



# Characteristics of collective operations



[www.ichec.ie](http://www.ichec.ie)



An Roinn Post, Fiontar agus Nuálaitheolaíochta  
Department of Jobs, Enterprise and Innovation



SCIENCE FUNDING INITIATIVE  
INDEPENDENTLY ADMINISTERED BY SFI



AN ROINN  
OIDEACHAIS AGUS SCILEANNA  
DEPARTMENT OF  
EDUCATION AND SKILLS

HEA

Higher Education Authority  
an tArd-Chomhairle

# Characteristics of collective operations

- involve all the processes grouped in a communicator

# Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator

# Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...



# Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...
- optimised communication for performance

# Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...
- optimised communication for performance
- code readability and maintainability

# Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...
- optimised communication for performance
- code readability and maintainability
- data synchronisation or temporal synchronisation if barriers are used

## Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...
- optimised communication for performance
- code readability and maintainability
- data synchronisation or temporal synchronisation if barriers are used
- up to MPI 3.0 they were only blocking



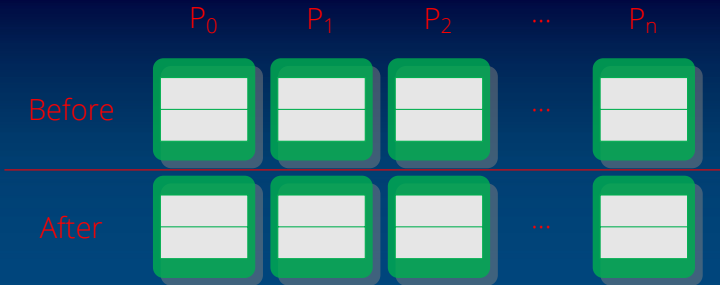
# Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...
- optimised communication for performance
- code readability and maintainability
- data synchronisation or temporal synchronisation if barriers are used
- up to MPI 3.0 they were only blocking
- receive buffers must have the same size as send buffers

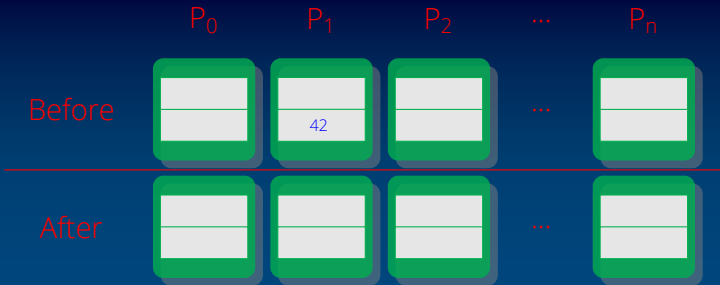
## Characteristics of collective operations

- involve all the processes grouped in a communicator
- must be called by all the processes in the communicator
- can be done with point-to-point operations but...
- optimised communication for performance
- code readability and maintainability
- data synchronisation or temporal synchronisation if barriers are used
- up to MPI 3.0 they were only blocking
- receive buffers must have the same size as send buffers
- **Examples:** broadcast, scatter, gather, global sum, global maximum, barrier synchronisation...

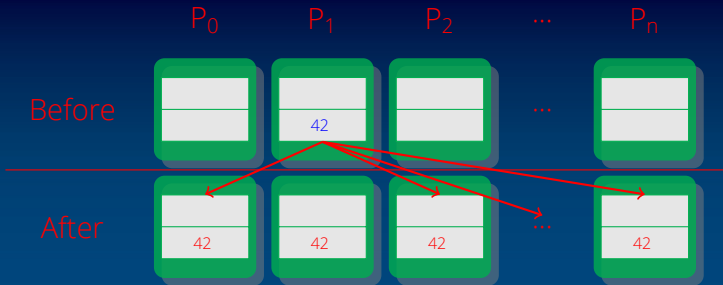
# broadcast



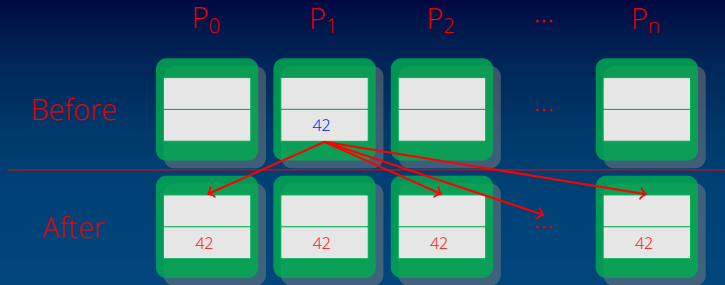
## broadcast



# broadcast

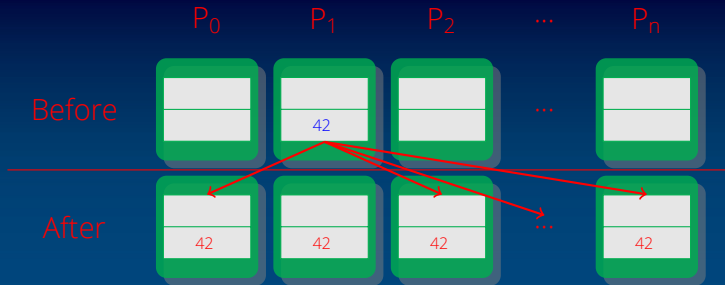


## broadcast



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

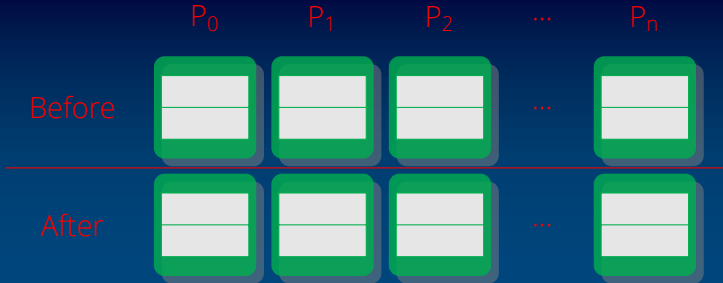
## broadcast



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

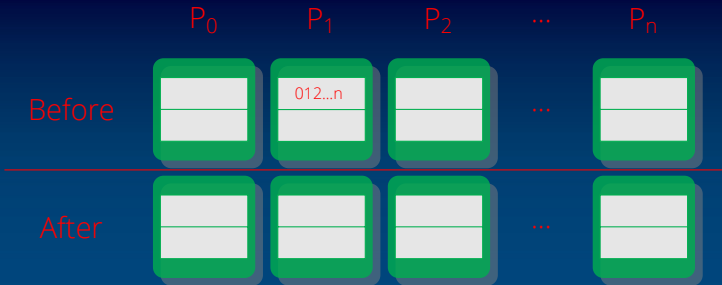
```
|| MPI_Bcast(buffer, count, datatype, root, comm, ierror)
|| <type>    :: buffer(:)
|| integer  :: count, datatype, root, comm, ierror
```

# scatter

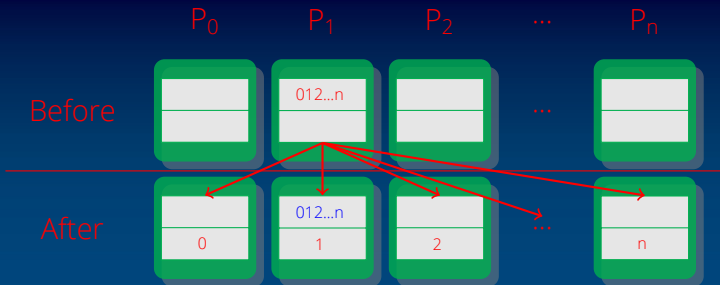




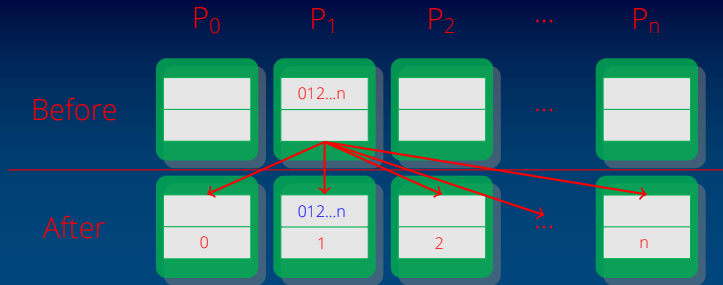
## scatter



## scatter



## scatter

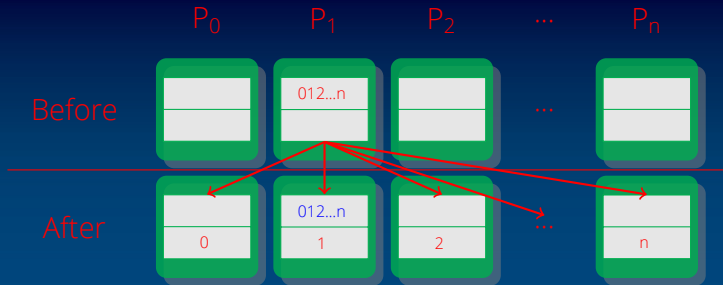


```

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

```

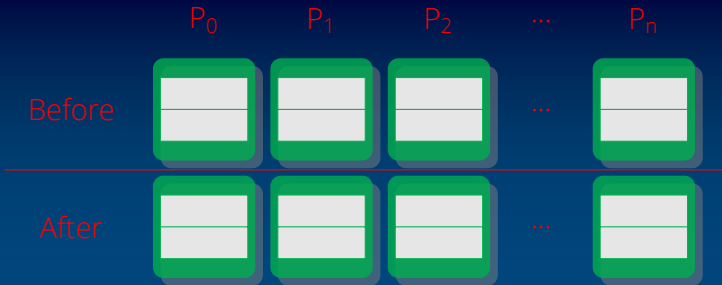
# scatter



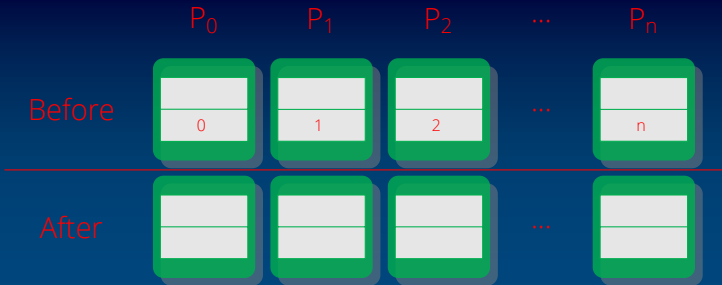
```
int MPI_Scatter(void *sendbuf, int sendcount, MPI_Datatype
sendtype, void *recvbuf, int recvcount, MPI_Datatype
recvtype, int root, MPI_Comm comm)
```

```
MPI_Scatter(sendbuf, sendcount, sendtype, recvbuf, recvcount,
recvtype, root, comm, ierror)
<type> :: sendbuf(:), recvbuf(:)
integer :: sendcount, sendtype, recvcount, recvtype, root
integer :: comm, ierror
```

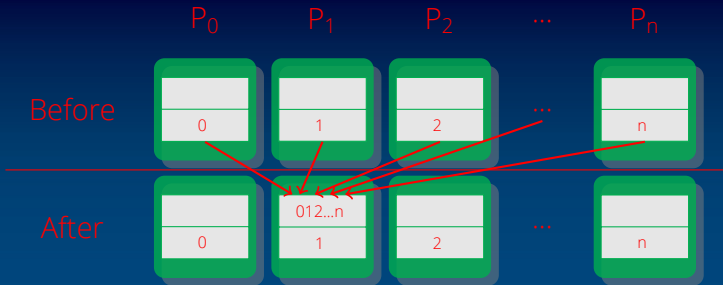
# gather



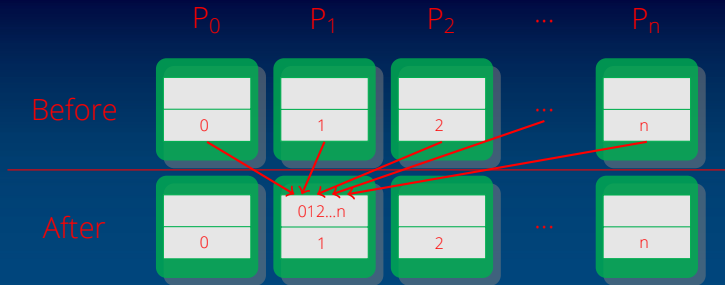
## gather



# gather



## gather

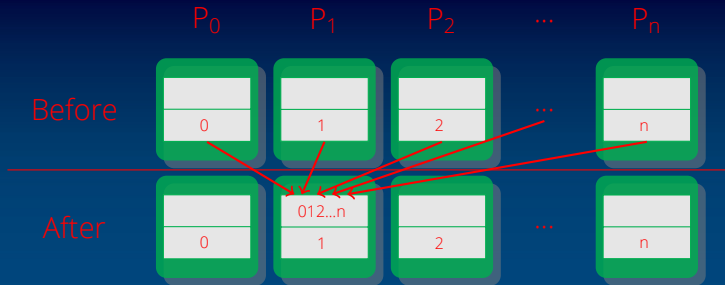


```

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



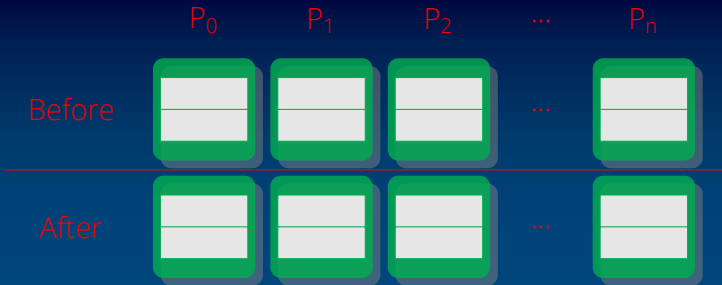
## gather



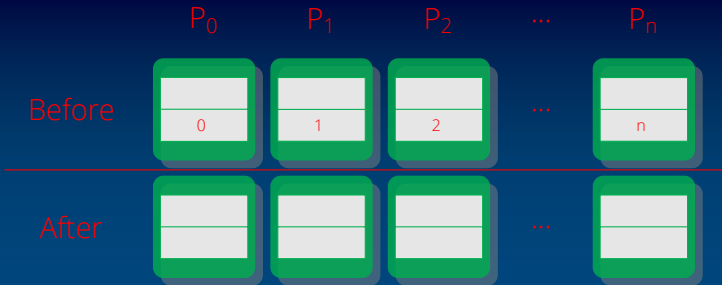
```
int MPI_Gather(void *sendbuf, int sendcnt, MPI_Datatype
sendtype, void *recvbuf, int recvcnt, MPI_Datatype recvtype
, int root, MPI_Comm comm)
```

```
MPI_Gather(sendbuf, sendcount, sendtype, recvbuf, recvcount,
recvtype, root, comm, ierror)
<type> :: sendbuf(:), recvbuf(:)
integer :: sendcount, sendtype, recvcount, recvtype, root
integer :: comm, ierror
```

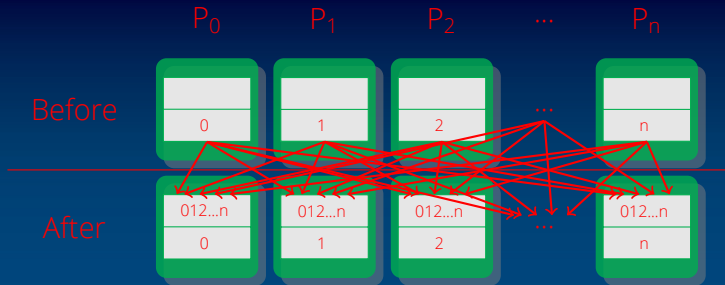
# AllGather



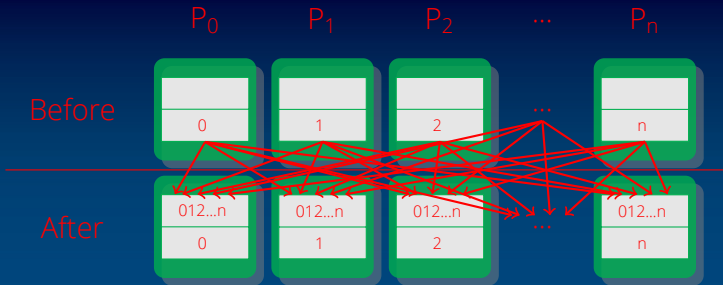
# AllGather



# AllGather

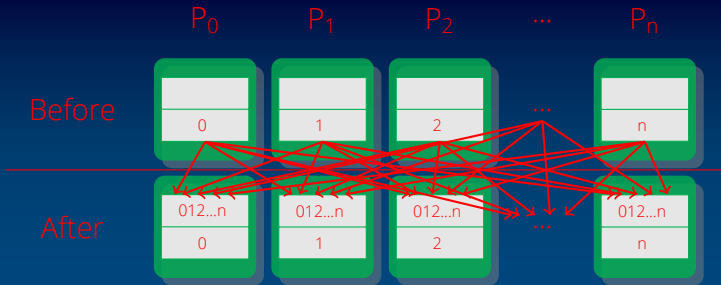


# AllGather



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

# AllGather



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

```
MPI_Allgather(sendbuf, sendcount, sendtype, recvbuf, recvcount,
recvtype, comm, ierror)
<type> :: sendbuf (:), recvbuf (:)
integer :: sendcount, sendtype, recvcount, recvtype, comm,
ierror
```



# AlltoAll



[www.ichec.ie](http://www.ichec.ie)



An Roinn Post, Fiontar agus Nuálachas  
Department of Jobs, Enterprise and Innovation



AN ROINN  
OIDEACHAIS AGUS SCILEANNA  
DEPARTMENT OF  
EDUCATION AND SKILLS

HEA

Higher Education Authority  
an tArd-Chomhairle

# AlltoAll

$P_0$	1	2	3	4	5	6	7	8
$P_1$	9	10	11	12	13	14	15	16
$P_2$	17	18	19	20	21	22	23	24
$P_3$	25	26	27	28	29	30	31	32



# AlltoAll

$P_0$	1	2	3	4	5	6	7	8
$P_1$	9	10	11	12	13	14	15	16
$P_2$	17	18	19	20	21	22	23	24
$P_3$	25	26	27	28	29	30	31	32

1	2
3	4
5	6
7	8

## AlltoAll

$P_0$	1	2	3	4	5	6	7	8
$P_1$	9	10	11	12	13	14	15	16
$P_2$	17	18	19	20	21	22	23	24
$P_3$	25	26	27	28	29	30	31	32

1	2	9	10
3	4	11	12
5	6	13	14
7	8	15	16

## AlltoAll

$P_0$	1	2	3	4	5	6	7	8
$P_1$	9	10	11	12	13	14	15	16
$P_2$	17	18	19	20	21	22	23	24
$P_3$	25	26	27	28	29	30	31	32

1	2	9	10	17	18
3	4	11	12	19	20
5	6	13	14	21	22
7	8	15	16	23	24

# AlltoAll

$P_0$	1	2	3	4	5	6	7	8
$P_1$	9	10	11	12	13	14	15	16
$P_2$	17	18	19	20	21	22	23	24
$P_3$	25	26	27	28	29	30	31	32

1	2	9	10	17	18	25	26
3	4	11	12	19	20	27	28
5	6	13	14	21	22	29	30
7	8	15	16	23	24	31	32

# AlltoAll

P <sub>0</sub>	1	2	3	4	5	6	7	8
P <sub>1</sub>	9	10	11	12	13	14	15	16
P <sub>2</sub>	17	18	19	20	21	22	23	24
P <sub>3</sub>	25	26	27	28	29	30	31	32

1	2	9	10	17	18	25	26
3	4	11	12	19	20	27	28
5	6	13	14	21	22	29	30
7	8	15	16	23	24	31	32

```

int MPI_Alltoall(void *sendbuf, int sendcount, MPI_Datatype
sendtype, void *recvbuf, int recvcount, MPI_Datatype
recvtype, MPI_Comm comm)

```

# AlltoAll

P <sub>0</sub>	1	2	3	4	5	6	7	8
P <sub>1</sub>	9	10	11	12	13	14	15	16
P <sub>2</sub>	17	18	19	20	21	22	23	24
P <sub>3</sub>	25	26	27	28	29	30	31	32

1	2	9	10	17	18	25	26
3	4	11	12	19	20	27	28
5	6	13	14	21	22	29	30
7	8	15	16	23	24	31	32

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

```
MPI_Alltoall(sendbuf, sendcount, sendtype, recvbuf, recvcount,
recvtype, comm, ierror)
<type>      :: sendbuf(:), recvbuf(:)
integer     :: sendcount, sendtype, recvcount, recvtype
integer     :: comm, ierror
```

# reduce

$P_0$

$P_1$

$P_2$

$P_3$

$P_4$

Before



+

After



## reduce

$P_0$

$P_1$

$P_2$

$P_3$

$P_4$

Before



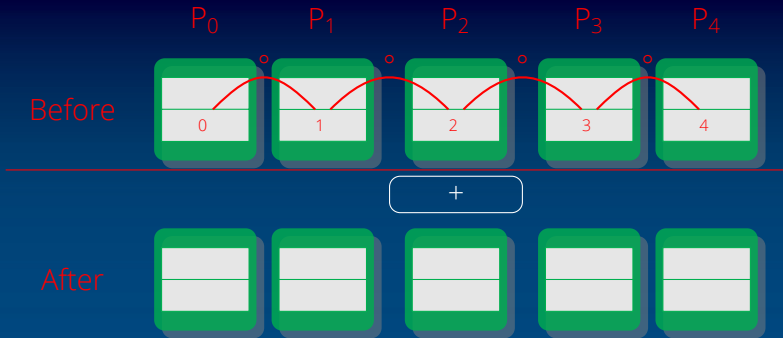
+

After

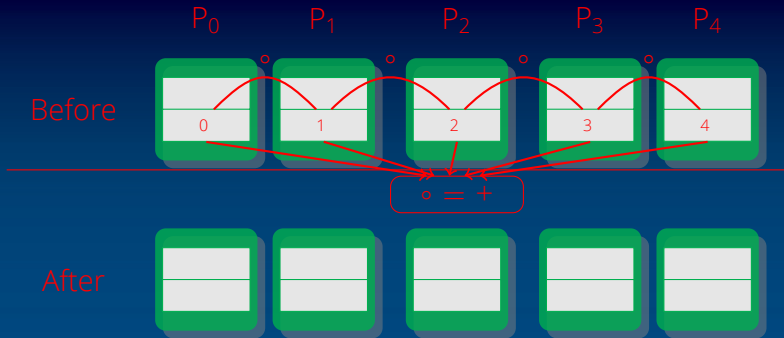




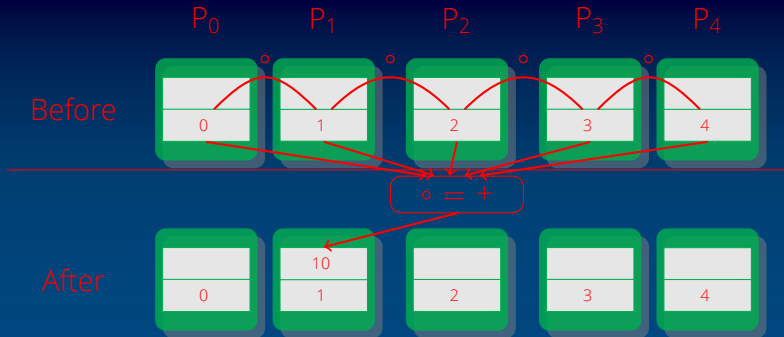
## reduce



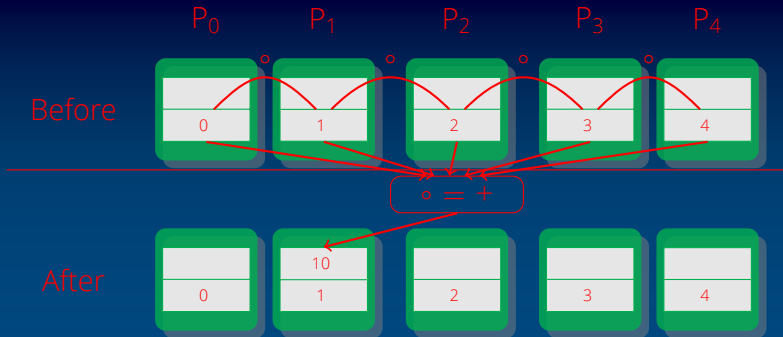
## reduce



## reduce

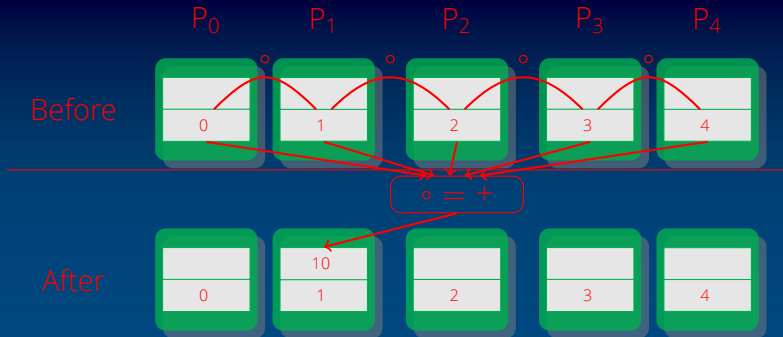


## reduce



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

## reduce



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

```
MPI_Reduce(sendbuf, recvbuf, count, datatype, op, root, comm,
ierror)
<type>    :: sendbuf(:), recvbuf(:)
integer   :: count, datatype, op, root, comm, ierror
```

# allreduce

 $P_0$  $P_1$  $P_2$  $P_3$  $P_4$ 

Before



+

After



# allreduce

 $P_0$  $P_1$  $P_2$  $P_3$  $P_4$ 

Before

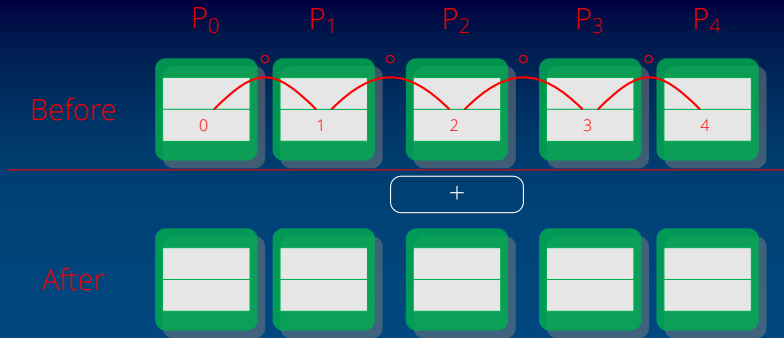


+

After

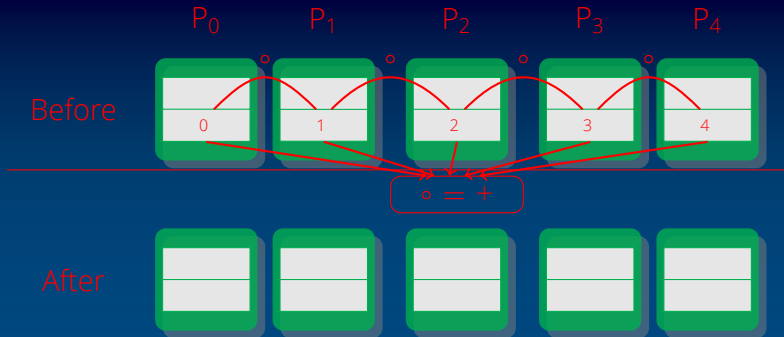


# allreduce

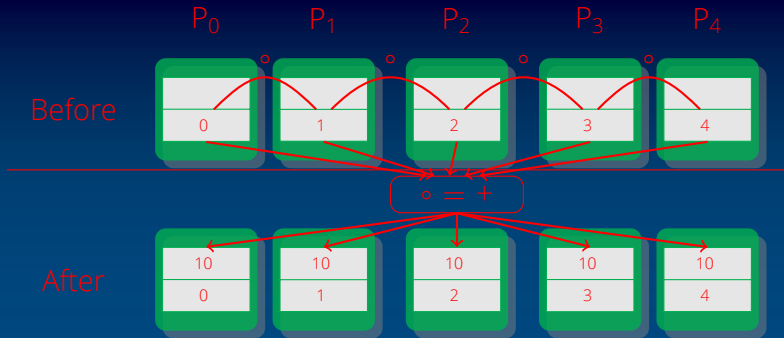




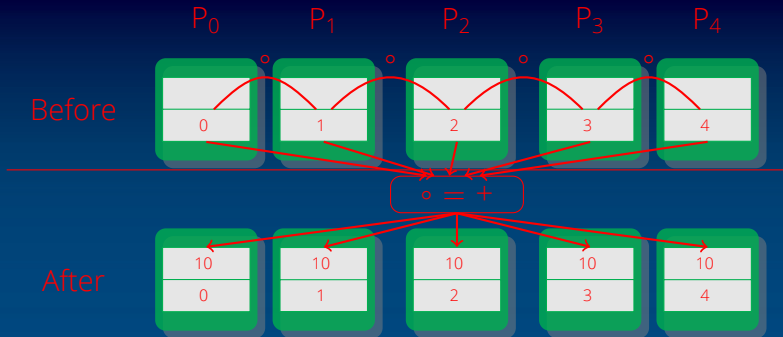
# allreduce



# allreduce



## allreduce

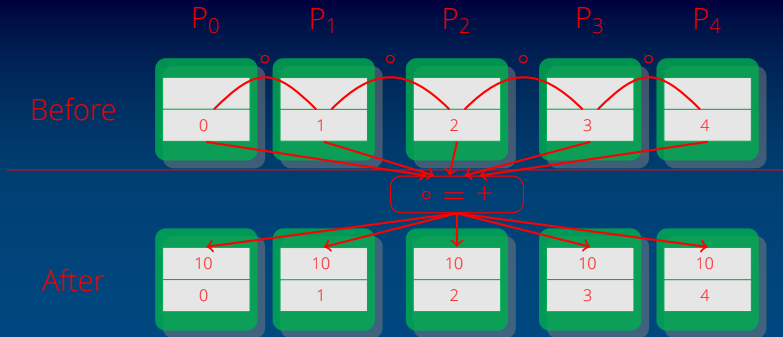


```

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

```

## allreduce



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

```
MPI_Allreduce(sendbuf, recvbuf, count, datatype, op, comm,
ierror)
<type>      :: sendbuf(*), recvbuf(*)
integer     :: count, datatype, op, comm, ierror
```

scan

 $P_0$  $P_1$  $P_2$  $P_3$  $P_4$ 

Before



+

After



scan

$P_0$

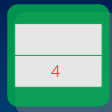
$P_1$

$P_2$

$P_3$

$P_4$

Before

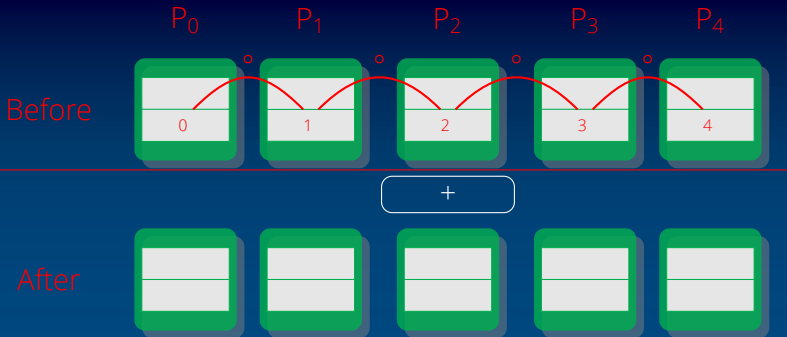


+

After



scan



scan

 $P_0$  $P_1$  $P_2$  $P_3$  $P_4$ 

Before

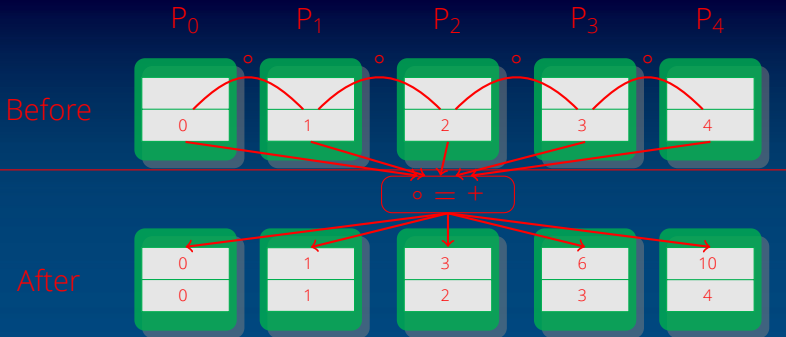


After



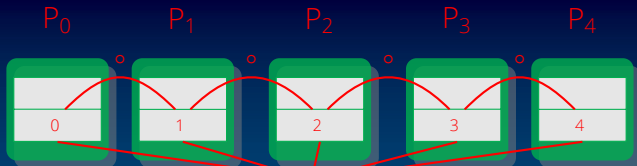


## scan

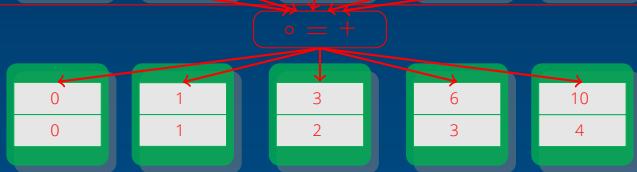


## scan

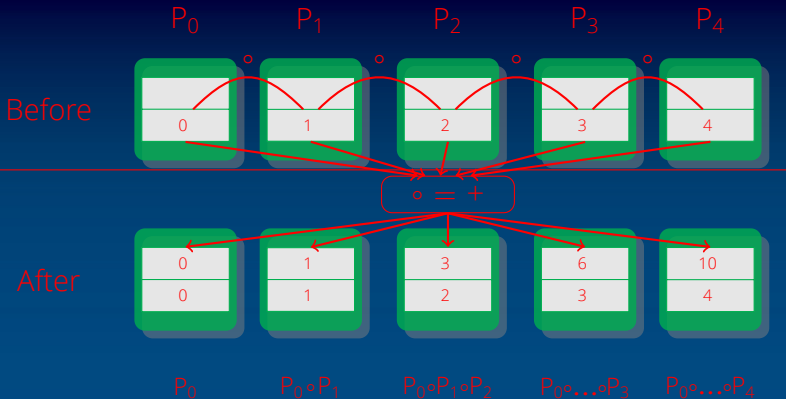
Before



After

 $P_0$  $P_0 \circ P_1$  $P_0 \circ P_1 \circ P_2$  $P_0 \circ \dots \circ P_3$  $P_0 \circ \dots \circ P_4$

## scan

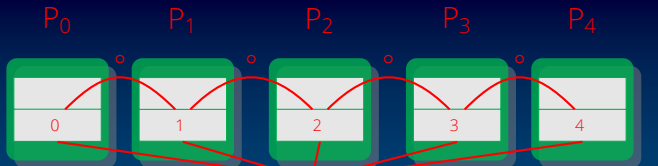


```

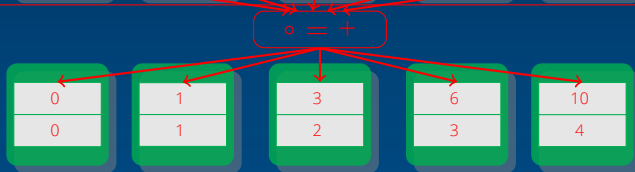
int MPI_Scan(void *sendbuf, void *recvbuf, int count,
             MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
  
```

## scan

Before



After

 $P_0$  $P_0 \circ P_1$  $P_0 \circ P_1 \circ P_2$  $P_0 \circ \dots \circ P_3$  $P_0 \circ \dots \circ P_4$ 

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

```
MPI_Scan(sendbuf, recvbuf, count, datatype, op, comm, ierror)
<type>   :: sendbuf(:), recvbuf(:)
integer  :: count, datatype, op, comm, ierror
```

# ExScan

 $P_0$  $P_1$  $P_2$  $P_3$  $P_4$ 

Before



+

After



# ExScan

$P_0$

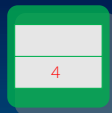
$P_1$

$P_2$

$P_3$

$P_4$

Before

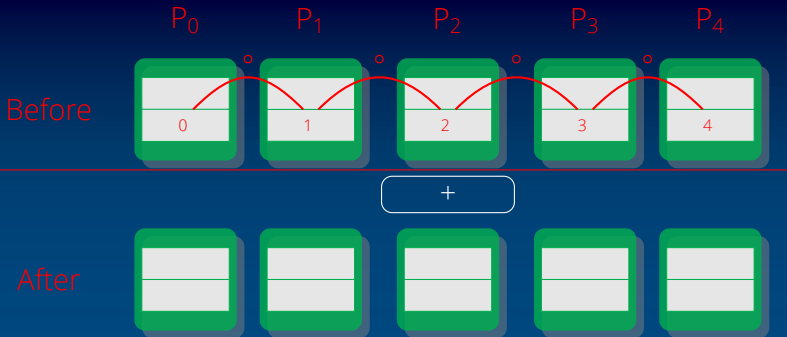


+

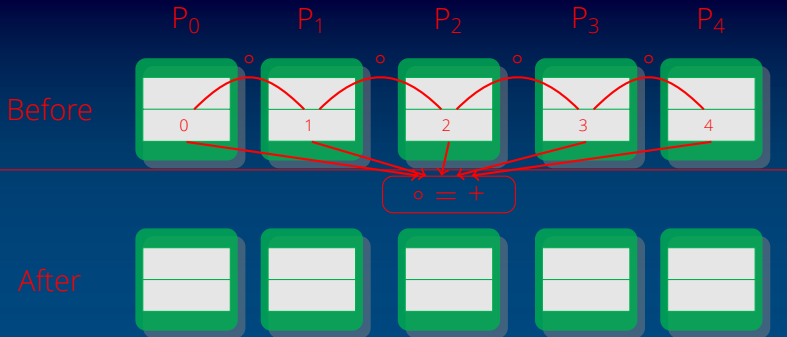
After



# ExScan

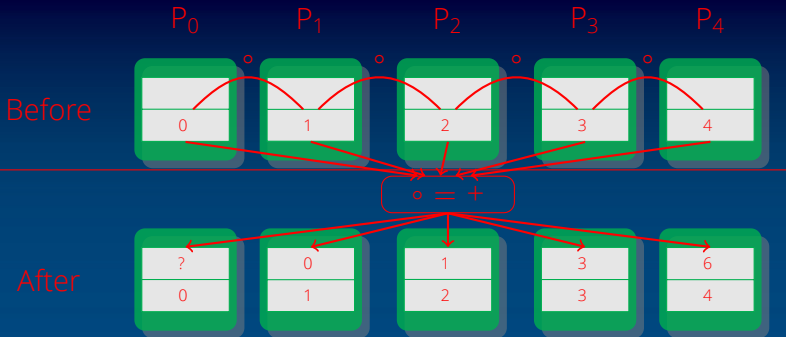


# ExScan

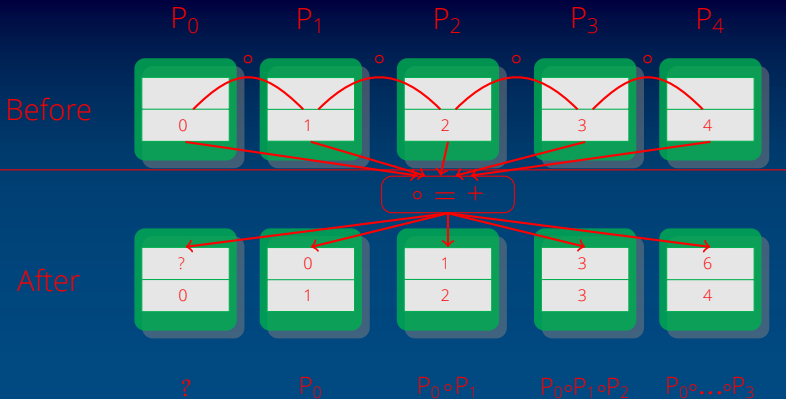




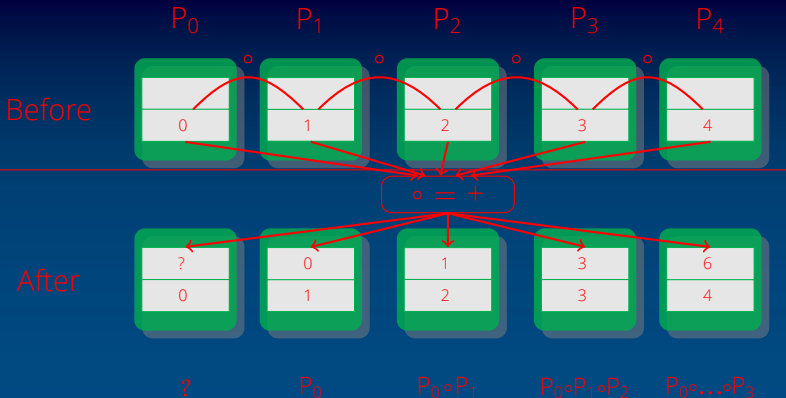
# ExScan



## ExScan

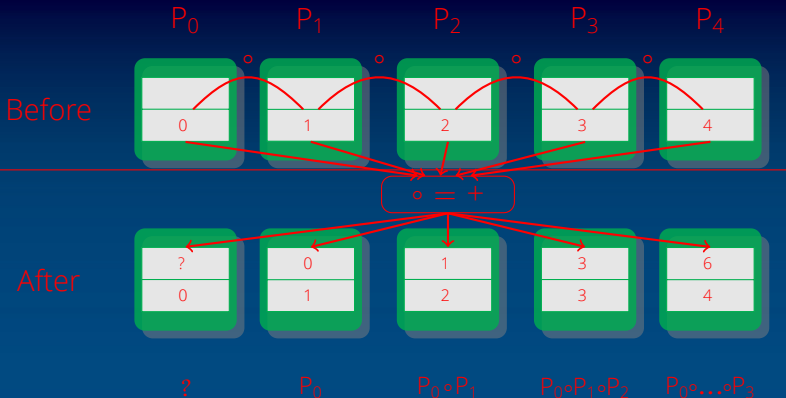


## ExScan



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

## ExScan



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

```
MPI_Exscan(sendbuf, recvbuf, count, datatype, op, comm, ierror)
<type>    :: sendbuf(:), recvbuf(:)
integer   :: count, datatype, op, comm, ierror
```



# MPI reduction operations



[www.ichec.ie](http://www.ichec.ie)



An Roinn Post, Fiontar agus Nuálaitheolaíochta  
Department of Jobs, Enterprise and Innovation



AN ROINN  
OIDEACHAIS AGUS SCILEANNA  
DEPARTMENT OF  
EDUCATION AND SKILLS

HEA  
Higher Education Authority  
an tArd-Chomhairle



# MPI reduction operations

● associative



## MPI reduction operations

- associative
- operates on single variables or vector elements



## MPI reduction operations

- associative
- operates on single variables or vector elements
- pre-defined or user-defined



## MPI reduction operations

- associative
- operates on single variables or vector elements
- pre-defined or user-defined
- results may differ when reduction is applied to floats.

## MPI reduction operations

- associative
- operates on single variables or vector elements
- pre-defined or user-defined
- results may differ when reduction is applied to floats.

## MPI reduction operations

- associative
- operates on single variables or vector elements
- pre-defined or user-defined
- results may differ when reduction is applied to floats.

Handle Name	Operation
MPI_MAX	maximum
MPI_MIN	minimum
MPI_SUM	sum
MPI_PROD	product
MPI_LAND	logical and
MPI_BAND	bit-wise and
MPI_LOR	logical or
MPI_BOR	bit-wise or
MPI_LXOR	logical xor
MPI_BXOR	bit-wise xor
MPI_MAXLOC	max value and location
MPI_MINLOC	min value and location



# user-defined MPI operation





associative

# user-defined MPI operation



An Roinn Post, Fiontar agus Nuálaitheolaíochta  
Department of Jobs, Enterprise and Innovation



AN ROINN  
OIDEACHAIS AGUS SCILEANNA  
DEPARTMENT OF  
EDUCATION AND SKILLS

HEA  
Higher Education Authority  
an tArd-Chomhairle

[www.ichec.ie](http://www.ichec.ie)



## user-defined MPI operation

- associative
- acts on vectors





## user-defined MPI operation

- associative
- acts on vectors
- prototype of the function specified in standard





## user-defined MPI operation

- associative
- acts on vectors
- prototype of the function specified in standard

```
|| void MPI_User_function(void *invec, void *inoutvec, int *len,  
|| MPI_Datatype *datatype);
```

```
|| subroutine user_function(invec, inoutvec, length, MPItype)  
|| <type> :: invec(length), inoutvec(length)  
|| integer :: length, MPItype
```







## user-defined MPI operation

- associative
- acts on vectors
- prototype of the function specified in standard

```
|| void MPI_User_function(void *invec, void *inoutvec, int *len,
    MPI_Datatype *datatype);
```

```
|| subroutine user_function(invec, inoutvec, length, MPItype)
    <type> :: invec(length), inoutvec(length)
    integer :: length, MPItype
```

```
|| int MPI_Op_create(MPI_User_function *function, int commute, MPI_Op *op)
    int MPI_op_free(MPI_Op *op)
```

```
|| MPI_Op_create(function, commute, op, ierror)
    external function
    logical :: commute
    integer :: op, ierror
    MPI_Op_free(op, ierror)
    integer :: op, ierror
```



# Generalized collectives



[www.ichec.ie](http://www.ichec.ie)



An Roinn Post, Fiontar agus Nuálachas  
Department of Jobs, Enterprise and Innovation



AN ROINN  
OIDEACHAIS AGUS SCILEANNA  
DEPARTMENT OF  
EDUCATION AND SKILLS

HEA

Higher Education Authority  
an tArd-Chomhairle



# Generalized collectives

MPI\_Allgatherv  
MPI\_Alltoallv  
MPI\_Alltoallw  
MPI\_Gatherv  
MPI\_Scatterv  
MPI\_Reduce\_scatter



## Generalized collectives

```
MPI_Allgatherv  
MPI_Alltoallv  
MPI_Alltoallw  
MPI_Gatherv  
MPI_Scatterv  
MPI_Reduce_scatter
```

```
int MPI_Alltoallv(void *sendbuf, int *sendcounts, int *sdispls,  
MPI_Datatype sendtype, void *recvbuf, int *recvcounts, int  
*rdispls, MPI_Datatype recvtype, MPI_Comm comm)
```

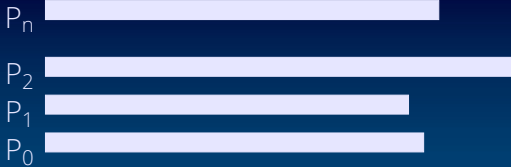
# Generalized collectives

```
MPI_Allgatherv  
MPI_Alltoallv  
MPI_Alltoallw  
MPI_Gatherv  
MPI_Scatterv  
MPI_Reduce_scatter
```

```
int MPI_Alltoallv(void *sendbuf, int *sendcounts, int *sdispls,  
MPI_Datatype sendtype, void *recvbuf, int *recvcounts, int  
*rdispls, MPI_Datatype recvtype, MPI_Comm comm)
```

```
MPI_Alltoallv(sendbuf, sendcounts, sdispls, sendtype, recvbuf,  
recvcounts, rdispls, recvtype, comm, ierror)  
<type> :: sendbuf(:), recvbuf(:)  
integer :: sendcounts(:), sdispls(:), sendtype  
integer :: recvcounts(:), rdispls(:), recvtype  
integer :: comm, ierror
```

# barrier



barrier



barrier





barrier



barrier



# barrier



## barrier



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

# barrier



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

```
|| MPI_Barrier(comm, ierror)
|| integer :: comm, ierror
```

## barrier



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

```
|| MPI_Barrier(comm, ierror)
|| integer :: comm, ierror
```

- never use in production useful only in debugging

## barrier



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

```
|| MPI_Barrier(comm, ierror)
|| integer :: comm, ierror
```

- **never use in production** useful only in debugging
- synchronisation, in general, does not need to be done by hand, data communication does it for you