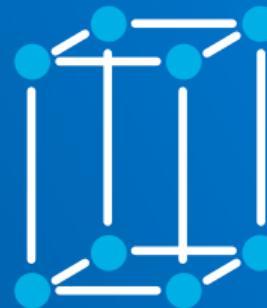


Fundamentals of Parallelism on Intel® Architecture

Week 5
Clusters and MPI

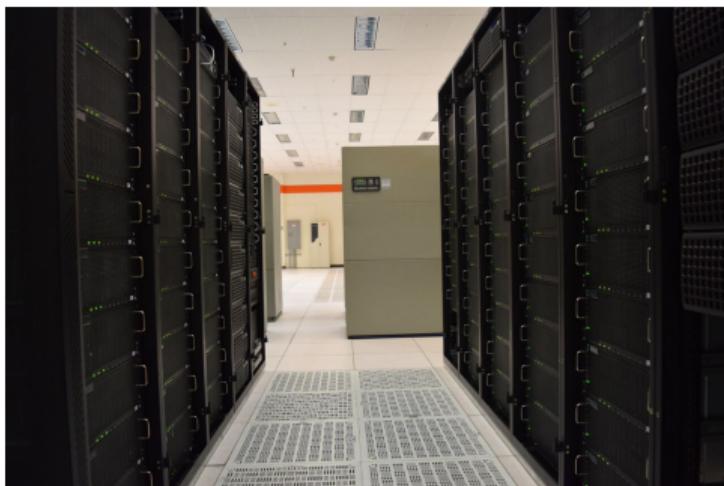


§1. Computing Clusters



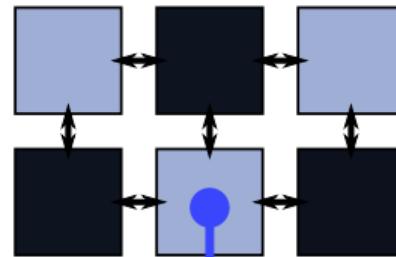
Clusters

Clusters often use Gigabit Ethernet for administration and InfiniBand or Intel Omni-Path for communication.



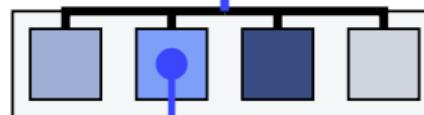
Parallel Programming Layers

CLUSTER COMPUTING
in distributed memory



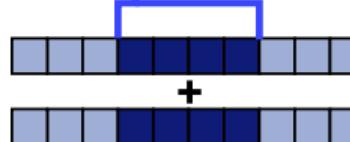
```
MPI_Sendrecv(data, k,  
MPI_DOUBLE, data2,  
... );
```

MULTITHREADING
in shared memory



```
#pragma omp parallel for  
for (j = 0; j < m; j++)  
    ComputeSubset(j);
```

VECTORIZATION
of floating-point math



```
#pragma omp simd  
for (i = 0; i < n; i++)  
    A[i] += B[i];
```

§2. Message Passing Interface, MPI



Message Passing Interface, MPI

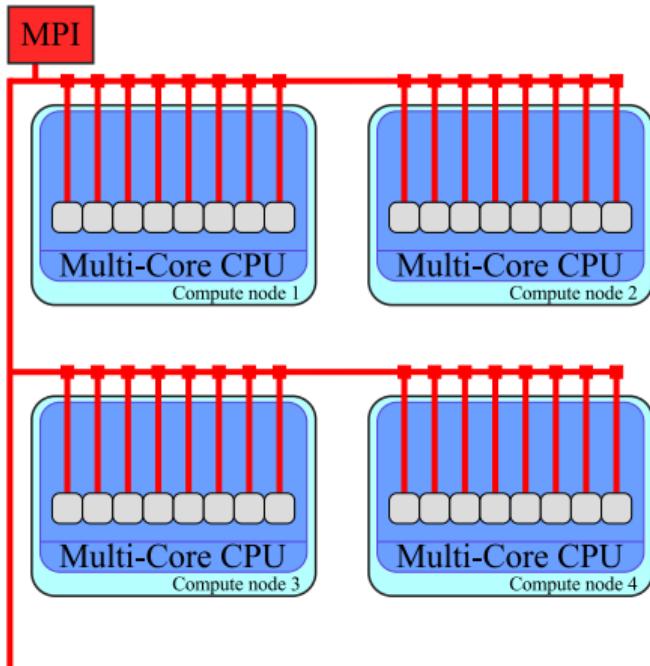


- Specification for message passing
- Multiple implementations exist

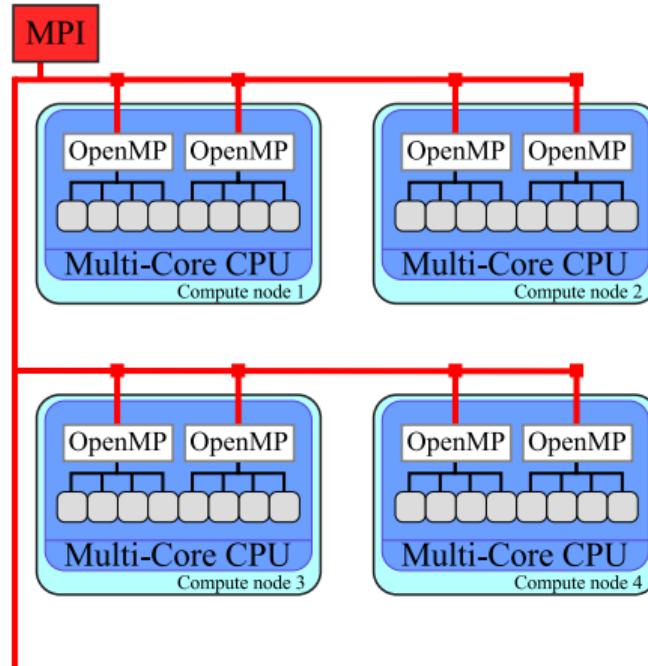
- Portable
- Efficient
- Designed for computing

- Distributed-memory computing
- Multiprocessing in shared memory

Hybrid MPI+OpenMP



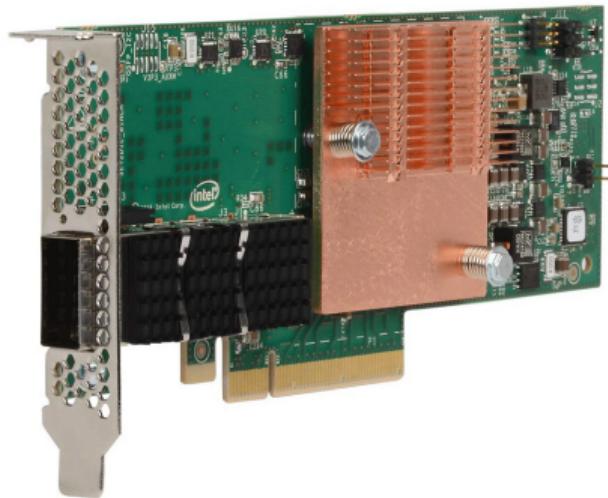
Works for low core counts



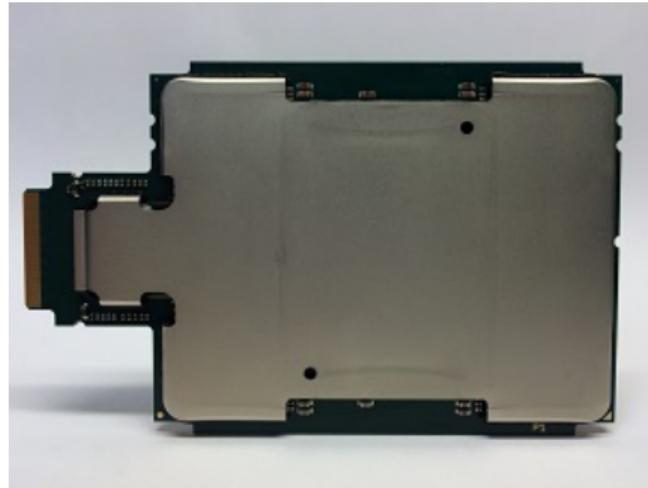
Necessary for multi-core CPUs

Intel's HPC Communication Fabric

Intel Omni-Path Architecture - low-latency, high-bandwidth, scalable communication fabric for HPC applications.



Discrete



Integrated

§3. Programming with MPI



Structure of MPI Applications: Hello World

```
1 #include "mpi.h"
2 #include <cstdio>
3 int main (int argc, char *argv[]) {
4     MPI_Init (&argc, &argv); // Initialize MPI environment
5     int rank, size, namelen;
6     char name[MPI_MAX_PROCESSOR_NAME];
7     MPI_Comm_rank (MPI_COMM_WORLD, &rank); // ID of current process
8     MPI_Get_processor_name (name, &namelen); // Hostname of node
9     MPI_Comm_size (MPI_COMM_WORLD, &size); // Number of processes
10    printf ("Hello World from rank %d running on %s!\n", rank, name);
11    if (rank == 0) printf("MPI World size = %d processes\n", size);
12    MPI_Finalize (); // Terminate MPI environment
13 }
```

MPICH site contains a list of MPI 3.2 routines



§4. Compiling and Running with MPI



Compiling and Running MPI Applications on Localhost

```
u100@c005% mpiicc -o HelloMPI HelloMPI.cc
```

Command file mympi:

```
#PBS -l nodes=1
cd $PBS_O_WORKDIR
mpirun -host localhost -np 2 ./HelloMPI
```

Results:

```
u100@c005% qsub mympi
2000
u100@c005% cat mympi.o2000
Hello World from rank 1 running on c005-n001!
Hello World from rank 0 running on c005-n001!
MPI World size = 2 processes
```



Running MPI Applications on Several Hosts

Command file mydistmpi:

```
#PBS -l nodes=2
cd $PBS_O_WORKDIR
cat $PBS_NODEFILE
mpirun -machinefile $PBS_NODEFILE ./HelloMPI
```

```
u100@c005% qsub mydistmpi
2001
u100@c005% cat mydistmpi.o2001
c005-n001
c005-n002
Hello World from rank 1 running on c005-n002!
Hello World from rank 0 running on c005-n001!
MPI World size = 2 processes
```



§5. Peer-to-Peer Messaging



Point to Point Communication

```
1 if (rank == sender) {  
2  
3     char outMsg[msgLen];  
4     strcpy(outMsg, "Hi There!");  
5     MPI_Send(&outMsg, msgLen, MPI_CHAR, receiver, tag, MPI_COMM_WORLD);  
6  
7 } else if (rank == receiver) {  
8  
9     char inMsg[msgLen];  
10    MPI_Recv (&inMsg, msgLen, MPI_CHAR, sender, tag, MPI_COMM_WORLD, &stat);  
11    printf ("Received message with tag %d: '%s'\n", tag, inMsg);  
12  
13 }
```

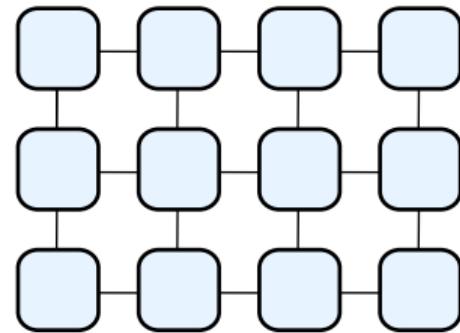
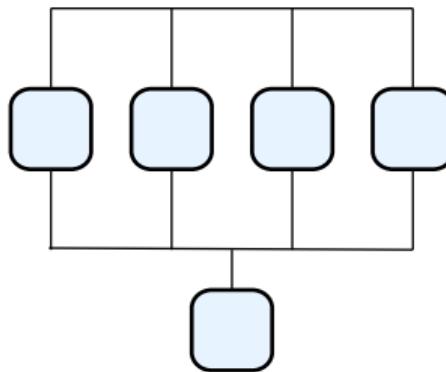
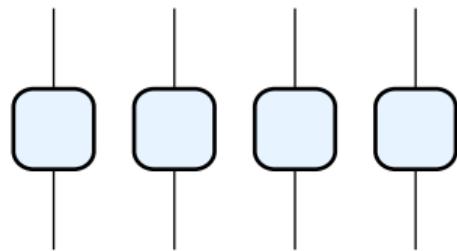


§6. Collective Communication



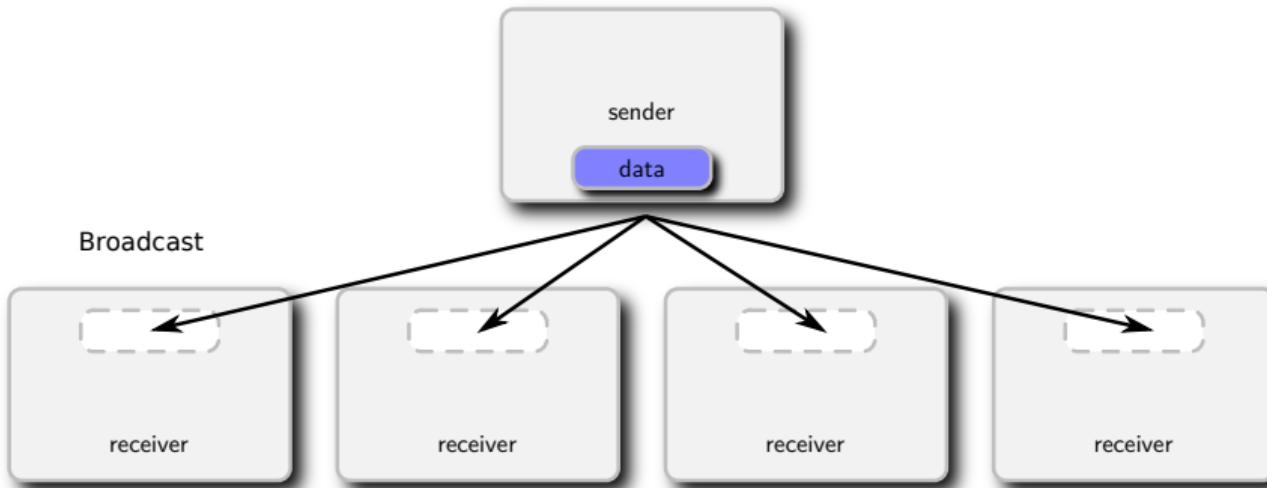
Parallel Patterns

Common parallel patterns call for collective communication



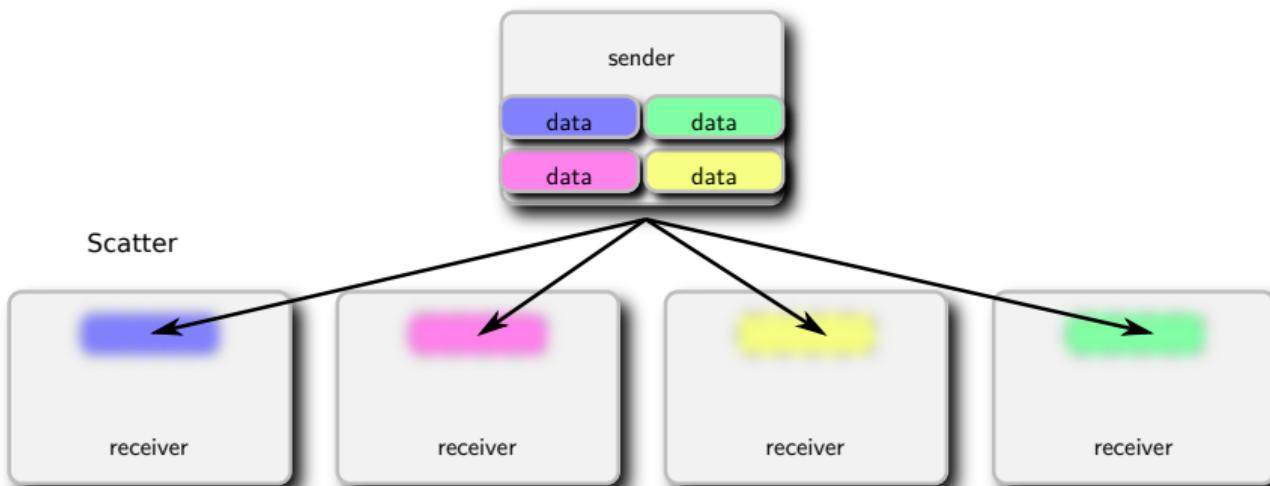
Collective Communication: Broadcast

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



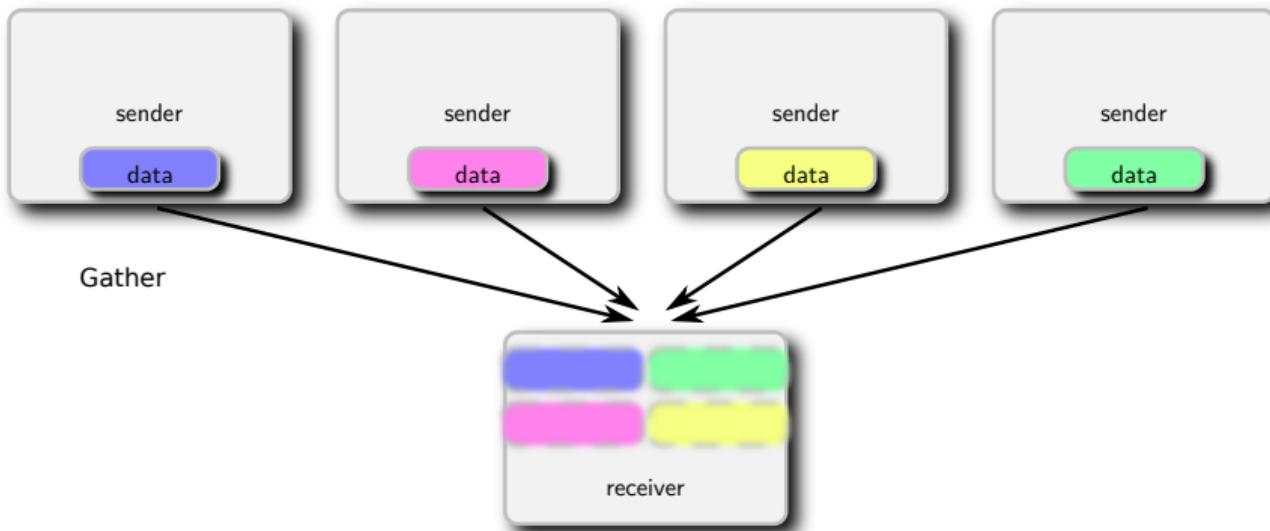
Collective Communication: Scatter

```
1 int MPI_Scatter(void *sendbuf, int sendcnt, MPI_Datatype sendtype, void *recvbuf,  
2                 int recvcnt, MPI_Datatype recvtype, int root, MPI_Comm comm);
```



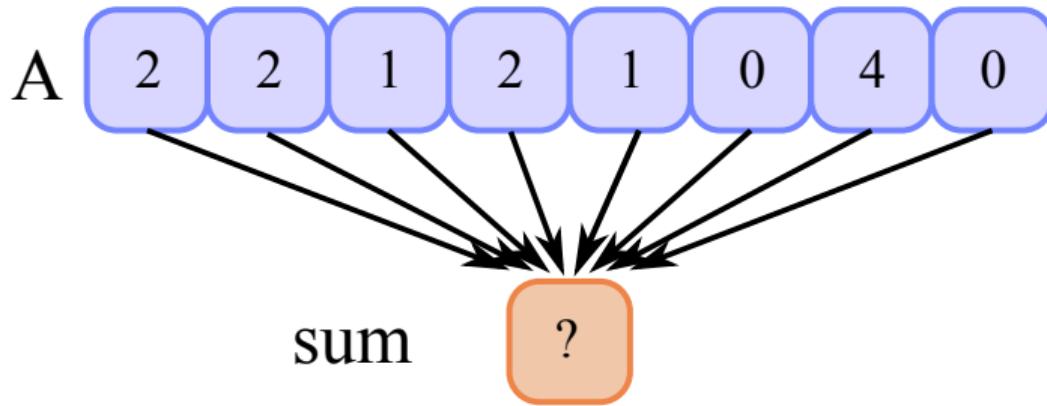
Collective Communication: Gather

```
1 int MPI_Gather(void *sendbuf, int sendcnt, MPI_Datatype sendtype,  
2                 void *recvbuf, int recvcnt, MPI_Datatype recvtype, int root, MPI_Comm com
```



Collective Communication: Reduction

```
1 int MPI_Reduce(void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype,  
2 MPI_Op op, int root, MPI_Comm comm);
```



Available reducers: max/min, minloc/maxloc, sum, product, AND, OR, XOR (logical or bitwise).

§7. Example: Stencil Code



Stencil Operators

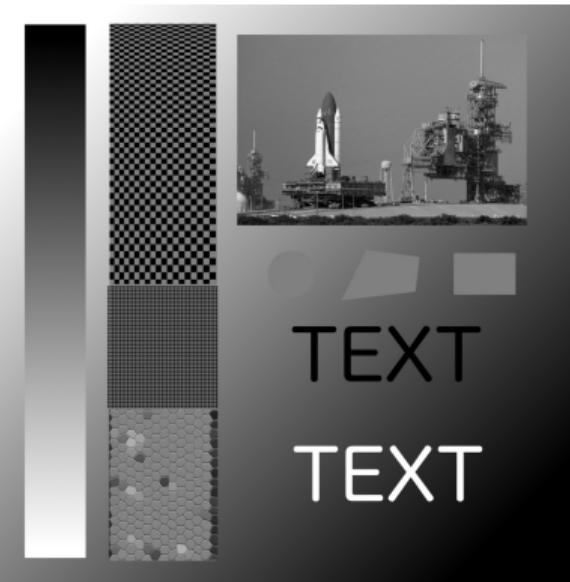
- ▷ Linear systems of equations
- ▷ Partial differential equations

$$Q_{x,y} = c_{00}P_{x-1,y-1} + c_{01}P_{x,y-1} + c_{02}P_{x+1,y-1} + \\ c_{10}P_{x-1,y} + c_{11}P_{x,y} + c_{12}P_{x+1,y} + \\ c_{20}P_{x-1,y+1} + c_{21}P_{x,y+1} + c_{22}P_{x+1,y+1}$$

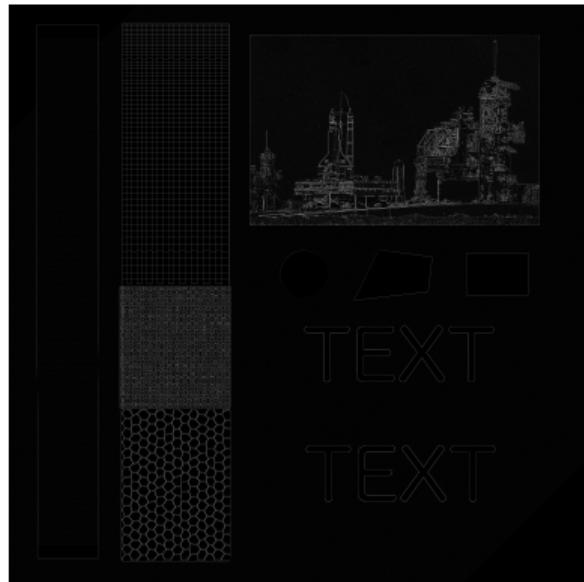
Fluid dynamics, heat transfer, image processing (convolution matrix),
cellular automata.



Edge Detection



$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \rightarrow$$



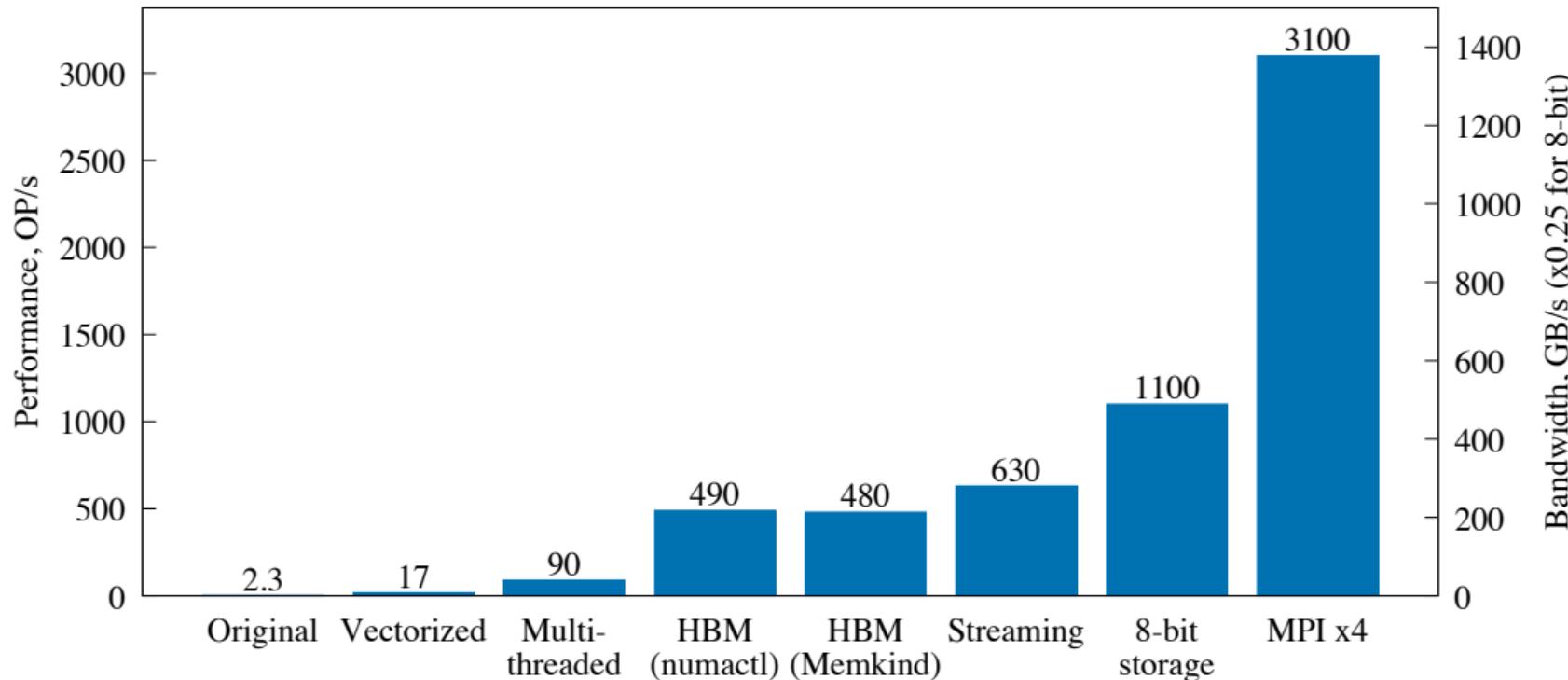
Clustering with MPI

```
#PBS -l nodes=4:flat
cd $PBS_O_WORKDIR
mpirun -machinefile $PBS_NODEFILE ./stencil test-image.png
```

```
1 MPI_Comm_rank(MPI_COMM_WORLD, &myRank);
2 MPI_Comm_size(MPI_COMM_WORLD, &nRanks);
3
4 const double rowsPerProcess = double(height-2)/double(nRanks);
5 const int myFirstRow = 1 + int(rowsPerProcess*myRank);
6 const int myLastRow = 1 + int(rowsPerProcess*(myRank+1));
7
8 #pragma omp parallel for
9 for (int i = myFirstRow; i < myLastRow; i++)
10 ...
```



Performance



§8. Example: Numerical Integration



Midpoint Rectangle Method

$$I(a, b) = \int_0^a f(x) dx \approx \sum_{i=0}^{n-1} f\left(x_{i+\frac{1}{2}}\right) \Delta x,$$

where

$$\Delta x = \frac{a}{n}, \quad x_{i+\frac{1}{2}} = \left(i + \frac{1}{2}\right) \Delta x.$$



Single-node Implementation

```
1 const double dx = a/(double)n;
2 double integral = 0.0;
3 #pragma omp parallel for simd reduction(+: integral)
4 for (int i = 0; i < n; i++) {
5     const double xip12 = dx*((double)i + 0.5);
6     const double dI = BlackBoxFunction(xip12)*dx;
7
8     integral += dI;
9 }
```

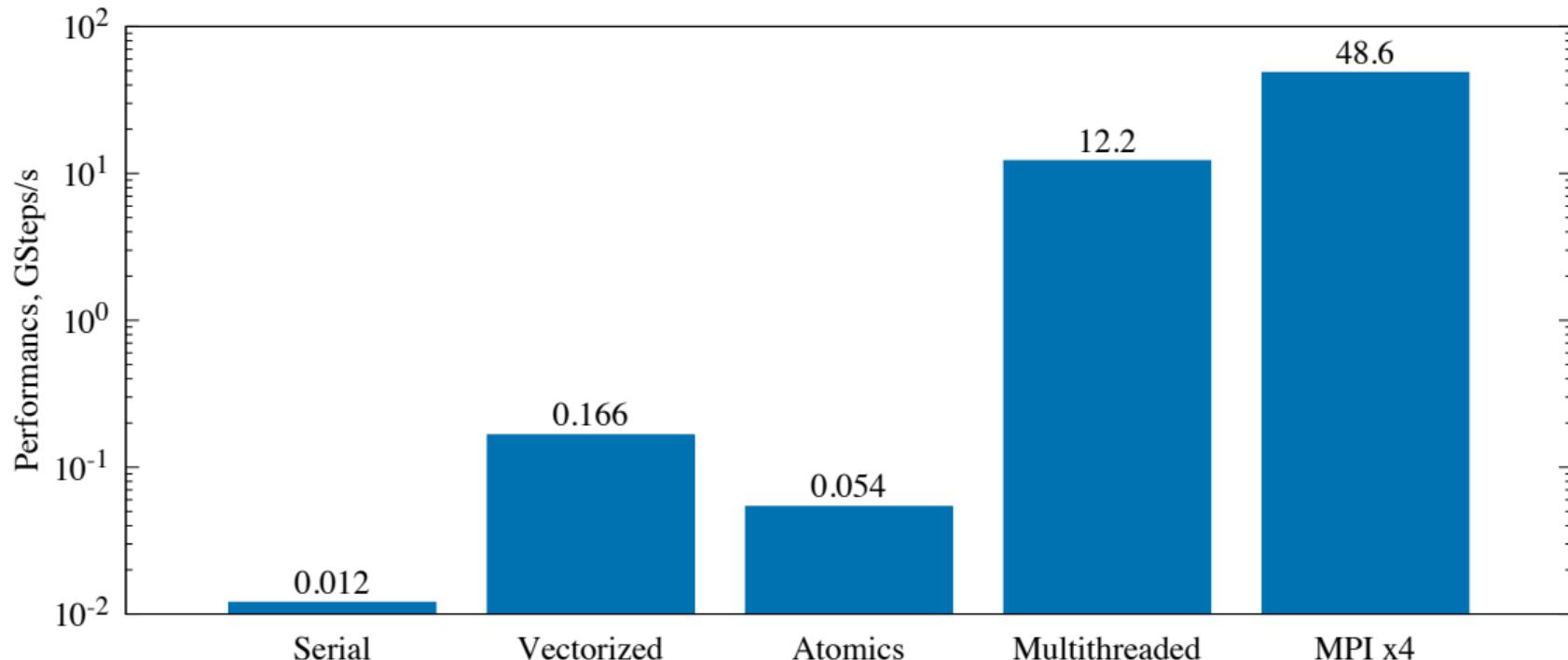


Multi-node Implementation

```
1 const int iStart = double(n)/double(nRanks)*double(rank);
2 const int iEnd   = double(n)/double(nRanks)*double(rank+1);
3
4 const double dx = a/(double)n;
5 double integral_partial = 0.0, integral = 0.0;
6 #pragma omp parallel for simd reduction(+: integral_partial)
7 for (int i = iStart; i < iEnd; i++) {
8     const double xip12 = dx*((double)i + 0.5);
9     const double dI = BlackBoxFunction(xip12)*dx;
10    integral_partial += dI;
11}
12
13 MPI_Allreduce(&integral_partial, &integral, 1, MPI_DOUBLE,
14                MPI_SUM, MPI_COMM_WORLD);
```



Performance



§9. Learn More



MPI Concepts and Functions

Communication modes – (non-)blocking, (a)synchronous, ready mode

Message buffers – application, system, user space

Communicators, Groups – creation, manipulation, usage

Data types – built-in, derived

Collective communication – patterns

Hybrid programming – co-existence with threads: safety, performance

One-sided communication – remote memory access

Click for links from the MPI Forum.

More information in our book and MPI tutorial from the LLNL



Summary on MPI

- ▷ Framework for distributed-memory programming
- ▷ Hides from developer complexity of programming a variety of fabrics
- ▷ Collective communication may use functions of the fabric
- ▷ Intel Omni-Path Architecture is a native solution for Xeon Phi
- ▷ Intel tool for tuning load balance, communication: ITAC

