

CS250: FOUNDATIONS OF COMPUTER SYSTEMS

[MEMORY]

Caching and Locality

Access a data item once

Likely, then, you are to access it again

Later in time (temporality)

Access a data item once

The one close to it will likely be accessed next

Spatial proximity

Caching's secret sauce? locality

Spatial, temporal or

Some combination thereof

SHRIDEEP PALLICKARA

Computer Science

Colorado State University

COMPUTER SCIENCE DEPARTMENT



COLORADO STATE UNIVERSITY

1

Frequently asked questions from the previous class survey

- Does each byte of memory have a physically numbered address?
- Would a 128-bit subsystem be faster or have more throughput?
- Why is the cache real estate so big on the chip?
- Why is it impossible to add more registers?



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.2

2

Topics covered in this lecture

- Data movements
- Caching
 - ▣ Direct mapped
 - ▣ Associative
 - ▣ N-way associativity
- A real-world example



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.3

3



4

Moving data between the registers and memory is strictly a program function

- The program loads data into registers and stores register data into memory using machine instructions like `mov`
- It is the programmer's or **compiler's responsibility** to keep heavily referenced data in the registers as long as possible
 - The CPU *will not* automatically place data in general-purpose registers in order to achieve higher performance



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.5

5

Program responsibility and obliviousness [1/2]

- Programs explicitly control access to registers, main memory, and memory-hierarchy subsystems that are at the file storage level or lower
- Programs are **largely unaware** of the memory *hierarchy between* the register level and main memory



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.6

6

Program responsibility and obliviousness [2/2]

- In particular, cache accesses are **transparent** to the program
 - Access to these levels of the memory hierarchy usually occurs *without any intervention* on a program's part
- Programs simply access main memory, and the hardware and operating system take care of the rest
 - Really? Yes!!!



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.7

7

Responsibility for data movements

- If a program always accesses main memory, it will run slowly
 - Modern DRAM main-memory subsystems are **much slower** than the CPU
- Cache memory subsystems and the CPU's cache controller **move data** between main memory and the L1, L2, and L3 caches
 - So that the CPU can quickly access oft-requested data
- Likewise, it is the **virtual memory** subsystem's responsibility to move oft-requested data from hard disk to main memory



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.8

8

Mechanics of transparent data accesses

- With few exceptions, most memory subsystem accesses take place **transparently between**
 - One level of the memory hierarchy and
 - the *level immediately below or above it*



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.9

9

For example, the CPU rarely accesses main memory directly

- Instead, when the CPU requests data from memory, the L1 cache subsystem takes over
- If the requested data is in the cache
 - **Cache hit**
 - The L1 cache subsystem returns the data to the CPU, and that concludes the memory access
- If the requested data is not in the cache?
 - **Cache miss!**



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.10

10

If the requested data isn't present in the L1 cache

...

- The L1 cache subsystem passes the request down to the L2 cache subsystem
- If the L2 cache subsystem has the data?
 - ▣ The L2 Cache returns this data to the L1 cache, which then returns the data to the CPU
- Requests for the same data *in the near future* will be fulfilled by the L1 cache rather than the L2 cache
 - ▣ Because the L1 cache **now has a copy** of the data



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.11

11

What if the L2 cache does not have it?

- After the L2 cache, the L3 cache kicks in
- If none of the L1, L2, or L3 cache subsystems have a copy of the data, the request goes to main memory
 - ▣ If the data is found in main memory, the main-memory subsystem passes it to the L3 cache
 - The L3 cache then passes it to the L2 cache, which then passes it to the L1 cache
 - The L1 cache then passes it to the CPU
- Once again, the data is now in the L1 cache, so any requests for this data *in the near future* will be fulfilled by the L1 cache



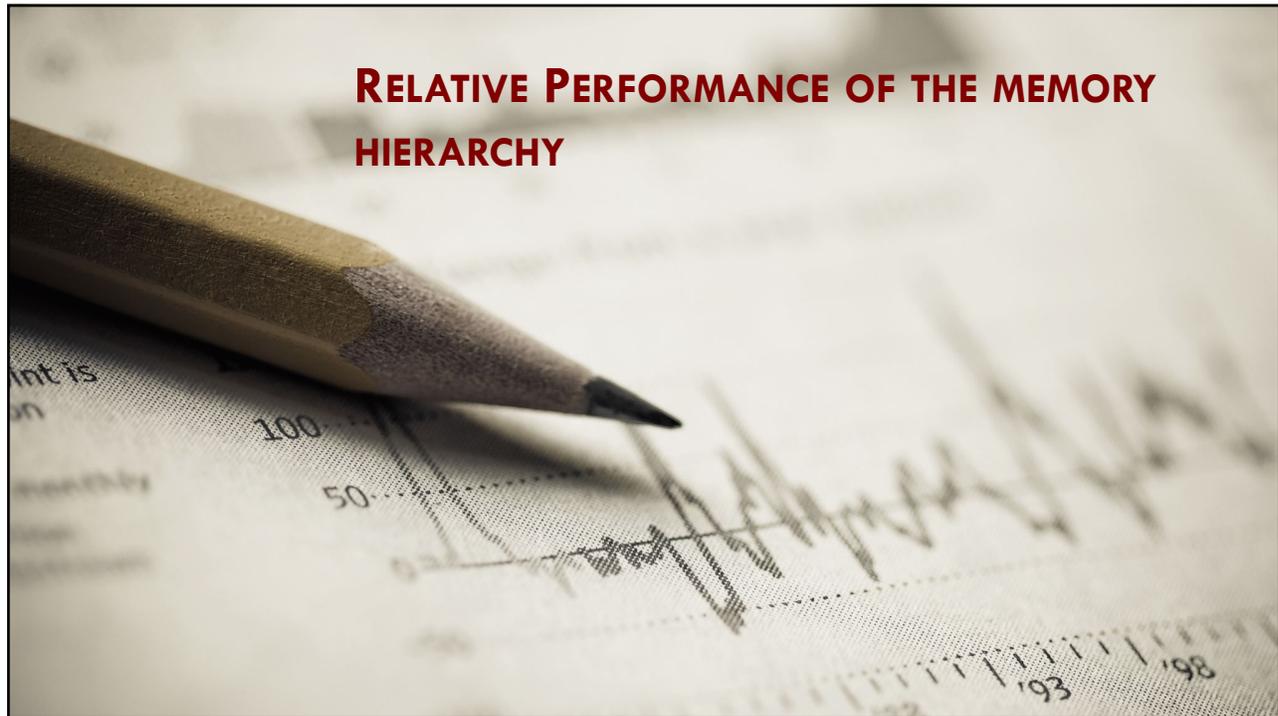
COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.12

12



13

Relative performance of the memory system

- Registers are, unquestionably, the best place to store data you need to access quickly
 - ▣ Accessing a register **never requires any extra time**, and
 - ▣ Most machine instructions that access data can access register data
- The difference in speed between the L1, L2, and L3 cache systems isn't so dramatic unless the secondary or tertiary cache is not packaged together on the CPU



14

There are several reasons why L2 cache accesses are slower than L1 accesses

- It takes the CPU time to **determine** that the data it's seeking is not in the L1 cache
 - By the time it does that, the memory access cycle is nearly complete, and there's no time to access the data in the L2 cache
- The circuitry of the L2 cache may be **slower than the circuitry of the L1** cache in order to make the L2 cache less expensive
- L2 caches are usually 16 to 64 times larger than L1 caches
 - Larger memory subsystems **tend to be slower** than smaller ones
 - All this amounts to additional wait states for accessing data in the L2 cache



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.15

15

A similar performance gulf separates the L2 and L3 caches and L3 and main memory

- Main memory is typically one order of magnitude slower than the L3 cache; L3 accesses are much slower than L2 access
- To speed up access to adjacent memory objects, the L3 cache reads data from main memory in **cache lines**
- Likewise, L2 cache reads cache lines from L3



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.16

16



17

Up until this point

- We have treated the cache as a magical place that
 - ▣ Automatically stores data when we need it, perhaps fetching new data as the CPU requires it
- But **how** exactly does the cache do this?
 - ▣ And **what happens when it is full**, and the CPU is requesting additional data that's not there?



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.18

18

Programs access only a small amount of data at a given time

- A cache that is sized accordingly will improve their performance
- Unfortunately, the data that programs want **rarely sits in contiguous memory location**
 - It's usually *spread out* all over the address space



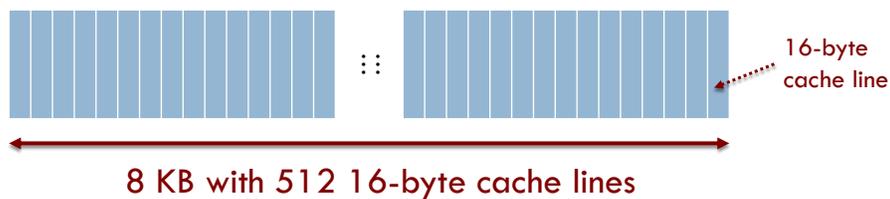
Cache design considerations ...

- Cache design must account for the fact that the cache must map data objects at **widely varying addresses** in memory
- Cache memory is not organized in a single group of bytes
 - Instead, it's usually organized in blocks of **cache lines**
 - Each line containing some number of bytes
 - Typically, a **small power of 2**: like 16, 32, or 64



Cache lines

- For example, an 8 KB cache line is often organized as a set of 512 cache lines of 16 bytes each



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.21

21

We can attach a different noncontiguous address to each of the cache lines

- Cache line 0 might correspond to addresses 0x10000 through 0x1000F
- Cache line 1 might correspond to addresses 0x21400 through 0x2140F



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALICKARA
COMPUTER SCIENCE DEPARTMENT

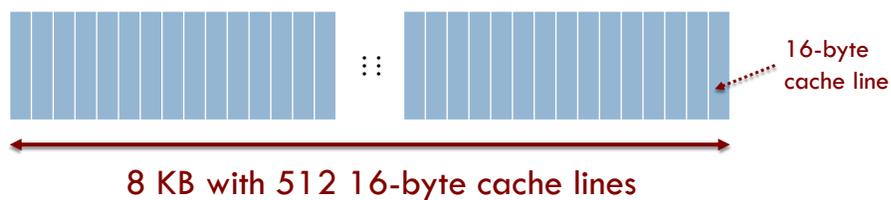
MEMORY

L10.22

22

Generally, if a cache line is n bytes long

- It will hold n bytes from main memory that fall on an **n -byte boundary**
- In our example of 16-byte cache lines, a cache line holds blocks of 16 bytes whose addresses fall on 16-byte boundaries in main memory
 - i.e., the least-significant 4 bits of the address of the first byte in the cache line are always 0



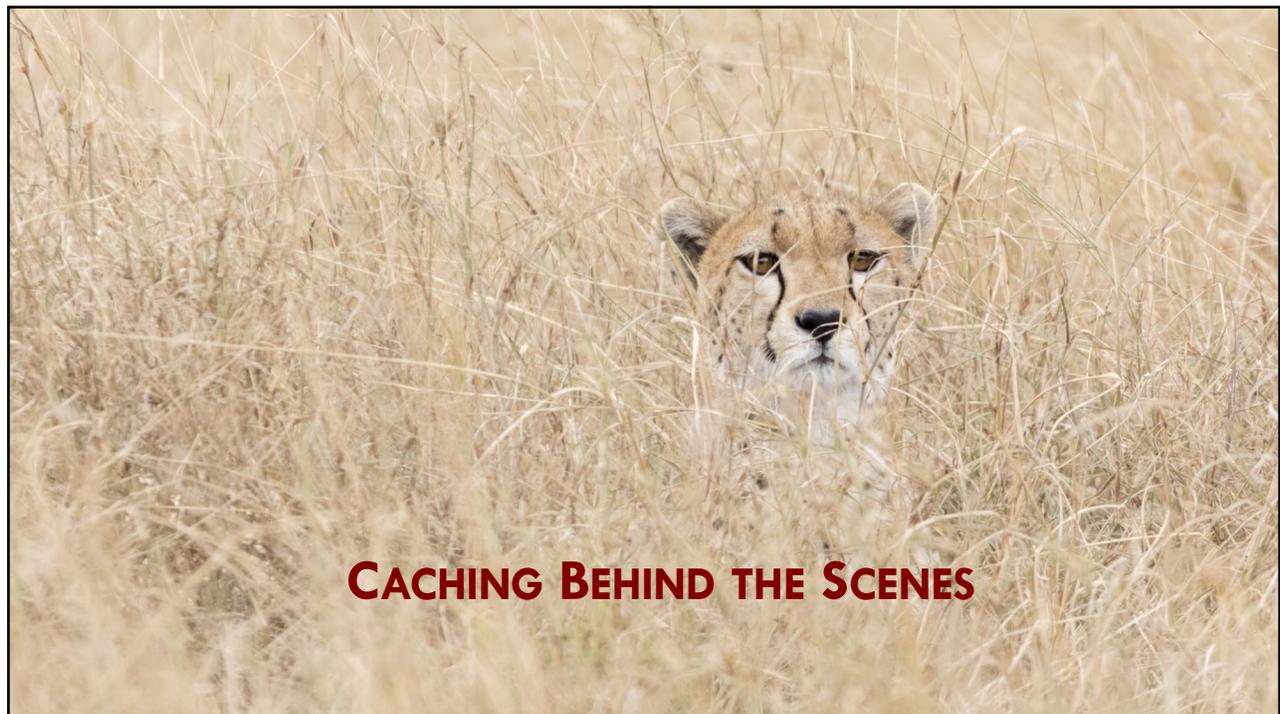
COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.23

23



24

Types of caches

- Direct mapped caches
- Fully associative caches
- N-way associative caches



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.25

25

A **direct mapped cache** is also known as a one-way associative cache

- In a direct-mapped cache, a particular block of main memory is always loaded into (*i.e.*, mapped to) the exact same cache line
- This mapping is **determined by a small number of bits** in the data block's memory address



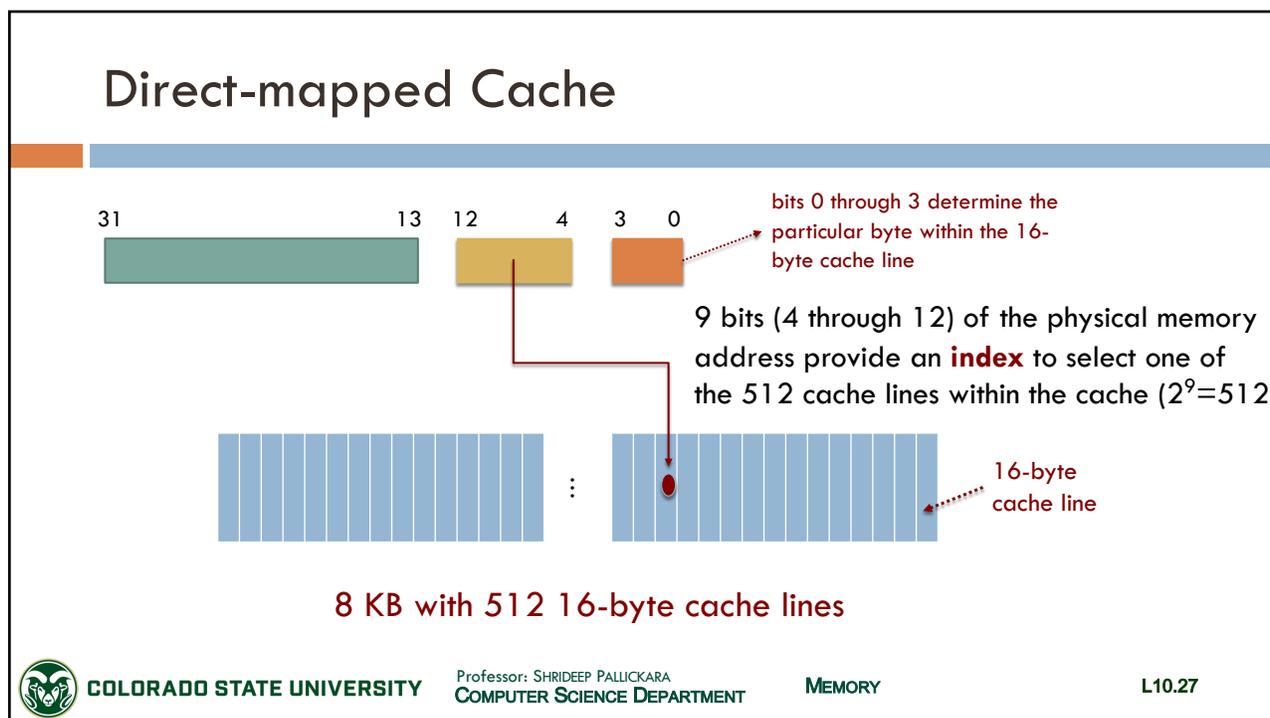
COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.26

26



27

Problems with a direct mapped cache

- Two different memory addresses located on 8KB boundaries **cannot both appear simultaneously** in the cache
- How many such addresses exist in our 32-bit system?
 - 2^{19} 8KB blocks exist in our system
 - 2^{19} 512 (2^9) blocks of 16-bytes (2^4) each
 - $2^{19} \cdot 2^9 \cdot 2^4 = 2^{32}$ (the size of the main memory in our example)

COLORADO STATE UNIVERSITY Professor: SHRIDEEP PALLICKARA COMPUTER SCIENCE DEPARTMENT MEMORY L10.28

28

The ideal world: A fully **associative** cache

- The cache controller can place a block of bytes in *any one* of the cache lines present in the cache memory
- While this is the most flexible cache system, the **extra circuitry** to achieve full associativity *is expensive* and, worse, *can slow down* the memory subsystem
- Most L1 and L2 caches are not fully associative for this reason



Trade-off space

- A fully associative cache is too complex, too slow, and too expensive to implement
- But a direct-mapped cache is too inefficient



A compromise: the **n-way associative cache**

- In an n-way set associative cache, the cache is broken up into sets of n cache lines
- The CPU determines the **particular set to use** based on
 - Some subset of the memory address bits, just as in the direct-mapping scheme, and ...
 - The cache controller uses a fully associative mapping algorithm to determine which one of the n cache lines within the set to use



31

For example, an 8KB two-way set associative cache subsystem with 16-byte cache lines [1/2]

- Organizes the cache into **256** cache-line sets with **two cache** lines each
- Eight bits from the memory address determine which one of these 256 different sets will contain the data
 - $2^8 = 256$
- Once the cache-line set is determined, the cache controller maps the block of bytes to one of the two cache lines within the set



32

For example, an 8KB two-way set associative cache subsystem with 16-byte cache lines [2/2]

- This means **two different memory addresses** located on 8KB boundaries (addresses having the same value in bits 4 through 11) can both appear simultaneously in the cache
- However, a **conflict** will occur if you attempt to access a **third** memory location at an address that is an even multiple of 8KB



COLORADO STATE UNIVERSITY

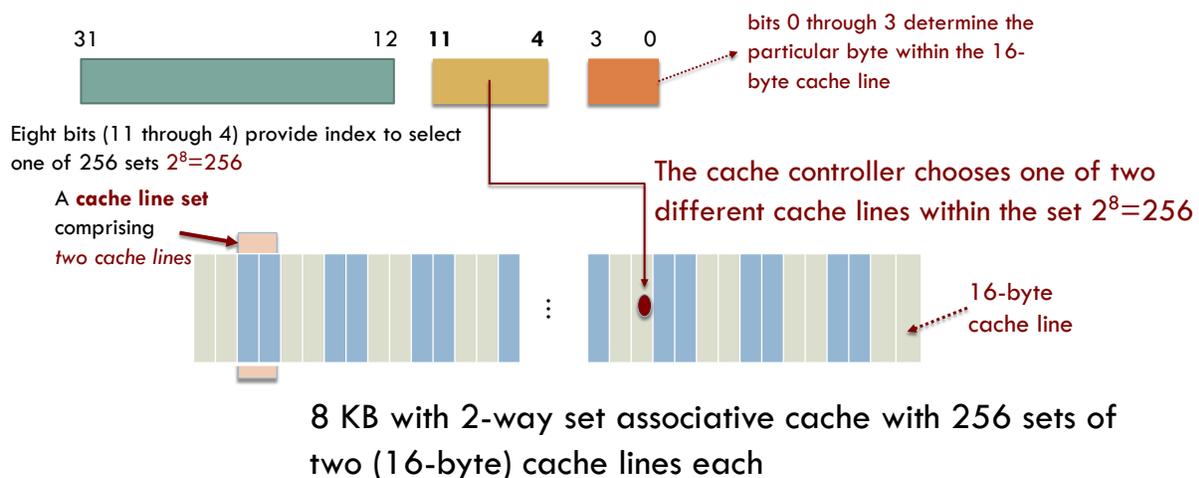
Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.33

33

An 8 KB two-way set associative cache subsystem with 16-byte cache lines



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.34

34

What if we have a 4-way associative cache

- A four-way set associative cache puts four associative cache lines in each cache-line set
- In our example, 8KB cache, a four-way set associative caching scheme would have **128 cache-line sets** with **four cache lines** each
- This would allow the cache to maintain up to four different blocks of data without a conflict, each of which would map to the same cache line in a direct-mapped cache



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.35

35

2-/4-way set associative vs direct mapped

- A 2- or 4-way set associative cache is
 - ▣ Much better than a direct-mapped cache and
 - ▣ Considerably less complex than a fully associative cache



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.36

36

Can we keep increasing the number of lines in each cache-line set?

- The more cache lines we have in each cache-line set, the closer we come to creating a fully associative cache
 - ▣ With all the attendant problems of complexity and speed
- Most cache designs are direct-mapped, two-way set associative, or four-way set associative
 - ▣ The various members of the 80x86 family make use of all three



WHY (MAIN) MEMORY MATTERS ...



An analogy: You at a government office

- Some interactions can be completed using your IDs/cards (in your wallet), documents (in your backpack), and documents (at home)
 - Items can be retrieved from the wallet in 2 seconds
 - The bag needs to be searched, and it takes about 120 seconds to do so
 - The trip home and back will take 36,000 seconds (or 10 hours)
- Average time to complete transaction if your wallet suffices 95% of the time but the backpack comes into play 5% of time?
 - $0.95 * (\text{wallet_time}) + 0.05 * (\text{backpack_time})$
 - $0.95 * 2 + 0.05 * 120 = 7.9$ seconds



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.39

39

An analogy: You at a government office

- Average time to complete transaction if your wallet suffices 95% of the time but the backpack comes into play 4% of time and you need to go home 1% of the time?
 - $0.95 * (\text{wallet_time}) + 0.04 * (\text{backpack_time}) + 0.01 * (\text{home_trip})$
 - $0.95 * 2 + 0.04 * 120 + 0.01 * 36,000 = 505.9$ seconds



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.40

40

Let's make it a little more real

- Cache access: 2 ns
- Main memory access: 50 ns
- Disk access: 8 milliseconds [8,000,000 ns]
- 97% cache and 3% main memory
 - $0.97 * 2 + 0.03 * 50 = 3.4$ ns
- 95% cache and 5% main memory
 - $0.95 * 2 + 0.05 * 50 = 4.4$ ns
- 95% cache, 4% main memory, and 1% disk
 - $0.95 * 2 + 0.04 * 50 + 0.01 * 8,000,000 = 80,003.9$ ns



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.41

41

The contents of this slide-set are based on the following references

- Jonathan E. Steinhart. *The Secret Life of Programs: Understand Computers -- Craft Better Code*. ISBN-10/ ISBN-13 : 1593279701/ 978-1593279707. No Starch Press. [Chapter 4]
- Randall Hyde. *Write Great Code, Volume 1, 2nd Edition: Understanding the Machine* 2nd Edition. ASIN: B07VSC1K8Z. No Starch Press. 2020. [Chapter 11]
- Matthew Justice. *How Computers Really Work: A Hands-On Guide to the Inner Workings of the Machine*. ISBN-10/ISBN-13 : 1718500661/ 978-1718500662. No Starch Press. 2020. [Chapter 7]



COLORADO STATE UNIVERSITY

Professor: SHRIDEEP PALLICKARA
COMPUTER SCIENCE DEPARTMENT

MEMORY

L10.42

42