



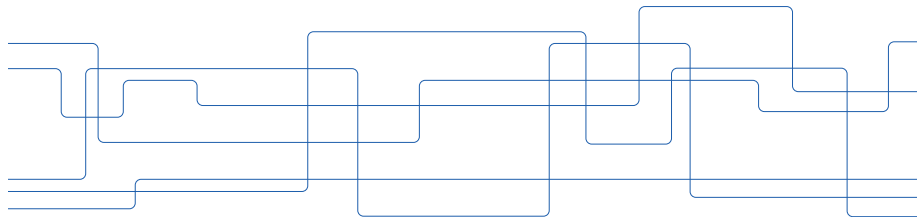
# MPI: Part I

## AQTIVATE Training Workshop I

**Dirk Pleiter**

CST | EECS | KTH

December 2023





# Overview

Introduction

First MPI Program

Blocking Send-Receive



# Content

Introduction

First MPI Program

Blocking Send-Receive

# MPI Introduction

- ▶ **MPI** = Message-Passing library Interface specification
  - ▶ Model: Data is moved from the address space of one process to that of another process through cooperative operations on each process
- ▶ Standard defined by MPI Forum
  - ▶ Members of forum: academic and commercial institutions
  - ▶ History:

May 1994:	Release of first version MPI-1.0
September 2009:	Release of version MPI-2.2
September 2012:	Release of version MPI-3.0
June 2015:	Release of version MPI-3.1
November 2020:	Publication of MPI-4.0 Release Candidate
June 2021:	Release of version MPI-4.0
November 2023:	Release of version MPI-4.1

# MPI Implementations

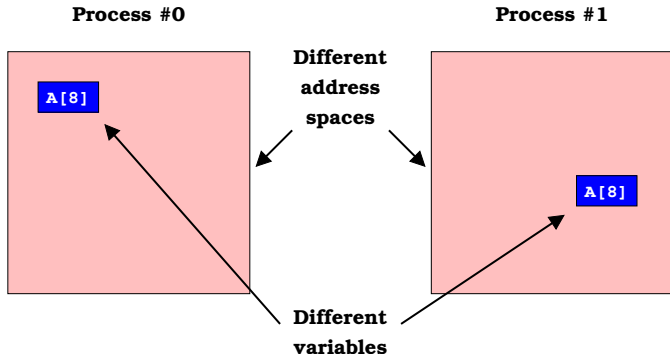
- ▶ MPI is only a standard, not an implementation
- ▶ Most popular implementations
  - ▶ OpenMPI (U Tennessee, LANL, Indiana U plus collaborators)
  - ▶ MPICH (ANL plus collaborators)
  - ▶ MVAPICH (Ohio State U plus collaborators)
- ▶ Different implementations have different advantages/disadvantages. Performance may differ on different platforms and depend on application.

# Message Passing Model

- ▶ MPI-based parallel programs consist of a set of separate processes
  - ▶ A process is an active instance of a program with its own program counter and address space
  - ▶ The processes can run on the same node or on different nodes
- ▶ Typically all processes execute the same program, but operate on different data
  - ▶ SPMD = Single Program Multiple Data
- ▶ Data has to be communicated explicitly between processes
  - ▶ Different communication patterns are supported:
    - ▶ Point-to-point communication operations
    - ▶ Collective communication operations
- ▶ Programmer has to take care of data distribution, communication and synchronisation

## Message Passing Model (2)

Beware: Variables of the same program with the same name in different processes are different variables stored at different memory addresses



# Why Using MPI? (1/2)

- ▶ **Observation:** MPI has become a standard parallel programming model and is supported on essentially all HPC systems
  - ▶ MPI replaced various other message passing solutions
- ▶ **Portability:** An application developed in a standard compliant manner on one system will run also on other systems that support an MPI implementation
  - ▶ In practice, some system specific optimisations may be required
- ▶ **Functionality:** Increasingly rich set of features
- ▶ **Wide uptake:** MPI is used for the vast majority of parallel applications leading to
  - ▶ Many training opportunities and widely available expertise
  - ▶ Broad eco-system



## Why Using MPI? (2/2)

- ▶ **Availability:** MPI runs on all Linux and most other systems
  - ▶ Open source implementations are available and distributed with popular Linux distributions
- ▶ **Performance opportunities:** MPI is co-developed with HPC solution providers that perform significant efforts optimising the software
  - ▶ Network technology providers provide optimised back-ends
  - ▶ HPC system providers provide their own implementations of MPI, e.g. Cray MPI, IBM Spectrum MPI

# References (and credits)

- ▶ MPI Standard documents:  
<http://mpi-forum.org/docs/docs.html>
- ▶ Current MPI Standard (MPI-4.1): <https://www.mpi-forum.org/docs/mpi-4.1/mpi41-report.pdf>
- ▶ W. Gropp, E. Lusk, A. Skjellum, "Using MPI: Portable Parallel Programming with the Message-passing Interface", Volume 1, MIT Press, 1999
- ▶ W. Gropp, E. Lusk, R. Thakur, "Using MPI-2: Advanced Features of the Message-passing Interface", Globe Pequot Press, 1999
- ▶ W. Gropp, T. Hoefler, R. Thakur, E. Lusk "Using Advanced MPI: Modern Features of the Message-Passing Interface", MIT Press, 2014
- ▶ M. Hava, "Parallel Programming with MPI", PRACE Autumn School 2016



# Content

Introduction

First MPI Program

Blocking Send-Receive

# First MPI Program: Hello World in C

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main (int argc, char *argv[])
5 {
6     int irank, isize;
7
8     MPI_Init(&argc, &argv);
9
10    MPI_Comm_rank(MPI_COMM_WORLD, &irank);
11    MPI_Comm_size(MPI_COMM_WORLD, &isize);
12
13    printf("Processor %d of %d: Hello World!\n",
14           irank, isize);
15
16    MPI_Finalize();
17
18    return 0;
19 }
```

# First MPI Program: Hello World in F90

```
1 PROGRAM MPI
2
3 IMPLICIT NONE
4
5 INCLUDE 'mpif.h'
6
7 INTEGER irank , isize , ierr
8
9 CALL MPI_Init(ierr)
10 CALL MPI_Comm_rank(MPI_COMM_WORLD, irank , ierr)
11 CALL MPI_Comm_size(MPI_COMM_WORLD, isize , ierr)
12
13 WRITE(*,*) "Processor " , irank , " of " , &
14           isize ,": Hello World!"
15
16 CALL MPI_Finalize(ierr)
17
18 END PROGRAM
```

# First MPI Program: Comments

- ▶ Need to include MPI header files or declare use of MPI module to make MPI functions, variable types and constants declared or defined
  - ▶ C or C++: `#include <mpi.h>`
  - ▶ Fortran 2008 or later: `USE mpi_f08`
  - ▶ Earlier Fortran: `USE mpi` or `INCLUDE 'mpif.h'`
- ▶ `MPI_Init` must be the first function that is called (few exceptions)
- ▶ `MPI_Finalize` is required to clean-up all MPI states
- ▶ C-language binding
  - ▶ MPI used through function calls
  - ▶ All functions return an error code of type `int`
- ▶ Fortran binding
  - ▶ MPI used mainly through subroutine but also function calls
  - ▶ Last argument is the error code

# MPI Initialisation (1/2)

- ▶ Single-threaded case

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

```
MPI_Init(ierror)
```

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

- ▶ Instead of the pointers argc and argv one may also provide NULL

```
int MPI_Init(NULL, NULL)
```

# MPI Initialisation (2/2)

## ► Multi-threaded case

```
int MPI_Init_thread(int *argc, char ***argv,  
                    int required, int *provided)
```

```
MPI_Init_thread(required, provided, ierror)  
    INTEGER, INTENT(IN)           :: required  
    INTEGER, INTENT(OUT)          :: provided  
    INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

with `required` taking the following values

- `MPI_THREAD_SINGLE`: Only one thread will execute
- `MPI_THREAD_FUNNELED`: Only main thread makes MPI calls
- `MPI_THREAD_SERIALIZED`: Only one thread at a time makes MPI calls
- `MPI_THREAD_MULTIPLE`: No restrictions



# MPI Finalisation

- ▶ Routine to clean up all MPI state

```
int MPI_Finalize(void)
```

```
MPI_Finalize(ierr)
```

```
INTEGER, OPTIONAL, INTENT(OUT) :: ierr
```

- ▶ Before calling for finalisation, a process must perform all MPI calls to complete its involvement in MPI communications

# MPI Communicators

- ▶ MPI **communicators** (more precisely: intra-communicators) allow to group MPI processes
  - ▶ Processes within a communicator can communicate among each other
- ▶ MPI communicators define a communication context
  - ▶ Communications within different communicators will never be mixed up
- ▶ Predefined Communicators
  - ▶ MPI\_COMM\_WORLD: Group including all processes
  - ▶ MPI\_COMM\_SELF: Group containing only the current process



# MPI Process Identification: Rank

- ▶ All processes within a communicator are assigned a unique identify called **rank**
  - ▶ Ranks are integer values starting from zero
- ▶ An MPI process may be part of multiple communicators and therefore have multiple ranks
  - ▶ Typically, rank refers to the identity of a process within the communicator `MPI_COMM_WORLD`

# Compiling MPI Programs

- ▶ MPI programs should be compiled using the wrapper compilers provided by the different MPI implementations, e.g.
  - ▶ C:
    - ▶ mpicc, mpiicc
  - ▶ C++:
    - ▶ mpicxx, mpiCC, mpic++, mpiicpc
  - ▶ Fortran:
    - ▶ mpifc, mpif77, mpif90, mpiifort
- ▶ Most wrappers allow to show what they are doing
  - ▶ Example OpenMPI: `mpicc --showme`

# Running MPI Programs

- ▶ Starting an MPI program requires a tool called `mpiexec`
  - ▶ Not standard compliant alternative: `mpirun`
- ▶ Typically, a single MPI is executed:  

```
mpiexec -n 8 ./hello_mpi.x
```

  - ▶ The option `-n <numprocs>` defines the number of processes that are started
  - ▶ Other options are MPI implementation specific
- ▶ MPI also allows to execute multiple programs, e.g.  

```
mpiexec -n 1 ./master.x : -n 23 ./worker.x
```



# Content

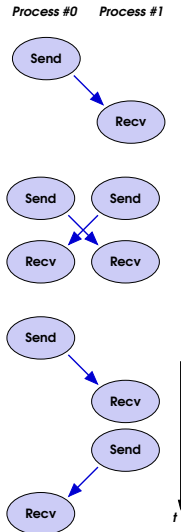
Introduction

First MPI Program

Blocking Send-Receive

# Digression: Communication Patterns

- ▶ Communication patterns involving a pair of processes
- ▶ **Ping:**
  - ▶ Process #0 sends a message to process #1 (Ping)
- ▶ **PingPing:**
  - ▶ Each process sends a message to the other process
- ▶ **PingPong:**
  - ▶ Process #0 sends a message to process #1 (Ping)
  - ▶ After receiving the message process #1 returns it to process #0 (Pong)



# Ping in C using MPI

```
1 #include <stdio.h>
2 #include <mpi.h>
3
4 int main(int argc ,char *argv[])
5 {
6     int irank;
7     int x, y;
8     int src, dst;
9     MPI_Status status;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &irank);
13
14     x = 100 + irank; y = -1;
15     src = 1; dst = 0;
16
17     if (irank == src) {
18         MPI_Send(&x, 1, MPI_INT, dst, 99, MPI_COMM_WORLD);
19         printf("#%d: have sent x=%d\n", irank, x);
20     }
21     else if (irank == dst) {
22         MPI_Recv(&y, 1, MPI_INT, src, 99, MPI_COMM_WORLD, &status);
23         printf("#%d: have received y=%d\n", irank, y);
24     }
25
26     MPI_Finalize();
27
28     return 0;
29 }
```



# Ping in Fortran 2008 using MPI

```
1 PROGRAM PING
2
3 USE mpi_f08
4
5 IMPLICIT NONE
6
7 INTEGER          :: irank, ierr
8 INTEGER          :: x, y
9 INTEGER          :: src, dst
10 TYPE(MPI_Status) :: xstat
11
12 CALL MPI_Init(ierr)
13 CALL MPI_Comm_rank(MPI_COMM_WORLD, irank);
14
15 x = 100 + irank; y = -1;
16 src = 1; dst = 0;
17
18 IF (irank == src) THEN
19     CALL MPI_Send(x, 1, MPI_INT, dst, 99, MPI_COMM_WORLD, ierr)
20     WRITE(*,*) "#", irank, ": have sent x=", x
21 ELSE IF (irank == dst) THEN
22     CALL MPI_Recv(y, 1, MPI_INT, src, 99, MPI_COMM_WORLD, xstat, ierr)
23     WRITE(*,*) "#", irank, ": have received y=", y
24 END IF
25
26 CALL MPI_Finalize();
27
28 END PROGRAM
```

# Standard MPI\_Send

- ▶ Syntax of a blocking send operation:

```
1 int MPI_Send(  
2     const void* buf,           /* Pointer to send buffer */  
3     int count,                 /* Number of elements */  
4     MPI_Datatype datatype,    /* Data type */  
5     int dest,                  /* Rank of destination */  
6     int tag,                   /* Message tag */  
7     MPI_Comm comm             /* Communicator */  
8 );
```

- ▶ The **message buffer** is described by `buf`, `count` and `datatype`
- ▶ The **target** is identified by `dest` and `comm`
- ▶ Messages are sent with an accompanying user-defined integer **tag** that is used by receiver for identification
  - ▶ Allow send-receiver pairs to perform concurrent communications

# MPI Data Types

- ▶ MPI is aware of intrinsic data types in the supported languages (C, Fortran, C++)
- ▶ Benefits
  - ▶ No need for buffer length in Bytes to be computed by programmer
  - ▶ Facilitate some support of type checking
  - ▶ Automatic data conversion, e.g.
    - ▶ Little versus big endian
    - ▶ Interfacing C and Fortran routines

# MPI Predefined Data Types

C:

MPI data type	C data type
MPI_CHAR	char
MPI_INT	signed int
MPI_FLOAT	float
MPI_DOUBLE	double

Fortran:

MPI data type	Fortran data type
MPI_CHAR	CHARACTER(1)
MPI_INTEGER	INTEGER
MPI_REAL	REAL
MPI_DOUBLE_PRECISION	DOUBLE PRECISION

# Standard MPI\_Receive

- ▶ Syntax of a blocking receive operation:

```
1 int MPI_Receive(  
2     const void* buf,           /* Pointer to receive buffer */  
3     int count,                 /* Number of elements */  
4     MPI_Datatype datatype,    /* Data type */  
5     int source,               /* Rank of source */  
6     int tag,                  /* Message tag */  
7     MPI_Comm comm,           /* Communicator */  
8     MPI_Status *status       /* Status object */  
9 );
```

- ▶ Waits until a matching (on source and tag) message is received
  - ▶ After function returned, receive buffer can be used
- ▶ Instead of providing a specific source or tag, MPI\_SOURCE\_ANY or MPI\_TAG\_ANY may be used
- ▶ The status object on return contains further information
- ▶ Receiving fewer than count occurrences of datatype is OK, but receiving more is an error

# Return Status Object

- ▶ Returned information
  - ▶ `MPI_SOURCE`: Source of received message
  - ▶ `MPI_TAG`: Tag of received message
  - ▶ `MPI_ERROR`: Error code
- ▶ Language binding
  - ▶ C: Structure with 3 elements, e.g. `status.MPI_SOURCE`
  - ▶ Fortran with `USE mpi` or `INCLUDE 'mpif.h'`: Array of length 3, e.g. `status(MPI_TAG)`
  - ▶ Fortran with `USE mpi_f08`: Structure with 3 elements, e.g. `status%MPI_ERROR`

# Error Codes

- ▶ The MPI standard leaves it to the implementations to define specific error codes, which must comply to the following rules
  - ▶ Predefined code 0 = `MPI_SUCCESS`
  - ▶ Other codes must be in the range `[1, MPI_ERR_LASTCODE]`
- ▶ The standard defines a set of error classes and a function to translate error code to error class:

```
int MPI_Error_class(int errorcode, int *errorclass)
```

```
MPI_ERROR_CLASS(ERRORCODE, ERRORCLASS, IERROR)
```

```
INTEGER ERRORCODE, ERRORCLASS, IERROR
```

- ▶ Error class examples:

<code>MPI_ERR_BUFFER</code>	Invalid buffer pointer
<code>MPI_ERR_RANK</code>	Invalid rank

# Standard MPI\_Sendrecv

- If all processes send and receive messages (PingPong pattern) send and receive can be initiated using a single MPI call:

```
1 int MPI_Sendrecv(  
2     const void *sendbuf,      /* Send buffer */  
3     int sendcount,            /* Number of elements */  
4     MPI_Datatype sendtype,    /* Data type */  
5     int dest,                 /* Rank of destination */  
6     int sendtag,              /* Send tag */  
7     void *recvbuf,            /* Receive buffer */  
8     int recvcnt,              /* Number of elements */  
9     MPI_Datatype recvtype,    /* Data type */  
10    int source,               /* Rank of source */  
11    int recvtag,              /* Receive tag */  
12    MPI_Comm comm,           /* Communicator */  
13    MPI_Status *status       /* Status object */  
14 );
```



# Probe Incoming Messages

- ▶ MPI allows incoming messages to be checked for, without actually receiving them, using the blocking function `MPI_Mprobe` (and similar variants)

```
1 int MPI_Mprobe(  
2     int source,           /* Source */  
3     int tag,              /* Tag */  
4     MPI_Comm comm,        /* Communicator */  
5     MPI_Message *message, /* Message handler */  
6     MPI_Status *status    /* Status object */  
7 );
```

- ▶ Instead of a specific source and/or tag, `MPI_SOURCE_ANY` or `MPI_TAG_ANY` may be used
- ▶ Information about the message length can be obtained using `MPI_Get_Count`

```
1 int MPI_Get_count(  
2     const MPI_Status *status, /* Status object */  
3     MPI_Datatype datatype,    /* Data type */  
4     int *count                /* Message length */  
5 );
```

# Probe Incoming Messages: Example

Function that allows to receive messages with an arbitrary number of integers:

```
1 #include <stdlib.h>
2 #include <stdio.h>
3 #include <mpi.h>
4
5 int recv_int(int* buf) {
6     int len;
7     MPI_Status status;
8     MPI_Message msg;
9
10    /* Blocking check for message, do not yet receive */
11    MPI_Mprobe(0, 0, MPI_COMM_WORLD, &msg, &status);
12
13    /* Get message length, allocate buffer and receive */
14    MPI_Get_count(&status, MPI_INT, &len);
15
16    buf = (int *) calloc(len, sizeof(int));
17    MPI_Mrecv(buf, len, MPI_INT, &msg, &status);
18
19    printf("Recieved %d integers from process 0.\n", len);
20
21    return len;
22 }
```

# Timing Measurement

- ▶ MPI defines a timer that returns time in units of seconds:

```
double MPI_Wtime(void)
```

```
DOUBLE PRECISION MPI_Wtime()
```

- ▶ Example program in C:

```
1 double tstart = MPI_Wtime();  
2 /* stuff to be timed */  
3 double tend = MPI_Wtime();  
4 printf("Execution time: %.2e seconds\n", tend - tstart);
```

- ▶ Time is measured as difference to some time in the past that is guaranteed not to change during process execution
- ▶ The times returned are local to the node that called them
- ▶ Time resolution can be checked using MPI\_Wtick()

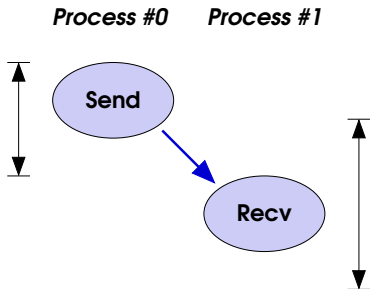
# Time Measurement: Pitfall (1/3)

- ▶ Will the modified version of the MPI Ping program correctly measure the time to communicate one integer?

```
1  src = 1; dst = 0;
2
3  double tstart = MPI_Wtime();
4  if (id == src) {
5      MPI_Send(&x, 1, MPI_INT, dst, 99, MPI_COMM_WORLD);
6      printf("#%d: have sent x=%d\n", id, x);
7  }
8  else if (id == dst) {
9      MPI_Recv(&y, 1, MPI_INT, src, 99, MPI_COMM_WORLD, &status);
10     printf("#%d: have received y=%d\n", id, y);
11 }
12 double tend = MPI_Wtime();
13 printf("Time for communication: %.2es\n", tend - tstart);
```

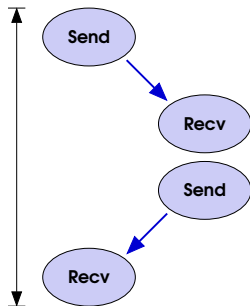
## Time Measurement: Pitfall (2/3)

- ▶ Answer: No
- ▶ Process #0 measures the time it takes to complete the send operation (plus printing a string)
- ▶ Process #1 measures the time from an arbitrary starting point until the receive operation completed



# Time Measurement: Pitfall (3/3)

- ▶ Alternative: Implement a PingPong communication pattern to measure **round-trip time**



# Finish with Past Networking Technology



[Wikipedia Commons]

Morse code receiver with paper tape recorder