

Application speed- up: an example with NAMD

FELICE PANTALEO - CERN

DANIELE CESINI - INFN-CNAF



Speedup of a real application

NAMD



- A parallel molecular dynamics code designed for high-performance simulation of large biomolecular systems
- Can be download here: <https://www.ks.uiuc.edu/Research/namd/>
- Installing...
 - https://github.com/inf-nesc/esc22/blob/main/hands-on/mpi/Make_and_Install_NAMD.sh
- Running on real molecule...
 - https://github.com/inf-nesc/esc22/blob/main/hands-on/mpi/HowTo_Launch_NAMD_on_APOA1.sh
- GPU Porting available on CUDA



Speedup of an application

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

$$S(P) = P$$



Parallel computation efficiency

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

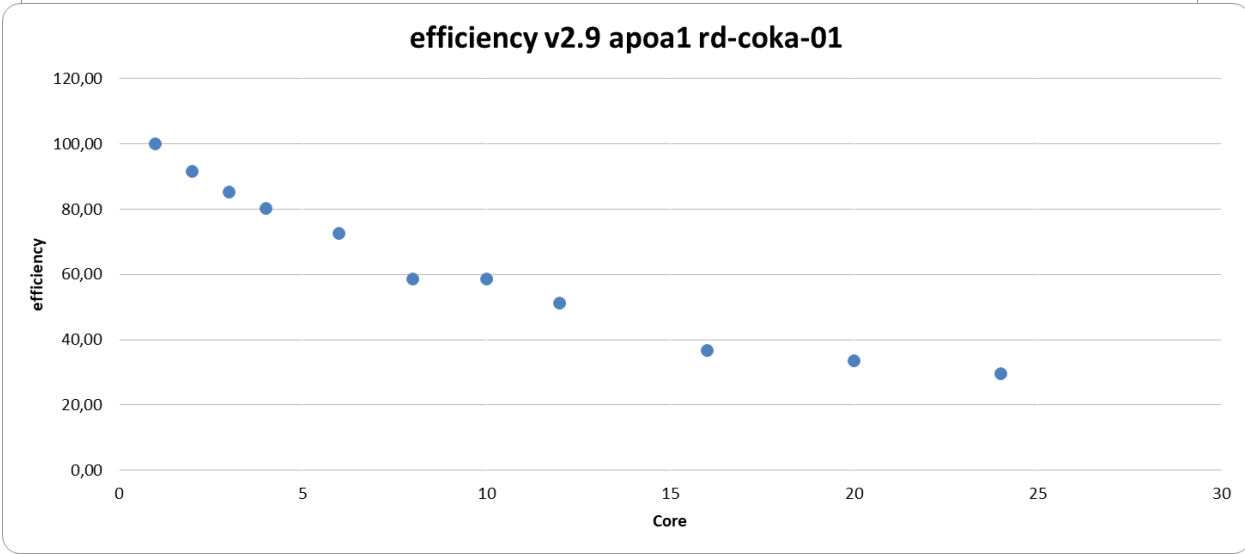
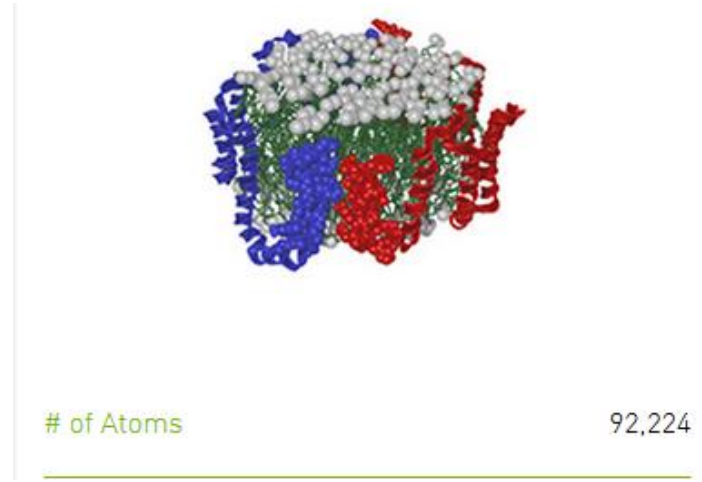
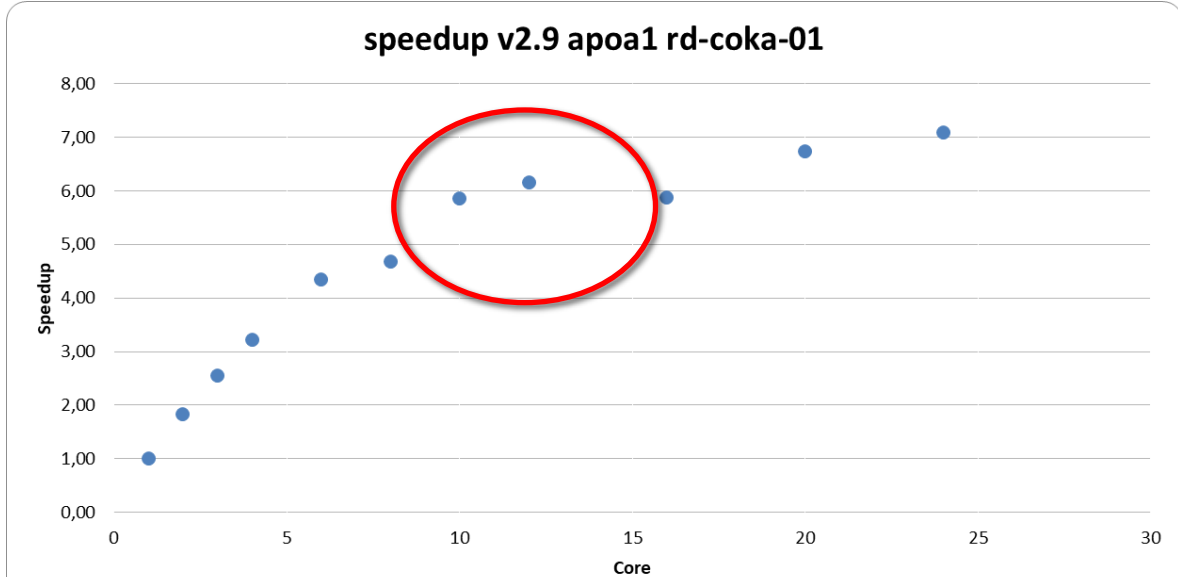
$$\varepsilon = \frac{T_{Seq}}{P * T_{Par}(P)} = \frac{S(P)}{P}$$

T = Elapsed Time

P = Number of processors Used



Speedup examples – NAMD APOA1

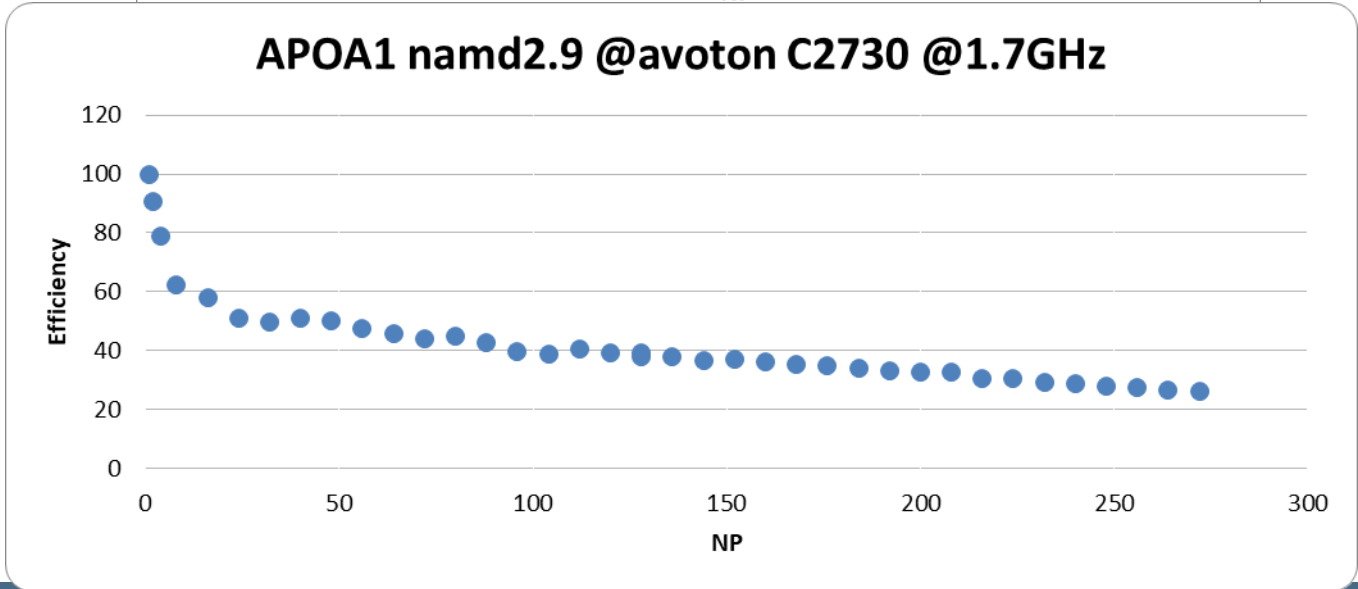
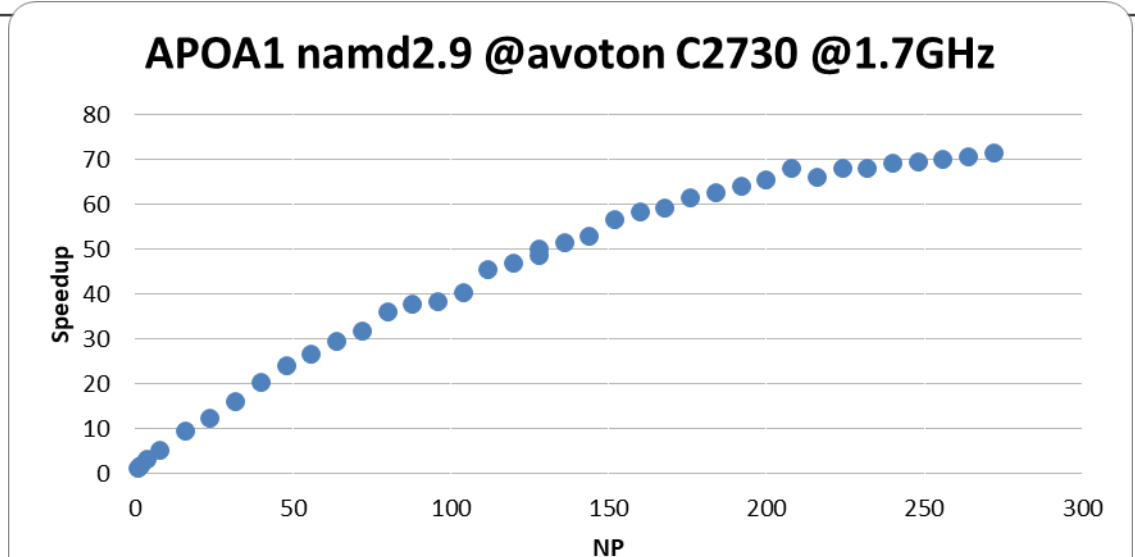


APOA1

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.



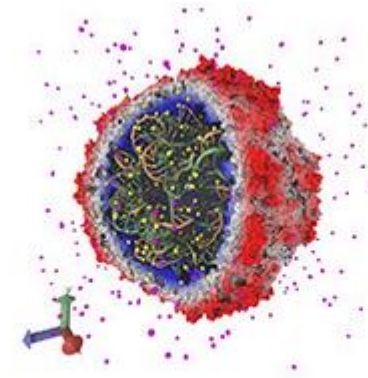
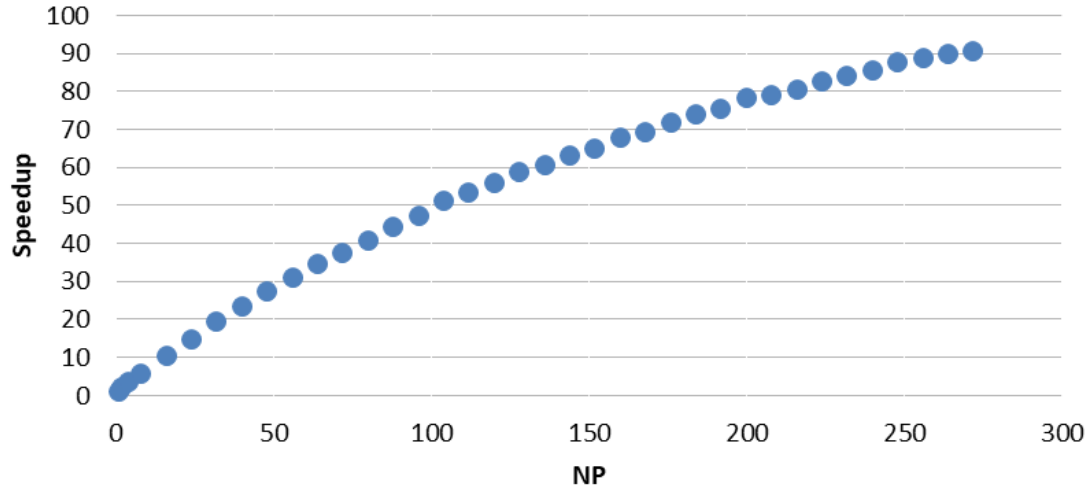
Speedup example – NAMD APOA1





Speedup – NAMD STMV

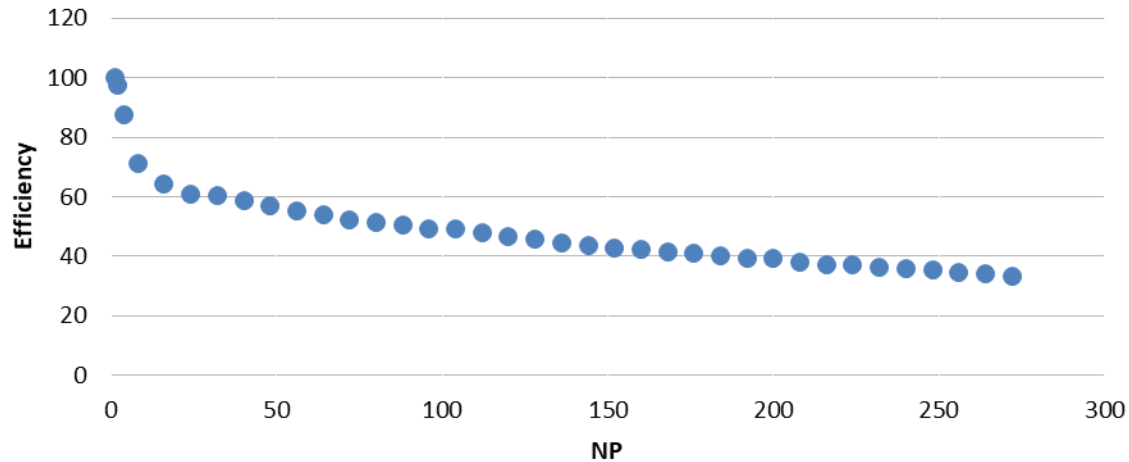
STMV namd2.9 @avoton C2730 @1.7GHz



of Atoms

1,066,628

STMV namd2.9 @avoton C2730 @1.7GHz



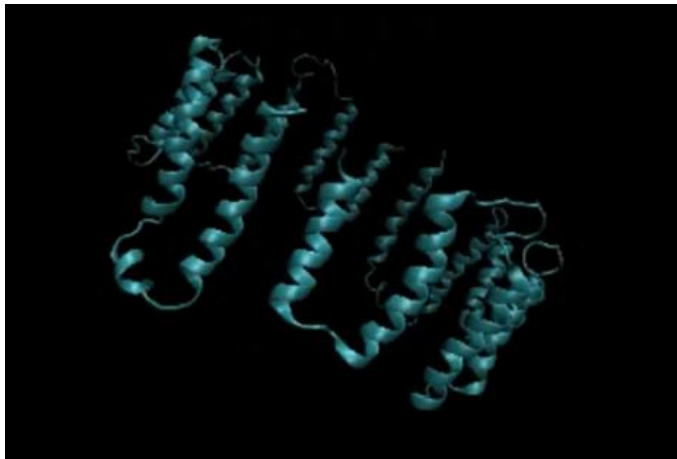
STMV

Satellite Tobacco Mosaic Virus (STMV) is a small, icosahedral plant virus that worsens the symptoms of infection by Tobacco Mosaic Virus (TMV). STMV 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.



Molecular Dynamics Simulation on three proteins

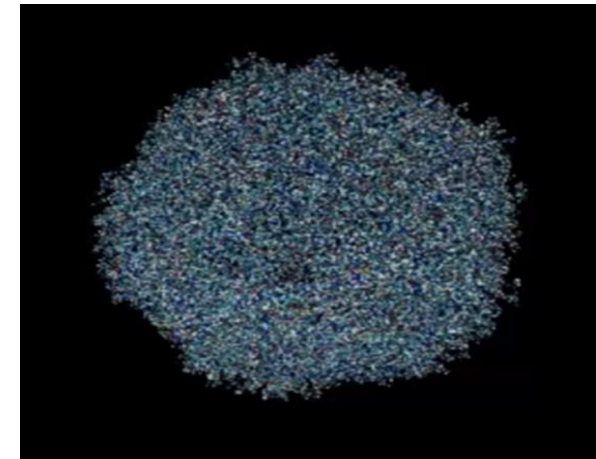
APOAI (92,224 atoms)
major protein component of high-density lipoprotein (HDL) in the bloodstream and plays a specific role in lipid metabolism.
standard NAMD cross-platform benchmark for years

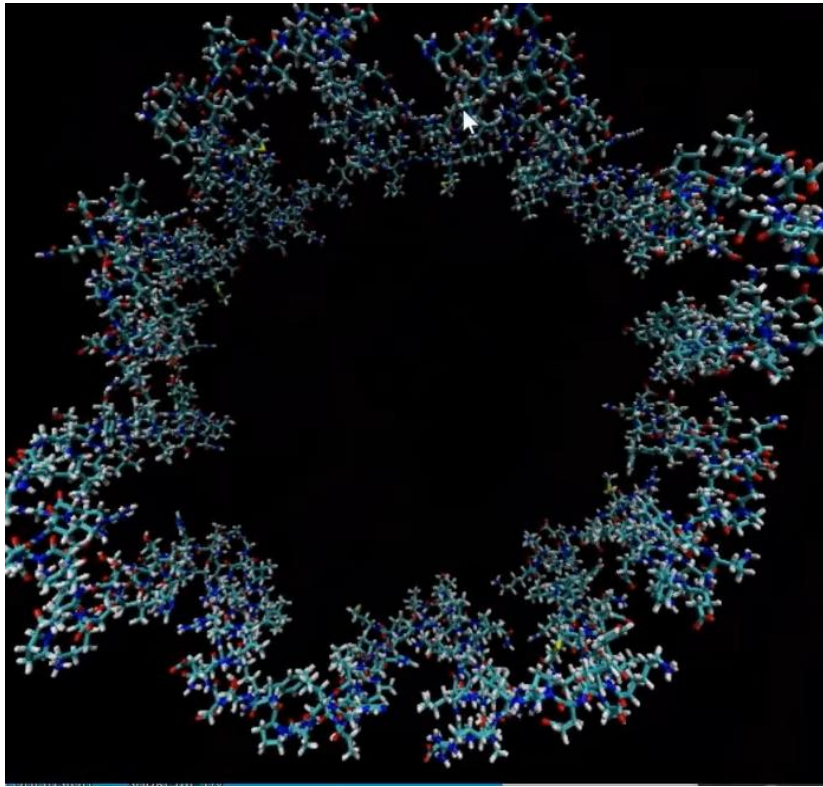


F1ATPASE (327,506 atoms):
enzyme that synthesizes adenosine triphosphate (ATP), the common molecular energy unit in cells.
The F1-ATPase benchmark is a model of the F1 subunit of ATP synthase.



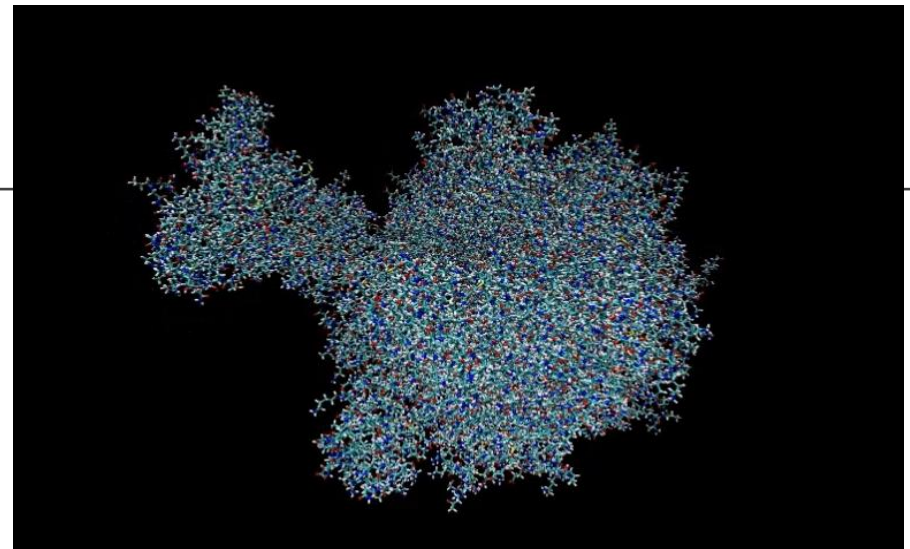
STMV(1,066,628 atoms):
is a small, icosahedral plant virus that worsens the symptoms of infection by Tobacco Mosaic Virus (TMV)



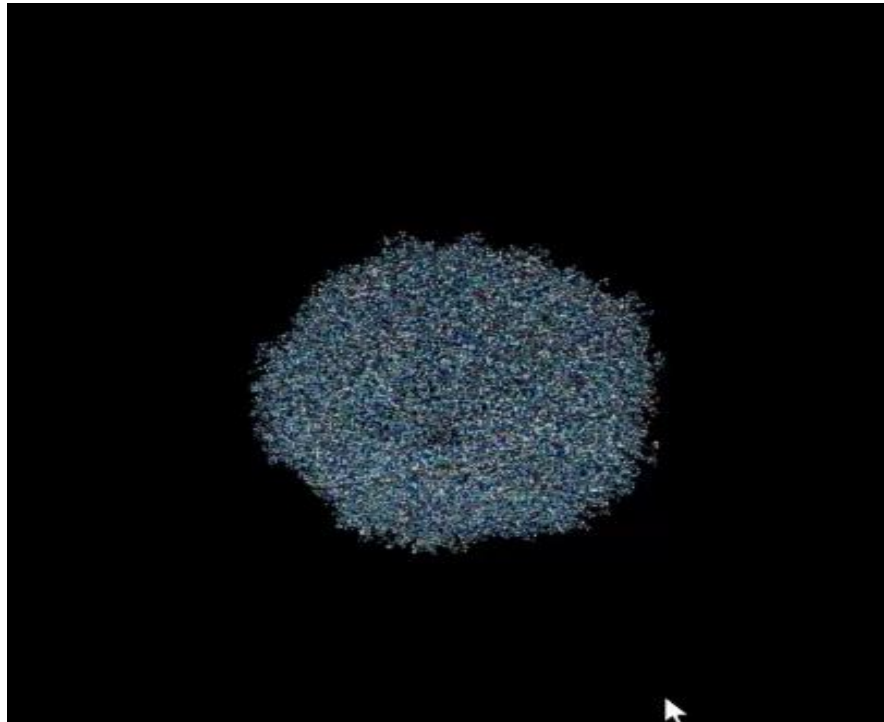


APOA1

VMD OUTPUT



F1ATPASE



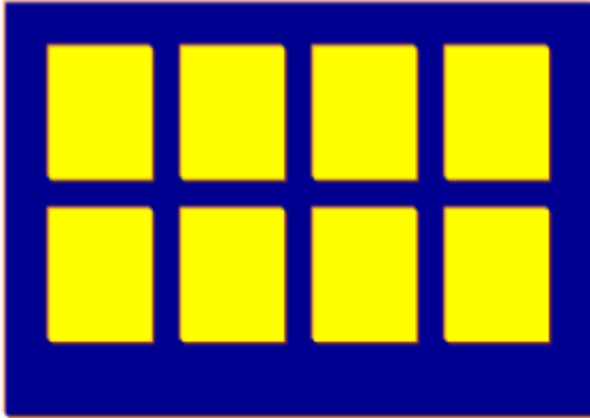
STMV



Let's add dimensions....accelerators

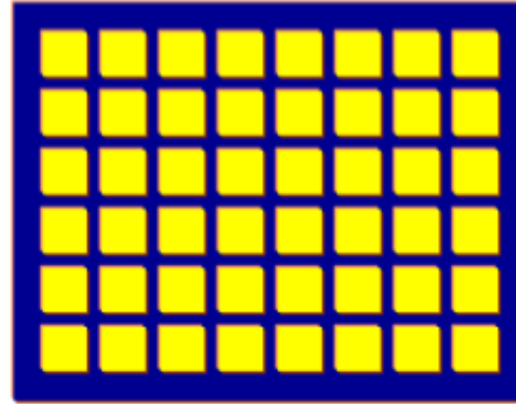
Reasonable Speedups - 1

• Chip A



$$Perf_A = Eff_A * Peak_A(Comp \text{ or } BW)$$

• Chip B



$$Perf_B = Eff_B * Peak_B(Comp \text{ or } BW)$$

$$Speedup \frac{B}{A} = \frac{Perf_B}{Perf_A} = \frac{Eff_B}{Eff_A} * \frac{Peak_A(Comp \text{ or } BW)}{Peak_B(Comp \text{ or } BW)}$$

© 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



Apply it to our system

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

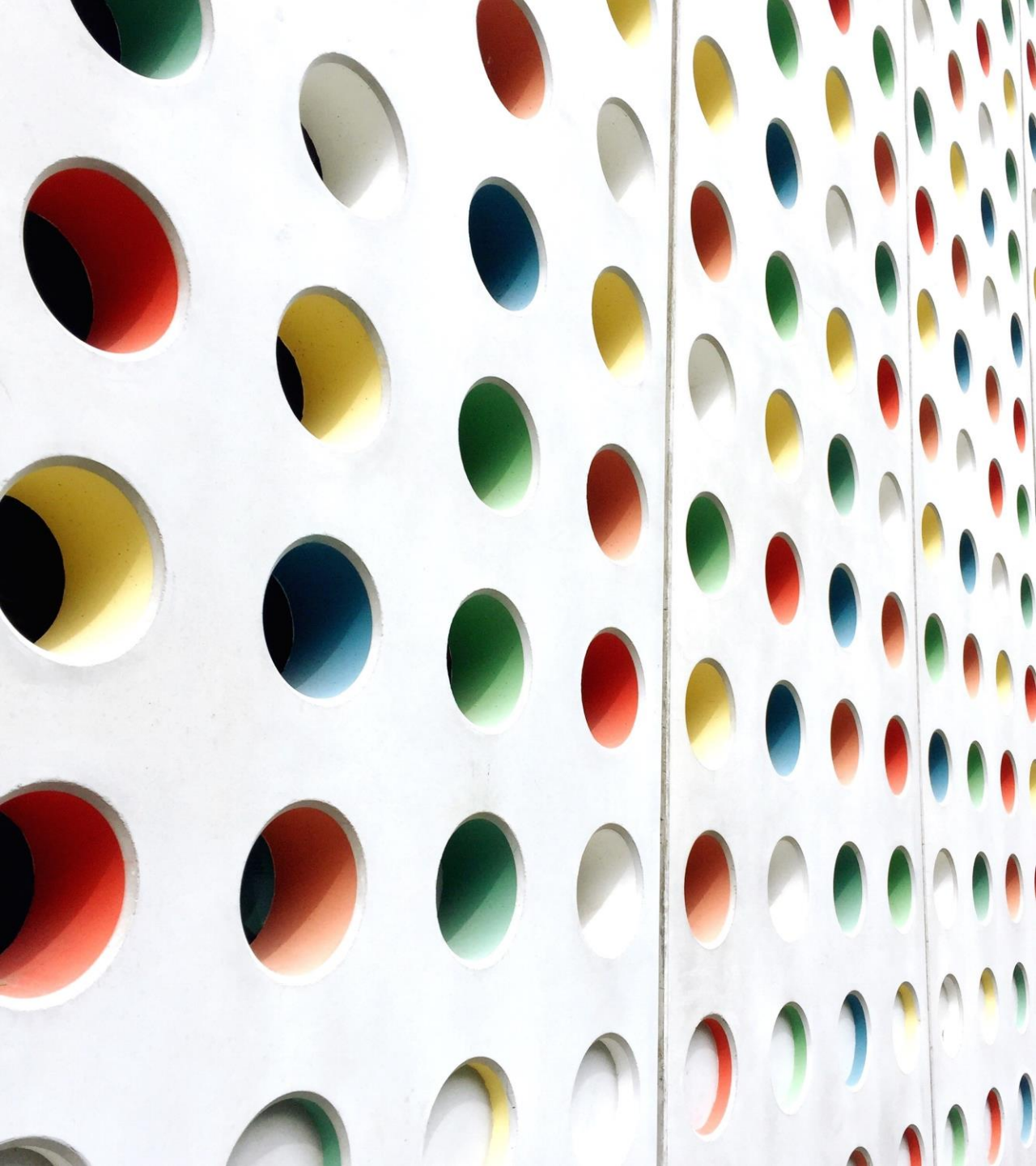
	Max Speed-Up CPU/GPU
Compute Bound App (sp)	$14.13/1.5 = 9.5$
Bandwidth Bound App	$572/240 = 2.4$

Common Mistakes in Comparing CPU and GPU performances



- 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: https://www.open-mpi.org/doc/v3.0/man3/MPI_Wtime.3.php
 - <https://www.open-mpi.org/faq/>
- MPICH.org: <https://www.mpich.org/>
- MPI Tutorial: <https://mpitutorial.com/>
- 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

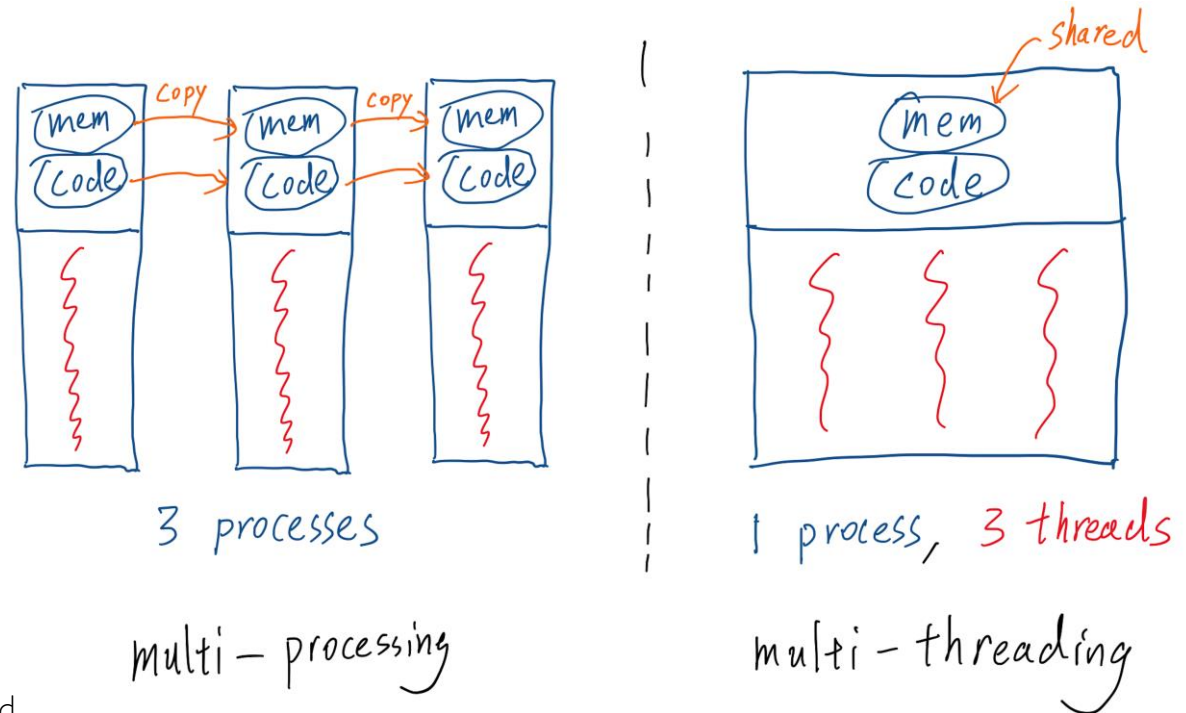
(Credits to Tim Mattson at IntelLab for his “**Hands-on** Introduction to MPI” at ESC15)

Multithread vs Multiprocess

- 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>

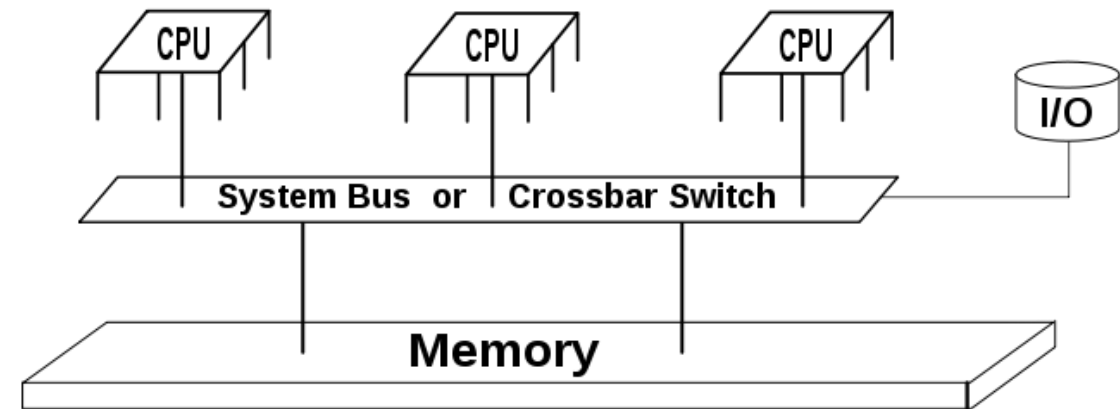
- Inter-process communication is slower due to isolated memory





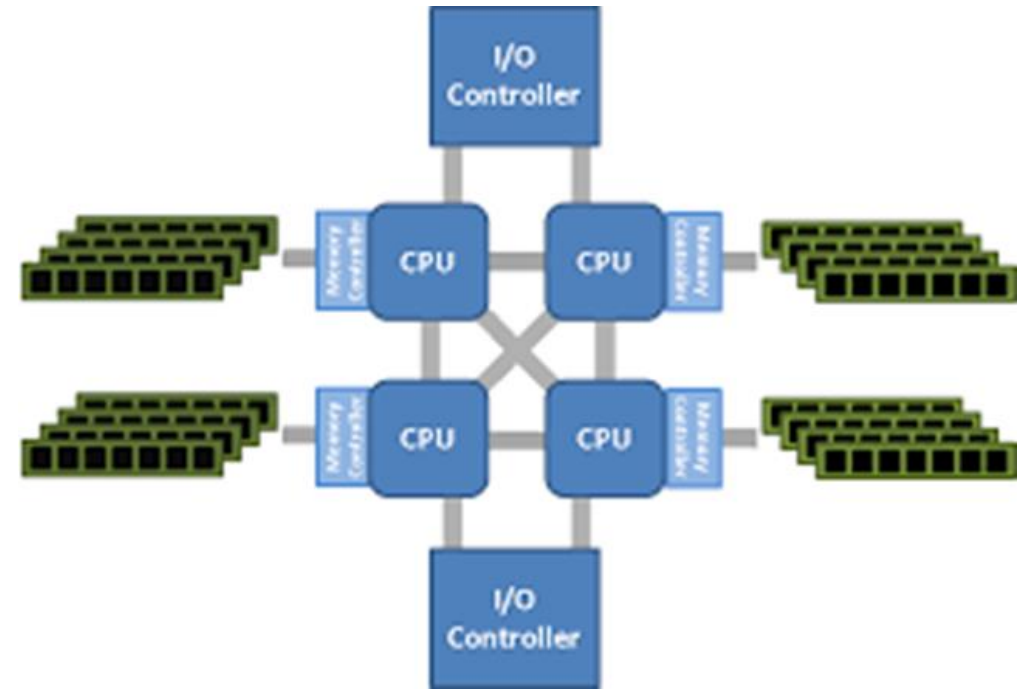
Shared Memory Systems

- 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

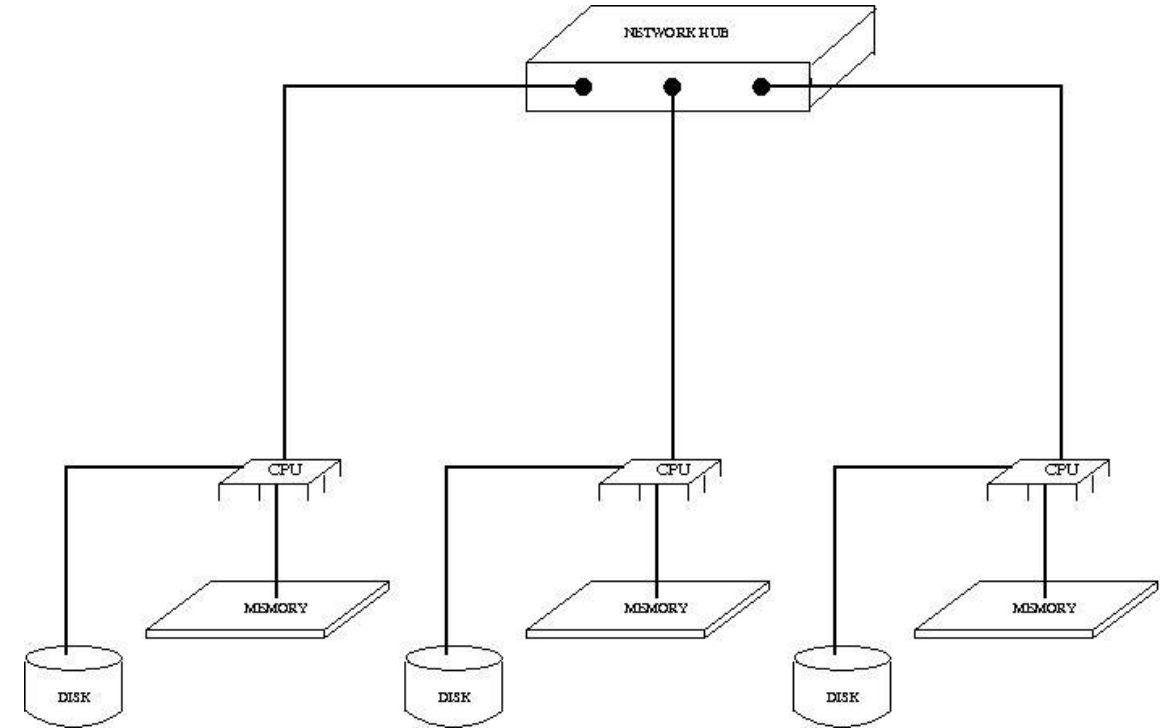
- A programmer can set an allocation policy for its program using a component of NUMA API called **libnuma**.
 - a user space shared library that can be linked to applications
 - provides explicit control of allocation policies to user programs.
- The NUMA execution environment for a process can also 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>

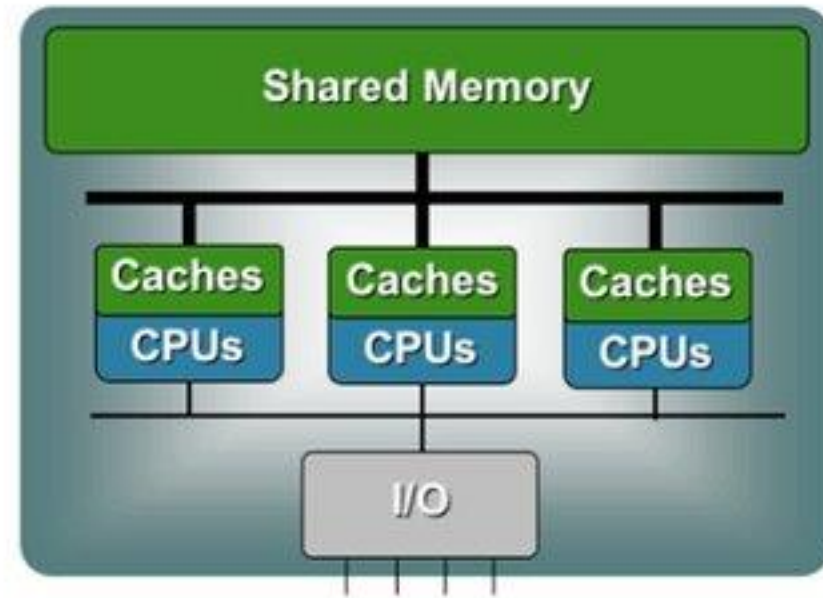
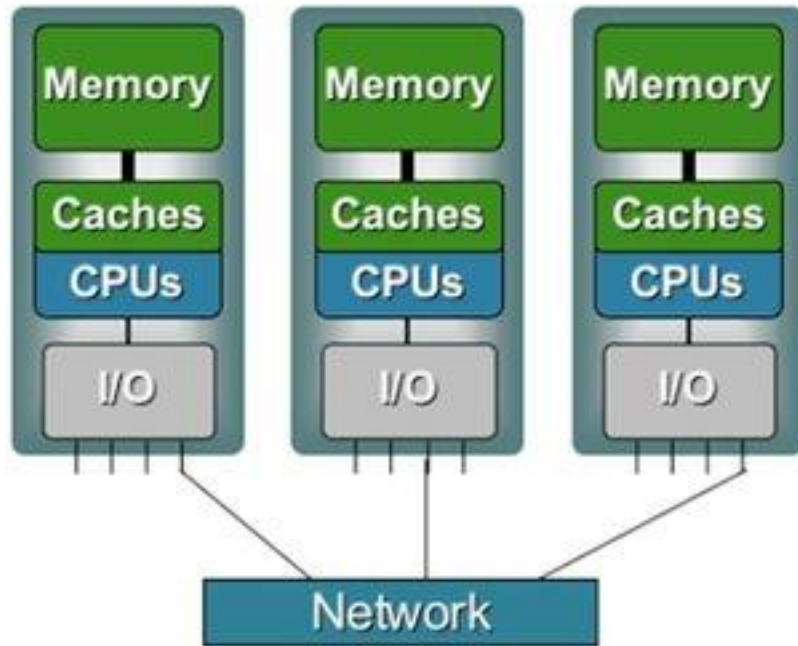
Distributed Memory Systems

- 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

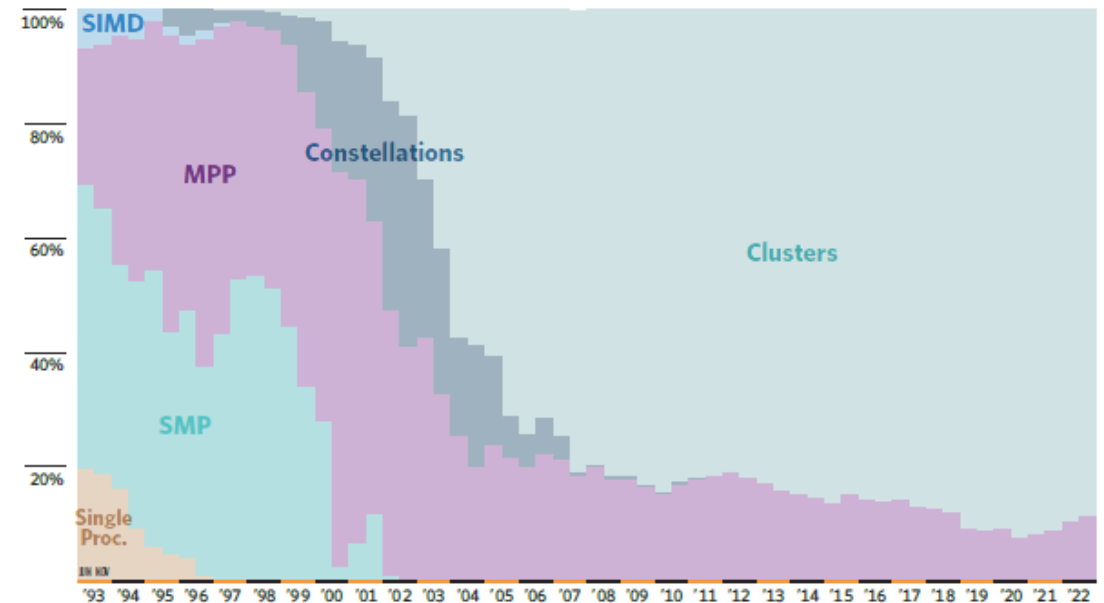
[a cluster is a] parallel computer system comprising an integrated collection of independent nodes, each of which is a system in its own right, capable of independent operation and derived from products developed and marketed for other stand-alone purposes

© 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

ARCHITECTURES Top500.org 2023 stats





System Topology

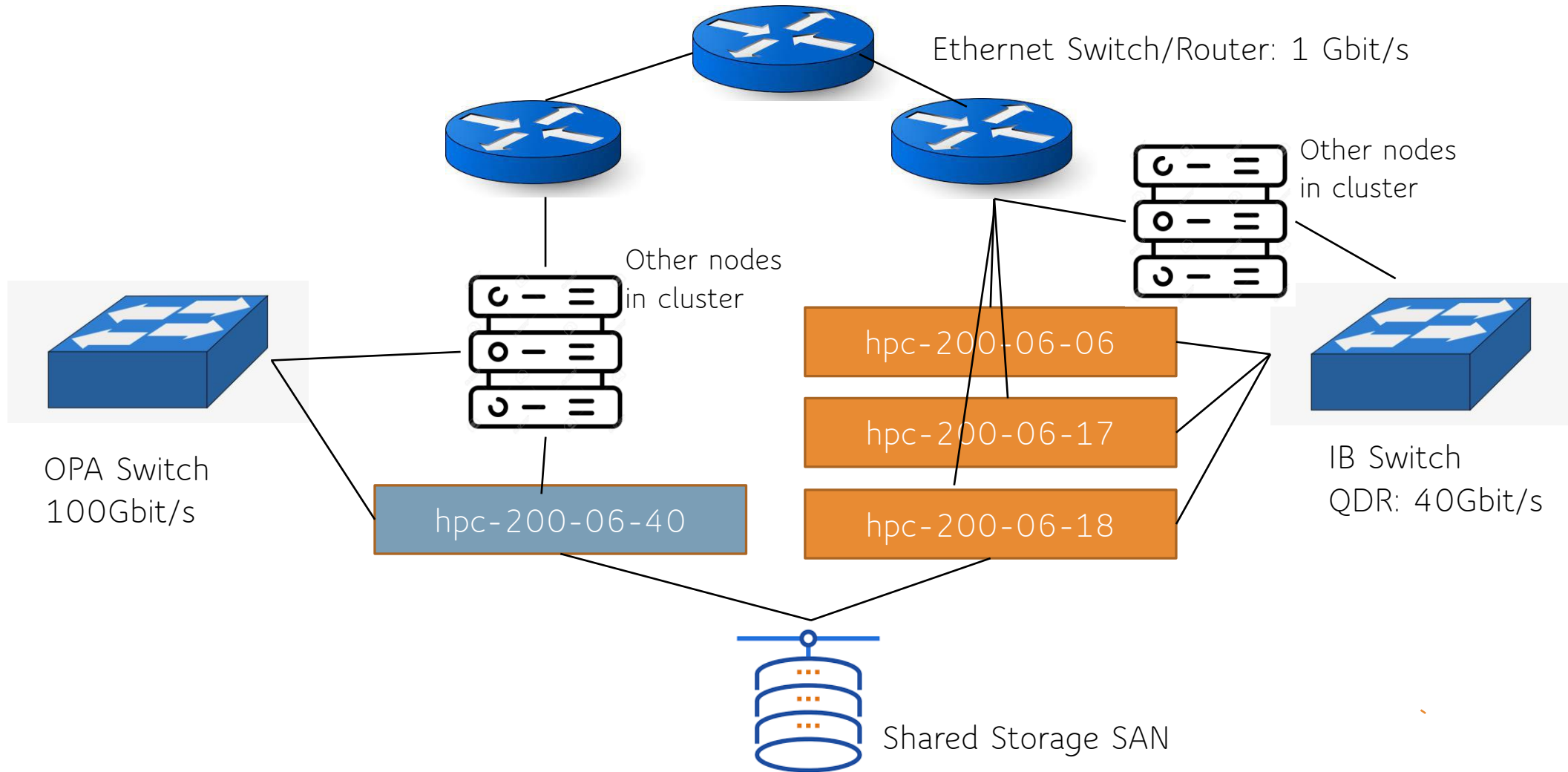
- Knowing where you are is important!!
 - Always try to understand the details of the system you are running on

lstopo -no-graphics
--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

ib0: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 65520
    inet 192.168.184.77 netmask 255.255.255.0 broadcast 192.168.184.255
Infiniband hardware address can be incorrect! Please read BUGS section in ifconfig(8).
    infiniband 80:00:00:02:FE:80:00:00:00:00:00:00:00:00:00:00:00:00:00:00 txqueuelen 256 (InfiniBand)
    RX packets 163 bytes 9128 (8.9 KiB)
    RX errors 0 dropped 0 overruns 0 frame 0
    TX packets 1 bytes 60 (60.0 B)
    TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

lo: flags=73<UP,LOOPBACK,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
```

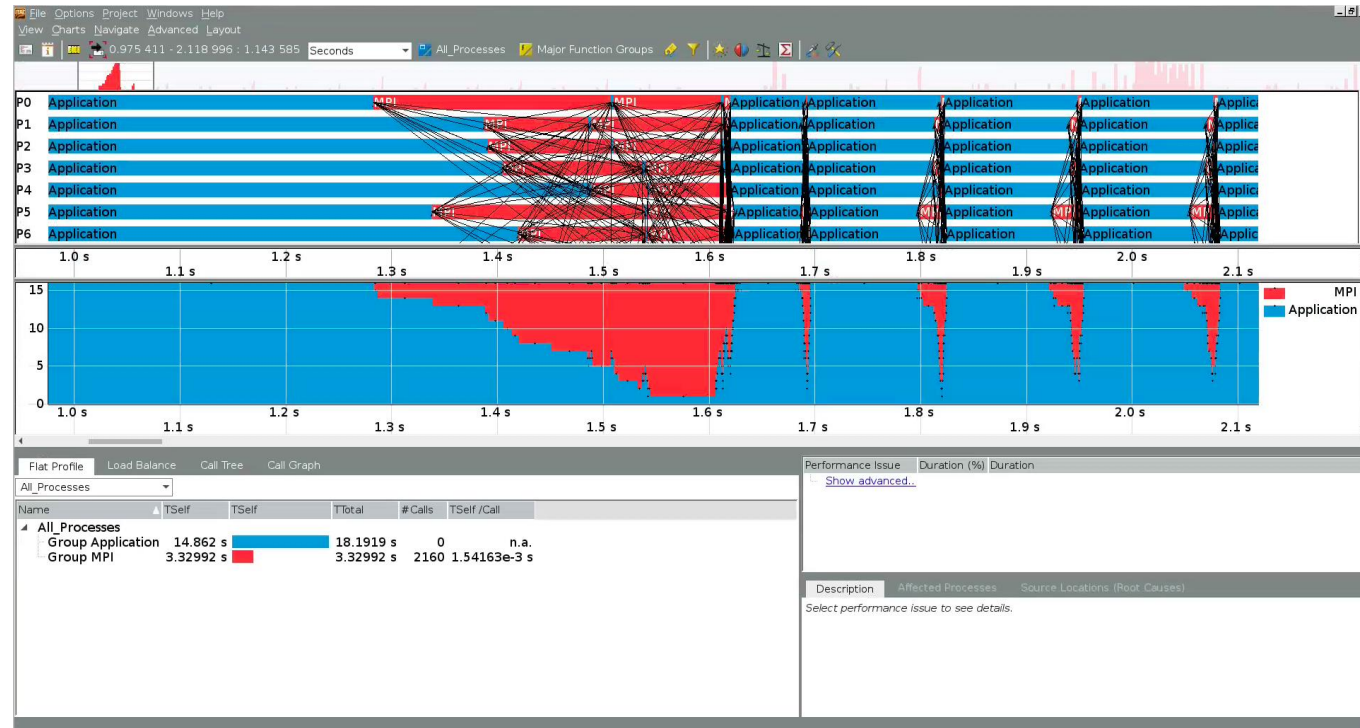
```
[cesinihpc@hpc-201-11-40 ~]$ ibstatus
Infiniband device 'hfi1_0' port 1 status:
    default gid:    fe80:0000:0000:0000:001
    base lid:      0x19
    sm lid:        0x1
    state:         4: ACTIVE
    phys state:    5: LinkUp
    rate:          100 Gb/sec (4X EDR)
    link_layer:    InfiniBand
```

```
[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 or 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 1Gbps latency $O(10.000ns)$
 - Infiniband HDR latency $O(200ns)$
 - DDR4-3600 latency $O(60ns)$
 - DDR5-5600 latency $O(10ns)$

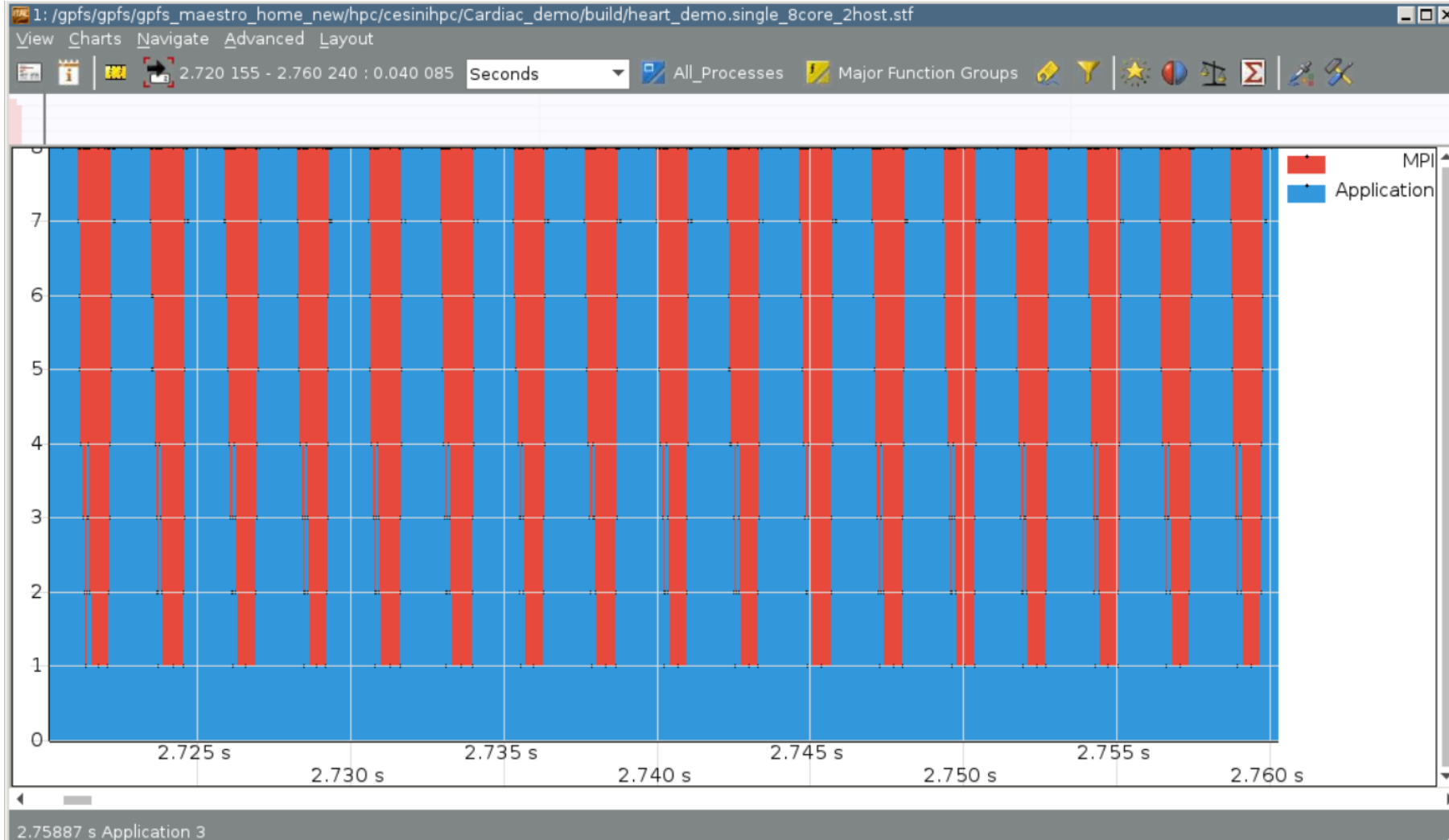
Latency Numbers every programmer should know

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		



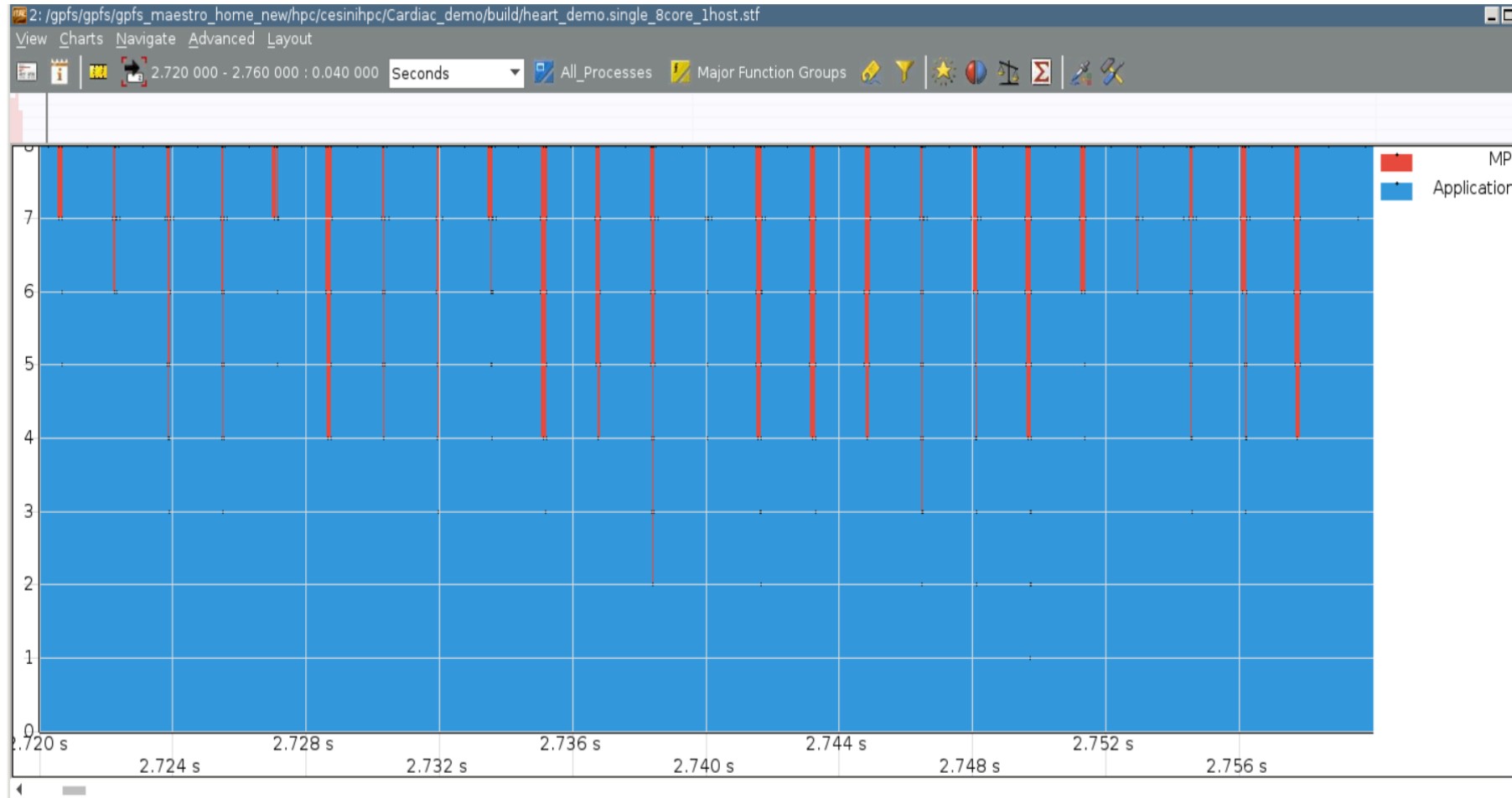
Communication performances in MPI Applications



8 processes
2 hosts
MPI send/receive
over ethernet



Communication performances in MPI Applications



8 processes
2 hosts
MPI send/receive via shared memory

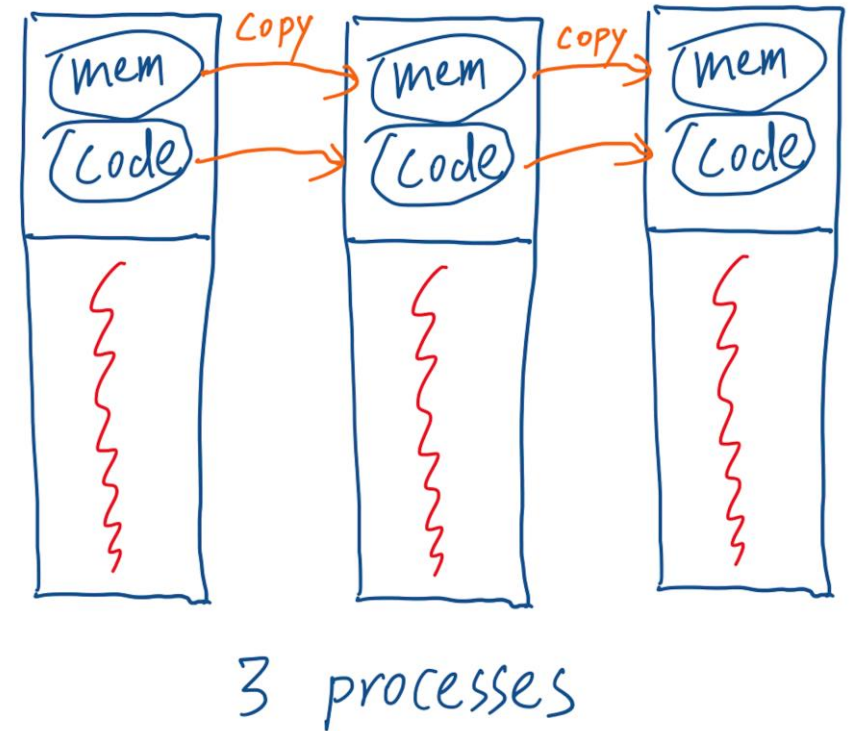


Communication performances in MPI Applications



MPI

- MPI is a standard : <http://www.mpi-forum.org/>
 - 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





SPMD – Single Program Multiple 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++

Process 2

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

https://github.com/inf-nesc/sesame23/blob/main/hands-on/mpi/MPI_Hello.cpp

```
#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, &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(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, &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;
```

```
}
```

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

Called to close any MPI program

- Frees memory allocated by MPI



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;

    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;

    MPI_Finalize();
    return 0;
}
```

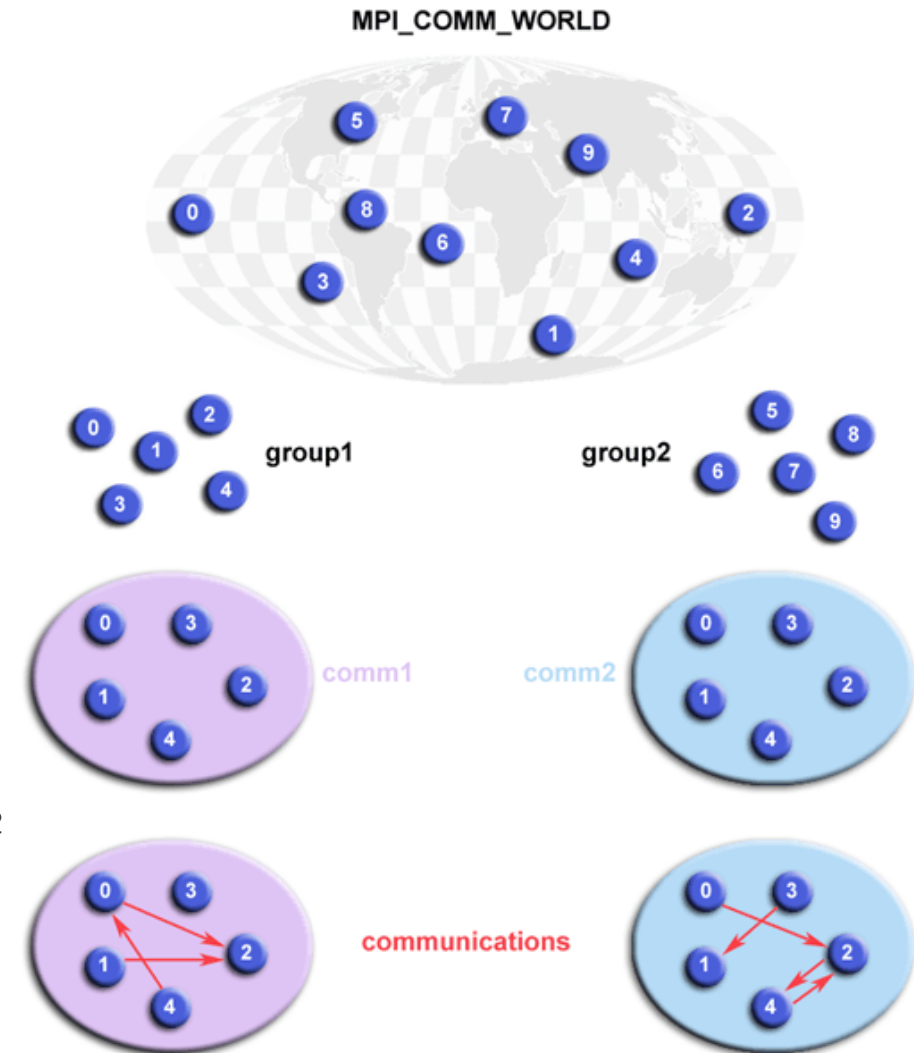
Note that other than `init()` and `finalize()`, every MPI function has a communicator which defines the context and group of processes that the MPI functions impact

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

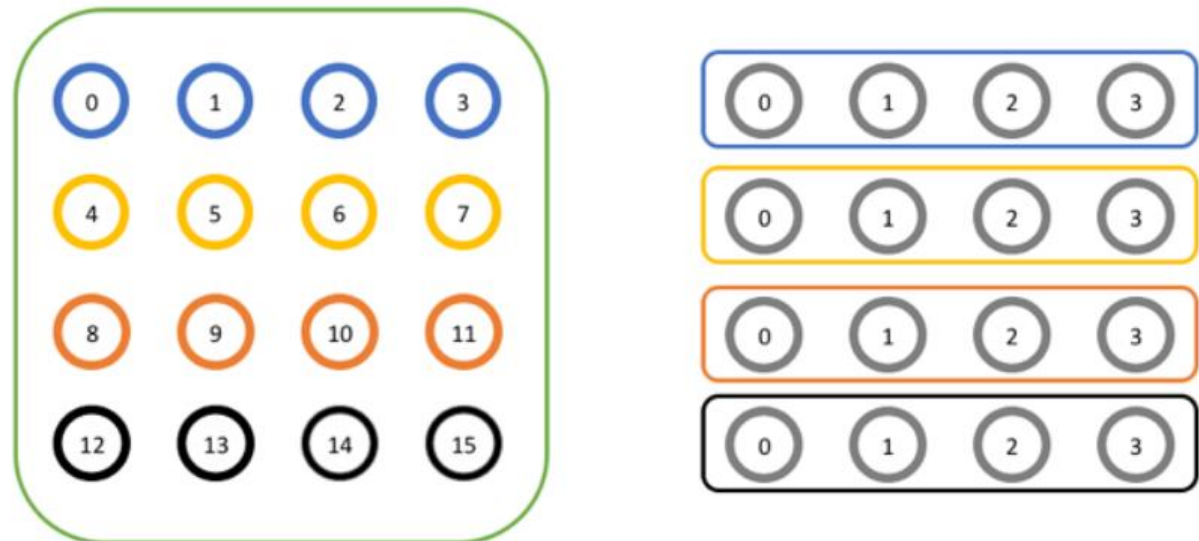
// Split the communicator based on the color and use the
// original rank for ordering
MPI_Comm row_comm;
MPI_Comm_split(MPI_COMM_WORLD, color, world_rank,
&row_comm);

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-groups-and-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
hpc-200-06-17 slots=2
hpc-200-06-02 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 warnings 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
```

```
[cesinihpc@hpc-200-06-18 mpi]$ mpirun --mca btl_openib_allow_ib 1 --hostfile machinefile.txt -np 6 MPI_Hello.out
Hello world from processor hpc-200-06-18.cr.cnaf.infn.it rank 0 of 6
Hello world from processor hpc-200-06-18.cr.cnaf.infn.it rank 1 of 6
Hello world from processor hpc-200-06-17.cr.cnaf.infn.it rank 2 of 6
Hello world from processor hpc-200-06-17.cr.cnaf.infn.it rank 3 of 6
Hello world from processor hpc-200-06-06.cr.cnaf.infn.it rank 4 of 6
Hello world from processor hpc-200-06-06.cr.cnaf.infn.it rank 5 of 6
```



A couple of notes

- The executable must be present in all the hosts used, in the same path
 - You are lucky in the school nodes - shared home directories!!
- OpenMPI in our cluster uses ssh to connect to the remote hosts
 - ssh should work passwordless (HostBasedAuthentication yes in sshd_config)
 - Or create an identity key pair and add the public part to the authorized_keys2 file in .ssh

```
[cesinihpc@hpc-200-06-17 .ssh]$ ll
total 515
-rw-r--r-- 1 cesinihpc hpc 414 May 26 17:09 authorized_keys2
-rw----- 1 cesinihpc hpc 1671 Jun 10 2014 id_rsa
-rw-r--r-- 1 cesinihpc hpc 414 Jun 10 2014 id_rsa.pub
-rw-r--r-- 1 cesinihpc hpc 7881 May 26 17:08 known_hosts
```

```
[cesinihpc@hpc-200-06-17 .ssh]$ cat id_rsa.pub >> authorized_keys2
```

- During login the OpenMPI environment should be loaded
 - Typically via the .basrc file



Point-to-Point Communication



Messages

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

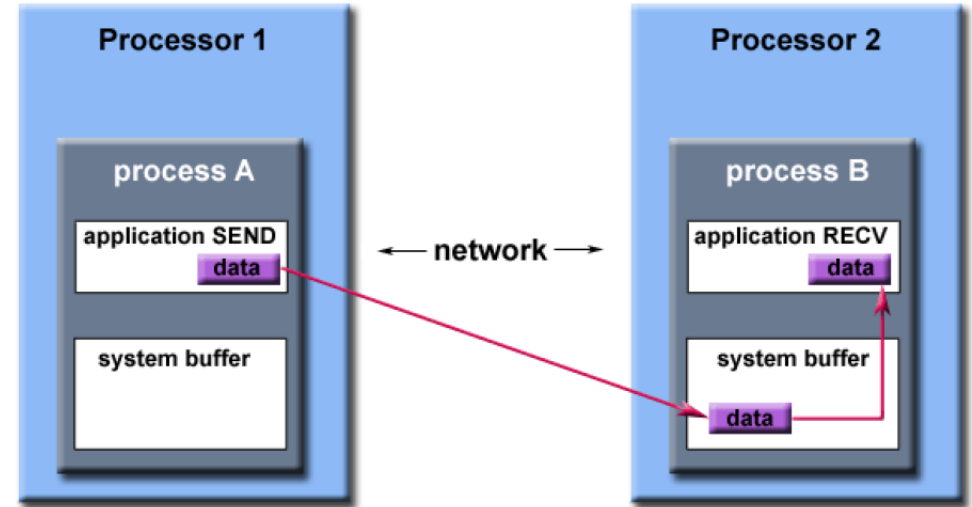
```
int MPI_Send (const void* buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)
```

```
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 cannot 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**

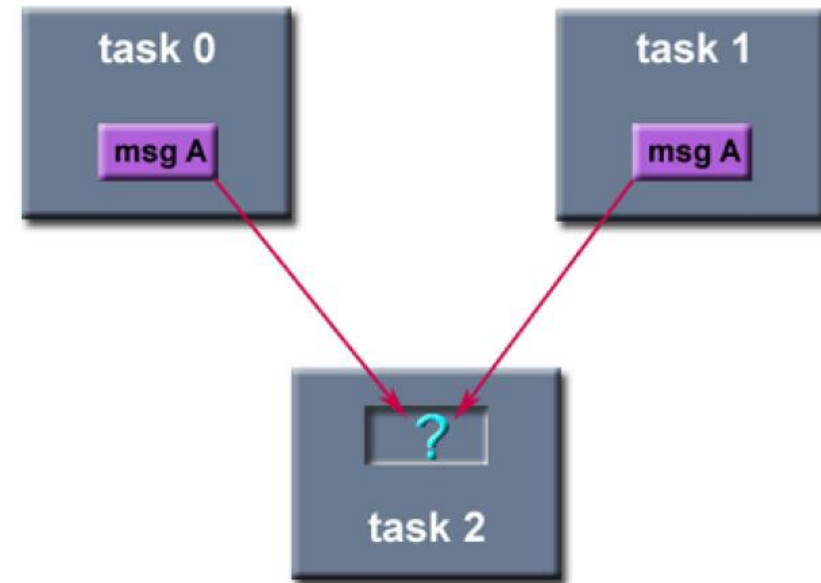


Order

- 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.

Fairness

- 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.



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/inf-nesc/sesame23/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;
    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;
}
// 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;
    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;
    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;
}
```



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**
 - **AT least up to version3**
 - **This changed in version4**
 - OpenMPI 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 20000000 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 enp6s0 --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
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 51.8564 seconds
Bandwidth = 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/inf-n-esc/sesame23/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;

    // 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;
}
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;

    // 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;
}

MPI_Finalize();
```



Collective Communication



Scope

- 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

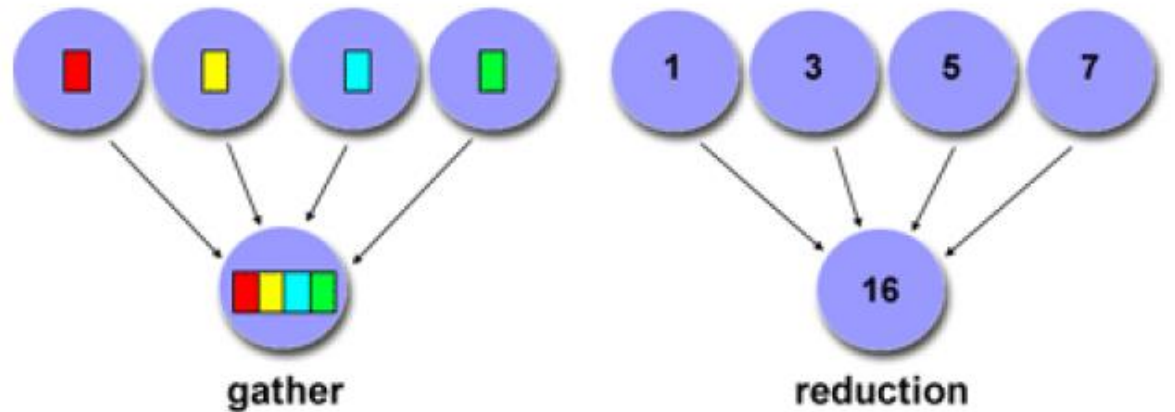
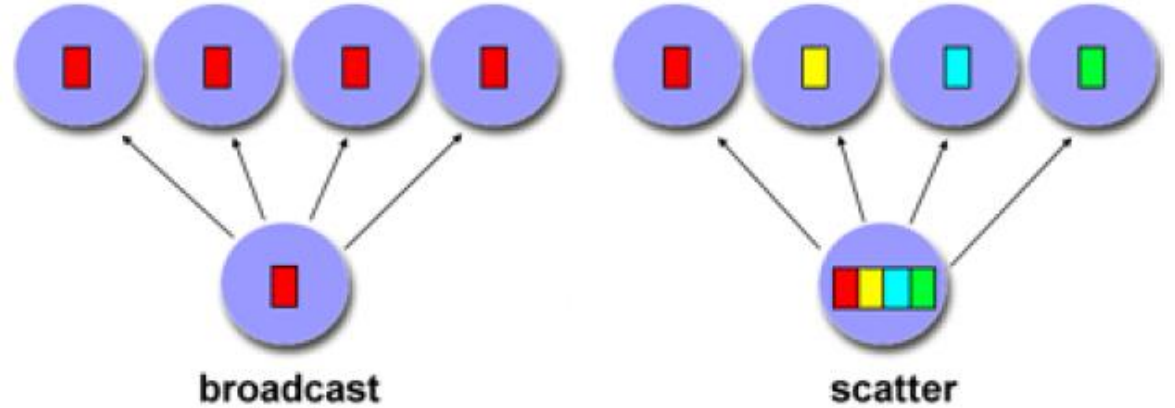
Types of Collective Operations

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

- Data Movement
- Broadcast
- Scatter/gather
- All-to-All.

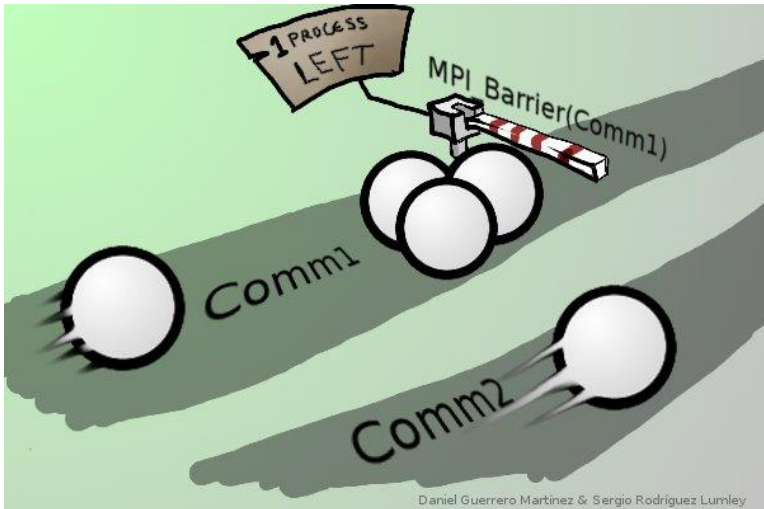
■ Collective Computation (reductions)

- one member of the group collects data from the other members and performs an operation on that data
 - Min
 - Max
 - Add
 - multiply



MPI_Barrier

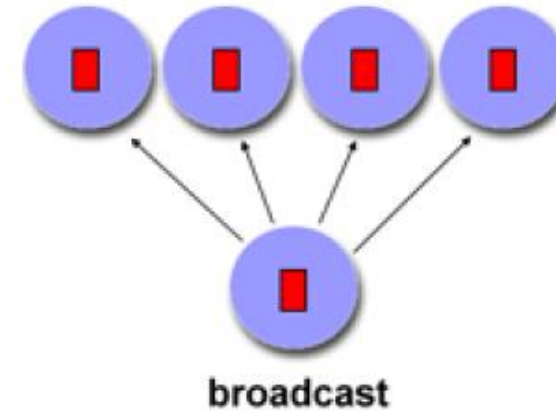
- MPI_Barrier(MPI_Comm Comm)
- Blocks until all processes have reached this routine
 - Blocks the caller until all group members have called it



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.

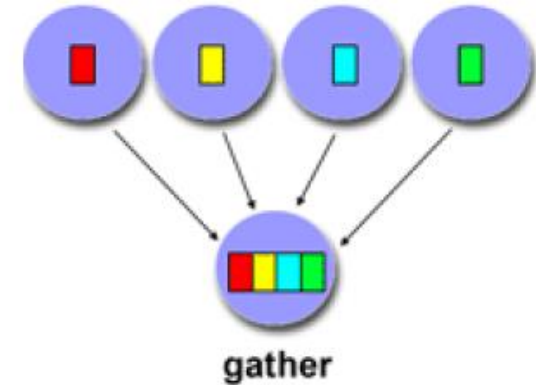


`MPI_BCAST(buffer, count, datatype, root, comm)`

INOUT	<code>buffer</code>	starting address of buffer (choice)
IN	<code>count</code>	number of entries in buffer (non-negative integer)
IN	<code>datatype</code>	datatype of buffer (handle)
IN	<code>root</code>	rank of broadcast root (integer)
IN	<code>comm</code>	communicator (handle)

Gather

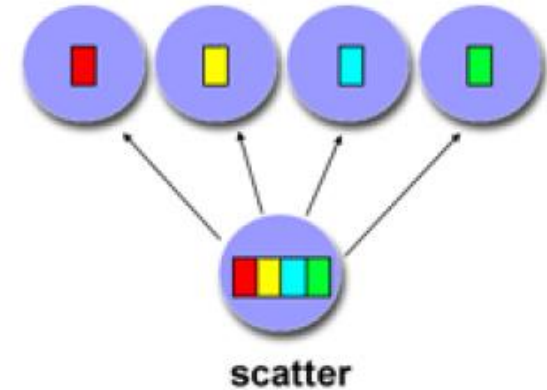
- `int MPI_Gather(const void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)`
 - each process (root process included) sends the contents of its send buffer to the root process.
 - The root process receives the messages and 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



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

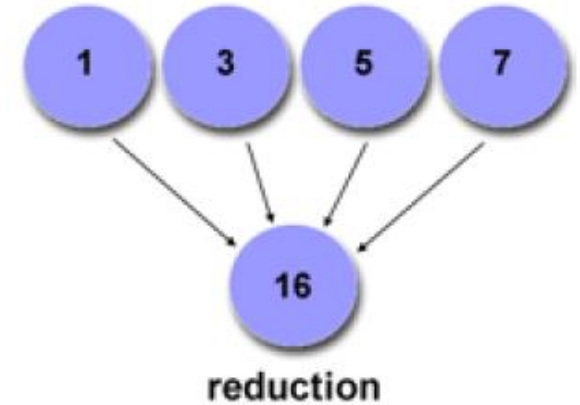


<code>MPI_SCATTER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)</code>		
IN	<code>sendbuf</code>	address of send buffer (choice, significant only at root)
IN	<code>sendcount</code>	number of elements sent to each process (non-negative integer, significant only at root)
IN	<code>sendtype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



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 buffer is defined by the arguments recvbuf, count and datatype;



MPI_REDUCE(sendbuf, recvbuf, count, datatype, op, root, comm)

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	op	reduce operation (handle)
IN	root	rank of root process (integer)
IN	comm	communicator (handle)

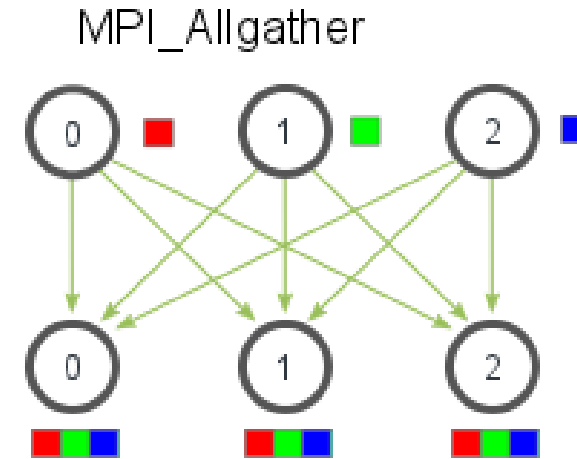


Reduce Operations

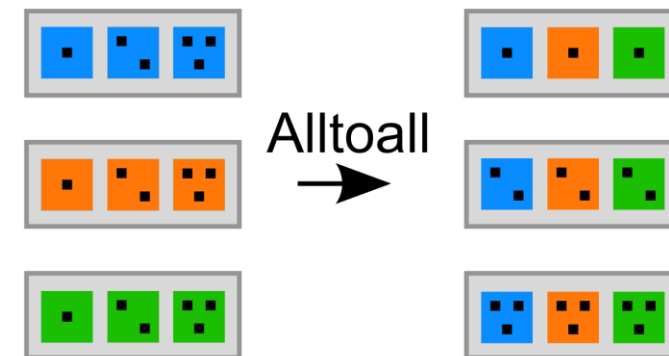
- MPI_MAX - Returns the maximum element
- MPI_MIN - Returns the minimum element
- MPI_SUM - Sums the elements.
- MPI_PROD - Multiplies all elements.
- MPI_LAND - Performs a logical and across the elements
- MPI_LOR - Performs a logical or across the elements
- MPI_BAND - Performs a bitwise and across 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 Operations

- MPI_ALLGATHER can be thought of as MPI_GATHER, but where all processes receive the result, instead of just the root



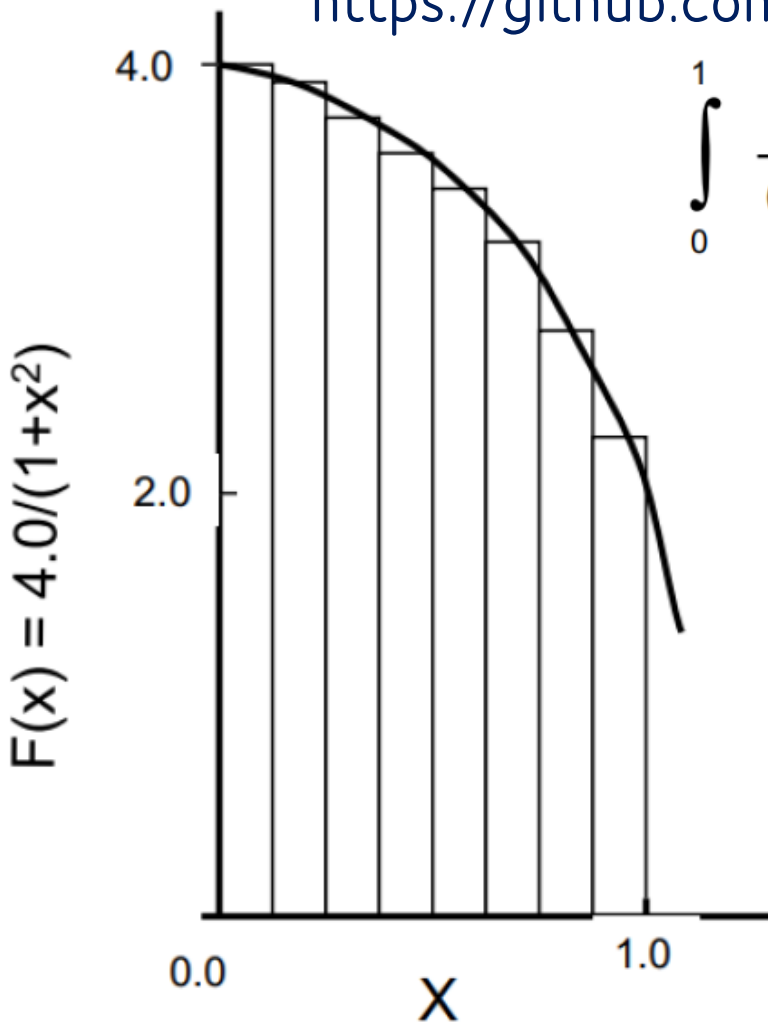
- MPI_ALLTOALL is an extension of MPI_ALLGATHER to the case where each process sends distinct data to each of the receivers. The j -th block sent from process i is received by process j and is placed in the i -th block of recvbuf.





The MPI_Pi

https://github.com/inf-n-esc/sesame23/blob/main/hands-on/mpi/MPI_Pi.cpp



$$\int_0^1 \frac{4.0}{(1+x^2)} dx = \pi \quad \longrightarrow \quad \sum_{i=0}^N F(x_i)\Delta x \approx \pi$$

```
double start_x = (myid * 1.0 / num_procs);

for (long long int i=0; i < steps_per_process; i++){
    auto x = start_x + (i + 0.5)*step;
    sum = sum + 4.0/(1.0 +x*x);
}

mypi = step * sum;

MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
if (myid == 0) {
    std::cout << "result: " << std::setprecision(15) << pi << std::endl;
}

MPI_Finalize();
return 0;
}
```



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

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.open-mpi.org/doc/v4.1/man1/mpirun.1.php>



Run-time Tuning - Binding

- 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
 - 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=1:0-2
```

Rank 0 runs on node aa, bound to logical socket 1, cores 0-2.

```
rank 1=bb slot=0:0,1
```

Rank 1 runs on node bb, bound to logical socket 0, cores 0 and 1.

```
rank 2=cc slot=1-2
```

Rank 2 runs on node cc, bound to logical cores 1 and 2.



Bind-to core example

```
[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 0 bound to socket 0[core 0[hwt 0-1]]: [BB/...../...../...../.....][...../...../...../.....]
[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]]: [..//..//BB/...../...../...../.....][...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 3 bound to socket 0[core 3[hwt 0-1]]: [..//...../BB/...../...../...../.....][...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 4 bound to socket 0[core 4[hwt 0-1]]: [..//...../...../BB/...../...../...../.....][...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 5 bound to socket 0[core 5[hwt 0-1]]: [..//...../...../...../BB/...../...../...../.....][...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 6 bound to socket 0[core 6[hwt 0-1]]: [..//...../...../...../...../BB/...../...../...../.....][...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 7 bound to socket 0[core 7[hwt 0-1]]: [..//...../...../...../...../...../BB][...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 8 bound to socket 1[core 8[hwt 0-1]]: [..//...../...../...../...../.....][BB/...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 9 bound to socket 1[core 9[hwt 0-1]]: [..//...../...../...../...../.....][..//BB/...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 10 bound to socket 1[core 10[hwt 0-1]]: [..//...../...../...../...../.....][...../BB/...../...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 11 bound to socket 1[core 11[hwt 0-1]]: [..//...../...../...../...../.....][...../...../BB/...../...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 12 bound to socket 1[core 12[hwt 0-1]]: [..//...../...../...../...../.....][...../...../...../BB/...../.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 13 bound to socket 1[core 13[hwt 0-1]]: [..//...../...../...../...../.....][...../...../...../...../BB/.....]
[hpc-200-06-17.cr.cnaf.infn.it:04770] MCW rank 14 bound to socket 1[core 14[hwt 0-1]]: [..//...../...../...../...../.....][...../...../...../...../...../BB/.....]
[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 = 4000000000. Step = 2.5e-11.
Numsteps per process = 250000000.
```

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




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

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

- Matrix Multiplication

- <https://www.hpc.cineca.it/content/exercise-16>
- <https://www.hpc.cineca.it/content/solution-16>

- 2D Laplace Equation

- <https://www.mcs.anl.gov/research/projects/mpi/tutorial/mpiexmpl/src/jacobi/C/main.html>