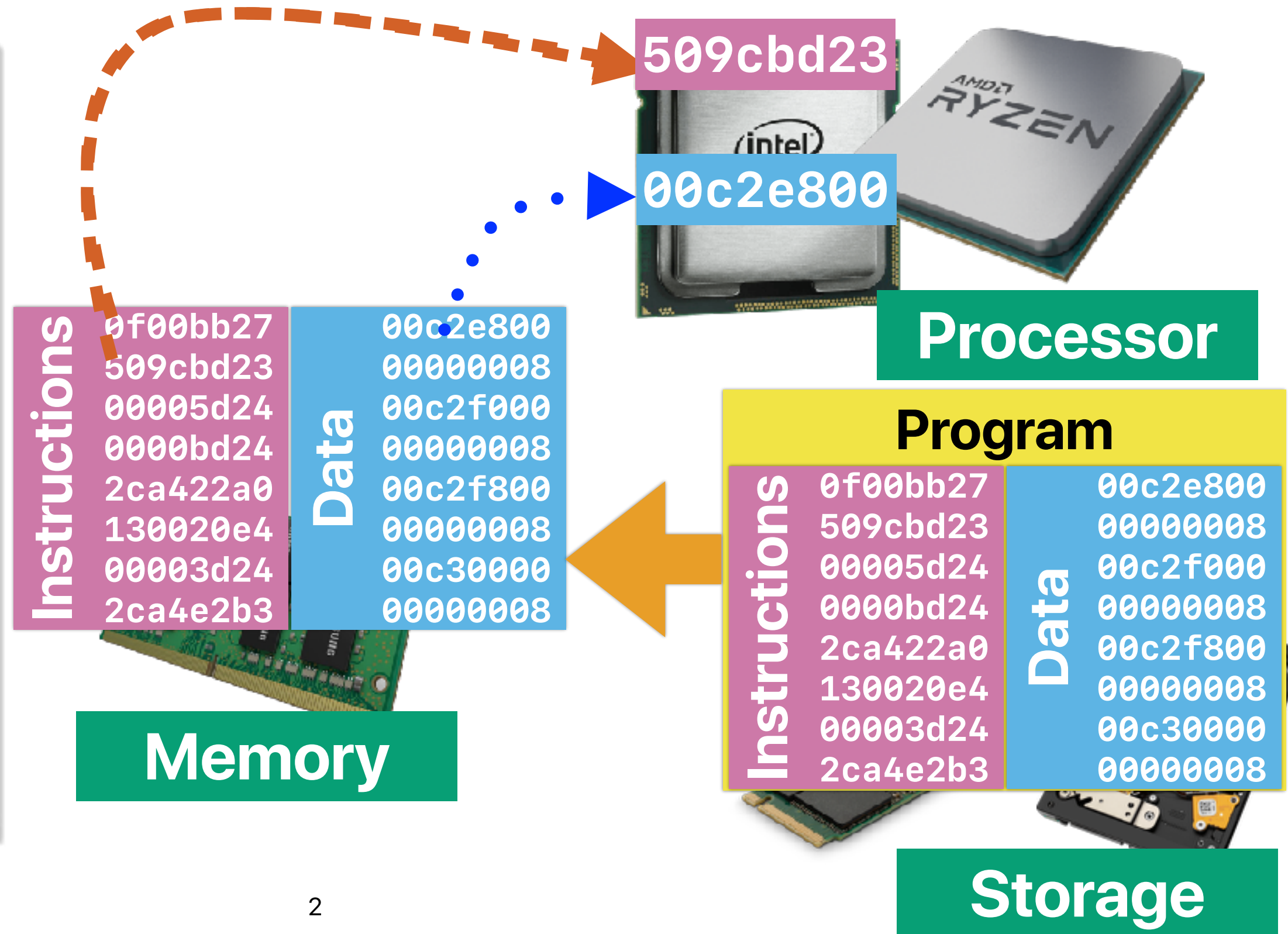


# Memory Hierarchy Inside Out:

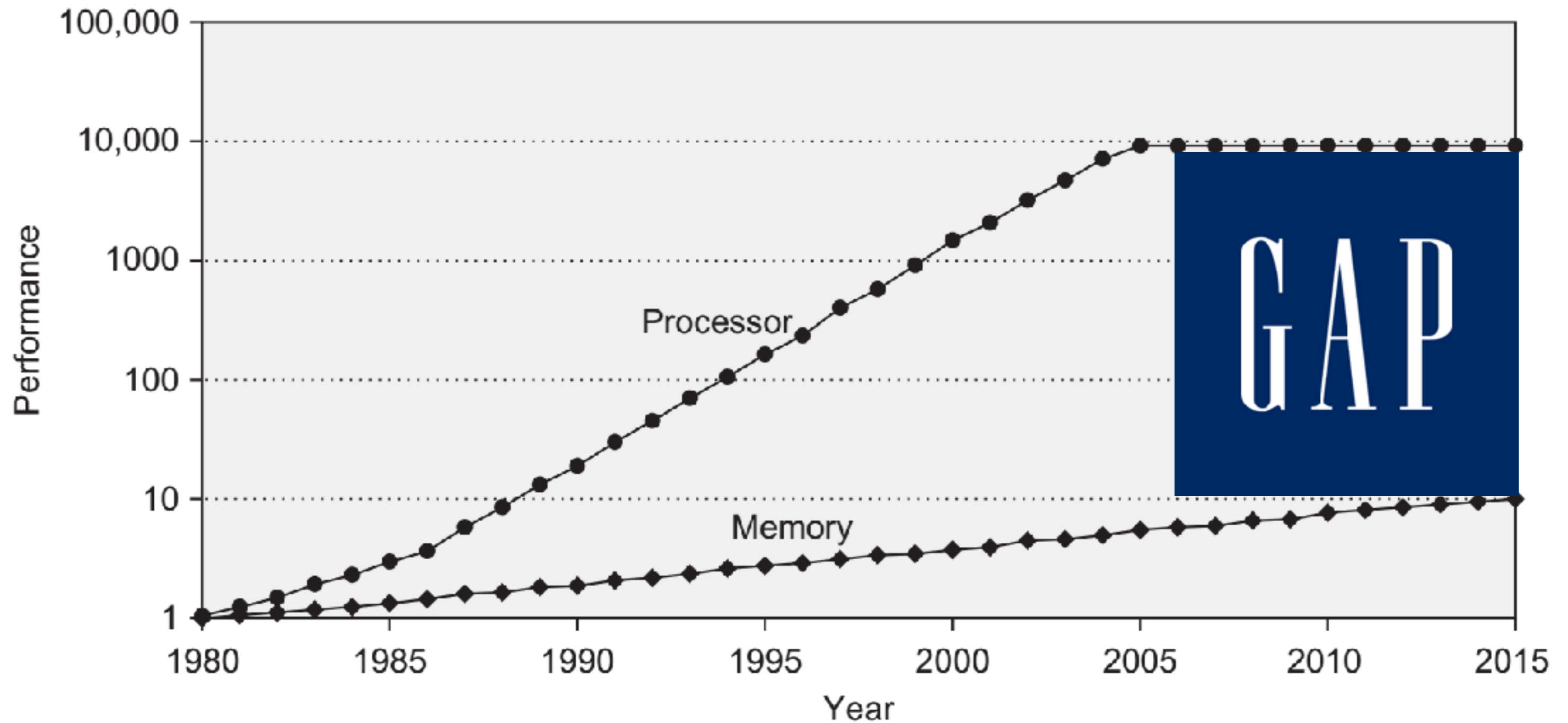
## (2) The A, B, C s of caches

Hung-Wei Tseng

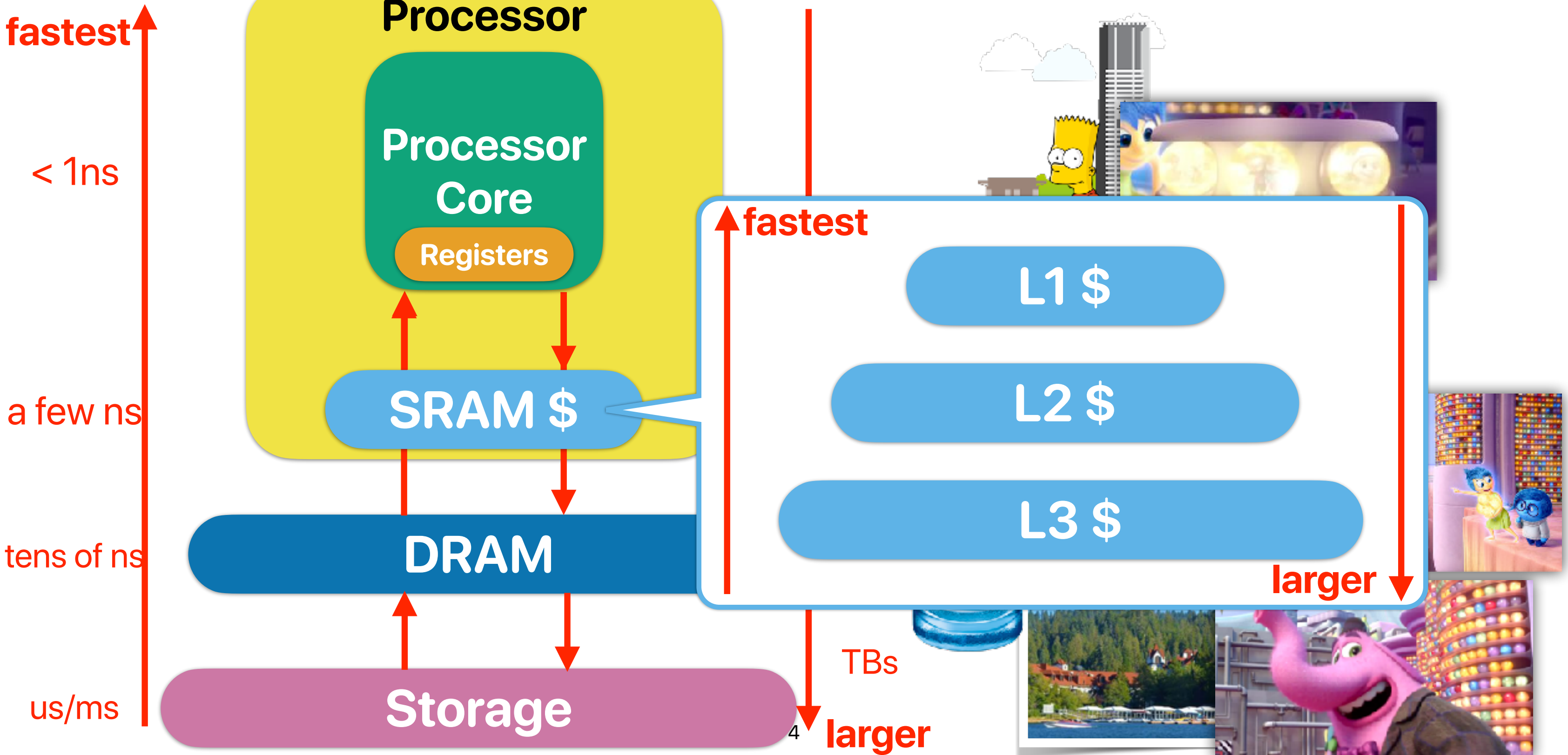
# von Neumann Architecture



# Recap: Performance gap between Processor/Memory



# Memory Hierarchy



# Outline

- Locality in applications (cont.)
- Architecting the cache
- The A, B, Cs of the cache

**Why adding small SRAMs would  
work?**

**Because of localities of memory references!**

# 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 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++;  
}
```



# How do you prepare closed-book exams?

- Review questions from prior years **Temporal locality**
- Review the whole chapter **Spatial locality**
- Practice similar questions **Spatial locality**
- Practice many times **Temporal locality**

# Locality

- Spatial locality — application tends to visit nearby stuffs in the memory
  - Code — the current instruction, and then the next

**Most of time, your program is just visiting a limited amount of data/instructions within a given timeframe**

- 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

# Locality and cache design

- The cache must be able to get chunks of near-by items every time to exploit spatial locality
- The cache must be able to keep a frequently used block for a while to exploit temporal locality

# **Architecting the Cache: capture the localities!**

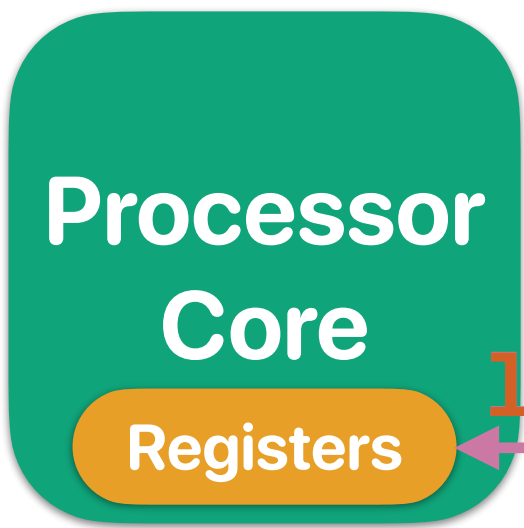
# Locality and cache design

- The cache must be able to get chunks of near-by items every time to exploit spatial locality

**We need to keep a block of data every time we put things in the cache**

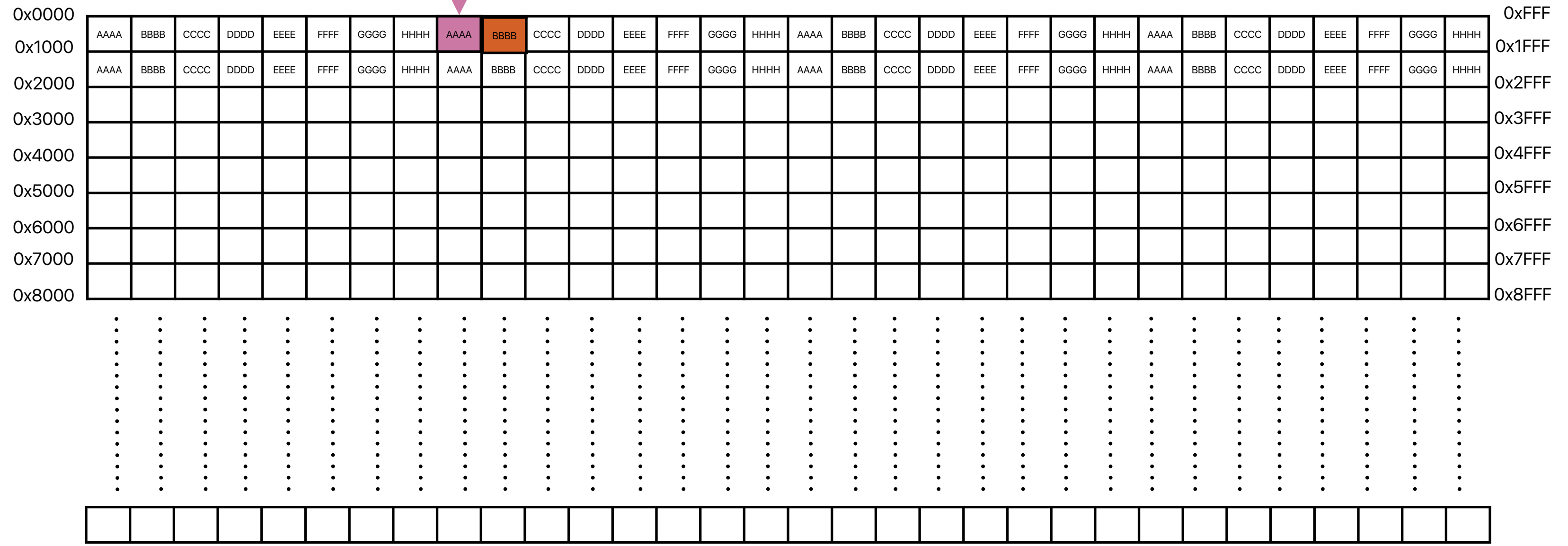
- The cache must be able to keep a frequently used block for a while to exploit temporal locality


**We need to keep multiple blocks of data in the cache**



# Load/store only access a "word" each time

load 0x000A






# Processor Core

## Registers

## Registers

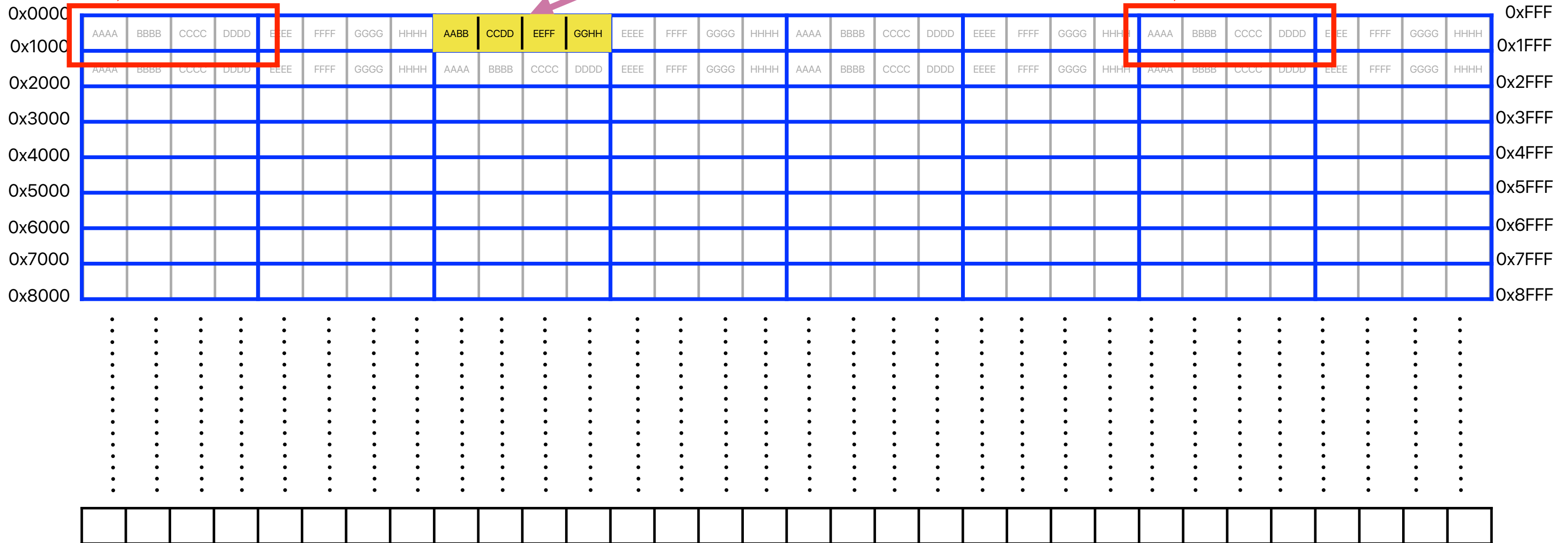
1w 0x0024



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"**





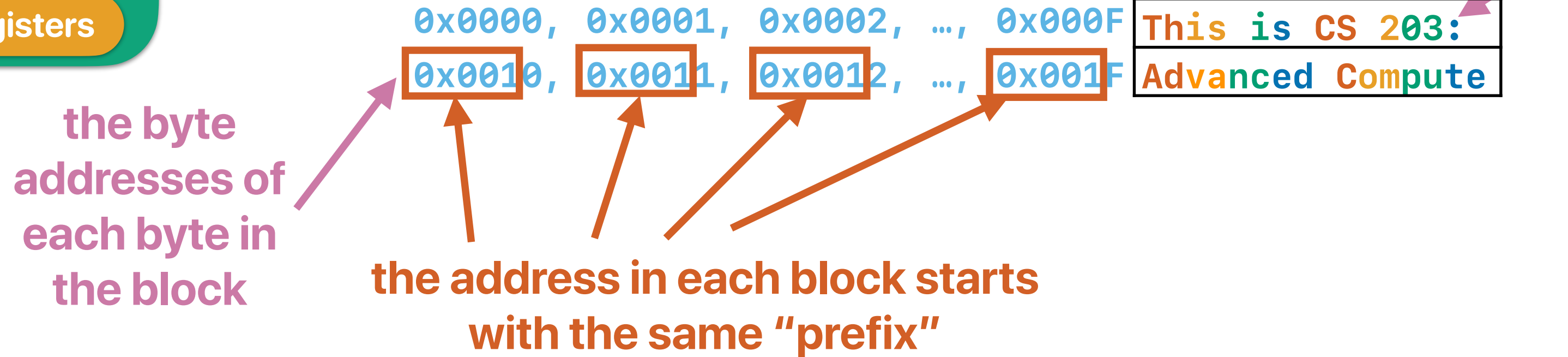
Processor  
Core

Registers

# What's a block?

the offset of the  
byte within a block

the data in  
memory





# How to tell who is there?

the common address  
prefix in each block

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

Processor  
Core

Registers

# 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

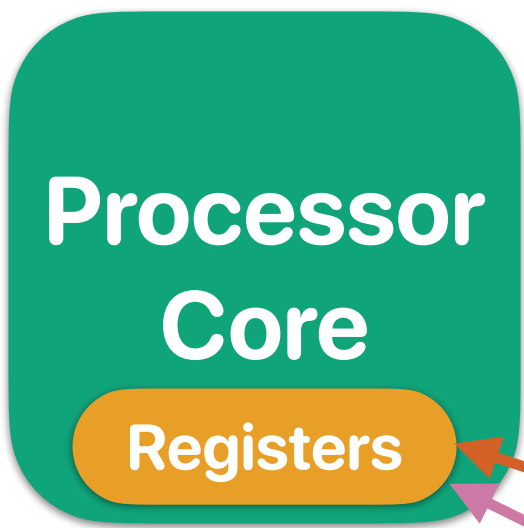
Tell if the block here is modified

Valid Bit  
Dirty Bit

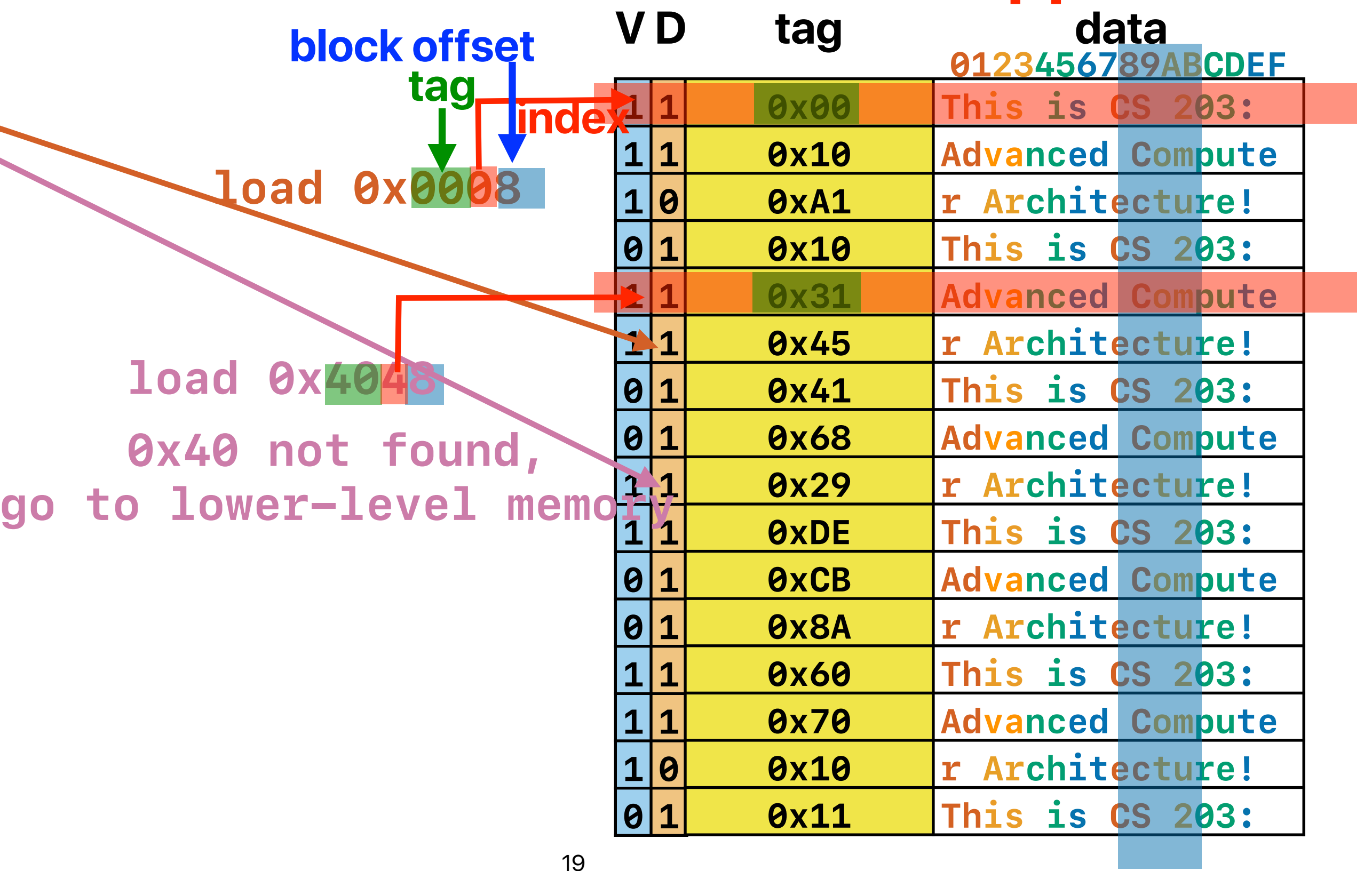
tag

data

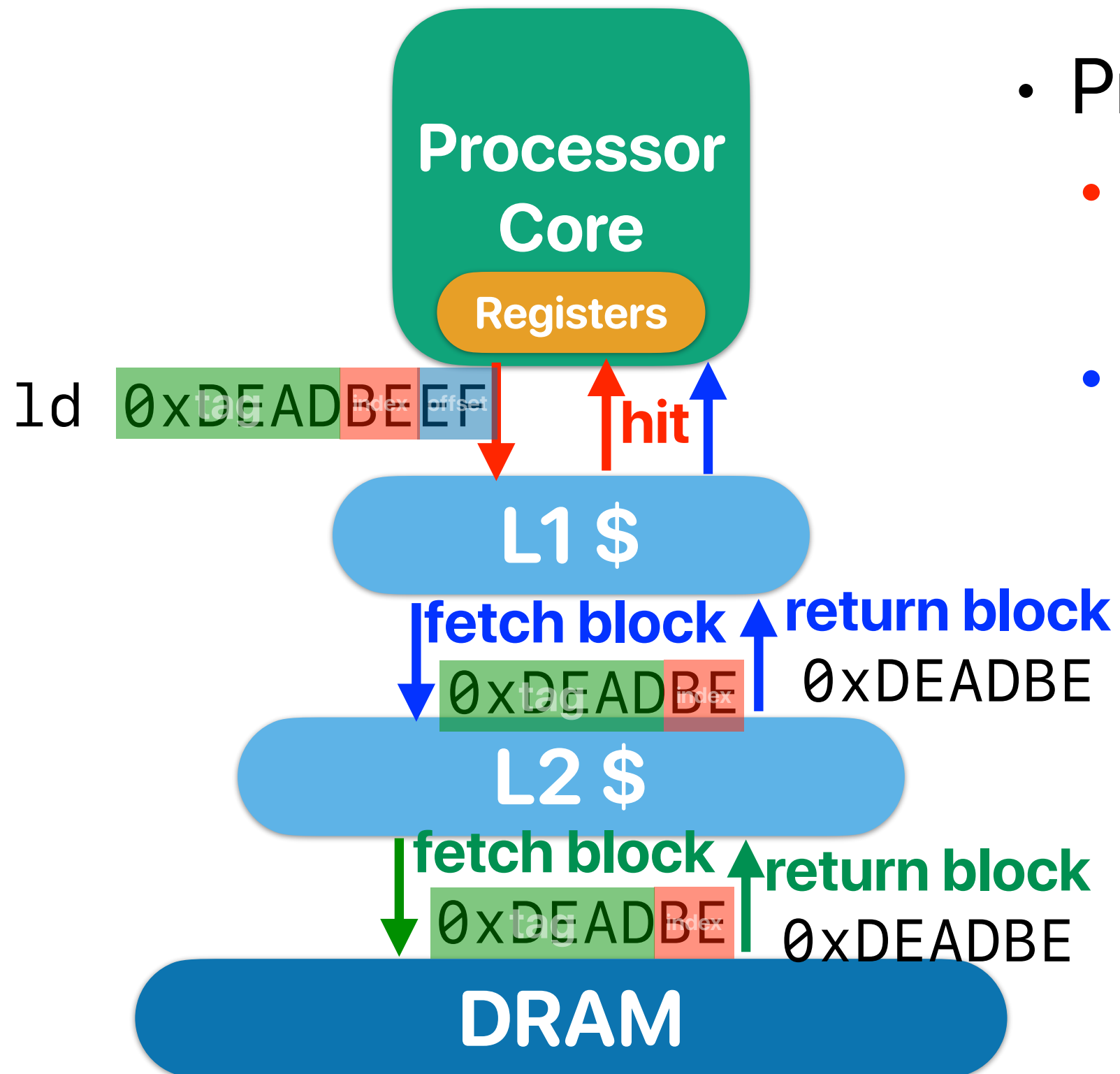
|   |   |       |                  | 0123456789ABCDEF |
|---|---|-------|------------------|------------------|
| 1 | 1 | 0x000 | This is CS 23:   |                  |
| 1 | 1 | 0x001 | Advanced Compute |                  |
| 1 | 0 | 0xF07 | r Architecture!  |                  |
| 0 | 1 | 0x100 | This 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



# What happens when we read data



- Processor sends load request to L1-\$
  - **if hit**
    - **return data**
  - **if miss**
    - Fetch a block
    - Select a victim block
      - If the target is not occupied — place the fetched block in the target location
      - If the target is full — select a victim block using some policy

**Let's simulate the simple cache!**

# Simulate a direct-mapped cache

- A direct mapped (1-way) cache with 256 bytes total capacity, a block size of 16 bytes

- # of blocks =  $\frac{256}{16} = 16$
- $\lg(16) = 4$  : 4 bits are used for the index
- $\lg(16) = 4$  : 4 bits are used for the byte offset
- The tag is  $64 - (4 + 4) = 56$  bits
- For example: 0x      8      0      0      0      0      0      8      0

= 0b1000 0000 0000 0000 0000 0000 0000 1000 0000





# Matrix vector revisited

```
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;  
}
```

# Matrix vector revisited

```
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;
}
```

tagindex

|          | Address (Hex)  | Address (Binary)                                   |
|----------|----------------|----------------------------------------------------|
| &a[0][0] | 0x558FE0A1D330 | 0b10101011000111111100000101000011101001100110000  |
| &b[0]    | 0x558FE0A1DC30 | 0b10101011000111111100000101000011101110000110000  |
| &a[0][1] | 0x558FE0A1D338 | 0b10101011000111111100000101000011101001100111000  |
| &b[1]    | 0x558FE0A1DC38 | 0b101010110001111111000001010000111011100000111000 |
| &a[0][2] | 0x558FE0A1D340 | 0b10101011000111111100000101000011101001101000000  |
| &b[2]    | 0x558FE0A1DC40 | 0b101010110001111111000001010000111011100001000000 |
| &a[0][3] | 0x558FE0A1D348 | 0b10101011000111111100000101000011101001101001000  |
| &b[3]    | 0x558FE0A1DC48 | 0b101010110001111111000001010000111011100001001000 |
| &a[0][4] | 0x558FE0A1D350 | 0b10101011000111111100000101000011101001101010000  |
| &b[4]    | 0x558FE0A1DC50 | 0b101010110001111111000001010000111011100001010000 |
| &a[0][5] | 0x558FE0A1D358 | 0b10101011000111111100000101000011101001101011000  |
| &b[5]    | 0x558FE0A1DC58 | 0b101010110001111111000001010000111011100001011000 |
| &a[0][6] | 0x558FE0A1D360 | 0b10101011000111111100000101000011101001101100000  |
| &b[6]    | 0x558FE0A1DC60 | 0b101010110001111111000001010000111011100001100000 |
| &a[0][7] | 0x558FE0A1D368 | 0b10101011000111111100000101000011101001101101000  |
| &b[7]    | 0x558FE0A1DC68 | 0b101010110001111111000001010000111011100001101000 |
| &a[0][8] | 0x558FE0A1D370 | 0b10101011000111111100000101000011101001101110000  |
| &b[8]    | 0x558FE0A1DC70 | 0b101010110001111111000001010000111011100001110000 |
| &a[0][9] | 0x558FE0A1D378 | 0b10101011000111111100000101000011101001101111000  |
| &b[9]    | 0x558FE0A1DC78 | 0b101010110001111111000001010000111011100001111000 |

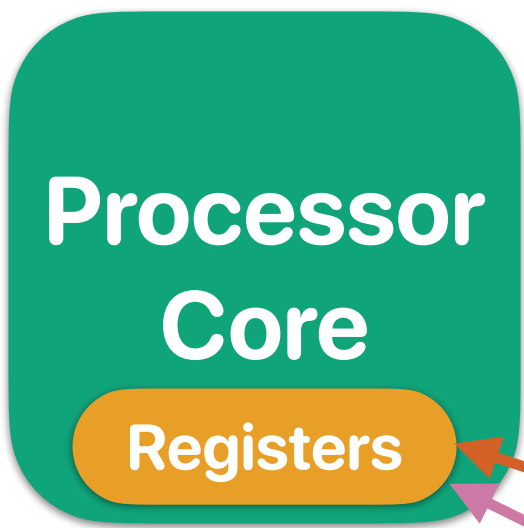
# Simulate a direct-mapped cache

tag index

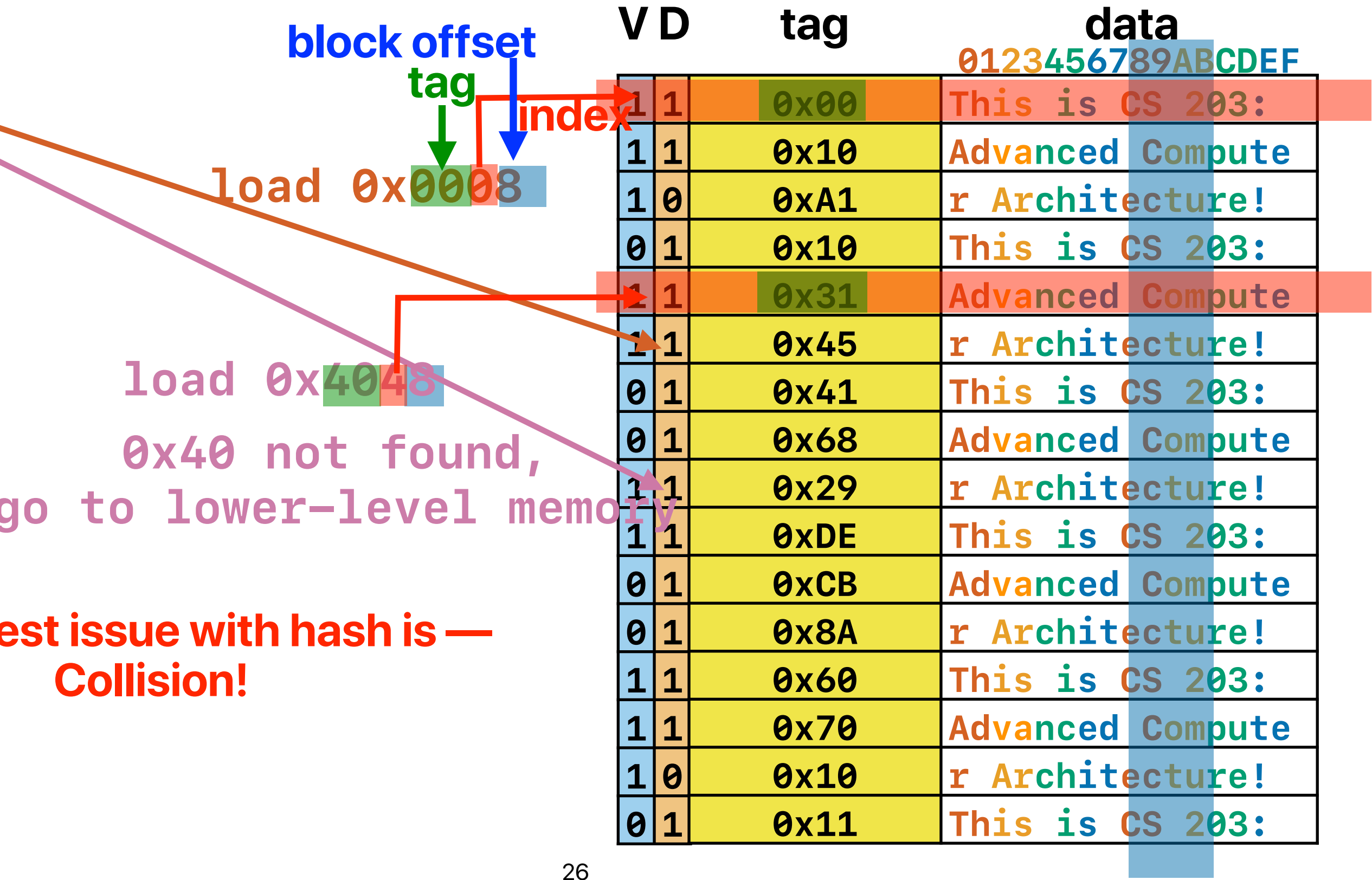
|    | V | D | Tag          | Data       |
|----|---|---|--------------|------------|
| 0  | 0 | 0 |              |            |
| 1  | 0 | 0 |              |            |
| 2  | 0 | 0 |              |            |
| 3  | 1 | 0 | 0x558FE0A1DC | b[0], b[1] |
| 4  | 1 | 0 | 0x558FE0A1DC | b[2], b[3] |
| 5  | 0 | 0 |              |            |
| 6  | 0 | 0 |              |            |
| 7  | 0 | 0 |              |            |
| 8  | 0 | 0 |              |            |
| 9  | 0 | 0 |              |            |
| 10 | 0 | 0 |              |            |
| 11 | 0 | 0 |              |            |
| 12 | 0 | 0 |              |            |
| 13 | 0 | 0 |              |            |
| 14 | 0 | 0 |              |            |
| 15 | 0 | 0 |              |            |

This cache doesn't work!!!  
— collisions!

| Address (Hex) |                |      |
|---------------|----------------|------|
| &a[0][0]      | 0x558FE0A1D330 | miss |
| &b[0]         | 0x558FE0A1DC30 | miss |
| &a[0][1]      | 0x558FE0A1D338 | miss |
| &b[1]         | 0x558FE0A1DC38 | miss |
| &a[0][2]      | 0x558FE0A1D340 | miss |
| &b[2]         | 0x558FE0A1DC40 | miss |
| &a[0][3]      | 0x558FE0A1D348 | miss |
| &b[3]         | 0x558FE0A1DC48 | miss |
| &a[0][4]      | 0x558FE0A1D350 | miss |
| &b[4]         | 0x558FE0A1DC50 | miss |
| &a[0][5]      | 0x558FE0A1D358 | miss |
| &b[5]         | 0x558FE0A1DC58 | miss |
| &a[0][6]      | 0x558FE0A1D360 | miss |
| &b[6]         | 0x558FE0A1DC60 | miss |
| &a[0][7]      | 0x558FE0A1D368 | miss |
| &b[7]         | 0x558FE0A1DC68 | miss |
| &a[0][8]      | 0x558FE0A1D370 | miss |
| &b[8]         | 0x558FE0A1DC70 | miss |
| &a[0][9]      | 0x558FE0A1D378 |      |
| &b[9]         | 0x558FE0A1DC78 |      |



# Hash-like structure — direct-mapped cache



The biggest issue with hash is — Collision!

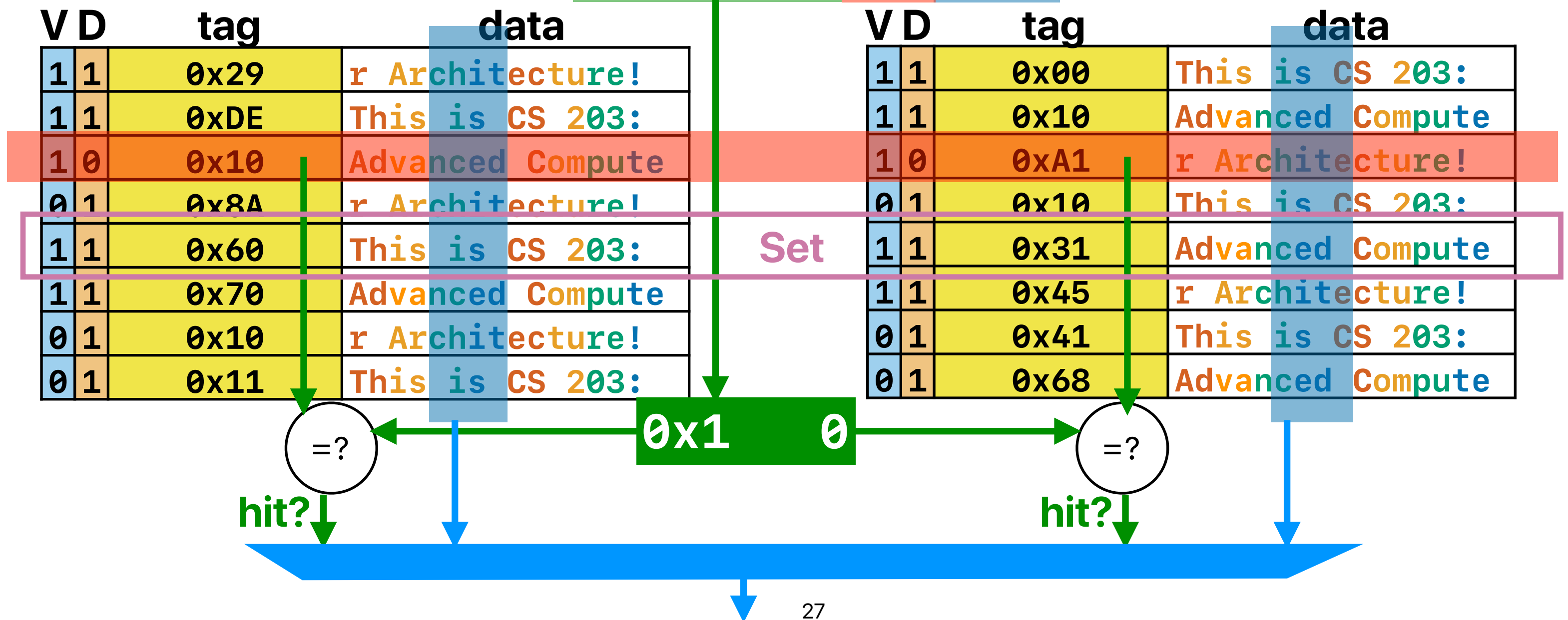
# Way-associative cache

memory address:  $0x0$  8 2 4

set block

tag index offset

memory address:  $0b00001000000100100$



# Now, 2-way, same-sized cache

- A 2-way cache with 256 bytes total capacity, a block size of 16 bytes

- # of blocks =  $\frac{256}{16} = 16$

- # of sets =  $\frac{16}{2} = 8$  (2-way: 2 blocks in a set)

- $\lg(8) = 3$  : 3 bits are used for the index

- $\lg(16) = 4$  : 4 bits are used for the byte offset

- The tag is  $64 - (4 + 4) = 56$  bits

- For example:  $0x \quad 8 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 8 \quad 0$   
=  $0b \quad 1000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 1000 \quad 0000$   

tag

index

offset

# Matrix vector revisited

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;  
}

tagindex

tag

index

|          | Address (Hex)  | Address (Binary)                                  |
|----------|----------------|---------------------------------------------------|
| &a[0][0] | 0x558FE0A1D330 | 0b10101011000111111100000101000011101001100110000 |
| &b[0]    | 0x558FE0A1DC30 | 0b10101011000111111100000101000011101110000110000 |
| &a[0][1] | 0x558FE0A1D338 | 0b10101011000111111100000101000011101001100111000 |
| &b[1]    | 0x558FE0A1DC38 | 0b10101011000111111100000101000011101110000111000 |
| &a[0][2] | 0x558FE0A1D340 | 0b10101011000111111100000101000011101001101000000 |
| &b[2]    | 0x558FE0A1DC40 | 0b10101011000111111100000101000011101110001000000 |
| &a[0][3] | 0x558FE0A1D348 | 0b10101011000111111100000101000011101001101001000 |
| &b[3]    | 0x558FE0A1DC48 | 0b10101011000111111100000101000011101110001001000 |
| &a[0][4] | 0x558FE0A1D350 | 0b10101011000111111100000101000011101001101010000 |
| &b[4]    | 0x558FE0A1DC50 | 0b10101011000111111100000101000011101110001010000 |
| &a[0][5] | 0x558FE0A1D358 | 0b10101011000111111100000101000011101001101011000 |
| &b[5]    | 0x558FE0A1DC58 | 0b10101011000111111100000101000011101110001011000 |
| &a[0][6] | 0x558FE0A1D360 | 0b10101011000111111100000101000011101001101100000 |
| &b[6]    | 0x558FE0A1DC60 | 0b10101011000111111100000101000011101110001100000 |
| &a[0][7] | 0x558FE0A1D368 | 0b10101011000111111100000101000011101001101101000 |
| &b[7]    | 0x558FE0A1DC68 | 0b10101011000111111100000101000011101110001101000 |
| &a[0][8] | 0x558FE0A1D370 | 0b10101011000111111100000101000011101001101110000 |
| &b[8]    | 0x558FE0A1DC70 | 0b10101011000111111100000101000011101110001110000 |
| &a[0][9] | 0x558FE0A1D378 | 0b10101011000111111100000101000011101001101111000 |
| &b[9]    | 0x558FE0A1DC78 | 0b10101011000111111100000101000011101110001111000 |



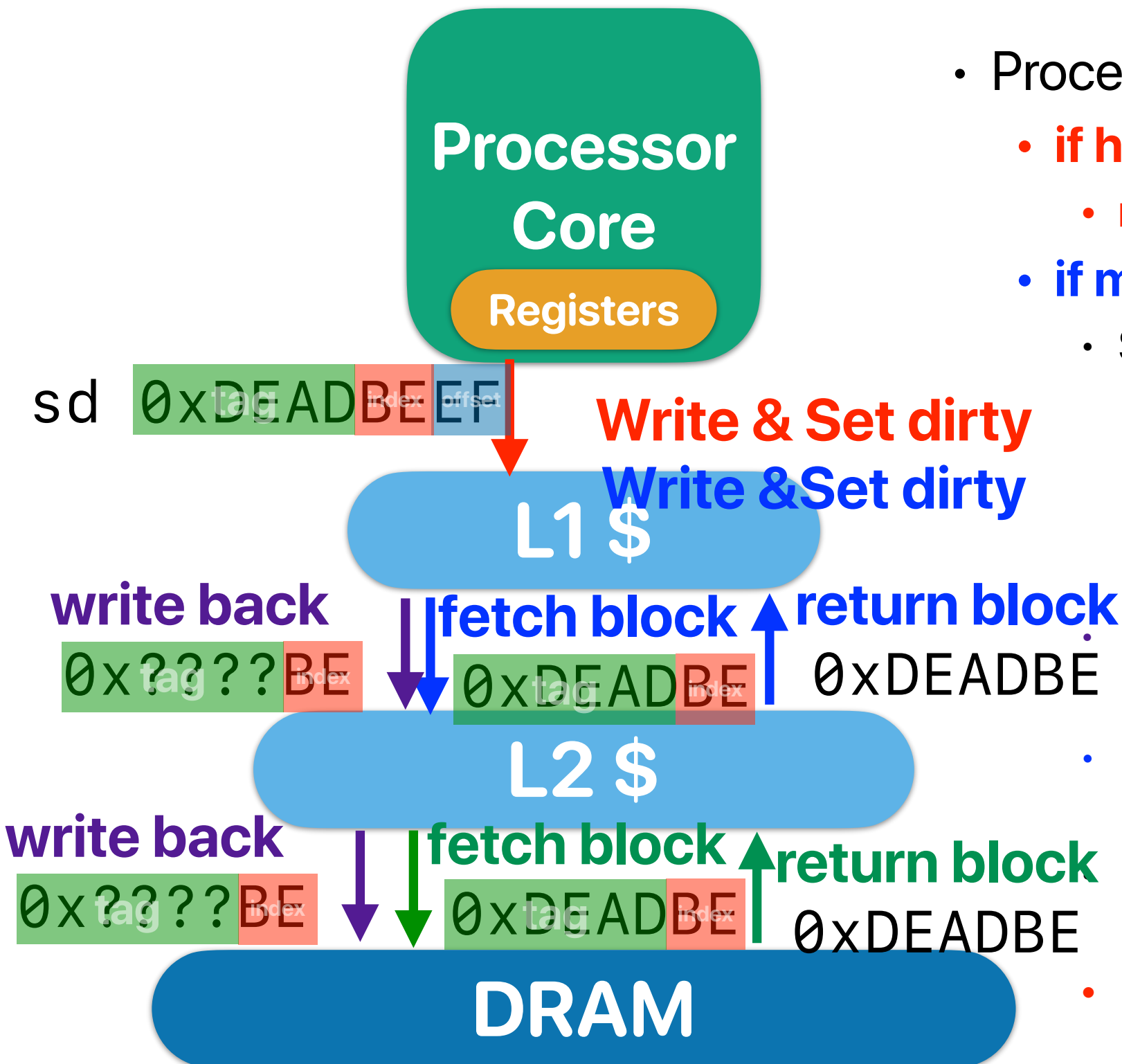
# Simulate a 2-way cache

| V | D | Tag          | Data             | V | D | Tag          | Data       |
|---|---|--------------|------------------|---|---|--------------|------------|
| 0 | 0 |              |                  | 0 | 0 |              |            |
| 0 | 0 |              |                  | 0 | 0 |              |            |
| 0 | 0 |              |                  | 0 | 0 |              |            |
| 1 | 0 | 0xAB1FC143A6 | a[0][0], a[0][1] | 1 | 0 | 0xAB1FC143B8 | b[0], b[1] |
| 1 | 0 | 0xAB1FC143A6 | a[0][2], a[0][3] | 1 | 0 | 0xAB1FC143B8 | b[2], b[3] |
| 0 | 0 |              |                  | 0 | 0 |              |            |
| 0 | 0 |              |                  | 0 | 0 |              |            |
| 0 | 0 |              |                  | 0 | 0 |              |            |

|          | Address (Hex)  | Tag          | Index |      |
|----------|----------------|--------------|-------|------|
| &a[0][0] | 0x558FE0A1D330 | 0xAB1FC143A6 | 0x3   | miss |
| &b[0]    | 0x558FE0A1DC30 | 0xAB1FC143B8 | 0x3   | miss |
| &a[0][1] | 0x558FE0A1D338 | 0xAB1FC143A6 | 0x3   | hit  |
| &b[1]    | 0x558FE0A1DC38 | 0xAB1FC143B8 | 0x3   | hit  |
| &a[0][2] | 0x558FE0A1D340 | 0xAB1FC143A6 | 0x4   | miss |
| &b[2]    | 0x558FE0A1DC40 | 0xAB1FC143B8 | 0x4   | miss |
| &a[0][3] | 0x558FE0A1D348 | 0xAB1FC143A6 | 0x4   | hit  |
| &b[3]    | 0x558FE0A1DC48 | 0xAB1FC143B8 | 0x4   | hit  |
| &a[0][4] | 0x558FE0A1D350 | 0xAB1FC143A6 | 0x5   | miss |
| &b[4]    | 0x558FE0A1DC50 | 0xAB1FC143B8 | 0x5   | miss |
| &a[0][5] | 0x558FE0A1D358 | 0xAB1FC143A6 | 0x5   | hit  |
| &b[5]    | 0x558FE0A1DC58 | 0xAB1FC143B8 | 0x5   | hit  |
| &a[0][6] | 0x558FE0A1D360 | 0xAB1FC143A6 | 0x6   | miss |
| &b[6]    | 0x558FE0A1DC60 | 0xAB1FC143B8 | 0x6   | miss |
| &a[0][7] | 0x558FE0A1D368 | 0xAB1FC143A6 | 0x6   | hit  |
| &b[7]    | 0x558FE0A1DC68 | 0xAB1FC143B8 | 0x6   | hit  |
| &a[0][8] | 0x558FE0A1D370 | 0xAB1FC143A6 | 0x7   | miss |
| &b[8]    | 0x558FE0A1DC70 | 0xAB1FC143B8 | 0x7   | miss |
| &a[0][9] | 0x558FE0A1D378 | 0xAB1FC143A6 | 0x7   | hit  |
| &b[9]    | 0x558FE0A1DC78 | 0xAB1FC143B8 | 0x7   | hit  |

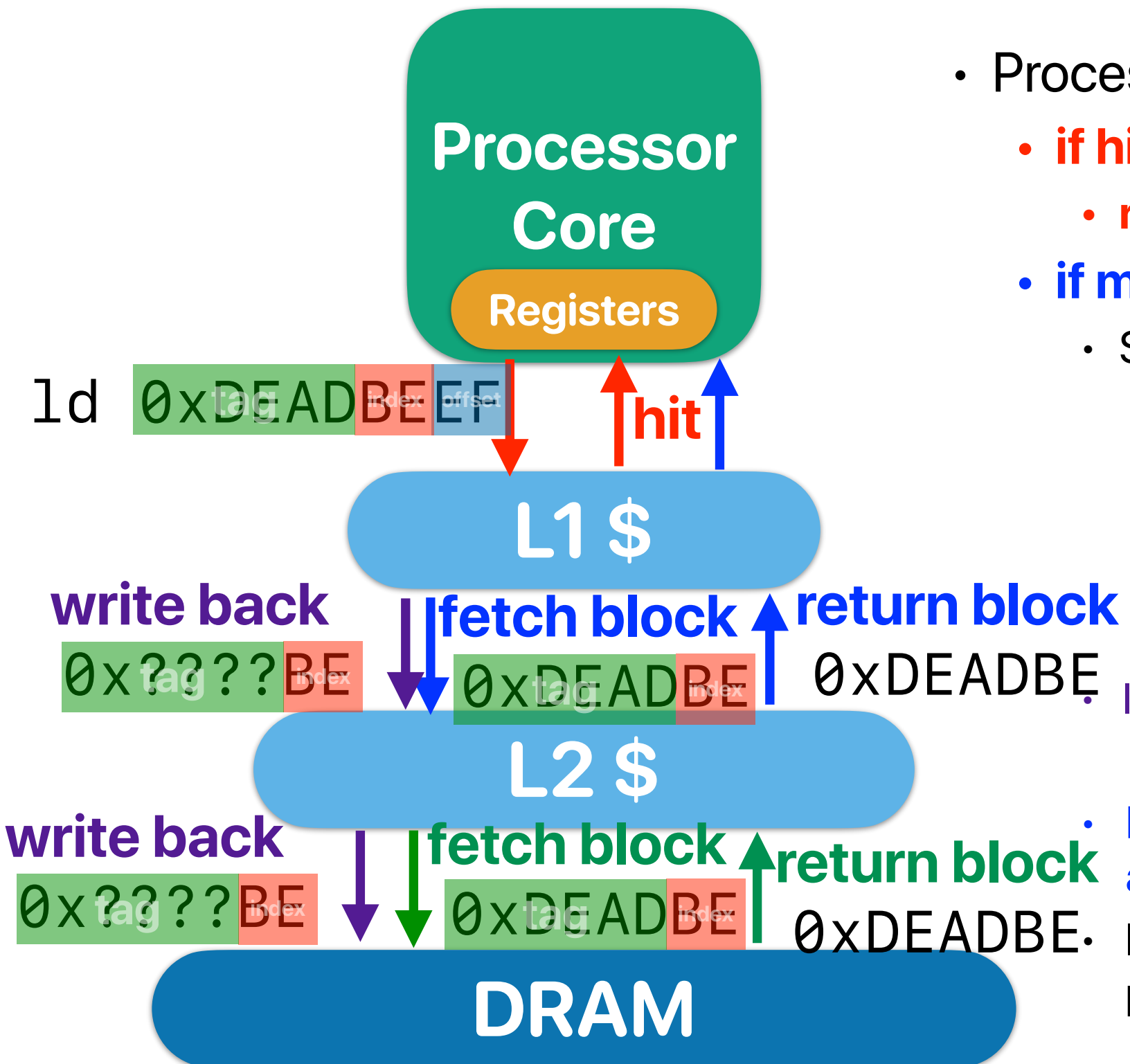
**Put everything all together:  
How cache interacts with CPU**

# What happens when we write data



- Processor sends load request to L1-\$
  - **if hit**
    - **return data — 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**

# What happens when we read data

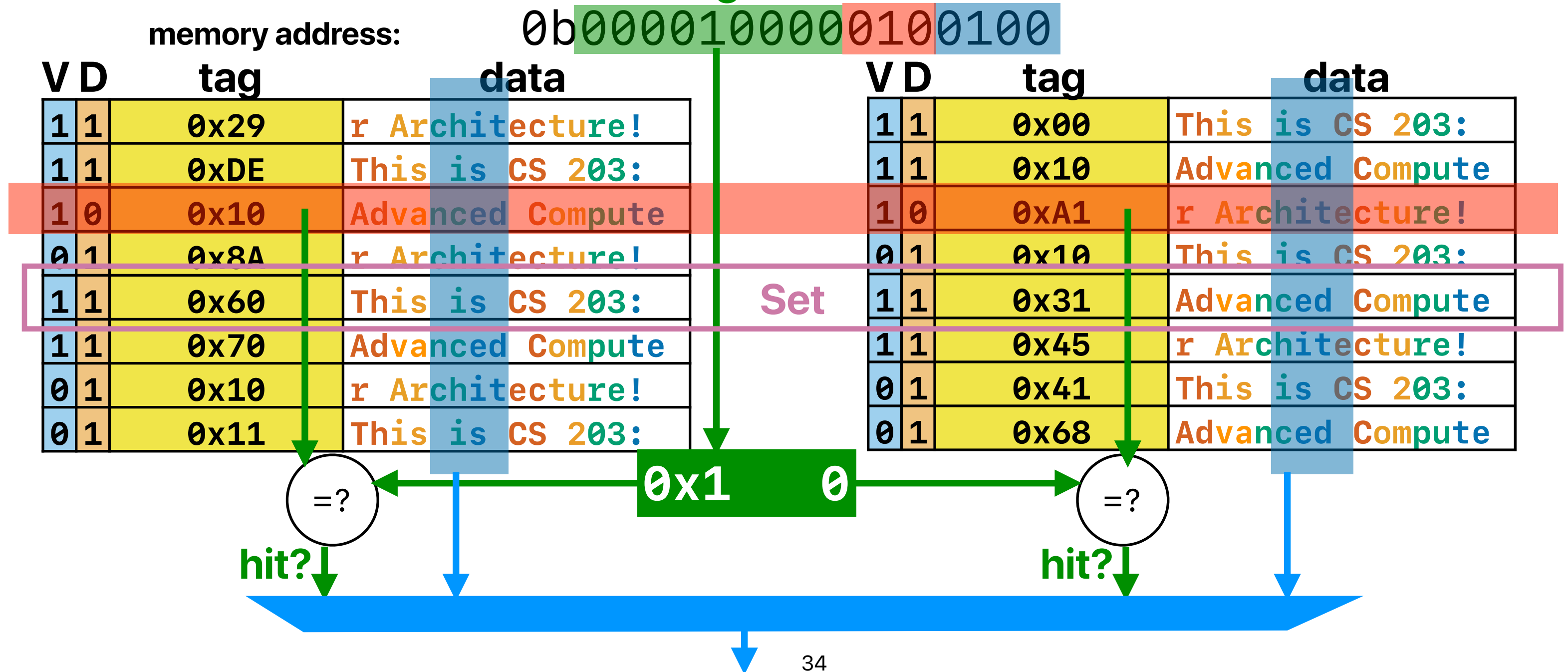


- Processor sends load request to L1-\$
  - **if hit**
    - **return data**
  - **if miss**
    - Select a victim block
      - If the target "set" is not full — select an empty/invalidated block as the victim block
      - If the target "set" is full — select a victim block using some policy
      - 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

# Way-associative cache

memory address: 0x0 8 2 4

tag set index block offset



# 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 (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 64 bits for 64-bit machine

- $$\frac{\text{address}}{\text{block\_size}} \pmod{S} = \text{set index}$$

# 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



# 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

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

$$32KB = 4 \times 64B \times S$$

$$S = \frac{32KB}{4 \times 64B} = 128$$

$$index = \log_2(128) = 7 \text{ bits}$$

$$offset = \log_2(64) = 6 \text{ bits}$$

$$tag = 64 - 7 - 6 = 51 \text{ bits}$$

# intel Core i7

- L1 data (D-L1) cache configuration of Core i7
  - Size 32KB, 8-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



# intel Core i7

- L1 data (D-L1) cache configuration of Core i7
  - Size 32KB, 8-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

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

$$32KB = 8 \times 64B \times S$$

$$S = \frac{32KB}{8 \times 64B} = 64$$

$$index = \log_2(64) = 6 \text{ bits}$$

$$offset = \log_2(64) = 6 \text{ bits}$$

$$tag = 64 - 6 - 6 = 52 \text{ bits}$$

# Announcement

- Regarding assignments
  - We will drop your lowest scored assignment
    - Please don't email to ask for extension — the dropping policy is to accommodate any potential reason for that
    - We don't accept late assignment
  - Please follow the EXACT instructions — any small thing you missed in the document can lead to undesirable outcome
  - Start early
    - We don't work 24/7 and we cannot help you last minute
    - Server could get busy last minute, too.
    - Gradescope has different test cases than released ones to prevent any shortcut of performance results — you have to test your code carefully to prevent failed execution on gradescope.
  - Assignment 2 is up tomorrow and please START EARLY
  - C++ programming — can't the demo regarding performance convince you to use C/C++?
  - Any kind of cheating is NOT ALLOWED and each assignment should be an individual work
    - Gradescope already identified several 100% similar ones — we will send those cases to Student Conduct & Academic Integrity Programs Office

| ↕ First & Last Name | ⇄ Swap | 📎 File  | ▼ Match Length | ↕ % Similarity | 📎 Top Source |
|---------------------|--------|---------|----------------|----------------|--------------|
| Si                  |        | sum.cpp | 96             | 100%           | S            |
| Si                  |        | sum.cpp | 96             | 98%            | S            |

# Announcement (cont.)

- Reading quiz due next Tuesday

# Computer Science & Engineering

# 203

# つづく

