

Lec 25: Parallel Processors

Kavita Bala
CS 3410, Fall 2008
Computer Science
Cornell University

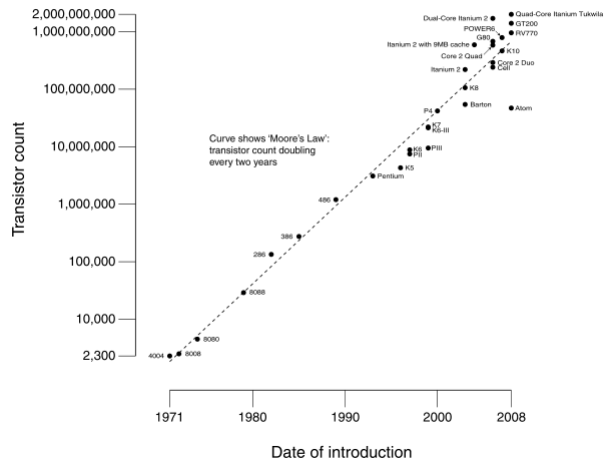
Announcements

- PA 3 out
 - Hack ‘n Seek
 - The goal is to have fun with it
 - Recitations today will talk about it
- Pizza party: Dec 2
- Final project (distributed ray tracer) out last week
 - Demos: Dec 16 and 17

Moore's Law

- Law about transistor count

CPU Transistor Counts 1971-2008 & Moore's Law



Parallel Processing

- Been around for long time in different forms
- Instruction Level Parallelism
- Data Parallelism
- Etc.

Instruction-Level Parallelism (ILP)

- Pipelining: executing multiple instructions in parallel
- To increase ILP
 - Deeper pipeline
 - Less work per stage \Rightarrow shorter clock cycle
 - Multiple issue
 - Replicate pipeline stages \Rightarrow multiple pipelines
 - Start multiple instructions per clock cycle
 - $CPI < 1$, so use Instructions Per Cycle (IPC)
 - E.g., 4GHz 4-way multiple-issue
 - 16 BIPS, peak $CPI = 0.25$, peak $IPC = 4$
 - But dependencies reduce this in practice

© Kavita Bala, Computer Science, Cornell University

Multiple Issue

- Static multiple issue
 - Compiler groups instructions to be issued together
 - Packages them into “issue slots”
 - Specifies multiple concurrent operations
 - \Rightarrow Very Long Instruction Word (VLIW)
 - Compiler detects and avoids hazards
- Dynamic multiple issue
 - CPU examines instruction stream and chooses instructions to issue each cycle
 - Compiler can help by reordering instructions
 - CPU resolves hazards using advanced techniques at runtime

© Kavita Bala, Computer Science, Cornell University

Speculation

- “Guess” what to do with an instruction
 - Start operation as soon as possible
 - Check whether guess was right
 - If so, complete the operation
 - If not, roll-back and do the right thing
- Common to static and dynamic multiple issue
- Examples
 - Speculate on branch outcome
 - Roll back if path taken is different
 - Speculate on load
 - Roll back if location is updated

© Kavita Bala, Computer Science, Cornell University

Compiler/Hardware Speculation

- Compiler can reorder instructions
 - e.g., move load before branch
 - Can include “fix-up” instructions to recover from incorrect guess
- Hardware can look ahead for instructions to execute
 - Buffer results until it determines they are actually needed
 - Flush buffers on incorrect speculation

© Kavita Bala, Computer Science, Cornell University

Issues

- Hazards
- Exceptions

© Kavita Bala, Computer Science, Cornell University

MIPS with Static Dual Issue

- Two-issue packets
 - One ALU/branch instruction
 - One load/store instruction
 - 64-bit aligned
 - ALU/branch, then load/store
 - Pad an unused instruction with nop

Address	Instruction type	Pipeline Stages						
		IF	ID	EX	MEM	WB		
n	ALU/branch	IF	ID	EX	MEM	WB		
n + 4	Load/store	IF	ID	EX	MEM	WB		
n + 8	ALU/branch		IF	ID	EX	MEM	WB	
n + 12	Load/store		IF	ID	EX	MEM	WB	
n + 16	ALU/branch			IF	ID	EX	MEM	WB
n + 20	Load/store			IF	ID	EX	MEM	WB

© Kavita Bala, Computer Science, Cornell University

Scheduling Example

- Schedule this for dual-issue MIPS

```
Loop: lw   $t0, 0($s1)    # $t0=array element
      addu $t0, $t0, $s2  # add scalar in $s2
      sw   $t0, 0($s1)    # store result
      addi $s1, $s1, -4   # decrement pointer
      bne $s1, $zero, Loop # branch $s1!=0
```

	ALU/branch	Load/store	cycle
Loop:	nop	lw \$t0, 0(\$s1)	1
	addi \$s1, \$s1, -4	nop	2
	addu \$t0, \$t0, \$s2	nop	3
	bne \$s1, \$zero, Loop	sw \$t0, 4(\$s1)	4

– $IPC = 5/4 = 1.25$ (c.f. peak $IPC = 2$)

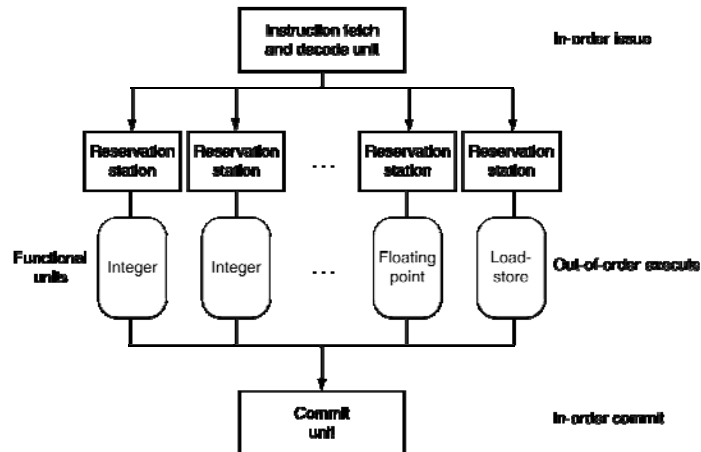
© Kavita Bala, Computer Science, Cornell University

Dynamic Multiple Issue

- “Superscalar” processors
- CPU decides whether to issue 0, 1, 2, ... each cycle
 - Avoiding structural and data hazards
- Avoids the need for compiler scheduling
 - Though it may still help
 - Code semantics ensured by the CPU

© Kavita Bala, Computer Science, Cornell University

Dynamic Multiple Issue



© Kavita Bala, Computer Science, Cornell University

Dynamic Pipeline Scheduling

- Allow the CPU to execute instructions out of order to avoid stalls
 - But commit result to registers in order
- Example

```
lw    $t0, 20($s2)
addu  $t1, $t0, $t2
sub   $s4, $s4, $t3
slli  $t5, $s4, 20
```

 - Can start sub while addu is waiting for lw

© Kavita Bala, Computer Science, Cornell University

Why Do Dynamic Scheduling?

- Why not just let the compiler schedule code?
- Not all stalls are predictable
 - e.g., cache misses
- Can't always schedule around branches
 - Branch outcome is dynamically determined
- Different implementations of an ISA have different latencies and hazards

© Kavita Bala, Computer Science, Cornell University

Does Multiple Issue Work?

- Yes, kind of
 - But not as much as we'd like
- Programs have real dependencies that limit ILP
- Some dependencies are hard to eliminate
 - e.g., pointer aliasing
- Some parallelism is hard to expose
 - Limited window size during instruction issue
- Memory delays and limited bandwidth
 - Hard to keep pipelines full
- Speculation can help if done well

© Kavita Bala, Computer Science, Cornell University

Power Efficiency

- Complexity of dynamic scheduling and speculations requires power
- Multiple simpler cores may be better

Microprocessor	Year	Clock Rate	Pipeline Stages	Issue width	Out-of-order/Speculation	Cores	Power
i486	1989	25MHz	5	1	No	1	5W
Pentium	1993	66MHz	5	2	No	1	10W
Pentium Pro	1997	200MHz	10	3	Yes	1	29W
P4 Willamette	2001	2000MHz	22	3	Yes	1	75W
P4 Prescott	2004	3600MHz	31	3	Yes	1	103W
Core	2006	2930MHz	14	4	Yes	2	75W
UltraSparc III	2003	1950MHz	14	4	No	1	90W
UltraSparc T1	2005	1200MHz	6	1	No	8	70W

© Kavita Bala, Computer Science, Cornell University

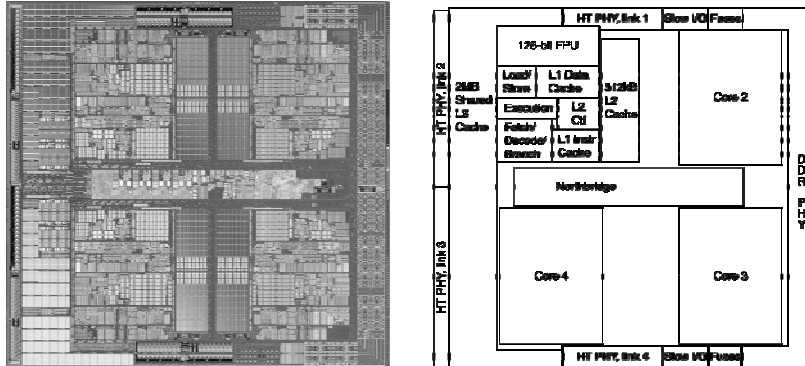
Why Multicore?

- Moore's law
 - A law about transistors
 - Smaller means faster transistors
- Power consumption growing with transistors
- The power wall
 - We can't reduce voltage further
 - We can't remove more heat
- How else can we improve performance?

© Kavita Bala, Computer Science, Cornell University

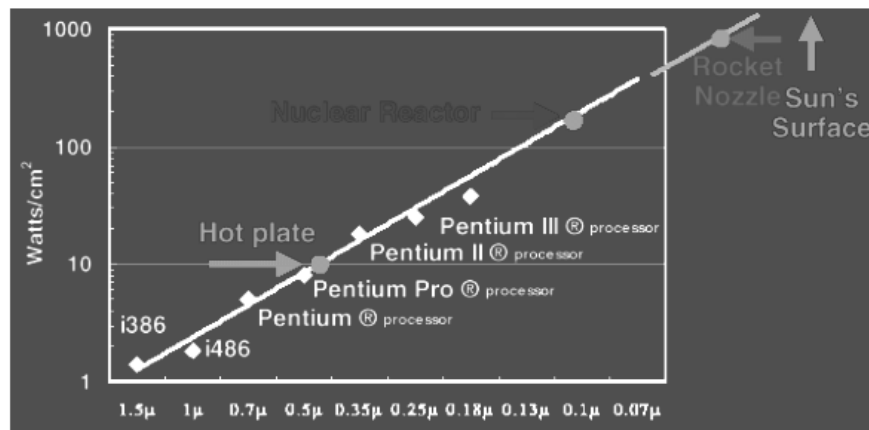
Inside the Processor

- AMD Barcelona: 4 processor cores



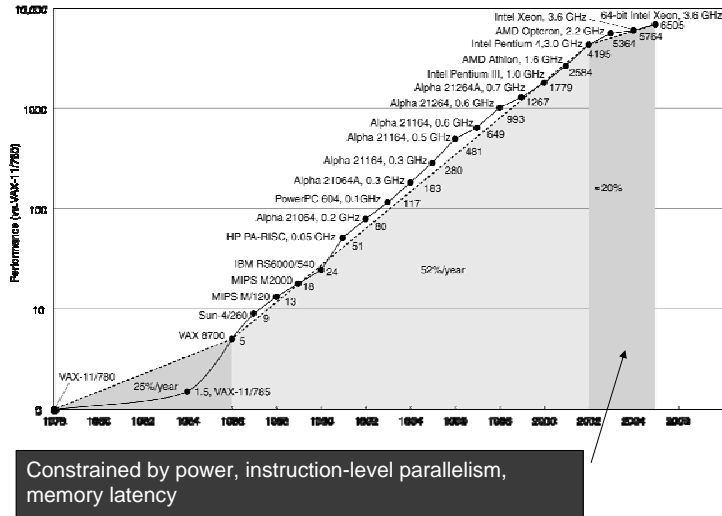
© Kavita Bala, Computer Science, Cornell University

Power Limits Performance



© Kavita Bala, Computer Science, Cornell University

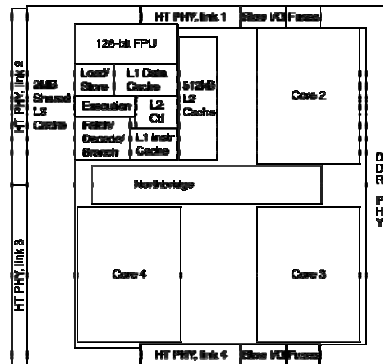
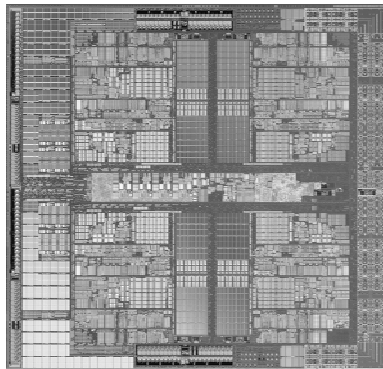
Uniprocessor Performance



© Kavita Bala, Computer Science, Cornell University

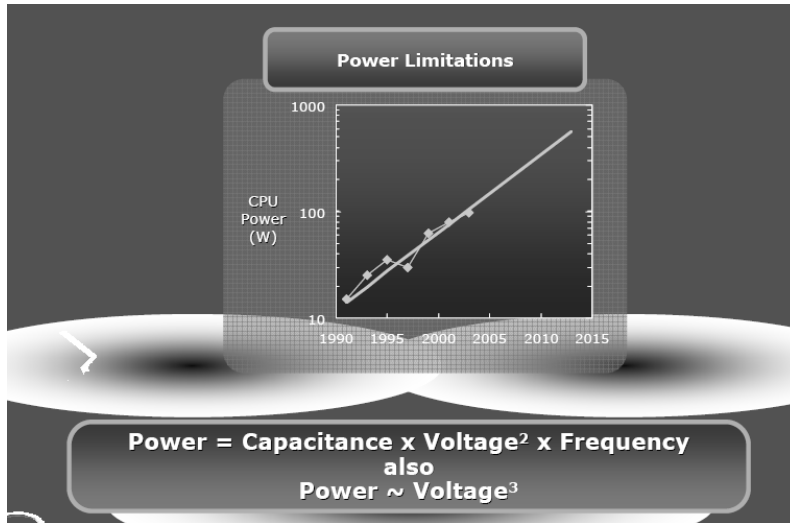
Inside the Processor

- AMD Barcelona: 4 processor cores



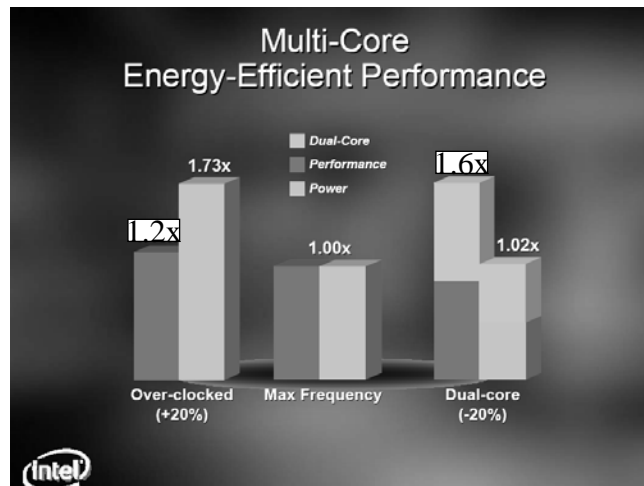
© Kavita Bala, Computer Science, Cornell University

Intel's argument



© Kavita Bala, Computer Science, Cornell University

Multi-Core Energy-Efficient Performance



© Kavita Bala, Computer Science, Cornell University

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
 - Balance load

© Kavita Bala, Computer Science, Cornell University

Load Balancing

- Need to manage work so all units are actually operating



© Kavita Bala, Computer Science, Cornell University

Amdahl's Law

- Task: serial part, parallel part
- As number of processors increases,
 - time to execute parallel part goes to zero
 - time to execute serial part remains the same
- *Serial part eventually dominates*
- Must parallelize ALL parts of task

$$\text{Speedup}(E) = \frac{\text{Execution Time without } E}{\text{Execution Time with } E}$$

© Kavita Bala, Computer Science, Cornell University

Amdahl's Law

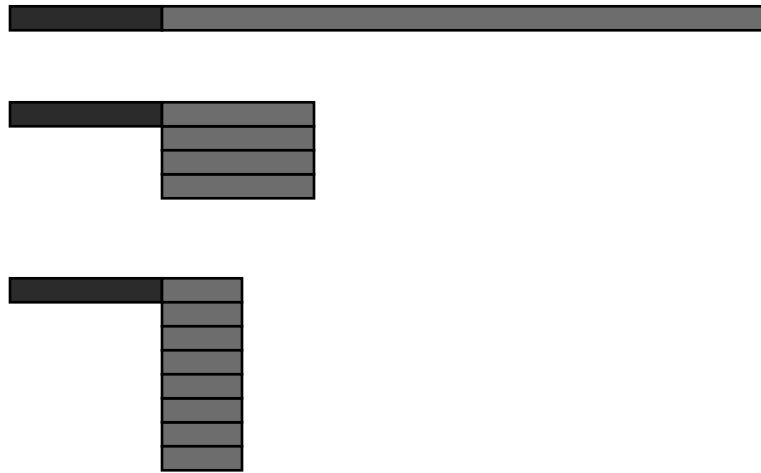
- Consider an improvement E
- F of the execution time is affected
- S is the speedup

Execution time (with E) = $((1 - F) + F/S) \cdot$ Execution time (without E)

$$\text{Speedup (with } E) = \frac{1}{(1 - F) + F/S}$$

© Kavita Bala, Computer Science, Cornell University

Amdahl's Law



© Kavita Bala, Computer Science, Cornell University

Pitfall: Amdahl's Law

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

$$T_{\text{improved}} = \frac{T_{\text{affected}}}{\text{improvement factor}} + T_{\text{unaffected}}$$

- Example: multiply accounts for 80s/100s
 - How much improvement in multiply performance to get 5x overall?

$$20 = \frac{80}{n} + 20 \quad \text{– Can't be done!}$$

© Kavita Bala, Computer Science, Cornell University

Amdahl's Law

- Sequential part can limit speedup
- Example: 100 processors, 90x speedup?
 - $T_{\text{new}} = T_{\text{parallelizable}}/100 + T_{\text{sequential}}$
 - $\text{Speedup} = \frac{1}{(1 - F_{\text{parallelizable}}) + F_{\text{parallelizable}}/100} = 90$
 - Solving: $F_{\text{parallelizable}} = 0.999$
- Need sequential part to be 0.1% of original time

© Kavita Bala, Computer Science, Cornell University

Scaling Example

- Workload: sum of 10 scalars, and 10×10 matrix sum
 - Speed up from 10 to 100 processors
- Single processor: Time = $(10 + 100) \times t_{\text{add}}$
- 10 processors
 - Time = $10 \times t_{\text{add}} + 100/10 \times t_{\text{add}} = 20 \times t_{\text{add}}$
 - Speedup = $110/20 = 5.5$ (55% of potential)
- 100 processors
 - Time = $10 \times t_{\text{add}} + 100/100 \times t_{\text{add}} = 11 \times t_{\text{add}}$
 - Speedup = $110/11 = 10$ (10% of potential)
- Assumes load can be balanced across processors

© Kavita Bala, Computer Science, Cornell University

Scaling Example (cont)

- What if matrix size is 100×100 ?
- Single processor: Time = $(10 + 10000) \times t_{\text{add}}$
- 10 processors
 - Time = $10 \times t_{\text{add}} + 10000/10 \times t_{\text{add}} = 1010 \times t_{\text{add}}$
 - Speedup = $10010/1010 = 9.9$ (99% of potential)
- 100 processors
 - Time = $10 \times t_{\text{add}} + 10000/100 \times t_{\text{add}} = 110 \times t_{\text{add}}$
 - Speedup = $10010/110 = 91$ (91% of potential)
- Assuming load balanced