

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

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

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

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?

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 $nprocs * recvcount$

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

