MPI, continued

Outline

- Finish MPI discussion
 - Review blocking and non-blocking communication
 - Introduce one-sided communication
- Sources for this lecture:
 - http://mpi.deino.net/mpi_functions/

Today's MPI Focus - Communication Primitives

- Collective communication
 - Reductions, Broadcast, Scatter, Gather
- Blocking communication
 - Overhead
 - Deadlock?
- Non-blocking
- One-sided communication

Quick MPI Review

- Six most common MPI Commands (aka, Six Command MPI)
 - MPI_Init
 - MPI_Finalize
 - MPI_Comm_size
 - MPI_Comm_rank
 - MPI_Send
 - MPI_Recv
- Send and Receive refer to "point-to-point" communication
- Last time we also showed collective communication
 - Reduce

Deadlock?

int a[10], b[10], myrank; MPI_Status status; ... MPI_Comm_rank(MPI_COMM_WORLD, &myrank);

```
if (myrank == 0) {
    MPI_Send(a, 10, MPI_INT, 1, 1, MPI_COMM_WORLD);
    MPI_Send(b, 10, MPI_INT, 1, 2, MPI_COMM_WORLD); }
```

```
else if (myrank == 1) {
	MPI_Recv(b, 10, MPI_INT, 0, 2, MPI_COMM_WORLD);
	MPI_Recv(a, 10, MPI_INT, 0, 1, MPI_COMM_WORLD);
}
```

Deadlock?

Consider the following piece of code:

```
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_Send(a, 10, MPI_INT, (myrank+1)%npes, 1,
MPI_COMM_WORLD);
MPI_Recv(b, 10, MPI_INT, (myrank-1+npes)%npes, 1,
MPI_COMM_WORLD); ...
```

Global Sum using Butterfly



Example: MPI Code for Butterfly Allreduce

• This is a class time exercise!

MPI Allreduce

MPI_Allreduce

Combines values from all processes and distributes the result back to all processes

Synopsis

Input Parameters

sendbuf
starting address of send buffer (choice)
count
number of elements in send buffer (integer)
datatype
data type of elements of send buffer (handle)
op
operation (handle)
comm

communicator (handle)

Output Parameters

recvbuf

starting address of receive buffer (choice)

Global Sum using Allreduce

• Write the same program!

Matrix-Vector Multiplication



```
void Mat_vect_mult(
1
        double A[] /* in */,
2
3
        double x[] /* in */,
4
      double y[] /* out */,
5
     int m /* in */,
     int n /* in */) {
6
7
    int i, j;
8
9
      for (i = 0; i < m; i++) {
        y[i] = 0.0;
10
        for (j = 0; j < n; j++)
11
           y[i] += A[i*n+j]*x[j];
12
13
     /* Mat_vect_mult */
14
   }
```

MPI Version

```
void Mat_vect_mult(
1
        double local_A[] /* in */,
2
3
      double local_x[] /* in */.
      double local_y[] /* out */,
4
5
      int local_m /* in */,
      int n /* in */,
6
     int local_n /* in */.
7
      MPI_Comm comm /* in */) {
8
9
      double* x:
     int local_i, j;
10
11
      int local_ok = 1;
12
      x = malloc(n*sizeof(double)):
13
      MPI_Allgather(local_x, local_n, MPI_DOUBLE,
14
15
           x, local_n, MPI_DOUBLE, comm);
16
17
      for (local_i = 0; local_i < local_m; local_i++) {
18
        local_y[local_i] = 0.0;
        for (j = 0; j < n; j++)
19
           local_y[local_i] += local_A[local_i*n+j]*x[j]:
20
21
22
  free(x):
23 } /* Mat_vect_mult */
```

Homework

 Include the parallel and serial algorithms of Matrix Mupltiply and compare their timing on a parallel machine. Use a 1000x1000 double matrix and a 1000 element vector.

More difficult p2p example: 2D relaxation

Replaces each interior value by the average of its four nearest neighbors.

```
Sequential code:
for (i=1; i<n-1; i++)
for (j=1; j<n-1; j++)
b[i,j] = (a[i-1][j]+a[i][j-1]+a[i+1][j]+a[i][j+1])/4.0;
```

1.0	0.0	0.0	0.0	0.0	0.0
1.0	0.0	0.0	0.0	0.0	0.0
1.0	0.0	0.0	0.0	0.0	0.0
1.0	0.0	0.0	0.0	0.0	0.0
1.0	0.0	0.0	0.0	0.0	0.0
1.0	0.0	0.0	0.0	0.0	0.0
Interior value					

Boundary value

MPI code, main loop of 2D SOR computation

```
#define Top
                       0
  1
  2
     #define Left
                       0
  3
     #define Right
                      (Cols-1)
  4
     #define Bottom (Rows-1)
  5
     #define NorthPE(i)
  6
                              ((i)-Cols)
     #define SouthPE(i)
  7
                             ((i)+Cols)
  8
     #define EastPE(i)
                              ((i)+1)
  9
     #define WestPE(i)
                               ((i)-1)
. . .
101
     do
     { /*
102
       * Send data to four neighbors
103
104
       */
       if(row !=Top)
105
                                                      /* Send North */
106
       {
107
         MPI_Send(&val[1][1], Width-2, MPI_FLOAT,
108
                   NorthPE(myID), tag, MPI COMM WORLD);
109
       }
110
111
       if(col !=Right)
                                                      /* Send East */
112
       {
         for(i=1; i<Height-1; i++)</pre>
113
114
          {
115
            buffer[i-1]=val[i][Width-2];
116
         }
         MPI Send(buffer, Height-2, MPI FLOAT,
117
118
                   EastPE(myID), tag, MPI COMM WORLD);
119
       }
120
121
                                                      /* Send South */
       if(row !=Bottom)
122
       {
123
         MPI Send(&val[Height-2][1], Width-2, MPI FLOAT,
124
                   SouthPE(myID), tag, MPI COMM WORLD);
125
       }
126
127
       if(col !=Left)
                                                      /* Send West */
128
       {
```

MPI code, main loop of 2D SOR computation, cont.

```
129
          for(i=1; i<Height-1; i++)</pre>
130
          {
131
            buffer[i-1]=val[i][1];
132
          }
133
          MPI Send(buffer, Height-2, MPI FLOAT,
134
                    WestPE(myID), tag, MPI COMM WORLD);
135
        }
136
        /*
137
138
         * Receive messages
139
         */
140
        if(row !=Top)
                                                 /* Receive from North */
141
        {
142
          MPI Recv(&val[0][1], Width-2, MPI FLOAT,
                    NorthPE(myID), tag, MPI COMM WORLD, &status);
143
144
        }
145
146
        if(col !=Right)
                                                 /* Receive from East */
147
        {
148
          MPI Recv(&buffer, Height-2, MPI FLOAT,
149
                    EastPE(myID), tag, MPI COMM WORLD, &status);
150
          for(i=1; i<Height-1; i++)</pre>
151
          {
            val[i][Width-1]=buffer[i-1];
152
153
          }
154
        }
155
156
                                                 /* Receive from South */
        if(row !=Bottom)
157
        {
158
          MPI Recv(&val[0][Height-1], Width-2, MPI FLOAT,
159
                    SouthPE(myID), tag, MPI COMM WORLD, &status);
160
        }
161
162
        if(col !=Left)
                                                 /* Receive from West */
163
        {
164
          MPI Recv(&buffer, Height-2, MPI FLOAT,
165
                    WestPE(myID), tag, MPI COMM WORLD, &status);
          for(i=1; i<Height-1; i++)</pre>
166
167
          {
168
            val[i][0]=buffer[i-1];
169
          }
170
        }
171
172
        delta=0.0; /* Calculate average, delta for all points */
173
        for(i=1; i<Height-1; i++)</pre>
174
        {
175
          for(j=1; j<Width-1; j++)</pre>
176
          {
```

```
177
            average=(val[i-1][j]+val[i][j+1]+
178
                      val[i+1][j]+val[i][j-1])/4;
179
            delta=Max(delta, Abs(average-val[i][j]));
180
            new[i][j]=average;
181
         }
182
       }
183
184
       /* Find maximum diff */
185
       MPI Reduce(&delta, &globalDelta, 1, MPI FLOAT, MPI MIN,
186
                   RootProcess, MPI COMM WORLD);
187
       Swap(val, new);
188
     } while(globalDelta < THRESHOLD);</pre>
```

MPI Scatter()

MPI_Scatter()

```
int MPI_Scatter(
  void *sendbuffer,
  int sendcount,
  MPI_Datatype sendtype,
  int destbuffer,
  int destcount,
  MPI_Datatype desttype,
  int root,
  MPI_Comm *comm
```

- // Scatter routine
- // Address of the data to send
- // Number of data elements to send
- // Type of data elements to send
- // Address of buffer to receive data
- // Number of data elements to receive
- // Type of data elements to receive
- // Rank of the root process
- // An MPI communicator

```
);
```

Arguments:

- The first three arguments specify the address, size, and type of the data elements to send to each process. These arguments only have meaning for the root process.
- The second three arguments specify the address, size, and type of the data elements for each receiving process. The size and type of the sending data and the receiving data may differ as a means of converting data types.
- The seventh argument specifies the root process that is the source of the data.
- The eighth argument specifies the MPI communicator to use.

Notes:

This routine distributes data from the root process to all other processes, including the root. A more sophisticated version of the routine, MPI_Scatterv(), allows the root process to send different amounts of data to the various processes. Details can be found in the MPI standard.

Return value:

An MPI error code.

MPI_Scatter Example

```
void Read_vector(
 1
         double local_a[] /* out */.
2
3
         int local_n /* in */,
         int
               n /* in */.
4
         char vec_name[] /* in */,
5
      int my_rank /* in */,
6
         MPI_Comm comm /* in */) {
7
8
9
      double* a = NULL:
      int i:
10
11
      if (my_rank == 0) {
12
         a = malloc(n*sizeof(double)):
13
         printf("Enter the vector %s\n", vec_name);
14
         for (i = 0; i < n; i++)
15
            scanf("%]f". &a[i]):
16
17
         MPI_Scatter(a, local_n, MPI_DOUBLE, local_a, local_n,
              MPI_DOUBLE, 0, comm);
18
19
        free(a):
20
      } else {
21
         MPI_Scatter(a, local_n, MPI_DOUBLE, local_a, local_n,
              MPI_DOUBLE, 0, comm);
22
23
      }
      /* Read_vector */
24
```

Distribute Data from input using a scatter operation

```
16
    length per process=length/size;
17
    myArray=(int *) malloc(length per process*sizeof(int));
18
    array=(int *) malloc(length*sizeof(int));
19
20
21
    /* Read the data, distribute it among the various processes */
22
    if(myID==RootProcess)
23
    {
24
      if((fp=fopen(*argv, "r"))==NULL)
25
      {
26
        printf("fopen failed on %s\n", filename);
27
        exit(0);
28
      }
29
      fscanf(fp,"%d", &length); /* read input size */
30
31
      for(i=0; i<length-1; i++)  /* read entire input file */</pre>
32
      {
33
        fscanf(fp,"%d", myArray+i);
34
      }
35
    }
36
37
    MPI Scatter(Array, length per process, MPI INT,
38
                myArray, length per process, MPI INT,
39
                RootProcess, MPI COMM WORLD);
```

MPI Gather

MPI Gather

Gathers together values from a group of processes

Synopsis

int MPI Gather(const void *sendbuf, int sendcount, MPI Datatype sendtype, void *recvbuf, int recvcount, MPI Datatype recvtype, int root, MPI Comm comm)

Input Parameters

sendbuf starting address of send buffer (choice) sendcount number of elements in send buffer (integer) sendtype data type of send buffer elements (handle) recvcount recvtype root rank of receiving process (integer) comm

number of elements for any single receive (integer, significant only at root)

data type of recv buffer elements (significant only at root) (handle)

communicator (handle)

Output Parameters

recvbuf

address of receive buffer (choice, significant only at root)

MPI_Gather Example

```
void Print_vector(
1
         double local_b[] /* in */.
2
3
         int local_n /* in */,
4
         int n /* in */.
        char title[] /* in */,
5
6
      int my_rank /* in */.
7
         MPI_Comm comm /* in */) {
8
9
      double* b = NULL:
10
      int i:
11
      if (my_rank == 0) {
12
         b = malloc(n*sizeof(double)):
13
14
         MPI_Gather(local_b, local_n, MPI_DOUBLE, b, local_n,
               MPI_DOUBLE. O. comm):
15
         printf("%s\n", title):
16
         for (i = 0; i < n; i++)
17
            printf("%f ". b[i]):
18
19
         printf("\n");
         free(b):
20
      } else {
21
         MPI_Gather(local_b, local_n, MPI_DOUBLE. b. local_n.
22
               MPI_DOUBLE. O. comm):
23
24
      /* Print_vector */
25
   }
```

Non-Blocking Communication

- The programmer must ensure semantics of the send and receive.
- This class of non-blocking protocols returns from the send or receive operation before it is semantically safe to do so.
- Non-blocking operations are generally accompanied by a check-status operation.
- When used correctly, these primitives are capable of overlapping communication overheads with useful computations.
- Message passing libraries typically provide both blocking and non-blocking primitives.

Non-Blocking Communication

- To overlap communication with computation, MPI provides a pair of functions for performing nonblocking send and receive operations ("I" stands for "Immediate"):
 - 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 waits for the operation to complete. int MPI_Wait(MPI_Request *request, MPI_Status *status) CS4230

Improving SOR with Non-Blocking Communication

```
if (row != Top) {
    MPI_Isend(&val[1][1], Width-
    2,MPI_FLOAT,NorthPE(myID),tag,MPI_COMM_WORLD, &requests[0]);
}
// analogous for South, East and West
...
if (row!=Top) {
    MPI_Irecv(&val[0][1],Width-2,MPI_FLOAT,NorthPE(myID),
        tag, MPI_COMM_WORLD, &requests[4]);
}
```

// Perform interior computation on local data

//Now wait for Recvs to complete

MPI_Waitall(8,requests,status);

//Then, perform computation on boundaries

Taking Time in MPI

• MPI_Wtime

```
start = MPI_Wtime():
/* Code to be timed */
finish = MPI_Wtime():
printf("Proc %d > Elapsed time = %e seconds\n"
      my_rank, finish-start);
double local_start, local_finish, local_elapsed, elapsed;
MPI_Barrier(comm);
local_start = MPI_Wtime():
/* Code to be timed */
. . .
local_finish = MPI_Wtime():
local_elapsed = local_finish - local_start;
MPI_Reduce(&local_elapsed, &elapsed, 1, MPI_DOUBLE,
   MPI_MAX, 0, comm);
if (my_rank == 0)
   printf("Elapsed time = %e seconds\n", elapsed);
```

One-Sided Communication



MPI Constructs supporting One-Sided Communication (RMA)

- MPI_Win_create exposes local memory to RMA operation by other processes in a communicator
 - Collective operation
 - Creates window object
- MPI_Win_free deallocates window object
- MPI_Put moves data from local memory to remote memory
- MPI_Get retrieves data from remote memory into local memory
- MPI_Accumulate updates remote memory using local values

MPI_Put and MPI_Get

int MPI_Put(void *origin_addr, int origin_count, MPI_Datatype origin_datatype, int target_rank, MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype, MPI_Win win);

int MPI_Get(void *origin_addr, int origin_count, MPI_Datatype origin_datatype, int target_rank, MPI_Aint target_disp, int target_count, MPI_Datatype target_datatype, MPI_Win win);

Specify address, count, datatype for origin and target, rank for target and MPI_win for 1-sided communication.

MPI Critique (Snyder)

- Message passing is a very simple model
- Extremely low level; heavy weight
 - Expense comes from ${\mbox{\sc h}}$ and lots of local code
 - Communication code is often more than half
 - Tough to make adaptable and flexible
 - Tough to get right and know it
 - Tough to make perform in some (Snyder says most) cases
- Programming model of choice for scalability
- Widespread adoption due to portability, although not completely true in practice