

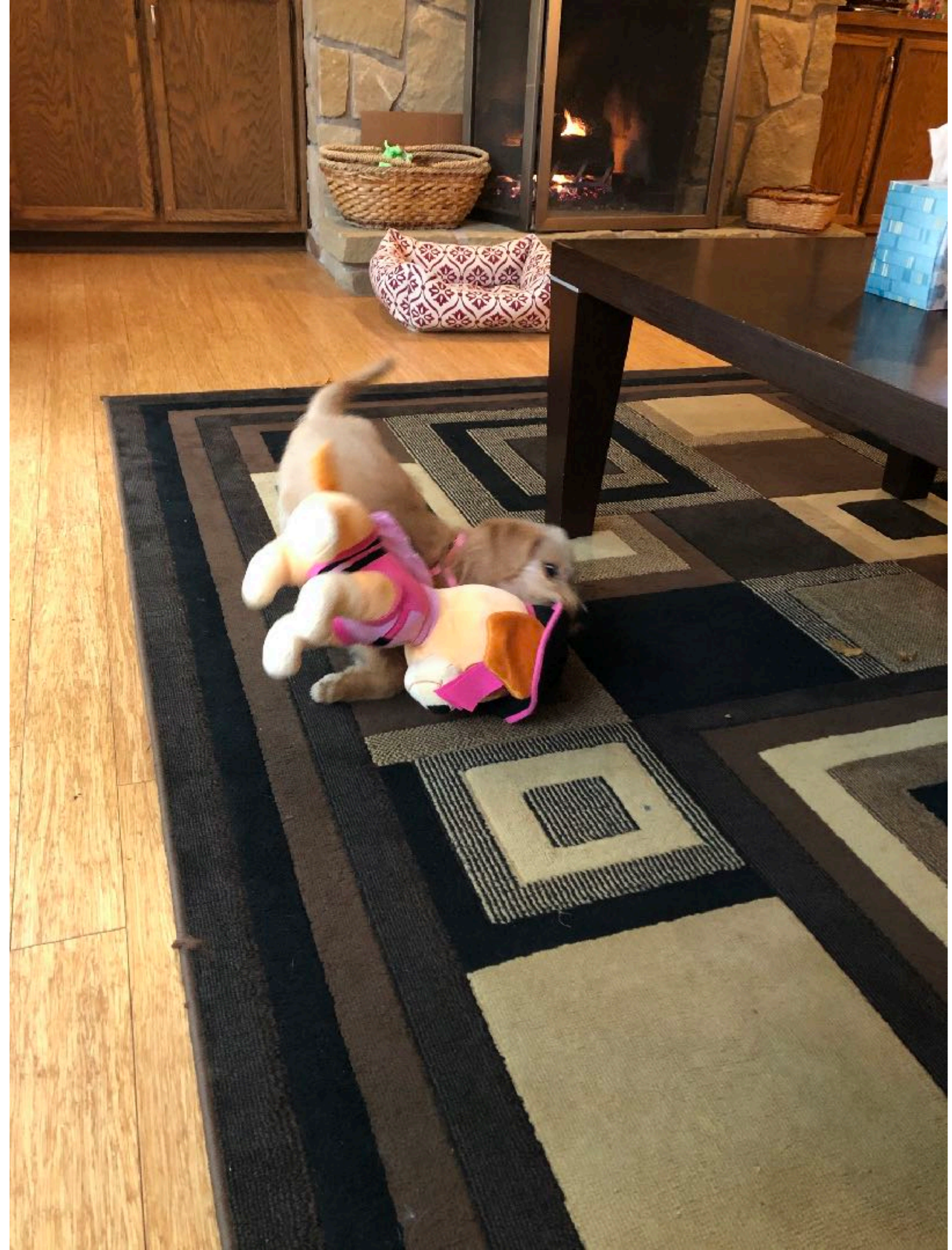


# Lecture 9: MPI Collectives

**CMSE 822: Parallel Computing**  
**Prof. Sean M. Couch**



# Puppy time

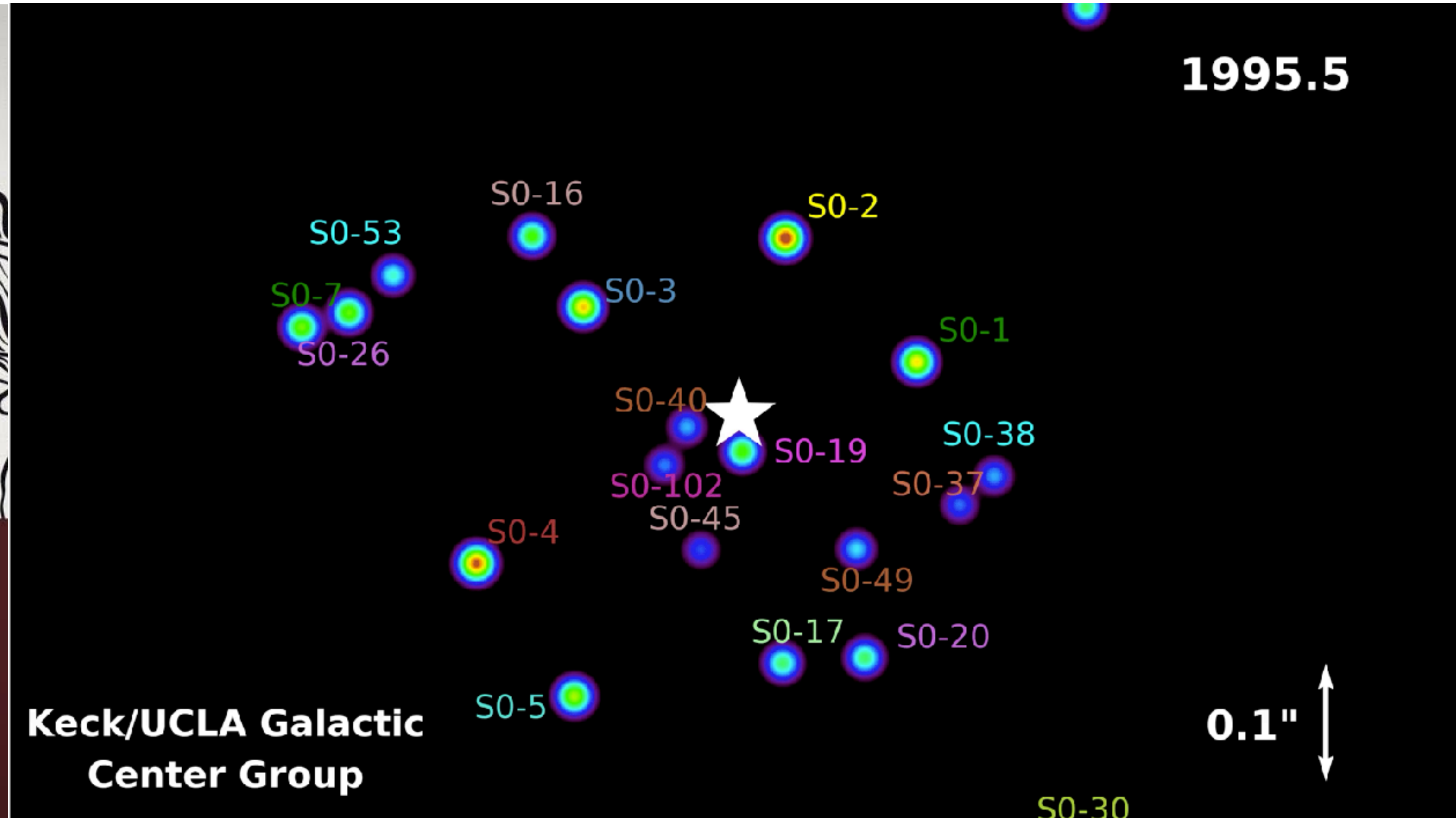
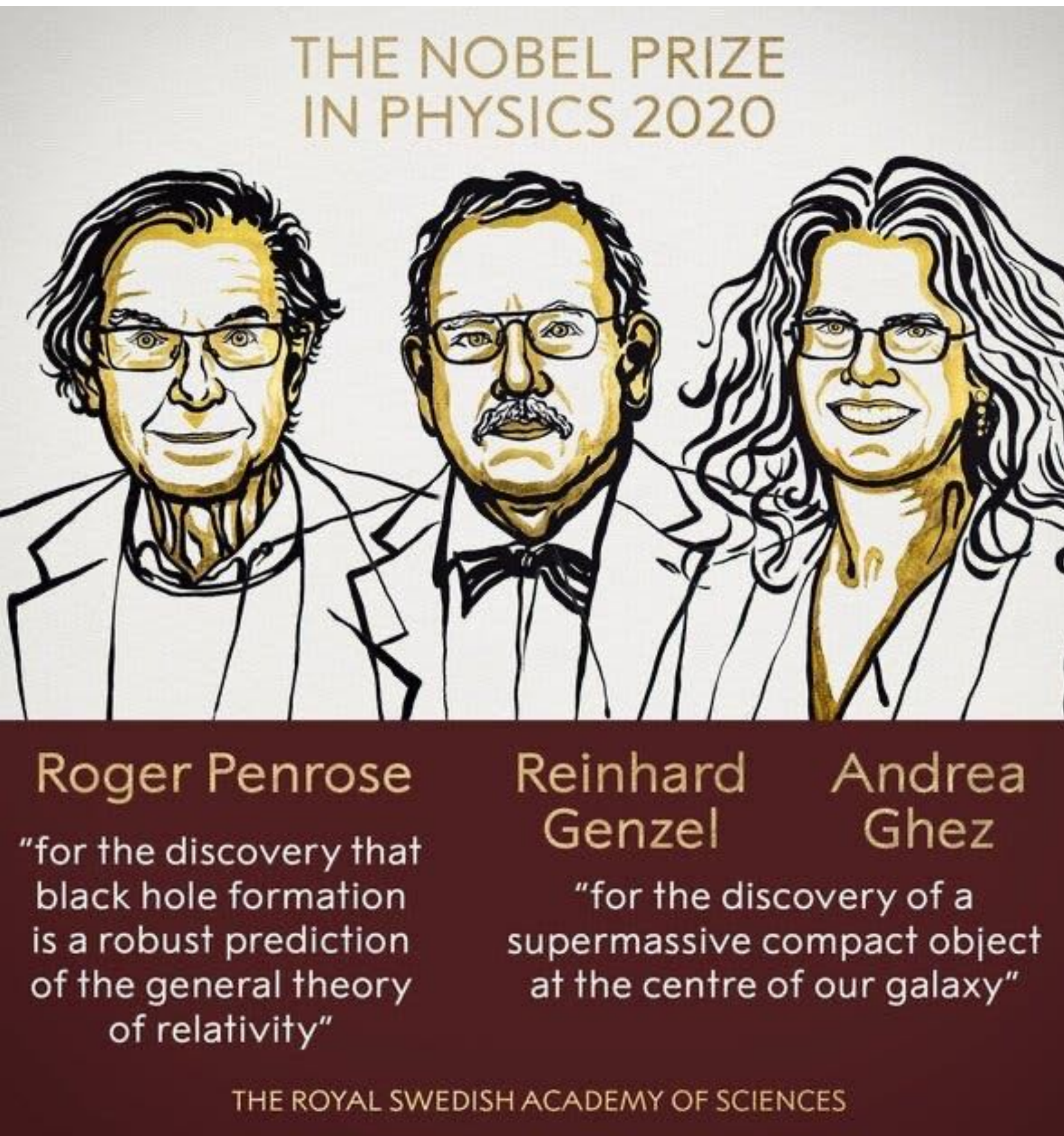






# Nobel Prize in Physics

## Black Holes!





# Brief MPI Tutorial

See <https://computing.llnl.gov/tutorials/mpi/>

also: <http://www.mpi-forum.org/docs/>

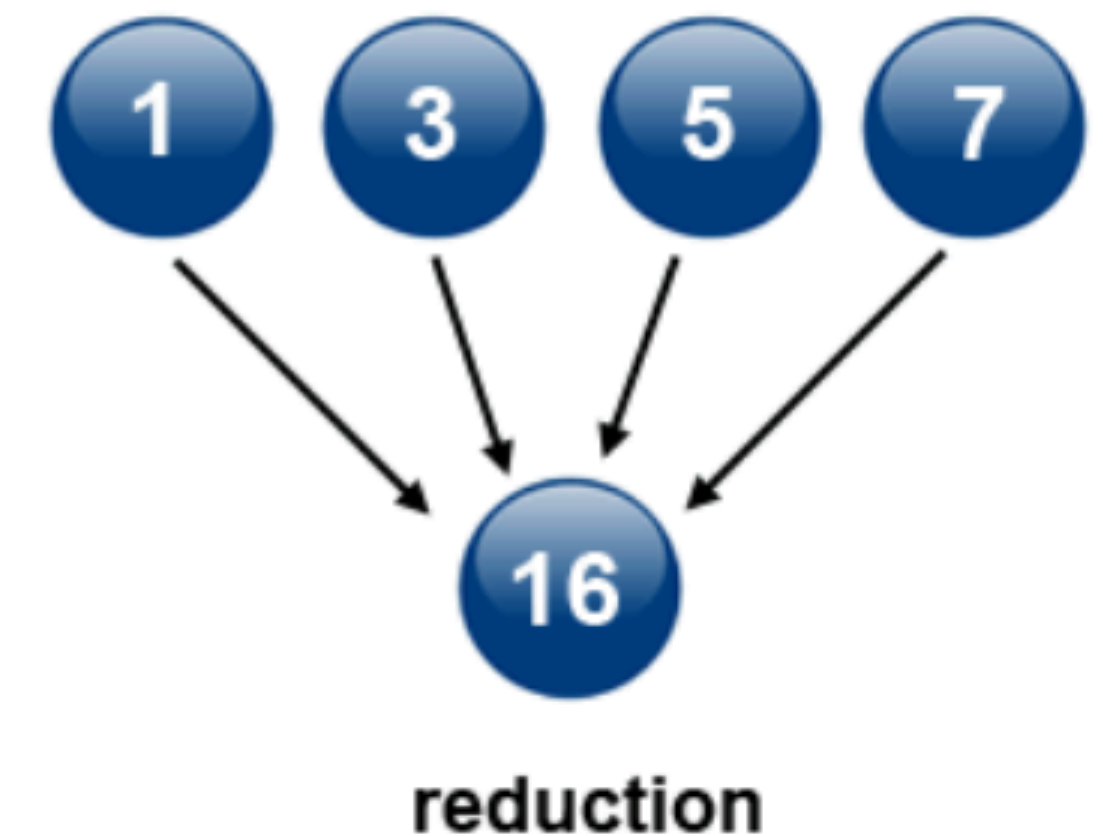
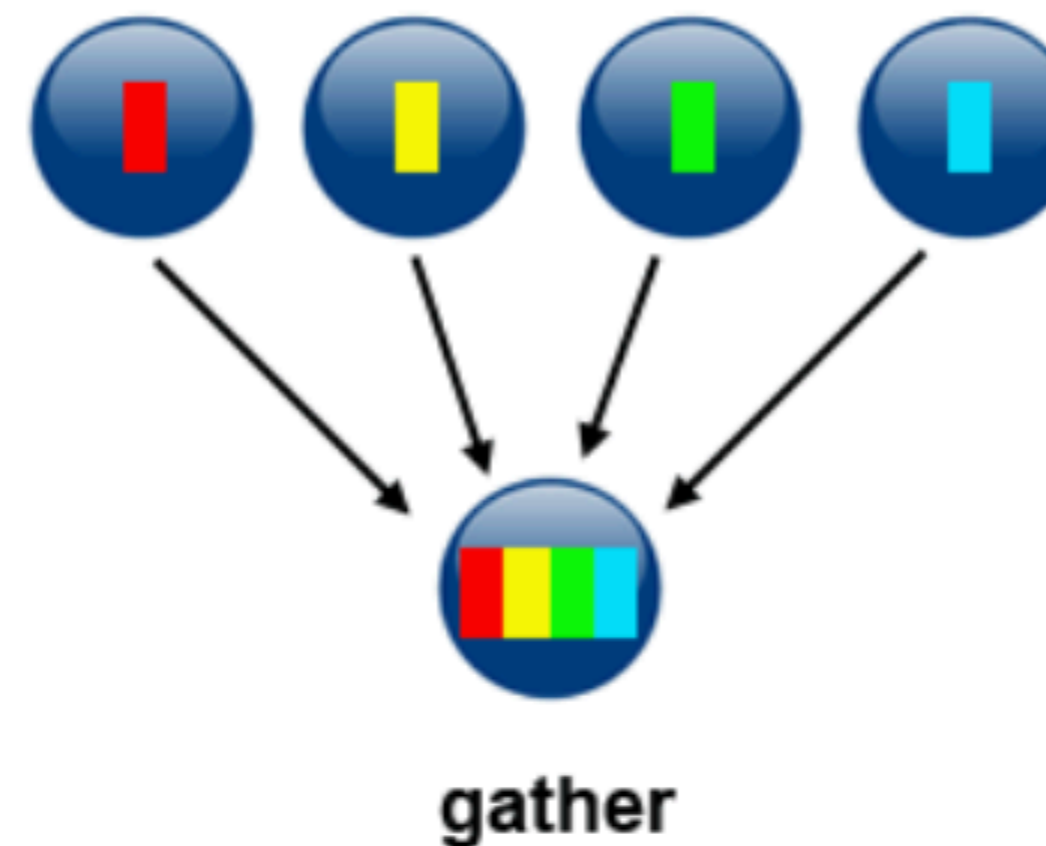
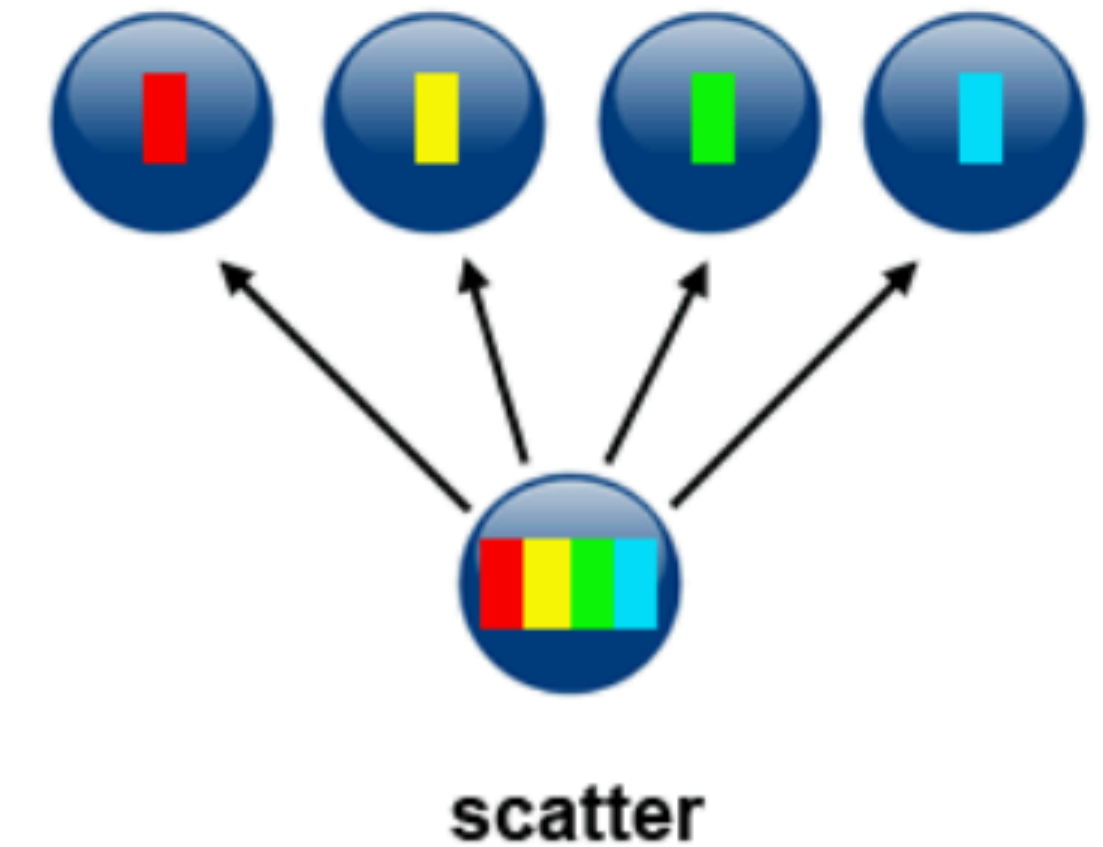
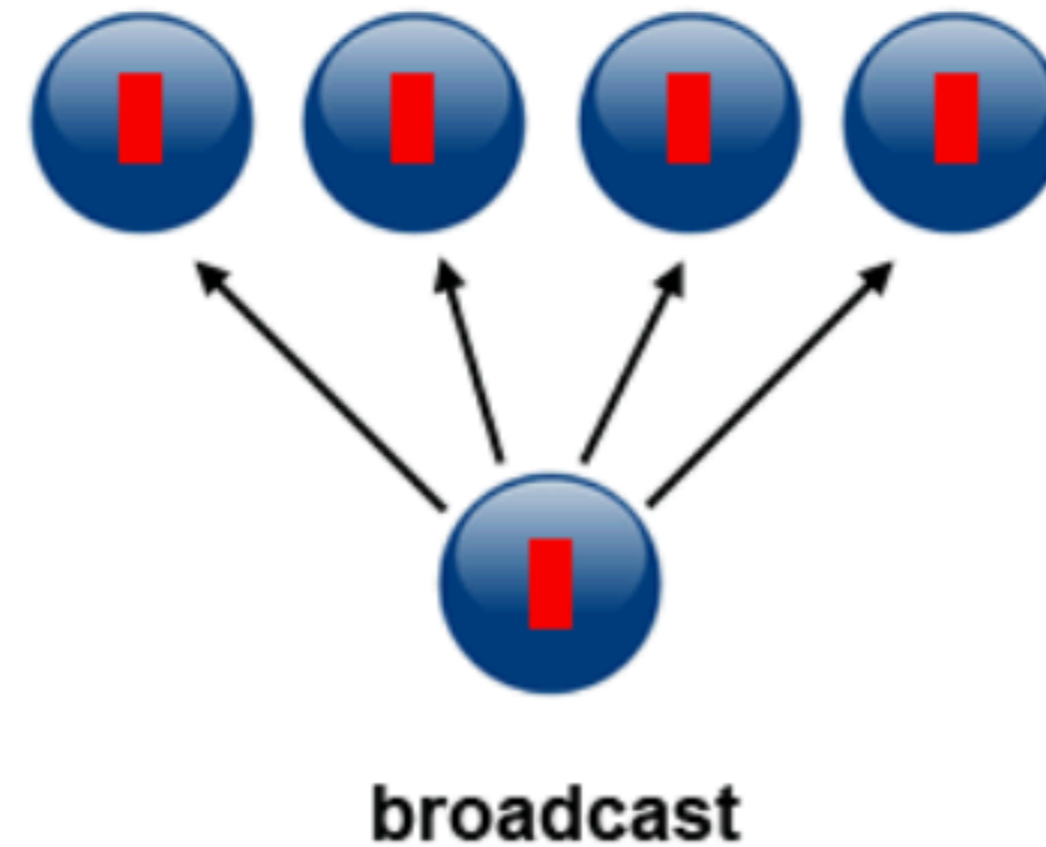




# MPI Collectives

## ► Types of Collective Operations:

- **Synchronization** - processes wait until all members of the group have reached the synchronization point.
- **Data Movement** - broadcast, scatter/gather, all to all.
- **Collective Computation** (reductions) - one member of the group collects data from the other members and performs an operation (min, max, add, multiply, etc.) on that data.





# MPI Collectives

## ► Scope:

- Collective communication routines must involve **all** processes within the scope of a communicator.
  - All processes are by default, members in the communicator MPI\_COMM\_WORLD.
  - Additional communicators can be defined by the programmer. See the [Group and Communicator Management Routines](#) section for details.
- Unexpected behavior, including program failure, can occur if even one task in the communicator doesn't participate.
- It is the programmer's responsibility to ensure that all processes within a communicator participate in any collective operations.

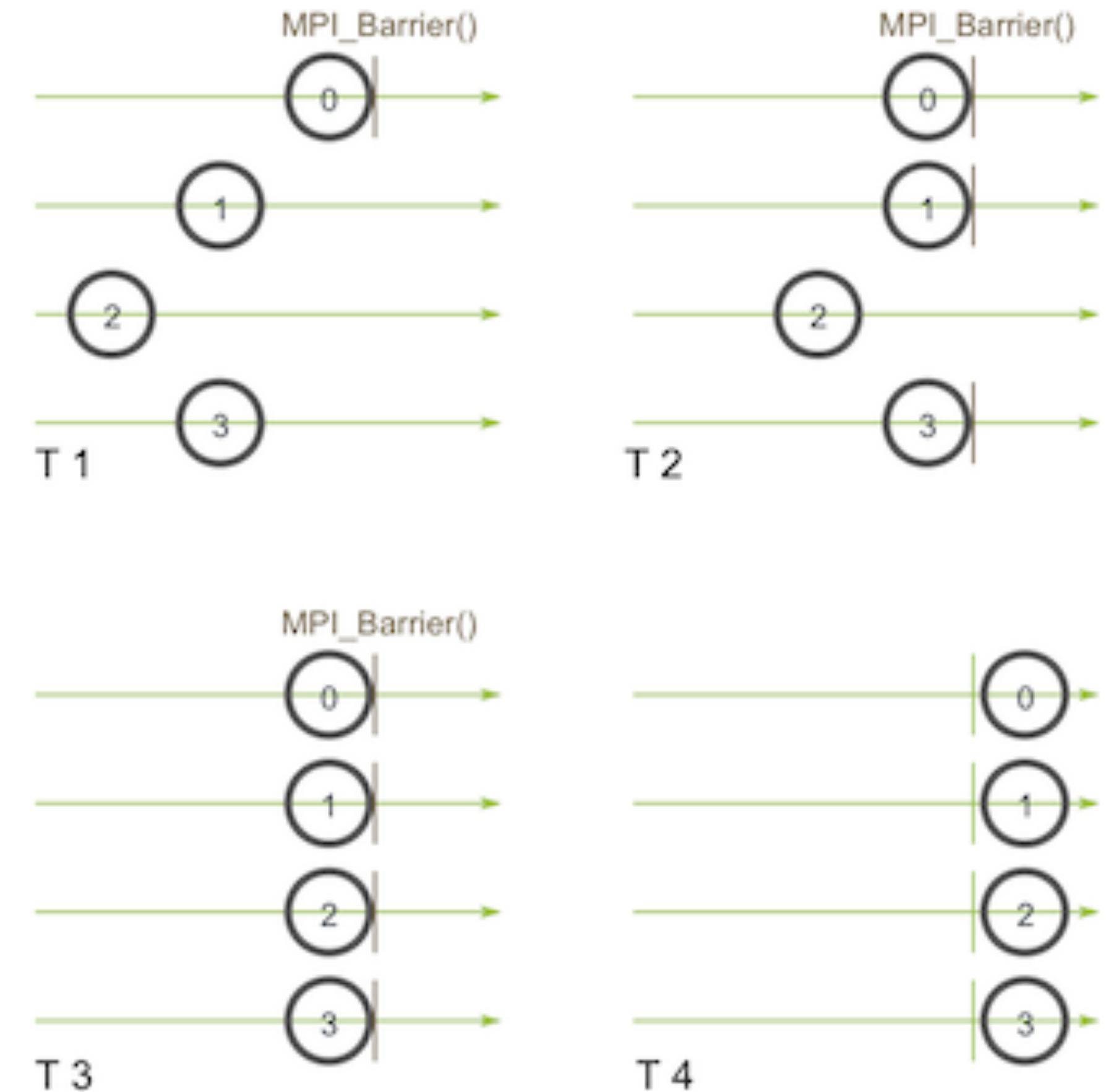


# MPI Collectives

## MPI\_Barrier

Synchronization operation. Creates a barrier synchronization in a group. Each task, when reaching the MPI\_Barrier call, blocks until all tasks in the group reach the same MPI\_Barrier call. Then all tasks are free to proceed.

```
MPI_Barrier (comm)  
MPI_BARRIER (comm,ierr)
```





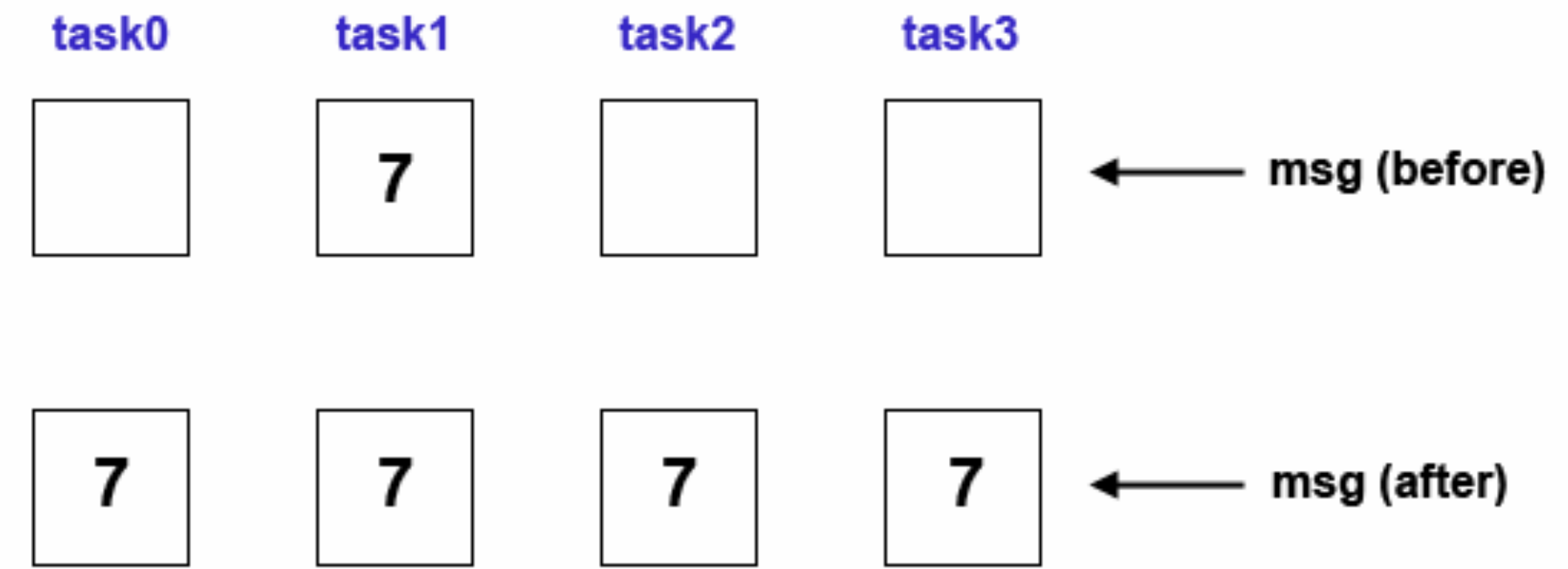


## MPI\_Bcast

Broadcasts a message from one task to all other tasks in communicator

```
count = 1;
source = 1;
MPI_Bcast(&msg, count, MPI_INT, source, MPI_COMM_WORLD);
```

task1 contains the message to be broadcast



### MPI\_Bcast

Data movement operation. Broadcasts (sends) a message from the process with rank "root" to all other processes in the group.

Diagram Here

```
MPI_Bcast (&buffer, count, datatype, root, comm)
MPI_BCAST (buffer, count, datatype, root, comm, ierr)
```





# MPI Collectives

## MPI\_Scatter

Data movement operation. Distributes distinct messages from a single source task to each task in the group.

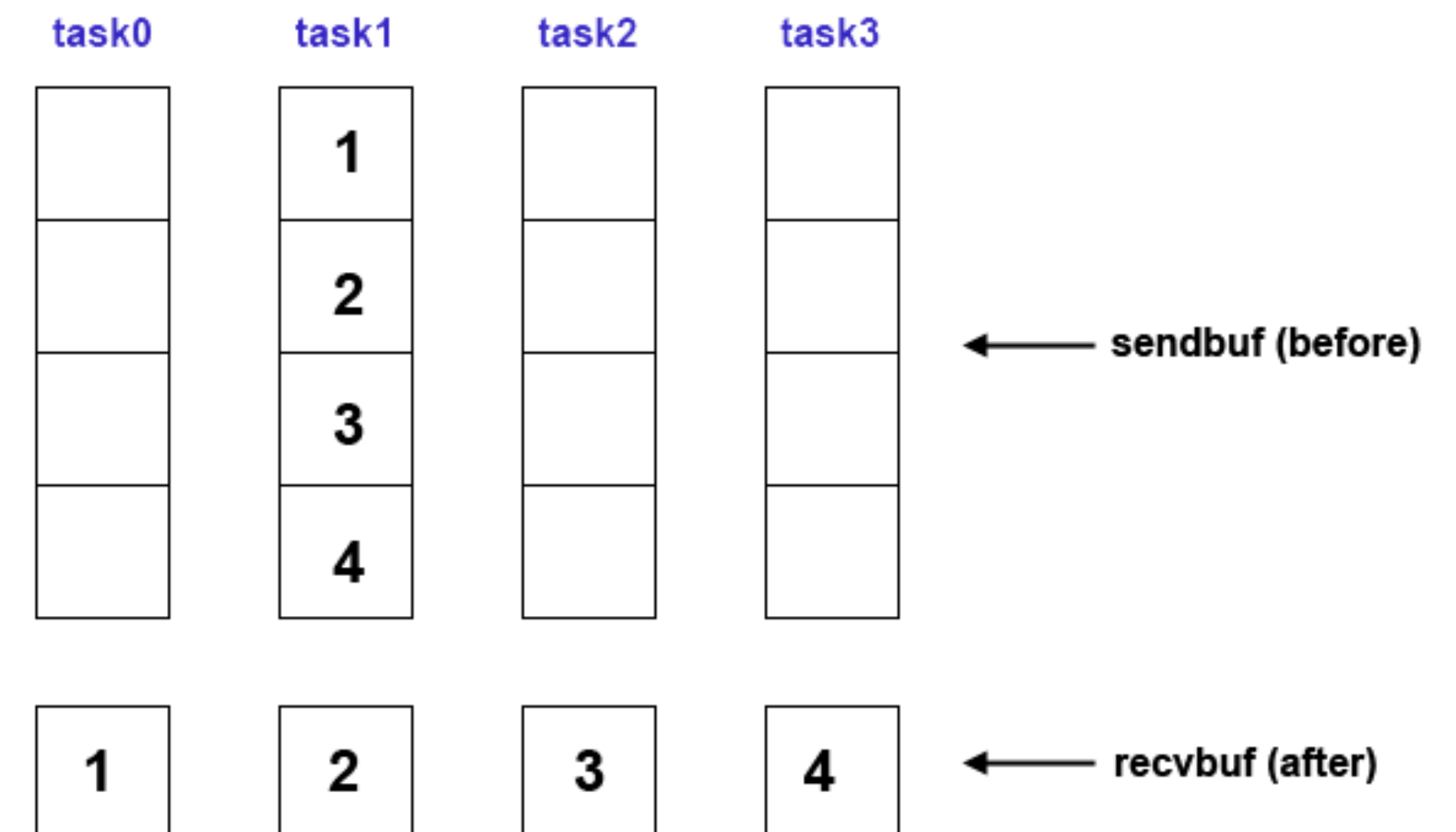
Diagram Here

```
MPI_Scatter (&sendbuf, sendcnt, sendtype, &recvbuf,  
             recvcnt, recvtype, root, comm)  
MPI_SCATTER (sendbuf, sendcnt, sendtype, recvbuf,  
             recvcnt, recvtype, root, comm, ierr)
```

## MPI\_Scatter

Sends data from one task to all other tasks in communicator

```
sendcnt = 1;  
recvcnt = 1;  
src = 1;  
task1 contains the data to be scattered  
MPI_Scatter(sendbuf, sendcnt, MPI_INT  
            recvbuf, recvcnt, MPI_INT  
            src, MPI_COMM_WORLD);
```







# MPI Collectives

## MPI\_Gather

### MPI\_Gather

Data movement operation. Gathers distinct messages from each task in the group to a single destination task. This routine is the reverse operation of MPI\_Scatter.

Diagram Here

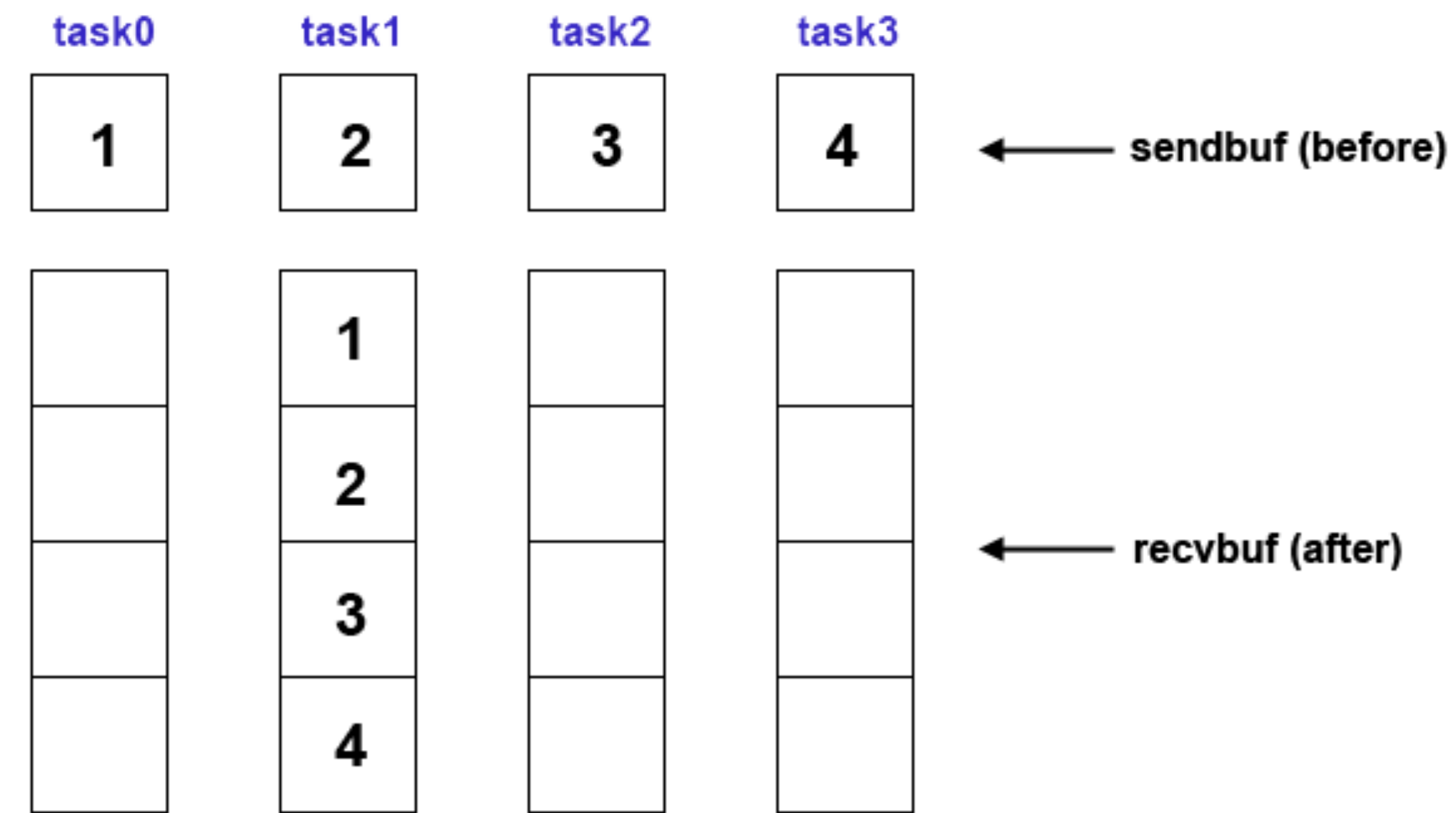
```
MPI_Gather (&sendbuf, sendcnt, sendtype, &recvbuf,  
           recvcount, recvtype, root, comm)  
MPI_GATHER (sendbuf, sendcnt, sendtype, recvbuf,  
            recvcount, recvtype, root, comm, ierr)
```

Gathers data from all tasks in communicator to a single task

```
sendcnt = 1;  
recvcnt = 1;  
src = 1;
```

message will be gathered into task1

```
MPI_Gather(sendbuf, sendcnt, MPI_INT  
           recvbuf, recvcnt, MPI_INT  
           src, MPI_COMM_WORLD);
```







# MPI Collectives

## MPI\_Allgather

Data movement operation. Concatenation of data to all tasks in a group. Each task in the group, in effect, performs a one-to-all broadcasting operation within the group.

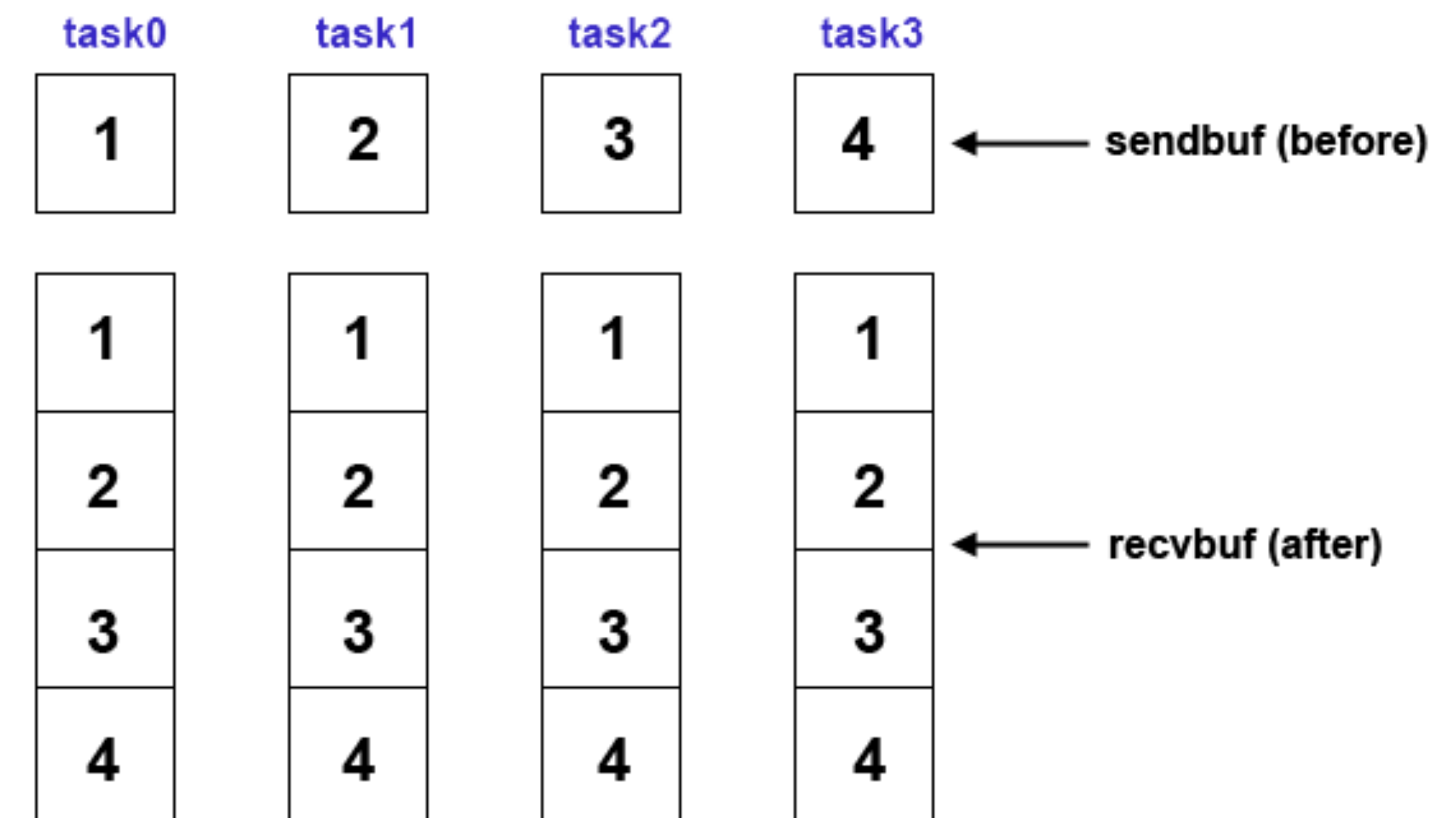
Diagram Here

```
MPI_Allgather (&sendbuf, sendcount, sendtype, &recvbuf,  
               recvcount, recvtype, comm)  
MPI_ALLGATHER (sendbuf, sendcount, sendtype, recvbuf,  
               recvcount, recvtype, comm, info)
```

## MPI\_Allgather

Gathers data from all tasks and then distributes to all tasks in communicator

```
sendcnt = 1;  
recvcnt = 1;  
MPI_Allgather(sendbuf, sendcnt, MPI_INT,  
              recvbuf, recvcnt, MPI_INT,  
              MPI_COMM_WORLD);
```







# MPI Collectives

## MPI\_Reduce

Collective computation operation. Applies a reduction operation on all tasks in the group and places the result in one task.

Diagram Here

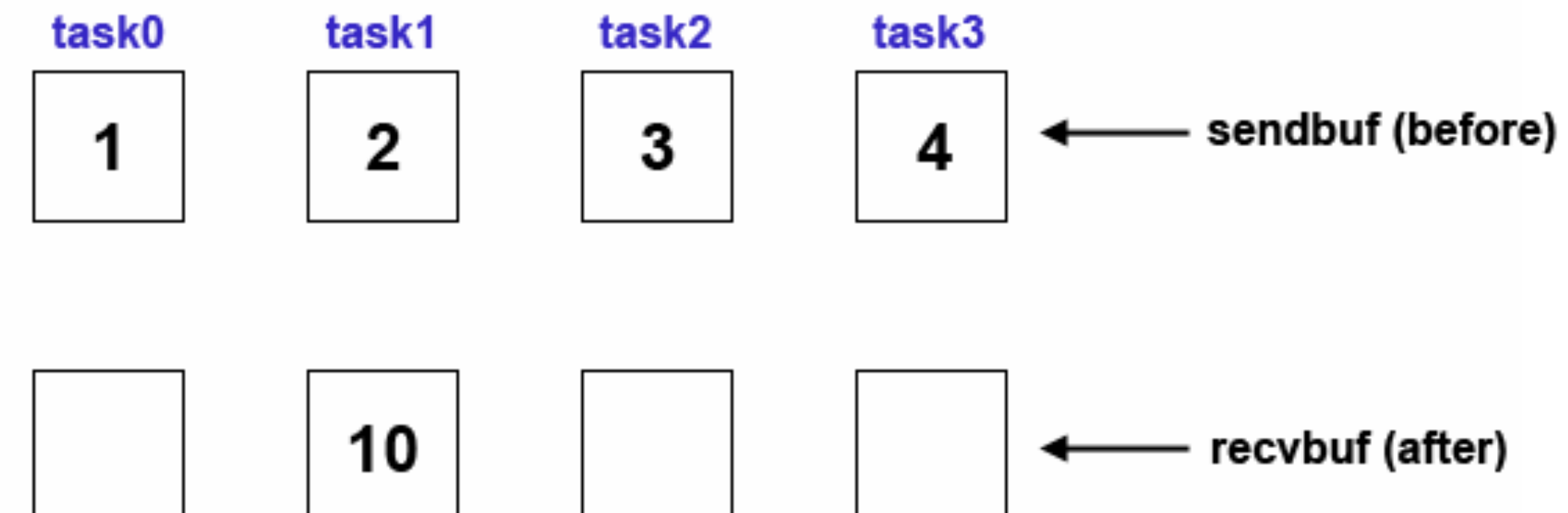
```
MPI_Reduce (&sendbuf, &recvbuf, count, datatype, op, root, comm)
MPI_REDUCE (sendbuf, recvbuf, count, datatype, op, root, comm, ierr)
```

## MPI\_Reduce

Perform reduction across all tasks in communicator and store result in 1 task

```
count = 1;
dest = 1;
MPI_Reduce(sendbuf, recvbuf, count, MPI_INT,
           MPI_SUM, dest, MPI_COMM_WORLD);
```

task1 will contain result







# MPI Collectives

The predefined MPI reduction operations appear below. Users can also define their own reduction functions by using the [MPI\\_Op\\_create](#) routine.

MPI Reduction Operation		C Data Types	Fortran Data Type
<b>MPI_MAX</b>	maximum	integer, float	integer, real, complex
<b>MPI_MIN</b>	minimum	integer, float	integer, real, complex
<b>MPI_SUM</b>	sum	integer, float	integer, real, complex
<b>MPI_PROD</b>	product	integer, float	integer, real, complex
<b>MPI_LAND</b>	logical AND	integer	logical
<b>MPI_BAND</b>	bit-wise AND	integer, MPI_BYTE	integer, MPI_BYTE
<b>MPI_LOR</b>	logical OR	integer	logical
<b>MPI_BOR</b>	bit-wise OR	integer, MPI_BYTE	integer, MPI_BYTE
<b>MPI_LXOR</b>	logical XOR	integer	logical
<b>MPI_BXOR</b>	bit-wise XOR	integer, MPI_BYTE	integer, MPI_BYTE
<b>MPI_MAXLOC</b>	max value and location	float, double and long double	real, complex, double precision
<b>MPI_MINLOC</b>	min value and location	float, double and long double	real, complex, double precision





# MPI Collectives

## MPI\_Allreduce

Collective computation operation + data movement. Applies a reduction operation and places the result in all tasks in the group. This is equivalent to an MPI\_Reduce followed by an MPI\_Bcast.

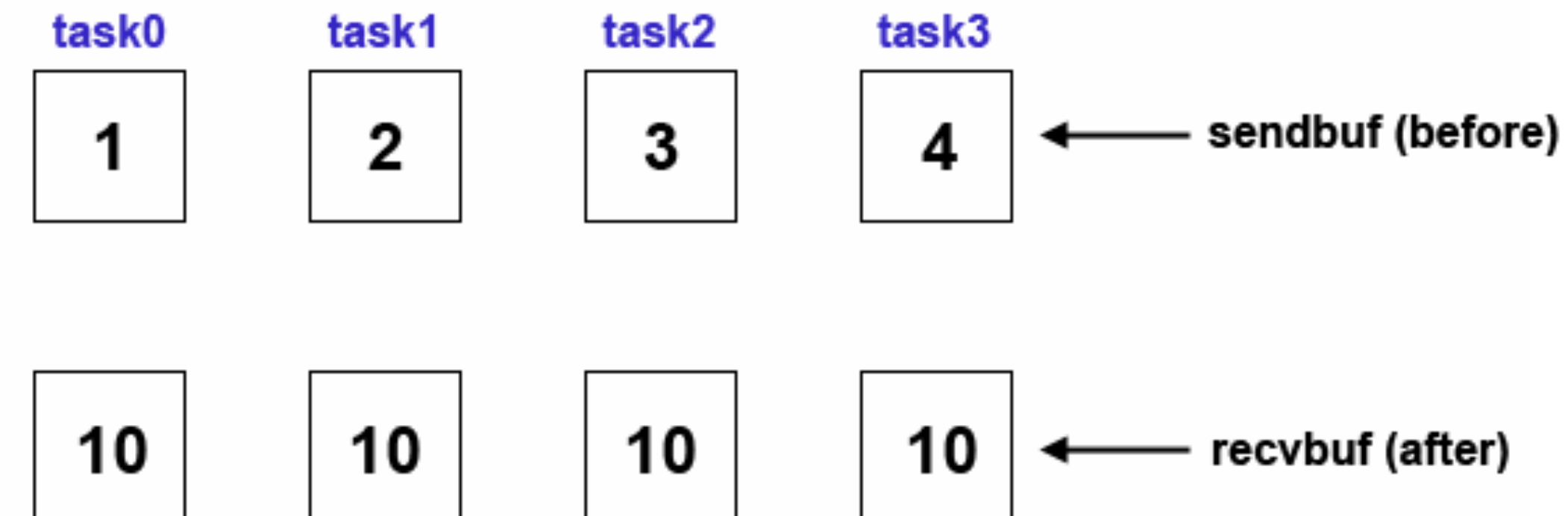
Diagram Here

```
MPI_Allreduce (&sendbuf,&recvbuf,count,datatype,op,comm)
MPI_ALLREDUCE (sendbuf,recvbuf,count,datatype,op,comm,ierr)
```

## MPI\_Allreduce

Perform reduction and store result across all tasks in communicator

```
count = 1;
MPI_Allreduce(sendbuf, recvbuf, count, MPI_INT,
              MPI_SUM, MPI_COMM_WORLD);
```







# MPI Collectives

## MPI\_Reduce\_scatter

Collective computation operation + data movement. First does an element-wise reduction on a vector across all tasks in the group. Next, the result vector is split into disjoint segments and distributed across the tasks. This is equivalent to an MPI\_Reduce followed by an MPI\_Scatter operation.

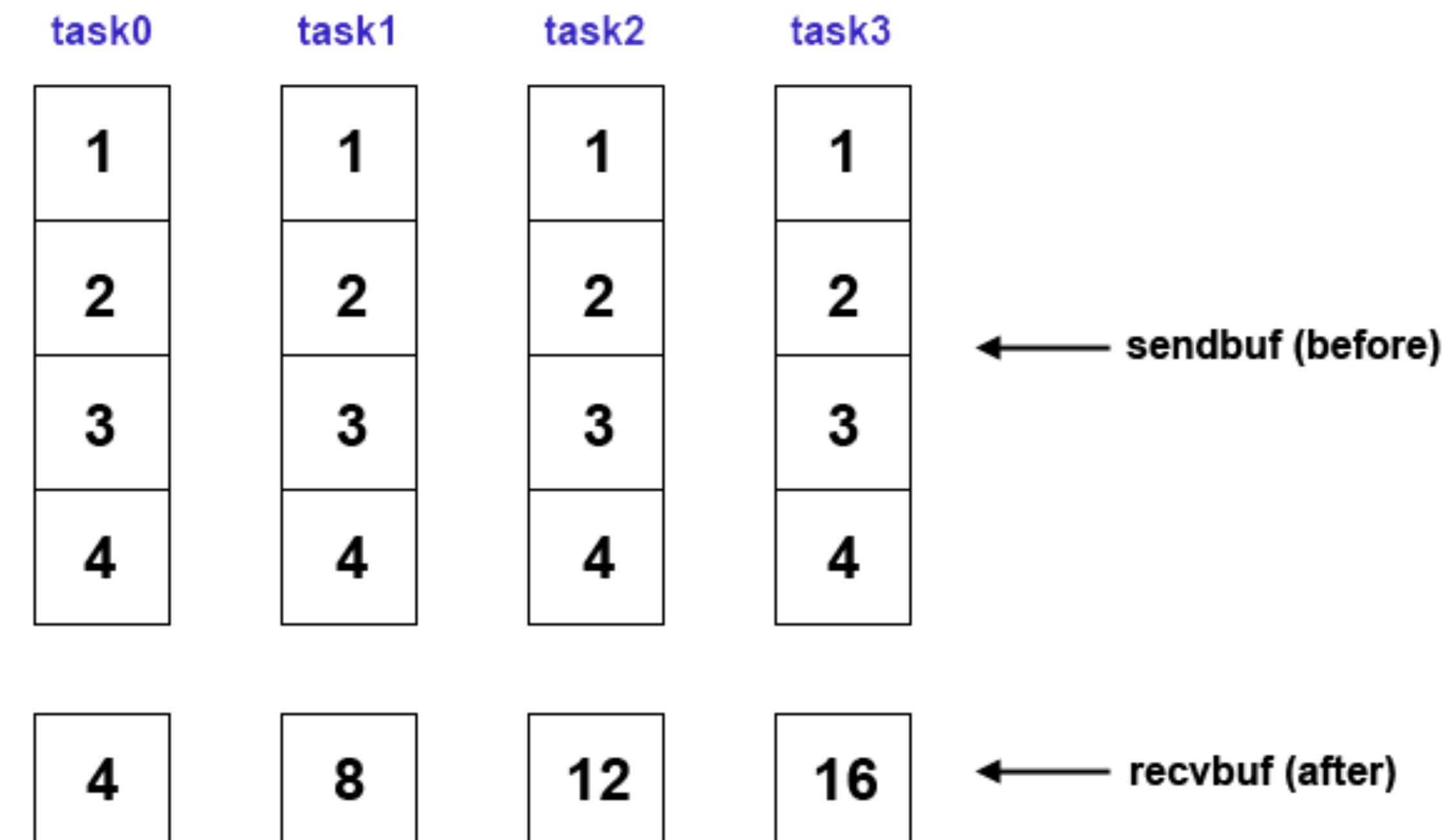
Diagram Here

```
MPI_Reduce_scatter (&sendbuf,&recvbuf,recvcount,datatype,  
                    op,comm)  
MPI_REDUCE_SCATTER (sendbuf,recvbuf,recvcount,datatype,  
                    op,comm,ierr)
```

## MPI\_Reduce\_scatter

Perform reduction on vector elements and distribute segments of result vector across all tasks in communicator

```
recvcnt = 1;  
MPI_Reduce_scatter(sendbuf, recvbuf, recvcount,  
                  MPI_INT, MPI_SUM, MPI_COMM_WORLD);
```





## MPI Collectives

Scatter data from all tasks to all tasks in communicator

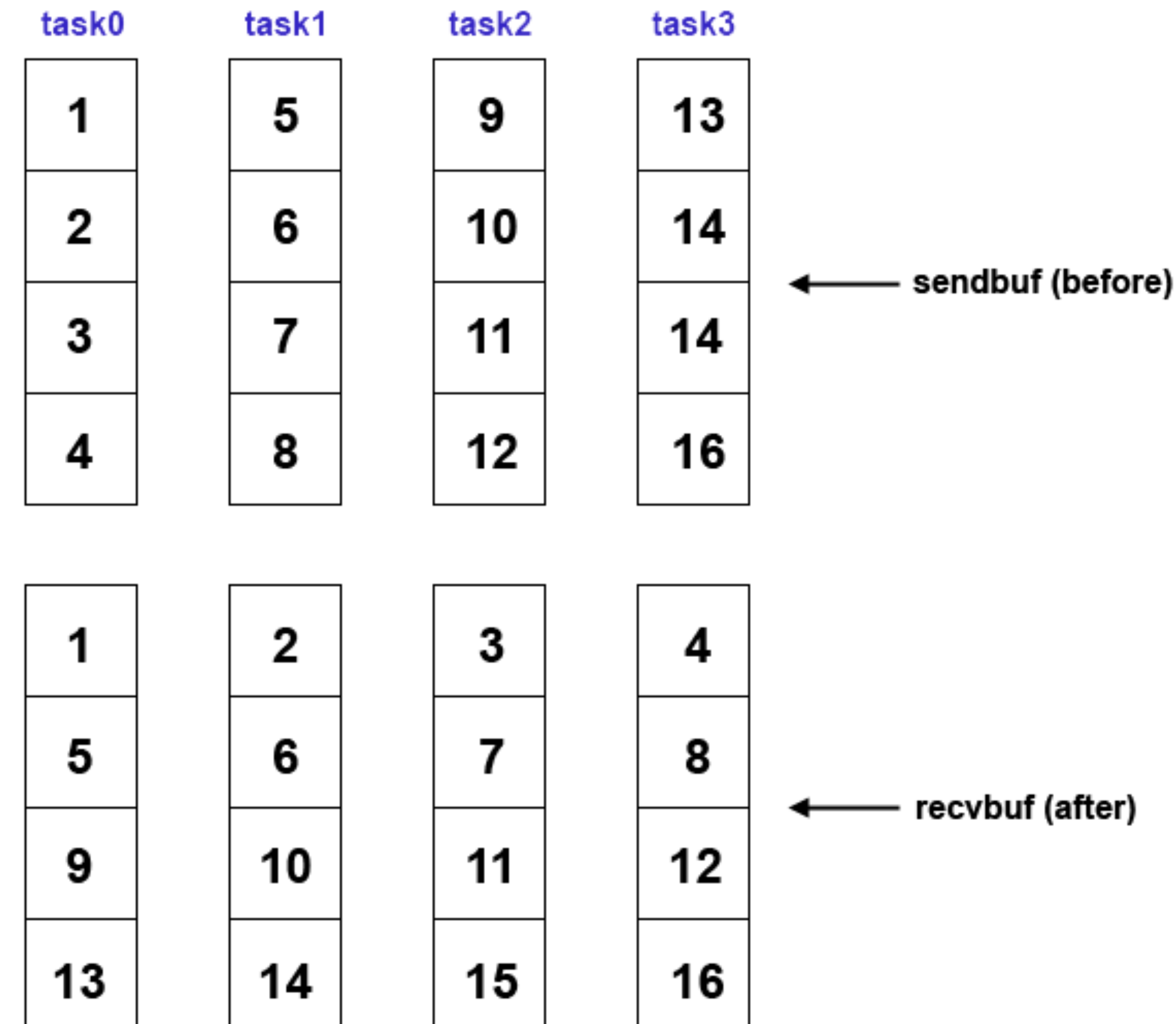
### MPI\_Alltoall

Data movement operation. Each task in a group performs a scatter operation, sending a distinct message to all the tasks in the group in order by index.

Diagram Here

```
MPI_Alltoall (&sendbuf, sendcount, sendtype, &recvbuf,
               recvcnt, recvtype, comm)
MPI_ALLTOALL (sendbuf, sendcount, sendtype, recvbuf,
               recvcnt, recvtype, comm, ierr)
```

```
sendcnt = 1;
recvcnt = 1;
MPI_Alltoall(sendbuf, sendcnt, MPI_INT,
              recvbuf, recvcnt, MPI_INT,
              MPI_COMM_WORLD);
```







# MPI Collectives

## MPI\_Scan

Performs a scan operation with respect to a reduction operation across a task group.

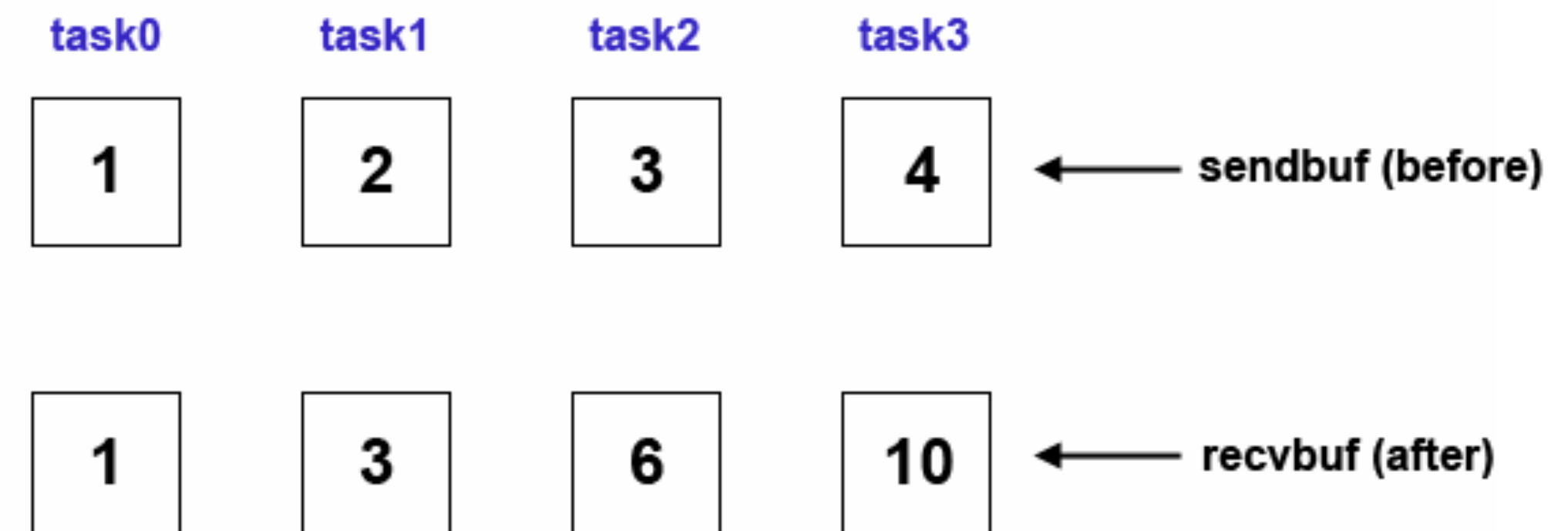
Diagram Here

```
MPI_Scan (&sendbuf,&recvbuf,count,datatype,op,comm)  
MPI_SCAN (sendbuf,recvbuf,count,datatype,op,comm,ierr)
```

## MPI\_Scan

Computes the scan (partial reductions) across all tasks in communicator

```
count = 1;  
MPI_Scan(sendbuf, recvbuf, count, MPI_INT,  
        MPI_SUM, MPI_COMM_WORLD);
```







# Example



## C Language - Collective Communications Example

```
1  #include "mpi.h"
2  #include <stdio.h>
3  #define SIZE 4
4
5  main(int argc, char *argv[]) {
6      int numtasks, rank, sendcount, recvcount, source;
7      float sendbuf[SIZE][SIZE] = {
8          {1.0, 2.0, 3.0, 4.0},
9          {5.0, 6.0, 7.0, 8.0},
10         {9.0, 10.0, 11.0, 12.0},
11         {13.0, 14.0, 15.0, 16.0} };
12     float recvbuf[SIZE];
13
14     MPI_Init(&argc,&argv);
15     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
16     MPI_Comm_size(MPI_COMM_WORLD, &numtasks);
17
18     if (numtasks == SIZE) {
19         // define source task and elements to send/receive, then perform collective scatter
20         source = 1;
21         sendcount = SIZE;
22         recvcount = SIZE;
23         MPI_Scatter(sendbuf,sendcount,MPI_FLOAT,recvbuf,recvcount,
24                     MPI_FLOAT,source,MPI_COMM_WORLD);
25
26         printf("rank= %d  Results: %f %f %f %f\n",rank,recvbuf[0],
27                recvbuf[1],recvbuf[2],recvbuf[3]);
28     }
29     else
30         printf("Must specify %d processors. Terminating.\n",SIZE);
31
32     MPI_Finalize();
33 }
```





# Group work

- What is the output of the example collective program?





# Homework 5

## Pi by MPI

