

CISC 372: Parallel Computing

MPI Collectives

Stephen F. Siegel

Department of Computer and Information Sciences
University of Delaware

The collective model of computation

- ▶ a collective operation is invoked by all processes in a communicator
- ▶ some processes contribute data, some receive data, some operations may be performed
- ▶ collective operations capture many commonly used parallel patterns
 - ▶ one process broadcasts data to all
 - ▶ add up values across processes and return the result to one process
 - ▶ gather data from all processes into one big array on one process
- ▶ many parallel algorithms can be expressed most easily using collectives
- ▶ collective programs are **easier to understand** than point-to-point
 - ▶ a higher-level, abstract, “big step” view of an algorithm
 - ▶ all procs do this, then all procs do that, then all procs do, ...
 - ▶ you can reason about them much like a sequential program
- ▶ MPI implementations have very optimized implementations of collectives
 - ▶ these can be quite complicated, but the implementation details are hidden from the user

MPI_Reduce

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

`sendbuf` address of send buffer (`void*`)
`recvbuf` address of recv buffer (root only, `void*`)
`count` number of elements in send buffer (`int`)
`datatype` data type of elements in send buffer (`MPI_Datatype`)
`op` reduce operation (`MPI_Op`)
`root` rank of root process (`int`)
`comm` communicator (`MPI_Comm`)

| | | | |
|----------------|----------|----------|----------|
| Rank 0 sendbuf | x_{00} | x_{01} | x_{02} |
| Rank 1 sendbuf | x_{10} | x_{11} | x_{12} |
| Rank 2 sendbuf | x_{20} | x_{21} | x_{22} |

| | | | |
|--------------|----------------------------|----------------------------|----------------------------|
| Root recvbuf | $x_{00} + x_{10} + x_{20}$ | $x_{01} + x_{11} + x_{21}$ | $x_{02} + x_{12} + x_{22}$ |
|--------------|----------------------------|----------------------------|----------------------------|

► see `reduce.c`

Creating global synchronization points: MPI_Barrier

`MPI_Barrier(comm)`

`comm` communicator (`MPI_Comm`)

- ▶ blocks calling process until all processes in `comm` call `MPI_Barrier`
- ▶ if one process in `comm` calls `MPI_Barrier(comm)`, all should
 - ▶ else **deadlock** ensues
- ▶ explicit barriers are rarely needed; some exceptions. . .
 - ▶ good practice: use barriers before calling `MPI_Wtime`
 - ▶ ensures all processes have reached that point
 - ▶ programs that use `MPI_ANY_SOURCE` (coming soon)
 - ▶ barriers may be necessary to control how sends are matched with received
- ▶ using barriers to control order of printing from different procs is **not** reliable
 - ▶ sometimes it works, sometimes it doesn't

Broadcast: MPI_Bcast

`MPI_Bcast(buffer, count, datatype, root, comm)`

| | |
|-----------------------|---------------------------------------------------------------|
| <code>buffer</code> | address of buffer (<code>void*</code>) |
| <code>count</code> | number of elements in buffer (<code>int</code>) |
| <code>datatype</code> | data type of elements in buffer (<code>MPI_Datatype</code>) |
| <code>root</code> | rank of root process (<code>int</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ broadcasts a message from a single process (`root`) to all other processes in `comm`
- ▶ on root, `buffer` acts as a send buffer; on non-root procs, `buffer` acts as a receive buffer
- ▶ after return, buffer will contain the same data as that on root
- ▶ is a **barrier** (global synchronization point) induced?
 - ▶ all procs must enter before any proc can exit ? **No**
- ▶ see `bcast.c`

MPI_Allreduce

`MPI_Allreduce(sendbuf, recvbuf, count, datatype, op, comm)`

| | |
|-----------------------|--------------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (<code>void*</code>) |
| <code>recvbuf</code> | address of recv buffer (<code>void*</code>) |
| <code>count</code> | number of elements in send buffer (<code>int</code>) |
| <code>datatype</code> | data type of elements in send buffer (<code>MPI_Datatype</code>) |
| <code>op</code> | reduce operation (<code>MPI_Op</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ just like `MPI_Reduce`, but no `root`
- ▶ instead, result is returned to `all processes` in `comm`
- ▶ equivalent to `MPI_Reduce` followed by `MPI_Bcast`
- ▶ barrier? `yes`

MPI_Scatter

```
MPI_Scatter(sendbuf, sendcount, sendtype,  
            recvbuf, recvcount, recvtype, root, comm)
```

| | |
|------------------------|--------------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (root only, <code>void*</code>) |
| <code>sendcount</code> | num. elements sent to each proc (root, <code>int</code>) |
| <code>sendtype</code> | data type of send buf. elements (root, <code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (<code>void*</code>) |
| <code>recvcount</code> | number of elements in recv buffer (<code>int</code>) |
| <code>recvtype</code> | type of data to receive (<code>MPI_Datatype</code>) |
| <code>root</code> | rank of sending process (<code>int</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ similar to broadcast: the root sends, everyone else receives
- ▶ but: root sends a **different block** of `sendbuf` to each proc
 - ▶ rank 0 gets the first `sendcount` elements, rank 1 gets the next `sendcount` elements ...
- ▶ number of elements in `sendbuf` is `nprocs*sendcount`; see `scatter.c`

MPI_Scatterv

```
MPI_Scatterv(sendbuf, sendcounts, displs, sendtype,  
             recvbuf, recvcount, recvtype, root, comm)
```

| | |
|-------------------------|-------------------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (root only, <code>void*</code>) |
| <code>sendcounts</code> | num. elements sent to each proc (root only, <code>int[nprocs]</code>) |
| <code>displs</code> | displacements for each proc (root only, <code>int[nprocs]</code>) |
| <code>sendtype</code> | data type of send buf. elements (root only, <code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (<code>void*</code>) |
| <code>recvcount</code> | number of elements in recv buffer (<code>int</code>) |
| <code>recvtype</code> | type of data to receive (<code>MPI_Datatype</code>) |
| <code>root</code> | rank of sending process (<code>int</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ generalizes `MPI_Scatter`: the amount of data sent to each process can vary
- ▶ `sendcounts[i]` = number of elements to send to proc *i*
- ▶ `displs[i]` = offset of send buffer for proc *i* relative to `sendbuf`

MPI_Gather

```
MPI_Gather(sendbuf, sendcount, sendtype,  
          recvbuf, recvcount, recvtype, root, comm)
```

| | |
|------------------------|---------------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (<code>void*</code>) |
| <code>sendcount</code> | num. elements to send (<code>int</code>) |
| <code>sendtype</code> | data type of send buf. elements (<code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (root only, <code>void*</code>) |
| <code>recvcount</code> | number of elements to recv from each proc (root, <code>int</code>) |
| <code>recvtype</code> | type of data to receive (root, <code>MPI_Datatype</code>) |
| <code>root</code> | rank of receiving process (<code>int</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ **inverse** of MPI_Scatter: everyone sends to root, root receives from everyone
- ▶ root receives into a **different block** of `recvbuf` for each proc
 - ▶ rank 0's message goes into first `recvcount` elements
 - ▶ rank 1's message goes into next `recvcount` elements ...
- ▶ number of elements in `recvbuf` is $nprocs * \text{recvcount}$

MPI_Allgather

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

| | |
|------------------------|----------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (<code>void*</code>) |
| <code>sendcount</code> | num. elements to send (<code>int</code>) |
| <code>sendtype</code> | data type of send buf. elements (<code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (<code>void*</code>) |
| <code>recvcount</code> | number of elements to recv from each proc (<code>int</code>) |
| <code>recvtype</code> | type of data to receive (<code>MPI_Datatype</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

► like MPI_Gather done once for each proc

MPI_Gatherv

```
MPI_Gatherv(sendbuf, sendcount, sendtype,  
            recvbuf, recvcunts, displs, recvtype, root, comm)
```

| | |
|------------------------|-----------------------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (<code>void*</code>) |
| <code>sendcount</code> | num. elements to send (<code>int</code>) |
| <code>sendtype</code> | data type of send buf. elements (<code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (root, <code>void*</code>) |
| <code>recvcunts</code> | number of elements to recv from each proc (root, <code>int[nprocs]</code>) |
| <code>displs</code> | displacements of receive buffers (root, <code>int[nprocs]</code>) |
| <code>recvtype</code> | type of data to receive (root, <code>MPI_Datatype</code>) |
| <code>root</code> | rank of receiving process (<code>int</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ generalizes `MPI_Gather`
- ▶ the amount sent by each process can vary
- ▶ see `gatherv.c`

MPI_Alltoall

```
MPI_Alltoall(sendbuf, sendcount, sendtype,  
             recvbuf, recvcount, recvtype, comm)
```

| | |
|------------------------|-----------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (<code>void*</code>) |
| <code>sendcount</code> | num. elements to send (<code>int</code>) |
| <code>sendtype</code> | data type of send buf. elements (<code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (<code>void*</code>) |
| <code>recvcount</code> | number of elements to receive on each proc (<code>int</code>) |
| <code>recvtype</code> | type of element to receive (<code>MPI_Datatype</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ every process sends distinct buffers to all others
- ▶ amount of data sent to each process is the same
- ▶ symmetric (no root)

MPI_Alltoallv

```
MPI_Alltoallv(sendbuf, sendcounts, sdispls, sendtype,  
              recvbuf, recvcunts, rdispls, recvtype, comm)
```

| | |
|-------------------------|------------------------------------------------------------------------|
| <code>sendbuf</code> | address of send buffer (<code>void*</code>) |
| <code>sendcounts</code> | num. elements to send to others (<code>int[nprocs]</code>) |
| <code>sdispls</code> | displacements of send buffers (<code>int[nprocs]</code>) |
| <code>sendtype</code> | data type of send buf. elements (<code>MPI_Datatype</code>) |
| <code>recvbuf</code> | address of receive buffer (<code>void*</code>) |
| <code>recvcunts</code> | number of elements to receive from others (<code>int[nprocs]</code>) |
| <code>rdispls</code> | displacements of receive buffers (<code>int[nprocs]</code>) |
| <code>recvtype</code> | type of element to receive (<code>MPI_Datatype</code>) |
| <code>comm</code> | communicator (<code>MPI_Comm</code>) |

- ▶ generalizes `MPI_Alltoall`
- ▶ see `emily.c`