



# An introduction to MPI



# MPI is a **Library** for Message-Passing

- Not built in to compiler
- Function calls that can be made from any compiler, many languages
- Just link to it
- Wrappers: `mpicc`, `mpif77`



## C

```
#include <stdio.h>
#include <mpi.h>

int main(int argc, char **argv) {
    int rank, size;
    int ierr;

    ierr = MPI_Init(&argc, &argv);

    ierr = MPI_Comm_size(MPI_COMM_WORLD, &size);
    ierr = MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    printf("Hello from task %d of %d, world!\n", rank, size);

    MPI_Finalize();

    return 0;
}
```

## Fortran

```
program hellompiworld
include "mpif.h"

integer rank, size
integer ierr

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

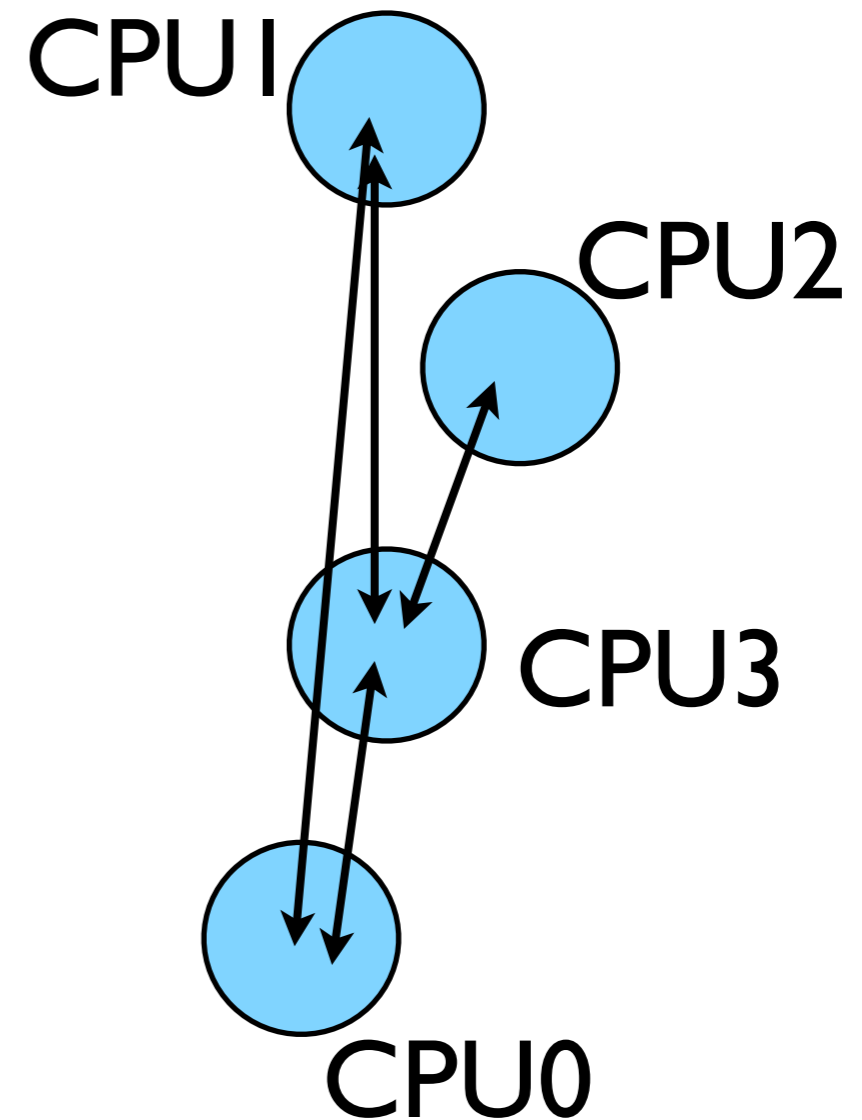
print *, "Hello from task ", rank, " of ", size, ", world!"

call MPI_FINALIZE(ierr)

return
end
```

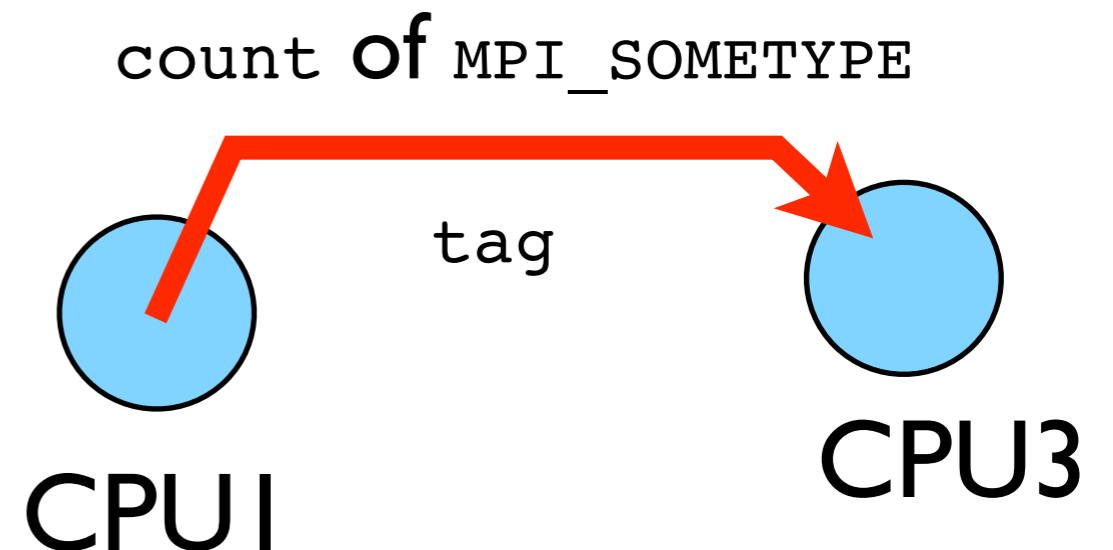
# MPI is a Library for **Message-Passing**

- Communication/coordination between tasks done by sending and receiving messages.
- Each message involves a function call from each of the programs.



# Messages

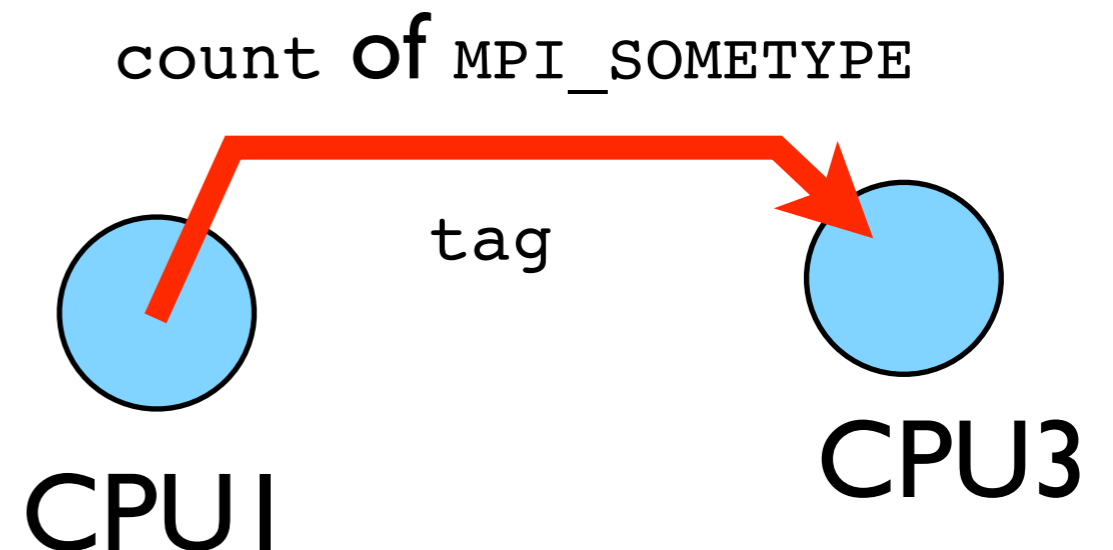
- Messages have a **sender** and a **receiver**
- When you are sending a message, don't need to specify sender (it's the current processor),
- A sent message has to be actively received by the receiving process





# Messages

- MPI messages are a string of length `count` all of some fixed MPI type
- MPI types exist for characters, integers, floating point numbers, etc.
- An arbitrary integer `tag` is also included - helps keep things straight if lots of messages are sent.



# Size of MPI Library

- Many, many functions (>200)
- Not nearly so many concepts
- We'll get started with just 6-10, use more as needed.

```
MPI_Init()  
MPI_Comm_size()  
MPI_Comm_rank()  
MPI_Send()  
MPI_Recv()  
MPI_Finalize()
```



# Hello World

- The obligatory starting point
- Let's compile and run it together

```
program hellompiworld
include "mpif.h"
```

```
integer rank, size
integer ierr
```

```
call MPI_INIT(ierr)
```

```
call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
```

```
print *, "Hello from task ", rank, " of ", size, ", world!"
```

```
call MPI_FINALIZE(ierr)
```

```
return
end
```

## Fortran

```
#include <stdio.h>
#include <mpi.h>
```

```
int main(int argc, char **argv) {
```

```
    int rank, size;
    int ierr;
```

```
    ierr = MPI_Init(&argc, &argv);
```

```
    ierr = MPI_Comm_size(MPI_COMM_WORLD, &size);
    ierr = MPI_Comm_rank(MPI_COMM_WORLD, &rank);
```

```
    printf("Hello from task %d of %d, world!\n", rank, size);
```

```
    MPI_Finalize();
```

```
    return 0;
```

## C

```
$ cd ~/pca/src/mpi-intro
```

```
edit hello-world.c or .f
```

```
$ mpif77 hello-world.f -o hello-
world
```

```
or
```

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

```
$ mpirun -np 1 hello-world
```

```
$ mpirun -np 2 hello-world
```



# What mpicc/ mpif77 do

- Just wrappers for the system C, Fortran compilers that have the various -I, -L clauses in there automatically
- --showme option shows (for this MPI version) which flags are being used

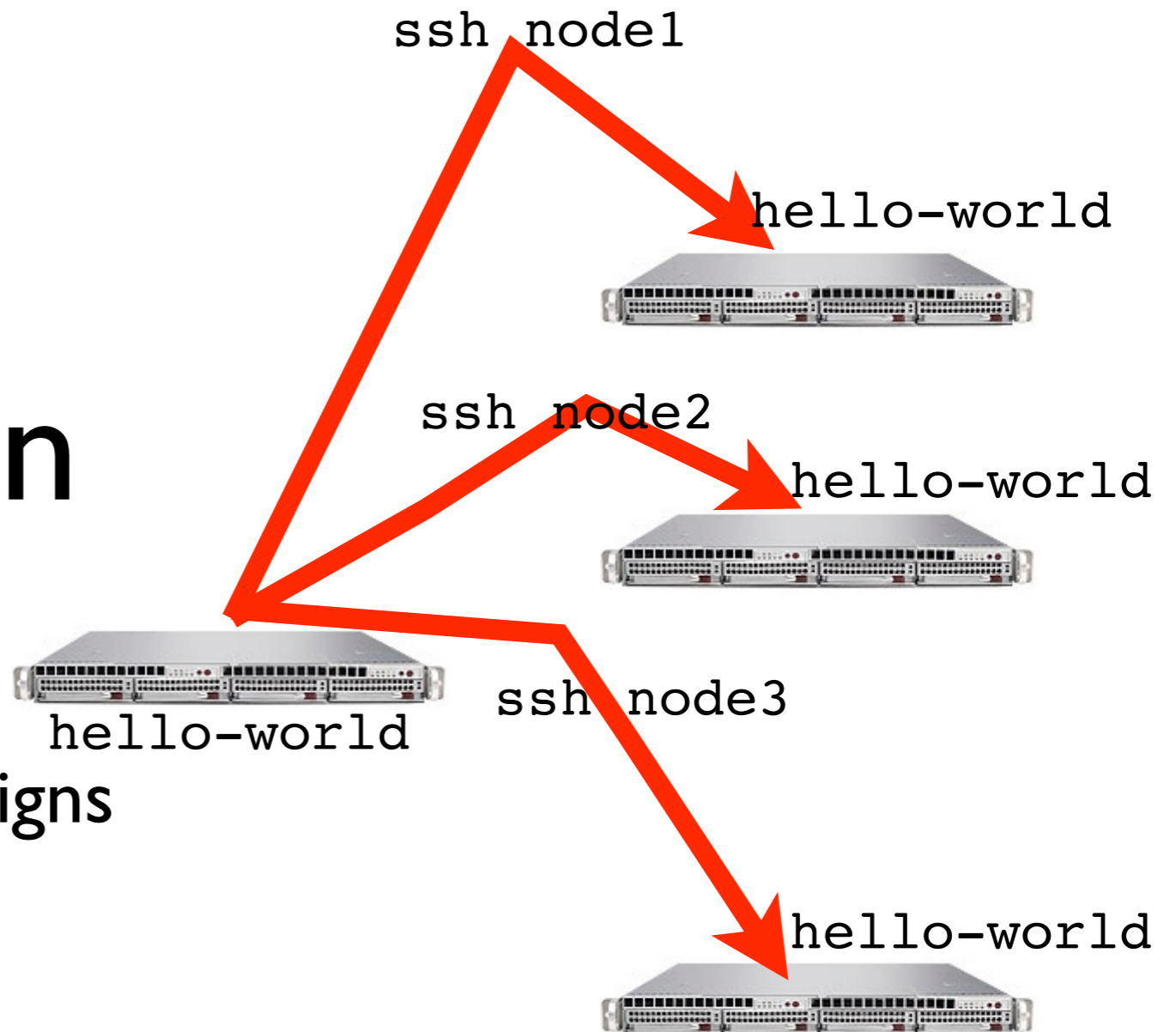
```
$ mpicc --showme hello-  
world.c -o hello-world
```

```
gcc -I/usr/local/include -  
pthread hello-world.c -o  
hello-world -L/usr/local/  
lib -lmpi -lopen-rte -  
lopen-pal -ldl -Wl,--  
export-dynamic -lnsl -lutil  
-lm -ldl
```



# What mpirun does

- Launches n processes, assigns each a PE and starts the program
- For multinode run, has a list of files, ssh's to each node and launches the program



# Number of Processes

- Number of processes to use is almost always equal to the number of processors
- But not necessarily.
- On your nodes, what happens when you run this?

```
$ mpirun -np 24 hello-world
```



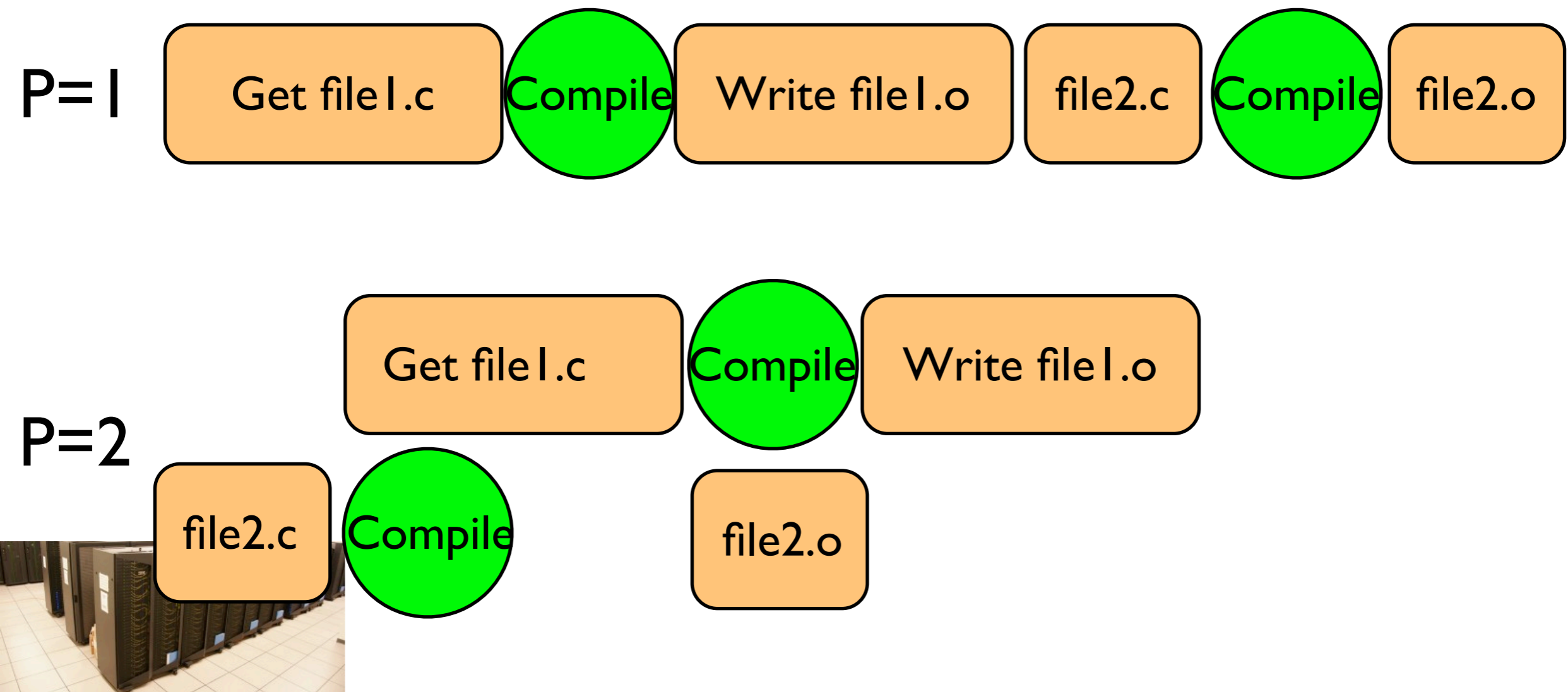
# make

- Make builds an executable from a list of source code files and rules
- Many files to do, of which order doesn't matter for most
- Parallelism!
- `make -j N` - launches N processes to do it
- `make -j 2` often shows speed increase even on single processor systems

```
$ make  
$ make -j 2  
$ make -j
```



# Overlapping Computation with I/O





# What the code does

```
program hellompiworld
include "mpif.h"

integer rank, size
integer ierr

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

print *, "Hello from task ", rank, " of ", size, ", world!"

call MPI_FINALIZE(ierr)

return
end
```

- (FORTRAN version; C is similar)



`include "mpif.h"`: imports declarations for MPI function calls

```
program hellompiworld
include "mpif.h"

integer rank, size
integer ierr

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

print *, "Hello from task ", rank, " of ", size, " MPI processes"

call MPI_FINALIZE(ierr)

return
end
```

`call MPI_INIT(ierr)`: initialization for MPI library. Must come first.

`ierr`: Returns any error code.

`call MPI_FINALIZE(ierr)`: close up MPI stuff. Must come last.  
`ierr`: Returns any error code.



```
program hellompiworld
include "mpif.h"

integer rank, size
integer ierr

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

print *, "Hello from task ", rank, " of ", size, ", world!"

call MPI_FINALIZE(ierr)

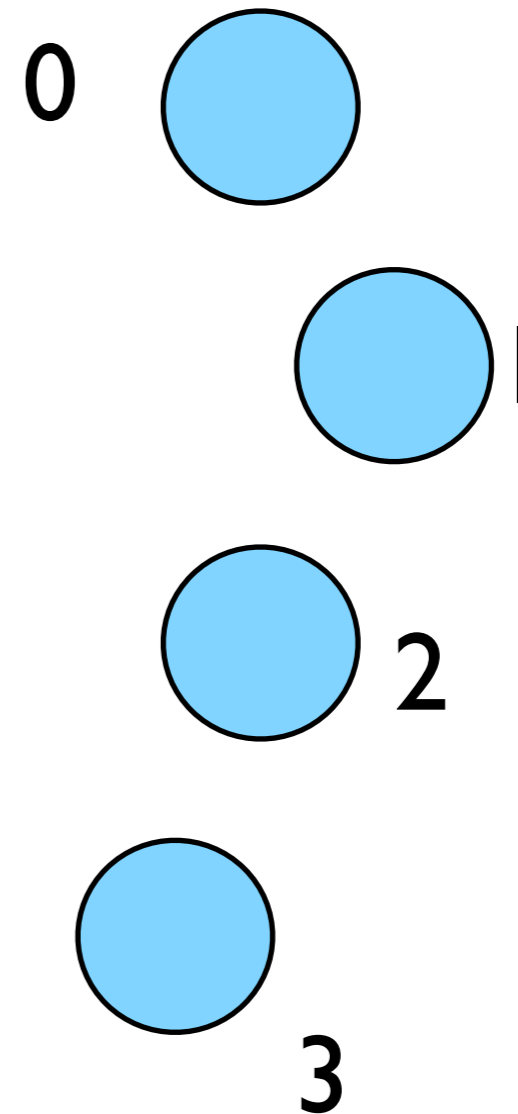
return
end
```

call MPI\_COMM\_RANK,  
call MPI\_COMM\_SIZE:  
requires a little more exposition.



# Communicators

- MPI Groups processes into communicators.
- Each communicator has some size -- number of tasks.
- Each task has a rank 0..size-1
- Every task belongs to `MPI_COMM_WORLD`

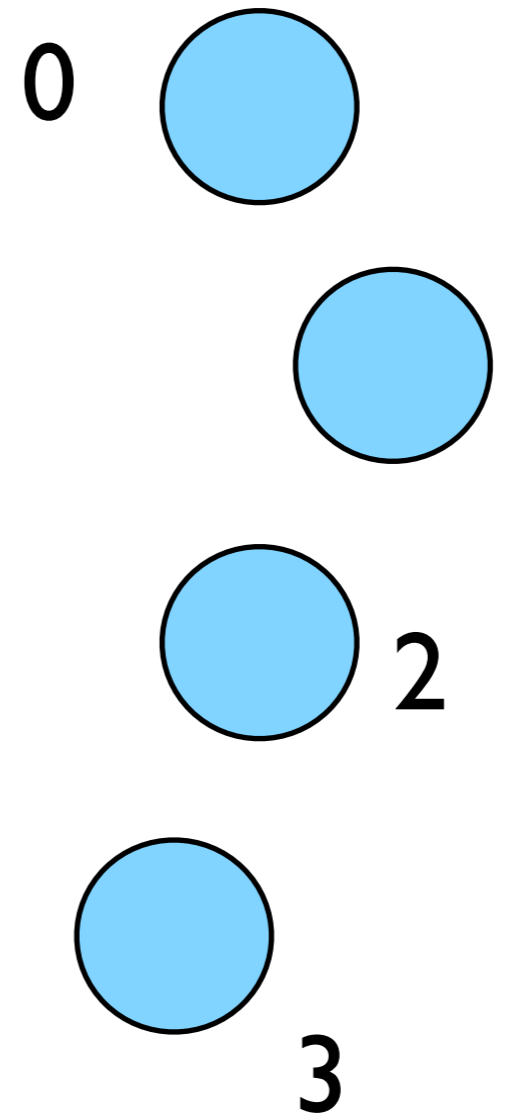


`MPI_COMM_WORLD:`  
`size=4, ranks=0..3`

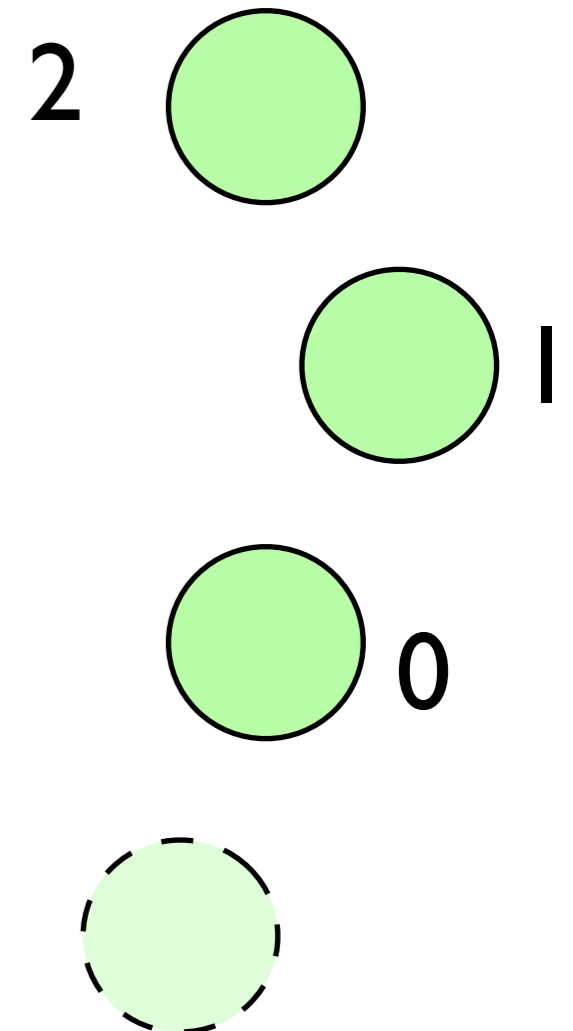
# Communicators

- Can create our own communicators over the same tasks
- May break the tasks up into subgroups
- May just re-order them for some reason

MPI\_COMM\_WORLD:  
size=4, ranks=0..3



new\_comm  
size=3, ranks=0..2





```
program hellompiworld
include "mpif.h"

integer rank, size
integer ierr

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

print *, "Hello from task ", rank, " of ", size, ", world!"

call MPI_FINALIZE(ierr)

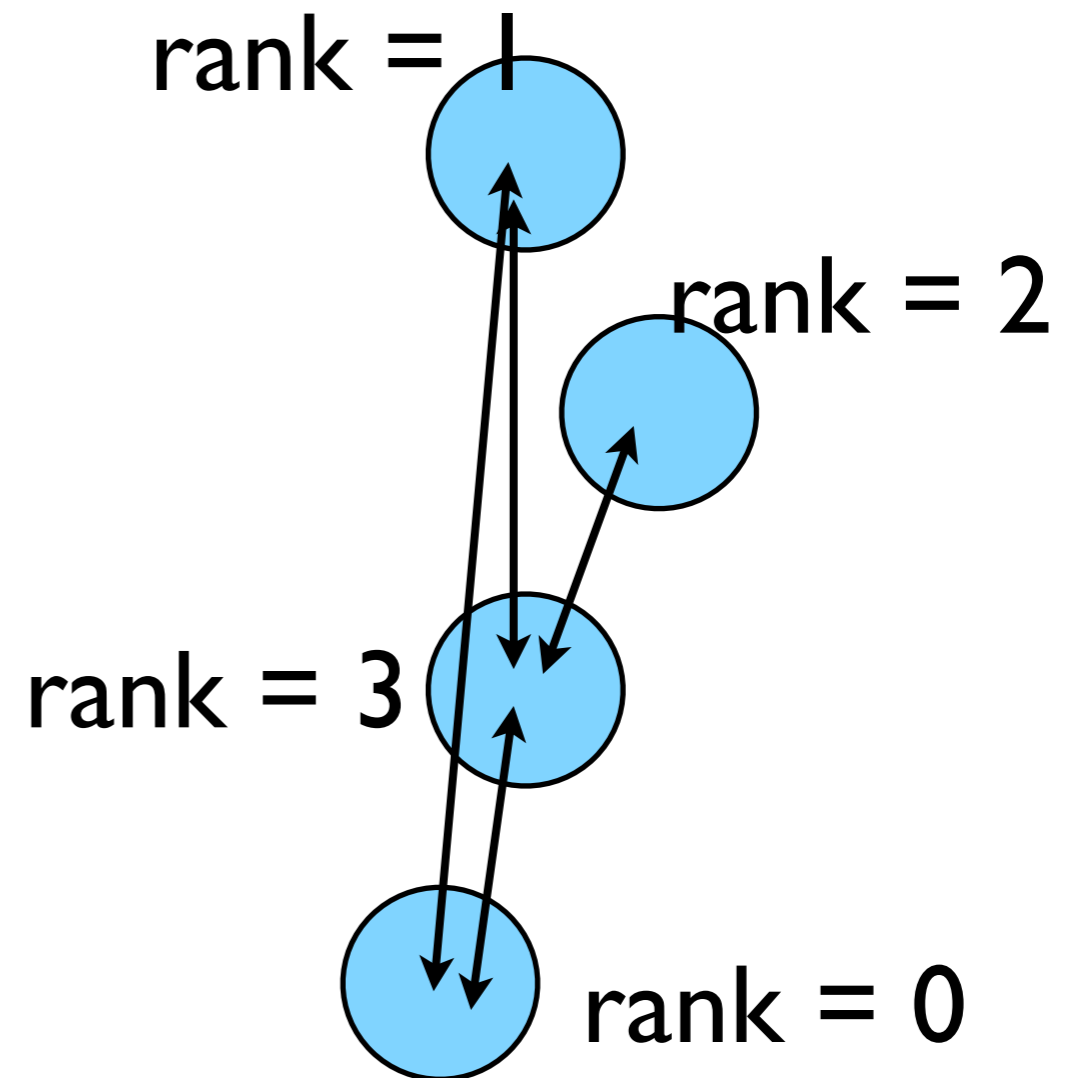
return
end
```

call MPI\_COMM\_RANK,  
call MPI\_COMM\_SIZE:  
get the size of  
MPI\_COMM\_WORLD and the  
current task's rank within  
MPI\_COMM\_WORLD  
put answers in rank and  
size



# Rank and Size much more important in MPI than OpenMP

- In OpenMP, compiler assigns jobs to each thread; don't need to know which one you are.
- MPI: processes determine amongst themselves which piece of puzzle to work on, then communicate with appropriate others.



# C

# Fortran

```
#include <stdio.h>
#include <mpi.h>

int main(int argc, char **argv) {
    int rank, size;
    int ierr;

    ierr = MPI_Init(&argc, &argv);

    ierr = MPI_Comm_size(MPI_COMM_WORLD, &size);
    ierr = MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    printf("Hello from task %d of %d, world!\n", rank, size);

    MPI_Finalize();

    return 0;
}
```

```
program hellompiworld
include "mpif.h"

integer rank, size
integer ierr

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

print *, "Hello from task ", rank, " of ", size, ", world!"

call MPI_FINALIZE(ierr)

return
end
```

- Fortran: All caps
- C - functions return ierr;
- Fortran - pass ierr
- MPI\_Init



# Our first real MPI program - but no Ms are P'ed!

- Let's fix this
- `cp hello-world.c firstmessage.c`
- `mpicc -o firstmessage firstmessage.c`
- `mpirun -np 2 ./firstmessage`

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <assert.h>
4 #include <mpi.h>
5
6 int main(int argc, char **argv) {
7
8     int rank, size;      /* the usual MPI stuff */
9     int ierr;
10    char hearmessage[6]; /* we recieve into here */
11    char sendmessage[]="Hello"; /* send from here*/
12    int sendto;          /* PE # we send to */
13    int recvfrom;        /* PE # we recv from */
14    const int OURTAG = 1; /* our shared tag */
15    MPI_Status status;   /* recieve status info */
16
17    ierr = MPI_Init(&argc, &argv);
18    ierr = MPI_Comm_size(MPI_COMM_WORLD, &size);
19    ierr = MPI_Comm_rank(MPI_COMM_WORLD, &rank);
20
21    assert(size > 1);    /* won't work otherwise */
22
23    if (rank == 0) {
24        sendto = 1;
25        ierr = MPI_Ssend(sendmessage, 6, MPI_CHAR, sendto,
26                          OURTAG, MPI_COMM_WORLD);
27        printf("%d: Sent message. <%s>\n", rank, sendmessage);
28    }
29
30    if (rank == 1) {
31        recvfrom = 0;
32        ierr = MPI_Recv(hearmessage, 6, MPI_CHAR, recvfrom,
33                          OURTAG, MPI_COMM_WORLD, &status);
34        printf("%d: Recieved message <%s>.\n", rank, hearmessage);
35    }
36
37    MPI_Finalize();
38
39    return 0;
40 }
```





# Fortran version

- Let's fix this
- cp hello-world.f firstmessage.f
- mpif77 -o firstmessage firstmessage.f
- mpirun -np 2 ./ firstmessage



```
1  program hellompiworld
2  implicit none
3  include "mpif.h"
4
5  integer rank, size
6  integer ierr
7  integer sendto
8  integer recvfrom
9  integer ourtag
10 parameter (ourtag=1)
11 character(5) hearmessage
12 character(5) sendmessage
13 integer status(MPI_STATUS_SIZE)
14
15 call MPI_INIT(ierr)
16 call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
17 call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
18
19 if (size .le. 1) then
20     print *, 'FAIL: only one task'
21 else
22     if (rank .eq. 0) then
23         sendmessage = 'Hello'
24         sendto = 1
25         call MPI_SSEND(sendmessage, 5, MPI_CHARACTER, sendto, &
26 & ourtag, MPI_COMM_WORLD, ierr)
27         print *, rank, ': sent message <',sendmessage,'>.'
28     else if (rank .eq. 1) then
29         recvfrom = 0
30         call MPI_RECV(hearmessage, 5, MPI_CHARACTER, recvfrom, &
31 & ourtag, MPI_COMM_WORLD, status, ierr)
32         print *, rank, ': got message <',hearmessage,'>.'
33     endif
34 endif
35 call MPI_FINALIZE(ierr)
36
37 return
38 end
```



# C - Send and Receive

```
MPI_Status status;
```

```
ierr = MPI_Ssend(sendptr, count, MPI_TYPE, destination,  
                tag, Communicator);
```

```
ierr = MPI_Recv(rcvptr, count, MPI_TYPE, destination, tag,  
               Communicator, status);
```



# Fortran - Send and Receive

```
integer status(MPI_STATUS_SIZE)
```

```
call MPI_SSEND(sendarr, count, MPI_TYPE, destination,  
              tag, Communicator)
```

```
call MPI_RECV(rcvarr, count, MPI_TYPE, destination, tag,  
             Communicator, status, ierr)
```



# More complicated example:

- make hello-world-send
- mpirun -np 4 hello-world-send



```
1 program helloworld
2
3 implicit none
4 include "mpif.h"
5
6 integer rank, size
7 integer ierr
8 integer leftneighbour
9 integer rightneighbour
10 integer tag
11 integer status(MPI_STATUS_SIZE)
12 character(15) greeting
13 character(15) received
14
15 call MPI_INIT(ierr)
16
17 call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
18 call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
19
20 leftneighbour = rank-1
21 rightneighbour = rank+1
22
23 if (leftneighbour.ge.0) print *,rank,': My left neighbor is ',
24 & leftneighbour
25 if (rightneighbour.lt.size)print *,rank,': My right neighbor is ',
26 & rightneighbour
27 print *, rank, ': lets say Hi to our right!'
28
29 write(greeting,'(a i4)'), 'Hello from', rank
30 tag = 1
31
32 if (rightneighbour .lt. size) then
33 call MPI_SSEND(greeting, 15, MPI_CHARACTER, rightneighbour,
34 & tag, MPI_COMM_WORLD, ierr)
35 endif
36
37 print *, rank, ": ...and see if we are greeted!"
38
39 if (leftneighbour .ge. 0) then
40 call MPI_RECV(received, 15, MPI_CHARACTER, leftneighbour,
41 & tag, MPI_COMM_WORLD, status, ierr)
42 print *, rank, ": leftneighbor says <",received,">!"
43 endif
44
45 call MPI_FINALIZE(ierr)
46
47 return
48 end
```

```

integer leftneighbour
integer rightneighbour
integer tag
integer status(MPI_STATUS_SIZE)
character(15) greeting
character(15) received

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

leftneighbour = rank-1
rightneighbour = rank+1

if (leftneighbour.ge.0) print *,rank,': My left neighbor is ',
& leftneighbour
if (rightneighbour.lt.size)print *,rank,': My right neighbor is ',
& rightneighbour
print *, rank, ': lets say Hi to our right!'

write(greeting,'(a i4)'), 'Hello from', rank
tag = 1

if (rightneighbour .lt. size) then
  call MPI_SSEND(greeting, 15, MPI_CHARACTER, rightneighbour,
& tag, MPI_COMM_WORLD, ierr)
endif

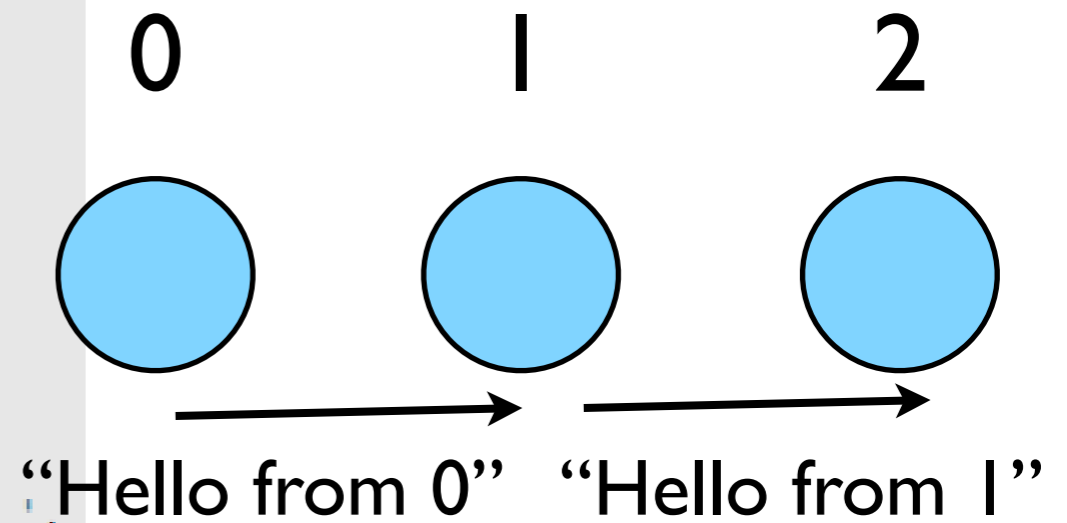
print *, rank, ": ...and see if we are greeted!"

if (leftneighbour .ge. 0) then
  call MPI_RECV(received, 15, MPI_CHARACTER, leftneighbour,
& tag, MPI_COMM_WORLD, status, ierr)
  print *, rank, ": leftneighbor says <",received,">!"
endif

call MPI_FINALIZE(ierr)

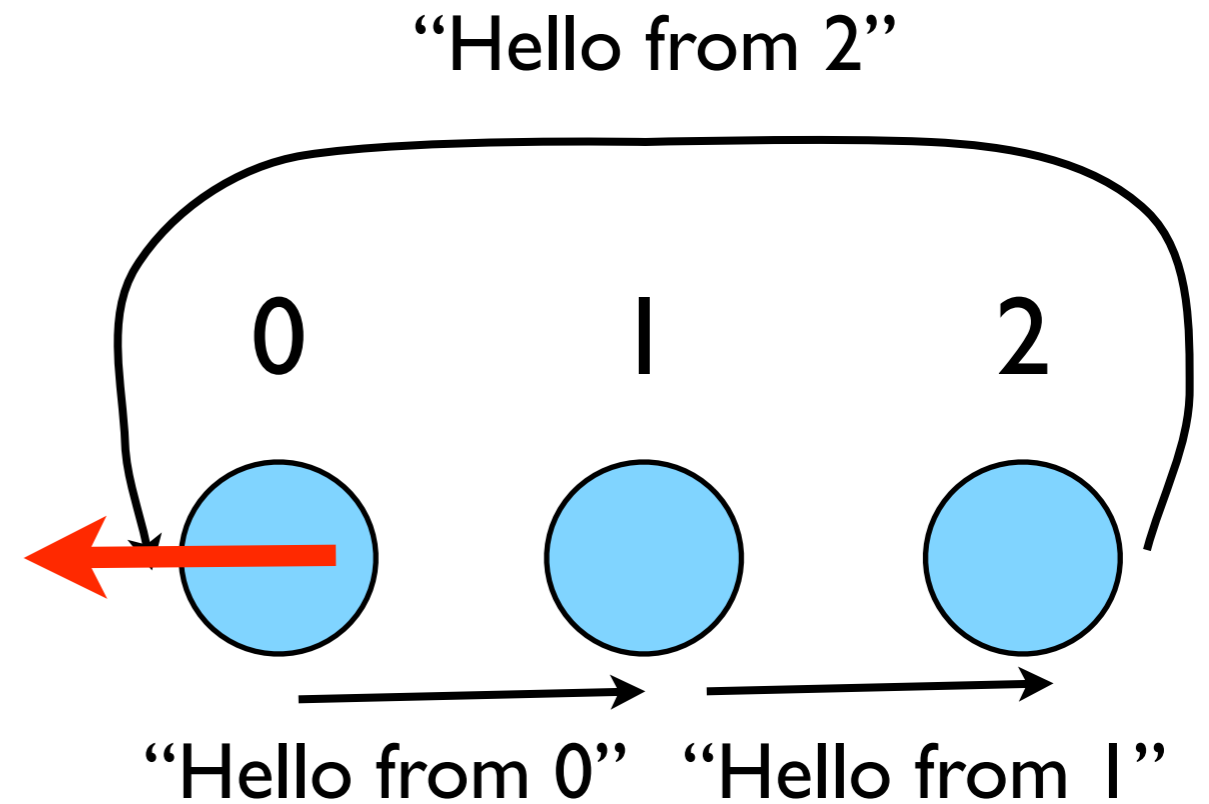
return
end

```



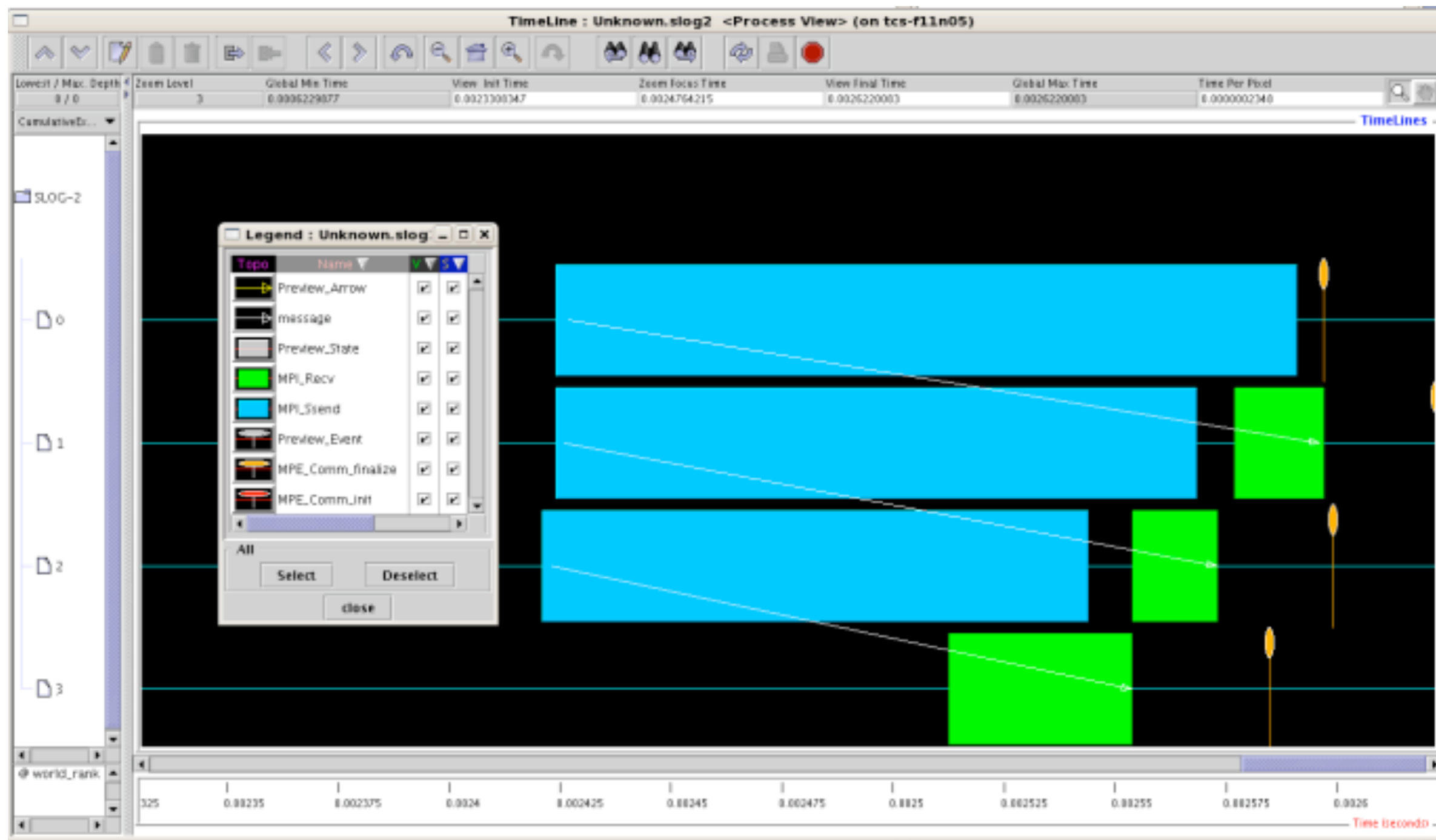
# Implement periodic boundary conditions

- `cp hello-world-send.c hello-world-send-around.c`
- edit so it `wraps around`
- make `hello-world-send-around`
- `mpirun -np 3 hello-world-send-around`

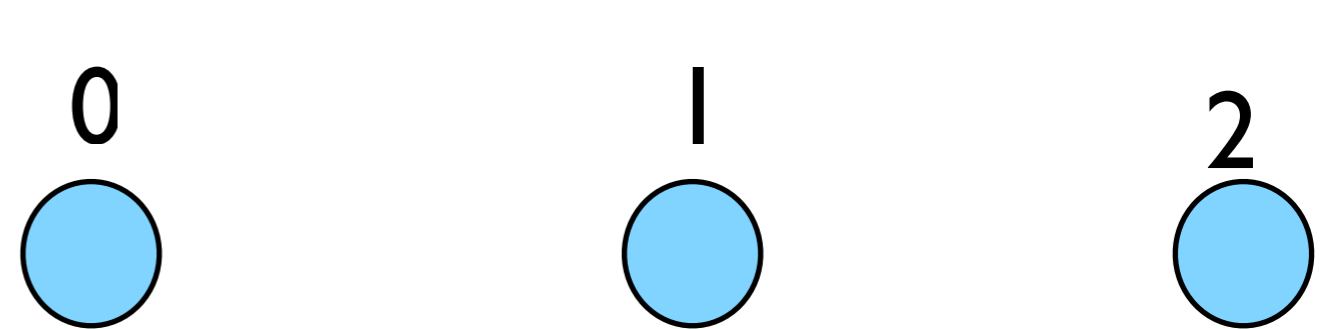
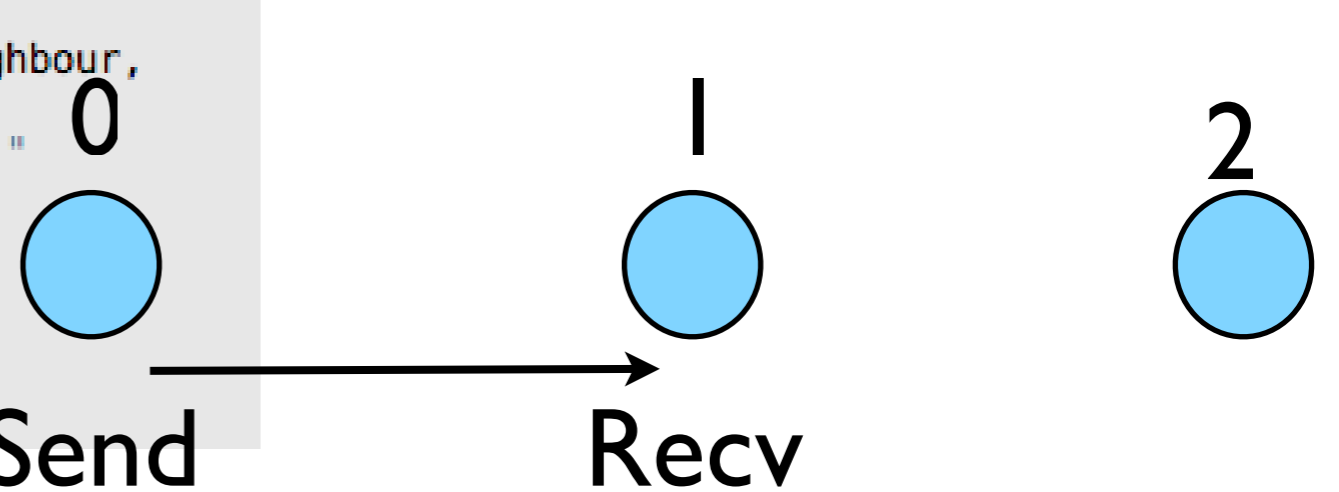
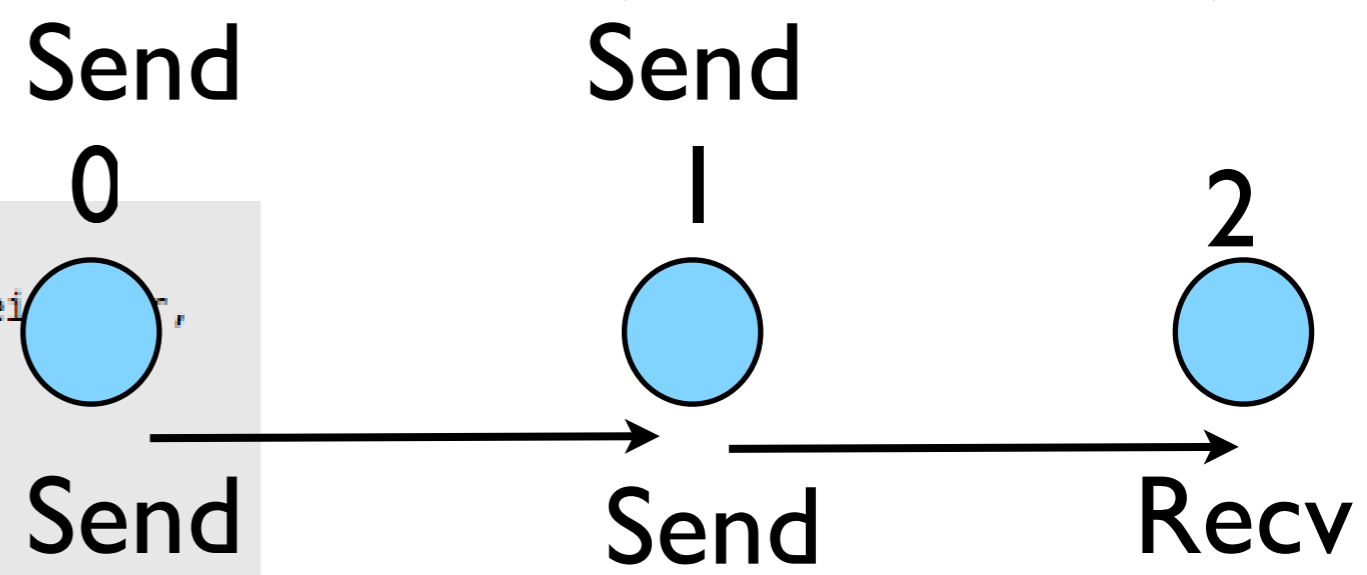
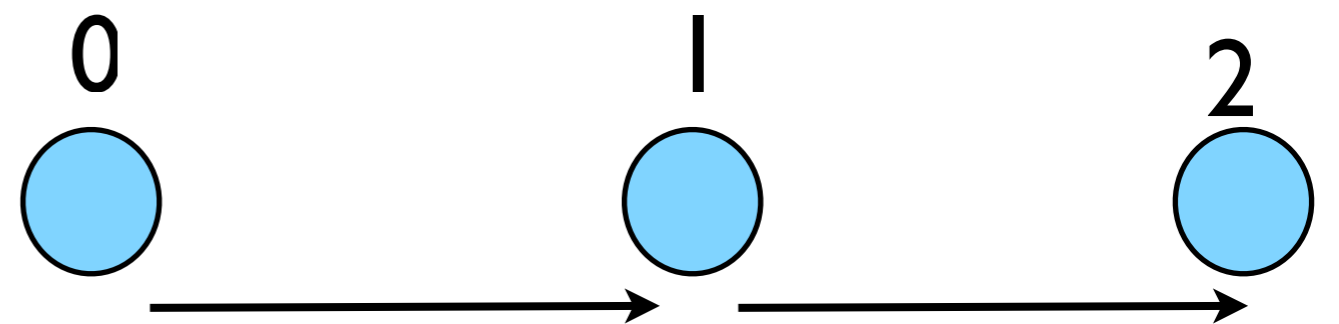


- make `hello-world-send-around`
- `mpirun -np 3 hello-world-send-around`









```

if (rightneighbour .lt. size) then
  call MPI_SSEND(greeting, 15, MPI_CHARACTER, rightneighbour,
    tag, MPI_COMM_WORLD, ierr)
endif
print *, rank, ": ...and see if we are greeted!"

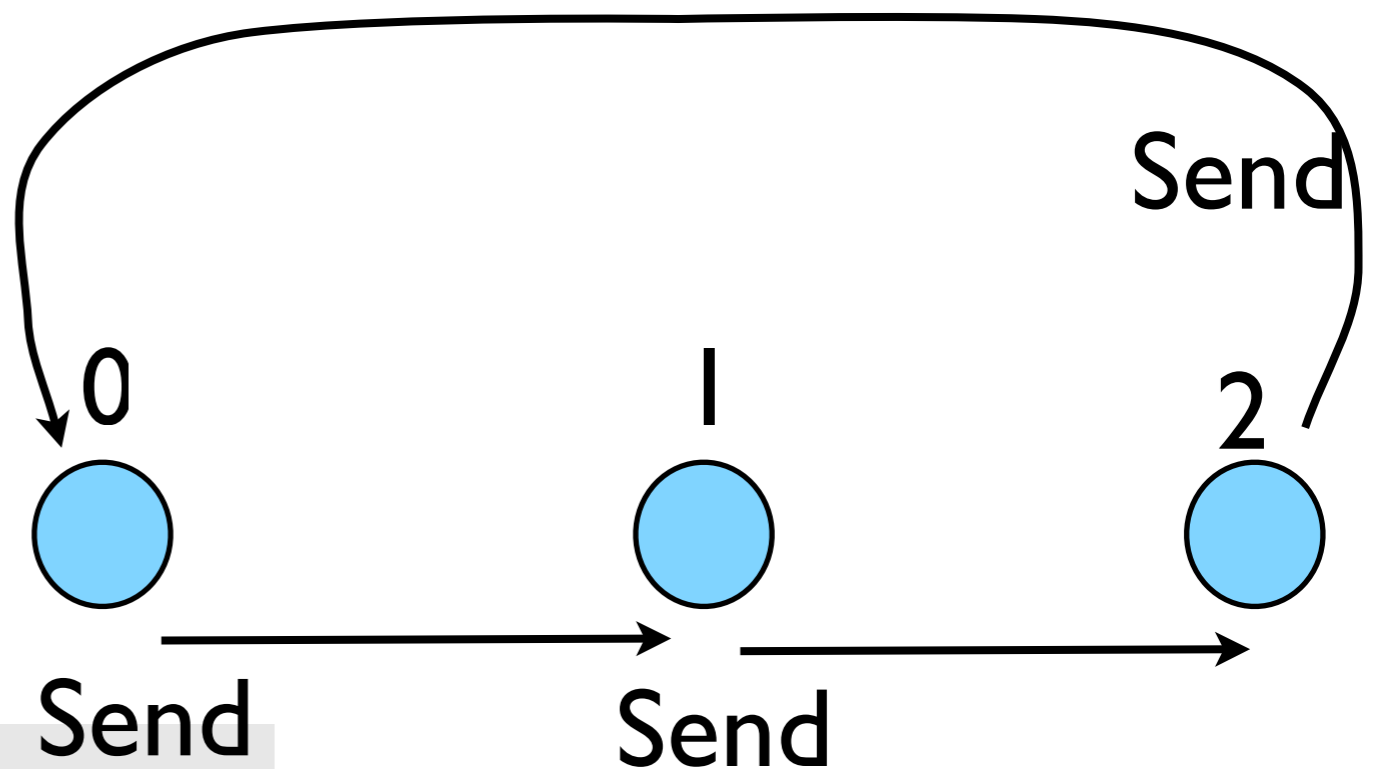
if (leftneighbour .ge. 0) then
  call MPI_RECV(received, 15, MPI_CHARACTER, leftneighbour,
    tag, MPI_COMM_WORLD, status, ierr)
  print *, rank, ": leftneighbor says <", received, ">!"
endif

call MPI_FINALIZE(ierr)

return
end

```





←  
0,1,2

```

if (rightneighbour .lt. size) then
  call MPI_SSEND(greeting, 15, MPI_CHARACTER, rightneighbour,
&
    tag, MPI_COMM_WORLD, ierr)
endif

print *, rank, ": ...and see if we are greeted!"

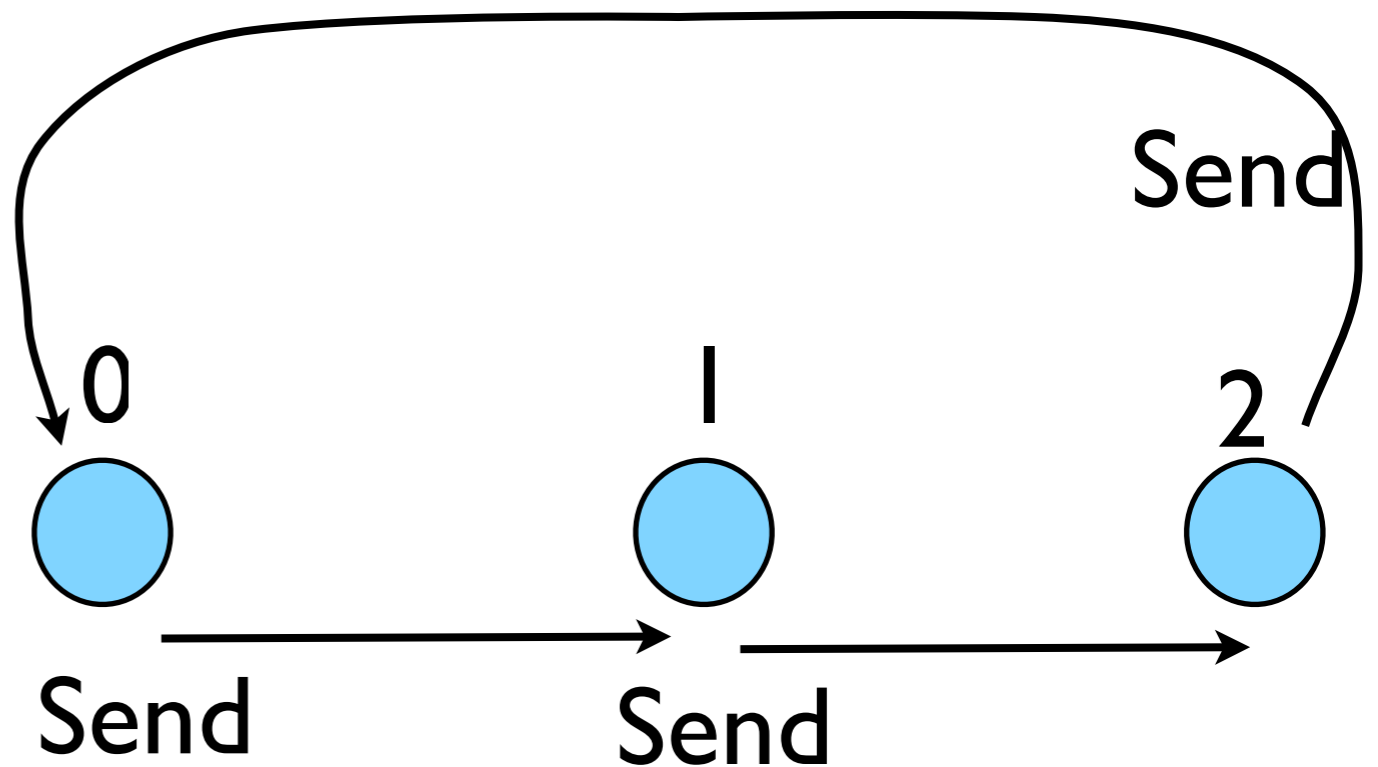
if (leftneighbour .ge. 0) then
  call MPI_RECV(received, 15, MPI_CHARACTER, leftneighbour,
&
    tag, MPI_COMM_WORLD, status, ierr)
  print *, rank, ": leftneighbor says <", received, ">!"
endif

call MPI_FINALIZE(ierr)

return
end
  
```

# Deadlock

- The other classic parallel bug
- Occurs when a cycle of tasks are for the others to finish.
- Whenever you see a closed cycle, you likely have (or risk) deadlock.



# Big MPI

## Lesson #1

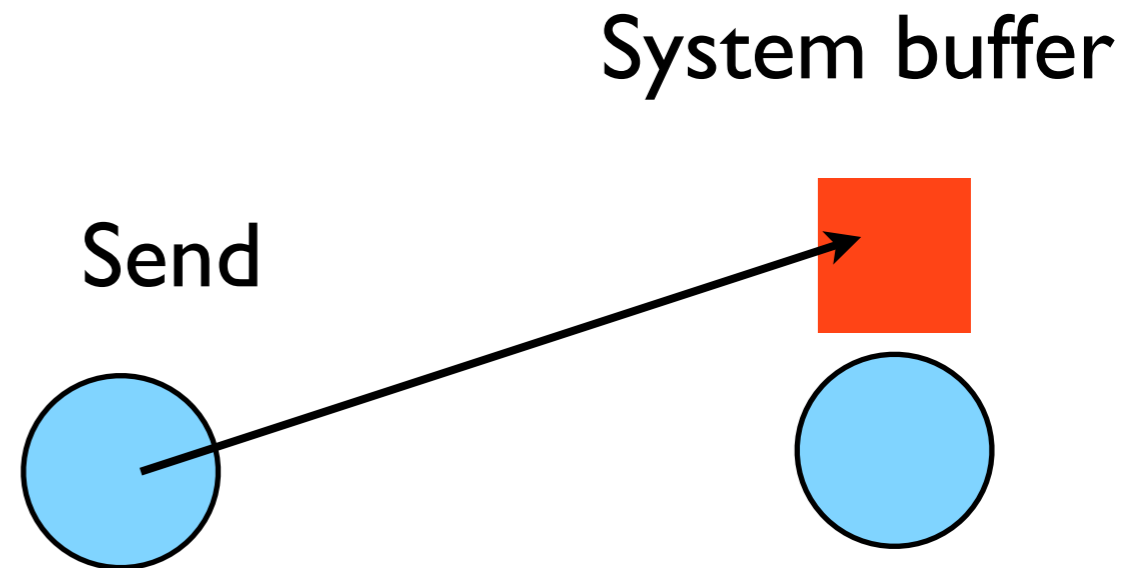
All sends and receives must be paired, **at time of sending**



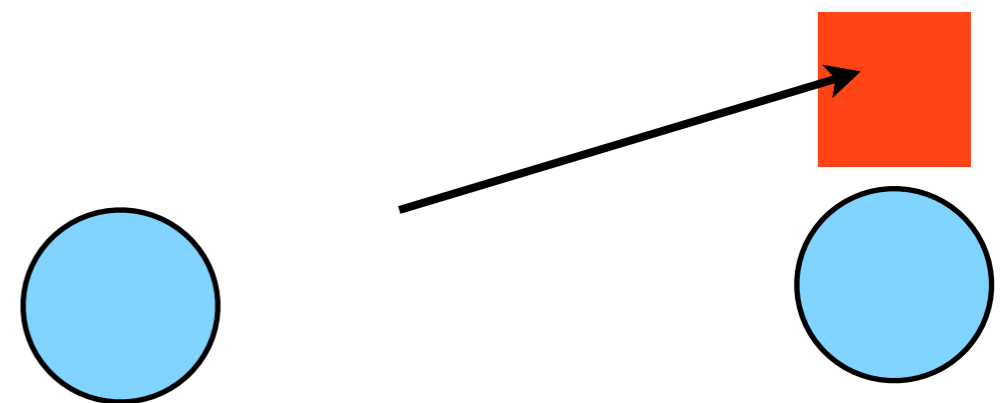
# Different versions of SEND

- SSEND: safe send; doesn't return until receive has started. Blocking, no buffering.
- SEND: Undefined. Blocking, probably buffering
- ISEND : Unblocking, no buffering
- IBSEND: Unblocking, buffering

## Buffering



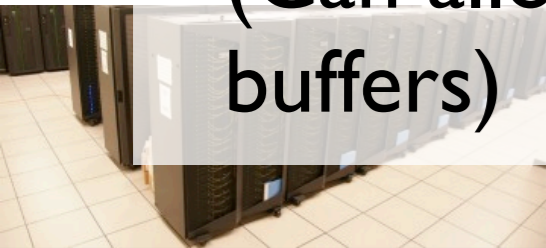
## (Non) Blocking



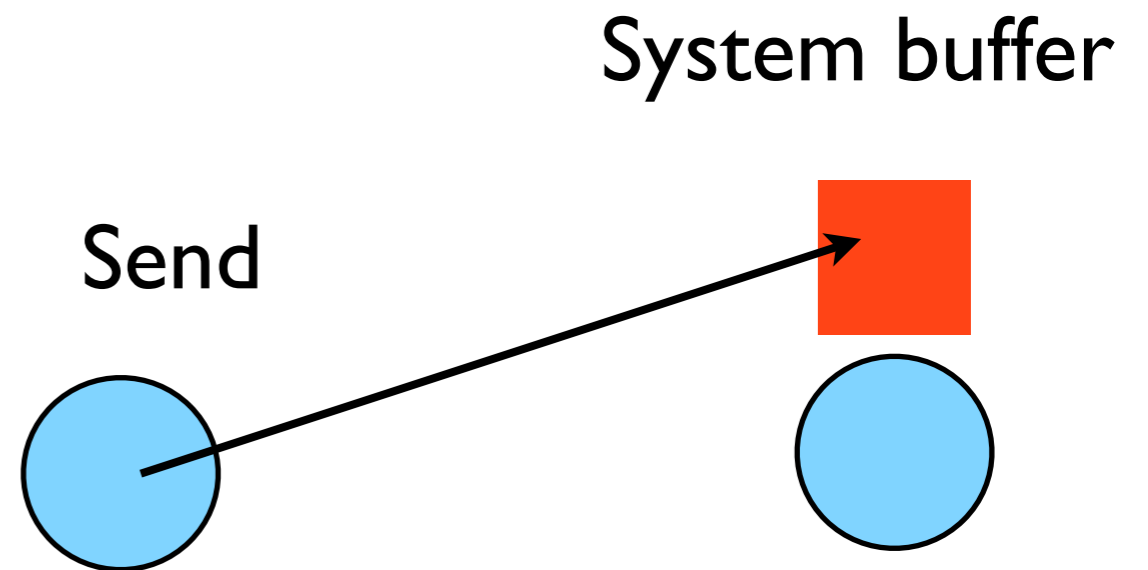


# Buffering is dangerous!

- Worst kind of danger: will usually work.
- Think voice mail; message sent, reader reads when ready
- But voice mail boxes do fill
- Message fails.
- Program fails/hangs mysteriously.
- (Can allocate your own buffers)

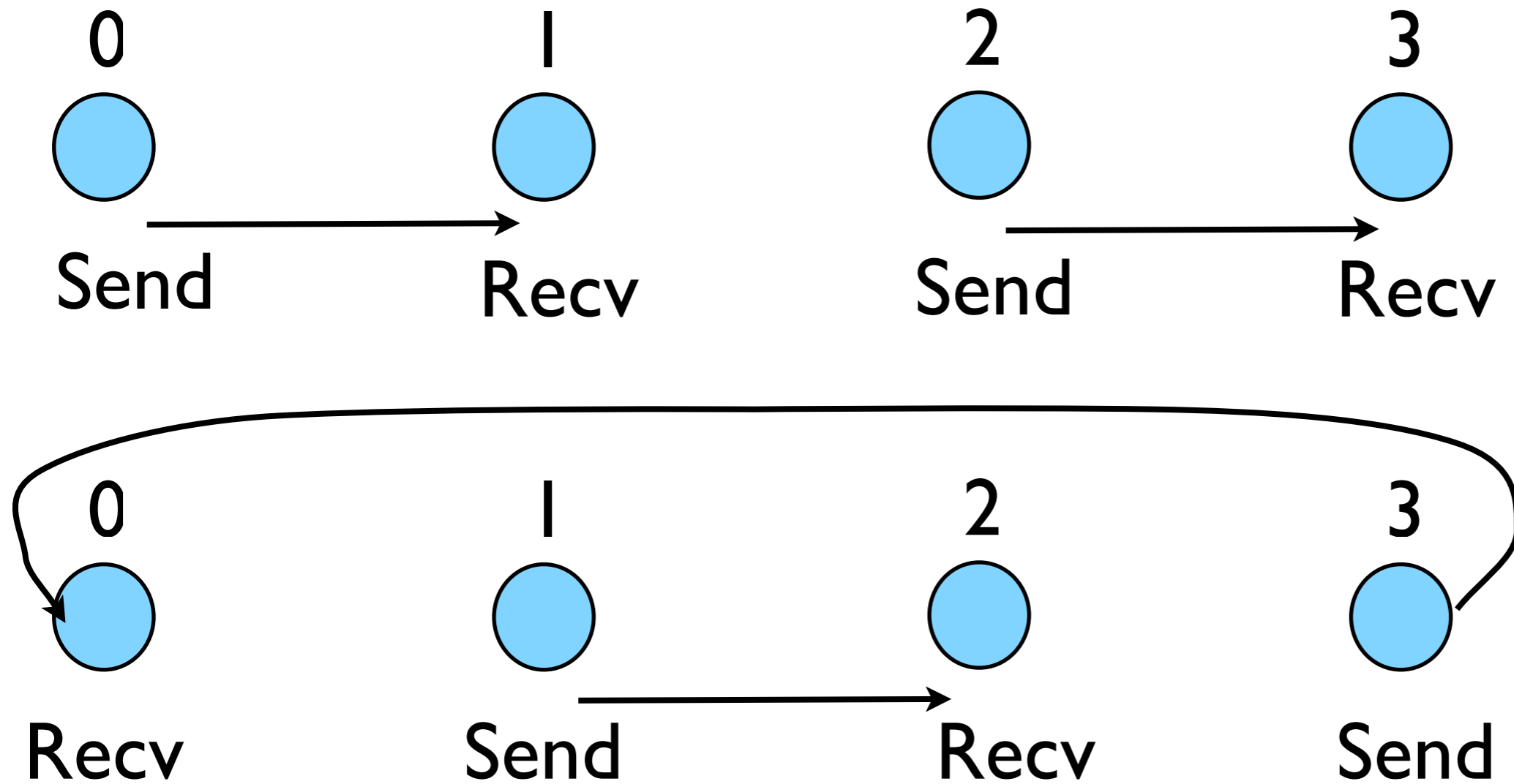


## Buffering



Without using new MPI  
routines, how can we fix  
this?





- First: evens send, odds receive
- Then: odds send, evens receive
- Will this work with an odd # of processes?
- How about 2? 1?



```

program hellompiworld

implicit none
include "mpif.h"

integer rank, size
integer ierr
integer leftneighbour
integer rightneighbour
integer tag
integer status(MPI_STATUS_SIZE)
character(15) greeting
character(15) received

call MPI_INIT(ierr)

call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)

leftneighbour = mod((rank-1)+size,size)
rightneighbour = mod((rank+1),size)

if (leftneighbour.ge.0)
&     print *,rank,': My left neighbour is ', leftneighbour
if (rightneighbour.lt.size)
&     print *,rank,': My right neighbour is ',rightneighbour

write(greeting,'(a i4)'), 'Hello from', rank
tag = 1

if (mod(rank,2) .eq. 0) then
    call MPI_SSEND(greeting, 15, MPI_CHARACTER,
&                 rightneighbour, tag, MPI_COMM_WORLD, ierr)
    call MPI_RECV(received, 15, MPI_CHARACTER,
&                 leftneighbour, tag, MPI_COMM_WORLD, status,
&                 ierr)
    print *, rank, ": leftneighbour says <",received,>!"
else
    call MPI_RECV(received, 15, MPI_CHARACTER,
&                 leftneighbour, tag, MPI_COMM_WORLD, status,
&                 ierr)
    print *, rank, ": leftneighbour says <",received,>!"
    call MPI_SSEND(greeting, 15, MPI_CHARACTER,
&                 rightneighbour, tag, MPI_COMM_WORLD, ierr)
endif

call MPI_FINALIZE(ierr)

return
end

```

← Evens send first

← Then odds

hello-world-send-fixed.f

# Something new: Sendrecv

- A blocking send and receive built in together
- Lets them happen simultaneously
- Can automatically pair the sends/recvs!
- dest, source does not have to be same; nor do types or size.



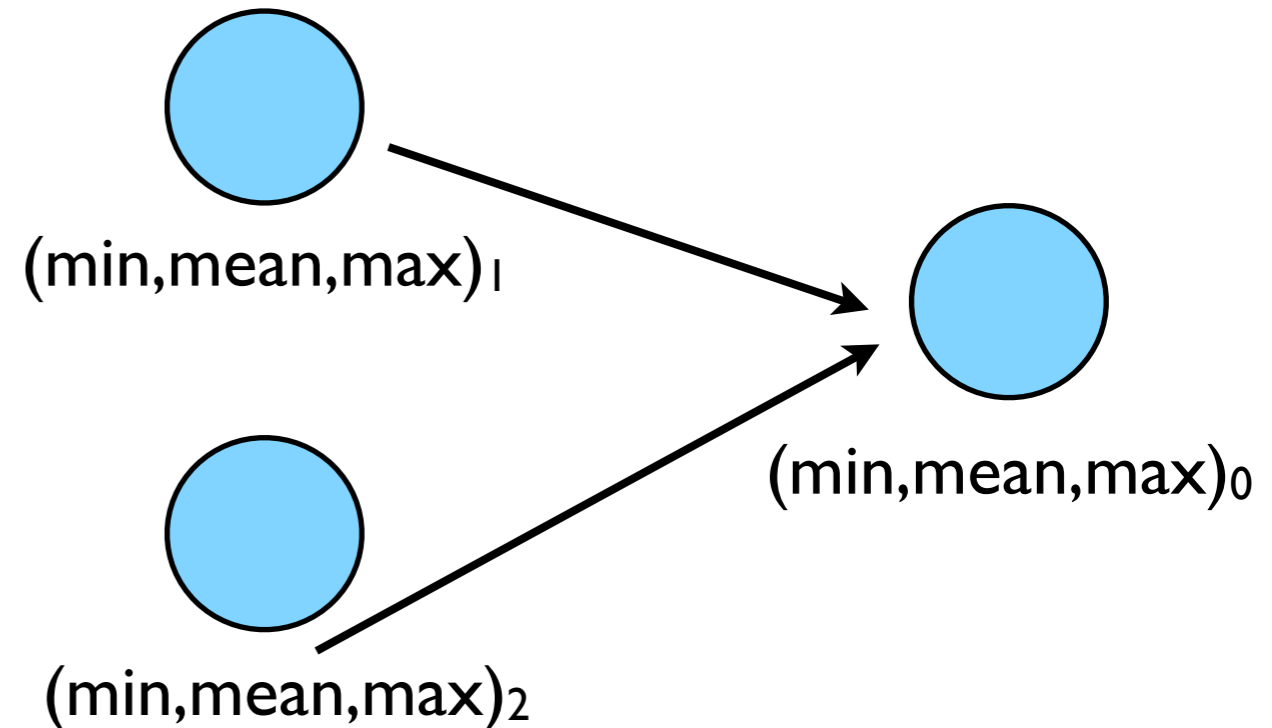
```
1  program hellompiworld
2
3  implicit none
4  include "mpif.h"
5
6  integer rank, size
7  integer ierr
8  integer leftneighbour
9  integer rightneighbour
10 integer tag
11 integer status(MPI_STATUS_SIZE)
12 character(15) greeting
13 character(15) received
14
15 call MPI_INIT(ierr)
16
17 call MPI_COMM_RANK(MPI_COMM_WORLD, rank, ierr)
18 call MPI_COMM_SIZE(MPI_COMM_WORLD, size, ierr)
19
20 leftneighbour = mod((rank-1)+size,size)
21 rightneighbour = mod((rank+1),size)
22
23 if (leftneighbour.ge.0)
24 &     print *,rank,': My left neighbour is ', leftneighbour
25 if (rightneighbour.lt.size)
26 &     print *,rank,': My right neighbour is ',rightneighbour
27
28 write(greeting,'(a i4)'), 'Hello from', rank
29 tag = 1
30
31 call MPI_SENDRECV(greeting, 15, MPI_CHARACTER, rightneighbour,tag,
32 &     received, 15, MPI_CHARACTER,leftneighbour, tag,
33 &     MPI_COMM_WORLD, status, ierr)
34
35 print *, rank, ": leftneighbour says <",received,">!"
36
37 call MPI_FINALIZE(ierr)
38
39 return
40 end
```

hello-world-sendrecv.f



# Min, Mean, Max of numbers

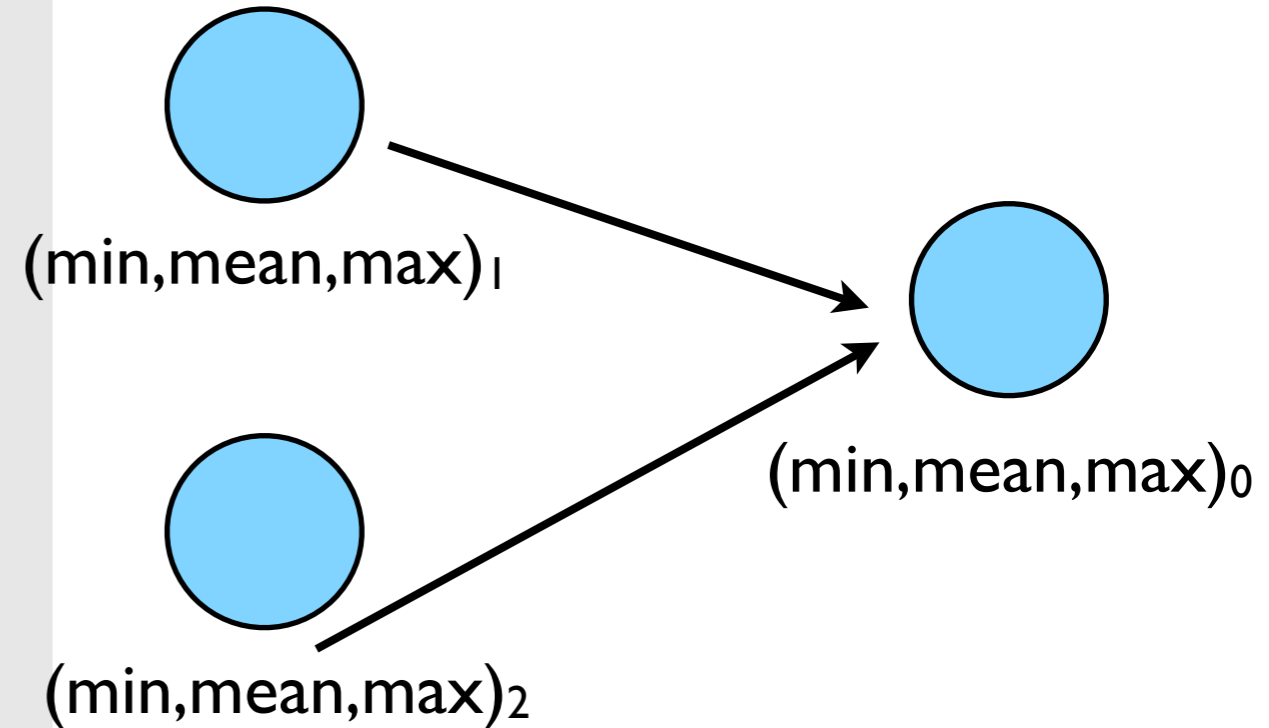
- Lets try some code that calculates the min/mean/max of a bunch of random numbers  $-1..1$ . Should go to  $-1, 0, +1$  for large  $N$ .
- Each gets their partial results and sends it to some node, say node 0 (why node 0?)



```

33 c
34 c find min/mean/max
35 c
36     datamin = 1e+19
37     datamax = -1e+19
38     datamean = 0
39
40     do i=1,nx
41         do j=1,ny
42             if (dat(i,j) .lt. datamin) datamin = dat(i,j)
43             if (dat(i,j) .gt. datamax) datamax = dat(i,j)
44             datamean = datamean + dat(i,j)
45         enddo
46     enddo
47     datamean = datamean/(1.*nx*ny)
48
49     print *,myid,': min/mean/max = ', datamin, datamean, datamax
50 c
51 c combine data
52 c
53     if (myid .ne. 0) then
54         datapack(1) = datamin
55         datapack(2) = datamean
56         datapack(3) = datamax
57         call MPI_SSEND(datapack,3,MPI_REAL,0,1,MPI_COMM_WORLD,ierr)
58     else
59         globmin = datamin
60         globmax = datamax
61         globmean = datamean
62         do proc=1,nprocs-1
63             call MPI_RECV(datapack, 3, MPI_REAL, MPI_ANY_SOURCE, 1,
64                 MPI_COMM_WORLD, status, ierr)
65             if (datapack(1) .lt. globmin) globmin=datapack(1)
66             globmean = globmean + datapack(2)
67             if (datapack(3) .gt. globmax) globmax=datapack(3)
68         enddo
69         globmean = globmean/nprocs
70         print *,'Global min/mean/max=',globmin,globmean,globmax
71     endif
72
73     call MPI_FINALIZE(ierr)
74     return
75     end
76
77
78

```



Q: are these sends/recvd adequately paired?

minmeanmax-sendrecv.f

```

C
C find min/mean/max
C
datamin = 1e+19
datamax = -1e+19
datamean = 0

do i=1,nx
  do j=1,ny
    if (dat(i,j) .lt. datamin) datamin = dat(i,j)
    if (dat(i,j) .gt. datamax) datamax = dat(i,j)
    datamean = datamean + dat(i,j)
  enddo
enddo
datamean = datamean/(1.*nx*ny)

print *,myid,': min/mean/max = ', datamin, datamean, datamax

C
C combine data
C
call MPI_ALLREDUCE(datamin, globmin, 1, MPI_REAL, MPI_MIN,
& MPI_COMM_WORLD, ierr)

C to just send to task 0:
call MPI_REDUCE(datamin, globmin, 1, MPI_REAL, MPI_MIN,
& 0, MPI_COMM_WORLD, ierr)
etc.

call MPI_ALLREDUCE(datamax, globmax, 1, MPI_REAL, MPI_MAX,
& MPI_COMM_WORLD, ierr)
call MPI_ALLREDUCE(datamean, globmean, 1, MPI_REAL, MPI_SUM,
& MPI_COMM_WORLD, ierr)
globmean = globmean/nprocs
print *, myid,': Global min/mean/max=',globmin,globmean,globmax

call MPI_FINALIZE(ierr)
return
end

```

## MPI\_Reduce and MPI\_Allreduce

Performs a reduction  
and sends it to one node or  
everywhere

minmeanmax-allreduce.f

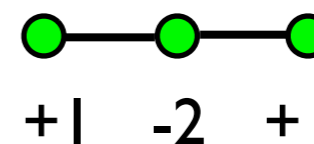
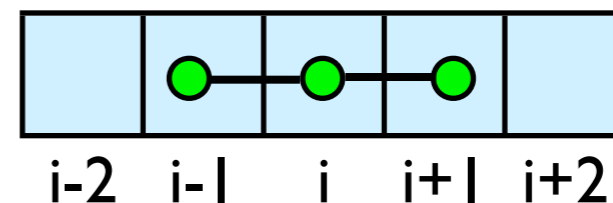


# Discretizing Derivatives

$$\left. \frac{dQ}{dx} \right|_i \approx \frac{Q_{i+1} - 2Q_i + Q_{i-1}}{\Delta x}$$

- Done by finite differencing the discretized values
- Implicitly or explicitly involves interpolating data and taking derivative of the interpolant
- More accuracy - larger

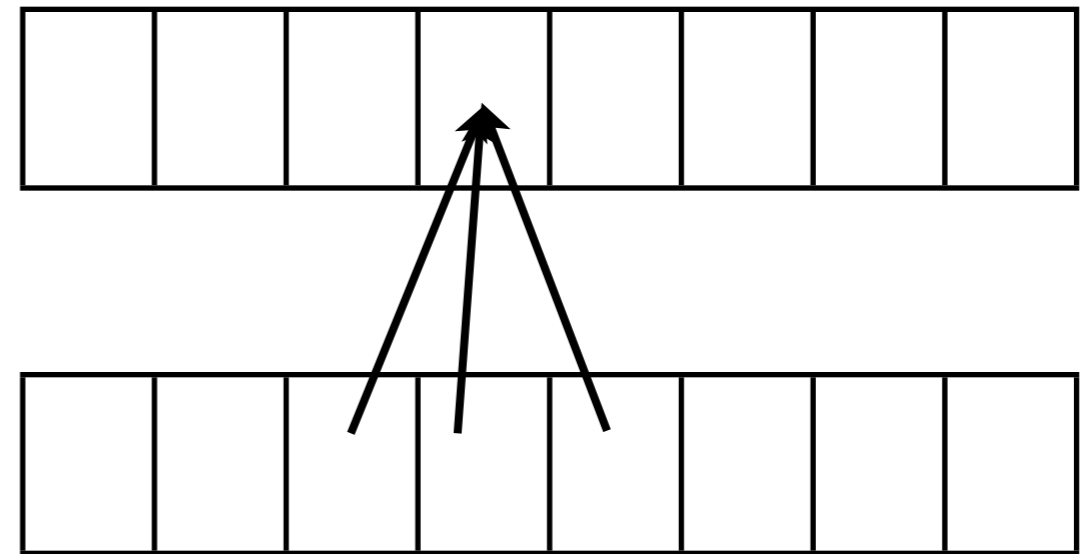
‘stencils’



# Diffusion Equation

$$\begin{aligned}\frac{\partial T}{\partial t} &= -D \frac{\partial^2 T}{\partial x^2} \\ \frac{\partial T_i^{(n)}}{\partial t} &\approx \frac{T_i^{(n)} - T_i^{(n-1)}}{\Delta t} \\ \frac{\partial T_i^{(n)}}{\partial x} &\approx \frac{T_{i+1}^{(n)} - 2T_i^{(n)} + T_{i-1}^{(n)}}{\Delta x^2} \\ T_i^{(n+1)} &\approx T_i^{(n)} - \frac{D\Delta t}{\Delta x^2} (T_{i+1}^{(n)} - 2T_i^{(n)} + T_{i-1}^{(n)})\end{aligned}$$

- Simple 1d ODE
- Each timestep, new data for  $T[i]$  requires old data for  $T[i+1], T[i], T[i-1]$



```
cd ~/pca/mipi-intro
make diffusion
./diffusion
```

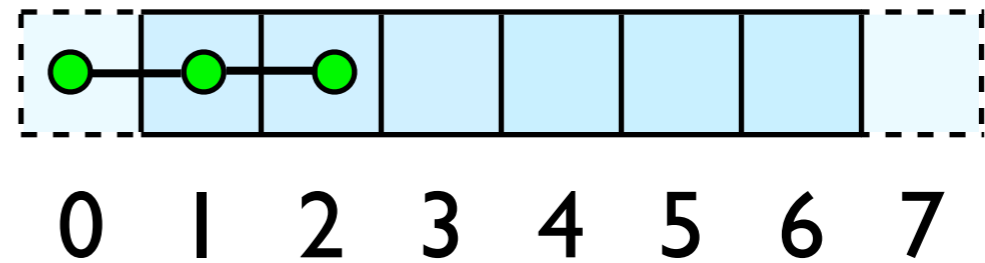




# Guardcells

- How to deal with boundaries?
- Because stencil juts out, need information on cells beyond those you are updating
- Pad domain with 'guard cells' so that stencil works even for the first point in domain
- Fill guard cells with values such that the required boundary conditions are met

## Global Domain



$$ng = 1$$

loop from  $ng$ ,  $N - 2 \cdot ng$



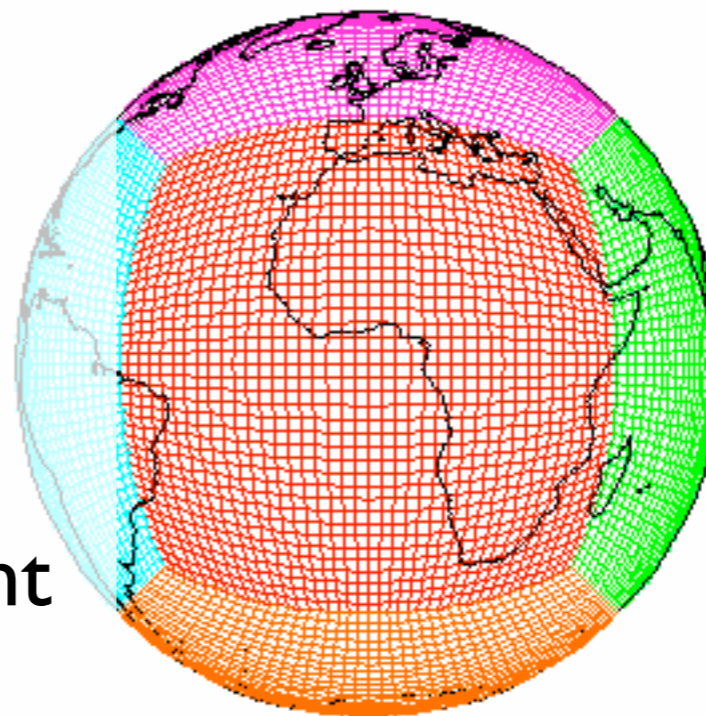
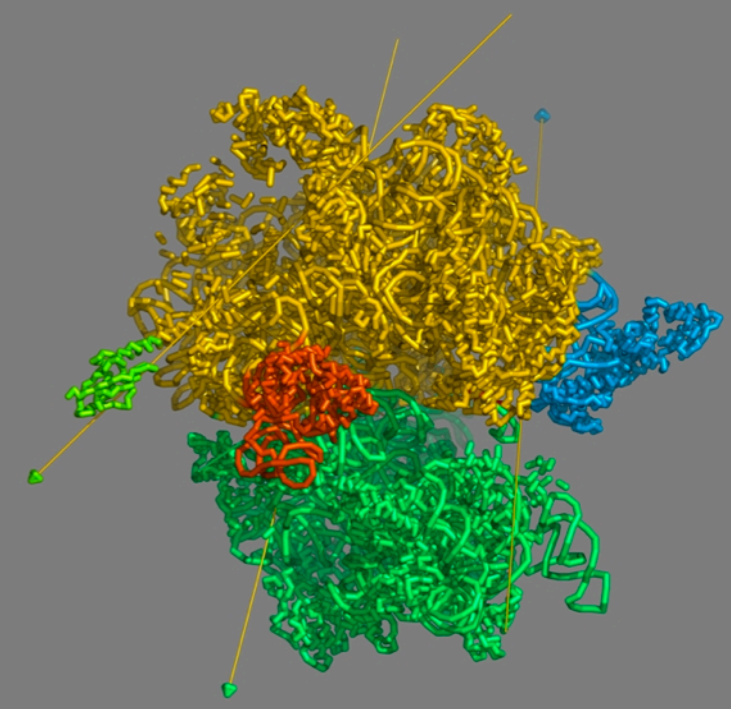
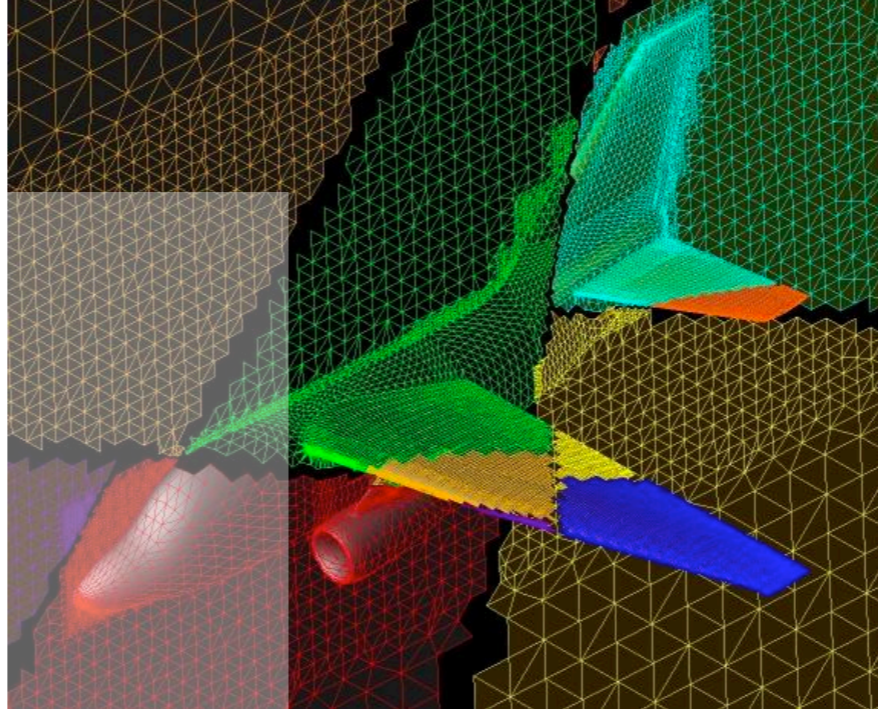
# Domain

# Decomposition

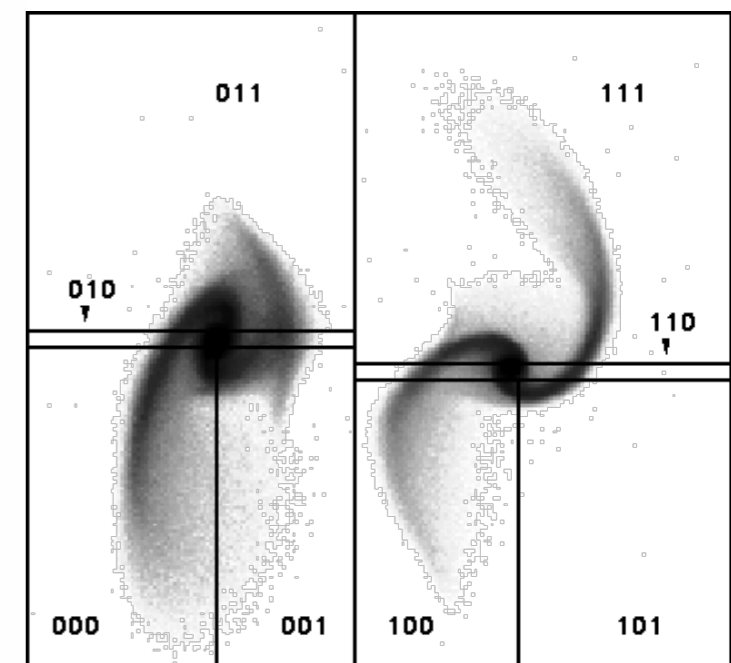
<http://adg.stanford.edu/aa241/design/compaero.html>

<http://www.uea.ac.uk/cmp/research/cmpbio/Protein+Dynamics,+Structure+and+Function>

- A very common approach to parallelizing on distributed memory computers
- Maintain Locality; need local data mostly, this means only surface data needs to be sent between processes.



[http://sivo.gsfc.nasa.gov/cubedsphere\\_comp.html](http://sivo.gsfc.nasa.gov/cubedsphere_comp.html)



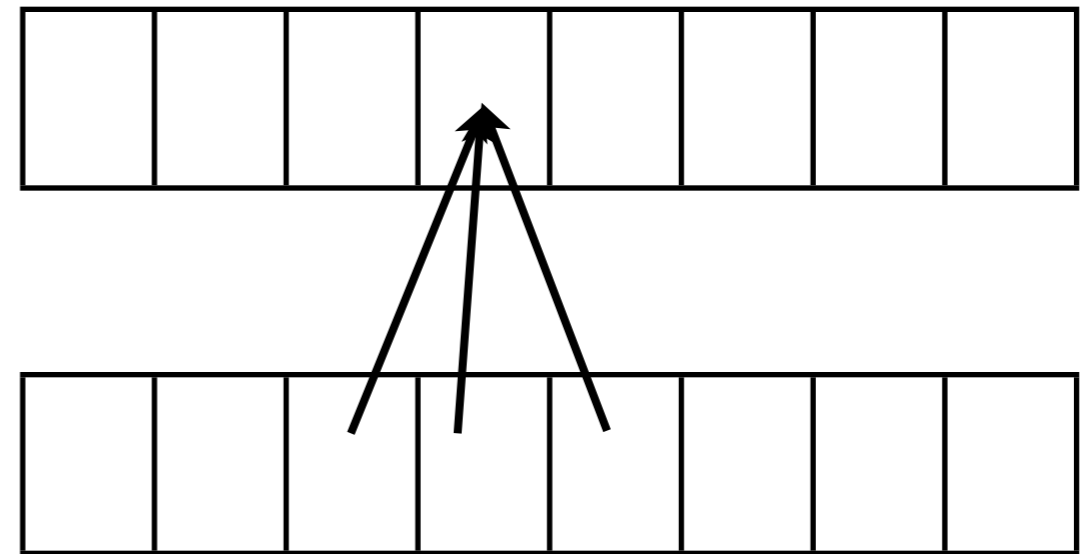
<http://www.cita.utoronto.ca/~dubinski/treecode/node8.html>



# Implement a diffusion equation in MPI

- Need one neighboring number per neighbor per timestep

$$\frac{dT}{dt} = D \frac{d^2T}{dx^2}$$
$$T_i^{n+1} = T_i^n + \frac{D\Delta t}{\Delta x^2} (T_{i+1}^n - 2T_i^n + T_{i-1}^n)$$



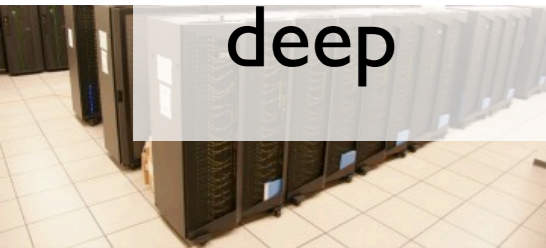
```
cd ~/pca/mpl-intro
make diffusion
./diffusion
```



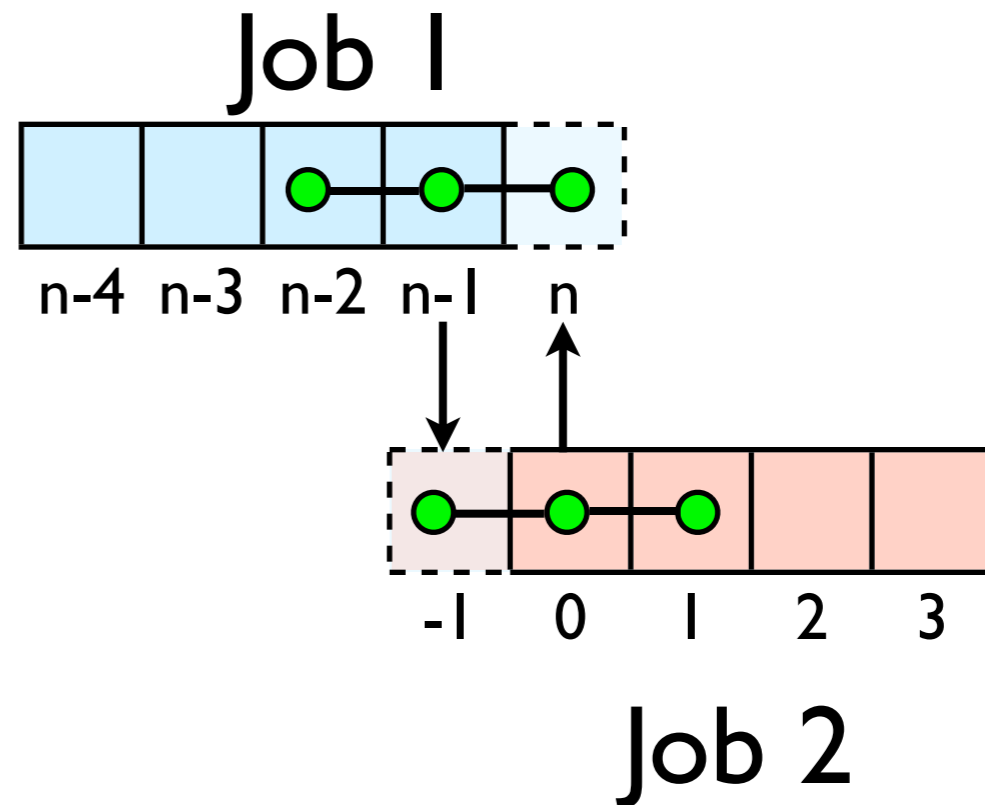
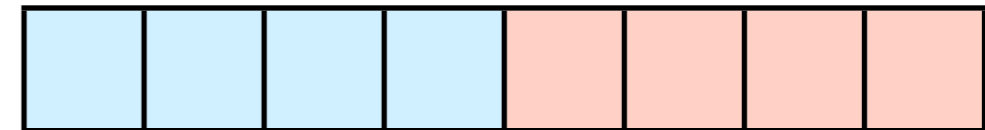


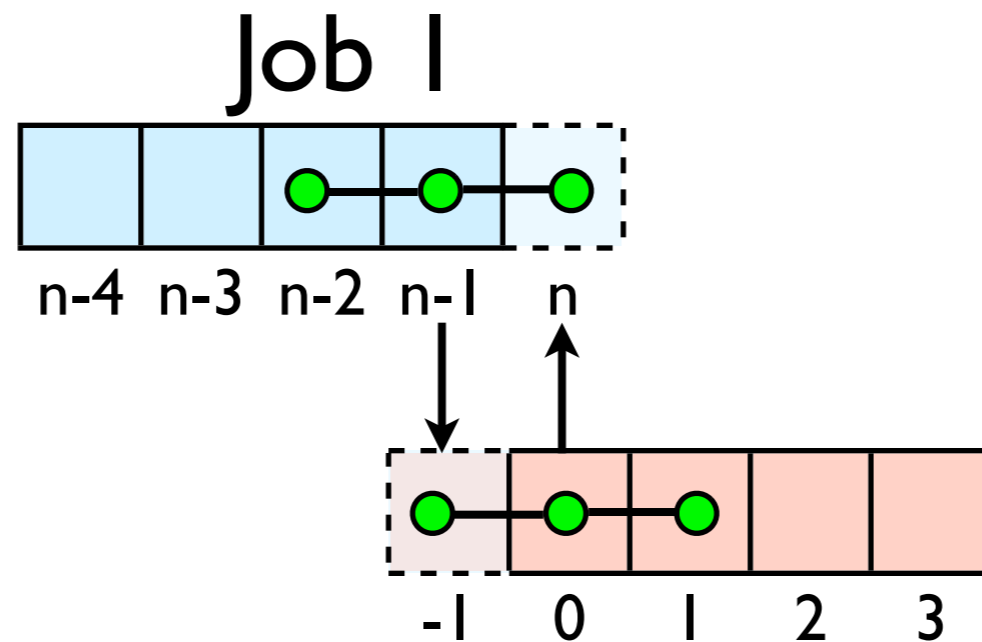
# Guardcells

- Works for parallel decomposition!
- Job 1 needs info on Job 2s 0th zone, Job 2 needs info on Job 1s last zone
- Pad array with 'guardcells' and fill them with the info from the appropriate node by message passing or shared memory
- Hydro code: need guardcells 2 deep



## Global Domain





- Do computation
- guardcell exchange: each cell has to do 2 sendrecv
  - its rightmost cell with neighbors leftmost
  - its leftmost cell with neighbors rightmost
- Use even/odd trick from before; if even do rightmost
- For simplicity, use periodic BCs -- task 0 is neighbors with task P-1



# Hands-on Assignment I: MPI diffusion

- `cp diffusion.f diffusion-mpi.f`
- Make an MPI-ed version of diffusion equation
- (Sample executable in completedexecutables)
- Test on 1..8 procs
- Copy file to `~/hw-intrompi/`

- add standard MPI calls: `init`, `finalize`, `comm_size`, `comm_rank`
- Figure out how many points PE is responsible for (`totpoints/size`)
- Figure out neighbors
- Start at 1, but end at `totpoints/size`
- `localxleft` and `localxright`
- At end of step, exchange guardcells; use `sendrecv`





# Hands-on

## Assignment II:

### MinMeanMax

- Run minmeanmax-sendrecv on 1,2,4 procs
- Test diffusion timing on 1,2,4,8 procs
- Put timing info in hw directory

```
$ cd mpi-intro
$ make clean
$ make minmeanmax-sendrecv
$ time -p mpirun -np 1 ./
minmeanmax-sendrecv
$ make diffusion-mpi
$ time -p mpirun -np 1 ./diffusion-
mpi
```



# Hands-on Assignment III: Use the Queue

- Copy files to assegai
- Log in
- Run minmeanmax-sendrecv on 1,2,4,8 nodes
- Test diffusion timing on 1,2,4,8 nodes
- Put results in hw directory

From gpc01..04

```
$ cp ~/pca/scripts/sample.pbs .  
$ make clean  
$ make diffusion-mpi USEPGPLOT=""  
PGPLIBS=""
```

```
$ qstat  
$ qsub sample.pbs  
$ qstat
```

to change number of nodes, change  
#PBS -l nodes=2:ppn=8...  
to 1,2,4,8  
and change  
-np 16  
as appropriate (8,16,32,64)



**sample.pbs:**

```
#!/bin/bash
```

```
# MOAB/Torque submission script for SciNet GPC (ethernet)
```

```
#
```

```
#PBS -l nodes=2:ppn=8,walltime=1:00:00
```

```
#PBS -N test
```

```
# DIRECTORY TO RUN - $PBS_O_WORKDIR is directory job was submitted from  
cd $PBS_O_WORKDIR
```

```
# EXECUTION COMMAND; -np nodes*ppn
```

```
mpirun -np 16 -hostfile $PBS_NODEFILE ./diffusion-mpi
```



## C syntax

```
MPI_Status status;
```

```
ierr = MPI_Init(&argc, &argv);
```

```
ierr = MPI_Comm_{size,rank}(Communicator, &{size,rank});
```

```
ierr = MPI_Send(sendptr, count, MPI_TYPE, destination,  
                tag, Communicator);
```

```
ierr = MPI_Recv(rcvptr, count, MPI_TYPE, source, tag,  
                Communicator, &status);
```

```
ierr = MPI_Sendrecv(sendptr, count, MPI_TYPE, destination, tag,  
                    rcvptr, count, MPI_TYPE, source, tag,  
                    Communicator, &status);
```

```
ierr = MPI_Allreduce(&mydata, &globaldata, count, MPI_TYPE,  
                    MPI_OP, Communicator);
```

Communicator -> MPI\_COMM\_WORLD

MPI\_Type -> MPI\_FLOAT, MPI\_DOUBLE, MPI\_INT, MPI\_CHAR...

MPI\_OP -> MPI\_SUM, MPI\_MIN, MPI\_MAX,...



# FORTRAN syntax

```
integer status(MPI_STATUS_SIZE)
```

```
call MPI_INIT(ierr)
```

```
call MPI_COMM_{SIZE,RANK}(Communicator, {size,rank},ierr)
```

```
call MPI_SSEND(sendarr, count, MPI_TYPE, destination,  
              tag, Communicator)
```

```
call MPI_RECV(rcvarr, count, MPI_TYPE, destination,tag,  
             Communicator, status, ierr)
```

```
call MPI_SENDRECV(sendptr, count, MPI_TYPE, destination,tag,  
                 recvptr, count, MPI_TYPE, source, tag,  
                 Communicator, status, ierr)
```

```
call MPI_ALLREDUCE(&mydata, &globaldata, count, MPI_TYPE,  
                 MPI_OP, Communicator, ierr)
```

Communicator -> MPI\_COMM\_WORLD

MPI\_Type -> MPI\_REAL, MPI\_DOUBLE\_PRECISION,  
 MPI\_INTEGER, MPI\_CHARACTER

MPI\_OP -> MPI\_SUM, MPI\_MIN, MPI\_MAX, ...

