

# Application speedup: an example with NAMD

FELICE PANTALEO – CERN DANIELE CESINI – INFN-CNAF



# Speedup of a real application





- A parallel molecular dynamics code designed for high-performance simulation of large biomolecular systems
- Can be download here: <u>https://www.ks.uiuc.edu/Research/namd/</u>
- Installing...
  - https://github.com/infn-esc/esc22/blob/main/hands-on/mpi/Make\_and\_Install\_NAMD.sh
- Running on real molecule...
  - https://github.com/infn-esc/esc22/blob/main/hands-on/mpi/HowTo\_Launch\_NAMD\_on\_APOA1.sh
- GPU Porting available on CUDA



Speedup: measures the increased performance in running in parallel on P processors

$$S(P) = \frac{T_{Seq}(1)}{T_{Par}(P)}$$

Perfect Linear Speedup: no overhead due to parallelism. Speedup equals the number of processors

$$S(P) = P$$



Efficiency: measures how well the hardware resources (processors) are utilized

$$\varepsilon = \frac{T_{seq}}{P*TPar(P)} = \frac{S(P)}{P}$$

T = Elapsed TimeP = Number of processors Used

#### Speedup examples – NAMD APOA1







#### APOA1

Apolipoprotein A1 (ApoA1) is the major protein component of high-density lipoprotein (HDL) in the bloodstream and plays a specific role in lipid metabolism. The ApoA1 benchmark consists of 92,224 atoms and has been a standard NAMD cross-platform benchmark for years.

### Speedup example – NAMD APOA1





#### Speedup – NAMD STMV

#### STMV namd2.9 @avoton C2730 @1.7GHz Speedup NP

STMV namd2.9 @avoton C2730 @1.7GHz Efficiency NP

#### # of Atoms

#### STMV

Satellite Tobacco Mosaic Virus (STMV) is a small, icosahedral plant virus that worsens the symptoms of infection by Tobacco Mosaic Virus (TMV). SMTV is an excellent candidate for research in molecular dynamics because it is relatively small for a virus and is on the medium to high end of what is feasible to simulate using traditional molecular dynamics in a workstation or a small server.

1,066,628







#### Let's add dimensions....accelerators

#### Reasonable Speedups - 1

Chip A





 $Perf_A = Eff_A * Peak_A(Comp or BW)$ 

 $Perf_B = Eff_B * Peak_B(Comp or BW)$ 

$$Speedup \frac{B}{A} = \frac{Perf_B}{Perf_A} = \frac{Eff_B}{Eff_A} * \frac{Peak_A(Comp\_or\_BW)}{Peak_B(Comp\_or\_BW)}$$

© Tim Matson @ESC school



### Reasonable Speedups - 3



#### Core i7 960

- Four OoO Superscalar Cores, 3.2GHz
- Peak SP Flop: 102GF/s
- Peak BW: 30 GB/s

#### GTX 280

- 30 SMs (w/ 8 In-order SP each), 1.3GHz
- Peak SP Flop: 933GF/s\*
- Peak BW: 141 GB/s

Assuming both Core i7 and GTX280 have the same efficiency:

	Max Speedup: GTX 280 over Core i7 960
Compute Bound Apps: (SP)	933/102 = 9.1x
Bandwidth Bound Apps:	141/30 = 4.7x

\* 933GF/s assumes mul-add and the use of SFU every cycle on GPU

Source: Victor Lee et. al. "Debunking the 100X GPU vs. CPU Myth", ISCA 2010

© Tim Matson @ESC school



	CPU	GPU
Chip	DUAL IntelXeon Gold 6148 CPU @ 2.40GHz	QUAD NVIDIA Tesla V100 SXM2 32GB
Compute Perf Peak (single precision)	2(socket)*20(core)*2.4(clock GHz)*512/32(avx) = 1.5TF (sp)	FP(32)(float performance) ==> 14.13TF
Bandwidth Peak	2 * socket * 6 (channel/sock) * 20GB/s =240GB/s	4x 143GB/s = 572GB/s

	Max Speed-Up CPU/GPU
Compute Bound App (sp)	14.13/1.5 = 9.5
Bandwidth Bound App	572/240 = 2.4



- Compare the latest GPU against an old CPU
- Highly optimized GPU code vs. unoptimized CPU code
- Compare optimized CUDA vs. Matlab or python
- Parallel GPU code vs. serial, unvectorized CPU code
- Ignore the GPU penalty of moving data across the PCI bus from the CPU to the GPU
- GPUs and other accelerators can be great but be suspicious when people claim speedups of 100+

#### © Tim Matson @ESC school



Parallelism beyond the node: Introduction to MPI Programming

FELICE PANTALEO – CERN DANIELE CESINI – INFN-CNAF

### Reference Material



- MPI Standard: https://www.mpi-forum.org/docs/
- Open-mpi.org: <u>https://www.open-mpi.org/doc/v3.0/man3/MPI\_Wtime.3.php</u>
  - https://www.open-mpi.org/faq/
- MPICH.org: <u>https://www.mpich.org/</u>
- MPI Tutorial: <u>https://mpitutorial.com/</u>
- Message Passing Interface (MPI). Author: Blaise Barney, Lawrence Livermore National
  - https://hpc-tutorials.llnl.gov/mpi/
- Tutorial and exercises @ Argonne National Laboratory: https://www.mcs.anl.gov/research/projects/mpi/tutorial/mpiexmpl/contents.html
- www.google.com

(Credits to Tim Mattson at IntelLab for his ""Hands-on" Introduction to MPI" at ESC15)

#### Multithread vs Multiprocess

men

1.00

- Multithreading and multiprocessing are two ways to achieve multitasking
- A process has its own memory
- •A thread shares the memory with the parent process and other threads within the process.
- pid is process identifier; tid is thread identifier
- (\*)But as it happens, the kernel doesn't make a real distinction between them: threads are just like processes but they share some things (memory, fds...) with other instances of the same group

(\*)https://stackoverflow.com/questions/4517301/difference-between-pid-and-tid

 Inter-process communication is slower due to isolated memory





#### Shared Memory Systems

- Shared memory is memory that may be simultaneously accessed by multiple programs with an intent to provide communication among them or avoid redundant copies
- Shared memory is an efficient means of passing data between programs
- Shared memory systems may use uniform memory access (UMA): all the processors share the physical memory uniformly
- Non-uniform memory access (NUMA): memory access time depends on the memory location relative to a processor





### NUMA Architecture Programming

- A programmer can set an allocation policy for its program using a component of NUMA API called libnuma.
  - a user space shared library that can be linked to applications
  - provides explicit control of allocation policies to user programs.
- The NUMA execution environment for a process can also be set up by using the numactl tool
- Numactl can be used to control process mapping to cpuset and restrict memory allocation to specific nodes without altering the program's source code



#### http://halobates.de/numaapi3.pdf



### Distributed Memory Systems

- Distributed memory refers to a multiprocessor computer system in which each processor has its own private memory
- Computational tasks can only operate on local data
- if remote data is required, the computational task must communicate with one or more remote processors
- In contrast, a shared memory multiprocessor offers a single memory space used by all processors







### Shared vs Distributed Memory Systems







#### Clusters



[a cluster is a] parallel computer system comprising an integrated collection of independent nodes, each of which is a system in its own right, capable of independent operation <u>and derived from products</u> <u>developed and marketed for other stand-alone purposes</u>

© Dongarra et al. : "High-performance computing: clusters, constellations, MPPs, and future directions", Computing in Science & Engineering (Volume:7, Issue: 2)



(\*) Picture from: http://en.wikipedia.org/wiki/Computer\_cluster



# System Topology

ESC

- Knowing where you are is important!!
  - Always try to understand the details of the system you are running on

adimperformance in the second product of the														
NUMANode P#0 (64GB)														
Socket P#0														
L3 (35M8)														
Ľ	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)
Ľ	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)
Ľ	Lli (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	Lli (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	Lli (32KB)
	Core P#0	Core P#1	Core P#2	Core P#3	Core P#4	Core P#5	Core P#6	Core P#8	Core P#9	Core P#10	Core P#11	Core P#12	Core P#13	Core P#14
	PU P#0	PU P#1	PU P#2	PU P#3	PU P#4	PU P#5	PU P#6	PU P#7	PU P#8	PU P#9	PU P#10	PU P#11	PU P#12	PU P#13
Ľ	PU P#28	PU P#29	PU P#30	PU P#31	PU P#32	PU P#33	PU P#34	PU P#35	PU P#36	PU P#37	PU P#38	PU P#39	PU P#40	PU P#41
NLMANode P#1 (6408)														
NUP	MANode P#1 (64G	8)												
Soc	1ANode P#1 (646) tket P#1	8)												
Soc	MANode P#1 (640 cket P#1 _3 (35MB)	8)												
Soc L	1ANode P#1 (64G cket P#1 _3 (35MB) _2 (256KB)	6) L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)	L2 (256KB)
500 [L	(ANode P#1 (640) cket P#1 .3 (35MB) .2 (256KB) .1d (32KB)	B) L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)	L2 (256KB) L1d (32KB)
	(ANode P#1 (640) cket P#1 .3 (35%B) .2 (256KB) .1d (32KB) .1i (32KB)	<pre>b) L2 (256KB) L1d (32KB) L1i (32KB) </pre>	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)	L2 (256KB) L1d (32KB) L1i (32KB)			
	(ANode P#1 (640) cket P#1 .3 (35%B) .2 (256KB) .1d (32KB) .1i (32KB) .1i (32KB)	<pre>b) L2 (256KB) L1d (32KB) L1i (32KB) Core P#1</pre>	L2 (256KB) L1d (32KB) L1i (32KB) Core P#2	L2 (256KB) L1d (32KB) L1i (32KB) Core P#3	L2 (256KB) L1d (32KB) L1i (32KB) Core P#4	L2 (256KB) L1d (32KB) L1i (32KB) Core P#5	L2 (256KB) L1d (32KB) L1i (32KB) Core P#5	L2 (256KB) L1d (32KB) L1i (32KB) Core P#8	L2 (256KB) L1d (32KB) L1i (32KB) Core P#9	L2 (256KB) L1d (32KB) L1i (32KB) Core P#10	L2 (256KB) L1d (32KB) L1i (32KB) Core F#11	L2 (256KB) L1d (32KB) L1i (32KB) Core P#12	L2 (256KB) L1d (32KB) L1i (32KB) Core P#13	L2 (256KB) L1d (32KB) L1i (32KB) Core P#14
	<pre>4ANode P#1 (6468 cket P#1 .3 (35MB) .2 (256KB) .1d (32KB) .1i (32KB) .1i (32KB) Core P#0 PU P#14 .01 E#42</pre>	<ul> <li>E)</li> <li>E2 (256KB)</li> <li>E1d (32KB)</li> <li>E1i (32KB)</li> <li>Core P#1</li> <li>PU P#15</li> <li>PU P#15</li> <li>PU P#42</li> </ul>	L2 (256KB) L1d (32KB) L1i (32KB) Core P#2 PU P#16 PU P#16	L2 (256KB) L1d (32KB) L1i (32KB) Core P#3 PU P#17 PL P#45	L2 (256KB) L1d (32KB) L1i (32KB) Core P## PU P#18 PU P#18	L2 (256KB) L1d (32KB) L1i (32KB) Core P#5 PU P#19 PU P#19 PU P#19	L2 (256KB) L1d (32KB) L1i (32KB) Core P#6 PU P#20 PU P#20	L2 (256KB) L1d (32KB) L1i (32KB) Core P#8 PU P#21 PU P#21 PU P#49	L2 (256KB) L1d (32KB) L1i (32KB) Core P#9 PU P#22 PU P#58	L2 (256KB) L1d (32KB) L1i (32KB) Core P#10 PU P#23 PU P#53	L2 (256KB) L1d (32KB) L1i (32KB) Core P#11 PU P#24 PU P#24 PU P#52	L2 (256KB) L1d (32KB) L1i (32KB) Core P#12 PU P#25 PU P#25 PU P#53	L2 (256KB) L1d (32KB) L1i (32KB) Core P#13 PU P#26 PU P#54	L2 (256KB) L1d (32KB) L1i (32KB) Core P#14 PU P#27 PU P#27 PU P#25
	<pre>4ANode P#1 (646) cket P#1 .3 (35%B) .2 (256KB) .1d (32KB) .1i (32KB) .1i (32KB) .0ore P#0 .PU P#14</pre>	<pre>b) L2 (256KB) L1d (32KB) L1i (32KB) Core P#1 PU P#15 PU P#43</pre>	L2 (256KB) L1d (32KB) L1i (32KB) Core P#2 PU P#16 PU P#44	L2 (256KB) L1d (32KB) L1i (32KB) Core P#3 PU P#17 PU P#45	L2 (256KB) L1d (32KB) L1i (32KB) Core P#1 PU P#18 PU P#46	L2 (256KB) L1d (32KB) L1i (32KB) Core P#5 PU P#19 PU P#47	L2 (256KB) L1d (32KB) L1i (32KB) Core P#6 PU P#20 PU P#48	L2 (256KB) L1d (32KB) L1i (32KB) Core P#S PU P#21 PU P#49	L2 (256KB) L1d (32KB) L1i (32KB) Core P#9 PU P#22 PU P#50	L2 (256KB) L1d (32KB) L1i (32KB) Core P#10 PU P#23 PU P#51	L2 (256KB) L1d (32KB) L1i (32KB) Core P#11 PU P#24 PU P#52	L2 (256KB) L1d (32KB) L1i (32KB) Core P#12 PU P#25 PU P#53	L2 (256KB) L1d (32KB) L1i (32KB) Core P#13 PU P#26 PU P#54	L2 (256KB) L1d (32KB) L1i (32KB) Core P#14 PU P#27 PU P#55

# lstopo --no-io -.txt

#### System Networking





#### System Networking



[cesinihpc@hpc-201-11-40 ~]\$ ifconfig

eth0: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 9000 inet 131.154.184.77 netmask 255.255.255.0 broadcast 131.154.184.255 ether ac:1f:6b:41:d3:00 txqueuelen 1000 (Ethernet) RX packets 126504834 bytes 19185206222 (17.8 GiB) RX errors 0 dropped 0 overruns 0 frame 0 TX packets 17873813 bytes 11835855740 (11.0 GiB) TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

.o: flags=73<UP,L00PBACK,RUNNING> mtu 65536 inet 127.0.0.1 netmask 255.0.0.0 loop txqueuelen 1000 (Local Loopback) RX packets 2722587 bytes 531228851 (506.6 MiB) RX errors 0 dropped 0 overruns 0 frame 0 TX packets 2722587 bytes 531228851 (506.6 MiB) TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

#### [cesinihpc@hpc-201-11-40 ~]\$ ibstatus

infiniband device 'hfil_(	0' port 1 status:
default gid:	fe80:0000:0000:0000:001
base lid:	0x19
sm lid:	0x1
state:	4: ACTIVE
phys state:	5: LinkUp
rate:	100 Gb/sec (4X EDR)
link_layer:	InfiniBand

[cesinihpc@hpc-200-06-18	~]\$ ibstatus
Infiniband device 'qib0'	port 1 status:
default gid:	fe80:0000:0000:0000:0011:7500
base lid:	0xd
sm lid:	0x1
state:	4: ACTIVE
phys state:	5: LinkUp
rate:	40 Gb/sec (4X QDR)
link_layer:	InfiniBand

# The Message Passing Programming Model

- Program consists of a collection of named processes
  - Number of processes almost always fixed at program startup time
  - Local address space per node -- NO physically shared memory.
  - Logically shared data is partitioned over local processes
- Communication happens by explicit send/receive statements

Message can be passed over a network infrastructure o via the main memory, "shared" memory





# Performance and Efficiency Loss?



- The latency of the DRAM can be measured in tens of nanoseconds
- Sending a byte to a networked computer can take 2-3 orders of magnitude longer than DRAM, depending on the interconnect technology
- In using Message Passing, try hard to minimize communication
- In any case, the interconnection technology greatly affects the program performances
  - Ethernet 1Gbs latency O(10.000ns)
  - Infiniband HDR latency O(200ns)
  - DDR4-3600 latency O(60ns)
  - DDR5-5600 latency O(10ns)

#### Latency Numbers every programmer should know

Latency Comparison Numbers (~2012)					
L1 cache reference	0.	5 ns			
Branch mispredict	5	ns			
L2 cache reference	7	ns			14x L1 cache
Mutex lock/unlock	25	ns			
Main memory reference	100	ns			20x L2 cache, 200x L1 cache
Compress 1K bytes with Zippy	3,000	ns	3 us		
Send 1K bytes over 1 Gbps network	10,000	ns	10 us		
Read 4K randomly from SSD*	150,000	ns	150 us		~1GB/sec SSD
Read 1 MB sequentially from memory	250,000	ns	250 us		
Round trip within same datacenter	500,000	ns	500 us		
Read 1 MB sequentially from SSD*	1,000,000	ns	1,000 us	1 ms	~1GB/sec SSD, 4X memory
Disk seek	10,000,000	ns	10,000 us	10 ms	20x datacenter roundtrip
Read 1 MB sequentially from disk	20,000,000	ns	20,000 us	20 ms	80x memory, 20X SSD
Send packet CA->Netherlands->CA	150,000,000	ns	150,000 us	150 ms	

### Communication performances in MPI Applications



## Communication performances in MPI Applications





## Communication performances in MPI Applications



#### ESC22 - BERTINORO

#### MPI is a standard : <u>http://www.mpiforum.org/</u>

- Defines API for C, C++, Fortran77, Fortran90
- Library with diverse functionalities:
  - Communication primitives (blocking, non-blocking)
  - Parallel I/O
  - RMA
  - Neighborhood collectives
- •When you run an MPI program, multiple processes all running the same program are launched working on their own block of data









• Every process runs the same program...

• .....on P processing elements where P can be arbitrarily large

Each process has a unique identifier and runs the version of the program with that particular identifier

the rank - an ID ranging from 0 to (P-1)

Each process access its own private data

•You usually run one process per socket/core depending on the parallelization strategy

And on the system topology

### SPMD – Single Program Multiple Data

Process 1 If pid == 1: a = 5 Send (a,2) Else: Recv(b,1) b++



### MPI Implementations



MPICH

- The initial implementation of the MPI 1.x standard, from Argonne National Laboratory (ANL) and Mississippi State University.
- ANL has continued developing MPICH for over a decade, and now offers MPICH-3.2, implementing the MPI-3.1 standard
- •IBM also was an early implementor, and most early 90s supercomputer companies either commercialized MPICH, or built their own implementation.
- LAM/MPI from Ohio Supercomputer Center
  - another early open implementation..
- Open MPI (not to be confused with OpenMP) was formed by the merging FT-MPI, LA-MPI, LAM/MPI, and PACX-MPI, and is found in many TOP-500 supercomputers.
  - We will use OpenMPI for our exercises!!
- Many other efforts are derivatives of MPICH, LAM, and other works, including, but not limited to, commercial implementations from HP, Intel, Microsoft, and NEC.

#### MPI HelloWorld

Esc

#include <iostream>
#include <mpi.h>

int main(intargc,char\*\*argv){
 int rank, world\_size;

MPI\_Init(&argc,&argv); MPI\_Comm\_rank(MPI\_COMM\_WORLD,&rank); MPI\_Comm\_size(MPI\_COMM\_WORLD,&world\_size); char processor\_name[MPI\_MAX\_PROCESSOR\_NAME]; int name\_len;

MPI\_Get\_processor\_name(processor\_name, &name\_len);

std::cout << "Hello world from processor " << processor\_name << " rank " << rank << " of " << world\_size << std::endl;

```
MPI_Finalize();
return 0;
```

### MPI\_Init and MPI\_Finalize



#include <iostream>
#include <mpi.h>

int main(intargc,char\*\*argv){
 int rank, world\_size;

MPI\_Init(&argc,&argv)<sup>.</sup>

MPI\_Comm\_rank(MPI\_COMM\_WORLD,&rank); MPI\_Comm\_size(MPI\_COMM\_WORLD,&world\_size); char processor\_name[MPI\_MAX\_PROCESSOR\_NAME]; int name\_len;

#### Called before any other MPI functions

- Initializes the library
- Argc and argv are the command line args passed to main
  - Open MPI accepts the C/C++ argc and argv arguments to main, but neither modifies, interprets, nor distributes them

```
MPI_Get_processor_name(processor_name, &name_len);
```

std::cout << "Hello world from processor " << processor\_name << " rank " << rank << " of " world\_size << std::endl; MPI\_Finalize();

- Frees memory allocated by MPI

return 0;

### How many processes?

#include <iostream>
#include <mpi.h>

int main(int argc, char\*\* argv){
 int rank, world\_size;

MPI\_Init(&argc, &argv);

MPI\_Comm\_rank(MPI\_COMM\_WORLD,&rank);

MPI\_Comm\_size(MPI\_COMM\_WORLD,&size);

char processor\_name[MPI\_MAX\_PROCESSOR\_NAME];
int name\_len;

MPI\_Get\_processor\_name(processor\_name, &name\_len);
std::cout << "Hello world from processor " << processor\_name
<< " rank " << rank << " of " << world\_size << std::endl;</pre>

```
MPI_Finalize();
return 0;
```

**Communicators** consist of two parts, a **context** and a **process group**. The communicator lets us control how groups of messages interact.

int MPI\_Comm\_size (MPI\_Comm comm, int\*
size)

- MPI\_Comm, an opaque data type called a communicator. Default context:
   MPI\_COMM\_WORLD (all processes)
- MPI\_Comm\_size returns the number of processes in the process group associated with the communicator


# Who am I? (which is my rank?)

#include <iostream>
#include <mpi.h>

int main(int argc, char\*\* argv){
 int rank, world\_size;

MPI\_Init(&argc, &argv):

MPI\_Comm\_rank(MPI\_COMM\_WORLD,&rank);

MPI\_Comm\_size(MPI\_COMM\_WORLD,&size);

char processor\_name[MPI\_MAX\_PROCESSOR\_NAME];
int name\_len;

MPI\_Get\_processor\_name(processor\_name, &name\_len);
std::cout << "Hello world from processor " << processor\_name
<< " rank " << rank << " of " << world\_size << std::endl;</pre>

```
MPI_Finalize();
return 0;
```

Note that other than init() and finalize(), every MPI function has a communicator which defines the context and group of processes that the MPI functions impact

int MPI\_Comm\_rank (MPI\_Comm comm, int\*
rank)

- MPI\_Comm, an opaque data type called a communicator. Default context:
   MPI\_COMM\_WORLD (all processes)
- MPI\_Comm\_rank returns an integer ranging from 0 to "(num of procs)-1"



# Communicators and Groups - 1



Internally, MPI has to keep up with (among other things) two major parts of a communicator

- the context (or ID) that differentiates one communicator from another
  - prevents an operation on one communicator from matching with a similar operation on another communicator
- the group of processes contained by the communicator

Communicators provides a separate communication space

- It's not unusual to do everything using MPI\_COMM\_WORLD, but for more complex use cases, it might be helpful to have more communicators.
  - MPI\_Comm\_split is the simplest way to create a new communicator
- A Group is a little simpler, since it is just the set of all processes in the communicator.
  - MPI offers function to manage Groups: Union or Intersection
  - Groups can be used to create Communicators



### Communicators and Groups - 2



// Get the rank and size in the original communicator int world\_rank, world\_size; MPI\_Comm\_rank(MPI\_COMM\_WORLD, &world\_rank); MPI\_Comm\_size(MPI\_COMM\_WORLD, &world\_size);

int color = world\_rank / 4; // Determine color based on row

// Split the communicator based on the color and use the // original rank for ordering MPI\_Comm row\_comm; MPI\_Comm\_split(MPI\_COMM\_WORLD, color, world\_rank, &row\_comm);

int row\_rank, row\_size; MPI\_Comm\_rank(row\_comm, &row\_rank); MPI\_Comm\_size(row\_comm, &row\_size);

printf("WORLD RANK/SIZE: %d/%d \t ROW RANK/SIZE: %d/%d\n",

world\_rank, world\_size, row\_rank, row\_size);

MPI\_Comm\_free(&row\_comm);

Split a Large Communicator into a Smaller ones





© https://mpitutorial.com/tutorials/introduction-to-groupsand-communicators/

# How do I run it?



#### Compile it:

mpic++ -o MPI\_Hello.out MPI\_Hello.cpp

#### Run it:

 mpirun -hostfile machinefile.txt -np <np> MPI\_Hello.out

the command is implementation dependent [cesinihpc@hpc-200-06-18 esc22]\$ cat machinefile.txt hpc-200-06-18 slots=2

The same of running this

Add to the .bashrc the following two lines: module load compilers/gcc-12.2\_sl7 module load compilers/openmpi-4-1-4\_gcc12.2

Use option --mca btl\_openib\_allow\_ib 1 To suppress worning on IB usage (for openMPI4.0 and later)

mpirun --mca btl\_openib\_allow\_ib 1 -H hpc-200-06-18:2,hpc-200-06-17:2,hpc-200-06-06:2 -np 6 MPI\_Hello.out

cesi	ihpc@h	pc - 20	00-06-18 m	pi]\$ mpiru	inmca	btl_	openib_	allow_i	b 1	hostfile	<pre>machinefile.txt</pre>	-np 6	MPI_Hello.	out
lello	world	from	processor	hpc-200-(	)6-18.cr	.cnaf	.infn.i	t rank	0 01	f 6				
lello	world	from	processor	hpc-200-(	)6-18.cr	.cnaf	.infn.i	t rank	1 01	f 6				
lello	world	from	processor	hpc-200-(	)6-17.cr	.cnaf	.infn.i	t rank	2 01	f 6				
lello	world	from	processor	hpc-200-(	)6-17.cr	.cnaf	.infn.i	t rank	3 01	f 6				
lello	world	from	processor	hpc-200-(	)6-06.cr	.cnaf	.infn.i	t rank	4 01	f 6				
lello	world	from	processor	hpc - 200 - (	)6-06.cr	.cnaf	.infn.i	t rank	5 01	f 6				

hpc-200-06-17 slots=2

hpc-200-06-02 slots=2



- The executable must be present in all the hosts used, in the same path
  - You are lucky in the school nodes shared homes!!
  - OpenMPI in our cluster uses ssh to connect to the remote hosts
    - ssh should work passwordless (HostBasedAuthentication yes in sshd\_config)
    - During login the OpenMPI environment should be loaded
      - Typically via the .basrc file



#### Point-to-Point Communication





In general, in order to be able to communicate using messages you need to fill in a header and a payload

Synchronous send: If the sender waits for the message to be received

Asynchronous send returns immediately after the message has been sent

Receiving is usually synchronous

Messages have to match, otherwise deadlocks can occur

#### Messages – Send and Receive



int MPI\_Send (const void\* buf, int count, MPI\_Datatype datatype, int dest, int tag, MPI\_Comm comm)

int MPI\_Recv (void\* buf, int count, MPI\_Datatype datatype, int source, int tag, MPI\_Comm comm, MPI\_Status\* status)

- Int MPI\_Send performs a blocking send of the specified data ("count" copies of type "datatype," stored in "buf") to the specified destination (rank "dest" within communicator "comm"), with message ID "tag"
- int MPI\_Recv (void\* buf, int count, MPI\_Datatype datatype, int source, int tag, MPI\_Comm comm, MPI\_Status\* status)
- "blocking" means the functions return as soon as the buffer, "buf", can be safely used.

#### MPI Message Buffer

- In a perfect world, every send operation would be perfectly synchronized with its matching receive. This is rarely the case.
- •The MPI implementation must be able to deal with storing data when the two tasks are out of sync.
- Consider the following two cases:
  - A send operation occurs 5 seconds before the receive is ready - where is the message while the receive is pending?
  - Multiple sends arrive at the same receiving task which can only accept one send at a time - what happens to the messages that are "backing up"?
- The MPI implementation (not the MPI standard) decides what happens to data in these types of cases.
- Typically, a system buffer area is reserved to hold data in transit



- Opaque to the programmer and managed entirely by the MPI library
- A finite resource that can be easy to exhaust
- •Often mysterious and not well documented
- Able to exist on the sending side, the receiving side, or both
- Something that may improve program performance because it allows send - receive operations to be asynchronous





# Blocking vs Non-Blocking



Blocking:

- A blocking send routine will only "return" after it is safe to modify the application buffer (your sent data) for reuse.
- Safe means that modifications will not affect the data intended for the receive task.
- Safe does not imply that the data was actually received it may very well be sitting in a system buffer

#### Non-Blocking

- Non-blocking send and receive routines behave similarly they will return almost immediately.
- They do not wait for any communication events to complete, such as message copying from user memory to system buffer space or the actual arrival of message
- Non-blocking operations simply "request" the MPI library to perform the operation when it is able. The user can not predict when that will happen.
- It is unsafe to modify the application buffer (your variable space) until you know for a fact the requested non-blocking operation was actually performed by the library. There are "wait" routines used to do this.
- Non-blocking communications are primarily used to overlap computation with communication and exploit possible performance gains





•MPI guarantees that messages will not overtake each other.

If a sender sends two messages (Message 1 and Message 2) in succession to the same destination, and both match the same receive, the receive operation will receive Message 1 before Message 2.

If a receiver posts two receives (Receive 1 and Receive 2), in succession, and both are looking for the same message, Receive 1 will receive the message before Receive 2.

#### ESC22 - BERTINORO

•MPI does not guarantee fairness - it's up to the programmer to prevent "operation starvation".

 Example: task 0 sends a message to task 2.
 However, task 1 sends a competing message that matches task 2's receive. Only one of the sends will complete

The scenario requires that the receive used the wildcard MPI\_ANY\_SOURCE as its source argument.

07/10/2022





#### Fairness

### Non-blocking Send and Receive



int MPI\_ISend (const void \*buf, int count, MPI\_Datatype datatype, int dest, int tag, MPI\_Comm comm, MPI\_Request \*request)

int MPI\_IRecv (void \*buf, int count, MPI\_Datatype datatype, int source, int tag, MPI\_Comm comm, MPI\_Request \*request)

- int MPI\_Isend begins a non-blocking send of the variable buf to destination dest.
- Int MPI\_Irecv begins a non-blocking receive
- Since non-blocking operations may return before the requested system buffer space is obtained, the system issues a unique "request number".
  - The programmer uses this system assigned "handle" later (in a WAIT type routine) to determine completion of the non-blocking operation
    - MPI\_Wait( request, status )
    - MPI\_Test( request, flag, status )

#### Anywhere you use MPI\_Send or MPI\_Recv, you can use the pair of MPI\_Isend/MPI\_Wait or MPI\_Irecv/MPI\_Wait

# Send and Receive exercise – the PingPong



https://github.com/infn-esc/esc22/blob/main/hands-on/mpi/PingPong.cpp

```
// rank 0 sends to rank 1 and waits to receive a return message
  if (rank == 0) {
     dest = 1;
     source = 1;
    MPI Ssend(&outmsg, 1, MPI CHAR, dest, tag, MPI COMM WORLD);
     std::cout << "Rank 0 successfully sent a message to Rank 1: " << outmsg << std::endl;</pre>
    MPI Recv(&inmsg, 1, MPI CHAR, source, tag, MPI COMM WORLD, &Stat);
     std::cout << "Rank 0 successfully received a message from Rank 1: " << inmsg << std::endl;</pre>
// rank 1 waits for rank 0 message then returns a message
  else if (rank == 1) {
     std::cout << "Rank 1 is waiting for a message from Rank 0" << std::endl;</pre>
     source = 0;
     dest = 0;
    MPI Recv(&inmsg, 1, MPI CHAR, source, tag, MPI COMM WORLD, &Stat);
     std::cout << "Rank 1 successfully received a message from Rank 0: " << inmsg << std::endl;</pre>
     MPI Ssend(&outmsg, 1, MPI CHAR, dest, tag, MPI COMM WORLD);
     std::cout << "Rank 1 successfully sent a message to Rank 0: " << outmsg << std::endl;</pre>
```

### Change the Network interface



Following the so-called "Law of Least Astonishment", Open MPI assumes that if you have both an IP network and at least one high-speed network (such InfiniBand), you will likely only want to use the high-speed network(s) for MPI message passing

 Open MPI may still use TCP for setup and teardown information – so you'll see traffic across your IP network during startup and shutdown of your MPI job. This is normal and does not affect the MPI message passing channels.

mpirun --mca btl\_openib\_allow\_ib 1 -np 2 --hostfile machinefile.txt BandWidth.out

[cesinihpc@hpc-200-06-17 esc22]\$ cat machinefile.txt hpc-200-06-17 slots=1 hpc-200-06-18 slots=1

### Effect of changing the network interface



mpirun --mca btl\_openib\_allow\_ib 1 -np 2 --hostfile machinefile.txt BandWidth.out

[cesinihpc@hpc-200-06-17 mpi]\$ mpirun --mca btl\_openib\_allow\_ib 1 -np 2 --hostfile machinefile.txt BandWidth.out Arrays created and initialized Arrays created and initialized Rank 1 is waiting for a message from Rank 0 Rank 0 Received 200000000 INTs from rank 1 with tag 1 10 Iterations took 0.493944 seconds 8e+08 Bytes sent in 0.493944 seconds Bandwidth = 1.61962e+09 B/s = 12.9569 Gbit/s Rank 1 Received 20000000 INTs from rank 0 with tag 1

mpirun --mca btl tcp,self,vader --mca pml ob1 --mca btl\_tcp\_if\_include eth1 --hostfile machinefile.txt np 2 BandWidth.out

[cesinihpc@hpc-200-06-17 mpi]\$ mpirun --mca btl tcp,self,vader --mca pml ob1 --mca btl\_tcp\_if\_include eth1 --hostfile machinefile.txt -np 2 BandWidth.out Arrays created and initialized Rank 1 is waiting for a message from Rank 0 Rank 1 Received 40000000 INTs from rank 0 with tag 1 Rank 0 Received 40000000 INTs from rank 1 with tag 1 20 Iterations took 51.8564 seconds 3.2e+09 Bytes sent in 51.6564 seconds 3.2e+09 Bytes sent in 51.6564 seconds 3.andwidth = 6.17089e+07 B/s = 0.493671 Gbit/s

Now try using the SH (shared memory) MCA...any improvement?

# Non Blocking PingPong



https://github.com/infn-esc/esc22/blob/main/hands-on/mpi/NoBloc\_PingPong.cpp

```
if(my rank == 0)
        int value sent = 9999;
        MPI Request request;
        // Launch the non-blocking send
        MPI_Isend(&value_sent, 1, MPI_INT, 1, 0, MPI_COMM_WORLD, &request);
        std::cout << "MPI process " << my_rank << " I launched the non-blocking send." << std::endl;</pre>
        // Wait for the non-blocking send to complete
        MPI_Wait(&request, MPI_STATUS_IGNORE);
        std::cout << "MPI process " << my_rank << " The wait completed, so I sent value " << value_sent << std::endl;</pre>
else
        int value received = 0;
        MPI Request request;
        // Launch the non-blocking receive
        MPI Irecv(&value received, 1, MPI INT, 0, MPI ANY TAG, MPI COMM WORLD, &request);
        std::cout << "MPI process " << my rank << " I launched the non-blocking receive." << std::endl;</pre>
        // Wait for the non-blocking send to complete
        MPI_Wait(&request, MPI_STATUS_IGNORE);
        std::cout << "MPI process " << my_rank << " The wait completed, so I received value " << value_received << std::endl;</pre>
```

#### MPI\_Finalize();



#### Collective Communication





A message can be sent to/received from a group of processes

 Collective communication routines must involve all processes within the scope of a communicator

It is the programmer's responsibility to ensure that all processes within a communicator participate in any collective operations.

- Use collective communication when possible
  - They are implemented more efficiently than the sum of their point-to-point equivalent calls

# Types of Collective Operations



Synchronization - processes wait until all members c
 the group have reached the synchronization point.

- Data Movement
- Broadcast
- Scatter/gather
- All-to-All.
- Collective Computation (reductions)
  - one member of the group collects data from the other members and performs an operation on that data
    - Min
    - Max
    - Add
    - multiply



#### MPI\_Barrier

- MPI\_Barrier(MPI\_Comm Comm)
- Blocks until all processes have reached this routine
  - Blocks the caller until all group members have called it



				_	Vampir	<ul> <li>Timeline</li> </ul>	ir		
			flash_modified_ne	ew_00256.otf	(4:00.633 - 4:	00.794 = 0.	161 s)		
15 S.M.	4:00.6	64 4:00.66 4	:00,68 4:00,	.7 4:00	.72 4:00.	74 4:0	0.76 4:0	0,78	
Process 0	user	MPI_Barrier	81 MPI_Ssend		MPI_Waitall	US	r	1	Application
Process 1	user 1	<u>1PI_</u> Barrier	81 MPI_Ssend		MPI_Waitall	user			MPI
Process 2	user	MPI_Barnier	81 MPI_Ssend		MPI.	<u>Waitall</u> us	er		NoSendRecvReplac
Process 3	user	MPI_Bannien	81 MPI_Sser	nd	MPI_Waitall	us	er 🛛		
Process 4	user	MP1_Barrier	81Sse	end		131	user	10	
Process 5	user	MPI_Barrier	81 MPI_Sse	and	MPI	_Waitall	user		
Process 6	user	MPI_Barrier	81 MPI_Ssend	d		131	user		
Process 7	user	MP)_Bannier	81 MPI_Sse	nd			user		
Process 8	user	MPI_Barrier	81 MPI_Ssend			MPI_Waita	ll user		
Process 9	user	MPI_Barrier	81 MPI_Ssend		MPI_Waital1		user		
Process 10	user	MPI_Barrier	81 MPI_Ssend	1	MPI_Wa	aitall	user		
Process 11	user	MPI_Barrier	MPI_Irecv MP	'I_Ssend			131 user		
Process 12	user	MPI_Barrier	81 MPI_S	send			131 use	9 <b>5</b>	
Process 13	user	MPI_Barrier	81 107	MPI_Ssend	¢		1 1	iser	
Process 14	user	MPI_Barrier	81 MPI_Ssend	1		MPI	_Waitall us	er	
Process 15	user	MPI_Barrier	81 MPI_See	nd	č		us	er -	
Process 16	usen	MPI_Barrier	81 MPI_Ssen	d			MPI_Waitall u	ser	
Process 17	user M	PI_Barrier	81 MP1_Ssend	D.	MP.	I_Waitall		iser	
1				ſ					
	_	_	_	Bisplayed 18	3 from 256 bars	i Timeline –	25		
<u>x</u>	_			Displayed 18	3 from 256 bars Vampir - 1	imeline <	2>	_	
X)		f1	ash_lessbarrier_0	Displayed 18	3 from 256 bars Vampir - 1 tf (4:01.584 -	imeline < 4:01.744 =	2>. 0.161 s)	1.04 7	
	1000	f1 4:01.6 4:01.62	ash_lessbarnier_0 4:01.64	Displayed 18 5_x1_00256.o 4:01.66	3 from 258 bars Vampir - 1 tf (4:01.584 - 4:01.68	; fimeline < 4:01.744 = 4:01.7	2> 0,161 s) 4:01,72	4:01.74	Opolication
X)	Jiser	f1 4:01.6 4:01.62 31MPX_Socie	ash_lessbarrier_0 4:01.64 nd	Displayed 18 5_x1_00256.o 4:01.65 HPI_Waita	3 from 256 bars Vampir - T tf (4:01.584 - 4:01.68 all user	; fimeline < 4:01.744 = 4:01.7	2> 0.161 s) 4:01.72	4:01.74	Application
Process 0 Process 1	user user	f1 4:01.5 4:01.52 31MP1_Ssend 31.MP1_Ssend 01.MP1_Ssend	ash_lessbarrier_0 4:01.64 end	Displayed 18 5_x1_00256.o 4:01.66 HPI_Waitz NPI_Waitz	3 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user all user	; fimeline < 4:01.744 = 4:01.7	2> 0.161 s) 4:01.72	4:01.74	Application NPI
Process 0 Process 1 Process 2	user User User	f1 4:01.6 4:01.62 81MP2_See 81 MP1_Ssend 91 MP1_Ssend 91 MP1_Ssend	ash_lessbarrier_0 4:01.64 end MDT Uni	Displayed 18 5_x1_00256.o 4:01.66 NPI_Maits NPI_Maits 1 +=11	B         From 256 bars           Vampir - I           tf (4:01.584 -           4:01.68           all         user           all         user           all         user	; fimeline < 4:01.744 = 4:01.7	2> 0.161 s) 4:01.72	4:01.74	Application NPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4	user user user user	f1 4:01.6 4:01.62 81MP1_See 81 MPI_Seend 81 MPI_Seend 81 MPI_Seend 81 MPI_Seend	ash_lessbarrier_0 4:01.64 end end MPI_Wai	Displayed 18 5_x1_00256.o 4:01.66 NPI_Maits NPI_Maits 1 tall	3 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user 31 user 131	imeline < 4:01,744 = 4:01.7	2> 0.161 s) 4:01.72	4:01.74	Application NPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5	user user user user	f1 4:01.6 4:01.62 81MP'_Sse 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend	ash_lessbarnien_0 4:01.64 end end MPI_Wai NPI Seend	Displayed 18 5_x1_00256.o 4:01.66 HPI_Maits NPI_Maits 1 tall	3 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user 31 user user 131 [Maital]	imeline < 4:01,744 = 4:01.7	2> 0.161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 5	user User User User User	f1 4:01.6 4:01.62 31MP1_See 31 MP1_Ssend 31 MP1_Ssend 31 MP1_Ssend 31 MP1_Ssend 31 MP1_Ssend 31 MP1_Ssend 31 MP1_Ssend 31 MP1_Ssend	ash_lessbarrier_0 4:01.64 end end I_Ssend MPI_Ssend end	Displayed 18 5 x1 00256.0 4:01.66 HPI_Waits NPI_Maits 1 tall HPI	3 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user 131 User 131 [_Waital]	imeline < 4:01,744 = 4:01.7 user user	2> 0.161 s) 4:01.72	4:01.74	Application NPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 5 Process 7	usen Usen Usen Usen Usen	f1 4:01,6 4:01,62 31MP2_See 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MP 31 31 MP[_Sse 31 MPI_Sse	ash_lessbarrier_0 4:01.64 end I_Ssend MPI_Wai MPI_Ssend MPI_Ssend	Displayed 18 15_x1_00256.o 4:01.66 HPI_Waitz NPI_Waitz 1 tall	8 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user 131 User 131 131 131	imeline < 4:01.74 = 4:01.7 user user	2> 0,161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 7 Process 7	usen usen usen usen usen usen usen	f1 4:01.6 4:01.62 31MP3_See 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Sse 31 MPI_Sse 31 MPI_Sse 31 MPI_Sse	ash_lessbarrier_0 4:01.64 end MPI_Wai I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend nd	Displayed 18 15_x1_00256.o 4:01.65 HPI_Waitz NPI_Waitz 1; tall HPI	B from 256 bars           Vampir - 1           tf (4:01.584 -           4:01.68           all user           31 user           31 user           131 I.Maitall           131	imeline < 4:01.74 = 4:01.7 user user user user	2> 0,161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 6 Process 7 Process 8 Process 8	usen usen usen usen usen usen usen usen	f1 4:01.6 4:01.62 31MP7_See 31.MPI_Ssend 31.MPI_Ssend 31.MPI_Ssend 31.MPI_Ssend 31.MPI_Ssend 31.MPI_Ssend 31.MPI_Ssend	ash_lessbarrier_0 4:01.64 end MPI_Wai I_Ssend MPI_Ssend end MPI_Ssend nd	Displayed 18 15_x1_00256.0 4:01.68 NPI_Waitz NPI_Waitz 1: tall NPI NPI NPI NPI NPI NPI NPI NPI	B from 256 bars           Vampir - 1           tf (4:01.584 -           4:01.68           all user           all user           31 user           131 user           131           131           User           131           Waittall           1	imeline < 4:01.744 = 4:01.7 user user user user user	2> 0,161 s) 4:01,72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 5 Process 5 Process 6 Process 6 Process 8 Process 9 Process 10	user User User User User User User User	f1 4:01.6 4:01.62 31MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend	ash_lessbarrier_0 4:01.64 end MPI_blai I_Ssend MPI_Ssend MPI_Ssend nd Ssend	Displayed 18 15_x1_00256.c 4:01.68 NPI_Maitz NPI_Maitz 1: tall HPI HPI_Waita	B from 256 bars           Vampir - 1           tf (4:01.584 -           4:01.68           all user           all user           31 user           11 user           131 user           131 user           131 user           131 user           131 user           131           Waitall           1           MPI Naital	imeline < 4:01.744 = 4:01.7 user user user user user user	2> 0.161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 5 Process 5 Process 6 Process 7 Process 8 Process 9 Process 10	usen usen usen usen usen usen usen usen	4101.6 4101.62 81MP'_See 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssen 81 MPI_Ssen 81 MPI_Ssen 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend 81 MPI_Ssend	ash_lessbarrier_0 4:01.64 end end MPI_Wai I_Ssend MPI_Ssend end MPI_Ssend nd Ssend MPI_Ssend	Displayed 18 5_x1_00256.o 4:01.66 HPI_Maits NPI_Maits 1 tall HPI <u>HPI_Maital</u>	3 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user 131 User 131 [Waitall MPI_Waital 131 131	imeline < 4:01.744 4:01.7 user user user user user user user user	2> 0.161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 5 Process 6 Process 6 Process 7 Process 8 Process 9 Process 10 Process 11	usen usen usen usen usen usen usen usen	f1 4:01,6 4:01,62 81MP1_See 31 MP1_Ssend 31 MP1_Ssend	ash_lessbarrier_0 4:01.64 end I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend nd MPI_Ssend MPI_Ssend MPI_Ssend	Displayed 18 5 x1 00256.o 4:01.66 HPI_Waits NPI_Maits 1: tall HPI HPI_Waital	3 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user 11 user 131 User 131 [Waital] 1 MPI_Waital MPI_Waital MPI Vaital	imeline < 4:01.744 = 4:01.7 user user user user user user user user	2> 0,161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Arocess 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 6 Process 6 Process 7 Process 9 Process 10 Process 11 Process 12 Process 12	usen usen usen usen usen usen usen usen	f1 4:01,6 4:01,62 31MP3_Soc 31 MP1_Ssend 31 MP1_Ssend	ash_lessbarrier_0 4:01.64 end I_Ssend MPI_Wai NPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend	Displayed 18 5_x1_00256.o 4:01.66 HPI_Waitz NPI_Waitz tall HPI_Waital	8 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user 31 user 131 1. Waitall 1 MP[_Waitall MP[_Waitall 131 MP[_Waitall	imeline < 4:01.74 = 4:01.7 user user user user user user user user	2> 0,161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 6 Process 7 Process 7 Process 9 Process 10 Process 11 Process 12 Process 13 Process 13	usen usen usen usen usen usen usen usen	f1 4:01,6 31MP2_See 31MP2_See 31 MPI_Ssend 31 MPI_Ssend 31 MPI_See 31 31 MPI_See 31 MPI_	ash_lessbarrier_0 4:01.64 end MPI_Wai I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend	Displayed 18 15_x1_00256.o 4:01.66 MPI_Waitz NPI_Waitz tall MPI_Waital	8 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user 131 user 131 131 14 MPI Waitall MPI Waitall MPI Waitall	imeline < 4:01.74 = 4:01.7 user user user user user user user user	2> 0,161 c) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 4 Process 5 Process 7 Process 7 Process 9 Process 10 Process 11 Process 12 Process 13 Process 15	usen usen usen usen usen usen usen usen	f1 4:01.6 4:01.62 31MP3_Soc 31 MPI_Ssend 31 MPI_Ssend	ash_lessbarrier_0 4:01.64 end NPI_Wai I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend	Displayed 18 15_x1_00256.o 4:01.65 HPI_Waitz NPI_Maitz 1; tall HPI_Maital	8 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user 31 user 131 131 Maitall MPI_Waitall MPI_Waitall MPI_Waitall	imeline < 4:01.74 = 4:01.7 user user user user user user user user	2> 0.161 s) 4:01.72	<b>4:01.74</b>	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 5 Process 5 Process 6 Process 7 Process 10 Process 10 Process 11 Process 11 Process 14 Process 15 Process 14 Process 15 Process 16	user user user user user user user user	f1 4:01,6 4:01,62 31MP/_See 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_See 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend 31 MPI_Ssend	ash_lessbarrier_0 4:01.64 end MPI_Wai I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend Ssend send	Displayed 18 5_x1_00256.0 4:01.68 HPI_Maitz NPI_Maitz 1: tall HPI_Maital	8 from 256 bars Vampir - 1 tf (4:01.584 - 4:01.68 all user all user 131 user 131 Maitall MPI_Waitall MPI_Waitall MPI_Waitall	imeline < 4:01.744 = 4:01.7 user user user user user user user user	2> 0,161 s) 4:01,72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 5 Process 5 Process 6 Process 7 Process 9 Process 10 Process 10 Process 11 Process 11 Process 13 Process 14 Process 16 Process 16 Process 17	usen usen usen usen usen usen usen usen	41 4:01.6 31MP/_See	ash_lessbarrier_0 4:01.64 end MPI_Wai I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend Ssend send	Displayed 18 5_x1_00256.0 4:01.68 HPI_Waitz NPI_Waitz 1: tall HPI_Waital	8 from 256 bars Vampir - T tf (4:01.584 - 4:01.68 all user 31 user 131 user 131 Meitall MPI_Waitall MPI_Waitall MPI_Waitall MPI_Waitall	imeline < 4:01.744 = 4:01.7 user user user user user user user user	2> 0.161 s) 4:01.72	4:01.74	Application MPI NoSendRecvReplac
Process 0 Process 1 Process 2 Process 3 Process 5 Process 6 Process 6 Process 7 Process 10 Process 10 Process 11 Process 11 Process 13 Process 14 Process 16 Process 16 Process 17	usen usen usen usen usen usen usen usen	41 4:01.6 4:01.62 81MP/_See 8	ash_lessbarrier_0 4:01.64 end MPI_Wai I_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend MPI_Ssend Ssend ssend	Displayed 18 5_x1_00256.0 4:01.68 HPI_Maitz NPI_Maitz 1: tall HPI_Maital HPI_Maital	8 from 256 bars Vampir - T tf (4:01.584 - 4:01.68 all user all user 131 user 131 Maitall MPI_Waitall MPI_Waitall MPI_Waitall MPI_Waitall	imeline < 4:01.744 = 4:01.7 user user user user user user user 10.5 user 13:0.05 user 13:05 user	2> 0,161 s) 4:01,72	4:01.74	Application MPI NoSendRecvReplac

An MPI Barrier call before a communication phase ensures a synchronized start of the communication calls (top). When removing the barrier there is an un-synchronized start (bottom)

#### Broadcast

Int MPI\_Bcast\_c(void \*buffer, MPI\_Count count, MPI\_Datatype datatype, int root, MPI\_Comm comm)

- broadcasts a message from the process with rank root to all processes of the group, itself included.
- It is called by all members of the group using the same arguments for comm and root.
- On return, the content of root's buffer is copied to all other processes.



#### MPI\_BCAST(buffer, count, datatype, root, comm)

INOUT

IN

IN

IN

IN





#### Gather

- int MPI\_Gather(const void \*sendbuf, int sendcount, MPI\_Datatype sendtype, void \*recvbuf, int recvcount, MPI\_Datatype recvtype, int root, MPI\_Comm comm)
  - each process (root process included) sends the contents of its send buffer to the root process.
  - The root process receives the messages and stores them in rank order
  - The receive buffer is ignored for all non-root processes
  - Note that the recvcount argument at the root indicates the number of items it receives from each process, not the total number of items it receives



MPI\_GATHER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)

IN	sendbuf	starting address of send buffer (choice)
IN	sendcount	number of elements in send buffer (non-negative integer)
IN	sendtype	datatype of send buffer elements (handle)
OUT	recvbuf	address of receive buffer (choice, significant only at root)
IN	recvcount	number of elements for any single receive (non-negative integer, significant only at root)
IN	recvtype	datatype of recv buffer elements (handle, significant only at root)
IN	root	rank of receiving process (integer)
IN	comm	communicator (handle)



#### Scatter

int MPI\_Scatter(const void \*sendbuf, int sendcount, MPI\_Datatype sendtype, void \*recvbuf, int recvcount, MPI\_Datatype recvtype, int root, MPI\_Comm comm)

the root sends a message with

MPI\_Send(sendbuf,sendcountn, sendtype,...). This message is split into n equal segments, the i-th segment is sent to the i-th process in the group, and each process receives this message as above.

The send buffer is ignored for all non-root processes



#### scatter

IN	sendbuf	address of send buffer (choice, significant only at root)
IN	sendcount	number of elements sent to each process (non-negative integer, significant only at root)
IN	sendtype	data type of send buffer elements (handle, significant only at root)
OUT	recvbuf	address of receive buffer (choice)
IN	recvcount	number of elements in receive buffer (non-negative integer)
IN	recvtype	datatype of receive buffer elements (handle)
IN	root	rank of sending process (integer)
IN	comm	communicator (handle)

MPI\_SCATTER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)

IN

IN

IN

IN

IN



#### ESC22 - BERTINORO

#### Reduce

 int MPI\_Reduce(const void \*sendbuf, void \*recvbuf, int count, MPI\_Datatype datatype, MPI\_Op op, int root, MPI\_Comm comm)

- combines the elements provided in the input buffer of each process in the group, using the operation op, and returns the combined value in the output buffer of the process with rank root.
- The input buffer is defined by the arguments sendbuf, count and datatype; the output bu er is defined by the arguments recvbuf, count and datatype;



IN	sendbuf	address of send buffer (choice)
OUT	recvbuf	address of receive buffer (choice, significant only at root)
IN	count	number of elements in send buffer (non-negative integer)
IN	datatype	datatype of elements of send buffer (handle)
IN	ор	reduce operation (handle)
IN	root	rank of root process (integer)
IN	comm	communicator (handle)





#### **Reduce** Operations

Esc

- MPI\_MAX Returns the maximum element
- MPI\_MIN Returns the minimum element
- •MPI\_SUM Sums the elements.
- •MPI\_PROD Multiplies all elements.
- •MPI\_LAND Performs a logical and across the elements
- MPI\_LOR Performs a logical or across the elements
- MPI\_BAND Performs a bitwise and across the bits of the elements
- MPI\_BOR Performs a bitwise or across the bits of the elements
- •MPI\_MAXLOC Returns the maximum value and the rank of the process that owns it
- •MPI\_MINLOC Returns the minimum value and the rank of the process that owns it

#### Other Collective Opertaions

 MPI\_ALLGATHER can be thought of as MPI\_GATHER, but where all processes receive the result, instead of just the root

 MPI\_ALLTOALL is an extension of MPI\_ALLGATHER to the case where each process sends distinct data to each of the receivers. The j-th block sent from process i is received by process j and is placed in the i-th block of recvbuf.









## The MPI\_Pi





### Run-time Tuning: Process Affinity



- •Open MPI supports processor affinity on a variety of systems through process binding
  - Each MPI process is "bound" to a specific subset of processing resources (cores, sockets, L\* cache, hwthread etc.).
  - The operating system will constrain that process to run on only that subset
- •Affinity can improve performance by inhibiting excessive process movement
  - for example, away from "hot" caches or NUMA memory.
- Judicious bindings can improve performance
  - by reducing resource contention (by spreading processes apart from one another)
  - improving interprocess communications (by placing processes close to one another).
- Binding can also improve performance reproducibility by eliminating variable process placement.
- •Unfortunately, binding can also degrade performance by inhibiting the OS capability to balance loads.
- Depending on how processing units on your node are numbered, the binding pattern may be good, bad, or even disastrous
  - If you want to control affinity you have to know what you are doing

# Mapping, Ranking, and Binding: Oh My!



- •Open MPI employs a three-phase procedure for assigning process locations and ranks:
  - Mapping
    - Assigns a default location to each process
  - Ranking
    - Assigns an MPI\_COMM\_WORLD rank value to each process
  - Binding
    - Constrains each process to run on specific processors
- •To control process mapping in the command line:
  - --map-by <foo>
  - Map to the specified object, defaults to socket.

Often a good choice is to let MPI decide for you. But if you want to master the MPI mapping, the mpirun manual is a good starting point: https://www.openmpi.org/doc/v4.1/man1/mpirun.1.php

<foo> can be: slot, hwthread, core, L1cache, L2cache, L3cache, socket, numa, board, node, sequential, distance, and ppr.



- In Open-MPI mpirun automatically binds processes as of the start of the v1.8 series
  - Two binding patterns are used in the absence of any further directives:
    - Bind to core: when the number of processes is <= 2</p>
    - Bind to socket: when the number of processes is > 2
- To control process binding in the command line:
  - --bind-to <foo>:
    - Bind processes to the specified object, defaults to core.
    - Supported options include slot, hwthread, core, l1cache, l2cache, l3cache, socket, numa, board, and none.
  - -report-bindings, --report-bindings: Report any bindings for launched processes.

# Fine binding: The rankfile



#### -rf, --rankfile <rankfile>

- Provide a rankfile file for fine control of the process allocation
- rank <N>=<hostname> slot=<slot list>

For example:

\$ cat myrankfile

rank 0=aa slot=1:0-2Rank 0 runs on node aa, bound to logical socket 1, cores 0-2.rank 1=bb slot=0:0,1Rank 1 runs on node bb, bound to logical socket 0, cores 0 and 1.rank 2=cc slot=1-2Rank 2 runs on node cc, bound to logical cores 1 and 2.





[cesinihpc@hpc-200-06-17 mpi]\$ time mpirun --mca btl\_openib\_allow\_ib 1 --map-by core --bind-to core --report-bindings -np 16 MPI\_Pi.o hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 1 bound to socket 0[core 1[hwt 0-1]]: [../BB/../../../ hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 2 bound to socket 0[core 2[hwt 0-1]]: [hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]: hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 4 bound to socket 0[core 4[hwt 0-1]]: hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 5 bound to socket 0[core 5[hwt 0-1]]: [hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 6 bound to socket 0[core 6[hwt 0-1]]: [hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 7 bound to socket 0[core 7[hwt 0-1]]: [hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 8 bound to socket 1[core 8[hwt 0-1]]: [hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 9 bound to socket 1[core 9[hwt 0-1]]: [../ hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 10 bound to socket 1[core 10[hwt 0-1]]: hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 11 bound to socket 1[core 11[hwt 0-1]]: hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 12 bound to socket 1[core 12[hwt 0-1]]: hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 13 bound to socket 1[core 13[hwt 0-1]]: [hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 14 bound to socket 1[core 14[hwt 0-1]]: [../../../ hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 15 bound to socket 1[core 15[hwt 0-1]]: [../../../../../../../../../BB] Integrating Pi with numsteps = 40000000000. Step = 2.5e-11. Numsteps per process = 2500000000.

result:	3.14159265358959
real	0m16.133s
user	4m6.876s
sys	0m4.320s

#### A disastrous binding example



cesinihpc@hpc-200-06-17 mpi]\$ time mpirun --mca btl openib allow ib 1 --map-by hwthread --bind-to hwthread --report-bindings -np 16 MPI Pi.o [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 1 bound to socket 0[core 0[hwt 1]]: [.B/../../../../../../.. hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 2 bound to socket 0[core 1[hwt 0]]: [../B./ [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 3 bound to socket 0[core 1[hwt 1]]: [../.B/.. [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 4 bound to socket 0[core 2[hwt 0]]: [../../B./ [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 5 bound to socket 0[core 2[hwt 1]]: [../../.B/.. hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 6 bound to socket 0[core 3[hwt 0]]: [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 7 bound to socket 0[core 3[hwt 1]]: [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 8 bound to socket 0[core 4[hwt 0]]: [../../ [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 9 bound to socket 0[core 4[hwt 1]]: [../... ′../.B/.. hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 10 bound to socket 0[core 5[hwt 0]]: [../ [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 11 bound to socket 0[core 5[hwt 1]]: [../.. [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 12 bound to socket 0[core 6[hwt 0]]: [../../../../../ [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 13 bound to socket 0[core 6[hwt 1]]: [../../ hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 14 bound to socket 0[core 7[hwt 0]]: [../... [hpc-200-06-17.cr.cnaf.infn.it:05178] MCW rank 15 bound to socket 0[core 7[hwt 1]]: [../../../../../../../..B][../../../../../../../ Integrating Pi with numsteps = 40000000000. Step = 2.5e-11. lumsteps per process = 2500000000.

result:	3.14159265358959
real	0m31.250s
user	8m11.443s
sys	0m2.772s

#### Run-time tuning: Memory Affinity



•Open MPI supports general and specific memory affinity,

- •it generally tries to allocate all memory local to the processor that asked for it.
- •When shared memory is used for communication, Open MPI uses memory affinity to make certain pages local to specific processes in order to minimize memory network/bus traffic.

#### Homework Exercise

Matrix transpose

- <u>https://www.hpc.cineca.it/content/exercise-15</u>
- Solution: <u>https://www.hpc.cineca.it/content/solution-15</u>

#### Matrix Multiplication

- <u>https://www.hpc.cineca.it/content/exercise-16</u>
- <u>https://www.hpc.cineca.it/content/solution-16</u>
- 2D Laplace Equation
  - https://www.mcs.anl.gov/research/projects/mpi/tutorial/mpiexmpl/src/jacobi/C/main.html

