

# Parallel Programming

## MPI – Part 1

Prof. **Paolo Bientinesi**

HPAC, RWTH Aachen  
pauldj@aices.rwth-aachen.de

WS18/19



**High Performance and  
Automatic Computing**

# Preliminaries

---

- Distributed-memory architecture
- Topologies → later on
- Back in the days: Node  $\equiv$  CPU  $\equiv$  process
- Nowadays: NODE  $\rightarrow$  CPUs  $\rightarrow$  multi cores  $\rightarrow$  many processes
- **Assumption: Fully connected topology**
- **Assumption: Each process can simultaneously send and receive**
- **Assumption: Messages in opposite directions do not cause a conflict**

# What is “MPI”?

---

- A **library**, not a language, not a program.
- De-facto standard for distributed-memory parallelism.
- In fact, it's the specification of a library, not the actual implementation.
- MPI defines the interface, the functionality and the semantics of functions that deliver a message passing mechanism.
- Idea: clear separation between data communication and application.
- Both open-source and proprietary implementations.
- [www.mpi-forum.org](http://www.mpi-forum.org)

# History – before 1994

---

- Before MPI, no standards
- Different computers, different needs  
⇒ **many** message passing environments
- N-cube, P4, PICL, PVM, ISIS, Express, Zipcode;  
Intel NX, IBM EUI, IBM CCL, ...
- **A lot of duplication!**
- **No portability whatsoever**

# History

---

- **[1992]** First “MPI Forum” meeting (Supercomputing '92)
- **[1993–94]** Seven “MPI Forum” meetings. Working on the MPI standard
- **[1994]** Release of the first MPI standard: MPI-1
- **[1995]** First implementations of MPI: MPICH, LAM MPI, ...
- **[1998]** Release of the second MPI standard: MPI-2  
More than 100 new functions!
- **[2002]** Complete implementations of MPI-2  
Dynamic process management, 1-sided communication, MPI-I/O
- **[2012]** Release of MPI-3  
Non-blocking collectives, sparse collectives, ...

Thanks to Jesper Larsson Träff (TU Wien).

## "Minimal" MPI

<code>MPI_Init(...)</code>	MPI Initialization
<code>MPI_Comm_size(...)</code>	How many processes are there?
<code>MPI_Comm_rank(...)</code>	What rank am I?
<code>MPI_Send(...)</code>	Send data to another process
<code>MPI_Recv(...)</code>	Receive data from another process
<code>MPI_Finalize( )</code>	MPI termination

## `int MPI_Init( ... )`

- `MPI_Init(&argc, &argv);`
- First MPI function
- Args not specified; an implementation might use them
- Query: `MPI_Initialized`

## `int MPI_Finalize()`

- Last MPI function
- No arguments
- Query: `MPI_Finalized`

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

- Returns the number of processes in the communicator `comm`
- Communicator: for now `MPI_COMM_WORLD`  $\equiv$  “everybody”

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

- Returns the rank of the calling process within the communicator
- The rank is THE unique process identifier!
- NOTE: each process (rank) can be multi-threaded



# Send ↔ Recv

---

- Objective: data movement
- MPI\_Send and MPI\_Recv must be matched
- Blocking communication

Necessary information:

Send	Recv
dest	source
*buffer	*target
size	size
datatype	datatype
tag	tag
comm	comm

```
int MPI_Send(*buffer, count, datatype, dest, tag, comm)
```

- \*buffer is an address!
- count is indispensable; so is datatype
- dest is a rank (in comm)
- tag is an integer

```
int MPI_Recv(*target, count, datatype, source, tag, comm, *status)
```

- \*target, datatype as for the Send
- count is the size of target. Actual size: MPI\_Get\_count
- source is either a rank (in comm) or MPI\_ANY\_SOURCE
- tag is either an integer or MPI\_ANY\_TAG
- \*status on exit, contains info about the message

# Communication modes

---

- Point-to-point                      Send, Receive, . . . , many variants
- Collective                              Broadcast, Reduce, . . . , many operations
- One-sided                                Put, Get.