# Introduction to MPI

Anne Cadiou <sup>1</sup> Fabrice Roy <sup>2</sup> (translation)

<sup>1</sup>Laboratoire de Mécanique des Fluides et d'Acoustique <sup>2</sup>Laboratoire Univers et Théories

Ecole Thématique PhyNuBE : Première rencontre de Physique Nucléaire de Basse Energie 2021 (december 9th 2021)

# Message passing programming model

#### Definition

- The program is written in a classic language (Fortran, C, C++, etc.).
- All the program variables are private and reside in the local memory of each process.
- Each process has the possibility of executing different parts of a program.
- A variable is exchanged between two or several processes via a programmed call to specific subroutines.

This slide comes from the TDRIS MPT course

# Why use this programming model?

#### Modern supercomputers

- Distributed memory computers composed of several nodes
- One or several processors in each node
- $\Rightarrow$  One or several processes of a distributed application can run on each node
- The nodes are connected with a high performance network
- $\Rightarrow$  Messages can be exchanged through this network to enable distribution of the workload, synchronization, etc.

#### Introduction to MPI

#### Definition

- Message Passing Interface
- Library and standard for communications between computing nodes

Usable on computers with shared or distributed memory Fast and portable

#### **MPI**

- Manages message passing between processes (data transfer, synchronization, global operations)
- Based on SPMD principle (Single Program Multiple Data)
- Each process has its own data
- Communications between processes are done in a communicator
- Each process is identified by its rank within the communicator

Introduction

# History

- Concept of standard discussed in 1991
- First standard presented at Supercomputing 93'
- MPI-1 release, 1994
  - library of functions usable with C, C++, Fortran
- MPI-2, 1998
  - MPI I/O,
  - dynamic processes management,
  - one-sided communications,
  - C++ interface becomes obsolete (external "C")
- MPI-3, 2012
  - collective non-blocking operations,
  - explicit functions for shared memory architectures,
  - modern Fortran interface (2003, 2008)
  - C++ support is dropped
  - interfacing with external tools (debugging, profiling)
- MPI-4?
  - accelerators support?
  - fault tolerance?
  - Discussion at EuroMPI, sept. 2019
  - Usage survey (feb. 2019), https://bosilca.github.io/MPIsurvey/

# Implementations

#### Open source

- MPICH (MPI 1.x), MPICH 2 (MPI-2)
- OpenMPI
- Boost (C++)

Some manufacturers also develop their own implementation.

Classic implementations for C, C++, Fortran. There are also implementations in Python, Julia, Java, OCaml, Perl, etc.

# Alternatives?

#### Other kind of parallelism

- OpenMP, TBB, Pthreads (shared memory)
- OpenACC, CUDA, OpenCL (accelerators)
- StarPU (task parallelism)
- kokkos, Alpaka (cross-platforms)

Other features

# General structure of a MPI program

```
// beginnng of the program
// include of the library
// initialization of the MPI environment
// calls to the library to exchange messages
// closing of the MPI environment
// end of the program
```

Usage

00000

# Include MPI

#### Fortran

include 'mpif.h'

use mpi

f95 prog.f -lmpi
mpif90 prog.f
mpirun -np 4 ./a.out

Fortran 2008 (MPI-3)

use mpi\_f08

C

#include <mpi.h>

gcc prog.c -lmpi
mpicc prog.c
mpirun -np 4 ./a.out

C++

#include <mpi.h>
extern "C" {
 #include <mpi.h>
}

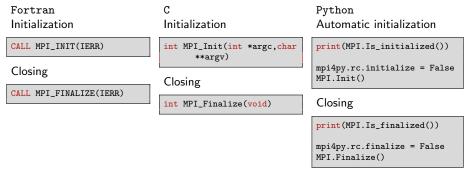
gcc prog.cpp -lmpi mpicxx prog.cpp mpirun -np 4 ./a.out

#### Python

import mpi4py.MPI as MPI

mpirun -np 4 python \
script.py

# MPI Environment



MPI creates a communicator that contains all processes. Its default name is MPI\_COMM\_WORLD (MPI.COMM\_WORLD for Python)

Introduction

Type MPI : MPI\_Comm

example:

Fortran 2008

C

 Python

comm = MPI.COMM\_WORLD

Older Fortran

- Defines a group of active processes
  - Can be created or destroyed during the execution ( $\geq$  MPI-2)
- A process has one or several identifiers (communicator, rank)
  - A process can belong to several communicators.
  - A process can have a different rank in each communicator.
- MPI communications must specify the communicator in which they take place.

A communicator can be decomposed, associated with a particular topology, etc.

Use the prefix MPI\_ mpi4py with Python is base on the obsolete C++ syntax. Class MPI.

#### Fortran

```
CALL MPI_XXX(parameter,...,
     ierr)
call mpi_xxx(parameter,...,
     ierr)
CALL MPI BSEND (buf.count.
     type, dest, tag, comm,
     ierr)
```

ierr returns the error code

MPI\_SUCCESS if ok.

#### C

```
rc = MPI Xxx(parameter....)
rc = MPI_Bsend(&buf,count,
     type.dest.tag.comm)
```

rc returns the error code. MPT SUCCESS if ok.

Warning: in C, you must respect the case. MPI is always capital, first letter of the second word in capital.

### Python based on C++ interface from MPT-2

```
comm = MPI.COMM WORLD
 # communication of
       Python objects (
       SLOW)
comm.xxx(data,parameters
     ...)
comm.bsend(buffer,dest=0,
     tag=1)
# communication of buffer-
     like objects
comm.Xxx([data,count,MPI.
     type],parameters
     ...)
comm.Bsend([buffer,
     buffsize, MPI.INT],
     dest=0.tag=1)
```

# Example of process identification

Initialize the MPI environment and returns the process rank for each active process in the default communicator.

#### Execution:

```
mpirun -np 4 ./a.out
```

```
mpirun -np 4 python script.py
```

### Output:

```
Hello! I am process 0 of 4 on plume.
Hello! I am process 1 of 4 on plume.
Hello! I am process 2 of 4 on plume.
Hello! I am process 3 of 4 on plume.
```

# Example of process identification in Fortran

```
PROGRAM hello
 USE mpi
 IMPLICIT NONE
 INTEGER :: numtasks, rank, reslen, ierr = 0
 CHARACTER(MPI_MAX_PROCESSOR_NAME) :: hostname
 CALL MPI_INIT(ierr)
 CALL MPI COMM SIZE(MPI COMM WORLD, numtasks, ierr)
 CALL MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
 CALL MPI_GET_PROCESSOR_NAME(hostname, reslen, ierr)
 WRITE(*, '(2A, I2, A, I2, 3A)') &
      'Hello! '. &
      'I am process ', rank, &
      ' of ', numtasks, &
      ' on ', hostname(1:reslen), '.'
 CALL MPI FINALIZE(ierr)
END PROGRAM hello
```

```
#include <mpi.h>
#include <stdio.h>
int main(int argc, char *argv[]) {
 int numtasks, rank, reslen, rc;
 char hostname[MPI_MAX_PROCESSOR_NAME];
 MPI_Init(&argc,&argv);
 MPI_Comm_size(MPI_COMM_WORLD,&numtasks);
 MPI_Comm_rank(MPI_COMM_WORLD,&rank);
 MPI_Get_processor_name(hostname, &reslen);
 printf("Hello! I am process %d of %d on %s.\n",rank,numtasks,hostname);
 MPI Finalize():
```

Introduction

00

```
#!/usr/bin/env python
"""
for python3
"""
import mpi4py.MPI as MPI

rank = MPI.COMM_WORLD.Get_rank()
numtasks = MPI.COMM_WORLD.Get_size()
hostname = MPI.Get_processor_name()

mess = "Hello! I am process %d of %d on %s."
print(mess % (rank, numtasks, hostname))
```

#### Communications exchange messages between at least 2 processes

- one sends
- the other receives

#### Different kinds of communications

- Point to point communications
- Global (or collective) communications
- One-sided communications, where data movement is decoupled from process synchronization (\geq MPI-2)

Communications can be blocking or non-blocking The messages exchanged are typed data.

Introduction

00

| MPI data type        | Fortran data type |
|----------------------|-------------------|
| MPI_INTEGER          | INTEGER           |
| MPI_REAL             | REAL              |
| MPI_DOUBLE_PRECISION | DOUBLE PRECISION  |
| MPI_COMPLEX          | COMPLEX           |
| MPI_LOGICAL          | LOGICAL           |
| MPI_CHARACTER        | CHARACTER         |
| MPI_BYTE             |                   |
| MPI_PACKED           |                   |

# Main data types in C

| MPI data type          | C data type            |
|------------------------|------------------------|
| MPI_CHAR               | char                   |
| MPI_SHORT              | short int              |
| MPI_INT                | int                    |
| MPI_LONG               | long int               |
| MPI_LONG_LONG          | long long int          |
| MPI_UNSIGNED_CHAR      | unsigned char          |
| MPI_UNSIGNED_SHORT     | unsigned short int     |
| MPI_UNSIGNED           | unsigned int           |
| MPI_UNSIGNED_LONG      | unsigned long int      |
| MPI_UNSIGNED_LONG_LONG | unsigned long long int |
| MPI_FLOAT              | float                  |
| MPI_DOUBLE             | double                 |
| MPI_LONG_DOUBLE        | long double            |
| MPI_BYTE               |                        |
| MPI_PACKED             |                        |

References

0

# Main data types in Python

Introduction

00

| MPI data type | NumPy data type |
|---------------|-----------------|
| MPI.INTEGER   | np.intc         |
| MPI.LONG      | np.int          |
| MPI.FLOAT     | np.float32      |
| MPI.DOUBLE    | np.float64      |

Other features

# Point to point communications

Communication between 2 processes : the sender and the receiver It includes 2 stages : sending and receiving

#### A message consists of its header, which contains

- the communicator
- the rank of the sending process
- the rank of the receiving process
- a (tag) which allows the program to distinguish between different messages

#### and its content, which contains

- the exchanged data
- its type
- its size

#### Fortran

Introduction

00

```
CALL MPI_SEND(message, size, type, destination, tag, MPI_COMM_WORLD, ierr)
```

```
CALL MPI_RECV(message, size, type, source, tag, MPI_COMM_WORLD, status, ierr)
```

MPI\_Send(void\* data, int count, datatype, int dest, int tag, communicator)

(

```
***
```

```
MPI_Recv(void* data, int count, datatype, int source, int tag, communicator, status)
```

#### Python

```
MPI.COMM_WORLD.send(data, dest=#, tag=#)
```

```
MPI.COMM_WORLD.Send(data, dest=#, tag=#)
```

```
data = MPI.COMM_WORLD.recv(source=#, tag=#)
```

```
MPI.COMM WORLD.Recv(data, source=#, tag=#)
```

(send/recv : generic Python objects, Send/Recv : NumPy arrays, faster)

# Example : One process reads data

Process 0 (always exists) reads an integer and sends it to the other processes, which send back this integer multiplied by their rank in the communicator. (master/slave setup)

```
$ mpirun -np 4 ./a.out
Enter integer number:
12
Slave 1 has received n=12 from 0
Slave 2 has received n=12 from 0
Slave 3 has received n=12 from 0
Master 0 received from slave 1: 12
Master 0 received from slave 2: 24
Master 0 received from slave 3: 36
```

```
PROGRAM point_a_point
 2
    USE mpi
    IMPLICIT NONE
 4
    INTEGER, DIMENSION(MPI STATUS SIZE) :: statut
6
    INTEGER, PARAMETER :: tagm=101, tagr=201, master=0
    INTEGER :: rang.nprocs.i.n.ierr
8
10
    CALL MPI_INIT(ierr)
    CALL MPI_COMM_RANK(MPI_COMM_WORLD, rang, ierr)
12
    CALL MPI_COMM_SIZE(MPI_COMM_WORLD, nprocs, ierr)
14
    IF (rang == master) THEN
      WRITE(*,*) 'Enter integer number:'
16
      READ(*,*) n
18
      DO i=1,nprocs-1
        CALL MPI SEND(n.1.MPI INTEGER.i.tagm.MPI COMM WORLD.ierr)
      END DO
20
      DO i=1,nprocs-1
        CALL MPI RECV(n.1.MPI INTEGER.i.tagr.MPI COMM WORLD.statut.ierr)
22
        WRITE(*,'(A,12,A,12,A,13)') "Master", rang, " received from slave ",i,": n=",n
        WRITE(*,*) n,statut(MPI_SOURCE),statut(MPI_TAG),statut(MPI_ERROR)
24
      END DO
26
    ELSE
```

References

0

Introduction

00

```
CALL MPI_RECV(n,1,MPI_INTEGER,master,tagm,MPI_COMM_WORLD,statut,ierr)

WRITE(*,*)(A,12,A,13,A,12)') "Slave ",rang," has received n=",n," from ",master WRITE(*,*) n,statut(MPI_SOURCE),statut(MPI_TAG),statut(MPI_ERROR)

n = n*rang
    CALL MPI_SEND(n,1,MPI_INTEGER,master,tagr,MPI_COMM_WORLD,ierr)

END IF

CALL MPI_FINALIZE(ierr)

END PROGRAM point_a_point
```

Other features

# Example in C

```
#include <mpi.h>
2 #include <stdio.h>
4 int main(int argc, char *argv[]) {
    int rang,nprocs;
6
    int master,tagm,tagr;
    int n,i;
8
    MPI_Status status;
10
    MPI_Init(&argc,&argv);
12
    master = 0:
    tagm = 101;
14
    tagr = 201;
16
    MPI_Comm_size(MPI_COMM_WORLD,&nprocs);
18
    MPI_Comm_rank(MPI_COMM_WORLD,&rang);
    if (rang == master) {
20
```

00

```
printf("Enter integer number:\n");
      scanf("%d",&n);
22
      for (i=1; i<nprocs; i++) {</pre>
        MPI_Send(&n,1,MPI_INT,i,tagm,MPI_COMM_WORLD);
24
      for (i=1; i<nprocs; i++) {</pre>
26
        MPI_Recv(&n,1,MPI_INT,i,tagr,MPI_COMM_WORLD,&status);
        printf("Master %d received from slave %d: %d\n".rang. i. n):
28
    } else {
30
      MPI_Recv(&n,1,MPI_INT,master,tagm,MPI_COMM_WORLD,&status);
      printf("Slave %d has received n=%d from %d\n",rang,n,master);
32
      n = n*rang:
      MPI_Send(&n,1,MPI_INT,master,tagr,MPI_COMM_WORLD);
34
36
    MPI_Finalize();
38 | }
```

```
#!/usr/bin/env python
  """ for python3 """
  import mpi4pv.MPI as MPI
  rang = MPI.COMM WORLD.Get rank()
6 nprocs = MPI.COMM_WORLD.Get_size()
8 \mid master = 0
  tagm = 101
10 tagr = 201
12 if rang == master:
      print('Enter integer number:')
      n = int(input())
14
      for i in range(1,nprocs):
          MPI.COMM_WORLD.send(n,dest=i,tag=tagm)
16
      for i in range(1,nprocs):
          n = MPI.COMM_WORLD.recv(source=i,tag=tagr)
18
          print("Master ",rang," received from slave ",i,": n=",n)
20 else:
      n = MPI.COMM_WORLD.recv(source=0,tag=tagm)
      print("Slave ",rang," has received n=",n," from ",master)
22
      n = n*rang
      MPI.COMM WORLD.send(n.dest=0.tag=tagr)
24
```

# Example in Python with NumPy structures

```
#!/usr/bin/env pvthon
   for python3
6 import mpi4py.MPI as MPI
  import numpy as np
8
  rang = MPI.COMM WORLD.Get rank()
10 nprocs = MPI.COMM_WORLD.Get_size()
12 \mid master = 0
  tagm = 101
14 tagr = 201
16 if rang == master:
      print('Enter integer number:')
      n = input()
18
      data = np.array([n],dtype='i')
      for i in range(1,nprocs):
20
          MPI.COMM_WORLD.Send([data,MPI.INT],dest=i,tag=tagm)
      for i in range(1,nprocs):
22
          MPI.COMM_WORLD.Recv([data,MPI.INT],source=i,tag=tagr)
          print("Master ",rang," received from slave ",i,": n=",data[0])
24
  else:
```

```
Introduction
00
```

```
data = np.empty(1,dtype='i')
26
      MPI.COMM_WORLD.Recv([data,MPI.INT],source=0,tag=tagm)
      print("Slave ",rang," has received n=",data[0]," from ",master)
28
      data[0] = data[0] *rang
      MPI.COMM_WORLD.Send([data, MPI.INT], dest=0, tag=tagr)
30
```

```
(mpi4py3) $ mpirun -np 4 python point_a_point_numpy.py
Enter integer number:
12
Slave 1 has received n= 12 from 0
Slave 3 has received n= 12 from 0
Master 0 received from slave 1 : n= 12
Slave 2 has received n= 12 from 0
Master O received from slave 2: n= 24
```

Master O received from slave 3: n= 36

# Blocking communications

#### MPI\_SEND and MPI\_RECV are blocking communications

- Execution stops until sending and receiving are completed
- Warning: headers and datas in send and receive functions must match (source, tag, etc.)

# Advantage

Data is completely sent or received before it can be accessed or modified by the program

### Disadvantage

Computations stop until the end of the communication

# Optimize point to point communications

Communications are blocking because of

- data are copied in temporary memory space (buffer)
- synchronization (the application waits for a matching receive begins before it continues the execution after a send)

Optimize consists in minimizing the time spent doing something other than **computations** (overhead)

#### Options:

- Overlap communications with computations
- Avoid using buffers
- Minimize overheads related to multiple calls to communications functions

Comunications can be standard, synchronous, buffered ou persistent.

# Asynchronous communications

The execution continues before the communication has completed : it is possible to compute during the communication.

#### Fortran

Introduction

```
CALL MPI_ISEND(message, size, type, destination, tag, comm, request, ierr)
```

```
CALL MPI_IRECV(message, size, type, source, tag, comm, request, ierr)
```

c

```
MPI_Isend(void* data, int count, datatype, int dest, int tag, comm, request)
```

```
MPI_Irecv(void* data, int count, datatype, int source, int tag, comm, status)
```

#### Python

```
MPI.COMM_WORLD.isend(data, dest=#, tag=#)
```

```
MPI.COMM_WORLD.Isend(data, dest=#, tag=#)
```

```
data = MPI.COMM_WORLD.irecv(source=#, tag=#)
```

```
MPI.COMM_WORLD.Irecv(data, source=#, tag=#)
```

# Fortran

Introduction

Wait until a request has completed

CALL MPI\_WAIT(request, status, ierr)

Check whether a request has completed

CALL MPI\_TEST(request,flag,status,ierr)

Check if a message has arrived

CALL MPI\_PROBE(source, tag, status, comm, ierr)

There are asynchronous versions of these functions.

00

Wait until a request has completed

MPI\_Wait(request, status)

Check whether a request has completed

MPI\_Test(request,flag,status)

Check if a message has arrived

MPI\_Probe(source,tag,comm,flag,status)

Python

00

Wait until a request has completed

req.wait()

MPI.Request.Wait(req)

Check whether a request has completed

MPI\_Test(request,flag,status)

Check if a message has arrived

MPI\_Probe(source,tag,comm,flag,status)

Fortran Send

MPI\_SEND

MPI\_ISEND

MPI\_SSEND

MPI\_ISSEND MPI\_BSEND

MPI\_IBSEND

Receive

MPT RECV

MPT TRECV

Check

MPI\_WAIT

С

Send

MPI\_Send

MPI\_Isend MPI\_Ssend

MPT\_Tssend

MPI\_Bsend

MPI\_Ibsend

Receive

MPT Recv

MPT Trecy

Check

MPT Wait

Python (Numpy)

Send

MPT.Send

MPT Tsend

MPT Ssend MPT.Issend

MPI.Bsend

MPT Thsend

Receive

 MPT.Recv MPT Trecv

Check

MPT.Wait

type

Send

blocking, standard

non blocking, standard

References

blocking, synchronous non blocking,

synchronous

blocking, buffered non blocking, buffered

Receive

blocking, standard

non blocking, standard

Check

wait until the communication has completed

# Some simple rules

- Initialize receptions before sends
- ⇒ Write calls to MPI\_IRECV before MPI\_SEND
- Avoid use of buffers
- ⇒ Use synchronous functions MPI\_SSEND
- Overlap communications with computations
- ⇒ Use non blocking communications MPI\_ISEND and MPI\_IRECV

- Communication that involves every processes in the communicator (point-to-point communications sequence).
- Processes call the same function with corresponding arguments.
- There is no tag.
- Some collective communication have only one sending process or one receiving process, usually called root process.

Introduction

## Send and receive

- Synchronization
- Broadcast of a data from the root process to every other processes in the communicator
- Scatter of data from the root process on each process
- Gather of data on the root process

## Operation on the data

Reduction (max, min, sum, product, etc.)

# Synchronization

#### Fortran

CALL MPI\_BARRIER(MPI\_COMM\_WORLD,ierr)

C

MPI\_Barrier(MPI\_COMM\_WORLD)

Python

comm = MPI.COMM\_WORLD

comm.barrier()

comm.Barrier()

#### **Broadcast**

Send data from one process to all others (broadcast).

#### Fortran

CALL MPI\_BCAST(message, size, type, source, MPI\_COMM\_WORLD, ierr)

C

MPI\_Bcast(message,size,type,source,MPI\_COMM\_WORLD)

#### Python

comm = MPI.COMM\_WORLD

data = comm.bcast(data,root=source)

comm.Bcast(data,root=source)

```
PROGRAM bcast
    USE mpi
2
    IMPLICIT NONE
    INTEGER :: n,nprocs,rang
    INTEGER :: ierr = 0
    INTEGER, PARAMETER :: idat=91
6
    CALL MPI_INIT(ierr)
8
    CALL MPI_COMM_SIZE(MPI_COMM_WORLD, nprocs, ierr)
    CALL MPI COMM RANK(MPI COMM WORLD, rang, ierr)
10
12
    IF (rang.eq.0) THEN
        OPEN(unit=idat.file='data.txt')
        READ(idat,*) n
14
        CLOSE(idat)
    END IF
16
    CALL MPI_BCAST(n,1,MPI_INT,0,MPI_COMM_WORLD,ierr)
18
20
    print *,'Process',rang,' has received ',n
    CALL MPI_FINALIZE(ierr)
22
24 END PROGRAM bcast
```

File:

00

cat data.txt

29

Execution:

mpirun -np 4 ./a.out

## Output:

Process 0 has received 29

Process 1 has received 29

Process 2 has received 29

Process 3 has received 29

Introduction

```
#include "mpi.h"
2 #include <stdio.h>
  #include <stdlib.h>
  int main(int argc, char *argv[]) {
6
    int nprocs,rang;
8
    int n;
    MPI_Init(&argc,&argv);
10
    MPI_Comm_size(MPI_COMM_WORLD, &nprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &rang);
12
    if (rang==0) {
14
      FILE *fp = fopen("data.txt", "r");
      fscanf(fp, "%d", &n);
16
      fclose(fp);
18
    MPI_Bcast(&n,1,MPI_INTEGER,0,MPI_COMM_WORLD);
20
    printf("Process %d has received %d \n",rang,n);
22
    MPI_Finalize();
24
    return 0;
26 }
```

# Example in Python

```
#!/usr/bin/env python
  for python3
6 import mpi4py.MPI as MPI
8 comm = MPI.COMM_WORLD
  rang = comm.Get_rank()
10 nprocs = comm.Get_size()
12 if rang == 0:
      f = open('data.txt','r')
      n = int(f.read())
14
  else:
      n = None
16
18 n = comm.bcast(n,root=0)
20 print("Process ",rang,"has received ",n)
```

```
#!/usr/bin/env python
  for python3
6 import mpi4py.MPI as MPI
  import numpy as np
8
  comm = MPI.COMM WORLD
10 rang = comm.Get_rank()
  nprocs = comm.Get_size()
12
  if rang == 0:
      f = open('data.txt','r')
      n = np.array([int(f.read())], dtype=int)
16 else:
      n = np.zeros(1, dtype=int)
18
  comm.Bcast(n,root=0)
20
  print("Process ",rang,"has received ",n[0])
```

## Selective distribution of data

Distribute data from one process to all processes.

#### Fortran

Introduction

```
CALL MPI_SCATTER(sendbuf, sendcount, sendtype, & recvbuf, recvcount, recvtype, & root, comm, ierr)
```

C

#### Python

```
comm = MPI.COMM_WORLD
```

```
data = comm.scatter(sendbuff, root)
```

```
comm.Scatter([senddata, data_size, data_type], [recvdata, data_size, data_type], root)
```

# Example: distribution of a 2D array on 4 processes

Example inspired by https://computing.llnl.gov/tutorials/mpi.

Goal : distribute  $4 \times 5$  array on 4 processes.

## Execution:

```
mpirun -np 4 ./a.out
```

```
mpirun -np 4 python script.py
```

#### Output:

```
rank= 0 Results: 1.00000000 2.00000000 3.00000000 4.00000000 4.50000000 rank= 1 Results: 5.00000000 6.00000000 7.00000000 8.00000000 8.50000000 rank= 2 Results: 9.00000000 10.0000000 11.0000000 12.0000000 12.5000000 rank= 3 Results: 13.0000000 14.0000000 15.0000000 16.0000000 16.5000000
```

Introduction

```
program scatter
     include 'mpif.h'
 2
     integer SIZE_X,SIZE_Y
4
     parameter(SIZE_X=5,SIZE_Y=4)
     integer numtasks, rank, sendcount, recycount, source, ierr
6
8
     real*4 sendbuf(SIZE X.SIZE Y), recvbuf(SIZE X)
     data sendbuf /1.0, 2.0, 3.0, 4.0, 4.5, &
                  5.0, 6.0, 7.0, 8.0, 8.5, &
10
                  9.0, 10.0, 11.0, 12.0, 12.5, &
                  13.0, 14.0, 15.0, 16.0, 16.5/
12
     call MPI INIT(ierr)
14
     call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
     call MPI COMM SIZE(MPI COMM WORLD, numtasks, ierr)
16
     if (numtasks .eq. SIZE_Y) then
18
        source = 1
20
        sendcount = SIZE_X
        recycount = SIZE X
22
```

Introduction

```
#include "mpi.h"
2 #include <stdio.h>
  #define SIZE_X 5
4 #define SIZE Y 4
6 int main(int argc, char *argv[]) {
    int numtasks, rank, sendcount, recvcount, source;
8
    float sendbuf[SIZE Y][SIZE X] = {
      \{1.0, 2.0, 3.0, 4.0, 4.5\},\
10
     {5.0, 6.0, 7.0, 8.0, 8.5},
12
      {9.0, 10.0, 11.0, 12.0, 12.5},
      {13.0, 14.0, 15.0, 16.0, 16.5} };
    float recybuf[SIZE X]:
14
    MPI_Init(&argc,&argv);
16
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI Comm size(MPI COMM WORLD, &numtasks):
18
20
    if (numtasks == SIZE Y) {
      // define source task and elements to send/receive, then perform collective
            scatter
      source = 1:
22
      sendcount = SIZE_X;
      recycount = SIZE X:
24
```

```
MPI_Scatter(sendbuf, sendcount, MPI_FLOAT, recvbuf, recvcount,
                  MPI_FLOAT, source, MPI_COMM_WORLD);
26
28
      printf("rank= %d Results: %f %f %f %f %f \n",rank,recvbuf[0],
             recvbuf[1],recvbuf[2],recvbuf[3],recvbuf[4]);
30
    else
      printf("Must specify %d processors. Terminating.\n",SIZE_Y);
32
    MPI_Finalize();
34
    return 0;
36 }
```

```
#!/usr/bin/env python
   for python3
6 import mpi4py.MPI as MPI
  comm = MPI.COMM_WORLD
  rang = comm.Get_rank()
10 nprocs = comm.Get_size()
12 if rang == 0:
     data = [[1.0, 2.0, 3.0, 4.0, 4.5], \
             [5.0, 6.0, 7.0, 8.0, 9.5], \
14
             [9.0, 10.0, 11.0, 12.0, 12.5], 
             [13.0, 14.0, 15.0, 16.0, 16.5]]
16
  else:
     data = None
18
20 data = comm.scatter(data, root=0)
  print('rank=',rang,'Results:',data)
```

```
#!/usr/bin/env python
  for python3
6 import mpi4py.MPI as MPI
  import numpy as np
8
  comm = MPI.COMM WORLD
10 rang = comm.Get_rank()
  nprocs = comm.Get size()
12
  mv_N = 5
14 N = mv N * nprocs
16 if rang == 0:
     data = np.array([[1.0, 2.0, 3.0, 4.0, 4.5], \
18
             [5.0, 6.0, 7.0, 8.0, 9.5], \
             [9.0, 10.0, 11.0, 12.0, 12.5], \
             [13.0, 14.0, 15.0, 16.0, 16.5]], dtype=np.float64)
20
  else:
     data = np.empty(N, dtype=np.float64)
24 recv_data = np.empty(my_N, dtype=np.float64)
  comm.Scatter([data,my_N,MPI.DOUBLE], [recv_data,my_N,MPI.DOUBLE], root=0)
26 print('rank=',rang,'Results:',recv_data)
```

# Aggregation of the data

Aggregate data from all the processes to one process (root), possibly with broadcast of the result

#### Fortran

```
CALL MPI_GATHER(sendbuf, sendcount, sendtype, & recvbuf, recvcount, recvtype, & root, comm, ierr)
```

С

```
MPI_Gather(sendbuf, sendcount, sendtype,
    recvbuf, recvcount, recvtype,
    root, comm)
```

#### Python

```
comm = MPI.COMM_WORLD
```

```
data = comm.gather(sendbuff, root)
```

comm.Gather([senddata, data\_size, data\_type], [recvdata, data\_size, data\_type], root)

#### Reduction

Some operations are carried out on the transferred data, possibly with broadcast of the result.

#### Fortran

```
CALL MPI REDUCE (SENDBUF, RECVBUF, COUNT, DATATYPE, OP, ROOT, COMM, ierr)
CALL MPI_ALLREDUCE(SENDBUF, RECVBUF, COUNT, DATATYPE, OP, COMM, ierr)
```

C

```
MPI_Reduce(sendbuf,recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root,
     MPI Comm comm)
MPI_Reduce(sendbuf, recybuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

#### Python

```
comm = MPI.COMM WORLD
```

```
comm.reduce(sendobj=None, recvobj=None, op=MPI.SUM, root=0)
comm.allreduce(sendobi=None, recvobi=None, op=MPI.SUM)
```

```
comm.Reduce(sendbuf, recvbuf, op=MPI.SUM, root=0)
comm.Allreduce(sendbuf, recvbuf, op=MPI.SUM)
```

MPI.MIN min
MPI.MAX max
MPI.SUM sum
MPI.PROD product

MPI.MAXLOC index of the max value MPI.MINLOC index of the min value

(Python syntax)

## Example of the use of the reduction

Sum the elements of an array of size N = 1001. The array is distributed in pchunks, p is the number of processes used.

Execution for p = 4:

```
mpirun -np 4 ./a.out
```

```
mpirun -np 4 python script.py
```

#### Output:

```
0 ] part : 251.000000
 2 1 part : 250.000000
 3 ] part : 250.000000
[ 1 ] part : 250.000000
Sum : 1001.000000
```

```
PROGRAM main
 USE mpi
 IMPLICIT NONE
 INTEGER, parameter :: N=1001
 INTEGER :: ierr, i, rank, nprocs
 INTEGER :: nstart, nstop, npart, ncount, nrem
 DOUBLE PRECISION, allocatable :: vec(:)
 DOUBLE PRECISION :: local sum, total sum
 CALL MPI_INIT(ierr)
 CALL MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
 CALL MPI_Comm_size(MPI_COMM_WORLD, nprocs, ierr)
 ncount = N/nprocs
 nrem = MOD(N,nprocs)
 IF (rank < nrem) THEN
   nstart = rank * (ncount + 1)
   nstop = nstart + ncount
 ELSE
   nstart = rank * ncount + nrem
   nstop = nstart + (ncount - 1)
 END IF
 npart = nstop-nstart+1
 allocate(vec(npart))
```

```
DO i=1,npart
   vec(i) = 1.0D0
 END DO
 local sum = 0.0D0
 DO i=1,npart
  local sum = local sum + vec(i)
 END DO
 WRITE(*,'(A,I3,A,F15.8)') "[",rank,"] part : ",local_sum
 CALL MPI BARRIER (MPI COMM WORLD.ierr)
 CALL MPI_REDUCE(local_sum, total_sum, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD
       . ierr)
 IF (rank.eq.0) THEN
   WRITE(*,*) "Sum : ", total_sum
 END IF
 deallocate(vec)
 CALL MPI FINALIZE(ierr)
END PROGRAM
```

# Example in C

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>
int main(int argc, char *argv[]) {
 int i, N;
 int nprocs, rank;
 MPI_Init(&argc,&argv);
 MPI_Comm_size(MPI_COMM_WORLD,&nprocs);
 MPI_Comm_rank(MPI_COMM_WORLD,&rank);
 double *vec;
 double local_sum,sum;
 N=1001;
 int count = N/nprocs;
 int remainder = N%nprocs;
 int start, stop;
```

```
if (rank < remainder) {
  start = rank * (count + 1):
  stop = start + count;
} else {
  start = rank * count + remainder;
  stop = start + (count - 1);
int npart = stop-start+1;
vec = malloc(sizeof(double)*npart);
for (i=0; i<npart; i++) {
  vec[i] = 1.0:
local_sum = 0.0;
for (i=0; i<npart; i++) {
 local_sum += vec[i];
printf("[ %d ] part : %g\n",rank,local_sum);
MPI_Barrier(MPI_COMM_WORLD);
MPI_Reduce(&local_sum,&sum,1,MPI_DOUBLE,MPI_SUM,0,MPI_COMM_WORLD);
if (rank == 0) {
 printf("Sum : %g\n",sum);
MPI Finalize():
return 0;
```

# Example in Python

```
#!/usr/bin/env python
""" for python3 """
import mpi4py.MPI as MPI
import numpy as np
import part as part
comm = MPI.COMM_WORLD
nprocs = comm.size
rank = comm.rank
# Define the size of the problem
N = 1001
start, end = part.partition(rank, nprocs, N)
vec = np.ones((end-start+1),dtype=np.float64)
# Calculate the local sum of local vectors
local sum = vec.sum()
print("[ %d ] part : %f"%(rank, local_sum))
comm.barrier()
# Get the global sum
global_sum = comm.reduce(local_sum, op=MPI.SUM, root=0)
if rank == 0:
   print("Sum : %f"%(global sum))
```

```
#!/usr/bin/env python
""" for python3 """
import mpi4py.MPI as MPI
import numpy as np
import part as part
comm = MPI.COMM WORLD
nprocs = comm.size
rank = comm.rank
# Define the size of the problem
N = 1001
start, end = part.partition(rank, nprocs, N)
vec = np.ones((end-start+1),dtype=np.float64)
# Calculate the local sum of local vectors
local_sum = vec.sum()
print("[ %d ] part : %f"%(rank, local_sum))
comm.Barrier()
# Get the global sum
global_sum = np.zeros(1, dtype='float64')
comm.Reduce(local_sum, global_sum, op=MPI.SUM, root=0)
if rank == 0:
   print("Sum : %f"%(global_sum[0]))
```

References

# Other features

The MPI library allows other kinds of features, for instance :

- Definition and use of derived data types
- Creation of communicators
- Use of topologies for communicators
- Parallel I/O

#### References

- Parallélisme sur machines à mémoire distribuée, Bastien Di Pierro, LyonCalcul, 2016 (in French)
- Calcul parallèle avec MPI, Guy Moebs, Univ. Nantes, 2010 (in French)
- Formations MPI, IDRIS (in French and English)
- https://computing.llnl.gov/tutorials/mpi
- https://www.mpi-forum.org/docs
- OpenMPI documentation
- mpi4py documentation