# Memory & Caches I

CSE 351 Autumn 2024

#### Instructor:

**Ruth Anderson** 

#### **Teaching Assistants:**

Alexandra Michael **Connie Chen** Chloe Fong Chendur Jayavelu Joshua Tan Nikolas McNamee Nahush Shrivatsa Naama Amiel Neela Kausik **Renee Ruan** Rubee Zhao Samantha Dreussi Sean Siddens Waleed Yagoub



**Alt text:** I looked at some of the data dumps from vulnerable sites, and it was ... bad. I saw emails, passwords, password hints. SSL keys and session cookies. Important servers brimming with visitor IPs. Attack ships on fire off the shoulder of Orion, c-beams glittering in the dark near the Tannhäuser Gate. I should probably patch OpenSSL.

http://xkcd.com/1353/

# **Relevant Course Information**

- HW14 due TONIGHT, Wednesday (10/30) @ 11:59 pm
- No Lecture on Fri 11/01 (No HW/Reading due)
  - HW15 not due until next Monday (11/04)
  - HW16 de next Wed (11/06)
- Midterm Exam: See Ed Post and exams page
  - Take home, on Gradescope
  - Open: Thursday 10/31 at 5pm; Due: Saturday 11/02 at 11:59pm
  - Review in section this week (10/31)
- Lab 3 due Mon 11/11
  - Encouraged to aim for Fri 11/08, actual deadline Mon 11/11
  - You have everything you need to do the lab as of last lecture!
  - Last part of HW15 (due Mon 11/04) is useful for Lab 3

# Aside: Units and Prefixes (Review)

- Here focusing on large numbers (exponents > 0)
- Note that  $10^3 \approx 2^{10}$
- SI prefixes are *ambiguous* if base 10 or 2
- IEC prefixes are *unambiguously* base 2

SI Size	Prefix	Symbol	IEC Size	Prefix	Symbol
10 <sup>3</sup>	Kilo-	K	2 <sup>10</sup>	Kibi-	Ki
10 <sup>6</sup>	Mega-	М	2 <sup>20</sup>	Mebi-	Mi
109	Giga-	G	2 <sup>30</sup>	Gibi-	Gi
1012	Tera-	Т	2 <sup>40</sup>	Tebi-	Ti
10 <sup>15</sup>	Peta-	Р	2 <sup>50</sup>	Pebi-	Pi
10 <sup>18</sup>	Exa-	E	2 <sup>60</sup>	Exbi-	Ei
10 <sup>21</sup>	Zetta-	Z	2 <sup>70</sup>	Zebi-	Zi
10 <sup>24</sup>	Yotta-	Y	2 <sup>80</sup>	Yobi-	Yi

SIZE PREFIXES (10<sup>x</sup> for Disk, Communication; 2<sup>x</sup> for Memory)

## How to Remember?

- Will be given to you on Final reference sheet
- Mnemonics
  - There unfortunately isn't one well-accepted mnemonic
    - But that shouldn't stop you from trying to come with one!
  - Killer Mechanical Giraffe Teaches Pet, Extinct Zebra to Yodel
  - Kirby Missed Ganondorf Terribly, Potentially Exterminating Zelda and Yoshi
  - xkcd: Karl Marx Gave The Proletariat Eleven Zeppelins, Yo
    - <u>https://xkcd.com/992/</u>
  - Post your best on Ed Discussion!

# **Reading Review**

- Terminology:
  - Caches: cache blocks, cache hit, cache miss
  - Principle of locality: temporal and spatial
  - Average memory access time (AMAT): hit time, miss penalty, hit rate, miss rate

# **Review Questions**

- Convert the following to or from IEC:
  - 512 Ki-books
  - 2<sup>27</sup> cats
- Compute the average memory access time (AMAT) for the following system properties:
  - Hit time of 1 ns
  - Miss rate of 1%
  - Miss penalty of 100 ns

#### How does execution time grow with SIZE?

```
int array[SIZE];
int sum = 0;
for (int i = 0; i < 200000; i++) {
  for (int j = 0; j < SIZE; j++) {
    sum += array[j];
  }
}</pre>
```



#### **Actual Data**



SIZE

# Making memory accesses fast!

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

#### **Processor-Memory Gap**



# **Problem: Processor-Memory Bottleneck**



#### **Problem: lots of waiting on memory**

cycle: single machine step (fixed-time)

# **Problem: Processor-Memory Bottleneck**



cycle: single machine step (fixed-time)



- Pronunciation: "cash"
  - We abbreviate this as "\$"
- <u>English</u>: A hidden storage space for provisions, weapons, and/or treasures
- <u>Computer</u>: Memory with short access time used for the storage of frequently or recently used instructions (i-cache/I\$) or data (d-cache/D\$)
  - More generally: Used to optimize data transfers between any system elements with different characteristics (network interface cache, I/O cache, etc.)

# **General Cache Mechanics (Review)**



## **General Cache Concepts: Hit (Review)**



Data in block b is needed

Block b is in cache: Hit!

Data is returned to CPU

## General Cache Concepts: Miss (Review)



Data in block b is needed

Block b is not in cache: Miss!

Block b is fetched from memory

#### Block b is stored in cache

- Placement policy: determines where b goes
- Replacement policy: determines which block gets evicted (victim)

Data is returned to CPU

# Why Caches Work (Review)

 Locality: Programs tend to use data and instructions with addresses <u>near or equal to those they have used</u> <u>recently</u>

# Why Caches Work (Review)

- Locality: Programs tend to use data and instructions with addresses <u>near or equal to those they have used</u> <u>recently</u>
- *Temporal* locality:



 Recently referenced items are *likely* to be referenced <u>again</u> in the near future

# Why Caches Work (Review)

- Locality: Programs tend to use data and instructions with addresses <u>near or equal to those they have used</u> <u>recently</u>
- Temporal locality:
  - Recently referenced items are *likely* to be referenced <u>again</u> in the near future
- Spatial locality:
  - Items with <u>nearby addresses</u> tend to be referenced <u>close together in time</u>
- How do caches take advantage of this?



# **Example: Any Locality?**

```
sum = 0;
for (i = 0; i < n; i++)
{
    sum += a[i];
}
return sum;</pre>
```

#### Data:

- Temporal: sum referenced in each iteration
- Spatial: consecutive elements of array a [] accessed

#### Instructions:

- <u>Temporal</u>: cycle through loop repeatedly
- Spatial: reference instructions in sequence

#### **Locality Example #1 Code**

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}</pre>
```

## **Locality Example #1**

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;
    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}</pre>
```

	Layout in Memory											
	a [0] [0]	a [0] [1]	a [0] [2]	a [0] [3]	a [1] [0]	a [1] [1]	a [1] [2]	a [1] [3]	a [2] [0]	a [2] [1]	a [2] [2]	a [2] [3]
7	6			9	2			1(	8			

Note: 76 is just one possible starting address of array a

M = 3, N=4								
a[0][0]	a[0][1]	a[0][2]	a[0][3]					
a[1][0]	a[1][1]	a[1][2]	a[1][3]					
a[2][0]	a[2][1]	a[2][2]	a[2][3]					

Access Pattern
stride = ?

1)	a[0]	[0]
2)	a[0]	[1]
3)	a[0]	[2]
4)	a[0]	[3]
5)	a[1]	[0]
6)	a[1]	[1]
7)	a[1]	[2]
8)	a[1]	[3]
9)	a[2]	[0]
10)	a[2]	[1]
11)	a[2]	[2]
12)	a[2]	[3]

#### **Locality Example #2 Code**

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

## **Locality Example #2**

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;
    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}</pre>
```

# M = 3, N=4a[0][0]a[0][1]a[0][2]a[0][3]a[1][0]a[1][1]a[1][2]a[1][3]a[2][0]a[2][1]a[2][2]a[2][3]

Access Pattern: stride = ?



#### Layout in Memory

	a	a	a	a	a	a	a	a	a	a	a	a
	[0]	[0]	[0]	[0]	[1]	[1]	[1]	[1]	[2]	[2]	[2]	[2]
	[0]	[1]	[2]	[3]	[0]	[1]	[2]	[3]	[0]	[1]	[2]	[3]
7	<b>† †</b> 76 92				2			10	8			

# **Locality Example #3 Code**

```
int sum_array_3D(int a[M][N][L])
{
    int i, j, k, sum = 0;
    for (i = 0; i < N; i++)
        for (j = 0; j < L; j++)
            for (k = 0; k < M; k++)
                sum += a[k][i][j];
    return sum;
}</pre>
```

- What is wrong with this code?
- How can it be fixed?



#### **Locality Example #3**

```
int sum_array_3D(int a[M][N][L])
{
    int i, j, k, sum = 0;
    for (i = 0; i < N; i++)
        for (j = 0; j < L; j++)
            for (k = 0; k < M; k++)
                sum += a[k][i][j];
    return sum;
}</pre>
```

- What is wrong with this code?
- How can it be fixed?

Layout in Memory (M = ?, N = 3, L = 4)



# **Cache Performance Metrics (Review)**

- Huge difference between a cache hit and a cache miss
  - Could be 100x speed difference between accessing cache and main memory (measured in *clock cycles*)
- Miss Rate (MR)
  - Fraction of memory references not found in cache (misses / accesses) = 1 - Hit Rate
- Hit Time (HT)

- Miss takes HI+MP
- 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

# **Cache Performance (Review)**

- Two things hurt the performance of a cache:
  - Miss rate and miss penalty
- Average Memory Access Time (AMAT): average time to access memory considering both hits and misses
   AMAT = Hit time + Miss rate × Miss penalty (abbreviated AMAT = HT + MR × MP)
- 99% hit rate twice as good as 97% hit rate!
  - Assume HT of 1 clock cycle and MP of 100 clock cycles
  - 97%: AMAT =
  - 99%: AMAT =

# **Practice Question**

 Processor specs: 200 ps clock, MP of 50 clock cycles, MR of 0.02 misses/instruction, and HT of 1 clock cycle

AMAT =

- Which improvement would be best?
   A. 190 ps clock
  - **B.** Miss penalty of 40 clock cycles
  - **C.** MR of 0.015 misses/instruction

# Can we have more than one cache?

- Why would we want to do that?
  - Avoid going to memory!
- Typical performance numbers:
  - Miss Rate
    - L1 MR = 3-10%
    - L2 MR = Quite small (*e.g.* < 1%), depending on parameters, etc.
  - Hit Time
    - L1 HT = 4 clock cycles
    - L2 HT = 10 clock cycles
  - Miss Penalty
    - P = 50-200 cycles for missing in L2 & going to main memory
    - Trend: increasing!

# **An Example Memory Hierarchy**



# Summary

- Memory Hierarchy
  - Successively higher levels contain "most used" data from lower levels
  - Makes use of temporal and spatial locality
  - Caches are intermediate storage levels used to optimize data transfers between any system elements with different characteristics
- Cache Performance
  - Ideal case: found in cache (hit)
  - Bad case: not found in cache (miss), search in next level
  - Average Memory Access Time (AMAT) = HT + MR × MP
    - Hurt by Miss Rate and Miss Penalty