



# Parallelization with MPI

Beginner

28-29 April 2025

ONLINE



# MPI

## Message Passing Interface

One of the most iconic tools in HPC.

No matter if it's large-scale HPC simulations or training large AI models, the backbone that makes this possible is MPI or one of its derivatives.



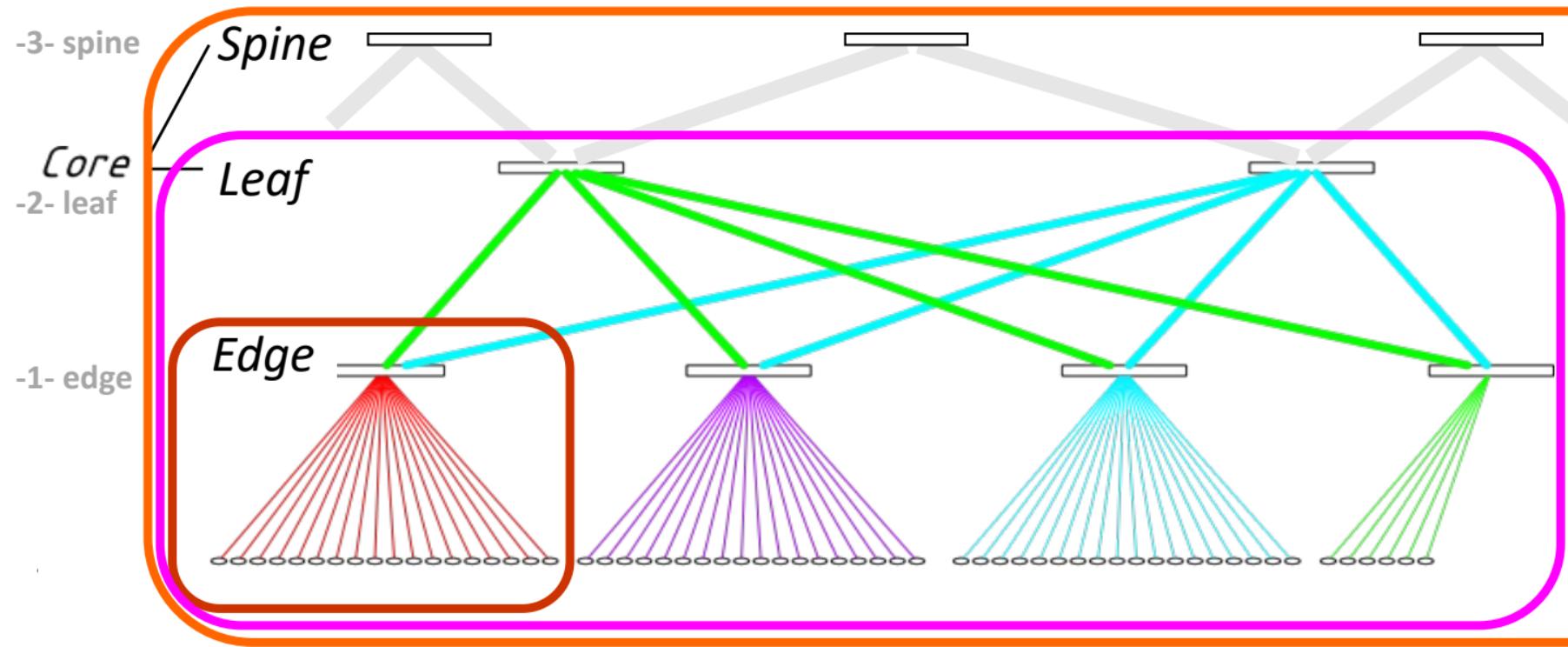
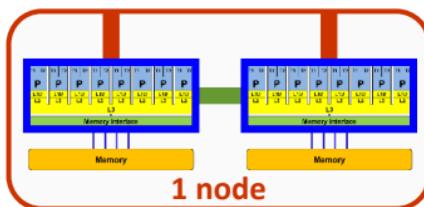
# MPI – communication is key



# VSC-3 ping pong benchmark

to be shown  
after exercises & homework  
of “messages and point-to-point communication”

# VSC-3 – fat-tree interconnect



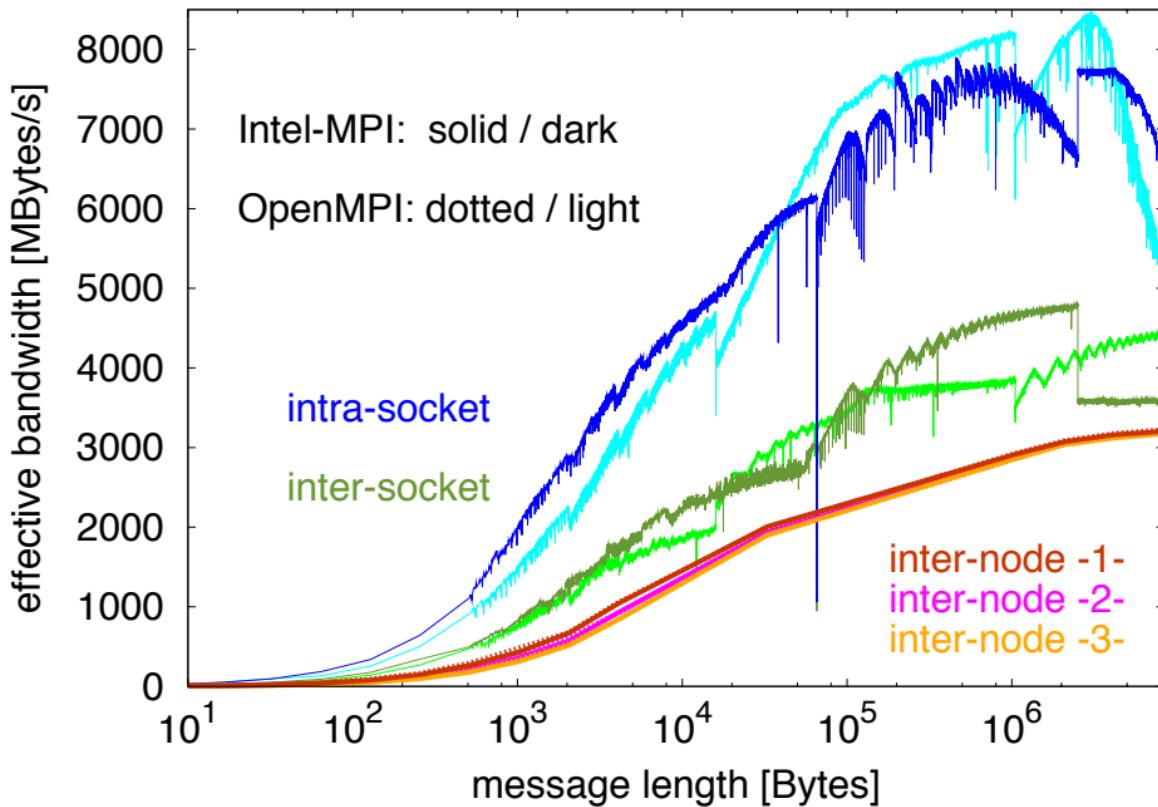
# VSC-3 ping pong benchmark – latency

Latency [μs]	MPI_Send(...)		MPI_Ssend(...)	
	OpenMPI	Intel-MPI	Intel-MPI	OpenMPI
	C~Fortran	C < Fortran	C~<F.	C~<F.
intra-socket	0.3	0.3	0.3	1.2
inter-socket	0.6	0.7	0.7	1.7
IB -1- edge	1.2	1.4	1.5	2.0
IB -2- leaf	1.6	1.8	1.9	2.5
IB -3- spine	2.1	2.3	2.4	3.0

module load intel/16.0.3 intel-mpi/5.1.3 | openmpi/1.10.2

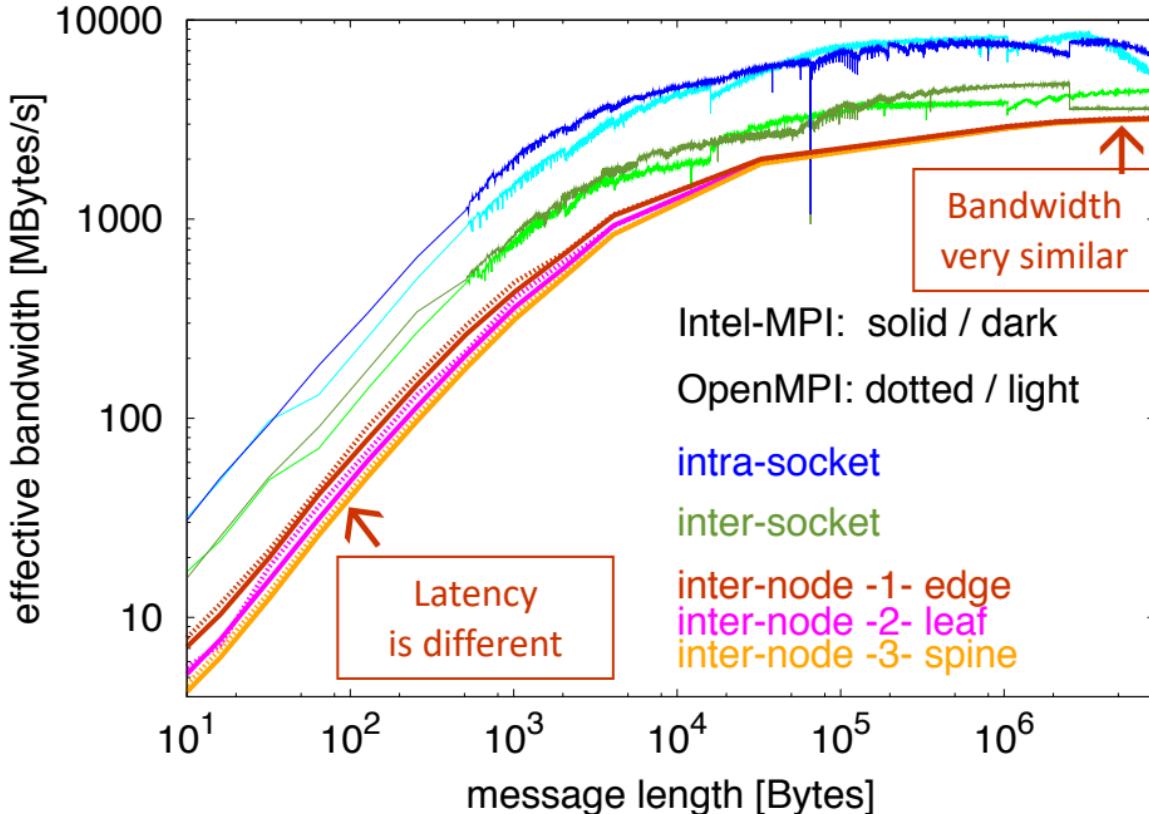
# ping pong benchmark – bandwidth

VSC-3



# ping pong benchmark – bandwidth

VSC-3



# Optimizing MPI communication

## a real-world example

to be shown  
after exercises  
of “collective communication”

# WRONG (ORIG) version – a real world example

s\_p\_e:

loop over irank except self:

**MPI\_ISEND(snd\_buf ... irank ...)**

**MPI\_BARRIER(...)**

loop over irank except self:

**MPI\_RECV(rcv\_buf ... irank ...)**

**MPI\_BARRIER(...)**

DO SOME WORK (snd\_buf is not used = okay)

**MPI\_BARRIER(...)**

problem size	6x6		8x8	
	1_s_p_e	total	1_s_p_e	total
WRONG (ORIG)	1500 s	14 h	21000 s	90 h

# REMOVE BARRIER – a real world example

s\_p\_e:

loop over irank except self:

**MPI\_ISEND(snd\_buf ... irank ...)**

~~MPI\_BARRIER(...)~~

loop over irank except self:

**MPI\_RECV(rcv\_buf ... irank ...)**

~~MPI\_BARRIER(...)~~

DO SOME WORK (snd\_buf is not used = okay)

~~MPI\_BARRIER(...)~~

problem size	6x6		8x8	
timing	1_s_p_e	total	1_s_p_e	total
WRONG (ORIG)	1500 s	14 h	21000 s	90 h
REMOVE BARRIER	960 s	12 h	3500 s	32 h

NEVER use MPI\_BARRIER  
in production code !!!!!

# CORRECT NB comm. – a real world example

s\_p\_e:

ASYNCHRONOUS :: snd\_buf

loop over irank except self:

`MPI_ISEND(snd_buf ... irank ...)`

loop over irank except self:

`MPI_RECV(rcv_buf ... irank ...)`

`MPI_WAIT(...)`

`if (.NOT. MPI_ASYNC_PROT...)`

DO SOME WORK (snd\_buf is not used = okay)

problem size	6x6		8x8	
timing	1_s_p_e	total	1_s_p_e	total
WRONG (ORIG)	1500 s	14 h	21000 s	90 h
REMOVE BARRIER	960 s	12 h	3500 s	32 h
CORRECT NB comm.	425 s	10 h	960 s	18 h

# BCAST version

## – a real world example

s\_p\_e:

~~ASYNCHRONOUS :: snd\_buf~~

~~loop over irank except self:~~

~~MPI\_ISEND(snd\_buf ... irank ...)~~

loop over irank except self:

rcv\_buf = snd\_buf

MPI\_BCAST(rcv\_buf ... irank ...)

problem size	6x6		8x8	
timing	1_s_p_e	total	1_s_p_e	total
WRONG (ORIG)	1500 s	14 h	21000 s	90 h
REMOVE BARRIER	960 s	12 h	3500 s	32 h
CORRECT NB comm.	425 s	10 h	960 s	18 h
BCAST	420 s	10 h	1025 s	20 h

DO SOME WORK (snd\_buf is not used = okay)