

A "Hands-on" Introduction to MPI

Tim Mattson

Intel Corp.

timothy.g.mattson@intel.com



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.

Our MPI progression

Topic	Exercise	concepts
I. MPI introduction	mpiexec hello world	Running programs on clusters
II. Group+context = communicator	MPI hello world	The organization of processes in an MPI program.
III. The Bulk Synchronous pattern.	MPI Pi Program	Collective communications
IV. Scaling	Running the MPI Pi program	Weak vs. Strong scaling
V. Message passing: basics	Ring program with send receive	How messages move through the MPI runtime
VI. Message passing diversity in MPI	Different versions of the ring program	Data environment details, software optimization
VII. Geometric decomposition	Matrix transpose	Using MPI for geometric decomposition problems
VIII. Wrap-up		The 12 fundamental MPI constructs

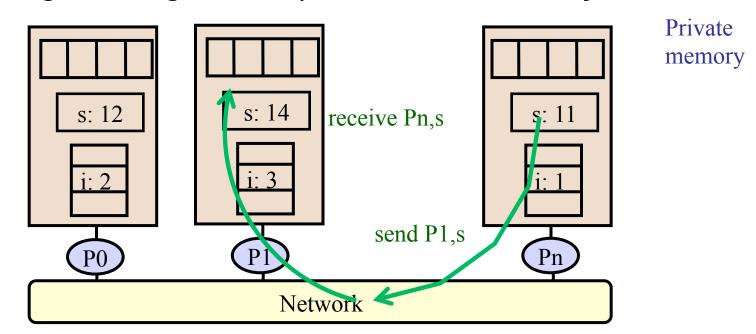
Outline



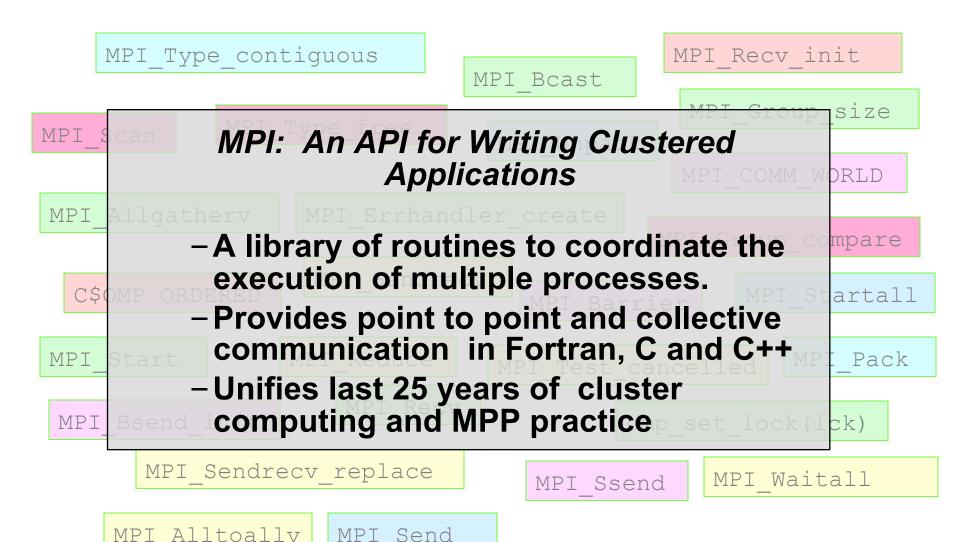
- The Distributed memory platform
 - MPI and the Bulk Synchronous Pattern
 - Scalability in Parallel Computing
 - Message Passing
 - Geometric Decomposition
 - Some closing thoughts

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 by explicit send/receive pairs
 - Coordination is implicit in every communication event.
 - MPI (Message Passing Interface) is the most commonly used SW

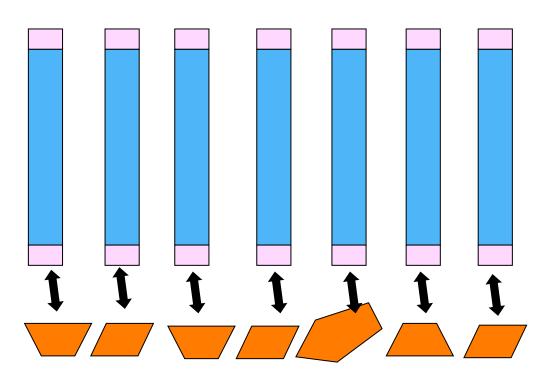


Parallel API's: MPI the Message Passing Interface



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.



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.

Running SPMD programs

- MPI implementations include a way to start "P processes" on the system.
- For MPIch (the most common MPI implementation, this is done with the mpiexec command:
 - > mpiexec -n P ./a.out

Run the program locally as P processes

- There are many options for mpiexec.
 - > mpiexec -f hostfile -n P ./a.out

> mpiexec -h

Run the program as P processes on the nodes from hostfile. A hostfile has node name on each line followed by a colon and the number of available poceeworsr

Ask mpieec for information about how to use the mpiexec commnds.

Exercise: Hello world part 1

- Goal
 - To confirm that you can run a program in parallel.
- Program
 - Add MPI to your path. In your ".bashrc file" add the line
 - PATH=\$PATH:/usr/lib64/mpich/bin
 - Write a program that prints "hello world" to the screen.
 - Use mpiexec to run multiple copies of the program.
 - Run them on your shared memory node
 - Run them across the nodes of a cluster (hint: you'll need a hostfile)
 - To run 3 processes on one node and 4 on another, my hostfile would be:

esc-33:3 esc-55:4

To run the executable hello on 2 processes on my local machine type:

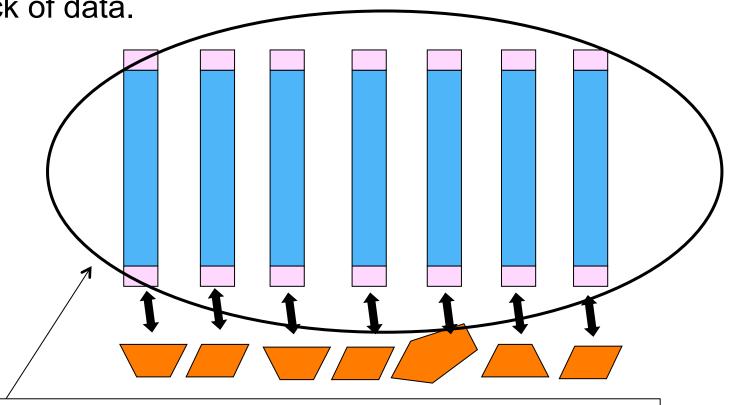
> mpiexec -n 4 ./a.out

To run the executable hello on 7 processes on my two node clusster:

> mpiexec -f hostfile -n 7 ./a.out

An MPI program at runtime

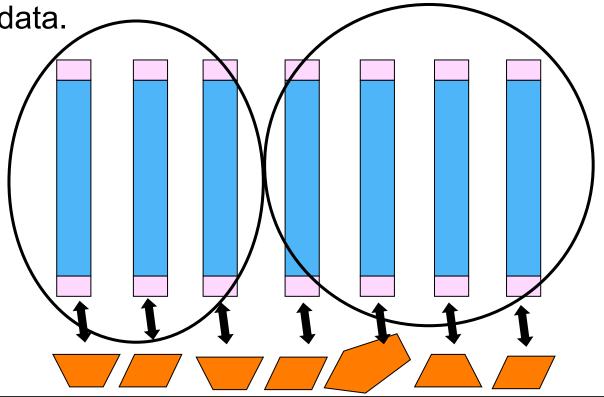
 Typically, when you run an MPI program, multiple processes all running the same program are launched ... working on their own block of data.



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 **process group** 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 argy are the command line args passed from main()

```
#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();
                    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

```
#inclu
#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"

```
#inclu
#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: Hello world

- Goal
 - To confirm that you can run an MPI 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 processes on my local node:

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

Running the program

```
#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;
```

```
    On a 4 node cluster, I'd run this program (hello) as:
    > mpirun -n 4 hello
    Hello from process 1 of 4
    Hello from process 2 of 4
    Hello from process 0 of 4
    Hello from process 3 of 4
```

Outline

- The Distributed memory platform
- - MPI and the Bulk Synchronous Pattern
 - Scalability in Parallel Computing
 - Message Passing
 - Geometric Decomposition
 - Some closing thoughts

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

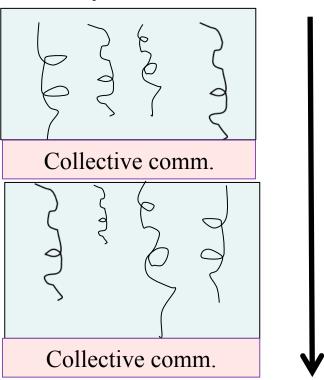
MPI Data Type	C Data Type
MPI_BYTE	
MPI_CHAR	signed char
MPI_DOUBLE	double
MPI_FLOAT	float
MPI_INT	int
MPI_LONG	long
MPI_LONG_DOUBLE	long double
MPI_PACKED	
MPI_SHORT	short
MPI_UNSIGNED_SHORT	unsigned short
MPI_UNSIGNED	unsigned int
MPI_UNSIGNED_LONG	unsigned long
MPI_UNSIGNED_CHAR	unsigned char

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

A typical pattern with MPI Programs

- Many MPI applications have few (if any) sends and receives. They use the following very common pattern:
- Time

- 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



 $P_1 P_2$

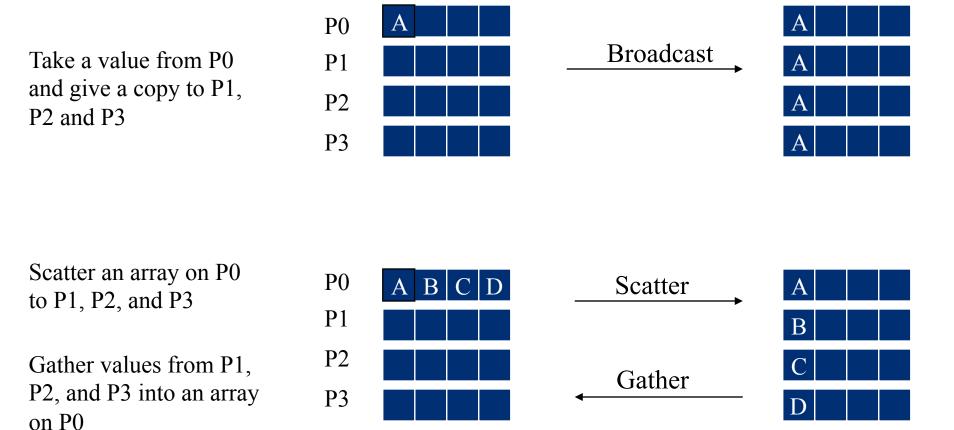
Processes

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

MPI Collective Routines

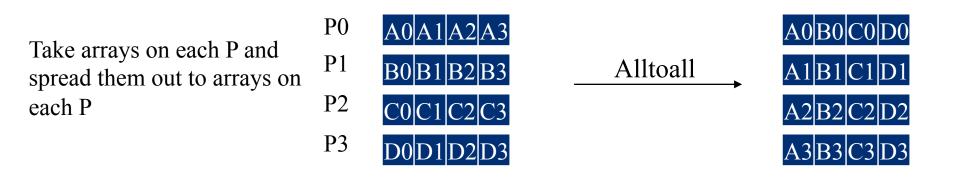
- Collective communications: called by all processes in the group to create a global result and share with all participating processes.
 - Allgather, Allgatherv, Allreduce, Alltoall,
 Alltoallv, Bcast, Gather, Gatherv, Reduce,
 Reduce scatter, Scan, Scatter, Scatterv
- Notes:
 - Allreduce, Reduce, Reduce_scatter, 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





Collective Computation

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

P0 В P1

P2

P3

Reduce

ABCD

Take values on each P and combine them with a scan operation and spread the scan array out among all P. P0 P1

P3

P2

Scan

AB ABC

Reduction

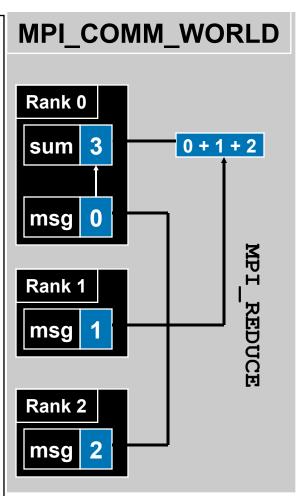
- MPI_Reduce performs specified reduction operation on specified data from all processes in communicator, places result in process "root" only.
- MPI_Allreduce places result in all processes (avoid unless necessary)

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

```
#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;
 msq = myrank;
 MPI Reduce (&msg, &sum, 1, MPI INT,
          MPI SUM, 0 MPI COMM WORLD);
 MPI Finalize();
```



Exercise 2: Pi Program

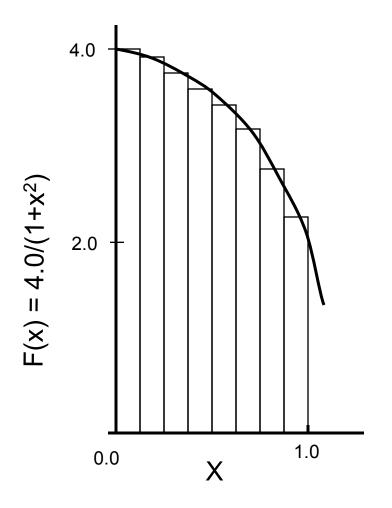
- Goal
 - To write a simple Bulk Synchronous, SPMD program
- Program
 - Start with the provided "pi program" and using an MPI reduction, write a parallel version of the program. Explore its scalability on your system.

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\leq 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

sum += 4.0/(1.0+x*x):

MPI COMM WORLD)

MPI Reduce(&sum, &pi, 1, MPI DOUBLE, MPI SUM, 0,

sum *= step;

Pi program in MPI #include <mpi.h> **MPI** Thread OpenMP OpenMP void main (int argc, char *argv∏) SPMD PI Loop or procs int i, my id, numprocs; double x, pi, step, sum critical step = 1.0/(double) num steps; MPI Init(&argc, &argy); 0.84 0.850.431 MPI Comm Rank(MPI COMM WORLD, MPI Comm Size(MPI COMM WORLD, & 0.480.230.483 0.23 0.470.46 for (i=my id; i<num steps; ; i=i+numprocs) 4 0.46 0.23 0.46 x = (i+0.5)*step;

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.

Outline

- The Distributed memory platform
- MPI and the Bulk Synchronous Pattern
- Scalability in Parallel Computing
 - Message Passing
 - Geometric Decomposition
 - Some closing thoughts

Amdahl's Law

- What is the maximum speedup you can expect from a parallel program?
- Approximate the runtime as a part that can be sped up with additional processors and a part that is fundamentally serial.

$$Time_{par}(P) = (serial_fraction + \frac{parallel_fraction}{P})*Time_{seq}$$

■ If serial_fraction is α and parallel_fraction is (1- α) then the speedup is:

$$S(P) = \frac{Time_{seq}}{Time_{par}(P)} = \frac{Time_{seq}}{(\alpha + \frac{1 - \alpha}{P}) * Time_{seq}} = \frac{1}{\alpha + \frac{1 - \alpha}{P}}$$

- If you had an unlimited number of processors: $P \rightarrow \infty$
- The maximum possible speedup is: $S = \frac{1}{\alpha}$ ← Amdahl's Law

What if the problem size grows

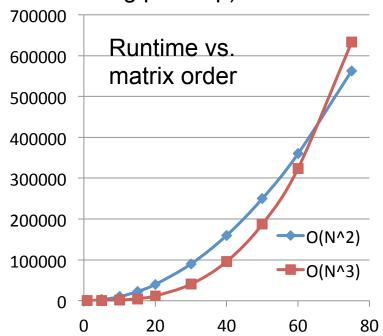
- Consider the dense linear algebra problems.
- A key feature of many of these operations between matrices (such as LU factorization or matrix multiplication)
 ... work scales as the cube of the order of the matrix.
- Assume we can parallelize the linear algebra operation
 (O(N³)) but not the loading of the matrices from memory
 (O(N²)). How does the serial fraction vary with matrix order
 (assume loading from memory is much slower than a
 floating point op).

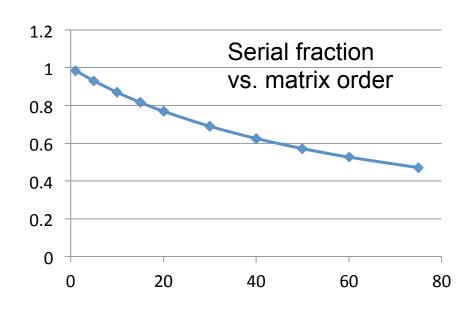
What would plots of runtime vs. problem size look like for the N squared and N cubed terms?

What would plots of serial fraction vs. problem size look like for the N squared and N cubed terms?

What if the problem size grows

- Consider the dense linear algebra problems.
- A key feature of many of these operations between matrices (such as LU factorization or matrix multiplication) ... work scales as the cube of the order of the matrix.
- Assume we can parallelize the linear algebra operation (O(N³)) but not the loading of the matrices from memory (O(N²)). How does the serial fraction vary with matrix order (assume loading from memory is much slower than a floating point op).

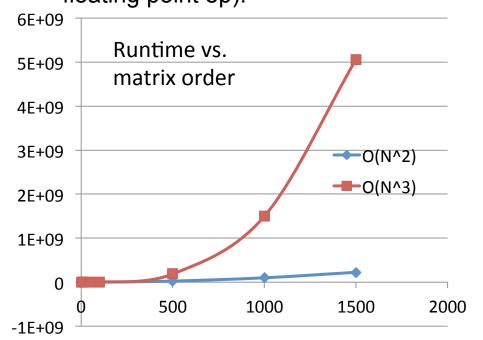


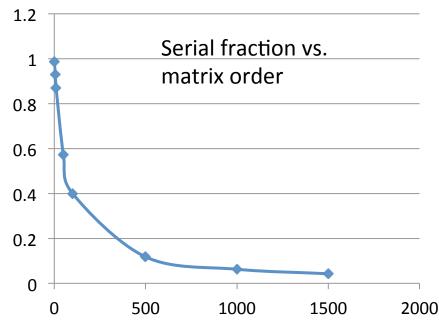


What if the problem size grows

- Consider the dense linear algebra problems.
- A key feature of many of these operations between matrices (such as LU factorization or matrix multiplication) ... work scales as the cube of the order of the matrix.

 Assume we can parallelize the linear algebra operation (O(N³)) but not the loading of the matrices from memory (O(N²)). How does the serial fraction vary with matrix order (assume loading from memory is much slower than a floating point op).





For much larger Matrix orders ...

Weak Scaling: a response to Amdhal

 Gary Montry and John Gustafson (1988, Sandia National Laboratories) observed that for many problems the serial fraction of a function of the problem size (N) decreases:

$$S(P,N) = \frac{T_{seq}(1)}{(\alpha(N) + \frac{1 - \alpha(N)}{P}) * T_{seq}(1)} \lim_{N \to N_{larg e}} \alpha(N) = 0$$

$$S(P,N_{larg e}) \to P$$

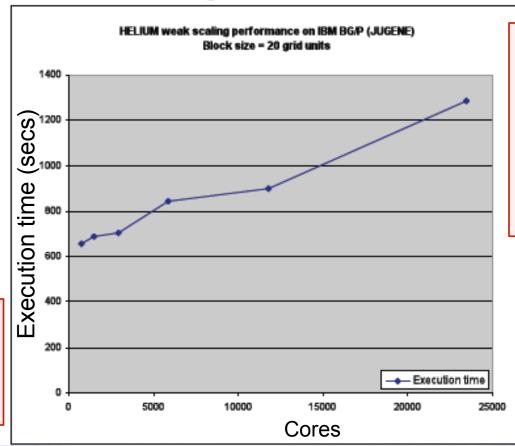
- In other words ... if parallelizable computations asymptotically dominate the runtime, then you can increase a problem size until limitations due to Amdahl's law can be ignored. This is an easier form of scalability for a programmer to meet ... so its called "weak scaling":
 - Weak Scaling: Performance of an application when the problem size increases with the number of processors (fixed size problem per node)

Example of weak scaling

HELIUM Weak Scaling Performance on BG/P

epcc

Local block size fixed to 20 grid units



A time dependent Quantum simulation of helium atoms with 20 grid units per processing element.

IBM Blue Gene P, 0.85 GHz, PowerPC 450, 4way processors

May 13, 10

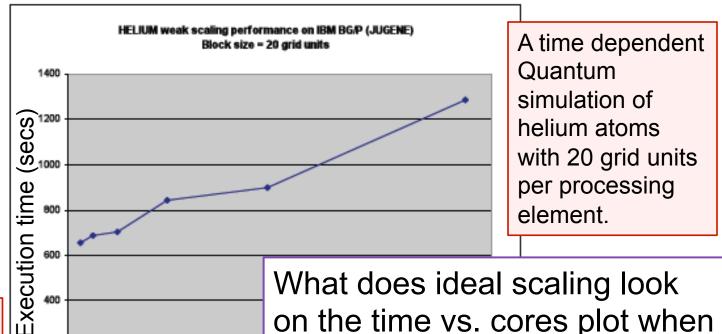
NAMD & HELIUM Enabling Work on the PRACE IBM Prototypes

Example of weak scaling

200

HELIUM Weak Scaling Performance on BG/P

Local block size fixed to 20 grid units



IBM Blue Gene P, 0.85 GHz, PowerPC 450, 4way processors

NAMD & HELIUM Enabling Work on the PRACE IBM Prototypes

10000

15000

Cores

34

on the time vs. cores plot when

25000

you have ideal weak scaling?

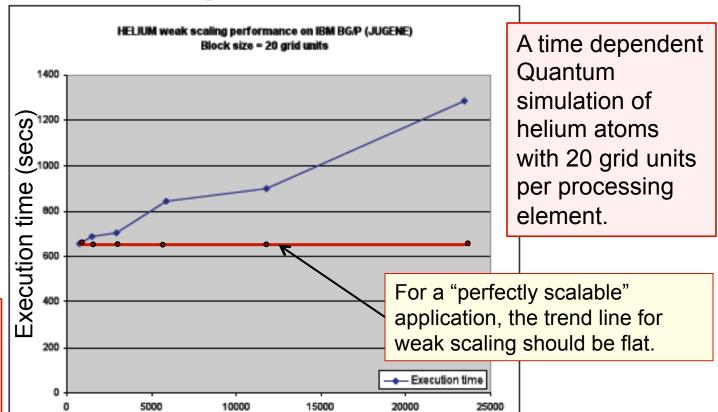
20000

Example of weak scaling

HELIUM Weak Scaling Performance on BG/P

epcc

Local block size fixed to 20 grid units



IBM Blue Gene P, 0.85 GHz, PowerPC 450, 4way processors

May 13, 10

NAMD & HELIUM Enabling Work on the PRACE IBM Prototypes

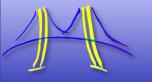
Cores

Outline

- The Distributed memory platform
- MPI and the Bulk Synchronous Pattern
- Scalability in Parallel Computing

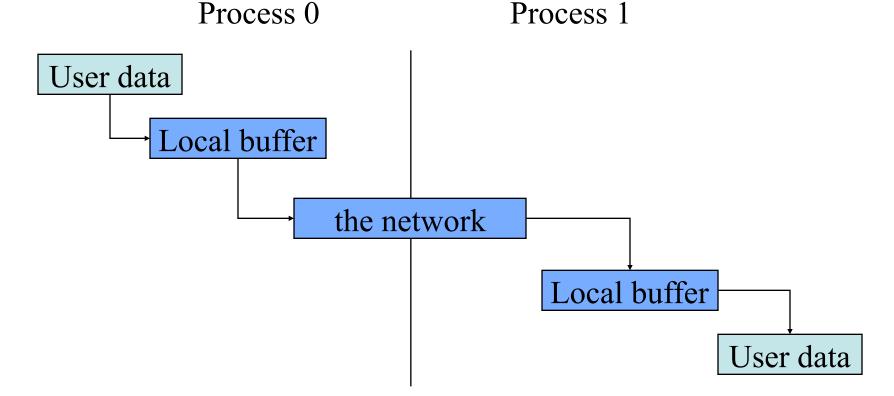


- Message Passing
- Geometric Decomposition
- Some closing thoughts



Buffers

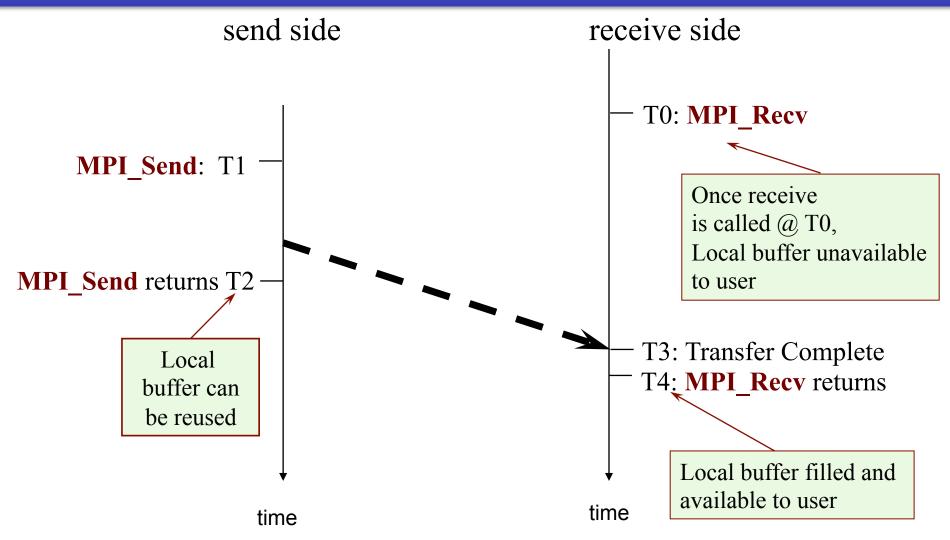
- Message passing has a small set of primitives, but there are subtleties
 - Buffering and deadlock
 - Deterministic execution
 - Performance
- When you send data, where does it go? One possibility is:



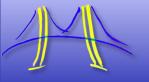
Derived from: Bill Gropp, UIUC

Blocking Send-Receive Timing Diagram

(Receive before Send)



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

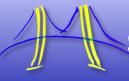


Sources of Deadlocks

- 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

Order the operations more carefully:

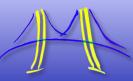
Process 0	Process 1
Send(1)	Recv(0)
Recv(1)	Send(0)

Supply receive buffer at same time as send:

Process 0 Process 1

Sendrecv(1) Sendrecv(0)

Slide source: Bill Gropp, UIUC



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)

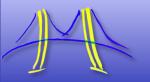
Isend(0)

Irecv(1)

Irecv(0)

Waitall

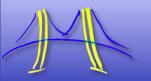
Waitall



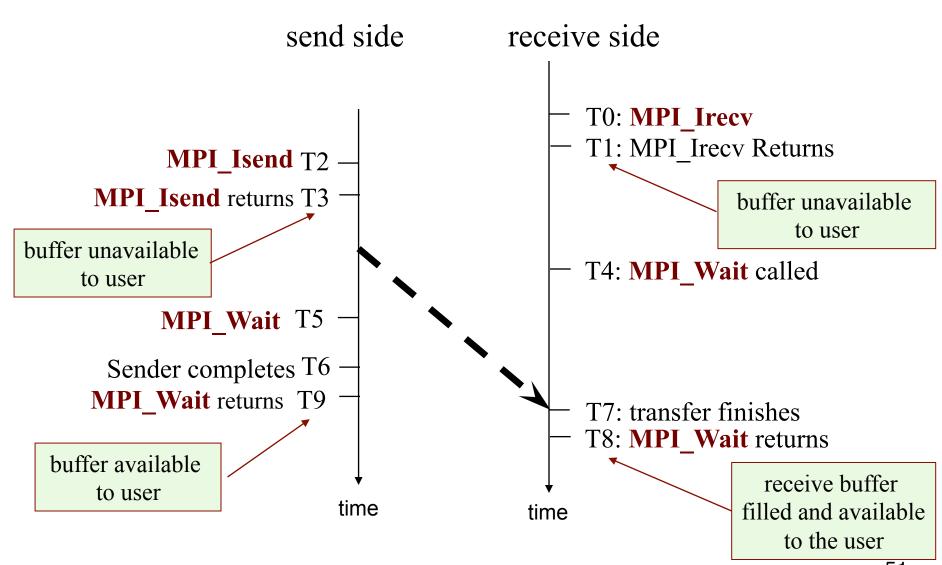
Non-Blocking Communication

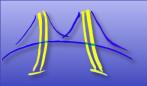
- Non-blocking operations return immediately and pass "request handles" that can be waited on and queried
 - MPI_Isend(start, count, datatype, dest, tag, comm, request)
 - MPI_Irecv(start, count, datatype, src, tag, comm, request)
 - MPI_Wait(request, status)
- One can also test without waiting using 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.



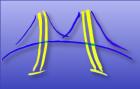
Non-Blocking Send-Receive Diagram





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(&req2, &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.
 - It may deadlock if the network stalls due to there being no place to put a message (i.e. no receives in place to the send that is blocking on when its buffer can be reused hangs).
 - Try to make it more stable for large messages by:
 - Split-phase ... have the nodes "send than receive" while the other half "receive then send".
 - Sendrecv ... a collective communication send/receive.

This is any integer ... it must be the same on matched sends and receives

double *buff; int buff count, to, from, tag=3; MPI Status stat;

MPI_Recv (buff, buff_count, MPI_DOUBLE, from, tag, MPI_COMM_WORLD, &stat); MPI_Send (buff, buff_count, MPI_DOUBLE, to, tag, MPI_COMM_WORLD);

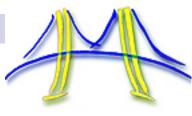
MPI_Sendrecv (snd_buff, buff_count, MPI_DOUBLE, to, tag,

rcv_buf, buff_count, MPI_DOUBLE, to, tag, MPI_COMM_WORLD, &stat);

Outline

- The Distributed memory platform
- MPI and the Bulk Synchronous Pattern
- Scalability in Parallel Computing
- Message Passing
- Geometric Decomposition
 - Some closing thoughts

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

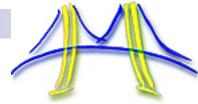
$$u_{j}^{n+1} = (1-2r)u_{j}^{n} + ru_{j-1}^{n} + ru_{j+1}^{n} \qquad r = \frac{k}{h^{2}}$$

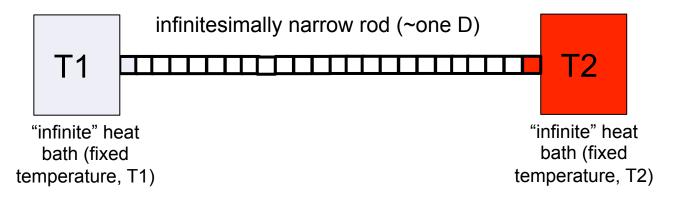
The solution at t = t_{n+1} is determined explicitly from the solution at t = t_n (assume u[t][0] = u[t][N] = Constant for all t).

```
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>
```

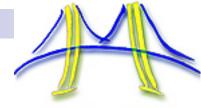
 Explicit methods are easy to compute ... each point updated based on nearest neighbors. Converges for r<1/2.

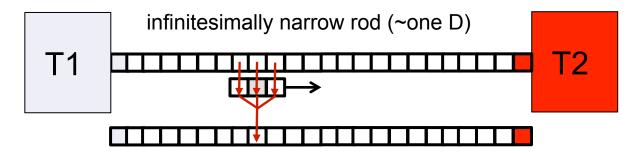




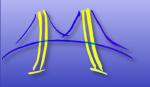






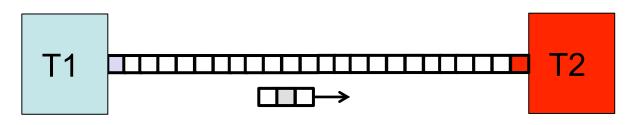


Pictorially, you are sliding a three point "stencil" across the domain (u[t]) and computing a new value of the center point (u[t+1]) at each stop.



return 0;

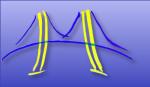
Heat Diffusion equation

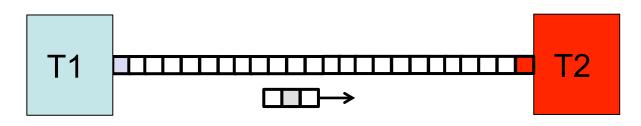


```
int main()
{
    double *u = malloc (sizeof(double) * (N));
    double *up1 = malloc (sizeof(double) * (N));

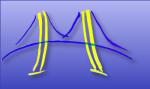
initialize_data(uk, ukp1, N, P); // init to zero, set end temperatures
for (int t = 0; t < N_STEPS; ++t){
    for (int x = 1; x < N-1; ++x)
        up1[x] = u[x] + (k / (h*h)) * (u[x+1] - 2*u[x] + u[x-1]);

    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>
```

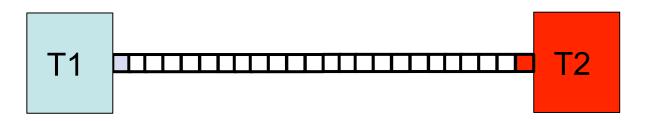


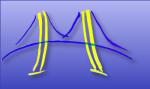


```
How would
int main()
                                                      you parallelize
  double *u = malloc (sizeof(double) * (N));
                                                      this program?
  double *up1 = malloc (sizeof(double) * (N));
   initialize data(uk, ukp1, N, P); // init to 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;
```



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



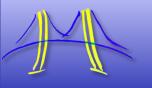


Seven strategies for parallelizing software

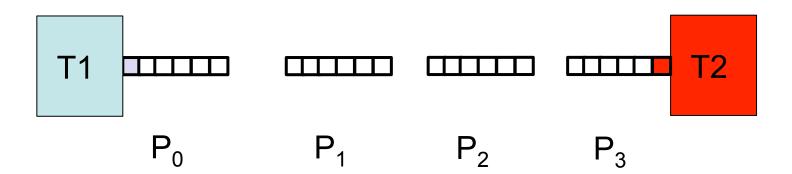
- These seven strategies for parallelizing software give us:
 - Names: so we can communicate better
 - Categories: so we can gather and share information
 - A palette (like an artist's palette) of approaches that is:
 - Necessary: we should consider them all and
 - Sufficient: once we have considered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

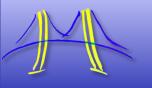
Parallel Algorithm Strategy Patterns

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

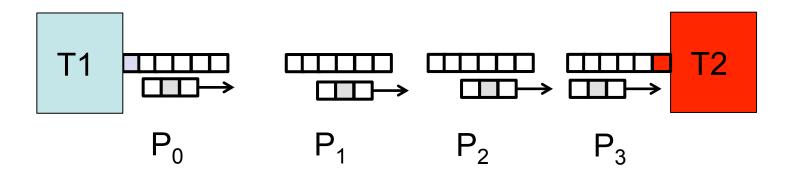


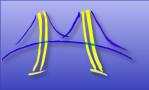
Break it into chunks assigning one chunk to each process.



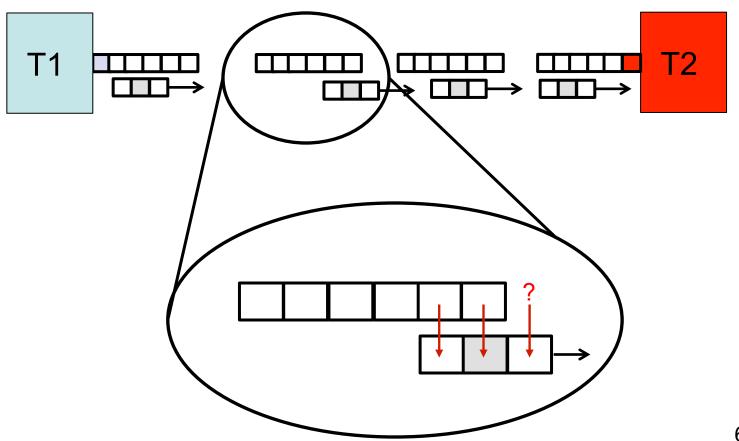


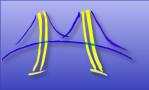
Each process works on it's own chunk ... sliding the stencil across the domain to updates its own data.



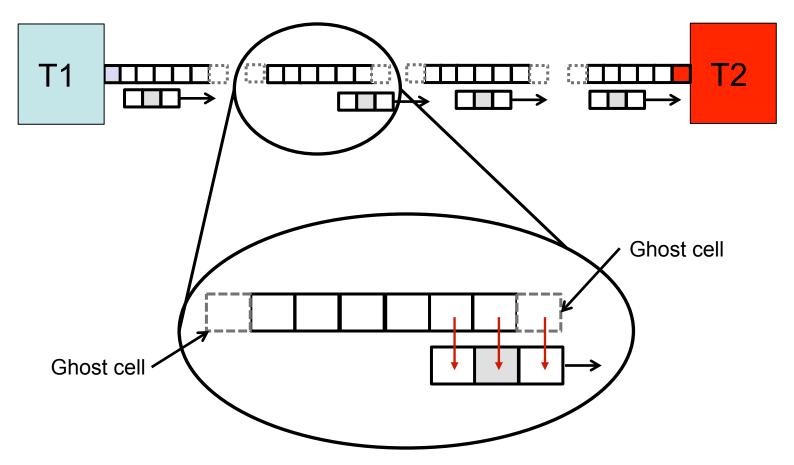


What about the ends of each chunk ... where the stencil will run off the end and hence have missing values for the computation?





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.



Geometric Decomposition

• Use when:

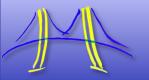
 The problem is organized around a central data structure that can be decomposed into smaller segments (chunks) that can be updated concurrently.

Solution

- Typically, the data structure is updated iteratively where a new value for one chunk depends on neighboring chunks.
- The computation breaks down into three components: (1) exchange boundary data, (2) update the interiors or each chunk, and (3) update boundary regions. The optimal size of the chunks is dictated by the properties of the memory hierarchy.

Note:

 This pattern is often used with the Structured Mesh and linear algebra computational strategy pattern.

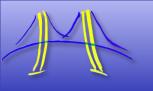


SPMD: Single Program Mulitple Data

- 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.

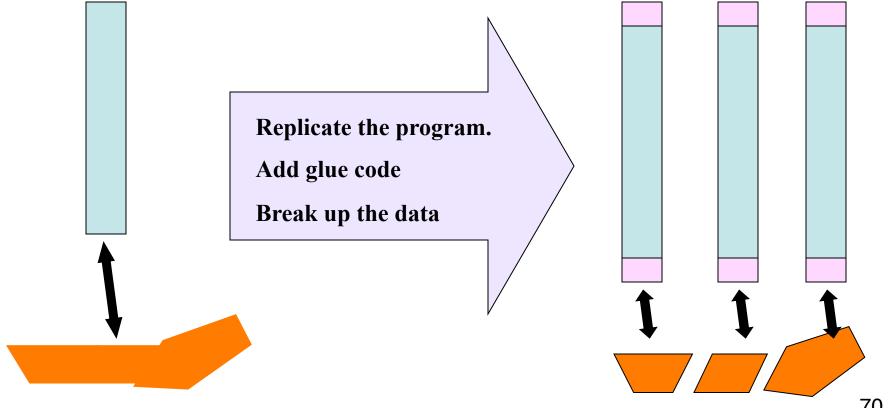
MPI programs almost always use this pattern ... it is probably the most commonly used pattern in the history of parallel programming.

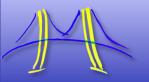


How do people use MPI? **The SPMD Design Pattern**

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.

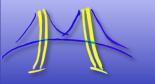




Heat Diffusion MPI Example

```
MPI Init (&argc, &argv);
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
MPI Finalize();
return 0;
```

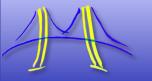
We write/explain this part first and then address the communication and data structures



return 0;

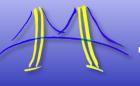
Heat Diffusion MPI Example

```
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();
```



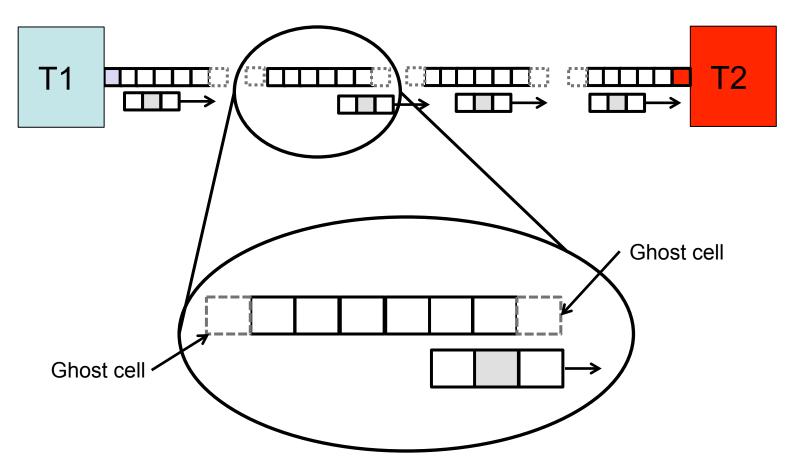
Heat Diffusion MPI Example

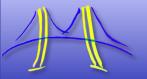
```
MPI_Init (&argc, &argv);
                                        1D PDE solver ... the simplest "real" message
                                        passing code I can think of. Note: edges of
MPI Comm size (MPI COMM WORLD, &P);
                                        domain held at a fixed temperature
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) 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 */
                                                                             73
```



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:

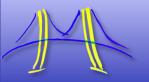
Arrays often need to be partitioned between multiple UEs. How can we do this so the resulting program is both readable and efficient?

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:

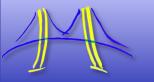
- Express algorithm in blocks
- Abstract indexing inside mapping functions ... programmer works in an index space natural to the domain, functions map into distribution needed for efficient execution.
- The text of the pattern defines some of these common mapping functions (which can get quite confusing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



Partitioned Arrays

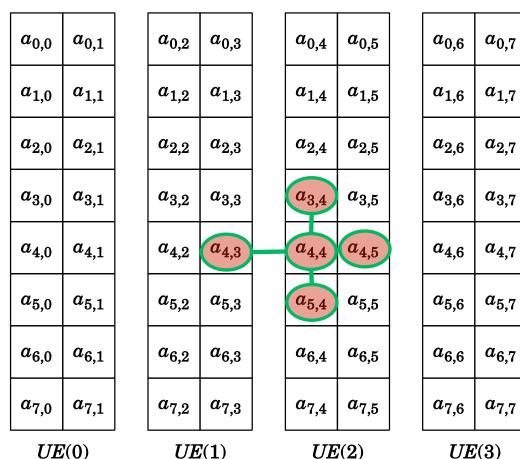
- Realistic problems are 2D or 3D; require move complex data distributions.
- We need to parallelize the computation by partitioning this index space
- Example: Consider a 2D domain over which we wish to solve a PDE using an explicit finite difference solver. The figure shows a five point stencil ...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

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



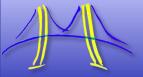
Partitioned Arrays: Column block distribution

- 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



P is the # of processors

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



Partitioned Arrays: Block distribution

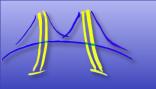
UE(1, 1)

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).

<i>UE</i> (0, 0)				<i>UE</i> (0, 1)				
$a_{0,0}$	$a_{0,1}$	$a_{0,2}$	$a_{0,3}$		$a_{0,4}$	$a_{0,5}$	$a_{0,6}$	$a_{0,7}$
$a_{1,0}$	$a_{1,1}$	$a_{1,2}$	$a_{1,3}$		$a_{1,4}$	$a_{1,5}$	$a_{1,6}$	$a_{1,7}$
$a_{2,0}$	$a_{2,1}$	$a_{2,2}$	$a_{2,3}$		$a_{2,4}$	$a_{2,5}$	$a_{2,6}$	$a_{2,7}$
$a_{3,0}$	$a_{3,1}$	$a_{3,2}$	$a_{3,3}$		$a_{3,4}$	$a_{3,5}$	$a_{3,6}$	$a_{3,7}$
$a_{4,0}$	$a_{4,1}$	$a_{4,2}$	$a_{4,3}$		$a_{4,4}$	$a_{4,5}$	$a_{4,6}$	$a_{4,7}$
$a_{5,0}$	$a_{5,1}$	$a_{5,2}$	$a_{5,3}$		$a_{5,4}$	$a_{5,5}$	$a_{5,6}$	$a_{5,7}$
$a_{6,0}$	$a_{6,1}$	$a_{6,2}$	$a_{6,3}$		$a_{6,4}$	$a_{6,5}$	$a_{6,6}$	$a_{6,7}$
$a_{7,0}$	$a_{7,1}$	$a_{7,2}$	$a_{7,3}$		$a_{7,4}$	$a_{7,5}$	$a_{7,6}$	$a_{7,7}$

UE(1, 0)

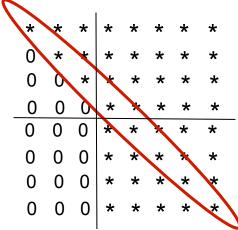
P is the # of processors



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).



 $A_{3.0}$

	_		
$a_{0,0} \mid a_{0,0}$	$a_{0,2} \mid a_{0,3}$	$oxed{a_{0,4} a_{0,5}}$	$a_{0,6} \mid a_{0,7}$
$\left[\begin{array}{c c}a_{1,0}\end{array}\middle a_{1,}$	$\begin{bmatrix} a_{1,2} & a_{1,3} \end{bmatrix}$	$oxed{a_{1,4} a_{1,5}}$	$oxed{a_{1,6} \mid a_{1,7}}$
$\overline{A_{0,0}}$	$A_{0,1}$	$\overline{A_{0,2}}$	$A_{0,3}$
$oxed{a_{2,0} \mid a_{2,0}}$	$a_{2,2} a_{2,3}$	$oxed{a_{2,4} \mid a_{2,5}}$	$\boxed{a_{2,6} \mid a_{2,7}}$
$a_{3,0} a_{3,0}$	$a_{3,2} a_{3,3}$	$a_{3,4} \mid a_{3,5} \mid$	$a_{3,6} a_{3,7}$
$\overline{A_{1,0}}$	$\overline{}$ $A_{1,1}$	$\overline{\hspace{1cm}}_{A_{1,2}}$	$\overline{}_{1,3}$
$a_{4,0} a_{4,0}$	$a_{4,2} \mid a_{4,3}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\boxed{a_{4,6} \mid a_{4,7}}$
$oxed{a_{5,0} a_{5,0}}$	$a_{5,2} \mid a_{5,3}$	$oxed{a_{5,4} a_{5,5}}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$
$oxed{A_{2,0}}$	$\overline{}$ $A_{2,1}$	$oxed{A_{2,2}}$	$A_{2,3}$
$a_{6,0} \mid a_{6,0}$	$a_{6,2} a_{6,3}$	$oxed{a_{6,4} a_{6,5}}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$
$oxed{a_{7,0} a_{7,0}}$	$a_{7,2} a_{7,3}$	$oxed{a_{7,4} a_{7,5}}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$

 $A_{3.1}$

 $A_{3,3}$

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

```
double *buff; int buff_count, to, from, tag=3; MPI_Status stat; MPI_Request s_req, r_req; MPI_Irecv (buff, buff_count, MPI_DOUBLE, from, tag, MPI_COMM_WORLD, &r_req); MPI_Isend (buff, buff_count, MPI_DOUBLE, to, tag, MPI_COMM_WORLD, &s_req); MPI_Wait(&recv_req, &stat) MPI_Wait(&send_req, &stat)
```

Outline

- The Distributed memory platform
- MPI and the Bulk Synchronous Pattern
- Scalability in Parallel Computing
- Message Passing
- Geometric Decomposition



The 12 core functions in MPI

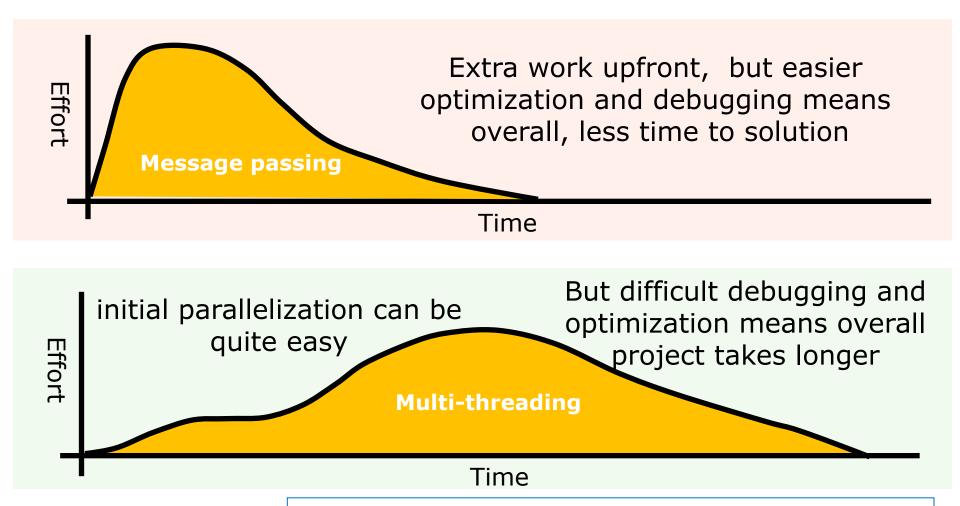
- MPI Init
- MPI Finish
- MPI_Comm_size
- MPI_Comm_rank
- MPI_Send
- MPI_Recv
- MPI_Reduce
- MPI_Isend
- MPI Irecv
- MPI_Wait
- MPI_Wtime
- MPI_Bcast

The 12 core functions in MPI

- MPI Init
- MPI Finish
- MPI_Comm_size
- MPI Comm rank
- MPI_Send
- MPI_Rocy
- MPI Reduce
- MPI Isend
- MPI Irecv
- MPI Wait
- MPI_Wtime
- MPI_Bcast

Real Programmers always try to overlap communication and computation .. Post your receives using MPI_Irecv() then where appropriate, MPI_Isend().

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, nonstandard 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

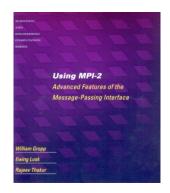
- The Standard itself:
 - at http://www.mpi-forum.org
 - All MPI official releases, in both postscript and HTML
- Other information on Web:
 - at http://www.mcs.anl.gov/mpi
 - pointers to lots of stuff, including other talks and tutorials, a FAQ, other MPI pages

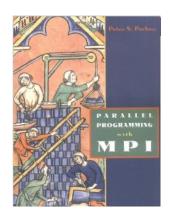
Books for learning MPI

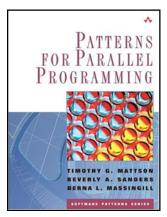
• Using MPI-2: Portable Parallel Programming with the Message-Passing Interface, by Gropp, Lusk, and Thakur, MIT Press, 1999..

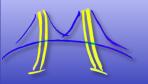
Parallel Programming with MPI, by Peter Pacheco, Morgan-Kaufmann, 1997.

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

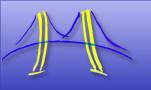








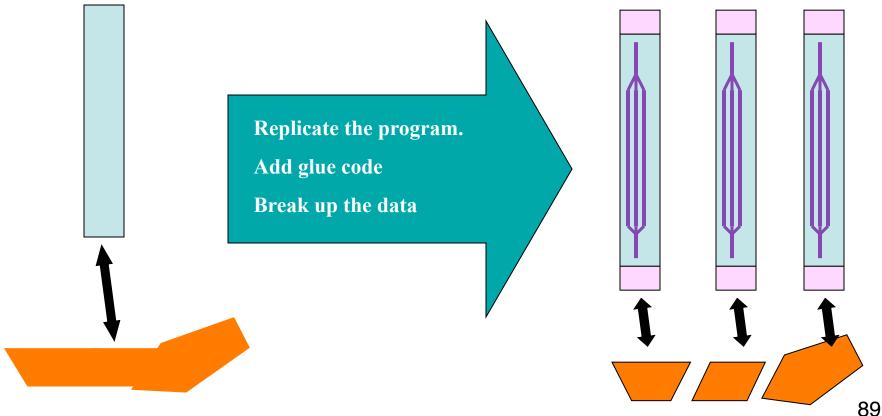
MIXING MPI AND OPENMP

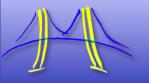


How do people mix MPI and OpenMP?

A sequential program working on a data set

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

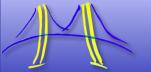




Pi program with MPI and OpenMP

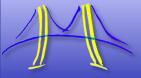
```
#include <mpi.h>
#include "omp.h"
void main (int argc, char *argv[])
        int i, my id, numprocs; double x, pi, step, sum = 0.0;
        step = 1.0/(double) num steps;
        MPI Init(&argc, &argv);
        MPI_Comm_Rank(MPI_COMM_WORLD, &my_id);
        MPI_Comm_Size(MPI_COMM_WORLD, &numprocs);
        my steps = num steps/numprocs;
#pragma omp parallel for reduction(+:sum) private(x)
        for (i=my id*my steps; i < (m id+1)*my steps; i++)
                 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);
```

Get the MPI part done first, then add OpenMP pragma where it makes sense to do so



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)
- Environment variables are not propagated by mpirun.
 You'll need to broadcast OpenMP parameters and set them with the library routines.

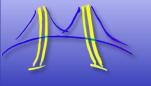


Dangerous Mixing of MPI and OpenMP

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

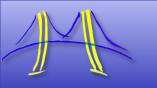
```
MPI_Comm_Rank(MPI COMM WORLD, &mpi id);
#pragma omp parallel
  int tag, swap_neigh, stat, omp_id = omp_thread_num();
  long buffer [BUFF SIZE], incoming [BUFF SIZE];
  big_ugly_calc1(omp_id, mpi_id, buffer);
                                               // Finds MPI id and tag
SO
  neighbor(omp 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

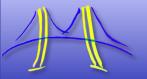
- Keep message passing and threaded sections of your program separate:
 - Setup message passing outside OpenMP parallel regions (MPI_Thread_funneled)
 - Surround with appropriate directives (e.g. critical section or master) (MPI Thread Serialized)
 - For certain applications depending on how it is designed 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.



Safe Mixing of MPI and OpenMP Put MPI in sequential regions

```
MPI Init(&argc, &argv);
                       MPI Comm Rank(MPI COMM WORLD, &mpi id);
// 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.

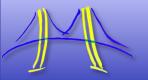


Safe Mixing of MPI and OpenMP

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
#pragma omp for
                                                have never had a problem with
  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);
```

95 **95**



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.