



BITS Pilani
Pilani Campus

Message Passing Interface (MPI)

K Hari Babu
Department of Computer Science & Information Systems



BITS Pilani
Pilani Campus



Programming with MPI: Derived Data Types

Sending Data in a Single Message

- Sending a, b and n in one message will take lesser time compared to sending a, b and n individually
 - MPI provides three basic approaches to consolidate data
 - Send in an array using count argument if data is of same type
 - Derived data types
 - MPI_Pack/MPI_Unpack are used to pack data into a contiguous buffer space
 - The pack/unpack routines are provided for compatibility with previous libraries.
- Derived data types is offered in later versions of MPI

```
09 if (my_rank == 0) {
10     printf("Enter a, b, and n\n");
11     scanf("%lf %lf %d", a,b,n);
12 for (dest = 1; dest < comm_sz; dest++) {
13     MPI_Send(a, 1, MPI_DOUBLE, dest, 0, MPI_COMM_WORLD);
14     MPI_Send(b, 1, MPI_DOUBLE, dest, 0, MPI_COMM_WORLD);
15     MPI_Send(n, 1, MPI_INT, dest, 0, MPI_COMM_WORLD);
16 }
17 } else { /* my rank != 0 */
18     MPI_Recv(a, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
20     MPI_Recv(b, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
22     MPI_Recv(n, 1, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
24 }
```

Create a New Data Type

- In MPI, a derived datatype can be used to represent any collection of data items in memory by storing both the types of the items and their relative locations in memory
- In our trapezoidal rule example, suppose that on process 0 the variables a, b, and n are stored in memory locations with the following addresses
- Then the following derived datatype could represent these data items
- `{(MPI_DOUBLE,0),(MPI_DOUBLE,16),(MPI_INT,24)}`
 - The first element of each pair corresponds to the type of the data, and the second element of each pair is the displacement of the data element from the beginning of the type

Variable	Address
a	24
b	40
n	48

Create a New Data Type

```
int MPI_Type_create_struct(  
    int count, //no of elements in the struct  
    int array_of_blocklengths[],  
    MPI_Aint array_of_displacements[],  
    MPI_Datatype array_of_types[],  
    MPI_Datatype *newtype  
);
```

- We can use MPI_Type_create_struct to build a derived datatype that consists of individual elements that have different basic types
- The argument count is the number of elements in the datatype
 - for our example, it should be three
- Each of the next 3 array arguments should have count elements
- The first array, array of block lengths, allows for the possibility that the individual data items might be arrays or subarrays
 - for example, the first element were an array containing five elements, we would have array_of_blocklengths[0] = 5;
 - in our case array_of_blocklengths[3] = {1, 1, 1};

Create a New Data Type

```
int MPI_Type_create_struct(  
    int count, //no of elements in the struct  
    int array_of_blocklengths[],  
    MPI_Aint array_of_displacements[],  
    MPI_Datatype array_of_types[],  
    MPI_Datatype *newtype  
);
```

array_of_blocklengths allows for the possibility that the individual data items might be arrays or subarrays.

- array_of_displacements, specifies the displacements, in bytes, from the start of the message
 - array_of_displacements[] = {0, 16, 24};
- We can use MPI_Get_address to calculate displacements

```
int MPI_Get_address(  
    void* location_p /* in */,  
    MPI_Aint* address_p /* out */);  
/*returns address of the memory  
location referenced by location_p*/
```

```
MPI_Aint a_addr, b_addr, n_addr;  
MPI_Get_address(&a, &a_addr);  
array_of_displacements[0] = 0;  
MPI_Get_address(&b, &b_addr);  
array_of_displacements[1] = b_addr - a_addr;  
MPI_Get_address(&n, &n_addr);  
array_of_displacements[2] = n_addr - a_addr;
```

Create a New Data Type

```
int MPI_Type_create_struct(  
    int count, //no of elements in the struct  
    int array_of_blocklengths[],  
    MPI_Aint array_of_displacements[],  
    MPI_Datatype array_of_types[],  
    MPI_Datatype *newtype  
);
```

- The array_of_types should store the MPI datatypes of the elements
 - array_of_types[3] = {MPI_DOUBLE, MPI_DOUBLE, MPI_INT}
- With these initializations, we can build the new datatype with the call

```
MPI_Datatype input_mpi_t;  
MPI_Type_create_struct(3, array_of_blocklengths,  
    array_of_displacements, array_of_types,  
    &input_mpi_t);
```

- Before we can use input_mpi_t in a communication function, we must first commit it with a call to

```
int MPI_Type_commit(MPI_Datatype* new_mpi_t p /*in/out*/);
```

Using Derived Data Type

- Now, in order to use `new_mpi_t`, we make the following call to `MPI_Bcast` on each process:
- `MPI_Bcast(&a, 1, input_mpi_t, 0, comm);`

Example



BITS Pilani

```
91  /*-----  
92  * Function:      Build_mpi_type  
93  * Purpose:       Build a derived datatype so that the three  
94  *               input values can be sent in a single message.  
95  * Input args:    a_p:  pointer to left endpoint  
96  *               b_p:  pointer to right endpoint  
97  *               n_p:  pointer to number of trapezoids  
98  * Output args:   input_mpi_t_p:  the new MPI datatype  
99  */  
100 void Build_mpi_type(  
101     double*      a_p          /* in */,  
102     double*      b_p          /* in */,  
103     int*         n_p          /* in */,  
104     MPI_Datatype* input_mpi_t_p /* out */) {  
105  
106     int array_of_blocklengths[3] = {1, 1, 1};  
107     MPI_Datatype array_of_types[3] = {MPI_DOUBLE, MPI_DOUBLE, MPI_INT};  
108     MPI_Aint a_addr, b_addr, n_addr;  
109     MPI_Aint array_of_displacements[3] = {0};  
110  
111     MPI_Get_address(a_p, &a_addr);  
112     MPI_Get_address(b_p, &b_addr);  
113     MPI_Get_address(n_p, &n_addr);  
114     array_of_displacements[1] = b_addr-a_addr;  
115     array_of_displacements[2] = n_addr-a_addr;  
116     MPI_Type_create_struct(3, array_of_blocklengths,  
117         array_of_displacements, array_of_types,  
118         input_mpi_t_p);  
119     MPI_Type_commit(input_mpi_t_p);  
120 } /* Build_mpi_type */
```

Example

```
132 void Get_input(  
133     int      my_rank  /* in */,  
134     int      comm_sz  /* in */,  
135     double*  a_p      /* out */,  
136     double*  b_p      /* out */,  
137     int*     n_p      /* out */) {  
138     MPI_Datatype input_mpi_t;  
139  
140     Build_mpi_type(a_p, b_p, n_p, &input_mpi_t);  
141  
142     if (my_rank == 0) {  
143         printf("Enter a, b, and n\n");  
144         scanf("%lf %lf %d", a_p, b_p, n_p);  
145     }  
146     MPI_Bcast(a_p, 1, input_mpi_t, 0, MPI_COMM_WORLD);  
147  
148     MPI_Type_free(&input_mpi_t);  
149 } /* Get_input */
```

Taking timings



- We're usually not interested in the time taken from the start of program execution to the end of program execution

```
78 MPI_Barrier(comm);
79 start = MPI_Wtime();
80 Mat_vect_mult(local_A, local_x, local_y, local_m, n, local_n, comm);
81 finish = MPI_Wtime();
82 loc_elapsed = finish-start;
83 MPI_Reduce(&loc_elapsed, &elapsed, 1, MPI_DOUBLE, MPI_MAX, 0, comm);
84
85 #ifdef DEBUG
86 Print_vector("y", local_y, m, local_m, my_rank, comm);
87 #endif
88
89 if (my_rank == 0)
90     printf("Elapsed time = %e\n", elapsed);
```

- Time taken by parallel algorithm is the time of the slowest process. We can't get exactly this time because we can't insure that all the processes start at the same instant. However, we can come reasonably close. The MPI collective communication function MPI Barrier insures that no process will return from calling it until every process in the communicator has started calling it.

Odd-even Transposition Sort



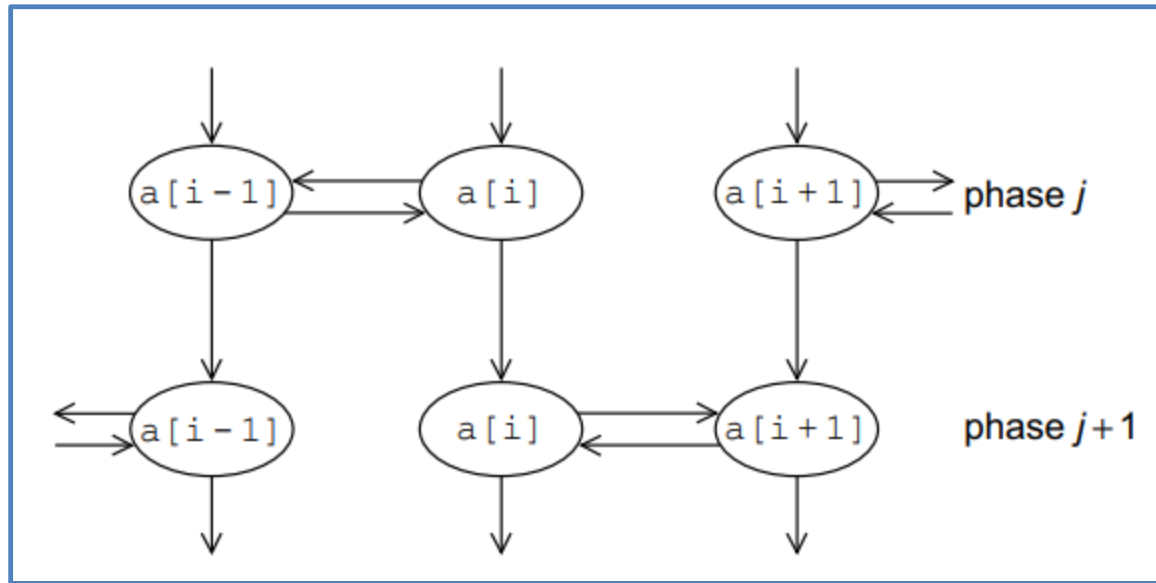
- Odd-even transposition sort is a sorting algorithm that's similar to bubble sort, but that has more opportunities for parallelism

```
137 void Odd_even_sort(  
138     int a[] /* in/out */,  
139     int n   /* in    */) {  
140     int phase, i, temp;  
141  
142     for (phase = 0; phase < n; phase++)  
143         if (phase % 2 == 0) { /* Even phase */  
144             for (i = 1; i < n; i += 2)  
145                 if (a[i-1] > a[i]) {  
146                     temp = a[i];  
147                     a[i] = a[i-1];  
148                     a[i-1] = temp;  
149                 }  
150             } else { /* Odd phase */  
151                 for (i = 1; i < n-1; i += 2)  
152                     if (a[i] > a[i+1]) {  
153                         temp = a[i];  
154                         a[i] = a[i+1];  
155                         a[i+1] = temp;  
156                     }  
157             }  
158     } /* Odd_even_sort */  
159 }
```

Suppose $a = \{9, 7, 8, 6\}$.

Phase	Subscript in Array						
	0		1		2		3
0	9	↔	7		8	↔	6
	7		9		6		8
1	7		9	↔	6		8
	7		6		9		8
2	7	↔	6		9	↔	8
	6		7		8		9
3	6		7	↔	8		9
	6		7		8		9

Parallel odd-even transposition sort



- When $n = p$: depending on the phase, process i can send its current value, $a[i]$, either to process $i - 1$ or process $i + 1$. At the same time, it should receive the value stored on process $i - 1$ or process $i + 1$, respectively, and then decide which of the two values it should store as $a[i]$ for the next phase.

Parallel odd-even transposition sort



- It's unlikely that $n = p$
- How should this be modified when each process is storing $n/p > 1$ elements?

At any time, each process should have exactly n/p elements.

```

1  Sort local keys;
2  for (phase = 0; phase < comm sz; phase++) {
3      partner = Compute partner(phase, my rank);
4      if (I'm not idle) {
5          Send my keys to partner;
6          Receive keys from partner;
7          if (my rank < partner)
8              Keep smaller keys;
9          else
10             Keep larger keys;
11     }
12 }
    
```

Time	Process			
	0	1	2	3
Start	15, 11, 9, 16	3, 14, 8, 7	4, 6, 12, 10	5, 2, 13, 1
After Local Sort	9, 11, 15, 16	3, 7, 8, 14	4, 6, 10, 12	1, 2, 5, 13
After Phase 0	3, 7, 8, 9	11, 14, 15, 16	1, 2, 4, 5	6, 10, 12, 13
After Phase 1	3, 7, 8, 9	1, 2, 4, 5	11, 14, 15, 16	6, 10, 12, 13
After Phase 2	1, 2, 3, 4	5, 7, 8, 9	6, 10, 11, 12	13, 14, 15, 16
After Phase 3	1, 2, 3, 4	5, 6, 7, 8	9, 10, 11, 12	13, 14, 15, 16

Parallel odd-even transposition sort



- how do we compute the partner rank?

```

1  if (phase % 2 == 0) /* Even phase */
2      if (my_rank % 2 != 0) /* Odd rank */
3          partner = my_rank - 1;
4      else /* Even rank */
5          partner = my_rank + 1;
6  else /* Odd phase */
7      if (my_rank % 2 != 0) /* Odd rank */
8          partner = my_rank + 1;
9      else /* Even rank */
10         partner = my_rank - 1;
11  if (partner == -1 || partner == comm sz)
12     partner = MPI_PROC_NULL;
    
```

MPI_PROC_NULL is a constant defined by MPI. When it's used as the source or destination rank in a point-to-point communication, no communication will take place and the call to the communication will simply return.

Time	Process			
	0	1	2	3
Start	15, 11, 9, 16	3, 14, 8, 7	4, 6, 12, 10	5, 2, 13, 1
After Local Sort	9, 11, 15, 16	3, 7, 8, 14	4, 6, 10, 12	1, 2, 5, 13
After Phase 0	3, 7, 8, 9	11, 14, 15, 16	1, 2, 4, 5	6, 10, 12, 13
After Phase 1	3, 7, 8, 9	1, 2, 4, 5	11, 14, 15, 16	6, 10, 12, 13
After Phase 2	1, 2, 3, 4	5, 7, 8, 9	6, 10, 11, 12	13, 14, 15, 16
After Phase 3	1, 2, 3, 4	5, 6, 7, 8	9, 10, 11, 12	13, 14, 15, 16

Parallel odd-even transposition sort



- how do we exchange elements?
 - If a process is not idle, we might try to implement the communication with a call to MPI_Send and a call to MPI_Recv

```
1 MPI_Send(my_keys, n/comm_sz, MPI_INT, partner, 0, comm);  
2 MPI_Recv(temp_keys, n/comm_sz, MPI_INT, partner, 0, comm, MPI STATUS IGNORE);  
_
```

- This might result in deadlock.
- MPI standard allows MPI_Send to behave in two different ways:
 - it can simply copy the message into an MPI-managed buffer and return,
 - or it can block until the matching call to MPI Recv starts
 - Furthermore, many implementations of MPI set a threshold at which the system switches from buffering to blocking. That is, messages that are relatively small will be buffered by MPI Send, but for larger messages, it will block.
- If the MPI_Send executed by each process blocks, no process will be able to start executing a call to MPI_Recv leading to deadlock

Parallel odd-even transposition sort



- A program that relies on MPI-provided buffering is said to be unsafe
 - Such a program may run without problems for various sets of input, but it may hang or crash with other sets. If we use `MPI_Send` and `MPI_Recv` in this way, our program will be unsafe, and it's likely that for small values of n the program will run without problems, while for larger values of n , it's likely that it will hang or crash
 - we can use an alternative to `MPI_Send` i.e. `MPI_Ssend`
 - `MPI_Ssend` is guaranteed to block until the matching receive starts

Parallel odd-even transposition sort



- How can we modify the communication in the parallel odd-even sort program so that it is safe?
 - Communication must be restructured. The most common cause of an unsafe program is multiple processes simultaneously first sending to each other and then receiving
 - need to restructure the communications so that some of the processes receive before sending

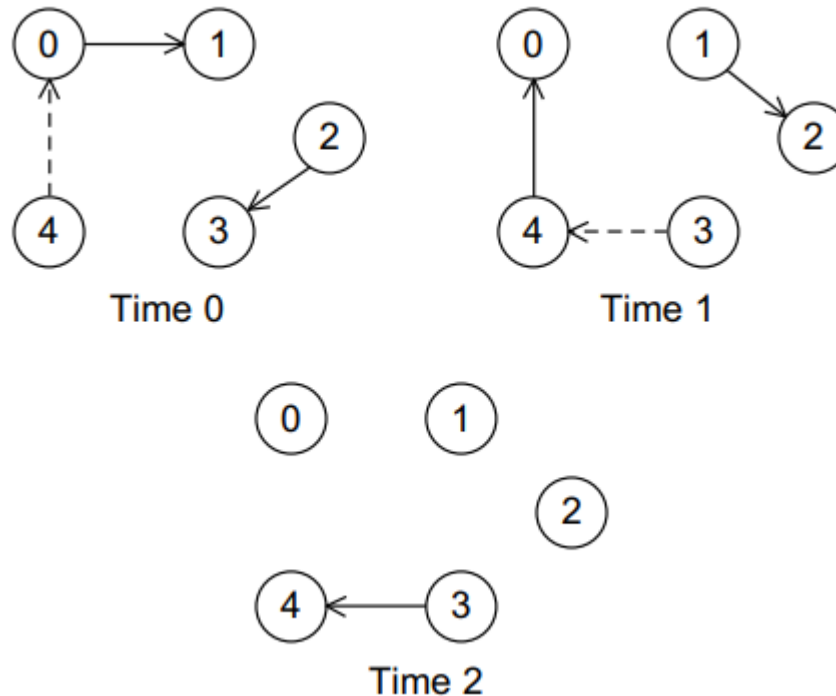
```
1  if (my_rank % 2 == 0) {  
2      MPI_Send(msg, size, MPI_INT, (my_rank+1) % comm_sz, 0, comm);  
3      MPI_Recv(new_msg, size, MPI_INT, (my_rank+comm_sz-1) % comm_sz,  
4          0, comm, MPI_STATUS_IGNORE, MPI_STATUS_IGNORE).  
5  } else {  
6      MPI_Recv(new_msg, size, MPI_INT, (my_rank+comm_sz-1) % comm_sz,  
7          0, comm, MPI_STATUS_IGNORE).  
8      MPI_Send(msg, size, MPI_INT, (my_rank+1) % comm_sz, 0, comm);  
9  }
```

- this will work if comm_sz is even. But when it is odd?

Parallel odd-even transposition sort



- For 5 processes, dashed arrows shows waiting communications though expected to be complete in the same phase.



Parallel odd-even transposition sort



- MPI provides an alternative to scheduling the communications ourselves— MPI_Sendrecv:

```
1 int MPI_Sendrecv(  
2     void* send_buf_p /* in */,  
3     int send_buf_size /* in */,  
4     MPI_Datatype send_buf_type /* in */,  
5     int dest /* in */,  
6     int send_tag /* in */,  
7     void* recv_buf_p /* out */,  
8     int recv_buf_size /* in */,  
9     MPI_Datatype recv_buf_type /* in */,  
10    int source /* in */,  
11    int recv_tag /* in */,  
12    MPI_Comm communicator /* in */,  
13    MPI_Status* status_p /* in */  
14 );
```

- This function carries out a blocking send and a receive in a single call. The dest and the source can be the same or different.
- MPI implementation schedules the communications so that the program won't hang or crash.

Parallel odd-even transposition sort

innovate

achieve

lead

```
306  * Purpose:      One iteration of Odd-even transposition sort
307  * In args:      local_n, phase, my_rank, p, comm
308  * In/out args:  local_A
309  * Scratch:      temp_B, temp_C
310  */
311  void Odd_even_iter(int local_A[], int temp_B[], int temp_C[],
312                    int local_n, int phase, int even_partner, int odd_partner,
313                    int my_rank, int p, MPI_Comm comm) {
314      MPI_Status status;
315
316      if (phase % 2 == 0) {
317          if (even_partner >= 0) {
318              MPI_Sendrecv(local_A, local_n, MPI_INT, even_partner, 0,
319                          temp_B, local_n, MPI_INT, even_partner, 0, comm,
320                          &status);
321              if (my_rank % 2 != 0)
322                  Merge_high(local_A, temp_B, temp_C, local_n);
323              else
324                  Merge_low(local_A, temp_B, temp_C, local_n);
325          }
326      } else { /* odd phase */
327          if (odd_partner >= 0) {
328              MPI_Sendrecv(local_A, local_n, MPI_INT, odd_partner, 0,
329                          temp_B, local_n, MPI_INT, odd_partner, 0, comm,
330                          &status);
331              if (my_rank % 2 != 0)
332                  Merge_low(local_A, temp_B, temp_C, local_n);
333              else
334                  Merge_high(local_A, temp_B, temp_C, local_n);
335          }
336      }
```

References

- <https://computing.llnl.gov/tutorials/mpi/>
- “An introduction to Parallel Programming”, Peter S. Pacheco (Chapter 3)



BITS Pilani
Pilani Campus



Thank You