

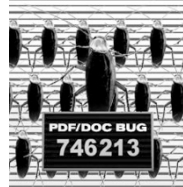
**Lecture 23
 Caches II
 2010-07-29**



Instructor Paul Pearce

**TOOLS THAT
 AUTOMATICALLY FIND
 SOFTWARE BUGS**

Black Hat (a security conference) is currently underway, and this year there are several talks focusing on automatically identifying and classifying software bugs. Noah Johnson (our TA) is currently at this conference giving a talk on BitBlaze, one of the tools discussed in this article. Check it out.



<http://www.technologyreview.com/computing/25869/?a=f>

And in Review...

- We would like to have the capacity of disk at the speed of the processor: unfortunately this is not feasible.
- So we create a memory hierarchy:
 - each successively lower level contains "most used" data from next higher level
 - exploits temporal & spatial locality
 - do the common case fast, worry less about the exceptions (design principle of MIPS)
- Locality of reference is a Big Idea



Where we left off: direct mapped cache example

- Ex.: 16KB of data, direct-mapped, 4 word blocks

- Offset: 4 bits
- Index: 10 bits
- Tag: 18 bits

- Read 4 addresses

1. 0x00000014
2. 0x0000001C
3. 0x00000034
4. 0x00008014

Address (hex)	Value of Word
00000010	a
00000014	b
00000018	c
0000001C	d
...	...
00000030	e
00000034	f
00000038	g
0000003C	h
...	...
00008010	i
00008014	j
00008018	k
0000801C	l
...	...



Memory vals here:

3. Load Word from 0x00000034

- 00000000000000000000 0000000011 0100

Valid	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				

1022	0				
1023	0				



So read block 3

- 00000000000000000000 0000000011 0100

Valid	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				

1022	0				
1023	0				



No valid data

- 00000000000000000000 0000000011 0100

Valid	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	0				
4	0				
5	0				
6	0				
7	0				

1022	0				
1023	0				



Load that cache block, return word f

- 00000000000000000000 0000000011 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	1	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				

Cal

4. Load word from 0x00008014

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	1	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				

Cal

So read Cache Block 1, Data is Valid

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	1	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				

Cal

Cache Block 1 Tag does not match (0 != 2)

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	d	c	b	a
2	0				
3	1	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				

Cal

Miss, so replace block 1 with new data & tag

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	2	l	k	i
2	0				
3	1	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				

Cal

And return word J

- 00000000000000000010 000000001 0100

Valid	Tag field	Index field	Offset		
Index	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	1	2	l	k	i
2	0				
3	1	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					
1022	0				
1023	0				

Cal

Do an example yourself. What happens?

- Chose from: Cache: Hit, Miss, Miss w. replace
Values returned: a, b, c, d, e, ..., k, l
- Load word from address 0x00000030 ?
000000000000000000 000000011 0000
- Load word from address 0x0000001c ?
000000000000000000 0000000001 1100

Cache

Valid	Tag	0xc-f	0x8-b	0x4-7	0x0-3
0	0				
1	2	l	k	n	i
2	0				
3	0	h	g	f	e
4	0				
5	0				
6	0				
7	0				
...					



CS61C L23 Caches II (13)

Pearce, Summer 2010 © UCB

What to do on a cache write?

- **First, we need to get data into the cache**
 - Follow same procedure as a read. If a cache miss, read from memory. Once there, we have to make a policy decision.
- **Write-through**
 - update the word in cache block and corresponding word in memory
- **Write-back**
 - update word in cache block
 - allow memory word to be "stale"
 - => add 'dirty' bit to each block indicating that memory needs to be updated when block is replaced
 - => OS flushes cache before I/O...



Performance trade-offs?

CS61C L23 Caches II (15)

Pearce, Summer 2010 © UCB

Block Size Tradeoff (1/3)

- **Benefits of Larger Block Size**
 - **Spatial Locality:** 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 array accesses too



CS61C L23 Caches II (16)

Pearce, Summer 2010 © UCB

Block Size Tradeoff (2/3)

- **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
 - Ping-pong effect, compromises temporal locality
 - Result: miss rate goes up
- **In general, minimize Average Memory Access Time (AMAT)**



CS61C L23 Caches II (17)

= Hit Time + Miss Penalty x Miss Rate

Pearce, Summer 2010 © UCB

Block Size Tradeoff (3/3)

- **Hit Time**
 - time to find and retrieve data from current level cache
- **Miss Penalty**
 - average time to retrieve data on a current level miss (includes the possibility of misses on successive levels of memory hierarchy)
- **Hit Rate**
 - % of requests that are found in current level cache
- **Miss Rate**
 - 1 - Hit Rate

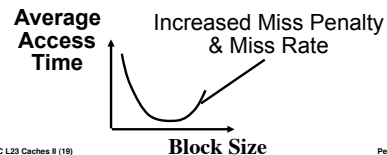
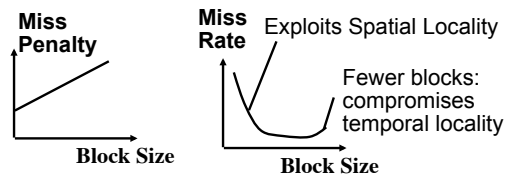


CS61C L23 Caches II (18)

Pearce, Summer 2010 © UCB

Block Size Tradeoff Conclusions

For a fixed cache size:



CS61C L23 Caches II (19)

Pearce, Summer 2010 © UCB

Administrivia

- Drop deadline Friday.
 - But you shouldn't drop. That would make me sad. =(
- My office hours are in 551 Soda today, not 511 Soda.
 - 551 is directly across the hall from 511. Just turn right instead of left when leaving elevator.
- Midterm regrades returned in lab today



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



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?
- Dealing with Conflict Misses
 - Solution 1: Make the cache size bigger
 - Falls at some point
 - Solution 2: Multiple distinct blocks can fit in the same cache index?



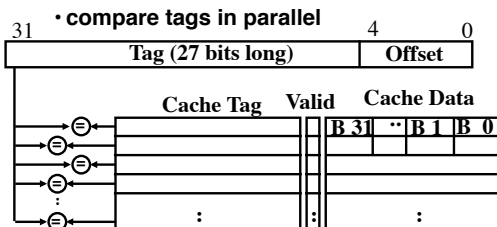
Fully Associative Cache (1/3)

- Memory address fields:
 - Tag: same as before
 - Offset: same as before
 - Index: non-existent
- What does this mean?
 - any block can go anywhere in the cache
 - must compare with all tags in entire cache to see if data is there



Fully Associative Cache (2/3)

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



Fully Associative Cache (3/3)

- Benefit of Fully Assoc Cache
 - No Conflict Misses (since data can go anywhere)
- Drawbacks of Fully Assoc Cache
 - Need hardware comparator for every single block: if we have a 64KB of data in cache with 4B blocks, we need 16K comparators: infeasible



Final Type of Cache Miss

• 3rd C: 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
- This is the primary type of miss for Fully Associative caches.



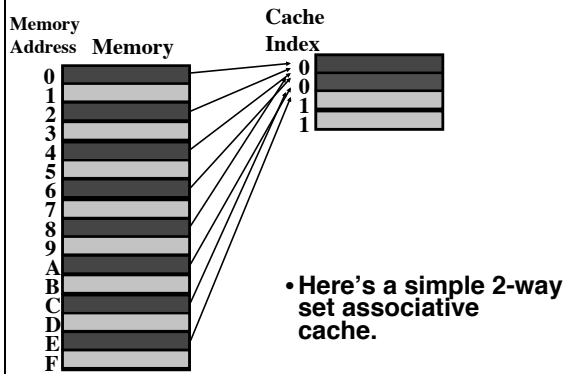
N-Way Set Associative Cache (1/3)

• Memory address fields:

- Tag: same as before
- Offset: same as before
- Index: points us to the correct "row" (called a set 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



Associative Cache Example



N-Way Set Associative Cache (2/3)

• Basic Idea

- cache is direct-mapped w/respect to sets
- each set is fully associative with N blocks in it
- Given 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.



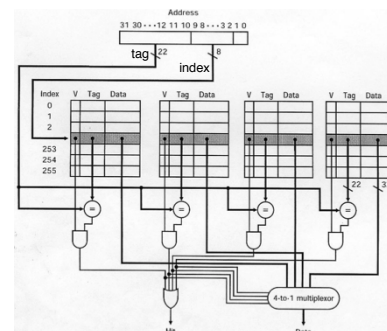
N-Way Set Associative Cache (3/3)

• What's so great about this?

- even a 2-way set assoc cache avoids a lot of conflict misses
- hardware cost isn't that bad: only need N comparators
- In fact, for a cache with M blocks,
 - it's Direct-Mapped if it's 1-way set assoc
 - it's Fully Assoc if it's M-way set assoc
 - so these two are just special cases of the more general set associative design



4-Way Set Associative Cache Circuit



Block Replacement Policy

- **Direct-Mapped Cache**
 - index completely specifies position which position a block can go in on a miss
- **N-Way Set Assoc**
 - index specifies a set, but block can occupy any position within the set on a miss
- **Fully Associative**
 - block can be written into any position
- **Question: if we have the choice, where should we write an incoming block?**
 - If there are any locations with valid bit off (empty), then usually write the new block into the first one.
 - If all possible locations already have a valid block, we must pick a replacement policy: rule by which we determine which block gets "cached out" on a miss.



CS61C L23 Caches II (32)

Pearce, Summer 2010 © UCB

Block Replacement Policy: LRU

- **LRU (Least Recently Used)**
 - Idea: cache out block which has been accessed (read or write) least recently
 - Pro: temporal locality \Rightarrow recent past use implies likely future use: in fact, this is a very effective policy
 - Con: with 2-way set assoc, easy to keep track (one LRU bit); with 4-way or greater, requires complicated hardware and much time to keep track of this



CS61C L23 Caches II (33)

Pearce, Summer 2010 © UCB

Block Replacement Example

- We have a 2-way set associative cache with a four word total capacity and one word blocks. We perform the following word accesses (ignore bytes for this problem):

0, 2, 0, 1, 4, 0, 2, 3, 5, 4

- How many hits and how many misses will there be for the LRU block replacement policy?



CS61C L23 Caches II (34)

Pearce, Summer 2010 © UCB

Block Replacement Example: LRU

0: miss, bring into set 0 (loc 0)

2: miss, bring into set 0 (loc 1)

0: hit

1: miss, bring into set 1 (loc 0)

4: miss, bring into set 0 (loc 1, replace 2)

Addresses 0, 2, 0, 1, 4, 0, ...

0: hit

	loc 0	loc 1
set 0	0	iru
set 1		
set 0	iru	2
set 1		
set 0	0	iru, 2
set 1		
set 0	0	iru, 2
set 1	1	iru
set 0	iru	4
set 1	1	iru
set 0	0	iru, 4
set 1	1	iru



CS61C L23 Caches II (35)

Pearce, Summer 2010 © UCB

Big Idea

- How to choose between associativity, block size, replacement & write policy?
- Design against a performance model
 - Minimize: Average Memory Access Time
 - = Hit Time
 - + Miss Penalty x Miss Rate
 - influenced by technology & program behavior, architecture, cache purpose
- Create the illusion of a memory that is large, cheap, and fast - on average

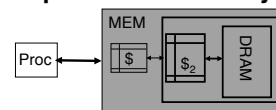


CS61C L23 Caches II (36)

Pearce, Summer 2010 © UCB

Improving Miss Penalty

- When caches first became popular, Miss Penalty \sim 10 processor clock cycles
- Today 2400 MHz Processor (0.4 ns per clock cycle) and 80 ns to go to DRAM \Rightarrow 200 processor clock cycles!



Solution: another cache between memory and the processor cache: Second Level (L2) Cache



CS61C L23 Caches II (37)

Pearce, Summer 2010 © UCB

Peer Instruction

1. A 2-way set-associative cache can be outperformed by a direct-mapped cache.
2. Larger block size \Rightarrow lower miss rate

	12
a)	FF
b)	FT
c)	TF
d)	TT



CS61C L23 Caches II (38)

Pearce, Summer 2010 © UCB

And in Conclusion...

- We've discussed memory caching in detail. Caching in general shows up over and over in computer systems
 - Filesystem cache, Web page cache, Game databases / tablebases, Software memoization, and many more!
- **Big idea:** if something is expensive but we want to do it repeatedly, do it once and cache the result.
- **Cache design choices:**
 - Size of cache: speed v. capacity
 - Block size (i.e., cache aspect ratio)
 - Write Policy (Write through v. write back)
 - Associativity choice of N (direct-mapped v. set v. fully associative)
 - Block replacement policy
 - 2nd level cache?
 - 3rd level cache?
- Use performance model to pick between choices, depending on programs, technology, budget, ...



CS61C L23 Caches II (40)

Pearce, Summer 2010 © UCB

Bonus slides

- These are extra slides that used to be included in lecture notes, but have been moved to this, the "bonus" area to serve as a supplement.
- The slides will appear in the order they would have in the normal presentation

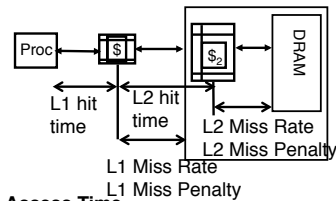
Bonus



CS61C L23 Caches II (41)

Pearce, Summer 2010 © UCB

Analyzing Multi-level cache hierarchy



$$\text{Avg Mem Access Time} = \text{L1 Hit Time} + \text{L1 Miss Rate} * \text{L1 Miss Penalty}$$

$$\text{L1 Miss Penalty} = \text{L2 Hit Time} + \text{L2 Miss Rate} * \text{L2 Miss Penalty}$$

$$\text{Avg Mem Access Time} = \text{L1 Hit Time} + \text{L1 Miss Rate} * (\text{L2 Hit Time} + \text{L2 Miss Rate} * \text{L2 Miss Penalty})$$



CS61C L23 Caches II (42)

Pearce, Summer 2010 © UCB

Example

- **Assume**
 - Hit Time = 1 cycle
 - Miss rate = 5%
 - Miss penalty = 20 cycles
 - Calculate AMAT...
- **Avg mem access time**

$$= 1 + 0.05 * 20$$

$$= 1 + 1 \text{ cycles}$$

$$= 2 \text{ cycles}$$



CS61C L23 Caches II (43)

Pearce, Summer 2010 © UCB

Ways to reduce miss rate

- **Larger cache**
 - limited by cost and technology
 - hit time of first level cache < cycle time (bigger caches are slower)
- **More places in the cache to put each block of memory – associativity**
 - fully-associative
 - any block any line
 - N-way set associated
 - N places for each block
 - direct map: N=1



CS61C L23 Caches II (44)

Pearce, Summer 2010 © UCB

Typical Scale

- L1
 - size: tens of KB
 - hit time: complete in one clock cycle
 - miss rates: 1-5%
- L2:
 - size: hundreds of KB
 - hit time: few clock cycles
 - miss rates: 10-20%
- L2 miss rate is fraction of L1 misses that also miss in L2
- why so high?



CS61C L23 Caches II (45)

Pearce, Summer 2010 © UCB

Example: with L2 cache

- Assume
 - L1 Hit Time = 1 cycle
 - L1 Miss rate = 5%
 - L2 Hit Time = 5 cycles
 - L2 Miss rate = 15% (% L1 misses that miss)
 - L2 Miss Penalty = 200 cycles
- L1 miss penalty = $5 + 0.15 * 200 = 35$
- Avg mem access time = $1 + 0.05 * 35 = 2.75$ cycles



CS61C L23 Caches II (46)

Pearce, Summer 2010 © UCB

Example: without L2 cache

- Assume
 - L1 Hit Time = 1 cycle
 - L1 Miss rate = 5%
 - L1 Miss Penalty = 200 cycles
- Avg mem access time = $1 + 0.05 * 200 = 11$ cycles
- 4x faster with L2 cache! (2.75 vs. 11)

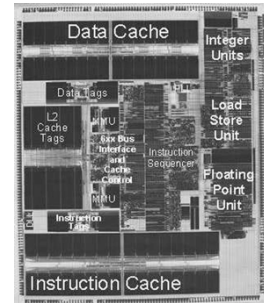


CS61C L23 Caches II (47)

Pearce, Summer 2010 © UCB

An actual CPU – Early PowerPC

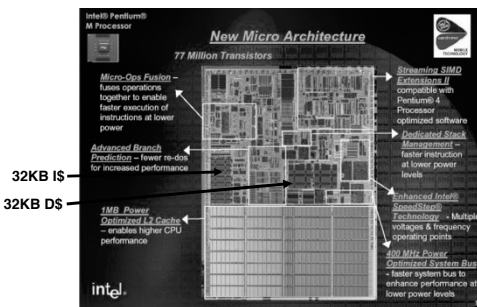
- Cache
 - 32 KB Instructions and 32 KB Data L1 caches
 - External L2 Cache interface with integrated controller and cache tags, supports up to 1 MByte external L2 cache
 - Dual Memory Management Units (MMU) with Translation Lookaside Buffers (TLB)
- Pipelining
 - Superscalar (3 inst/cycle)
 - 6 execution units (2 integer and 1 double precision IEEE floating point)



CS61C L23 Caches II (48)

Pearce, Summer 2010 © UCB

An Actual CPU – Pentium M



CS61C L23 Caches II (49)

Pearce, Summer 2010 © UCB