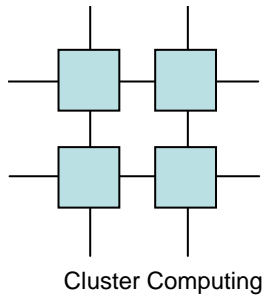


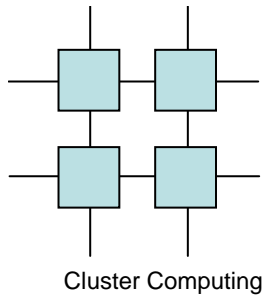
MPI

Industrial Standard Message Passing



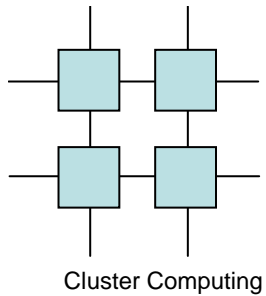
MPI Features

- Industrial Standard
 - Highly portable
 - Widely available
- SPMD programming model
 - Synchronous execution



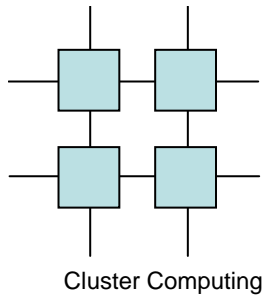
MPI Outer scope

- `int MPI_Init(int *argc, char ** argv)`
- `int MPI_Finalize()`



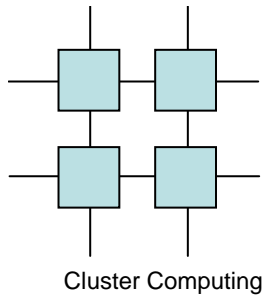
Basic Communication with MPI

- Recieve
 - Blocking
 - `int MPI_Recv(void *buf, int count, MPI_Datatype, int source, int tag, MPI_comm comm, MPI_Status *status)`
 - Nonblocking
 - `int MPI_Irecv(void *buf, int count, MPI_Datatype, int source, int tag, MPI_comm comm, MPI_Request *request)`



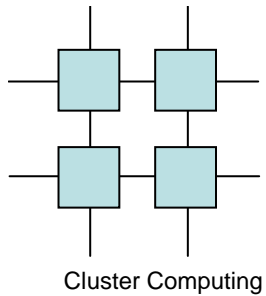
Basic Communication with MPI

- Send
 - Standard
 - Blocking
 - `int MPI_Send(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_Comm comm)`
 - Nonblocking
 - `int MPI_Isend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_Request *request)`



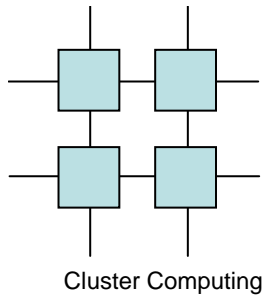
Basic Communication with MPI

- Send
 - Ready
 - Blocking
 - `int MPI_Rsend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_Comm comm)`
 - Nonblocking
 - `int MPI_IrSend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_Request *request)`



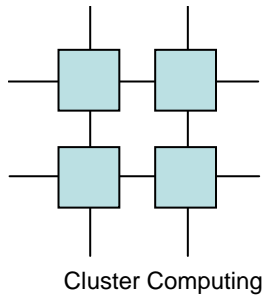
Basic Communication with MPI

- Send
 - Synchronous
 - Blocking
 - `int MPI_Ssend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_comm comm)`
 - Nonblocking
 - `int MPI_Issend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_Request *request)`



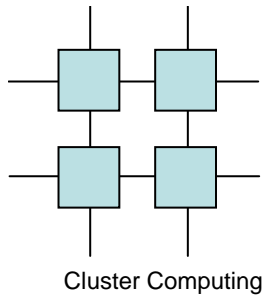
Basic Communication with MPI

- Send
 - Buffered
 - Blocking
 - `int MPI_Bsend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_comm comm)`
 - Nonblocking
 - `int MPI_Ibsend(void *buf, int count, MPI_Datatype, int dest, int tag, MPI_Request *request)`



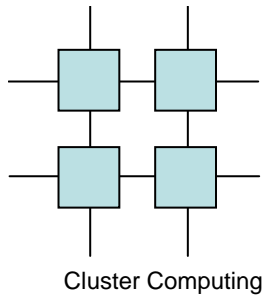
Finalizing nonblocking operations

- Wait for operation to complete
 - `Int MPI_Wait(MPI_Request *request, MPI_Status *status)`
- Test if the operation has completed
 - `Int MPI_Test(MPI_Request *request, int *flag, MPI_Status *status)`



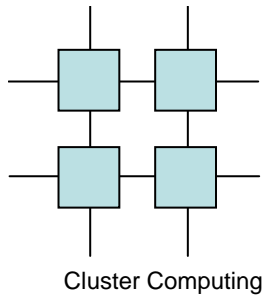
Wait variations

- `MPI_Waitall`
- `MPI_Waitany`
- `MPI_Waitsome`



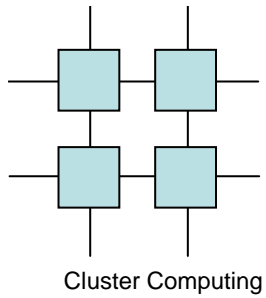
Peeking for messages

- Blocking test for a message
 - `MPI_Probe(int source, int tag,
MPI_Comm comm,
MPI_Status *status);`
- Nonblocking probe
 - `MPI_Probe(int source, int tag,
MPI_Comm comm,
int *flag;
MPI_Status *status);`



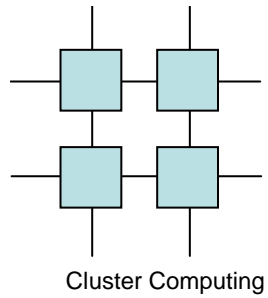
Addressing in MPI

- All processes are named continuously from 0 within one or more process groups
- Process address
 - (id, communicator)



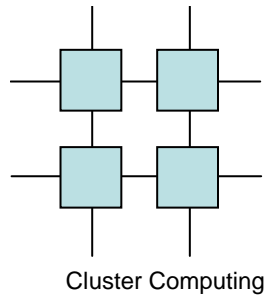
Data Type Support in MPI

C Datatype	MPI Datatype
int	MPI_INT
long	MPI_LONG
float	MPI_FLOAT
double	MPI_DOUBLE
char	MPI_CHAR
(struct)	MPI_PACKED
	MPI_BYTE

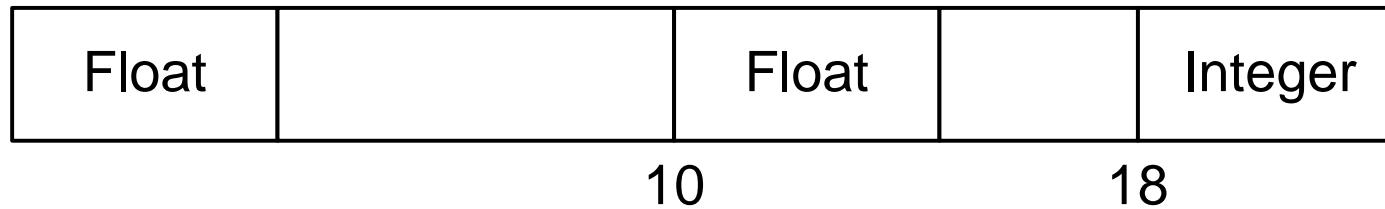


Application defined data-types

- MPI can handle userdefined datatypes
- A datatype is defined as
 - A vector of types
 - A vector on instances
 - A vector of displacements

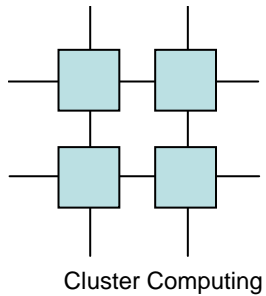


Application defined data- types



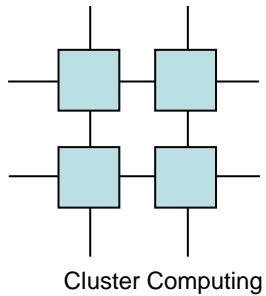
(float, float, int)

 $(1,1,1)$
$$(0, 10, 18)$$



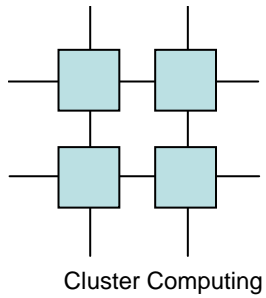
Application defined data-types

```
int blk_len[ 3] = { 1,1,1 };
MPI_Aint displ[ 3], start_addr, addr;
MPI_Datatype typel[ 3]={ MPI_FLOAT, MPI_FLOAT, MPI_INT };
displ[ 0] = 0;
MPI_Address (a,& start_addr);
MPI_Address (b,& addr);
displ[1] = addr - start_addr;
MPI_Address (n,& addr);
displ[2] = addr - start_addr;
MPI_Type_struct (3, blk_len, displ, typel, msg_ptr);
MPI_Type_commit (msg_ptr);
```

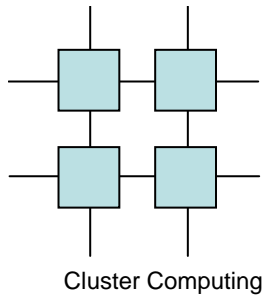
Communicator

- MPI processes belong to one or more communicators
- Addressing MPI processes always include a communicator and a absolute address within the communicator
- All processes initially belong to the communicator `MPI_COMM_WORLD`



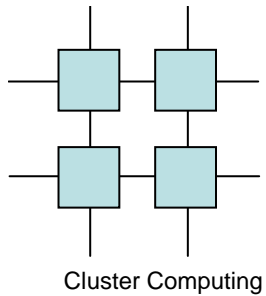
Communicators

- `int MPI_Comm_Size(MPI_comm comm, int* size)`
- `int MPI_Comm_Rank(MPI_comm comm, int* rank)`



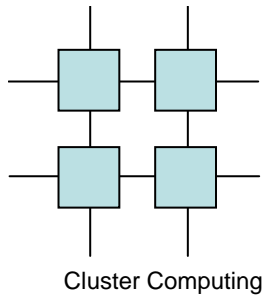
All for one and one for all

- Barriers
- Broadcasting
- Multicasting
- Scattering data
- Gathering data
- All to All
- Reductions
- Scanning



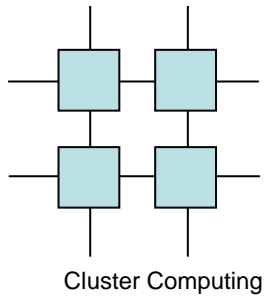
Barriers

- MPI provides a general barrier mechanism
- `int MPI_Barrier(MPI_COMM_WORLD);`
 - + Very easy to use
 - + Can be optimized for the specific topology



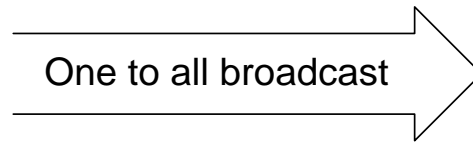
Broadcasting

- Send a message to all participants
- `int MPI_Bcast(void *buf, int count,
MPI_Datatype datatype,
int root,
MPI_comm comm)`

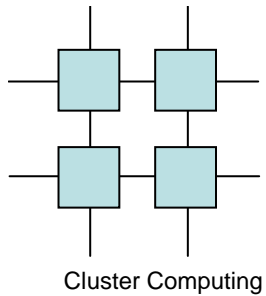


Broadcast

A_0			

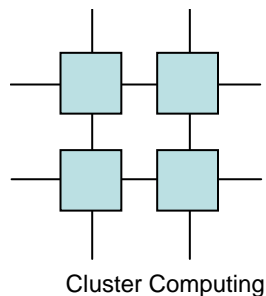


A_0			
A_0			
A_0			
A_0			



All to All

- Efficient way to distribute partial results to all participants
- `int MPI_Alltoall(void *sendbuf, int scount, MPI_Datatype sdatatype, void *recvbuf, int rcount, MPI_Datatype rdatatype, MPI_comm comm)`



All to All

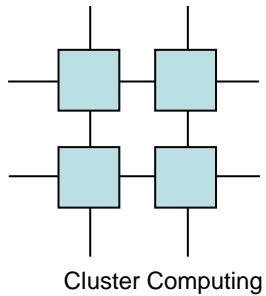
Send Buffer

P r o c e s s o r	Data ➤					
	A₀	B₀	C₀	D₀	E₀	F₀
	A₁	B₁	C₁	D₁	E₁	F₁
	A₂	B₂	C₂	D₂	E₂	F₂
	A₃	B₃	C₃	D₃	E₃	F₃
	A₄	B₄	C₄	D₄	E₄	F₄
	A₅	B₅	C₅	D₅	E₅	F₅



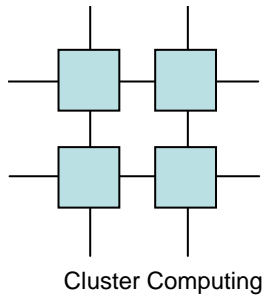
Receive Buffer

P r o c e s s o r	Data ➤					
	A₀	A₁	A₂	A₃	A₄	A₅
	B₀	B₁	B₂	B₃	B₄	B₅
	C₀	C₁	C₂	C₃	C₄	C₅
	D₀	D₁	D₂	D₃	D₄	D₅
	E₀	E₁	E₂	E₃	E₄	E₅
	F₀	F₁	F₂	F₃	F₄	F₅



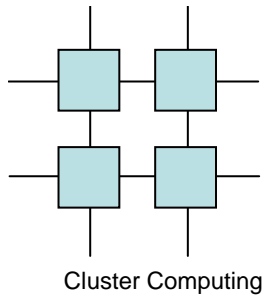
Multicasting

- Multicasting does not make sense in MPI
- Instead a processgroup can be created and broadcast within the group can be used



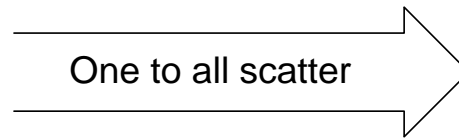
Scattering data

- Allows easy distribution of data amongst processors
- `int MPI_Scatter(void *sendbuf, int scount, MPI_Datatype sdatatype, void *recvbuf, int rcount, MPI_Datatype rdatatype, int root, MPI_comm comm)`

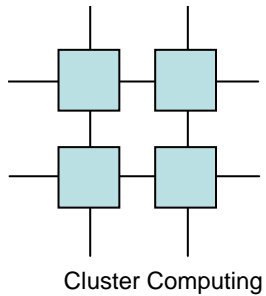


Scattering data

A_0	A_1	A_2	A_3

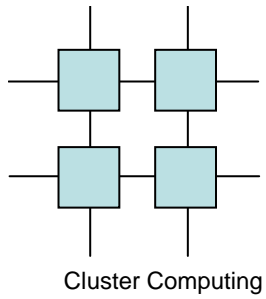


A_0			
A_1			
A_2			
A_3			



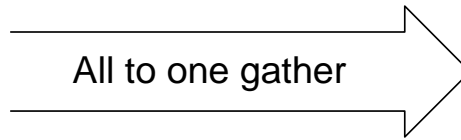
Gathering data

- Easy collection of results
- `int MPI_Gather(void *sendbuf, int scount, MPI_Datatype sdatatype, void *recvbuf, int rcount, MPI_Datatype rdatatype, int root, MPI_comm comm)`

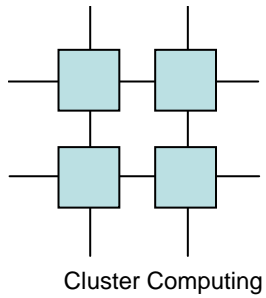


Gathering Data

A_0			
A_1			
A_2			
A_3			

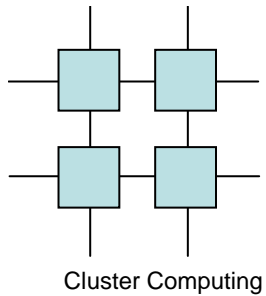


A_0	A_1	A_2	A_3



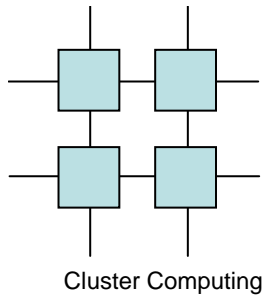
Reductions

- $D = D_0 \oplus D_1 \oplus D_2 \oplus D_n$
- Builtin reductions
 - MPI_MAX
 - MPI_MIN
 - MPI_SUM
 - MPI_PROD



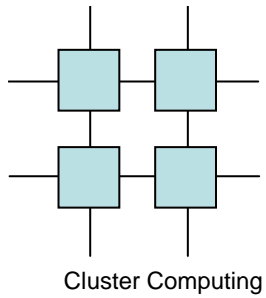
MPI_Reduce

- `int MPI_Reduce(void *sndbuf,
void *recvbuf,
int count,
MPI_Datatype
datatype,
MPI_Op op,
int root,
MPI_comm comm)`
- Similar with `MPI_Allreduce`



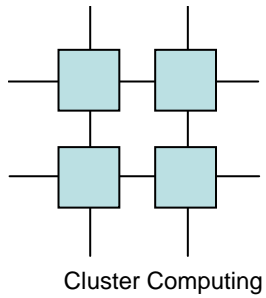
Scanning

- $D_i = D_0 \oplus D_1 \oplus D_2 \oplus \dots \oplus D_{i-1}$
- Same builtins as with Reductions
- `int MPI_Scan (void *sndbuf,
void *recvbuf,
int count,
MPI_Datatype
datatype,
MPI_Op op,
MPI_comm comm)`



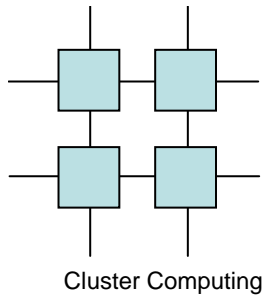
Mapping topologies to MPI

- MPI applications all have a topology
 - Graph
 - Cartesian
 - No topology



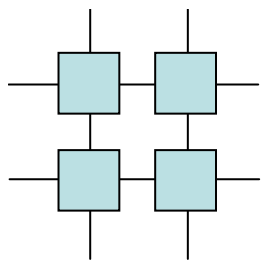
Mapping topologies to MPI

- MPI_Graphcreate
- MPI_Cartcreate
 - MPI_Carttrank
 - MPI_Cartcoords



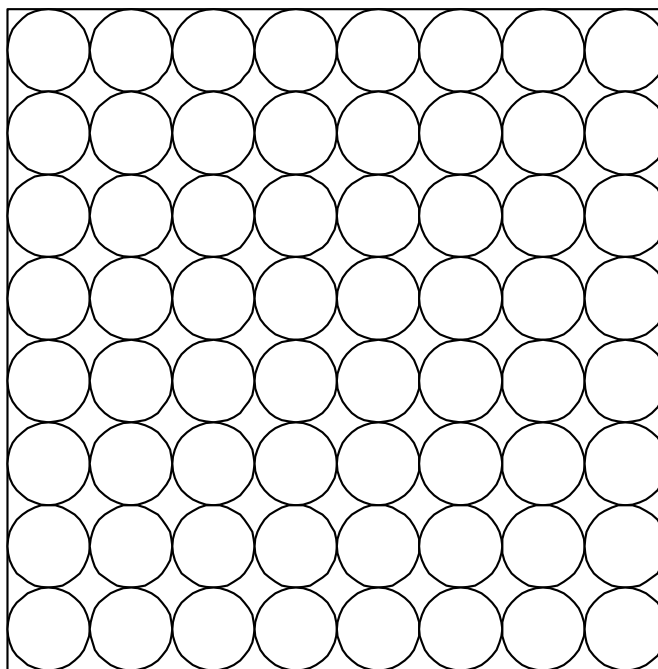
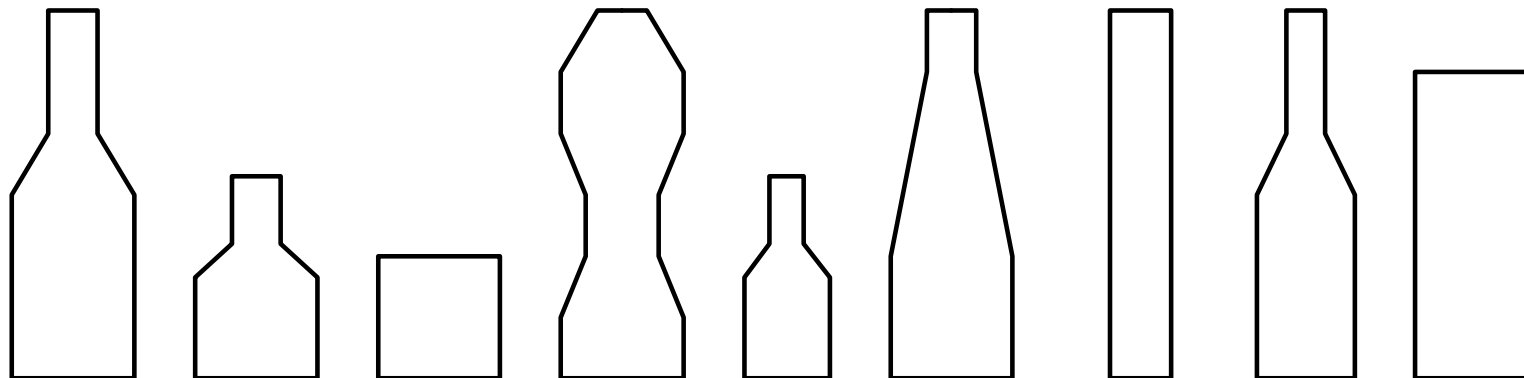
Example Neural Network training

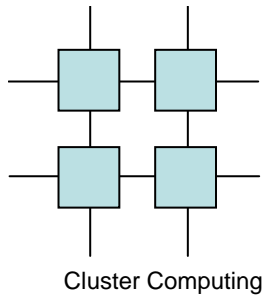
- Training of large neural networks is very hard



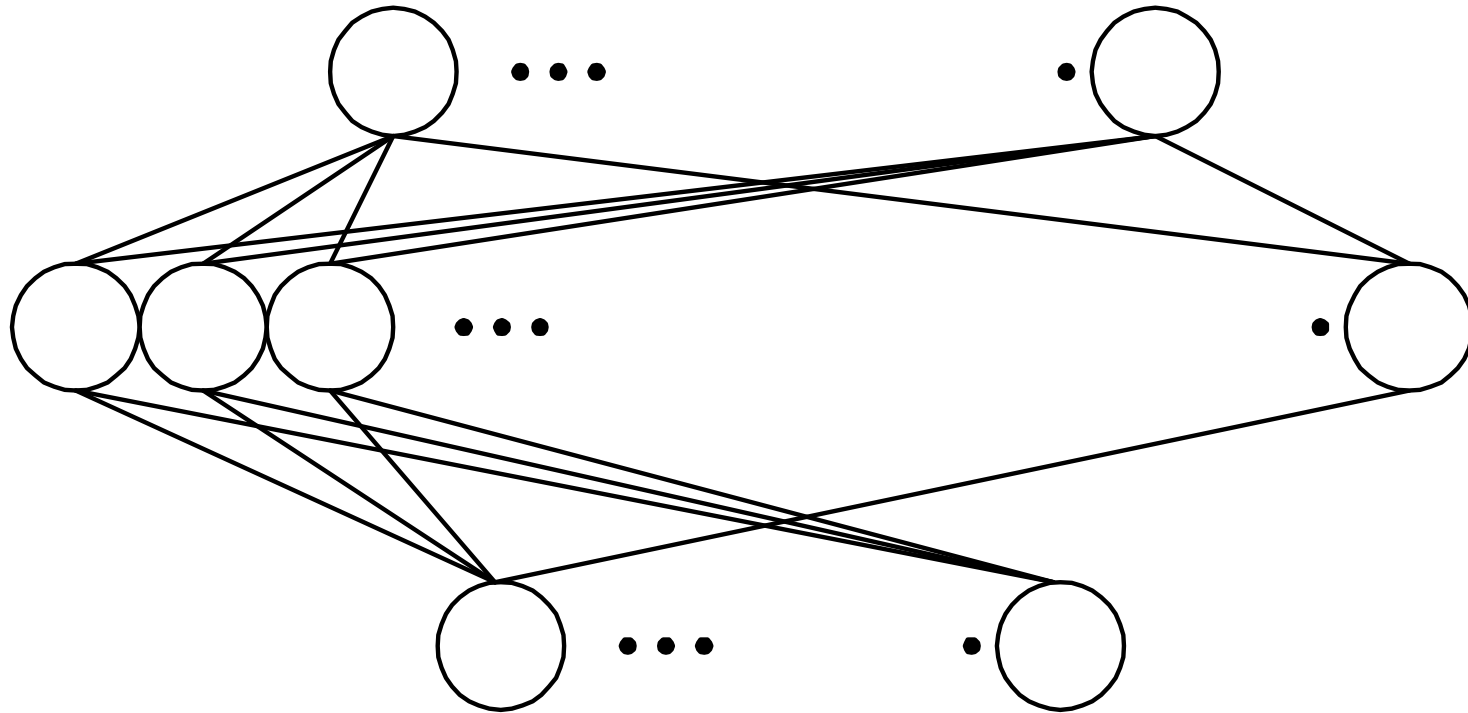
Cluster Computing

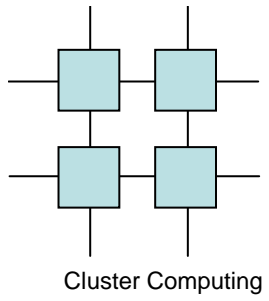
NN Training



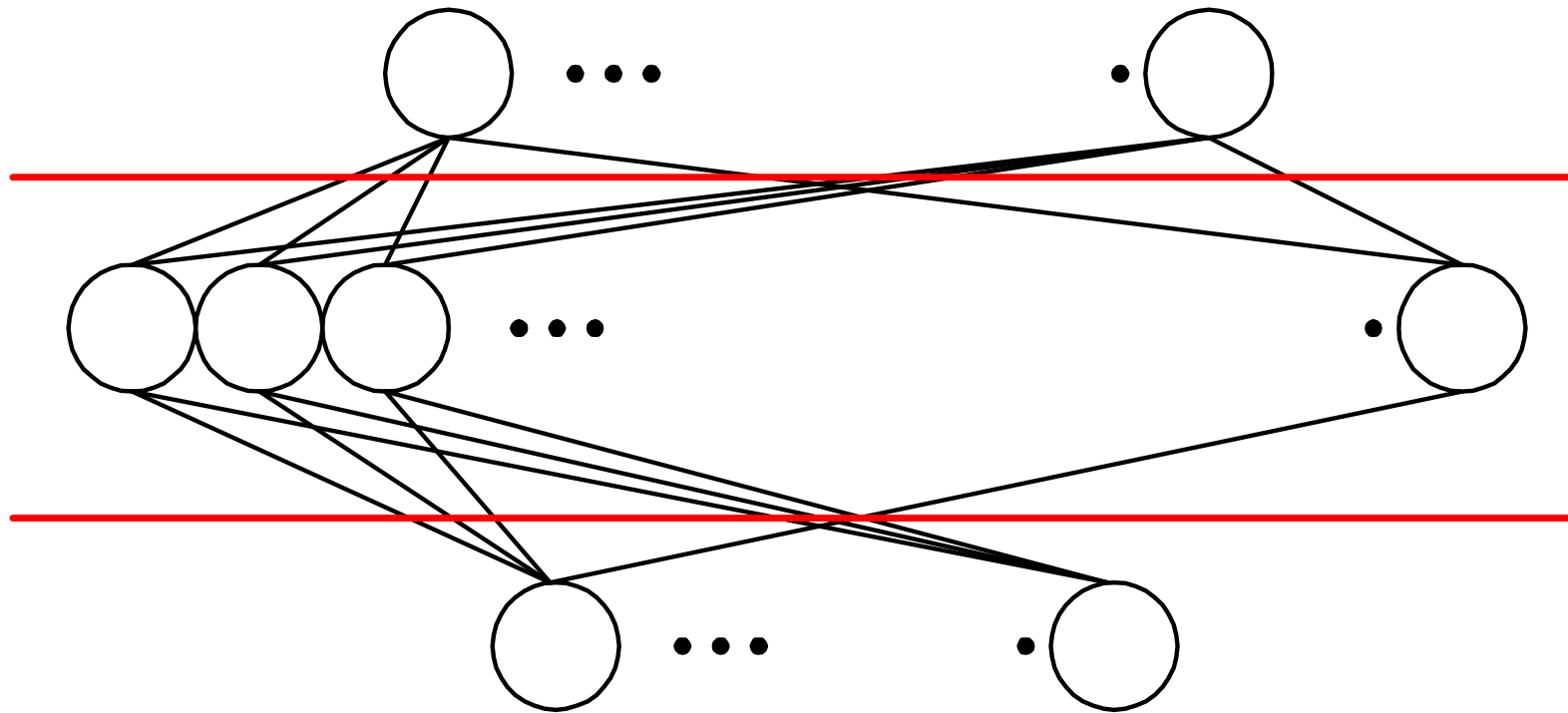


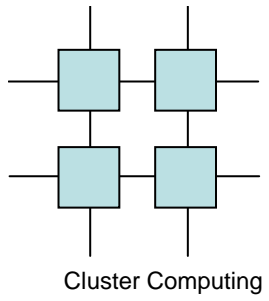
NN Training



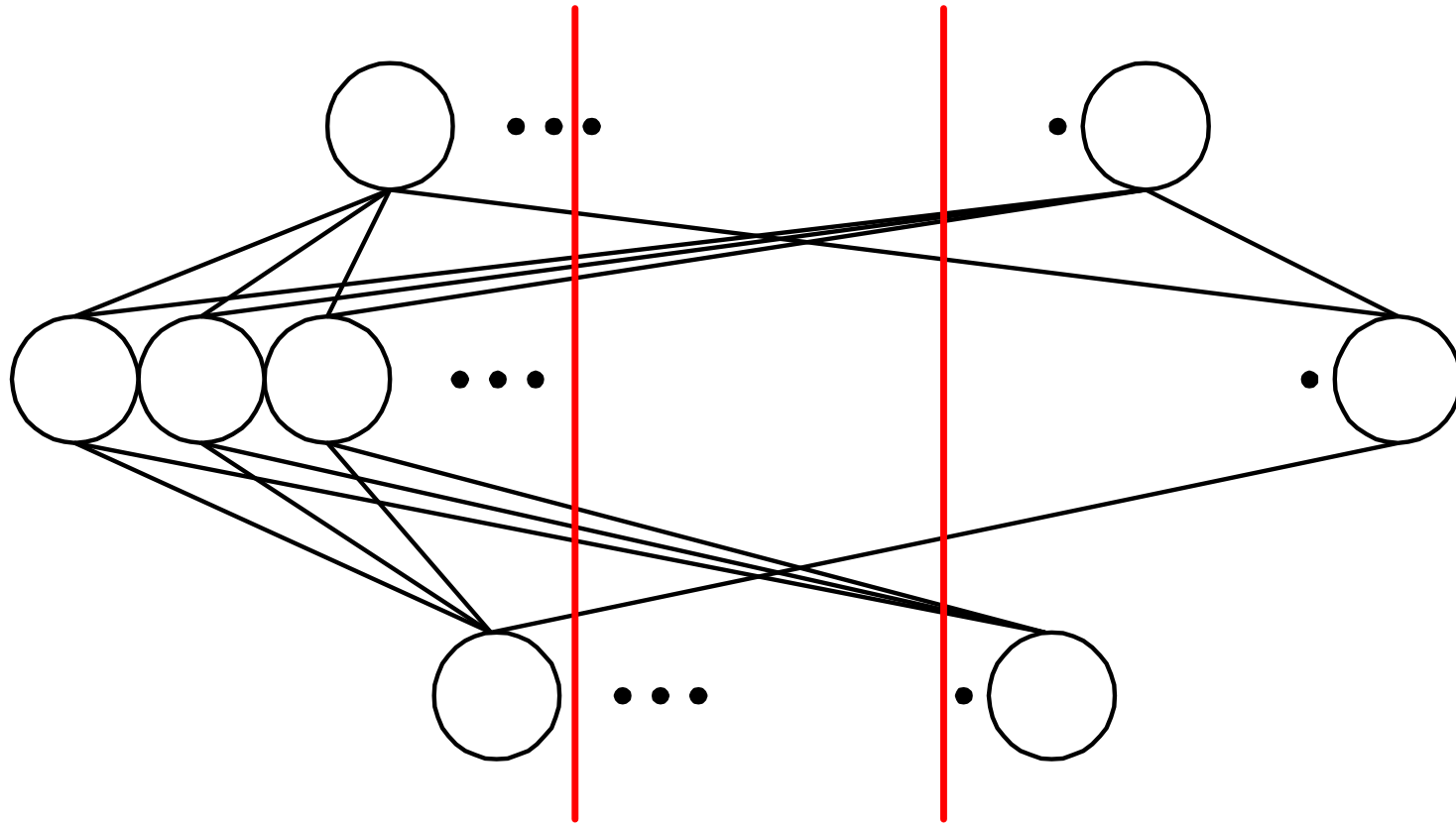


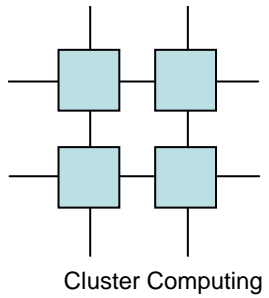
NN Training





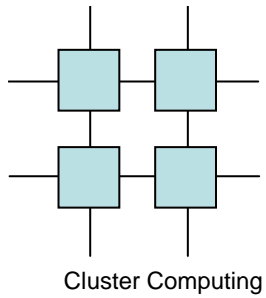
NN Training





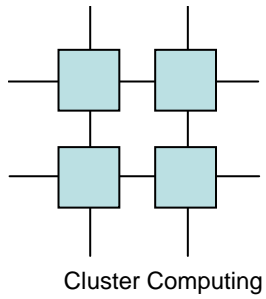
NN Training

- Can be handled by an MPI_Alltoall after each iteration

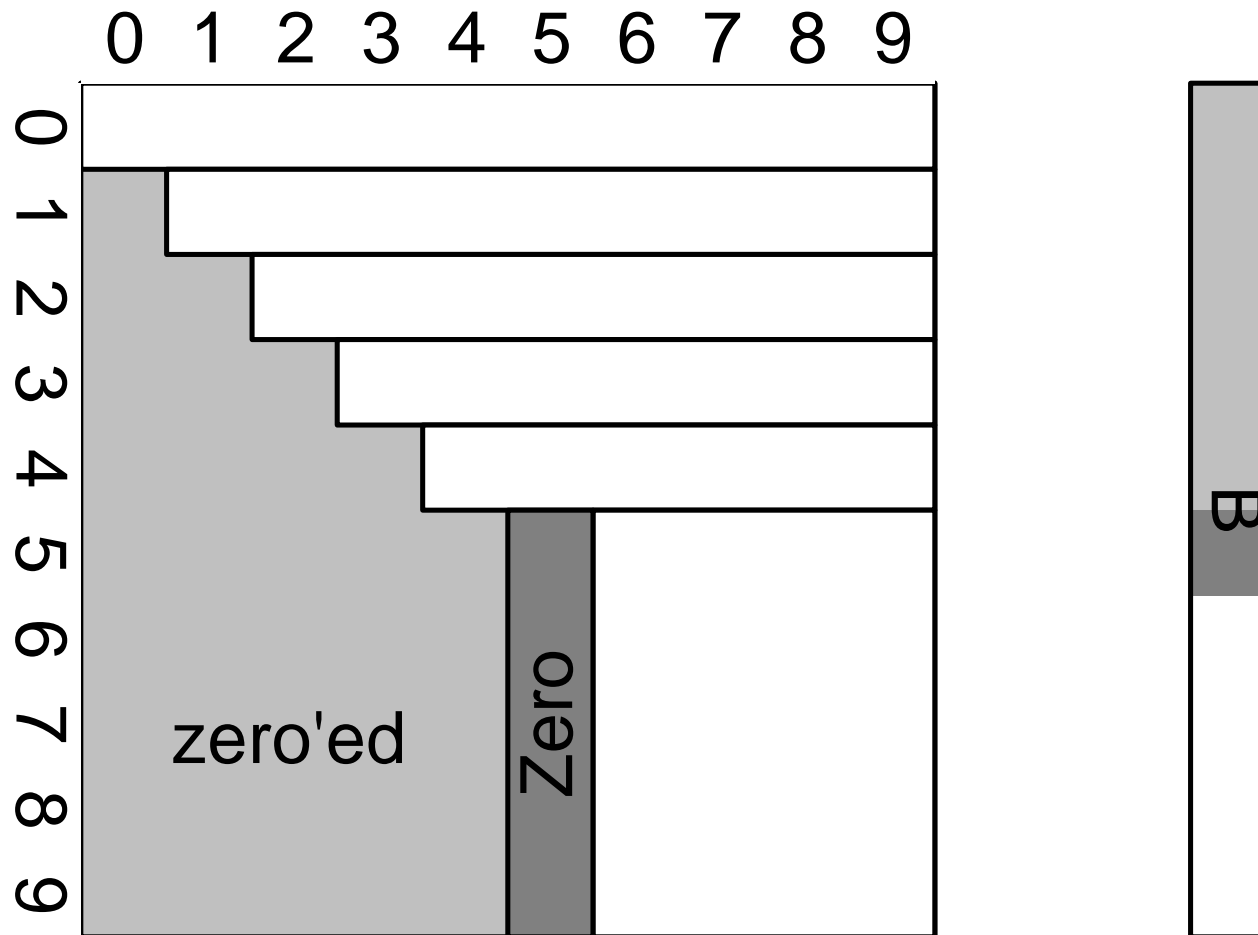


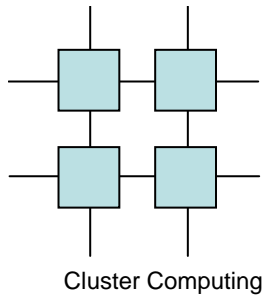
Example Gaussian Elimination

- Solve n equations with n unknown
- Frequently used in scientific applications



Gaussian Elimination





Summary

- MPI is the definite industry leader in MPP programming
- Source code portability
- Loosely synchronous execution is very different from PVM