

|epcc|

Collective Communications

- ▶ Communications involving a group of processes.
 - ▶ Called by all processes in a communicator.
 - ▶ Examples:
 - Barrier synchronisation.
 - Broadcast, scatter, gather.
 - Global sum, global maximum, etc.
-

- ▶ Collective action over a communicator.
 - ▶ All processes must communicate.
 - ▶ Synchronisation may or may not occur.
 - ▶ All collective operations are blocking.
 - ▶ No tags.
 - ▶ Receive buffers must be exactly the right size.
-

▶ C:

```
int MPI_Barrier (MPI_Comm comm)
```

▶ Fortran:

```
MPI_BARRIER (COMM, IERROR)  
INTEGER COMM, IERROR
```

▶ C:

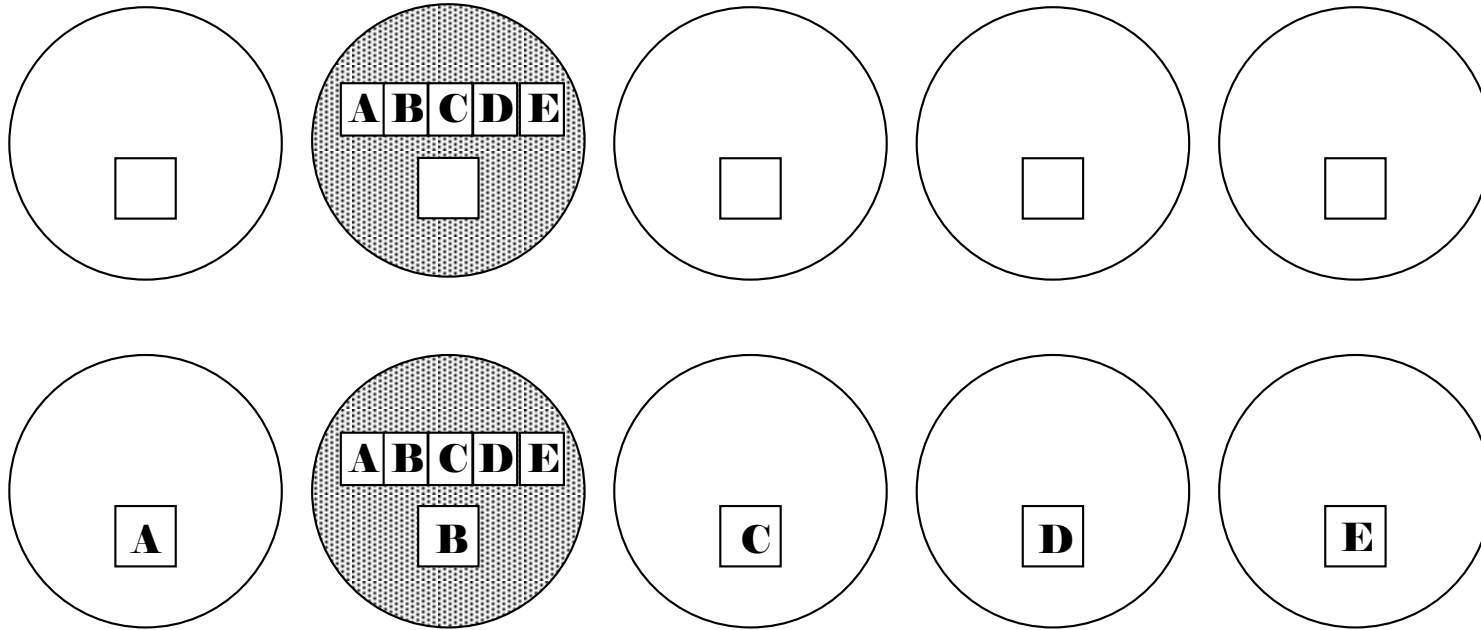
```
int MPI_Bcast (void *buffer, int count,  
              MPI_Datatype datatype, int root,  
              MPI_Comm comm)
```

▶ Fortran:

```
MPI_BCAST (BUFFER, COUNT, DATATYPE, ROOT,  
          COMM, IERROR)
```

```
<type> BUFFER(*)
```

```
INTEGER COUNT, DATATYPE, ROOT, COMM, IERROR
```



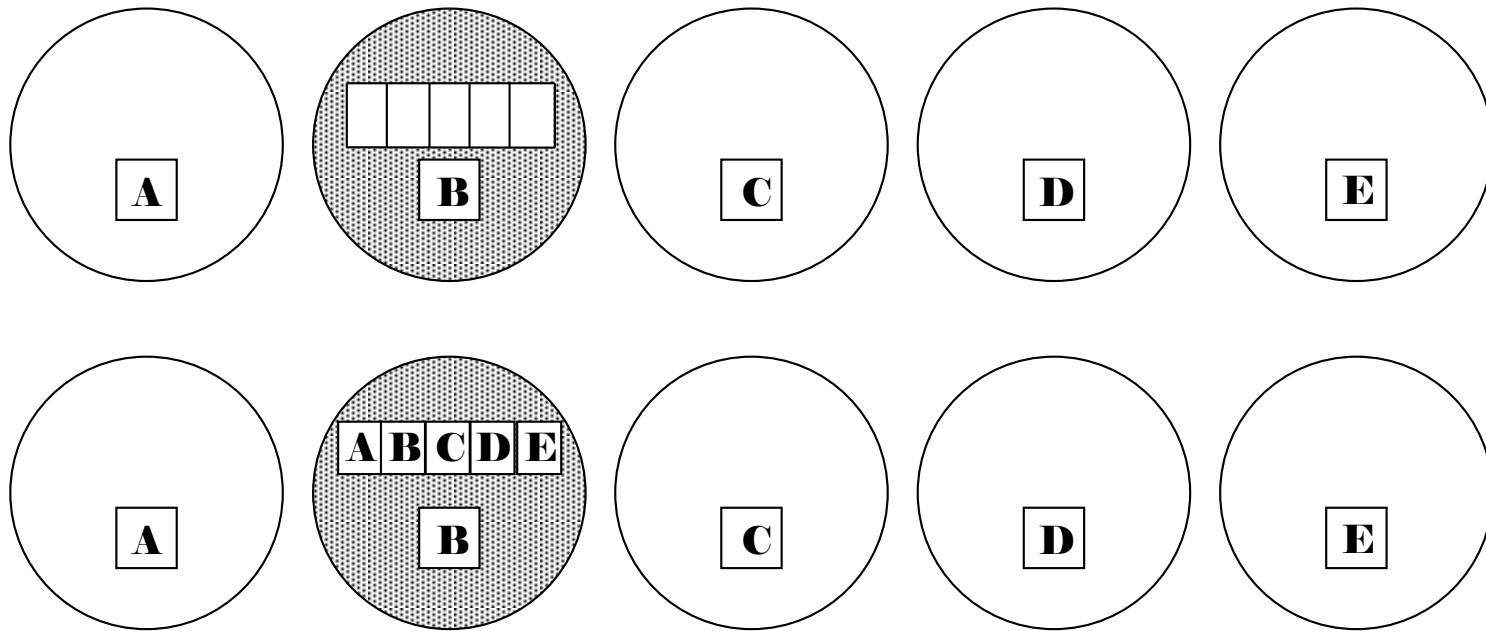
▶ C:

```
int MPI_Scatter(void *sendbuf,  
               int sendcount, MPI_Datatype sendtype,  
               void *recvbuf, int recvcount,  
               MPI_Datatype recvtype, int root,  
               MPI_Comm comm)
```

▶ Fortran:

```
MPI_SCATTER(SENDBUF, SENDCOUNT, SENDTYPE,  
            RECVBUF, RECVCOUNT, RECVTYPE,  
            ROOT, COMM, IERROR)
```

```
<type> SENDBUF, RECVBUF  
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT  
INTEGER RECVTYPE, ROOT, COMM, IERROR
```



▶ C:

```
int MPI_Gather(void *sendbuf, int sendcount,  
              MPI_Datatype sendtype, void *recvbuf,  
              int recvcount, MPI_Datatype recvtype,  
              int root, MPI_Comm comm)
```

▶ Fortran:

```
MPI_GATHER(SENDBUF, SENDCOUNT, SENDTYPE,  
           RECVBUF, RECVCOUNT, RECVTYPE,  
           ROOT, COMM, IERROR)
```

```
<type> SENDBUF, RECVBUF  
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT  
INTEGER RECVTYPE, ROOT, COMM, IERROR
```

- ▶ Used to compute a result involving data distributed over a group of processes.
 - ▶ Examples:
 - global sum or product
 - global maximum or minimum
 - global user-defined operation
-

MPI Name	Function
MPI_MAX	Maximum
MPI_MIN	Minimum
MPI_SUM	Sum
MPI_PROD	Product
MPI_LAND	Logical AND
MPI_BAND	Bitwise AND
MPI_LOR	Logical OR
MPI_BOR	Bitwise OR
MPI_LXOR	Logical exclusive OR
MPI_BXOR	Bitwise exclusive OR
MPI_MAXLOC	Maximum and location
MPI_MINLOC	Minimum and location

▶ C:

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

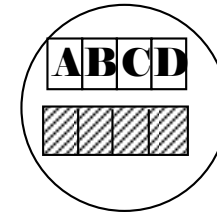
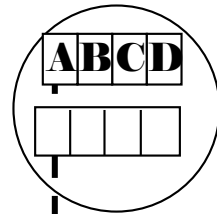
▶ Fortran:

```
MPI_REDUCE(SENDBUF, RECVBUF, COUNT,  
           DATATYPE, OP, ROOT, COMM, IERROR)
```

```
<type> SENDBUF, RECVBUF  
INTEGER SENDCOUNT, SENDTYPE, RECVCOUNT  
INTEGER RECVMODE, ROOT, COMM, IERROR
```

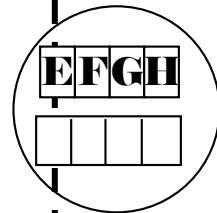
Rank

0

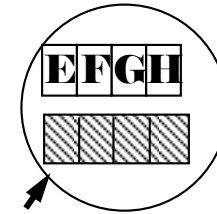
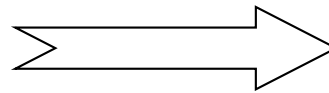


1

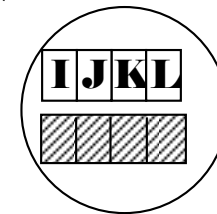
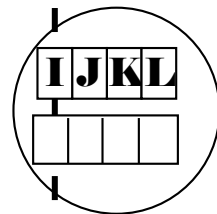
Root



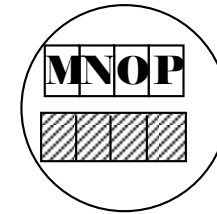
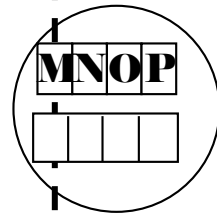
MPI_REDUCE



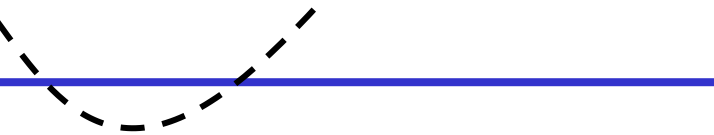
2



3



AoEoIoM



Integer global sum

▶ C:

```
MPI_Reduce(&x, &result, 1, MPI_INT,  
           MPI_SUM, 0, MPI_COMM_WORLD)
```

▶ Fortran:

```
CALL MPI_REDUCE(x, result, 1, MPI_INTEGER,  
               MPI_SUM, 0,  
               MPI_COMM_WORLD, IERROR)
```

- ▶ Sum of all the x values is placed in *result*.
 - ▶ The result is only placed there on processor 0.
-

- ▶ Reducing using an arbitrary operator, ○
- ▶ C - function of type MPI_User_function:

```
void my_op (void *invec,  
            void *inoutvec, int *len,  
            MPI_Datatype *datatype)
```

- ▶ Fortran - external subprogram of type

```
SUBROUTINE MY_OP (INVEC (*), INOUTVEC (*),  
                 LEN, DATATYPE)  
<type> INVEC (LEN), INOUTVEC (LEN)  
INTEGER LEN, DATATYPE
```

- ▶ Operator function for \circ must act as:

```
for (i = 1 to len)
```

```
    inoutvec(i) = inoutvec(i)  $\circ$  invec(i)
```

- ▶ Operator \circ need not commute but must be associative.
-

- ▶ Operator handles have type `MPI_Op` or `INTEGER`
- ▶ C:

```
int MPI_Op_create(MPI_User_function *my_op,  
                 int commute, MPI_Op *op)
```

- ▶ Fortran:

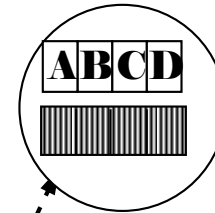
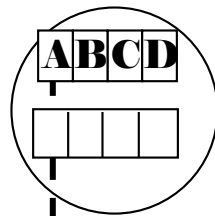
```
MPI_OP_CREATE (MY_OP, COMMUTE, OP, IERROR)
```

```
EXTERNAL MY_OP  
LOGICAL COMMUTE  
INTEGER OP, IERROR
```

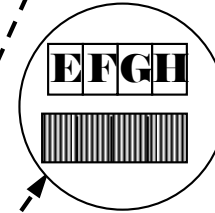
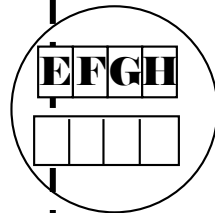
- ▶ MPI_Allreduce no root process
 - ▶ MPI_Reduce_scatter result is scattered
 - ▶ MPI_Scan “parallel prefix”
-

Rank

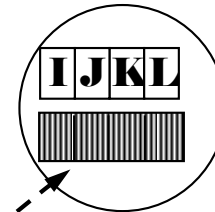
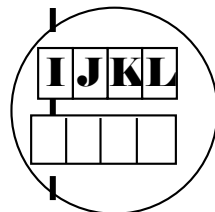
0



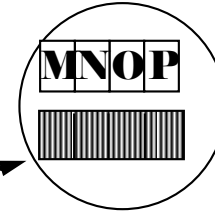
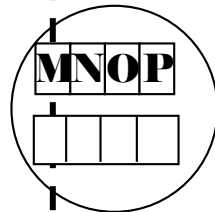
1



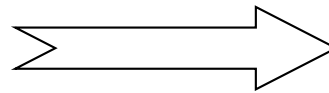
2



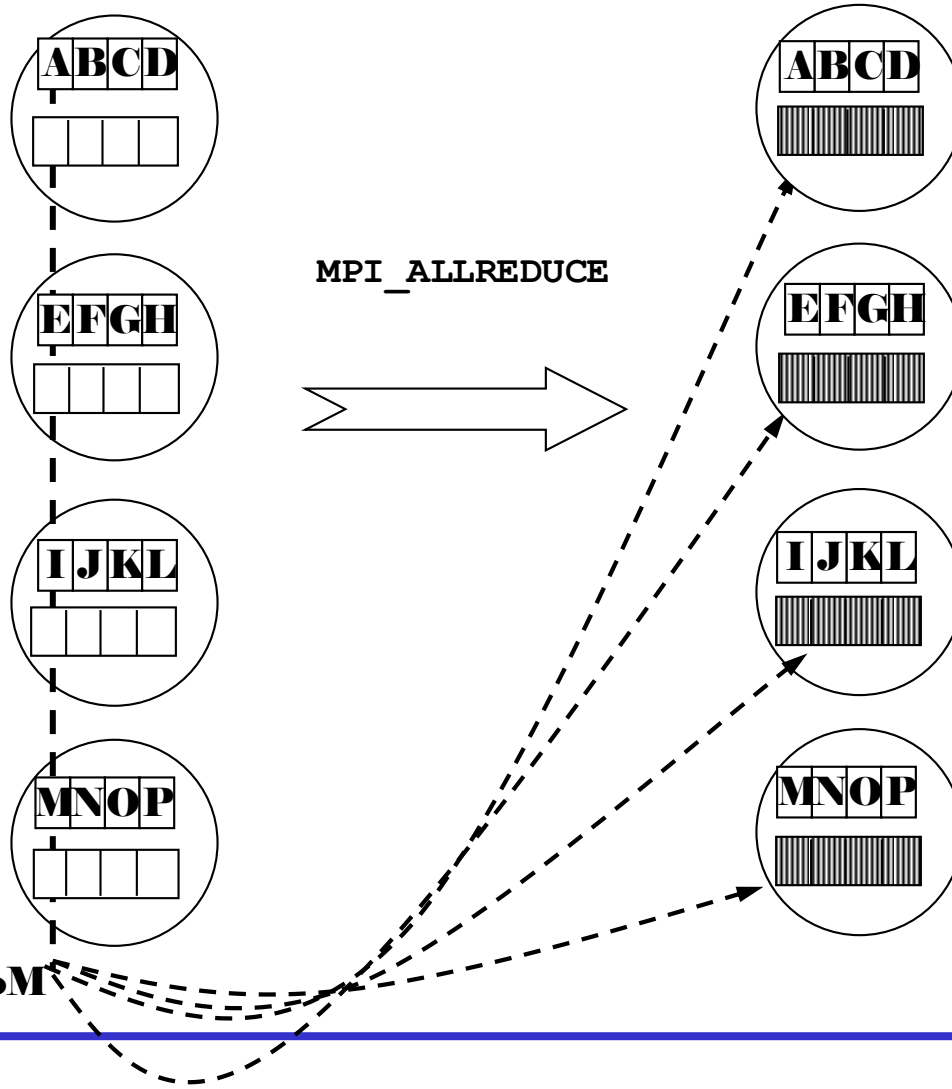
3



MPI_ALLREDUCE



AoEoIoM



Integer global sum

▶ C:

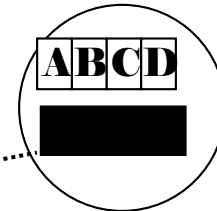
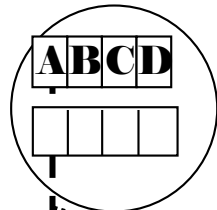
```
int MPI_Allreduce(void* sendbuf,  
                 void* recvbuf, int count,  
                 MPI_Datatype datatype,  
                 MPI_Op op, MPI_Comm comm)
```

▶ Fortran:

```
MPI_ALLREDUCE(SENDBUF, RECVBUF, COUNT,  
              DATATYPE, OP, COMM, IERROR)
```

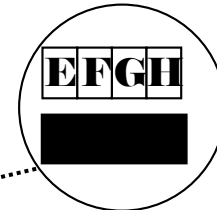
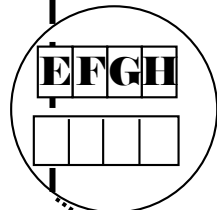
Rank

0

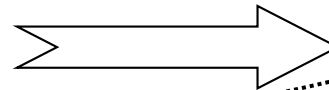


A

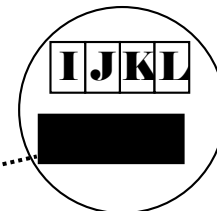
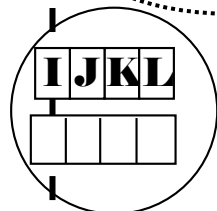
1



MPI_SCAN



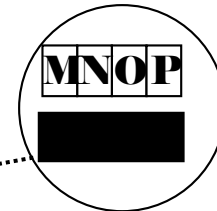
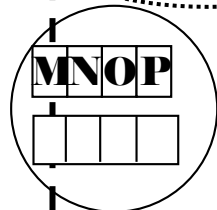
2



AoE

AoEoI

3



AoEoIoM



Integer partial sum

▶ C:

```
int MPI_Scan(void* sendbuf, void* recvbuf,  
            int count, MPI_Datatype datatype,  
            MPI_Op op, MPI_Comm comm)
```

▶ Fortran:

```
MPI_SCAN(SENDBUF, RECVBUF, COUNT,  
         DATATYPE, OP, COMM, IERROR)
```

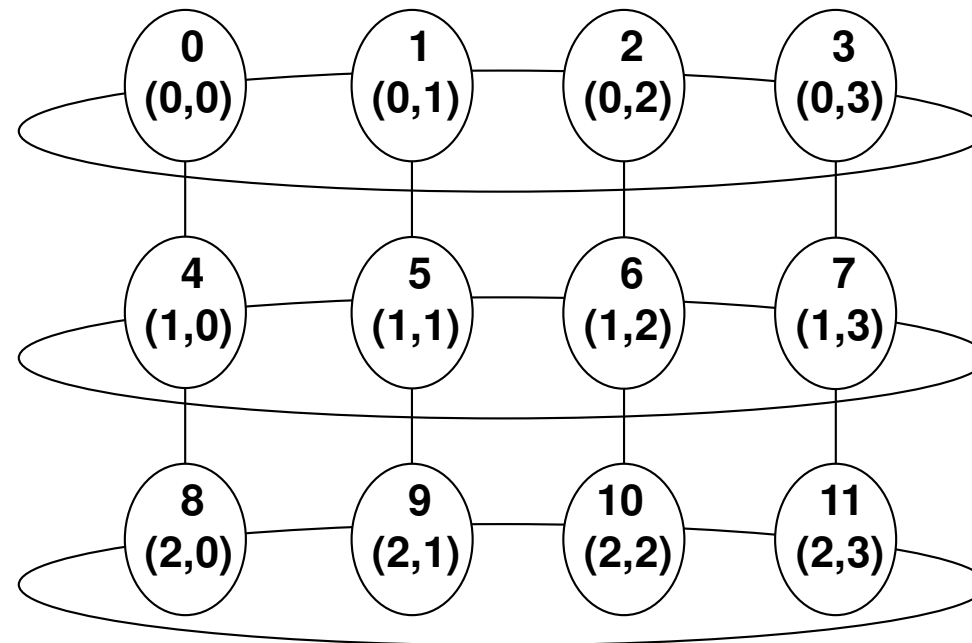
|epcc|

Virtual Topologies

- ▶ Convenient process naming.
 - ▶ Naming scheme to fit the communication pattern.
 - ▶ Simplifies writing of code.
 - ▶ Can allow MPI to optimise communications.
-

- ▶ Creating a topology produces a new communicator.
 - ▶ MPI provides ``mapping functions".
 - ▶ Mapping functions compute processor ranks, based on the topology naming scheme.
-

A 2-dimensional Cylinder



▶ Cartesian topologies

- each process is “connected” to its neighbours in a virtual grid.
 - boundaries can be cyclic, or not.
 - optionally re-order ranks to allow MPI implementation to optimise for underlying network interconnectivity.
- processes are identified by cartesian coordinates.

▶ Graph topologies

- general graphs
 - not covered here
-

▶ C:

```
int MPI_Cart_create(MPI_Comm comm_old,  
                   int ndims, int *dims, int *periods,  
                   int reorder, MPI_Comm *comm_cart)
```

▶ Fortran:

```
MPI_CART_CREATE(COMM_OLD, NDIMS, DIMS,  
                PERIODS, REORDER, COMM_CART, IERROR)
```

```
INTEGER COMM_OLD, NDIMS, DIMS(*), COMM_CART, IERROR  
LOGICAL PERIODS(*), REORDER
```

▶ C:

```
int MPI_Dims_create(int nnodes, int ndims,  
                   int *dims)
```

▶ Fortran:

```
MPI_DIMS_CREATE(NNODES, NDIMS, DIMS, IERROR)
```

```
INTEGER NNODES, NDIMS, DIMS(*), IERROR
```

- ▶ Call tries to set dimensions as close to each other as possible

dims before the call	function call	dims on return
(0, 0)	MPI_DIMS_CREATE(6, 2, dims)	(3, 2)
(0, 0)	MPI_DIMS_CREATE(7, 2, dims)	(7, 1)
(0, 3, 0)	MPI_DIMS_CREATE(6, 3, dims)	(2, 3, 1)
(0, 3, 0)	MPI_DIMS_CREATE(7, 3, dims)	erroneous call

- ▶ Non zero values in dims sets the number of processors required in that direction.
 - **WARNING:- make sure dims is set to 0 before the call!**
-

Mapping process grid coordinates to ranks

▶ C:

```
int MPI_Cart_rank(MPI_Comm comm,  
                 int *coords, int *rank)
```

▶ Fortran:

```
MPI_CART_RANK (COMM, COORDS, RANK, IERROR)
```

```
INTEGER COMM, COORDS (*), RANK, IERROR
```

Mapping ranks to process grid coordinates

▶ C:

```
int MPI_Cart_coords(MPI_Comm comm, int rank,  
                   int maxdims, int *coords)
```

▶ Fortran:

```
MPI_CART_COORDS(COMM, RANK, MAXDIMS, COORDS,  
                IERROR)
```

```
INTEGER COMM, RANK, MAXDIMS, COORDS(*), IERROR
```

Computing ranks of my neighbouring processes
Following conventions of MPI_SendRecv

▶ C:

```
int MPI_Cart_shift(MPI_Comm comm,  
                  int direction, int disp,  
                  int *rank_source, int *rank_dest)
```

▶ Fortran:

```
MPI_CART_SHIFT(COMM, DIRECTION, DISP,  
               RANK_SOURCE, RANK_DEST, IERROR)
```

```
INTEGER COMM, DIRECTION, DISP,  
        RANK_SOURCE, RANK_DEST, IERROR
```

- ▶ What if you ask for the rank of a non-existent process?
 - or look off the edge of a non-periodic grid?
 - ▶ MPI returns a NULL processor
 - rank is `MPI_PROC_NULL`
 - ▶ `MPI_PROC_NULL` is a black hole
 - sends and receives complete immediately
 - send buffer disappears, receive buffer isn't touched
 - like UNIX `/dev/null`
-

- ▶ Cut a grid up into “slices”.
 - ▶ A new communicator is produced for each slice.
 - ▶ Each slice can then perform its own collective communications.
 - ▶ `MPI_Cart_sub` and `MPI_CART_SUB` generate new communicators for the slices.
 - Use array to specify which dimensions should be retained in the new communicator.
-

▶ C:

```
int MPI_Cart_sub (MPI_Comm comm,  
                 int *remain_dims,  
                 MPI_Comm *newcomm)
```

▶ Fortran:

```
MPI_CART_SUB (COMM, REMAIN_DIMS,  
              NEWCOMM, IERROR)
```

```
INTEGER COMM, NEWCOMM, IERROR  
LOGICAL REMAIN_DIMS (*)
```

```

1 #include <stdio.h>
2 #define SIZE 16
3 #define UP 0
4 #define DOWN 1
5 #define LEFT 2
6 #define RIGHT 3
7
8
9 main(int argc, char *argv[]) {
10 int numtasks, rank, source, dest, outbuf, i, tag=1,
11     inbuf[4]={MPI_PROC_NULL,MPI_PROC_NULL,MPI_PROC_NULL,MPI_PROC_NULL},
12     nbrs[4], dims[2]={4,4},
13     periods[2]={0,0}, reorder=0, coords[2];
14
15 MPI_Request reqs[8];
16 MPI_Status stats[8];
17 MPI_Comm cartcomm; // required variable
18
19 MPI_Init(&argc,&argv);
20 MPI_Comm_size(MPI_COMM_WORLD, &numtasks);
21
22 if (numtasks == SIZE) {
23     // create cartesian virtual topology, get rank, coordinates, neighbor ranks
24     MPI_Cart_create(MPI_COMM_WORLD, 2, dims, periods, reorder, &cartcomm);
25     MPI_Comm_rank(cartcomm, &rank);
26     MPI_Cart_coords(cartcomm, rank, 2, coords);
27     MPI_Cart_shift(cartcomm, 0, 1, &nbrs[UP], &nbrs[DOWN]);
28     MPI_Cart_shift(cartcomm, 1, 1, &nbrs[LEFT], &nbrs[RIGHT]);
29
30     printf("rank= %d coords= %d %d neighbors(u,d,l,r)= %d %d %d %d\n",
31           rank, coords[0], coords[1], nbrs[UP], nbrs[DOWN], nbrs[LEFT],
32           nbrs[RIGHT]);
33
34     outbuf = rank;
35
36     // exchange data (rank) with 4 neighbors
37     for (i=0; i<4; i++) {
38         dest = nbrs[i];
39         source = nbrs[i];
40         MPI_Isend(&outbuf, 1, MPI_INT, dest, tag,
41                 MPI_COMM_WORLD, &reqs[i]);
42         MPI_Irecv(&inbuf[i], 1, MPI_INT, source, tag,
43                 MPI_COMM_WORLD, &reqs[i+4]);
44     }
45
46     MPI_Waitall(8, reqs, stats);
47
48     printf("rank= %d inbuf(u,d,l,r)= %d %d %d %d\n",
49           rank, inbuf[UP], inbuf[DOWN], inbuf[LEFT], inbuf[RIGHT]); }
50 else
51     printf("Must specify %d processors. Terminating.\n", SIZE);
52
53 MPI_Finalize();
54 }

```

Sample program output: (partial)

```

rank= 0 coords= 0 0 neighbors(u,d,l,r)= -1 4 -1 1
rank= 0 inbuf(u,d,l,r)= -1 4 -1 1
rank= 8 coords= 2 0 neighbors(u,d,l,r)= 4 12 -1 9
rank= 8 inbuf(u,d,l,r)= 4 12 -1 9
rank= 1 coords= 0 1 neighbors(u,d,l,r)= -1 5 0 2
rank= 1 inbuf(u,d,l,r)= -1 5 0 2
rank= 13 coords= 3 1 neighbors(u,d,l,r)= 9 -1 12 14
rank= 13 inbuf(u,d,l,r)= 9 -1 12 14
...
rank= 3 coords= 0 3 neighbors(u,d,l,r)= -1 7 2 -1
rank= 3 inbuf(u,d,l,r)= -1 7 2 -1
rank= 11 coords= 2 3 neighbors(u,d,l,r)= 7 15 10 -1
rank= 11 inbuf(u,d,l,r)= 7 15 10 -1
rank= 10 coords= 2 2 neighbors(u,d,l,r)= 6 14 9 11
rank= 10 inbuf(u,d,l,r)= 6 14 9 11
rank= 9 coords= 2 1 neighbors(u,d,l,r)= 5 13 8 10
rank= 9 inbuf(u,d,l,r)= 5 13 8 10

```

Calendar till Minor 2

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
Interconnection, MPI	10	11 assign2 released	12	13	14	15	16
Synchronization	17	18 assign1, minor1	19	20	21	22	23
Memory consistency	24	25	26	27	28	29	1
	2	3 assign2 due	4	5	6	7	8
Problem solving	9	10	11	12	13	14	15
	16	17	18	19	20	21	22