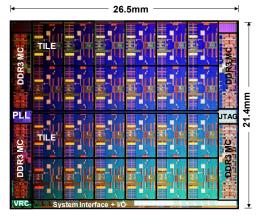
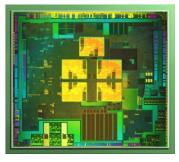


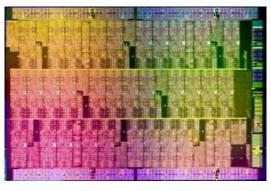
NVIDIA GTX 480 processor



Intel labs 48 core SCC processor



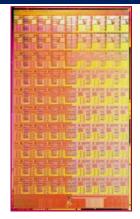
NVIDIA Tegra 3 (quad Arm Corex A9 cores + GPU)



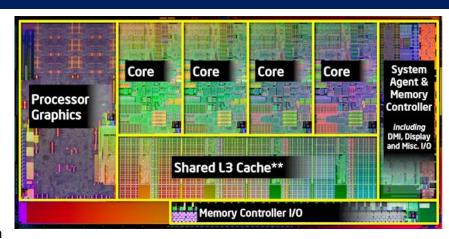
An Intel MIC processor

A hand's on introduction to Cluster Computing

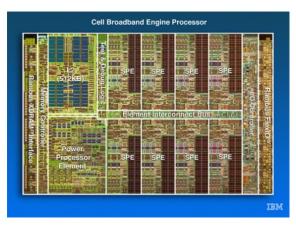
Tim Mattson (Intel Labs)



Intel Labs 80 core Research processor



Intel "Sandybridge" processor



IBM Cell Broadband engine processor

Third party names are the property of their owners

Other than the Intel lab's research processors. Die photos from UC Berkeley CS194 lecture notes

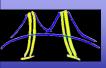
DisclaimerREAD THIS ... its very important



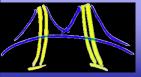
- The views expressed in this talk are those of the speakers and not their employer.
- This is an academic style talk and does not address details of any particular Intel product. You will learn nothing about Intel products from this presentation.
- This was a team effort, but if we say anything really stupid, it's our fault ... don't blame our collaborators.



Slides marked with this symbol were produced-with Kurt Keutzer and his team for CS194 ... A UC Berkeley course on Architecting parallel applications with Design Patterns.

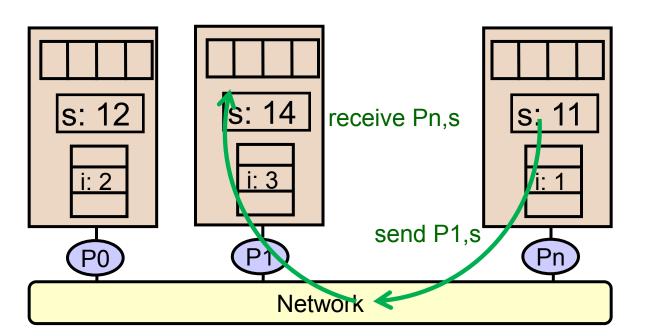


INTRODUCTION TO MPI

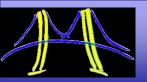


Programming Model: Message Passing

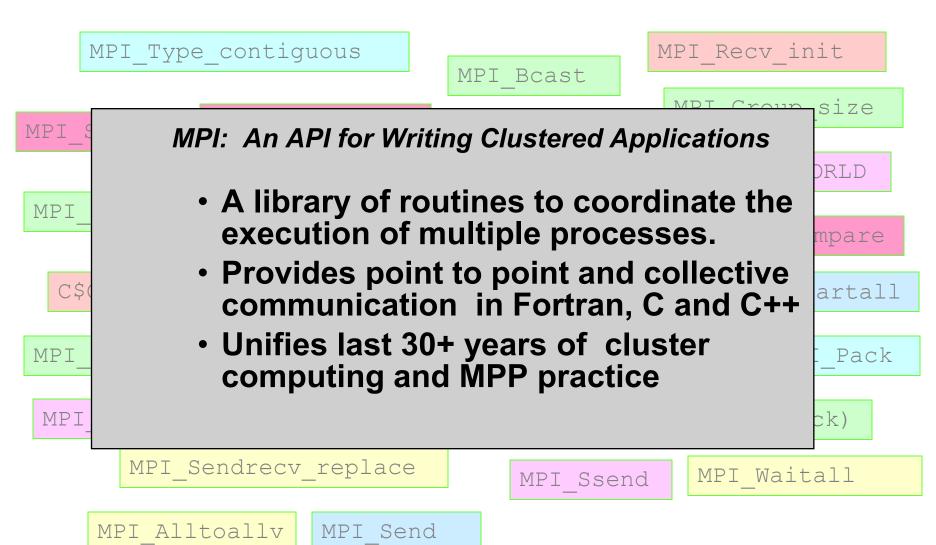
- Program consists of a collection of named processes.
 - Number of processes almost always fixed at program startup time
 - Local address space per node -- NO physically shared memory.
 - Logically shared data is partitioned over local processes.
- Processes communicate through explicit communication events.
 - Coordination is implicit in every communication event.
 - MPI (Message Passing Interface) is the most commonly used SW

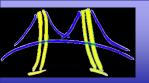


Private memory



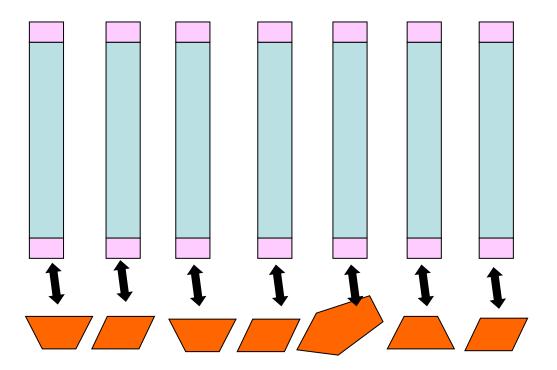
MPI: Message Passing Interface

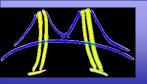




An MPI program at runtime

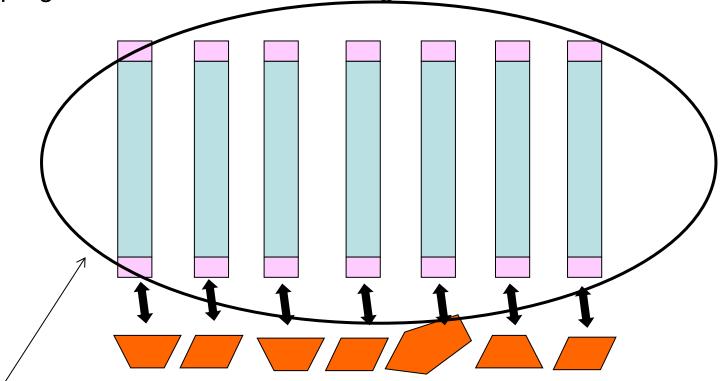
Typically, when you run an MPI program, multiple processes all running the same program are launched ... working on their own block of data.



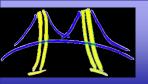


An MPI program at runtime

Typically, when you run an MPI program, multiple processes all running the same program are launched ... working on their own block of data.

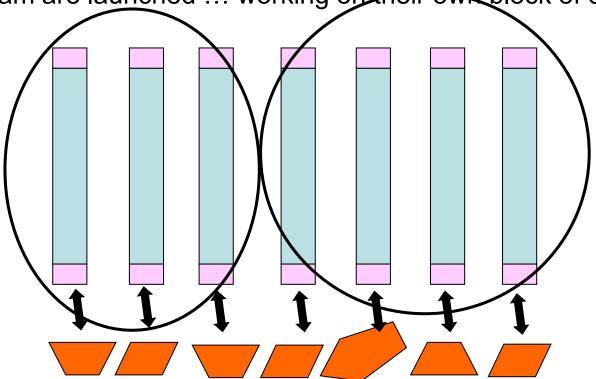


The collection of processes involved in a computation is called "a **process group**"



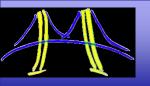
An MPI program at runtime

Typically, when you run an MPI program, multiple processes all running the same program are launched ... working on their own block of data.



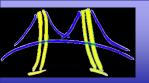
You can dynamically split a <u>process group</u> into multiple subgroups to manage how processes are mapped onto different tasks

MPI functions work within a "**context**" ... events in different contexts ... even if they share a process group ... cannot interfere with each other.



MPI Hello World

```
#include <stdio.h>
#include <mpi.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( "Hello from process %d of %d\n",
                                 rank, size );
    MPI_Finalize();
    return 0;
```



Initializing and finalizing MPI

```
int MPI Init (int* argc, char* argv[])
```

- Initializes the MPI library ... called before any other MPI functions.
- agrc and argv are the command line args passed

```
#include <stdio.h>
                          from main()
#include <mpi.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( "Hello from process %d of %d\n",
                                rank, size );
    MPI_Finalize();
                    int MPI Finalize (void)
    return 0;
```

Frees memory allocated by the MPI library ... close every MPI program with a call to MPI_Finalize

How many processes are involved?

```
int MPI Comm size (MPI Comm comm, int* size)
```

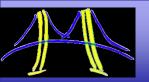
- MPI_Comm, an opaque data type called a communicator. Default context: MPI_COMM_WORLD (all processes)
- MPI_Comm_size returns the number of processes in the process group associated with the communicator

```
#include <mpi.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( "Hello from process %d of %d\n",
                                 rank, size );
    MPI_Finalize();
    return 0;
```

Communicators consist of two parts, a context and a process group.

The communicator lets me control how groups of messages interact.

The communicator lets me write modular SW ... i.e. I can give a library module its own communicator and know that it's messages can't collide with messages originating from outside the module



Which process "am I" (the rank)

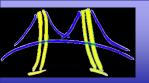
```
int MPI_Comm_rank (MPI_Comm comm, int* rank)
```

- MPI_Comm, an opaque data type, a communicator. Default context:
 MPI_COMM_WORLD (all processes)
- MPI Comm rank An integer ranging from 0 to "(num of procs)-1"

```
#incl___
#include <mpi.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( "Hello from process %d of %d\n",
                                 rank, size );
    MPI_Finalize();
    return 0;
```

Note that other than init() and finalize(), every MPI function has a communicator.

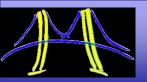
This makes sense .. You need a context and group of processes that the MPI functions impact ... and those come from the communicator.



Running the program

- On a 4 node cluster, I'd run this program (hello) as:> mpiexec –n 4 hello
- What would this program would output?

```
#include <stdio.h>
#include <mpi.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( "Hello from process %d of %d\n",
                                 rank, size );
    MPI_Finalize();
    return 0;
```



Exercise 1: Hello world

- Goal
 - To confirm that you can run a program on our cluster
- Program
 - Write a program that prints "hello world" to the screen.
 - Modify it to run as an MPI program ... with each process in the process group printing "hello world" and its rank

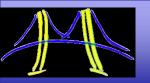
```
#include <mpi.h>
int size, rank, argc; char **argv;
MPI_Init (&argc, &argv);
MPI_Comm_rank (MPI_COMM_WORLD, &rank);
MPI_Comm_size (MPI_COMM_WORLD, &size);
MPI_Finalize();
```

To run the executable hello on 2 nodes of the cluster, type:

```
> mpiexec -n 4 a.out
```

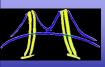
If you need to start the MPI deamon (mpd) go to your home directory and type:

- > touch .mpd.conf
- > chmod 600 .mpd.conf

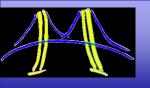


Running the program

```
#include <stdio.h>
#include <mpi.h>
int main (int argc, char **a
    int rank, size;
    MPI Init (&argc, &argv);
    MPI Comm rank (MPI COMM WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);
    printf( "Hello from process %d of %d\n",
                                 rank, size );
    MPI_Finalize();
    return 0;
```



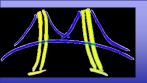
MPI FOR BULK SYNCHRONOUS PROGRAMS



Sending and Receiving Data

- MPI_Send performs a blocking send of the specified data ("count" copies of type "datatype," stored in "buf") to the specified destination (rank "dest" within communicator "comm"), with message ID "tag"
- MPI_Recv performs a blocking receive of specified data from specified source whose parameters match the send; information about transfer is stored in "status"

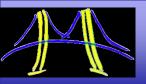
By "blocking" we mean the functions return as soon as the buffer, "buf", can be safely used.



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.

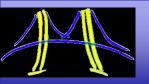


The data in a message: datatypes

- The data in a message to send or receive is described by a triple:
 - (address, count, datatype)
- An MPI datatype is defined as:
 - Predefined, simple data type from the language (e.g., MPI_DOUBLE)
 - Complex data types (contiguous blocks or even custom types.
- E.g. ... A particle's state is defined by its 3 coordinates and 3 velocities
 MPI_Datatype PART;
 MPI_Type_contiguous(6, MPI_DOUBLE, &PART);
 - MPI_Type_commit(&PART);
- You can use this data type in MPI functions, for example, to send data for a single particle:

MPI_Send (buff, 1, PART, Dest, tag, MPI_COMM_WORLD);
address

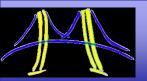
Datatype



Receiving the right message

- The receiving process identifies messages with the double :
 - (source, tag)
- Where:
 - Source is the rank of the sending process
 - Tag is a user-defined integer to help the receiver keep track of different messages from a single source

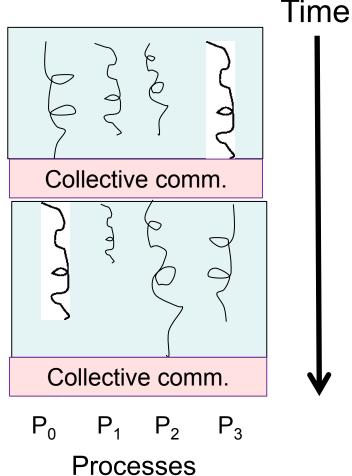
- Can relax tag checking by specifying MPI_ANY_TAG as the tag in a receive.
- Can relax source checking by specifying MPI_ANY_SOURCE MPI_Recv (buff, 1, PART, MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &status);
- This is a useful way to insert race conditions into an MPI program

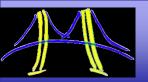


A typical pattern with MPI Programs

- Many MPI applications have few (if any) sends and receives. They use the following very common pattern:
 - Use the Single Program Multiple Data pattern
 - Each process maintains a local view of the global data
 - A problem broken down into phases each of which is composed of two subphases:
 - Compute on local view of data
 - Communicate to update global view on all processes (collective communication).
 - Continue phases until complete

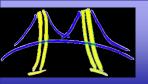
This is a subset or the SPMD pattern sometimes referred to as the Bulk Synchronous pattern.



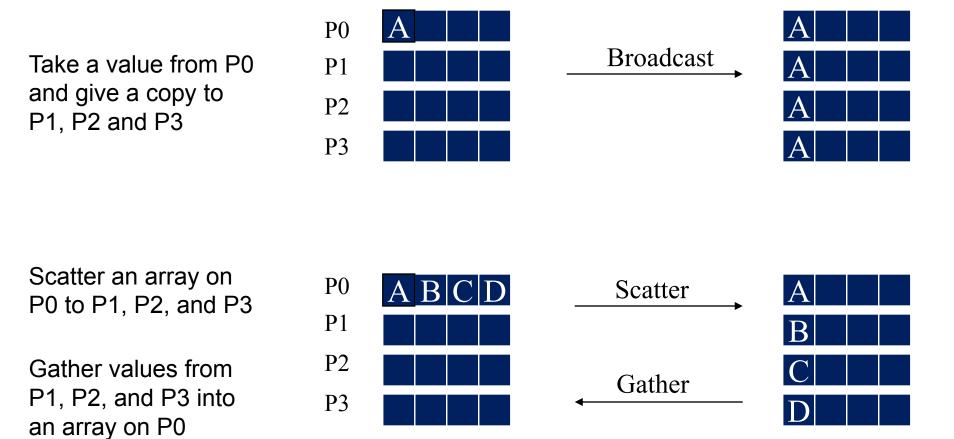


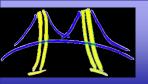
MPI Collective Routines

- Collective communications: called by all processes in the group to create a global result and share with all participating processes.
 - Allgather, Allgatherv, Allreduce, Alltoall, Alltoallv, Bcast, Gather, Gatherv, Reduce, Reduce scatter, Scan, Scatter, Scatterv
- Notes:
 - Allreduce, Reduce, Reduce_scatter, and Scan use the same set of built-in or user-defined combiner functions.
 - Routines with the "All" prefix deliver results to all participating processes
 - Routines with the "v" suffix allow chunks to have different sizes
- Global synchronization is available in MPI
 - MPI_Barrier(comm)
- Blocks until all processes in the group of the communicator comm call it.

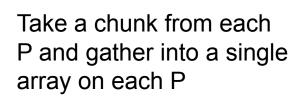


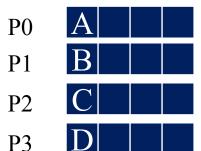
Collective Data Movement

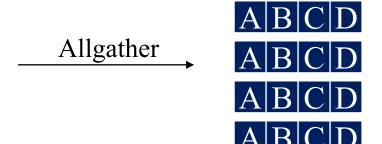




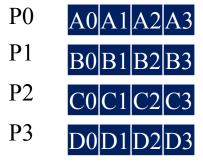
More Collective Data Movement

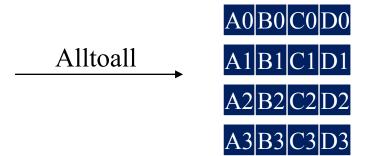


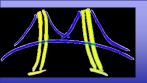




Take arrays on each P and spread them out to arrays on each P







Collective Computation

Take values on each P and combine them with an op (such as add) into a single value on one P.

A P0

P1

P2

P3

Reduce

Take values on each P and combine them with a scan operation and spread the scan array out among all P.

P0

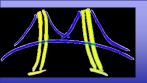
P1

P2

P3

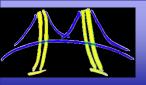
Scan

AB ABC



MPI_BCAST Example

```
#include <mpi.h>
                                                 MPI COMM WORLD
                                                  Rank 0
int main(int argc, char *argv[]) {
 int nprocs, myrank, msg[4] = \{0,0,0,0,0\};
                                                  msg 1
MPI Init(&argc, &argv);
                                                  Rank 1
 MPI Comm size (MPI COMM WORLD, &nprocs);
MPI Comm rank(MPI COMM WORLD, &myrank);
                                                  msg ()
 if (myrank == 0) msg[0] = 1;
                                                                  M
                                                  msg
 else
                  msg[0] = 0;
MPI_Bcast(msg, 4, MPI_INT, 0, MPI_COMM_WORLD);
                                                  Rank 2
                                                  msg 0
MPI Finalize();
                                                  msg 1
```

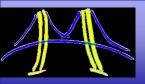


Reduction

- MPI_Reduce performs specified reduction operation on specified data from all processes in communicator, places result in process "root" only.
- MPI_Allreduce places result in all processes (avoid unless necessary)

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

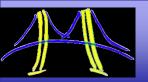
Operation	Function
MPI_BAND	Bitwise AND
MPI_LOR	Logical OR
MPI_BOR	Bitwise OR
MPI_LXOR	Logical exclusive OR
MPI_BXOR	Bitwise exclusive OR
User-defined	It is possible to define new reduction operations



MPI_REDUCE Example MPI_COMM_WORLD

```
#include <mpi.h>
int main(int argc, char* argv[]) {
  int msg, sum, nprocs, myrank;
 MPI Init(&argc, &argv);
 MPI Comm size(MPI COMM WORLD, &nprocs);
 MPI Comm rank(MPI COMM WORLD, &myrank);
  sum = 0;
 msg = myrank;
 MPI Reduce (&msg, &sum, 1, MPI INT, MPI SUM, 0,
             MPI COMM WORLD);
 MPI Finalize();
```

```
Rank 0
 sum 3
msg ()
Rank 1
                  REDUCE
msg
Rank 2
msg
```



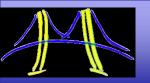
Exercise 2: Pi Program

- Goal
 - To write a simple Bulk Synchronous, SPMD program
- Program
 - Start with the provided "pi program" and using an MPI reduction, write a parallel version of the program.

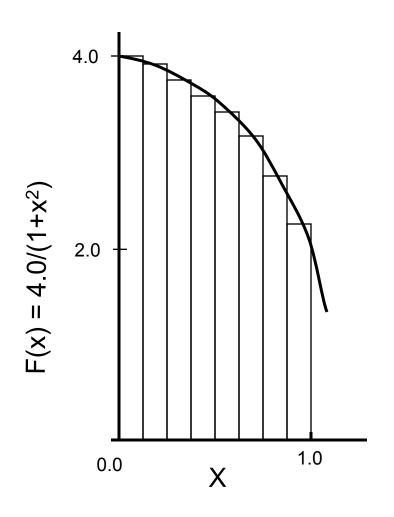
Operation	Function	
MPI_SUM	Summation	
MPI_PROD	Product	

```
#include <mpi.h>
int size, rank, argc; char **argv;
MPI_Init (&argc, &argv);
MPI_Comm_rank (MPI_COMM_WORLD, &rank);
MPI_Comm_size (MPI_COMM_WORLD, &size);
MPI_Finalize();
```

MPI Data Type	C Data Type
MPI_DOUBLE	double
MPI_FLOAT	float
MPI_INT	int
MPI_LONG	long



Example Problem: Numerical Integration



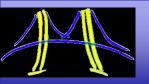
Mathematically, we know that:

$$\int_{0}^{1} \frac{4.0}{(1+x^2)} dx = \pi$$

We can approximate the integral as a sum of rectangles:

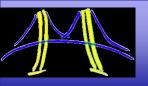
$$\sum_{i=0}^{N} F(x_i) \Delta x \approx \pi$$

Where each rectangle has width Δx and height $F(x_i)$ at the middle of interval i.



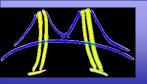
PI Program: an example

```
static long num_steps = 100000;
double step;
void main ()
        int i; double x, pi, sum = 0.0;
        step = 1.0/(double) num steps;
        x = 0.5 * step;
        for (i=0;i<= num_steps; i++){
               x+=step;
               sum += 4.0/(1.0+x*x);
        pi = step * sum;
```



Pi program in MPI

```
#include <mpi.h>
void main (int argc, char *argv[])
       int i, my id, numprocs; double x, pi, step, sum = 0.0;
       step = 1.0/(double) num steps;
       MPI Init(&argc, &argv);
       MPI Comm Rank(MPI COMM WORLD, &my id);
       MPI Comm Size(MPI COMM WORLD, &numprocs);
       my steps = num steps/numprocs;
       for (i=my id*my steps; i<(my id+1)*my steps; i++)
                x = (i+0.5)*step;
                                              Sum values in "sum" from
                sum += 4.0/(1.0+x*x);
                                              each process and place it
                                                 in "pi" on process 0
       sum *= step;
       MPI_Reduce(&sum, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,
               MPI COMM WORLD);
```



MPI Pi program performance

A

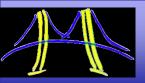
Pi program in MPI

```
#include <mpi.h>
void main (int argc, char *argv∏)
       int i, my id, numprocs; double x, pi, step, sum
       step = 1.0/(double) num steps;
       MPI Init(&argc, &argy);
       MPI Comm Rank(MPI COMM WORLD,
       MPI Comm Size(MPI COMM WORLD, &
       for (i=my_id; i<num_steps; ; i=i+numprocs)
                x = (i+0.5)*step;
                sum += 4.0/(1.0+x*x):
       sum *= step ;
       MPI Reduce(&sum, &pi, 1, MPI DOUBLE, MPI SUM, 0,
               MPI COMM WORLD)
```

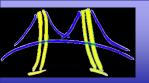
Thread	OpenMP	OpenMP	MPI
or procs	SPMD	PI Loop	
	critical		
1	0.85	0.43	0.84
2	0.48	0.23	0.48
3	0.47	0.23	0.46
4	0.46	0.23	0.46

Note: OMP loop used a Blocked loop distribution. The others used a cyclic distribution. Serial .. 0.43.

^{*}Intel compiler (icpc) with −O3 on Apple OS X 10.7.3 with a dual core (four HW thread) Intel® CoreTM i5 processor at 1.7 Ghz and 4 Gbyte DDR3 memory at 1.333 Ghz.

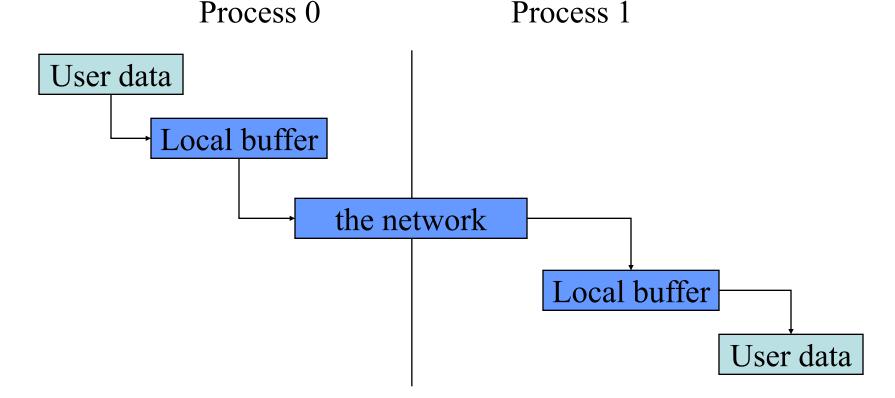


UNDERSTANDING MESSAGE PASSING WITH MPI



Buffers

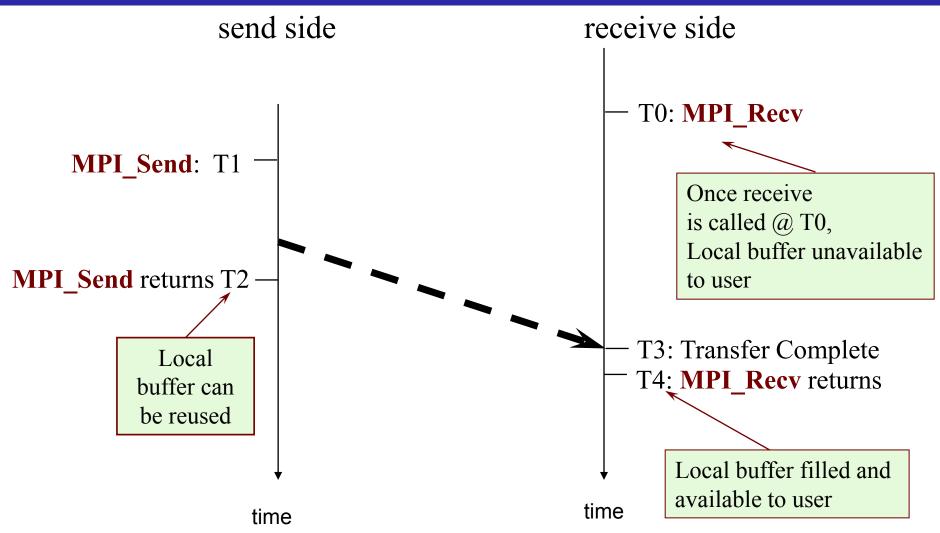
- Message passing has a small set of primitives, but there are subtleties
 - Buffering and deadlock
 - Deterministic execution
 - Performance
- When you send data, where does it go? One possibility is:



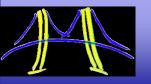
Derived from: Bill Gropp, UIUC

Blocking Send-Receive Timing Diagram

(Receive before Send)



It is important to post the receive before sending, for highest performance.



Sources of Deadlocks

- Send a large message from process 0 to process 1
 - If there is insufficient storage at the destination, the send must wait for the user to provide the memory space (through a receive)
- What happens with this code?

Process 0	Process 1		
Send(1)	Send(0)		
Recv(1)	Recv(0)		

 This code could deadlock ... it depends on the availability of system buffers in which to store the data sent until it can be received

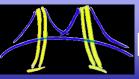
Some Solutions to the "deadlock" Problem



Process 0	Process 1
Send(1)	Recv(0)
Recv(1)	Send(0)

• Supply receive buffer at same time as send:

Sendrecv(1)	Sendrecv(0)		
Process 0	Process 1		



More Solutions to the "unsafe" Problem

Supply a sufficiently large buffer in the send function

Process 0

Process 1

Bsend(1)

Bsend(0)

Recv(1)

Recv(0)

Use non-blocking operations:

Process 0

Process 1

Isend(1)

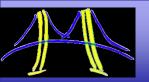
Isend(0)

Irecv(1)

Irecv(0)

Waitall

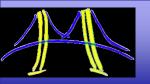
Waitall



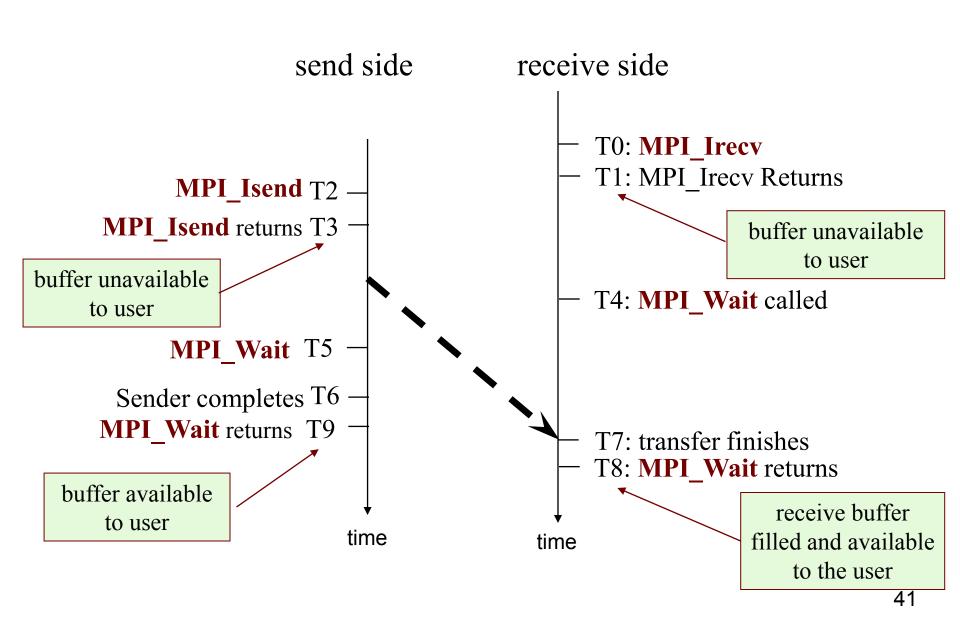
Non-Blocking Communication

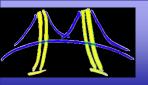
- Non-blocking operations return immediately and pass "request handles" that can be waited on and queried
 - MPI_ISEND(start, count, datatype, dest, tag, comm, request)
 - MPI_IRECV(start, count, datatype, src, tag, comm, request)
 - MPI_WAIT(request, status)
- One can also test without waiting using MPI_TEST
 - MPI_TEST(request, flag, status)
- Anywhere you use MPI_Send or MPI_Recv, you can use the pair of MPI_Isend/MPI_Wait or MPI_Irecv/MPI_Wait

Non-blocking operations are extremely important ... they allow you to overlap computation and communication.



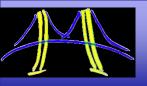
Non-Blocking Send-Receive Diagram





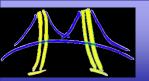
Example: shift messages around a ring (part 1 of 2)

```
#include <stdio.h>
#include <mpi.h>
int main(int argc, char **argv)
 int num, rank, size, tag, next, from;
 MPI Status status1, status2;
 MPI Request req1, req2;
 MPI Init(&argc, &argv);
 MPI Comm rank( MPI COMM WORLD, &rank);
 MPI Comm size(MPI COMM WORLD, &size);
 tag = 201;
 next = (rank+1) % size;
 from = (rank + size - 1) % size;
 if (rank == 0) {
  printf("Enter the number of times around the ring: ");
  scanf("%d", &num);
  printf("Process %d sending %d to %d\n", rank, num, next);
  MPI Isend(&num, 1, MPI INT, next, tag, MPI COMM WORLD,&req1);
  MPI Wait(&req1, &status1);
```



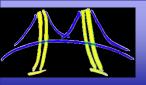
Example: shift messages around a ring (part 2 of 2)

```
do {
 MPI Irecv(&num, 1, MPI INT, from, tag, MPI COMM WORLD, &req2);
 MPI Wait(&req2, &status2);
 printf("Process %d received %d from process %d\n", rank, num, from);
 if (rank == 0) {
  num--;
   printf("Process 0 decremented number\n");
 printf("Process %d sending %d to %d\n", rank, num, next);
 MPI Isend(&num, 1, MPI INT, next, tag, MPI COMM WORLD, &req1);
 MPI Wait(&req1, &status1);
} while (num != 0);
if (rank == 0) {
 MPI Irecv(&num, 1, MPI INT, from, tag, MPI COMM WORLD, &req2);
 MPI Wait(&reg2, &status2);
MPI_Finalize();
return 0;
```



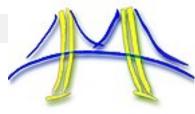
Exercise 3: Ring program

- Goal
 - Explore other modes of message passing in MPI
- Program
 - Start with the basic ring program we provide. Run it for a range of message sizes and notes what happens for large messages.
 - If the program deadlocks (and it should) figure out why and how to fix it.
 - Try a range of message passing functions to understand how they work.



MPI AND THE GEOMETRIC DECOMPOSITION PATTERN

Example: finite difference methods



- Solve the heat diffusion equation in 1 D:
 - \square u(x,t) describes the temperature field
 - □ We set the heat diffusion constant to one
 - □ Boundary conditions, constant u at endpoints.

$$\frac{\partial^2 u}{\partial x^2} = \frac{\partial u}{\partial t}$$

map onto a mesh with stepsize h and k

$$x_i = x_0 + ih \qquad t_i = t_0 + ik$$

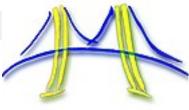
 Central difference approximation for spatial derivative (at fixed time)

$$\frac{\partial^2 u}{\partial x^2} = \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2}$$

■ Time derivative at t = tⁿ⁺¹

$$\frac{du}{dt} = \frac{u^{n+1} - u^n}{k}$$

Example: Explicit finite differences



Combining time derivative expression using spatial derivative at t = tⁿ

$$\frac{u_j^{n+1} - u_j^n}{k} = \frac{u_{j+1}^n - 2u_j^n + u_{j-1}^n}{h^2}$$

Solve for u at time n+1 and step j

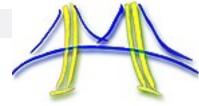
$$u_{j}^{n+1} = (1-2r)u_{j}^{n} + ru_{j-1}^{n} + ru_{j+1}^{n} \qquad r = \frac{k}{h^{2}}$$

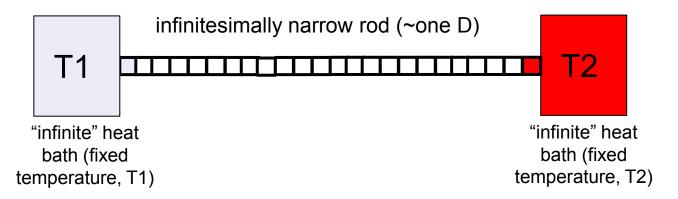
The solution at t = t_{n+1} is determined explicitly from the solution at t = t_n (assume u[t][0] = u[t][N] = Constant for all t).

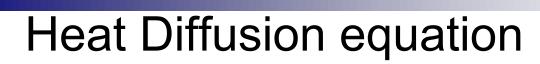
```
for (int t = 0; t < N_STEPS-1; ++t)
  for (int x = 1; x < N-1; ++x)
      u[t+1][x] = u[t][x] + r*(u[t][x+1] - 2*u[t][x] + u[t][x-1]);</pre>
```

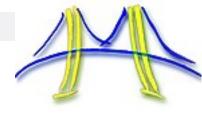
 Explicit methods are easy to compute ... each point updated based on nearest neighbors. Converges for r<1/2.

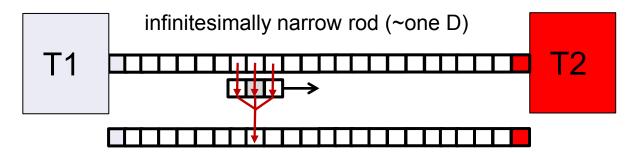




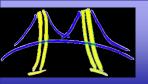






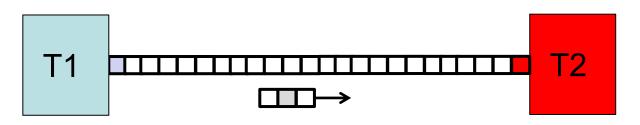


Pictorially, you are sliding a three point "stencil" across the domain (u[t]) and computing a new value of the center point (u[t+1]) at each stop.



return 0;

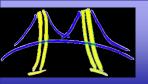
Heat Diffusion equation

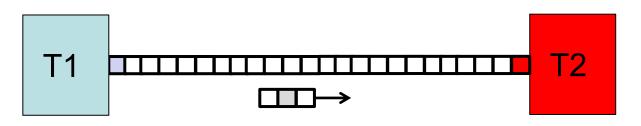


```
int main()
{
    double *u = malloc (sizeof(double) * (N));
    double *up1 = malloc (sizeof(double) * (N));

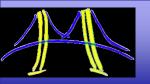
initialize_data(uk, ukp1, N, P); // init to zero, set end temperatures
for (int t = 0; t < N_STEPS; ++t){
    for (int x = 1; x < N-1; ++x)
        up1[x] = u[x] + (k / (h*h)) * (u[x+1] - 2*u[x] + u[x-1]);

    A well known trick with 2 arrays so I
    don't overwrite values from step k-1
    as I fill in for step k</pre>
```





```
How would
int main()
                                                      you parallelize
  double *u = malloc (sizeof(double) * (N));
                                                      this program?
  double *up1 = malloc (sizeof(double) * (N));
  initialize data(uk, ukp1, N, P); // init to zero, set end temperatures
  for (int t = 0; t < N STEPS; ++t){
     for (int x = 1; x < N-1; ++x)
         up1[x] = u[x] + (k / (h*h)) * (u[x+1] - 2*u[x] + u[x-1]);
     temp = up1; up1 = u; u = temp;
return 0;
```

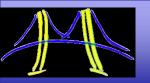


Seven strategies for parallelizing software

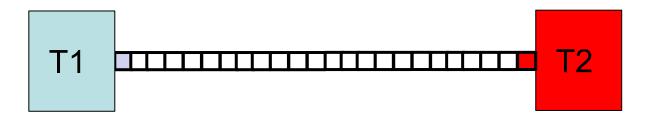
- These seven strategies for parallelizing software give us:
 - Names: so we can communicate better
 - Categories: so we can gather and share information
 - A palette (like an artist's palette) of approaches that is:
 - Necessary: we should consider them all and
 - Sufficient: once we have considered them all then we don't' have to worry that we forgot something

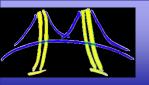
Parallel Algorithm Strategy Patterns

Task-Parallelism Divide and Conquer Data-Parallelism Pipeline Discrete-Event Geometric-Decomposition Speculation



Start with our original picture of the problem ... a one dimensional domain with end points set at a fixed temperature.



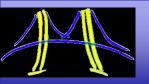


Seven strategies for parallelizing software

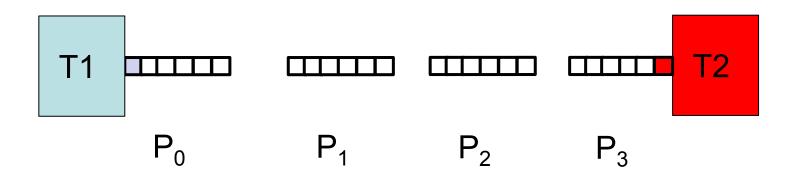
- These seven strategies for parallelizing software give us:
 - Names: so we can communicate better
 - Categories: so we can gather and share information
 - A palette (like an artist's palette) of approaches that is:
 - Necessary: we should consider them all and
 - Sufficient: once we have considered them all then we don't' have to worry that we forgot something

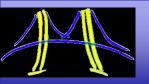
Parallel Algorithm Strategy Patterns

Task-Parallelism Divide and Conquer Data-Parallelism Pipeline Discrete-Event
Geometric-Decomposition
Speculation

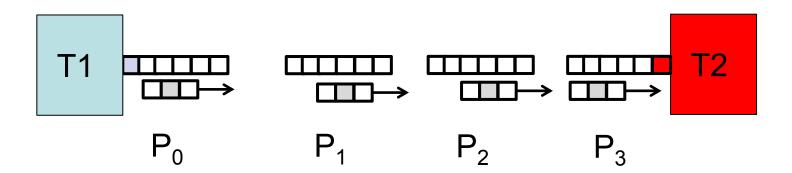


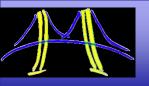
Break it into chunks assigning one chunk to each process.



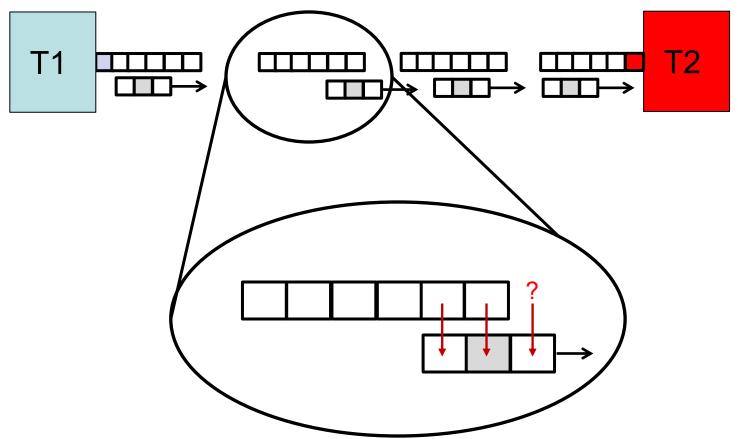


Each process works on it's own chunk ... sliding the stencil across the domain to updates its own data.

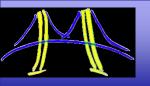




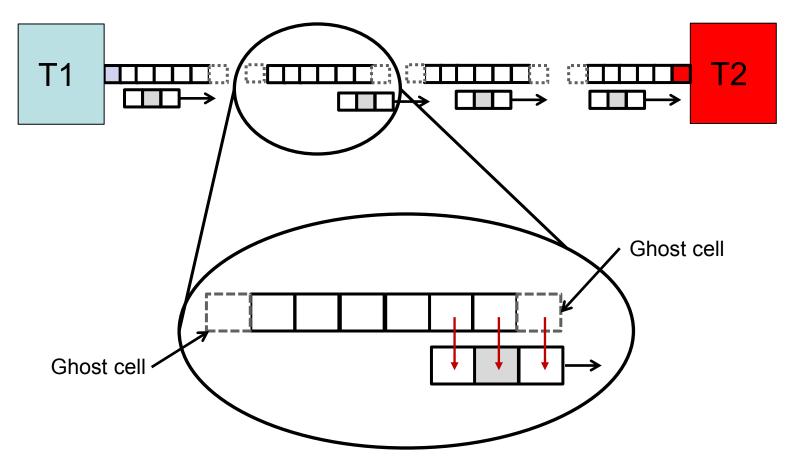
What about the ends of each chunk ... where the stencil will run off the end and hence have missing values for the computation?

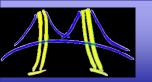


57



We add ghost cells to the ends of each chunk, update them with the required values from neighbor chunks at each time step ... hence giving the stencil everything it needs on any given chunk to update all of its values.



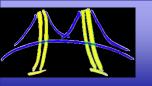


SPMD: Single Program Mulitple Data

- Run the same program on P processing elements where P can be arbitrarily large.
- Use the rank ... an ID ranging from 0 to (P-1) ... to select between a set of tasks and to manage any shared data structures.

This pattern is very general and has been used to support most (if not all) the algorithm strategy patterns.

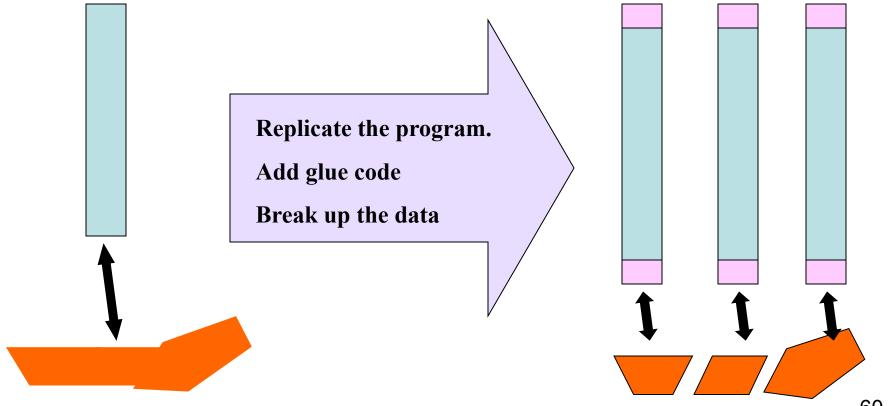
MPI programs almost always use this pattern ... it is probably the most commonly used pattern in the history of parallel programming.

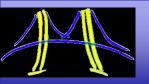


How do people use MPI? The SPMD Design Pattern

A sequential program working on a data set

- •A single program working on a decomposed data set.
- •Use Node ID and numb of nodes to split up work between processes
- Coordination by passing messages.





MPI_Finalize();

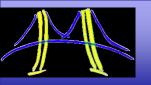
return 0;

MPI Init (&argc, &argv);

Heat Diffusion MPI Example

```
MPI Comm size (MPI COMM WORLD, &P);
MPI Comm rank (MPI COMM WORLD, &myID);
double *u = malloc (sizeof(double) * (2 + N/P)) // include "Ghost Cells"
double *up1 = malloc (sizeof(double) * (2 + N/P)); // to hold values
                                                    // from my neighbors
initialize_data(uk, ukp1, N, P);
for (int t = 0; t < N STEPS; ++t){
  if (myID != 0) MPI Send (&u[1], 1, MPI DOUBLE, myID-1, 0, MPI COMM WORLD);
  if (myID != P-1) MPI_Recv (&u[N/P+1], 1, MPI_DOUBLE, myID+1, 0, MPI_COMM_WORLD, &status);
  if (myID != P-1) MPI Send (\&u[N/P], 1, MPI DOUBLE, myID+1, 0, MPI COMM WORLD);
  if (myID != 0) MPI Recv (&u[0], 1, MPI DOUBLE, myID-1, 0, MPI COMM WORLD, &status);
 for (int x = 2; x < N/P; ++x)
    up1[x] = u[x] + (k / (h*h)) * (u[x+1] - 2*u[x] + u[x-1]);
 if (myID != 0)
    up1[1] = u[1] + (k / (h*h)) * (u[1+1] - 2*u[1] + u[1-1]);
 if (myID != P-1)
    up1[N/P] = u[N/P] + (k/(h*h)) * (u[N/P+1] - 2*u[N/P] + u[N/P-1]);
 temp = up1; up1 = u; u = temp;
} // End of for (int t ...) loop
```

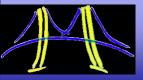
We write/explain this part first and then address the communication and data structures



return 0;

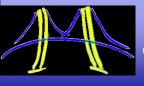
Heat Diffusion MPI Example

```
Update array values using local data
// Compute interior of each "chunk"
                                               and values from ghost cells.
  for (int x = 2; x < N/P; ++x)
    up1[x] = u[x] + (k / (h*h)) * (u[x+1] - 2*u[x] + u[x-1]);
// update edges of each chunk keeping the two far ends fixed
// (first element on Process 0 and the last element on process P-1).
  if (myID != 0)
                                                                   u[0] and u[N/P+1]
    up1[1] = u[1] + (k / (h*h)) * (u[1+1] - 2*u[1] + u[1-1]);
                                                                     are the ghost
                                                                         cells
  if (myID != P-1)
    up1[N/P] = u[N/P] + (k/(h*h)) * (u[N/P+1] - 2*u[N/P] + u[N/P-1]);
// Swap pointers to prepare for next iterations
  temp = up1; up1 = u; u = temp;
} // End of for (int t ...) loop
                                       Note I was lazy and assumed N was evenly
                                       divided by P. Clearly, I'd never do this in a
                                       "real" program.
MPI Finalize();
```



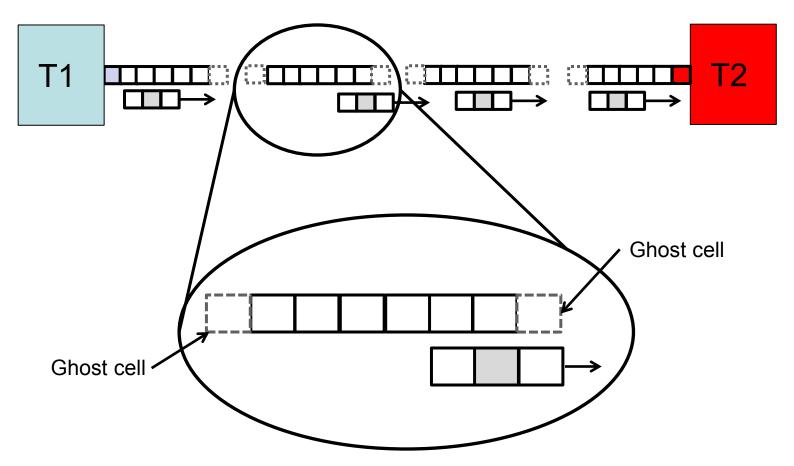
Heat Diffusion MPI Example

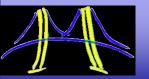
```
MPI Init (&argc, &argv);
                                         1D PDE solver ... the simplest "real" message
                                        passing code I can think of. Note: edges of
MPI Comm size (MPI COMM WORLD, &P);
MPI_Comm_rank (MPI_COMM_WORLD, &myID); domain held at a fixed temperature
double *u = malloc (sizeof(double) * (2 + N/P)) // include "Ghost Cells"
double *up1 = malloc (sizeof(double) * (2 + N/P)); // to hold values
                                                      // from my neighbors
initialize_data(uk, ukp1, N, P);
for (int t = 0; t < N_STEPS; ++t){
  if (myID != 0) Send my "left" boundary value to the neighbor on my "left"
    MPI Send (&u[1], 1, MPI DOUBLE, myID-1, 0, MPI COMM WORLD);
                         Receive my "right" ghost cell from the neighbor to my "right"
  if (myID != P-1)
    MPI Recv (&u[N/P+1], 1, MPI DOUBLE, myID+1, 0, MPI COMM WORLD, &status);
  if (myID != P-1) Send my "right" boundary value to the neighbor to my "right"
    MPI Send (&u[N/P], 1, MPI DOUBLE, myID+1, 0, MPI COMM WORLD);
                          Receive my "left" ghost cell from the neighbor to my "left"
  if (myID != 0)
    MPI_Recv (&u[0], 1, MPI_DOUBLE, myID-1, 0,MPI_COMM_WORLD, &status);
/* continued on previous slide */
                                                                             63
```



The Geometric Decomposition Pattern

- This is an instance of a very important design pattern ... the Geometric decomposition pattern.
- We will cover this pattern in more detail in a later lecture.





Partitioned Array Pattern

Problem:

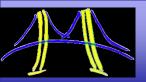
Arrays often need to be partitioned between multiple UEs. How can we do this so the resulting program is both readable and efficient?

Forces

- Large number of small blocks organized to balance load.
- Able to specialize organization to different platforms/problems.
- Understandable indexing to make programming easier.

Solution:

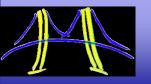
- Express algorithm in blocks
- Abstract indexing inside mapping functions ... programmer works in an index space natural to the domain, functions map into distribution needed for efficient execution.
- The text of the pattern defines some of these common mapping functions (which can get quite confusing ... and in the literature are usually left as "an exercise for the reader").



Partitioned Arrays

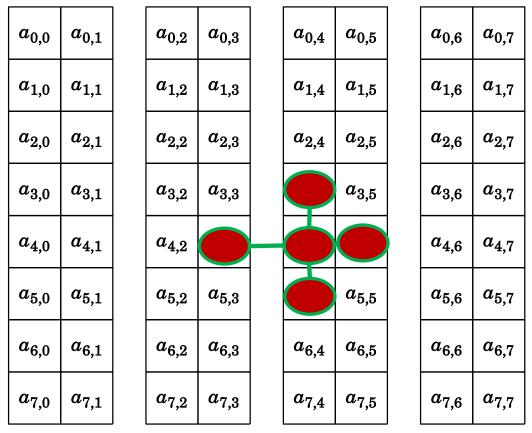
- Realistic problems are 2D or 3D; require move complex data distributions.
- We need to parallelize the computation by partitioning this index space
- Example: Consider a 2D domain over which we wish to solve a PDE using an explicit finite difference solver. The figure shows a five point stencil ... update a value based on its value and its 4 neighbors.
- Start with an array →

$a_{0,0}$	$a_{0,1}$	$a_{0,2}$	$a_{0,3}$	$a_{0,4}$	$a_{0,5}$	$a_{0,6}$	$a_{0,7}$
$a_{1,0}$	$a_{1,1}$	$a_{1,2}$	$a_{1,3}$	$a_{1,4}$	$a_{1,5}$	$a_{1,6}$	$a_{1,7}$
$\boxed{a_{2,0}}$	$a_{2,1}$	$a_{2,2}$	$a_{2,3}$	$a_{2,4}$	$a_{2,5}$	$a_{2,6}$	$a_{2,7}$
$a_{3,0}$	$a_{3,1}$	$a_{3,2}$	$a_{3,3}$		$a_{3,5}$	$a_{3,6}$	$a_{3,7}$
$a_{4,0}$	$a_{4,1}$	$a_{4,2}$				$a_{4,6}$	$a_{4,7}$
$a_{5,0}$	$a_{5,1}$	$a_{5,2}$	$a_{5,3}$		$a_{5,5}$	$a_{5,6}$	$a_{5,7}$
$a_{6,0}$	$a_{6,1}$	$a_{6,2}$	$a_{6,3}$	$a_{6,4}$	$a_{6,5}$	$a_{6,6}$	$a_{6,7}$
$a_{7,0}$	$a_{7,1}$	$a_{7,2}$	$a_{7,3}$	$a_{7,4}$	$a_{7,5}$	$a_{7,6}$	$a_{7,7}$



Partitioned Arrays: Column block distribution

- Split the non-unit-stride dimension (P-1) times to produce P chunks, assign the ith chunk to P_i WIth N = n * n, P = p * p
- In a 2D finite-differencing program (exchange edges), how much do we have to communicate? 2*n = 2*sqrt(N) messages per processor



UE(2)

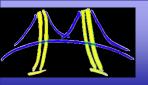
UE(1)

P is the # of processors

UE(0)

UE = unit of execution ... think of it as a generic term for "process or thread"

UE(3)



Partitioned Arrays: Block distribution

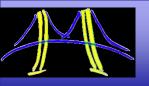
If we parallelize in both dimensions, then we have (n/p)² elements per processor, and we need to send 4*(n/p) = 4 *sqrt(N/P) messages from each processor. Asymptotically better than 2*sqrt(N).

UE(0, 0)				UE(0, 1)			
$a_{0,0}$	$a_{0,1}$	$a_{0,2}$	$a_{0,3}$	$a_{0,4}$	$a_{0,5}$	$a_{0,6}$	$a_{0,7}$
$a_{1,0}$	$a_{1,1}$	$a_{1,2}$	$a_{1,3}$	$a_{1,4}$	$a_{1,5}$	$a_{1,6}$	$a_{1,7}$
$a_{2,0}$	$a_{2,1}$	$a_{2,2}$	$a_{2,3}$	$a_{2,4}$	$a_{2,5}$	$a_{2,6}$	$a_{2,7}$
$a_{3,0}$	$a_{3,1}$	$a_{3,2}$	$a_{3,3}$		$a_{3,5}$	$a_{3,6}$	$a_{3,7}$
$\boxed{a_{4,0}}$	$a_{4,1}$	$a_{4,2}$				$a_{4,6}$	$a_{4,7}$
$a_{5,0}$	$a_{5,1}$	$a_{5,2}$	$a_{5,3}$		$a_{5,5}$	$a_{5,6}$	$a_{5,7}$
$a_{6,0}$	$a_{6,1}$	$a_{6,2}$	$a_{6,3}$	$a_{6,4}$	$a_{6,5}$	$a_{6,6}$	$a_{6,7}$
$a_{7,0}$	$a_{7,1}$	$a_{7,2}$	$a_{7,3}$	$a_{7,4}$	$a_{7,5}$	$a_{7,6}$	$a_{7,7}$

UE(1, 1)

UE(1, 0)

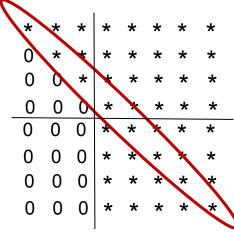
P is the # of processors



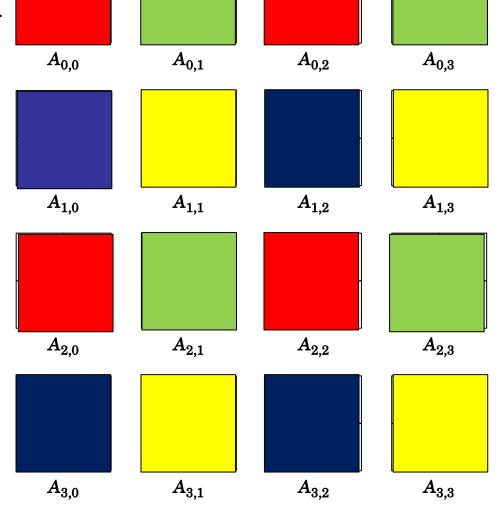
Partitioned Arrays:

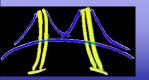
block cyclic distribution

LU decomposition (A= LU) .. Move down the diagonal transform rows to "zero the column" below the diagonal.



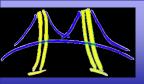
- Zeros fill in the right lower triangle of the matrix ... less work to do.
- Balance load with cyclic distribution of blocks of A mapped onto a grid of nodes (2x2 in this case ... colors show the mapping to nodes).



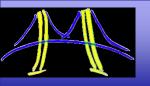


Exercise 4: Transpose

- Goal
 - Explore interaction of partitioned arrays and message passing
- Program
 - We provide a matrix transposition program ... which is one of the simplest examples of a program based on partitioned arrays.
 - Notice how the SPMD pattern interacts with the partitioned array pattern.
 - Modify the program to use isend/irecv and overlap communication with local transpose to maximize aggregate bandwidth



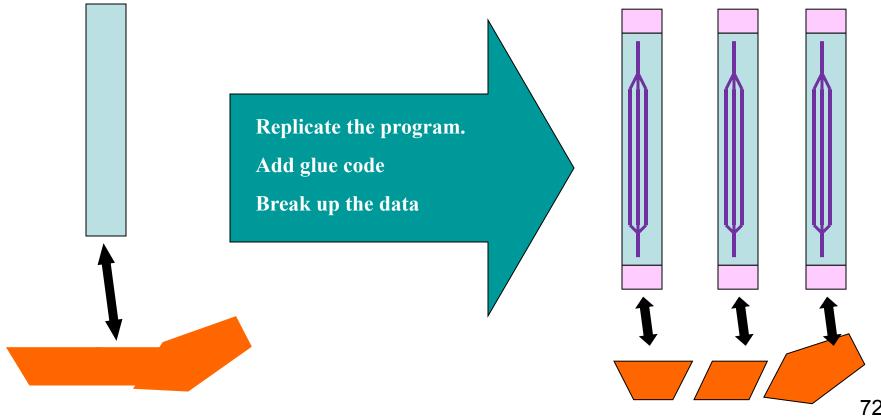
MIXING MPI AND OPENMP

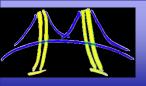


How do people mix MPI and OpenMP?

A sequential program working on a data set

- •Create the MPI program with its data decomposition.
- Use OpenMP inside each MPI process.





Pi program with MPI and OpenMP

```
int i, my id, numprocs; double x, pi, step, sum = 0.0;
                       step = 1.0/(double) num steps;
                       MPI_Init(&argc, &argv);
                       MPI Comm Rank(MPI COMM WORLD, &my id);
Get the MPI
                       MPI_Comm_Size(MPI_COMM_WORLD, &numprocs);
part done
                       my_steps = num_steps/numprocs;
first, then add
               #pragma omp parallel for reduction(+:sum) private(x)
OpenMP
                       for (i=my id*my steps; i < (m id+1)*my steps; i++)
pragma
where it
                                 x = (i+0.5)*step;
makes sense
                                 sum += 4.0/(1.0+x*x);
to do so
                       sum *= step;
```

#include <mpi.h>

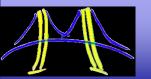
#include "omp.h"

void main (int argc, char *argv[])

73

MPI COMM WORLD);

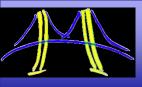
MPI_Reduce(&sum, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,



Key issues when mixing OpenMP and MPI

- 1. Messages are sent to a process not to a particular thread.
 - Not all MPIs are threadsafe. MPI 2.0 defines threading modes:
 - MPI_Thread_Single: no support for multiple threads
 - MPI_Thread_Funneled: Mult threads, only master calls MPI
 - MPI_Thread_Serialized: Mult threads each calling MPI, but they
 do it one at a time.
 - MPI_Thread_Multiple: Multiple threads without any restrictions
 - Request and test thread modes with the function:
 MPI_init_thread(desired_mode, delivered_mode, ierr)
- 2. Environment variables are not propagated by mpirun. You'll need to broadcast OpenMP parameters and set them with the library routines.

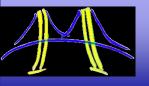
74



Dangerous Mixing of MPI and OpenMP

The following will work only if MPI_Thread_Multiple is supported ... a level of support I wouldn't depend on.

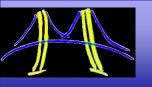
```
MPI Comm Rank(MPI COMM WORLD, &mpi id);
#pragma omp parallel
  int tag, swap_neigh, stat, omp_id = omp_thread_num();
  long buffer [BUFF_SIZE], incoming [BUFF_SIZE];
  big_ugly_calc1(omp_id, mpi_id, buffer);
                                              // Finds MPI id and tag so
  neighbor(omp_id, mpi_id, &swap_neigh, &tag); // messages don't conflict
  MPI_Send (buffer, BUFF_SIZE, MPI_LONG, swap_neigh,
           tag, MPI_COMM_WORLD);
  MPI Recv (incoming, buffer_count, MPI_LONG, swap_neigh,
           tag, MPI COMM WORLD, &stat);
  big ugly calc2(omp id, mpi id, incoming, buffer);
#pragma critical
  consume(buffer, omp id, mpi id);
```



Messages and threads

- Keep message passing and threaded sections of your program separate:
 - Setup message passing outside OpenMP parallel regions (MPI Thread funneled)
 - Surround with appropriate directives (e.g. critical section or master) (MPI Thread Serialized)
 - For certain applications depending on how it is designed it may not matter which thread handles a message. (MPI Thread Multiple)
 - Beware of race conditions though if two threads are probing on the same message and then racing to receive it.

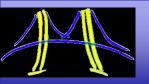
76 76



Safe Mixing of MPI and OpenMP Put MPI in sequential regions

```
MPI Comm Rank(MPI COMM WORLD, &mpi id);
MPI Init(&argc, &argv);
// a whole bunch of initializations
#pragma omp parallel for
for (I=0;I<N;I++) {
   U[I] = big calc(I);
   MPI Send (U, BUFF SIZE, MPI DOUBLE, swap neigh,
           tag, MPI COMM WORLD);
    MPI Recv (incoming, buffer count, MPI DOUBLE, swap neigh,
           tag, MPI COMM WORLD, &stat);
#pragma omp parallel for
for (I=0;I<N;I++) {
   U[I] = other big calc(I, incoming);
consume(U, mpi id);
```

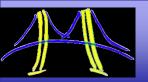
Technically Requires
MPI_Thread_funneled, but I
have never had a problem with
this approach ... even with preMPI-2.0 libraries.



Safe Mixing of MPI and OpenMP

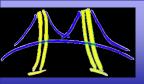
Protect MPI calls inside a parallel region

```
MPI Init(&argc, &argv); MPI Comm Rank(MPI COMM WORLD, &mpi id);
// a whole bunch of initializations
                                                Technically Requires
#pragma omp parallel
                                                MPI_Thread_funneled, but I
                                                have never had a problem with
#pragma omp for
  for (I=0;I<N;I++) U[I] = big_calc(I);
                                                this approach ... even with pre-
                                                MPI-2.0 libraries.
#pragma master
  MPI Send (U, BUFF SIZE, MPI DOUBLE, neigh, tag, MPI COMM WORLD);
   MPI Recv (incoming, count, MPI DOUBLE, neigh, tag, MPI COMM WORLD,
                                                                    &stat);
#pragma omp barrier
#pragma omp for
  for (I=0;I<N;I++) U[I] = other big calc(I, incoming);
#pragma omp master
  consume(U, mpi id);
```



Hybrid OpenMP/MPI works, but is it worth it?

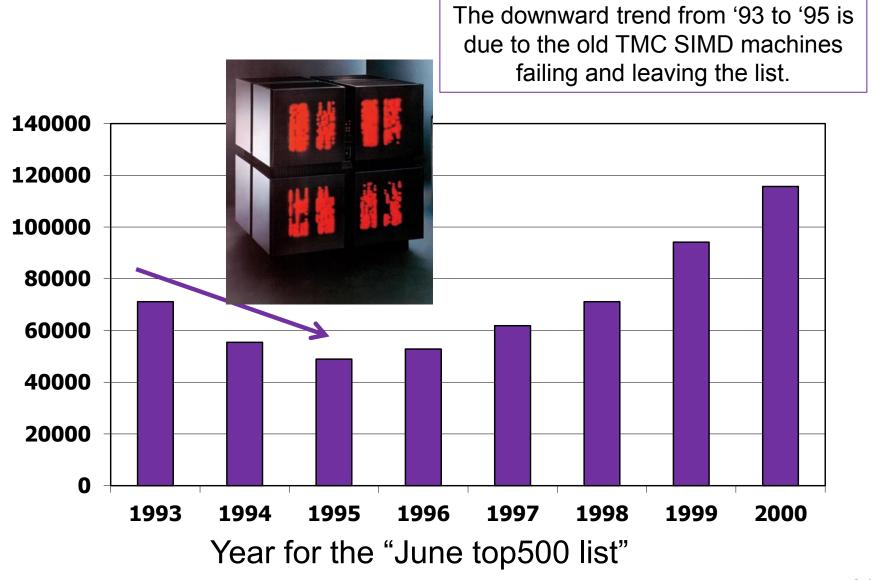
- Literature* is mixed on the hybrid model: sometimes its better, sometimes MPI alone is best.
- There is potential for benefit to the hybrid model
 - MPI algorithms often require replicated data making them less memory efficient.
 - Fewer total MPI communicating agents means fewer messages and less overhead from message conflicts.
 - Algorithms with good cache efficiency should benefit from shared caches of multi-threaded programs.
 - The model maps perfectly with clusters of SMP nodes.
- But really, it's a case by case basis and to large extent depends on the particular application.

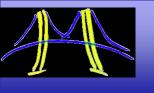


CLOSING COMMENTS

Parallel hardware trends

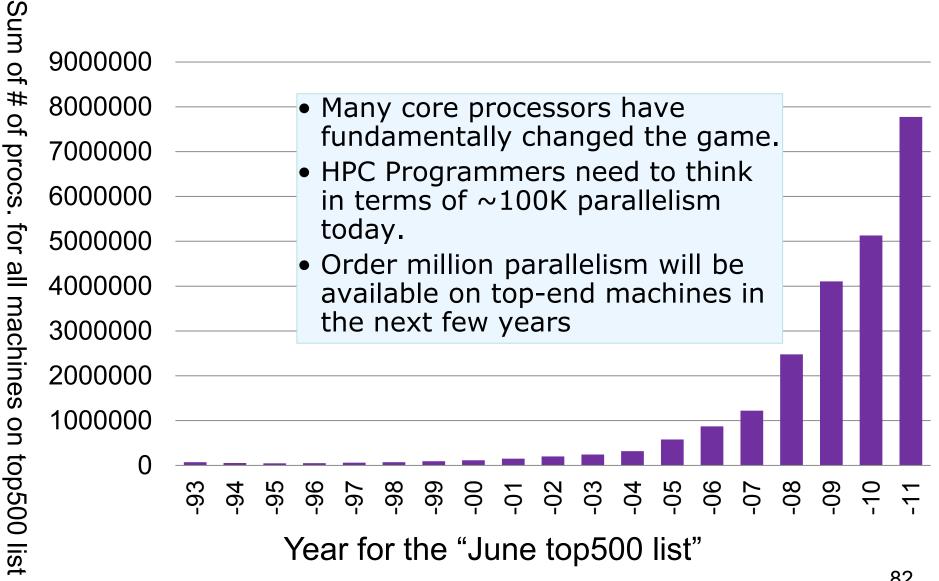
Top 500: total number of processors (1993-2000)

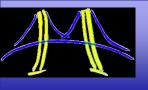




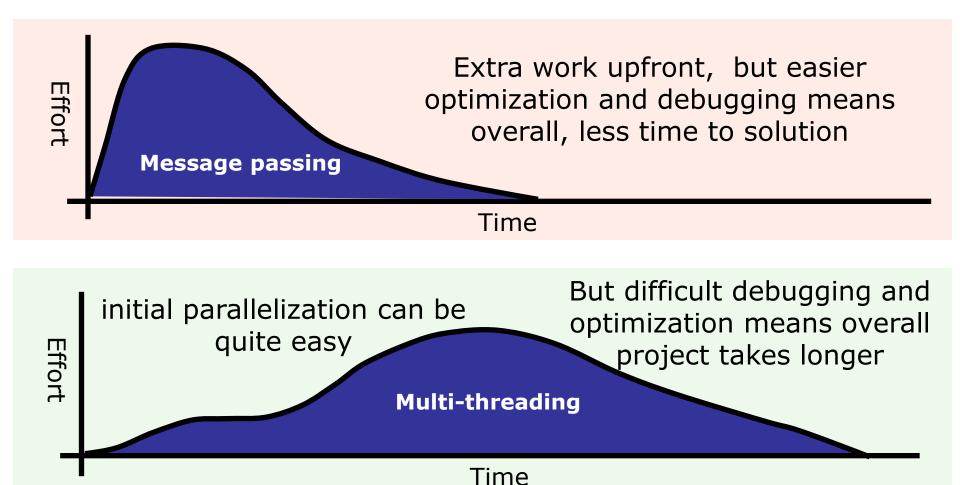
Parallel hardware trends

Top 500: total number of processors (1993-2011)

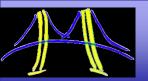




Does a shared address space make programming easier?

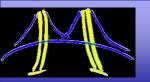


Proving that a shared address space program using semaphores is race free is an NP-complete problem*



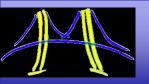
Closing comments

- Question conventional wisdom.
 - Do we really need cache coherence? If the memory hierarchy can't be hidden, isn't it better to expose the hierarchy so I can control it?
 - Debugging and Maintenance costs more than coding. So extra work up front to organize a problem to exploit the concurrency (e.g. decomposing and distributing data structures) shouldn't be such a big deal.
 - SW lives longer than HW. So why would anyone use a non-portable, non-standard programming model? That's just nuts!!
- As you move forward through the course
 - Notice that the patterns used in creating parallel code only weakly depend on the programming model. I can do loop parallelism with MPI, message passing with pthreads, kernel parallelism with OpenMP.
 - So learn multiple programming models and enjoy them ... but don't obsess about them. Ultimately, it's the design patterns and learning how to apply them to different problems that matter.



MPI References

- The Standard itself:
 - at http://www.mpi-forum.org
 - All MPI official releases, in both postscript and HTML
- Other information on Web:
 - at http://www.mcs.anl.gov/mpi
 - pointers to lots of stuff, including other talks and tutorials, a FAQ, other MPI pages

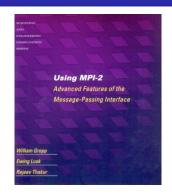


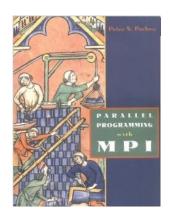
Books for learning MPI

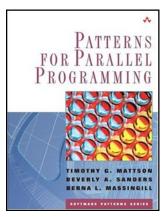
Using MPI-2: Portable Parallel Programming with the Message-Passing Interface, by Gropp, Lusk, and Thakur, MIT Press, 1999..

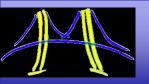
Parallel Programming with MPI, by Peter Pacheco, Morgan-Kaufmann, 1997.

Patterns for Parallel Programing, by Tim Mattson, Beverly Sanders, and Berna Massingill.





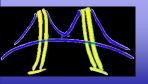




MPI core functions

```
#include <mpi.h>
int size, rank, argc, err_code=0;
char **argv;

MPI_Init (&argc, &argv);
MPI_Comm_rank (MPI_COMM_WORLD, &rank);
MPI_Comm_size (MPI_COMM_WORLD, &size);
MPI_Abort(MPI_COMM_WORLD, err_code);
MPI_Finalize();
```



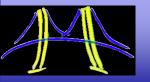
Basic MPI communication

Blocking send/receive

- MPI_Send (void* buf, int count, MPI_Datatype datatype,
 int dest, int tag, MPI_Comm comm)
- MPI_Recv (void* buf, int count, MPI_Datatype datatype, int source, int tag, MPI Comm comm, MPI Status* status)

Nonblocking send/receive

- MPI_Isend(const void *buf, int count, MPI_Datatype datatype,
 int dest, int tag, MPI_Comm comm, MPI_Request *request)
- MPI_Wait(MPI_Request *request, MPI_Status *status)



Collective Communication

- MPI_Barrier(MPI_Datatype datatype)