

# CSE 490/590 Computer Architecture

## Multithreading I

Steve Ko  
Computer Sciences and Engineering  
University at Buffalo

CSE 490/590, Spring 2011

### Last time...

- Superscalar suffers from the sequential nature of the ISA
- VLIW instructions consists of multiple operations
- Techniques such as loop unrolling, software pipelining, and trace scheduling gives an opportunity to extract ILP necessary in VLIW

CSE 490/590, Spring 2011

2

### Rotating Register Files

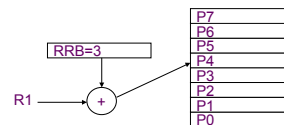
Problems: Scheduled loops require lots of registers,  
Lots of duplicated code in prolog, epilog

Solution: Allocate new set of registers for each loop iteration

CSE 490/590, Spring 2011

3

### Rotating Register File



Rotating Register Base (RRB) register points to base of current register set. Value added on to logical register specifier to give physical register number. Usually, split into rotating and non-rotating registers.

CSE 490/590, Spring 2011

4

### Rotating Register File (Previous Loop Example)

Three cycle load latency encoded as difference of 3 in register specifier number ( $f4 - f1 = 3$ )

Four cycle fadd latency encoded as difference of 4 in register specifier number ( $f9 - f5 = 4$ )

ld f1, ()	fadd f5, f4, ...	sd f9, ()	bloop	
ld P9, ()	fadd P13, P12,	sd P17, ()	bloop	RRB=8
ld P8, ()	fadd P12, P11,	sd P16, ()	bloop	RRB=7
ld P7, ()	fadd P11, P10,	sd P15, ()	bloop	RRB=6
ld P6, ()	fadd P10, P9,	sd P14, ()	bloop	RRB=5
ld P5, ()	fadd P9, P8,	sd P13, ()	bloop	RRB=4
ld P4, ()	fadd P8, P7,	sd P12, ()	bloop	RRB=3
ld P3, ()	fadd P7, P6,	sd P11, ()	bloop	RRB=2
ld P2, ()	fadd P6, P5,	sd P10, ()	bloop	RRB=1

CSE 490/590, Spring 2011

5

### Cydra-5: Memory Latency Register (MLR)

Problem: Loads have variable latency  
Solution: Let software choose desired memory latency

- Compiler schedules code for maximum load-use distance
- Software sets MLR to latency that matches code schedule
- Hardware ensures that loads take exactly MLR cycles to return values into processor pipeline
  - Hardware buffers loads that return early
  - Hardware stalls processor if loads return late

CSE 490/590, Spring 2011

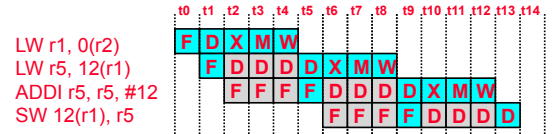
6

## Multithreading

- Difficult to continue to extract instruction-level parallelism (ILP) or data-level parallelism (DLP) from a single sequential thread of control
- Many workloads can make use of thread-level parallelism (TLP)
  - TLP from multiprogramming (run independent sequential jobs)
  - TLP from multithreaded applications (run one job faster using parallel threads)
- Multithreading uses TLP to improve utilization of a single processor

CSE 490/590, Spring 2011

## Pipeline Hazards



- Each instruction may depend on the next

*What is usually done to cope with this?*

CSE 490/590, Spring 2011

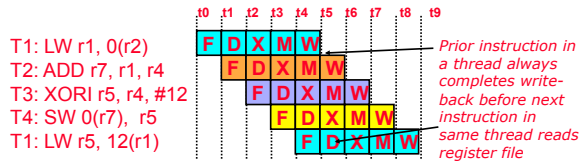
8

## Multithreading

How can we guarantee no dependencies between instructions in a pipeline?

- One way is to interleave execution of instructions from different program threads on same pipeline

*Interleave 4 threads, T1-T4, on non-bypassed 5-stage pipe*



CSE 490/590, Spring 2011

9

## CDC 6600 Peripheral Processors (Cray, 1964)

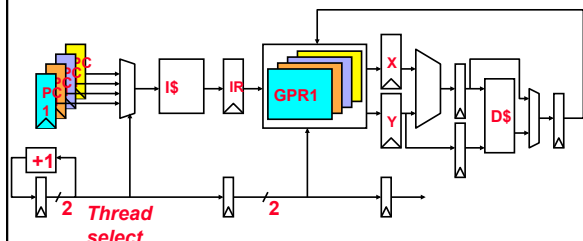


- First multithreaded hardware
- 10 "virtual" I/O processors
- Fixed interleave on simple pipeline
- Pipeline has 100ns cycle time
- Each virtual processor executes one instruction every 1000ns
- Accumulator-based instruction set to reduce processor state

CSE 490/590, Spring 2011

10

## Simple Multithreaded Pipeline



- Have to carry thread select down pipeline to ensure correct state bits read/written at each pipe stage
- Appears to software (including OS) as multiple, albeit slower, CPUs

CSE 490/590, Spring 2011

11

## Multithreading Costs

- Each thread requires its own user state
  - PC
  - GPRs
- Also, needs its own system state
  - virtual memory page table base register
  - exception handling registers
- *Other overheads:*
  - Additional cache/TLB conflicts from competing threads
  - (or add larger cache/TLB capacity)
  - More OS overhead to schedule more threads (where do all these threads come from?)

CSE 490/590, Spring 2011

12

## Thread Scheduling Policies

- Fixed interleave (*CDC 6600 PPU's, 1964*)
  - Each of N threads executes one instruction every N cycles
  - If thread not ready to go in its slot, insert pipeline bubble
- Software-controlled interleave (*TI ASC PPU's, 1971*)
  - OS allocates S pipeline slots amongst N threads
  - Hardware performs fixed interleave over S slots, executing whichever thread is in that slot
- Hardware-controlled thread scheduling (*HEP, 1982*)
  - Hardware keeps track of which threads are ready to go
  - Picks next thread to execute based on hardware priority scheme



CSE 490/590, Spring 2011

13

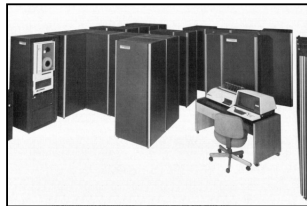
## CSE 490/590 Administrivia

- HW2 & midterm solution out
- Quiz 2 (next Friday 4/8): After midterm until next Monday

CSE 490/590, Spring 2011

14

## Denelcor HEP (Burton Smith, 1982)



- First commercial machine to use hardware threading in main CPU
- 120 threads per processor
  - 10 MHz clock rate
  - Up to 8 processors
  - precursor to Tera MTA (Multithreaded Architecture)

CSE 490/590, Spring 2011

15

## Tera MTA (1990-)

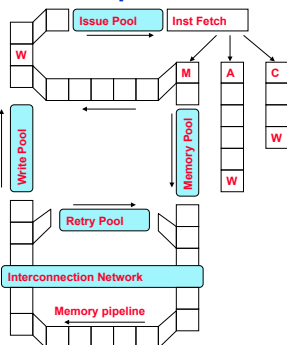
- Up to 256 processors
- Up to 128 active threads per processor
- Processors and memory modules populate a sparse 3D torus interconnection fabric
- Flat, shared main memory
  - No data cache
  - Sustains one main memory access per cycle per processor
- GaAs logic in prototype, 1KW/processor @ 260MHz
  - Second version CMOS, MTA-2, 50W/processor
  - New version, XMT, fits into AMD Opteron socket, runs at 500MHz



CSE 490/590, Spring 2011

16

## MTA Pipeline



- Every cycle, one VLIW instruction from one active thread is launched into pipeline
- Instruction pipeline is 21 cycles long
- Memory operations incur ~150 cycles of latency

CSE 490/590, Spring 2011

17

## Coarse-Grain Multithreading

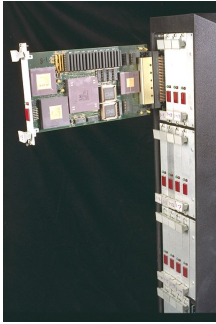
- Tera MTA designed for supercomputing applications with large data sets and low locality
- No data cache
  - Many parallel threads needed to hide large memory latency

- Other applications are more cache friendly
- Few pipeline bubbles if cache mostly has hits
  - Just add a few threads to hide occasional cache miss latencies
  - Swap threads on cache misses

CSE 490/590, Spring 2011

18

## MIT Alewife (1990)



- Modified SPARC chips
  - register windows hold different thread contexts
- Up to four threads per node
- Thread switch on local cache miss

CSE 490/590, Spring 2011

19

## IBM PowerPC RS64-IV (2000)

- Commercial coarse-grain multithreading CPU
- Based on PowerPC with quad-issue in-order five-stage pipeline
- Each physical CPU supports two virtual CPUs
- On L2 cache miss, pipeline is flushed and execution switches to second thread
  - short pipeline minimizes flush penalty (4 cycles), small compared to memory access latency
  - flush pipeline to simplify exception handling

CSE 490/590, Spring 2011

20

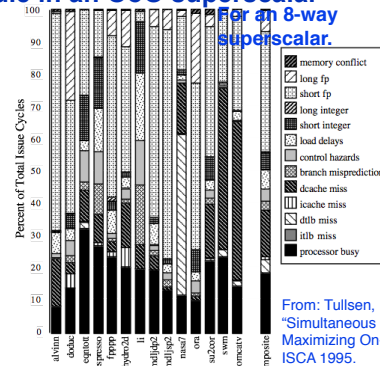
## Simultaneous Multithreading (SMT) for OoO Superscalars

- Techniques presented so far have all been "vertical" multithreading where each pipeline stage works on one thread at a time
- SMT uses fine-grain control already present inside an OoO superscalar to allow instructions from multiple threads to enter execution on same clock cycle. Gives better utilization of machine resources.

CSE 490/590, Spring 2011

21

## For most apps, most execution units lie idle in an OoO superscalar



From: Tullsen, Eggers, and Levy, "Simultaneous Multithreading: Maximizing On-chip Parallelism", ISCA 1995.

CSE 490/590, Spring 2011

22

## Acknowledgements

- These slides heavily contain material developed and copyright by
  - Krste Asanovic (MIT/UCB)
  - David Patterson (UCB)
- And also by:
  - Arvind (MIT)
  - Joel Emer (Intel/MIT)
  - James Hoe (CMU)
  - John Kubiatowicz (UCB)
- MIT material derived from course 6.823
- UCB material derived from course CS252

CSE 490/590, Spring 2011

23