

# **Message Passing Interface**

## **Distributed-Memory Parallel Programming**

Orian Louant

# Motivations for Parallel Computing

Most applications today are parallel applications

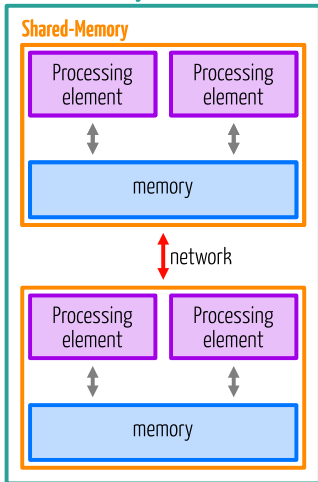
- in the years 2000's the CPU manufacturers have run out of room for boosting CPU performance
- instead of driving clock speeds and straight-line instruction throughput higher, they turn to hyperthreading and multicore architectures

In HPC in particular, it's crucial to be able to run your application in parallel

- in HPC, recent high-end CPUs have high core count: 36-64 cores
- high core count but lower clock, single core performance of an HPC CPU can be worse than your laptop CPU with turbo boost

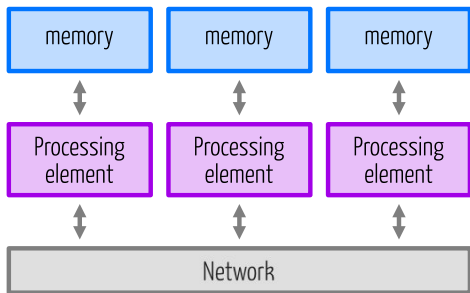
# Types of Parallel Systems

## Distributed-memory



- multiple compute nodes: distributed memory
- in HPC the dominant model for distributed memory programming is MPI, a standard that was introduced by the MPI Forum in May 1994 and updated in June 1995, last version 4.0 (2021).
- single compute node: shared memory
- in HPC the dominant model for shared-memory programming is OpenMP, a standard that was introduced in 1997 (Fortran) and 1998 (C/C++), last version is 5.1 (2020)

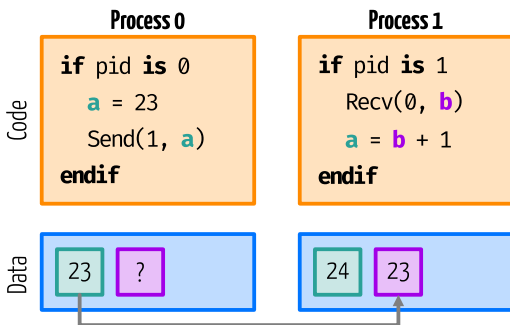
# Distributed-Memory



- MPI use processes to parallelize applications for distributed-memory systems
- the system consists of processing elements (nodes, CPUs, cores) and memory
- the processing elements have their own private memory
- the processing elements cannot have a direct access to the memory of the other processing elements
- the processing elements can communicate via a network

# Single Program Multiple Data

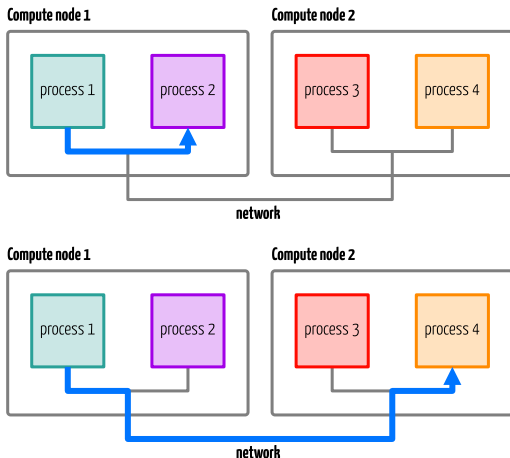
- The Single Program Multiple Data (SPMD) model is the main model for application message passing
- Multiple parallel processes run the same executable, they are identified by a unique identifier
- Each process has its own separate memory space and copy of the data
- Depending on their identifier, processes can follow different control paths



# Message Passing

The message passing model allows for inter-node and intra-node communication to co-exists

- Intra-node communication can leverage interconnect loopback capabilities or make use of a shared memory region
- In the case of inter-node communication, the messages are sent through a network



## What is MPI?

In scientific computing, the dominant paradigm for process parallelism is the single program multiple data model using MPI for interprocess communication

- MPI, which stands for Message Passing Interface, is a communication protocol for programming parallel computers
- The MPI standard was introduced by the MPI Forum in May 1994 and updated in June 1995. Version 3.1 of the standard has been approved in June 2015, and the last version (4.0) in June 2021
- It is a portable standard which defines the syntax and semantics of a core of library functions allowing the user to write message-passing programs
- It allows for both point-to-point and collective communication between processes

## What MPI provides

- **Management:** managing processes, communicators, creating datatypes, topologies, ...
- **Point-to-Point:** blocking and non-blocking communication based on send and receive
- **Collectives:** communications that involve a group or groups of processes
- **One-Sided:** where one process specifies all communication parameters, both for the sending side and for the receiving side (including shared memory)
- **I/O:** to write and read files to disk in parallel



# MPI Basics

# The Six-Functions MPI

The basics of MPI can be limited to six functions:

- `MPI_Init`: Initialize MPI
- `MPI_Comm_size`: Find out how many processes there are
- `MPI_Comm_rank`: Find out which process I am
- `MPI_Send`: Send a message
- `MPI_Recv`: Receive a message
- `MPI_Finalize`: Terminate MPI

These functions are accessible by including `mpi.h` (C/C++) or by using the `mpi` or `mpi_f08` module (Fortran).

## Initializing and Finalizing the MPI Environment

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

```
MPI_Init(ierror)
```

```
integer, optional, intent(out) :: ierror
```

Initialize the MPI environment. It must be called by **each MPI process, once and before any other MPI function call**. In C you can pass NULL for both the argc and argv arguments

```
MPI_Finalize()
```

```
MPI_Finalize(ierror)
```

```
integer, optional, intent(out) :: ierror
```

Terminates MPI execution environment. Once this function has been called, **no MPI function may be called afterward**.

# Initializing and Finalizing the MPI Environment

- the MPI standard does not say what a program can do before an `MPI_Init` or after an `MPI_Finalize`
- in general you should do as little as possible before calling `MPI_Init`
- these functions should be called by one thread only (the main thread)
- `MPI_Init` and `MPI_Finalize` should be called by the same thread

## Return Value of MPI Function/Routine

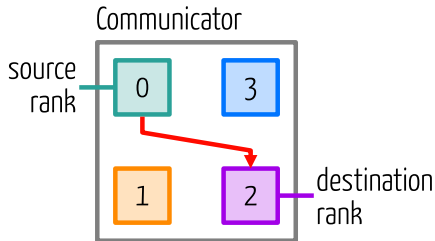
All MPI routines return an error value.

- In C this is the return value of the function
- In Fortran this is returned in the last argument (optional for the Fortran 2008 binding)

<code>MPI_SUCCESS</code>	Successful return code
<code>MPI_ERR_BUFFER</code>	Invalid buffer pointer
<code>MPI_ERR_COUNT</code>	Invalid count argument
<code>MPI_ERR_TYPE</code>	Invalid datatype argument
<code>MPI_ERR_TAG</code>	Invalid tag argument
<code>MPI_ERR_COMM</code>	Invalid communicator

## Communicator and Rank

- A communicator represents a group of processes that can communicate with each other
- Inside a communicator, each process is identified by an integer which is called a rank
- Each process has a unique rank inside a communicator but if a process is present in multiple communicators, it may have different ranks for each communicator.



## Communicator Size

To get the size of a communicator, i.e. the number of processes, use the `MPI_Comm_size` function.

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

```
MPI_Comm_size(communicator, size, ierror)  
  type(MPI_Comm), intent(in)      :: comm  
  integer, intent(out)            :: size  
  integer, optional, intent(out) :: ierror
```

After this function call, the size of the group associated with a communicator is stored in the variable `size`.

## Rank in a Communicator

The rank of a process in a communicator is obtained using the `MPI_Comm_rank` function.

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

```
MPI_Comm_rank(communicator, rank, ierror)  
  type(MPI_Comm), intent(in)      :: comm  
  integer, intent(out)            :: rank  
  integer, optional, intent(out) :: ierror
```

After this function call, the rank of the calling process inside the communicator is stored in the variable `rank`.



# The World Communicator

The MPI specification provides a predefined communicator: `MPI_COMM_WORLD`

- it allows communication with all processes that are accessible after MPI initialization
- this communicator act as the default communicator as most applications use a flat name space for processes as well as a single communication context

Using the `MPI_COMM_WORLD` communicator, you can retrieve

- the total number of processes with `MPI_Comm_size`
- the global rank of a process with `MPI_Comm_rank`

# MPI Hello Worlds

```
#include <stdio.h>
#include <mpi.h>

int main(int argc, char** argv) {
    int world_size, rank, name_len;
    char proc_name[MPI_MAX_PROCESSOR_NAME];

    MPI_Init(&argc, &argv);

    MPI_Comm_size(MPI_COMM_WORLD, &world_size);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Get_processor_name(proc_name, &name_len);

    printf("Hello world from rank %d "
           "out of %d on node %s.\n",
           rank, world_size, proc_name);

    MPI_Finalize();

    return 0;
}
```

```
program hello_world
    use mpi_f08

    implicit none

    integer :: rank, size, namelen, ierror
    character*(MPI_MAX_PROCESSOR_NAME) ::
        procname

    call MPI_Init(ierror)
    call MPI_Comm_size(MPI_COMM_WORLD, size)
    call MPI_Comm_rank(MPI_COMM_WORLD, rank)
    call MPI_Get_processor_name(procname,
        namelen)

    print 100, rank, size, procname(1:namelen)
100 format('Hello world from rank ', i0, &
&         ' out of ', i0, &
&         ' on node ', a, '.')

    call MPI_Finalize(ierror)
end
```

## Compile a MPI Application

To have access to the MPI tools, first load the **OpenMPI** module:

```
$ module load OpenMPI
```

Then you can compile your application using the **mpicc (mpifort)** compiler wrapper

```
$ mpicc -o hello_world hello_world.c  
$ mpifort -o hello_world hello_world.f90
```

**mpicc (mpifort)** is not a compiler, it is a wrapper: it adds all the relevant compiler or linker flags and then invoke the underlying compiler or linker. In our case it invokes **gcc (gfortran)**

## Running an MPI Application on a CÉCI Cluster

Create your job submission script. **n tasks** indicates the number of processes that we want to run.

```
#!/bin/bash
# Submission script for NIC5
#SBATCH --job-name=mpi-job
#SBATCH --time=00:01:00 # hh:mm:ss
#SBATCH --ntasks=4
#SBATCH --cpus-per-task=1
#SBATCH --mem-per-cpu=1024 # megabytes
```

```
module load OpenMPI
```

```
cd $SLURM_SUBMIT_DIR
```

```
# alternative: srun ./hello_world
```

```
mpirun ./hello_world
```

## Running an MPI application on a CÉCI Cluster

Save your job submission script and run it with **sbatch**

```
$ sbatch submit_mpi.job  
Submitted batch job <jobid>
```

[...]

```
$ cat mpi_hello.out  
Hello world from rank 0 out of 4 on node nic5-w030.  
Hello world from rank 2 out of 4 on node nic5-w035.  
Hello world from rank 3 out of 4 on node nic5-w037.  
Hello world from rank 1 out of 4 on node nic5-w034.
```

# MPI Programming on the CÉCI Cluster

MPI implementations are available on all the CÉCI. For example, for OpenMPI:

	Module	Execution
<b>Nic5</b>	<code>module load OpenMPI/4.1.2-GCC-11.2.0</code>	multi-node
<b>Lemaitre3</b>	<code>module load OpenMPI/4.1.2-GCC-11.2.0</code>	multi-node
<b>Hercules2</b>	<code>module load OpenMPI/4.1.2-GCC-11.2.0</code>	Restricted to one node
<b>Dragon2</b>	<code>module load OpenMPI/4.1.1-GCC-11.2.0</code>	Restricted to one node

Other implementations are available and can be used by loading the appropriate environment module: **impi** (Intel MPI) and **MPICH**

## Determine Where you are Running

The name of the processor on which a process is running can be queried by a call to the `MPI_Get_processor_name` function. Most of the time this function return the hostname.

```
MPI_Get_processor_name(char *name, int *namelen)
```

```
MPI_Get_processor_name(name, namelen, ierror)  
character(len=MPI_MAX_PROCESSOR_NAME), intent(out) :: name  
integer, intent(out) :: namelen  
integer, optional, intent(out) :: ierror
```

The maximum length of the name is defined by `MPI_MAX_PROCESSOR_NAME`. You should ensure that name storage space is at least this value. The actual length of the name is returned in `namelen`.

## The MPI\_Wtime function

An exception to the “all MPI functions return an error” is the `MPI_Wtime` function: it returns a double which represents a time stamp. The difference between two values obtained from this function allow you to compute the elapsed walltime between these two calls.

```
double start = MPI_Wtime();  
// [...]  
double end = MPI_Wtime();  
  
double elapsed = end - start;
```



# **Point to Point Communication**

## Point to Point Communication

So far, we didn't exchange messages between the processes that we spawned. The simplest type of communication is called point-to-point communication where

- the sender calls the send function and specifies the data to be sent
- the receiver calls the receive function and specify the data to be received
- the data is transferred as well as the associated metadata

In addition, we have to specify

- the communicator to use
- on the send side, the rank of the receiver in the communicator
- on the receive side, the rank of the sender in the communicator
- a tag to identify the message

## Sending a Message

Sending a message with MPI is done with the `MPI_Send` function.

<code>MPI_Send</code>	<code>(void*</code>	<code>buf,</code>	= address of the data you want to send
	<code>int</code>	<code>count,</code>	= number of elements to send
	<code>MPI_Datatype</code>	<code>datatype,</code>	= the type of data we want to send
	<code>int</code>	<code>dest,</code>	= the recipient of the message (rank)
	<code>int</code>	<code>tag,</code>	= identify the type of the message
	<code>MPI_Comm</code>	<code>comm)</code>	= the communicator used for this message

```
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
```

## Receiving a Message

Receiving a message with MPI is done with the `MPI_Recv` function.

<code>MPI_Recv</code>	<code>(void*</code> buf,	= where to receive the data
	<code>int</code> count,	= the receive buffer capacity
	<code>MPI_Datatype</code> datatype,	= the type of data we want to receive
	<code>int</code> source,	= the sender of the message (rank)
	<code>int</code> tag,	= identify the type of the message
	<code>MPI_Comm</code> comm,	= the communicator used for this message
	<code>MPI_Status*</code> status)	= informations about the message

```
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
```

## MPI Data types (C/C++)

MPI has a number of elementary data types, corresponding to the simple data types of the C programming language.

MPI Data Type	C Data Type	MPI Data Type	C Data Type
<code>MPI_CHAR</code>	<code>char</code>	<code>MPI_UNSIGNED_CHAR</code>	<code>unsigned char</code>
<code>MPI_INT</code>	<code>int</code>	<code>MPI_UNSIGNED</code>	<code>unsigned int</code>
<code>MPI_LONG</code>	<code>long int</code>	<code>MPI_UNSIGNED_LONG</code>	<code>unsigned long int</code>
<code>MPI_FLOAT</code>	<code>float</code>	<code>MPI_BYTE</code>	<code>unsigned char</code>
<code>MPI_DOUBLE</code>	<code>double</code>		

In addition you can create your own custom data type.

## MPI Data types (Fortran)

MPI has a number of elementary data types, corresponding to the simple data types of the Fortran programming language.

MPI Data Type	Fortran Data Type	MPI Data Type	Fortran Data Type
<code>MPI_INTEGER</code>	<code>integer</code>	<code>MPI_REAL4</code>	<code>real*4</code>
<code>MPI_REAL</code>	<code>real</code>	<code>MPI_REAL8</code>	<code>real*8</code>
<code>MPI_DOUBLE_PRECISION</code>	<code>double precision</code>	<code>MPI_INTEGER4</code>	<code>integer*4</code>
<code>MPI_COMPLEX</code>	<code>complex</code>	<code>MPI_INTEGER8</code>	<code>integer*8</code>
<code>MPI_LOGICAL</code>	<code>logical</code>	<code>MPI_DOUBLE_COMPLEX</code>	<code>double complex</code>
<code>MPI_CHARACTER</code>	<code>character(1)</code>		

In addition you can create your own custom data type.

# Simple Send and Receive Example

```
#include <stdio.h>
#include <mpi.h>

#define BUFSZ 5

int main(int argc, char* argv[]) {
    double buf[BUFSZ];

    int rank;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    if (rank == 0) {
        for(unsigned int i = 0; i < BUFSZ; i++) buf[i] = (double)i;
        MPI_Send(buf, BUFSZ, MPI_DOUBLE, 1, 99, MPI_COMM_WORLD);
    } else if(rank == 1) {
        MPI_Recv(buf, BUFSZ, MPI_DOUBLE, 0, 99, MPI_COMM_WORLD,
                MPI_STATUS_IGNORE);
    }

    printf("Process with rank %d: ", rank);
    for(unsigned int i = 0; i < BUFSZ; i++)
        printf("%5.11f", buf[i]);
    printf("\n");

    MPI_Finalize();

    return 0;
}
```

```
program main
    use mpi_f08

    implicit none

    integer, parameter :: bufsz = 5

    integer :: rank, i
    real    :: buf(bufsz)

    call MPI_Init()
    call MPI_Comm_rank(MPI_COMM_WORLD, rank)

    if (rank .eq. 0) then
        buf = (/ (real(i), i=1,bufsz) /)
        call MPI_Send(buf, bufsz, MPI_REAL, 1, 99, MPI_COMM_WORLD)
    else if (rank .eq. 1) then
        call MPI_Recv(buf, bufsz, MPI_REAL, 0, 99, MPI_COMM_WORLD,
                    MPI_STATUS_IGNORE)
    end if

    print 100, rank, buf
100  format('Process with rank ', i0, ': ', *(f5.1));

    call MPI_Finalize()
end
```

## Simple Send and Receive Example

```
$ mpicc -o send_recv_array send_recv_array.c
$ mpirun -np 2 ./send_recv_array
Process with rank 0:  0.0  1.0  2.0  3.0  4.0
Process with rank 1:  0.0  1.0  2.0  3.0  4.0

$ mpifort -o send_recv_array send_recv_array.f90
$ mpirun -np 2 ./send_recv_array
Process with rank 0:  0.0  1.0  2.0  3.0  4.0
Process with rank 1:  0.0  1.0  2.0  3.0  4.0
```



## About Communication Parameters

For the communication to succeed

- the source and destination ranks must be valid and use the same communicator
- the tags must match
- the receive buffer should be large enough to hold the message but the buffer may be larger than the data received

The **count** argument passed to the `MPI_Recv` function is the **maximum** number of elements of a certain MPI datatype that the buffer can contain. The number of data elements actually received may be lower

# The MPI Status Structure

Information (metadata) about the message can be obtained through the `MPI_Status` structure which contains the following fields

```
int MPI_SOURCE;    = source of the message
int MPI_TAG;      = tag of the message
int MPI_ERROR;    = error associated with the message
```

In Fortran, the `MPI_Status` type can be an integer array or a derived datatype depending on the binding you choose to use

```
use mpi
integer :: status(MPI_STATUS_SIZE)
...
src = status(MPI_SOURCE)
```

```
use mpi_f08
type(MPI_Status) :: status
...
src = status % MPI_SOURCE
```

## Get the Size of a message

The `MPI_Get_count` function gets the actual number elements received.

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

```
MPI_Get_count(status, datatype, count, ierror)
    type(MPI_Status), intent(in)    :: status
    type(MPI_Datatype), intent(in)  :: datatype
    integer, intent(out)           :: count
    integer, optional, intent(out) :: ierror
```

- The datatype argument should match the argument provided by the receive call that set the `status` variable.
- If the number of elements received exceeded the limits of the `count` parameter, then `MPI_Get_count` sets the value of count to `MPI_UNDEFINED`.

## Send and Receive with Get Count

```
char msg[20];
int recv_count;
MPI_Status status;

if (rank == 0) {
  strcpy(msg, "Hello mate!");
  MPI_Send(msg, strlen(msg)+1, MPI_CHAR, 1, tag,
           MPI_COMM_WORLD);

  printf("Process %d send: %s\n", rank, msg);
} else if (rank == 1) {
  MPI_Recv(msg, 20, MPI_CHAR, 0, tag,
           MPI_COMM_WORLD, &status);
  MPI_Get_count(&status, MPI_CHAR, &recv_count);

  printf("Process %d received: %s (size = %d)\n",
         rank, msg, recv_count);
}
```

```
character(len = 20) :: msg
integer :: rank, recv_count
type(MPI_Status) :: status

if (rank .eq. 0) then
  msg = 'Hello mate!'
  call MPI_Send(msg, 11, MPI_CHARACTER, 1, tag, &
               & MPI_COMM_WORLD)

  print 100, rank, msg
100  format('Process ', i0, ' send: ', a);
else if (rank .eq. 1) then
  call MPI_Recv(msg, 20, MPI_CHARACTER, 0, tag, &
               & MPI_COMM_WORLD, status)
  call MPI_Get_count(status, MPI_CHARACTER,
                    recv_count)

  print 200, rank, msg(1:recv_count), recv_count
200  format('Process ', i0, ' received: ', a, &
           & ' (size = ', i0, ')')
end if
```

## Simple Send and Receive Example

```
$ mpicc -o send_recv_get_count send_recv_get_count.c
$ mpirun -np 2 ./send_recv_get_count
Process 0 send: Hello mate!
Process 1 received: Hello mate! (size = 12)
```

In this example we set the maximum allowed length of the message to 20. This is the value we use on the receiver side, but the value used on the sender side was the actual size of the message. At the end we retrieve the actual length of the message using the `MPI_Get_count` function.

## Get the Size of a Message Before Receiving It

You can determine the size of a message before receiving it using a combination of the `MPI_Probe` and `MPI_Get_count` functions.

<code>MPI_Probe</code>	<code>(int source,</code>	<code>= the sender of the message (rank)</code>
	<code>int tag,</code>	<code>= identify the type of the message</code>
	<code>MPI_Comm comm)</code>	<code>= the communicator used for this message</code>
	<code>MPI_Status* status)</code>	<code>= informations about the message</code>

```
MPI_Probe(source, tag, comm, status, ierror)  
integer, intent(in)           :: source, tag  
type(MPI_Comm), intent(in)    :: comm  
type(MPI_Status)             :: status  
integer, optional, intent(out) :: ierror
```

# Get the Size of a Message Before Receiving It Example

```
MPI_Comm_size(MPI_COMM_WORLD, &size);
MPI_Comm_rank(MPI_COMM_WORLD, &rank);

if (rank == 0) {
    int buffer[3] = {123, 456, 789};
    printf("Process %d: sending 3 ints: %d, %d, %d\n",
           rank, buffer[0], buffer[1], buffer[2]);
    MPI_Send(buffer, 3, MPI_INT, 1, 10, MPI_COMM_WORLD);
} else if (rank == 1) {
    MPI_Status status;
    int count;

    MPI_Probe(0, 10, MPI_COMM_WORLD, &status);
    MPI_Get_count(&status, MPI_INT, &count);
    printf("Process %d retrieved the size of the message: %d.\n",
           rank, count);

    int* buffer = (int*)malloc(sizeof(int) * count);
    MPI_Recv(buffer, count, MPI_INT, 0, 10, MPI_COMM_WORLD, &status);
    ;
    printf("Process %d received message:", rank, count);
    for(int i = 0; i < count; ++i) printf(" %d", buffer[i]);
    printf(".\n");

    free(buffer);
}
```

```
call MPI_Init()
call MPI_Comm_size(MPI_COMM_WORLD, size)
call MPI_Comm_rank(MPI_COMM_WORLD, rank)

if (rank .eq. 0) then
    allocate(buffer(count))
    buffer = (/123, 456, 789/)

    print 100, rank, buffer
100    format('Process ', i0, ': sending 3 ints:', 3(1x,i0), '.')

    call MPI_Send(buffer, 3, MPI_INTEGER, 1, 10,
                  MPI_COMM_WORLD)
else if (rank .eq. 1) then
    call MPI_Probe(0, 10, MPI_COMM_WORLD, status)
    call MPI_Get_count(status, MPI_INTEGER, count)

    print 200, rank, count
200    format('Process ', i0, &
&        ' retrieved the size of the message: ', i0, '.')

    allocate(buffer(count))
    call MPI_Recv(buffer, count, MPI_INTEGER, 0, 10, &
&                MPI_COMM_WORLD, status)

    print 300, rank, buffer
300    format('Process ', i0, ' received message:', *(1x,i0), '.')
end if
```

## If You Know Nothing About The Message

You can receive a message with no prior information about the source, the tag or the size.

In order to receive such message, you can use a combination of `MPI_Probe` and `MPI_Status` as well as the `MPI_ANY_SOURCE` and `MPI_ANY_TAG` wildcards



## If You Know Nothing About The Message

```
if(rank == 0) {
    strcpy(sendbuf, "Hello Mate!");
    MPI_Send(sendbuf, strlen(sendbuf)+1, MPI_CHAR,
             1, 10, MPI_COMM_WORLD);
} else {
    MPI_Probe(MPI_ANY_SOURCE, MPI_ANY_TAG,
              MPI_COMM_WORLD, &status);
    MPI_Get_count(&status, MPI_CHAR, &msgsize);

    printf("Message incoming from process %d"
           " with tag %d and size %d\n"
           status.MPI_SOURCE, status.MPI_TAG, msgsize);

    if (msgsize != MPI_UNDEFINED)
        recvbuf = (char *)malloc(msgsize*sizeof(char));

    MPI_Recv(recvbuf, msgsize, MPI_CHAR,
             status.MPI_SOURCE,
             status.MPI_TAG,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);

    printf("Received message: %s\n", recvbuf);
}
```

```
if(rank .eq. 0) then
    sendbuf = 'Hello Mate!'
    call MPI_Send(sendbuf, LEN_TRIM(sendbuf),
                  MPI_CHAR, 1, 10, MPI_COMM_WORLD)
&
else
    call MPI_Probe(MPI_ANY_SOURCE, MPI_ANY_TAG,
                  MPI_COMM_WORLD, status)
&
    call MPI_Get_count(status, MPI_CHAR, msgsize)

    print 100, status % MPI_SOURCE,
           status % MPI_TAG, msgsize
100 format('Message incoming from process ', i0,
&
           ' with tag ', i0, ' and size ', i0)

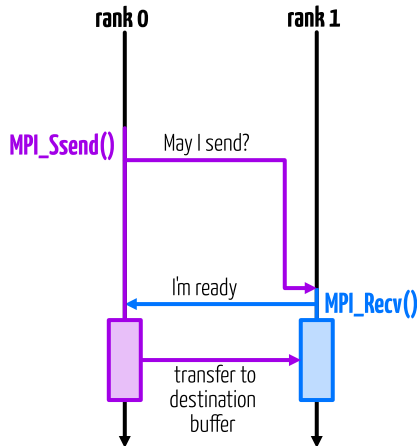
    call MPI_Recv(recvbuf, msgsize, MPI_CHAR,
                  status % MPI_SOURCE,
&
                  status % MPI_TAG,
&
                  MPI_COMM_WORLD, MPI_STATUS_IGNORE)

    print 200, recvbuf
200 format('Received message: ', a)
end if
```

# **Communication mode and Message Matching**

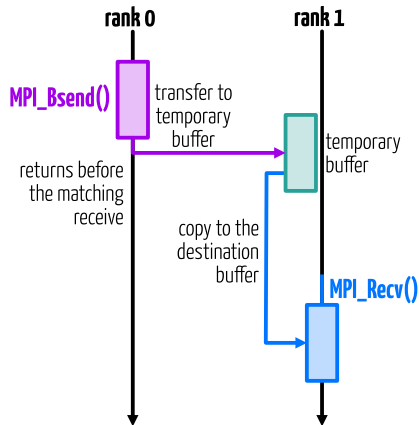
## Communications Mode: Rendezvous Protocol

- Using the rendezvous Protocol, communication does not complete at either end before both processes rendezvous at the communication
- At first, the sender sends the envelope of the message and the actual transfer occurs when a matching buffer is available on the receiving side
- This communication mode is also called synchronous



## Communications Mode: Eager Protocol

- Using the eager protocol, data is transferred to the receiver before a matching receive is posted
- This communication mode assumes that the destination can store the data
- This communication mode is also called buffered



## Communications Mode

- `MPI_Send` is blocking, it only returns when the send buffer is safe to reuse but this does not mean that the data has reached the receive buffer
- depending on the implementation and the size of the message, `MPI_Send` use the eager (buffered) or rendezvous (synchronous) protocol

Function	Description	Mode
<code>MPI_Send</code>	Blocking send	synchronous or buffered
<code>MPI_Ssend</code>	Synchronous send	synchronous
<code>MPI_Bsend</code>	Buffered send	buffered
<code>MPI_Rsend</code>	Ready Send	

## Message Matching

```
if (rank == 0) {  
    MPI_Ssend(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);  
    MPI_Ssend(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);  
} else if(rank == 1) {  
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
    MPI_Recv(&val2, 1, MPI_INT, 0, tag2, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

First Recv with tag 99 gets 12345

Second Recv with tag 1 gets 67890

First Ssend with tag 99 sent 12345

Second Ssend with tag 1 sent 67890

## Message Matching

```
if (rank == 0) {
    MPI_Ssend(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);
    MPI_Ssend(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);
} else if(rank == 1) {
    MPI_Recv(&val2, 1, MPI_INT, 0, tag2, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
```

With a synchronous send, this piece of code will produce a deadlock as there is no matching receive

- the first `MPI_Recv` is blocking until a message with `tag2` is received
- the first `MPI_Ssend` is blocking until its message with `tag1` received

## Message Matching

```
if (rank == 0) {  
    MPI_Bsend(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);  
    MPI_Bsend(&val2, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);  
} else if(rank == 1) {  
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
    MPI_Recv(&val2, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

First Bsend with tag 99 sent 12345  
Second Bsend with tag 99 sent 67890

First Recv with tag 99 gets 12345  
Second Recv with tag 99 gets 67890

**Message with the same tags are matched in order: the first message issued is received first, even in non-synchronous mode**



## Message Matching

```
if (rank == 0) {  
    MPI_Bsend(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);  
    MPI_Bsend(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);  
} else if(rank == 1) {  
    MPI_Recv(&val2, 1, MPI_INT, 0, tag2, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

First Bsend with tag 99 sent 12345  
Second Bsend with tag 1 sent 67890

First Recv with tag 1 gets 67890  
Second Recv with tag 99 gets 12345

**Message are matched by tags but they are not received in the same order as they were sent**

## Message Matching

```
if (rank == 0) {  
    MPI_Bsend(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);  
    MPI_Bsend(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);  
} else if(rank == 1) {  
    MPI_Recv(&val1, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
    MPI_Recv(&val2, 1, MPI_INT, 0, MPI_ANY_TAG, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

First Bsend with tag 99 sent 12345  
Second Bsend with tag 1 sent 67890

First Recv gets 12345  
Second Recv gets 67890

**We did not specify a tag at the receiving end: message are match in the order they were sent**

## Back to the Blocking-Send

```
if (rank == 0) {  
    MPI_Send(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);  
    MPI_Send(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);  
} else if(rank == 1) {  
    MPI_Recv(&val2, 1, MPI_INT, 0, tag2, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

In the previous examples, we have seen that the for the piece of code above

- a deadlock will occur if the synchronous mode is used
- the message will be matched by tag at the receiving end if the buffered mode is used

## Back to the Blocking-Send

The fact that the blocking send (`MPI_Send`) can change the communication protocol means that we might get different results depending on the implementation and the size of the message

```
if (rank == 0) {
    MPI_Send(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);
    MPI_Send(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);
} else if(rank == 1) {
    MPI_Recv(&val2, 1, MPI_INT, 0, tag2, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
```

A safer code will be

```
if (rank == 0) {
    MPI_Send(&val1, 1, MPI_INT, 1, tag1, MPI_COMM_WORLD);
    MPI_Send(&val2, 1, MPI_INT, 1, tag2, MPI_COMM_WORLD);
} else if(rank == 1) {
    MPI_Recv(&val1, 1, MPI_INT, 0, tag1, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    MPI_Recv(&val2, 1, MPI_INT, 0, tag2, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
```

# Communication Cost

# Communication Performance

There is a cost to communication: nothing is free.

$$T_{comm} = T_{latency} + \frac{n}{B_{peak}}$$

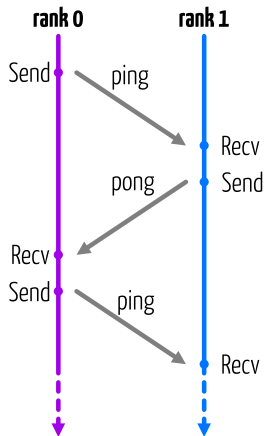
- $T_{latency}$ : inherent cost of communication (in s)
- $n$ : number of bytes to transfer
- $B_{peak}$ : asymptotic bandwidth of the network (in bytes/s)

From this we can compute the effective bandwidth (in bytes/s), i.e. the transfer rate for a given message size.

$$B_{eff} = \frac{n}{T_{latency} + \frac{n}{B_{peak}}}$$

# The Ping-Pong

We can visualize what is described theoretically in a real case: an MPI ping-pong application.



- We have two processes with respective rank 0 and 1
- process 0 sends a message to process 1 (ping)
- process 1 sends a message back to process 0 (pong)

We repeat this ping-pong 50 times and measure the time to determine the transfer time of one message with increasing message size.

# Ping Pong Code

```
for (int i = 0; i <= 27; i++) {
    // Actual code has a warmup loop
    double elapsed_time = -1.0 * MPI_Wtime();
    for (int i = 1; i <= 50; ++i) {
        if (rank == 0) {
            MPI_Send(A, N, MPI_DOUBLE, 1, 10, MPI_COMM_WORLD);
            MPI_Recv(A, N, MPI_DOUBLE, 1, 20, MPI_COMM_WORLD, &status);
        } else if (rank == 1) {
            MPI_Recv(A, N, MPI_DOUBLE, 0, 10, MPI_COMM_WORLD, &status);
            MPI_Send(A, N, MPI_DOUBLE, 0, 20, MPI_COMM_WORLD);
        }
    }
    elapsed_time += MPI_Wtime();

    long int num_bytes = 8 * N;
    double num_gbytes = (double)num_bytes / (double)bytes_to_gbytes;
    double avg_time_per_transfer = elapsed_time / (2.0 * 50);

    if(rank == 0)
        printf("Transfer size (B): %10li, Transfer Time (s): %15.9f, "
            "Bandwidth (GB/s): %15.9f\n",
            num_bytes, avg_time_per_transfer,
            num_gbytes/avg_time_per_transfer );
    free(A);
}
```

```
do i = 0,27
!   Actual code has a warmup loop
    elapsed_time = -1.0 * MPI_Wtime()
    do j = 1,50
        if (rank .eq. 0) then
            call MPI_Send(buffer, length, MPI_DOUBLE, 1, 10, &
                MPI_COMM_WORLD, ierror)
            call MPI_Recv(buffer, length, MPI_DOUBLE, 1, 20, &
                MPI_COMM_WORLD, status, ierror)
        else if (rank .eq. 1) then
            call MPI_Recv(buffer, length, MPI_DOUBLE, 0, 10, &
                MPI_COMM_WORLD, status, ierror)
            call MPI_Send(buffer, length, MPI_DOUBLE, 0, 20, &
                MPI_COMM_WORLD, ierror)
        end if
    end do
    elapsed_time = elapsed_time + MPI_Wtime()
    num_bytes = 8*length
    num_gbytes = dble(num_bytes) / bytes_to_gbytes;
    avg_time_per_transfer = elapsed_time / (2.0 * 50)

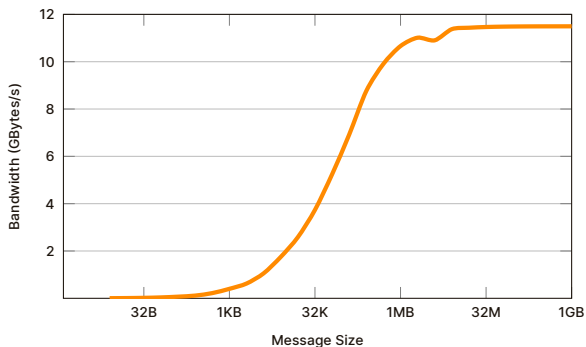
    if (rank .eq. 0) then
        print 100, num_bytes, avg_time_per_transfer, &
            num_gbytes/avg_time_per_transfer
    &
    100    format('Transfer size (B): ', i10, &
    &        ', Transfer Time (s): ', f15.9, &
    &        ', Bandwidth (GB/s): ', f15.9)
    end if

    length = length * 2
end do
```



# Ping Pong Result

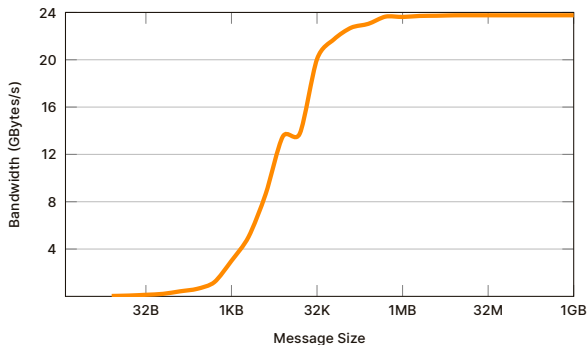
Ping-Pong Application running on NIC5



- With message size <1KB, the communication is dominated by the latency
- Transfer rate close to the theoretical performance of the network is observed for message size >10MB

## Ping Pong Result (bidirectional)

Ping-Pong Application running on NIC5  
(bidirectional)



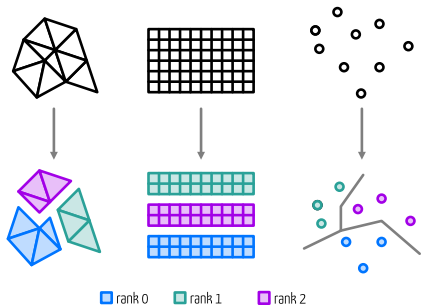
- Bidirectional bandwidth can be measured by issuing the send and receive at the same time
- The bandwidth is double compared to the mono-directional case
- General behavior with respect to the message size is the same as for the mono-directional case

# Parallelism with MPI

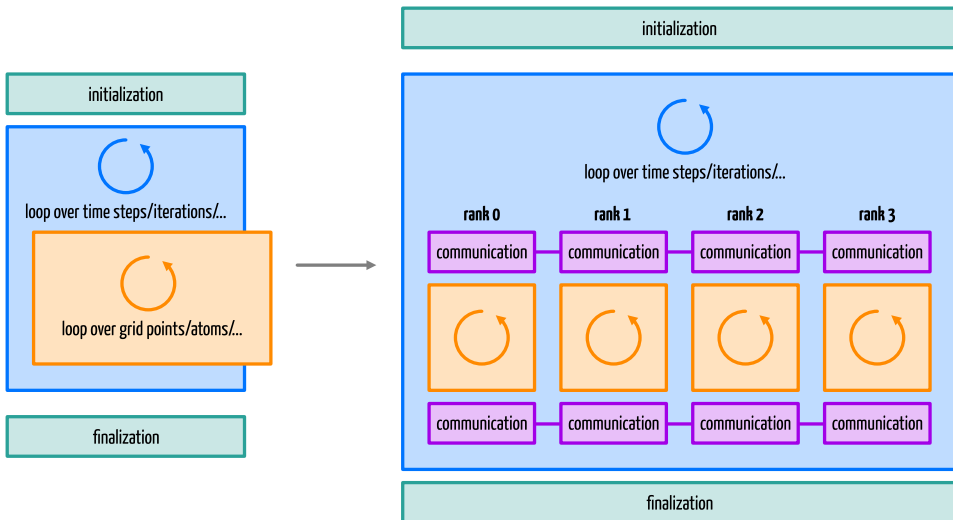
## Dividing the Problem

So far, we have focused our attention on the communication between the processes (ranks) but did not really address the question the parallelization.

- divide the problem in smaller problems that are processed by different ranks
- usually this subdivision is based on the rank and world size
- MPI also includes convenience routines to create virtual topologies



# Your Typical Scientific Application



# Sum of Integers: Serial

As a first example, we will consider an application that sums the integer between 1 and N. The serial implementation of such an application is

```
#include <stdio.h>

int main(int argc, char* argv[]) {
    const unsigned int N = 100000000;

    unsigned long sum = 0;
    for (unsigned int i = 1; i <= N; ++i)
        sum += i;

    printf("The sum of 1 to %d is %lu.\n", N, sum);

    return 0;
}
```

```
program main
    implicit none

    integer, parameter :: N = 100000000

    integer(kind=8) :: sum
    integer :: i

    sum = 0
    do i = 1, N
        sum = sum + i
    end do

    print 100, N, sum
100 format('The sum of 1 to ', i0, ' is ' i0 '.')

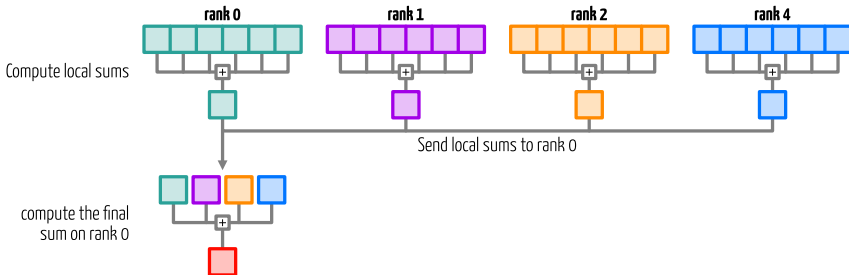
end
```

## Sum of Integers: Parallelization Strategy

The lower and upper bound for a rank can be computed using the following equations (assuming integer division):

$$start_i = \frac{N \cdot rank_i}{worldsize} + 1$$

$$end_i = \frac{N \cdot (rank_i + 1)}{worldsize}$$



# Sum of Integers: Computation

```
const unsigned int N = 10000000;  
  
int rank, size;  
MPI_Init(&argc, &argv);  
MPI_Comm_size(MPI_COMM_WORLD, &size);  
MPI_Comm_rank(MPI_COMM_WORLD, &rank);  
  
if (rank == 0)  
    printf("Running with %d processes.\n", size);  
  
unsigned int startidx = (N * rank / size) + 1;  
unsigned int  endidx = N * (rank+1) / size;  
  
unsigned long lsum = 0;  
for (unsigned int i = startidx; i <= endidx; ++i)  
    lsum += i;  
  
printf("Process %d has local sum %lu\n", rank, lsum);
```

```
integer, parameter :: N = 10000000  
integer(kind=8) :: lsum = 0, rsum = 0  
integer :: startidx, endidx, src, i, rank, wsize  
  
call MPI_Init()  
call MPI_Comm_size(MPI_COMM_WORLD, wsize)  
call MPI_Comm_rank(MPI_COMM_WORLD, rank)  
  
if (rank .eq. 0) then  
    print 100, wsize  
100    format('Running with ', i0, ' processes.')end if  
  
startidx = (N * rank / wsize) + 1  
endidx = N * (rank+1) / wsize  
  
do i = startidx, endidx  
    lsum = lsum + i  
end do  
  
print 200, rank, lsum  
200    format('Process ' i0, ' has local sum ', i0, '.')
```



# Sum of Integers: Communication

```
if (rank > 0) {
    MPI_Send(&lsum, 1, MPI_UNSIGNED_LONG, 0, 1,
             MPI_COMM_WORLD);
} else {
    unsigned long rsum;
    for(int src = 1; src < size; ++src) {
        MPI_Recv(&rsum, 1, MPI_UNSIGNED_LONG, src, 1,
                 MPI_COMM_WORLD, MPI_STATUS_IGNORE);
        lsum += rsum;
    }
}

if(rank == 0)
    printf("The sum of 1 to %d is %lu.\n", N, lsum);
```

```
if (rank .gt. 0) then
    call MPI_Send(lsum, 1, MPI_INTEGER8, 0, 1, &
& MPI_COMM_WORLD)
else
    do src = 1, wsize-1
        call MPI_Recv(rsum, 1, MPI_INTEGER8, src, 1, &
& MPI_COMM_WORLD, MPI_STATUS_IGNORE)
        lsum = lsum + rsum
    end do
end if

if (rank .eq. 0) then
    print 300, N, lsum
300 format('The sum of 1 to ', i0, ' is ' i0 '.')
end if
```

## What About Interfaces?

The previous example is almost an embarrassingly parallel problem

- little is needed to separate the problem into a number of parallel tasks
- the main loop iterations (the sum) are independent of each other
- communication is only needed at the end to compute the final sum

What about problems where we have dependence on values computed by the other processes?

## Diffusion Equation in 1D

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2}$$

After discretization and by using a forward difference in time and a central difference in space, the previous equation reads

$$\frac{\partial}{\partial t} u(x_i, t_n) = \alpha \frac{\partial^2}{\partial x^2} u(x_i, t_n) \rightarrow \frac{u_i^{n+1} - u_i^n}{\Delta t} = \alpha \frac{u_{i+1}^n - 2 \cdot u_i^n + u_{i-1}^n}{\Delta x^2}$$

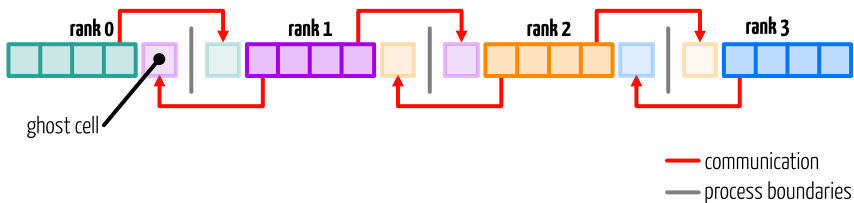
so that for each time step,  $u_i^{n+1}$  can be computed as

$$u_i^{n+1} = u_i^n + \alpha \frac{\Delta t}{\Delta x^2} (u_{i+1}^n - 2 \cdot u_i^n + u_{i-1}^n)$$

## Diffusion Equation in 1D

$$u_i^{n+1} = u_i^n + \alpha \frac{\Delta t}{\Delta x^2} (u_{i+1}^n - 2 \cdot u_i^n + u_{i-1}^n)$$

Each point depends on the values on the left and right: use ghost cells that are updated (via communication) at each time step



## Diffusion 1D (C)

The communication part is written so that we do the left-right communication first and then right-left.

```
const int left_rank = my_rank > 0           ? my_rank - 1 : MPI_PROC_NULL;
const int right_rank = my_rank < world_size-1 ? my_rank + 1 : MPI_PROC_NULL;

for (unsigned int iter = 1; iter <= NITERS; iter++) {
    MPI_Send(&uold[my_size], 1, MPI_DOUBLE, right_rank, 0, /* ... */);
    MPI_Recv(&uold[0],          1, MPI_DOUBLE, left_rank, 0, /* ... */);

    MPI_Send(&uold[1],          1, MPI_DOUBLE, left_rank, 1, /* ... */);
    MPI_Recv(&uold[my_size+1], 1, MPI_DOUBLE, right_rank, 1, /* ... */);

    for (unsigned int i = 1; i < my_size+1; i++)
        unew[i] = uold[i] + DIFF_COEF * dtdx2 * (uold[i+1] - 2.0 * uold[i] + uold[i-1]);

    // ...
}
```

## Diffusion 1D (Fortran)

The communication part is written so that we do the left-right communication first and then right-left.

```
left_rank = merge(my_rank - 1, MPI_PROC_NULL, my_rank > 0)
right_rank = merge(my_rank + 1, MPI_PROC_NULL, my_rank < world_size-1)

do iter = 1, niters
  call MPI_Send(uold(my_size), 1, MPI_DOUBLE_PRECISION, right_rank, 0, ...)
  call MPI_Recv(uold(0), 1, MPI_DOUBLE_PRECISION, left_rank, 0, ...)

  call MPI_Send(uold(1), 1, MPI_DOUBLE_PRECISION, left_rank, 1, ...)
  call MPI_Recv(uold(my_size+1), 1, MPI_DOUBLE_PRECISION, right_rank, 1, ...)

  do i = 1, my_size
    unew(i) = uold(i) + DIFF_COEF * dtdx2 * (uold(i+1) - 2.0 * uold(i) + uold(i-1))
  end do

! ...
end do
```

## Combined Sendrecv (C)

An alternative implementation could use combined send-receive

The send-receive operations combine in one operation the sending of a message to one destination and the receiving of another message, from another process

```
int MPI_Sendrecv(const void* sendbuf,           = address of the data you want to send
                 int sendcount,                 = number of elements to send
                 MPI_Datatype sendtype,         = the type of data we want to send
                 int dest,                      = the recipient of the message (rank)
                 int sendtag,                  = identify the type of the message
                 void* recvbuf,                = where to receive the data
                 int recvcount,                = the receive buffer capacity
                 MPI_Datatype recvtype,         = the type of data we want to receive
                 int source,                   = the sender of the message (rank)
                 int recvtag,                  = identify the type of the message
                 MPI_Comm comm,                = the communicator used for this message
                 MPI_Status* status);          = informations about the message
```

## Combined Sendrecv (Fortran)

The send-receive operations combine in one operation the sending of a message to one destination and the receiving of another message, from another process

```
MPI_Sendrecv(sendbuf, sendcount, sendtype, dest, sendtag, recvbuf,  
              recvcount, recvtype, source, recvtag, comm, status, ierror)  
type(*), dimension(..), intent(in) :: sendbuf  
integer, intent(in) :: sendcount, dest, sendtag, recvcount, source, recvtag  
type(MPI_Datatype), intent(in) :: sendtype, recvtype  
type(*), dimension(..) :: recvbuf  
type(MPI_Comm), intent(in) :: comm  
type(mpi_status) :: status  
integer, optional, intent(out) :: ierror
```



## Diffusion 1D: Sendrecv version

The 1D diffusion communication code can be rewritten using two `MPI_Sendrecv`: one for left-right communication and a second call for the right-left communication.

```
MPI_Sendrecv(&uold[my_size], 1, MPI_DOUBLE, right_rank, 0,  
            &uold[      0], 1, MPI_DOUBLE, left_rank,  0, /* ... */);
```

```
MPI_Sendrecv(&uold[      1], 1, MPI_DOUBLE, left_rank,  1,  
            &uold[my_size+1], 1, MPI_DOUBLE, right_rank, 1, /* ... */);
```

```
call MPI_Sendrecv(uold(my_size), 1, MPI_DOUBLE_PRECISION, right_rank, 0, &  
&uold(      0), 1, MPI_DOUBLE_PRECISION, left_rank,  0, ...)
```

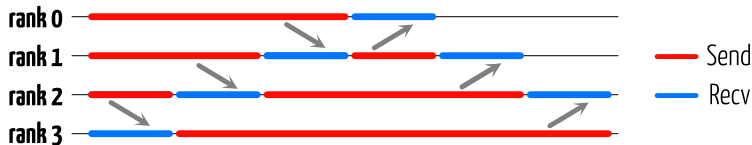
```
call MPI_Sendrecv(uold(      1), 1, MPI_DOUBLE_PRECISION, left_rank,  1, &  
&uold(my_size+1), 1, MPI_DOUBLE_PRECISION, right_rank, 1, ...)
```

# Diffusion 1D: Communication Serialization

```
MPI_Send(&uold[my_size], 1, ..., right_rank, 0, ... );  
MPI_Recv(&uold[0], 1, ..., left_rank, 0, ...);  
  
MPI_Send(&uold[1], 1, ..., left_rank, 1, ...);  
MPI_Recv(&uold[my_size+1], 1, ..., right_rank, 1, ...);
```

```
call MPI_Send(uold(my_size), 1, ..., right_rank, 0, ...)  
call MPI_Recv(uold(0), 1, ..., left_rank, 0, ...)  
  
call MPI_Send(uold(1), 1, ..., left_rank, 1, ...)  
call MPI_Recv(uold(my_size+1), 1, ..., right_rank, 1, ...)
```

Communication performance is poor: the communication is serial if the MPI implementation choose to use the synchronous mode



# The Serialization of Communication Problem

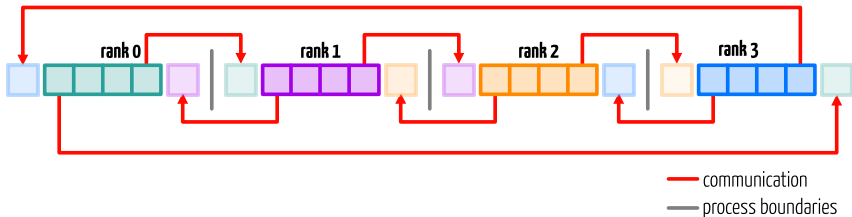
The problem with the serialization of the communication is that

- the processes we add, the longer the communication time
- it leads to a significant impact on the parallel efficiency

The best way to improve the performance is to use **non-blocking communication** where the send and/or receive calls return immediately without waiting for the communication to be completed

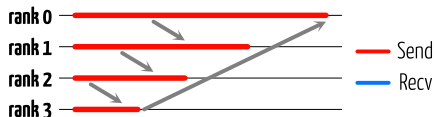
## Back to Diffusion Equation

Let's go back to the diffusion equation but this time with periodic boundary conditions



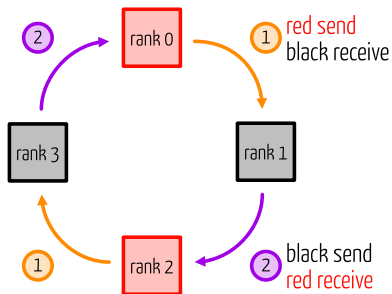
If we use the same communication pattern that we use for the non-periodic code, we end up with a deadlock

```
MPI_Send(..., left_rank, ...);  
MPI_Recv(..., right_rank, ...);
```



## Ring Communication: Odd-Even

One way to enable communication through a ring is to selectively reorder the send and receive calls



The communication is performed in two steps:

- even ranks send first to an odd rank process and then receive
- odd ranks receive first to an even rank process and then send

## Diffusion Equation: Odd-Even

With the odd-even pattern for communication, we can write a left to right transfer of the periodic diffusion equation that do not deadlock when using synchronous mode

```
if (my_rank%2 == 0) {  
    MPI_Send(&uold[1],          1, MPI_DOUBLE, left_rank, 1, ...);  
    MPI_Recv(&uold[my_size+1], 1, MPI_DOUBLE, right_rank, 1, ...);  
} else {  
    MPI_Recv(&uold[my_size+1], 1, MPI_DOUBLE, right_rank, 1, ...);  
    MPI_Send(&uold[1],          1, MPI_DOUBLE, left_rank, 1, ...);  
}
```

```
if (modulo(my_rank, 2) .eq. 0) then  
    call MPI_Send(uold(1),          1, MPI_DOUBLE_PRECISION, left_rank, 1, ...)  
    call MPI_Recv(uold(my_size+1), 1, MPI_DOUBLE_PRECISION, right_rank, 1, ...)  
else  
    call MPI_Recv(uold(my_size+1), 1, MPI_DOUBLE_PRECISION, right_rank, 1, ...)  
    call MPI_Send(uold(1),          1, MPI_DOUBLE_PRECISION, left_rank, 1, ...)  
end if
```

## Diffusion Equation: Odd-Even

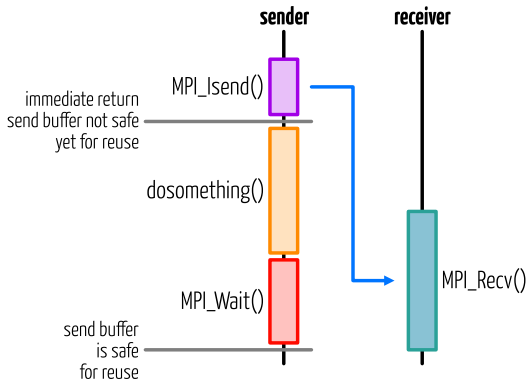
- with the odd-even communication pattern, our code is deadlock free but the communication is done in two steps introducing unnecessary latency and synchronization
- on our simple diffusion example, the impact on performance is limited but for more complex application this can have a larger negative impact

Fortunately, MPI provides an alternative to the blocking (or buffered) send and receive:

- non-blocking functions allow to perform asynchronous communication

# Non-Blocking Communication

The other way in which these send and receive operations can be done is by using the "I" functions. The "I" stands for Immediate returns and allow to perform the communication in three phases:



- call the non-blocking function `MPI_Isend` or `MPI_Irecv`
- Do some work
- Wait for the non-blocking communication to complete



## Non-Blocking Send

A non-blocking send is executed with the `MPI_Isend` function

<code>MPI_Isend</code>	<code>(void*</code>	<code>buf,</code>	= address of the data you want to send
	<code>int</code>	<code>count,</code>	= number of elements to send
	<code>MPI_Datatype</code>	<code>datatype,</code>	= the type of data we want to send
	<code>int</code>	<code>dest,</code>	= the recipient of the message (rank)
	<code>int</code>	<code>tag,</code>	= identify the type of the message
	<code>MPI_Comm</code>	<code>comm,</code>	= the communicator used for this message
	<code>MPI_Request*</code>	<code>request)</code>	= the handle on the non-blocking communication

```
MPI_Isend(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
```

## Non-Blocking Receive

A non-blocking receive is executed with the `MPI_Irecv` function

<code>MPI_Irecv</code>	<code>(void*</code> buf,	= where to receive the data
	<code>int</code> count,	= number of elements to send
	<code>MPI_Datatype</code> datatype,	= the type of data we want to receive
	<code>int</code> source,	= the sender of the message (rank)
	<code>int</code> tag,	= identify the type of the message
	<code>MPI_Comm</code> comm	= the communicator used for this message
	<code>MPI_Request*</code> request)	= the handle on the non-blocking communication

```
MPI_Irecv(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
```

## Waiting for a Non-Blocking Communication

An important feature of the non-blocking communication is the `MPI_Request` handle.

The value of this handle

- is generated by the non-blocking communication function
- is used by the `MPI_Wait` or `MPI_Test` functions

When using non-blocking communication, you have to be careful and avoid to

- modify the send buffer before the send operation completes
- read the receive buffer before the receive operation completes

## Waiting for a Non-Blocking Communication

The `MPI_Wait` returns when the operation identified by request is complete. This is a blocking function.

```
MPI_Wait(MPI_Request* request, = handle of the non-blocking communication  
        MPI_Status* status) = status of the completed communication
```

```
MPI_Wait(request, status, ierror)  
  type(MPI_Request), intent(inout) :: request  
  type(MPI_Status)           :: status  
  integer, optional, intent(out) :: ierror
```

## Wait on Multiple Requests

You can wait on multiple requests in one call with the `MPI_Waitall` function

```
MPI_Waitall(int count,                = number of request handlers to wait on
            MPI_Request array_of_requests[], = request handlers to wait on
            MPI_Status* array_of_statuses[]) = array in which write the statuses
```

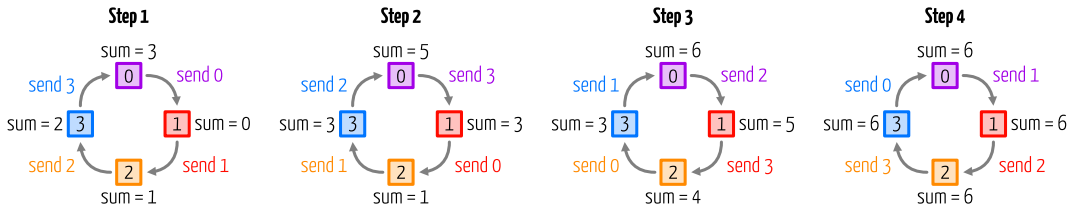
```
MPI_Waitall(count, array_of_requests, array_of_statuses, ierror)
integer, intent(in)                :: count
type(MPI_Request), intent(inout)    :: array_of_requests(count)
type(MPI_Status)                   :: array_of_statuses(*)
integer, optional, intent(out)     :: ierror
```

The *i*-th entry in `array_of_statuses` is set to the return status of the *i*-th operation

## Example: Ring Communication

As an example, let's consider the rotation of information inside a ring

- Each process stores its rank in the send buffer
- Each process sends its rank to its neighbour on the right
- Each process adds the received value to the sum
- Repeat



## Example: Ring Communication (blocking)

The blocking version of the ring communication is similar to the diffusion code that we considered earlier

```
sum = 0;
send_buf = rank;

for(int i = 0; i < world_size; ++i) {
    if (rank % 2 == 0) {
        MPI_Send(&send_buf, 1, MPI_INT, right_rank, ...);
        MPI_Recv(&recv_buf, 1, MPI_INT, left_rank, ...);
    } else {
        MPI_Recv(&recv_buf, 1, MPI_INT, left_rank, ...);
        MPI_Send(&send_buf, 1, MPI_INT, right_rank, ...);
    }

    send_buf = recv_buf;
    sum += recv_buf;
}
```

```
sum = 0
send_buf = rank

do i = 1,world_size
    if (modulo(rank, 2) .eq. 0) then
        call MPI_Send(send_buf, 1, MPI_INTEGER, right_rank, ...)
        call MPI_Recv(recv_buf, 1, MPI_INTEGER, left_rank, ...)
    else
        call MPI_Recv(recv_buf, 1, MPI_INTEGER, left_rank, ...)
        call MPI_Send(send_buf, 1, MPI_INTEGER, right_rank, ...)
    end if

    send_buf = recv_buf
    sum = sum + recv_buf
end do
```

## Example: Ring Communication (non-blocking, C)

For the non-blocking version, we use `MPI_Isend` instead of `MPI_Send`

- no need to separate the processes in two groups for the code to be correct:  
`MPI_Isend` returns immediately so the blocking receive is called
- we use `MPI_Wait` to make sure that the send buffer is safe to be reused

```
sum = 0;
send_buf = rank;

for(int i = 0; i < world_size; ++i) {
    MPI_Isend(&send_buf, 1, MPI_INT, right_rank, ..., &request);
    MPI_Recv(&recv_buf, 1, MPI_INT, left_rank, ...);

    MPI_Wait(&request, MPI_STATUS_IGNORE);

    send_buf = recv_buf;
    sum += recv_buf;
}
```



## Example: Ring Communication (non-blocking, Fortran)

For the non-blocking version, we use `MPI_Isend` instead of `MPI_Send`

- no need to separate the processes in two groups for the code to be correct:  
`MPI_Isend` returns immediately so the blocking receive is called
- we use `MPI_Wait` to make sure that the send buffer is safe to be reused

```
sum = 0
send_buf = rank

do i = 1, world_size
  call MPI_Isend(send_buf, 1, MPI_INTEGER, right_rank, ..., request)
  call MPI_Recv(recv_buf, 1, MPI_INTEGER, left_rank, ...)

  call MPI_Wait(request, MPI_STATUS_IGNORE)

  send_buf = recv_buf
  sum = sum + recv_buf
end do
```

## Diffusion Using Non-Blocking Communication (C)

Like the communication around a ring, rewriting the periodic diffusion code with non-blocking avoid to do the communication in two steps

```
for (int iter = 1; iter <= NITERS; iter++) {
    MPI_Irecv(&uold[0],          1, MPI_DOUBLE, left_rank,  0, MPI_COMM_WORLD, &reqs[0]);
    MPI_Irecv(&uold[my_size+1], 1, MPI_DOUBLE, right_rank, 1, MPI_COMM_WORLD, &reqs[1]);

    MPI_Isend(&uold[my_size], 1, MPI_DOUBLE, right_rank, 0, MPI_COMM_WORLD, &reqs[2]);
    MPI_Isend(&uold[1],        1, MPI_DOUBLE, left_rank,  1, MPI_COMM_WORLD, &reqs[3]);

    MPI_Waitall(4, reqs, MPI_STATUSES_IGNORE);

    for (int i = 1; i <= my_size; i++) {
        unew[i] = uold[i] + alpha * (uold[i+1] - 2.0 * uold[i] + uold[i-1]);
    }

    SWAP_PTR(uold, unew);
}
```

## Diffusion Using Non-Blocking Communication (Fortran)

Like the communication around a ring, rewriting the periodic diffusion code with non-blocking avoid to do the communication in two steps

```
do iter = 1, niters
  call MPI_Isend(uold(my_size), 1, MPI_DOUBLE_PRECISION, right_rank, 0, MPI_COMM_WORLD, reqs(0))
  call MPI_Isend(uold(1), 1, MPI_DOUBLE_PRECISION, left_rank, 1, MPI_COMM_WORLD, reqs(1))

  call MPI_Irecv(uold(0), 1, MPI_DOUBLE_PRECISION, left_rank, 0, MPI_COMM_WORLD, reqs(2))
  call MPI_Irecv(uold(my_size+1), 1, MPI_DOUBLE_PRECISION, right_rank, 1, MPI_COMM_WORLD, reqs(3))

  call MPI_Waitall(4, reqs, MPI_STATUSES_IGNORE)

  do i = 1, my_size
    unew(i) = uold(i) + alpha * (uold(i+1) - 2.0 * uold(i) + uold(i-1))
  end do

  uold(1:my_size) = unew(1:my_size)
end do
```

## Better Hide the Communications

As discussed before, communication comes at a cost. You should try to hide this cost as much as possible:

- if you want to use a large number of processes, you should try to hide this cost as much as possible
- achieved by overlapping communication and computation with non-blocking communication
- let the CPU do useful science instead of waiting for a communication to complete

The basic idea is to

- start the communication as soon as possible
- only wait at a point where you need to use the result of the communication or reuse a buffer involved in a communication

# Overlap Computation and Communication (C)

The 1D diffusion can be restructured so that overlap computation and communication

```
for (int iter = 1; iter <= NITERS; iter++) {
    MPI_Irecv(&uold[0],          1, MPI_DOUBLE, left_rank,  0, MPI_COMM_WORLD, &recv_reqs[0]);
    MPI_Irecv(&uold[my_size+1], 1, MPI_DOUBLE, right_rank, 1, MPI_COMM_WORLD, &recv_reqs[1]);

    MPI_Isend(&uold[my_size], 1, MPI_DOUBLE, right_rank, 0, MPI_COMM_WORLD, &send_reqs[0]);
    MPI_Isend(&uold[1],        1, MPI_DOUBLE, left_rank,  1, MPI_COMM_WORLD, &send_reqs[1]);

    for (int i = 2; i <= my_size-1; i++) {
        unew[i] = uold[i] + alpha * (uold[i+1] - 2.0 * uold[i] + uold[i-1]);
    }

    MPI_Waitall(2, recv_reqs, MPI_STATUSES_IGNORE);

    unew[    1] = uold[    1] + alpha * (uold[    2] - 2.0 * uold[    1] + uold[    0]);
    unew[my_size] = uold[my_size] + alpha * (uold[my_size+1] - 2.0 * uold[my_size] + uold[my_size-1]);

    MPI_Waitall(2, send_reqs, MPI_STATUSES_IGNORE);
    SWAP_PTR(uold, unew);
}
```

# Overlap Computation and Communication (Fortran)

The 1D diffusion can be restructured so that overlap computation and communication

```
do iter = 1, niters
  call MPI_Isend(uold(my_size), 1, MPI_DOUBLE_PRECISION, right_rank, 0, MPI_COMM_WORLD, send_reqs(1))
  call MPI_Isend(uold(1), 1, MPI_DOUBLE_PRECISION, left_rank, 1, MPI_COMM_WORLD, send_reqs(2))

  call MPI_Irecv(uold(0), 1, MPI_DOUBLE_PRECISION, left_rank, 0, MPI_COMM_WORLD, recv_reqs(1))
  call MPI_Irecv(uold(my_size+1), 1, MPI_DOUBLE_PRECISION, right_rank, 1, MPI_COMM_WORLD, recv_reqs(2))

  do i = 2, my_size-1
    unew(i) = uold(i) + alpha * (uold(i+1) - 2.0 * uold(i) + uold(i-1))
  end do

  call MPI_Waitall(2, recv_reqs, MPI_STATUSES_IGNORE)

  unew(1) = uold(1) + alpha * (uold(2) - 2.0 * uold(1) + uold(0))
  unew(my_size) = uold(my_size) + alpha * (uold(my_size+1) - 2.0 * uold(my_size) + uold(my_size-1))

  call MPI_Waitall(2, send_reqs, MPI_STATUSES_IGNORE)
  uold(1:my_size) = unew(1:my_size)
end do
```

# **Collective Communication**

## Collective Communication

So far, we have covered the topic of point-to-point communication: with a message that is exchanged between a sender and a receiver. However, in a lot of applications, collective communication may be required:

- **All-To-All:** All processes contribute to the result and all processes receive the result
- **All-To-One:** all processes contribute to the result and one process receives the result
- **One-To-All:** one process contributes to the result and all processes receive the result

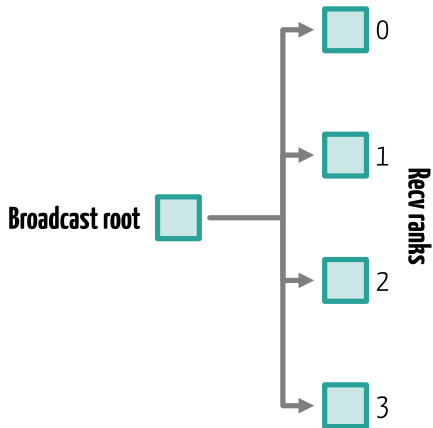


# Collective Communication

- the key argument of the collective communication routine is the communicator that define the group of participating processes
- the amount of data sent must exactly match the amount of data specified by the receiver
  
- **Broadcast:** Send data to all the processes
- **Scatter:** Distribute data between the processes
- **Gather:** Collect data from multiple processes to one process
- **Reduce:** Perform a reduction

## Broadcast

During a broadcast, one process (the root) sends the same data to all processes in a communicator.



# Broadcast

<code>MPI_Bcast</code>	<code>(void*</code>	<code>buffer,</code>	= address of the data you want to broadcast
	<code>int</code>	<code>count,</code>	= number of elements to broadcast
	<code>MPI_Datatype</code>	<code>datatype,</code>	= the type of data we want to broadcast
	<code>int</code>	<code>root,</code>	= rank of the broadcast root
	<code>MPI_Comm</code>	<code>comm)</code>	= the communicator used for this broadcast

```
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
```

# Broadcast Example

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

int rank;
MPI_Comm_rank(MPI_COMM_WORLD, &rank);

int bcast_root = 0;

int value;
if(rank == bcast_root) {
    value = 12345;
    printf("I am the broadcast root with rank %d "
           "and I send value %d.\n", rank, value);
}

MPI_Bcast(&value, 1, MPI_INT, bcast_root, MPI_COMM_WORLD);

if(rank != bcast_root) {
    printf("I am a broadcast receiver with rank %d "
           "and I obtained value %d.\n", rank, value);
}

MPI_Finalize();
```

```
integer :: i, value, ierror
integer :: rank, bcast_root

call MPI_Init(ierror)
call MPI_Comm_rank(MPI_COMM_WORLD, rank, ierror)

bcast_root = 0

if (rank .eq. bcast_root) then
    value = 12345
    print 100, rank, value
    format('I am the broadcast root with rank ', i0, &
           ' and I send value ', i0)
end if

call MPI_Bcast(value, 1, MPI_INT, bcast_root, &
               MPI_COMM_WORLD, ierror)

if (rank .ne. bcast_root) then
    print 200, rank, value
    format('I am a broadcast receiver with rank ', i0, &
           ' and I obtained value ', i0)
end if

call MPI_Finalize(ierror)
```

## Broadcast Example

```
$ mpicc -o broadcast broadcast.c
```

```
$ mpirun -np 4 ./broadcast
```

```
I am the broadcast root with rank 0 and I send value 12345.
```

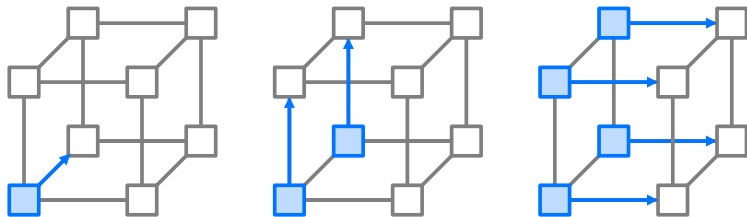
```
I am a broadcast receiver with rank 2 and I obtained value 12345.
```

```
I am a broadcast receiver with rank 1 and I obtained value 12345.
```

```
I am a broadcast receiver with rank 3 and I obtained value 12345.
```

## Broadcast hypercube

A one to all broadcast can be visualized on a hypercube of  $d$  dimensions with  $d = \log_2 p$ .

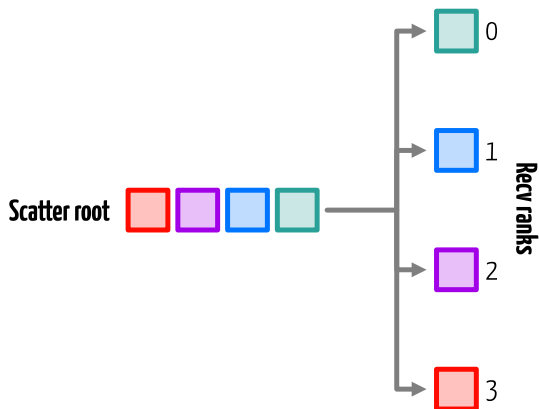


The broadcast procedure involves  $\log_2 p$  point-to-point simple message transfers.

$$T_{broadcast} = \left( T_{latency} + \frac{n}{B_{peak}} \right) \log_2 p$$

## MPI Scatter

During a scatter, the elements of an array are distributed in the order of process rank.



## MPI Scatter

<code>MPI_Scatter</code>	<code>(void*</code>	<code>sendbuf,</code>	<code>= address of the data you want to scatter</code>
	<code>int</code>	<code>sendcount,</code>	<code>= number of elements sent to each process</code>
	<code>MPI_Datatype</code>	<code>sendtype,</code>	<code>= the type of data we want to scatter</code>
	<code>void*</code>	<code>recvbuf,</code>	<code>= where to receive the data</code>
	<code>int</code>	<code>recvcount,</code>	<code>= number of elements to receive</code>
	<code>MPI_Datatype</code>	<code>recvtype,</code>	<code>= the type of data we want to receive</code>
	<code>int</code>	<code>root,</code>	<code>= rank of the scatter root</code>
	<code>MPI_Comm</code>	<code>comm)</code>	<code>= the communicator used for this scatter</code>

```
MPI_Scatter(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype,  
            root, comm, ierror)  
type(*), dimension(..), intent(in) :: sendbuf  
integer, intent(in)                :: sendcount, recvcount, root  
type(MPI_Datatype), intent(in)     :: sendtype, recvtype  
type(*), dimension(..)             :: recvbuf  
type(MPI_Comm), intent(in)         :: comm  
integer, optional, intent(out)     :: ierror
```



# Scatter Example

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

int size, rank, value, scatt_root = 0;
MPI_Comm_size(MPI_COMM_WORLD, &size);
MPI_Comm_rank(MPI_COMM_WORLD, &rank);

int* data = NULL;
if(rank == scatt_root) {
    data = (int*)malloc(sizeof(int)*size);

    printf("Values to scatter from process %d:", rank);
    for (int i = 0; i < size; i++) {
        data[i] = 100 * i;
        printf(" %d", data[i]);
    }
    printf("\n");
}

MPI_Scatter(data, 1, MPI_INT, &value, 1, MPI_INT,
           scatt_root, MPI_COMM_WORLD);
printf("Process %d received value %d.\n", rank, value);

if(rank == scatt_root) free(data);

MPI_Finalize();
```

```
integer :: size, rank, ierror
integer :: value, scatt_root, i
integer, dimension(:), allocatable :: buffer

call MPI_Init(ierror)
call MPI_Comm_size(MPI_COMM_WORLD, size, ierror)
call MPI_Comm_rank(MPI_COMM_WORLD, rank, ierror)

if (rank .eq. scatt_root) then
    allocate(buffer(size))
    do i = 1,size
        buffer(i) = 100 * i
    end do

    print 100, rank, data
    format('Values to scatter from process ', i0, &
          ':', *(1x,i0))
end if

call MPI_Scatter(buffer, 1, MPI_INTEGER, value, 1,
                MPI_INTEGER, &
                scatt_root, MPI_COMM_WORLD, ierror)

print 200, rank, value
format('Process ', i0, ' received value ', i0)

if (rank .eq. scatt_root) deallocate(buffer)

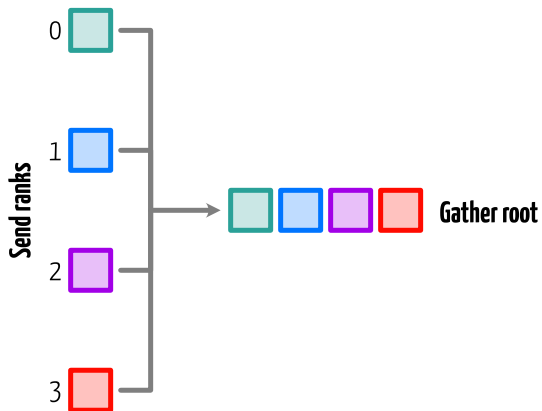
call MPI_Finalize(ierror)
```

## Scatter Example

```
$ mpicc -o scatter scatter.c
$ mpirun -np 4 ./scatter
Values to scatter from process 0: 0 100 200 300
Process 1 received value 100.
Process 2 received value 200.
Process 0 received value 0.
Process 3 received value 300.
```

## MPI Gather

A gathering is taking elements from each process and gathers them to the root process.



# MPI Gather

<code>MPI_Gather</code>	<code>(void*</code>	<code>sendbuf,</code>	<code>= address of the data you want to gather</code>
	<code>int</code>	<code>sendcount,</code>	<code>= number of elements to gather</code>
	<code>MPI_Datatype</code>	<code>sendtype,</code>	<code>= the type of data we want to gather</code>
	<code>void*</code>	<code>recvbuf,</code>	<code>= where to receive the data</code>
	<code>int</code>	<code>recvcount,</code>	<code>= number of elements to receive</code>
	<code>MPI_Datatype</code>	<code>recvtype,</code>	<code>= the type of data we want to receive</code>
	<code>int</code>	<code>root,</code>	<code>= rank of the gather root</code>
	<code>MPI_Comm</code>	<code>comm)</code>	<code>= the communicator used for this gather</code>

```
MPI_Gather(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype,  
           root, comm, ierror)  
type(*), dimension(..), intent(in) :: sendbuf  
integer, intent(in)                :: sendcount, recvcount, root  
type(MPI_Datatype), intent(in)     :: sendtype, recvtype  
type(*), dimension(..)             :: recvbuf  
type(MPI_Comm), intent(in)         :: comm  
integer, optional, intent(out)     :: ierror
```

# Gather Example

```
int gath_root = 0;

int size, rank, ierror, value;
int *buffer;

MPI_Init(&argc, &argv);
MPI_Comm_size(MPI_COMM_WORLD, &size);
MPI_Comm_rank(MPI_COMM_WORLD, &rank);

value = rank * 100;
printf("Process %d has value %d.\n", rank, value);

if (rank == gath_root) buffer = (int*)malloc(sizeof(int)*size);
MPI_Gather(&value, 1, MPI_INT, buffer, 1, MPI_INT,
          gath_root, MPI_COMM_WORLD);

if (rank == gath_root) {
    printf("Values collected on process %d:", rank);
    for (int i = 0; i < size; ++i) printf(" %d", buffer[i]);
    printf(".\n");

    free(buffer);
}

MPI_Finalize();
```

```
integer, parameter :: gath_root = 0

integer :: size, rank, ierror, value
integer, dimension(:), allocatable :: buffer

call MPI_Init(ierror)
call MPI_Comm_size(MPI_COMM_WORLD, size, ierror)
call MPI_Comm_rank(MPI_COMM_WORLD, rank, ierror)

value = rank * 100

print 100, rank, value
100 format('Process ', i0, ' has value ', i0, '.')

if (rank .eq. gath_root) allocate(buffer(size))
call MPI_Gather(value, 1, MPI_INTEGER, buffer, 1, &
& MPI_INTEGER, gath_root, MPI_COMM_WORLD, ierror)

if (rank .eq. gath_root) then
    print 200, rank, buffer
200 format('Values collected on process ', i0, &
& ' : ', *(1x,i0), '.')

    deallocate(buffer)
end if

call MPI_Finalize(ierror)
```

## Gather Example

```
$ mpicc -o gather gather.c
$ mpirun -np 4 ./gather
Process 2 has value 200.
Process 0 has value 0.
Process 3 has value 300.
Process 1 has value 100.
Values collected on process 0: 0 100 200 300.
```

## Back to the Sum of Integer

If we go back to the communication part of the sum of integers.

```
if (rank > 0) {
    MPI_Send(&proc_sum, 1, MPI_UNSIGNED, 0, 1, MPI_COMM_WORLD);
} else {
    unsigned int remote_sum;
    for(int src = 1; src < world_size; ++src) {
        MPI_Recv(&remote_sum, 1, MPI_UNSIGNED, src, 1, MPI_COMM_WORLD, &status);
        proc_sum += remote_sum;
    }
}
```

We can rewrite this part of the code with a gather

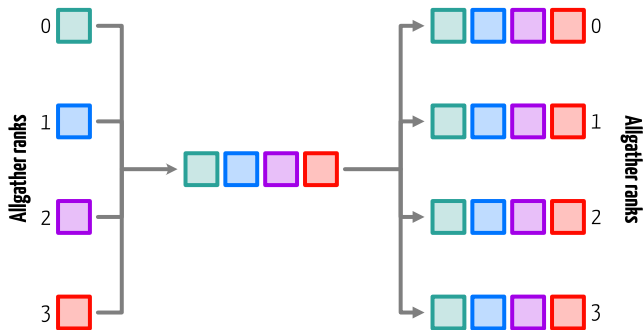
```
unsigned int* remote_sums;
if(rank == 0) remote_sums = (unsigned int*)malloc(sizeof(int)*world_size);

MPI_Gather(&proc_sum, 1, MPI_UNSIGNED, remote_sums, 1, MPI_UNSIGNED, 0, MPI_COMM_WORLD);

if(rank == 0) {
    unsigned int sum = 0;
    for(int i = 0; i < world_size; ++i)
        sum += remote_sums[i];
}
```

## MPI All Gather

You can perform an all gather operation where all the pieces of data are gathered in the order of the ranks and then, the result is broadcast to all the processes in the communicator.





## MPI All Gather

<code>MPI_Allgather</code>	<code>(void*</code>	<code>sendbuf,</code>	<code>= address of the data you want to gather</code>
	<code>int</code>	<code>sendcount,</code>	<code>= number of elements to gather</code>
	<code>MPI_Datatype</code>	<code>sendtype,</code>	<code>= the type of data we want to gather</code>
	<code>void*</code>	<code>recvbuf,</code>	<code>= where to receive the data</code>
	<code>int</code>	<code>recvcount,</code>	<code>= number of elements to receive</code>
	<code>MPI_Datatype</code>	<code>recvtype,</code>	<code>= the type of data we want to receive</code>
	<code>MPI_Comm</code>	<code>comm)</code>	<code>= the communicator used for this gather</code>

```
MPI_Allgather(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype,  
              comm, ierror)  
type(*), dimension(..), intent(in) :: sendbuf  
integer, intent(in)                :: sendcount, recvcount  
type(MPI_Datatype), intent(in)     :: sendtype, recvtype  
type(*), dimension(..)            :: recvbuf  
type(MPI_Comm), intent(in)        :: comm  
integer, optional, intent(out)    :: ierror
```

# All Gather Example

```
int gath_root = 0;
int size, rank, ierror, value;
int *buffer;

MPI_Init(&argc, &argv);
MPI_Comm_size(MPI_COMM_WORLD, &size);
MPI_Comm_rank(MPI_COMM_WORLD, &rank);

value = rank * 100;

printf("Process %d has value %d.\n", rank, value);

buffer = (int*)malloc(sizeof(int)*size);
MPI_Allgather(&value, 1, MPI_INT, buffer, 1,
             MPI_INT, MPI_COMM_WORLD);

printf("Values collected on process %d:", rank);
for (int i = 0; i < size; ++i) printf(" %d", buffer[i]);
printf(".\n");

free(buffer);

MPI_Finalize();
```

```
integer, parameter :: gath_root = 0

integer :: size, rank, ierror, value
integer, dimension(:), allocatable :: buffer

call MPI_Init(ierror)
call MPI_Comm_size(MPI_COMM_WORLD, size, ierror)
call MPI_Comm_rank(MPI_COMM_WORLD, rank, ierror)

value = rank * 100
print 100, rank, value
100 format('Process ', i0, ' has value ', i0, '.')

allocate(buffer(size))

call MPI_Allgather(value, 1, MPI_INTEGER, buffer, 1, &
                  & MPI_INTEGER, MPI_COMM_WORLD, ierror)

print 200, rank, buffer
200 format('Values collected on process ', i0, &
          & ': ', *(1x,i0), '.')

deallocate(buffer)

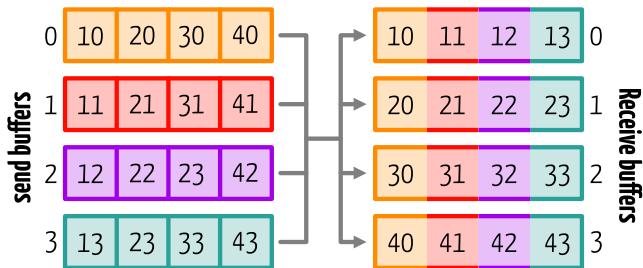
call MPI_Finalize(ierror)
```

## All Gather Example

```
$ mpicc -o allgather allgather.c
$ mpirun -np 4 ./allgather
Process 2 has value 200.
Process 0 has value 0.
Process 3 has value 300.
Process 1 has value 100.
Values collected on process 1: 0 100 200 300.
Values collected on process 3: 0 100 200 300.
Values collected on process 0: 0 100 200 300.
Values collected on process 2: 0 100 200 300.
```

## MPI All to All

All to all Scatter/Gather communication allows for data distribution to all processes: the  $j$ -th block sent from rank  $i$  is received by rank  $j$  and is placed in the  $i$ -th block of the receive buffer.



## MPI All to All

<code>MPI_Alltoall</code>	<code>(void*</code>	<code>sendbuf,</code>	<code>= address of the data you want to send</code>
	<code>int</code>	<code>sendcount,</code>	<code>= number of elements to send</code>
	<code>MPI_Datatype</code>	<code>sendtype,</code>	<code>= the type of data we want to send</code>
	<code>void*</code>	<code>recvbuf,</code>	<code>= where to receive the data</code>
	<code>int</code>	<code>recvcount,</code>	<code>= number of elements to receive</code>
	<code>MPI_Datatype</code>	<code>recvtype,</code>	<code>= the type of data we want to receive</code>
	<code>MPI_Comm</code>	<code>comm)</code>	<code>= the communicator used</code>

```
MPI_Alltoall(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype,  
             comm, ierror)  
type(*), dimension(..), intent(in) :: sendbuf  
integer, intent(in)                :: sendcount, recvcount  
type(MPI_Datatype), intent(in)     :: sendtype, recvtype  
type(*), dimension(..)            :: recvbuf  
type(MPI_Comm), intent(in)         :: comm  
integer, optional, intent(out)     :: ierror
```

## Vector Variants

Ok, but what if I do not want to transfer the same number of elements from each process?

- all the functions presented previously have a "v" variant
- these variants allow the messages received to have different lengths and be stored at arbitrary locations

On the root process, instead of specifying the number of elements to send, the vector variants take:

- an array where the entry  $i$  specifies the number of elements to send to rank  $i$
- an array where the entry  $i$  specifies the displacement from which to take the data to send to rank  $i$

## Vector Variants: Scatterv

For example, the vector variant of the `MPI_Scatter` function is `MPI_Scatterv`

<code>MPI_Scatterv</code>	<code>(void*</code>	= address of the data you want to send
	<code> sendbuf,</code>	
	<code> int sendcounts[],</code>	= the number of elements to send to each process
	<code> int displs[],</code>	= the displacement to the message sent to each process
	<code> MPI_Datatype sendtype,</code>	= the type of data we want to send
	<code> void* recvbuf,</code>	= where to receive the data
	<code> int recvcount,</code>	= number of elements to receive
	<code> MPI_Datatype recvtype,</code>	= the type of data we want to receive
	<code> int root,</code>	= rank of the root proces
	<code> MPI_Comm comm)</code>	= the communicator used

## Vector Variants: Scatterv

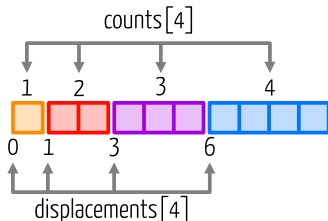
For example, the vector variant of the `MPI_Scatter` function is `MPI_Scatterv`

```
MPI_Scatterv(sendbuf, sendcounts, displs, sendtype, recvbuf, recvcount,  
             recvtype, root, comm, ierror)  
type(*), dimension(..), intent(in) :: sendbuf  
integer, intent(in)                :: sendcounts(*), displs(*)  
integer, intent(in)                :: recvcount, root  
type(MPI_Datatype), intent(in)     :: sendtype, recvtype  
type(*), dimension(..)             :: recvbuf  
type(MPI_Comm), intent(in)         :: comm  
integer, optional, intent(out)     :: ierror
```



## Scatterv Example (C)

- process with rank 0 is the root, it fills an array and dispatches the values to all the processes
- the processes receive  $\text{rank} + 1$  elements



- to rank 0
- to rank 1
- to rank 3
- to rank 4

```
int* displs = NULL;
int* nelems = NULL;
int* sendbuf = NULL;
int* recvbuf = (int*)malloc(sizeof(int)*(rank+1));

if (rank == 0) {
    int n = world_size*(world_size+1)/2;

    sendbuf = (int*)malloc(sizeof(int)*n);
    displs = (int*)malloc(sizeof(int)*world_size);
    nelems = (int*)malloc(sizeof(int)*world_size);

    for (int i = 0; i < n; ++i) sendbuf[i] = 100*(i+1);

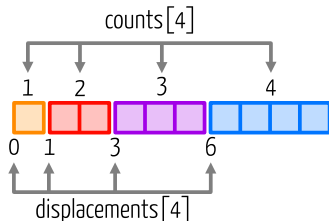
    for (int i = 0; i < world_size; ++i) {
        displs[i] = i*(i+1)/2;
        nelems[i] = i+1;
    }
}

MPI_Scatterv(sendbuf, nelems, displs, MPI_INT,
            recvbuf, rank+1, MPI_INT, ...);

printf("Process %d received values:", rank);
for(int i = 0; i < rank+1; i++) printf(" %d", recvbuf[i]);
printf("\n");
```

## Scatterv Example (Fortran)

- process with rank 0 is the root, it fills an array and dispatches the values to all the processes
- the processes receive  $\text{rank} + 1$  elements



■ to rank 0   ■ to rank 1  
■ to rank 3   ■ to rank 4

```
integer, dimension(:), allocatable :: sendbuf, recvbuf
integer, dimension(:), allocatable :: displs, nelems

allocate(recvbuf(rank+1))

if (rank .eq. 0) then
  n = world_size*(world_size+1)/2

  allocate(sendbuf(n))
  allocate(displs(world_size), nelems(world_size))

  do i = 1,n
    sendbuf(i) = 100*i
  end do

  do i = 1,world_size
    displs(i) = i*(i-1)/2
    nelems(i) = i
  end do
end if

call MPI_Scatterv(sendbuf, nelems, displs, MPI_INTEGER, &
  & recvbuf, rank+1, MPI_INTEGER, ...)

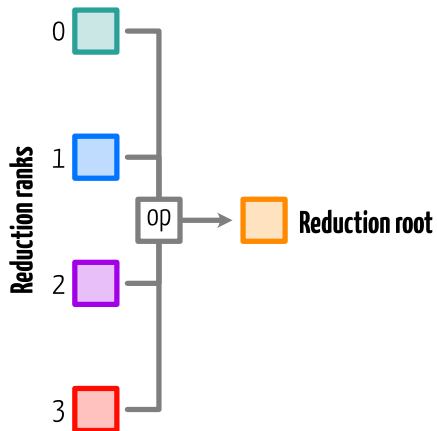
print 100, rank, recvbuf
100 format('Process ', i0, ' received values:', *(1x,i0))
```

## Scatterv Example

```
$ mpicc -o scatterv scatterv.c
$ mpirun -np 4 ./scatterv
Process 0 received values: 100
Process 1 received values: 200 300
Process 3 received values: 700 800 900 1000
Process 2 received values: 400 500 600
```

## MPI Reduce

Data reduction is reducing a set of numbers into a smaller set of numbers. For example, summing elements of an array or find the min/max value in an array.



## MPI Reduce

<code>MPI_Reduce</code>	<code>(void*</code>	<code>sendbuf,</code>	<code>= address of the data you want to reduce</code>
	<code>void*</code>	<code>recvbuf,</code>	<code>= address of where to store the result</code>
	<code>int</code>	<code>count,</code>	<code>= the number of data elements</code>
	<code>MPI_Datatype</code>	<code>datatype,</code>	<code>= the type of data we want to reduce</code>
	<code>MPI_Op</code>	<code>op,</code>	<code>= the type operation to perform</code>
	<code>int</code>	<code>root,</code>	<code>= rank of the reduction root</code>
	<code>MPI_Comm</code>	<code>comm)</code>	<code>= the communicator used for this reduction</code>

```
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
```

## MPI Reduction Operators

MPI has a number of elementary reduction operators, corresponding to the operators of the C programming language.

MPI Op	Operation	MPI Op	Operation
<code>MPI_MIN</code>	<code>min</code>	<code>MPI_LAND</code>	<code>&amp;&amp;</code>
<code>MPI_MAX</code>	<code>max</code>	<code>MPI_LOR</code>	<code>  </code>
<code>MPI_SUM</code>	<code>+</code>	<code>MPI_BAND</code>	<code>&amp;</code>
<code>MPI_PROD</code>	<code>*</code>	<code>MPI_BOR</code>	<code> </code>

In addition you can create your own custom operator type.

## Back to the Sum of Integers

If we go back to the communication part of the sum of integers.

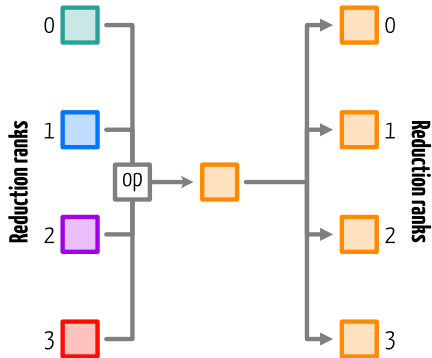
```
if (rank > 0) {
    MPI_Send(&proc_sum, 1, MPI_UNSIGNED, 0, 1, MPI_COMM_WORLD);
} else {
    unsigned int remote_sum;
    for(int src = 1; src < world_size; ++src) {
        MPI_Recv(&remote_sum, 1, MPI_UNSIGNED, src, 1, MPI_COMM_WORLD, &status);
        proc_sum += remote_sum;
    }
}
```

We can rewrite this part of the code with a reduction

```
unsigned int final_sum;
MPI_Reduce(&proc_sum, &final_sum, 1, MPI_UNSIGNED, MPI_SUM, 0, MPI_COMM_WORLD);
```

## MPI All Reduce

You can also use an all reduce operation so that the result is available to all the processes in the communicator.





## MPI All Reduce

<code>MPI_Allreduce</code>	<code>(void*</code>	<code>sendbuf,</code>	= address of the data you want to reduce
	<code>void*</code>	<code>recvbuf,</code>	= address of where to store the result
	<code>int</code>	<code>count,</code>	= the number of data elements
	<code>MPI_Datatype</code>	<code>datatype,</code>	= the type of data we want to reduce
	<code>MPI_Op</code>	<code>op,</code>	= the type operation to perform
	<code>MPI_Comm</code>	<code>comm)</code>	= the communicator used for this reduction

```
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
```

## Reduce and Allreduce Example

As an example of the `MPI_Allreduce` and `MPI_Reduce` functions, we will consider the computation of standard deviation

$$\sigma = \sqrt{\frac{\sum_i (x_i - \mu)^2}{N}}$$

- $x_i$ : value of the population
- $\mu$ : mean of the population
- $N$ : size of the population

## Reduce and Allreduce Example (C)

```
for (int i = 0; i < nelems_per_rank; ++i)
    local_sum += values[i];

MPI_Allreduce(&local_sum, &global_sum, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
mean = global_sum / (nelems_per_rank * world_size);

for (int i = 0; i < nelems_per_rank; ++i)
    local_sq_diff += (values[i] - mean) * (values[i] - mean);

MPI_Reduce(&local_sq_diff, &global_sq_diff, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);

if (rank == 0) {
    double stddev = sqrt(global_sq_diff / (nelems_per_rank * world_size));
    printf("Mean = %lf, Standard deviation = %lf\n", mean, stddev);
}
```

- for the mean as all ranks need it for the subsequent step we use **MPI\_Allreduce**
- we use **MPI\_Reduce** to sum the squared difference on rank 0 and compute the final result

## Reduce and Allreduce Example (Fortran)

```
do i = 1, nelems_per_rank
  local_sum = local_sum + values(i)
enddo

call MPI_Allreduce(local_sum, global_sum, 1, MPI_DOUBLE_PRECISION, MPI_SUM, MPI_COMM_WORLD)
mean = global_sum / (nelems_per_rank * size)

do i = 1, nelems_per_rank
  local_sq_diff = local_sq_diff + (values(i) - mean) * (values(i) - mean)
end do

call MPI_Reduce(local_sq_diff, global_sq_diff, 1, MPI_DOUBLE_PRECISION, MPI_SUM, 0, MPI_COMM_WORLD)

if (rank .eq. 0) then
  stddev = dsqrt(global_sq_diff / dble(nelems_per_rank * size))
  print 100, mean, stddev
100  format('Mean = ', f12.6, ', Standard deviation = ', f12.6, f12.6)
end if
```

- for the mean as all ranks need it for the subsequent step we use **MPI\_Allreduce**
- we use **MPI\_Reduce** to sum the squared difference on rank 0 and compute the final result

# MPI Barrier

A barrier can be used to synchronize all processes in a communicator. Each process wait until all processes reach this point before proceeding.

`MPI_Barrier`(MPI\_Comm communicator)

For example:

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

int rank;
MPI_Comm_rank(MPI_COMM_WORLD, &rank);

printf("Process %d: I start waiting at the barrier.\n",
       rank);

MPI_Barrier(MPI_COMM_WORLD);

printf("Process %d: I'm on the other side of the barrier.\n",
       rank);

MPI_Finalize();
```

```
integer :: rank, ierror

call MPI_Init(ierror)
call MPI_Comm_rank(MPI_COMM_WORLD, rank, ierror)

print 100, rank
100 format('Process ', i0, ': I start waiting at the barrier.')

call MPI_Barrier(MPI_COMM_WORLD, ierror);

print 200, rank
200 format('Process ', i0, &
&       ': I am on the other side of the barrier.')

call MPI_Finalize(ierror)
```

**Wrapping up**

# Summary

Today, we covered the following topics:

- Point-to-point communication
- Non-blocking communication
- Collective communication

But we only scratch the surface: the possibilities offered by MPI are much broader than what we have discussed.

## Going further

The possibilities offered by MPI are much broader than what we have discussed.

- User-defined datatype
- Persistent communication
- One-sided communication
- File I/O
- Topologies