

Introduction to Parallel Programming with MPI

Dr. Alireza Ghasemi, Dr. Georg Hager

Erlangen National High Performance Computing Center

Collective Operations

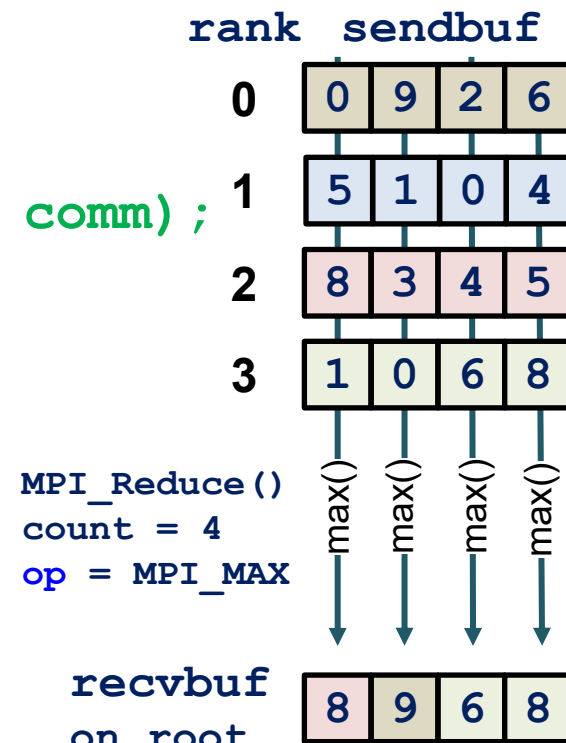


Global operations: reduction

- Compute results over distributed data

```
MPI_Reduce(sendbuf, recvbuf, count,  
           datatype, MPI_Op op, root, comm);
```

- Result in `recvbuf` only available on root process
- Perform operation on all `count` elements of an array
- If all ranks need the result, then use `MPI_Allreduce()`
- If the 12 predefined ops are not enough use `MPI_Op_create/MPI_Op_free` to create own ones



Global operations – predefined operators

Name	Operation	Name	Operation
<code>MPI_SUM</code>	Sum	<code>MPI_PROD</code>	Product
<code>MPI_MAX</code>	Maximum	<code>MPI_MIN</code>	Minimum
<code>MPI_LAND</code>	Logical AND	<code>MPI_BAND</code>	Bit-AND
<code>MPI_LOR</code>	Logical OR	<code>MPI_BOR</code>	Bit-OR
<code>MPI_LXOR</code>	Logical XOR	<code>MPI_BXOR</code>	Bit-XOR
<code>MPI_MAXLOC</code>	Maximum+Position	<code>MPI_MINLOC</code>	Minimum+Position

- Define own operations with `MPI_Op_create`/`MPI_Op_free`
- MPI assumes that the operations are associative
→ be careful with floating-point operations

“In-place” buffer specification

Override local input buffer with a result

MPI_Reduce

```
int partial_sum = ..., total_sum;
MPI_Reduce(&partial_sum, &total_sum,
          1, MPI_INT,
          MPI_SUM, root, comm);

int partial_sum = ..., total_sum;
if (rank == root) {
    total_sum = partial_sum;
    MPI_Reduce(MPI_IN_PLACE, &total_sum,
              1, MPI_INT,
              MPI_SUM, root, comm);
}
else {
    MPI_Reduce(&partial_sum, &total_sum,
              1, MPI_INT,
              MPI_SUM, root, comm);
}
```

MPI_Allreduce

```
int partial_sum = ..., total_sum;
MPI_AllReduce(&partial_sum, &total_sum,
             1, MPI_INT,
             MPI_SUM, comm);

int partial_sum = ..., total_sum;

total_sum = partial_sum;
MPI_AllReduce(MPI_IN_PLACE, &total_sum,
             1, MPI_INT,
             MPI_SUM, comm);
```

MPI_IN_PLACE cheat sheet

Function	MPI_IN_PLACE argument	@rank(s)	Comment [MPI 3.0]
<code>MPI_GATHER</code>	send buffer	root	Recv value at root already in the correct place in receive buffer.
<code>MPI_GATHERV</code>	send buffer	root	Recv value at root already in the correct place in receive buffer.
<code>MPI_SCATTER</code>	receive buffer	root	Root-th segment of send buffer is not moved.
<code>MPI_SCATTERV</code>	receive buffer	root	Root-th segment of send buffer is not moved.
<code>MPI_ALLGATHER</code>	send buffer	all	Input data at the correct place were process would receive its own contribution.
<code>MPI_ALLGATHERV</code>	send buffer	all	Input data at the correct place were process would receive its own contribution.
<code>MPI_ALLTOALL</code>	send buffer	all	Data to be sent is taken from receive buffer and replaced by received data, data sent/received must be of the same type map specified in receive count/receive type.
<code>MPI_ALLTOALLV</code>	send buffer	all	Data to be sent is taken from receive buffer and replaced by received data. Data sent/received must be of the same type map specified in receive count/receive type. The same amount of data and data type is exchanged between two processes.
<code>MPI_REDUCE</code>	send buffer	root	Data taken from receive buffer, replaced with output data.
<code>MPI_ALLREDUCE</code>	send buffer	all	Data taken from receive buffer, replaced with output data.

Summary of MPI collective communication

- MPI (blocking) **collectives**
 - **All ranks** in communicator **must call** the function
- **Communication** and synchronization
 - Barrier, broadcast, scatter, gather, and combinations thereof
- **Global operations**
 - **Reduce**, **allreduce**, some more...
- **In-place buffer** specification **MPI_IN_PLACE**
 - Save some space if need be

