#### DISTRIBUTED MEMORY PROGRAMMING WITH MPI

Carlos Jaime Barrios Hernández, PhD.



### **Remember Special Features of Architecture**

- Remember "concurrency": it exploits better the resources (shared) within a computer.
- Exploit SIMD and MIMD Architectures





SIMD

### **Cluster Computing Architecture**



**Parallel Applications** 

**Parallel Applications** 

Parallel Programming Environment

**Middleware** 

(Single System Image and Availability Infrastructure)



## **Distributed Computing Paradigms**

- Communication Models:
  - Message Passing
  - Shared Memory
- Computation Models:
  - Functional Parallel
  - Data Parallel

## Message Passing

- A process is a program counter and address space.
- Message passing is used for communication among processes.
- Inter-process communication:
  - Type: Synchronous / Asynchronous
  - Movement of data from one process's address space to another's

### 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. )



## What is message passing?

- Data transfer.
- Requires cooperation of sender and receiver
- Cooperation not always apparent in code

#### • MPI in a nutshell

- It is a library specification
- Works natively with C and Fortran
- Not a specific implementation or product Passing
- -Scalable
  - Must handle multiple machines
- -Portable
  - Sockets API change from one OS to another
  - Handles Big-endian/little-endian architectures

- Efficient

- Optimized communication algorithms
- Allow communication and computation overlap





#### MPI – Message Passing Interface

#### • MPI References

- Books

– Using MPI: Portable Parallel Programming with the Message Passing Interface, by Gropp, Lusk, and Skejellum, MIT Press, 1994.

 MPI: The Complete Reference, by Snir, Otto, Huss-Lederman, Walker, and Dongarra, MIT Press, 1996.

– Parallel Programming with MPI, by Peter Pacheco, Morgan Kaufmann, 1997.

- The standard:

- at <u>http://www.mpi-forum.org</u>

## **MPI History**

- 1990 PVM: Parallel Virtual Machine (Oak Ridge Nat' I Lab)
  - Message-passing routines
  - Execution environment (spawn + control parallel processes)
  - No an industry standard
- 1992 meetings (Workshop, Supercomputing'92)
- 1993 MPI draft
- 1994 MPI Forum (debates)
- 1994 MPI-1.0 release (C & Fortran bindings) + standardization
- 1995 MPI-1.1 release
- 1997 MPI-1.2 release (errata) + MPI-2 release (new features, C++ & Fortran 90 bindings)
- ???? MPI-3 release (new: FT, hybrid, p2p, RMA, ...)
- 2000 MPI (ch), Madeline, V4....
- 2005 OpenMPI...

#### • MPI

- Use of a single program, on multiple data
- What does it do?
  - way of identifying process
  - Independent of low-level API
  - Optimized communication
  - Allow communication and computation overlap
- What does it do not?
  - gain performance of application for free
  - application must be adapted

### Features of MPI

- General
  - Communications combine context and group for message security.
  - Thread safety can't be assumed for MPI programs.

### Features that are NOT part of MPI

- Process Management
- Remote memory transfer
- Threads
- Virtual shared memory

## 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.
- Good way to learn about subtle issues in parallel computing

## How big is the MPI library?

• Huge (125 Functions).

• Basic ( 6 Functions ).

### Group and Context



## Group and Context (cont.)

- Are two important and indivisible concepts of MPI.
- Group: is the set of processes that communicate with one another.
- Context: it is somehow similar to the frequency in radio communications.
- Communicator: is the central object for communication in MPI. Each communicator is associated with a group and a context.

#### **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.

## 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.

#### • Possible Programming Workflow



- Master/Slave
  - Master is one process that centrilizes all tasks
  - •Slaves starve for work



- Master/Slave
  - Master is often the bottleneck
  - Scalability is limited due to centralization
  - Possible to use replication to improve performance
  - It is adatable to heterogenous platforms

Task 1

Task 2

Task 3

Task 4

#### • Pipeline

- Each process plays a specific role, pipeline stages
- Data follows in a single direction
- Parallelism is achieved when the pipeline is full



#### • Pipeline

- Scalabillity is limited by the number of stages
- Synchronization may lead to bubbles
  - Slow sender
  - Fast receiver
- Difficult to use on heterogenous platforms

Result(60)

#### • Divide and Conquer

- Recursevely partion task on roughly equal sized tasks
- Or process the taks if it is small





• Divide and Conquer

- More scalable
- Possible to use replicated branches
- In practice is difficult to split tasks
- Suitable for branch and bound algorithms

#### • Installing

- Some common MPI implementations, all free:
  - OpenMPI

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

 MPICH-2 http://www.mcs.anl.gov/research/projects/ mpich2/

- LAM/MPI
http://www.lam-mpi.org/

#### • Installing

- I'm using MPICH-2
- Installed in Ubuntu 10.04 Lucid Lynx with
  - \$ sudo apt-get install mpich2
- Should work for most Debian based distributions
- Must create a local configuration file
  - \$ echo "MPD\_SECRET\_WORD=ChangeMe" > ~/.mpd.conf

#### • Test program

```
#include <mpi.h>
#include <stdio.h>
int main(int argc, char **argv){
   /* Initialize MPI */
   MPI_Init(&argc, &argv);
   printf("Test Program\n");
   /* Finalize MPI */
   return MPI_Finalize();
```

}

### **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
 partitioning and the parallelization
 architecture
 */
    MPI Finalize();
}
```

## 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;
}
```

#### Compiling

 Compiled with gcc, but a mpicc script is provided to invoke gcc with specific MPI options enabled

\$ mpicc mpi\_program.c -o my\_mpi\_executable

- Executed with a specital script

- \$ mpirun -np 1 my\_mpi\_executable
- \$ mpirun -np 2 my\_mpi\_executable
- \$ mpirun -np 3 my\_mpi\_executable

#### Running

 Compiled with gcc, but a mpicc script is provided to invoke gcc with specific mpi functions

\$ mpicc mpi\_program.c -o my\_mpi\_executable

#### - For a complete list of parameters try

\$ man mpicc

- Executed with a specital scrip

\$ mpirun -np 2 my\_mpi\_executable

### A minimal MPI program(c) (cont.)

- #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
- Note: 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.

### Improved 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;
```

#### How many processing units are available?

int MPI\_Comm\_size(MPI\_Comm comm, int \*psize)

– Group of process to communicate

Default Communicator: For grouping all process use
 MPI\_COMM\_WORLD

- psize

 Passed as reference will return the total amoung of proccess in this communicator

## 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

### **Basic MPI types**

#### MPI datatype

#### <u>C datatype</u>

MPI CHAR signed char MPI SIGNED CHAR signed char MPI UNSIGNED\_CHAR unsigned char signed short **MPI SHORT** MPI UNSIGNED SHORT unsigned short MPI INT signed int MPI UNSIGNED unsigned int MPI LONG signed long MPI UNSIGNED\_LONG unsigned long MPI FLOAT float MPI DOUBLE double MPI LONG DOUBLE long double

• Exercise 1 – Hello World

• Create program that prints hello world and the total number of available process on the screen

• Use –np with a variable number to verify that your program is working

- Exercise 2 Who am I
  - •If I am process 0
    - Prints: "hello world"
  - else
    - Prints: "I'm process <ID>"
    - Replacing <ID> by the process rank

# 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.

## 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 blocking receive

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.

## MPI\_STATUS

Status is a data structure

•In C:

```
int recvd_tag, recvd_from, recvd_count;
MPI_Status status;
MPI_Recv(..., MPI_ANY_SOURCE, MPI_ANY_TAG, ...,
&status)
recvd_tag = status.MPI_TAG;
recvd_from = status.MPI_SOURCE;
```

MPI\_Get\_count(&status, datatype, &recvd\_count);

### More info

- 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.

## 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

### **Collective Communications**

- Point-to-point communications involve pairs of processes.
- Many message passing systems provide operations which allow larger numbers of processes to participate

## **Types of Collective Transfers**

- Barrier
  - Synchronizes processors
  - No data is exchanged but the barrier blocks until all processes have called the barrier routine
- Broadcast (sometimes multicast)
  - A broadcast is a one-to-many communication
  - One processor sends one message to several destinations
- Reduction
  - Often useful in a many-to-one communication

#### Barrier



#### **Broadcast and Multicast**

**Broadcast** 

**Multicast** 



#### All-to-All



#### Reduction



#### Introduction to Collective Operations in MPI

- Collective ops are called by all processes in a communicator.
  - No tags
  - Blocking
- MPI\_BCAST distributes data from one process (the root) to all others in a communicator.
- MPI\_REDUCE/ALLREDUCE combines data from all processes in communicator and returns it to one process.
- In many numerical algorithms, SEND/RECEIVE can be replaced by BCAST/REDUCE, improving both simplicity and efficiency.
- Others:
  - MPI\_[ALL] SCATTER[V] / [ALL] GATHER[V]





## Collectives at Work (2)

• Reduce: RANK



Predefined Ops (assocociative & commutative) / user ops (assoc.)

MPI Name	Function
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_BXOR	Bitwise exclusive OR
MPI_MAXLOC	Maximum & location
MPI_MINLOC	Minimum & location

#### Collectives at Work (3)

• Allreduce:



### Simple full example

```
#include <stdio.h>
#include <mpi.h>
int main(int argc, char *argv[])
{
 const int tag = 42; /* Message tag */
 int id, ntasks, source id, dest id, err, i;
 MPI Status status;
 int msg[2]; /* Message array */
 err = MPI Init(&argc, &argv); /* Initialize MPI */
  if (err != MPI SUCCESS) {
   printf("MPI initialization failed!\n");
   exit(1);
  }
 err = MPI Comm size(MPI COMM WORLD, &ntasks); /* Get nr of tasks */
 err = MPI Comm rank(MPI COMM WORLD, &id); /* Get id of this process */
  if (ntasks < 2) {
   printf("You have to use at least 2 processors to run this program\n");
   MPI Finalize(); /* Quit if there is only one processor */
   exit(0);
  }
```

### Simple full example (Cont.)

```
if (id == 0) { /* Process 0 (the receiver) does this */
   for (i=1; i<ntasks; i++) {</pre>
     err = MPI Recv(msg, 2, MPI INT, MPI ANY SOURCE, tag, MPI COMM WORLD, \
                   &status); /* Receive a message */
     source id = status.MPI SOURCE; /* Get id of sender */
     printf("Received message %d %d from process %d\n", msg[0], msg[1], \
           source id);
 else { /* Processes 1 to N-1 (the senders) do this */
   msg[0] = id; /* Put own identifier in the message */
   msg[1] = ntasks; /* and total number of processes */
   dest id = 0; /* Destination address */
   err = MPI Send(msg, 2, MPI INT, dest id, tag, MPI COMM WORLD);
 if (id==0) printf("Ready\n");
 exit(0);
 return 0;
```

#### MPI One-to-one Communication

#### Assynchronous/Non-Blocking

- Process signs it is waiting for a message
- Continue working meanwhile



Time

#### MPI Collective Communication

#### Proccess master wants to send a message to everybody

- First solution, process master send N-1 messages
- Optimized collective communication send in parallel



Finishes in 2 slices of time

## Work@class



## Example: Compute PI (1)

#include "mpi.h"
#include <math.h>

```
int main(int argc, char *argv[])
 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)
 }
```

## Example: Compute PI (2)

```
h = 1.0 / (double)n;
sum = 0.0;
for (i = myid + 1; i <= 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;
```

}

## Profiling Support: PMPI

- Profiling layer of MPI
- Implemented via additional API in MPI library
  - Different name: PMPI\_Init()
  - Same functionality as MPI\_Init()
- Allows user to:
  - define own MPI\_Init()
  - Need to call PMPI\_Init():

```
MPI_Init(...) {
   collect pre stats;
   PMPI_Init(...);
   collect post stats;
}
```

- User may choose subset of MPI routines to be profiled
- Useful for building performance analysis tools
  - Vampir: Timeline of MPI traffic (Etnus, Inc.)
  - Paradyn: Performance analysis (U. Wisconsin)
  - mpiP: J. Vetter (LLNL)
  - ScalaTrace: F. Mueller et al. (NCSU)

### When to use MPI

- Portability and Performance
- Irregular data structure
- Building tools for others
- Need to manage memory on a per processor basis

## Summary

- The parallel computing community has cooperated on the development of a standard for message-passing libraries.
- There are many implementations, on nearly all platforms.
- MPI subsets are easy to learn and use.
- Lots of MPI material is available.

## Para Observar y Ejecutar

- <u>http://people.sc.fsu.edu/~jburkardt/cpp\_src/mpi/mpi.html</u>
- <u>http://www.slac.stanford.edu/comp/unix/farm/mpi.html</u>
- <u>http://www.mcs.anl.gov/research/projects/mpi/usingmpi/</u> <u>examples/main.htm</u>