Thread-Level Parallelism — Simultaneous MultiThreading (SMT) & Chip Multi-Processors (CMP)

Hung-Wei



Recap: What about "linked list"

Static instructions

LOOP: ld X10, 8(X10) addi X7, X7, 1 bne X10, X0, LOOP

Dynamic instructions



3

5

7

9

ene

6

8



Demo: ILP within a program

 perf is a tool that captures performance counters of your processors and can generate results like branch mis-prediction rate, cache miss rates and ILP.



Simultaneous multithreading: maximizing on-chip parallelism Dean M. Tullsen, Susan J. Eggers, Henry M. Levy **Department of Computer Science and Engineering, University of Washington**



Simultaneous multithreading

- The processor can schedule instructions from different threads/processes/programs
- Fetch instructions from different threads/processes to fill the not utilized part of pipeline
 - Exploit "thread level parallelism" (TLP) to solve the problem of insufficient ILP in a single thread
 - You need to create an illusion of multiple processors for OSs



Simultaneous multithreading





ld add bne ld add bne ld add

X1, 0(X10)addi X10, X10, 8 X20, X20, X1 X10, X2, LOOP X1, 0(X10)addi X10, X10, 8 X20, X20, X1 X10, X2, LOOP X1, 0(X10)addi X10, X10, 8 X20, X20, X1 bne X10, X2, LOOP





SMT

- Improve the throughput of execution
 - May increase the latency of a single thread
- Less branch penalty per thread
- Increase hardware utilization
- Simple hardware design: Only need to duplicate PC/Register **Files**
- Real Case:
 - Intel HyperThreading (supports up to two threads per core)
 - Intel Pentium 4, Intel Atom, Intel Core i7
 - AMD RyZen



Wider-issue processors won't give you much more

Program	IPC	BP Rate %	I cache %MPCI	D cache %MPCl	L2 cache %MPCI	Program	IPC	BP Rate %	l cache %MPCI	D cache %MPCI	L2 cache %MPCI
compress	0.9	85.9	0.0	3.5	1.0	compress	1.2	86.4	0.0	3.9	1.1
eqntott	1.3	79.8	0.0	0.8	0.7	eqntott	1.8	80.0	0.0	1.1	1.1
m88ksim	1.4	91.7	2.2	0.4	0.0	m88ksim	2.3	92.6	0.1	0.0	0.0
MPsim	0.8	78.7	5.1	2,3	2.3	MPsim	1.2	81.6	3.4	1.7	2.3
applu	0.9	79.2	0.0	2.0	1.7	applu	1.7	79.7	0.0	2.8	2.8
apsi	0.6	95.1	1.0	4.1	2.1	apsi	1.2	95.6	0.2	3.1	2.6
swim	0.9	99.7	0.0	1.2	1.2	swim	2.2	99.8	0.0	2.3	2.5
tomcatv	0.8	99.6	0.0	7.7	2.2	tomcatv	1.3	99,7	0.0	4.2	4.3
pmake	1.0	86.2	2.3	2.1	0.4	pmake	1.4	82.7	0.7	1.0	0.6

Table 5. Performance of a single 2-issue superscalar processor. Table 6. Performance of the 6-issue superscalar processor.

The case for a Single-Chip Multiprocessor Kunle Olukotun, Basem A. Nayfeh, Lance Hammond, Ken Wilson, and Kunyung

Kunle Olukotun, Basem A. Nayfeh, Lance Hammond, Ken W Chang Stanford University

Wide-issue SS processor v.s. multiple narrower-issue SS processors



Intel SkyLake





AMDL

RYZEN

Concept of CMP





Architectural Support for Parallel Programming

Parallel programming

- To exploit parallelism you need to break your computation into multiple "processes" or multiple "threads"
- Processes (in OS/software systems)
 - Separate programs actually running (not sitting idle) on your computer at the same time.
 - Each process will have its own virtual memory space and you need explicitly exchange data using inter-process communication APIs
- Threads (in OS/software systems)
 - Independent portions of your program that can run in parallel
 - All threads share the same virtual memory space
- We will refer to these collectively as "threads"
 - A typical user system might have 1-8 actively running threads.
 - Servers can have more if needed (the sysadmins will hopefully configure it that way)



What software thinks about "multiprogramming" hardware



What software thinks about "multiprogramming" hardware



Coherency & Consistency

- Coherency Guarantees all processors see the same value for a variable/memory address in the system when the processors need the value at the same time
 - What value should be seen
- Consistency All threads see the change of data in the same order
 - When the memory operation should be done



Simple cache coherency protocol

- Snooping protocol
 - Each processor broadcasts / listens to cache misses
- State associate with each block (cacheline)
 - Invalid
 - The data in the current block is invalid
 - Shared
 - The processor can read the data
 - The data may also exist on other processors
 - Exclusive
 - The processor has full permission on the data
 - The processor is the only one that has up-to-date data







Snooping Protocol



read miss/hit

What happens when we write in coherent caches?



Cache coherency



4Cs of cache misses

- 3Cs:
 - Compulsory, Conflict, Capacity
- Coherency miss:
 - A "block" invalidated because of the sharing among processors.



False sharing

- True sharing
 - Processor A modifies X, processor B also want to access X.
- False sharing
 - Processor A modifies X, processor B also want to access Y. However, Y is invalidated because X and Y are in the same block!

fence instructions

- x86 provides an "mfence" instruction to prevent reordering across the fence instruction
- x86 only supports this kind of "relaxed consistency" model. You still have to be careful enough to make sure that your code behaves as you expected



thread 2

mfenceb=1 must occur/update before mfence

Take-aways of parallel programming

- Processor behaviors are non-deterministic
 - You cannot predict which processor is going faster
 - You cannot predict when OS is going to schedule your thread
- Cache coherency only guarantees that everyone would eventually have a coherent view of data, but not when
- Cache consistency is hard to support

