

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Caches II

CSE 351 Autumn 2017

Instructor:
Justin Hsia

Teaching Assistants:
Lucas Wotton
Michael Zhang
Parker DeWilde
Ryan Wong
Sam Gehman
Sam Wolfson
Savanna Yee
Vinny Palaniappan

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Administrivia

- ❖ Homework 4 released tomorrow (Structs, Caches)
- ❖ Midterm Regrade Requests due Wednesday (11/8)
- ❖ Lab 3 due *Friday* (11/10)
- ❖ **Mid-Quarter Survey Feedback**
 - Pace is “moderate” to “a bit too fast”
 - You talk too fast in lecture (or rush at the end) and I wish there were more peer instruction questions
 - Canvas quiz answer keys are annoying, but instant homework feedback is great

2

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Memory Hierarchies

- ❖ Some fundamental and enduring properties of hardware and software systems:
 - Faster storage technologies almost always cost more per byte and have lower capacity
 - The gaps between memory technology speeds are widening
 - True for: registers ↔ cache, cache ↔ DRAM, DRAM ↔ disk, etc.
 - Well-written programs tend to exhibit good locality
- ❖ These properties complement each other beautifully
 - They suggest an approach for organizing memory and storage systems known as a [memory hierarchy](#)

3

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

An Example Memory Hierarchy

4

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

An Example Memory Hierarchy

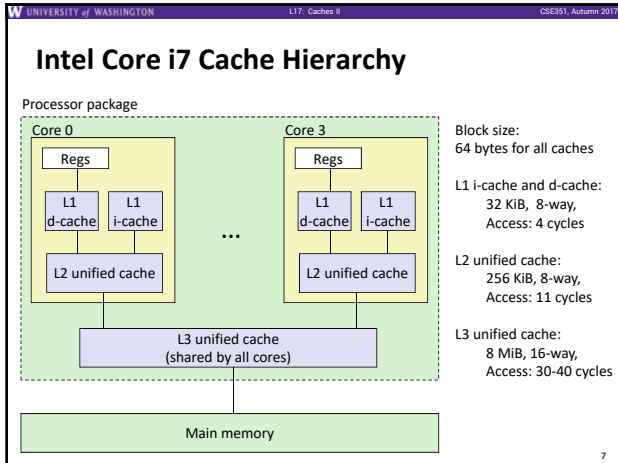
5

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Memory Hierarchies

- ❖ Fundamental idea of a memory hierarchy:
 - For each level k , the faster, smaller device at level k serves as a cache for the larger, slower device at level $k+1$
- ❖ Why do memory hierarchies work?
 - Because of locality, programs tend to access the data at level k more often than they access the data at level $k+1$
 - Thus, the storage at level $k+1$ can be slower, and thus larger and cheaper per bit
- ❖ **Big Idea:** The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top

6



UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Making memory accesses fast!

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

8

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Cache Organization (1)

Note: The textbook uses "B" for block size

- ❖ Block Size (K): unit of transfer between \$ and Mem
 - Given in bytes and always a power of 2 (e.g. 64 B)
 - Blocks consist of adjacent bytes (differ in address by 1)
 - Spatial locality!

9

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Cache Organization (1)

Note: The textbook uses "b" for offset bits

- ❖ Block Size (K): unit of transfer between \$ and Mem
 - Given in bytes and always a power of 2 (e.g. 64 B)
 - Blocks consist of adjacent bytes (differ in address by 1)
 - Spatial locality!
- ❖ Offset field
 - Low-order $\log_2(K) = k$ bits of address tell you which byte within a block
 - $(\text{address}) \bmod 2^k = n$ lowest bits of address
 - $(\text{address}) \bmod (\# \text{ of bytes in a block})$

m -bit address:

$m - k$ bits	k bits
Block Number	Block Offset

(refers to byte in memory)

10

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Cache Organization (2)

- ❖ Cache Size (C): amount of *data* the \$ can store
 - Cache can only hold so much data (subset of next level)
 - Given in bytes (C) or number of blocks (C/K)
 - Example: $C = 32 \text{ KiB} = 512 \text{ blocks}$ if using 64-B blocks
- ❖ Where should data go in the cache?
 - We need a mapping from memory addresses to specific locations in the cache to make checking the cache for an address **fast**
- ❖ What is a data structure that provides fast lookup?
 - Hash table!

11

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Review: Hash Tables for Fast Lookup

Insert:

5	
1	
27	
34	
102	
119	

Apply hash function to map data to "buckets"

12

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Place Data in Cache by Hashing Address

Memory: Block Num, Block Data

Cache: Index, Block Data

Here $K = 4$ B and $C/K = 4$

- Map to *cache index* from block address
 - Use next $\log_2(C/K) = s$ bits
 - $(\text{block address}) \bmod (\# \text{ blocks in cache})$

13

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Place Data in Cache by Hashing Address

Memory: Block Addr, Block Data

Cache: Index, Block Data

Here $K = 4$ B and $C/K = 4$

- Map to *cache index* from block address
 - Lets adjacent blocks fit in cache simultaneously!
 - Consecutive blocks go in consecutive cache indices

14

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Place Data in Cache by Hashing Address

Memory: Block Addr, Block Data

Cache: Index, Block Data

Here $K = 4$ B and $C/K = 4$

Collision!

- This might confuse the cache later when we access the data
- Solution?

15

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Tags Differentiate Blocks in Same Index

Memory: Block Addr, Block Data

Cache: Index, Tag, Block Data

Here $K = 4$ B and $C/K = 4$

- Tag = rest of address bits
 - t bits = $m - s - k$
 - Check this during a cache lookup

16

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Checking for a Requested Address

- CPU sends address request for chunk of data
 - Address and requested data are not the same thing!
 - Analogy: your friend \neq his or her phone number
- TIO address breakdown:

Tag (t)	Index (s)	Offset (k)
Block Number		

 - Index field tells you where to look in cache
 - Tag field lets you check that data is the block you want
 - Offset field selects specified start byte within block
- Note: t and s sizes will change based on hash function

17

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Cache Puzzle #1

Vote at <http://PollEv.com/justin>

- Based on the following behavior, which of the following block sizes is NOT possible for our cache?
 - Cache starts *empty*, also known as a *cold cache*
 - Access (addr: hit/miss) stream:
 - (14: miss), (15: hit), (16: miss)

A. 4 bytes
 B. 8 bytes
 C. 16 bytes
 D. 32 bytes
 E. We're lost...

18

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Direct-Mapped Cache

Block Addr | Block Data

00	00
00	01
00	10
00	11
01	00
01	01
01	10
01	11
10	00
10	01
10	10
10	11
11	00
11	01
11	10
11	11

Index | Tag | Block Data

00	00	
01	11	
10	01	
11	01	

Here $K = 4$ B and $C/K = 4$

- Hash function: $(\text{block address}) \bmod (\# \text{ of blocks in cache})$
 - Each memory address maps to *exactly* one index in the cache
 - Fast (and simpler) to find an address

19

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Direct-Mapped Cache Problem

Block Addr | Block Data

00	00
00	01
00	10
00	11
01	00
01	01
01	10
01	11
10	00
10	01
10	10
10	11
11	00
11	01
11	10
11	11

Index | Tag | Block Data

00	??	
01	??	
10		
11	??	

Here $K = 4$ B and $C/K = 4$

- What happens if we access the following addresses?
 - 8, 24, 8, 24, 8, ...?
 - Conflict in cache (misses!)
 - Rest of cache goes *unused*
- Solution?

20

UNIVERSITY of WASHINGTON L17: Caches II CSE351, Autumn 2017

Associativity

- What if we could store data in any place in the cache?
 - More complicated hardware = more power consumed, slower
- So we *combine* the two ideas:
 - Each address maps to exactly one **set**
 - Each set can store block in more than one **way**

1-way: 8 sets, 1 block each

2-way: 4 sets, 2 blocks each

4-way: 2 sets, 4 blocks each

8-way: 1 set, 8 blocks

direct mapped

fully associative

21