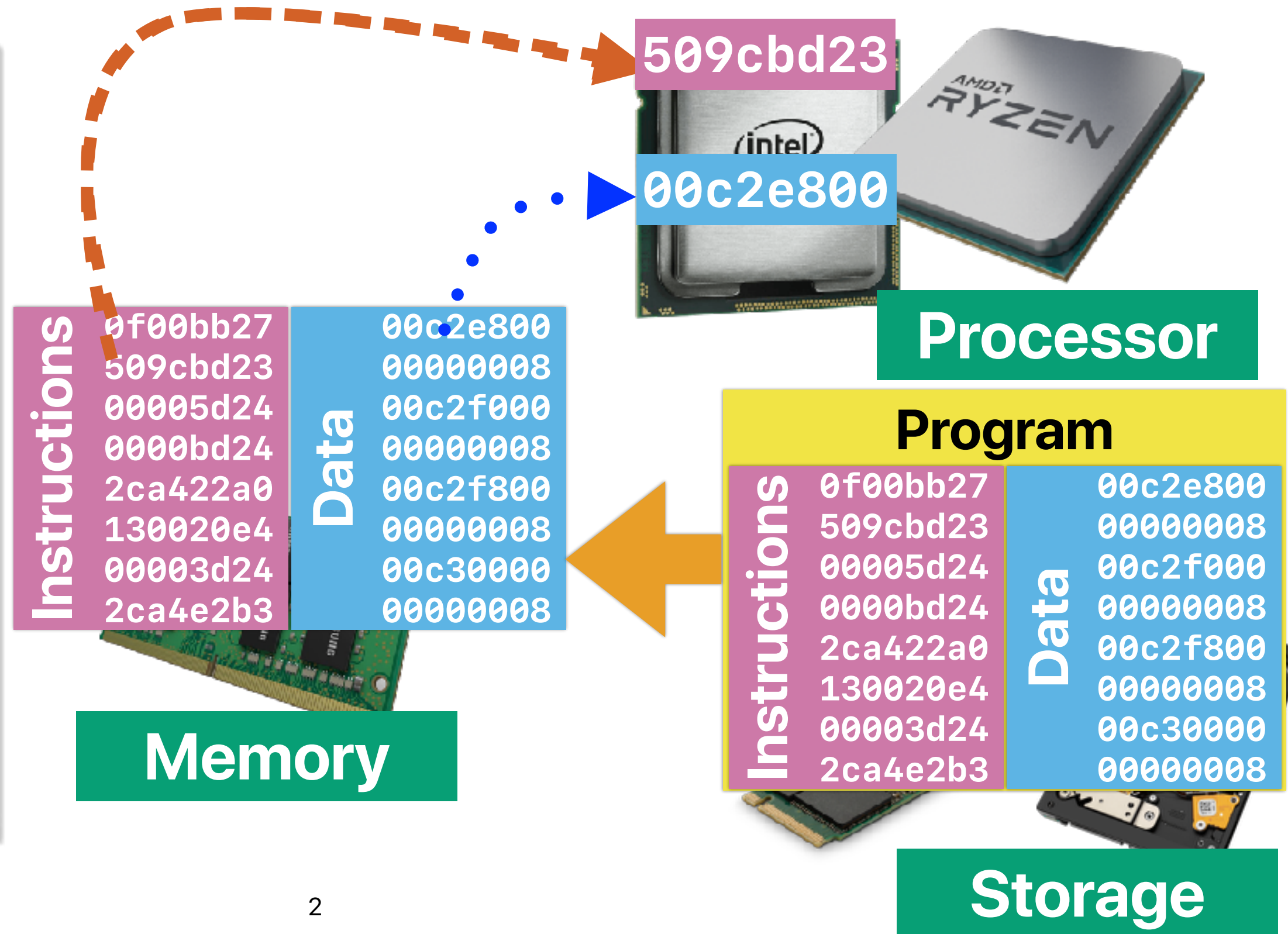


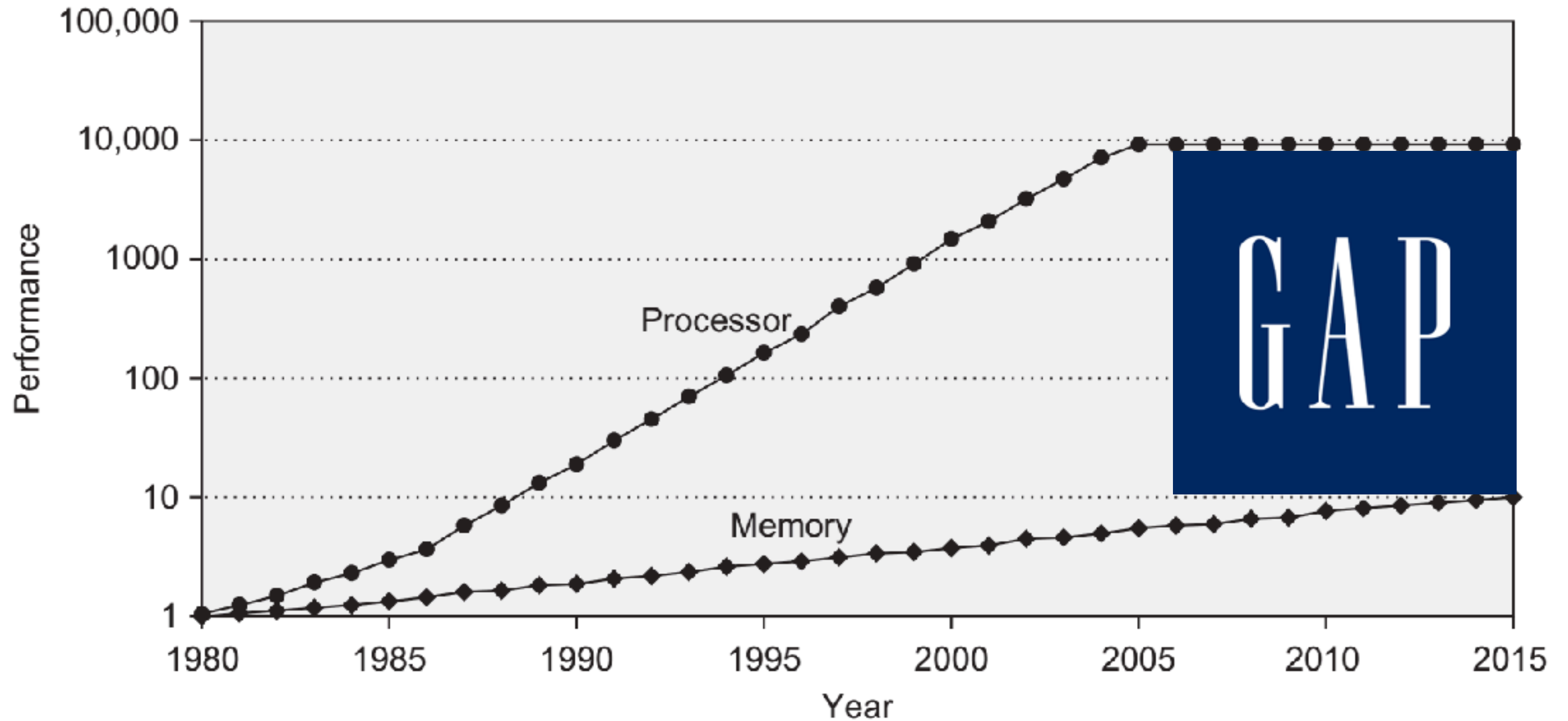
Memory Hierarchy (2): The A, B, Cs of Your Cache

Hung-Wei Tseng

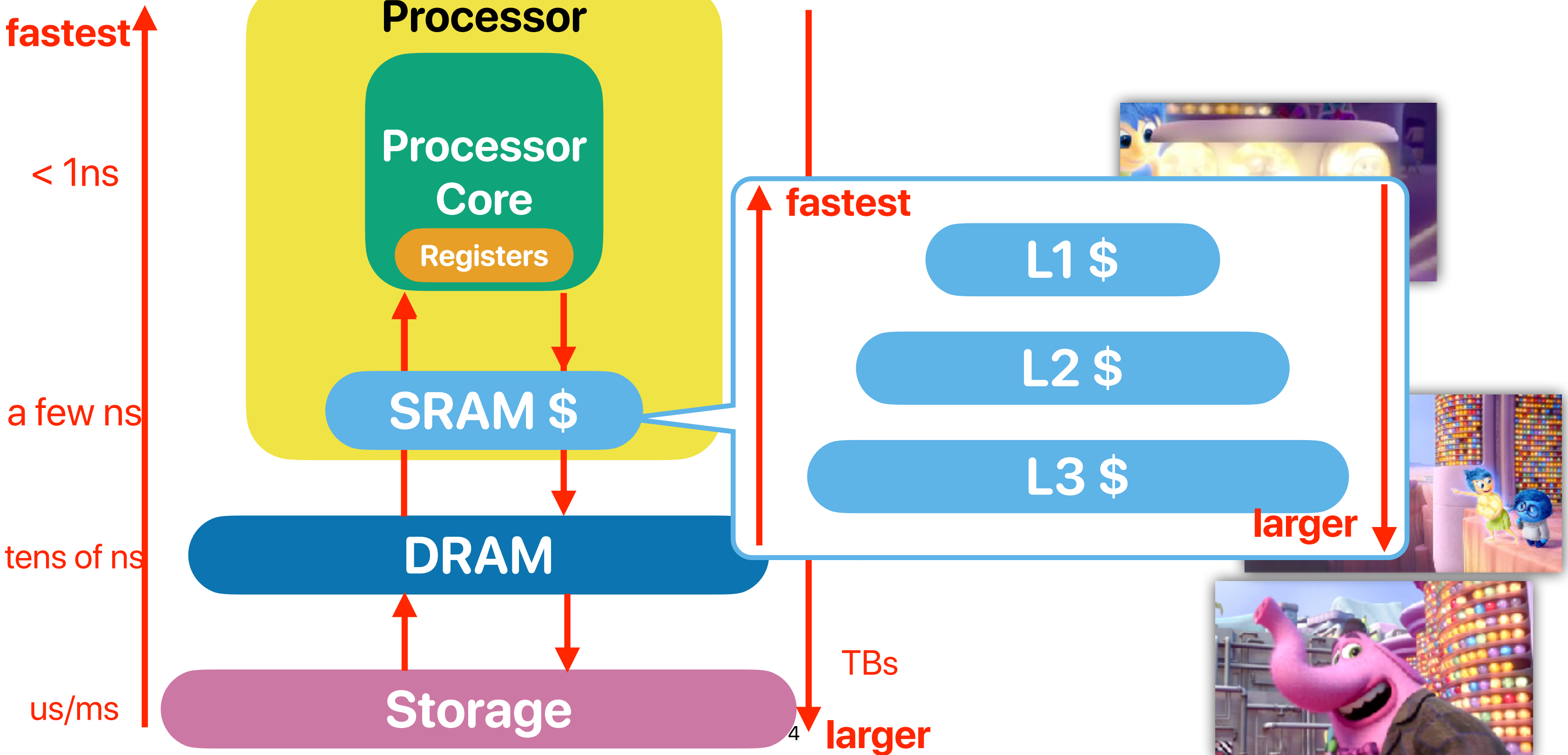
Recap: von Neumann Architecture



Recap: Performance gap between Processor/Memory



Memory Hierarchy



How can "memory hierarchy" help in performance?

- Assume that we have a processor running @ 4 GHz and a program with 20% of load/store instructions. If the instruction has no memory access, the CPI is just 1. Now, in addition to we DDR5, whose latency 13.75 ns, we also got an SRAM cache with latency of just at 0.5 ns and can capture 90% of the desired data/instructions. what's the average CPI (pick the closest one)?

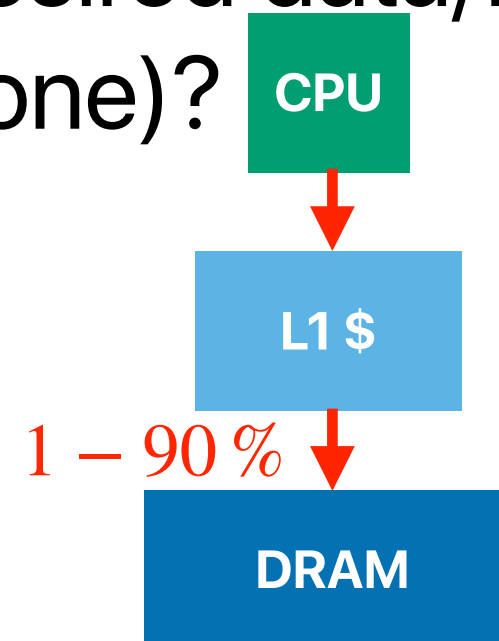
A. 6

B. 8

C. 10

D. 12

E. 67



$$CPU \text{ cycle time} = \frac{1}{4 \times 10^9} = 0.25ns$$

$$Each \$ \text{ access} = \frac{0.5}{0.25} = 2 \text{ cycles}$$

$$Each DRAM \text{ access} = \frac{13.75}{0.25} = 55 \text{ cycles}$$

$$CPI_{average} = 1 + 100\% \times [2 + (1 - 90\%) \times 55] + 20\% \times [2 + (1 - 90\%) \times 55] = 10 \text{ cycles}$$

L1? L2? L3?

Can we really “predict” upcoming data accurately (e.g., 90%) with such small caches?

CPU-Z - ID : vfljgr	
CPU Mainboard Memory SPD Graphics Bench About	
Processor	
Name	AMD Ryzen 7 7700X
Code Name	Raphael
Max TDP	105 W
Package	AM5
Technology	5 nm
Core Voltage	1.2 V
Specification	
AMD Ryzen 7 7700X 8-Core Processor	
Ext. Family	19
Ext. Model	1
Revision	RPL/2
Instructions	
MMX(+), SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, SSE4A, x86-64, AMD-V, AES, AVX, AVX2, AVX512, FMA3, SHA	
Clocks (Core #0)	
Core Speed	5188.99 MHz
Multiplier	x 52.0 (4 - 55.5)
Bus Speed	99.79 MHz
Rated FSB	
Cache	
L1 Data	8 x 32 KB
L1 Inst.	8 x 32 KB
Level 2	8 x 1024 KB
Level 3	32 MBytes
Selection Socket #1	
Cores	8
Threads	16
CPU-Z Ver. 2.09.0.x64 Tools Validate Close	

CPU-Z - ID : pk15b	
CPU Mainboard Memory SPD Graphics Bench About	
Processor	
Name	Intel Core i7 14700K
Code Name	Raptor Lake
Max TDP	125 W
Package	LGA1700
Technology	14 nm
Core Voltage	1.1 V
Specification	
Intel(R) Core(TM) i7-14700K	
Ext. Family	6
Ext. Model	b7
Revision	B0
Instructions	
MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, EM64T, VT-x, AES, AVX, AVX2, FMA3, SHA	
Clocks (Core #0)	
Core Speed	5287.07 MHz
Multiplier	x 53.0 (8 - 55)
Bus Speed	99.76 MHz
Rated FSB	
Cache	
L1 Data	8 x 48 KB + 12 x 32 KB
L1 Inst.	8 x 32 KB + 12 x 64 KB
Level 2	8 x 2 MB + 3 x 4 MB
Level 3	33 MBytes
Selection Socket #1	
Cores	8 + 12
Threads	28
CPU-Z Ver. 2.08.0.x64 Tools Validate Close	

Outline

- The “predictable” code behavior
- Designing a cache that captures the predictability
- Estimating the code performance on cache

The predictability of your code

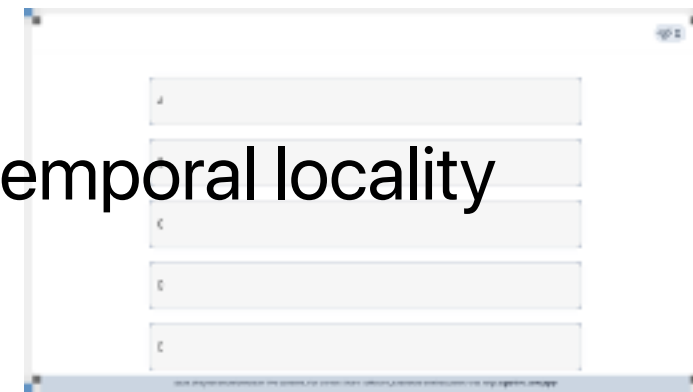


Locality of data

- Which description about locality of arrays `matrix` and `vector` in the following code is the **most accurate**?

```
for(uint32_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint32_t j = 0; j < n; j++) {  
        result += matrix[i][j]*vector[j];  
    }  
    output[i] = result;  
}
```

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



Data locality

- Which description about locality of arrays `matrix` and `vector` in the following code is the **most accurate**?

```
for(uint32_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint32_t j = 0; j < n; j++) {  
        result += matrix[i][j]*vector[j];  
    }  
    output[i] = result;  
}
```

spatial locality:

`matrix[0][0], matrix[0][1], matrix[0][2], ...`

`vector[0], vector[1], ..., vector[n]`

temporal locality:

reuse of `vector[0], vector[1], ...,`

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

Code also has locality

keep going to the
next instruction —
spatial locality

```
for(uint32_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint32_t j = 0; j < n; j++) {  
        result += matrix[i][j]*vector[j];  
    }  
    output[i] = result;  
}
```

**repeat many times —
temporal locality!**

```
i = 0;  
while(i < m) {  
    result = 0;  
    j = 0;  
    while(j < n) {  
        a = matrix[i][j];  
        b = vector[j];  
        temp = a*b;  
        result = result + temp;  
    }  
    output[i] = result;  
    i++;  
}
```

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

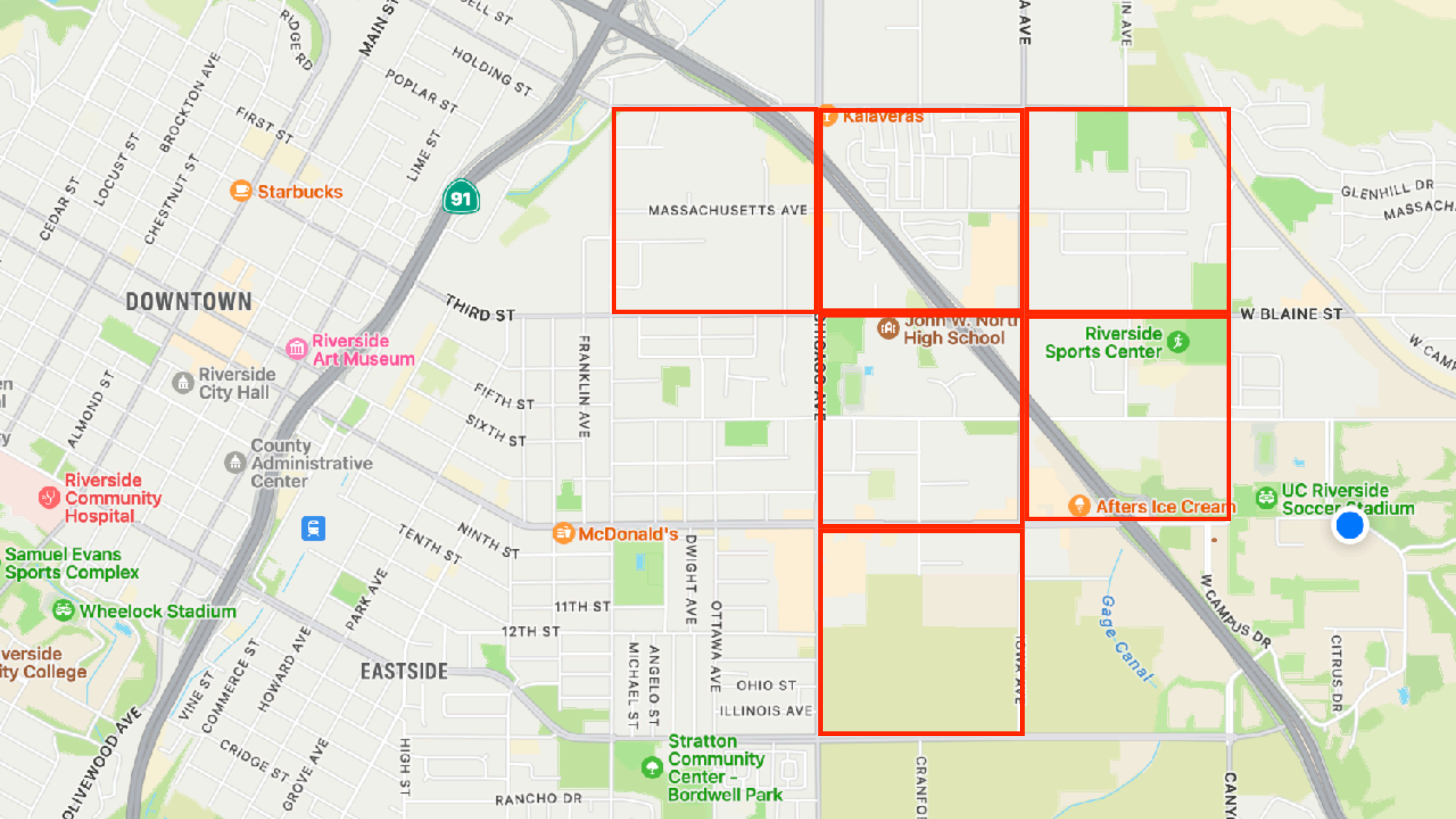
- Code — loops, frequently invoked functions
 - Typically tens of static instructions — at most several KBs
 - Data — program can read/write the same data many times (e.g., vectors in matrix-vector product)

Designing a hardware to exploit locality

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

We need to "cache consecutive memory locations" every time — the cache should store a "block" of code/data

- Temporal locality — application revisit the same thing again and again
 - Code — loops, frequently invoked functions
 - Typically tens of static instructions — at most several KBs
 - Data — program can read/write the same data many times (e.g., vectors in matrix-vector product)



DOWNTOWN

EASTSIDE

MASSACHUSETTS AVE

91

Starbucks

Kalaveras

Riverside Art Museum

Riverside City Hall

County Administrative Center

Riverside Community Hospital

Samuel Evans Sports Complex

Wheelock Stadium

Riverside City College

McDonald's

John W. North High School

Riverside Sports Center

Afters Ice Cream

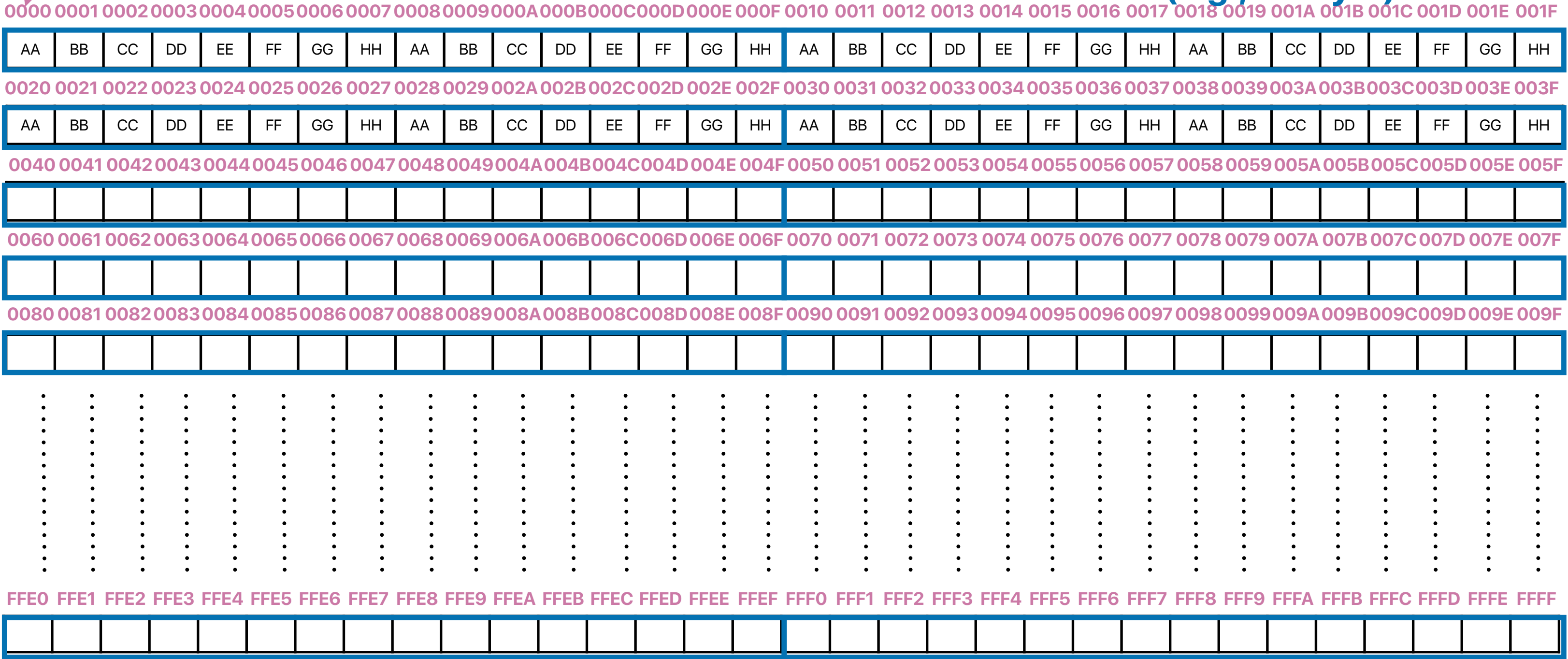
UC Riverside Soccer Stadium

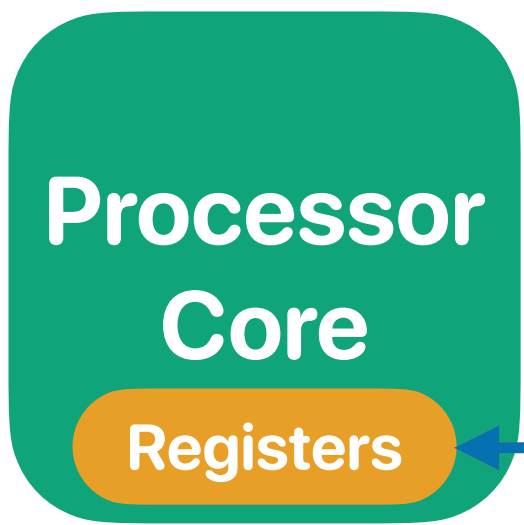
Stratton Community Center - Bordwell Park

Gage Canal

Block and the memory space

Each byte of memory location has an "address" Partition the space with fixed-size "blocks" (e.g., 16-byte)



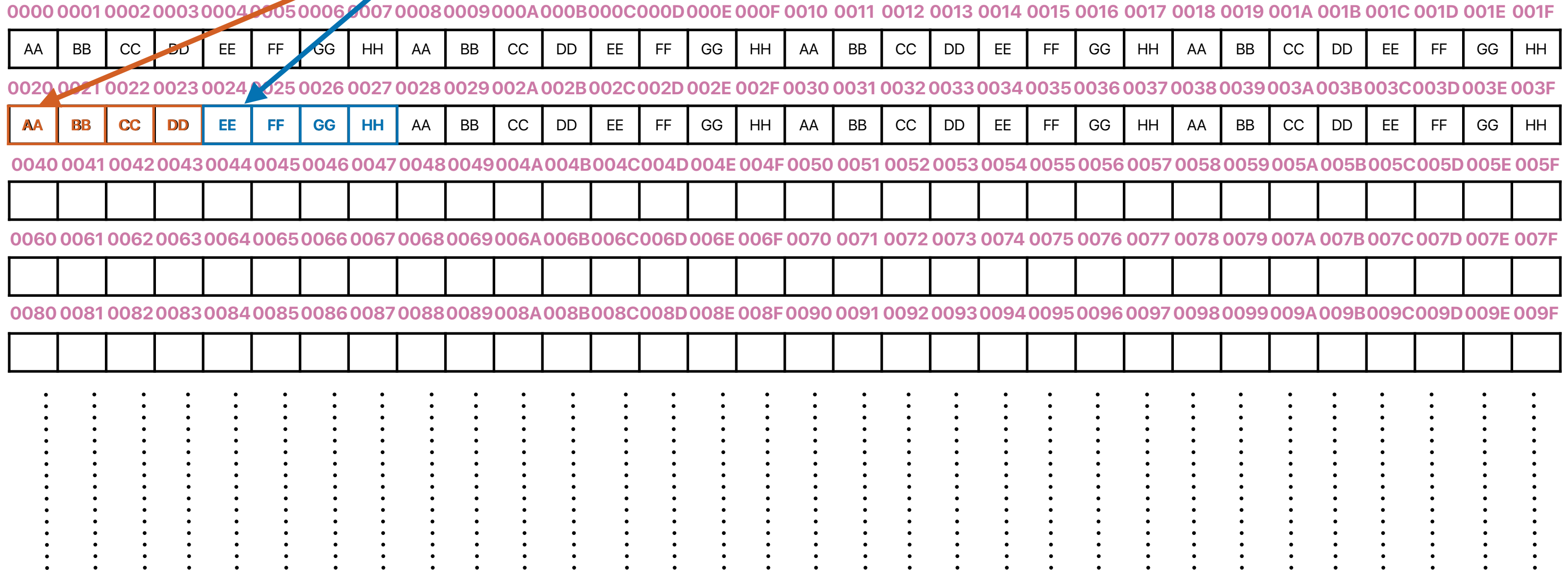


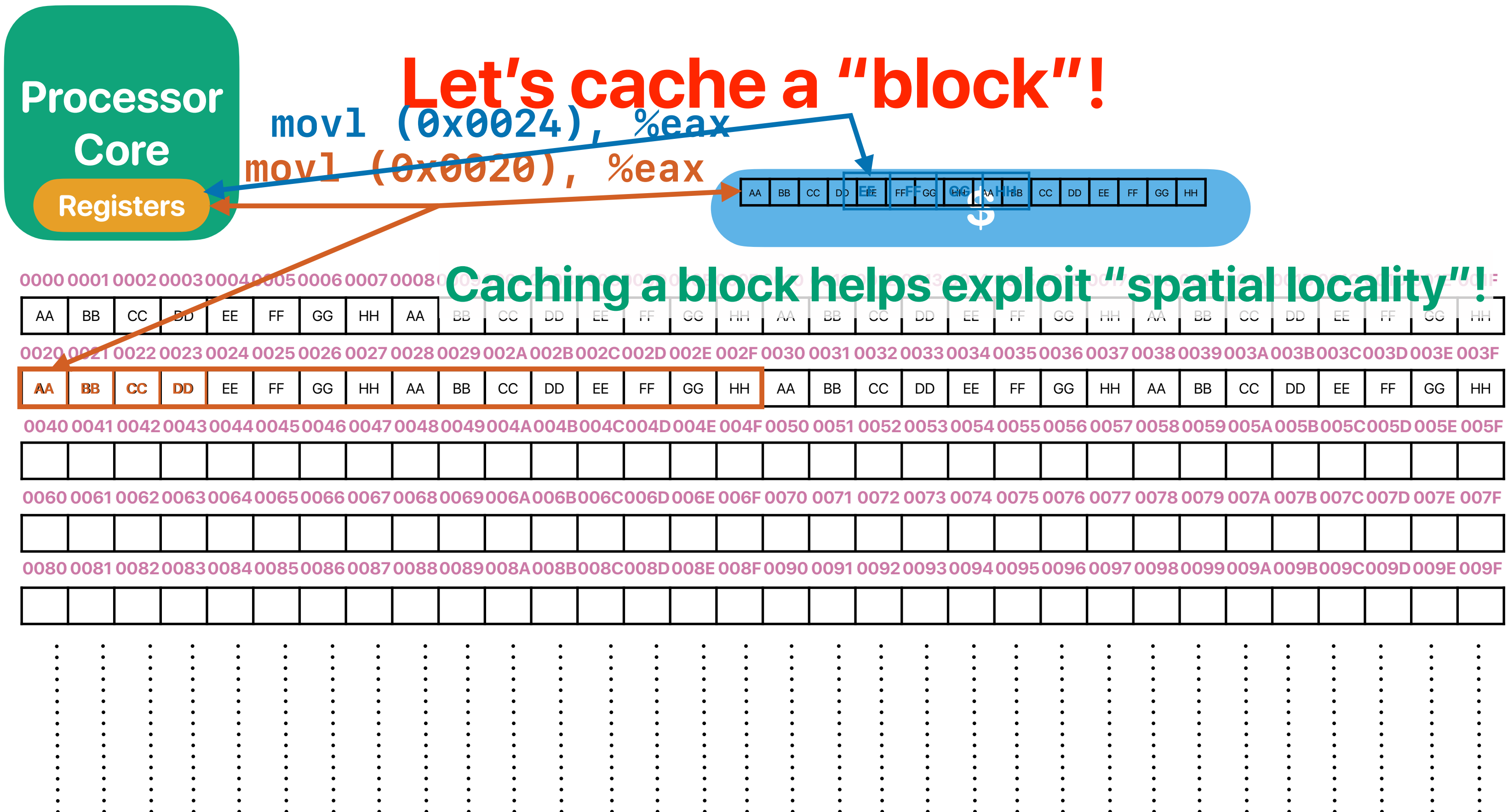
When there is a "movl"

movl (0x0024), %eax

movl (0x0020), %eax

Every "movl" has to visit the slow memory!





Recap: Locality

- Which description about locality of arrays `matrix` and `vector` in the following code is the **most accurate**?

```
for(uint32_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint32_t j = 0; j < n; j++) {  
        result += matrix[i][j]*vector[j];  
    }  
    output[i] = result;  
}
```

**Simply caching one block
isn't enough**

Designing a hardware to exploit locality

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

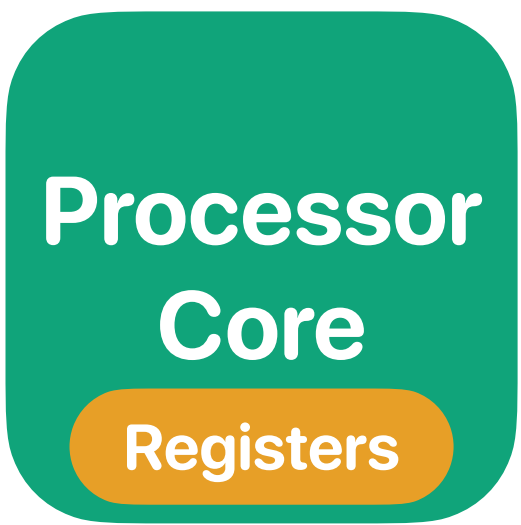
We need to “cache consecutive memory locations” every time — the cache should store a “block” of code/data

- Temporal locality — application revisit the same thing again and again

We need to “cache frequently used memory blocks”

— the cache should store a few blocks several KBs

— the cache must be able to distinguish blocks • Data — program can read/write the same data many times (e.g., vectors in matrix-vector product)



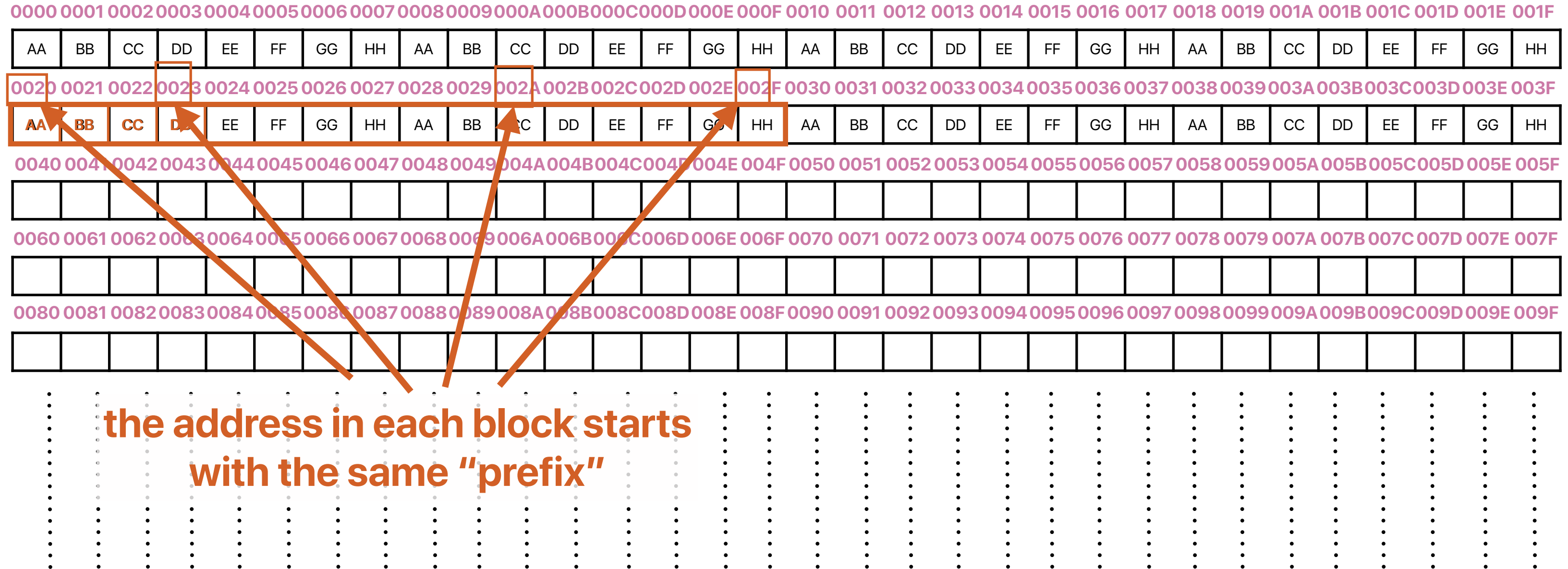
How to tell who is there?

?															
0123456789ABCDEF															
?															
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:															

Registers

Let's cache a "block"!

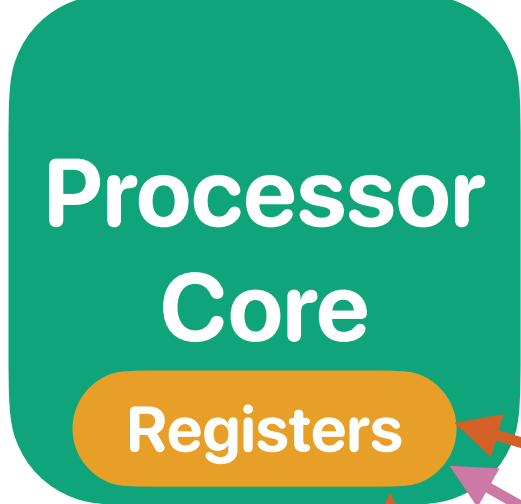
```
movl  (0x0024), %eax
```

~~movl (0x0020), %eax~~

Registers



tag array



How to tell w

block offset

tag

1w 0x0008

1w 0x4048

0x404 not found,
go to lower-level memory

Tell if the block here can be used

Tell if the block here is modified

Valid Bit

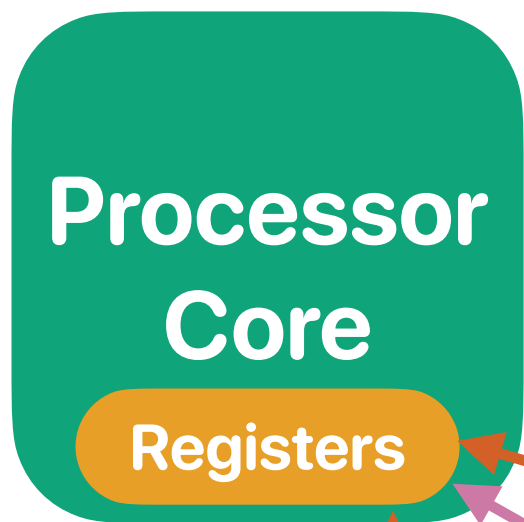
Dirty Bit

tag

data

		0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
	Valid Bit	Dirty Bit	tag		data												
	1	1	0x000		This is CSE13:												
	1	1	0x001		Advanced Computer Architecture!												
	1	0	0xF07		Advanced Computer Architecture!												
	0	1	0x100		This is CS 203:												
	1	1	0x310		Advanced Computer Architecture!												
	1	1	0x450		Advanced Computer Architecture!												
	0	1	0x006		This is CS 203:												
	0	1	0x537		Advanced Computer Architecture!												
	1	1	0x266		Advanced Computer Architecture!												
	1	1	0x307		This is CS 203:												
	0	1	0x265		Advanced Computer Architecture!												
	0	1	0x80A		Advanced Computer Architecture!												
	1	1	0x620		This is CS 203:												
	1	1	0x630		Advanced Computer Architecture!												
	1	0	0x705		Advanced Computer Architecture!												
	0	1	0x216		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

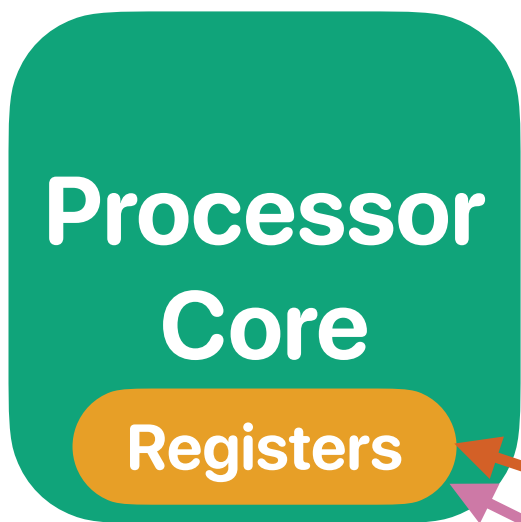
Valid Bit

Dirty Bit

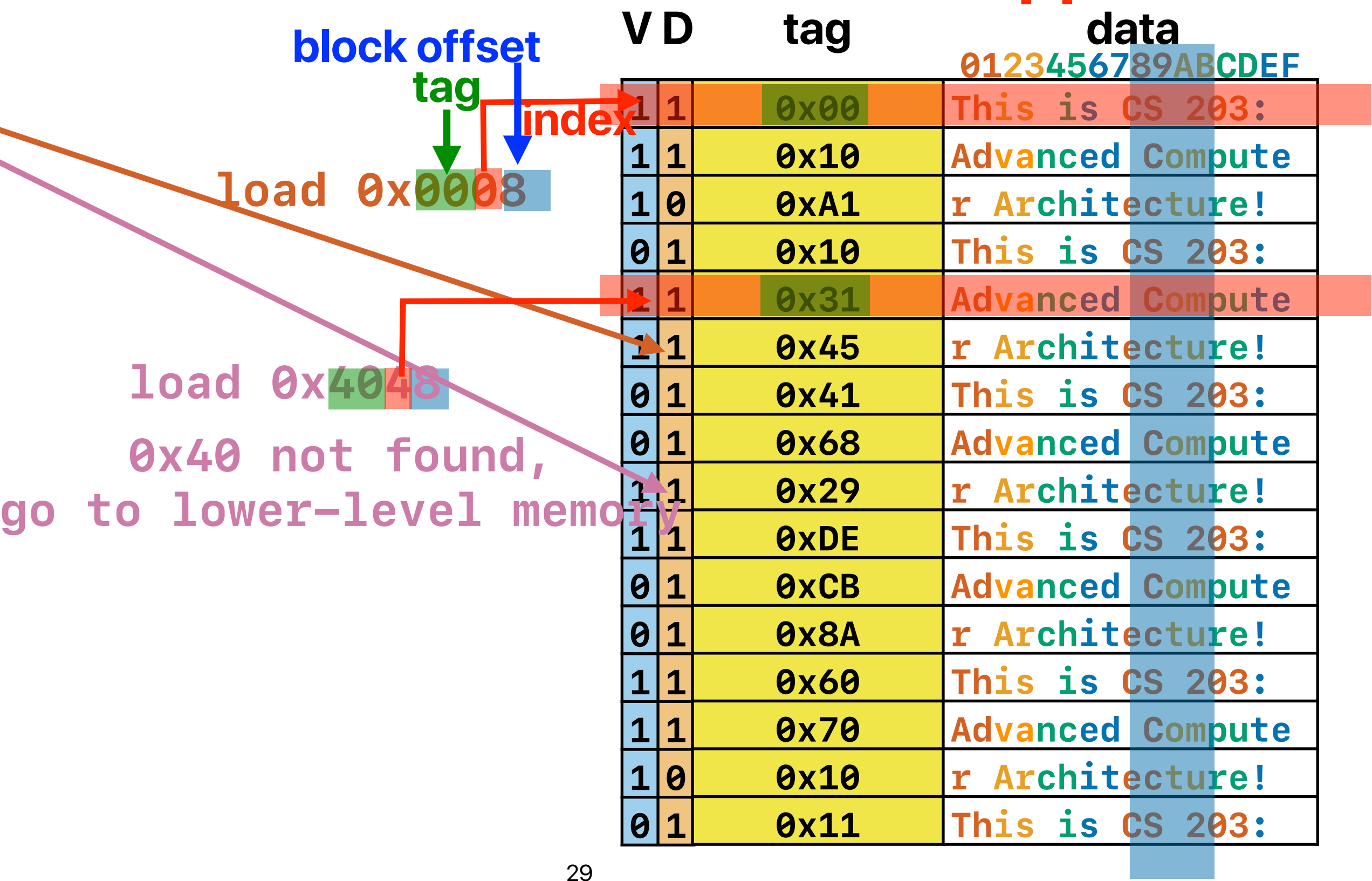
tag

data

				0123456789ABCDEF
1	1	0x000	This is CSE13:	
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:	



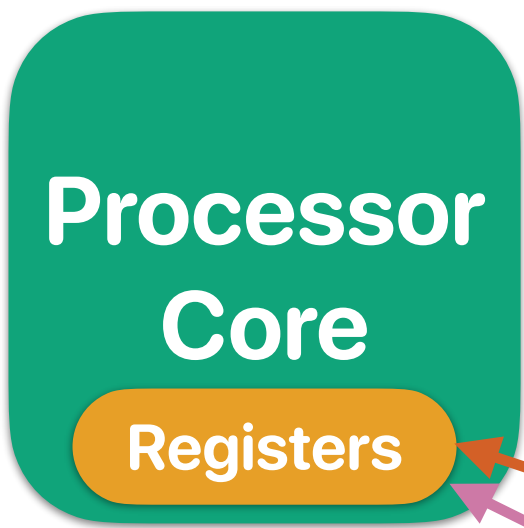
Hash-like structure — direct-mapped cache



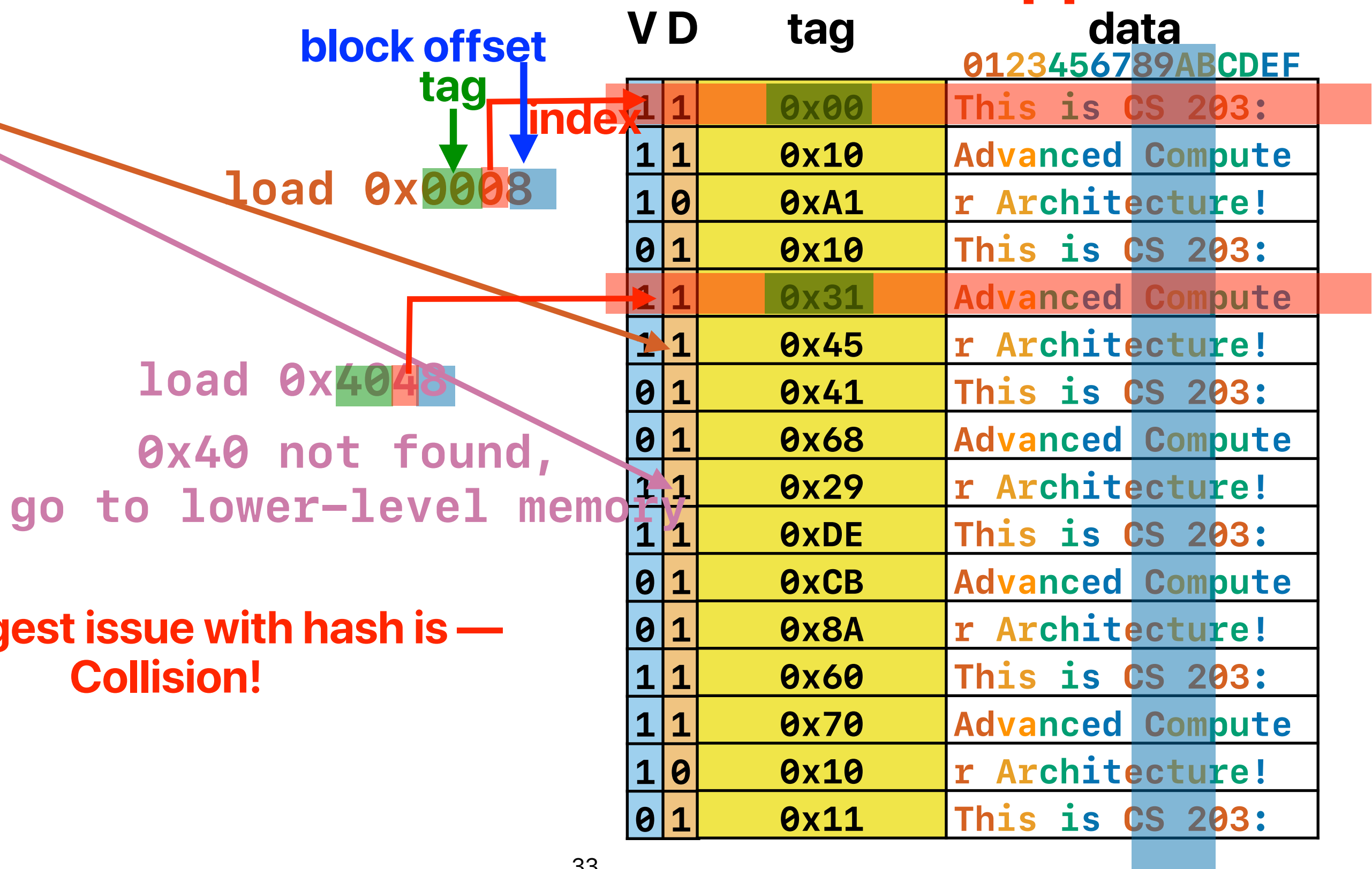
Blocksize == Linesize

```
[5]: # Your CS203 Cluster
! cs203 demo "lscpu | grep 'Model name'; getconf -a | grep CACHE"

ssh htseng@horsea " srun -N1 -p datahub lscpu | grep 'Model name'"
Model name:                  12th Gen Intel(R) Core(TM) i3-12100F
ssh htseng@horsea " srun -N1 -p datahub getconf -a | grep CACHE"
LEVEL1_ICACHE_SIZE           32768
LEVEL1_ICACHE_ASSOC           8
LEVEL1_ICACHE_LINESIZE       64
LEVEL1_DCACHE_SIZE           49152
LEVEL1_DCACHE_ASSOC           12
LEVEL1_DCACHE_LINESIZE       64
LEVEL2_CACHE_SIZE             1310720
LEVEL2_CACHE_ASSOC            10
LEVEL2_CACHE_LINESIZE        64
LEVEL3_CACHE_SIZE             12582912
LEVEL3_CACHE_ASSOC            12
LEVEL3_CACHE_LINESIZE        64
LEVEL4_CACHE_SIZE             0
LEVEL4_CACHE_ASSOC            0
LEVEL4_CACHE_LINESIZE        0
```



Hash-like structure — direct-mapped cache



The biggest issue with hash is — Collision!

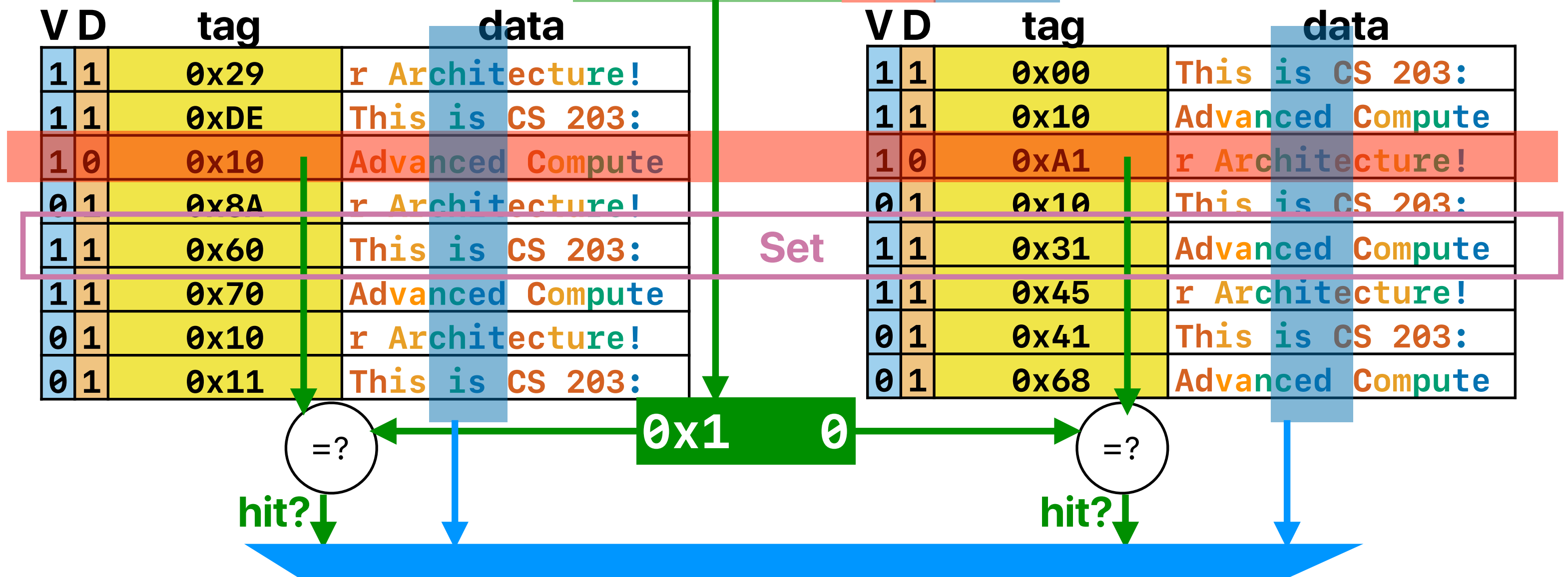
Way-associative cache

memory address: 0x0 8 2 4

 set block

 tag index offset

memory address: 0b00001000000100100



What is Associativity?

```
[5]: # Your CS203 Cluster
! cs203 demo "lscpu | grep 'Model name'; getconf -a | grep CACHE"

ssh htseng@horsea " srun -N1 -p datahub lscpu | grep 'Model name'"
Model name:                  12th Gen Intel(R) Core(TM) i3-12100F
ssh htseng@horsea " srun -N1 -p datahub getconf -a | grep CACHE"
LEVEL1_ICACHE_SIZE           32768
LEVEL1_ICACHE_ASSOC           8
LEVEL1_ICACHE_LINESIZE        64
LEVEL1_DCACHE_SIZE           49152
LEVEL1_DCACHE_ASSOC           12
LEVEL1_DCACHE_LINESIZE        64
LEVEL2_CACHE_SIZE             1310720
LEVEL2_CACHE_ASSOC            10
LEVEL2_CACHE_LINESIZE         64
LEVEL3_CACHE_SIZE             12582912
LEVEL3_CACHE_ASSOC            12
LEVEL3_CACHE_LINESIZE         64
LEVEL4_CACHE_SIZE             0
LEVEL4_CACHE_ASSOC            0
LEVEL4_CACHE_LINESIZE         0
```


The A, B, Cs of your cache

$$C = ABS$$

- **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 (Linesize)**
 - 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 N bits for N-bit machines (e.g., 64-bit for 64-bit machines)
- $(\text{address} / \text{block_size}) \% \mathbf{S} = \text{set index}$



NVIDIA Tegra X1

- L1 data (D-L1) cache configuration of NVIDIA Tegra X1 (used by Nintendo Switch and Jetson Nano)
 - Size 32KB, 4-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

A screenshot of a Pollev poll interface. It shows a list of five input fields, each preceded by a letter (A, B, C, D, E) in a small font. The input fields are empty, suggesting a multiple-choice or open-ended poll where users can select or enter an answer.

NVIDIA Tegra X1

- L1 data (D-L1) cache configuration of NVIDIA Tegra X1 (used by Nintendo Switch and Jetson Nano)
 - Size 32KB, 4-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

$$C = A \times B \times S$$

$$32 \times 1024 = 4 \times 64 \times S$$

$$S = 128$$

$$\text{Offset} = \log_2(64) = 6$$

$$\text{Index} = \log_2(128) = 7$$

$$\text{Tag} = 64 - 7 - 6 = 51$$



intel Core i7

- L1 data (D-L1) cache configuration of Core i7
 - Size 48KB, 12-way set associativity, 64B block
 - Assume 64-bit memory address
 - Which of the following is **NOT** correct?
 - A. Tag is 52 bits
 - B. Index is 6 bits
 - C. Offset is 6 bits
 - D. The cache has 128 sets
 - E. All of the above are correct

A screenshot of a Pollev poll interface. It shows a list of five input boxes, each preceded by a letter (A, B, C, D, E). The boxes are currently empty, indicating that no answers have been submitted yet.

intel Core i7

- L1 data (D-L1) cache configuration of Core i7
 - Size 48KB, 12-way set associativity, 64B block
 - Assume 64-bit memory address
 - Which of the following is **NOT** correct?
 - A. Tag is 52 bits
 - B. Index is 6 bits
 - C. Offset is 6 bits
 - D. The cache has 128 sets**
 - E. All of the above are correct

$$C = A \times B \times S$$

$$48 \times 1024 = 12 \times 64 \times S$$

$$S = 64$$

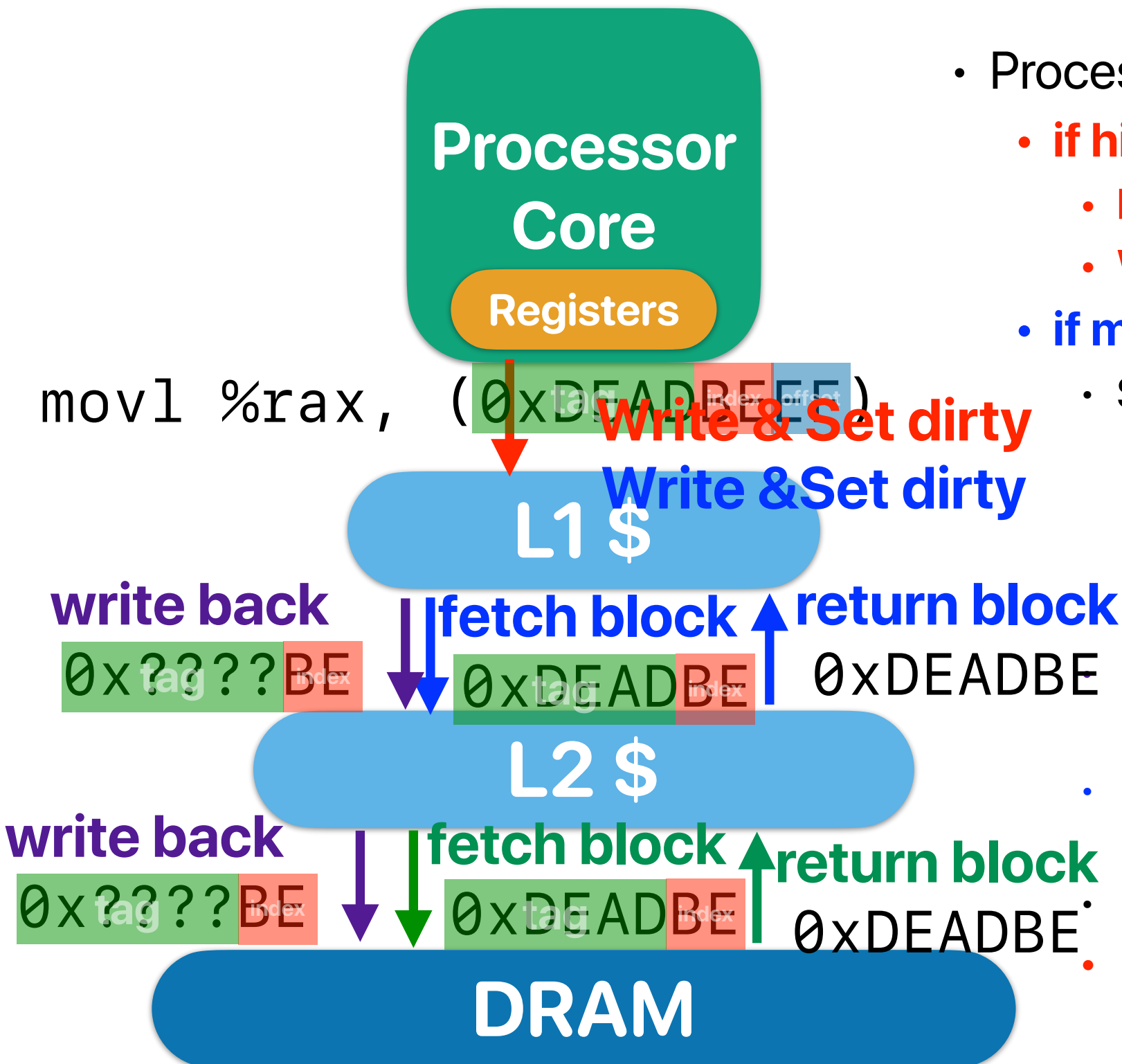
$$Offset = \log_2(64) = 6$$

$$Index = \log_2(12) = 3.58$$

$$Tag = 64 - 6 - 6 = 52$$

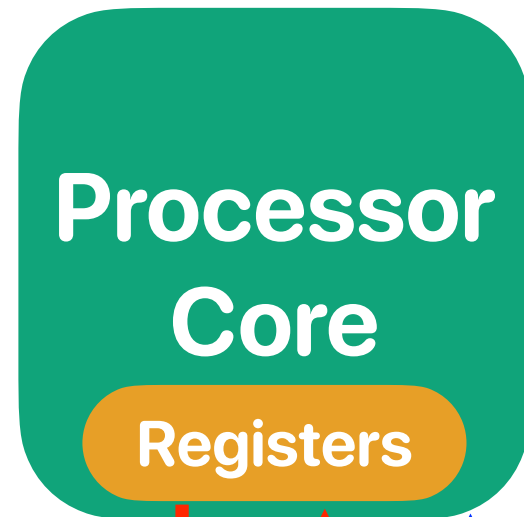
**Put everything all together:
How cache interacts with CPU**

The complete picture



- Processor sends memory access request to L1-\$
 - **if hit**
 - **Read - return data**
 - **Write - update & set DIRTY**
 - **if miss**
 - Select a victim block
 - If the target "set" is not full — select an empty/invalidated block as the victim block
 - If the target "set" is full — select a victim block using some policy
 - LRU is preferred — to exploit temporal locality!
 - If the victim block is "dirty" & "valid"
 - **Write back** the block to lower-level memory hierarchy
 - Fetch the requesting block from lower-level memory hierarchy and place in the victim block
 - If write-back or fetching causes any miss, repeat the same process
 - **Present the write "ONLY" in L1 and set DIRTY**

Processor/cache interaction



- Processor sends memory access request to L1-\$

- **if hit**
 - **return data**

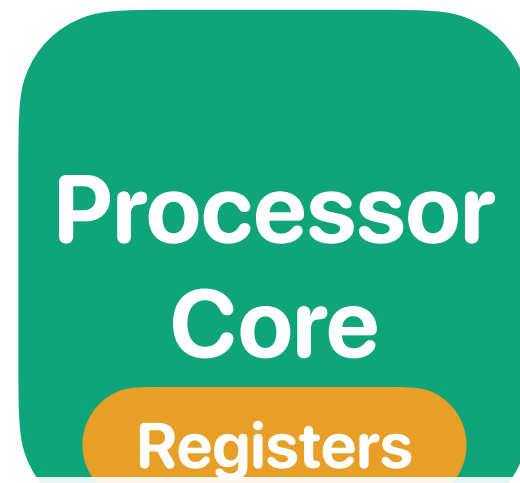
- **if miss**

- **fetch the requesting block from lower-level memory hierarchy and place in the cache**

What if we run out of \$ blocks?



Considering we have limited space in \$



- Processor sends memory access request to L1-\$
- **if hit**

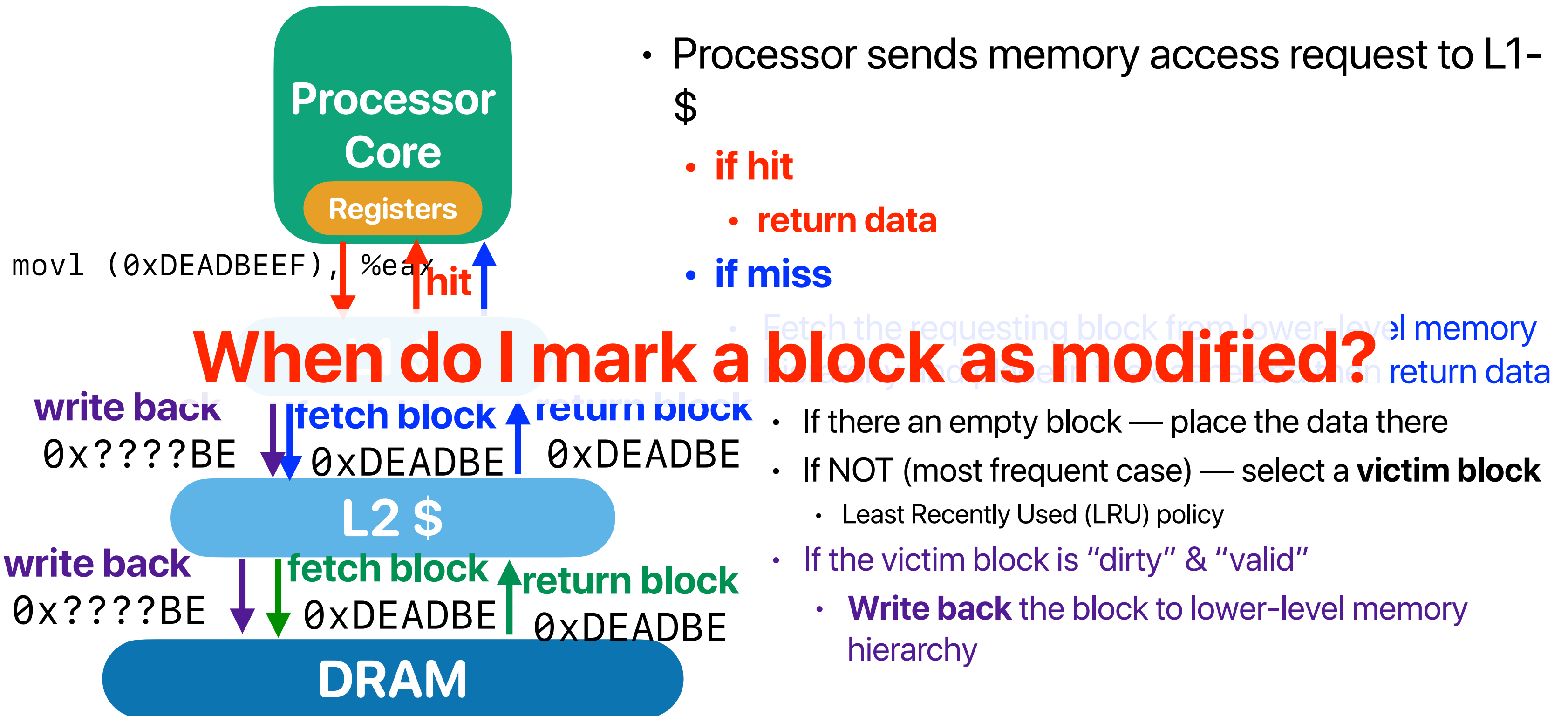
What if the victim block is modified?

— ignoring the update is not acceptable!

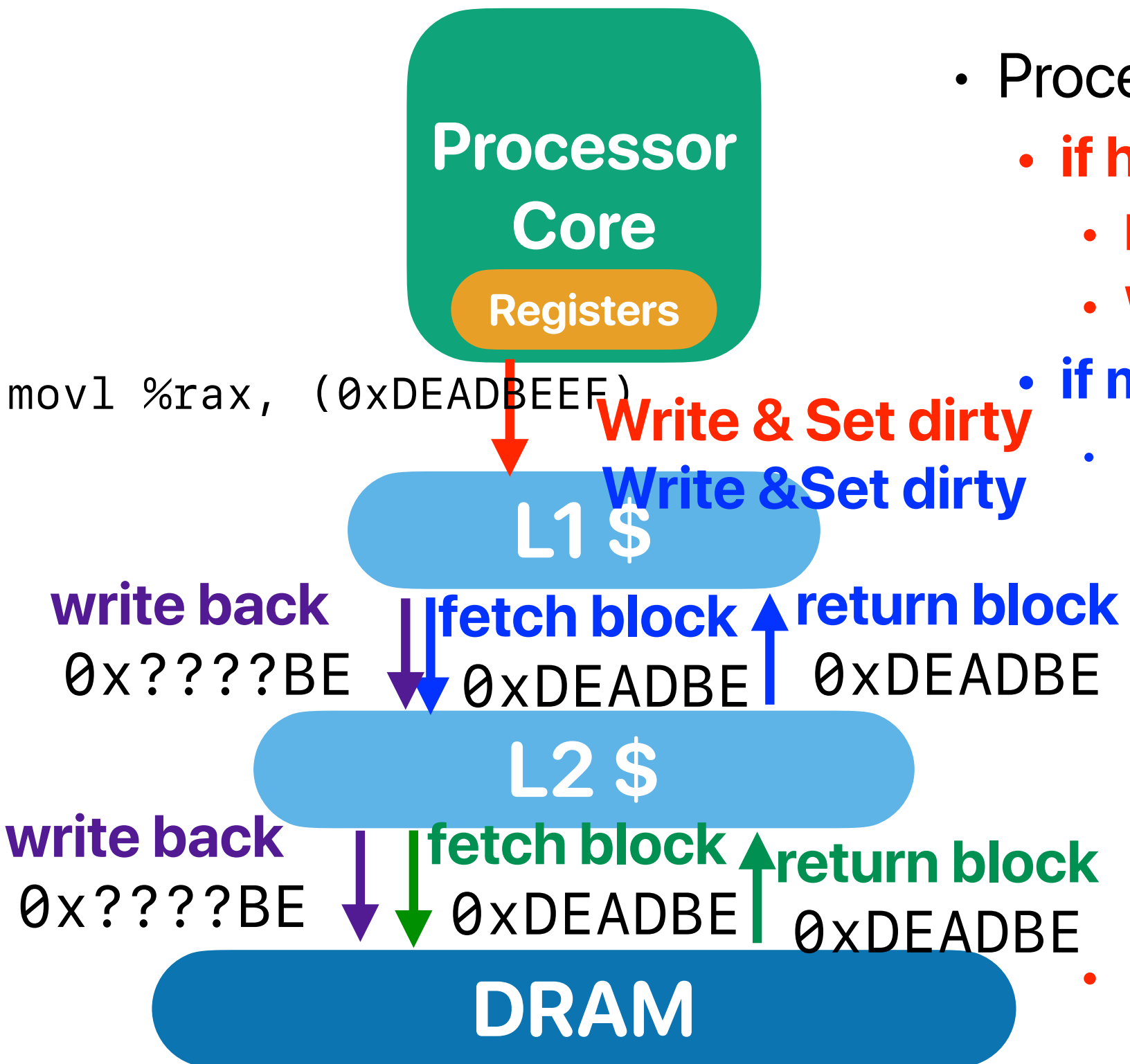


- If there an empty block — place the data there
- If NOT (most frequent case) — select a **victim block**
 - Least Recently Used (LRU) policy

Considering a victim block may be modified



Considering we have "writes"



- Processor sends memory access request to L1-\$
 - **if hit**
 - **Read: Return data**
 - **Write: Update "ONLY" in L1 and set DIRTY**
 - **if miss**
 - Fetch the requesting block from lower-level memory hierarchy and place in the cache and then return data
 - If there an empty block — place the data there
 - If NOT (most frequent case) — select a **victim block**
 - Least Recently Used (LRU) policy
 - If the victim block is "dirty" & "valid"
 - **Write back** the block to lower-level memory hierarchy
 - **Write: Update "ONLY" in L1 and set DIRTY**

Announcement

- Assignment #2 due **tonight**
 - Review the “demo”s of previous lectures if you need inspiration for programming assignment
 - You should run the experiments yourself and calculate results based on that
 - Everyone should have a different answer
 - The autograder won’t credit you if your answer does not contain working progress, equations in LaTeX format (if appropriate) and numbers from your experiments
 - We won’t give you credits if you don’t show your equations but only answers in the midterm either
 - If you consult your classmates, you need to put their names in the cell where you state your name
- Reading quiz #3 due **next Tuesday** before the lecture
- Assignment #3 will be ready this weekend
 - Please check our course website
<https://www.escalab.org/classes/cs203-2024sp/>

Computer Science & Engineering

203

つづく

