



MPI Programming

Intel Software College



Course Objectives

After successfully completing this module, you will be able to:

- Compile and run MPI programs using the MPICH implementation
- Write MPI programs using the core library

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

Advanced MPI

- Useful Features
- Factors Affecting MPI Performance

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

Advanced MPI

- Useful Features
- Factors Affecting MPI Performance

MPI

A library - not a language

A library for inter-process communication and data exchange

Function categories

- Initialization/finalization
- Point-to-point communication
- Collective communication
- Communicator topologies
- User-defined data types
- Utilities (e.g.- timing and initialization)

Parallel Models Compared

	MPI	Threads	OpenMP*
Portable	✓		✓
Scalable	✓	✓	✓
Performance Oriented	✓		✓
Supports Data Parallel	✓	✓	✓
Incremental Parallelism			✓
High Level			✓
Serial Code Intact			✓
Verifiable Correctness			✓
Distributed Memory	✓		



Common MPI Implementations

MPICH* (Argonne National Laboratory)

- Most common MPI implementation
- Derivatives
 - MPICH GM* – Myrinet* support (available from Myricom)
 - MVAPICH* – InfiniBand* support (available from Ohio State University)
 - Intel® MPI – version tuned to Intel Architecture systems
 - LAM/MPI* (Indiana University)
 - Contains many MPI 2.0 features
 - Daemon-based MPI implementation
 - MPI/Pro* (MPI Software Technology)
 - Scali MPI Connect*
 - Provides native support for most high-end interconnects

Getting Started

OSCAR? NPACI Rocks?

- MPI Implementation(s) pre-configured
- For a manual build see material in backup to these slides

What do you need to run?

- Program – supplied or build with
`mpicc <program.c>` or
`mpif77 <program.f>`
- Which machines to use? (machines file)
- A way to launch programs on remote machines (ssh for both OSCAR and Rocks)

Run with:

- `mpirun -n xx <program name>`
where `xx` is the number of processes (ranks) to launch

Absolute Minimum MPI Program Skeleton

Initialize MPI Lib on my node – required

Do work – optional, but important 😊

Wrap it up once we are finished – required

Compiling an MPI Program

Most MPI implementations supply compilation scripts:

```
mpif77 mpi_prog.f  
mpicc mpi_proc.c  
mpif90 mpi_prof.f90  
mpiCC mpi_prof.C
```

Manual compilation/linking also possible:

```
ifc mpi_prog.f -L/usr/local/mpich-1.2.5/lib -lfmpich -lmpich
```

MPI Machine File

A text file telling MPI where to launch processes

Put separate machine name on each line

Example:

```
compute-0-0  
compute-0-1  
compute-0-2  
compute-0-3
```

Check implementation for multi-processor node formats

Default file found in MPI installation

Parallel Program Execution

Launch scenario for MPIRUN

- Find machine file (to know where to launch)
- Use SSH or RSH to execute a copy of program on each node in machine file
- Once launched, each copy establishes communication with local MPI lib (MP_Init)
- Each copy ends MPI interaction with MPI_Finalize

Starting an MPI Program

1 Execute on node1:
`mpirun -np 4 mpi_prog`

2 Check the MPICH machine file:
node1
node2
node3
node4
node5
node6

3

node1
mpi_prog (rank 0)

node2
mpi_prog (rank 1)

node3
mpi_prog (rank 2)

node4
mpi_prog (rank 3)

Activity 1: "Hello, World" in MPI

Demonstrates how to create, compile and run a simple MPI program on the lab cluster using the Intel MPI implementation

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

int main (int argc, char* argv[])
{
    MPI_Init (&argc, &argv); /* Initialize the library */

    printf ("Hello world\n");

    MPI_Finalize (); /* Wrap it up. */
}
```

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

Advanced MPI

- Useful Features
- Factors Affecting MPI Performance

The Whole Library

MPI is large and complex

- MPI 1.0 – 125 functions
- MPI 2.0 is even larger

But, many MPI features are rarely used

- Inter-communicators
- Topologies
- Persistent communication
- Functions designed for library developers

The Absolute Minimum

Six MPI functions:

- `MPI_Init` ★
- `MPI_Comm_size`
- `MPI_Comm_rank`
- `MPI_Send`
- `MPI_Recv`
- `MPI_Finalize` ★

Many parallel algorithms can be implemented efficiently with only these functions

Initializing the MPI Library

Fortran:

```
MPI_INIT (ierr)
```

C:

```
int MPI_Init (int* argc, char** argv)
```

`MPI_Init` prepares the system for MPI execution

Call to `MPI_Init` may update arguments in C

- Implementation dependent

No MPI functions may be called before `MPI_Init`

Shutting Down MPI

Fortran:

```
MPI_FINALIZE (ierr)
```

C:

```
int MPI_Finalize (void)
```

`MPI_Finalize` frees any memory allocated by the MPI library

No MPI functions may be called after calling `MPI_Finalize`

- Exception: `MPI_Init`

Sizing the MPI Communicator

Fortran:

```
MPI_COMM_SIZE (comm, size, ierr)
integer :: comm, size, ierr
```

C:

```
int MPI_Comm_size (
    MPI_Comm comm, int* size)
```

`MPI_Comm_size` returns the number of processes in the specified communicator

The communicator structure, `MPI_Comm`, is defined in `mpi.h`

Determining MPI Process Rank

Fortran:

```
MPI_COMM_RANK (comm, rank, ierr)  
integer :: comm, rank, ierr
```

C:

```
int MPI_Comm_rank (  
    MPI_Comm comm, int* rank)
```

`MPI_Comm_rank` returns the rank of calling process within the specified communicator

Processes are numbered from 0 to N-1

Activity 2: "Hello, World" with IDs

Demonstrates how to modify, compile and run a simple MPI program on the lab cluster using the Intel MPI implementation

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

int main (int argc, char* argv[])
{
    int numProc, myRank;

    MPI_Init (&argc, &argv); /* Initialize the library */
    MPI_Comm_rank (MPI_COMM_WORLD, &myRank); /* Who am I?" */
    MPI_Comm_size (MPI_COMM_WORLD, &numProc); /*How many? */

    printf ("Hello. Process %d of %d here.\n", myRank, numProc);

    MPI_Finalize (); /* Wrap it up. */
}
```

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

Advanced MPI

- Useful Features
- Factors Affecting MPI Performance

Sending Data

Fortran:

```
MPI_SEND (buf, count, datatype,  
          dest, tag, comm, ierr)  
  
<type> buf (*)  
integer :: count, datatype, ierr,  
          dest, tag, comm
```

C:

```
int MPI_Send (void* buf, int count,  
             MPI_Datatype datatype, int dest,  
             int tag, MPI_Comm comm)
```

`MPI_Send` performs a blocking send of the specified data to the specified destination

Receiving Data

Fortran:

```
MPI_RECV (buf, count, datatype, source,  
          tag, comm, status, ierr)  
<type> buf (*)  
integer :: count, datatype, ierr, source,  
          tag, comm,  
          status(MPI_STATUS_SIZE)
```

C:

```
int MPI_Recv (void* buf, int count,  
             MPI_Datatype datatype, int source,  
             int tag, MPI_Comm comm,  
             MPI_Status* status)
```

`MPI_Recv` performs a blocking receive of the specified data from the specified source

MPI Data Types for C

MPI Data Type	C Data Type
MPI_BYTE	
MPI_CHAR	signed char
MPI_DOUBLE	double
MPI_FLOAT	float
MPI_INT	int
MPI_LONG	long
MPI_LONG_DOUBLE	long double
MPI_PACKED	
MPI_SHORT	short
MPI_UNSIGNED_SHORT	unsigned short
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long
MPI_UNSIGNED_CHAR	unsigned char

MPI provides predefined data types that must be specified when passing messages.

MPI Data Types for Fortran

MPI Data Type	Fortran Data Type
MPI_BYTE	
MPI_CHARACTER	CHARACTER
MPI_COMPLEX	COMPLEX
MPI_DOUBLE_PRECISION	DOUBLE PRECISION
MPI_INTEGER	INTEGER
MPI_LOGICAL	LOGICAL
MPI_PACKED	
MPI_REAL	REAL

MPI provides predefined data types that must be specified when passing messages.

MPI Status Details

C

type `MPI_Status` has three fields

- `MPI_SOURCE`
- `MPI_TAG`
- `MPI_ERROR`

```
status.MPI_SOURCE
```

```
status.MPI_TAG
```

```
status.MPI_ERROR
```

Fortran

Status is an array indexed by

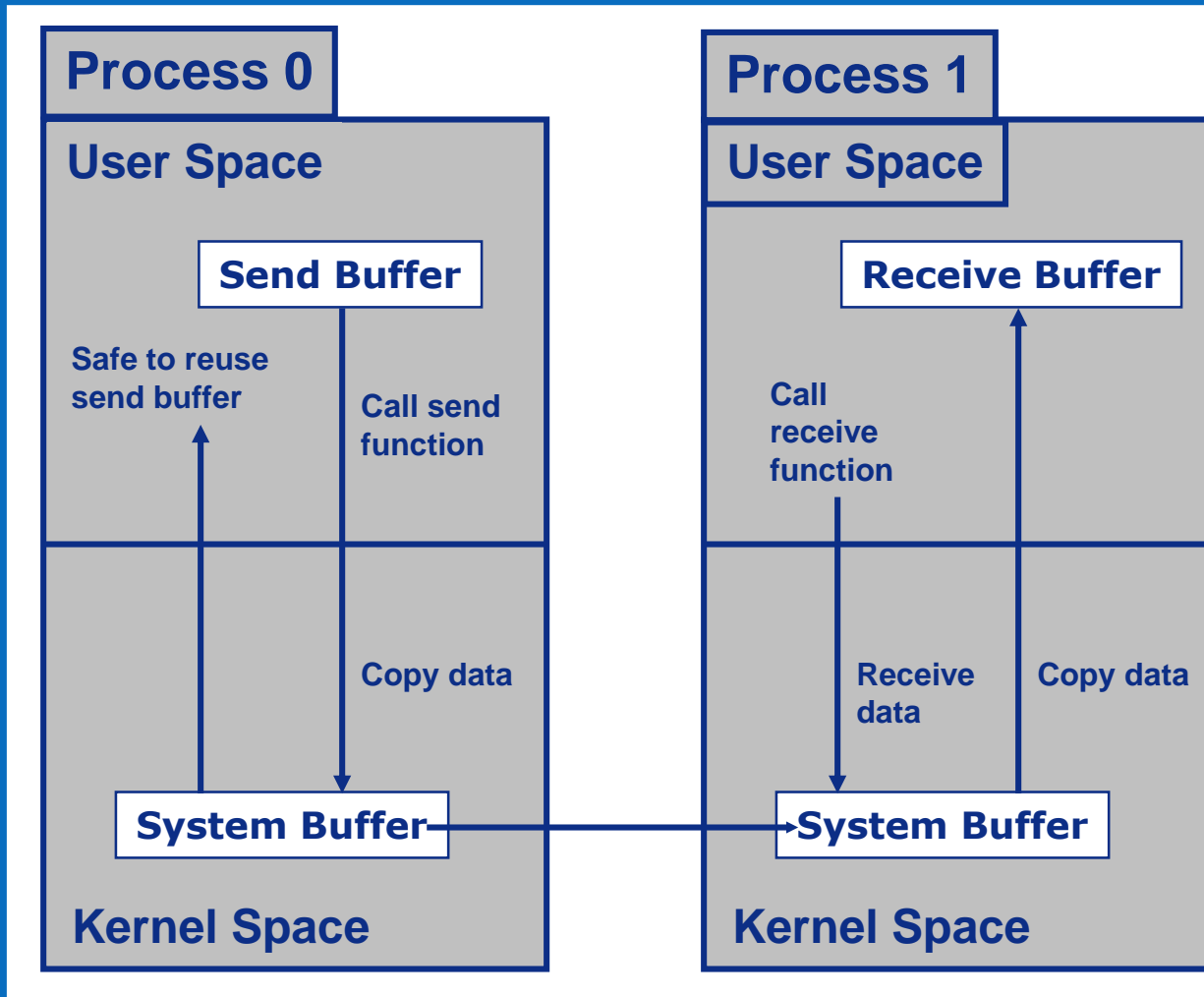
- `MPI_SOURCE`
- `MPI_TAG`
- `MPI_ERROR`

```
status(MPI_SOURCE)
```

```
status(MPI_TAG)
```

```
status(MPI_ERROR)
```

Blocking Communication



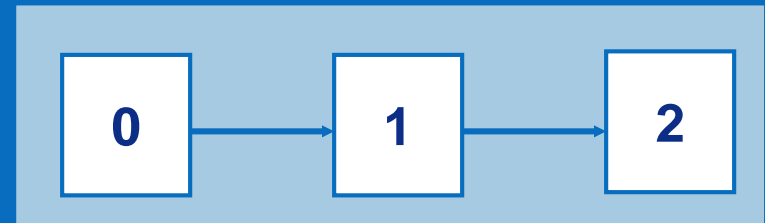
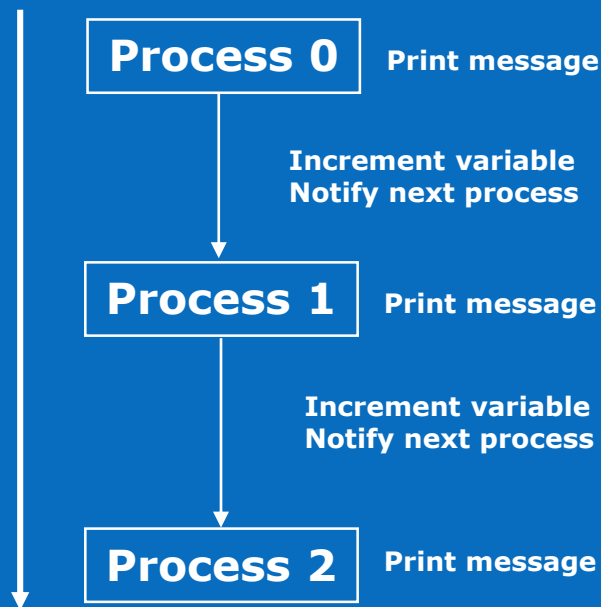
Key Points

- Sending and receiving data is a paired operation. Ignoring this principle can result in deadlock.
- Communication and synchronization are coupled.

Activity 3: "Hello, World" with Message

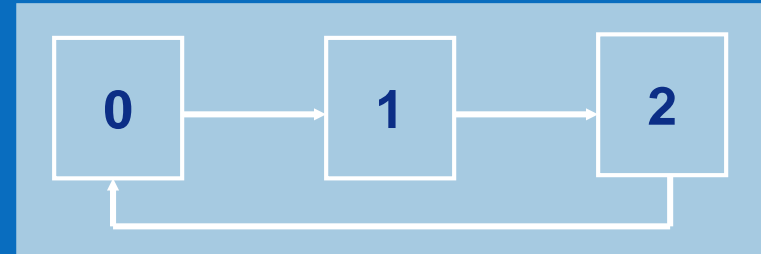
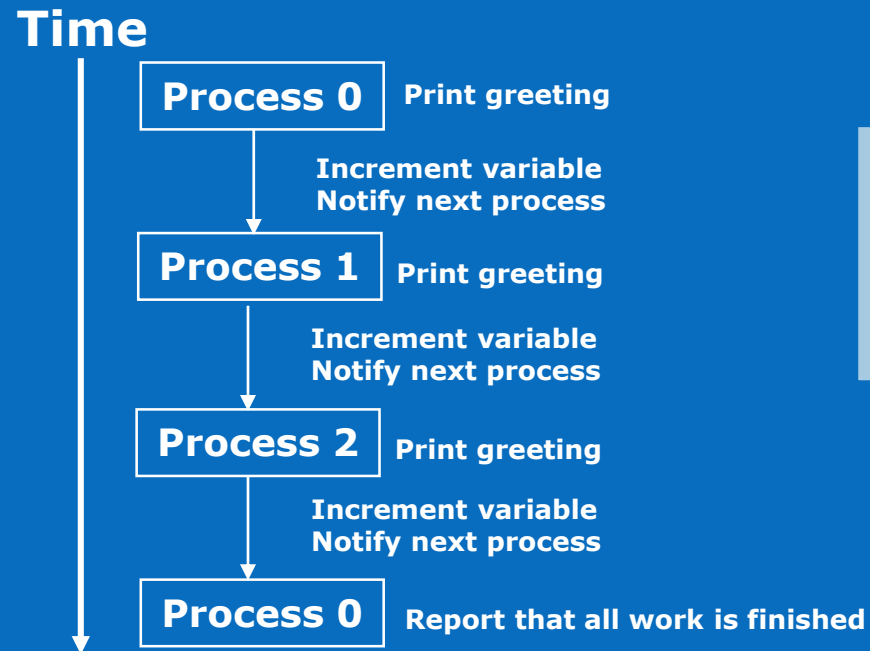
Modify the "Hello, World" example so that the MPI processes pass an integer variable from left to right, reporting and incrementing its value at each step.

Time



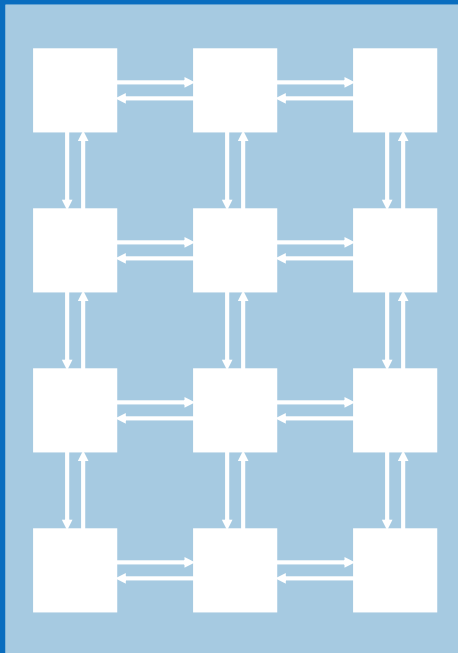
Activity 4: "Hello, World" with Finish

Modify the "Hello, World" example so that the MPI processes pass an integer variable from left to right, reporting and incrementing its value at each step. The master process should report when all processes are finished.



Thought Problem

Consider a 2D mesh with nearest-neighbor communication, coded in MPI:



```
do i = 1, neighbors
    call MPI_RECV (edge(1,i), len, MPI_REAL,
nbr(i),
                                tag, comm, requests(i),
ierr)
enddo

do i = 1, neighbors
    call MPI_SEND (edge(1,i), len, MPI_REAL,
nbr(i),
                                tag, comm, ierr)
enddo
```

What is wrong with this code?

Definite Deadlock

```
do i = 1, neighbors
  call MPI_RECV (edge(1,i), len, MPI_REAL, nbr(i),
                tag, comm, requests(i), ierr)
enddo

do i = 1, neighbors
  call MPI_SEND (edge(1,i), len, MPI_REAL, nbr(i),
                tag, comm, ierr)
enddo
```

This program should always deadlock because each MPI process executes a blocking 'receive' before the corresponding 'send'.

Possible Deadlock

```
do i = 1, neighbors
  call MPI_SEND (edge(1,i), len, MPI_REAL, nbr(i),
                tag, comm, ierr)
enddo

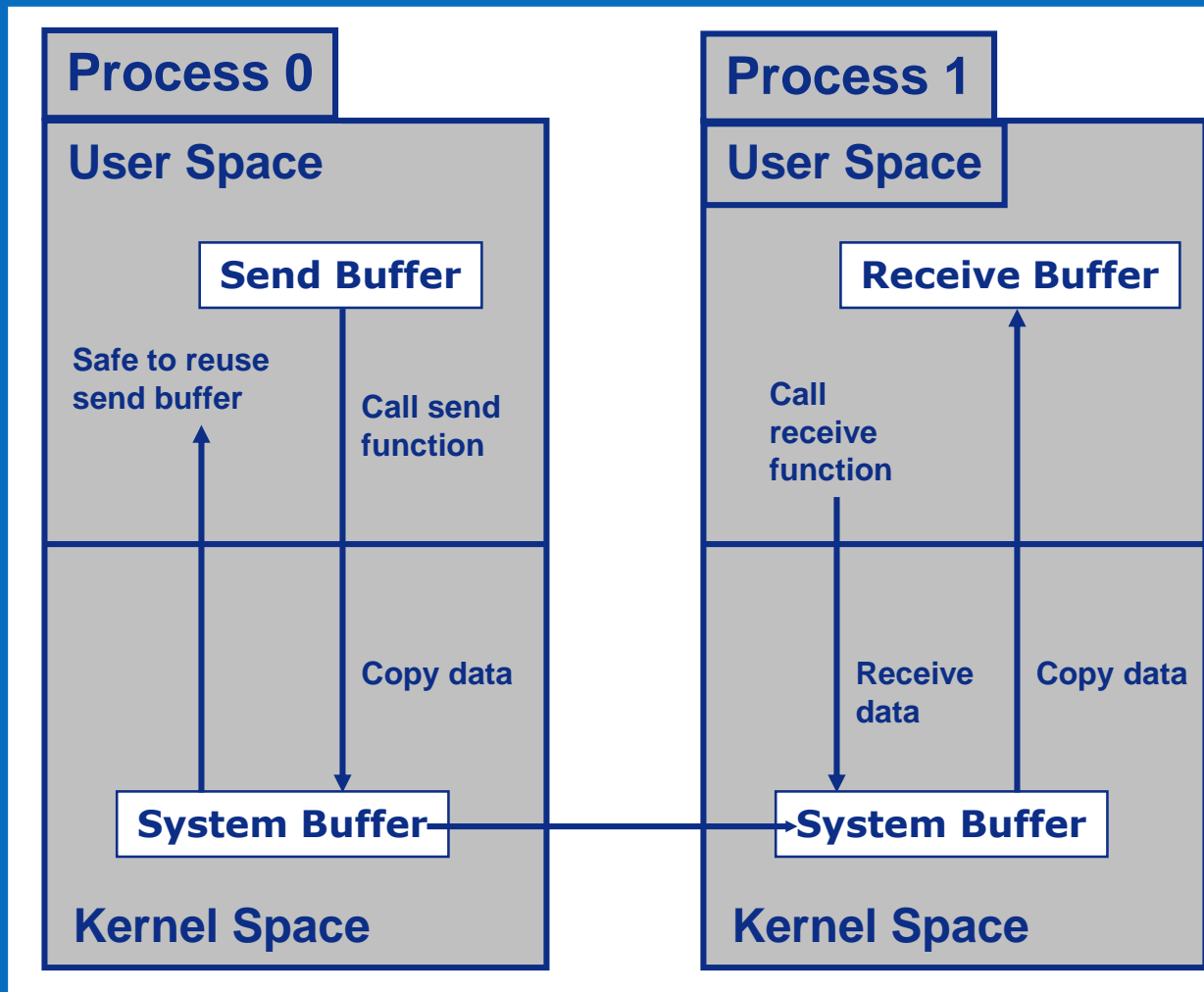
do i = 1, neighbors
  call MPI_RECV (edge(1,i), len, MPI_REAL, nbr(i),
                tag, comm, requests(i), ierr)
enddo
```

This program may deadlock.

In some MPI implementations, MPI_SEND returns when the send buffer is copied to the system buffer. Result:

- Program might work if messages are small, but would deadlock when messages are larger than the system buffer

Worth Repeating

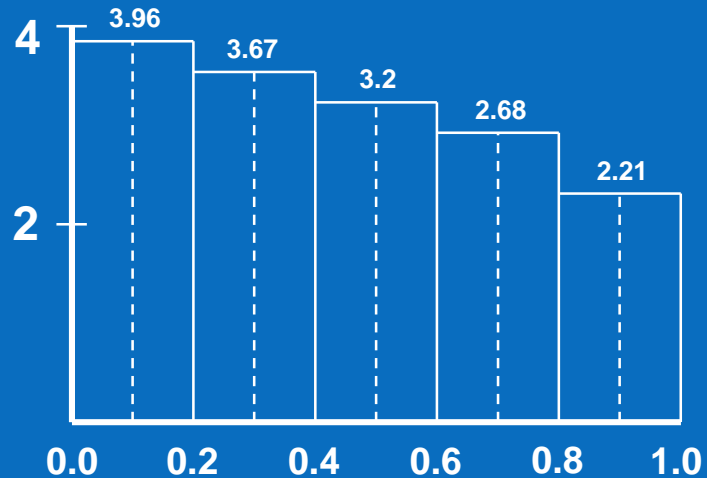


Key Points

- Sending and receiving data is a paired operation.
- **Ignoring this principle can result in deadlock.**
- Communication and synchronization are coupled.

Activity 5: Parallelize Pi Integration

$$f(x) = 4 / (1 + x^2)$$



$$\frac{3.96 + 3.67 + 3.2 + 2.68 + 2.21}{5 \text{ intervals}} = 3.144$$

The serial version

```
#include <stdio.h>
#define INTERVALS 100000

int main (int argc, char* argv[])
{
    int i;
    double h, x, pi = 0.0;

    h = 1.0 / (double)INTERVALS

    for (i = 0; i < INTERVALS; i++)
    {
        x = h * ((double)i - 0.5);
        pi += 4.0 / (1.0 + x * x);
    }
    pi *= h;
    printf ("Pi = %f\n", pi);
}
```

No Deadlock

```

Process 0 code
do
<Prepare data>
    call MPI_SEND (data, len, MPI_REAL, 1,
                  tag, comm, ierr)
    Call MPI_RECV (buff, len, MPI_REAL, 1,
                  tag, comm, ierr)
enddo
Process 1 code
do
    call MPI_RECV (buff, len, MPI_REAL, 0,
                  tag, comm, ierr)
<Prepare data>
    call MPI_SEND (buff, len, MPI_REAL, 0,
                  tag, comm, ierr)
enddo

```

Though this program will not deadlock, there is another issue . . .

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

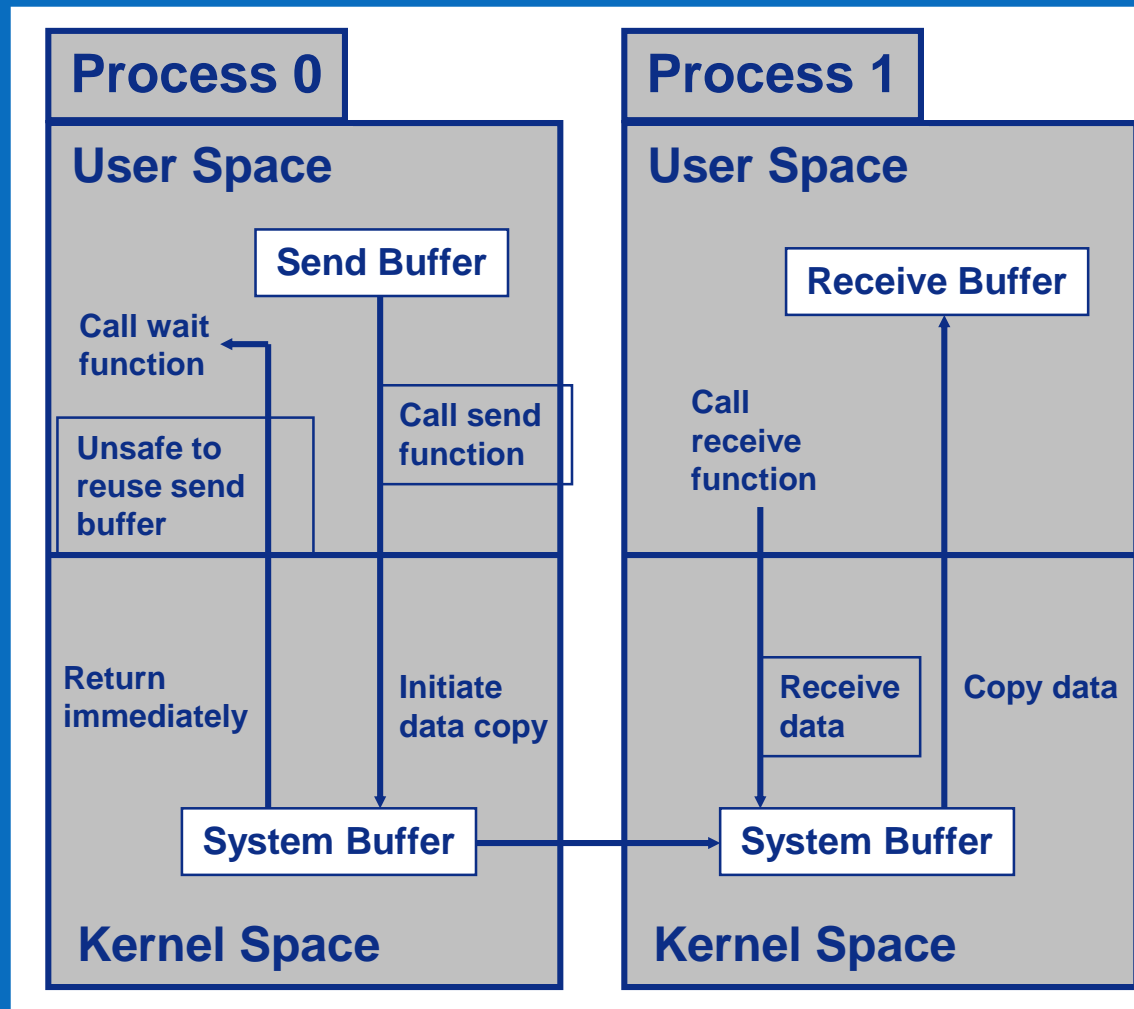
Advanced MPI

- Useful Features
- Factors Affecting MPI Performance

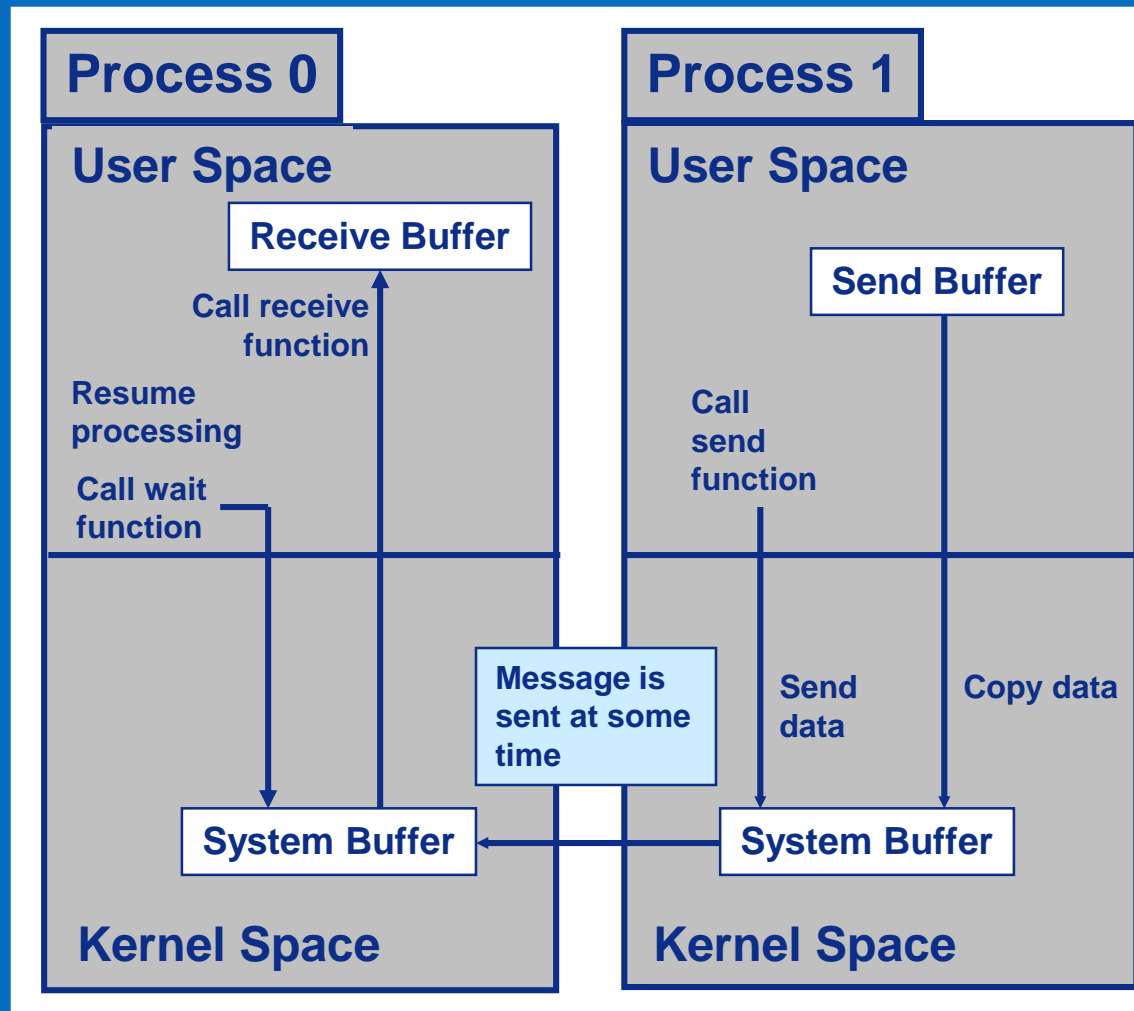
Non-Blocking Send

Key Points

1. The sender does not wait for message buffering to complete.
2. The programmer must guarantee that the send buffer is not modified before the message is delivered.
3. It is difficult to debug MPI programs that ignore the previous point.



Non-Blocking Receive



Key Points

1. The receiver does not wait for message to be received.
2. The program may continue with other processing until message is needed.
3. Call *wait* function to block until message arrives. Cannot guarantee buffer contents until *wait* function returns.

Immediate Send

Fortran:

```
MPI_ISEND (buf, count, datatype, dest,  
           tag, comm, request, ierr)  
<type> buf(*)  
integer :: count, datatype, dest, tag,  
           comm, request, ierr
```

C:

```
int MPI_Isend (void* buf, int count,  
              MPI_Datatype datatype, int dest,  
              int tag, MPI_Comm comm,  
              MPI_Request* request)
```

`MPI_Isend` performs a non-blocking send of the specified data to the specified destination

Immediate Receive

Fortran:

```
MPI_Irecv (buf, count, datatype, source,  
           tag, comm, request, ierr)  
<type> buf (*)  
integer :: count, datatype, ierr, source,  
           tag, comm, request
```

C:

```
int MPI_Irecv (void* buf, int count,  
              MPI_Datatype datatype, int source,  
              int tag, MPI_Comm comm,  
              MPI_Request* request)
```

`MPI_Irecv` performs a non-blocking receive of the specified data from the specified source

Confirming Delivery

Fortran:

```
MPI_WAIT (request, status, ierr)
integer :: request, ierr,
          status(MPI_STATUS_SIZE)
```

C:

```
int MPI_Wait (MPI_Request* request,
              MPI_Status* status)
```

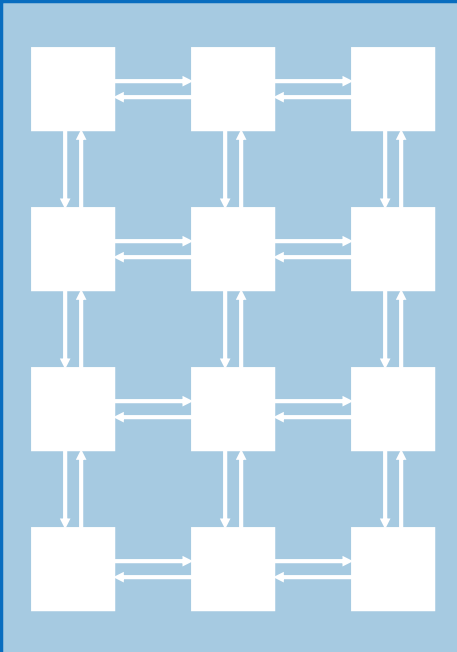
`MPI_Wait` returns when the operation that started the specified request gets completed

`MPI_Send` or `MPI_Recv` can also be used to confirm delivery

`MPI_Waitall` and `MPI_Waitany` can be used to wait on groups of messages without imposing order

Thought Problem

Consider a two-dimensional mesh with nearest-neighbor communication, coded in MPI



```
do i = 1, neighbors
    call MPI_Irecv (edge(1,i), len, MPI_REAL,
nbr(i),
                                tag, comm, requests(i),
ierr)
enddo

do i = 1, neighbors
    call MPI_Send (edge(1,i), len, MPI_REAL,
nbr(i),
                                tag, comm, requests(i),
ierr)
enddo
```

Why does this code perform poorly?

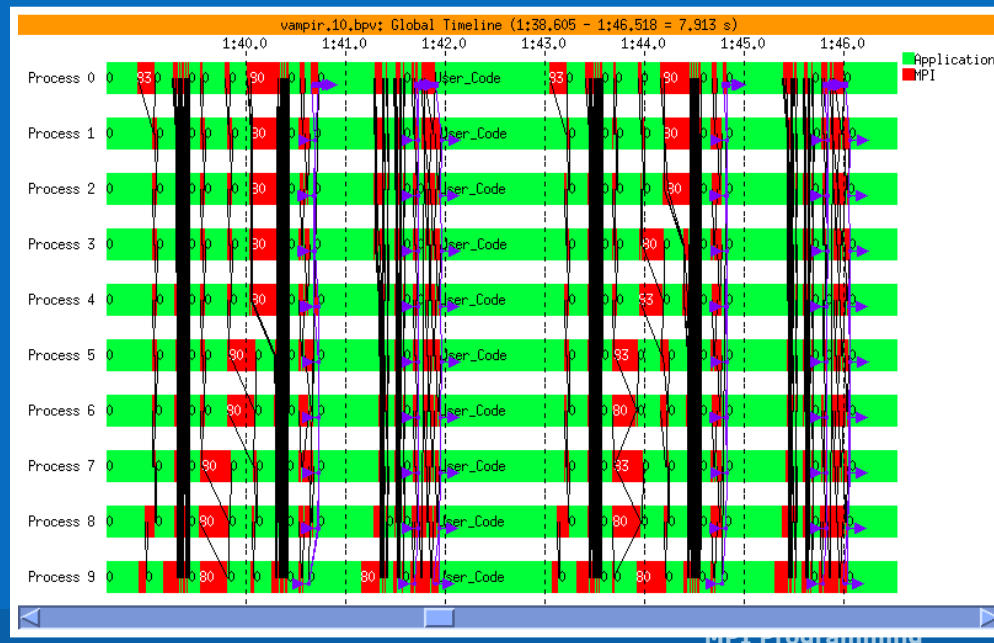
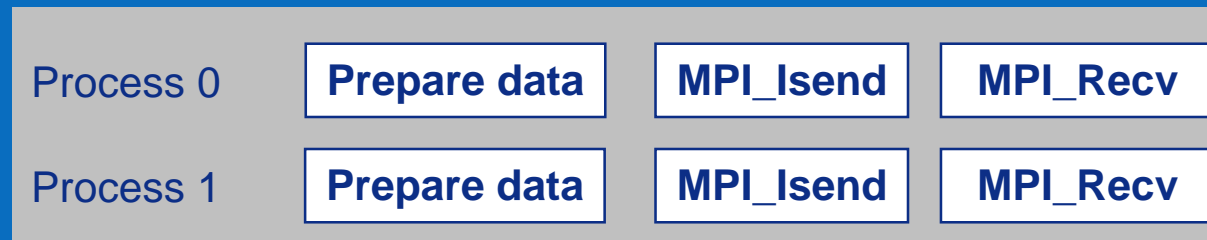
Increased Parallelism

```
Process 0 code
do
<Prepare data>
  call MPI_ISEND (data, len, MPI_REAL, 1,
                 tag, comm, request, ierr)
  call MPI_RECV (buff, len, MPI_REAL, 1,
                tag, comm, ierr)
enddo
Process 1 code
Do
<Prepare data>
  call MPI_ISEND (data, len, MPI_REAL, 0,
                 tag, comm, request, ierr)
  call MPI_RECV (buff, len, MPI_REAL, 0,
                tag, comm, ierr)
enddo
```

How does this run?

Optimizing Communication

Parallel data preparation alleviates the bottleneck



Intel® Trace Analyzer
MPI timeline showing
improved computation
to communication ratio

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

Advanced MPI

- Useful Features
- Factors Affecting MPI Performance

Collective Communication

Functions to facilitate communication among all MPI processes in a communicator

Categories:

- Barrier and broadcast
- Gather and scatter
- Reduction

Collective communication functions can be optimized for the underlying network topology

Broadcasting Data

Fortran:

```
MPI_BCAST (buf, count, datatype, root,  
           comm, ierr)  
<type> buf (*)  
integer :: count, datatype, root, comm,  
         ierr
```

C:

```
int MPI_Bcast (void* buf, int count,  
              MPI_Datatype datatype, int root,  
              MPI_Comm comm)
```

`MPI_Bcast` sends the specified data to all processes in the communicator

MPI_BCAST Example

```

program BroadcastData
include 'mpif.h'

integer, dimension(4) :: msg

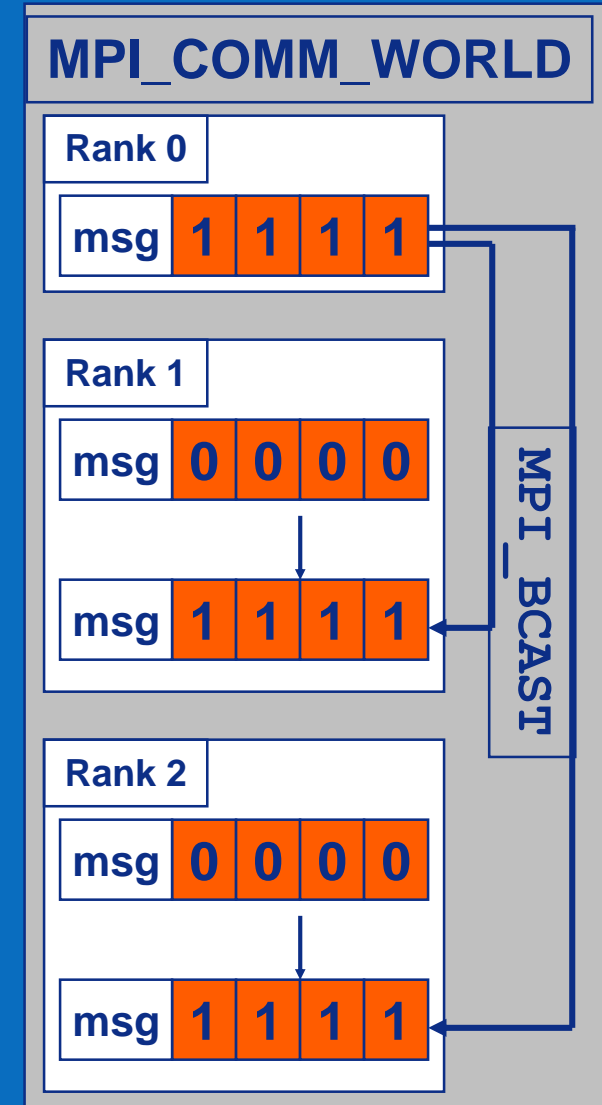
call MPI_INIT (ierr)
call MPI_COMM_SIZE (MPI_COMM_WORLD, nprocs, ierr)
call MPI_COMM_RANK (MPI_COMM_WORLD, myrank, ierr)

if (myrank == 0) then
    msg = 1
else
    msg = 0
endif

call MPI_BCAST &
    (msg, 4, MPI_INTEGER, 0, MPI_COMM_WORLD, ierr)

call MPI_FINALIZE (ierr)
end program BroadcastData

```



Data Reduction

Fortran:

```
MPI_REDUCE (sendbuf, recvbuf, count,  
            datatype, operation, root,  
            comm, ierr)  
  
<type> sendbuf(*), recvbuf(*)  
integer :: count, datatype, operation,  
            root, comm, ierr
```

C:

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

`MPI_Reduce` performs the specified reduction operation on the specified data from all processes in the communicator

MPI Reduction Operations

Operation	Function
MPI_SUM	Summation
MPI_PROD	Product
MPI_MIN	Minimum value
MPI_MINLOC	Minimum value and location
MPI_MAX	Maximum value
MPI_MAXLOC	Maximum value and location
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
User-defined	It is possible to define new reduction operations

MPI_REDUCE Example

```

program ReduceData
include 'mpif.h'

integer :: msg, sum

call MPI_INIT (ierr)
call MPI_COMM_SIZE (MPI_COMM_WORLD, nprocs, ierr)
call MPI_COMM_RANK (MPI_COMM_WORLD, myrank, ierr)

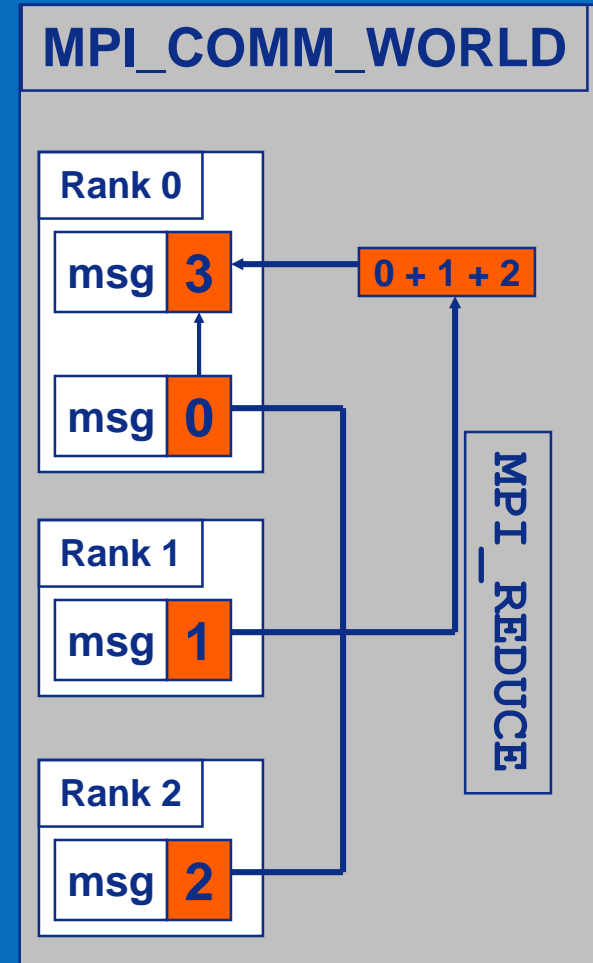
sum = 0
msg = myrank

call MPI_REDUCE
    (msg, sum, 1, MPI_INTEGER, MPI_SUM, 0,
     MPI_COMM_WORLD, ierr)

call MPI_FINALIZE (ierr)

end program ReduceData

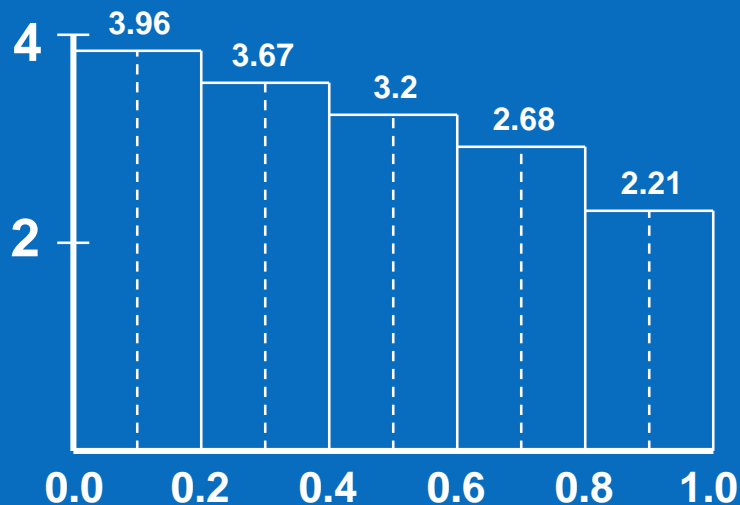
```



Activity 6: Pi Revisited

Simplify your parallel Pi integration using collective communication routines

$$f(x) = 4 / (1 + x^2)$$



$$\frac{3.96 + 3.67 + 3.2 + 2.68 + 2.21}{5 \text{ intervals}} = 3.144$$

```
#include <stdio.h>
#define INTERVALS 100000

int main (int argc, char* argv[])
{
    int i;
    double h, x, pi = 0.0;

    h = 1.0 / (double)INTERVALS

    for (i = 0; i < INTERVALS; i++)
    {
        x = h * ((double)i - 0.5);
        pi += 4.0 / (1.0 + x * x);
    }
    pi *= h;
    printf ("Pi = %f\n", pi);
}
```

Activity 7: Master-worker

Review

MPI is best for distributed-memory parallel systems

There are several MPI implementations

The MPI library is large but the core is small

MPI communication:

- Point-to-point, blocking
- Point-to-point, non-blocking
- Collective

Core MPI

11 Out of 125 Total Functions

Program startup and shutdown

- `MPI_Init`, `MPI_Finalize`
- `MPI_Comm_size`, `MPI_Comm_rank`

Point-to-point communication

- `MPI_Send`, `MPI_Recv`

Non-blocking communication

- `MPI_Isend`, `MPI_Irecv`, `MPI_Wait`

Collective communication

- `MPI_Bcast`, `MPI_Reduce`

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

Collective Communications

Advanced MPI

- **Useful Features**
- Factors Affecting MPI Performance

Useful Features

Communicators

Timing routines

User-defined data types

Inter-communicators

Topologies

Persistent communication

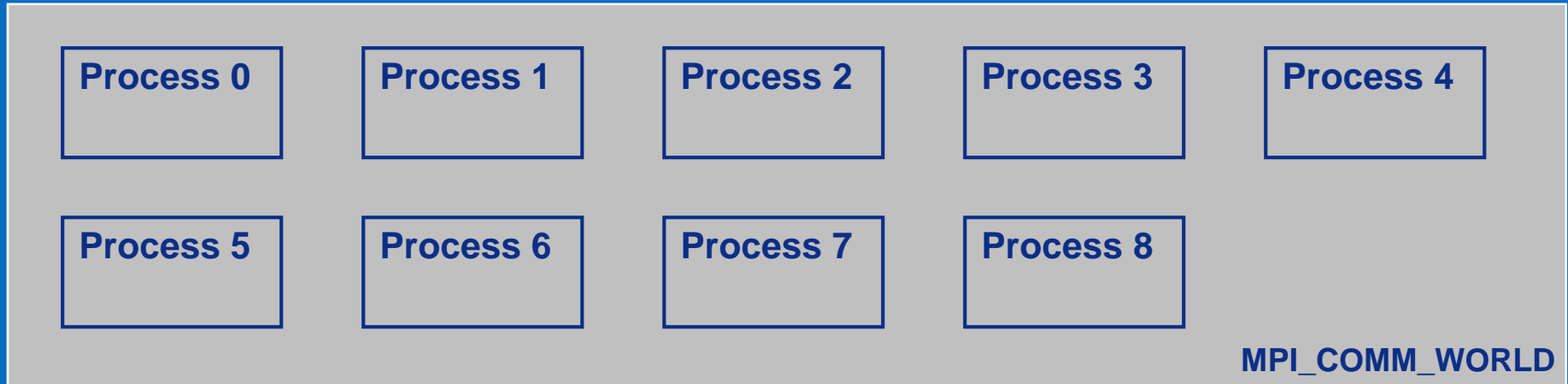
New features in MPI 2.0

Communicators

MPI uses communicators to direct message traffic.

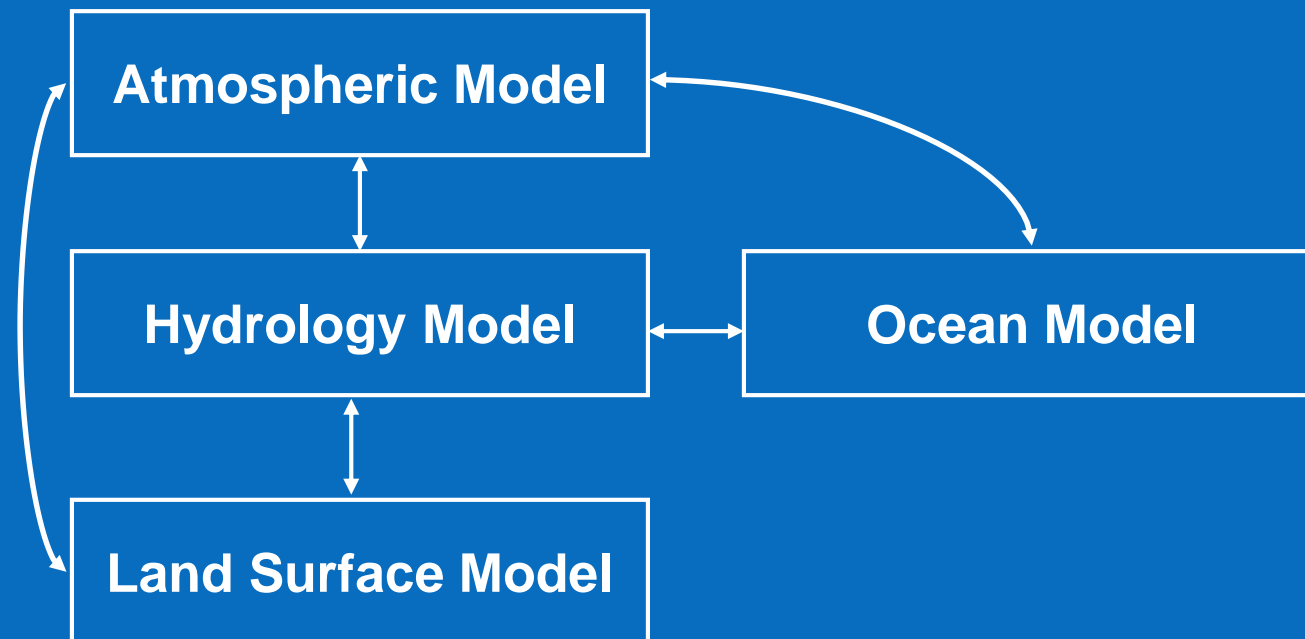
MPI_COMM_WORLD is default communicator of all processes.

Program may define other communicators.



Communicators: Example

Each model conveniently maps to its own communicator.



Timing Routines

Performance is a major design goal of MPI

Therefore, MPI provides functions to measure runtime:

- `MPI_Wtime` – returns the number of seconds from some arbitrary starting point
- `MPI_Wtick` – used to determine timer resolution

```
int start = MPI_Wtime ();  
  
// Do some work  
  
int end = MPI_Wtime ();  
int runtime = end - start;
```

Handling Non-Contiguous Data

How to send only the red elements of V?



MPI_Send

This method requires inefficient copy of non-contiguous data

One possibility, copy these elements to a temporary array before sending.

Derived Data Types

How to send only the red elements of V, while avoiding copying non-contiguous data to a temporary array?



Define a new data type, in this case a vector with stride of two from original.



```
MPI_Type_vector (3, 1, 2, MPI_REAL, &vType);
MPI_Send (V[2], 1, vType, dest, tag, MPI_COMM_WORLD);
```

MPI provides several functions to define new data types. MPICH contains optimizations for user-defined type.

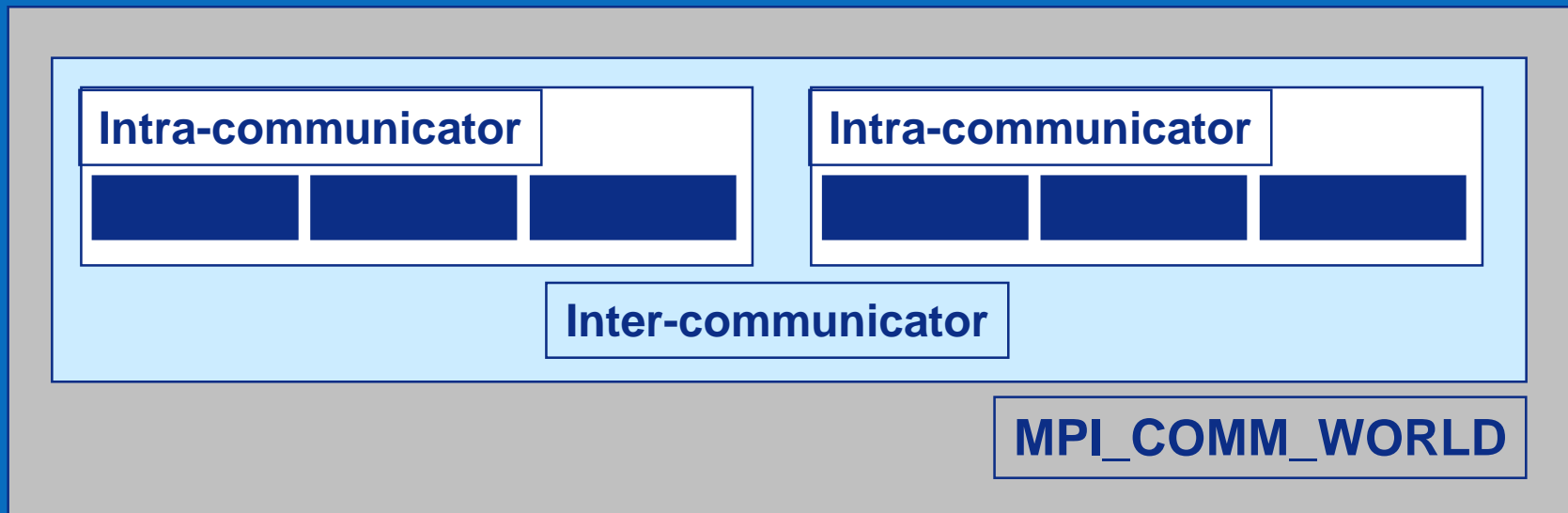
Inter-communicators

Useful for library development to hide library messages from calling program

Can be used to divide the MPI processes among different parallel computations (for example, multi-block CFD)

Functions:

- `MPI_Comm_create`, `MPI_Comm_dup`, `MPI_Comm_split`



Topologies

MPI provides functions to map application topology to processor topology



For example, consider an application that maps naturally onto a 4 x 3 Cartesian grid. Nearest neighbors are not obvious from process rank

Use `MPI_Cart_create` function to create a new communicator with the specified topology

Use `MPI_Cart_shift` to identify nearest neighbors, even for periodic grids

MPI topologies can be user-defined graphs

MPI 2.0 – More Features

One-sided communication

- MPI_Put, MPI_Get
- Uses Remote Memory Access (RMA)
- Separates communication from synchronization

Dynamic process spawning

- MPI_Spawn function
- Creates a new inter-communicator

Collective communication can be used across disjoint intra-communicators

Parallel I/O

MPI 2.0 – Better Language Support

Fortran 90/95

- Array sections
- Modules

C++

- MPI class library
- MPI functions are methods within MPI classes

Profiling MPI Communications

MPICH* contains profiling tools that are not installed by default

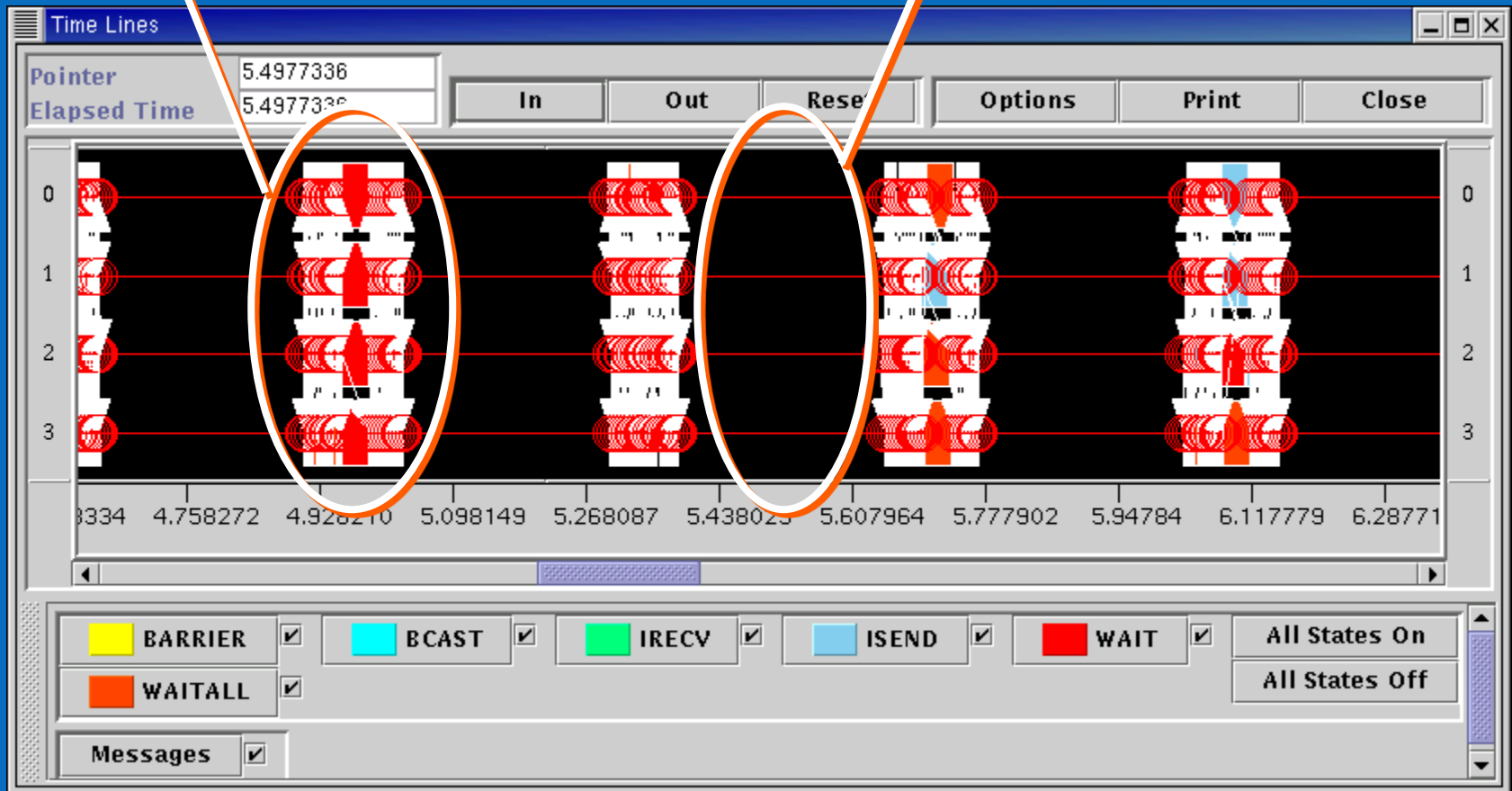
Jumpshot-2

- Installation
 - Install Java Development Kit
 - Additional MPICH configuration option:
 - `-mpe-opts=--with-java=/usr/local/jdk`
- Usage
 - Add `-mpilog` option to MPICH compilation, e.g.:
 - `mpicc -mpilog -o prog.exe prog.c`

Communications Timeline

Communication

Computation



MPI Programming

Agenda

What is MPI?

Core MPI Library

Basic Communication

Non-Blocking Communications

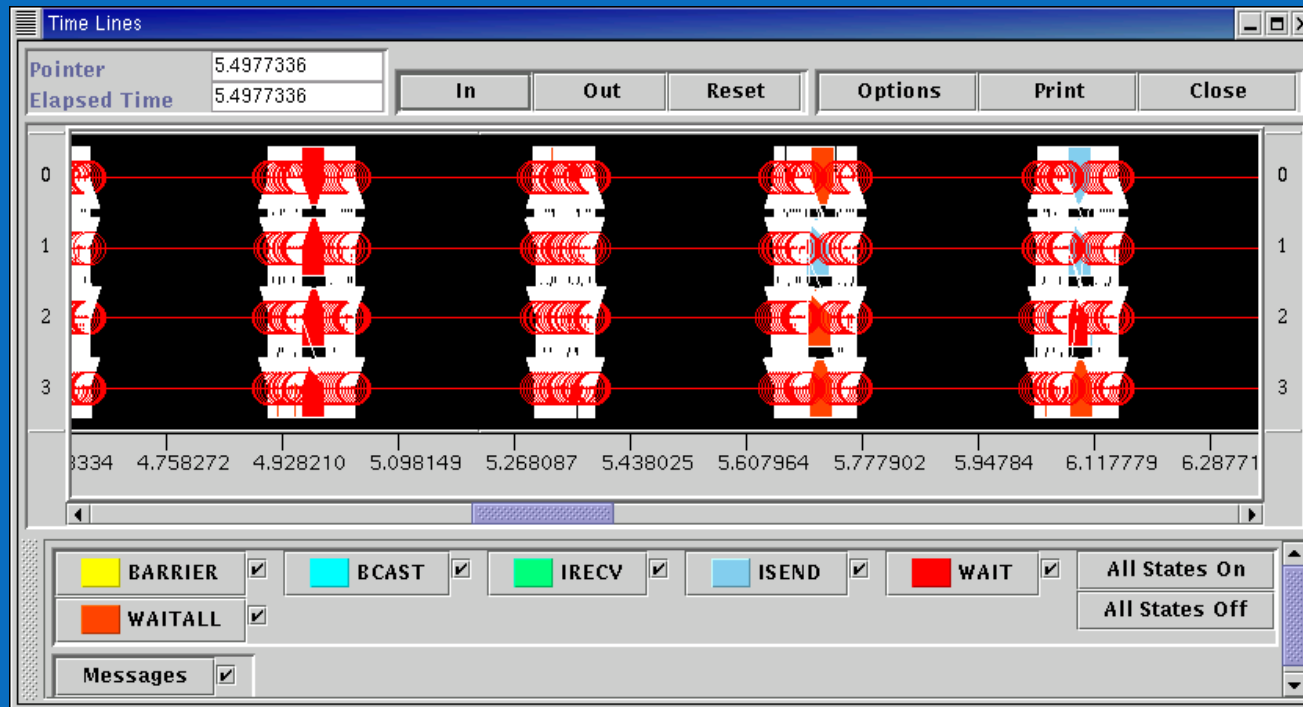
Collective Communications

Advanced MPI

- Useful Features
- **Factors Affecting MPI Performance**

Network Interconnect

- Timeline shows a repeating pattern of computation and communication but does not indicate a communications bottleneck



- The only way to improve the performance of this application is to use a faster network interconnect

Network Interconnect

Fast Ethernet

Gigabit Ethernet

Myrinet*

InfiniBand*

Quadrics*

Message Size

In general, given the same amount of data, it is better to send fewer large messages than many small messages

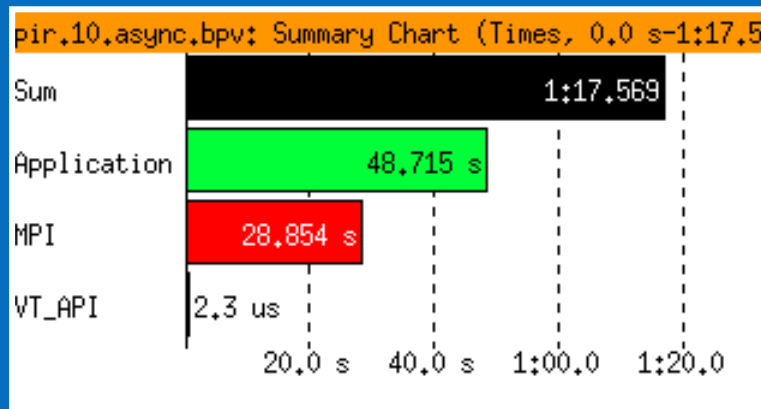
- Fewer MPI function calls, less system overhead
- Fewer messages, less TCP/IP overhead

User-defined data types in MPI can help pack several messages into one

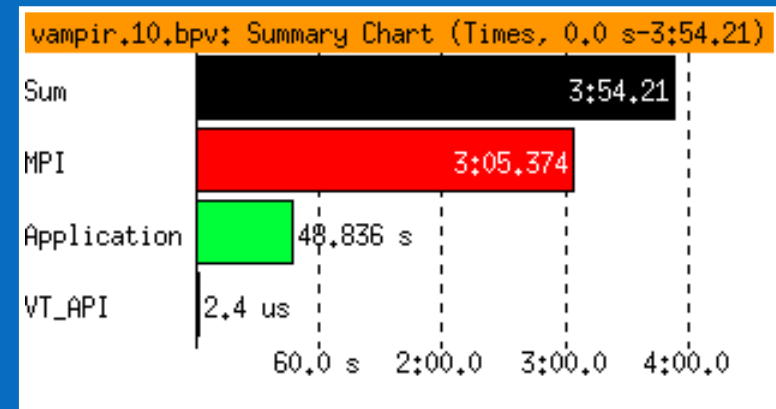
Computation: Communication

Which of these applications is likely to show better parallel speedup and scalability?

Application 1



Application 2



Obviously, Application 1 has a better ratio of computation to communication!

Overlapping Communication & Computation

```
do i = 1, neighbors
  call MPI_Irecv (edge(1,i), len, MPI_REAL, nbr(i),
                 tag, comm, requests(i), ierr)
enddo

do i = 1, neighbors
  call MPI_Isend (edge(1,i), len, MPI_REAL, nbr(i),
                 tag, comm, requests(i), ierr)
enddo

! Perform computation not requiring message data

call MPI_Waitall ((2 * neighbors), requests, status, ierr)
```

MPICH Communication Devices

ch_p4

- General communication device
- Support for SMP nodes
- Supports heterogeneous nodes
- MPICH default communication device

ch_p4mpd

- Only supports homogeneous, uni-processor clusters
- Provides faster, more scalable job startup

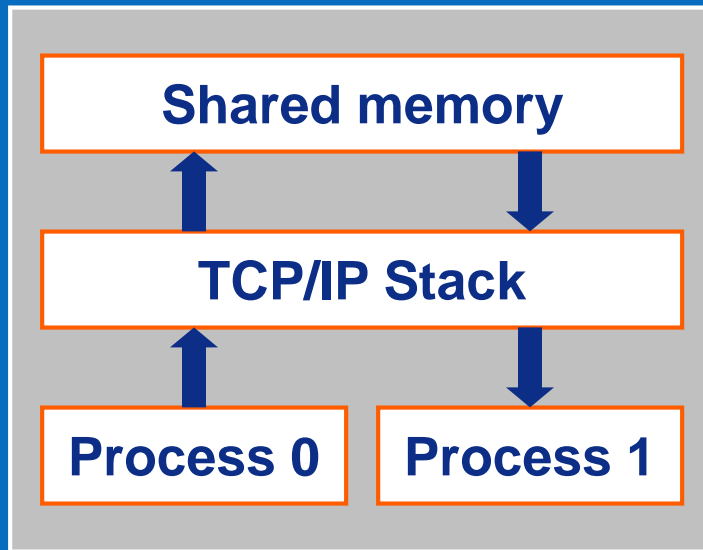
ch_shmem – Best for machines supporting shmem

ch_globus2 – Grid-enabled MPICH for Globus

ch_nt – MPICH* for Microsoft Windows*

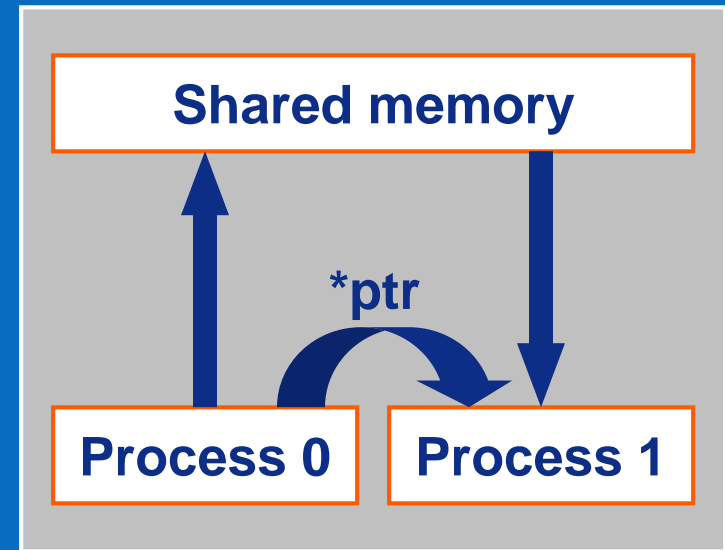
MPICH SMP Support

**Intra-node
communication through
TCP/IP**



Dual-CPU node

**Intra-node
communication through
memcpy**



Dual-CPU node

**Configure MPICH with the `-comm=shared` option
Set the `MPI_MAX_CLUSTER_SIZE` environment variable**

Backup



Getting Started with MPICH

Step 1: Download most recent package

- <ftp://ftp.mcs.anl.gov/pub/mpi/mpich.tar.gz>

Step 2: Unpack the archive in /usr/local

Step 3: Patch MPICH

- Download the latest cumulative patch
 - <ftp://ftp.mcs.anl.gov/pub/mpi/patch/1.2.5/patch.all>
- Copy patch file(s) to /usr/local/mpich-1.2.5

- Execute:

```
patch -p0 < patch.all
```

Configuring MPICH

Step 4: Configure

- Important configuration options
 - Communication device
 - Compilers
 - SMP support

- Execute:

```
configure -with-device=ch_p4 \  
-cc=icc -fc=ifc -f90=ifc -c++=icc \  
-prefix=/usr/local/mpich-1.2.5
```

- Check the configuration log for errors

Installing MPICH

Step 5: Build

- Execute: `make >& make.log`
- Check the make log for errors

Step 6: Install

- Execute: `make install`

At this point, anyone logging into the system should have access to the MPI compiler scripts, libraries, mpirun, and machine files.

MPI_GET_COUNT and Wild Card

Fortran:

```
MPI_GET_COUNT (status, datatype, count, ierr)
integer :: status(MPI_STATUS_SIZE), count,
           datatype, ierr
```

C:

```
int MPI_Get_count (MPI_Status* status,
                  MPI_Datatype datatype, int* count)
```

MPI_Get_count returns number of elements received.

MPI_Recv may use “wild card” values:

MPI_ANY_TAG

MPI_ANY_SOURCE

