Caches III

CSE 351 Summer 2024

Instructor: Ellis Haker

Teaching Assistants:

Naama Amiel Micah Chang Shananda Dokka Nikolas McNamee Jiawei Huang



Async (Paris Arc II) @0xAsync

No mom it's not a "messy pile of clothes on my chair" it's an L1 cache for fast random access to my frequently used clothes in O(1) time. It needs to be big to avoid expensive cache misses (looking in my closet). I NEED to be minimizing latency, this is important to me. Please.

Administrivia

- Today
 - HW14 due (11:59pm)
- Friday, 7/26
 - RD17 (1pm)
 - No HW due!
- Sunday, 7/28
 - · Lab3 due (11:59pm) late Lve Late Tre 7/30
- Monday, 7/29
 - No RD due!
 - HW15 + 16 due (11:59pm)

Caches

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
 - Direct-mapped (sets; index + tag)
 - o Associativity (ways)
 - Replacement policy
 - Handling writes
- Program optimizations that consider caches

Recap: Direct-Mapped Cache



- Each block in memory can only go into one set
 - Fast and simple
- Hash function
 - o (Block #) mod (# of sets)

Direct-Mapped: a Problem!



- What if we have the following access pattern:
 - 0, 16, 0, 16, ...
 - Each access will miss
 - Thrashing!
 - Rest of the cache unused

One Solution: Fully Associative

- Can store any data block anywhere in the cache
 - When loading in a new block, store it in the first unused space
 - Only evict old blocks when the entire cache is full!
- Problems with this?
 - Requires complicated hardware to check each set
 - More power consumed, slower



Memory

tog= linchest

Cache

Full Solution: Associativity

- Rather than just one block per set, we can have multiple
 - Cache line = block + relevant metadata (i.e. tag)
 - Associativity (E) = number of lines in a set
 - Cache is "*E*-way set associative"
- If any block can go in any set, the cache is Fully Associative



Cache Organization: Associativity

K = block size (in bytes), **C** = cache size (in bytes), **E** = associativity

- **S** = number of sets = **C**÷**K**÷**E**
- Use lowest $s = \log_2(S)$ bits of block number to find the set
 - Direct-Mapped: E = 1, so $s = \log_2(C \div K)$, same as we saw previously
 - Fully Associative: E = C:K, so log₂(1) = 0
 as E increases, 5 Lecreases -> fewer in dex bits



Example Placement #1



K=4

Block size K: 16 B

Block Placement and Replacement

- If there is an empty line in the set, store data there
 - How to tell if the line is empty? Valid bit (0 = empty, 1 = used)
 - Stored in cache line
- If there are no empty lines, which one do we replace?
 - No choice for direct-mapped, easy!
 - For other caches, need a replacement policy
 - Ideal is least recently used (LRU)
 - Most real caches approximate LRU

Polling Questions

1. We have a cache of size 2 KiB with block size of 128 B. If our cache has 2 sets, what is its associativity?

7. 2 = 2 K



2 5

2. If addresses are 16 bits wide, how wide is the Tag field? $k = \log_2(k) = 7 \qquad s = \log_2(s)^2 \cdot 3$ $toslife = 1(-7 - 3 = 5 - \log_2(s)^2 \cdot 3)$



Notation Review

| | Parameter | Variable | Formulas |
|--------------------------------------|--------------------|------------------------------|---|
| | Block size | K (B in book) | |
| # Chile | Cache size | С | $C = K^* E^* S \leftrightarrow S = C/K/E$ $k = \log_2(K) \leftrightarrow K = 2^k$ |
| obbresses in memory ((we prem | Associativity | E | |
| | Number of sets | S | |
| | Address space | М | |
| & Dote | Address width | т | $\mathbf{S} = \log_2(\mathbf{S}) \leftrightarrow \mathbf{S} = 2^\circ$ |
| | Offset field width | k (b in book) | $\boldsymbol{m} = \log_2(\boldsymbol{M}) \leftrightarrow \boldsymbol{M} = 2^m$ |
| | Index field width | S | <i>m</i> = <i>t</i> + <i>s</i> + <i>k</i> |
| | Tag field width | t | |

Example Cache Problem 1KiB address space, 125 cycles to go to memory. Fill in the following table: 2^{13} $= 2^{12}$ $= 2^{12}$ = 12 (cost)

| 2 - 2 - 2 - 2 - 2 | Cache size (C) | 64B | |
|----------------------------------|----------------------------|-----------|----------|
| - | Block size (<i>K</i>) | 8B | |
| S=C+K+E | Associativity (E) | 2-way | |
| $= 2^{10} \div 2^3 \div 2 = 2^6$ | Hit Time | 3 cycles | |
| | Miss Rate | 20% | |
| _ | Address width (<i>m</i>) | lo hites | logz (M) |
| | Offset bits (<i>k</i>) | 3 Lits | 1.52 (K) |
| AMAT = 6+1 + MK - M | Index bits (s) | 6 Lits | Log 2 (5 |
| = 78 | Tag bits (<i>t</i>) | 1 617 | m-5 - k |
| | ΑΜΑΤ | 28 cycles | |

Read: General Steps

- 1. Break up address into offset, index, and tag
- 2. Locate the set using index
- 3. Check if any line in the set has our data
 - a. Valid bit = 1 and tag matches
 - b. If yes hit!
 - i. Otherwise, load in from memory
- 4. Read out data from block starting at offset

Read: Direct-Mapped Cache

One line per set (*E*=1), *K*=4B

1. Locate set



Read: Direct-Mapped Cache (pt 2)

One line per set (*E*=1), *K*=4B

- 1. Locate set
- 2. Check if valid and compare tag
 - a. No match? Old block gets evicted, replaced with new one

model if V==1 and Tag == 0x abcAddress of short0xabc 0...01 01valid?V Tag

Read: Direct-Mapped Cache (pt 3)

One line per set (*E*=1), *K*=4B

- 1. Locate set
- 2. Check if valid and compare tag
- 3. Get data starting at offset



Read: Set Associative Cache

Two lines per set (*E*=2), *K*=4B

1. Locate set



Read: Set Associative Cache (pt 2)

Two lines per set (*E*=2), *K*=4B

- 1. Locate set
- 2. Check if valid and compare tag for every line in the set
 - a. No match? One line is selected to get evicted and replaced by new one



Read: Set Associative Cache (pt 3)

Two lines per set (*E*=2), *K*=4B

- 1. Locate set
- 2. Check if valid and compare tag for every line in the set
- 3. Get data starting at offset



Types of Cache Misses: 3 C's

- Compulsory (cold-start) miss
 - Occurs on the first access to a block
 - Smaller block size = more compulsory misses
- Conflict miss
 - Occurs when cache is large enough to hold multiple data blocks, but they cannot be in the cache at the same time because they conflict
 - Lower associativity = more conflict misses
 - Does not occur in fully associative caches
- Capacity miss
 - Occurs when the set of active blocks (the working set) is too big to fit in the cache
 - Smaller cache size = more capacity misses

Code Analysis - more detailed explanation in vext

• Assuming cache starts **cold** (i.e. all blocks invalid), and sum, i, and j are all

stored in registers, calculate the **miss rate**. $2^{B} e^{-L_{1}} \circ m = 10$ bits, C = 64B, K = 8B, E = 2 $5 = C4 \div 3 \div 2 = 4$ k= 3 hits 5 = 2 hits محصا المحالي each block is #define SIZE 8 exicted short ar[SIZE][SIZE], sum = 0; // &ar=0x200 for (int i = 0; i < SIZE; i++) {</pre> Lefore its - [-]...-[5] for (int j = 0; j < SIZE; j++)</pre> vset again. 2 sum += ar[j][i]; MQ=100.1 } ~[0][0], ~[2] [0], ~[4][0], ~[(][0]= 5et 0 access a Lo J [e] ~[1][0].~[3][0], ~[5][1], ~[7][0] = set 2 c[1][0] ~[v] [o]

~[3][0]...

Summary

- Associativity: number of cache lines in a set
 - Direct-Mapped caches have associativity 1
 - Fully Associative caches have associativity equal to # of blocks (all in one set)
 - Most caches are somewhere in between
- 3 types of misses
 - **Compulsory**: first time accessing block
 - Conflict: block was evicted by another when the cache was not full
 - Capacity: block was evicted because the cache was full