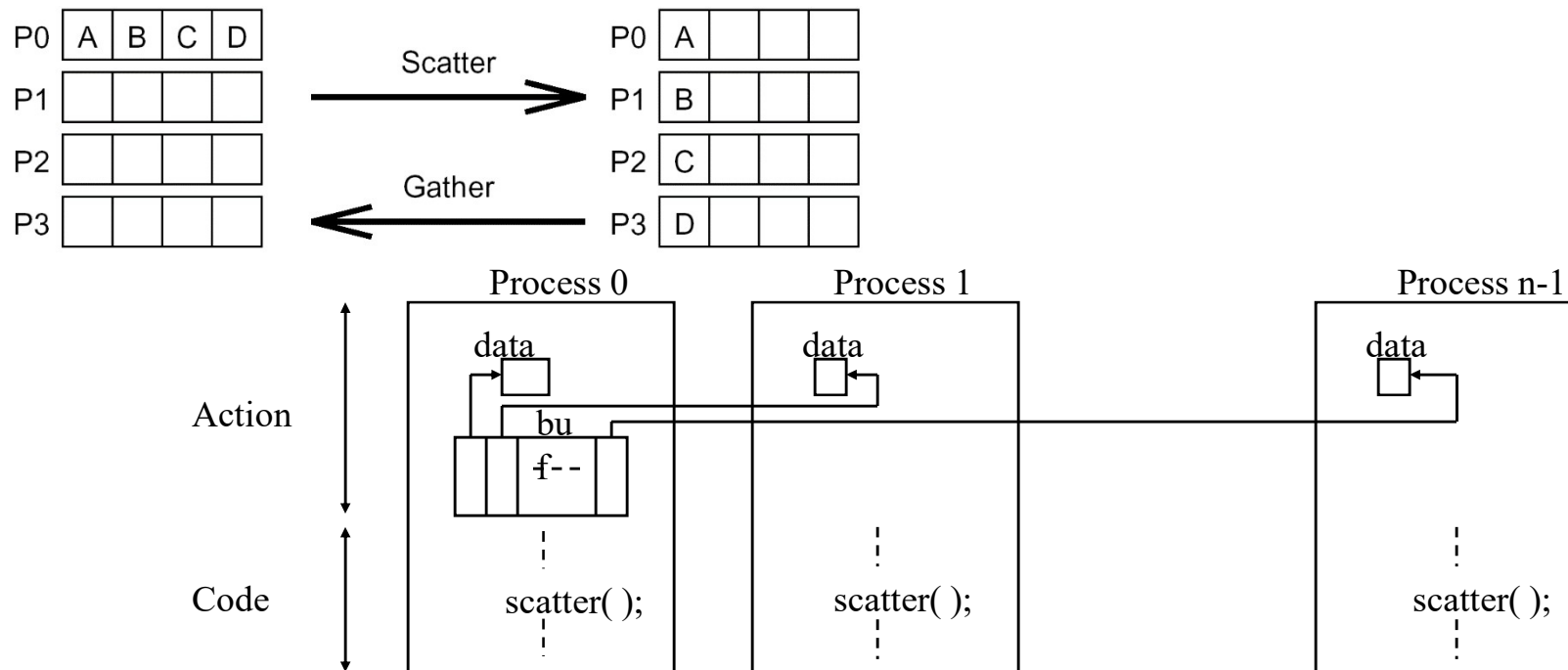
Collective Communication and Computation Operations

- MPI provides an extensive set of functions for performing common collective communication operations.
- Each of these operations is defined over a group corresponding to the communicator.
- All processors in a communicator must call these operations.

Collective Communication Operations

- The corresponding scatter operation is:

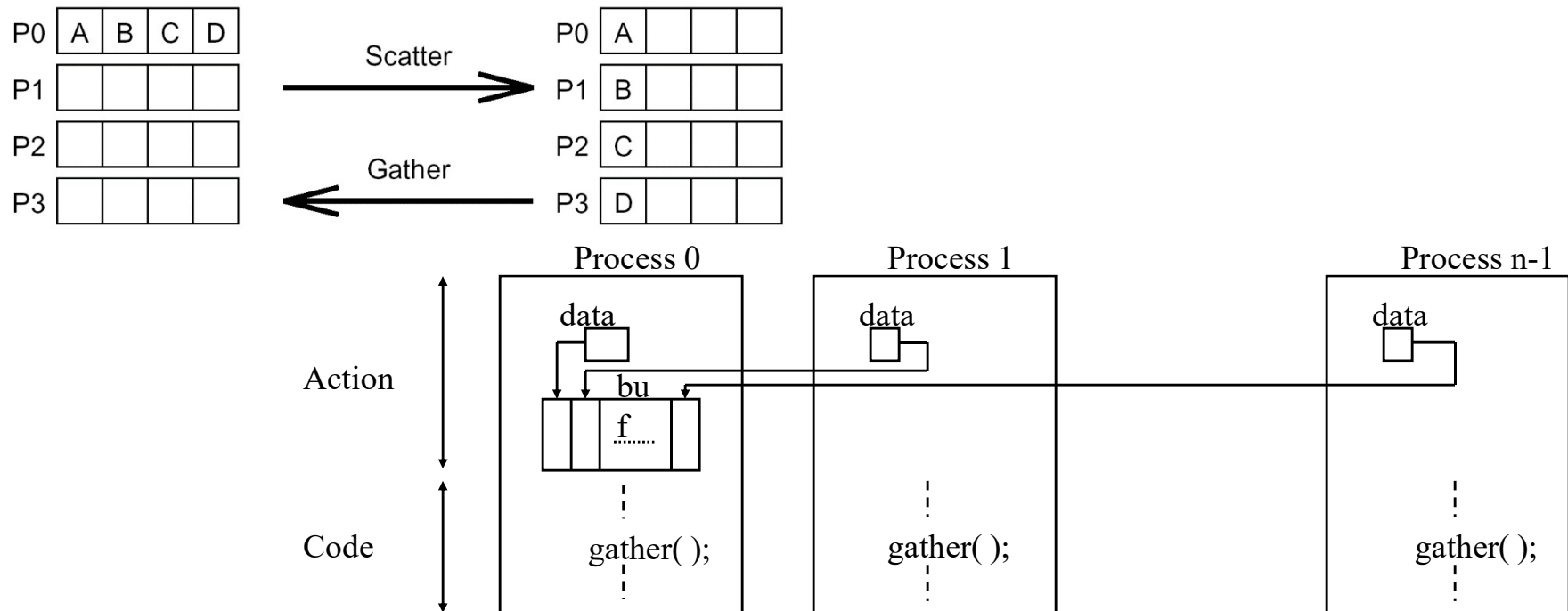
```
int MPI_Scatter(void *sendbuf, int sendcount,  
               MPI_Datatype senddatatype, void *recvbuf,  
               int recvcount, MPI_Datatype recvdatatype,  
               int source, MPI_Comm comm)
```



Collective Communication Operations

- The gather operation is performed in MPI using:

```
int MPI_Gather(void *sendbuf, int sendcount,  
              MPI_Datatype senddatatype, void *recvbuf,  
              int recvcount, MPI_Datatype recvdatatype,  
              int target, MPI_Comm comm)
```

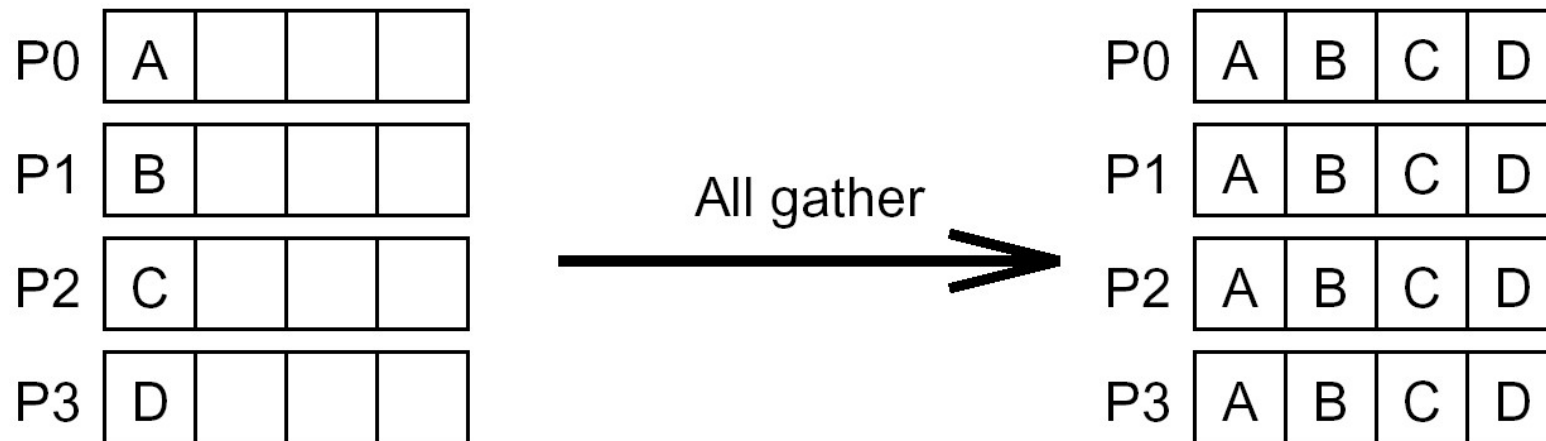


Gather Operation

Collective Communication Operations

- MPI also provides the MPI_Allgather function in which the data are gathered at all the processes.

```
int MPI_Allgather(void *sendbuf, int sendcount,  
                 MPI_Datatype senddatatype, void *recvbuf,  
                 int recvcount, MPI_Datatype recvdatatype,  
                 MPI_Comm comm)
```

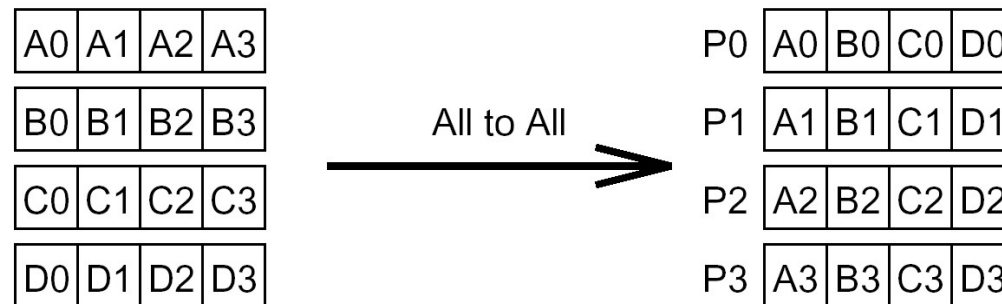


Collective Communication Operations

- The all-to-all personalized communication operation is performed by:

```
int MPI_Alltoall(void *sendbuf, int sendcount,  
                 MPI_Datatype senddatatype, void *recvbuf,  
                 int recvcount, MPI_Datatype recvdatatype,  
                 MPI_Comm comm)
```

- Using this core set of collective operations, a number of programs can be greatly simplified.



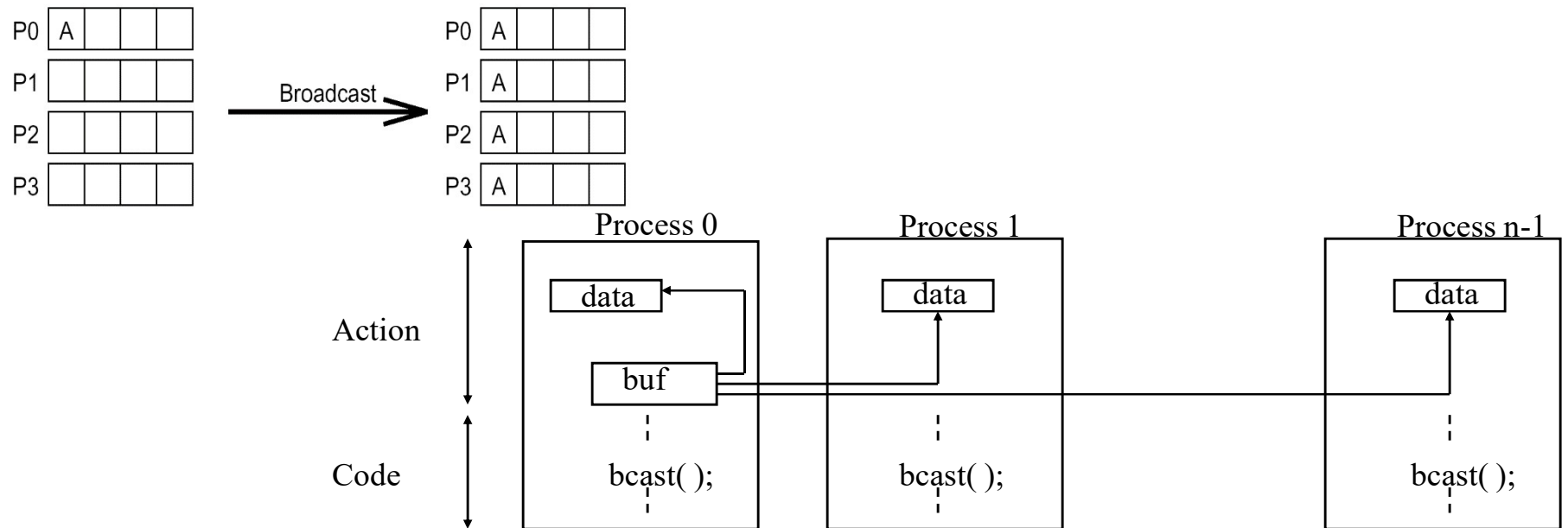
Collective Communication Operations

- The barrier synchronization operation is performed in MPI using:

```
int MPI_Barrier(MPI_Comm comm)
```

The one-to-all broadcast operation is:

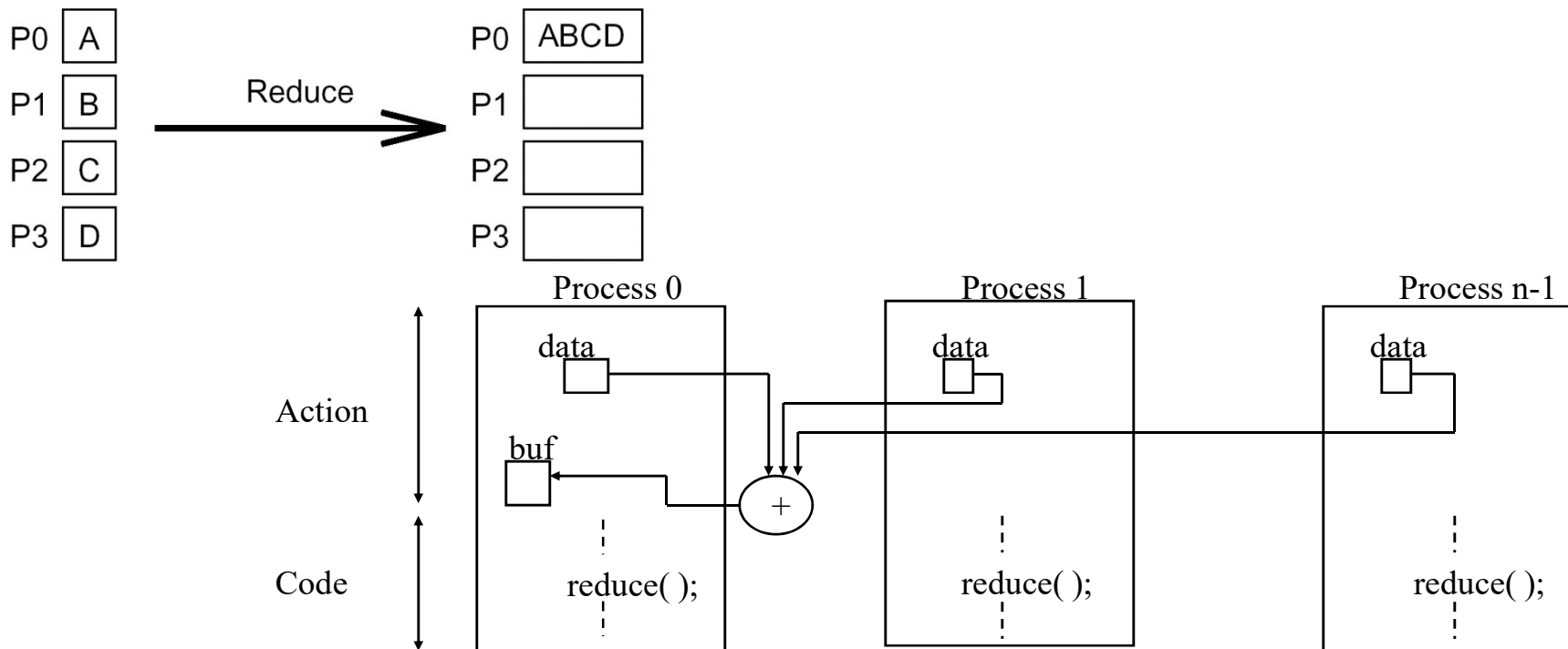
```
int MPI_Bcast(void *buf, int count, MPI_Datatype datatype,  
             int source, MPI_Comm comm)
```



Collective Communication Operations

- The all-to-one reduction operation is:

```
int MPI_Reduce(void *sendbuf, void *recvbuf, int count,
               MPI_Datatype datatype, MPI_Op op, int target,
               MPI_Comm comm)
```



Predefined Reduction Operations

Operation	Meaning	Datatypes
MPI_MAX	Maximum	C integers and floating point
MPI_MIN	Minimum	C integers and floating point
MPI_SUM	Sum	C integers and floating point
MPI_PROD	Product	C integers and floating point
MPI LAND	Logical AND	C integers
MPI_BAND	Bit-wise AND	C integers and byte

Predefined Reduction Operations

Operation	Meaning	Datatypes
MPI_LOR	Logical OR	C integers
MPI_BOR	Bit-wise OR	C integers and byte
MPI_LXOR	Logical XOR	C integers
MPI_BXOR	Bit-wise XOR	C integers and byte
MPI_MAXLOC	Max-min value location	Data-pairs
MPI_MINLOC	Min-min value location	Data-pairs

Collective Communication Operations

- The operation `MPI_MAXLOC` combines pairs of values (v_i, l_i) and returns the pair (v, l) such that v is the maximum among all v_i 's and l is the corresponding l_i (if there are more than one, it is the smallest among all these l_i 's).
- `MPI_MINLOC` does the same, except for minimum value of v_i .

Value	15	17	11	12	17	11
Process	0	1	2	3	4	5

`MinLoc(Value, Process) = (11, 2)`

`MaxLoc(Value, Process) = (17, 1)`

An example use of the `MPI_MINLOC` and `MPI_MAXLOC` operators.

Collective Communication Operations

MPI datatypes for data-pairs used with the `MPI_MAXLOC` and `MPI_MINLOC` reduction operations.

MPI Datatype	C Datatype
<code>MPI_2INT</code>	pair of ints
<code>MPI_SHORT_INT</code>	short and int
<code>MPI_LONG_INT</code>	long and int
<code>MPI_LONG_DOUBLE_INT</code>	long double and int
<code>MPI_FLOAT_INT</code>	float and int
<code>MPI_DOUBLE_INT</code>	double and int

Collective Communication Operations

- If the result of the reduction operation is needed by all processes, MPI provides:

```
int MPI_Allreduce(void *sendbuf, void *recvbuf,  
                  int count, MPI_Datatype datatype, MPI_Op op,  
                  MPI_Comm comm)
```

- To compute prefix-sums, MPI provides:

```
int MPI_Scan(void *sendbuf, void *recvbuf, int count,  
             MPI_Datatype datatype, MPI_Op op,  
             MPI_Comm comm)
```

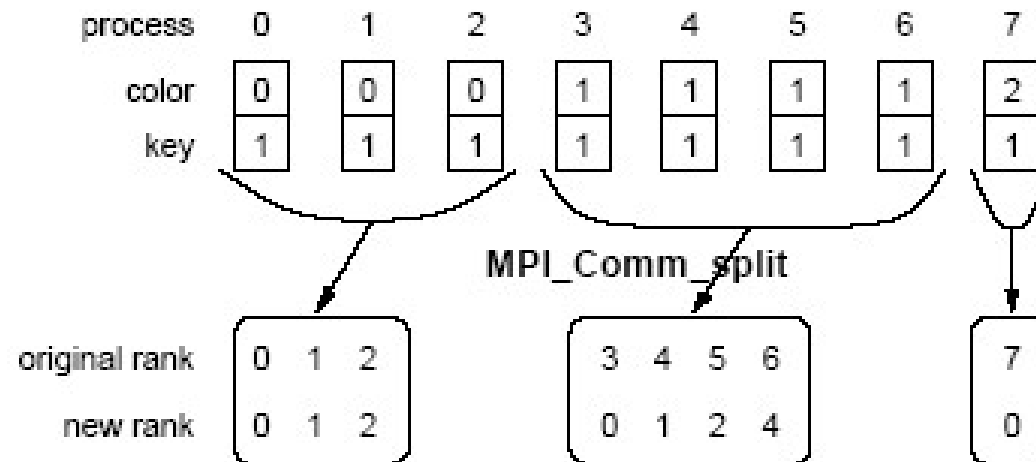
Groups and Communicators

- In many parallel algorithms, communication operations need to be restricted to certain subsets of processes.
- MPI provides mechanisms for partitioning the group of processes that belong to a communicator into subgroups each corresponding to a different communicator.
- The simplest such mechanism is:

```
int MPI_Comm_split(MPI_Comm comm, int color, int key,  
                  MPI_Comm *newcomm)
```

- This operation groups processors by color and sorts resulting groups on the key.

Groups and Communicators



Using `MPI_Comm_split` to split a group of processes in a communicator into subgroups.

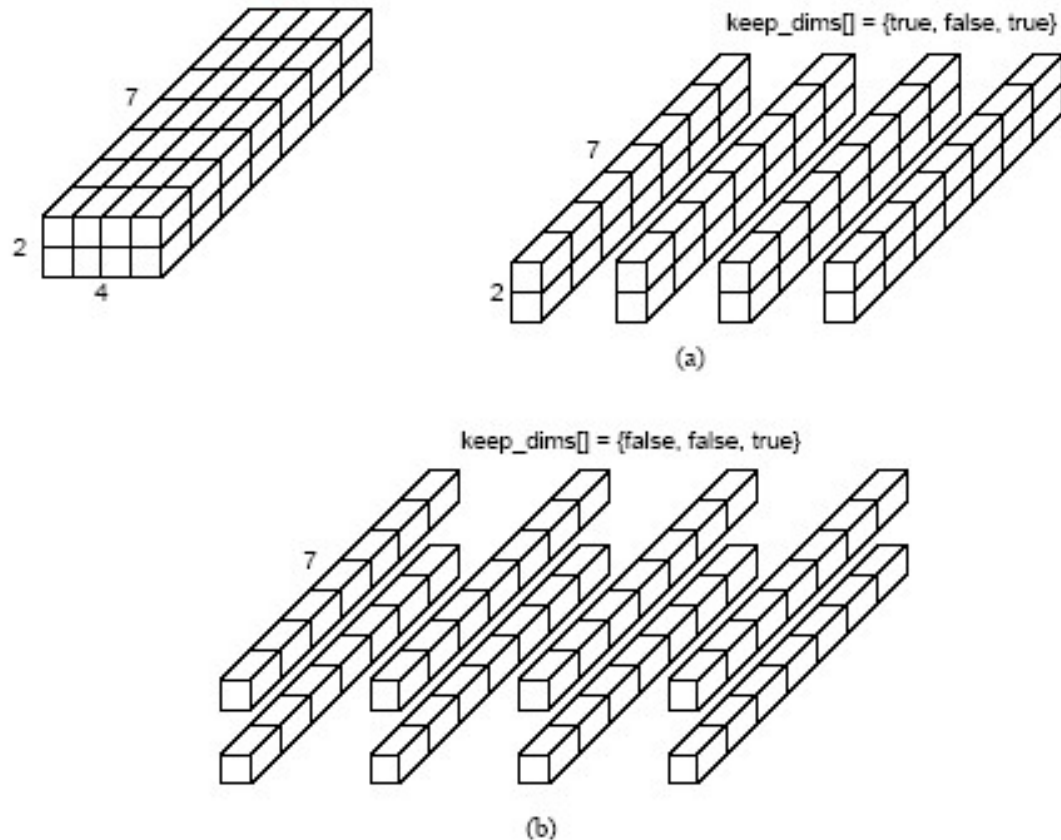
Groups and Communicators

- In many parallel algorithms, processes are arranged in a virtual grid, and in different steps of the algorithm, communication needs to be restricted to a different subset of the grid.
- MPI provides a convenient way to partition a Cartesian topology to form lower-dimensional grids:

```
int MPI_Cart_sub(MPI_Comm comm_cart, int *keep_dims,  
                MPI_Comm *comm_subcart)
```

- If `keep_dims[i]` is true (non-zero value in C) then the `i`th dimension is retained in the new sub-topology.
- The coordinate of a process in a sub-topology created by `MPI_Cart_sub` can be obtained from its coordinate in the original topology by disregarding the coordinates that correspond to the dimensions that were not retained.

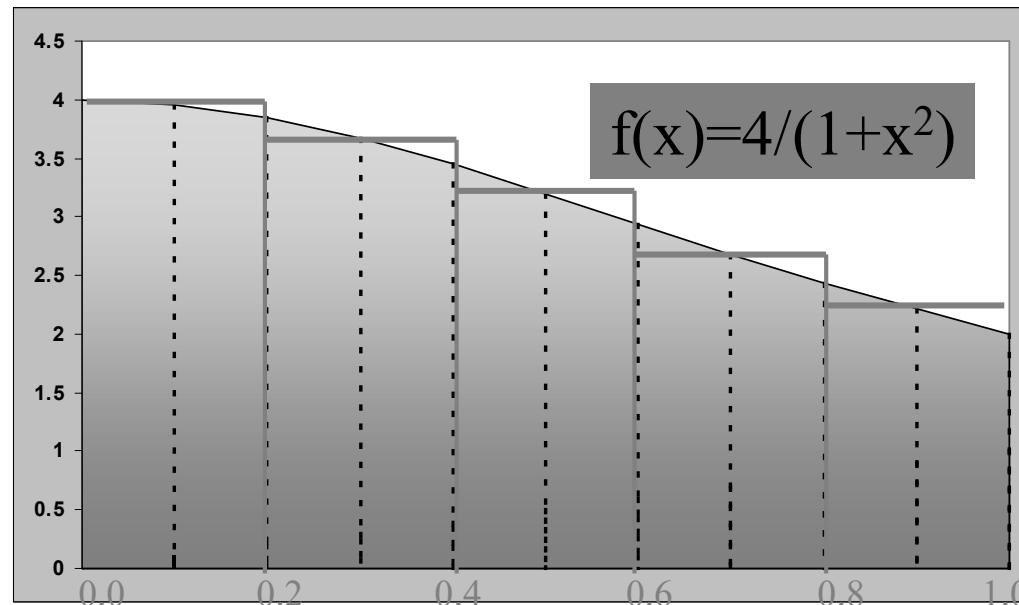
Groups and Communicators



Splitting a Cartesian topology of size $2 \times 4 \times 7$ into (a) four subgroups of size $2 \times 1 \times 7$, and (b) eight subgroups of size $1 \times 1 \times 7$.

Example using collective operations: Numerical Integration (computation of π)

- Integrate the function f (which equals π):



$$4 \int_0^1 \frac{1}{1+x^2}$$

- A perfect parallel program
 - it can be expressed with a minimum of communication
 - Load balancing is automatic
 - We can verify the answer

Implementation in C

```
#include "mpi.h"
#include <math.h>

int main (int argc, char *argv[]) {
    int n, myid, nprocs, i;
    double PI25DT = 3.141592653589793238462643
    double mypi, pi, h, sum, x;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &nprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &myid);
    while (1) {
        if (myid == 0) {
            printf("Enter the number of intervals: (0 quits) ");
            scanf("%d",&n);
        }

        MPI_Bcast(&n,1,MPI_INT,0,MPI_COMM_WORLD);
        if (n == 0)
            break;

        else{
            h = 1.0 / (double) n;
            sum = 0.0;
            for (i = myid + 1; i <= n ; i+= nprocs) {
                x = h * ((double) i - 0.5);
                sum += (4.0 / (1.0 + x*x));
            }
            mypi = h * sum;
            MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE,
                MPI_SUM, 0, MPI_COMM_WORLD);
            if (myid == 0) printf("pi is approximately %.16f,
                Error is %.16f\n",pi, fabs(pi-PI25DT));
        }
    }

    MPI_Finalize();
    return 0;
}
```