

COMS3008A: Parallel Computing

Introduction to MPI II

Hairong Wang

School of Computer Science & Applied Mathematics
University of the Witwatersrand, Johannesburg

2021-10-7

1 Point to Point Communication

- Blocking vs. Non-blocking
- MPI Message Passing Function Arguments
- Avoiding Deadlocks
- Sending and Receiving Messages Simultaneously
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

1 Point to Point Communication

- Blocking vs. Non-blocking
- MPI Message Passing Function Arguments
- Avoiding Deadlocks
- Sending and Receiving Messages Simultaneously
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

1 Point to Point Communication

- **Blocking vs. Non-blocking**
- MPI Message Passing Function Arguments
- Avoiding Deadlocks
- Sending and Receiving Messages Simultaneously
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

Blocking vs. Non-blocking

- Most of the MPI point-to-point functions can be used in either blocking or non-blocking mode.
- Blocking:
 - A blocking send function will only “return” after it is safe to modify the application buffer (your send data) for reuse.
 - A blocking send can be synchronous which means there is handshaking occurring with the receive task to confirm a safe send.
 - A blocking send can be asynchronous if a system buffer is used to hold the data for eventual delivery to the receive.
 - A blocking receive only “returns” after the data has arrived and is ready for use by the program.

Blocking vs. Non-blocking cont.

- Non-blocking:
 - Non-blocking send and receive functions behave similarly - they will return almost immediately. They do not wait for any communication events to complete.
 - Non-blocking operations simply “request” the MPI library to perform the operation when it is able.
 - It is unsafe to modify the application buffer (your variable space) until you know for a fact the requested non-blocking operation was actually performed by the library.
 - Non-blocking communications are primarily used to overlap computation with communication and exploit possible performance gains.

1 Point to Point Communication

- Blocking vs. Non-blocking
- **MPI Message Passing Function Arguments**
- Avoiding Deadlocks
- Sending and Receiving Messages Simultaneously
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

MPI Message Passing Function Arguments

MPI point-to-point communication functions generally have an argument list that takes one of the following formats:

Blocking sends	<code>MPI_Send(buffer, count, type, dest, tag, comm)</code>
Non-blocking sends	<code>MPI_Isend(buffer, count, type, dest, tag, comm, request)</code>
Blocking receive	<code>MPI_Recv(buffer, count, type, source, tag, comm, status)</code>
Non-blocking receive	<code>MPI_Irecv(buffer, count, type, source, tag, comm, request)</code>

1 Point to Point Communication

- Blocking vs. Non-blocking
- MPI Message Passing Function Arguments
- **Avoiding Deadlocks**
- Sending and Receiving Messages Simultaneously
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

Avoiding Deadlocks

- The semantics of `MPI_Send` and `MPI_Recv` place some restrictions on how we can mix and match send and receive operations.
- Sources of deadlocks:
 - Send a large message from one process to another process
 - If the receive buffer is not large enough at the destination, the send must wait for the user to provide the memory space (through a receive)
 - Mismatched send and receive – unsafe.
- What happens with
- Order the operations.

Process 0	Process 1
Send(1)	Send(0)
Recv(0)	Recv(1)

Process 0	Process 1
Send(1)	Recv(1)
Recv(0)	Send(0)

Avoiding Deadlocks

Example 1

Process 0 sends two messages with different tags to process 1, and process 1 receives them in reverse order.

```
1 int a[10], b[10], myrank;
2 MPI_Status status;
3 ...
4 MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
5 if (myrank == 0) {
6     MPI_Send(a, 10, MPI_INT, 1, 1, MPI_COMM_WORLD);
7     MPI_Send(b, 10, MPI_INT, 1, 2, MPI_COMM_WORLD);
8 } else if (myrank == 1) {
9     MPI_Recv(b, 10, MPI_INT, 0, 2, MPI_COMM_WORLD,
10             MPI_STATUS_IGNORE);
11     MPI_Recv(a, 10, MPI_INT, 0, 1, MPI_COMM_WORLD,
12             MPI_STATUS_IGNORE);
13 }
14 ...
```

Avoiding Deadlocks

Example 2

Consider the following piece of code, in which process i sends a message to process $i + 1$ (modulo the number of processes) and receives a message from process $i - 1$ (modulo the number of processes). (Note that this example is different from the token ring example in Lec7.)

```
1 int a[10], b[10], npes, myrank;  
2 MPI_Status status;  
3 ...  
4 MPI_Comm_size(MPI_COMM_WORLD, &npes);  
5 MPI_Comm_rank(MPI_COMM_WORLD, &myrank);  
6 MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,  
7          MPI_COMM_WORLD);  
8 MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,  
9          MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
10 ...
```

Avoiding Deadlocks

Example 2 cont.

We can break the circular wait to avoid deadlocks as follows:

```
1 int a[10], b[10], npes, myrank;  
2 MPI_Status status;  
3 ...  
4 MPI_Comm_size(MPI_COMM_WORLD, &npes);  
5 MPI_Comm_rank(MPI_COMM_WORLD, &myrank);  
6 if (myrank%2 == 0) {  
7     MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,  
8         MPI_COMM_WORLD);  
9     MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,  
10        MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
11 } else {  
12     MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,  
13        MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
14     MPI_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,  
15        MPI_COMM_WORLD);  
16 }  
17
```

1 Point to Point Communication

- Blocking vs. Non-blocking
- MPI Message Passing Function Arguments
- Avoiding Deadlocks
- **Sending and Receiving Messages Simultaneously**
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

Sending and Receiving Messages Simultaneously

To exchange messages, MPI provides the following function that both sends and receives a message:

```
1 int MPI_Sendrecv(void *sendbuf, int sendcount,  
2 MPI_Datatype senddatatype, int dest, int sendtag,  
3 void *recvbuf, int recvcount, MPI_Datatype  
4 recvdatatype,  
5 int source, int recvtag, MPI_Comm comm,  
MPI_Status *status)
```

The arguments include arguments to the send and receive functions.

Using MPI_Sendrecv in Example 2

Example 2 can be made “safe” by using MPI_Sendrecv:

```
int a[10], b[10], npes, myrank;  
MPI_Status status;  
...  
MPI_Comm_size(MPI_COMM_WORLD, &npes);  
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);  
MPI_Sendrecv(a, 10, MPI_INT, (myrank+1)%npes, 1, b, 10,  
             MPI_INT, (myrank-1+npes)%npes, 1,  
             MPI_COMM_WORLD, &status);  
...
```


1 Point to Point Communication

- Blocking vs. Non-blocking
- MPI Message Passing Function Arguments
- Avoiding Deadlocks
- Sending and Receiving Messages Simultaneously
- **Overlapping Communication with Computation**

2 Collective Communications

- MPI Collective Communication Illustrations

3 Examples

4 Summary

Overlapping Communication with Computation

- In order to overlap communication with computation, MPI provides a pair of functions for performing non-blocking send and receive operations.

```
int MPI_Isend(void *buf, int count, MPI_Datatype datatype,
              int dest, int tag, MPI_Comm comm,
              MPI_Request *request)
int MPI_Irecv(void *buf, int count, MPI_Datatype datatype,
              int source, int tag, MPI_Comm comm,
              MPI_Request *request)
```

- These operations return before the operations have been completed. Function `MPI_Test` tests whether or not the non-blocking send or receive operation identified by its `request` has finished.
- ```
int MPI_Test(MPI_Request *request, int *flag,
 MPI_Status *status)
```

- **MPI\_Wait** blocks until a specified non-blocking send or receive operation has completed. For multiple non-blocking operations, the programmer can specify all completions.

```
int MPI_Wait(MPI_Request *request, MPI_Status *status)
```

- **MPI\_Request** handle is used to determine whether an operations has completed.
  - Non-blocking wait: **MPI\_Test**
  - Blocking wait: **MPI\_Wait**

- Anywhere you use **MPI\_Send** or **MPI\_Recv**, you can use the pair of **MPI\_Isend/MPI\_Wait** or **MPI\_Irecv/MPI\_Wait**.

- It is sometimes desirable to wait on multiple requests:

- ```
MPI_Waitall(int count,  
MPI_Request array_of_requests[],  
MPI_Status array_of_statuses[])
```

- The corresponding version of **MPI_Test**

```
int MPI_Testall(int count,  
MPI_Request array_of_requests[], int *flag,  
MPI_Status array_of_statuses[])
```

flag: true if all the requests are completed, otherwise false

Example 3

```
1  int main(int argc, char *argv[]){
2      int myid, numprocs, left, right, flag=0;
3      int buffer1[10], buffer2[10];
4      MPI_Request request; MPI_Status status;
5      MPI_Init(&argc,&argv);
6      MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
7      MPI_Comm_rank(MPI_COMM_WORLD, &myid);
8      /* initialize buffer2 */
9      .....
10     right = (myid + 1) % numprocs;
11     left = myid - 1;
12     if (left < 0)
13         left = numprocs - 1;
14     MPI_Irecv(buffer1, 10, MPI_INT, left, 123,
15              MPI_COMM_WORLD,
16              &request);
17     MPI_Send(buffer2, 10, MPI_INT, right, 123,
18              MPI_COMM_WORLD);
19     MPI_Test(&request, &flag, &status);
20     while (!flag){
21         /* Do some work ... */
22         MPI_Test(&request, &flag, &status);
23     }
```

Example 4

```
1 int main(int argc, char *argv){
2     int numtasks, rank, next, prev, buf[2], tag1=1, tag2=2;
3     MPI_Request reqs[4]; MPI_Status stats[4];
4     MPI_Init(&argc,&argv);
5     MPI_Comm_size(MPI_COMM_WORLD, &numtasks);
6     MPI_Comm_rank(MPI_COMM_WORLD, &rank);

8     prev = rank-1; next = rank+1;
9     if (rank == 0) prev = numtasks - 1;
10    if (rank == (numtasks - 1)) next = 0;
11    MPI_Irecv(&buf[0],1,MPI_INT,prev,tag1,MPI_COMM_WORLD,
12             &reqs[0]);
13    MPI_Irecv(&buf[1],1,MPI_INT,next,tag2,MPI_COMM_WORLD,
14             &reqs[1]);

16    MPI_Isend(&rank,1,MPI_INT,prev,tag2,MPI_COMM_WORLD,
17             &reqs[2]);
18    MPI_Isend(&rank,1,MPI_INT,next,tag1,MPI_COMM_WORLD,
19             &reqs[3]);
20    MPI_Waitall(4, reqs, stats);
21    MPI_Finalize();
22 }
```

Example 5

- We can use **the trapezoidal rule** to approximate the area between the graph of a function, $y = f(x)$, two vertical lines, and the x -axis.

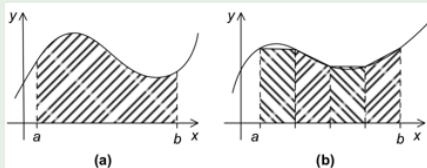


Figure: The trapezoidal rule: (a) area to be estimated, (b) estimate area using trapezoids

Example 5 cont.

- If the endpoints of the subinterval are x_i and x_{i+1} , then the length of the subinterval is $h = x_{i+1} - x_i$. Also, if the lengths of the two vertical segments are $f(x_i)$ and $f(x_{i+1})$, then the area of the trapezoid is

$$\text{Area of one trapezoid} = \frac{h}{2}(f(x_i) + f(x_{i+1})).$$

- Since we chose the N subintervals, we also know that the bounds of the region are $x = a$ and $x = b$ then

$$h = \frac{b - a}{N}$$

Example 5 cont.

- The pseudo code for a serial program:

```
h = (b-a) / N;  
approx = (f(a) + f(b)) / 2.0;  
for(i=1; i<=n-1; i++){  
    x_i = a + i * h;  
    approx += f(x_i);  
}  
approx = h * approx;
```

Recall we can design a parallel program using four basic steps:

- 1 Partition the problem solution into tasks.
- 2 Identify the communication between the tasks.
- 3 Aggregate the tasks into composite tasks.
- 4 Map the composite tasks to cores.

Example 5: Parallel Algorithm for the Trapezoidal Rule

Assuming `comm_sz` evenly divides n , the pseudo-code for the parallel program looks like the following:

```
1  Get a, b, n;
2  h = (b - a)/n;
3  local_n = n/comm_sz;
4  local_a = a + my_rank * local_n * h;
5  local_b = local_a + local_n * h;
6  local_integral = Trap(local_a, local_b, local_n, h);
7  if (my_rank != 0)
8      Send local integral to process 0;
9  else { /* my_rank == 0 */
10     total_integral = local_integral;
11     for (proc = 1; proc < comm_sz; proc++) {
12         Receive local_integral from proc;
13         total_integral += local_integral;
14     }
15 }
16 if (my_rank == 0)
17     print result;
```

Dealing with I/O

- In most cases, all the processes in `MPI_COMM_WORLD` have access to `stdout` and `stderr`.
- The order in which the processes' output appears is indeterministic.
- For the input, i.e., `stdin`, usually, only process 0 has access to.
- If an MPI program uses `scanf` function, then process 0 reads in the data, and sends it to the other processes.

Outline

- 1 Point to Point Communication
 - Blocking vs. Non-blocking
 - MPI Message Passing Function Arguments
 - Avoiding Deadlocks
 - Sending and Receiving Messages Simultaneously
 - Overlapping Communication with Computation
- 2 Collective Communications
 - MPI Collective Communication Illustrations
- 3 Examples
- 4 Summary

Collective Communications

- Communication is coordinated among a group of processes, as specified by the communicator.
- All collective operations are blocking and no message tags are used.
- All processes in the communicator must call the collective operation.
- Three classes of collective operations
 - Data movement
 - Collective computation
 - Synchronization

Outline

- 1 Point to Point Communication
 - Blocking vs. Non-blocking
 - MPI Message Passing Function Arguments
 - Avoiding Deadlocks
 - Sending and Receiving Messages Simultaneously
 - Overlapping Communication with Computation
- 2 Collective Communications
 - MPI Collective Communication Illustrations
- 3 Examples
- 4 Summary

MPI_Bcast

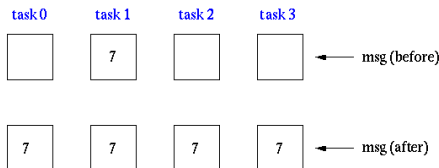
- A collective communication in which data belonging to a single process is sent to all of the processes in the communicator is called a broadcast — `MPI_Bcast`.

MPI_Bcast

Broadcasts a message to all other processes of that group

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

broadcast originates in task 1



- The process with rank `source` sends the contents of the memory referenced by `msg` to all the processes in the communicator `MPI_COMM_WORLD`.

Example 5 Cont.

- ❶ In our example `mpi_trapezoid_1.c`, we are using

```
1 if (my_rank == 0) {
2     for (dest = 1; dest < comm_sz; dest++) {
3         MPI_Send(a_p, 1, MPI_DOUBLE, dest, 0,
4                 MPI_COMM_WORLD);
5         MPI_Send(b_p, 1, MPI_DOUBLE, dest, 0,
6                 MPI_COMM_WORLD);
7         MPI_Send(n_p, 1, MPI_INT, dest, 0, MPI_COMM_WORLD
8                 );
9     }
10 } else { /* my rank != 0 */
11     MPI_Recv(a_p, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD
12             , MPI_STATUS_IGNORE);
13     MPI_Recv(b_p, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD
14             , MPI_STATUS_IGNORE);
15     MPI_Recv(n_p, 1, MPI_INT, 0, 0, MPI_COMM_WORLD,
16             MPI_STATUS_IGNORE);
17 }
```

- ❷ Instead of using point-to-point communications, you can use collective communications here. **Write another function to implement this part using** — `MPI_Bcast()`.

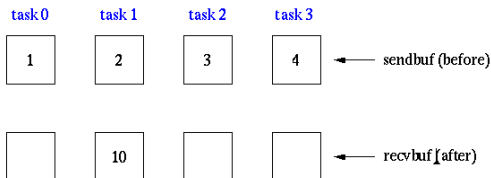
MPI_Reduce

- `MPI_Reduce` combines data from all processes in the communicator and returns it to one process.
- In many numerical algorithms, `Send/Receive` can be replaced by `Bcast/Reduce`, improving both simplicity and efficiency.

MPI_Reduce

Perform and associate reduction operation across all tasks in the group and place the result in one task

```
count = 1;  
dest = 1;           result will be placed in task 1  
MPI_Reduce(sendbuf, recvbuf, count, MPI_INT, MPI_SUM,  
            dest, MPI_COMM_WORLD);
```



MPI_Reduce cont.

- When the `count` is greater 1, `MPI_Reduce` operates on arrays instead of scalars.

```
1  double local_x[N], sum[N];  
2  ...  
3  MPI_Reduce(local_x, sum, N, MPI_DOUBLE, MPI_SUM, 0, \  
4  MPI_COMM_WORLD);
```

Example 5 cont.

- ❶ In our example `mpi_trapezoid_1.c`, we are using

```
1 /* Add up the integrals calculated by each process */
2 if (my_rank != 0) {
3     MPI_Send(&local_int, 1, MPI_DOUBLE, 0, 0,
4             MPI_COMM_WORLD);
5 } else {
6     total_int = local_int;
7     for (source = 1; source < comm_sz; source++) {
8         MPI_Recv(&local_int, 1, MPI_DOUBLE, source, 0,
9                 MPI_COMM_WORLD, MPI_STATUS_IGNORE);
10        total_int += local_int;
11    }
```

- ❷ Instead of using point-to-point communications, you can also use collective communications here. Rewrite this part using appropriate collective communication.

MPI_Reduce

- Suppose that each process calls `MPI_Reduce` with operator `MPI_SUM`, and **destination process 0**. What happens with the following multiple calls of `MPI_Reduce`? What are the values for `b` and `d`?

Time	Process 0	Process 1	Process 2
0	<code>a=1; c = 2;</code>	<code>a=1; c = 2;</code>	<code>a=1; c = 2;</code>
1	<code>MPI_Reduce (&a, &b, 1, ...)</code>	<code>MPI_Reduce (&c, &d, 1, ...)</code>	<code>MPI_Reduce (&a, &b, 1, ...)</code>
2	<code>MPI_Reduce (&c, &d, 1, ...)</code>	<code>MPI_Reduce (&a, &b, 1, ...)</code>	<code>MPI_Reduce (&c, &d, 1, ...)</code>

- The order of the calls will determine the matching.
- What will happen with the following code?

```
MPI_Reduce (&x, &x, 1, MPI_DOUBLE, MPI_SUM, 0, comm);
```

MPI_Allreduce

- If the result of the reduction operation is needed by all processes, MPI provides:

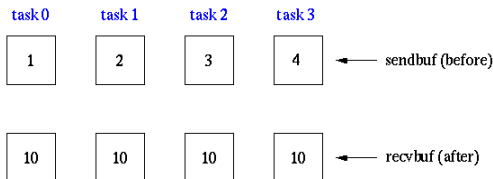
```
int MPI_Allreduce(void *sendbuf, void *recvbuf, int count,  
MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

- This is equivalent to an MPI_Reduce followed by an MPI_Bcast.

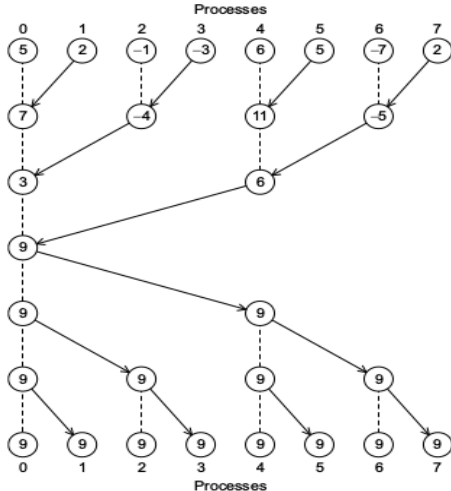
MPI_Allreduce

Perform and associate reduction operation across all tasks in the group and place the result in all tasks

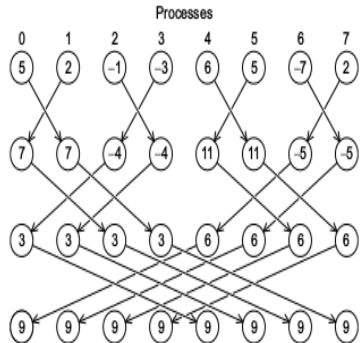
```
count = 1;  
MPI_Allreduce(sendbuf, recvbuf, count, MPI_INT, MPI_SUM,  
MPI_COMM_WORLD);
```



MPI_Allreduce cont.



(a)



(b)

Figure: (a) A global sum followed by a broadcasting; (2) A butterfly structured global sum.

- The scatter operation is to distribute **distinct messages** from a single source task (or process) to each task in the group.

```
int MPI_Scatter(void *sendbuf, int sendcount,  
               MPI_Datatype senddatatype, void *recvbuf,  
               int recvcount, MPI_Datatype recvdatatype,  
               int source, MPI_Comm comm)
```

MPI_Scatter

Sends data from one task to all other tasks in a group

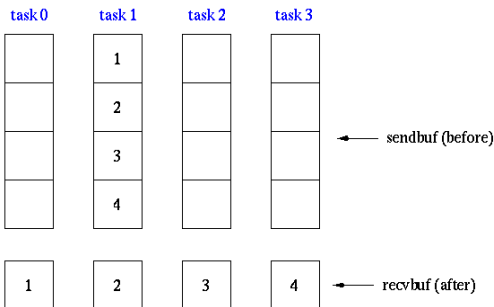
```
sendcnt = 1;
```

```
recvnt = 1;
```

```
src = 1;
```

task 1 contains the message to be scattered

```
MPI_Scatter(sendbuf, sendcnt, MPI_INT,  
            recvbuf, recvnt, MPI_INT,  
            src, MPI_COMM_WORLD);
```



- The gather operation is performed in MPI using `MPI_Gather`.
 - Gathers **distinct messages** from each task in the group to a single destination task.
 - Reverse operation of `MPI_Scatter`.

```
int MPI_Gather(void *sendbuf, int sendcount,  
              MPI_Datatype senddatatype, void *recvbuf,  
              int recvcnt, MPI_Datatype recvdtype,  
              int target, MPI_Comm comm)
```


MPI_Gather

Gathers together values from a group of processes

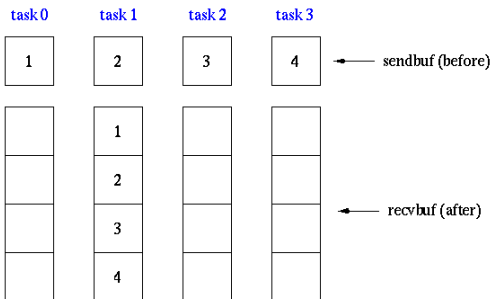
```
sendcnt = 1;
```

```
recvcnt = 1;
```

```
src = 1;
```

messages will be gathered in task 1

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



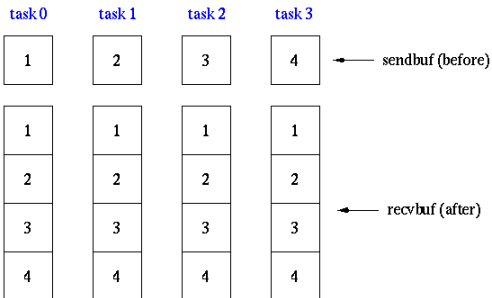
- MPI also provides the `MPI_Allgather` function in which the data are gathered at all the processes.

```
int MPI_Allgather(void *sendbuf, int sendcount,  
                 MPI_Datatype senddatatype, void *recvbuf,  
                 int recvcount, MPI_Datatype recvdatatype,  
                 MPI_Comm comm)
```

MPI_Allgather

Gathers together values from a group of processes and distributes to all

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



- The all-to-all communication operation is performed by:

```
int MPI_Alltoall(void *sendbuf, int sendcount,  
                MPI_Datatype senddatatype, void *recvbuf,  
                int recvcount, MPI_Datatype recvdatatype,  
                MPI_Comm comm)
```

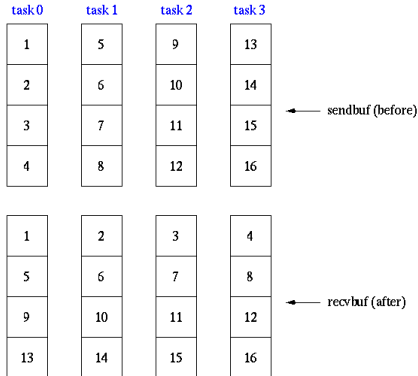
- Each task in a group performs a scatter operation, sending a distinct message to all the tasks in the group in order by index.

MPI_Alltoall cont.

MPI_Alltoall

Sends data from all to all processes. Each process performs a scatter operation.

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



Example 6 (Matrix vector multiplication)

If $A = (a_{ij})$ is an $m \times n$ matrix and \mathbf{x} is a vector with n components, then $\mathbf{y} = A\mathbf{x}$ is a vector with m components. Furthermore,

$$y_i = a_{i0}x_0 + a_{i1}x_1 + a_{i2}x_2 + \dots + a_{i,n-1}x_{n-1}.$$

A serial code can be as simple as

```
1  for (i = 0; i < m; i++) {  
2      y[i] = 0.0;  
3      for (j = 0; j < n; j++)  
4          y[i] += A[i*n+j]*x[j];  
5  }
```

Example 6 cont.

Process 0 reads in the matrix and distributes row blocks to all the processes in communicator `comm`.

```
1  if (my_rank == 0) {
2      A = malloc(m*n*sizeof(double));
3      if (A == NULL) local_ok = 0;
4      Check_for_error(local_ok, "Random_matrix", "Can't
        allocate temporary matrix", comm);
5      srand(2021);
6      for (i = 0; i < m; i++)
7          for (j = 0; j < n; j++)
8              A[i*n+j] = (double)rand( ) / RAND_MAX;
9      MPI_Scatter(A, local_m*n, MPI_DOUBLE, local_A,
        local_m*n, MPI_DOUBLE, 0, comm);
10     free(A);
11 } else {
12     Check_for_error(local_ok, "Random_matrix", "Can't
        allocate temporary matrix", comm);
13     MPI_Scatter(A, local_m*n, MPI_DOUBLE, local_A,
        local_m*n, MPI_DOUBLE, 0, comm);
14 }
```

Example 6 cont.

Each process gathers the entire vector, then proceeds to compute its share of sub-matrix and vector multiplication.

```
1  MPI_Allgather(local_x, local_n, MPI_DOUBLE,  
2      x, local_n, MPI_DOUBLE, comm);  
  
4  for (local_i = 0; local_i < local_m; local_i++) {  
5      local_y[local_i] = 0.0;  
6      for (j = 0; j < n; j++)  
7          local_y[local_i] += local_A[local_i*n+j]*x[j];  
8  }
```


MPI_Scan

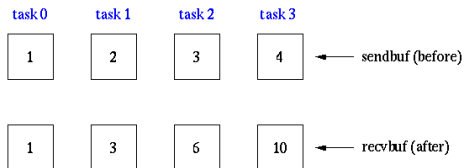
- To compute prefix-sums, MPI provides:

```
int MPI_Scan(void *sendbuf, void *recvbuf, int count,  
            MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
```

MPI_Scan

Computes the scan (partial reductions) of data
on a collection of processes

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



- Using this core set of collective operations, MPI communications can be greatly simplified.

Outline

1 Point to Point Communication

- Blocking vs. Non-blocking
- MPI Message Passing Function Arguments
- Avoiding Deadlocks
- Sending and Receiving Messages Simultaneously
- Overlapping Communication with Computation

2 Collective Communications

- MPI Collective Communication Illustrations

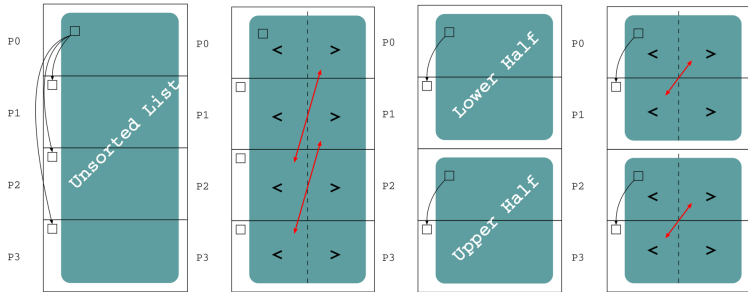
3 Examples

4 Summary

Parallel quicksort

- one process broadcast initial pivot to all processes;
- each process in the upper half swaps with a partner in the lower half
- recurse on each half
- swap among partners in each half
- each process uses quicksort on local elements

Parallel quicksort cont.



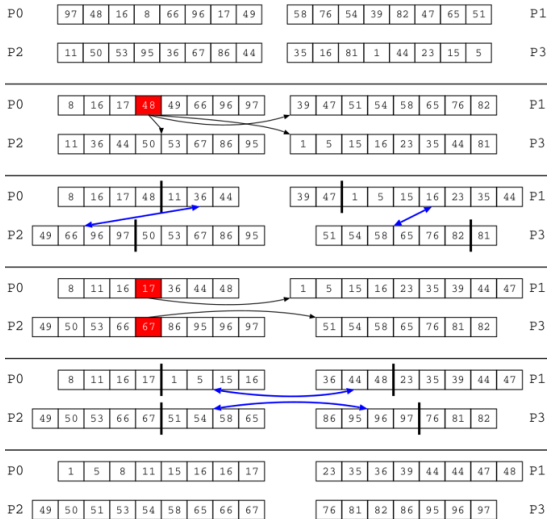
Hyperquicksort

Limitation of parallel quicksort: poor balancing of list sizes.

Hyperquicksort: sort elements before broadcasting pivot.

- sort elements in each process
- select median as pivot element and broadcast it
- each process in the upper half swaps with a partner in the lower half
- recurse on each half

Hyperquicksort cont.



Example 7

Task 0 pings task 1 and awaits return ping

```
1 #include "mpi.h"
2 #include <stdio.h>
3 main(int argc, char *argv[]){
4     int numtasks, rank, dest, source, rc, count, tag=1;
5     char inmsg, outmsg='x';
6     MPI_Status Stat;
7
8     MPI_Init(&argc,&argv);
9     MPI_Comm_size(MPI_COMM_WORLD, &numtasks);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11
12    if (rank == 0){
13        dest = 1;
14        source = 1;
15        rc = MPI_Send(&outmsg, 1, MPI_CHAR, dest, tag,
16                     MPI_COMM_WORLD);
17        rc = MPI_Recv(&inmsg, 1, MPI_CHAR, source, tag,
18                    MPI_COMM_WORLD, &Stat);
19    }
```

Example 7 cont.

```
1  else if (rank == 1) {
2      dest = 0;
3      source = 0;
4      rc = MPI_Recv(&inmsg, 1, MPI_CHAR, source, tag,
5                   MPI_COMM_WORLD, &Stat);
6      rc = MPI_Send(&outmsg, 1, MPI_CHAR, dest, tag,
7                   MPI_COMM_WORLD);
8  }
9
10 rc = MPI_Get_count(&Stat, MPI_CHAR, &count);
11 printf("Task %d: Received %d char(s) from task %d with
12        tag %d \n",
13        rank, count, Stat.MPI_SOURCE, Stat.MPI_TAG);
14
15 MPI_Finalize();
16 }
```


Example 8

Perform a scatter operation on the rows of an array

```
1 #include "mpi.h"
2 #include <stdio.h>
3 #define SIZE 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];

14    MPI_Init(&argc, &argv);
15    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
16    MPI_Comm_size(MPI_COMM_WORLD, &numtasks);
```

Example 8 cont.

```
1  if (numtasks == SIZE) {  
2      source = 1;  
3      sendcount = SIZE;  
4      recvcnt = SIZE;  
5      MPI_Scatter(sendbuf, sendcount, MPI_FLOAT, recvbuf,  
6                  recvcnt,  
7                  MPI_FLOAT, source, MPI_COMM_WORLD);  
8      printf("rank= %d Results: %f %f %f %f\n", rank,  
9             recvbuf[0], recvbuf[1], recvbuf[2], recvbuf[3]);  
10 } else  
    printf("Must specify %d processors. Terminating.\n",  
           SIZE);  
    MPI_Finalize();
```

Example 9

The Odd-Even Transposition Sort

- Sorts n elements in n phases (n is even), each of which requires $n/2$ compare-exchange operations.
- The algorithm alternates between two phases — odd and even phases.
- Let $\langle a_0, a_1, \dots, a_{n-1} \rangle$ be the sequence to be sorted.
 - During the odd phase, elements with odd indices are compared with their right neighbours, and if they are out of sequence they are exchanged; thus, the pairs $(a_1, a_2), (a_3, a_4), \dots, (a_{n-3}, a_{n-2})$ are compare exchanged.
 - During the even phase, elements with even indices are compared with their right neighbours, and if they are out of sequence they are exchanged; $(a_0, a_1), (a_2, a_3), \dots, (a_{n-2}, a_{n-1})$.
- After n phases of odd-even exchanges, the sequence is sorted. Each phase requires $n/2$ compare-exchange operations (sequential complexity $O(n^2)$).

Example 9 cont. – The serial algorithm

```
1  for i = 0 to n-1 do
2    if i is even then
3      for j = 0 to n/2 - 1 do
4        compare-exchange(a(2j), a(2j+1));
5    if i is odd then
6      for j = 0 to n/2 - 1 do
7        compare-exchange(a(2j+1), a(2j+2));
```

Example 9 cont. – The parallel algorithm

```
1  void oddevensort(int n)
2      id = process's label;
3      for i = 0 to n-1 do
4          if i is odd then
5              if id is odd then
6                  compare-exchange_min(id, id + 1); //increasing
                    comparator
7              else
8                  compare-exchange_max(id, id - 1); //decreasing
                    comparator
9          if i is even then
10             if id is even then
11                 compare-exchange_min(id, id + 1);
12             else
13                 compare-exchange_max(id, id - 1);
```

Outline

- 1 Point to Point Communication
 - Blocking vs. Non-blocking
 - MPI Message Passing Function Arguments
 - Avoiding Deadlocks
 - Sending and Receiving Messages Simultaneously
 - Overlapping Communication with Computation
- 2 Collective Communications
 - MPI Collective Communication Illustrations
- 3 Examples
- 4 Summary

Summary

- Point-to-point communication
 - Blocking vs non-blocking
 - Safety in MPI programs
- Collective communication
 - Collective communications involve all the processes in a communicator.
 - All the processes in the communicator must call the same collective function.
 - Collective communications do not use **tags**, the message is matched on the order in which they are called within the communicator.
 - The meanings of *local variable* and *global variable* in MPI
 - Some important MPI collective communications we learned:
MPI_Reduce, MPI_Allreduce, MPI_Bcast, MPI_Gather,
MPI_Scatter, MPI_Allgather, MPI_Alltoall, MPI_Scan
etc.

Useful references:

- Trobec et al., *Introduction to Parallel Computing: From Algorithms to Programming on State-of-the-Art Platforms*
- Quinn, *Parallel Programming in C with MPI and OpenMP*
- Grama et al., *Introduction to Parallel Computing*
- Barney, *Message Passing Interface*
- MPI, *MPI Forum*

Bibliography I

Barney, Blaise. *Message Passing Interface*.

<https://hpc-tutorials.llnl.gov/mpi/>. Accessed 2021-10-6.

Grama, Ananth et al. *Introduction to Parallel Computing*. Addison Wesley, 2003.

MPI. *MPI Forum*. <https://www.mpi-forum.org/docs/>. Accessed 2021-10-6.

Quinn, Michael J. *Parallel Programming in C with MPI and OpenMP*. McGraw-Hill Education Group, 2003.

Trobec, Roman et al. *Introduction to Parallel Computing: From Algorithms to Programming on State-of-the-Art Platforms*. Springer, 2018.