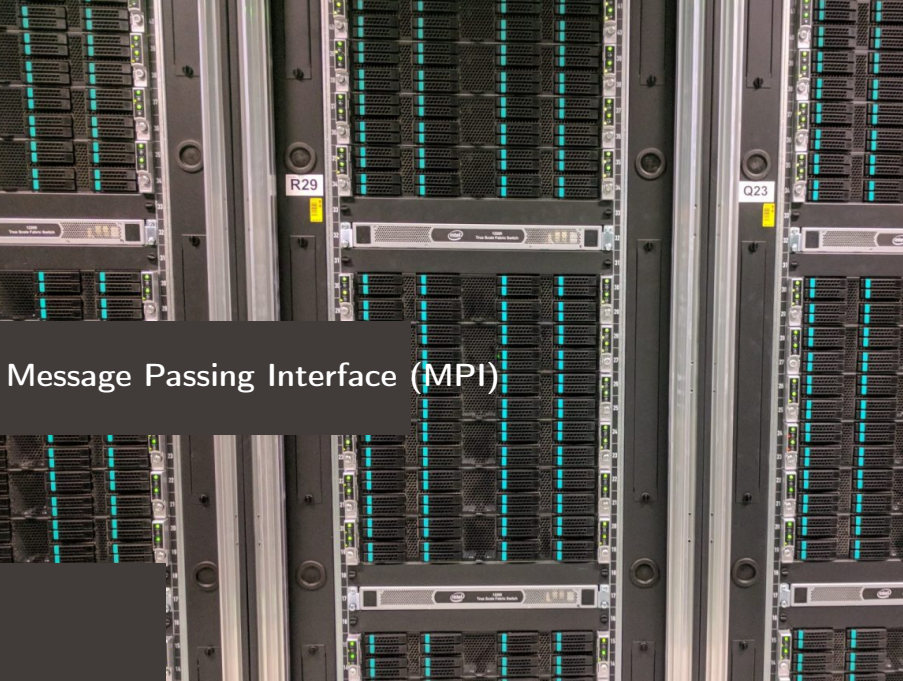# MATH-454 Parallel and High Performance Computing
# Lecture 4: MPI basics

**Pablo Antolin**

**Slides of N. Richart, E. Lanti, V. Keller's lecture notes**
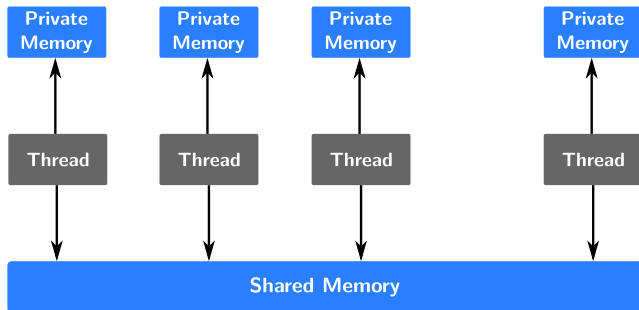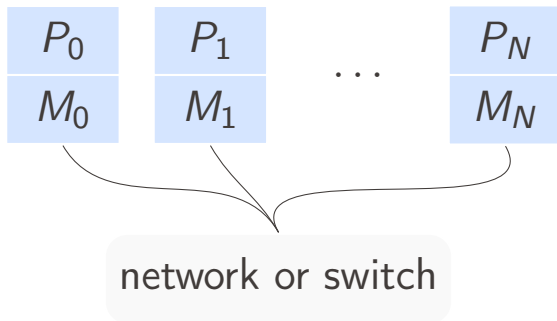
March 20 2025

SCITAS

EPFL

# Message Passing Interface (MPI)

EPFL

SCITAS

- Introduce distributed memory programming paradigm.
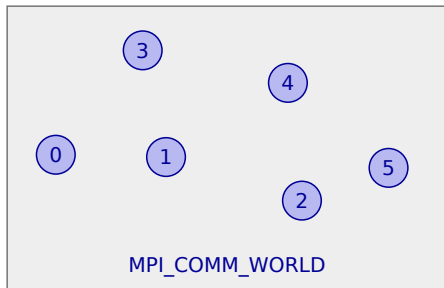- Point-to-point communications.
- Collective communications.

- MPI is a *Message-Passing Interface* specification.
- There are many implementations (Intel MPI, OpenMPI, MPICH, MVAPICH, etc.)
- Library interface, not a programming language.
- It is standarized:
  - ▶ Defined by the MPI forum
  - ▶ Current version is MPI 4.1
- As such, it is portable, flexible, and efficient.
- Interface to C and Fortran in standard. C interface usable with C++. Also MPI bindings for other languages (Python, Julia, MATLAB, etc.).

mpi/hello_mpi.cc

```cpp
1  #include <iostream>
2  #include <mpi.h>
3
4  int main(int argc, char *argv[]) {
5    MPI_Init(&argc, &argv);
6
7    int size, rank;
8    MPI_Comm_size(MPI_COMM_WORLD, &size);
9    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
10
11   std::cout << "I am process " << rank << " out of " << size <<
       ↪  std::endl;
12
13   MPI_Finalize();
14   return 0;
15 }
```

MPI_COMM_WORLD

- MPI code is bounded by a `MPI_Init` and a `MPI_Finalize`.
- MPI starts $N$ processes numbered $0, 1, \ldots, N - 1$.
  The number of every process is usually denoted as *rank*.
- They are grouped in a *communicator* of *size N*.
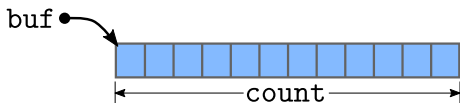- After init, MPI provides a default communicator called **MPI_COMM_WORLD**.

- Point-to-Point (One-to-One).
- Collectives (One-to-All, All-to-One, All-to-All).
- One-sided/Shared memory (One-to . . . ).
- Blocking and Non-Blocking of all types.

### Syntax

```
1 int MPI_Ssend(const void *buf, int count, MPI_Datatype datatype,
  ↪ int dest, int tag, MPI_Comm comm);
2
3 int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int
  ↪ source, int tag, MPI_Comm comm, MPI_Status *status);
```

- **buf**: pointer to the data to send/receive.
- **count**: number of elements to send/receive.
- **datatype**: datatype of the data to send/receive.
- **dest**, **source**: the rank of the destination/source of the communication.
- **tag**: a message tag to differentiate the communications.
- **comm**: communicator in which the communication happens.
- **status**: object containing information on the communication.

- Buffer is a pointer to the first data (`buf`), a size (`count`), and a `datatype`.
- Datatypes (extract):
  - ▶ `MPI_INT`
  - ▶ `MPI_UNSIGNED`
  - ▶ `MPI_FLOAT`
  - ▶ `MPI_DOUBLE`
- For `std::vector<double> vect`:
  - ▶ `buf = vect.data()`
  - ▶ `count = vect.size()`
  - ▶ `datatype = MPI_DOUBLE`

Constants:

- `MPI_STATUS_IGNORE`: to state that the status is ignored.
- `MPI_PROC_NULL`: placeholder for the source or destination.
- `MPI_ANY_SOURCE`: is a wildcard for the source of a receive.
- `MPI_ANY_TAG`: is a wildcard for the tag of a receive.

`MPI_Status`:

- Structure containing `tag` and `source`

```
1 MPI_Status status;
2 std::cout << "Tag: " << status.tag << " - Source: " <<
   ↪  status.source << std::endl;
```

- Size of the received buffer can be asked using the status:

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

mpi/send_recv.cc

```
16 MPI_Init(NULL, NULL);
17 MPI_Comm_rank(MPI_COMM_WORLD, &rank);
18 MPI_Comm_size(MPI_COMM_WORLD, &size);
19
20 assert(size == 2 && "Works only with 2 procs");
21
22 int tag = 0;
23
24 if (rank == 0) {
25   fill_buffer(buf);
26   MPI_Ssend(buf.data(), buf.size(), MPI_INT, 1, tag, MPI_COMM_WORLD);
27 } else if (rank == 1) {
28   MPI_Recv(buf.data(), buf.size(), MPI_INT, 0, tag, MPI_COMM_WORLD,
     ↪  MPI_STATUS_IGNORE);
29 }
```

<u>Rank 0</u>     <u>Rank 1</u>

send(1,0)  $\longrightarrow$  recv(0,0)

finalize()     finalize()

- `MPI_Ssend`: (S for Synchronous) function returns when other end posted matching recv and the buffer can be safely reused.
- `MPI_Bsend`: (B for Buffer) function returns immediately, send buffer can be reused immediately.
- `MPI_Rsend`: (R for Ready) can be used only when a receive is already listening (waiting).
- `MPI_Send`: acts like `MPI_Bsend` on small arrays, and like `MPI_Ssend` on bigger ones.

- `MPI_Ssend`: (S for Synchronous) function returns when other end posted matching recv and the buffer can be safely reused.
- `MPI_Bsend`: (B for Buffer) function returns immediately, send buffer can be reused immediately.
- `MPI_Rsend`: (R for Ready) can be used only when a receive is already listening (waiting).
- `MPI_Send`: acts like `MPI_Bsend` on small arrays, and like `MPI_Ssend` on bigger ones.

Be careful with deadlocks and race conditions!!

mpi/send_recv_deadlock.cc

```
25 assert(size == 2 && "Works only with 2 procs");
26 int tag = 0;
27
28 if (rank == 0)
29 {
30   fill_buffer(buf);
31   MPI_Ssend(buf.data(), buf.size(), MPI_INT, 1, tag, MPI_COMM_WORLD);
32   MPI_Recv(buf.data(), buf.size(), MPI_INT, 1, tag, MPI_COMM_WORLD,
   ↪   MPI_STATUS_IGNORE);
33 }
34 else if (rank == 1)
35 {
36   MPI_Ssend(buf.data(), buf.size(), MPI_INT, 0, tag, MPI_COMM_WORLD);
37   MPI_Recv(buf.data(), buf.size(), MPI_INT, 0, tag, MPI_COMM_WORLD,
   ↪   MPI_STATUS_IGNORE);
38 }
```

Rank 0          Rank 1

send(1,0)      send(0,0)

barrier!        barrier!

mpi/send_recv_deadlock_2.cc

```
16  MPI_Init(NULL, NULL);
17  MPI_Comm_rank(MPI_COMM_WORLD, &rank);
18  MPI_Comm_size(MPI_COMM_WORLD, &size);
19
20  // Let's now assume that we may have 2 or more processes.
21  int tag = 0;
22
23  if (rank == 0) {
24    fill_buffer(buf);
25    MPI_Ssend(buf.data(), buf.size(), MPI_INT, 1, tag, MPI_COMM_WORLD);
26  } else {
27    MPI_Recv(buf.data(), buf.size(), MPI_INT, 0, tag, MPI_COMM_WORLD,
       ↪  MPI_STATUS_IGNORE);
28  }
29  MPI_Finalize();
```

With 2 processes:

| Rank 0 | | Rank 1 |
|--------|-----|--------|
| send(1,0) | $\longrightarrow$ | recv(0,0) |
| finalize() | | finalize() |

With 3 processes:

| Rank 0 | | Rank 1 | Rank 2 |
| --- | --- | --- | --- |
| send(1,0) | $\longrightarrow$ | recv(0,0) | recv(0,0) |
| finalize() | | finalize() | barrier! |

## Syntax

```c
int MPI_Sendrecv(const void *sendbuf, int sendcount, MPI_Datatype
  sendtype, int dest, int sendtag, void *recvbuf, int recvcount,
  MPI_Datatype recvtype, int source, int recvtag, MPI_Comm comm,
  MPI_Status *status);
```

- Combines a send and a receive, to help mitigate deadlocks.
- Has a in-place variant `MPI_Sendrecv_replace`.

### Syntax

```
1 int MPI_Isend(const void *buf, int count, MPI_Datatype datatype,
  ↪  int dest, int tag, MPI_Comm comm, MPI_Request *request);
2
3 int MPI_Irecv(void *buf, int count, MPI_Datatype datatype, int
  ↪  source, int tag, MPI_Comm comm, MPI_Request *request);
```

- **I** for *immediate.*
- `request` in addition to parameters from blocking version.
  It is an object attached to the communication.
- Receive does not have a status.
- The communication starts but is not completed.
- **S**, **B**, and **R** variants are also defined.

### Syntax

```
1  int MPI_Wait(MPI_Request *request, MPI_Status *status);
2
3  int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status);
```

- For non-blocking communications, completion should be "manually" checked.
- `MPI_Test` or `MPI_Wait`.
- Send completed means the buffer can be reused.
- Receive completed means the buffer can be read.
- `status` is set at completion.
- `flag` is **true** if completed, **false** otherwise.

**EPFL**

mpi/isend_recv.cc

```cpp
22 MPI_Request request;
23 if (rank == 0) {
24   fill_buffer(buf);
25   MPI_Isend(buf.data(), buf.size(), MPI_INT, 1, 0, MPI_COMM_WORLD,
     ↪ &request);
26 } else if (rank == 1) {
27   MPI_Recv(buf.data(), buf.size(), MPI_INT, 0, 0, MPI_COMM_WORLD,
     ↪ MPI_STATUS_IGNORE);
28 }
29
30 // While MPI communications are performed here I can do computations as
   ↪ long as buf is not modified in rank 0, or read/modified in rank 1.
31
32 MPI_Wait(&request, MPI_STATUS_IGNORE);
33 MPI_Finalize();
```

- `MPI_Waitall`, `MPI_Testall` waits or tests completion of all the pending requests.
- `MPI_Waitany`, `MPI_Testany` waits or tests completion of one out of many.
- `MPI_Waitsome`, `MPI_Testsome` waits or tests completion of all the specified requests.
- For arrays of statuses can use `MPI_STATUSES_IGNORE`.
- `MPI_Request_get_status` equivalent to `MPI_Test` but does not free completed requests.

### Syntax

```
1  int MPI_Iprobe(int source, int tag, MPI_Comm comm, int *flag,
   ↪  MPI_Status *status);
2
3  int MPI_Probe(int source, int tag, MPI_Comm comm, MPI_Status
   ↪  *status);
```

- Check incoming messages without receiving.
- Immediate variant returns **true** if matching message exists.
- Can be used in combination with a successive `MPI_Get_count` for deducing the size of an incoming message before actually recieving it. Thus, we can allocate a buffer for holding the message.

### Syntax

```
1 int MPI_Barrier(MPI_Comm comm);
```

- Collective communications **must** be called by all processes in the communicator.
- Barrier is hard synchronization.
- Avoid as much as possible.

## Syntax

```
1 int MPI_Bcast(void *buffer, int count, MPI_Datatype datatype, int
  ↪   sender, MPI_Comm comm);
```

- The **sender** process sends data to every other process.

### Syntax

```
1 int MPI_Bcast(void *buffer, int count, MPI_Datatype datatype, int
  ↪  sender, MPI_Comm comm);
```

- The **sender** process sends data to every other process.

**EPFL**

### Syntax

```
1 int MPI_Scatter(const void *sendbuf, int sendcount, MPI_Datatype
↪    sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
↪    int sender, MPI_Comm comm);
```

- The **sender** process sends a piece of the data to all processes.
- The sendbuf, sendcount, and sendtype are only relevant on the **sender**.

$P_0$

$P_1$

$P_2$

$P_3$

**EPFL**

### Syntax

```
1 int MPI_Scatter(const void *sendbuf, int sendcount, MPI_Datatype
  ↪  sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
  ↪  int sender, MPI_Comm comm);
```

- The **sender** process sends a piece of the data to all processes.
- The sendbuf, sendcount, and sendtype are only relevant on the **sender**.

### Syntax

```
1  int MPI_Gather(const void *sendbuf, int sendcount, MPI_Datatype
   ↪   sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
   ↪   int receiver, MPI_Comm comm);
```

- Every process sends its data to the **receiver** process.
- The `recvbuf`, `recvcount`, and `recvtype` are only relevant on the **receiver**.
- `recvcount` is the size per process, not the total size.

## Syntax

```
1 int MPI_Gather(const void *sendbuf, int sendcount, MPI_Datatype
  ↪  sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
  ↪  int receiver, MPI_Comm comm);
```

- Every process sends its data to the **receiver** process.
- The `recvbuf`, `recvcount`, and `recvtype` are only relevant on the **receiver**.
- `recvcount` is the size per process, not the total size.

$P_0$

$P_1$

$P_2$

$P_3$

## Syntax

```
1 int MPI_Allgather(const void *sendbuf, int sendcount, MPI_Datatype
  ↪   sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
  ↪   MPI_Comm comm);
```

- Every process sends its data to all the other processes.

**EPFL**

### Syntax

```
1 int MPI_Allgather(const void *sendbuf, int sendcount, MPI_Datatype
  ↪  sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
  ↪  MPI_Comm comm);
```
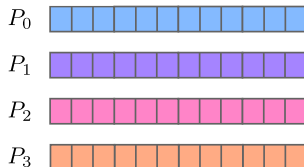
- Every process sends its data to all the other processes.

### Syntax

```
1 int MPI_Alltoall(const void *sendbuf, int sendcount, MPI_Datatype
  ↪   sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
  ↪   MPI_Comm comm);
```
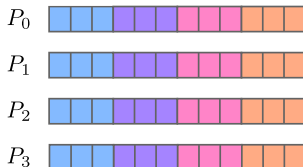
- Every process sends a piece of its data to all the other processes.

### Syntax

```
1 int MPI_Alltoall(const void *sendbuf, int sendcount, MPI_Datatype
  ↪  sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype,
  ↪  MPI_Comm comm);
```

- Every process sends a piece of its data to all the other processes.

### Syntax

```
1 int MPI_Reduce(const void *sendbuf, void *recvbuf, int count,
  ↪  MPI_Datatype datatype, MPI_Op op, int receiver, MPI_Comm comm);
```
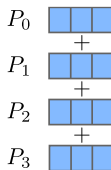
- Data from all processes are reduced on the **receiver** process.
- Common operations being `MPI_SUM`, `MPI_MAX`, `MPI_MIN`, `MPI_PROD`.
- A `MPI_Allreduce` variant exists where all processes get the results.
- `MPI_IN_PLACE` can be passed to `sendbuf` of **receiver**, for *reduce*, and of all processes, for *allreduce*.

$P_0$

$P_1$

$P_2$

$P_3$

**EPFL**

### Syntax

```
1 int MPI_Reduce(const void *sendbuf, void *recvbuf, int count,
  ↪   MPI_Datatype datatype, MPI_Op op, int receiver, MPI_Comm comm);
```

- Data from all processes are reduced on the **receiver** process.
- Common operations being `MPI_SUM`, `MPI_MAX`, `MPI_MIN`, `MPI_PROD`.
- A `MPI_Allreduce` variant exists where all processes get the results.
- `MPI_IN_PLACE` can be passed to `sendbuf` of **receiver**, for *reduce*, and of all processes, for *allreduce*.

### Syntax

```
1 int MPI_Reduce(const void *sendbuf, void *recvbuf, int count,
  ↪  MPI_Datatype datatype, MPI_Op op, int receiver, MPI_Comm comm);
```

- Data from all processes are reduced on the **receiver** process.
- Common operations being MPI_SUM, MPI_MAX, MPI_MIN, MPI_PROD.
- A MPI_Allreduce variant exists where all processes get the results.
- MPI_IN_PLACE can be passed to sendbuf of **receiver**, for *reduce*, and of all processes, for *allreduce*.

$P_0$

$P_1$

$P_2$

$P_3$

### Syntax

```
1 int MPI_Scan(const void *sendbuf, void *recvbuf, int count,
  ↪ MPI_Datatype datatype, MPI_Op op, MPI_Comm comm);
2
3 int MPI_Exscan(const void *sendbuf, void *recvbuf, int count,
  ↪ MPI_Datatype datatype, MPI_Op op, MPI_Comm comm);
```

- Performs the prefix reduction on data.
- `MPI_Scan`: on process $i$ contains the reduction of values from processes $[0, i]$
- `MPI_Exscan`: on process $i$ contains the reduction of values from processes $[0, i[$.
- `MPI_IN_PLACE` can be passed to `sendbuf`.