

CS4961 Parallel Programming

Lecture 16: Introduction to Message Passing

Mary Hall
October 29, 2009

10/29/2009

CS4961

1

Administrative

- Homework assignment 3 will be posted today (after class)
 - Due, Thursday, November 5 before class
 - Use the "handin" program on the CADE machines
 - Use the following command:


```
"handin cs4961 hw3 <gzipped tar file>"
```
- NEW: VTUNE PORTION IS EXTRA CREDIT!
- Mailing list set up: cs4961@list.eng.utah.edu
 - Next week we'll start discussing final project
 - Optional CUDA or MPI programming assignment part of this

10/27/2009

CS4961

2



A Few Words About Final Project

- Purpose:
 - A chance to dig in deeper into a parallel programming model and explore concepts.
 - Present results to work on communication of technical ideas
- Write a non-trivial parallel program that combines two parallel programming languages/models. In some cases, just do two separate implementations.
 - OpenMP + SSE-3
 - OpenMP + CUDA (but need to do this in separate parts of the code)
 - TBB + SSE-3
 - MPI + OpenMP
 - MPI + SSE-3
 - MPI + CUDA
- Present results in a poster session on the last day of class

10/29/2009

CS4961

3



Example Projects

- Look in the textbook or on-line
 - Recall Red/Blue from Ch. 4
 - Implement in MPI (+ SSE-3)
 - Implement main computation in CUDA
 - Algorithms from Ch. 5
 - SOR from Ch. 7
 - CUDA implementation?
 - FFT from Ch. 10
 - Jacobi from Ch. 10
 - Graph algorithms
 - Image and signal processing algorithms
 - Other domains...

10/29/2009

CS4961

4



Today's Lecture

- Message Passing, largely for distributed memory
- Message Passing Interface (MPI): a Local View language
- Sources for this lecture
 - Larry Snyder, <http://www.cs.washington.edu/education/courses/524/08wi/>
 - Online MPI tutorial <http://www-unix.mcs.anl.gov/mpi/tutorial/gropp/Talk.html>

10/29/2009

CS4961

5



Message Passing

- Message passing is the principle alternative to shared memory parallel programming
 - Based on Single Program, Multiple Data (SPMD)
 - Model with send() and recv() primitives
 - Message passing is universal, but low-level
 - More even than threading, message passing is locally focused -- what does each processor do?
 - Isolation of separate address spaces
 - + no data races
 - + forces programmer to think about locality, so good for performance
 - + architecture model exposed, so good for performance
 - low level
 - complexity
 - code growth!

10/29/2009

CS4961

6



Message Passing Libraries (1)

- Many "message passing libraries" were once available
 - Chameleon, from ANL.
 - CMMD, from Thinking Machines.
 - Express, commercial.
 - MPL, native library on IBM SP-2.
 - NX, native library on Intel Paragon.
 - Zipcode, from LLL.
 - PVM, Parallel Virtual Machine, public, from ORNL/UTK.
 - Others...
 - MPI, Message Passing Interface, now the industry standard.
- Need standards to write portable code.



Message Passing Libraries (2)

- All communication, synchronization require subroutine calls
 - No shared variables
 - Program run on a single processor just like any uniprocessor program, except for calls to message passing library
- Subroutines for
 - Communication
 - Pairwise or point-to-point: Send and Receive
 - Collectives all processor get together to
 - Move data: Broadcast, Scatter/gather
 - Compute and move: sum, product, max, ... of data on many processors
 - Synchronization
 - Barrier
 - No locks because there are no shared variables to protect
 - Queries
 - How many processes? Which one am I? Any messages waiting?



Novel Features of MPI

- Communicators encapsulate communication spaces for library safety
- Datatypes reduce copying costs and permit heterogeneity
- Multiple communication modes allow precise buffer management
- Extensive collective operations for scalable global communication
- Process topologies permit efficient process placement, user views of process layout
- Profiling interface encourages portable tools

Slide source: Bill Gropp



MPI References

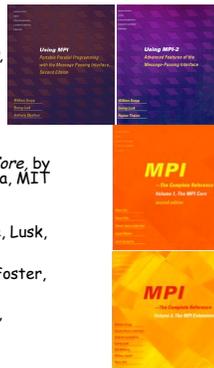
- The Standard itself:
 - at <http://www.mpi-forum.org>
 - All MPI official releases, in both postscript and HTML
- Other information on Web:
 - at <http://www.mcs.anl.gov/mpi>
 - pointers to lots of stuff, including other talks and tutorials, a FAQ, other MPI pages

Slide source: Bill Gropp



Books on MPI

- *Using MPI: Portable Parallel Programming with the Message-Passing Interface (2nd edition)*, by Gropp, Lusk, and Skjellum, MIT Press, 1999.
- *Using MPI-2: Portable Parallel Programming with the Message-Passing Interface*, by Gropp, Lusk, and Thakur, MIT Press, 1999.
- *MPI: The Complete Reference - Vol 1 The MPI Core*, by Snir, Otto, Huss-Lederman, Walker, and Dongarra, MIT Press, 1998.
- *MPI: The Complete Reference - Vol 2 The MPI Extensions*, by Gropp, Huss-Lederman, Lumsdaine, Lusk, Nitzberg, Saphir, and Snir, MIT Press, 1998.
- *Designing and Building Parallel Programs*, by Ian Foster, Addison-Wesley, 1995.
- *Parallel Programming with MPI*, by Peter Pacheco, Morgan-Kaufmann, 1997.



Slide source: Bill Gropp



Working through an example

- We'll write some message-passing pseudo code for Count3 (from Lecture 4)

```

1 int array[length];
2 int i;
3 int total;
4 forall(j in(0..p-1))
5 {
6   int size=mySize(array,0);
7   int myData[size]=localize(array[]);
8   int i, priv_count=0;
9   for(i=0; i<size; i++)
10  {
11    if(myData[i]==3)
12    {
13      priv_count++;
14    }
15  }
16  total +=priv_count;
17 }

```

The data is global
Number of desired threads
Result of computation, grand total

Figure size of local part of global data

Associate my part of global data with local variable
Local accumulation

compute grand total

10/29/2009

CS4961

12



Finding Out About the Environment

- Two important questions that arise early in a parallel program are:
 - How many processes are participating in this computation?
 - Which one am I?
- MPI provides functions to answer these questions:
 - MPI_Comm_size reports the number of processes.
 - MPI_Comm_rank reports the *rank*, a number between 0 and size-1, identifying the calling process

Slide source: Bill Gropp



Hello (C)

```
#include "mpi.h"
#include <stdio.h>

int main( int argc, char *argv[] )
{
    int rank, size;
    MPI_Init( &argc, &argv );
    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
    MPI_Comm_size( MPI_COMM_WORLD, &size );
    printf( "I am %d of %d\n", rank, size );
    MPI_Finalize();
    return 0;
}
```

Slide source: Bill Gropp



Hello (Fortran)

```
program main
include 'mpif.h'
integer ierr, rank, size

call MPI_INIT( ierr )
call MPI_COMM_RANK( MPI_COMM_WORLD, rank, ierr )
call MPI_COMM_SIZE( MPI_COMM_WORLD, size, ierr )
print *, 'I am ', rank, ' of ', size
call MPI_FINALIZE( ierr )
end
```

Slide source: Bill Gropp



Hello (C++)

```
#include "mpi.h"
#include <iostream>

int main( int argc, char *argv[] )
{
    int rank, size;
    MPI::Init(argc, argv);
    rank = MPI::COMM_WORLD.Get_rank();
    size = MPI::COMM_WORLD.Get_size();
    std::cout << "I am " << rank << " of " << size <<
        "\n";
    MPI::Finalize();
    return 0;
}
```

Slide source: Bill Gropp



Notes on Hello World

- All MPI programs begin with `MPI_Init` and end with `MPI_Finalize`
- `MPI_COMM_WORLD` is defined by `mpi.h` (in C) or `mpif.h` (in Fortran) and designates all processes in the MPI "job"
- Each statement executes independently in each process
 - including the `printf/print` statements
- I/O not part of MPI-1 but is in MPI-2
 - print and write to standard output or error not part of either MPI-1 or MPI-2
 - output order is undefined (may be interleaved by character, line, or blocks of characters),
- The MPI-1 Standard does not specify how to run an MPI program, but many implementations provide


```
mpirun -np 4 a.out
```

Slide source: Bill Gropp



MPI Basic Send/Receive

- We need to fill in the details in



- Things that need specifying:
 - How will "data" be described?
 - How will processes be identified?
 - How will the receiver recognize/screen messages?
 - What will it mean for these operations to complete?

Slide source: Bill Gropp



Some Basic Concepts

- Processes can be collected into groups
- Each message is sent in a context, and must be received in the same context
 - Provides necessary support for libraries
- A group and context together form a communicator
- A process is identified by its rank in the group associated with a communicator
- There is a default communicator whose group contains all initial processes, called `MPI_COMM_WORLD`

Slide source: Bill Gropp



MPI Datatypes

- The data in a message to send or receive is described by a triple (address, count, datatype), where
- An MPI datatype is recursively defined as:
 - predefined, corresponding to a data type from the language (e.g., `MPI_INT`, `MPI_DOUBLE`)
 - a contiguous array of MPI datatypes
 - a strided block of datatypes
 - an indexed array of blocks of datatypes
 - an arbitrary structure of datatypes
- There are MPI functions to construct custom datatypes, in particular ones for subarrays

Slide source: Bill Gropp



MPI Tags

- Messages are sent with an accompanying user-defined integer tag, to assist the receiving process in identifying the message
- Messages can be screened at the receiving end by specifying a specific tag, or not screened by specifying `MPI_ANY_TAG` as the tag in a receive
- Some non-MPI message-passing systems have called tags "message types". MPI calls them tags to avoid confusion with datatypes

Slide source: Bill Gropp



MPI Basic (Blocking) Send



MPI_Send(A, 10, MPI_DOUBLE, 1, ...)

MPI_Recv(B, 20, MPI_DOUBLE, 0, ...)

`MPI_SEND(start, count, datatype, dest, tag, comm)`

- The message buffer is described by (`start`, `count`, `datatype`).
- The target process is specified by `dest`, which is the rank of the target process in the communicator specified by `comm`.
- When this function returns, the data has been delivered to the system and the buffer can be reused. The message may not have been received by the target process.

Slide source: Bill Gropp



MPI Basic (Blocking) Receive



MPI_Send(A, 10, MPI_DOUBLE, 1, ...)

MPI_Recv(B, 20, MPI_DOUBLE, 0, ...)

`MPI_RECV(start, count, datatype, source, tag, comm, status)`

- Waits until a matching (both `source` and `tag`) message is received from the system, and the buffer can be used
- `source` is rank in communicator specified by `comm`, or `MPI_ANY_SOURCE`
- `tag` is a tag to be matched on or `MPI_ANY_TAG`
- receiving fewer than `count` occurrences of `datatype` is OK, but receiving more is an error
- `status` contains further information (e.g. size of message)

Slide source: Bill Gropp



A Simple MPI Program

```
#include "mpi.h"
#include <stdio.h>
int main( int argc, char *argv[] )
{
    int rank, buf;
    MPI_Status status;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank( MPI_COMM_WORLD, &rank );

    /* Process 0 sends and Process 1 receives */
    if (rank == 0) {
        buf = 123456;
        MPI_Send( &buf, 1, MPI_INT, 1, 0, MPI_COMM_WORLD );
    }
    else if (rank == 1) {
        MPI_Recv( &buf, 1, MPI_INT, 0, 0, MPI_COMM_WORLD,
                &status );
        printf( "Received %d\n", buf );
    }

    MPI_Finalize();
    return 0;
}
```

Slide source: Bill Gropp



Figure 7.1 An MPI solution to the Count 3s problem.

```

1 #include <stdio.h>
2 #include "mpi.h"
3 #include "globals.h"
4
5 int main(argc, argv)
6 int argc;
7 char **argv;
8 {
9     int myID, value, numProcs;
10    MPI_Status status;
11
12    MPI_Init(&argc, &argv);
13    MPI_Comm_size(MPI_COMM_WORLD, &numProcs);
14    MPI_Comm_rank(MPI_COMM_WORLD, &myID);
15
16    length_per_process=length/numProcs;
17    myArray=(int *) malloc(length_per_process*sizeof(int));
18
19    /* Read the data, distribute it among the various processes */
20    if(myID==RootProcess)
21    {
22        if((fp=fopen(argv, "r"))==NULL )
23        {
24            printf("fopen failed on %s\n", filename);
25            exit(0);
26        }
27        fscanf(fp, "%d", &length); /* read input size */
28
29        for(p=0; p<numProcs-1; p++) /* read data on behalf of each */
30        { /* of the other processes */
31            for(i=0; i<length_per_process; i++)
32            {
33                fscanf(fp, "%d", &myArray[i]);
34            }
35            MPI_Send(myArray, length_per_process, MPI_INT, p+1,
36                    tag, MPI_COMM_WORLD);

```

7-25



Figure 7.1 An MPI solution to the Count 3s problem. (cont.)

```

37     }
38
39     for(i=0; i<length_per_process; i++) /* Now read my data */
40     {
41         fscanf(fp, "%d", &myArray[i]);
42     }
43
44     else
45     {
46         MPI_Recv(myArray, length_per_process, MPI_INT, RootProcess,
47                 tag, MPI_COMM_WORLD, &status);
48     }
49
50     /* Do the actual work */
51     for(i=0; i<length_per_process; i++)
52     {
53         if(myArray[i]==3)
54         {
55             myCount++; /* Update local count */
56         }
57     }
58
59     MPI_Reduce(&myCount, &globalCount, 1, MPI_INT, MPI_SUM,
60              RootProcess, MPI_COMM_WORLD);
61
62     if(myID==RootProcess)
63     {
64         printf("Number of 3's: %d\n", globalCount);
65     }
66     MPI_Finalize();
67     return 0;
68 }

```

7-26



Code Spec 7.8 MPI Scatter().

```

MPI_Scatter()
int MPI_Scatter(
void *sendbuffer,          // Scatter routine
int sendcount,            // Address of the data to send
MPI_Datatype sendtype,    // Number of data elements to send
int destbuffer,          // Type of data elements to send
int destcount,           // Address of buffer to receive data
MPI_Datatype desttype,   // Number of data elements to receive
int root,                // Type of data elements to receive
MPI_Comm *comm           // Rank of the root process
);
MPI_Comm *comm           // 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.

(continued)

7-27



Code Spec 7.8 MPI Scatter(). (cont.)

- 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.

7-28



Figure 7.2

Replacement code (for lines 16–48 of Figure 7.1) to distribute data using a scatter operation.

```

16 length_per_process=length/size;
17 myArray=(int *) malloc(length_per_process*sizeof(int));
18
19 array=(int *) malloc(length*sizeof(int));
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 */
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);

```

7-29



Other Basic Features of MPI

- **MPI_Gather**
 - Analogous to MPI_Scatter
- **Scans and reductions**
- **Groups, communicators, tags**
 - Mechanisms for identifying which processes participate in a communication
- **MPI_Bcast**
 - Broadcast to all other processes in a "group"

10/29/2009

CS4961

30



Figure 7.4 Example of collective communication within a group.

```

1 int numCols; /* initialized elsewhere */
2
3 void broadcast_example()
4 {
5     int **ranks; /* the ranks that belong to each group */
6     int myRank; /* row number of this process */
7     int rowNumber; /* row number of this process */
8     int random; /* value that we would like to broadcast */
9     rowNumber=myRank/numCols;
10    MPI_Group globalGroup, newGroup;
11    MPI_Comm rowComm(numCols);
12
13    /* initialize ranks[] array */
14    ranks[0]={0,1,2,3}; /* not legal C */
15    ranks[1]={4,5,6,7};
16    ranks[2]={8,9,10,11};
17    ranks[3]={12,13,14,15};
18
19    /* Extract the original group handle */
20    MPI_Comm_group(MPI_COMM_WORLD, &globalGroup);
21
22    /* Define the new group */
23    MPI_Group_incl(globalGroup, P/numCols, ranks[rowNumber], &newGroup);
24
25    /* Create new communicator */
26    MPI_Comm_create(MPI_COMM_WORLD, newGroup, &newComm);
27
28    random=rand();
29
30    /* Broadcast 'random' across rows */
31    MPI_Bcast(&random, 1, MPI_INT, rowNumber*numCols, newComm);
32 }

```

7-31



Figure 7.5 A 2D relaxation replaces—on each iteration—all interior values by the average of their four nearest neighbors.

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

7-32



Figure 7.6 MPI code for the main loop of the 2D SOR computation.

```

1 #define Top 0
2 #define Left 0
3 #define Right (Cols-1)
4 #define Bottom (Rows-1)
5
6 #define NorthPE(i) ((i)-Cols)
7 #define SouthPE(i) ((i)+Cols)
8 #define EastPE(i) ((i)+1)
9 #define WestPE(i) ((i)-1)
10
11 do
12 {
13     /* Send data to four neighbors
14     */
15     if(row !=Top) /* Send North */
16     {
17         MPI_Send(val[i][j], Width-2, MPI_FLOAT,
18             NorthPE(myID), tag, MPI_COMM_WORLD);
19     }
20     if(col !=Right) /* Send East */
21     {
22         for(i=1; i<Height-1; i++)
23         {
24             buffer[i-1]=val[i][Width-2];
25         }
26         MPI_Send(buffer, Height-2, MPI_FLOAT,
27             EastPE(myID), tag, MPI_COMM_WORLD);
28     }
29     if(row !=Bottom) /* Send South */
30     {
31         MPI_Send(val[Height-2][j], Width-2, MPI_FLOAT,
32             SouthPE(myID), tag, MPI_COMM_WORLD);
33     }
34     if(col !=Left) /* Send West */
35     {

```

7-33



Figure 7.6 MPI code for the main loop of the 2D SOR computation. (cont.)

```

36     for(i=1; i<Width-1; i++)
37     {
38         buffer[i-1]=val[i][j];
39     }
40     MPI_Send(buffer, Height-2, MPI_FLOAT,
41         WestPE(myID), tag, MPI_COMM_WORLD);
42 }
43
44 /*
45 * Receive messages
46 */
47 if(row !=Top) /* Receive from North */
48 {
49     MPI_Recv(recv[i][j], Width-2, MPI_FLOAT,
50         NorthPE(myID), tag, MPI_COMM_WORLD, &status);
51 }
52 if(col !=Right) /* Receive from East */
53 {
54     MPI_Recv(buffer, Height-2, MPI_FLOAT,
55         EastPE(myID), tag, MPI_COMM_WORLD, &status);
56     for(i=1; i<Height-1; i++)
57     {
58         val[i][Width-1]=buffer[i-1];
59     }
60 }
61 if(row !=Bottom) /* Receive from South */
62 {
63     MPI_Recv(recv[Height-1][j], Width-2, MPI_FLOAT,
64         SouthPE(myID), tag, MPI_COMM_WORLD, &status);
65 }
66 if(col !=Left) /* Receive from West */
67 {
68     MPI_Recv(buffer, Height-2, MPI_FLOAT,
69         WestPE(myID), tag, MPI_COMM_WORLD, &status);
70     for(i=1; i<Height-1; i++)
71     {
72         val[i][0]=buffer[i-1];
73     }
74 }
75
76 /* Calculate average, delta for all points */
77 for(i=1; i<Height-1; i++)
78 {
79     for(j=1; j<Width-1; j++)

```

7-34



Figure 7.6 MPI code for the main loop of the 2D SOR computation. (cont.)

```

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);

```

7-35



MPI Critique (Snyder)

- Message passing is a very simple model
- Extremely low level; heavy weight
 - Expense comes from λ 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

10/29/2009

CS4961

36

