

# MPI Tutorial

Shao-Ching Huang

**Modified by**

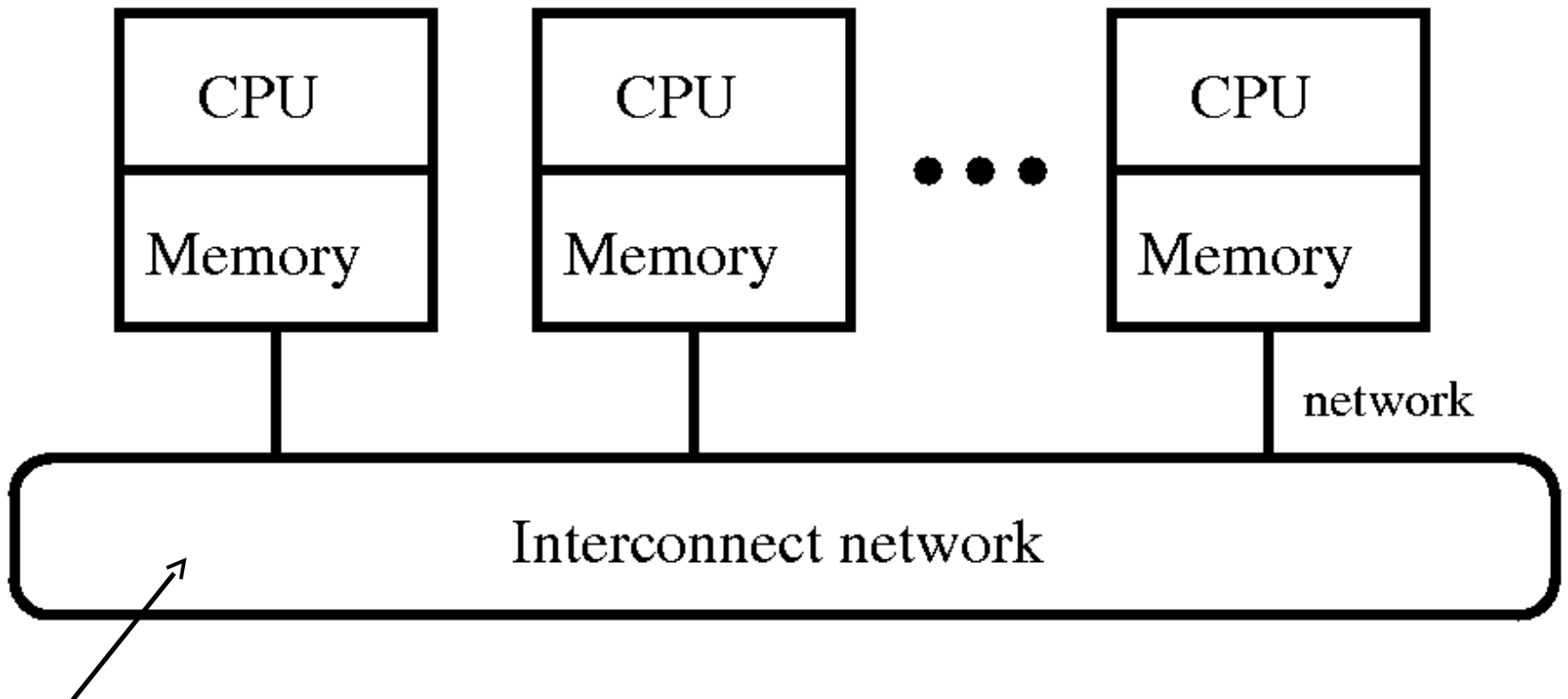
**V Krishna Nandivada, IIT Madras**

IDRE High Performance Computing Workshop

2013-02-13

# Distributed Memory

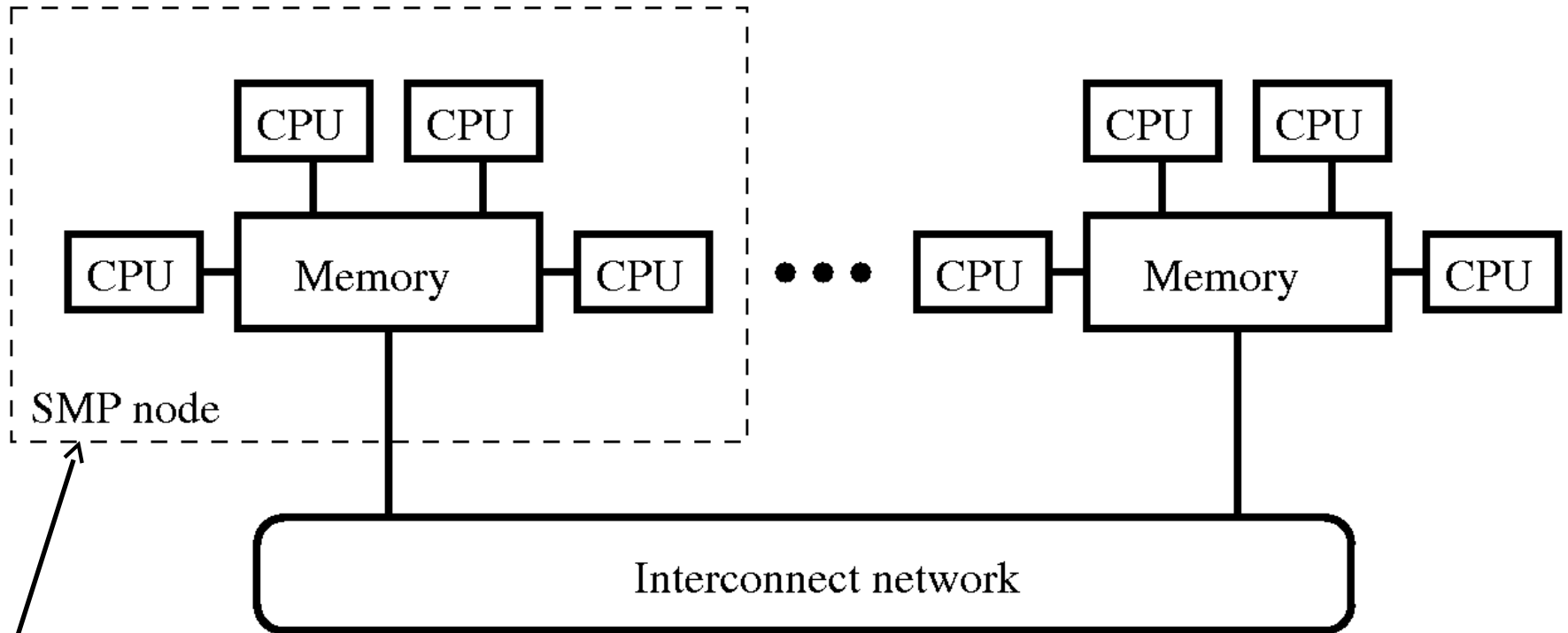
- Each CPU has its own (local) memory



This needs to be fast for parallel scalability (e.g. Infiniband, Myrinet, etc.)

# Hybrid Model

- Shared-memory within a node
- Distributed-memory across nodes



e.g. a compute node of the Hoffman2 cluster

# Today's Topics

- What is MPI
- Message passing basics
- Point to point communication
- Collective communication
- Derived data types
- Examples

# MPI = Message Passing Interface

- API for distributed-memory programming
  - parallel code that runs across multiple computers (nodes)
  - <http://www.mpi-forum.org/>
- De facto industry standard
  - available on (almost) every parallel computer for scientific computing
- Use from C/C++, Fortran, Python, R, ...
- More than 200 routines
- Using only 10 routines are enough in many cases
  - Problem dependent

# Clarification

- You can mix MPI and OpenMP in one program
- You *could* run multiple MPI processes on a single CPU
  - e.g. debug MPI codes on your laptop
  - An MPI job can span across multiple computer nodes (distributed memory)
- You *could* run multiple OpenMP threads on a single CPU
  - e.g. debug OpenMP codes on your laptop

# MPI Facts

- High-quality implementation available for free
  - Easy to install one on your desktop/laptop
  - OpenMPI: <http://www.open-mpi.org/>
  - MPICH2: <http://www.mcs.anl.gov/research/projects/mpich2/>
- Installation Steps
  - download the software
  - (assuming you already have C/C++/Fortran compilers)
  - On Mac or Linux: “configure, make, make install”

# Communicator

- A group of processes
  - processes are numbered 0,1,.. to N-1
- Default communicator
  - MPI\_COMM\_WORLD
  - contains all processes
- Query functions:
  - How many processes in total?  
MPI\_Comm\_size(MPI\_COMM\_WORLD, &nproc)
  - What is my process ID?  
MPI\_Comm\_rank(MPI\_COMM\_WORLD, &rank)
  - ...



# Hello world (C)

```
#include "mpi.h" // MPI header file
#include <stdio.h>
main(int argc, char *argv[])
{
    int np, pid;
    MPI_Init(&argc, &argv); // initialize MPI

    MPI_Comm_size(MPI_COMM_WORLD, &np);
    MPI_Comm_rank(MPI_COMM_WORLD, &pid);
    printf("N. of procs = %d, proc ID = %d\n", np, pid);

    MPI_Finalize(); // clean up
}
```

# Hello world (Fortran)

```
program hello
  Use mpi
  integer :: ierr,np,pid
  call mpi_init(ierr)
  call mpi_comm_size(MPI_COMM_WORLD,np,ierr)
  call mpi_comm_rank(MPI_COMM_WORLD,pid,ierr)
  write(*,('np = ",i2,2x,"id = ",i2)') np,pid
  call mpi_finalize(ierr)
end program hello
```

- ☞ When possible, use “use mpi”, instead of “include ‘mpif.h’”

# Error checking

- Most MPI routines returns an error code
  - C routines as the function value
  - Fortran routines in the last argument
- Examples
  - Fortran  
`MPI_Comm_rank(MPI_COMM_WORLD, myid, ierr)`
  - C/C++  
`int ierr = MPI_Comm_rank(MPI_COMM_WORLD, &myid);`

```
#include <mpi.h>
#include <stdio.h>
int main(int argc, char** argv) {
    // Initialize the MPI environment
    MPI_Init(NULL, NULL);

    // Get the number of processes
    int world_size;
    MPI_Comm_size(MPI_COMM_WORLD, &world_size);

    // Get the rank of the process
    int world_rank;
    MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);

    // Get the name of the processor
    char processor_name[MPI_MAX_PROCESSOR_NAME];
    int name_len;
    MPI_Get_processor_name(processor_name, &name_len);

    // Print off a hello world message
    printf("Hello world from processor %s, rank %d"
           " out of %d processors\n",
           processor_name, world_rank, world_size);

    // Finalize the MPI environment.
    MPI_Finalize();
}
```

```
$ mpicc hello.c -o hello
```

```
$ mpicc hello.c -o hello --showme
```

```
gcc hello.c -o hello -I/Users/nvk/software/include -L/Users/nvk/software/lib -Impi
```

```
$ mpirun -np 2 ./hello
```

```
Hello world from processor prerana.local, rank 0 out of 2 processors
```

```
Hello world from processor prerana.local, rank 1 out of 2 processors
```

```
$ mpirun -np 6 ./hello -- hostfile hosts.txt
```

# MPI built-in data types

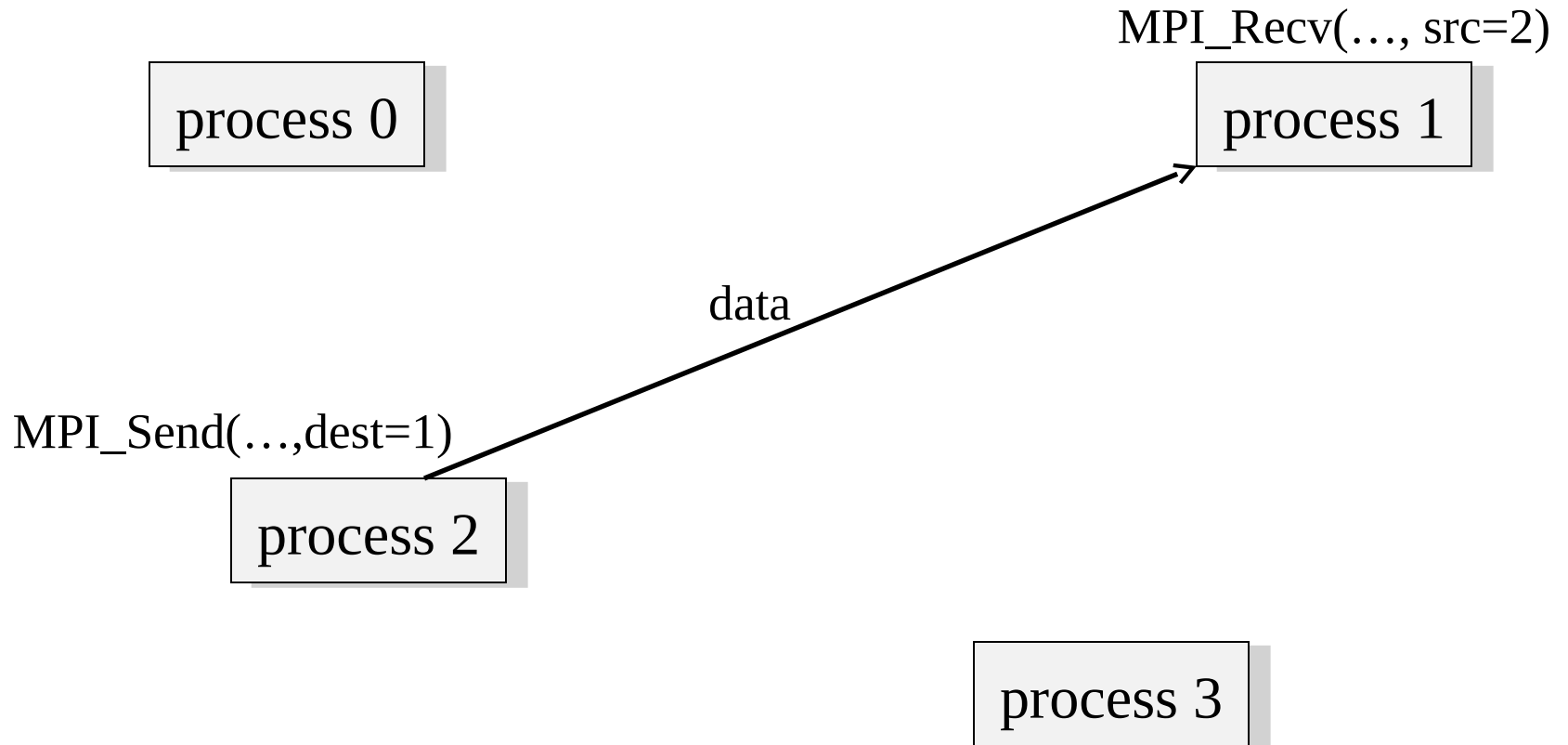
C/C++	Fortran
MPI_CHAR	MPI_CHARACTER
MPI_INT	MPI_INTEGER
MPI_FLOAT	MPI_REAL
MPI_DOUBLE	MPI_DOUBLE_PRECISION
...	...

- See MPI standard for a complete list
- New types can be (recursively) created/defined
  - based on existing types
  - called “derived data type”
  - discussed later

# Today's Topics

- Message passing basics
- Point to point communication
- Collective communication
- Derived data types
- Examples

# Point to point communication





# MPI\_Send: send data to another process

MPI\_Send(buf, count, data\_type, dest, tag, comm)

Arguments	Meanings
buf	starting address of send buffer
count	# of elements
data_type	data type of each send buffer element
dest	processor ID (rank) destination
tag	message tag
comm	communicator

Examples:

```
C/C++: MPI_Send(&x,1,MPI_INT,5,0,MPI_COMM_WORLD);  
Fortran: MPI_Send(x,1,MPI_INTEGER,5,0,MPI_COMM_WORLD,ierr)
```

# MPI\_Recv: receive data from another process

MPI\_Recv(buf, count, datatype, src, tag, comm, status)

Arguments	Meanings
buf	starting address of send buffer
count	# of elements
datatype	data type of each send buffer element
src	processor ID (rank) destination
tag	message tag
comm	communicator
status	status object (an integer array in Fortran)

## Examples:

```
C/C++: MPI_Recv(&x,1,MPI_INT,5,0,MPI_COMM_WORLD,&stat);  
Fortran: MPI_Recv(x,1,MPI_INTEGER,5,0,MPI_COMM_WORLD,stat,ierr)
```

# Notes on MPI\_Recv

- A message is received when the followings are matched:
  - Source (sending process ID/rank)
  - Tag
  - Communicator (e.g. MPI\_COMM\_WORLD)
- Wildcard values may be used:
  - MPI\_ANY\_TAG  
(don't care what the tag value is)
  - MPI\_ANY\_SOURCE  
(don't care where it comes from; always receive)

# Send/recv example (C)

- Send an integer array  $f[N]$  from process 0 to process 1

```
int f[N], src=0, dest=1;
MPI_Status status;
// ...
MPI_Comm_rank( MPI_COMM_WORLD, &rank);

if (rank == src)           // process "dest" ignores this
    MPI_Send(f, N, MPI_INT, dest, 0, MPI_COMM_WORLD);

if (rank == dest)        // process "src" ignores this
    MPI_Recv(f, N, MPI_INT, src, 0, MPI_COMM_WORLD, &status);
//...
```

# Send/recv example (F90)

- Send an integer array  $f(1:N)$  from process 0 to process 1

```
integer f(N), status(MPI_STATUS_SIZE), rank, src=0, dest=1,ierr
// ...
call MPI_Comm_rank( MPI_COMM_WORLD, rank,ierr);

if (rank == src) then                                !process "dest" ignores this
    call MPI_Send(f, N, MPI_INT, dest, 0, MPI_COMM_WORLD,ierr)
end if

if (rank == dest) then                               !process "src" ignores this
    call MPI_Recv(f, N, MPI_INT, src, 0, MPI_COMM_WORLD,
status,ierr)
end if
//...
```

# Send/Recv example (cont'd)

- Before

process 0 (send)	process 1 (recv)
f[0]=0	f[0]=0
f[1]=1	f[1]=0
f[2]=2	f[2]=0

- After

process 0 (send)	process 1 (recv)
f[0]=0	f[0]=0
f[1]=1	f[1]=1
f[2]=2	f[2]=2

## Ping-Pong

```
int ping_pong_count = 0;
int partner_rank = (my_rank + 1) % 2;
while (ping_pong_count < PING_PONG_LIMIT) {
    if (my_rank == ping_pong_count % 2) {
        // Increment the ping pong count before you send it
        ping_pong_count++;
        MPI_Send(&ping_pong_count, 1, MPI_INT, partner_rank, 0,
                MPI_COMM_WORLD);
        printf("%d sent and incremented ping_pong_count "
               "%d to %d\n", my_rank, ping_pong_count,
               partner_rank);
    } else {
        MPI_Recv(&ping_pong_count, 1, MPI_INT, partner_rank, 0,
                MPI_COMM_WORLD, MPI_STATUS_IGNORE);
        printf("%d received ping_pong_count %d from %d\n",
               my_rank, ping_pong_count, partner_rank);
    }
}
```

```
$ mpirun -np 2 ./ping_pong
0 sent and incremented ping_pong_count 1 to 1
0 received ping_pong_count 2 from 1
0 sent and incremented ping_pong_count 3 to 1
0 received ping_pong_count 4 from 1
0 sent and incremented ping_pong_count 5 to 1
0 received ping_pong_count 6 from 1
0 sent and incremented ping_pong_count 7 to 1
1 received ping_pong_count 1 from 0
1 sent and incremented ping_pong_count 2 to 0
1 received ping_pong_count 3 from 0
1 sent and incremented ping_pong_count 4 to 0
1 received ping_pong_count 5 from 0
1 sent and incremented ping_pong_count 6 to 0
1 received ping_pong_count 7 from 0
1 sent and incremented ping_pong_count 8 to 0
1 received ping_pong_count 9 from 0
1 sent and incremented ping_pong_count 10 to 0
0 received ping_pong_count 8 from 1
0 sent and incremented ping_pong_count 9 to 1
0 received ping_pong_count 10 from 1
```



## Ringa Ringa Roses!

```
int token;
if (world_rank != 0) {
    MPI_Recv(&token, 1, MPI_INT, world_rank - 1, 0,
            MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    printf("Process %d received token %d from process %d\n",
        world_rank, token, world_rank - 1);
} else {
    // Set the token's value if you are process 0
    token = -1;
}
MPI_Send(&token, 1, MPI_INT, (world_rank + 1) % world_size,
        0, MPI_COMM_WORLD);

// Now process 0 can receive from the last process.
if (world_rank == 0) {
    MPI_Recv(&token, 1, MPI_INT, world_size - 1, 0,
            MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    printf("Process %d received token %d from process %d\n",
        world_rank, token, world_size - 1);
}
```

```
$ mpirun -np 4 ./ring
```

```
Process 1 received token -1 from process 0
```

```
Process 2 received token -1 from process 1
```

```
Process 3 received token -1 from process 2
```

```
Process 0 received token -1 from process 3
```

# Blocking

- Function call does not return until the communication is complete
- MPI\_Send and MPI\_Recv are blocking calls
- Calling order matters
  - it is possible to wait indefinitely, called “deadlock”
  - improper ordering results in serialization (loss of performance)

# Deadlock

- This code always works:

```
MPI_Comm_rank(comm, &rank);

if (rank == 0) {
    MPI_Send(sendbuf, cnt, MPI_INT, 1, tag, comm);
    MPI_Recv(recvbuf, cnt, MPI_INT, 1, tag, comm, &stat);
} else { // rank==1
    MPI_Recv(recvbuf, cnt, MPI_INT, 0, tag, comm, &stat);
    MPI_Send(sendbuf, cnt, MPI_INT, 0, tag, comm);
}
```

# Deadlock

- This code deadlocks:

```
MPI_Comm_rank(comm, &rank);

if (rank == 0) {
    MPI_Recv(recvbuf, cnt, MPI_INT, 1, tag, comm, &stat);
    MPI_Send(sendbuf, cnt, MPI_INT, 1, tag, comm);
} else { /* rank==1 */
    MPI_Recv(recvbuf, cnt, MPI_INT, 0, tag, comm, &stat);
    MPI_Send(sendbuf, cnt, MPI_INT, 0, tag, comm);
}
```

reason: MPI\_Recv on process 0 waits indefinitely and never returns.

# Non-blocking

- Function call returns immediately, without completing data transfer
  - Only “starts” the communication (without finishing)
  - MPI\_Isend and MPI\_Irecv
  - Need an additional mechanism to ensure transfer completion (MPI\_Wait)
- Avoid deadlock
- Possibly higher performance
- Examples: MPI\_Isend & MPI\_Irecv

# MPI\_Isend

MPI\_Isend(buf, count, datatype, dest, tag, comm, request )

- Similar to MPI\_Send, except the last argument “request”
- Typical usage:

```
MPI_Request request_X, request_Y;  
MPI_Isend(..., &request_X);  
MPI_Isend(..., &request_Y);  
  
//... some ground-breaking computations ...  
  
MPI_Wait(&request_X, ...);  
MPI_Wait(&request_Y, ...);
```

# MPI\_Irecv

MPI\_Irecv(buf, count, datatype, src, tag, comm, request )

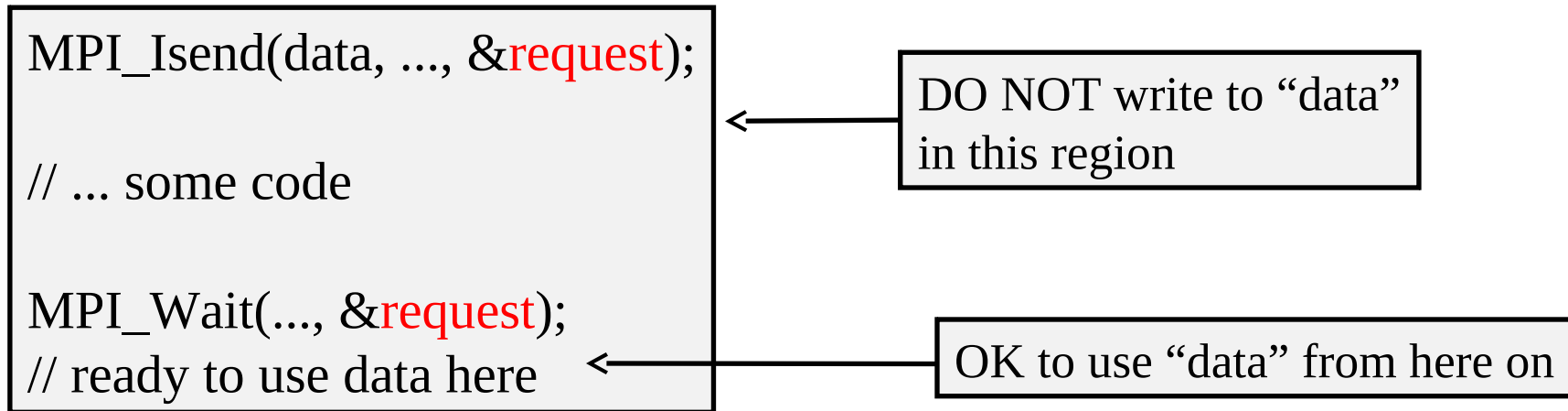
- Similar to MPI\_Recv, except the last argument “request”
- Typical usage:

```
MPI_Request request_X, request_Y;  
MPI_Irecv(..., &request_X);  
MPI_Irecv(..., &request_Y);  
  
//... more ground-breaking computations ...  
  
MPI_Wait(&request_X, ...);  
MPI_Wait(&request_Y, ...);
```



# Caution about MPI\_Isend and MPI\_Irecv

- The sending process should not access the send buffer until the send completes



# MPI\_Wait

`MPI_Wait(MPI_Request, MPI_Status)`

- Wait for an `MPI_Isend/recv` to complete
- Use the same “request” used in an earlier `MPI_Isend` or `MPI_Irecv`
- If they are multiple requests, one can use  
`MPI_Waitall(count, request[], status[]);`  
`request[]` and `status[]` are arrays.

# Other variants of MPI Send/Recv

- MPI\_Sendrecv
  - send and receive in one call
- Mixing blocking and non-blocking calls
  - e.g. MPI\_Isend + MPI\_Recv
- MPI\_Bsend
  - buffered send
- MPI\_Ibsend
- ... (see MPI standard for more)

## MPI\_Status

- o- Can ask MPI to read data from any source (MPI\_ANY\_SOURCE).
- o- Can ask MPI to read data of any tag (MPI\_ANY\_TAG)
- o- MPI\_Status structure has three fields:
  - rank of the sender (stat.MPI\_SOURCE)
  - tag of the message (stat.MPI\_TAG)
  - Length of the message.
    - Use MPI\_Get\_count (MPI\_Status\* status, MPI\_Datatype datatype, int\* count)
- o- Use MPI\_Probe (int source, int tag, MPI\_Comm comm, MPI\_Status\* status) to read the status. It is a pre-receive query.
- o- MPI\_IProbe — the non-blocking version
- o- Use MPI\_Get\_count to know the exact length.
- o- Advantage?

# Today's Topics

- Message passing basics
  - communicators
  - data types
- Point to point communication
- **Collective communication**
- Derived data types
- Examples

## Different Send/Recv functions

MPI\_Send - blocks till the buffer is ready to use.

MPI\_Isend - No blocking.

MPI\_Ssend - The recv request has been posted at the target  
- Can be used for synchronization.

MPI\_Recv

MPI\_Irecv

MPI\_Srecv

MPI\_Test (MPI\_Request \*req, int \*flag, MPI\_Status \*status)

- Tests for the completion of the request.
- \*flag is set to true if the operation is completed.

## Blocked Send/Receive

MPI\_Send - blocks till the message is added to internal buffer.

- What if buffer is full? - Blocks.

Solution? Create a buffer and use buffered send.

MPI\_Buffer\_attach:

MPI\_BSend:

```
MPI_Buffer_attach( b, n*sizeof(double) + MPI_BSEND_OVERHEAD );  
for (i=0; i<m; i++) {  
    MPI_Bsend( buf, n, MPI_DOUBLE, ... );  
}
```

No MPI\_BRecv. Why?

## Code Walkthrough

```
while (*){  
    compute data element to send;  
  
    for (each element e){  
        MPI_Send(e, ..., sending-neighbor)  
    }  
  
    MPI_Probe(receiving-neighbor 0, MPI_COMM_WORLD, &status);  
    MPI_Get_Count(&status, MPI_BYTE, &len);  
    MPI_Recv(buf, len, MPI_BYTE, ..., r);  
}
```

Any Issues with this code?

Suggested reading:

<http://mpitutorial.com/tutorials/point-to-point-communication-application-random-walk/>



# Synchronization

MPI\_Barrier(MPI\_Comm communicator) -

no processes in the communicator can pass the barrier until all of them call the function.

- Make sure that either no process calls MPI\_Barrier, or every process calls MPI\_Barrier.

Q: How to implement MPI\_Barrier?

# Collective communication

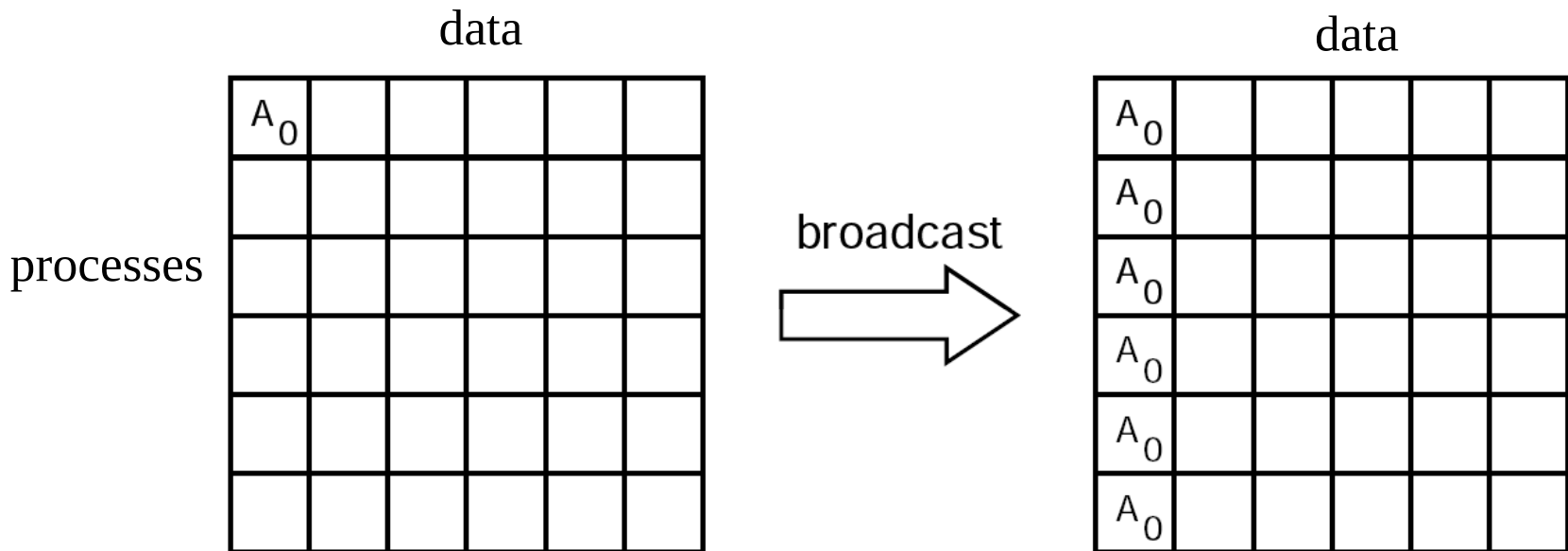
- One to all
  - MPI\_Bcast, MPI\_Scatter
- All to one
  - MPI\_Reduce, MPI\_Gather
- All to all
  - MPI\_Alltoall

Implicit Synchronization

# MPI\_Bcast

MPI\_Bcast(buffer, count, datatype, root, comm )

Broadcasts a message from “root” process to all other processes in the same communicator



# MPI\_Bcast Example

- Broadcast 100 integers from process “3” to all other processes

C/C++

```
MPI_Comm comm;  
int array[100];  
//...  
MPI_Bcast( array, 100, MPI_INT, 3, comm);
```

Fortran

```
INTEGER comm  
integer array(100)  
//...  
call MPI_Bcast( array, 100, MPI_INTEGER, 3, comm,ierr)
```

## MPI\_Bcast Vs MPI\_Send+MPI\_Recv

```
if (world_rank == root) {  
    // If we are the root process, send our data to everyone  
    for (int i = 0; i < world_size; i++) {  
        if (i != world_rank) {MPI_Send(data, count, datatype, i, 0, communicator);}  
    }  
} else { // If we are a receiver process, receive the data from the root  
    MPI_Recv(data, count, datatype, root, 0, communicator, MPI_STATUS_IGNORE);  
}
```

Vs

```
MPI_Bcast(data, num_elements, MPI_INT, 0, MPI_COMM_WORLD);
```

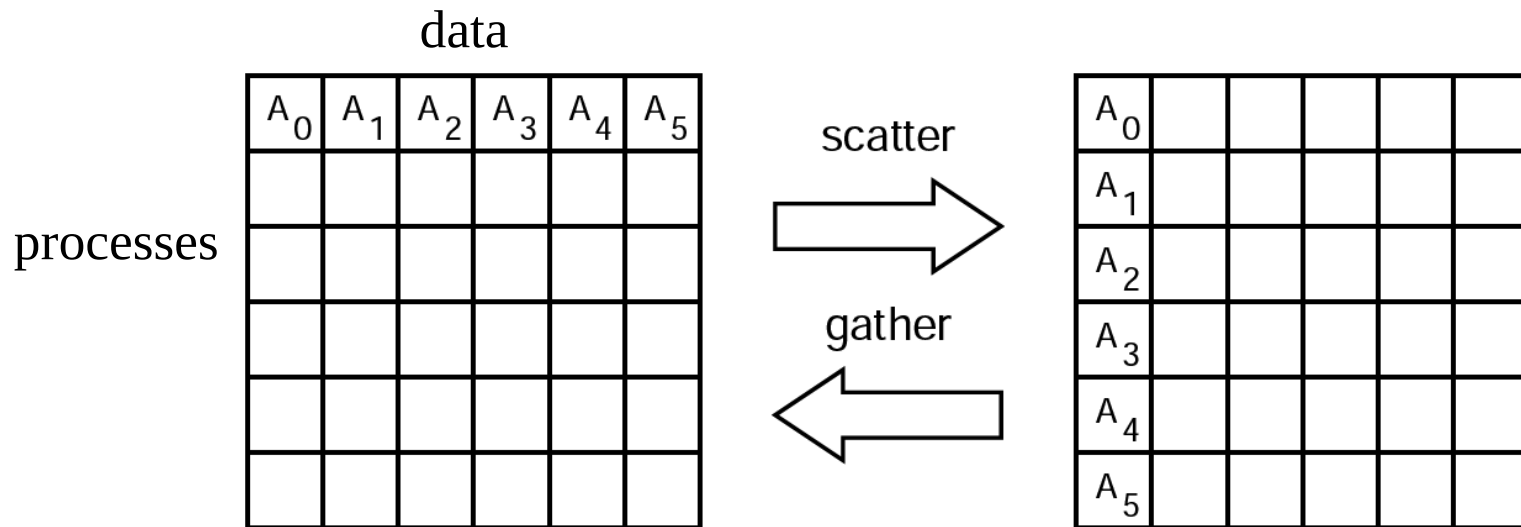
*A simple way to speedup:*

*- Use tree based communication.*

# MPI\_Gather & MPI\_Scatter

MPI\_Gather (sbuf, scnt, stype, rbuf, rcnt, rtype, root, comm )

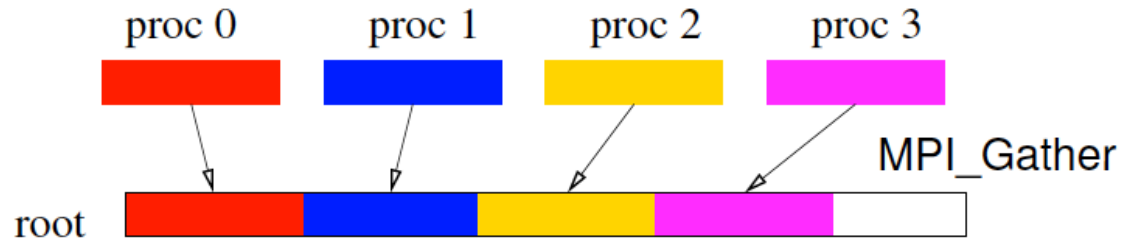
MPI\_Scatter(sbuf, scnt, stype, rbuf, rcnt, rtype, root, comm )



☞ When gathering, make sure the root process has big enough memory to hold the data (especially when you scale up the problem size).

Non blocking versions: MPI\_IGather and MPI\_IScatter

# MPI\_Gather Example



```
MPI_Comm comm;
int np, myid, sendarray[N], root;
double *rbuf;
MPI_Comm_size( comm, &np);    // # of processes
MPI_Comm_rank( comm, &myid); // process ID
if (myid == root)             // allocate space on process root
    rbuf = new double [np*N];

MPI_Gather( sendarray, N, MPI_INT, rbuf, N, MPI_INT,
            root, comm);
```

## Scatter Gather Example

```
if (world_rank == 0) {rand_nums = create_rand_nums(elements_per_proc * world_size);}

// Create a buffer that will hold a subset of the random numbers
float *sub_rand_nums = malloc(sizeof(float) * elements_per_proc);

// Scatter the random numbers to all processes
MPI_Scatter(rand_nums, elements_per_proc, MPI_FLOAT, sub_rand_nums,
            elements_per_proc, MPI_FLOAT, 0, MPI_COMM_WORLD);

// Compute the average of your subset
float sub_avg = compute_avg(sub_rand_nums, elements_per_proc);
// Gather all partial averages down to the root process
float *sub_avgs = NULL;

if (world_rank == 0) {sub_avgs = malloc(sizeof(float) * world_size);}

MPI_Gather(&sub_avg, 1, MPI_FLOAT, sub_avgs, 1, MPI_FLOAT, 0, MPI_COMM_WORLD);

// Compute the total average of all numbers.
if (world_rank == 0) {float avg = compute_avg(sub_avgs, world_size);}
```



# Variations of MPI\_Gather/Scatter

- Variable data size
  - MPI\_Gatherv
  - MPI\_Scatterv
- Gather + broadcast (in one call)
  - MPI\_Allgather All processes send same amount of data.
  - MPI\_Allgatherv Processes may send variable amount of data.

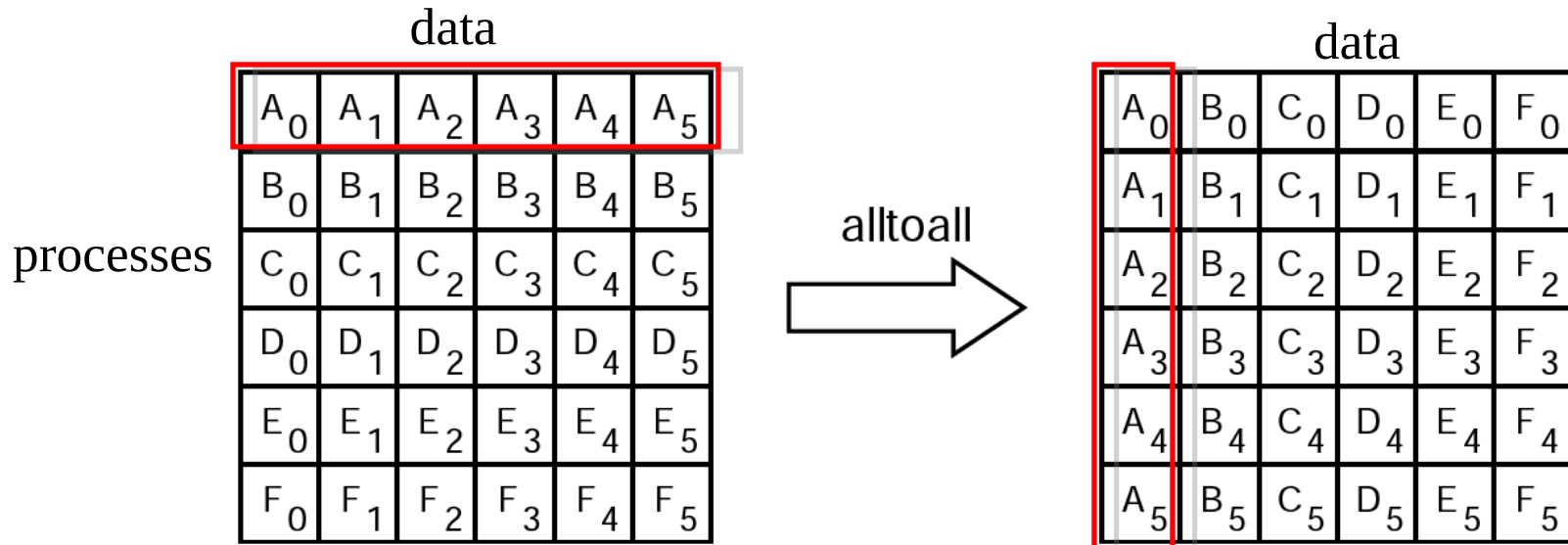
```
// Gather all partial averages down to all the processes
float *sub_avgs = (float *)malloc(sizeof(float) * world_size);
MPI_Allgather(&sub_avg, 1, MPI_FLOAT, sub_avgs, 1, MPI_FLOAT,
             MPI_COMM_WORLD);
```

```
// Everyone: compute the total average of all numbers.
float avg = compute_avg(sub_avgs, world_size);
```

# MPI\_Alltoall

MPI\_Alltoall( send\_buf, send\_count, send\_data\_type,  
recv\_buf, recv\_count, recv\_data\_type, comm)

The j-th block send\_buf from process i is received by process j and is placed in the i-th block of rbuf:



# MPI\_Reduce

count

MPI\_Reduce (send\_buf, recv\_buf, ██████████, OP, root, comm)

- Apply operation OP to send\_buf from all processes and return result in the recv\_buf on process “root”.
- Some predefined operations:

Operations (OP)	Meaning
MPI_MAX	maximum value
MPI_MIN	minimum value
MPI_SUM	sum
MPI_PROD	products
...	

(see MPI standard for more predefined reduce operations)

# MPI\_Reduce example

- Parallel vector inner product:

$$a \leftarrow x \cdot y$$

```
// loc_sum = local sum
float loc_sum = 0.0;           // probably should use double
for (i = 0; i < N; i++)
    loc_sum += x[i] * y[i];

// sum = global sum
MPI_Reduce(&loc_sum, &sum, 1, MPI_FLOAT, MPI_SUM,
           root, MPI_COMM_WORLD);
```

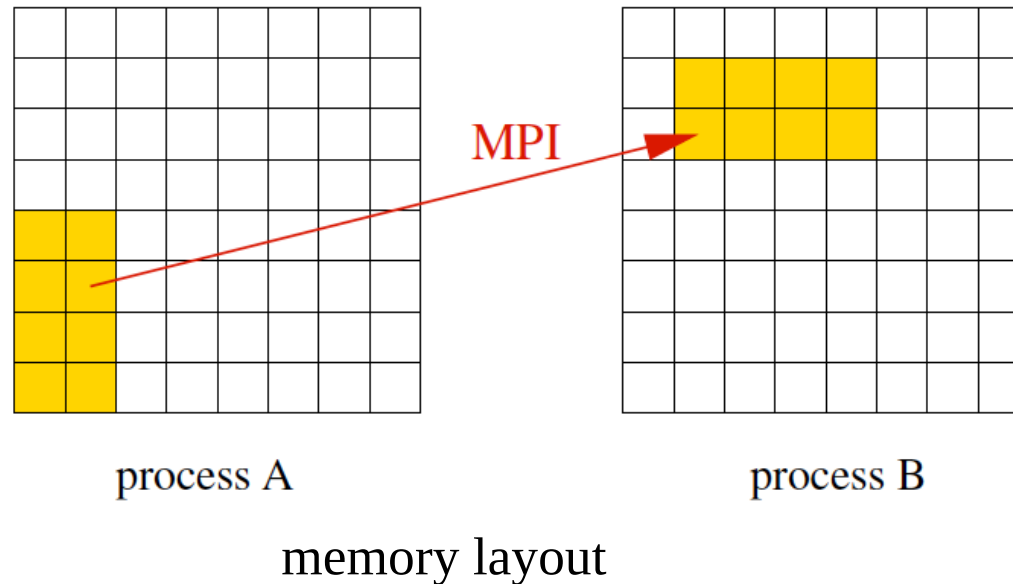
Variation: MPI\_Allreduce — results sent to all processes.

# Today's Topics

- Message passing basics
  - communicators
  - data types
- Point to point communication
- Collective communication
- Derived data types
- Examples

# Derived Data Type

- Define data objects of various sizes and shapes (memory layout)
- Example
  - The send and recv ends have same data size but different memory layouts



# Data Type Constructors

<b>Constructors</b>	<b>Usage</b>
Contiguous	contiguous chunk of memory
Vector	strided vector
Hvector	strided vector in bytes
Indexed	variable displacement
Hindexed	variable displacement in bytes
Struct	fully general data type

# MPI\_Type\_contiguous

MPI\_Type\_contiguous(count, old\_type, newtype)

- Define a contiguous chunk of memory
- Example – a memory block of 10 integers

```
int a[10];  
MPI_Datatype intvec;  
MPI_Type_contiguous(10, MPI_INT, &intvec);  
MPI_Type_commit(&intvec);  
MPI_Send(a, 1, intvec, ...); /* send 1 10-int vector */
```

new type



is equivalent to

```
MPI_Send(a, 10, MPI_INT,...); /* send 10 ints */
```



# MPI\_Type\_vector

MPI\_Type\_vector(count, blocklen, stride, old\_type, newtype )

To create a strided vector (i.e. with “holes”):



```
MPI_Datatype yellow_vec;  
MPI_Type_vector(3, 4, 6, MPI_FLOAT, &yellow_vec);  
MPI_Type_commit(&yellow_vec);
```

# Commit and Free

- A new type needs to be committed before use  
`MPI_Type_commit(datatype)`
- Once committed, it can be used many times
- To destroy a data type, freeing the memory:  
`MPI_Type_free(datatype)`

☞ If you repeatedly (e.g. in iterations) create MPI types, make sure you free them when they are no longer in use. Otherwise you may have memory leak.

# Rank Computaion

Problem: Each process has a key. Goal: rank order the processes based on their key.

Two options: gather keys from all processes to the root.

- Root sorts the keys.

- Root informs each process of its rank.

- Each process shares its keys to every other process.

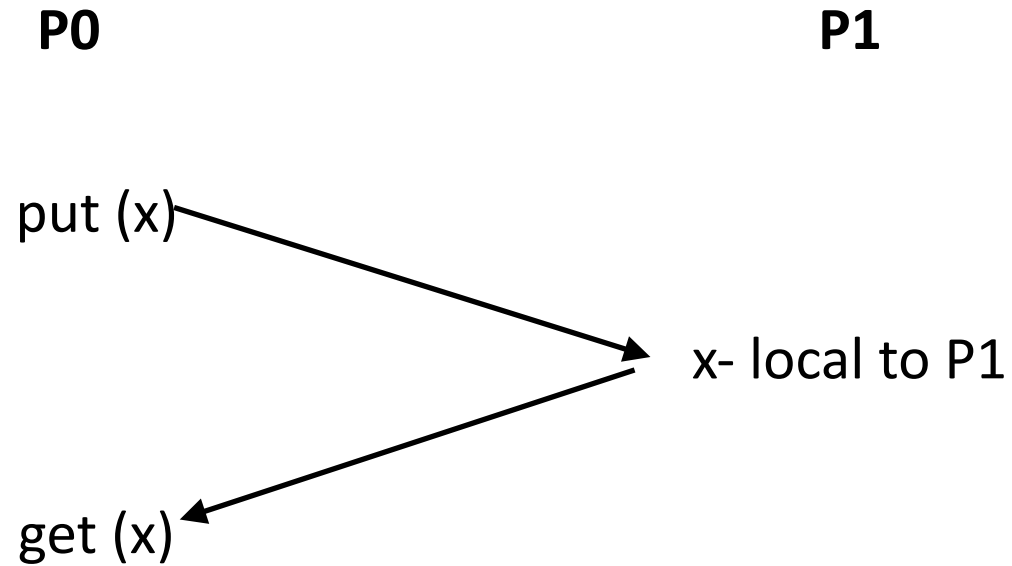
- Each process sorts the keys and finds own rank.

# Motivating Remote Memory Access

Consider the following two scenarios, in the context of a graph  $G$  (assume that each node know the rank of its neighbor)

- Under some condition neighbors exchange data
  - Handle multiple neighbors.
- Under some condition, each node  $n_1$  wants some info from its neighbor  $n_2$ .
  - the condition is evaluated locally at  $n_1$ .
  - the condition evaluation needs data from  $n_2$ .

## One side Communication (MPI 2)



P1 need not participate!

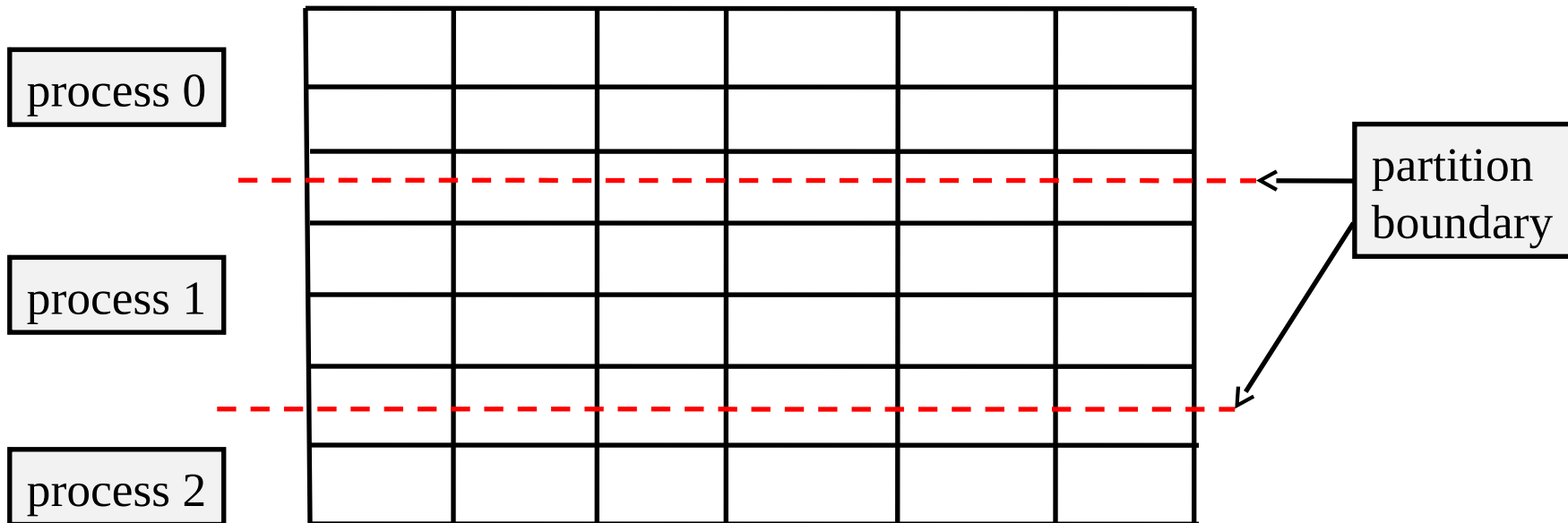
# Examples

- Poisson equation
- Fast Fourier Transform (FFT)

# Poisson equation (or any elliptic PDE)

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = R(x, y)$$

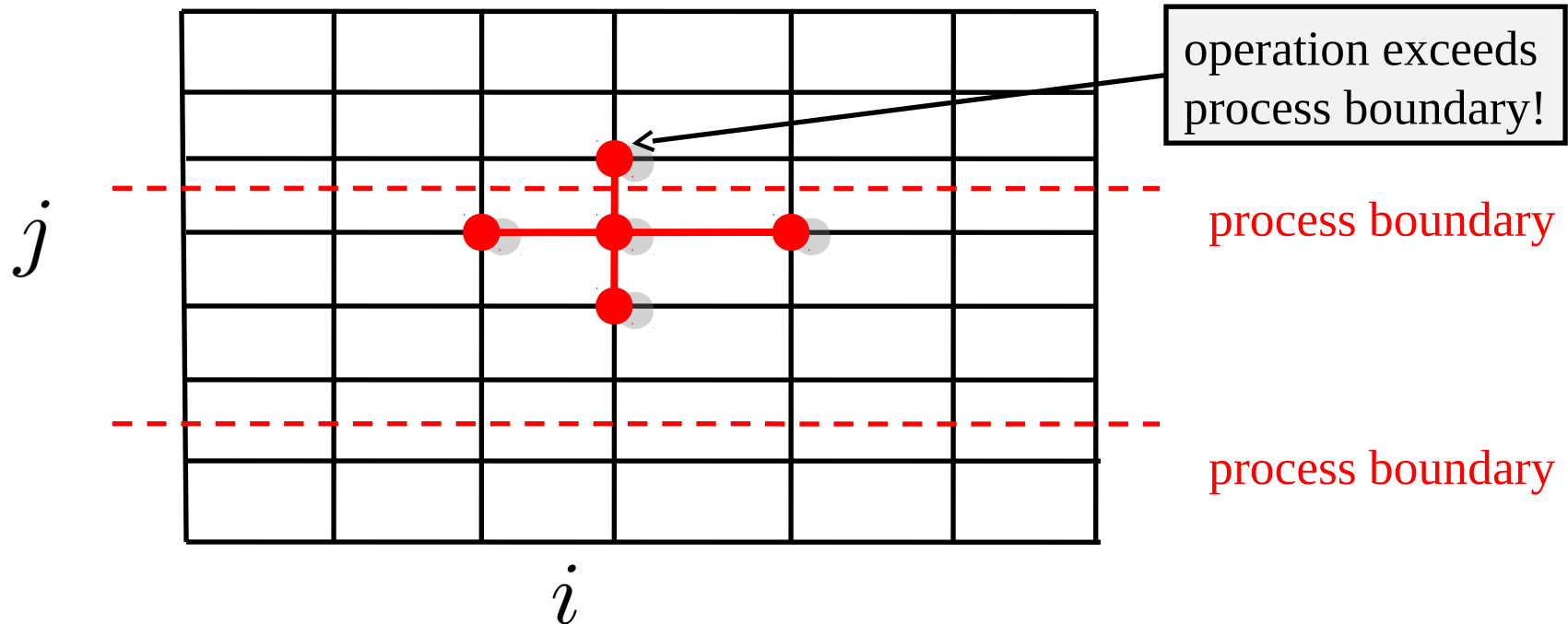
Computational grid:



# Poisson equation

Jacobi iterations (as an example)

$$f_{i,j}^{k+1} = \frac{1}{4} (f_{i+1,j}^k + f_{i-1,j}^k + f_{i,j+1}^k + f_{i,j-1}^k)$$

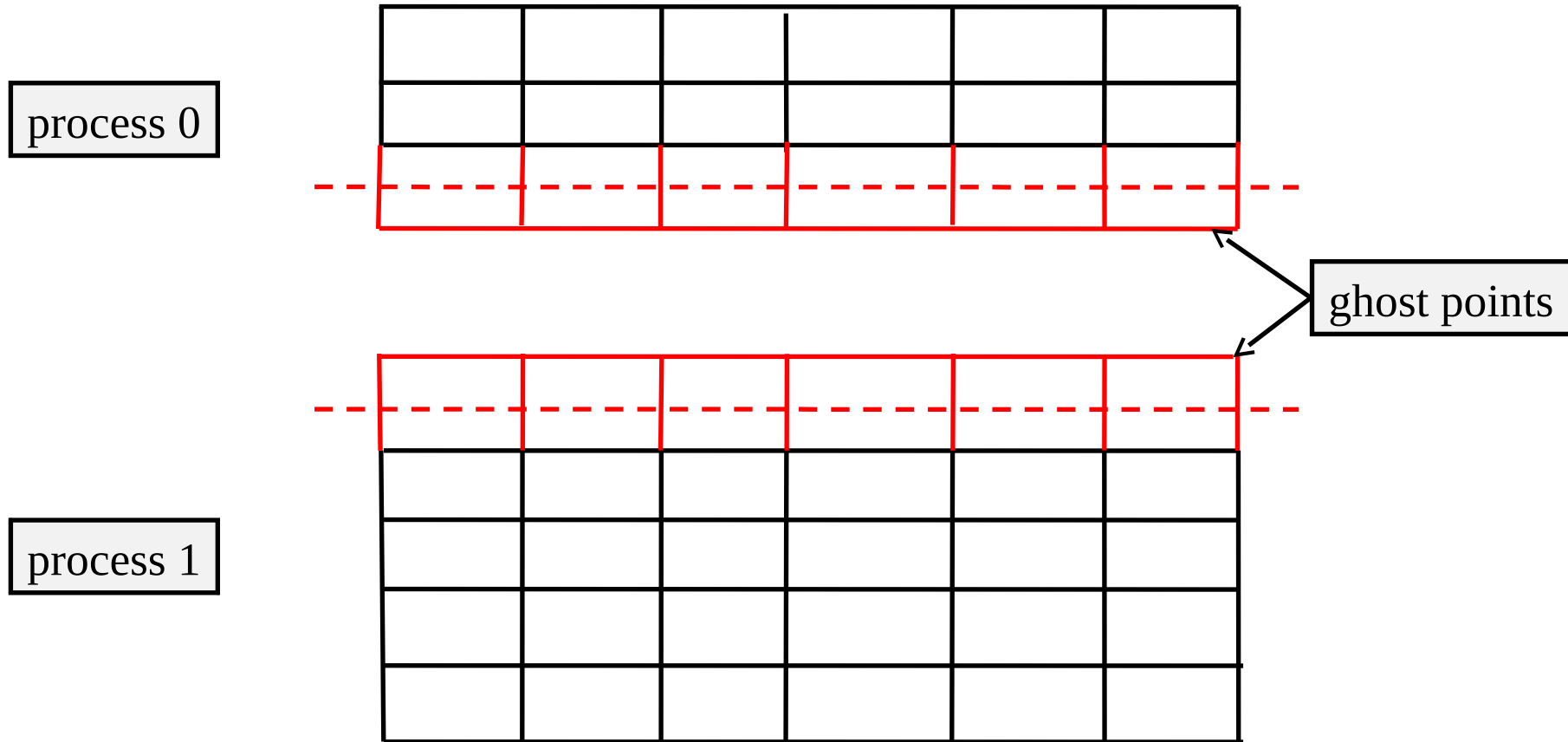


One solution is to introduce “ghost points” (see next slide)



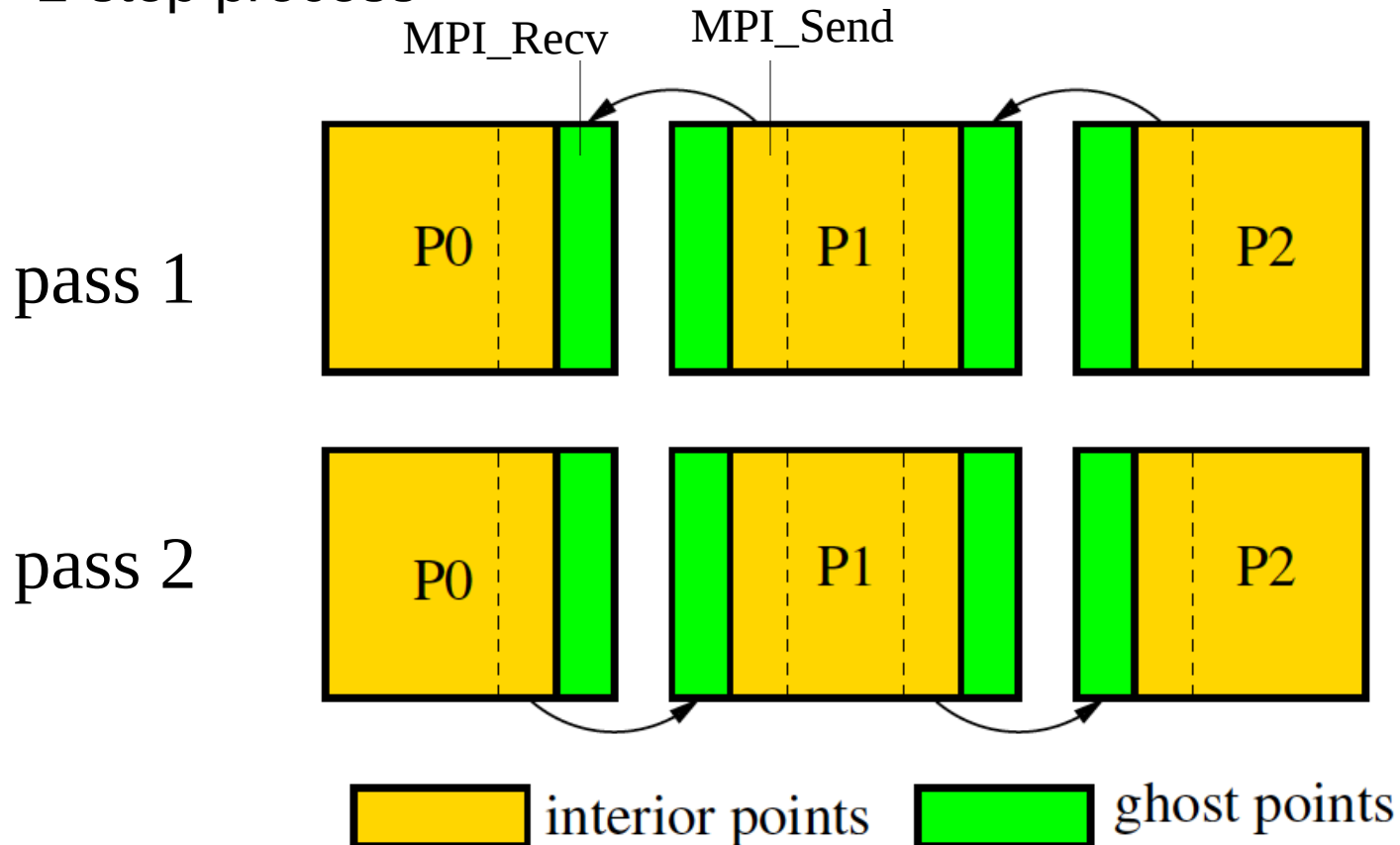
# Ghost points

Redundant copy of data held on neighboring processes



# Update ghost points in one iteration

- 2-step process



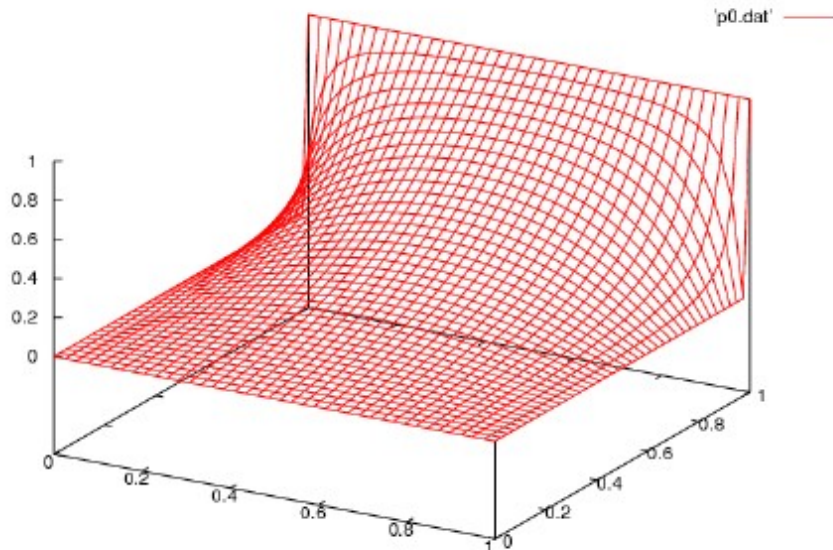
- Repeat for many iterations until convergence

# Poisson solution

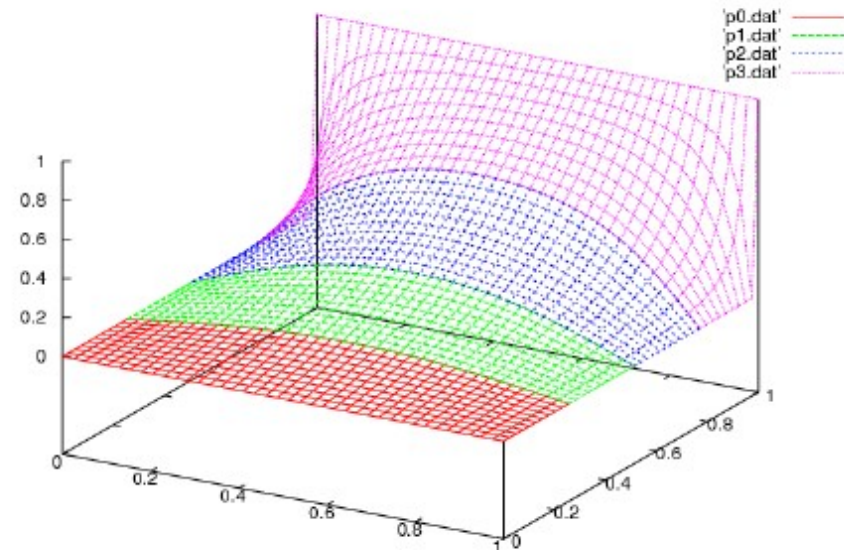
Dirichlet boundary conditions

$$\phi(x, 1) = 1, \phi(x, 0) = \phi(0, y) = \phi(1, y) = 0$$

1 process

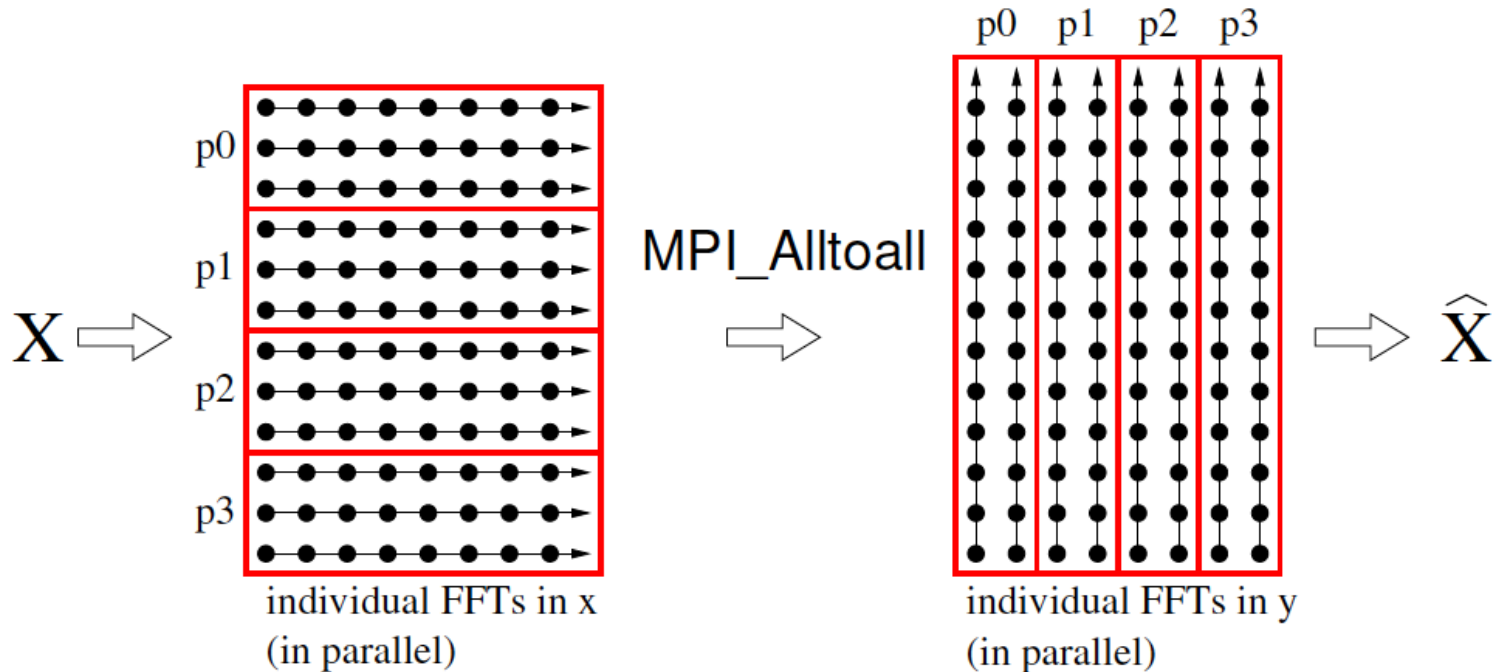


4 processes



# “Parallel” FFT

$$\hat{X}(k_x, k_y) = \sum \sum X(x, y) \exp^{-i(k_x x + k_y y)}$$



Doing multiple (sequential) FFT in parallel

# Timing

- MPI\_Wtime
  - elapsed wall-clock time in seconds
  - Note: wall-clock time is not CPU time
- Example

```
double t1,t2;  
t1 = MPI_Wtime();  
//... some heavy work ...  
t2 = MPI_Wtime();  
printf("elapsed time = %f seconds\n", t2-t1);  
Parallel
```

# How to run an MPI program

- Compile

C:        **mpicc** foo.c

C++:     **mpicxx** foo.cpp

F90:     **mpif90** foo.f90

☞ mpicc, mpicxx and mpif90 are sometimes called the MPI compilers (wrappers)

- Run

**mpiexec** -n 4 [options] a.out

- The options in mpiexec are implementation dependent
- Check out the user's manual

# Summary

- MPI for distributed-memory programming
  - works on shared-memory parallel computers too
- Communicator
  - a group of processes, numbered 0,1,...,to N-1
- Data Types
  - derived types can be defined based on built-in ones
- Point-to-point Communication
  - blocking (Send/Recv) and non-blocking (Isend/Irecv)
- Collective Communication
  - gather, scatter, alltoall

# Online Resources

- MPI-1 standard

<http://www.mpi-forum.org/docs/mpi-11-html/mpi-report.html>

- MPI-2 standard

<http://www.mpi-forum.org/docs/mpi-20-html/mpi2-report.html>

- MPI-3 standard

<http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf>

MPI Tutorial

[www.mpitutorial.com](http://www.mpitutorial.com)