

# CS184b: Computer Architecture (Abstractions and Optimizations)

Day 19: May13, 2005  
Multithreading



Caltech CS184 Spring2005 -- DeHon

## Today

- Multitasking/Multithreading model
- Fine-Grained Multithreading
- SMT (Symmetric Multi-Threading)

2

Caltech CS184 Spring2005 -- DeHon

## Problem

- Long latency of operations
  - IO or page-fault
  - Non-local memory fetch
    - Main memory, L3, remote node in distributed memory
  - Long latency operations (mpy, fp)
- Wastes processor cycles while stalled
- If processor stalls on return
  - Latency problem turns into a throughput (utilization) problem
  - CPU sits idle

3

Caltech CS184 Spring2005 -- DeHon

## Idea

- Run something else useful while stalled
- In particular, another process/thread
  - Another PC
  
- Use parallelism to “tolerate” latency

4

Caltech CS184 Spring2005 -- DeHon

## Old Idea

- Share expensive machine among multiple users (jobs)
- When one user task must wait on IO
  - Run another one
- Time multiplex machine among users

5

Caltech CS184 Spring2005 -- DeHon

## Mandatory Concurrency

- Some tasks must be run “concurrently” (interleaved) with user tasks
  - DRAM Refresh
  - IO
    - Keyboard, network, ...
  - Window system (xclock...)
  - Autosave ☺
  - Clippy ☹

6

Caltech CS184 Spring2005 -- DeHon

## Other Useful Concurrency

- Print spooler
- Web browser
  - Download images in parallel
- Instant Messenger/Zephyr (Gale)
- biff/xbiff/xfaces...

Caltech CS184 Spring2005 -- DeHon

7

## Multitasking

- Single machine run multiple tasks
- Machine provides same ISA/sequential semantics to each task
  - Task believes it own machines
  - Same as if other tasks running on different machines
- Tasks isolated from one another
  - Cannot affect each other functionally
  - (may impact each other's performance)

Caltech CS184 Spring2005 -- DeHon

8

## Each task/process

- **Process** – virtualization of the CPU
  - Has own unique set of state:
    - PC
    - Registers
    - VM Page Table (hence memory image)

Caltech CS184 Spring2005 -- DeHon

9

## Sharing the CPU

- Save/Restore
  - PC/Registers/Page Table
- Virtual Memory Isolation
- Privileged system software
  - User/System mode execution
- Functionally, task not notice that it gave up the CPU for period of time

Caltech CS184 Spring2005 -- DeHon

10

## Threads

- **Threads** – separate PC, but **shares and address space**
- Has own processor state:
  - PC
  - Registers
- Shares
  - Memory
  - VM Page Table
- Process may have multiple threads

Caltech CS184 Spring2005 -- DeHon

11

## Multitasking/Multithreading

- Gives us an initial model for parallelism
- So far, parallelism of unrelated tasks
- Eventually, cooperating
  - Have to address concurrent memory model
  - (next time)

Caltech CS184 Spring2005 -- DeHon

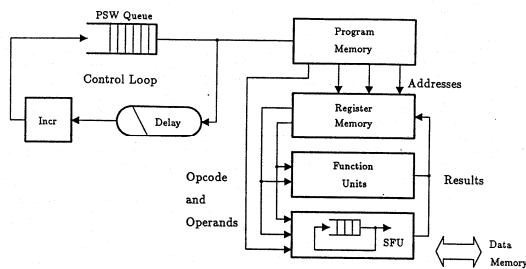
12

## Fine Grained

## HEP/ $\mu$ Unity/Tera

- Provide a number of contexts
  - Separate PCs, register files, ...
- Number of contexts  $\geq$  operation latency
  - Pipeline depth
  - Roundtrip time to main memory
- Run each context in round-robin fashion

## HEP Pipeline



[figure: Arvind+Innucci, DFVLR'87]

## Strict Interleaved Threading

- Uses parallelism to get throughput
  - Avoid interlocking, bypass...
  - Cover memory latency
  - Essentially **C-slow** transformation of processor
- Potentially poor single-threaded performance
  - Increases end-to-end latency of thread

## Compare Graph Machine

- How does this compare to our Graph Machine Model?
  - What's a thread?
  - What latency are we hiding?

## SMT

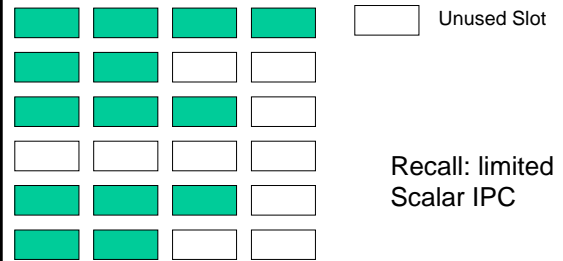
## Superscalar and Multithreading?

- Do both?
- Issue from multiple threads into pipeline
- No worse than (super)scalar on single thread
- More throughput with multiple threads
  - Fill in what would have been empty issue slots with instructions from different threads

Caltech CS184 Spring2005 -- DeHon

19

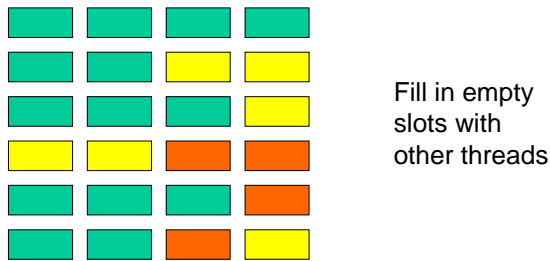
## SuperScalar Inefficiency



Caltech CS184 Spring2005 -- DeHon

20

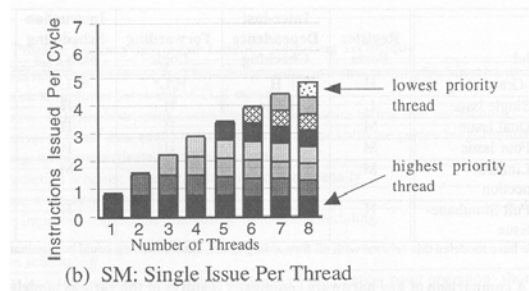
## SMT Promise



Caltech CS184 Spring2005 -- DeHon

21

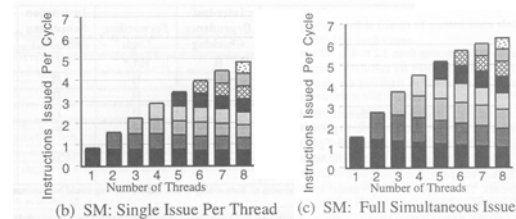
## SMT Estimates (ideal)



Caltech CS184 Spring2005 -- DeHon

22

## SMT Estimates (ideal)



[Tullsen et al. ISCA '95]

Caltech CS184 Spring2005 -- DeHon

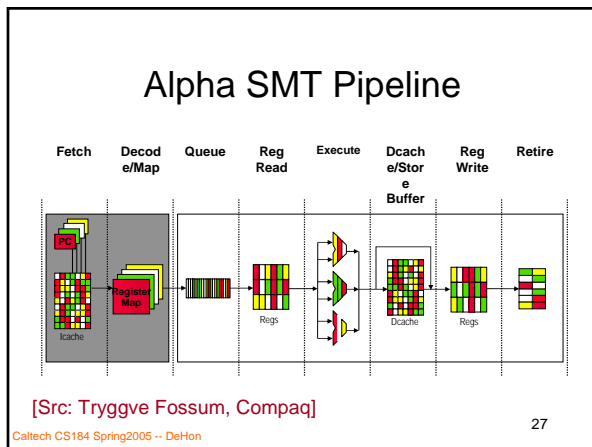
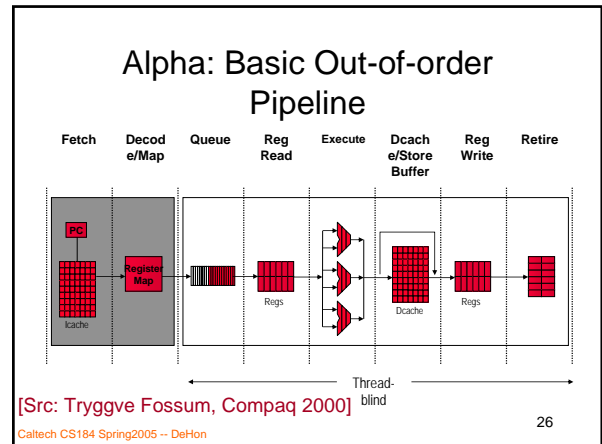
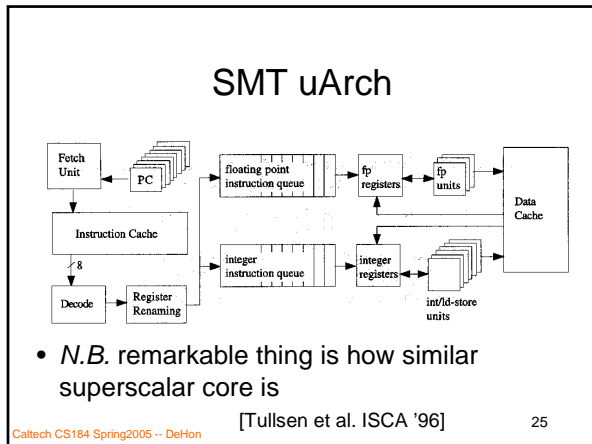
23

## SMT uArch

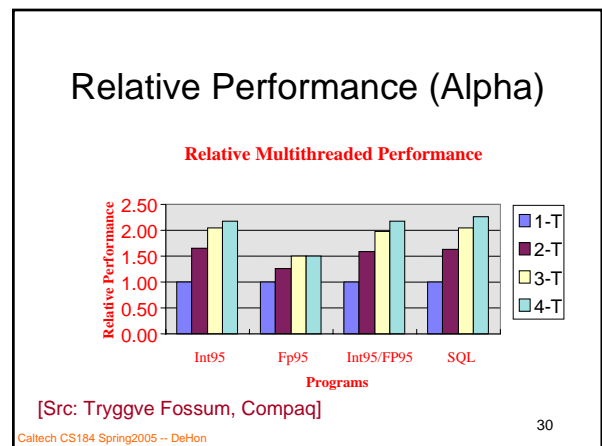
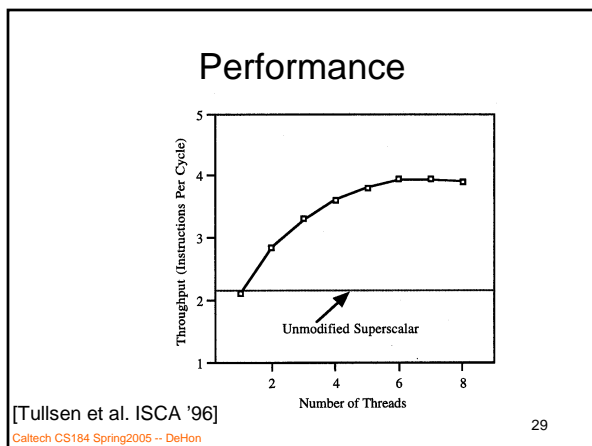
- **Observation:** exploit register renaming
  - Get small modifications to existing superscalar architecture
  - **Key trick:** different threads (processes) get distinct physical register assignments

Caltech CS184 Spring2005 -- DeHon

24



- ### SMT uArch
- Changes:
    - Multiple PCs
    - Control to decide how to fetch from
    - Separate return stacks per thread
    - Per-thread reorder/commit/flush/trap
    - Thread id w/ BTB
    - Larger register file
      - More things outstanding
- Caltech CS184 Spring2005 -- DeHon 28



## Alpha SMT

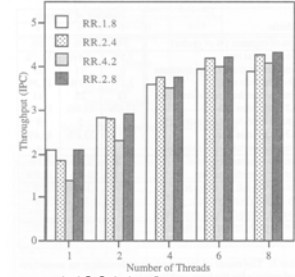
- Cost-effective Multiprocessing--increased throughput
- 4 X architectural registers
- 2 X performance gain with little additional cost and complexity

Caltech CS184 Spring2005 -- DeHon

31

## Optimizing: fetch freedom

- RR=Round Robin
- RR.X.Y
  - X - threads do fetch in cycle
  - Y - instructions fetched/thread



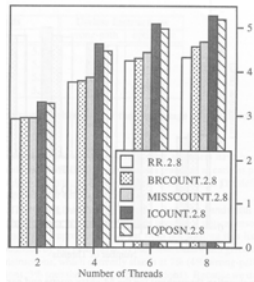
[Tullsen et al. ISCA '96]

Caltech CS184 Spring2005 -- DeHon

32

## Optimizing: Fetch Alg.

- ICOUNT - priority to thread w/ fewest pending instrs
- BRCCOUNT
- MISSCOUNT
- IQPOSN - penalize threads w/ old instrs (at front of queues)



[Tullsen et al. ISCA '96]

Caltech CS184 Spring2005 -- DeHon

33

## Throughput Improvement

- 8-issue superscalar
  - Achieves little over 2 instructions per cycle
- Optimized SMT
  - Achieves 5.4 instructions per cycle on 8 threads
- 2.5x throughput increase

Caltech CS184 Spring2005 -- DeHon

34

## Costs

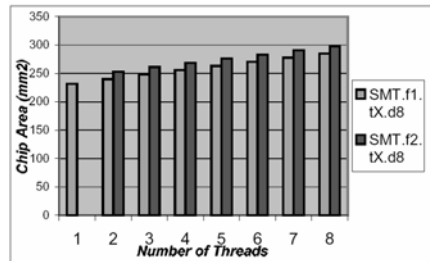
Function	Chip Block	MIPS		Relative area increase of adding SMT		
		R10K.2x w/0.19uB SMT (mm <sup>2</sup> )	R10K.2x 0.19uB (mm <sup>2</sup> )	% increase of adding SMT to R10K.2x Block	% increase versus core area	% increase versus chip area excluding L2 cache
Dcache	Dcache	11.4	11.4	0	0.0	0.0
	Dtag	1.6	1.6	0	0.0	0.0
Icache	Icache	6.3	13.7	50	2.9	2.2
	Tag	1.3	1.3	0	0.0	0.0
ILB	ILB	2.4	5.7	30	0.7	0.6
	Tag	0.0	0.0	0	0.0	0.0
Fetch	Fetch	1.0	4.2	137	5.0	4.1
	Branch	1.6	1.6	0	0.0	0.0
Decode	Decode	2.3	4.5	96	1.7	1.2
	Branch/Logic	2.5	3.3	33	10.4	13.9
Overhead/Order execution	Branch/Order	1.4	16.2	108	18.9	10.6
	Fetch/Order	2.5	2.9	16	0.0	0.0
	IQ	7.0	8.9	27	0.0	0.0
	L20	0.4	11.8	29	0.0	0.0
	FP0	6.3	9.0	43	0.0	0.0
	Branch	6.1	7.9	29	0.0	0.0
	RAS	0.3	2.1	60	0.0	0.0
	Int0	3.7	18.4	231	20.0	14.3
	FP0	5.5	12.6	128	0.0	0.0
	FP1	7.6	7.6	0	0.0	0.0
Arithmetic Units	FPALU	4.0	4.0	0	0.0	0.0
	FPALU	5.3	5.3	0	0.0	0.0
Miscellaneous	Int0	5.2	5.2	0	0.0	0.0
	FP0	0.9	0.9	0	0.0	0.0
	Misc.	2.4	2.4	0	0.0	0.0
	Int0	18.7	18.7	0	0.0	0.0
Heating	Heating	12.9	12.9	0	0.0	0.0
	Heating	5.5	5.5	0	0.0	0.0
SMT L2 Cache	Cache	126.7	183.8	46.7	31.1	31.1
	Cache w/ L2 cache	176.7	257.7	46.7	31.1	31.1
	Cache w/ L2 cache	281.7	287.7	0	0.0	0.0
	Cache w/ L2 cache	281.7	287.7	0	0.0	0.0

[Burns+Gaudiot HPCA'2000]

Caltech CS184 Spring2005 -- DeHon

35

## Costs



Single and double cache line fetches

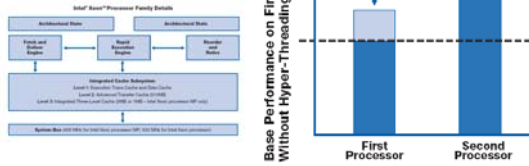
Caltech CS184 Spring2005 -- DeHon

[Burns+Gaudiot HPCA'2000]

36

## Intel: Hyperthreading

- Supports 2 threads



[http://www.intel.com/business/bss/products/hyperthreading/server/ht\\_server.pdf](http://www.intel.com/business/bss/products/hyperthreading/server/ht_server.pdf)

Caltech CS184 Spring2005 -- DeHon

## Admin

- No class on Monday
- Class meet next on Wednesday
  - ...and will meet on Friday

Caltech CS184 Spring2005 -- DeHon

38

## Big Ideas

- 0, 1, Infinity → virtualize resources
  - Processes virtualize CPU
- Latency Hiding
  - Processes, Threads
  - Find something else useful to do while wait...
    - C-Slow

Caltech CS184 Spring2005 -- DeHon

39