CS61C – Machine Structures

Lecture 32 - Caches II

4/12/2006

John Wawrzynek

(www.cs.berkeley.edu/~johnw)

www-inst.eecs.berkeley.edu/~cs61c/

CS 61C L32 Caches II (1)

Wawrzynek Spring 2006 © UCB



address	hex	binary	
ſ	00	00000	maps to row 0
_	02	00010	maps to row 1
Ē	04	00100	maps to row 2
-	06	00110	maps to row 3
F	07	01000	maps to row 0
F	09 0A	01001	inaps to row o
-	0B 0C	01011 01100	maps to row 1
-	0D OE	01101	maps to row 2
-	<u>0</u> F	01111	maps to row 3
-	11	10001	maps to row 0
	13	10010	maps to row 1
	14 15	10100 10101	maps to row 2
	16 17	$10110 \\ 10111$	maps to row 3
	18 19	11000 11001	maps to row 0
	1A 1B	$11010 \\ 11011$	maps to row 1
Ī	1C 1D	11100 11101	maps to row 2
Ī	1E 1F	11110 11111	maps to row 3

DM Cache Review

Assume: 32 Byte Main Memory (addresses shown to left)

2 Bytes/block \Rightarrow 1 bit for "offset"

4 blocks in cache \Rightarrow 2 bits for "index" Slots in cache called: row 0, 1, 2, & 3

CS 61C L32 Caches II (3)

address	: hex	binary		
	00	00000		
_	01	00001	maps to row 0	0 tag
	02 03	00010 00011	maps to row 1	0 tag
	04 05	00100 00101	maps to row 2	0 tag
	06 07	00110 00111	maps to row 3	0 tag
-	08 09	01000 01001	maps to row 0	1 tag
	0A 0B	01010 01011	maps to row 1	1 tag
	0C 0D	$01100 \\ 01101$	maps to row 2	1 tag
	OE OF	$\begin{array}{c} 01110\\ 01111 \end{array}$	maps to row 3	1 tag
	10 11	10000 10001	maps to row 0	2 tag
	12 13	10010 10011	maps to row 1	2 tag
	14 15	$10100 \\ 10101$	maps to row 2	2 tag
	16 17	$10110 \\ 10111$	maps to row 3	2 tag
	18 19	11000 11001	maps to row 0	3 tag
	1A 1B	$11010 \\ 11011$	maps to row 1	3 tag
	1C 1D	$11100 \\ 11101$	maps to row 2	3 tag
	1E 1F	11110 11111	maps to row 3	3 tag

DM Cache Review

Wawrzynek Spring 2006 © UCB

Assume: 32 Byte Main Memory (addresses shown to left)

2 Bytes/block \Rightarrow 1 bit for "offset"

4 blocks in cache \Rightarrow 2 bits for "index" Slots in cache called: row 0, 1, 2, & 3

5 - 1 - 3 \Rightarrow 2 bits for "tag"

Wawrzynek Spring 2006 © UCB

CS 61C L32 Caches II (4)

Accessing data in a direct mapped cache



CS 61C L32 Caches II (5)

Do an example yourself. What happens?



CS 61C L32 Caches II (6)

Answers



Block Size Tradeoff (1/3)

^o Benefits of Larger Block Size

- <u>Spatial Locality</u>: if we access a given word, we're likely to access other nearby words soon
- Very applicable with Stored-Program Concept: if we execute a given instruction, it's likely that we'll execute the next few as well
- Works nicely in sequential access to arrays and structs too

CS 61C L32 Caches II (8)

Block Size Tradeoff (2/3)

^o Drawbacks of Larger Block Size

- · Larger block size means larger miss penalty
 - on a miss, takes longer time to load a new block from next level
- If block size is too big relative to cache size, then there are too few blocks
 - Result: miss rate goes up
- ° In general, try to minimize

Average Memory Access Time (AMAT)

= Hit Time + Miss Penalty x Miss Rate

CS 61C L32 Caches II (9)

Wawrzynek Spring 2006 © UCB

Block Size Tradeoff (3/3)

•<u>Hit Time</u> = time to find and retrieve data from current level cache

^o <u>Miss Penalty</u> = average time to retrieve data on a current level miss (includes the possibility of misses on successive levels of memory hierarchy)

^o<u>Hit Rate</u> = % of requests that are found in current level cache

^o<u>Miss Rate</u> = 1 - Hit Rate

CS 61C L32 Caches II (10)



CS 61C L32 Caches II (11)

Wawrzynek Spring 2006 © UCB



Block Size Tradeoff Conclusions

Administrivia

^oRegraded Midterms handout after class.

°Exam 2 next Wednesday

- ·7-9pm, Pimentel.
- Covers through pipelining (last week).
- Special review session Monday evening.

CS 61C L32 Caches II (13)

Wawrzynek Spring 2006 © UCB

Types of Cache Misses (1/2)

° "Three Cs" Model of Misses

°1st C: Compulsory Misses

- occur when a program is first started
- cache does not contain any of that program's data yet, so misses are bound to occur
- can't be avoided easily, so won't focus on these in this course

CS 61C L32 Caches II (14)

Types of Cache Misses (2/2)

° 2nd C: Conflict Misses

- miss that occurs because two distinct memory addresses map to the same cache location
- two blocks (which happen to map to the same location) can keep overwriting each other
- big problem in direct-mapped caches
- how do we lessen the effect of these?

^o Dealing with Conflict Misses

- Solution 1: Make the cache size bigger
 Fails at some point
- Solution 2: Multiple distinct blocks can fit in the same cache Index?

CS 61C L32 Caches II (15)

Wawrzynek Spring 2006 © UCB

Fully Associative Cache (1/3)

^o Memory address fields:

- Tag: same as before
- Offset: same as before
- Index: non-existant

°What does this mean?

- no "rows": any block can go anywhere in the cache
- must compare with all tags in entire cache to see if data is there

CS 61C L32 Caches II (16)

Fully Associative Cache (2/3)

° Fully Associative Cache (e.g., 32 B block)

compare tags in parallel



CS 61C L32 Caches II (17)

Wawrzynek Spring 2006 © UCB

Fully Associative Cache (3/3)

^oBenefit of Fully Assoc Cache

 No Conflict Misses (since data can go anywhere)

^o Drawbacks of Fully Assoc Cache

 Need hardware comparator for every single entry: if we have a 64KB of data in cache with 4B entries, we need 16K comparators: very expensive!

CS 61C L32 Caches II (18)

Third Type of Cache Miss

° Capacity Misses

- miss that occurs because the cache has a limited size
- miss that would not occur if we increase the size of the cache
- sketchy definition, so just get the general idea

^o This is the primary type of miss for Fully Associative caches.

CS 61C L32 Caches II (19)

Wawrzynek Spring 2006 © UCB

N-Way Set Associative Cache (1/4)

^oMemory address fields:

- Tag: same as before
- Offset: same as before
- Index: points us to the correct "row" (called a <u>set</u> in this case)

°So what's the difference?

- · each set contains multiple blocks
- once we've found correct set, must compare with all tags in that set to find our data

CS 61C L32 Caches II (20)

N-Way Set Associative Cache (2/4)

°Summary:

- · cache is direct-mapped w/respect to sets
- · each set is fully associative
- basically N direct-mapped caches working in parallel: each has its own valid bit and data

CS 61C L32 Caches II (21)

Wawrzynek Spring 2006 © UCB

N-Way Set Associative Cache (3/4)

^oGiven memory address:

- Find correct set using Index value.
- Compare Tag with all Tag values in the determined set.
- If a match occurs, hit!, otherwise a miss.
- Finally, use the offset field as usual to find the desired data within the block.

CS 61C L32 Caches II (22)



CS 61C L32 Caches II (23)





Cache Things to Remember

^o Caches are NOT mandatory:

- Processor performs arithmetic
- Memory stores data
- Caches simply make data transfers go faster
- ^o Each level of Memory Hiererarchy subset of next higher level
- ^o Caches speed up due to temporal locality: store data used recently
- ^o Block size > 1 wd spatial locality speedup: Store words next to the ones used recently
- ^o Cache design choices:
 - size of cache: speed v. capacity
 - N-way set assoc: choice of N (direct-mapped, fully-associative just special cases for N)

CS 61C L32 Caches II (26)