ECE 677: Distributed Computing Systems

Salim Hariri

High Performance Distributed Computing Laboratory

> University of Arizona Tele: (520) 621-4378

Fall 2012

Distributed Systems Design Framework (Cont)

Distributed Computing Paradigms (DCP)				
Computation Models		Communication Models		
Functional Parallel	Data Parallel	Message Passing	Shared <u>Memory</u>	
System Architecture and Services (SAS)				
Architecture Models		System Level Services		
Computer Networks and Protocols (CNP)				
Computer Networks		Communication Protocols		

Brief overview

- What: standard for a message passing library (C, C++ and Fortran) to be used for message-passing parallel computing.
- When: 92-94 MPI1; 95-97 MPI2
- Size: MPI1: 127 calls; MPI2: ~150 calls.
 - Many parallel programs can be written with 6 basic functions.
 - Functions are orthogonal.
 - Support for many different communication paradigms.
 - Support for different communication modes.
 - Options offered via different function names, rather than parameters.
- Where:
 - Parallel computers and clusters (distributed or shared memory)
 - NOWs (Network of workstations, heterogeneous systems)
- Find more: <u>http://www.mcs.anl.gov/research/project_detail.php?id=</u> <u>2</u>

Companion Material

- Online examples available at <u>http://www.mcs.anl.gov/mpi/tutorial</u>
- ftp://ftp.mcs.anl.gov/pub/mpi/mpiexmple.tar.
 gz contains source code and run scripts that allows you to evaluate your own MPI implementation

The Message-Passing Model

- A process is (traditionally) a program counter and address space
- Processes may have multiple threads(program counters and associated stacks) sharing a single address space. MPI is for communication among processes, which have separate address spaces.
- Interprocess communication consists of
 - Synchronization/Asynchronization
 - Movement of data from one process's address space to another's

What is message passing?

Data transfer.

Requires cooperation of sender and receiver

 Cooperation not always apparent in code

Communication Modes

- Based on the type of send:
 - Synchronous: Completes once the acknowledgement is received by the sender.
 - Buffered send: completes immediately, unless if an error occurs.
 - Standard send: completes once the message has been sent, which may or may not imply that the message has arrived at its destination.
 - Ready send: completes immediately, if the receiver is ready for the message it will get it, otherwise the message is dropped silently.

Synchronous Vs. Asynchronous

 A synchronous communication is not complete until the message has been received.

 An asynchronous communication completes as soon as the message is on the way.

Synchronous Vs. Asynchronous (cont.)









Blocking vs. Non-Blocking

 Blocking, means the program will not continue until the communication is completed.

Non-Blocking, means the program will continue, without waiting for the communication to be completed.

What is MPI?

- A message-passing library specifications:
 - Extended message-passing model
 - Not a language or compiler specification
 - Not a specific implementation or product
- For parallel computers, clusters, and heterogeneous networks.
- Communication modes: standard, synchronous, buffered, and ready.
- Designed to permit the development of parallel software libraries.
- Designed to provide access to advanced parallel hardware for
 - End users
 - Library writers
 - Tool developers

Why to use MPI?

- MPI provides a powerful, efficient, and portable way to express parallel programs.
- MPI was explicitly designed to enable libraries which may eliminate the need for many users to learn (much of) MPI.
- Portable !!!!!!!!!!!!!!!!!!!!!!!!!!!!!!

Is MPI large or small?

MPI is large(125 functions)

- MPI's extensive functionality requires many functions.
- Number of functions not necessarily a measure of complexity.
- MPI is small(6 functions)
 - Many parallel programs can be written with just 6 basic functions.
- MPI is just right
 - One can access flexibility when it is required.
 - One need not master all parts of MPI to use it.
 - MPI is whatever size you like

Features that are NOT part of MPI

Process Management

Remote memory transfer

Threads

Virtual shared memory

Why MPI is simple?

- Many parallel programs can be written using just these six functions, only two of which are non-trivial;
 - MPI_INIT
 - MPI_FINALIZE
 - MPI_COMM_SIZE
 - MPI_COMM_RANK
 - MPI_SEND
 - MPI_RECV

Skeleton MPI Program

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

```
main( int argc, char** argv ) {
    MPI_Init( &argc, &argv );
    /* main part of the program */
    Use MPI function call depend on your data partition
    and parallization architecture
    MPI_Finalize();
}
```

Initializing MPI

- The first MPI routine called in any MPI program must be the initialization routine MPI_INIT
- MPI_INIT is called once by every process, before any other MPI routines

int mpi_Init(int *argc, char **argv);

Startup and endup

- int MPI_Init(int *argc, char ***argv)
 - The first MPI call in any MPI process
 - Establishes MPI environment
 - One and only one call to MPI_INIT per process
- int MPI_Finalize(void)
 - Exiting from MPI
 - Cleans up state of MPI
 - The last call of an MPI process

A minimal MPI program(c)

```
#include ``mpi.h"
#include <stdio.h>
int main(int argc, char *argv[])
{
    MPI_Init(&argc, &argv);
    printf(``Hello, world!\n");
    MPI_Finalize();
    Return 0;
}
```

Commentary

- #include "mpi.h" provides basic MPI definitions and types.
- MPI_Init starts MPI
- MPI_Finalize exits MPI
- Note that all non-MPI routines are local; thus printf run on each process

Notes on C

- In C:
 - mpi.h must be included by using #include mpi.h
 - MPI functions return error codes or MPI_SUCCESS

Error handling

By default, an error causes all processes to abort.

- The user can have his/her own error handling routines.
- Some custom error handlers are available for downloading from the net.

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.

Better 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\n", rank, size);
    MPI_Finalize();
    return 0;
```

```
}
```

Some basic concepts

- Processes can be collected into groups.
- Each message is sent in a context, and must be received in the same context.
- 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.

Compiling and running

Head file

- Fortran -- mpif.h
- C -- mpi.h (*we use C in this presentation)

Compile:

- implementation dependent. Typically requires specification of header file directory and MPI library.
- mpiCC –o destination-filename source-file.c
- mpiCC filename
- Run:
 - mpirun -np <# proc> <executable>

```
#include <stdio.h>
#include <string.h> // this allows us to manipulate text strings
#include "mpi.h" // this adds the MPI header files to the program
int main(int argc, char* argv[]) {
 int my rank; // process rank
 int p; // number of processes
 int source; // rank of sender
 int dest; // rank of receiving process
 int tag = 0; // tag for messages
 char message[100]; // storage for message
 MPI Status status; // stores status for MPI Recv statements
 // starts up MPI
 MPI Init(&argc, &argv);
 // finds out rank of each process
 MPI Comm rank (MPI COMM WORLD, &my rank);
 // finds out number of processes
 MPI Comm size (MPI COMM WORLD, &p);
 if (my rank!=0) {
   sprintf(message, "Greetings from process %d!", my rank);
   dest = 0; // sets destination for MPI Send to process 0
   // sends the string to process 0
   MPI Send(message, strlen(message)+1, MPI CHAR, dest, tag, MPI COMM WORLD);
 } else {
   for (source = 1; source < p; source++) {
     // receives greeting from each process
     MPI Recv(message, 100, MPI CHAR, source, tag, MPI COMM WORLD, &status);
     printf("%s\n", message); // prints out greeting to screen
   }
  }
 MPI Finalize(); // shuts down MPI
 return 0;
}
```

Result

- mpicc hello.c
- mpirun -np 6 a.out
 - Greetings from process 1! Greetings from process 2! Greetings from process 3! Greetings from process 4!
 - Greetings from process 5!

MPI blocking send

- MPI_SEND(void *start, int count,MPI_DATATYPE datatype, int dest, int tag, MPI_COMM comm)
- The message buffer is described by (start, count, datatype).
- dest is the rank of the target process in the defined communicator.
- tag is the message identification number.

MPI_RECV(void *start, int count, MPI_DATATYPE datatype, int source, int tag, MPI_COMM comm, MPI_STATUS *status)

- Source is the rank of the sender in the communicator.
- The receiver can specify a wildcard value for souce (MPI_ANY_SOURCE) and/or a wildcard value for tag (MPI_ANY_TAG), indicating that any source and/or tag are acceptable
- Status is used for exrtra information about the received message if a wildcard receive mode is used.
- If the count of the message received is less than or equal to that described by the MPI receive command, then the message is successfully received. Else it is considered as a buffer overflow error.

More comment on send and receive

- A receive operation may accept messages from an arbitrary sender, but a send operation must specify a unique receiver.
- Source equals destination is allowed, that is, a process can send a message to itself.

Review of Basic MPI routines

- MPI is used to create parallel programs based on message passing
- Usually the same program is run on multiple processors
- The 6 basic calls in MPI are:
 - MPI_INIT(ierr)
 - MPI_COMM_RANK(MPI_COMM_WORLD, myid, ierr)
 - MPI_COMM_SIZE(MPI_COMM_WORLD, numprocs, ierr)
 - MPI_Send(buffer, count, MPI_INTEGER, destination, tag, MPI_COMM_WORLD, ierr)
 - MPI_Recv(buffer, count, MPI_INTEGER, source, tag, MPI_COMM_WORLD, status, ierr)
 - MPI_FINALIZE(ierr)

Communication Primitives

Communications on distributed memory computers:

- Point to Point
- One to All Broadcast
- All to All Broadcast
- One to All Personalized
- All to All Personalized
- Shifts
- Collective Computation

MPI basic send/receive

We need to fill in the details in

Process 0 Process 1 Send(data) Receive(d ata)

- 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 operation to complete?

Data Types

- The data message which is sent or received is described by a triple (address, count, datatype).
- The following data types are supported by MPI:
 - Predefined data types that are corresponding to data types from the programming language.
 - Arrays.
 - Sub blocks of a matrix
 - User defined data structure.
 - A set of predefined data types

MPI Data Types in C

C MPI Types			
MPI_CHAR	signed char		
MPI_SHORT	signed short int		
MPI_INT	signed int		
MPI_LONG	signed long int		
MPI_UNSIGNED_CHAR	unsigned char		
MPI_UNSIGNED_SHORT	unsigned short int		
MPI_UNSIGNED	unsigned int		
MPI_UNSIGNED_LONG	unsigned long int		
MPI_FLOAT	float		
MPI_DOUBLE	double		
MPI_LONG_DOUBLE	long double		
MPI_BYTE	-		
MPI_PACKED	-		
Why defining the data types during the send of a message?

Because communications take place between heterogeneous machines. Which may have different data representation and length in the memory.

Broadcast and reduce

- MPI_Bcast(buffer, count, datatype, root, comm)
 - Broadcast the message of length count in buffer from the process root to all other processes in the group. All processes must call with same arguments.
- MPI_Reduce(sbuf, rbuf, count, stype, op, root, comm)
 - Apply the reduction function op to the data of each process in the group (type stype in sbuf) and store the result in rbuf on the root process. op can be a pre-defined function, or defined by the user.

Global Communications in MPI: Broadcast

- All nodes call MPI_Bcast
- One node (root) sends a message all others receive the message
- C
 - MPI_Bcast(&buffer, count, datatype, root, communicator);
- Fortran
 - call MPI_Bcast(buffer, count, datatype, root, communicator, ierr)

Root is node that sends the message

Global Communications in MPI: Broadcast

- broadcast.c is a parallel program to broadcast data using MPI_Bcast
 - Initialize MPI
 - Have processor 0 broadcast an integer
 - Have all processors print the data
 - Quit MPI

Global Communications in MPI: Broadcast

```
This is a simple broadcast program in MPI
#include <stdio.h>
#include "mpi.h"
int main(argc,argv)
int argc;
char *argv[];
  int i, myid, numprocs;
  int source, count;
  int buffer[4];
  MPI Status status;
  MPI Request request;
  MPI_Init(&argc,&argv);
  MPI Comm size(MPI COMM WORLD,&numprocs);
  MPI Comm rank(MPI COMM WORLD,&myid);
  source=0;
  count=4;
  if(myid == source)
   for(i=0;i<count;i++)</pre>
    buffer[i]=i;
  }
  MPI_Bcast(buffer,count,MPI_INT,source,MPI_COMM_WORLD);
  for(i=0;i<count;i++)</pre>
   printf("%d ",buffer[i]);
  printf("\n");
  MPI Finalize();
}
```

Global Communications in MPI: Reduction

- Used to combine partial results from all processors
- Result returned to root processor
- Several types of operations available. For example summation, maximum etc
- Works on single elements and arrays

Global Communications in MPI: MPI_Reduce

- C
 - int MPI_Reduce(&sendbuf, &recvbuf, count, datatype, operation,root, communicator)
- Fortran
 - call MPI_Reduce(sendbuf, recvbuf, count, datatype, operation,root, communicator, ierr)
- Parameters
 - Like MPI_Bcast, a <u>root</u> MPI process is specified.
 - <u>Operation</u> is mathematical operation

Global Communications in MPI: MPI_Reduce

MPI_MAX MPI MIN MPI_PROD MPI SUM MPI_LAND MPI_LOR MPI_LXOR MPI_BAND MPI BOR MPI_BXOR MPI MAXLOC MPI MINLOC

Maximum Minimum Product Sum Logical and Logical or Logical exclusive or Bitwise and Bitwise or Bitwise exclusive or Maximum value and location Minimum value and location

Example: PI in C - 1

```
#include "mpi.h"
#include <math.h>
int main(int argc, char *argv[])
ł
  int done = 0, n, myid, numprocs, i, rc;
  double PI25DT = 3.141592653589793238462643;
  double mypi, pi, h, sum, x, a;
  MPI Init(&argc,&argv);
  MPI Comm size (MPI COMM WORLD, & numprocs);
  MPI Comm rank (MPI COMM WORLD, & myid);
  while (!done) {
    if (myid == 0) {
      printf("Enter the number of intervals: (0 quits) ");
      scanf("%d", &n);
    }
    MPI Bcast(&n, 1, MPI INT, 0, MPI COMM WORLD);
    if (n == 0) break;
```

Example: PI in C - 2

```
h = 1.0 / (double) n;
  sum = 0.0;
  for (i = myid + 1; i \le n; i += numprocs) {
    x = h * ((double)i - 0.5);
    sum += 4.0 / (1.0 + x*x);
  }
  mypi = h * sum;
  MPI Reduce (& mypi, & pi, 1, MPI DOUBLE, MPI SUM, 0,
             MPI COMM WORLD);
  if (myid == 0)
    printf("pi is approximately %.16f, Error is
%.16f\n",
            pi, fabs(pi - PI25DT));
}
MPI Finalize();
return 0;
```

Point to Point Communications in MPI

- Basic operations of Point to Point (PtoP) communication in MPI
- Several steps are involved in the PtoP communication
- Sending process
 - data is copied to the user buffer by the user
 - User calls one of the MPI send routines
 - System copies the data from the user buffer to the system buffer
 - System sends the data from the system buffer to the destination processor

Point to Point Communications in MPI

Receiving process

- User calls one of the MPI receive subroutines
- System receives the data from the source process, and copies it to the system buffer
- System copies the data from the system buffer to the user buffer
- User uses the data in the user buffer

Point to Point Communications in MPI



More information of point to point communication are in the Appendixes

MPI tags

- Message are sent with an accompanying user-defined integer tag, to assist the receiving process in identifying the message.
- Message 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 datatype.

MPI_Barrier

- Blocks the caller until all members in the communicator have called it.
- Used as a synchronization tool.
- C
 - MPI_Barrier(comm)
- Fortran
 - Call MPI_BARRIER(COMM, IERROR)
- Parameter
 - Comm: communicator (often MPI_COMM_WORLD)

Overview of Some Advanced MPI Routines

- Can split MPI communicators (MPI_Comm_split)
- Probe incoming messages (MPI_Probe)
- Asynchronous communication (MPI_Isend, MPI_Irecv, MPI_Wait, MPI_Test etc)
- Scatter different data to different processors (MPI_Scatter), Gather (MPI_Gather)
- MPI_AllReduce, MPI_Alltoall
- Derived data types (MPI_TYPE_STRUCT etc)
- MPI I/O

Group routines

- MPI_Group_size returns number of processes in group
- MPI_Group_rank returns rank of calling process in group
- MPI_Group_compare compares group members and group order
- MPI_Group_translate_ranks translates ranks of processes in one group to those in another group
- MPI_Comm_group returns the group associated with a communicator
- MPI_Group_union creates a group by combining two groups
- MPI_Group_intersection creates a group from the intersection of two groups

Group routines ...

- MPI_Group_difference creates a group from the difference between two groups
- MPI_Group_incl creates a group from listed members of an existing group
- MPI_Group_excl creates a group excluding listed members of an existing group
- MPI_Group_range_incl creates a group according to first rank, stride, last rank
- MPI_Group_range_excl creates a group by deleting according to first rank, stride, last rank
- MPI_Group_free marks a group for deallocation

Communicator routines

- MPI_Comm_size returns number of processes in communicator's group
- MPI_Comm_rank returns rank of calling process in communicator's group
- MPI_Comm_compare compares two communicators
- MPI_Comm_dup duplicates a communicator
- MPI_Comm_create creates a new communicator for a group
- MPI_Comm_split splits a communicator into multiple, nonoverlapping communicators
- MPI_Comm_free marks a communicator for deallocation

Collective communication

- MPI_Allgather All processes gather messages
- MPI_Allreduce Reduce to all processes
- MPI_Alltoall All processes gather distinct messages
- MPI_Bcast Broadcast a message
- MPI_Gather Gather a message to root
- MPI_Reduce Global reduce operation
- MPI_ReduceScatter Reduce and scatter results
- MPI_Scatter Scatter a message from root
- MPI_Scan Global prefix reduction

Collective Data Movement





More Collective Data Movement





Collective Computation



Timing

 MPI Wtime() returns the wall-clock time.

```
double start, finish, time;
MPI_Barrier(MPI_COMM_WORLD);
start = MPI_Wtime();
```

```
...
MPI_Barrier(MPI_COMM_WORLD);
finish = MPI_Wtime();
time = finish - start;
```

MPI Trace Output

MPI Routine	#ca.	lls	avg. by	rtes	time(sec)
MPI_Comm_size MPI_Comm_rank MPI_Send MPI_Recv MPI_Barrier		1 1 500 500 500	102 102	0.0 0.0 24.0 24.0 0.0	0.000 0.000 0.001 0.008 0.013
total communication time total elapsed time user cpu time system time maximum memory size	= 0.022 = 3.510 = 3.500 = 0.010 = 15856	seconds seconds seconds seconds KBytes.			
Message size distribution	is:				
MPI_Send	#calls 500	avg.	bytes 1024.0	time	e(sec) 0.001
MPI_Recv	#calls 500	avg.	bytes 1024.0	time	e(sec) 0.008
Call Graph Section:					
<pre>communication time = 0.02 MPI Routine MPI_Send MPI_Recv MPI_Barrier</pre>	2 sec,]	parent = #calls 500 500 500	poisso t	on time(sec) 0.001 0.008 0.013	
<pre>communication time = 0.00 MPI Routine MPI_Comm_size MPI Comm rank</pre>	00 sec,]	parent = #calls 1 1	dot t	ime(sec) 0.000 0.000	

MPI-2

- MPI-2 new topics:
 - process creation and management, including client/server routines
 - one-sided communications (put/get, active messages)
 - extended collective operations
 - external interfaces
 - **I/O**

Designing MPI programs

- Partitioning
 - Before tackling MPI
- Communication
 - Many point to collective operations
- Agglomeration
 - Needed to produce MPI processes
- Mapping
 - Handled by MPI



MPI

Pros:

- Very portable
- Requires no special compiler
- Requires no special hardware but can make use of high performance hardware
- Very flexible -- can handle just about any model of parallelism
- No shared data! (You don't have to worry about processes "treading on each other's data" by mistake.)
- Can download free libraries for your Linux/PC!
- Forces you to do things the "right way" in terms of decomposing your problem.

Cons:

- All-or-nothing parallelism (difficult to incrementally parallelize existing serial codes)
- No shared data! Requires distributed data structures
- Could be thought of assembler for parallel computing -- you generally have to write more code
- Partitioning operations on distributed arrays can be messy.

MPI v.s. OpenMP

- Message passing v.s. shared data
- Processes v.s. Threads
- MPI has no work sharing structure.

OpenMP

Pros:

- Incremental parallelism -- can parallelize existing serial codes one bit at a time
- Quite simple set of directives
- Shared data!
- Partitioning operations on arrays is very simple.
- Cons:
 - Requires proprietary compilers
 - Requires shared memory multiprocessors
 - Shared data!
 - Having to think about what data is shared and what data is private
 - Cannot handle models like master/slave work allocation (yet)
 - Generally not as scalable (more synchronization points)
 - Not well-suited for non-trivial data structures like linked lists, trees etc

Homework #1 (programming) will be posted.

Due: September 24 before the class

Appendix

Unidirectional Communication

Blocking send and blocking receive

 if (myrank == 0) then call MPI_Send(...) elseif (myrank == 1) then call MPI_Recv(....) endif

Non-blocking send and blocking receive

 if (myrank == 0) then call MPI_ISend(...) call MPI_Wait(...) else if (myrank == 1) then call MPI_Recv(....) endif **Unidirectional Communication**

Blocking send and non-blocking recv

endif

Non-blocking send and non-blocking recv

Bidirectional Communication

- Need to be careful about deadlock when two processes exchange data with each other
- Deadlock can occur due to incorrect order of send and recv or due to limited size of the system buffer


Case 1 : both processes call send first, then recv if (myrank == 0) then call MPI_Send(....) call MPI_Recv (...) elseif (myrank == 1) then call MPI_Send(....) call MPI_Send(....)

endif

- No deadlock as long as system buffer is larger than send buffer
- Deadlock if system buffer is smaller than send buf
- If you replace MPI_Send with MPI_Isend and MPI_Wait, it is still the same
- Moral : there may be error in coding that only shows up for larger problem size

 Case 2 : both processes call recv first, then send if (myrank == 0) then call MPI_Recv(....) call MPI_Send (...) elseif (myrank == 1) then call MPI_Recv(....) call MPI_Send(....) endif

 The above will always lead to deadlock (even if you replace MPI_Send with MPI_Isend and MPI_Wait)

The following code can be safely executed

if (myrank == 0) then
 call MPI_Irecv(....)
 call MPI_Send (...)
 call MPI_Wait(...)
elseif (myrank == 1) then
 call MPI_Irecv(....)
 call MPI_Send(....)
 call MPI_Send(....)
 call MPI_Wait(....)
endif

 Case 3 : one process call send and recv in this order, and the other calls in the opposite order

```
if (myrank == 0 ) then
            call MPI_Send(....)
            call MPI_Recv(...)
elseif (myrank == 1) then
            call MPI_Recv(....)
            call MPI_Send(....)
endif
```

- The above is always safe
- You can replace both send and recv on both processor with Isend and Irecv

Where to get MPI library?

MPICH (WINDOWS / UNICES)

http://www-unix.mcs.anl.gov/mpi/mpich/

Open MPI (UNICES)

http://www.open-mpi.org/

SHARED MEMORY OPENMP

Different types of parallel platforms: Distributed Memory

M0, M1, ... Mn are memories associated with processors P0, P1, ..., Pn.



Distributed Memory Architecture

Different types of parallel platforms: Shared Memory



Pittsburgh Supercomputing

Different types of parallel platforms: Shared Memory

- SMP: Symmetric Multiprocessing
 - Identical processing units working from the same main memory
 - SMP machines are becoming more common in the everyday workplace
 - Dual-socket motherboards are very common, and quad-sockets are not uncommon
 - 2 and 4 core CPUs are now commonplace
 - Intel Larabee: 12-48 cores in 2009-2010
- ASMP: Asymmetric Multiprocessing
 - Not all processing units are identical
 - Cell processor of PS3

Parallel Programming Models

- Shared Memory
 - Multiple processors sharing the same memory space
- Message Passing
 - Users make calls that explicitly share information between execution entities
- Remote Memory Access
 - Processors can directly access memory on another processor
- These models are then used to build more sophisticated models
 - Loop Driven
 - Function Driven Parallel (Task-Level)

Introduction

- OpenMP is designed for shared memory systems.
- OpenMP is easy to use
 - achieve parallelism through compiler directives
 - or the occasional function call
- OpenMP is a "quick and dirty" way of parallelizing a program.
- OpenMP is usually used on existing serial programs to achieve moderate parallelism with relatively little effort

Composition assignmented assign



Pittsburgh Supercomputing

OpenMP Execution Model In MPI, all threads are active all the time

- In OpenMP, execution begins only on the master thread. Child threads are spawned and released as needed.
 - Threads are spawned when program enters a parallel region.
 - Threads are released when program exits a parallel region

OpenMP Execution Model



Parallel Region Example: For loop

Fortran:

```
!$omp parallel do
do i = 1, n
        a(i) = b(i) + c(i)
enddo
```

This comment or pragma tells openmp compiler to spawn threads *and* distribute work among those threads

```
C/C++:

#pragma omp parallel for (i=1; i<=n; i++)

a[i] = b[i] + c[i]; between the threads
```

Pros of OpenMP

- Because it takes advantage of shared memory, the programmer does not need to worry (that much) about data placement
- Programming model is "serial-like" and thus conceptually simpler than message passing
- Compiler directives are generally simple and easy to use
- Legacy serial code does not need to be rewritten

Cons of OpenMP

- Codes can only be run in shared memory environments!
 - In general, shared memory machines beyond ~8 CPUs are much more expensive than distributed memory ones, so finding a shared memory system to run on may be difficult
- Compiler must support OpenMP
 - whereas MPI can be installed anywhere
 - However, gcc 4.2 now supports OpenMP

Cons of OpenMP

- In general, only moderate speedups can be achieved.
 - Because OpenMP codes tend to have serial-only portions, Amdahl's Law prohibits substantial speedups
- Amdahl's Law:
 - F = Fraction of serial execution time that cannot be
 - parallelized
 - N = Number of processors





If you have big loops that dominate execution time, these are ideal targets for

Compiling and Running OpenMP

- True64:
- SGI IRIX:
- IBM AIX:
- Portland Group:
- Intel:
- gcc (4.2)

-mp -mp -qsmp=omp -mp -openmp -fopenmp

Compiling and Running OpenMP

- OMP_NUM_THREADS environment variable sets the number of processors the OpenMP program will have at its disposal.
- Example script

#!/bin/tcsh
setenv OMP_NUM_THREADS 4
mycode < my.in > my.out

OpenMP Basics: 2 Approaches to Parallelism

Divide loop iterations among threads: We will focus mainly on loop level parallelism in this lecture



Divide various sections of code between threads

Sections: Functional parallelism

#pragma omp parallel master thread ł *#pragma omp sections* FORK #pragma omp section SECTIONS team block1 OIN #pragma omp section block2 master thread

Philip Blood (Scientific Specialist) Pittsburgh Supercomputing Image from: https://computing.llnl.gov/tutorials/openM P



Pitfall #1: Data dependencies

Consider the following code:

```
a[0] = 1;
for(i=1; i<5; i++)
a[i] = i + a[i-1];
```

- There are dependencies between loop iterations.
- Sections of loops split between threads will not necessarily execute in order
- Out of order loop execution will result in undefined behavior

Pitfall #1: Data dependencies

- 3 simple rules for data dependencies
 - 1. All assignments are performed on arrays.
 - 2. Each element of an array is assigned to by at most one iteration.
 - 3. No loop iteration reads array elements modified by any other iteration.

Avoiding dependencies by using Private Variables (Pitfall #1.5)

Consider the following loop:

```
#pragma omp parallel for
{
   for(i=0; i<n; i++) {
      temp = 2.0*a[i];
      a[i] = temp;
      b[i] = c[i]/temp;
   }
}</pre>
```

}

 By default, all threads share a common address space. Therefore, all threads will be modifying *temp* simultaneously Avoiding dependencies by using Private Variables (Pitfall #1.5)

The solution is to make *temp* a threadprivate variable by using the "private" clause:

```
#pragma omp parallel for private(temp)
{
```

```
for(i=0; i<n; i++) {
   temp = 2.0*a[i];
   a[i] = temp;
   b[i] = c[i]/temp;
}</pre>
```

Philip Blood (Scientific Specialist)

Avoiding dependencies by using Private Variables (Pitfall #1.5)

Default OpenMP behavior is for variables to be shared. However, sometimes you may wish to make the default private and explicitly declare your shared variables (but only in Fortran!):

```
!$omp parallel do default(private) shared(n,a,b,c)
    do i=1,n
        temp = 2.0*a(i)
        a(i) = temp
        b(i) = c(i)/temp;
    enddo
!$omp end parallel do
```

Private variables

- Note that the loop iteration variable (e.g. *i* in previous example) is *private* by default
- Caution: The value of any variable specified as *private* is *undefined* both upon entering and leaving the construct in which it is specified
- Use *firstprivate* and *lastprivate* clauses to retain values of variables declared as *private*

Use of function calls within parallel loops

- In general, the compiler will not parallelize a loop that involves a function call unless is can guarantee that there are no dependencies between iterations.
 - sin(x) is OK, for example, if x is private.
- A good strategy is to inline function calls within loops. If the compiler can inline the function, it can usually verify lack of dependencies.
- System calls do not parallelize!!!

Pitfall #2: Updating shared variables simultaneously

Consider the following serial code:

```
the_max = 0;
for (i=0;i<n; i++)
  the max = max(myfunc(a[i]), the max);
```

- This loop can be executed in any order, however the_max is modified every loop iteration.
- Use "critical" clause to specify code segments that can only be executed by one thread at a time:

```
#pragma omp parallel for private(temp)
{
   for(i=0; i<n; i++) {
      temp = myfunc(a[i]);
      #pragma omp critical
      the_max = max(temp, the_max);
   }
}</pre>
```

Reduction operations

Now consider a global sum:

```
for(i=0; i<n; i++)
```

sum = sum + a[i];

 This can be done by defining "critical" sections, but for convenience, OpenMP also provides a reduction clause:

```
#pragma omp parallel for reduction(+:sum)
{
   for(i=0; i<n; i++)
      sum = sum + a[i];
}</pre>
```

```
Philip Blood ( Scientific Specialist
)
Pittsburgh Supercomputing
```

Reduction operations

- C/C++ reduction-able operators (and initial values):
 - (0)• + (0)(1)* **&** (~0) (0)(0)Λ **&&** (1)(0)• ||

Pitfall #3: Parallel overhead

Spawning and releasing threads results in *significant* overhead.



Pittsburgh Supercomputing

Pitfall #3: Parallel overhead



Pitfall #3: Parallel Overhead

- Spawning and releasing threads results in *significant* overhead.
- Therefore, you want to make your parallel regions as large as possible
 - Parallelize over the largest loop that you can (even though it will involve more work to declare all of the private variables and eliminate dependencies)
 - Coarse granularity is your friend!
Separating "Parallel" and "For" directives to reduce overhead

In the following example, threads are spawned only once, not once per loop:

```
#pragma omp parallel {
    #pragma omp for
    for(i=0; i<maxi; i++)
        a[i] = b[i];
    #pragma omp for
    for(j=0; j<maxj; j++)
        c[j] = d[j];
}</pre>
```

Use "nowait" to avoid barriers

- At the end of every loop is an implied barrier.
- Use "nowait" to remove the barrier at the end of the first loop:

```
#pragma omp parallel {
    #pragma omp for nowait
    for(i=0; i<maxi; i++)
        a[i] = b[i];
    #pragma omp for
    for(j=0; j<maxj; j++)
        c[j] = d[j];
}</pre>
Barrier removed by
"nowait" clause
```

Thread stack

- Each thread has its own memory region called the thread stack
- This can grow to be quite large, so default size may not be enough
- This can be increased (e.g. to 16 MB): csh:

limit stacksize 16000; setenv KMP_STACKSIZE 16000000 **bash**:

ulimit -s 16000; export KMP_STACKSIZE=16000000

Useful OpenMP Functions

- void omp_set_num_threads(int num_threads)
 - Sets the number of OpenMP threads (overrides OMP_NUM_THREADS)
- int omp_get_thread_num()
 - Returns the number of the current thread
- int omp_get_num_threads()
 - Returns the total number of threads currently participating in a parallel region
 - Returns "1" if executed in a serial region
- For portability, surround these functions with #ifdef _OPENMP
- #include <omp.h>

- OpenMP partitions workload into "chunks" for distribution among threads
- Default strategy is static:



- This strategy has the least amount of overhead
- However, if not all iterations take the same amount of time, this simple strategy will lead to load imbalance.



- OpenMP offers a variety of scheduling strategies:
 - schedule(static,[chunksize])
 - Divides workload into equal-sized chunks
 - Default *chunksize* is Nwork/Nthreads
 - Setting *chunksize* to less than this will result in chunks being assigned in an interleaved manner
 - Lowest overhead
 - Least optimal workload distribution

- schedule(dynamic,[chunksize])
 - Dynamically assigned chunks to threads
 - Default *chunksize* is 1
 - Highest overhead
 - Optimal workload distribution
- schedule(guided, [chunksize])
 - Starts with big chunks proportional to (number of unassigned iterations)/(number of threads), then makes them progressively smaller until chunksize is reached
 - Attempts to seek a balance between overhead and workload optimization

schedule(runtime)

- Scheduling can be selected at runtime using OMP_SCHEDULE
- **e.g.** setenv OMP_SCHEDULE "guided, 100"
- In practice, often use:
 - Default scheduling (static, large chunks)
 - Guided with default chunksize
- Experiment with your code to determine optimal strategy

What we have learned

- How to compile and run OpenMP progs
- Private vs. shared variables
- Critical sections and reductions for updating scalar shared variables
- Techniques for minimizing thread spawning/exiting overhead
- Different scheduling strategies

Summary

- OpenMP is often the easiest way to achieve moderate parallelism on shared memory machines
- In practice, to achieve decent scaling, will probably need to invest some amount of effort in tuning your application.
- More information available at:
 - <u>https://computing.llnl.gov/tutorials/openMP/</u>
 - <u>http://www.openmp.org</u>
 - Using OpenMP, MIT Press, 2008

Processes, Threads, and Parallel Programming

- Process: is a program counter and address space.
- Thread: Threads have their own program counter and a memory stack, but share the other resources within the process.
- Shared memory uses thread as their milestone for parallel programming.

OpenMP

- Portable API for shared memory thread-based parallelism
- C/C++ and Fortran Master
 "Fork-Join" model Fork
 parall el
 join
 n Master

Greg Wolffe / Christian Trefftz Grand Valley State University

OpenMP Fundamentals
 Environment variables
 export OMP_NUM_THREADS="4"

- Library functions omp_set_num_threads (4);
- Compiler directives (#pragmas)
 #pragma omp parallel num_threads (4)

```
OpenMP Programming
Source code
  #include ``omp.h"
  #pragma omp parallel
  ł
     ...parallel region
  }
Compiling
  q++ -fopenmp filename.cc -o filename
```

Greg Wolffe / Christian Trefftz Grand Valley State University

Directive Responsibility

- Work-sharing
- Data scoping
- Synchronization
- Scheduling

OpenMP: Work Sharing

- Parallel region: partition work
 Each thread executes same code
- Parallel for loop: partition iterations
 Threads share iterations of loop
- Parallel section: functional parallelism
 Threads perform different tasks

Greg Wolffe / Christian Trefftz Grand Valley State University

OpenMP: Data

- Shared: threads access a single copy of the data object
- Private: each thread gets volatile copy
 - Firstprivate: initialized from master
 - Lastprivate: master's copy updated with last value of last thread

Critical Section Problem

- Shared memory system ⇒ shared data
- Shared data ⇒ concurrent access
- Concurrent access ⇒ corrupted variables

 Critical section problem: ensure "correct" access to shared data

Greg Wolffe / Christian Trefftz Grand Valley State University



Concurrency Control

- Synchronization
 - <u>Mutex</u> ensures exclusive access to critical section of code
 - <u>Barrier</u> causes a group of threads to pause until all have reached a defined point
- Signalling
 - <u>Conditional Wait</u> waits for some event; signals when it occurs
 - <u>Broadcasting</u> signals a group of waiting threads



// program code

Greg Wolffe / Christian Trefftz Grand Valley State University







Greg Wolffe / Christian Trefftz Grand Valley State University





Broadcast



Greg Wolffe / Christian Trefftz Grand Valley State University

OpenMP: Scheduling

- Static: splits iteration space into blocks of size *chunk*
- Dynamic: assign blocks to threads as they become idle (uneven workloads)
- Guided: adjusts *chunk*-size exponentially until all assigned

Home Work