

UltraSPARC T2
(Niagara-2)
multicore chip layout

(courtesy of T. Okazaki, Flickr)

- course web site: <http://cs.anu.edu.au/student/comp8320>
- course coordinator & lecturer/tutor:
Peter Strazdins, CSIT N219, 6125-5140, comp8320@cs
- discussion forum accessible by StReAMS:
- course schedule: 10 'modules'

Overview

- course contact
- review: instruction-level parallelism, memory hierarchy, cache memory, shared memory multiprocessor and programming model, OpenMP
- assumed knowledge and assessment
- intermission!
- discussion: your views on multicore
- what is multicore; four multicore designs
- advent of multicore (parallel programming decline and rebirth)
- the multicore challenge: why parallel programming is hard

credits: some material for this lecture is from CS194 course by Prof Kathy Yelick, UCB

Review: Instruction Level Parallelism: Pipelining

- pipelining: achieve fast cycle times by breaking instr'n execution into k stages
 - if complete one instr'n per cycle. get $\leq k$ -way parallelism
 - note: generally, the circuitry for each stage is independent
- e.g. ($k = 5$): stages FI = Fetch Instr'n., DI = Decode Instr'n., FO = Fetch Operand, EX = Execute Instr'n., WB = Write Back

(branch):	FI	DI	FO	EX	WB				
(delay slot):		FI	DI	FO	EX	WB			
(guess)			FI	DI	FO	EX	WB		
(guess)				FI	DI	FO	EX	WB	
(sure)					FI	DI	FO	EX	WB

- note: EX & WB stages may involve memory accesses (and may possibly stall the pipeline)

Review: Instruction Level Parallelism: Superscalar

- multiple instruction issue (superscalar):
 - a small number (w) of instr'n's (a group) are scheduled by the H/W to *execute together*
 - groups must have an appropriate 'instruction mix'
 - e.g. UltraSPARC III ($w = 4$): $\left. \begin{array}{l} \leq 2 \text{ different floating point} \\ \leq 1 \text{ load / store ; } \leq 1 \text{ branch} \\ \leq 2 \text{ integer / logical} \end{array} \right\} \text{instr'n's per group}$
 - have $\leq w$ -way parallelism over *different types of instr'n's*
 - generally requires:
 - ◆ multiple ($\geq w$) instr'n fetches; extra grouping stage (G) in the pipeline
 - **problem:** amplifies all problems with pipelining
 - **issue:** the instruction mix must be balanced for **maximum performance!**
- out-of-order execution
 - hardware dynamically re-orders instructions to obtain better grouping and resolve dependencies (nice, but very complex and hence expensive!)
 - e.g. swap the order of the **fadds** in


```
fmuld %f0,%f1,%f2 ; fadd %f2,%f4,%f4 ; ... ; fadd %f7,%f0,%f1
```

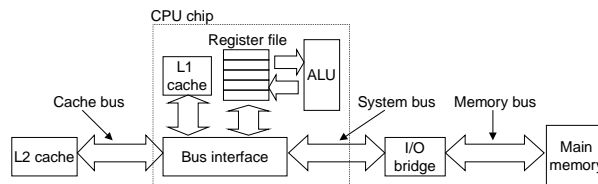
Review: Cache Memory – Details

- idea: data that is “*currently most needed*” is brought into a (smaller) faster memory
- observation: memory accesses in *most* programs exhibit:
 - temporal locality: if access address X , likely to access X again soon
 - spatial locality: if access address X , likely to access $X+1$ soon
 - ⇒ cache organized into lines (blocks) of B bytes ($B = 2^b$, e.g. $4 \leq b \leq 9$)
 - ✓ blocked memory accesses (faster) & less control info needed (per byte)
 - × redundant memory traffic if only ever use 1 byte per line
- if so, yields good cost-speed tradeoff
- **problem:** keeping consistency of data cache & main memory
 - when a store instruction is executed, the relevant line is updated 1st; when does main memory get updated?
 - ◆ 2 strategies: write-through and copy-back
- organization: S -way set associative caches ([O'H.&Bryant, fig 6.33])
 - middle bits of address X is used to select a set of S lines
 - typically $S = 1, 2, 4, 6, 8$ (direct-mapped cache corresponds to $S = 1$)
 - reduces chance of conflicts by factor of S

Review: the Memory Hierarchy

- in a computer system, there is a memory hierarchy, because:
 - many different mediums for the storage of data
 - generally, there is a *trade-off* between speed and capacity
 - ◆ fast memories tend to be small; large memories tend to be slow

medium	access time	typical size
registers	~1 ns	< 1 KB
cache mem.	~20 ns	< 2 MB
main mem.	~100 ns	< 2 GB
disk	~10 ⁷ ns	> 10 GB



[O'H.&Bryant, fig 6.24]

- idea of cache (\$) memory: store recently accessed data from main memory in a small, faster memory (closer to / inside the CPU)
 - Q: why might this be useful?
- issue: must maintain address (tag) of each data item in the caches
 - must be able to efficiently lookup whether data item for address x is in cache

Review: Shared Memory Multiprocessors

- a number of processors connected to a globally addressable memory
 - through a bus ([Lyn&Snyder, fig 2.3]) or, better,
 - memory is organized into modules, all connected by an interconnect
- need caches! memory consistency problem is now exacerbated!
 - 2 strategies: invalidate and update
 - question: how *soon* do changes need to be propagated?
- consider programs (threads) on processors 0 and 1 attempting to acquire a lock
 - hardware must support this via some atomic instructions, e.g.


```

                    ! %o0 has address of the lock
                    mov 0xff, %o1 ! 0xff is value for acquiring lock
                    loop: ! loop to acquire lock
                    brnz %o1, loop ! exit if %o1 = 0 (lock value was 0)
                    ldstub [%o0], %o1 ! atomic swap %o1 and value at lock
                    ... ! safely update shared data structure
                    stb %g0, [%o0] ! release lock
                    
```
 - for locking to work properly, any change to a memory location by 1 processor must be known at the next atomic instruction

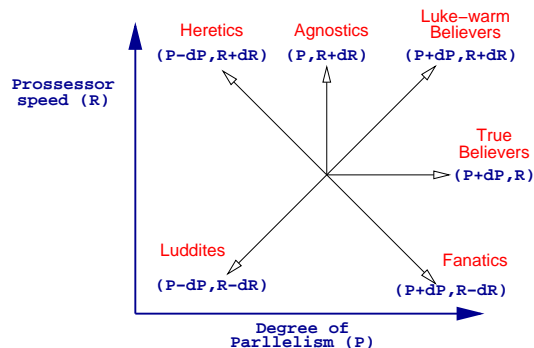
What is Multicore?

- (also known as chip multiprocessing, CMP): multiple shared memory processors ('cores') on a single chip
 - why? because we can! also because we must ...
 - can run multiple applications (processes) in parallel
 - can run a single (threaded) application in parallel; **herein lies the challenge!**
- (large-scale) parallelism is now cheap, mainstream
- some multicore designs:
 - Intel Core 2 Duo ([Hughes&Hughes, fig 2-8]): shared L2\$; coherency needed on separate L1\$s
 - AMD Opteron ([Hughes&Hughes, fig 2-3]): separate L2\$; coherency needed on L2\$s (both within and across chips ('sockets'))
 - UltraSPARC T1 ([Hughes&Hughes, fig 2-5]): note crossbar for the L2\$ and multiple memory controllers; also uses hardware threading (why?)
 - Cell Broadband Engine ([Hughes&Hughes, fig 2-6]): heterogeneous cores (SPEs) of limited capability \Rightarrow greater programming issues!

Memory (data access) is a key consideration!

(Pre-) Advent: Moore's Law undermines Parallel Computing

- looked promising in the 90's, but many companies failed due to Moore's Law: # transistors on a chip doubles every 18 months
 - for a long time, this also permitted an exponential increase in clock speed ...
 - *On several recent occasions, I have been asked whether parallel computing will soon be relegated to the trash heap ... Ken Kennedy, CRPC Directory, 1994*
- demography of Parallel Computing (mid 90's, origin unknown)

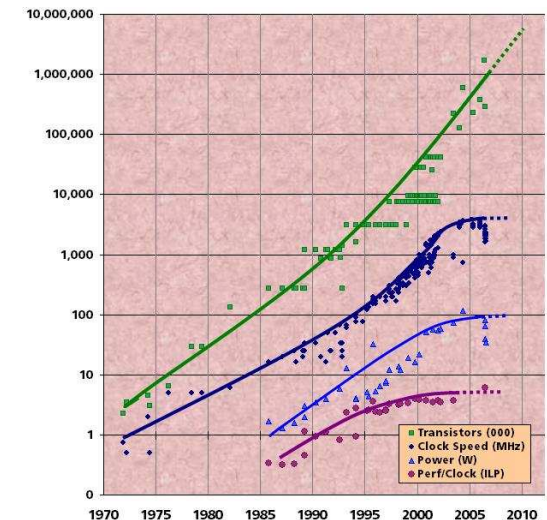


Advent: End of Clock Scaling and Instruction-level Parallelism

- extrapolation of exponential power density increase 1985–2000 indicates we are at the limit!
 - 2000 Intel chip equivalent to a hotplate, would have \Rightarrow a rocket nozzle by 2010!
- dissipated power is given by: $P \propto CV^2f \propto Cf^3$, V is the voltage, C is the chip's capacitance and f is the clock frequency
- (*ideal*) parallel performance for p cores is given by $R = pf$, but $C \propto p$
- double $p \Rightarrow$ double R , but also P
- double p , halve $f \Rightarrow$ maintain R , but quarter P
- instruction level parallelism has reached its limits
 - instructions per cycle (IPC) of a 6-way superscalar not much better than a 2-way (Olukotun et al, ASPLOS, 1996)
 - by utilizing technique of hardware threading, can hide effects of cache misses
 - ◆ for each core, have a number of 'virtual CPUs', each with own register set
 - ◆ the core's control unit can flexibly select instructions ready to execute in each of these

Advent: The Multicore Revolution

- transistor density is expected to increase (till 2015)
- instead of doubling clock speed, double the number of cores
- cores are much simpler (can utilize less ILP), but there can be many
- chip design and testing costs are significantly reduced
- parallelism still must be exposed by software



(source: CS194 lecture 1 \Rightarrow)

