

CS4230 Parallel Programming

Lecture 13: Introduction to Message Passing

Mary Hall
October 23, 2012

10/23/2012

CS4230

1

Administrative

- Preview of next programming assignment
 - due 11:59PM, Friday, November 2
 - SVD contains several reductions
 - We will strip out the Jacobi rotations. Your mission is to implement just the reductions in MPI, using point-to-point communication and then collective communication
- Subsequent assignment
 - Scalable MPI implementation of SVD, due Friday, Nov. 9

10/23/2012

CS4230

2



Today's Lecture

- Message Passing, largely for distributed memory
- Message Passing Interface (MPI):
 - The most commonly-used distributed-memory programming language for large-scale computation
- Chapter 3 in textbook
- Sources for this lecture
 - Textbook slides
 - Online MPI tutorial
<http://www-unix.mcs.anl.gov/mpi/tutorial/gropp/talk.html>

10/23/2012

CS4230

3



Recall from L3: Two main classes of parallel architecture organizations

- Shared memory multiprocessor architectures
 - A collection of autonomous processors connected to a memory system.
 - Supports a global address space where each processor can access each memory location.
- Distributed memory architectures
 - A collection of autonomous systems connected by an interconnect.
 - Each system has its own distinct address space, and processors must explicitly communicate to share data.
 - Clusters of PCs connected by commodity interconnect is the most common example.

08/28/2012

CS4230

4



Message Passing and MPI

- Message passing is the predominant programming model for supercomputers and clusters
 - Portable
 - Low-level, but universal and matches earlier hardware execution model
- What it is
 - A library used within conventional sequential languages (Fortran, C, C++)
 - Based on Single Program, Multiple Data (SPMD)
 - Isolation of separate address spaces
 - + no data races, but communication errors possible
 - + exposes execution model and forces programmer to think about locality, both good for performance
 - Complexity and code growth!

Like OpenMP, MPI arose as a standard to replace a large number of proprietary message passing libraries.

10/23/2012

CS4230

5



Message Passing Library Features

- All communication, synchronization require subroutine calls
 - No shared variables
 - Program runs on a single processor just like any uniprocessor program, except for calls to message passing library
- Subroutines for
 - Communication
 - Pairwise or point-to-point: A message is sent from a specific sending process (point a) to a specific receiving process (point b).
 - Collectives involving multiple processors
 - Move data: Broadcast, Scatter/gather
 - Compute and move: Reduce, AllReduce
 - Synchronization
 - Barrier
 - No locks because there are no shared variables to protect
 - Queries
 - How many processes? Which one am I? Any messages waiting?

10/23/2012

CS4230

6



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

10/23/2012

CS4230

Slide source: Bill Gropp

7



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

10/23/2012

CS4230

8



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( "Greetings from process %d of
           %d\n", rank, size );
    MPI_Finalize();
    return 0;
}
```

11/03/2011

CS4961

Slide source: Bill Gropp

9



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 << "Greetings from process " << rank << "
               of " << size << "\n";
    MPI::Finalize();
    return 0;
}
```

11/03/2011

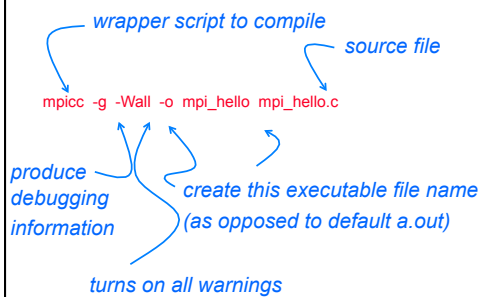
CS4961

Slide source: Bill Gropp

10



Compilation



Copyright © 2010, Elsevier Inc. All rights Reserved

10/23/2012

CS4230

11



Execution

mpiexec -n <number of processes> <executable>

mpiexec -n 1 ./mpi_hello

run with 1 process

mpiexec -n 4 ./mpi_hello

run with 4 processes

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

12



Execution

```
mpiexec -n 1 ./mpi_hello
```

Greetings from process 0 of 1 !

```
mpiexec -n 4 ./mpi_hello
```

Greetings from process 0 of 4 !
Greetings from process 1 of 4 !
Greetings from process 2 of 4 !
Greetings from process 3 of 4 !

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

13



MPI Components

• MPI_Init

- Tells MPI to do all the necessary setup.

```
int MPI_Init(
    int*    argc_p /* in/out */,
    char*** argv_p /* in/out */);
```

• MPI_Finalize

- Tells MPI we're done, so clean up anything allocated for this program.

```
int MPI_Finalize(void);
```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

14



Basic Outline

```
...
#include <mpi.h>
...
int main(int argc, char* argv[]) {
    ...
    /* No MPI calls before this */
    MPI_Init(&argc, &argv);
    ...
    MPI_Finalize();
    /* No MPI calls after this */
    ...
    return 0;
}
```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

15



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

10/23/2012

CS4230

16



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.

10/23/2012

CS4230

Slide source: Bill Gropp

17



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)

10/23/2012

CS4230

Slide source: Bill Gropp

18



Some Basic Clarifying Concepts

- How to organize processes
 - 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`

10/23/2012

CS4230

Slide source: Bill Gropp

19



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

10/23/2012

CS4230

Slide source: Bill Gropp

20



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

10/23/2012

CS4230

Slide source: Bill Gropp

21



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;
}
```

10/23/2012

CS4230

Slide source: Bill Gropp

22



Trapezoidal Rule: Serial algorithm

```
/* Input:  a, b, n */
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;
```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

23



Parallel pseudo-code (naïve)

```
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;
```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

24



First version (1)

```

1 int main(void) {
2     int my_rank, comm_sz, n = 1024, local_n;
3     double a = 0.0, b = 3.0, h, local_a, local_b;
4     double local_int, total_int;
5     int source;
6
7     MPI_Init(NULL, NULL);
8     MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
9     MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);
10
11     h = (b-a)/n;      /* h is the same for all processes */
12     local_n = n/comm_sz; /* So is the number of trapezoids */
13
14     local_a = a + my_rank*local_n*h;
15     local_b = local_a + local_n*h;
16     local_int = Trap(local_a, local_b, local_n, h);
17
18     if (my_rank != 0) {
19         MPI_Send(&local_int, 1, MPI_DOUBLE, 0, 0,
20                 MPI_COMM_WORLD);

```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

25



First version (2)

```

21 } else {
22     total_int = local_int;
23     for (source = 1; source < comm_sz; source++) {
24         MPI_Recv(&local_int, 1, MPI_DOUBLE, source, 0,
25                 MPI_COMM_WORLD, MPI_STATUS_IGNORE);
26         total_int += local_int;
27     }
28 }
29
30 if (my_rank == 0) {
31     printf("With n = %d trapezoids, our estimate\n", n);
32     printf("of the integral from %f to %f = %.15e\n",
33           a, b, total_int);
34 }
35 MPI_Finalize();
36 return 0;
37 } /* main */

```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

26



First version (3)

```

1 double Trap(
2     double left_endpt /* in */,
3     double right_endpt /* in */,
4     int trap_count /* in */,
5     double base_len /* in */) {
6     double estimate, x;
7     int i;
8
9     estimate = (f(left_endpt) + f(right_endpt))/2.0;
10    for (i = 1; i <= trap_count-1; i++) {
11        x = left_endpt + i*base_len;
12        estimate += f(x);
13    }
14    estimate = estimate*base_len;
15
16    return estimate;
17 } /* Trap */

```

CS4230

10/23/2012

27



MPI_Reduce

```

int MPI_Reduce(
    void*      input_data_p /* in */,
    void*      output_data_p /* out */,
    int        count /* in */,
    MPI_Datatype datatype /* in */,
    MPI_Op      operator /* in */,
    int        dest_process /* in */,
    MPI_Comm    comm /* in */);

```

```

MPI_Reduce(&local_int, &total_int, 1, MPI_DOUBLE, MPI_SUM, 0,
           MPI_COMM_WORLD);

```

```

double local_x[N], sum[N];
...
MPI_Reduce(local_x, sum, N, MPI_DOUBLE, MPI_SUM, 0,
           MPI_COMM_WORLD);

```

CS4230

10/23/2012

28



Replace with reduction: OpenMP version

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



```
h = (b-a)/n;
approx = (f(a) + f(b))/2.0;
# pragma omp parallel for num_threads(thread_count) \
    reduction(+: approx)
for (i = 1; i <= n-1; i++)
    approx += f(a + i*h);
approx = h*approx;
```

Copyright © 2010, Elsevier Inc. All rights Reserved

CS4230

10/23/2012

29



MPI also has reduction

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

```
MPI_REDUCE( sendbuf, recvbuf, count, datatype, op, root, comm)
[ IN sendbuf] address of send buffer (choice)
[ OUT recvbuf] address of receive buffer (choice, significant only at root)
[ IN count] number of elements in send buffer (integer)
[ IN datatype] data type of elements of send buffer (handle)
[ IN op] reduce operation (handle)
[ IN root] rank of root process (integer)
[ IN comm] communicator (handle)
```

10/23/2012

CS4230

30



Predefined reduction operators in MPI

Operation Value	Meaning
MPI_MAX	Maximum
MPI_MIN	Minimum
MPI_SUM	Sum
MPI_PROD	Product
MPI_LAND	Logical and
MPI_BAND	Bitwise and
MPI_LOR	Logical or
MPI_BOR	Bitwise or
MPI_LXOR	Logical exclusive or
MPI_EXOR	Bitwise exclusive or
MPI_MAXLOC	Maximum and location of maximum
MPI_MINLOC	Minimum and location of minimum

CS4230/12

31



Collective vs. Point-to-Point Communications

- All the processes in the communicator must call the same collective function.
- For example, a program that attempts to match a call to `MPI_Reduce` on one process with a call to `MPI_Recv` on another process is erroneous, and, in all likelihood, the program will hang or crash.

CS4230

10/23/2012

32



Collective vs. Point-to-Point Communications

- The arguments passed by each process to an MPI collective communication must be "compatible."
- For example, if one process passes in 0 as the `dest_process` and another passes in 1, then the outcome of a call to `MPI_Reduce` is erroneous, and, once again, the program is likely to hang or crash.

CS4230

10/23/2012

33



Collective vs. Point-to-Point Communications

- The `output_data_p` argument is only used on `dest_process`.
- However, all of the processes still need to pass in an actual argument corresponding to `output_data_p`, even if it's just `NULL`.

CS4230

10/23/2012

34



Collective vs. Point-to-Point Communications

- Point-to-point communications are matched on the basis of tags and communicators.
- Collective communications don't use tags.
- They're matched solely on the basis of the communicator and the order in which they're called.

CS4230/12

35



Next Time

- More detail on communication constructs
 - Blocking vs. non-blocking
 - One-sided communication
- Support for data and task parallelism

10/23/2012

CS4230

36

