Caches I

CSE 351 Summer 2024

Instructor: Ellis Haker

Teaching Assistants:

Naama Amiel Micah Chang Shananda Dokka Nikolas McNamee Jiawei Huang



Administrivia

- Today
 - HW12 due (11:59pm)
 - Lab2 due (11:59pm)
- Monday, 7/22
 - RD15 due (1pm)
 - HW13 due (11:59pm)
 - Quiz 2 released (11:59pm)
 - Same rules apply as Quiz 1
- Wednesday, 7/24
 - RD16 due (1pm)
 - HW14 due (11:59pm)

Quiz 1 Grades Released!

- Regrade requests will open tonight at 11:59pm
- If you received a message about possible academic misconduct, please email me
- And even though it's cheesy...
 - Your success in life is not defined by grades
 - You are not defined by grades
 - We know it seems critically important right now, but we promise, the numbers on a transcript will fade with time.

Topic Group 3: Scale & Coherence

- How do we make memory accesses faster?
- How do programs manage large amounts of memory?
- How does your computer run multiple programs at once?

Starting with caches, which are implemented in hardware.



Vil: - 2"= 1024 Aside: Units and Prefixes \Li\o= Louc

- Traditional prefixes represent powers of 10, we define new ones for base 2 Ex: 1 **Kibi**byte = 2^{10} bytes $\approx 10^3$ bytes = 1 **Kilo**byte Ο
- SI prefixes are *ambiguous* if base 10 or base 2 (does 'k' stand for kilo or kibi?) (people often son "kild" when they "kild"
- IEC prefixes are *unambiguously* base 2

SI Symbol	SI Prefix	SI size	IEC symbol	IEC Prefix	IEC Size
К	Kilo-	10 ³	Ki	Kibi-	2 ¹⁰
Μ	Mega-	10 ⁶	Mi	Mebi-	2 ²⁰
G	Giga-	10 ⁹	Gi	Gibi-	2 ³⁰
Т	Tera-	10 ¹²	Ti	Tebi-	240
Р	Peta-	10 ¹⁵	Pi	Pebi-	2 ⁵⁰
E	Exa-	10 ¹⁸	Ei	Exbi-	2 ⁶⁰
Z	Zetta-	10 ²¹	Zi	Zebi-	2 ⁷⁰
Y	Yotta-	10 ²⁴	Yi	Yobi-	2 ⁸⁰

How to Remember?

- You can always look it up :)
 - It's on the midterm reference sheet

- Mnemonics
 - Killer Mechanical Giraffe Teaches Pet Extinct Zebra to Yodel
 - Kirby Missed Ganondorf Terribly, Potentially Exterminating Zelda and Yoshi
 - From xkcd: Karl Marx Gave The Proletariat Eleven Zeppelins, Yo
 - https://xkcd.com/992/

Review Questions

- 1. Convert the following to or from IEC: a. 512 Mi-students $N_{i=2}^{2^{n}} \le \sqrt{2^{n}} \cdot \sqrt{2^{n}$
- Compute the average memory access time (AMAT) for a system with the following properties:
 Am At = HT + M+·M?
 - a. Hit time of 2 ns
 - b. Miss rate of 1%
 - c. Miss penalty of 300 ns

: ZNS + 0.01-30 Gus = ZNS + 3 NS

How does execution time grow with SIZE?





Actual Data



An Analogy





Caches

- Cache basics
- Principle of locality
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches



A Very Silly Analogy



A Very Silly Analogy (pt 2)





- Pronounced "cash"
 - Often abbreviated to '\$"
- <u>English</u>: hidden storage space for provisions, weapons, or treasures
- <u>Computer:</u> Memory with short access time used for the storage of frequently or recently used instructions or data
 - I-cache for instructions
 - o d-cache for data
 - More generally: Used to optimize data transfers between any system elements with different characteristics (network interface cache, I/O cache, etc.)

If caches are so much faster, why do we need memory?

- Two common memory technologies others exist too, but these are the most common right now
 - DRAM: high-capacity, cheap, energy efficient, but slow
 - **SRAM**: much faster, but less energy efficient and *expensive*
- We can't afford to have all our computer's memory be SRAM
 - Use DRAM to provide large amounts of memory for cheap
 - Have a small SRAM cache for speed



DRAM

SRAM

General Cache Mechanics

- Memory
 - Slower, larger, cheaper
 - Partition into "blocks"
- Cache
 - Smaller, faster, more expensive
 - Stores a subset of blocks from memory
 - Data is copied in *block-sized chunks*



Cache Mechanics: Hit

- Data we need is in block b
- Block *b* is in the cache
 - *Hit!*
- Data is returned to the CPU



Cache Mechanics: Miss

- Data we need is in block b
- Block b is not in the cache already
 - Miss!
- Block *b* is fetched from memory
- Block b is written into the cache



Cache Mechanics: Miss (pt 2)

- Data we need is in block *b*
- Block b is not in the cache already
 - Miss!
- Block *b* is fetched from memory
- Block b is written into the cache
 - Placement policy decides where it goes
 - Replacement policy decides what we kick out
- Data is returned to the CPU



Why Caches Work

- Takes advantage of locality
 - Common patterns in how programs access data
- Temporal locality
 - If a program accesses data once, it's likely to access it again
- Spatial locality
 - If a program accesses some data, it's likely to access other data that's *nearby* in memory

How do caches take advantage of this?





Data

- **Temporal**: sum, i, and n are accessed every loop iteration
- Spacial: consecutive elements of a

Instructions

- Temporal: loop body code
- Instructions: instructions executed in sequence



Layout in Memory

a	a	a	а	а	a	a	а	а	a	а	а
[0]	[0]	[0]	[0]	[1]	[1]	[1]	[1]	[2]	[2]	[2]	[2]
[0]	[1]	[2]	[3]	[0]	[1]	[2]	[3]	[0]	[1]	[2]	[3]

	1)	a[0][0]
Access pattern:	2)	a[0][1]
	3)	a[0][2]
	4)	a[0][3]
Stride? 1	5)	a[1][0]
-	6)	a[1][1]
6002 course performance.		a[1][2]
moving on to the next	8)	a[1][3]
	9)	a[2][0]
a	10)	a[2][1]
[2]	11)	a[2][2]
[3]	12)	a[2][3]



Layout in Memory

a	a	a	a	а	a	a	a	а	a	a	а
[0]	[0]	[0]	[0]	[1]	[1]	[1]	[1]	[2]	[2]	[2]	[2]
[0]	[1]	[2]	[3]	[0]	[1]	[2]	[3]	[0]	[1]	[2]	[3]
_ ر			_		トノ	_		2	ア		

Access nottorn.	1)	a[0][0]
Access patient:	2)	a[1][0]
	3)	a[2][0]
	4)	a[0][1]
Stride?	5)	a[1][1]
	6)	a[2][1]
Pri rolumence X	7)	a[0][2]
	8)	a[1][2]
	9)	a[2][2]
	10)	a[0][3]
	11)	a[1][3]
	12)	a[2][3]

Back cache performance i



really back cache preformance !

Cache Performance Metrics

- Miss Rate (MR)
 - Fraction of memory references not found in cache
 - o (misses / accesses) = 1 Hit Rate
- Hit Time (HT)
 - Time to deliver a block in the cache to the processor
 - Includes time to determine whether the block is in the cache
- Miss Penalty (MP)
 - Additional time required because of a miss
 - Total miss time = Hit Time + Miss Penalty

Cache Performance

- Average Memory Access Time (AMAT): average time to access data, considering both cache hits and misses
 - AMAT = Hit time + Miss Rate × Miss Penalty (abbreviated AMAT = HT + MR×MP)

Practice Questions

- Processor specs: 200 ps clock, MP of 50 clock cycles, MR of 0.02 misses/instruction, and HT of 1 clock cycle
 A) What is the AMAT (in clock cycles)? HT + MP · MP = 1 + 0.02 · 50 = 1+1 = 2-404
- 2. Which of these improvements would result in the lowest AMAT?
 - A) 190 ps clock some # of yoles = 2.190 = 380 ps
 - B) Miss penalty of 40 clock cycles $1+c_0 \cdot 2 \cdot 40 = 1 \cdot 8$ cycles $200 \frac{9}{3}$ cycle = 360 pc C) Miss rate of 0.015 misses/instruction

1+0.015.50 = 1.75 cycles - 200 PS/cycle = 350ps

personnt of time for CPU to do I took

Cache Performance (pt 2)

- Misses have a *much* larger effect on AMAT than hits!
 - Going to memory could be 100x slower than accessing the cache (measured in clock cycles)
- High miss rate or miss penalty hurt AMAT the most

<u>Ex</u>: Assume HT of 1 clock cycle and MP of 100 clock cycles MP = 1 - HPIf HR = 99%, AMAT = <u>Zycles</u> MP = 0.01, $1 + 0.01 \cdot 100 = Z$

If HR = 97%, AMAT = 4 agens MP: 6 03, 1+0.03:100 = 4

Increasing hit rate by only 2% cut acces time in half!

Can we have more than one cache? Yes!

- Why would we want to do that?
 - Avoid going to memory at all costs!
- Typical performance numbers
 - Miss Rate
 - L1: 3-10%
 - L2: very small (likely <1%)
 - Hit Time
 - L1: 4 clock cycles
 - L2: 10 clock cycles
 - Miss Penalty
 - 50-200 cycles for missing in L2 going to main memory)
 - Trend: increasing!

Here workers are outdeted! -trade secreta

An Example Memory Hierarchy



Learning About Your Machine

- Linux:
 - lscpu
 - o ls /sys/devices/system/cpu/cpu0/cache/index0/
 - Example: cat /sys/devices/system/cpu/cpu0/cache/index*/size

• Windows:

- wmic memcache get <query> (all values in KB)
- Example: wmic memcache get MaxCacheSize
- Modern processor specs: <u>http://www.7-cpu.com/</u>

Summary

Memory Hierarchy

- Higher levels are faster, smaller, and more expensive
 - Contain the most used data from lower levels
- Exploits temporal and spatial locality
- Caches are intermediate levels between memory and the CPU
- Cache Performance
 - Ideal case: data found in cache (hit)
 - Bad case: not found in cache (miss), go to next level in hierarchy
 - Average Memory Access Time (AMAT) = HT + MR × MP
 - Hurt my high miss rate and miss penalty