

Lecture 23

Introduction to Parallel Processing Hardware

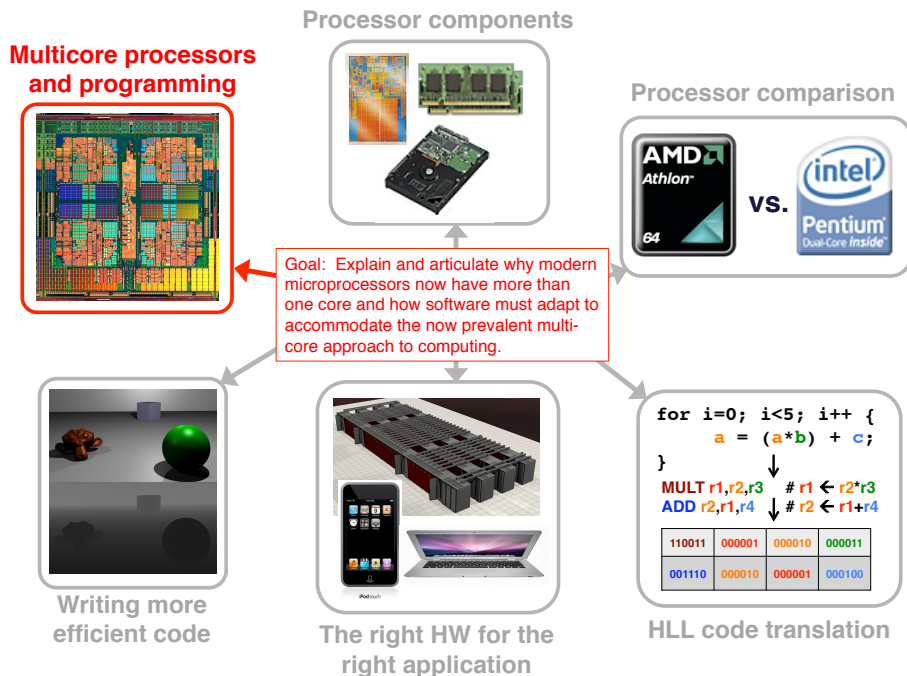
Suggested reading:
(HP Chapter 7 – over the next few weeks)

Important Dates

- Tuesday 11/15 Lab 5 out
- Thursday 11/17 Project summary due (email)
- Thursday 11/17 Lab 4 due
- Tuesday 11/22 HW 7 due
- Tuesday 11/22 HW 8 out
- Thursday 12/1 Lab 5 due
- Thursday 12/6 HW 8 due
- Thursday 12/8 Last day of class
- Monday 12/13 Final Exam (10:30-12:30)
- Thursday 12/15 Last day to turn in project
- Review session TBD

1

2



3

4

Fundamental lesson(s)

- This lecture will highlight:
 1. Different types of parallel systems
 2. What things count as "non-parallelizable" and (significantly) degrade parallel performance
 3. That parallel systems that once existed at the "room-level" now exist on-chip
 - Thus, "room-level" issues must now be dealt with "on-chip".

Why it's important...

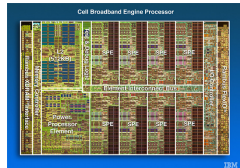
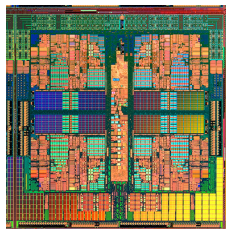
- It's a multi-core world out there...
- Understanding what degrades performance is **ESSENTIAL** for writing efficient software.

General context: Multiprocessors

- Multiprocessor is any computer with several processors



Lemieux cluster, Pittsburgh supercomputing center



Important idea:

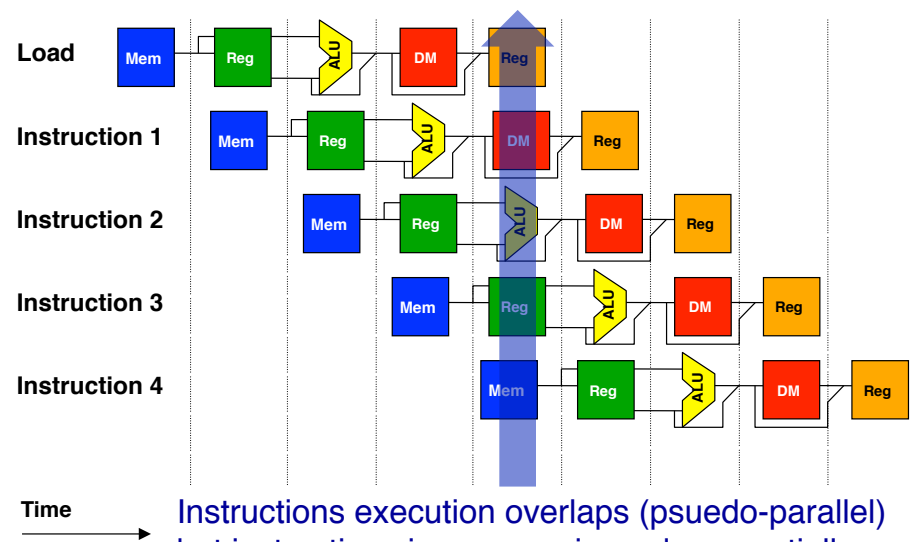
- For a long time, “parallel processing” = connecting multiple chips
- It still does, but now done on single chip too; challenges of connecting multiple chips now exist “on-chip”

INTRODUCTION

5

6

Pipelining and “Parallelism”



Important idea:

We've already done some parallel processing!

7

8

What comes after pipelining?

Dynamic Scheduling: Motivation

	1	2	3	4	5	6	7	8	9	10
<code>divf f0, f2, f4</code>	F	D	E/	E/	E/	E/	W			
<code>addf f6, f0, f2</code>		F	D	d*	d*	d*	E+	E+	W	
<code>mulf f8, f2, f4</code>			F	p*	p*	p*	D	E*	E*	W

- cycle4: `addf` stalls due to **RAW hazard**
 - OK, fundamental problem
- also cycle4: `mulf` stalls due to **pipeline hazard** (`addf` stalls)
 - why? `mulf` can't proceed into ID because `addf` is there
 - but that's the only reason \Rightarrow not good enough!
- why can't we decode `mulf` in cycle 4 and execute it in c5?
 - no fundamental reason why we can't do this!

© 2007 by Sorin, Roth, Hill, Wood, Sohi, Smith, Vijaykumar, Lipasti

ECE 252 / CPS 220 Lecture Notes
Dynamic Scheduling I

9

What comes after pipelining?

Dynamic Scheduling

dynamic scheduling (out-of-order execution)

- execute instructions in non-sequential (non-vonNeumann) order
 - + reduce stalls
 - + improve functional unit utilization
 - + enable parallel execution (not in-order \Rightarrow can be in parallel)
- make it *appear* like sequential execution: precise interrupts
 - very important
 - but hard

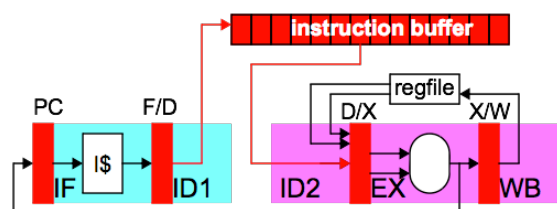
© 2007 by Sorin, Roth, Hill, Wood, Sohi, Smith, Vijaykumar, Lipasti

ECE 252 / CPS 220 Lecture Notes
Dynamic Scheduling I

10

Superscalar machines

Instruction Buffer



trick: **instruction buffer** (many names for this buffer)

- basically: a bunch of latches for holding instructions
 - this is the scope of instructions that the scheduler can see
- split ID into two pieces
 - accumulate decoded instructions in buffer **in-order**
 - buffer sends instructions down rest of pipe **out-of-order**

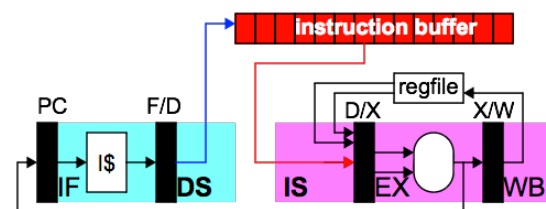
© 2007 by Sorin, Roth, Hill, Wood, Sohi, Smith, Vijaykumar, Lipasti

ECE 252 / CPS 220 Lecture Notes
Dynamic Scheduling I

11

How Superscalar Machines Work

Dispatch and Issue



- **dispatch (DS)**: first part of ID
 - allocate resources in instruction buffer
 - new kind of structural hazard (instruction buffer could be full)
 - *dispatch is in-order, and stall propagates to younger instructions*
- **issue (IS)**: second part of ID
 - send instructions from instruction buffer to execution units
 - *out-of-order, wait does NOT propagate to younger instructions*

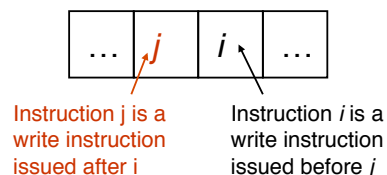
© 2007 by Sorin, Roth, Hill, Wood, Sohi, Smith, Vijaykumar, Lipasti

ECE 252 / CPS 220 Lecture Notes
Dynamic Scheduling I

12

Superscalar Introduces New Hazards

- With **WAW hazard**, instruction *j* tries to write an operand before instruction *i* writes it.
- The writes are performed in wrong order leaving the value written by earlier instruction
- Graphically/Example:



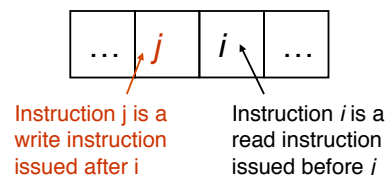
i: DIV F1, F2, F3
j: SUB F1, F4, F6

Problem: If *j* after *i*, and *j* finishes first, and then *i* writes to F1, the wrong value is in F1.

13

Superscalar Introduces New Hazards

- With **WAR hazard**, instruction *j* tries to write an operand before instruction *i* reads it.
- Instruction *i* would incorrectly receive newer value of its operand;
 - Instead of getting old value, it could receive some newer, undesired value:
- Graphically/Example:



i: DIV F7, F1, F3
j: SUB F1, F4, F6

Problem: what if instruction *j* completes before instruction *i* reads F1?

14

Solution: Register Renaming

Freeing Registers in R10K

raw instruction	map table	free locations	renamed instruction
	r1 r2 r3		
add r1, r2, r3	11 12 13	14, 15, 16, 17	add 14, 12, 13
sub r3, r2, r1	14 12 15	16, 17	sub 15, 12, 14
mul r1, r2, r3	16 12 15	17	mul 16, 12, 15
div r2, r1, r3	16 17 15		div 17, 16, 15

- when **add** commits: free 11
- when **sub** commits: free 13
- when **mul** commits: free ?
- when **div** commits: free ?
- see the pattern?

Refer to temporary variable.

Board Example: Superscalar Trace



CLASSIFYING PARALLEL MACHINES

Flynn's Taxonomy

- **MISD: Multiple I, Single D Stream**
 - Not used much, no example today
- **MIMD: Multiple I, Multiple D Streams**
 - Each processor executes its own instructions and operates on its own data
 - Pre multi-core, typical off-the-shelf multiprocessor (made using a bunch of “normal” processors)
 - Each node could also be superscalar
 - Historical lessons, challenges now apply to multi-core too!
 - Example: Intel Xeon e5345

Multiprocessing (Parallel) Machines

- **Flynn's Taxonomy of Parallel Machines**
 - How many Instruction streams?
 - How many Data streams?
- **SISD: Single I Stream, Single D Stream**
 - A uni-processor
 - Could have psuedo-parallel execution here
 - (These are the techniques we just talked about/reviewed...)
 - Example: Intel Pentium 4
- **SIMD: Single I, Multiple D Streams**
 - Each “processor” works on its own data
 - But all execute the same instructions in lockstep
 - Example: SSE instructions in Intel x86
 - (e.g. for vector addition in graphics processing)

17

18

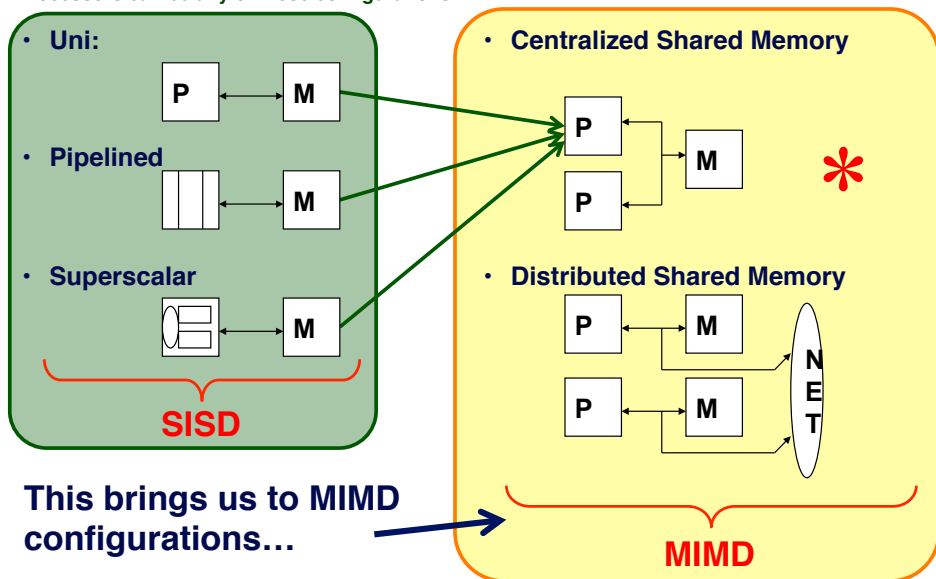
MIMD MACHINE TYPES

19

20

What's Next?

Processors can be any of these configurations



This brings us to MIMD configurations...

21

Multiprocessors

- Why did/do we need multiprocessors?
 - Uni-processor speed improved fast
 - But there are problems that needed even more speed
 - Before: Wait for a few years for Moore's law to catch up?
 - Or use multiple processors and do it sooner?
 - Now: Moore's Law still catching up.
- Multiprocessor software problem
 - Most code is sequential (for uni-processors)
 - MUCH easier to write and debug
 - Correct parallel code very, very difficult to write
 - Efficient and correct is much more difficult
 - Debugging even more difficult

Let's look more at example MIMD configurations...

22

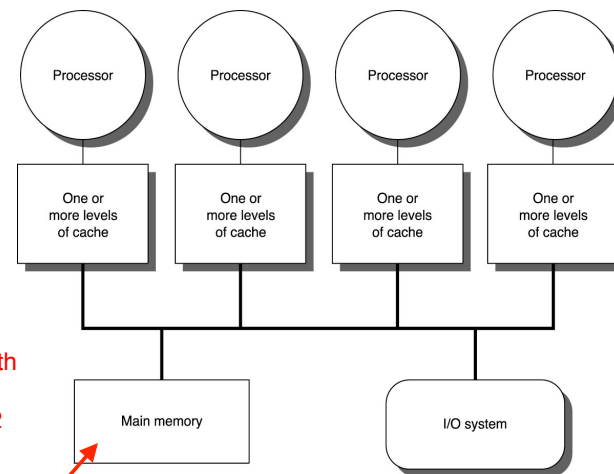
Multiprocessor Memory Types

- Shared Memory:
 - In this model, there is one (large) common shared memory for all processors
- Distributed memory:
 - In this model, each processor has its own (small) local memory, and its content is not replicated anywhere else

But processors can (and do) cache data from memory (wherever that may be).

MIMD Multiprocessors

Centralized Shared Memory



Actually, quite representative of quad core chip, with "main memory" being a shared, L2 cache.

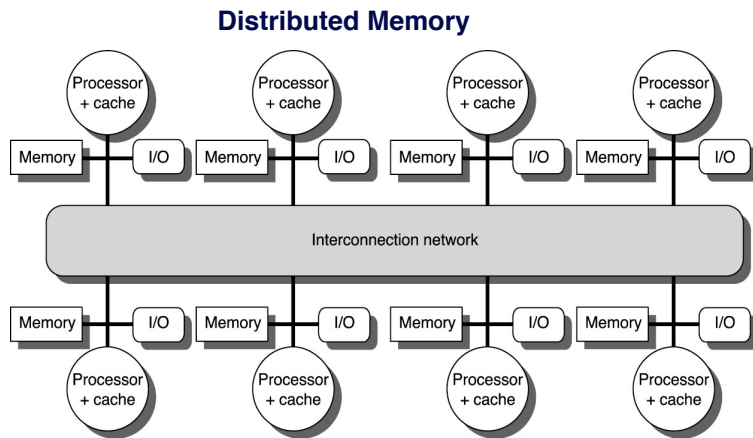
Note: just 1 memory

© 2003 Elsevier Science (USA). All rights reserved.

23

24

MIMD Multiprocessors



Multiple, distributed memories here.

© 2003 Elsevier Science (USA). All rights reserved.

25

More on Centralized-Memory Machines

- “Uniform Memory Access” (UMA)
 - All memory locations have similar latencies
 - Data sharing through memory reads/writes
 - P1 can write data to a physical address A, P2 can then read physical address A to get that data
- Problem: Memory Contention
 - All processor share the one memory
 - Memory bandwidth becomes bottleneck
 - Used only for smaller machines
 - Most often 2,4, or 8 processors
 - Up to 16-32 processors

26

More on Distributed-Memory Machines

- Two kinds
 - Distributed Shared-Memory (DSM)
 - All processors can address all memory locations
 - Data sharing possible
 - Also called NUMA (non-uniform memory access)
 - Latencies of different memory locations can differ
 - (local access faster than remote access)
 - Message-Passing
 - A processor can directly address only local memory
 - To communicate with other processors, must explicitly send/receive messages
- Most accesses local, so less memory contention (can scale to well over 1000s of processors)

27

Message-Passing Machines

- A cluster of computers
 - Each with its own processor and memory
 - An interconnect to pass messages between them
 - Producer-Consumer Scenario:
 - P1 produces data D, uses a SEND to send it to P2
 - The network routes the message to P2
 - P2 then calls a RECEIVE to get the message
 - Two types of send primitives
 - Synchronous: P1 stops until P2 confirms receipt of message
 - Asynchronous: P1 sends its message and continues
 - Standard libraries for message passing:
Most common is MPI – Message Passing Interface

28

Message Passing Pros and Cons

- Pros
 - Simpler and cheaper hardware
 - Explicit communication makes programmers aware of costly (communication) operations
- Cons
 - Explicit communication is painful to program
 - Requires manual optimization
 - If you want a variable to be local and accessible via LD/ST, you must declare it as such
 - If other processes need to read or write this variable, you must explicitly code the needed sends and receives to do this

29

Message Passing: A Program

- Calculating the sum of array elements

```

#define ASIZE 1024
#define NUMPROC 4
double myArray[ASIZE/NUMPROC];
double mySum=0;
for(int i=0;i<ASIZE/NUMPROC;i++)
    mySum+=myArray[i];
if(myPID=0){
    for(int p=1;p<NUMPROC;p++){
        int pSum;
        recv(p,pSum);
        mySum+=pSum;
    }
    printf("Sum: %lf\n",mySum);
}else
    send(0,mySum);
    
```

Must manually split the array

“Main” processor adds up partial sums and prints the result

Other processors send their partial results to main

30

Shared Memory Pros and Cons

- Pros
 - Communication happens automatically
 - More natural way of programming
 - Easier to write correct programs and gradually optimize them
 - No need to manually distribute data (but can help if you do)
- Cons
 - Needs more hardware support
 - Easy to write correct, but inefficient programs (remote accesses look the same as local ones)

More on this in the next few days.
(i.e. where inefficiencies come from...)

31

Shared Memory: A Program

- Calculating the sum of array elements

```

#define ASIZE 1024
#define NUMPROC 4
shared double array[ASIZE];
shared double allSum=0;
shared mutex sumLock;
double mySum=0;
for(int i=myPID*ASIZE/NUMPROC;i<(myPID+1)*ASIZE/NUMPROC;i++)
    mySum+=array[i];
lock(sumLock);
allSum+=mySum;
unlock(sumLock);
if(myPID=0)
    printf("Sum: %lf\n",allSum);
    
```

Array is shared

Each processor sums up “its” part of the array

$i < (2 * 256); i < 512$

Each processor adds its partial sums to the final result

Note need for synchronization support

Main processor prints the result

32

Real Issues and Problems

- Motivation is 2 fold
 - 50s ,60s, 70s, ... today – have rooms filled with machines working to solve a common problem
 - Could consider multi-core chips as “CSM”
 - i.e. 4-core chip, each core has L1 cache, share L2 cache
- Problems are often the same, just on different scales
- As technology scales, problems once seen at room level can be found “on chip”

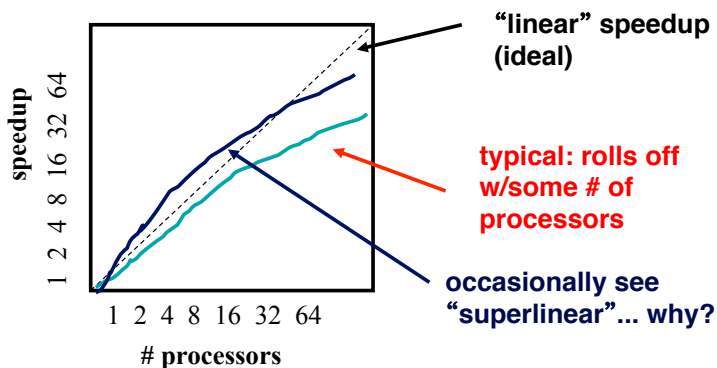
THE CHALLENGE(S) OF PARALLEL COMPUTING

33

Speedup

metric for performance on latency-sensitive applications

- $\text{Time}(1) / \text{Time}(P)$ for P processors
 - note: must use the best *sequential* algorithm for $\text{Time}(1)$; the parallel algorithm may be different.



34

Parallel Programming

- Parallel software is the problem
- Need to get significant performance improvement
 - Otherwise, just use a faster uniprocessor, since it’s easier!
- Difficulties
 - Partitioning
 - Coordination
 - Communications overhead



Challenges: Cache Coherency

- Shared memory easy with no caches
 - P1 writes, P2 can read
 - Only one copy of data exists (in memory)
- Caches store their own copies of the data
 - Those copies can easily get inconsistent
 - Classical example: adding to a sum
 - P1 loads allSum, adds its mySum, stores new allSum
 - P1's cache now has dirty data, but memory not updated
 - P2 loads allSum from memory, adds its mySum, stores allSum
 - P2's cache also has dirty data
 - Eventually P1 and P2's cached data will go to memory
 - Regardless of write-back order, final value ends up wrong

If moderate # of nodes, write-through not practical.

37

Challenges: Latency

- ...is already a major source of performance degradation
 - Architecture charged with hiding local latency
 - (that's why we talked about registers, caches, IC, etc.)
 - Hiding global latency is task of programmer
 - (i.e. manual resource allocation)
- Today:
 - multiple clock cycles to cross chip
 - access to DRAM in 100s of CCs
 - round trip remote access in 1000s of CCs
- In spite of progress in NW technology, *increases* in clock rate may cause delays to reach 1,000,000s of CCs in worst cases

39

Challenges: Contention

- Contention for access to shared resources - esp. memory banks or remote elements - may dominate overall system scalability
 - The problem:
 - Neither network technology nor chip memory bandwidth has grown at the same rate as the processor execution rate or data access demands

38

Challenges: Reliability

- Reliability:
 - Improving yield and achieving high “up time”
 - (think about how performance might suffer if one of the 1 million nodes fails every x number of seconds or minutes...)
 - Solve with checkpointing, other techniques



Part D

Challenges: Languages

- Programming languages, environments, & methodologies:
 - Need simple semantics and syntax that can also expose computational properties to be exploited by large-scale architectures

41

Challenges: Latency ★

- ...is already a major source of performance degradation
 - Architecture charged with hiding local latency
 - (that's why we talked about registers & caches)
 - Hiding global latency is also task of programmer
 - (i.e. manual resource allocation)

- ★ Overhead where no actual processing is done.

Impediments to Parallel Performance

- ★ Reliability:
 - Want to achieve high “up time” – especially in non-CMPs
- ★ Contention for access to shared resources
 - i.e. multiple accesses to limited # of memory banks may dominate system scalability
- Programming languages, environments, & methods:
 - Need simple semantics that can expose computational properties to be exploited by large-scale architectures
- Algorithms
 - Not all problems are parallelizable } $\text{Speedup} = \frac{1}{[1 - \text{Fraction}_{\text{parallelizable}}] + \frac{\text{Fraction}_{\text{parallelizable}}}{N}}$
- ★ Cache coherency
 - P1 writes, P2 can read
 - Protocols can enable \$ coherency but add overhead
- ★ Overhead where no actual processing is done.

What if you write good code for 4-core chip and then get an 8-core chip?

42

Impediments to Parallel Performance

- All ★'ed items also affect speedup that could be obtained...

$$\text{Speedup} = \frac{1}{[1 - \text{Fraction}_{\text{parallelizable}}] + \frac{\text{Fraction}_{\text{parallelizable}}}{N}}$$

44