

Memory Hierarchy (I): The Basics

Hung-Wei Tseng



Disney • PIXAR INSIDE OUT

GET DISNEY+

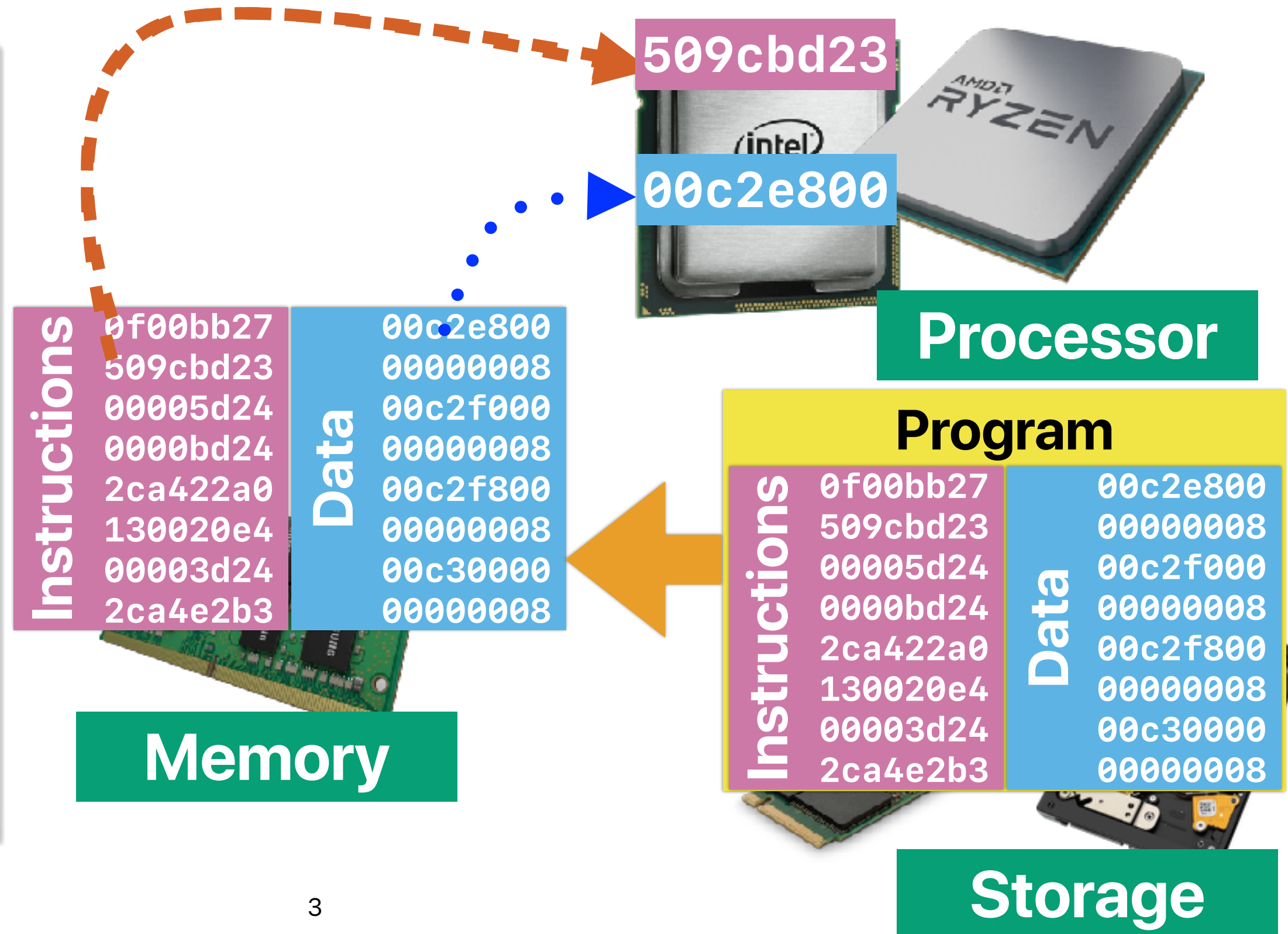
▶ TRAILER

PG 2015 • 1h 35m • Coming of age, Family, Animation

When 11-year-old Riley moves to a new city, her Emotions team up to help her through the transition. Joy, Fear, Anger, Disgust and Sadness work together, but when Joy and Sadness get lost, they must journey through unfamiliar places to get back home.



von Neumann Architecture



Recap: Summary of CPU Performance Equation

$$Performance = \frac{1}{Execution\ Time}$$

$$Execution\ Time = \frac{Instructions}{Program} \times \frac{Cycles}{Instruction} \times \frac{Seconds}{Cycle}$$

$$ET = IC \times CPI \times CT$$

$$Speedup = \frac{Execution\ Time_X}{Execution\ Time_Y}$$

- IC (Instruction Count)
 - ISA, Compiler, algorithm, programming language, **programmer**
- CPI (Cycles Per Instruction)
 - Machine Implementation, microarchitecture, compiler, application, algorithm, programming language, **programmer**
- Cycle Time (Seconds Per Cycle)
 - Process Technology, microarchitecture, **programmer**

Is TFLOPS (Tera Floating-point Operations Per Second) a good metric?

$$\begin{aligned} TFLOPS &= \frac{\# \text{ of floating point instructions} \times 10^{-12}}{\text{Execution Time}} \\ &= \frac{IC \times \% \text{ of floating point instructions} \times 10^{-12}}{IC \times CPI \times CT} \\ &= \frac{\% \text{ of floating point instructions} \times 10^{-12}}{CPI \times CT} \quad \text{IC is gone!} \end{aligned}$$

A good performance metric must cover IC, CPI, CT!

- Cannot compare different ISA/compiler
 - What if the compiler can generate code with fewer instructions?
 - What if new architecture has more IC but also lower CPI?
- Does not make sense if the application is not floating point intensive

Amdahl's Law

$$Speedup_{enhanced}(f, s) = \frac{1}{(1 - f) + \frac{f}{s}}$$

- Corollary #1: Maximum speedup
- Corollary #2: Make the common case fast
 - Common case changes all the time
- Corollary #3: Single-core performance still matters
- Corollary #4: Exploiting more parallelism from a program is the key to performance gain in modern architectures

$$Speedup_{max}(f, \infty) = \frac{1}{1 - f}$$
$$Speedup_{max}(f_1, \infty) = \frac{1}{1 - f_1}$$

$$Speedup_{max}(f_2, \infty) = \frac{1}{1 - f_2}$$

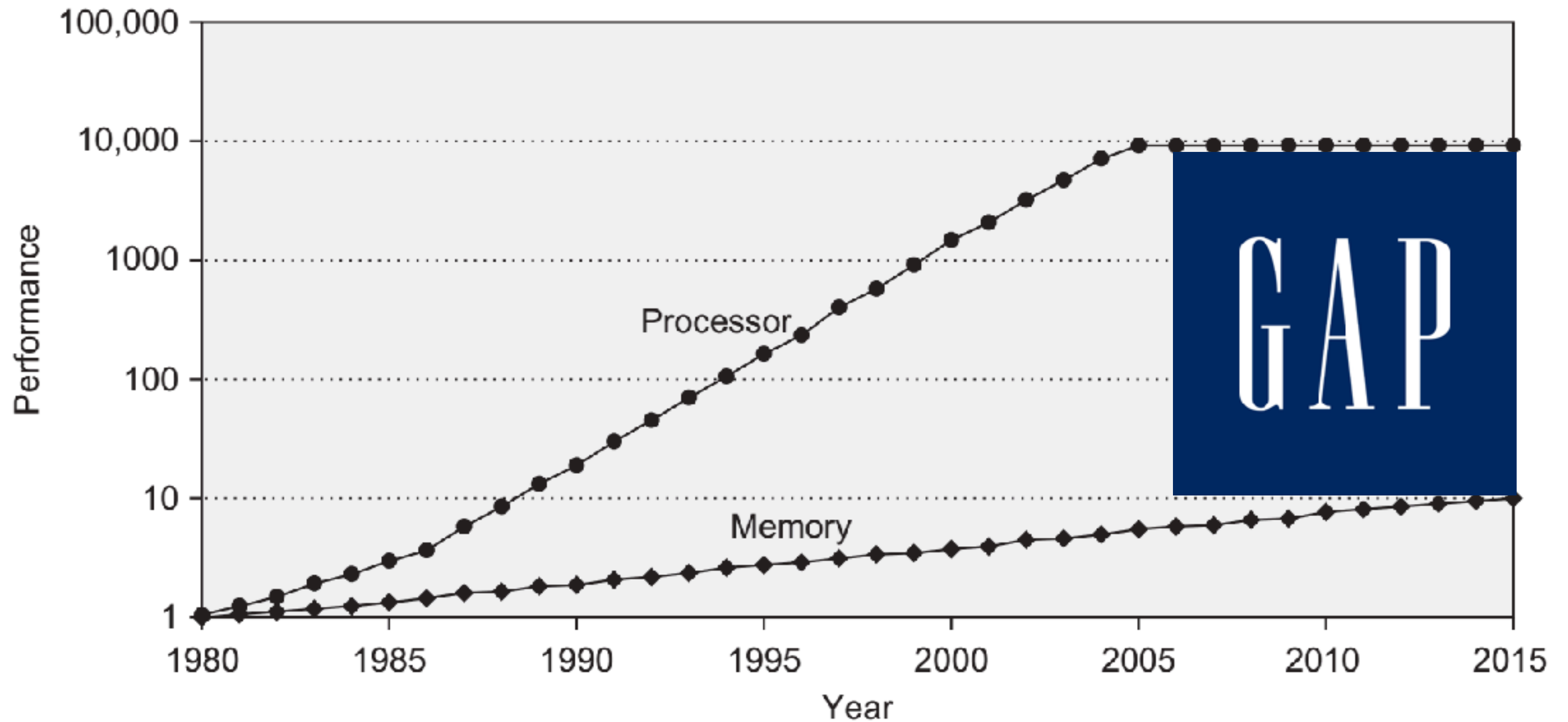
$$Speedup_{max}(f_3, \infty) = \frac{1}{1 - f_3}$$

$$Speedup_{max}(f_4, \infty) = \frac{1}{1 - f_4}$$

$$Speedup_{parallel}(f_{parallelizable}, \infty) = \frac{1}{1 - f_{parallelizable}}$$

$$Speedup_{parallel}(f_{parallelizable}, \infty) = \frac{1}{1 - f_{parallelizable}}$$

Performance gap between Processor/Memory



The impact of "slow" memory

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has "perfect" memory, the CPI is just 1. Now, consider we have DDR4 and the program is well-behaved that precharge is never necessary — the access latency is simply 26 ns. What's the average CPI (pick the most close one)?

A. 9

B. 17

C. 27

D. 35

E. 69

$$1 + 100\% \times (52) + 30\% \times 52 = 68.6 \text{ cycles}$$



#MAGA



BIDEN

Unite for a Better Architecture

APPROVED BY JOE BIDEN. PAID FOR BY BIDEN FOR PRESIDENT.

Team scores



0.5



2



2.5



3

Outline

- The Basic Idea behind Memory Hierarchy
- How cache works

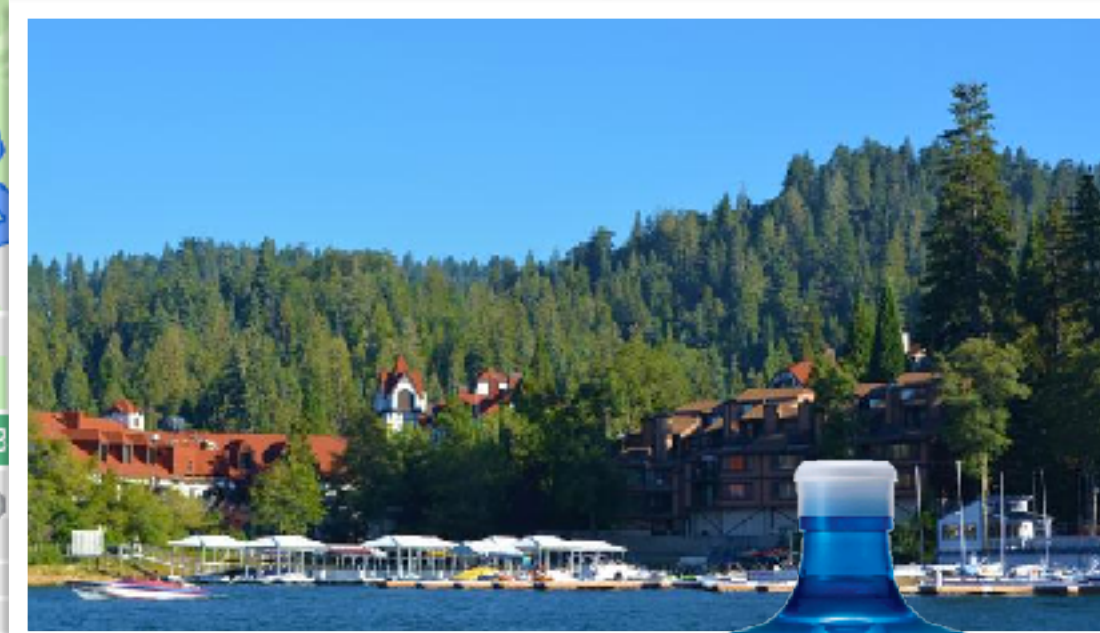
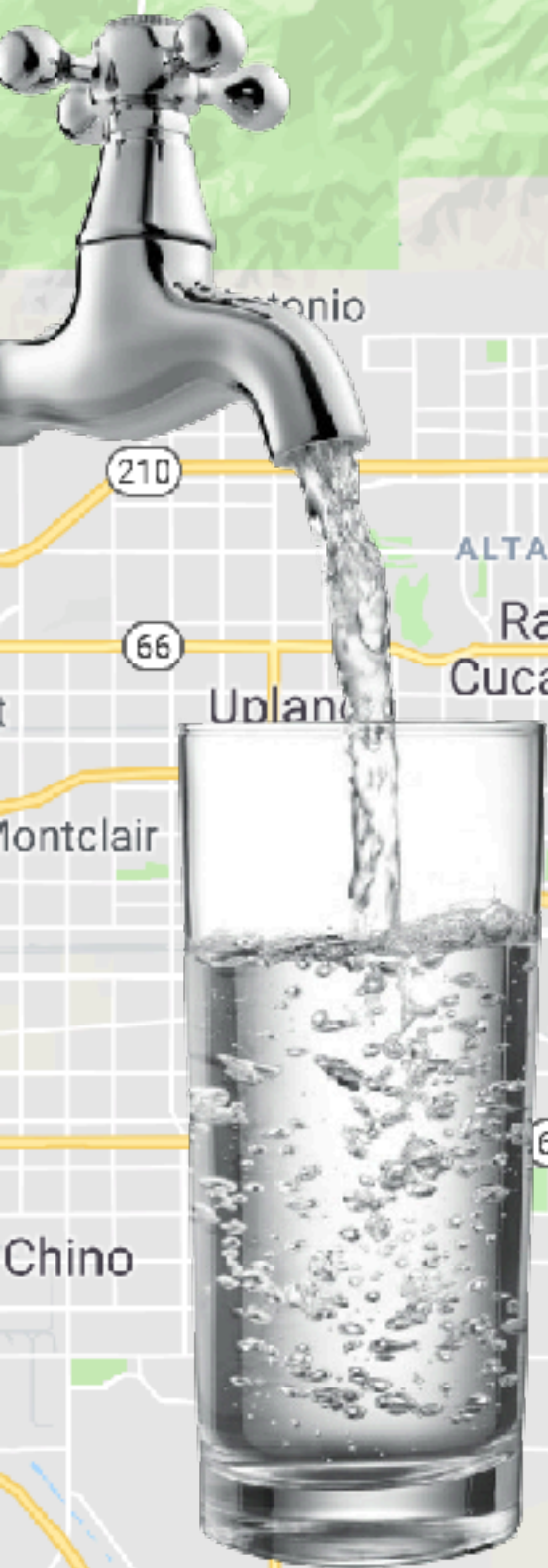
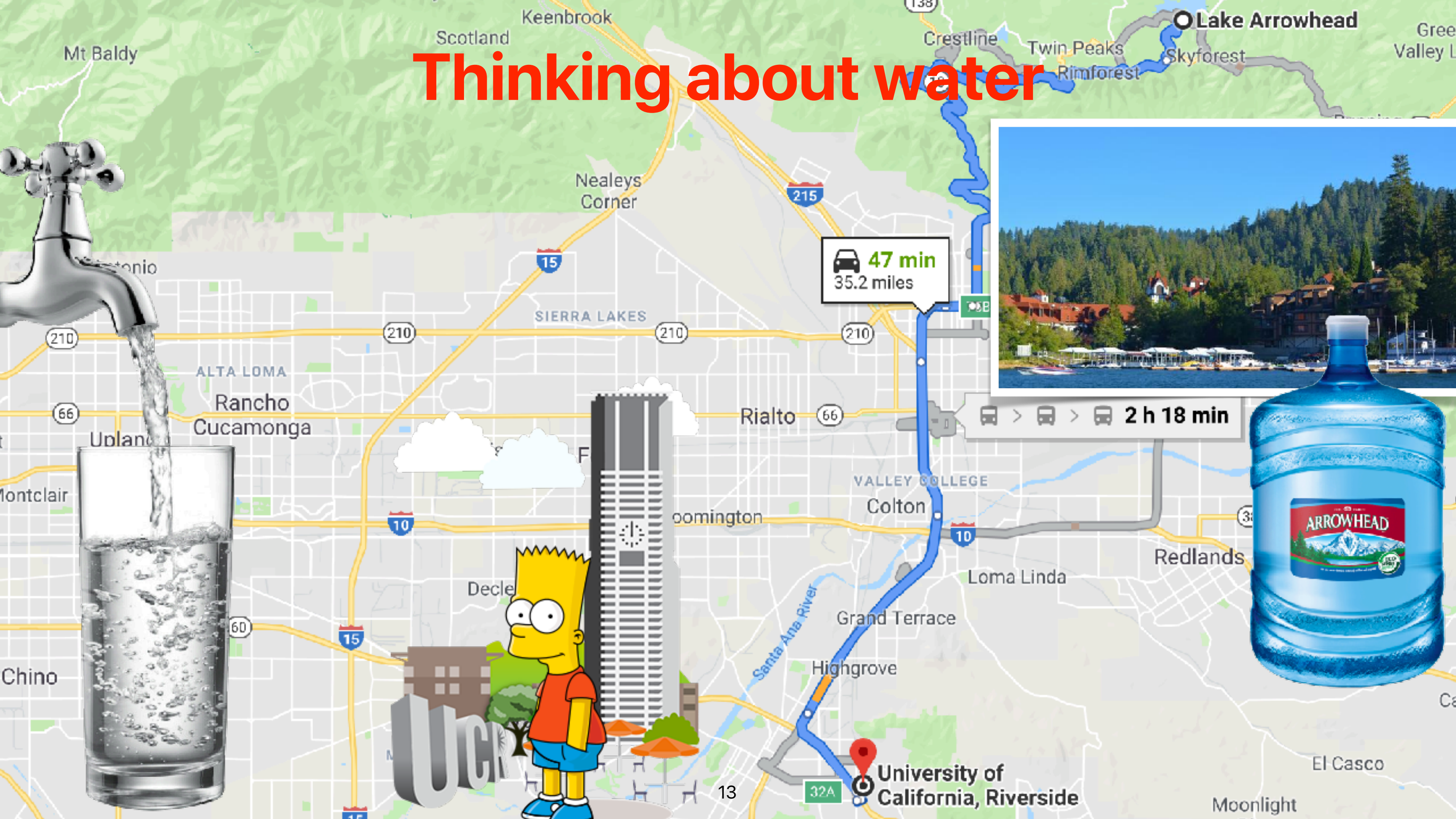
Alternatives?

Memory technology	Typical access time	\$ per GiB in 2012
SRAM semiconductor memory	0.5–2.5 ns	\$500–\$1000
DRAM semiconductor memory	50–70 ns	\$10–\$20
Flash semiconductor memory	5,000–50,000 ns	\$0.75–\$1.00
Magnetic disk	5,000,000–20,000,000 ns	\$0.05–\$0.10



Fast, but expensive \$\$\$

Thinking about water

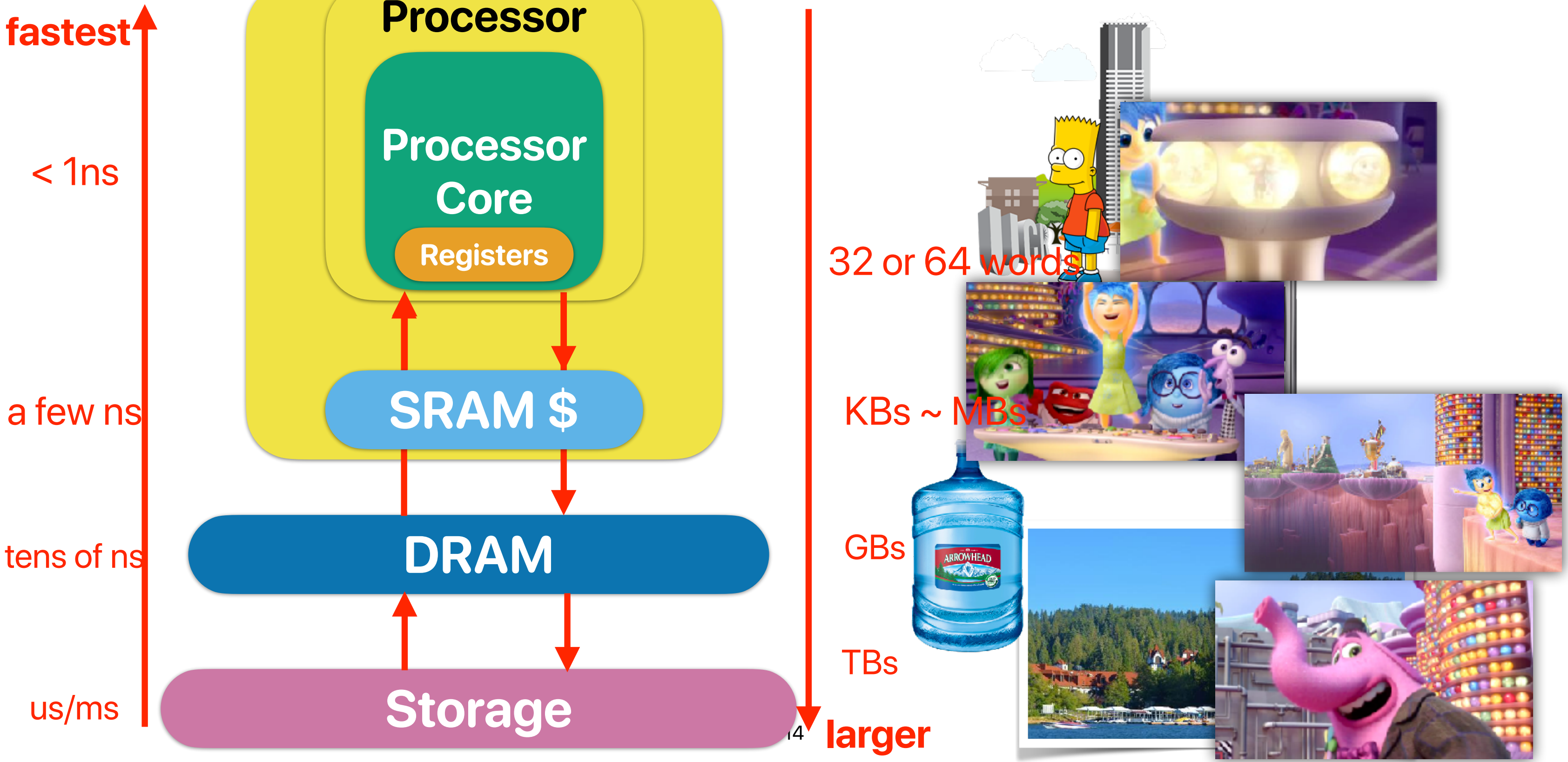


47 min
35.2 miles

2 h 18 min

University of
California, Riverside

Memory Hierarchy



How can memory hierarchy help in performance?

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has “perfect” memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got an SRAM cache with latency of just at 0.5ns and can capture 90% of the desired data/instructions. what’s the average CPI (pick the most close one)?

- A. 2
- B. 4
- C. 8
- D. 16
- E. 32

How can memory hierarchy help in performance



- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has “perfect” memory, the CPI is just 1. Now, in addition to DDR4, whose latency is 26 ns, we also got an SRAM cache with latency of just at 0.5ns and can capture 90% of the desired data/instructions. what’s the average CPI (pick the most close one)?

- A. 2
- B. 4
- C. 8
- D. 16
- E. 32

How can memory hierarchy help in performance?

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has “perfect” memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got an SRAM cache with latency of just at 0.5ns and can capture 90% of the desired data/instructions. what’s the average CPI (pick the most close one)?
 - A. 2
 - B. 4
 - C. 8

 $1 + (1 - 90\%) \times [100\% \times (52) + 30\% \times 52] = 7.76 \text{ cycles}$
 - D. 16
 - E. 32

L1? L2? L3?

CPU-Z - ID : wswpbb

CPU | Caches | Mainboard | Memory | SPD | Graphics | Bench | About

Processor

Name	AMD Ryzen 7 2700X		
Code Name	Pinnacle Ridge	Max TDP	105 W
Package	Socket AM4 (1331)		
Technology	12 nm	Core Voltage	1.36 V
Specification	AMD Ryzen 7 2700X Eight-Core Processor		
Family	F	Model	8
Ext. Family	17	Ext. Model	8
Instructions	MMX(+), SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, SSE4A, x86-64, AMD-V, AES, AVX, AVX2, FMA3, SHA		

Clocks (Core #0)

Core Speed	4290.73 MHz
Multiplier	x 43.0
Bus Speed	99.78 MHz
Rated FSB	

Cache

L1 Data	8 x 32 KBytes	8-way
L1 Inst.	8 x 64 KBytes	4-way
Level 2	8 x 512 KBytes	8-way
Level 3	2 x 8192 KBytes	16-way

Selection: Processor #1 | Cores: 8 | Threads: 16

CPU-Z Ver. 1.86.0.x64 | Tools | Validate | Close

CPU | Caches | Mainboard | Memory | SPD | Graphics | Bench | About

Processor

Name	Intel Core i7 9700K		
Code Name	Coffee Lake	Max TDP	95.0 W
Package	Socket 1151 LGA		
Technology	14 nm	Core Voltage	0.737 V
Specification	Intel® Core™ i7-9700K CPU @ 3.60GHz (ES)		
Family	6	Model	E
Ext. Family	6	Ext. Model	9E
Instructions	MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, EM64T, VT-x, AES, AVX, AVX2, FMA3, TSX		

Clocks (Core #0)

Core Speed	4798.85 MHz
Multiplier	x 48.0 (8 - 49)
Bus Speed	99.98 MHz
Rated FSB	

Cache

L1 Data	8 x 32 KBytes	8-way
L1 Inst.	8 x 32 KBytes	8-way
Level 2	8 x 256 KBytes	4-way
Level 3	12 MBytes	12-way

Selection: Socket #1 | Cores: 8 | Threads: 8

How can deeper memory hierarchy help in performance?

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has “perfect” memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got a 2-level SRAM caches with
 - it's 1st-level one at latency of 0.5ns and can capture 90% of the desired data/instructions.
 - the 2nd-level at latency of 5ns and can capture 60% of the desired data/instructions

What's the average CPI (pick the most close one)?

- A. 2
- B. 4
- C. 8
- D. 16
- E. 32



How can deeper memory hierarchy help in performance

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has “perfect” memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got a 2-level SRAM caches with
 - it's 1st-level one at latency of 0.5ns and can capture 90% of the desired data/instructions.
 - the 2nd-level at latency of 5ns and can capture 60% of the desired data/instructions

What's the average CPI (pick the most close one)?

- A. 2
- B. 4
- C. 8
- D. 16
- E. 32

How can deeper memory hierarchy help in performance?

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has “perfect” memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got a 2-level SRAM caches with
 - it's 1st-level one at latency of 0.5ns and can capture 90% of the desired data/instructions.
 - the 2nd-level at latency of 5ns and can capture 60% of the desired data/instructions

What's the average CPI (pick the most close one)?

A. 2

B. 4

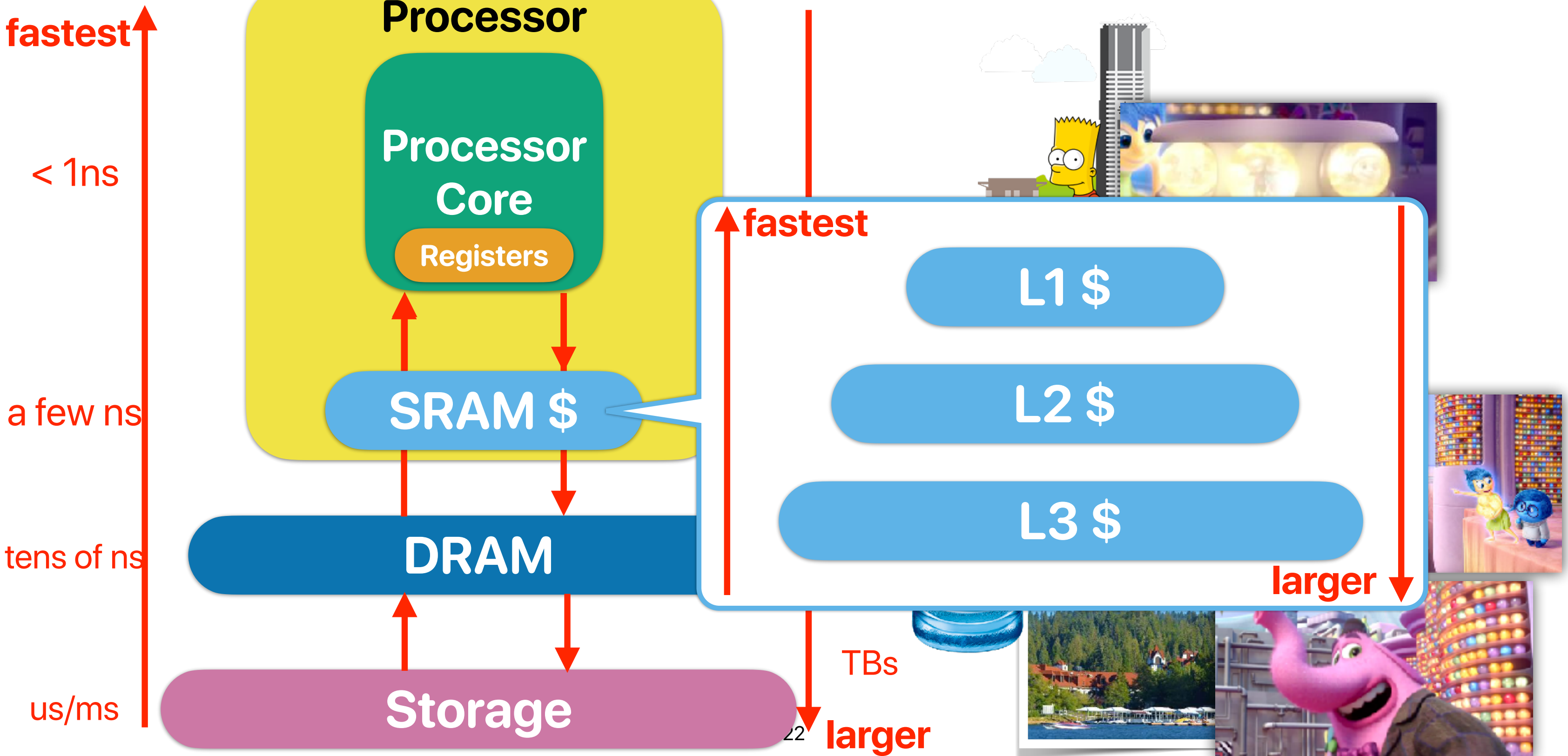
C. 8

D. 16

E. 32

$$1 + (1 - 90\%) \times [10 + (1 - 60\%) \times 52 + 30\% \times (10 + (1 - 60\%) \times 52)] = 5 \text{ cycles}$$

Memory Hierarchy



**Why adding small SRAMs would
work?**

Locality

- Which description about locality of arrays `sum` and `A` in the following code is the most accurate?

```
for(i = 0; i < 100000; i++)  
{  
    sum[i%10] += A[i];  
}
```

- A. Access of `A` has temporal locality, `sum` has spatial locality
- B. Both `A` and `sum` have temporal locality, and `sum` also has spatial locality
- C. Access of `A` has spatial locality, `sum` has temporal locality
- D. Both `A` and `sum` have spatial locality
- E. Both `A` and `sum` have spatial locality, and `sum` also has temporal locality

Locality

- Which description about locality of arrays `sum` and `A` in the following code is the most accurate?

```
for(i = 0; i < 100000; i++)  
{  
    sum[i%10] += A[i];  
}
```

- A. Access of `A` has temporal locality, `sum` has spatial locality
- B. Both `A` and `sum` have temporal locality, and `sum` also has spatial locality
- C. Access of `A` has spatial locality, `sum` has temporal locality
- D. Both `A` and `sum` have spatial locality
- E. Both `A` and `sum` have spatial locality, and `sum` also has temporal locality

Locality

- Which description about locality of arrays `sum` and `A` in the following code is the most accurate?

```
for(i = 0; i < 100000; i++)  
{  
    sum[i%10] += A[i];  
}
```

spatial locality:

`A[0], A[1], A[2], A[3],`

`sum[0], sum[1], ..., sum[9]`

temporal locality:

reuse of `sum[0], sum[1], ..., sum[9]`

- A. Access of `A` has temporal locality, `sum` has spatial locality
- B. Both `A` and `sum` have temporal locality, and `sum` also has spatial locality
- C. Access of `A` has spatial locality, `sum` has temporal locality
- D. Both `A` and `sum` have spatial locality
- E. Both `A` and `sum` have spatial locality, and `sum` also has temporal locality

Locality

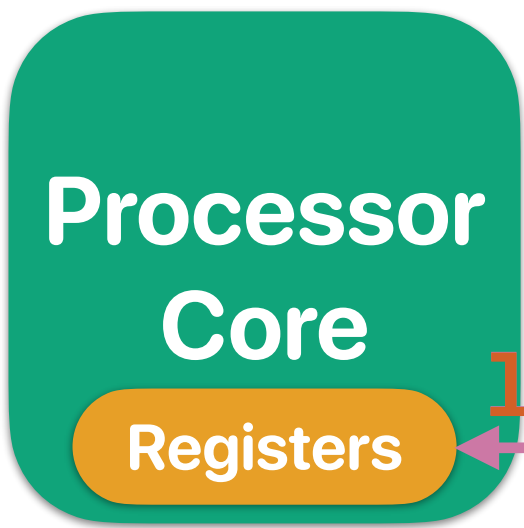
- Spatial locality — application tends to visit nearby stuffs in the memory

- Code — the current instruction, and then $PC + 4$

Most of time, your program is just visiting a very small amount of data/instructions within a given window

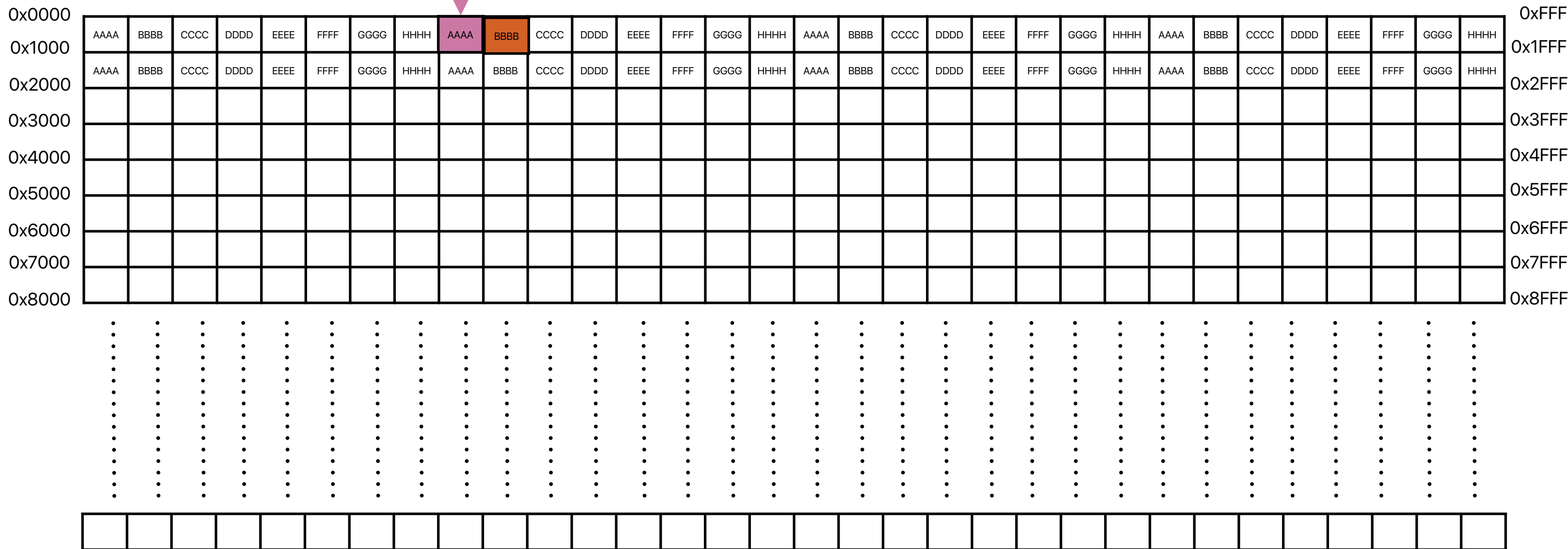
- Temporal locality — application revisit the same thing again and again
 - Code — loops, frequently invoked functions
 - Data — the same data can be read/write many times


Architecting the Cache



Load/store only access a "word" each time

load 0x000A






Processor Core

Registers

Registers

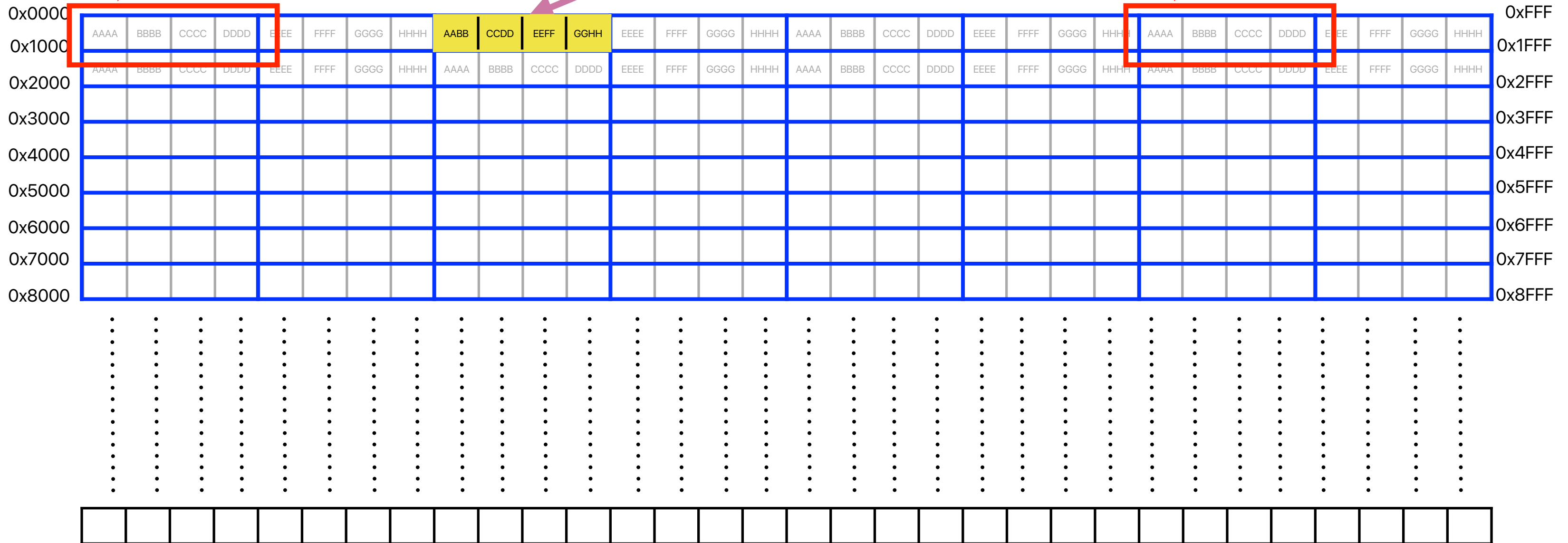
1w 0x0024



A diagram showing a blue rounded rectangle containing the text "SRAM \$". Below the text, there are two colored squares: a purple square labeled "AABB" and an orange square labeled "CCDD".

Assume each block is 16 bytes

**"Logically" partition
memory space into
↓
"blocks"**





How to tell who is there?

tag

00000000000000000000000000000000
0123456789ABCDEF
0x00000000000000000000000000000000

0x000	This is CS 203: Advanced Computer Architecture! This is CS 203: Advanced Computer Architecture! This is CS 203: Advanced Computer Architecture! This is CS 203: Advanced Computer Architecture! This is CS 203: Advanced Computer Architecture! This is CS 203:
-------	---



How to tell w

block offset

tag

1w 0x0008

1w 0x4048

0x404 not found,
go to lower-level memory

The complexity of search the matching tag—
 $O(n)$ —will be slow if our cache size grows!

Can we search things faster?
—hash table! $O(1)$

Tell if the block here can be used

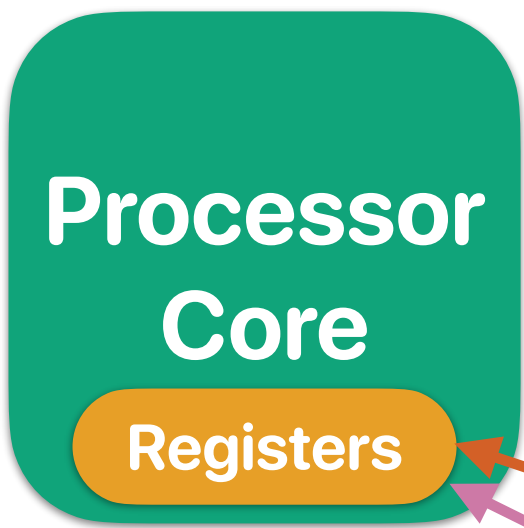
Tell if the block here is modified

Valid Bit
Dirty Bit

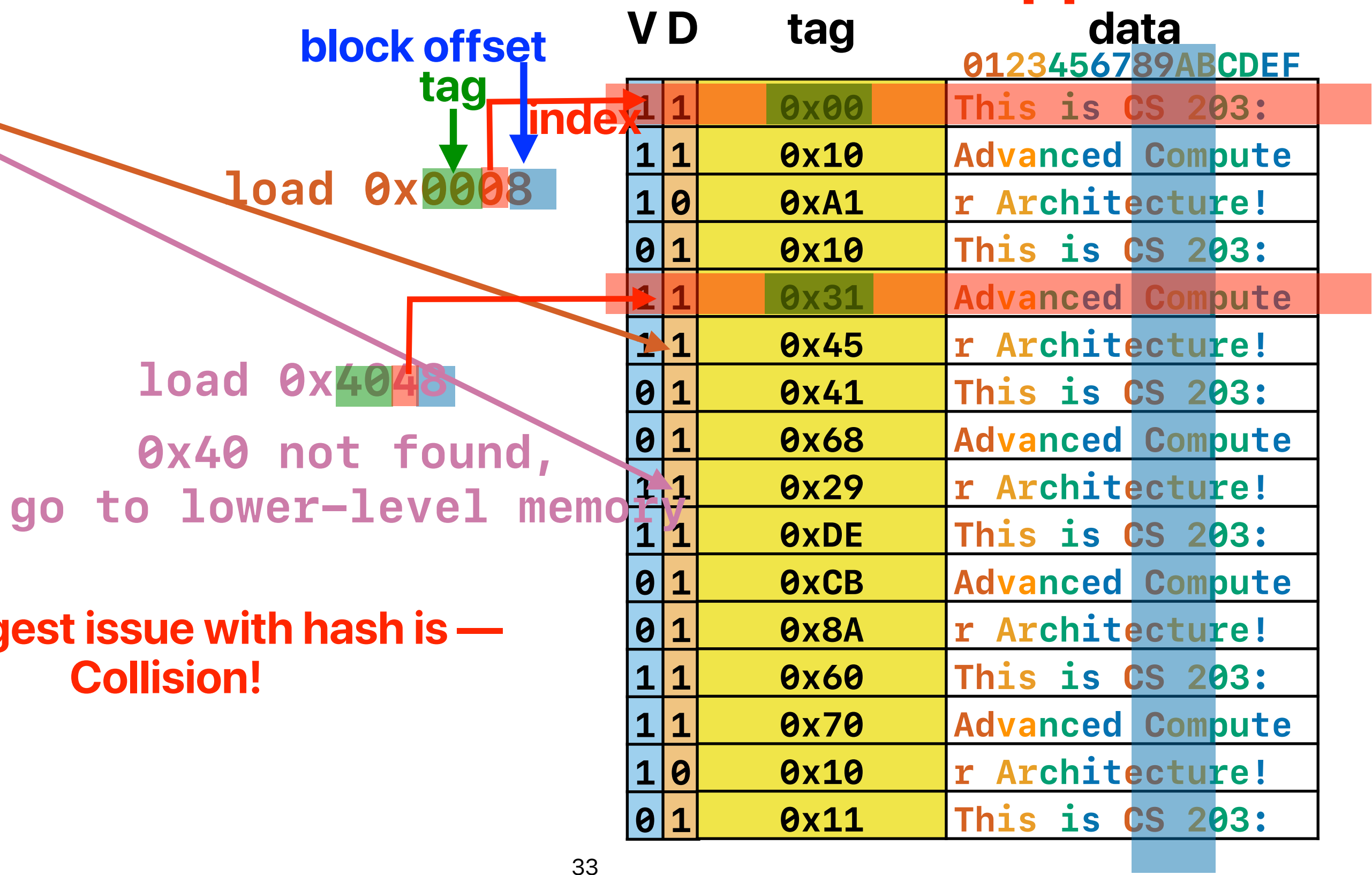
tag

data

				0123456789ABCDEF
1	1	0x000	This is CS 23:	
1	1	0x001	Advanced Compute	
1	0	0xF07	r Architecture!	
0	1	0x100	This is CS 203:	
1	1	0x310	Advanced Compute	
1	1	0x450	r Architecture!	
0	1	0x006	This is CS 203:	
0	1	0x537	Advanced Compute	
1	1	0x266	r Architecture!	
1	1	0x307	This is CS 203:	
0	1	0x265	Advanced Compute	
0	1	0x80A	r Architecture!	
1	1	0x620	This is CS 203:	
1	1	0x630	Advanced Compute	
1	0	0x705	r Architecture!	
0	1	0x216	This is CS 203:	



Hash-like structure — direct-mapped cache



The biggest issue with hash is — Collision!

Way-associative cache

memory address: $0x0$ 8 2 4

set block

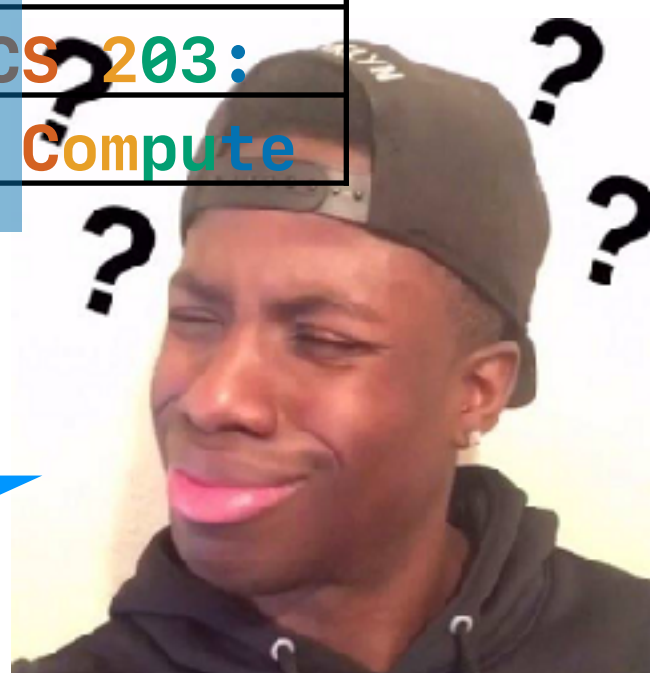
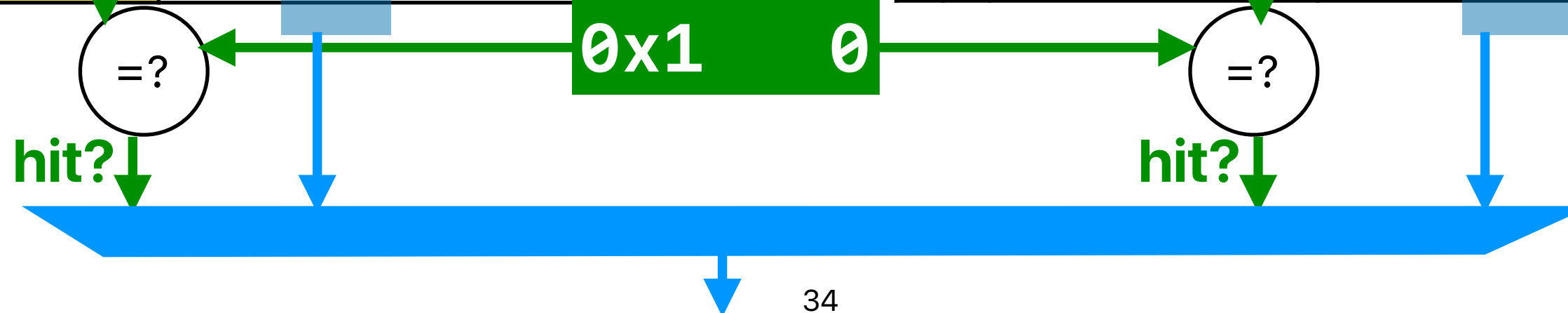
tag index offset

memory address: $0b00001000000100100$

V	D	tag	data
1	1	$0x29$	r Architecture!
1	1	$0xDE$	This is CS 203:
1	0	$0x10$	Advanced Compute
0	1	$0x8A$	r Architecture!
1	1	$0x60$	This is CS 203:
1	1	$0x70$	Advanced Compute
0	1	$0x10$	r Architecture!
0	1	$0x11$	This is CS 203:

V	D	tag	data
1	1	$0x00$	This is CS 203:
1	1	$0x10$	Advanced Compute
1	0	$0xA1$	r Architecture!
0	1	$0x10$	This is CS 203:
1	1	$0x31$	Advanced Compute
1	1	$0x45$	r Architecture!
0	1	$0x41$	This is CS 203:
0	1	$0x68$	Advanced Compute

Set



$$C = ABS$$

- **C: Capacity** in data arrays
- **A: Way-Associativity** — how many blocks within a set
 - N-way: N blocks in a set, $A = N$
 - 1 for direct-mapped cache
- **B: Block Size (Cacheline)**
 - How many bytes in a block
- **S: Number of Sets:**
 - A set contains blocks sharing the same index
 - 1 for fully associate cache



Corollary of $C = ABS$

memory address: 0b 000010000 010 0100

tag set index block offset

- number of bits in **block** offset — $\lg(\mathbf{B})$
- number of bits in **set** index: $\lg(\mathbf{S})$
- tag bits: $\text{address_length} - \lg(\mathbf{S}) - \lg(\mathbf{B})$
 - address_length is 32 bits for 32-bit machine
- $(\text{address} / \text{block_size}) \% \mathbf{S} = \text{set index}$

AMD Phenom II

- L1 data (D-L1) cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block
 - Assume 64-bit memory address

Which of the following is correct?

- A. Tag is 49 bits
- B. Index is 8 bits
- C. Offset is 7 bits
- D. The cache has 1024 sets
- E. None of the above

AMD Phenom II



- L1 data (D-L1) cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block
 - Assume 64-bit memory address

Which of the following is correct?

- A. Tag is 49 bits
- B. Index is 8 bits
- C. Offset is 7 bits
- D. The cache has 1024 sets
- E. None of the above

AMD Phenom II

- L1 data (D-L1) cache configuration of AMD Phenom II
 - Size 64KB, 2-way set associativity, 64B block
 - Assume 64-bit memory address

Which of the following is correct?

A. Tag is 49 bits

B. Index is 8 bits

C. Offset is 7 bits

D. The cache has 1024 sets

E. None of the above

$$C = ABS$$

$$64\text{KB} = 2 * 64 * S$$

$$S = 512$$

$$\text{offset} = \lg(64) = 6 \text{ bits}$$

$$\text{index} = \lg(512) = 9 \text{ bits}$$

$$\text{tag} = 64 - \lg(512) - \lg(64) = 49 \text{ bits}$$

Announcement

- Assignment #1 due tonight
 - Assignments SHOULD BE done individually — if discussed with others, make sure their names on your submission
 - We will drop your least performing assignment as well
 - Attendance counts as one assignment
- Reading quiz due Wednesday before the lecture
 - We will drop two of your least performing reading quizzes
 - You have two shots, both unlimited time
- Joel Emer's Talk next Monday @ 11am
 - We will not have a lecture next Monday to encourage you attend Joel Emer's talk
 - If you capture a screen shot and submit it through iLearn, you will receive a full credit reading quiz
 - The talk cannot be broadcasted on YouTube due to the license constraint
- Office Hours on Zoom (the office hour link, not the lecture one)
 - Walk-in, no appointment is necessary
 - Hung-Wei/Prof. Usagi: M 8p-9p, W 2p-3p
 - Quan Fan: F 1p-3p

Computer Science & Engineering

203

つづく

