



 Learning Objectives
 Understand the problem of race conditions in concurrent systems,
 Learn how to decompose a program for parallel execution,
 Be able to write simple parallel programs in the important programming models,
 Understand the operation of common cachecoherence algorithms, both bus-based and network-based, and
 Learn about common memory-consistency models, and appreciate the advantages and disadvantages of each.

CSC/ECE 506: Architecture of Parallel Computer

2

NC STATE UNIVERSITY



	4 programs: 24%
Homework 50%	3 problem sets: 18%
	1 peer-reviewed exercise: 8%
	Test 1: 10%
Tests 50%	Test 2: 15%
	Final exam: 25%

































Current trends: multicore and manycore IBM Cell Aspect Intel Cloverto AMD Barcelona # cores 4 4 8+1 Clock 2.66 GHz 2.3 GHz 3.2 GHz frequency Core type 000 000 2-issue SIMD Superscalar Superscalar 2x4MB L2 512KB L2 256KB local Caches (private), 2MB L3 (sh'd) store Chip power 120 watts 95 watts 100 watts Exercise: Browse the Web (or the textbook O) for information on more recent processors, and for each processor, fill out this form. (You can view the submissions.) NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computer









































38





Exercise

- Go to http://www.top500.org and look at the Lists and Statistics menus in the top menu bar.
 - choose either List Statistics or Development over time, - then select one of the statistics, e.g., Vendors, Processor
 - examine what kind of systems are prevalent. Then do the same for earlier lists, and report on the trend.
- · You can go all the way back to the first list from 1993.

Three parallel-programming models

- *Shared-memory* programming is like using a "bulletin board" where you can communicate with colleagues.
- *Message-passing* is like communicating via e-mail or telephone calls. There is a well defined event when a message is sent or received.
- Data-parallel programming is a "regimented" form of cooperation. Many processors perform an action separately on different sets of data, then exchange information globally before continuing en masse.

User-level communication primitives are provided to realize the programming model

There is a mapping between language primitives of the programming model and these primitives

These primitives are supported directly by hardware, or via OS, or via user software.

In the early days, the kind of programming model that could be used was closely tied to the architecture.

Today—

- · Compilers and software play important roles as bridges
- Technology trends exert a strong influence

The result is convergence in organizational structure, and relatively simple, general-purpose communication primitives.

A shared address space

In the shared-memory model, processes can access the same memory locations.

Communication occurs implicitly as result of loads and stores

This is convenient.

• Wide range of granularities supported.

Lecture 2

Architecture of Parallel Computers

The interconnect in a shared-memory multiprocessor can take several forms.

It may be a *crossbar switch*.

Each processor has a direct connection to each memory and I/O controller.



1

Bandwidth scales with the number of processors.

Unfortunately, cost scales with the square of the number of processors.

This is sometimes called the "mainframe approach."

At the other end of the spectrum is a *shared-bus* architecture.



All processors, memories, and I/O controllers are connected to the bus.

Such a multiprocessor is called a symmetric multiprocessor (SMP).

What are some advantages and disadvantages of organizing a multiprocessor this way? List them <u>here</u>.

•

A compromise between these two organizations is a *multistage interconnection network*.

- Similar programming model to time-sharing on uniprocessors, except that processes run on different processors
- · Wide range of scale: few to hundreds of processors

Good throughput on multiprogrammed workloads.

This is popularly known as the *shared memory* model, even though memory may be physically distributed among processors.

The shared-memory model

A process is a virtual address space plus one or more threads of control.

Portions of the address spaces of tasks are shared.



What does the private region of the virtual address space usually contain? Stack and private data, incl. register save areas and control flags.

Conventional memory operations can be used for communication.

Special atomic operations are used for synchronization.

The interconnection structure

Lecture 2

Architecture of Parallel Computers

The processors are on one side, and the memories and controllers are on the other.

Each memory reference has to traverse the stages of the network.

Why is this called a compromise between the other two strategies?



2

Like a bus, it doesn't require a connection from each processor to each memory. Like a crossbar, it can handle multiple accesses simultaneously. But, it has more connections than a bus, and can handle fewer simultaneous accesses than a crossbar.

For small configurations, however, a shared bus is quite viable.

Message passing

In a message-passing architecture, a complete computer, including the I/O, is used as a building block.

Communication is via explicit I/O operations, instead of loads and stores.

- A program can directly access only its private address space (in local memory).
- It communicates via explicit messages (send and receive).
- It is like a network of workstations (clusters), but more tightly integrated.

Easier to build than a scalable shared-memory machine.

Send-receive primitives

The programming model is further removed from basic hardware operations.



Library or OS intervention is required to do communication.

- · send specifies a buffer to be transmitted, and the receiving process.
- · receive specifies sending process, and a storage area to receive into. · A memory-to-memory copy is performed, from the address space of one process to the address space of the other.
- · There are several possible variants, including whether send completes-

when the receive has been executed, synchronous when the send buffer is available for reuse, or

- when the message has been sent. asynchronous
- · Similarly, a receive can wait for a matching send to execute, or simply fail if one has not occurred.

There are many overheads: copying, buffer management, protection. Let's describe each of these. Submit your descriptions here.

· Why is there an overhead to copying, compared to a share-memory machine?

Lecture 2	Arch

hitecture of Parallel Computers

```
5
```

Lecture 2

3

Architecture of Parallel Computers

// spawn child thread

// tell child to work

· Describe the overhead of buffer management.

Here's an example from the textbook of the difference between shared address-space and message

while (signal == 0) {}; // wait until instructed to work

printf("child thread> sum is d'', a + b);

while (signal == 1) {} // wait until child done

signal = 0; // my work is done

printf("all done, exiting\n");

· What is the overhead for protection?

A shared-memory system uses the thread model:

passing programming

int a, b, signal;

void main() {

a = 5, b = 3;

clone(&dosum,...)

signal = 0;

signal = 1;

void dosum(<args>) {

Message-passing uses the process model:

```
int a, b;
```

void dosum() {

```
recvMsg(mainID, &a, &b);
```

```
printf("child process> sum is %d", a + b);
```

```
}
```

void main() {

```
if (fork() == 0) // I am the child process
 dosum();
                  // I am the parent process
else {
 a = 5, b = 3;
 sendMsg(childID, a, b);
 wait(childID);
 printf("all done, exiting\n");
}
```

```
}
```

Here's the relevant section of documentation on the fork () function:

"Upon successful completion, fork () and fork1 () return 0 to the child process and return the process ID of the child process to the parent process."

Interconnection topologies

Early message-passing designs provided hardware primitives that were very close to the message-passing model.

Each node was connected to a fixed set of neighbors in a regular pattern by point-to-point links that behaved as FIFOs.

A common design was a hypercube, which had $2 \times n$ links per node, where n was the number of dimensions.

The diagram shows a 3D cube.

One problem with hypercubes was that they were difficult to lay out on silicon.

Because of that, 2D meshes eventually supplanted hypercubes.



Here is an example of a 16-node mesh. Note that the last element in one row is connected to the first element in the next.

If the last element in each row were connected to the first element in the same row, we would have a torus instead.

FIFO on each link.

Lecture 2



Thus, software sends were implemented as synchronous hardware operations at each node.

What was the problem with this, for multi-hop messages? Interrupts are required at intermediate nodes.

- Synchronous ops were replaced by DMA, enabling non-blocking operations
- A DMA device is a special-purpose controller that transfers data between memory and an I/O device without processor intervention.
- Messages were buffered by the message layer of the system at the destination until a receive took place.
- When a receive took place, the data was copied into the address space of the receiving process.

The diminishing role of topology.

- With store-and-forward routing, topology was important.
 Parallel algorithms were often changed to conform to the topology of the machine on which they would be run.
- Introduction of pipelined ("wormhole") routing made topology less important.

In current machines, it makes less difference how far the data travels.

This simplifies programming; cost of interprocessor communication is essentially independent of which processor is receiving the data.

Toward architectural convergence

In 1990, there was a clear distinction between message-passing and shared-memory machines. Today, there isn't a distinct boundary.

- Message-passing operations are supported on most shared-memory machines.
- A shared virtual address space can be constructed on a messagepassing machine, by sharing *pages* between processors.
 - When a missing page is accessed, a page fault occurs.

Lecture 2	Architecture of Parallel Computers

- Synchronization overhead ... one thread needs to wait until another thread arrives at a certain point.
- Communication overhead ... message-passing is not instantaneous.
- Load imbalance ... some threads may have work to do than others, and the fast ones need to wait for the slow ones.
- Serial programs don't have to use thread-safe libraries, and that helps them run faster.

[§4.3.1] If some portions of the problem don't have much concurrency, the speedup on those portions will be low, lowering the average speedup of the whole program.

Exercise: Submit your answers to the questions below.

Suppose that a program is composed of a serial phase and a parallel phase.

- The whole program runs for 1 time unit.
- The serial phase runs for time s, and the parallel phase for time 1-s.

Then regardless of how many processors $\it N$ are used, the execution time of the program will be at least $\it s$

and the speedup will be no more than 1/s. This is known as Amdahl's law.

For example, if 25% of the program's execution time is serial, then regardless of how many processors are used, we can achieve a speedup of no more than 4.

Efficiency is defined as

speedup

number of processors

Let us normalize computation time so that

 The OS fetches the page from the remote node via messagepassing.

At the machine-organization level, the designs have converged too.

The block diagrams for shared-memory and message-passing machines look essentially like this:

\square	Network		
M S	M S	000	M S
P	P		P

In shared memory, the network interface is integrated with the memory controller.

It initiates a transaction to access memory at a remote node.

In message-passing, the network interface is essentially an I/O device.

What does Solihin say about the ease of writing shared-memory and message-passing programs on these architectures?

- · Which model is easier to program for initially?
- · Why doesn't it make much difference in the long run?

The limits of parallelism: Amdahl's law

Speedup is defined as

time for serial execution

time for parallel execution

or, more precisely, as

time for serial execution of best serial algorithm time for parallel execution of our algorithm

Give two reasons why it is better to define it the second way than the first.

- Bad algorithms often have good speedup.
- Lecture 2

2

Architecture of Parallel Computers

10

• the serial phase takes time 1, and

• the parallel phase takes time *p* if run on a single processor.

Then if run on a machine with N processors, the parallel phase takes p/N.

Let $\boldsymbol{\alpha}$ be the ratio of serial time to total execution time. Thus

$$\alpha = \frac{1}{1+p/N} = \frac{N}{N+p}.$$

For large N, α approaches 1, so efficiency approaches 0.

Does it help to add processors? No ...

Gustafson's law: But this is a pessimistic way of looking at the situation.

In 1988, Gustafson et al. noted that as computers become more powerful, people run larger and larger programs.

Therefore, as *N* increases, *p* tends to increase too. Thus, the fraction of time $1-\alpha$ does not necessarily shrink with increasing *N*, and efficiency remains reasonable.

There may be a maximum to the amount of speedup for a given problem size, but since the problem is "scaled" to match the processing power of the computer, there is no clear maximum to "scaled speedup."

Gustafson's law states that any sufficiently large problem can be efficiently parallelized.

Shared-Memory Parallel Programming



seen this done in the last lecture?

What considerations are important in mapping threads to processors?

Solihin says that there are three levels of parallelism:

- program level
- algorithm level
- code level

Identifying loop-level parallelism

[\$3.2] Goal: given a code, without knowledge of the algorithm, find parallel tasks.

Focus on loop-dependence analysis.

Notations:

• S is a statement in the source code

Lecture 2	Architecture of Parallel Computers	13	Lecture 2	Architecture of Parallel Computers	14
				· · · · · · · · · · · · · · · · · · ·	

<pre>for (i=1; i<n; i++)="" pre="" {<=""></n;></pre>	$S1[i] \rightarrow TS1[i+1]: loop-carried$
S1: a[i] = a[i-1] + 1;	$S1[i] \rightarrow T S2[i]$: loop-independent
S2: b[i] = a[i];	S3[i,j] →TS3[i,j+1]:
}	 loop-carried on for j loop no loop-carried dependence in for i loop
<pre>for (i=1; i<n; i++)<="" pre=""></n;></pre>	
<pre>for (j=1; j< n; j++)</pre>	S4[1,]]→ S4[1+1,]]:
S3: a[i][j] = a[i][j-1] + 1;	• no loop-carried dependence in for j loop
<pre>for (i=1; i<n; i++)<="" pre=""></n;></pre>	• loop-carried on for i loop
for (j=1; j< n; j++)	Iteration-space Traversal Graph (ITG)
S4: a[i][j] = a[i-1][j] + 1;	[§3.2.1] The ITG shows graphically the order of traversal in the iteration space. This is

sometimes called the happens-before relationship. In an ITG,

- A node represents a point in the iteration space
- A directed edge indicates the next point that will be encountered after the current point is traversed

Example:

for (i=1; i<4; i++)		
<pre>for (j=1; j<4; j++)</pre>		
S3: a[i][j] = a[i][j-1]	+	1;

- S[i, j, ...] denotes a statement in the loop iteration [i, j, ...]
- "S1 then S2" means that S1 happens before S2
- If S1 then S2:

S1 \rightarrow T S2 denotes true dependence, i.e., S1 writes to a location that is read by S2 (RAW hazard)

S1 \rightarrow A S2 denotes anti-dependence, i.e., S1 reads a location written by S2 (WAR hazard)

S1 ${\rightarrow}{\rm O}$ S2 denotes output dependence, i.e., S1 writes to the same location written by S2 (WAW hazard)

Example:



Exercise: Identify the dependences in the above code.

Loop-independent vs. loop-carried dependences

[§3.2] Loop-carried dependence: dependence exists across iterations; i.e., if the loop is removed, the dependence *no longer exists*.

Loop-independent dependence: dependence exists within an iteration; i.e., if the loop is removed, the dependence still exists.

Example:



Loop-carried Dependence Graph (LDG)

- · LDG shows the true/anti/output dependence relationship graphically.
- A node is a point in the iteration space.
- A directed edge represents the dependence.

Example:

for (i=1; i<4; i++)
for (j=1; j<4; j++)
S3: a[i][j] = a[i][j-1] + 1;</pre>

What do we know about the ITG for these nested loops?



Another example:

for (i=1; i<=n; i++)
for (j=1; j<=n; j++)
S1: a[i][j] = a[i][j-1] + a[i][j+1] + a[i-1][j] + a[i+1][j];
for (i=1; i<=n; i++)
for (j=1; j<=n; j++) {
 S2: a[i][j] = b[i][j] + c[i][j];
 S3: b[i][j] = a[i][j-1] * d[i][j];
}</pre>

• Draw the ITG

• List all the dependence relationships

Note that there are two "loop nests" in the code.

- The first involves S1.
- The other involves S2 and S3.

Lecture 2

Architecture of Parallel Computers

17



Draw the LDG for Loop Nest 1.



Dependence relationships for Loop Nest 2

- True dependences:
- S2[i,j] →T S3[i,j+1]
- Output dependences:
 - None
- Anti-dependences:
 - \circ s2[i,j] →A s3[i,j] (loop-independent dependence)



Dependence relationships for Loop Nest 1

• True dependences:

$$\circ$$
 S1[i,j] →T S1[i,j+1]
 \circ S1[i,j] →T S1[i+1,j]

- Anti-dependences:

Exercise: Suppose we dropped off the first half of S1, so we had

S1: a[i][j] = a[i-1][j] + a[i+1][j];

or the last half, so we had

S1: a[i][j] = a[i][j-1] + a[i][j+1];

Which of the dependences would still exist?

Lecture	2
Lecture	2

Architecture of Parallel Computers

18





Why are there no vertical edges in this graph? Answer here.

Why is the anti-dependence not shown on the graph?

Exercise: Consider this code sequence.

for (i = 3; i < n; i++) {
 for (j = 0; j < n - 3; j++) {
 S1: A[i][j] = A[i - 3][j] + A[i][j + 3];
 S2: B[i][j] = A[i][j] / 2;
 }
}</pre>

List the dependences, and say whether they are loop independent or loop carried. Then draw the ITG and LDG (you don't need to submit these).

Finding parallel tasks across iterations

[§3.2.2] Analyze loop-carried dependences:

- Dependences must be enforced (especially true dependences; other • dependences can be removed by privatization)
- There are opportunities for parallelism when some dependences are not present.

Example 1







How many parallel tasks are there here? *n*, one per iteration of the *i* loop.

Example 3



Lecture 2

Architecture of Parallel Computers

21

In each anti-diagonal, the nodes are independent of each other



We

need to rewrite the code to iterate over anti-diagonals:

Calculate number of anti-diagonals for each anti-diagonal do Calculate the number of points in the current anti-diagonal for_all points in the current anti-diagonal do Compute the value of the current point in the matrix

Parallelize the loops highlighted above.

<pre>for (i=1; i <= 2*n-1; i+</pre>	+) {// 2n-1 anti-diagonals
if (i <= n) {	
points = i;	// number of points in anti-diag
row = i;	// first pt (row,col) in anti-diag $% \left(\left({{{\left({{{\left({r_{i}} \right)}} \right)}_{i}}} \right)$
col = 1;	// note that row+col = i+1 always $% \left($
}	
else {	
points = 2*n - i;	
row = n;	
col = i-n+1;	// note that row+col = i+1 always $% \left($

DOACROSS Parallelism

[§3.2.3] Suppose we have this code:

```
Can we execute anything
                                                                      in
                             for (i=1; i<=N; i++) {</pre>
parallel?
                               S: a[i] = a[i-1] + b[i] * c[i];
Well, we can't run the
```

iterations of the for loop in parallel, because ...

 $s[i] \rightarrow T s[i+1]$ (There is a loop-carried dependence.)

But, notice that the b[i] *c[i] part has no loop-carried dependence.

This suggests breaking up the loop into two:

<pre>for (i=1; i<=N; i++) { S1: temp[i] = b[i] * c[i]; }</pre>	The first loop is izable. The second is not. Execution time: $N \times (T_{S1}+T_{S2})$		
<pre>for (i=1; i<=N; i++) { S2: a[i] = a[i-1] + temp[i];</pre>	What is a approact	a disadvan n? <mark>Storag</mark> o	tage of this <mark>e o'head</mark> .
}	Here's he problem:	ow to solve	e this
post(0);	parallel task 1	parallel task 2	parallel task 3
for (i=1; i<=N; i++) {	(1=1)	(1=2)	(1=3)
S1: temp = b[i] * c[i];	S1	S1	S1
wait(i-1);	82 post(1)	wait(1)	
<pre>S2: a[i] = a[i-1] + temp;</pre>		S2	



post(i);

pc

Function parallelism

- [§3.2.4] Identify dependences in a loop body.
- If there are independent statements, can split/distribute the loops.

Example:

	Loop-carried dependences:
for (i=0; i <n; i++)="" th="" {<=""><td></td></n;>	
S1: a[i] = b[i+1] * a[i-1];	
S2: b[i] = b[i] * coef;	Loop-Indep. dependences:
S3: c[i] = 0.5 * (c[i] + a[i]);	Note that S4 has no
S4: d[i] = d[i-1] * d[i];	dependences with other
}	statements
	After loop distribution:
for (i=0; i <n; i++)="" th="" {<=""><th>Each loop is a parallel task.</th></n;>	Each loop is a parallel task.
S1: a[i] = b[i+1] * a[i-1];	This is called function
S2: b[i] = b[i] * coef;	parallelism.
S3: c[i] = 0.5 * (c[i] + a[i]);	It can be distinguished from
}	<i>data parallelism</i> , which we saw
	III DOALE and DOACKOSS.
for (i=0; i <n; i++)="" th="" {<=""><th>Further transformations can be</th></n;>	Further transformations can be
$S4 \cdot d[i] = d[i-1] * d[i]:$	performed (see p. 44 of text).
	"s1[i] →A s2[i+1]" implies
1	that S2 at iteration <i>i</i> +1 must be
	executed after S1 at iteration i

Hence, the dependence is not violated if all S2s execute after all S1s.

Characteristics of function parallelism:

- Only gives modest ||ism, does not grow with input size.
- Difficult to balance the load

Can use function parallelism along with data parallelism when data parallelism is limited.

lecture 2	Architecture of Parallel Computers	25

Intuitively, why are these cases different? Read only needs to be in shared memory, one copy

R/W non-conflicting: Only one copy, but can be in local memory.

R/W conflicting: One copy, in shared memory.

Example 1			
	for (i	=1; i<=n; i++)	
Let's assume each iteration	for	(j=1; j<=n; j++) {	
of the for i	S2	: a[i][j] = b[i][j] +	c[i][j];
loop is a parallel task.	\$3	: b[i][j] = a[i][j-1]	* d[i][j];
Fill in the	}		
tableaus <u>here</u> .			
Read-only	y	R/W non-conflicting	R/W conflicting
n, c, d		a, b	i,j

Now, let's assume that each for *j* iteration is a separate task.

Read-only	R/W non-conflicting	R/W conflicting
n, c, d, i	<mark>b</mark>	<mark>a, j</mark>

= Corrected after class

Do these two decompositions create the same number of tasks? No, the first creates n tasks, and the second creates n^2

DOPIPE Parallelism

[§3.2.5] Another strategy for loop-carried dependences is pipelining the statements in the loop.

Consider this situation:	for (i=2; i<=N; i++) {
Loop-carried dependences:	S1: a[i] = a[i-1] + b[i];
Loop-indep. dependences:	<pre>S2: c[i] = c[i] + a[i]; }</pre>
To parallelize, we just need to make sure the two statements are executed in sync:	parallel parallel task 1 task 2
for (i=2; i<=N; i++) {	
a[i] = a[i-1] + b[i];	S1 post(1) weit(1)
<pre>post(i);</pre>	S1 post(2)
}	s_2 wait(2)
	S2
<pre>for (i=2; i<=N; i++) {</pre>	S1 post(n)
<pre>wait(i);</pre>	wait(n)
c[i] = c[i] + a[i];	S2
} G	Question: What's the difference between OACROSS and DOPIPE?

Determining variable scope

[§3.4] This step is specific to the shared-memory programming model. For each variable, we need to decide how it is used. There are three possibilities:

- Read-only: variable is only read by multiple tasks
- R/W non-conflicting: variable is read, written, or both by only one task
- R/W conflicting: variable is written by one task and may be read by another

Architecture of Parallel Computers Lecture 2

26

Example 2 for (i=1; i<=n; i++)</pre> Let's assume that each for j iteration is а for (j=1; j<=n; j++) {</pre> separate task. S1: a[i][j] = b[i][j] + c[i][j]; S2: b[i][j] = a[i-1][j] * d[i][j]; S3: e[i][j] = a[i][j]; } Read-only R/W non-conflicting R/W conflicting n, c, d, i a, b, e

Exercise: Suppose each for *i* iteration were a separate task ...

Read-only	R/W non-conflicting	R/W conflicting
n, c, d	b, e	a, i, j

Privatization

Privatization means making private copies of a shared variable.

What is the advantage of privatization?

It removes read-write conflicts, so tasks can run in parallel, without paying attention to which other task is reading or writing a variable.

Of the three kinds of variables in the table above, which kind(s) does it make sense to privatize? R/W conflicting

Under what conditions is a variable privatizable?

- · If it is always defined (=written) in program order by a task before use (=read) by the same task (Case 1).
- · If its values in different parallel tasks are known ahead of time, allowing private copies to be initialized to the known values (Case 2).

When a variable is privatized, one private copy is made for each thread (not each task).

Result of privatization: R/W conflicting \rightarrow R/W non-conflicting

Let's revisit the examples.

Example 1 for (i=1; i<=n; i++)</pre> With each for i for (j=1; j<=n; j++) {</pre> iteration a separate task, which of the S2: a[i][j] = b[i][j] + c[i][j]; R/W conflicting S3: b[i][j] = a[i][j-1] * d[i][j]; variables are privatizable? }

i (Case 2), *j* (Case 1)

Which case does each such variable fall into?

We can think of privatized variables as arrays, indexed by process ID: i[id]

Example 2

Parallel tasks: each for j loop iteration.

Is the R/W conflicting variable j privatizable? If so, which case does it represent? Yes, Case 2

Reduction

Reduction is another way to remove conflicts. It is based

Suppose we have a large matrix, and need to perform some operation on all of the elementslet's say, a sum of products-to produce a single result.



We could have a single processor undertake this, this is slow and does not make good use of the parallel machine.

Lecture 2

Architecture of Parallel Computers

29

So, we divide the matrix into portions, and have one processor work on each portion.

Then after the partial sums are complete, they are combined into a global sum. Thus, the array has been "reduced" to a single element.

Examples:

- addition (+), multiplication (*)
- Logical (and, or, ...)

The reduction variable is the scalar variable that is the result of a reduction operation.

Criteria for reducibility:

- Reduction variable is updated by each task, and the order of update does not matter
- · Hence, the reduction operation must be associative and commutative

Goal: Compute

 $y = y_{init}$ op x1 op x2 op $x3 \dots$ op x_n

op is a reduction operator if it is commutative

 $u \mathbf{op} v = v \mathbf{op} u$

and associative

 $(u \mathbf{op} v) \mathbf{op} w = u \mathbf{op} (v \mathbf{op} w)$

Lecture 2

Architecture of Parallel Computers

30

Summary of scope criteria

Should be declared private	Should be declared shared	Should be de- clared reduction	Non-privatizable R/W conflicting
Privatizable R/W conflicting	Read-only R/W non-conflicting	Reduction	Declare as shared and protect by synchronization

Example 1	<pre>for (i=1; i<=n; i++)</pre>
with for <i>i</i> parallel tasks	for (j=1; j<=n; j++) {
Fill in the answers	<pre>S2: a[i][j] = b[i][j] + c[i][j];</pre>
here.	00. MIRTIAL = AIRLA IL * AIRLIAL.

Read-only	R/W non-conflicting	R/W conflicting
c, d, n	a, b	i, j

Declare as shared	Declare as private
a, b, c, d, n	i, j

Example 2
with for <i>j</i> parallel tasks
Fill in the answers <u>here</u> .

Read-only	R/W non-conflicting	R/W conflicting
<mark>C</mark> , <i>d</i> , <i>I</i> , n	a, b, e	j

Declare as shared	Declare as private
<i>a, b, c, d</i> , i, n	j

Example 3	<pre>for (i=0; i<n; i++)<="" pre=""></n;></pre>
Consider matrix	<pre>for (j=0; j<n; j++)="" pre="" {<=""></n;></pre>
multiplication.	C[i][j] = 0.0;
	for (k=0; k <n; k++)="" th="" {<=""></n;>
Exercise: Suppose the parallel tasks are for k	C[i][j] = C[i][j] + A[i][k]*B[k][j];

iterations. Determine which variables are conflicting, which should be declared as private, and which need to be protected against concurrent access by using a critical section

Read-only	R/W non-conflicting	R/W conflicting
A, B, i, j, n		C, k

Declare as shared	Declare as private	
A, B, [C], i, j, n	k	

Which variables, if any, need to be protected by a critical section? C

Now, suppose the parallel tasks are for *i* iterations. Determine which variables are conflicting, which should be declared as private, and which need to be protected against concurrent access by using a critical section

Read-only	R/W non-conflicting	R/W conflicting
A,B,n	С	i, j, k

Declare as shared	Declare as private
<i>A,B,</i> C, n	i, j, k

Architecture of Parallel Computers

Synchronization

Synchronization is how programmers control the sequence of operations that are performed by parallel threads.

Three types of synchronization are in widespread use.

- Point-to-point:
 - a pair of *post()* and *wait()*
 - o a pair of send() and recv() in message passing
- Lock
 - a pair of *lock*() and *unlock*()
 - $\circ\;$ only one thread is allowed to be in a locked region at a given time
 - o ensures mutual exclusion
 - used, for example, to serialize accesses to R/W concurrent variables.
- Barrier
 - a point past which a thread is allowed to proceed only when all threads have reached that point.

Lock

What problem may arise here?

	// inside a parallel region	
	<pre>for (i=start_iter; i<end_iter; i++)<="" pre=""></end_iter;></pre>	
Two may fetch a v back late	<pre>sum = sum + a[i]; rariable and increment it concurrently. Then the task that st er overwrites the contribution of the task that wrote it earlier</pre>	tasks each ores it r.

A lock prevents more than one thread from being inside the locked region.

Lecture 2	Architecture of Parallel Computers	33

Goal: Simulate the motion of water currents in the ocean. Important to climate modeling.

Motion depends on atmospheric forces, friction with ocean floor, and "friction" with ocean walls.



(a) Cross sections

To predict the state of the ocean at any instant, we need to solve complex systems of equations.

The problem is *continuous* in both space and time. But to solve it, we *discretize* it over both dimensions.

Every important variable, e.g.,

 pressure 	 velocity 	 currents
------------------------------	------------------------------	------------------------------

has a value at each grid point.

This model uses a set of 2D horizontal cross-sections through the ocean basin.

Equations of motion are solved at all the grid points in one time-step.

- The state of the variables is updated, based on this solution.
- The equations of motion are solved for the next time-step.

Tasks

The first step is to divide the work into tasks.

// inside a parallel region

```
for (i=start_iter; i<end_iter; i++) {
    lock(x);</pre>
```

sum = sum + a[i]; unlock(x);

granularity to lock?

· How to build a lock that is correct and fast.

Barrier: Global event synchronization



A barrier is used when the code that follows requires that all threads have gotten to this point. Example: Simulation that works in terms of timesteps.

Load balance is important.

Execution time is dependent on the slowest thread.

This is one reason for gang scheduling and avoiding time sharing and context switching.

Simulating ocean currents

We will study a parallel application that simulates ocean currents.

Lecture 2

Architecture of Parallel Computers

34

- A task is an arbitrarily defined portion of work.
- · It is the smallest unit of concurrency that the program can exploit.

Example: In the ocean simulation, a task can be computations on-

- a single grid point,
- · a row of grid points, or
- any arbitrary subset of the grid.

Tasks are chosen to match some natural granularity in the work.

- If the grain is small, the decomposition is called
- If it is large, the decomposition is called ______

Threads

Lecture 2

A *thread* is an abstract entity that performs tasks.

- A program is composed of cooperating threads.
- · Each thread is assigned to a processor.
- Threads need not correspond 1-to-1 with processors!

Example: In the ocean simulation, an equal number of rows may be assigned to each thread.

Four steps in parallelizing a program:

- Decomposition of the computation into tasks.
- Assignment of tasks to threads.
- Orchestration of the necessary data access, communication, and synchronization among threads.
- Mapping of threads to processors.

lssues:

What



Together, decomposition and assignment are called partitioning.

They break up the computation into tasks to be divided among threads.

The number of tasks available at a time is an upper bound on the achievable parallelism.

Table 2.1 Steps in the Parallelization Process and Their Goals

Step	Architecture- Dependent?	Major Performance Goals
Decomposition	Mostly no	Expose enough concurrency but not too much
Assignment	Mostly no	Balance workload Reduce communication volume
Orchestration	Yes	Reduce noninherent communication via data locality Reduce communication and synchronization cost as seen by the processor Reduce serialization at shared resources Schedule tasks to satisfy dependences early
Mapping	Yes .	Put related processes on the same processor if necessary Exploit locality in network topology

Architecture of Parallel Computers

Lecture 2

_

- · •
- .
- •
- If this difference is less than a "tolerance" parameter, the solution has converged.
- If so, we exit solver; if not, we do another sweep.

Here is the code for the solver.

 int n; double **A, diff = 0; 	/*size of matrix: (n + 2-by-n + 2) elements*/
 main() begin read(n); A ← malloc (a 2-d array of siz initialize(A); Solve (A); end main 	/*read input parameter: matrix size*/ re n + 2 by n + 2 doubles); /*initialize the matrix A somehow*/ /*call the routine to solve equation*/
10.procedure Solve (A) 11. double **A; 12.begin	/*solve the equation system*/ /*A is an (n + 2)-by-(n + 2) array*/
 int i, j, done = 0; float diff = 0, temp; while (!done) do diff = 0; for i < 1 to n do 	/*outermost loop over sweeps*/ /*initialize maximum difference to 0*/ /*sweep over nonborder points of grid*/
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	<pre>/*save old value of element*/ + A[i,j-1] + A[i-1,j] + /*compute average*/ mp);</pre>
24. end for 25. if $(diff/(n*n) < TOL)$ then d 26. end while 27.end procedure	ione = 1;

Decomposition

A simple way to identify concurrency is to look at loop iterations.

Is there much concurrency in this example? Does the algorithm let us perform more than one sweep concurrently?

Note that-

· Computation proceeds from left to right and top to bottom.

37

Parallelization of an Example Program

[\$2.3] In this lecture, we will consider a parallelization of the kernel of the Ocean application.

The serial program

The equation solver solves a PDE on a grid.

It operates on a regular 2D grid of (n+2) by (n+2) elements.

- The boundary elements in the border rows and columns do not change.
- The interior n-by-n points are updated, starting from their initial values.

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	Ŷ	0	0	0	0	0
0	0	0	0-	→Ě	-0	0	0	0	0
0	0	0	0	Ó	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

Expression for updating each interior point: $A[ij] = 0.2 \times (A[ij] + A[ij - 1] + A[i - 1, j] + A[ij + 1] + A[i + 1, j])$

- The old value at each point is replaced by the weighted average of itself and its 4 nearest-neighbor points.
- · Updates are done from left to right, top to bottom.
 - $^\circ\,$ The update computation for a point sees the new values of points above and to the left, and
 - · the old values of points below and to the right.
 - This form of update is called the Gauss-Seidel method.

During each sweep, the solver computes how much each element has changed since the last sweep.

Architecture of Parallel Computers

```
Lecture 2
```

38

- · Thus, to compute a point, we use
 - the updated values from the point above and the point to the left, but
 - the "old" values of the point itself and its neighbors below and to the right.

Here is a diagram that illustrates the dependences.

ې مېخ	÷γ∕-	»ợ-	×	»ơ-	×ó	×	»ŏ-	»∕o′
_ŏ→ŏ	-≫-	×∛-	>0 -)	×	×	÷0́-	ж
_∛→∛-	-∛-	>∛-	≫	≫	×	≫	≫	≫Ó
_ğ→ğ	-×∛-	>∛-	×	×	≫	×	≫	»∕∕
ð ð	ð	ğ	ø	ø	ø	ø	ø	ø
ð ð	ð	Å.	ø	ø	ø	ø	ø	ø
ð ð	\X	Þ	ø	ø	ø	ø	ø	ø
ðğ	Ķ	×.	ø	ø	ø	,Ó	ø	ø
ð ð	Å.	ğ	ø	ø	ø	ø	ø	ø
ð ð	ð	ð	ø	ø.	Ø	ø	0	0

The horizontal and vertical lines with arrows indicate dependences.

The dashed lines along the antidiagonal connect points with no dependences that can be computed in parallel.

Of the O(___) work in each sweep, ∃ concurrency proportional to _____ along antidiagonals.

How could we exploit this parallelism?

 We can *leave loop structure alone* and let loops run in parallel, inserting *synchronization ops* to make sure a value is computed before it is used.

Why isn't this a good idea?

- We can change the loop structure, making
 - the outer for loop (line 17) iterate over anti-diagonals, and
 - the inner **for** loop (line 18) iterate over elements within an antidiagonal.

Why isn't this a good idea?

Lecture 2

The Gauss-Seidel algorithm doesn't *require* us to update the points from left to right and top to bottom.

It is just a convenient way to program on a uniprocessor.

We can compute the points in another order, as long as we use updated values frequently enough (if we don't, the solution will converge, but more slowly).

Red-black ordering

Let's divide the points into alternating "red" and "black" points:



To compute a red point, we don't need the updated value of any other red point. But we need the updated values of 2 black points.

And similarly for computing black points.

Thus, we can divide each sweep into two phases.

- First we compute all red points.
- Then we compute all black points.

True, we don't use any updated black values in computing red points.

But we use all updated red values in computing black points.

Whether this converges more slowly or faster than the original ordering depends on the problem.

Lecture 2	Architecture of Parallel Computers

The only difference is that for has been replaced by for_all.

A for_all just tells the system that all iterations can be executed in parallel.

With **for_all** in both loops, all n^2 iterations of the nested loop can be executed in parallel.

We could write the program so that the computation of one row of grid points must be assigned to a single processor. How would we do this?

With each row assigned to a different processor, each task has to access about 2n grid points that were computed by other processors; meanwhile, it computes n grid points itself.

So the communication-to-computation ratio is O(1).

Assignment

How can we statically assign elements to processes?

 One option is "block assignment"—Row *i* is assigned to process *[i/p]*.



- Another option is "cyclic assignment—Process *i* is assigned rows *i*, *i*+*p*, *i*+2*p*, etc.
- · Another option is 2D contiguous block partitioning.

We could instead use dynamic assignment, where a process gets an index, works on the row, then gets a new index, etc. Is there any advantage to this?

What are advantages and disadvantages of these partitionings?

41

But it does have important advantages for parallelism.

- How many red points can be computed in parallel? $n^{2}/2$
- How many black points can be computed in parallel? n²/2

Red-black ordering is effective, but it doesn't produce code that can fit on a single display screen.

A simpler decomposition

Another ordering that is simpler but still works reasonably well is just to ignore dependences between grid points within a sweep.

A sweep just updates points based on their nearest neighbors, regardless of whether the neighbors have been updated yet.

Global synchronization is still used between sweeps, however.

Now execution is no longer deterministic.

The number of sweeps needed, and the results, may depend on the number of processors used.

But for most reasonable assignments of processors, the number of sweeps will not vary much.

Let's look at the code for this.

j] +

Lecture 2 Architecture of Parallel Computers

Static assignment of rows to processes reduces concurrency

But block assignment reduces communication, by assigning adjacent rows to the same processor.

How many rows now need to be accessed from other processors?

So the communication-to-computation ratio is now only O(____)

Orchestration

Once we move on to the orchestration phase, the computation model affects our decisions.

Data-parallel model

In the code below, we assume that global declarations are used for shared data, and that any data declared within a procedure is private.

Global data is allocated with g_malloc.

Differences from sequential program:

- for_all loops
- · decomp statement
- mydiff variable, private to each process
- reduce statement

1. 2.	<pre>int n, nprocs; double **A, diff = 0;</pre>	/*grid size (n+2×n+2) and # of processes*/
з.	main()	
4.	begin	
5.	read(n); <pre>read(nprocs);</pre>	;/*read input grid size and # of processes*/
6.	A ← <mark>G_MALLOC</mark> (a 2-d array o	f size n+2 by n+2 doubles);
7.	initialize(A);	/*initialize the matrix A somehow*/
8.	Solve (A);	/*call the routine to solve equation*/
9.	end main	
10.	procedure Solve(A)	/*solve the equation system*/
11.	double **A;	/* A is an (n+2×n+2) array*/
12.	begin	
13.	int i, j, done = 0;	
14.	<pre>float mydiff = 0, temp;</pre>	
14a.	<pre></pre>	
15.	while (!done) do	/*outermost loop over sweeps*/
16.	mydiff = 0;	/*initialize maximum difference to 0 */
17.	for_all i ← 1 to n do	/*sweep over non-border points of grid*/
18.	<mark>for_all</mark> j ← 1 to n do	
19.	<pre>temp = A[i,j];</pre>	/*save old value of element*/
20.	A[i,j] ← 0.2 * (A[i,j] + A[i,j-1] + A[i-1,j] +
21.	A[i, j+1] + A[i+1, j])	; /* compute average*/
22.	<pre>mydiff += abs(A[i,j] -</pre>	- temp);
23.	end <mark>for_all</mark>	
24.	end for_all	
24a.	REDUCE (mydiff, diff, Al	; ;
25.	if (diff/(n*n) < TOL) ther	done = 1;
∠0.	end whitte	

The **decomp** statement has a twofold purpose.

· It specifies the assignment of iterations to processes.

The first dimension (rows) is partitioned into *nprocs* contiguous blocks. The second dimension is not partitioned at all.

Specifying [CYCLIC, *, nprocs] would have caused a cyclic partitioning of rows among nprocs processes.

Specifying [*, CYCLIC, nprocs] would have caused a cyclic partitioning of columns among nprocs processes.

Lecture 2	Architecture of Parallel Computers	45



What are the main differences between the serial program and this program?

• The first process creates *nprocs*-1 worker processes. All *n* processes execute *Solve*.

All processes execute the same code.

But all do *not* execute the same instructions at the same time.

Specifying [BLOCK, BLOCK, nprocs] would have implied a 2D contiguous block partitioning.

• It specifies the assignment of grid data to memories on a distributedmemory machine. (Follows the *owner-computes* rule.)

The mydiff variable allows local sums to be computed.

The **reduce** statement tells the system to add together all the *mydiff* variables into the shared *diff* variable.

Shared-memory model

In this model, we need mechanisms to create processes and manage them.

After we create the processes, they interact as shown on the right.



Lecture 2

2

Architecture of Parallel Computers

46

- Private variables like *mymin* and *mymax* are used to control loop bounds.
- All processors need to—
 - complete an iteration before any process tests for convergence. <u>Why</u>?
 - test for convergence before any process starts the next iteration. Why?

Notice the use of barrier synchronization to achieve this.

 Locks must be placed around updates to *diff*, so that no two processors update it at once. Otherwise, inconsistent results could ensue.

<u>p1</u>	<u>p2</u>	
r1 \leftarrow diff		$\{p_1 \text{ gets 0 in its } r1\}$
	rl ← diff	{ p ₂ also gets 0}
r1 \leftarrow r1+r2		{ p1 sets its r1 to 1}
	rl ← rl+r2	{ p ₂ sets its rl to 1}
diff \leftarrow rl		$\{ p_1 \text{ sets } diff \text{ to } 1 \}$
	diff ← r1	{ p ₂ also sets <i>diff</i> to 1}

If we allow only one processor at a time to access *diff*, we can avoid this race condition.

What is one performance problem with using locks?

Note that at least some processors need to access *diff* as a non-local variable.

What is one technique that our shared-memory program uses to diminish this problem of serialization?

Message-passing model

The program for the message-passing model is also similar, but again there are several differences.

- There's no shared address space, so we can't declare array A to be shared.
- Instead, each processor holds the rows of A that it is working on.
- The subarrays are of size (*n*/*nprocs* + 2) × (*n* + 2). This allows each subarray to have a copy of the boundary rows from neighboring processors. Why is this done?

These *ghost* rows must be copied explicitly, via **send** and **receive** operations.

Note that **send** is not synchronous; that is, it doesn't make the process wait until a corresponding **receive** has been executed.

What problem would occur if it did?

 Since the rows are copied and then not updated by the processors they have been copied from, the boundary values are more out-ofdate than they are in the sequential version of the program.

This may or may not cause more sweeps to be needed for convergence.

• The indexes used to reference variables are *local* indexes, not the "real" indexes that would be used if array *A* were a single shared array.

Architecture of Parallel Computers

1. int pid, n, b;	/*process id, matrix dimension and number of processors to be used*/
<pre>2.float **myA; 3.main() 4 begin</pre>	
5. read(n); read(nprocs); 8a. CREATE (nprocs-1, Solve	/*read input matrix size and number of processes*/
<pre>8b. Solve(); 8c. WAIT_FOR_END (nprocs-1) 9.end main</pre>	<pre>/*main process becomes a worker too*/ /*wait for all child processes created to terminate*/</pre>
<pre>10. procedure Solve() 11. begin 13. int i.i. pid n' = n/nn</pre>	racs done = 0.
14. float temp, tempdiff, m	ydiff = 0; /*private variables*/
6. myA ← malloc(a 2-d array	of size [n/nprocs + 2] by n+2);
<pre>7. initialize (myA);</pre>	/*my assigned rows of A*/ /*initialize my rows of A, in an unspecified way*/
15. while (!done) do	
16. $mydiff = 0;$	/*set local diff to 0*/
16b. if (pid != nprocs-1) th	en
SEND (&myA[n',0], n*siz	<pre>ceof(float),pid+1,ROW);</pre>
16d. if (pid != 0) then RECE	en
RECEIVE(&myA[n'+1,0],	<pre>n*sizeof(float), pid+1,ROW);</pre>
	/*border rows of neighbors have now been copied into myA[0,*] and myA[n'+1,*]*/
17. for i ← 1 to n' do	/*for each of my (nonghost) rows*/
 18. IOT] ← 1 to n do 19. temp = mvA[i,j]; 	/*for all honborder elements in that row*/
20. myA [i,j] = 0.2 * (myA [i,j] + myA [i,j-1] + myA [i-1,j] +
21. myA[i,j+1] + myA 22 mydiff += abs(myA[[i+1,j]); i.il - temp):
23. endfor	r,jj cemp/,
24. endfor	
	/*communicate local diff values and determine if done; can be replaced by reduction and broadcast*/
25a. II (pid != 0) then 25b. SEND (mydiff, sizeof	(float), 0, DIFF);
25c. RECEIVE (done, sizeo	<pre>f(int),0,DONE);</pre>
25d. <mark>else</mark>	/*pid 0 does this*/
25e. for i ← 1 to nproc 25f PECETVE (tempdiff	cs-1 do /*for each other process*/
25q. mydiff += tempdiff	; /*accumulate into total*/
25h. endfor	
25i if (mydiff/(n*n) < TO	DL) then done = 1;
25]. for 1 ← 1 to nproc 25k SEND (done, sizeof	(int) i DONE)
251. endfor	
25m. endif	
26. endwhile 27. end procedure	
the procedure	
Lecture 2 Archite	ecture of Parallel Computers 50

Parallel access to linked data structures

[Solihin Ch. 4] Answer the questions below.

Name some linked data structures.

Lecture 2

What operations can be performed on all of these structures?

Why is it hard to parallelize these operations?

Explain how the following code illustrates such a dependence.

```
void addValue(pIntList pList, int key, int x) {
```

```
pIntListNode p = pList->head;
while (p != NULL) {
    if (p->key == key)
        S1: p->data = p->data + x;
    S2: p = p->next;
```

In the notation introduced in Lecture 5, how would the dependence be written?

If we just look at the loops in an "LDS" program, we won't find any parallelism to be exploited.

So, where can we find the opportunity to execute anything in parallel? The "algorithm level"—parallelism between the operations that are performed on the LDS.

Lecture 2

} } 49

Conceptually, we can allow several operations to be performed (partially) in parallel. What kind of operations? search, insertion, deletion

But how do we decide which operations can be performed in parallel?

Correctness of parallel LDS operations

Serializability: A parallel execution of a group of operations (or primitives) is said to be *serializable* if there is some sequence of operations (or primitives) that produce an identical result.

Suppose a node insertion i_1 and a node deletion d_1 are performed in parallel. The outcome must be equivalent to either

*i*₁ followed by *d*₁, or *d*₁ followed by *i*₁, or

Conflict between two insertions







Conflict between an insertion and a deletion

Let's look at the simple case of a singly-linked list.

Suppose two items are inserted in parallel: insert both 4 and 5. Serializable outcomes:

insert 4, then insert 5

or insert 5, then insert 4 In any case,

both 4 and 5 must be in the list at the end of execution

What could happen if the operations are not parallelized correctly? Node 4



Thread 1 wants to delete node "5"

prev points to node "3" p points to node "5"

Thread 0 wants to insert node "4" prev points to node "3" 4 •

Thread 0: newNode->next = p; prev->next = newNode;

(c)

Serializable outcome:

delete 5, insert 4

or insert 4, delete 5

in both cases, at the end of execution, node 4 is in the list, but node 5 is not in the list

In the case shown, node 4 is lost.

What would be a sequence that produces another incorrect result? What would happen with this sequence? (You may use this worksheet.)

Lecture 2

Architecture of Parallel Computers

53

Lecture 2

Architecture of Parallel Computers

54

Depending on when the insertion code is executed,

- node 6 will be found, or
- node 6 may not be found, and an uninitialized link may be followed.

Conflict between a deletion and a search

- Deletion and search
 - o delete 5, then search for 5
 - \circ search for 5, then delete 5
- Possible outcomes
 - Node 5 may be found or not found Node 5 is deleted from the list
 - What, if anything, is the problem with these outcomes? Neither; both are serializable

Conflict between a deletion and a search operations



Conflict between an insertion and a search

Conflict between an insertion and a search operations



we attempt

insert 5, then search 6

or search 6, then insert 5

in both cases, at the end of execution,

- 5 must be in the list, and
- 6 must be found •

- Fine-grain lock approach
 - o A lock is associated with each node
 - Each operation locks only nodes that need to be accessed exclusively. Complex: Deadlock can occur; memory allocation and deallocation become more
 - complex

Parallelization among readers

- Basic idea
 - o (Read-only) operations that do not modify the list can execute in parallel. o (Write) operations that modify the list execute sequentially
- How to enforce
 - A read-only operation acquires a read lock
 - A write operation acquires a write lock

Construct a lock-compatibility table

Already-granted lock	Read lock requested	Write lock requested
Read lock	Yes	No
Write lock	No	No

Example

<pre>IntListNode_Search(int x)</pre>	<pre>IntListNode_Insert(node *p)</pre>
{	{
<pre>acq_read_lock();</pre>	<pre>acq_write_lock();</pre>
<pre>rel_read_lock();</pre>	<pre>rel_write_lock();</pre>
}	}
Global-lock approach	

Each operation logically has two steps

Traversal

Node insertion: Find the correct location for the node

Lecture 2	Architecture of Parallel Computers

Architecture of Parallel Computers

• Parallel execution of two operations that affect a common node, in which at least one operation involves writing to the node, can produce conflicts that lead to non-serializable outcome Under some circumstances, a serializable outcome may still be achieved, despite the conflicts

Conflicts can also occur between LDS operations and memory-management functions such as

- Node search: Find the sought-for node
- o List modification
- structure, then modify the list, then release the lock.

- so the assumptions must be re-validated.

Example

....

Lecture 2

IntListNode_Insert(node *p)

{

/* perform traversal */

acq write lock(); /* then check validity:

nodes still there? link still valid? */

/* if not valid, repeat traversal */

/* if valid, modify list */

rel_write_lock();

}

Fine-grain locking approach

- Associate each node with a lock (read, write).
- Each operation locks only needed nodes.
- (Read and write) operations execute in parallel except when they conflict on some nodes. Fill in the blanks below Nodes that will be modified are write-locked. 0
 - o Nodes that are read and must remain unchanged are read-locked.
- Pitfall: Deadlock becomes possible.
 - Suppose one operation locks node 1 and then needs to lock node 2, while another operation locks node 2 and then needs to lock node 1.

Architecture of Parallel Computers

• Then neither operation can complete before the other operation frees the lock it is holding.

 Deadlocks can be prevented by imposing a global ordering on the order in which nodes are acquired

Example

void insert(pIntList pList, int x) {

int succeed;

- .../* traversal code to find where to insert */
- /* insert the node at head or between prev & p */

succeed = 0:

- do {
- acq write lock(prev);
- acq_read_lock(p);
- if (prev->next != p || prev->deleted || p->deleted)
- rel_write_lock(prev);
- rel_read_lock(p);
- .../* repeat traversal */
- }
- else
 - succeed = 1;
- } while (!succeed);
- /* prev and p are now valid, so insert node */
- newNode->next = p;
- if (prev != NULL)
- prev->next = newNode;
- else
- pList->head = newNode;
- rel write lock(prev);
- rel_read_lock(p);

ł

Questions

59

58

57

Works well if structure is modified infrequently

Parallel traversal, followed by sequential list modifications

- Node deletion: Find the node to delete

Basic idea: perform the traversal in parallel, but modify the list in a critical section, i.e., lock the

Pitfall

Main Observations

mentioned above.

Global lock approach

allocation and deallocation. Parallelization strategies Parallelization among readers Very simple

Relatively simple

- The list may have changed by the time the write-lock is acquired,

What do the tests prev->de	leted and p->deleted mean?		else /* delete non-	nead node */	
			prev->next = p->ne	ext;	
Why is garbage collection used	I, rather than explicit deletion?		p->deleted = 1; /*de	on'tdeallocate;mark deleted*/	
			rel_write_lock(prev)	;	
			<pre>rel_write_lock(p);</pre>		
The delete operation is similar	; code that is the same is shown in green.		}		
void delete (pIntList	<pre>pList, int x) {</pre>		-		
int succeed;					
/* traversal code	to find node to delete */				
/* node has been fo	ound; perform the deletion */				
<pre>succeed = 0;</pre>					
do {					
acq_write_lock(pr	cev);				
<pre>acq_write_lock(p)</pre>	;				
if (prev->next != p	> prev->deleted p->deleted)				
{					
rel_write_loc	k (prev) ;				
rel_write_loc	k (p) ;				
/* repeat tr	<pre>raversal; return if not found */</pre>				
}					
else					
<pre>succeed = 1;</pre>					
<pre>} while (!succeed);</pre>	:				
<pre>/* prev and p are r</pre>	now valid, so delete node */				
if (prev == NULL)	/* delete head node */				
acq_write_lock(pI	ist);				
pList->head = p->	<pre>>next;</pre>				
rel_write_lock(pI	.ist);				
}					
Lecture 2	Architecture of Parallel Computers	61	Lecture 2	Architecture of Parallel Computers	62

Data parallel algorithms¹

(Guy Steele): The data-parallel programming style is an approach to organizing programs suitable for execution on massively parallel computers.

In this lecture, we will-

- characterize the programming style,
- examine the building blocks used to construct data-parallel programs, and
- see how to fit these building blocks together to make useful algorithms.

All programs consist of code and data put together. If you have more than one processor, there are various ways to organize parallelism.

- Control parallelism: Emphasis is on extracting parallelism by orienting the program's organization around the parallelism in the code.
- parallelism: Emphasis is on organizing programs to extract parallelism from the organization of the data.

With data parallelism, typically all the processors are at roughly the same point in the program.

Control and data parallelism vs. SIMD and MIMD.

- · You may write a data-parallel program for a MIMD computer, or
- a control-parallel program which is executed on a SIMD computer.

Emphasis in this talk will be on styles of organizing programs. It becomes an engineering issue whether it is appropriate to organize the hardware to match the program.

The sequential programming style, typified by C and Pascal, has building blocks like-

- scalar arithmetic operators,
 control structures like if ... then ... else, and subscripted array references. .

¹Video © 1991, Thinking Machines Corporation. This video is available from University Video Communications, http://www.uvc.co

Lecture 9

Architecture of Parallel Computers

The programmer knows essentially how much these operations cost. E.g., addition and subtraction have similar costs; multiplication may be more expensive

To write data-parallel programs effectively, we need to understand the cost of data-parallel operations

- · Elementwise operations (carried on independently by processors; typically operations and tests).
- · Conditional operations (also elementwise, but some processors may not participate, or act in various ways).
- Replication
- Permutation
- Parallel prefix (scan)

An example of an elementwise operation:





The results can be used to "conditionalize" future operations:

if(A > B)C = A + B

© 2014 Edward F. Gehringer

CSC/ECE 506 Lecture Notes, Spring 2014

2

• • 0 0 0 • 0 • 4 5 2 1 3 2 3 1 1 1 6 2 1 3 0 5 0 0 5 8 2 0 4 0

The set of bits that is used to conditionalize the operations is frequently called a condition mask or a context. Each processor can perform different computations based on the data it contains.

Building blocks

Communications operations:

- : Get a single value out to all processors. This operation happens so frequently that is worthwhile to support in hardware. It is not unusual to see a hardware bus of some kind.
- Spreading (nearest-neighbor grid). One way is to have each row copied to its nearest neighbor.

4	3	6	2	5	3	4	9	2
Я	3	6	2	5	3	4	9	2
Я	3	6	2	5	3	4	9	2
Я	3	6	2	5	3	4	9	2
Я	3	6	2	5	3	4	9	2
Я	3	6	2	5	3	4	9	2
Я	3	6	2	5	3	4	9	2
Ч	3	6	2	5	3	4	9	2

A better way is to use a copy-scan:

- · On the first step, the data is copied to the row that is directly below
- · On the second step, data is copied from each row that has the data to the row that is two rows below
- · On the third step, data is copied from each row to the row that is four rows below.

In this way, the row can be copied in logarithmic time, if we have the necessary interconnections.

essentially the inverse of broadcasting. Each processor has an element, and you are trying to combine them in some way to produce a single result.



Summing a vector in logarithmic time:

x 0	X ₁	X 2	Х3	X 4	X5	x ₆	X	7
_							$\overline{\}$	
X 0	Σ01	X 2	Σ2	X 4	Σ4	X 6	Σ	7 6
		/	/			/	/	
x 0	Σ10	X 2	Σ30	X 4	Σ ⁵ 4	X 6	Σ	7 4
						_	_	
\mathbf{x}_0	Σ1	X 2	Σ 3	X_	Σ 3	X ₆	Σ	7

Most of the time during the course of this algorithm, most processors have not been busy.

So while it is fast, we haven't made use of all the processors.

Suppose you don't turn off processors; what do you get? Vector sum-prefix (sum-scan).



Each processor has received the sum of what it contained, plus all the processors preceding it.

We have computed the sums of all prefixes-initial segments-of the array.

This can be called the checkbook operation; if the numbers are a set of credits and debits, then the prefixes are the set of running balances that should appear in your checkbook.

_____. We wish to assign a different number to each processor.



· Regular permutation.



Of course, one can do shifting on two-dimensional arrays too; you might shift it one position to the north.

Another kind of permutation is an odd-even swap:



Distance 2^k swap:



Some algorithms call for performing irregular permutations on the data.

Lecture 9

Architecture of Parallel Computers

5

7

CBEAHDFG ABCDEFGH

The permutation depends on the data. Here we have performed a sort. (Real sorting algorithms have a number of intermediate steps.)

Example: image processing

Suppose we have a rocket ship and need to figure out where it is.

Some of the operations are strictly local. We might focus in on a particular region, and have each processor look at its values and those of its neighbor.

This is a local operation; we shift the data back and forth and have each processor determine whether it is on a boundary.

When we assemble this data and put it into a global object, the communication patterns are dependent on the data; it depends on where the object happened to be in the image.

Irregularly organized data

Most of our operations so far were on arrays, regularly organized data.

We may also have operations where the data are connected by pointers.

In this diagram, imagine the processors as being in completely different parts of the machine, known to each other only by an address.

doubling:

I originally thought that nothing could be more essentially sequential than processing a linked list. You just can't find the third one without going through the second one. But I forgot that there is processing power at each node.

The most important technique is *pointer doubling*. This is the pointer analogue of the spreading operation we looked at earlier to make a copy of a vector into a matrix in a logarithmic number of steps.

In the first step, each processor makes a copy of the pointer it has to its neighbor.

© 2014 Edward F. Gehringer

CSC/ECE 506 Lecture Notes, Spring 2014

6

In the rest of the steps, each processor looks at the processor it is pointing to with its extra pointer, and gets a copy of *its* pointer.

In the first step, each processor has a pointer to the next processor. But in the next step, each processor has a pointer to the processor two steps away in the linked list.

ŢĿŢĿŢĿŢĿŢĿ

In the next step, each processor has a pointer to the pointer four processors away (except that if you fall off the end of the chain, you don't update the pointer).

Eventually, in a logarithmic number of steps, each processor has a pointer to the end of the chain.

How can this be used? In partial sums of a linked list.

x <u>1 x 2 x 3 x 4 x 5 x 6 x 7</u>

At the first step, each processor takes the pointer to its neighbor.

At the next step, each processor takes the value that it holds, and adds it into the value in the place pointed to:

$\nabla 0 \nabla 1$	52	x 3	$\nabla 4$	5 5	~	5 5	7
14 ŏ 1 4 ó	1 4 6	4 2	<u> </u>		4	5 4	6
° - °	• ≚		- 4	► ∸	+ ·	<u>∼</u>	~
	ス	へ	ー	ー	7	_	_





And after the third step, you will find that each processor has gotten the sum of its own number plus all the preceding ones in the list.

$\Sigma_{0}^{0} \Sigma_{0}^{1} \Sigma_{0}^{2} \Sigma_{0}^{3} \Sigma_{0}^{4} \Sigma_{0}^{5} \Sigma_{0}^{6} \Sigma_{0}^{7}$

Speed vs. efficiency: In sequential programming, these terms are considered to be synonymous. But this coincidence of terms comes about only because you have a single processor.

In the parallel case, you may be able to get it to go fast by doing extra work.

Let's take a look at the serial vs. parallel algorithm for summing an array.

-Reduction

	Serial	Parallel
Processors	1	Ν
Time steps	<i>N</i> –1	log N
Additions	<i>N</i> –1	<i>N</i> –1
Cost	<i>N</i> –1	N log N
Efficiency	1	$-\frac{1}{\log N}$
	Sur	n – Prefix
	Serial	Parallel
Processors	1	n
Time steps	<i>n</i> –1	log n
Additions	<i>n</i> –1	n (log n-1)
Cost	<i>n</i> –1	n log n
Efficiency	1	$\frac{\log n - 1}{\log n}$

The serial version of sum-prefix is similar to the serial version of sum-reduction, but you save the partial sums. You don't need to do any more additions, though.

In the parallel version, the number of additions is much greater. You use n processors, and commit log n time steps, and nearly all of them were busy.

As *n* gets large, the efficiency is very close to 1. So this is a very efficient algorithm. But in some sense, the efficiency is bogus; we've kept the processors

busy doing more work than they had to do. Only n-1 additions are really required to compute sum-prefix. But n(log n-1) additions are required to do it

Thus, the business of measuring the speed and efficiency of a parallel algorithm is tricky. The measures I used are a bit naïve. We need to develop better measures.

Exercise: Submit your answers here.

Calculate the speedup of summing a vector using copy-scan (turning off the processors that are not in use)

- How long does it take to sum the vector serially?
- How long does it take to sum it using copy-scan? •
- What is the speedup?

What is the efficiency (speedup ÷ # of processors) of summing a vector with copy-scan?

In the parallel version of summing an array via sum-prefix, a "bogus" efficiency is mentioned. What would be the "non-bogus" efficiency of the same algorithm

Putting the building blocks together

Let's consider matrix multiply



One way of doing this with a brute-force approach is to use n^3 processors.



Lecture 9

Architecture of Parallel Computers



The columns of the second array are skewed



The two arrays are overlaid, and they then look like this

This is a systolic algorithm; it rotates both of the source matrices at the same time.

> The first source matrix is rotated horizontally The second source matrix is rotated vertically



At the first time step, the 2nd element of the first row and the 2nd element of the first column meet in the upper left corner. They are then multiplied and accumulated

At the second time step, the 3rd element of the first row and the 3rd element of the first column meet in the upper left corner. They are then multiplied and accumulated

At the third time step, the 4th element of the first row and the 4th element of the first column meet in the upper left corner. They are then multiplied and accumulated.

At the fourth time step, the 1st element of the first row and the 1st element of the first column meet in the upper left corner. They are then multiplied and accumulated

Architecture of Parallel Computers

The same thing is going on at all the other points of the matrix.

1. Replicate. The first step is to make copies of the first source array, using a spread operation.

2. Replicate. Then we will do the same thing with the second source, spreading those down the cube.

So far, we have used O(log n) time.

3. Elementwise multiply. n³ operations are performed, one by each processor.

4. Perform a parallel sum operation, using the doubling-reduction method.

We have multiplied two matrices in $O(\log n)$ time, but at the cost of using n^3 processors.

> Brute force: n^3 processors O(log n) time

Also, if we wanted to add the sum to one of the matrices, it's in the wrong place, and we would incur an additional cost to move it.

Cannon's method

There's another method that only requires n^2 processors. We take the two source arrays and put them in the same n^2 processors. The result will also show up in the same n^2 processors

We will prethe two source arrays.

· The first array has its rows skewed by different amounts.

© 2014 Edward F. Gehringer

CSC/ECE 506 Lecture Notes, Spring 2014

10

-sum-

The serves to cause the correct elements of each row and

O(n) time

Instead of the rocket ship earlier in the lecture, we'll consider a smaller region. (This is one of the problems in talking about data-parallel algorithms. They're useful for really large amounts of data, but it's difficult to show that on the screen.)

We have a number of regions in this image. There's a large central "green" region, and a "red-orange" region in the upper right-hand corner. Some disjoint regions have the same color.

We would like to compute a result in which each region gets assigned a distinct number.

We don't care which number gets assigned, as long as the numbers are distinct (even for regions of the same color

0	0	2	2	2	5	5	5
8	0	0	2	2	2	2	5
8	8	0	19	2	2	2	23
8	8	19	19	19	19	23	23
8	19	19	19	19	19	23	23
8	19	19	19	19	23	23	23
8	49	49	19	19	23	23	23
19	49	49	49	60	60	60	60

For example, here the central green region has had all its pixels assigned the value 19

The squiggly region in the upper left corner has received 0 in all its pixels.

The region in the upper right, even though the same color as the central green region. has received a different value

Let's see how all the building blocks we have discussed can fit together to make an interesting algorithm.

9

Cannon's method: n² processors

column to meet at the right time.

An additional benefit is that the matrix ends up in the right place.

Labeling regions in an image

Let's consider a really big example.

First, let's assign each processor a different number.

Here I've assigned the numbers sequentially across the rows, but any distinct numbering would do.

We've seen how the enumeration technique can do this in a logarithmic number of time steps.

					_	
		V				
	Ж	+				
/	1		\backslash		1	
			-	Ж	+	
			7			

 0
 1
 2
 3
 4
 5
 6
 7

 8
 9
 10
 11
 12
 13
 14
 15

 16
 17
 18
 19
 20
 21
 22
 23

 24
 25
 26
 27
 28
 29
 30
 31

 32
 33
 34
 5
 63
 37
 38
 39

 40
 41
 42
 43
 44
 45
 46
 47

 48
 49
 50
 51
 52
 53
 54
 55

 56
 57
 58
 59
 60
 61
 62
 63

Next, we have each of the pixels examine the values of its eight neighbors.

We shift it up, down, left, right, to the northeast, northwest, southeast, and southwest.

This is enough for each processor to do elementwise computation and decide whether it is on the border.

(There are messy details, but we won't discuss them here, since they have little to do with parallelism.)

The next computation will be carried out only by processors that are on the borders (an example of conditional operation).

We have each of the processors again consider the pixel values that came from its neighbors, and

inquire again, using shifting, if each of its neighbors are border elements.

This is enough information to figure out which of your neighbors are border elements in the same region, so you can construct pointers to them.



Lecture 9

Architecture of Parallel Computers

13

Data-parallel programming makes it easy to organize operations on large quantities of data in massively parallel computers.

It differs from sequential programming in that its emphasis is on operations on entire sets of data instead of one element at a time.

You typically find fewer loops, and fewer array subscripts.

On the other hand, data-parallel programs are like sequential programs, in that they have a single thread of control.

In order to write good data-parallel programs, we must become familiar with the necessary building blocks for the construction of data-parallel algorithms.

With one processor per element, there are a lot of interesting operations which can be performed in constant time, and other operations which take logarithmic time, or perhaps a linear amount of time.

This also depends on the connections between the processors. If the hardware doesn't support sufficient connectivity among the processors, a communication operation may take more time than would otherwise be required.

Once you become familiar with the building blocks, writing a data-parallel program is just as easy (and just as hard) as writing a sequential program. And, with suitable hardware, your programs may run much faster.

Exercise: Run through Lim's algorithm on the grid given here.

Questions and answers: [not shown during class] Question: (Bert Halstead): Do you ever get into problems when you have highly data-dependent computations, and it's hard to keep more than a small fraction of the processors doing the same operation at the same time?

Answer: Yes. That's one reason for making the distinction between the dataparallel style and _______hardware. The best way to design a system to give you the most flexibility without making it overly difficult to control is, I think, still an open research question.

Question (Franklin Turback): Your algorithms seem to be based on the assumption that you actually have enough processors to match the size of your problem. If you have more data than processors, it seems that the logarithmic time growth is no longer justified.

Answer: There's no such thing as a free lunch. Making the problem bigger makes it run slower. If you have a much larger problem that won't fit, you're going to have to buy a larger computer.

0	1	2		4	5	6	
8	9	10	11		13	14	15
	17	18	19	20	21	22	23
	25	26		28	29	30	
32	33				37	38	
40	41	42		44	45		
48	49	50	51	52	53	54	55
56	1		59	60	61	62	63

0	0	2		2	5	5	
8	0	0	2	2	2	2	5
	8	0	19	2	2	2	23
	8	19		19	19	23	
8	19				19	23	
8	19	19		19	23		
8	49	49	19	19	23	23	23
49			49	60	60	60	60

Now we have stitched together the borders in a linked list.

We now use the pointerdoubling algorithm. Each pixel on the borders considers the number that it was assigned in the enumeration step.

We use the pointer-doubling algorithm to do a reduction step using the **min** operation.

Each linked list performs pointer-doubling around that list, and determines which number is the smallest in the list.

Then another pointer-doubling algorithm makes copies of that minimum all around the list.

0 0 2 2 2 5 5 5

8 0 0 2 2 2 2 5

8 8 0 19 2 2 2 23

8 8 19 19 19 19 23 23

8 19 19 19 19 19 23 23

 8
 19
 19
 19
 23
 23
 23

 8
 49
 49
 19
 19
 23
 23
 23

49 49 49 49 60 60 60 60

Finally, we can use ______ operation, not on linked lists, but by operating on the columns (or the rows) to copy the processor labels from the borders to the rows.

Other items, particularly those on the edge, may need the numbers propagated up instead of down. So you do a scan in both directions.

The operation used is a noncommutative operation that copies the old number from the neighbor, unless it comes across a new number.

This is known as Lim's algorithm.

Region labeling: $O(n^2)$ processors. $O(\log n)$ time

(Each of the steps was either constant time or O(log n) time.)

© 2014 Edward F. Gehringer

CSC/ECE 506 Lecture Notes, Spring 2014

14

Question: How about portability of programs to different machines?

Answer: Right now it's very difficult, because so far, we haven't agreed on standards for the right building blocks to support. Some architectures support some building blocks but not others. This is why you end up with nonportabilities of efficiencies of running times.

Question: For dealing with large sparse matrices, there are methods that we use to reduce complexity. If this is true, how do you justify the overhead cost of parallel processing?

Answer: Yes, that is true. It would not be appropriate to use that kind of algorithm on a sparse matrix, just as you don't use the usual sequential triplynested loop.

_____ processing on a data-parallel computer calls for very different approaches. They typically call for the irregular communication and permutation techniques that I illustrated.

Question: What about non-linear programming and algorithms like branch-andbound?

Answer: It is sometimes possible to use data-parallel algorithms to do seemingly unstructured searches, as on a game tree, by maintaining a work queue, like you might do in a more control-parallel, and at every step, taking a large number of task items off the queue by using an enumeration step and using the results of that enumeration to assign them to the processors.

This may depend on whether the rest of the work to be done is sufficiently similar. If it's not, then control parallelism may be more appropriate.

Question: With the current software expertise in 4GLs for sequential machines, do you think that developing data-parallel programming languages will end up at least at 4GL level?

Answer: I think we are now at the point where we know how to design dataparallel languages at about the level of expressiveness as C, Fortran, and possibly Lisp. I think it will take awhile before we can raise our level of understanding to the level we need to design 4GLs.