Multicore & Parallel Processing

J

Kevin Walsh CS 3410, Spring 2010 Computer Science Cornell University

P&H Chapter 4.10-11, 7.1-6

Q: How to improve system performance?

- \rightarrow Increase CPU clock rate?
 - \rightarrow But I/O speeds are limited
 - Disk, Memory, Networks, etc.

Recall: Amdahl's Law

Solution: Parallelism

Pipelining: execute multiple instructions in parallel

- Q: How to get more instruction level parallelism?
- A: Deeper pipeline
 - Less work per stage \Rightarrow shorter clock cycle
- A: Multiple issue pipeline
 - Start multiple instructions per clock cycle in duplicate stages
 - Example: 1GHz 4-way multiple-issue

Peak CPI = 0.25 \rightarrow 16 billion instructions per second

Static Multiple Issue

a.k.a. Very Long Instruction Word (VLIW)

Compiler groups instructions to be issued together

- Packages them into "issue slots"
- Simple HW: Compiler detects and avoids hazards

Example: Static Dual-Issue MIPS

Instructions come in pairs (64-bit aligned)
– One ALU/branch instruction (or nop)

- One load/store instruction (or nop)

- (J ctra

Compiler scheduling for dual-issue MIPS...



ALU/branch slot

TOP:	nop			
	addu	\$s1,	\$s1,	+8
	addu	\$t0,	\$t0,	\$s2
	addu	\$t1,	\$t1,	\$s2
	bne	\$s1,	\$s3,	TOP

Load	l/store	e slot	cycle
lw \$	5t0, 0	(\$s1)	1
lw	\$t1,	4(\$s1)	2
nop	e		3
SW	\$t0,	-8(\$s1)	4
SW	\$t1,	-4(\$s1)	5

Dynamic Multiple Issue

a.k.a. SuperScalar Processor

- CPU examines instruction stream and chooses multiple instructions to issue each cycle
- Compiler can help by reordering instructions....
- ... but CPU is responsible for resolving hazards
- Even better: Speculation/Out-of-order Execution
 - Guess results of branches, loads, etc.
 - Execute instructions as early as possible
 - Roll back if guesses were wrong
 - Don't commit results until all previous insts. are retired

- Q: Does multiple issue / ILP work?
- A: Kind of... but not as much as we'd like Limiting factors?
 - Programs dependencies
 - Hard to detect dependencies → be conservative
 e.g. Pointer Aliasing: A[0] += 1; B[0] *= 2;
 - Hard to expose parallelism
 - Can only issue a few instructions ahead of PC
 - Structural limits
 - Memory delays and limited bandwidth
 - Hard to keep pipelines full

Q: Does multiple issue / ILP cost much?

A: Yes.

\rightarrow Dynamic issue and speculation requires power

CPU	Year	Clock Rate	Pi S	ipeline tages	lssue 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
UltraSparc III	2003	1950MHz		14	4	No	1	90W
P4 Prescott	2004	3600MHz		31	3	Yes	1	103W
Core	2006	2930MHz		14	4	Yes	2	75W
UltraSparc T1	2005	1200MHz		6	1	No	8	70W

\rightarrow Multiple simpler cores may be better?



Moore's law

- A law about transistors
- Smaller means more transistors per die
- And smaller means faster too

But: Power consumption growing too...

Power Limits Performance



Power = capacitance * voltage² * frequency In practice: Power ~ voltage³ \checkmark \checkmark

Reducing voltage helps (a lot) Sund Sole of So

The power wall

- We can't reduce voltage further
- We can't remove more heat

Why Multicore?



AMD Barcelona: 4 processor cores



Q: So lets just all use multicore from now on!A: Software must be written as parallel program

Multicore difficulties

- Partitioning work
- Coordination & synchronization
- Communications overhead
- Balancing load over cores
- How do you write parallel programs?
 - … without knowing exact underlying architecture?

Partition work so all cores have something to do



Load Balancing

Need to partition so all cores are actually working

		_
		_



If tasks have a serial part and a parallel part... Example:

- step 1: divide input data into *n* pieces
- step 2: do work on each piece
- step 3: combine all results
- Recall: Amdahl's Law
- As number of cores increases ...
 - time to execute parallel part?
 - time to execute serial part?

