## CS 7172 Parallel and Distributed Computation

#### Message Passing Interface (MPI)

#### **Kun Suo**

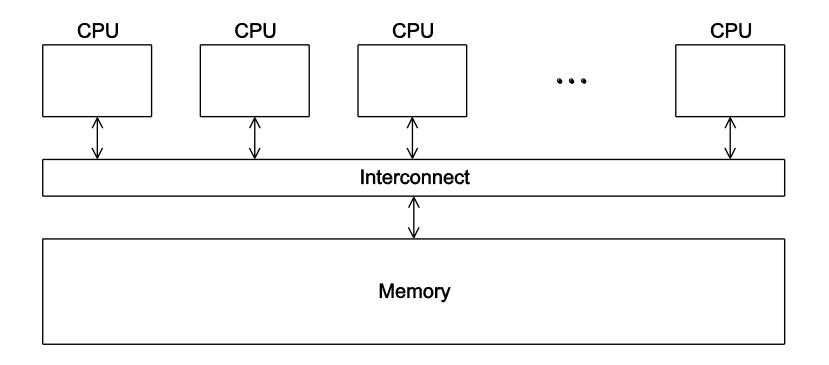
Computer Science, Kennesaw State University

https://kevinsuo.github.io/

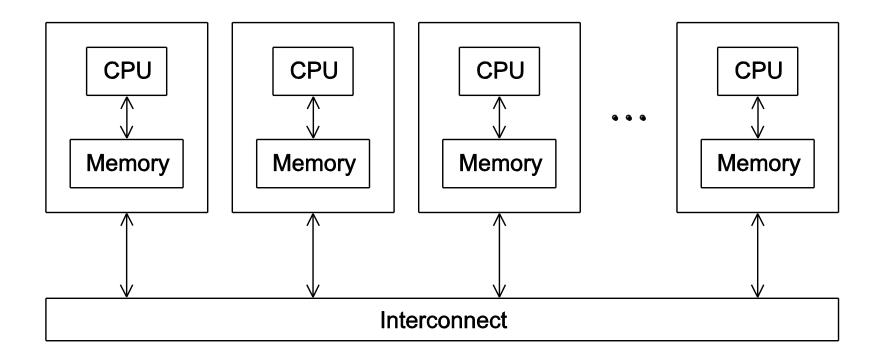
#### **Outline**

- Computer networks, primarily from an application perspective
- Protocol layering
- Client-server architecture
- End-to-end principle
- TCP
- Socket programming

## A shared memory system



## A distributed memory system



#### **Hello World!**

```
#include <stdio.h>
int main(void) {
   printf("hello, world\n");
   return 0;
}
```

## helloworld-mpi.c

```
#include <stdio.h>
int main(int argc, char** argv) {
   // Initialize the MPI environment
   MPI_Init(NULL, NULL);
   // Get the number of processes
   int world_size;
   MPI_Comm_size(MPI_COMM_WORLD, &world_size);
   // Get the rank of the process
   int world_rank;
   MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
   // Get the name of the processor
   char processor_name[MPI_MAX_PROCESSOR_NAME];
   int name_len;
   MPI_Get_processor_name(processor_name, &name_len);
   // Print off a hello world message
   printf("Hello world from processor %s, rank %d out of %d processors\n",
           processor_name, world_rank, world_size);
   // Finalize the MPI environment.
   MPI_Finalize();
```

#### Compile

wrapper script to compile source file mpicc -g -Wall -o mpi\_hello mpi\_hello.c produce create this executable file name debugging (as opposed to default a.out) information turns on all warnings

#### **Execution**

mpiexec -n <number of processes> <executable>

mpiexec -n 1 ./mpi\_hello

run with 1 process

mpiexec -n 4 ./mpi\_hello

run with 4 processes

#### **Execution**

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

Greetings from process 0 of 1!

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

```
Greetings from process 0 of 4!
```

Greetings from process 1 of 4!

Greetings from process 2 of 4!

Greetings from process 3 of 4!

#### **Execution**

```
ksuo@ksuo-VirtualBox ~/cs7172> mpiexec -n 4 ./helloworld-mpi.o

Hello world from processor ksuo-VirtualBox, rank 0 out of 4 processors

Hello world from processor ksuo-VirtualBox, rank 1 out of 4 processors

Hello world from processor ksuo-VirtualBox, rank 2 out of 4 processors

Hello world from processor ksuo-VirtualBox, rank 3 out of 4 processors

ksuo@ksuo-VirtualBox ~/cs7172> mpiexec -n 2 ./helloworld-mpi.o

Hello world from processor ksuo-VirtualBox, rank 1 out of 2 processors

Hello world from processor ksuo-VirtualBox, rank 0 out of 2 processors

ksuo@ksuo-VirtualBox ~/cs7172> mpiexec -n 1 ./helloworld-mpi.o

Hello world from processor ksuo-VirtualBox, rank 0 out of 1 processors
```

#### **MPI Programs**

- Written in C.
  - Has main.
  - Uses stdio.h, string.h, etc.
- Need to add mpi.h header file.
- Identifiers defined by MPI start with "MPI\_".
- First letter following underscore is uppercase.
  - For function names and MPI-defined types.
  - Helps to avoid confusion.

#### **MPI Components**

- MPI\_Init
  - Tells MPI to do all the necessary setup.

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

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

#### **Basic Outline**

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

#### **Communicators**

 A collection of processes that can send messages to each other.

 MPI\_Init defines a communicator that consists of all the processes created when the program is started.

Called MPI COMM WORLD.

#### **Communicators**

```
int MPI_Comm_size(
     MPI_Comm comm /* in */,
     int* comm_sz_p /* out */);
```

number of processes in the communicator

### Example 2

```
const int MAX_STRING = 100;
int main(int argc, char** argv) {
   char greeting[MAX_STRING];
   int comm_sz; //number of processes
   int my_rank; //my process rank
   // Initialize the MPI environment
   MPI_Init(NULL, NULL);
   MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);
   MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
   char processor_name[MPI_MAX_PROCESSOR_NAME];
   int name_len;
   MPI_Get_processor_name(processor_name, &name_len);
   if (mv_rank != 0) {
           sprintf(greeting, "Greetings from processor %s, rank %d out of %d processors\n",
                  processor_name, my_rank, comm_sz);
           MPI_Send(greeting, strlen(greeting)+1, MPI_CHAR, 0, 0, MPI_COMM_WORLD);
   } else {
           printf("This is process: %d, Greetings from processor %s\n",
                  my_rank, processor_name);
           for (int q = 1; q < comm_sz; q++) {
               MPI_Recv(greeting, MAX_STRING, MPI_CHAR, q, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
               printf("%s", greeting);
   // Finalize the MPI environment.
   MPI_Finalize();
```

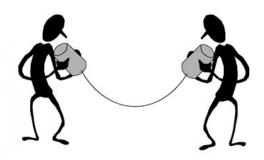
ksuo@ksuo-VirtualBox ~/cs7172> mpiexec -n 4 ./helloworld-mpi.o This is process: 0, Greetings from processor ksuo-VirtualBox Greetings from processor ksuo-VirtualBox, rank 1 out of 4 processors Greetings from processor ksuo-VirtualBox, rank 2 out of 4 processors Greetings from processor ksuo-VirtualBox, rank 3 out of 4 processors

#### **SPMD**

- Single-Program Multiple-Data
- We compile <u>one</u> program.
- Process 0 does something different.
  - Receives messages and prints them while the other processes do the work.

 The if-else construct makes our program SPMD.

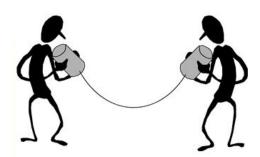
#### **Communication**



## **Data types**

MPI datatype	C datatype
MPI_CHAR	signed char
MPI_SHORT	signed short int
MPI_INT	signed int
MPI_LONG	signed long int
MPI_LONG_LONG	signed long 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	

#### **Communication**



## Message matching

```
MPI_Send(send_buf_p, send_buf_sz, send_type, dest, send_tag,
        send_comm);
                MPI_Send
                src = q
                                       MPI_Recv
MPI_Recv(recv_buf_p, recv_buf_sz, recv_type, src, recv_tag
         recv_comm, &status);
```

#### Receiving messages

- A receiver can get a message without knowing:
  - the amount of data in the message,
  - the sender of the message,
  - or the tag of the message.





#### status\_p argument

**MPI\_Status\*** 



MPI\_Status\* status;

status.MPI\_SOURCE status.MPI TAG

MPI\_SOURCE
MPI\_TAG
MPI\_ERROR

## How much data am I receiving?



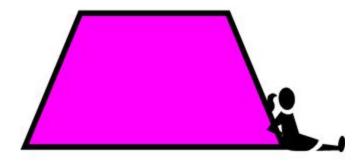
#### Issues with send and receive

- Exact behavior is determined by the MPI implementation.
- MPI\_Send may behave differently with regard to buffer size, cutoffs and blocking.
- MPI\_Recv always blocks until a matching message is received.

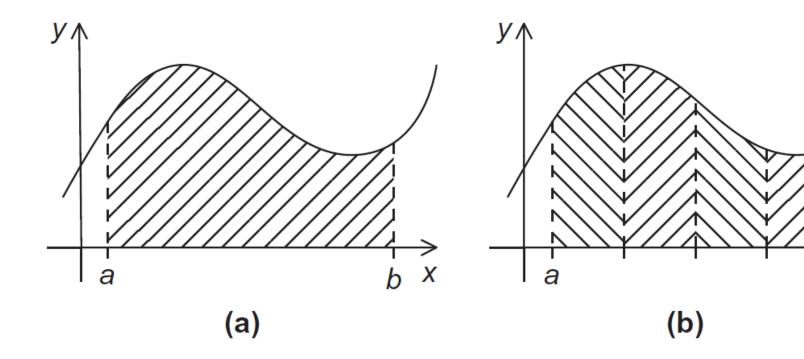
 Know your implementation; don't make assumptions!



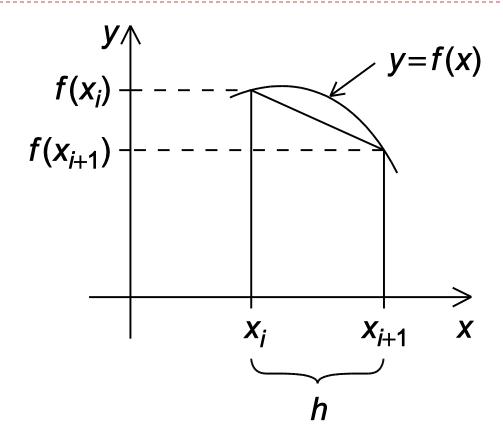
## Trapezoidal rule in mpi



## The Trapezoidal Rule



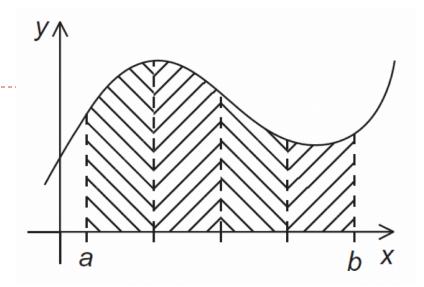
## One trapezoid



Area of one trapezoid 
$$=\frac{h}{2}[f(x_i) + f(x_{i+1})]$$

## The Trapezoidal Rule

Area of one trapezoid 
$$=\frac{h}{2}[f(x_i) + f(x_{i+1})]$$



$$h = \frac{b-a}{n}$$

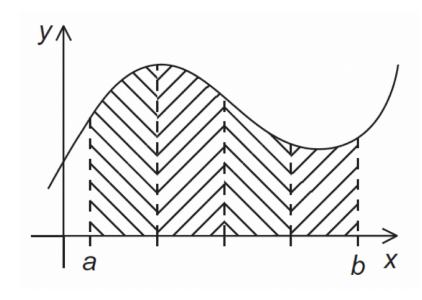
$$x_0 = a$$
,  $x_1 = a + h$ ,  $x_2 = a + 2h$ , ...,  $x_{n-1} = a + (n-1)h$ ,  $x_n = b$ 

Sum of trapezoid areas =  $h[f(x_0)/2 + f(x_1) + f(x_2) + \cdots + f(x_{n-1}) + f(x_n)/2]$ 

### Pseudo-code for a serial program

Sum of trapezoid areas =  $h[f(x_0)/2 + f(x_1) + f(x_2) + \cdots + f(x_{n-1}) + f(x_n)/2]$ 

```
/* Input: a, b, n */
h = (b-a)/n;
approx = (f(a) + f(b))/2.0;
for (i = 1; i <= n-1; i++) {
    x_i = a + i*h;
    approx += f(x_i);
}
approx = h*approx;</pre>
```



## Parallelizing the Trapezoidal Rule

1. Partition problem solution into tasks.

Identify communication channels between tasks.

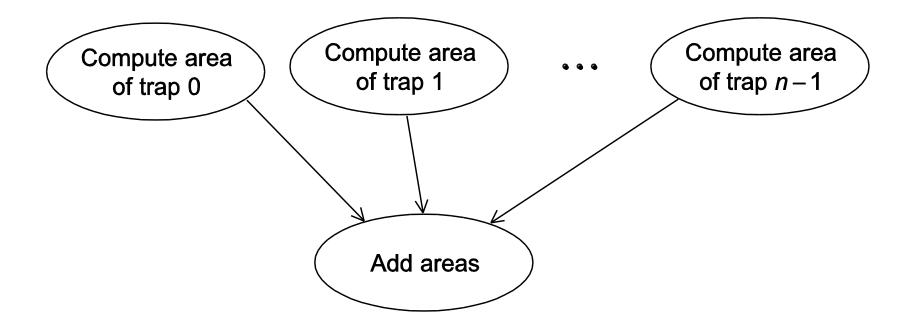
3. Aggregate tasks into composite tasks.

4. Map composite tasks to cores.

### Parallel pseudo-code

```
Get a, b, n;
      h = (b-a)/n;
      local n = n/comm sz;
      local_a = a + my_rank*local_n*h;
4
      local_b = local_a + local_n*h;
      local_integral = Trap(local_a, local_b, local_n, h);
      if (mv rank != 0)
         Send local integral to process 0;
9
      else /* my\_rank == 0 */
10
         total integral = local integral;
11
         for (proc = 1; proc < comm_sz; proc++) {</pre>
12
             Receive local_integral from proc;
13
             total integral += local integral;
14
15
16
      if (my_rank == 0)
17
         print result;
```

# Tasks and communications for Trapezoidal Rule



### First version (1)

```
int main(void) {
      int my rank, comm sz, n = 1024, local n;
      double a = 0.0, b = 3.0, h, local_a, local_b;
      double local int, total int;
      int source:
      MPI Init(NULL, NULL);
      MPI Comm rank(MPI_COMM_WORLD, &my_rank);
9
      MPI Comm size (MPI COMM WORLD, &comm sz);
10
      h = (b-a)/n; /* h is the same for all processes */
11
      local n = n/comm sz; /* So is the number of trapezoids */
12
13
14
      local_a = a + my rank*local n*h;
15
      local_b = local_a + local n*h;
16
      local int = Trap(local a, local b, local n, h);
17
18
      if (my rank != 0) {
19
         MPI Send(&local int, 1, MPI DOUBLE, 0, 0,
20
               MPI COMM WORLD);
```

## First version (2)

```
21
      } else {
22
         total int = local int;
23
         for (source = 1; source < comm_sz; source++) {</pre>
24
             MPI Recv(&local int, 1, MPI DOUBLE, source, 0,
25
                   MPI_COMM_WORLD , MPI_STATUS_IGNORE );
26
             total int += local int;
27
28
29
30
      if (my rank == 0) {
31
         printf("With n = %d trapezoids, our estimate \n", n);
32
         printf("of the integral from f to f = .15e\n",
33
              a, b, total int);
34
35
      MPI Finalize();
      return 0:
36
37
     /* main */
```

## First version (3)

```
double Trap(
         double left endpt /* in */.
         double right_endpt /* in */,
         int trap_count /* in */,
5
         double base_len /* in */) {
      double estimate, x;
6
      int i;
9
      estimate = (f(left endpt) + f(right endpt))/2.0;
10
      for (i = 1; i \le trap_count - 1; i++)
         x = left_endpt + i*base len;
11
12
         estimate += f(x);
13
14
      estimate = estimate * base len;
15
      return estimate;
16
17
     /* Trap */
```

## Dealing with I/O

```
#include < stdio.h>
#include <mpi.h>
                                   Each process just
                                   prints a message.
int main(void) {
   int my_rank, comm_sz;
   MPI Init(NULL, NULL);
   MPI Comm size (MPI COMM WORLD, &comm sz);
   MPI Comm rank (MPI COMM WORLD, &my rank);
   printf("Proc %d of %d > Does anyone have a toothpick?\n",
         my rank, comm sz);
   MPI Finalize();
   return 0:
  /* main */
```

### Running with 6 processes

```
Proc 0 of 6 > Does anyone have a toothpick?

Proc 1 of 6 > Does anyone have a toothpick?

Proc 2 of 6 > Does anyone have a toothpick?

Proc 4 of 6 > Does anyone have a toothpick?

Proc 3 of 6 > Does anyone have a toothpick?

Proc 5 of 6 > Does anyone have a toothpick?
```

unpredictable output



### Input

- Most MPI implementations only allow process
   0 in MPI\_COMM\_WORLD access to stdin.
- Process 0 must read the data (scanf) and send to the other processes.

```
MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);
Get_data(my_rank, comm_sz, &a, &b, &n);
h = (b-a)/n;
. . .
```

## Function for reading user input

```
void Get input(
         int my rank /*in */,
         int comm_sz /* in */,
         double* a_p /* out */,
         double* b_p /* out */,
         int * n p /* out */) {
       int dest;
       if (mv rank == 0) {
         printf("Enter a, b, and n\n");
         scanf("%lf %lf %d", a_p, b_p, n_p);
         for (dest = 1; dest < comm_sz; dest++) {
            MPI Send(a p, 1, MPI DOUBLE, dest, 0, MPI COMM WORLD);
            MPI Send(b p, 1, MPI DOUBLE, dest, 0, MPI COMM WORLD);
            MPI Send(n_p, 1, MPI_INT, dest, 0, MPI_COMM_WORLD);
      else \{ /* my\_rank != 0 */
         MPI Recv(a p, 1, MPI DOUBLE, 0, 0, MPI COMM WORLD,
               MPI STATUS IGNORE);
         MPI_Recv(b_p, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD,
               MPI STATUS IGNORE);
         MPI_Recv(n_p, 1, MPI_INT, 0, 0, MPI_COMM_WORLD,
               MPI STATUS IGNORE);
|} /* Get_input */
```

outation

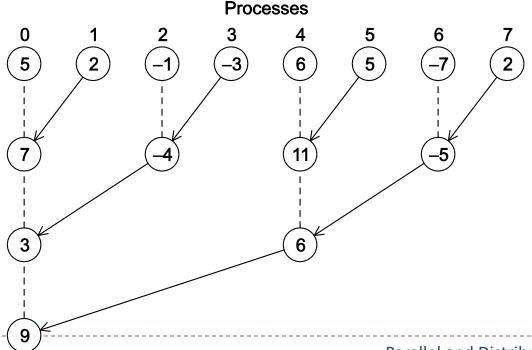
### **Collective communication**



#### Tree-structured communication

#### In the first phase:

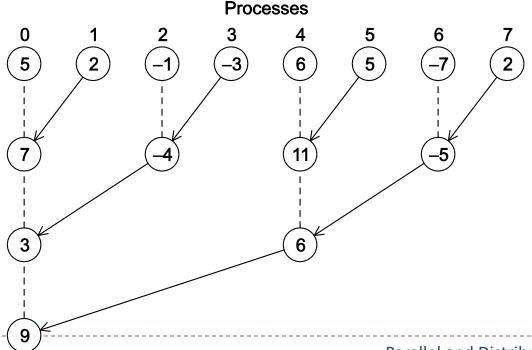
- o (a) Process 1 sends to 0; 3 sends to 2; 5 sends to 4; and 7 sends to 6.
- (b) Processes 0, 2, 4, and 6 add in the received values.
- (c) Processes 2 and 6 send their new values to processes 0 and 4, respectively.
- (d) Processes 0 and 4 add the received values into their new values.



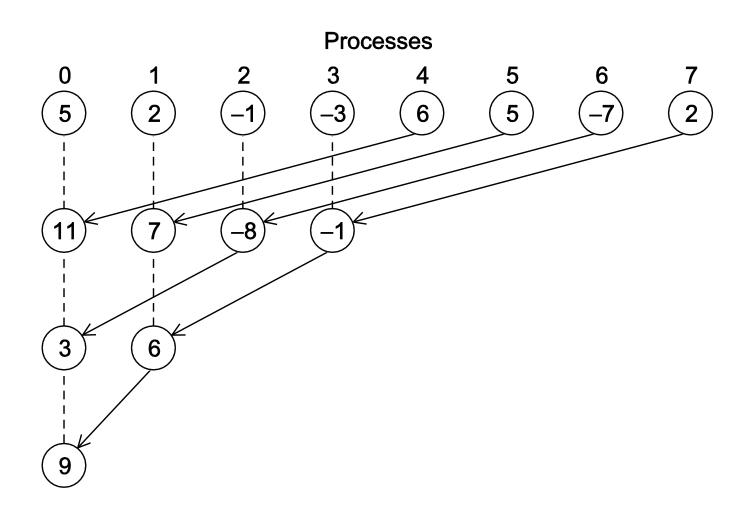
#### Tree-structured communication

#### In the second phase:

- (a) Process 4 sends its newest value to process 0.
- (b) Process 0 adds the received value to its newest value.



### An alternative tree-structured global sum



### MPI\_Reduce

```
\label{eq:mpi_reduce} \begin{split} \texttt{MPI\_Reduce}(\&\texttt{local\_int}\,,\;\&\texttt{total\_int}\,,\;\;1\,,\;\;\texttt{MPI\_DOUBLE}\,,\;\;\texttt{MPI\_SUM}\,,\;\;0\,,\\ \texttt{MPI\_COMM\_WORLD}\,); \end{split}
```

### Predefined reduction operators in MPI

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

• <u>All</u> the processes in the communicator must call the same collective function.

 For example, a program that attempts to match a call to MPI\_Reduce on one process with a call to MPI\_Recv on another process is erroneous, and, in all likelihood, the program will hang or crash.

 The arguments passed by each process to an MPI collective communication must be "compatible."

 For example, if one process passes in 0 as the dest\_process and another passes in 1, then the outcome of a call to MPI\_Reduce is erroneous, and, once again, the program is likely to hang or crash.

 The output\_data\_p argument is only used on dest\_process.

 However, all of the processes still need to pass in an actual argument corresponding to output\_data\_p, even if it's just NULL.

 Point-to-point communications are matched on the basis of tags and communicators.

Collective communications don't use tags.

 They're matched solely on the basis of the communicator and the order in which they're called.

## Example (1)

Tim	e Process 0	Process 1	Process 2
0	a = 1; c = 2	a = 1; c = 2	a = 1; c = 2
1	MPI_Reduce(&a, &b,)	MPI_Reduce(&c, &d,)	MPI_Reduce(&a, &b,)
2	MPI_Reduce(&c, &d,)	MPI_Reduce(&a, &b,)	MPI_Reduce(&c, &d,)

### Multiple calls to MPI\_Reduce

## Example (2)

 Suppose that each process calls MPI\_Reduce with operator MPI\_SUM, and destination process 0.

 At first glance, it might seem that after the two calls to MPI\_Reduce, the value of b will be 3, and the value of d will be 6.

## Example (3)

 However, the names of the memory locations are irrelevant to the matching of the calls to MPI\_Reduce.

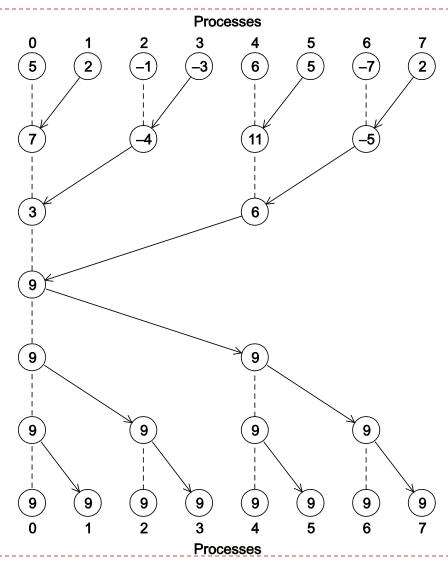
• The order of the calls will determine the matching so the value stored in b will be 1+2+1 = 4, and the value stored in d will be 2+1+2 = 5.

## MPI\_Allreduce

 Useful in a situation in which all of the processes need the result of a global sum in order to complete some larger computation.

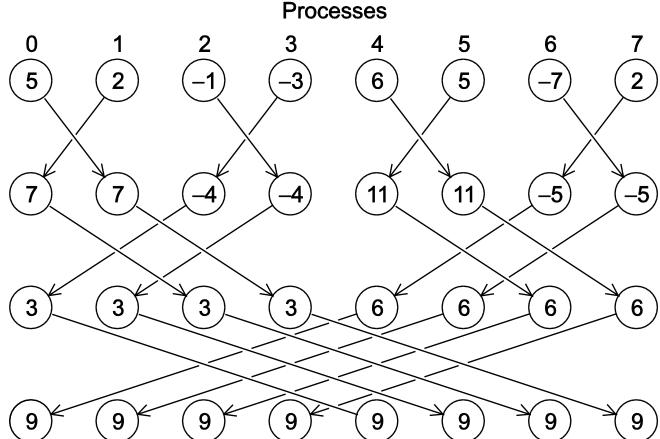
```
int MPI_Allreduce(
         void*
                      input_data_p /* in
         void*
                      output_data_p /* out */,
         int
                                   /* in */,
                     count
                                    /* in */,
        MPI_Datatype datatype
        MPI_Op
                                    /* in */.
                     operator
                                     /* in
                                           */);
        MPI Comm
                     comm
```

# A global sum followed by distribution of the result



## A butterfly-structured global sum

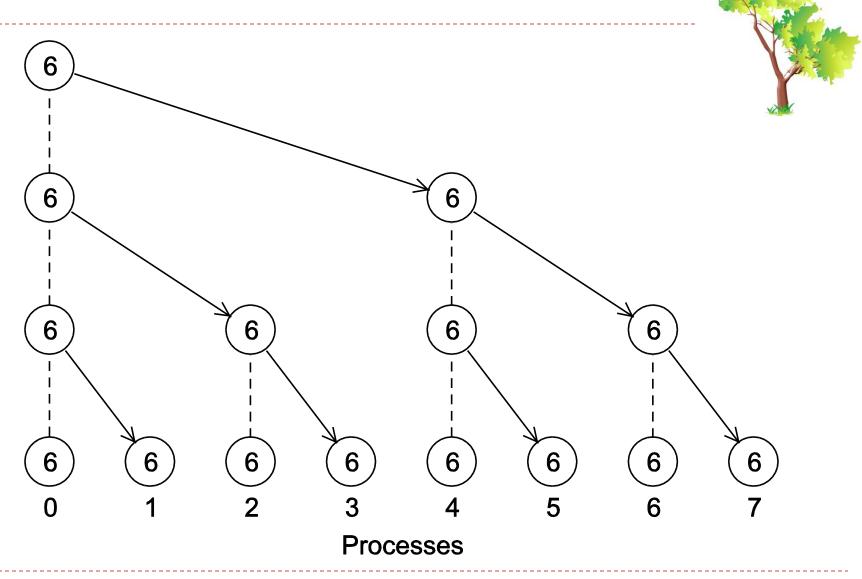




#### **Broadcast**

 Data belonging to a single process is sent to all of the processes in the communicator.

### A tree-structured broadcast



# A version of Get\_input that uses MPI Bcast

```
void Get input(
     int my_rank /* in */,
     int comm_sz /* in */,
     double * a_p /* out */,
     double* b_p /* out */,
     int * n_p /* out */) {
  if (my_rank == 0) {
     printf("Enter a, b, and n\n");
     scanf("%lf %lf %d", a_p, b_p, n_p);
  MPI_Bcast(a_p, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);
  MPI_Bcast(b_p, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);
  MPI_Bcast(n_p, 1, MPI_INT, 0, MPI_COMM_WORLD);
\} /* Get_input */
```

### **Data distributions**

$$\mathbf{x} + \mathbf{y} = (x_0, x_1, \dots, x_{n-1}) + (y_0, y_1, \dots, y_{n-1})$$

$$= (x_0 + y_0, x_1 + y_1, \dots, x_{n-1} + y_{n-1})$$

$$= (z_0, z_1, \dots, z_{n-1})$$

$$= \mathbf{z}$$

Compute a vector sum.

### Serial implementation of vector addition

```
void Vector_sum(double x[], double y[], double z[], int n) {
  int i;

for (i = 0; i < n; i++)
    z[i] = x[i] + y[i];
} /* Vector_sum */</pre>
```

# Different partitions of a 12-component vector among 3 processes

					С	omp	one	ents				
									B	Bloc	k-cyc	lic
Process		В	lock			Cy	clic	;	В	lock	size:	= 2
0	0	1	2	3	0	3	6	9	0	1	6	7
1	4	5	6	7	1	4	7	10	2	3	8	9
2	8	9	10	11	2	5	8	11	4	5	10	11

## **Partitioning options**

- Block partitioning
  - Assign blocks of consecutive components to each process.

- Cyclic partitioning
  - Assign components in a round robin fashion.

- Block-cyclic partitioning
  - Use a cyclic distribution of blocks of components.

# Parallel implementation of vector addition

```
void Parallel_vector_sum(
    double local_x[] /* in */,
    double local_y[] /* in */,
    double local_z[] /* out */,
    int local_n /* in */) {
    int local_i;

    for (local_i = 0; local_i < local_n; local_i++)
        local_z[local_i] = local_x[local_i] + local_y[local_i];
} /* Parallel_vector_sum */</pre>
```

#### Scatter

 MPI\_Scatter can be used in a function that reads in an entire vector on process 0 but only sends the needed components to each of the other processes.

```
int MPI_Scatter(
    void*
                send_buf_p /* in */,
    int
                send_count /* in */,
                send_type /* in */,
    MPI_Datatype
    void*
                recv buf p /* out */,
    int
                recv_count /*in */.
                recv_type /* in */,
    MPI_Datatype
    int
                src proc /* in */,
    MPI Comm
                comm /* in */):
```

## Reading and distributing a vector

```
void Read vector(
     double local_a[] /* out */,
     int local n /* in */,
                     /* in */.
     int n
     char vec name [] /* in */,
     int     my_rank  /* in */,
     MPI_Comm comm /* in */) {
  double * a = NULL;
  int i;
  if (my rank == 0) {
     a = malloc(n*sizeof(double));
     printf("Enter the vector %s\n", vec_name);
     for (i = 0; i < n; i++)
        scanf("%lf", &a[i]);
     MPI_Scatter(a, local_n, MPI_DOUBLE, local_a, local_n, MPI_DOUBLE,
           0. \text{comm}):
     free(a);
  } else {
     MPI Scatter(a, local n, MPI DOUBLE, local a, local n, MPI DOUBLE,
           0. \text{comm}):
  /* Read_vector */
```

utation

#### **Gather**

 Collect all of the components of the vector onto process 0, and then process 0 can process all of the components.

```
int MPI Gather(
     void*
                   send_buf_p /* in */,
                   send_count /*in */,
     int
                   send_type /*in */,
     MPI Datatype
     void*
                   recv_buf_p /* out */,
     int
                   recv_count /*in */,
                   recv_type /* in */,
     MPI Datatype
     int
                   dest_proc /* in */,
     MPI_Comm
                              /* in */);
                   comm
```

## Print a distributed vector (1)

```
void Print_vector(
     double local_b[] /* in */,
              local_n /* in */,
     int
     int
                    /* in */
     char
            title[] /* in */,
     int
          my_rank /* in */,
                       /* in */) {
     MPI_Comm comm
  double * b = NULL;
  int i;
```

## Print a distributed vector (2)

## **Allgather**

 Concatenates the contents of each process' send\_buf\_p and stores this in each process' recv\_buf\_p.

 As usual, recv\_count is the amount of data being received from each process.

## **Matrix-vector multiplication**

$$A = (a_{ij})$$
 is an  $m \times n$  matrix

 $\mathbf{x}$  is a vector with n components

y = Ax is a vector with m components

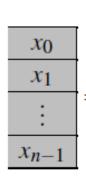
$$y_i = a_{i0}x_0 + a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{i,n-1}x_{n-1}$$

i-th component of y

Dot product of the ith row of A with x.

## Matrix-vector multiplication

<i>a</i> <sub>00</sub>	<i>a</i> <sub>01</sub>	• • •	$a_{0,n-1}$
<i>a</i> <sub>10</sub>	$a_{11}$	:	$a_{1,n-1}$
:	:		:
$a_{i0}$	$a_{i1}$	• • •	$a_{i,n-1}$
<i>a</i> <sub>i0</sub> :	<i>a</i> <sub>i1</sub> :	•••	$a_{i,n-1}$ :



уо
<i>y</i> <sub>1</sub>
<b>:</b>
$y_i = a_{i0}x_0 + a_{i1}x_1 + \cdots + a_{i,n-1}x_{n-1}$
$y_i = a_{i0}x_0 + a_{i1}x_1 + \cdots + a_{i,n-1}x_{n-1}$ :

#### Multiply a matrix by a vector

```
/* For each row of A */
for (i = 0; i < m; i++) {
    /* Form dot product of ith row with x */
    v[i] = 0.0;

for (j = 0; j < n; j++)
    v[i] += A[i][j]*x[j];
}</pre>
```

Serial pseudo-code

#### C style arrays

$$\begin{pmatrix} 0 & 1 & 2 & 3 \\ 4 & 5 & 6 & 7 \\ 8 & 9 & 10 & 11 \end{pmatrix}$$
 stored as

01234567891011

#### Serial matrix-vector multiplication

```
void Mat_vect_mult(
     double A[] /* in */,
     double x[] /* in */,
     double y[] /* out */,
     int m /*in */,
     int n /* in */) {
  int i, j;
  for (i = 0; i < m; i++) {
     y[i] = 0.0;
     for (j = 0; j < n; j++)
        y[i] += A[i*n+j]*x[j];
 /* Mat_vect_mult */
```

# An MPI matrix-vector multiplication function (1)

```
void Mat_vect_mult(
     double local_A[] /* in */,
     double local_x[] /* in */,
     double local_y[] /* out */,
     int
             local_m /* in */,
                     /* in */,
     int
             n
        local_n /*in */,
     int
     MPI_Comm comm /* in */) {
  double * x;
  int local_i, j;
  int local_ok = 1;
```

# An MPI matrix-vector multiplication function (2)

### **MPI** derived datatypes



#### **Derived datatypes**

- Used to represent any collection of data items in memory by storing both the types of the items and their relative locations in memory.
- The idea is that if a function that sends data knows this information about a collection of data items, it can collect the items from memory before they are sent.

 Similarly, a function that receives data can distribute the items into their correct destinations in memory when they're received.

#### **Derived datatypes**

 Formally, consists of a sequence of basic MPI data types together with a displacement for each of the data types.

Trapezoidal Rule example:

Variable	Address
a	24
b	40
n	48

 $\{(MPI\_DOUBLE, 0), (MPI\_DOUBLE, 16), (MPI\_INT, 24)\}$ 

#### MPI\_Type create\_struct

 Builds a derived datatype that consists of individual elements that have different basic types.

#### MPI\_Get\_address

 Returns the address of the memory location referenced by location\_p.

 The special type MPI\_Aint is an integer type that is big enough to store an address on the system.

```
int MPI_Get_address(
    void* location_p /* in */,
    MPI_Aint* address_p /* out */);
```

#### MPI\_Type\_commit

 Allows the MPI implementation to optimize its internal representation of the datatype for use in communication functions.

```
int MPI_Type_commit(MPI_Datatype* new_mpi_t_p /* in/out */);
```

#### MPI\_Type\_free

 When we're finished with our new type, this frees any additional storage used.

```
int MPI_Type_free(MPI_Datatype* old_mpi_t_p /* in/out */);
```

# Get input function with a derived datatype (1)

# Get input function with a derived datatype (2)

# Get input function with a derived datatype (3)

```
void Get_input(int my_rank, int comm_sz, double* a_p, double* b_p,
     int * n p) {
  MPI Datatype input mpi t;
  Build_mpi_type(a_p, b_p, n_p, &input_mpi_t);
   if (my rank == 0) {
     printf("Enter a, b, and n\n");
      scanf("%lf %lf %d", a_p, b_p, n_p);
  MPI_Bcast(a_p, 1, input_mpi_t, 0, MPI_COMM_WORLD);
  MPI Type free(&input mpi t);
 /* Get_input */
```

#### **Performance evaluation**



#### Elapsed parallel time

 Returns the number of seconds that have elapsed since some time in the past.

#### **Elapsed serial time**

In this case, you don't need to link in the MPI libraries.

 Returns time in microseconds elapsed from some point in the past.

```
#include "timer.h"
. . .
double now;
. . .
GET_TIME(now);
```



#### **Elapsed serial time**

```
#include "timer.h"
. . .
double start, finish;
. . .
GET_TIME(start);
/* Code to be timed */
. . .
GET_TIME(finish);
printf("Elapsed time = %e seconds\n", finish-start);
```

#### **MPI\_Barrier**

 Ensures that no process will return from calling it until every process in the communicator has started calling it.



#### **MPI\_Barrier**

```
double local_start, local_finish, local_elapsed, elapsed;
MPI_Barrier(comm);
local start = MPI Wtime();
/* Code to be timed */
local finish = MPI Wtime();
local_elapsed = local_finish - local_start;
MPI_Reduce(&local_elapsed, &elapsed, 1, MPI_DOUBLE,
  MPI MAX, 0, comm);
if (my rank == 0)
  printf("Elapsed time = %e seconds\n", elapsed);
```

### Run-times of serial and parallel matrixvector multiplication

	Order of Matrix				
comm_sz	1024	2048	4096	8192	16,384
1	4.1	16.0	64.0	270	1100
2	2.3	8.5	33.0	140	560
4	2.0	5.1	18.0	70	280
8	1.7	3.3	9.8	36	140
16	1.7	2.6	5.9	19	71

(Seconds)

#### Speedup

$$S(n, p) = \frac{T_{\text{serial}}(n)}{T_{\text{parallel}}(n, p)}$$

### **Efficiency**

$$E(n,p) = \frac{S(n,p)}{p} = \frac{T_{\text{serial}}(n)}{p \times T_{\text{parallel}}(n,p)}$$

# **Speedups of Parallel Matrix-Vector Multiplication**

	Order of Matrix				
comm_sz	1024	2048	4096	8192	16,384
1	1.0	1.0	1.0	1.0	1.0
2	1.8	1.9	1.9	1.9	2.0
4	2.1	3.1	3.6	3.9	3.9
8	2.4	4.8	6.5	7.5	7.9
16	2.4	6.2	10.8	14.2	15.5

# **Efficiencies of Parallel Matrix-Vector Multiplication**

	Order of Matrix				
comm_sz	1024	2048	4096	8192	16,384
1	1.00	1.00	1.00	1.00	1.00
2	0.89	0.94	0.97	0.96	0.98
4	0.51	0.78	0.89	0.96	0.98
8	0.30	0.61	0.82	0.94	0.98
16	0.15	0.39	0.68	0.89	0.97

### **Scalability**

 A program is scalable if the problem size can be increased at a rate so that the efficiency doesn't decrease as the number of processes increase.

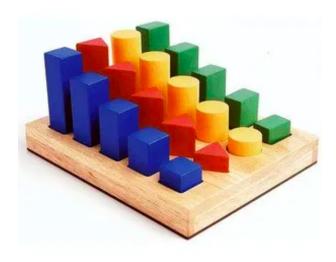


### **Scalability**

Programs that can maintain a constant
 efficiency without increasing the problem size
 are sometimes said to be strongly scalable.

 Programs that can maintain a constant efficiency if the problem size increases at the same rate as the number of processes are sometimes said to be weakly scalable.

### A parallel sorting algorithm



#### Sorting

- n keys and p = comm sz processes.
- n/p keys assigned to each process.
- No restrictions on which keys are assigned to which processes.
- When the algorithm terminates:
  - The keys assigned to each process should be sorted in (say) increasing order.
  - If 0 ≤ q < r < p, then each key assigned to process q should be less than or equal to every key assigned to process r.

#### Serial bubble sort

```
void Bubble_sort(
     int a[] /* in/out */,
     int n /* in */) {
   int list_length, i, temp;
  for (list_length = n; list_length \geq 2; list_length--)
     for (i = 0; i < list_length -1; i++)
         if (a[i] > a[i+1]) {
           temp = a[i];
           a[i] = a[i+1];
           a[i+1] = temp;
  /* Bubble_sort */
```

### **Odd-even transposition sort**

- A sequence of phases.
- Even phases, compare swaps:

$$(a[0], a[1]), (a[2], a[3]), (a[4], a[5]), \dots$$

Odd phases, compare swaps:

$$(a[1], a[2]), (a[3], a[4]), (a[5], a[6]), \dots$$

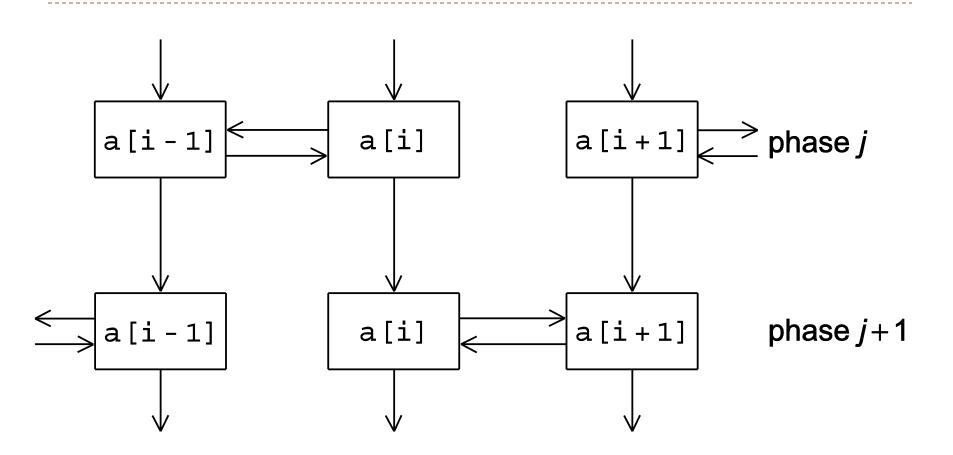
#### **Example**

```
Start: 5, 9, 4, 3
Even phase: compare-swap (5,9) and (4,3)
  getting the list 5, 9, 3, 4
Odd phase: compare-swap (9,3)
 getting the list 5, 3, 9, 4
Even phase: compare-swap (5,3) and (9,4)
  getting the list 3, 5, 4, 9
Odd phase: compare-swap (5,4)
 getting the list 3, 4, 5, 9
```

### Serial odd-even transposition sort

```
void Odd_even_sort(
      int a[] /* in/out */,
     int n /* in */) {
   int phase, i, temp;
  for (phase = 0; phase < n; phase++)
      if (phase % 2 == 0) { /* Even phase */
        for (i = 1; i < n; i += 2)
            if (a[i-1] > a[i]) {
              temp = a[i];
              a[i] = a[i-1];
              a[i-1] = temp;
     } else { /* Odd phase */
        for (i = 1; i < n-1; i += 2)
            if (a[i] > a[i+1]) {
              temp = a[i];
              a[i] = a[i+1];
              a[i+1] = temp;
  /* Odd_even_sort */
```

## Communications among tasks in odd-even sort



Tasks determining a[i] are labeled with a[i].

### Parallel odd-even transposition sort

	Process				
Time	0	1	2	3	
Start	15, 11, 9, 16	3, 14, 8, 7	4, 6, 12, 10	5, 2, 13, 1	
After Local Sort	9, 11, 15, 16	3, 7, 8, 14	4, 6, 10, 12	1, 2, 5, 13	
After Phase 0	3, 7, 8, 9	11, 14, 15, 16	1, 2, 4, 5	6, 10, 12, 13	
After Phase 1	3, 7, 8, 9	1, 2, 4, 5	11, 14, 15, 16	6, 10, 12, 13	
After Phase 2	1, 2, 3, 4	5, 7, 8, 9	6, 10, 11, 12	13, 14, 15, 16	
After Phase 3	1, 2, 3, 4	5, 6, 7, 8	9, 10, 11, 12	13, 14, 15, 16	

#### Pseudo-code

```
Sort local keys;
for (phase = 0; phase < comm_sz; phase++) {
   partner = Compute_partner(phase, my_rank);
   if (I'm not idle) {
      Send my keys to partner;
      Receive keys from partner;
      if (my_rank < partner)
         Keep smaller keys;
      else
         Keep larger keys;
```

#### Compute\_partner

```
if (phase % 2 == 0) /* Even phase */
  if (my_rank % 2 != 0) /* Odd rank */
     partner = my_rank - 1;
  else
                            /* Even rank */
     partner = my_rank + 1;
else
                       /* Odd phase */
  if (my_rank % 2 != 0) /* Odd rank */
     partner = my_rank + 1;
  else
                            /* Even rank */
     partner = my_rank - 1;
if (partner == -1 || partner == comm_sz)
  partner = MPI_PROC_NULL;
```

- The MPI standard allows MPI\_Send to behave in two different ways:
  - it can simply copy the message into an MPI managed buffer and return,
  - or it can block until the matching call to MPI\_Recv starts.

 Many implementations of MPI set a threshold at which the system switches from buffering to blocking.

 Relatively small messages will be buffered by MPI\_Send.

Larger messages, will cause it to block.

 If the MPI\_Send executed by each process blocks, no process will be able to start executing a call to MPI\_Recv, and the program will hang or deadlock.

 Each process is blocked waiting for an event that will never happen.

(see pseudo-code)

 A program that relies on MPI provided buffering is said to be unsafe.

 Such a program may run without problems for various sets of input, but it may hang or crash with other sets.

# MPI\_Ssend

- An alternative to MPI\_Send defined by the MPI standard.
- The extra "s" stands for synchronous and MPI\_Ssend is guaranteed to block until the matching receive starts.

### Restructuring communication

```
\label{eq:mpi_send} \begin{split} \text{MPI\_Send(msg, size, MPI\_INT, (my\_rank+1) \% comm\_sz, 0, comm);} \\ \text{MPI\_Recv(new\_msg, size, MPI\_INT, (my\_rank+comm\_sz-1) \% comm\_sz,} \\ 0, comm, MPI\_STATUS\_IGNORE. \end{split}
```



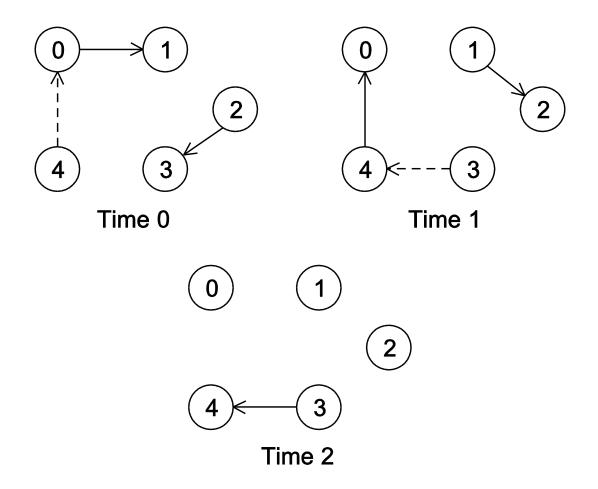
### MPI\_Sendrecv

- An alternative to scheduling the communications ourselves.
- Carries out a blocking send and a receive in a single call.
- The dest and the source can be the same or different.
- Especially useful because MPI schedules the communications so that the program won't hang or crash.

#### MPI\_Sendrecv

```
int MPI_Sendrecv(
    void*
           send_buf_p /*in */,
             send_buf_size /*in */,
    int
    MPI_Datatype send_buf_type /*in */,
    int
             dest /*in */.
    int
                send_tag /*in */,
    void*
                recv_buf_p /* out */.
    int
                recv_buf_size /*in */,
    MPI_Datatype recv_buf_type /* in */,
     int
                source /*in */.
                recv_tag /* in */,
     int
    MPI_Comm communicator /*in */,
    MPI_Status* status_p /*in */);
```

#### Safe communication with five processes



### Parallel odd-even transposition sort

```
void Merge_low(
     int my_keys[], /* in/out */
     int recv_keys[], /* in */
     int temp_keys[], /* scratch */
     int local n /* = n/p, in */) {
  int m_i, r_i, t_i;
  m i = r i = t i = 0;
  while (t_i < local_n) {
     if (my_keys[m_i] \le recv_keys[r_i]) 
        temp_keys[t_i] = my_keys[m_i];
        t i++; m i++;
     } else {
        temp_keys[t_i] = recv_keys[r_i];
       t i++; r i++;
  for (m i = 0; m i < local n; m i++)
     my_keys[m_i] = temp_keys[m_i];
  /* Merge_low */
```

# Run-times of parallel odd-even sort

	Number of Keys (in thousands)				
Processes	200	400	800	1600	3200
1	88	190	390	830	1800
2	43	91	190	410	860
4	22	46	96	200	430
8	12	24	51	110	220
16	7.5	14	29	60	130

(times are in milliseconds)

 MPI or the Message-Passing Interface is a library of functions that can be called from C, C++, or Fortran programs.

 A communicator is a collection of processes that can send messages to each other.

 Many parallel programs use the single-program multiple data or SPMD approach.

 Most serial programs are deterministic: if we run the same program with the same input we'll get the same output.

Parallel programs often don't possess this property.

 Collective communications involve all the processes in a communicator.

 When we time parallel programs, we're usually interested in elapsed time or "wall clock time".

 Speedup is the ratio of the serial run-time to the parallel run-time.

 Efficiency is the speedup divided by the number of parallel processes.

 If it's possible to increase the problem size (n) so that the efficiency doesn't decrease as p is increased, a parallel program is said to be scalable.

 An MPI program is unsafe if its correct behavior depends on the fact that MPI\_Send is buffering its input.