


MPI



# Collective Operations

Collective/Global Communication Operations

# Collective/Global Communication Operations

- A communication operation is called collective or global if all or a subset of the processes of a parallel program are involved.

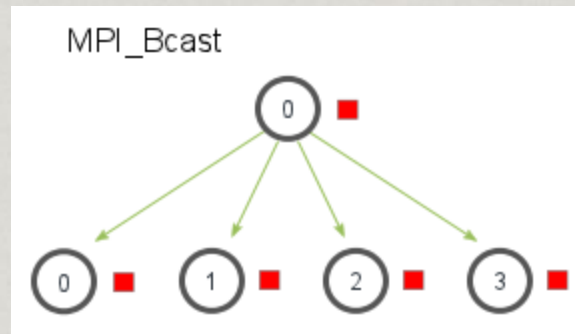
Global communication operation	MPI function
Broadcast operation	<code>MPI_Bcast()</code>
Accumulation operation	<code>MPI_Reduce()</code>
Gather operation	<code>MPI_Gather()</code>
Scatter operation	<code>MPI_Scatter()</code>
Multi-broadcast	<code>MPI_Allgather()</code>
Multi-accumulation operation	<code>MPI_Allreduce()</code>

- Collective MPI communication operations are always blocking; no non-blocking versions are provided as is the case for point-to-point operations

# Broadcast operation

- For a broadcast operation, one specific process of a group of processes sends the same data block to all other processes of the group.

- ```
int MPI_Bcast (void *message,  
              int count,  
              MPI_Datatype type,  
              int root,  
              MPI_Comm comm)
```



- Where **root** denotes the process which sends the data block.
- Each process must specify the same root process and must use the same communicator.
- Data blocks sent by `MPI_Bcast ()` cannot be received by an `MPI_Recv ()`.

# Broadcast operation cont.

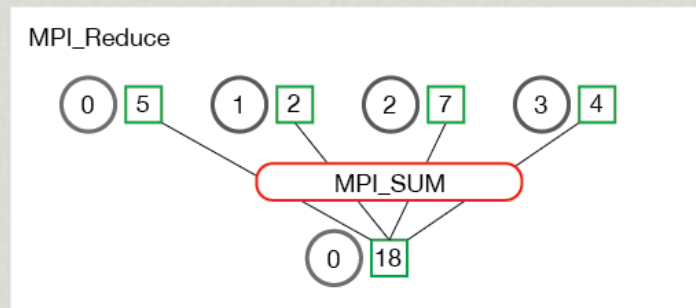
- The MPI runtime system guarantees that broadcast messages are received in the same order in which they have been sent by the root process, even if the corresponding broadcast operations are not executed at the same time.

```
if (my_rank == 0) {
    MPI_Bcast (&x, 1, MPI_INT, 0, comm);
    MPI_Bcast (&y, 1, MPI_INT, 0, comm);
    local_work ();
}
else if (my_rank == 1) {
    local_work ();
    MPI_Bcast (&y, 1, MPI_INT, 0, comm);
    MPI_Bcast (&x, 1, MPI_INT, 0, comm);
}
else if (my_rank == 2) {
    local_work ();
    MPI_Bcast (&x, 1, MPI_INT, 0, comm);
    MPI_Bcast (&y, 1, MPI_INT, 0, comm);
}
```

# Reduction Operation

- For such an operation, each participating process provides a block of data that is combined with the other blocks using a binary reduction operation. The accumulated result is collected at a root process.

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



- The parameter **recvbuf** specifies the receive buffer which is provided by the root process.

# Reduction Operation Cont.

Representation	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 exclusive or
MPI_BXOR	Bit-wise exclusive or
MPI_MAXLOC	Maximum value and corresponding index
MPI_MINLOC	Minimum value and corresponding index



# Reduction Operation Cont.

- For an **MPI\_Reduce()** operation, all participating processes must specify the same values for the parameters **count**, **type**, **op**, and **root**. The send buffers **sendbuf** and the receive buffer **recvbuf** must have the same size.
- MPI supports the definition of user-defined reduction operations using the following MPI function:

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



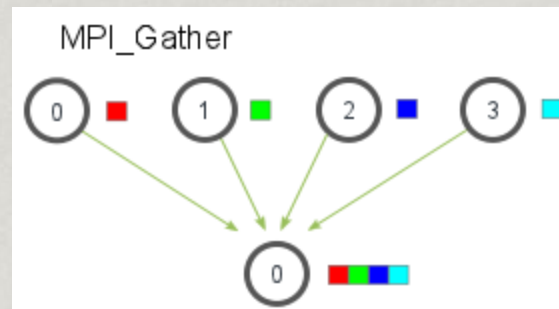
# Reduction Operation Cont.

```
double ain[30], aout[30];
int ind[30];
struct {double val; int rank;} in[30], out[30];
int i, my_rank, root=0;

MPI_Comm_rank (MPI_COMM_WORLD, &my_rank);
for (i=0; i<30; i++) {
    in[i].val = ain[i];
    in[i].rank = my_rank;
}
MPI_Reduce(in,out,30,MPI_DOUBLE_INT,MPI_MAXLOC,root,MPI_COMM_WORLD);
if (my_rank == root)
    for (i=0; i<30; i++) {
        aout[i] = out[i].val;
        ind[i] = out[i].rank;
    }
```

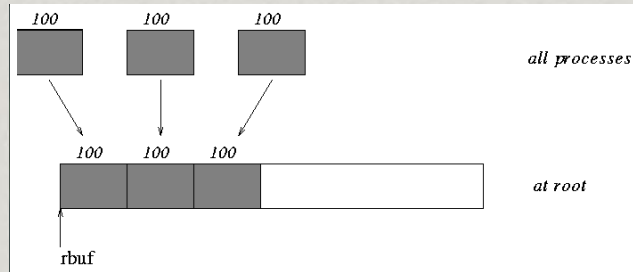
# Gather Operation

- For a gather operation, each process provides a block of data collected at a root process. Like reduce but with no operation which make the receive buffer bigger.
- ```
int MPI_Gather(void *sendbuf,  
              int sendcount,  
              MPI_Datatype sendtype,  
              void *recvbuf,  
              int recvcount,  
              MPI_Datatype recvtype,  
              int root,  
              MPI_Comm comm).
```
- They are stored in the order of the ranks of the processes according to comm.



# Gather Operation cont.

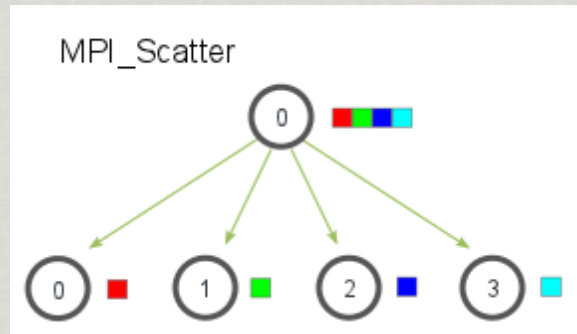
```
MPI_Comm comm;  
int sendbuf[100], my_rank, root = 0, gsize, *rbuf;  
MPI_Comm_rank (comm, &my_rank);  
if (my_rank == root) {  
    MPI_Comm_size (comm, &gsize);  
    rbuf = (int *) malloc (gsize*100*sizeof(int));  
}  
MPI_Gather(sendbuf, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);
```



# Scatter Operation

- For a scatter operation, a root process provides a different data block for each participating process.

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

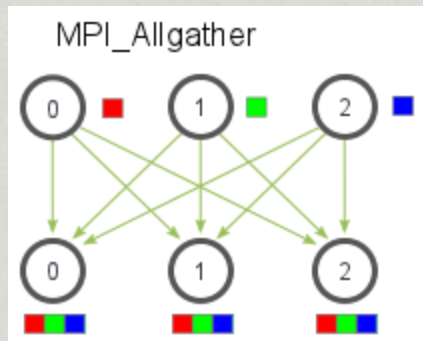


- They are scattered in the order of the ranks of the processes according to comm.

# Multi-broadcast Operation (MPI\_Allgather)

- Each participating process contributes a block of data which could, for example, be a partial result from a local computation

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



- They are gathered in the order of the ranks of the processes according to comm.

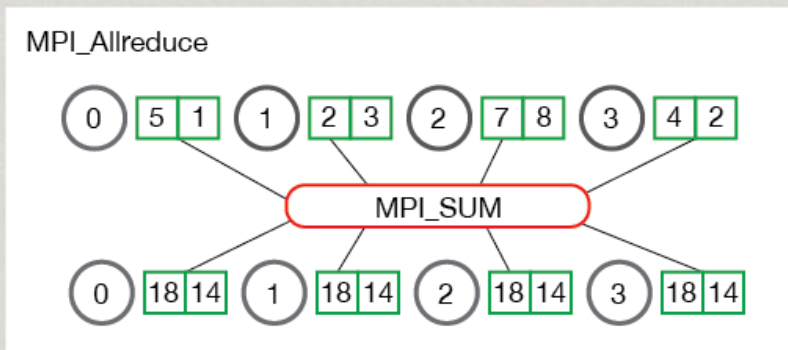
# Multi-broadcast Operation (MPI\_Allgather) cont.

```
int sbuf[100], gsize, *rbuf;  
MPI_Comm_size (comm, &gsize);  
rbuf = (int*) malloc (gsize*100*sizeof(int));  
MPI_Allgather (sbuf, 100, MPI_INT, rbuf, 100, MPI_INT, comm);
```

# Multi-accumulation Operation (MPI\_Allreduce)

- Each participating process performs a separate single-accumulation operation for which each process provides a different block of data.

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







QUESTIONS

# THANK YOU!

