**ECE-451/ECE-566 - Introduction to Parallel and Distributed Programming**

# Message Passing Computing

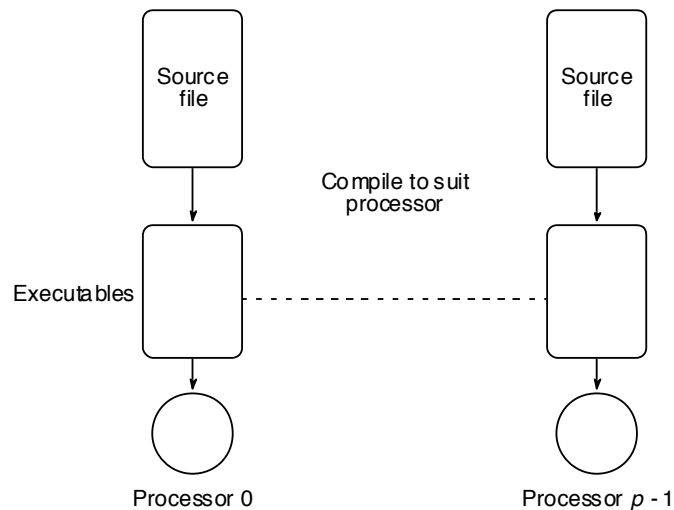Department of Electrical & Computer Engineering

Rutgers University

---

**Message-Passing Programming using User-level Message Passing Libraries**

Two primary mechanisms needed:

1. A method of creating separate processes for execution on different computers

2. A method of sending and receiving messages

2

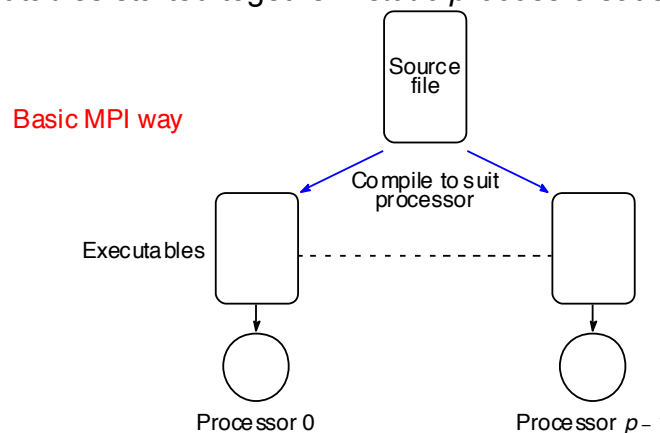# Multiple program, multiple data (MPMD) model

Source
file

Source
file

Compile to suit
processor

Executables

Processor 0

Processor $p$ - 1

3

# Single Program Multiple Data (SPMD) model

Different processes merged into one program. Control
statements select different parts for each processor to execute.
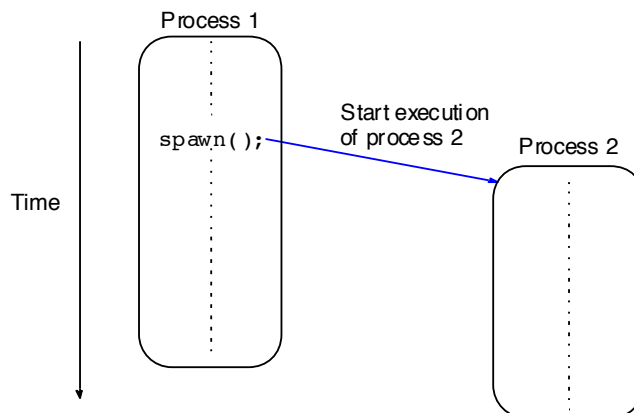All executables started together - *static process* creation

Basic MPI way

Source
file

Compile to suit
processor

Executables

Processor 0

Processor $p - 1$

4

# Multiple Program Multiple Data (MPMD) Model

Separate programs for each processor. One processor executes master process. Other processes started from within master process - *dynamic process* creation.
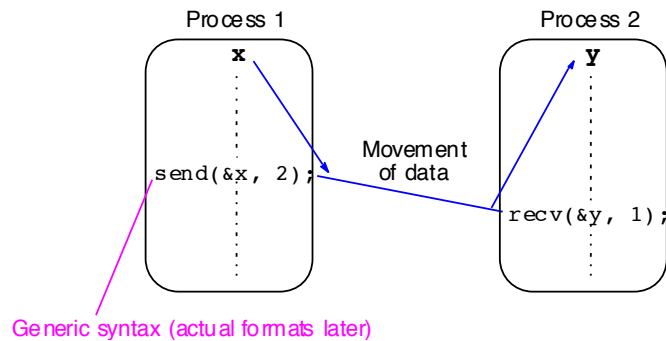
Process 1

```
spawn();
```

Start execution of process 2

Process 2

Time

5

# Basic "point-to-point" Send and Receive Routines

Passing a message between processes using send() and recv() library calls:

Process 1

**x**

```
send(&x, 2);
```

Movement of data

Process 2

**y**

```
recv(&y, 1);
```

Generic syntax (actual formats later)

6

# Synchronous Message Passing

Routines that actually return when message transfer completed.

## Synchronous send routine

• Waits until complete message can be accepted by the receiving process before sending the message.
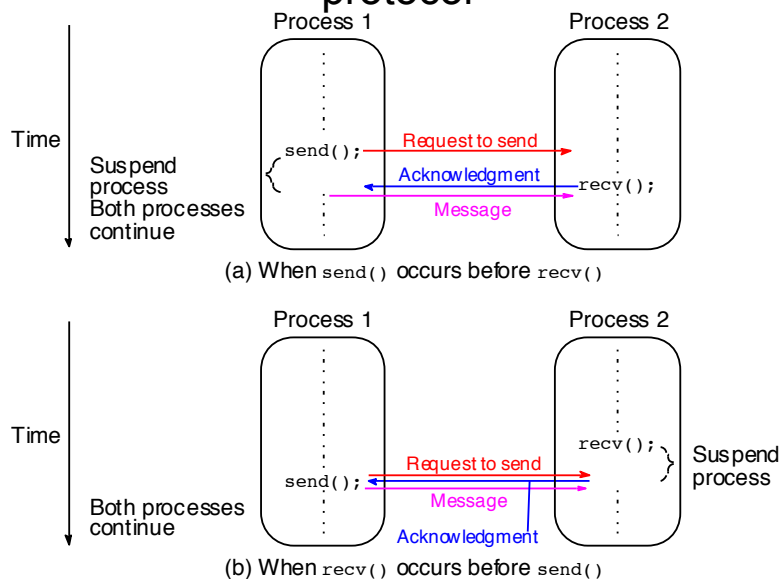
## Synchronous receive routine

• Waits until the message it is expecting arrives.

Synchronous routines intrinsically perform two actions:
They transfer data and they synchronize processes.

7

# Synchronous send() and recv() using 3-way protocol



(a) When send() occurs before recv()

(b) When recv() occurs before send()

8

4

# Asynchronous Message Passing

- Routines that do not wait for actions to complete before returning. Usually require local storage for messages.

- More than one version depending upon the actual semantics for returning.

- In general, they do not synchronize processes but allow processes to move forward sooner. Must be used with care.

9

# MPI Definitions of Blocking and Non-Blocking

- Blocking - return after their local actions complete, though the message transfer may not have been completed.
- Non-blocking - return immediately.

  Assumes that data storage used for transfer not modified by subsequent statements prior to being used for transfer, and it is left to the programmer to ensure this.
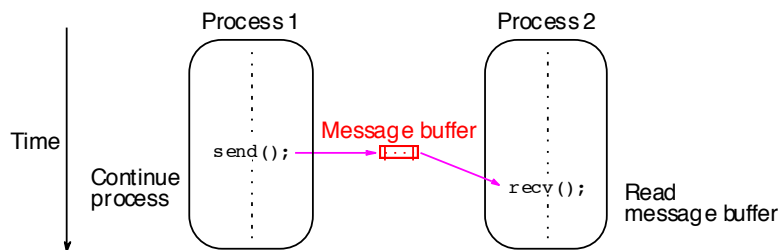
  These terms may have different interpretations in other systems.

10

# How message-passing routines return before message transfer completed

Message buffer needed between source and destination to hold message:

11

# Asynchronous (blocking) routines changing to synchronous routines

- Once local actions completed and message is safely on its way, sending process can continue with subsequent work.
- Buffers only of finite length and a point could be reached when send routine held up because all available buffer space exhausted.
- Then, send routine will wait until storage becomes re-available - i.e then routine behaves as a synchronous routine.
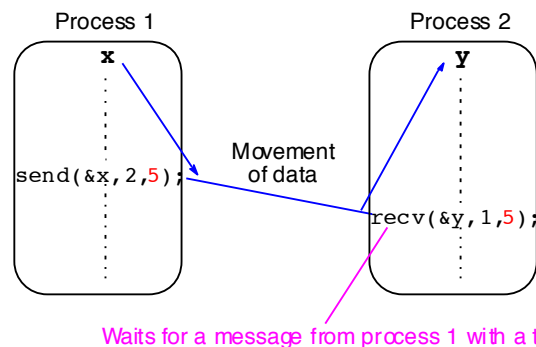
12

# Message Tag

- Used to differentiate between different types of messages being sent.

- Message tag is carried within message.

- If special type matching is not required, a wild card message tag is used, so that the recv() will match with any send().

13

# Message Tag Example

To send a message, x, with message tag 5 from a source process, 1, to a destination process, 2, and assign to y:



Process 1

Process 2

x

y

Movement
of data

send(&x,2,5);

recv(&y,1,5);

Waits for a message from process 1 with a tag of 5

14

# "Group" message passing routines

Have routines that send message(s) to a group of processes or receive message(s) from a group of processes

Higher efficiency than separate point-to-point routines although not absolutely necessary.
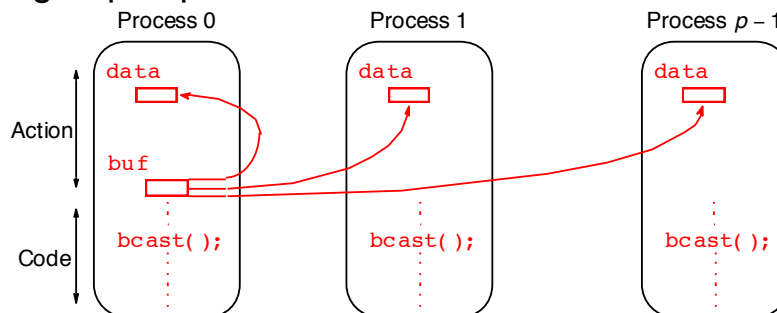
15

# Broadcast

Sending same message to all processes concerned with problem.
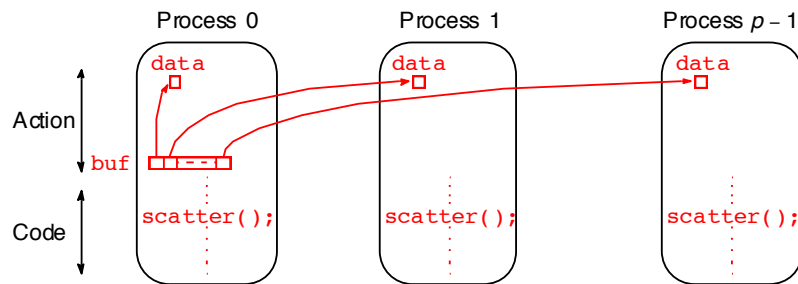Multicast - sending same message to defined group of processes.



MPI form

16

8

# Scatter

Sending each element of an array in root process to a separate process. Contents of *i*th location of array sent to *i*th process.

Process 0     Process 1     Process *p* − 1

data          data          data

Action

buf

Code

scatter();    scatter();    scatter();

MPI form

17

# Gather

Having one process collect individual values from set of processes.

Process 0     Process 1     Process *p* − 1

data          data          data

Action

buf

Code

gather();     gather();     gather();

MPI form

18

# Reduce

Gather operation combined with specified arithmetic/logical operation.

Example: Values could be gathered and then added together by root

19

# MPI
# (Message Passing Interface)

- Message passing library standard developed by group of academics and industrial partners to foster more widespread use and portability.

- Defines routines, not implementation.

- Several free implementations exist.

22

# MPI
## Process Creation and Execution

- Purposely not defined - Will depend upon implementation.
- Only static process creation supported in MPI version 1. All processes must be defined prior to execution and started together.
- Originally SPMD model of computation.
- MPMD also possible with static creation - each program to be started together specified.

23

# Communicators

- Defines scope of a communication operation.
- Processes have ranks associated with communicator.
- Initially, all processes enrolled in a "universe" called MPI_COMM_WORLD, and each process is given a unique rank, a number from 0 to $p$ - 1, with $p$ processes.
- Other communicators can be established for groups of processes.

24

# Using SPMD Computational Model

```
main (int argc, char *argv[])
{
MPI_Init(&argc, &argv);
        .
        .
MPI_Comm_rank(MPI_COMM_WORLD, &myrank); /*find process rank */

if (myrank == 0)
        master();
else
        slave();
        .
        .
MPI_Finalize();
}
```
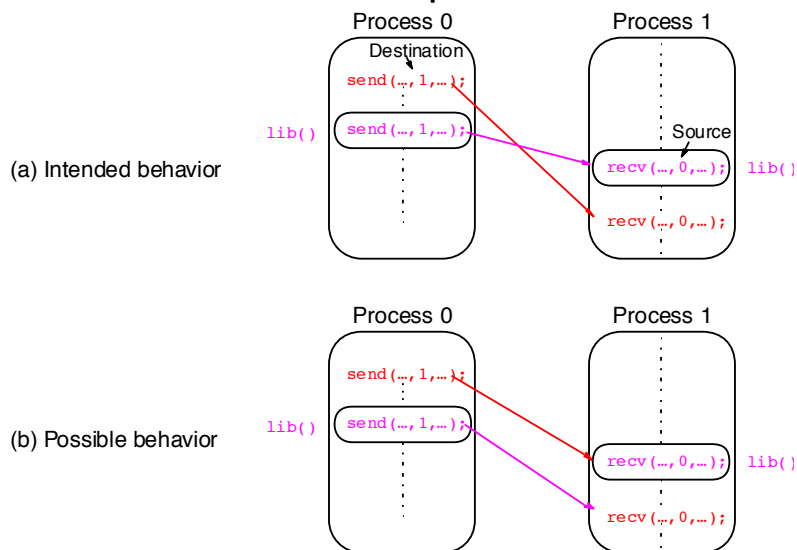
where master() and slave() are to be executed by master
process and slave process, respectively.

25

# Unsafe message passing
## Example

26

12

# MPI Solution
## "Communicators"

- Defines a communication domain - a set of processes that are allowed to communicate between themselves.

- Communication domains of libraries can be separated from that of a user program.

- Used in all point-to-point and collective MPI message-passing communications.

27

# Default Communicator
## MPI_COMM_WORLD

- Exists as first communicator for all processes existing in the application.

- A set of MPI routines exists for forming communicators.

- Processes have a "rank" in a communicator.

28

# MPI Point-to-Point Communication

- Uses send and receive routines with message tags (and communicator).
-  Wild card message tags available

29

# MPI Blocking Routines

- Return when "locally complete" - when location used to hold message can be used again or altered without affecting message being sent.

- Blocking send will send message and return - does not mean that message has been received, just that process free to move on without adversely affecting message.

30

# Parameters of blocking send

`MPI_Send(buf, count, datatype, dest, tag, comm)`

Address of
send buffer

Datatype of
each item

Message tag

Number of items
to send

Rank of destination
process

Communicator

31

# Parameters of blocking receive

`MPI_Recv(buf, count, datatype, src, tag, comm, status)`

Address of
receive buffer

Datatype of
each item

Message tag

Status
after operation

Maximum number
of items to receive

Rank of source
process

Communicator

32

15

# Example

To send an integer x from process 0 to process 1,

```
MPI_Comm_rank(MPI_COMM_WORLD,&myrank); /* find rank */

if (myrank == 0) {
  int x;
  MPI_Send(&x, 1, MPI_INT, 1, msgtag, MPI_COMM_WORLD);
} else if (myrank == 1) {
  int x;
  MPI_Recv(&x, 1, MPI_INT,
  0,msgtag,MPI_COMM_WORLD,status);
}
```

33

# MPI Nonblocking Routines

- Nonblocking send - MPI_Isend() - will return "immediately" even before source location is safe to be altered.

- Nonblocking receive - MPI_Irecv() - will return even if no message to accept.

34

16

# Nonblocking Routine Formats

`MPI_Isend(buf,count,datatype,dest,tag,comm,request)`

`MPI_Irecv(buf,count,datatype,source,tag,comm, request)`

Completion detected by `MPI_Wait()` and `MPI_Test()`.

`MPI_Wait()` waits until operation completed and returns then.

`MPI_Test()` returns with flag set indicating whether operation completed at that time.

Need to know whether particular operation completed.

Determined by accessing `request` parameter.

35

# Example

To send an integer x from process 0 to process 1 and allow process 0 to continue,

```
MPI_Comm_rank(MPI_COMM_WORLD, &myrank); /* find rank */
if (myrank == 0) {
    int x;
    MPI_Isend(&x,1,MPI_INT, 1, msgtag, MPI_COMM_WORLD, req1);
    compute();
    MPI_Wait(req1, status);
} else if (myrank == 1) {
    int x;
    MPI_Recv(&x,1,MPI_INT,0,msgtag, MPI_COMM_WORLD, status);
}
```

36

# Four Send Communication Modes

- **Standard Mode Send** - Not assumed that corresponding receive routine has started. Amount of buffering not defined by MPI. If buffering provided, send could complete before receive reached.
- **Buffered Mode** - Send may start and return before a matching receive. Necessary to specify buffer space via routine MPI_Buffer_attach().
- **Synchronous Mode** - Send and receive can start before each other but can only complete together.
- **Ready Mode** - Send can only start if matching receive already reached, otherwise error. Use with care.

37

# Communication Modes

- Each of the four modes can be applied to both blocking and nonblocking send routines.

- Only the standard mode is available for the blocking and nonblocking receive routines.

- Any type of send routine can be used with any type of receive routine.

38

# Collective Communication

Involves set of processes, defined by an intra-communicator.
Message tags not present. Principal collective operations:

- **MPI_Bcast()**       - Broadcast from root to all other processes
- **MPI_Gather()**       - Gather values for group of processes
- **MPI_Scatter()**      - Scatters buffer in parts to group of processes
- **MPI_Alltoall()**     - Sends data from all processes to all processes
- **MPI_Reduce()**       - Combine values on all processes to single value
- **MPI_Reduce_scatter()**    - Combine values and scatter results
- **MPI_Scan()**        - Compute prefix reductions of data on processes

39

# Example

To gather items from group of processes into process 0, using
dynamically allocated memory in root process:

```
int data[10];                    /*data to be gathered from processes*/
MPI_Comm_rank(MPI_COMM_WORLD, &myrank);      /* find rank */
if (myrank == 0) {
   MPI_Comm_size(MPI_COMM_WORLD, &grp_size); /*find group size*/
   buf = (int *)malloc(grp_size*10*sizeof (int)); /*allocate memory*/
}
MPI_Gather(data,10,MPI_INT,buf,grp_size*10,MPI_INT,0,MPI_COMM_WORLD;
```

**MPI_Gather()** gathers from all processes, including root.

40

# Barrier

- As in all message-passing systems, MPI provides a means of synchronizing processes by stopping each one until they all have reached a specific "barrier" call.

41

```
#include "mpi.h"
#include <stdio.h>
#include <math.h>
#define MAXSIZE 1000
void main(int argc, char *argv)                     Sample MPI program
{
    int myid, numprocs;
    int data[MAXSIZE], i, x, low, high, myresult, result;
    char fn[255];
    char *fp;
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);
    if (myid == 0) {  /* Open input file and initialize data */
            strcpy(fn,getenv("HOME"));
            strcat(fn,"/MPI/rand_data.txt");
            if ((fp = fopen(fn,"r")) == NULL) {
                        printf("Can't open the input file: %s\n\n", fn);
                        exit(1);
            }
            for(i = 0; i < MAXSIZE; i++) fscanf(fp,"%d", &data[i]);
    }
    MPI_Bcast(data, MAXSIZE, MPI_INT, 0, MPI_COMM_WORLD); /* broadcast data */
    x = n/nproc; /* Add my portion Of data */
    low = myid * x;
    high = low + x;
    for(i = low; i < high; i++)
            myresult += data[i];
    printf("I got %d from %d\n", myresult, myid); /* Compute global sum */
    MPI_Reduce(&myresult, &result, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
    if (myid == 0) printf("The sum is %d.\n", result);
    MPI_Finalize();
}
```

42

20