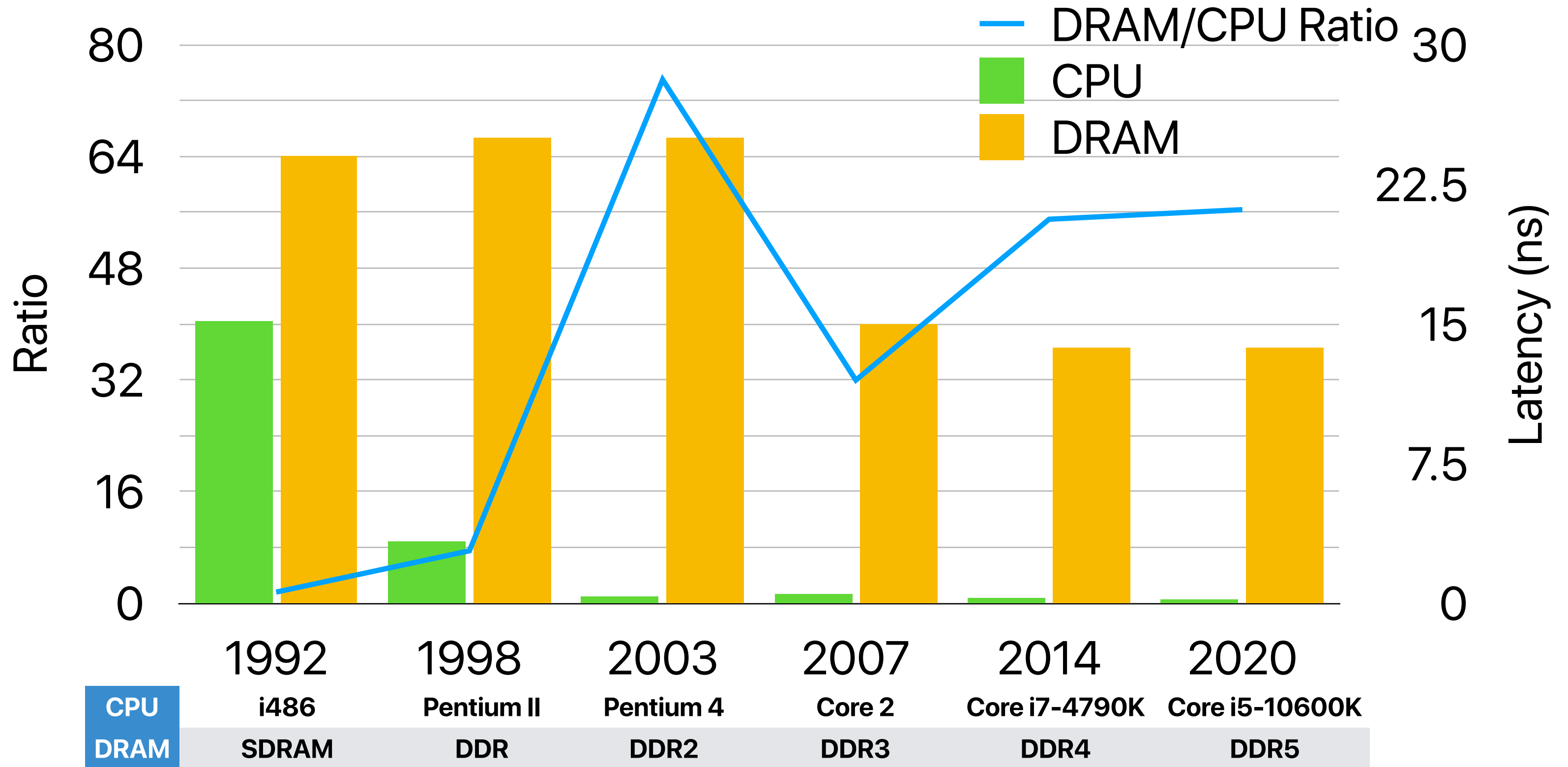


# Memory Hierarchy

Hung-Wei Tseng

# The "latency" gap between CPU and DRAM



# 20% is under-estimating ...

| Instruction class  | MIPS examples                     | HLL correspondence                                   | Frequency |         |
|--------------------|-----------------------------------|--|-----------|---------|
|                    |                                   |  | Integer   | Ft. pt. |
| Arithmetic         | add, sub, addi                    | Operations in assignment statements                  | 16%       | 48%     |
| Data transfer      | lw, sw, lb, lbu, lh, lhu, sb, lui | References to data structures, such as arrays        | 35%       | 36%     |
| Logical            | and, or, nor, andi, ori, sll, srl | Operations in assignment statements                  | 12%       | 4%      |
| Conditional branch | beq, bne, slt, slti, sltiu        | If statements and loops                              | 34%       | 8%      |
| Jump               | j, jr, jal                        | Procedure calls, returns, and case/switch statements | 2%        | 0%      |

**FIGURE 2.48** MIPS instruction classes, examples, correspondence to high-level program language constructs, and percentage of MIPS instructions executed by category for the average integer and floating point SPEC CPU2006 benchmarks.

Figure 3.24 in Chapter 3 shows average percentage of the individual MIPS instructions executed.

# Recap: Speedup and Amdahl's Law?

- Definition of "Speedup of Y over X" or say Y is n times faster than X:  $speedup_{Y\_over\_X} = n = \frac{Execution\ Time_X}{Execution\ Time_Y}$

- Amdahl's Law —  $Speedup_{enhanced}(f, s) = \frac{1}{(1-f) + \frac{f}{s}}$   $Speedup_{max}(f, \infty) = \frac{1}{(1-f)}$

- Corollary 1 — each optimization has an upper bound

- Corollary 2 — make the common case (the most time consuming case) fast!

$$\begin{aligned} Speedup_{max}(f_1, \infty) &= \frac{1}{(1-f_1)} \\ Speedup_{max}(f_2, \infty) &= \frac{1}{(1-f_2)} \\ Speedup_{max}(f_3, \infty) &= \frac{1}{(1-f_3)} \\ Speedup_{max}(f_4, \infty) &= \frac{1}{(1-f_4)} \end{aligned}$$

- Corollary 3: Optimization has a moving target

- Corollary 4: Exploiting more parallelism from a program is the key to performance gain in modern architectures

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

- Corollary 5: Single-core performance still matters

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

# Take-aways: inside out our memory hierarchy

- Memory access time is the most critical performance problem
  - One memory operation is as expensive as 50 arithmetic operations
  - Processor has to fetch instructions from memory
  - We have an average of 33% of data memory access instructions!

# Alternatives?

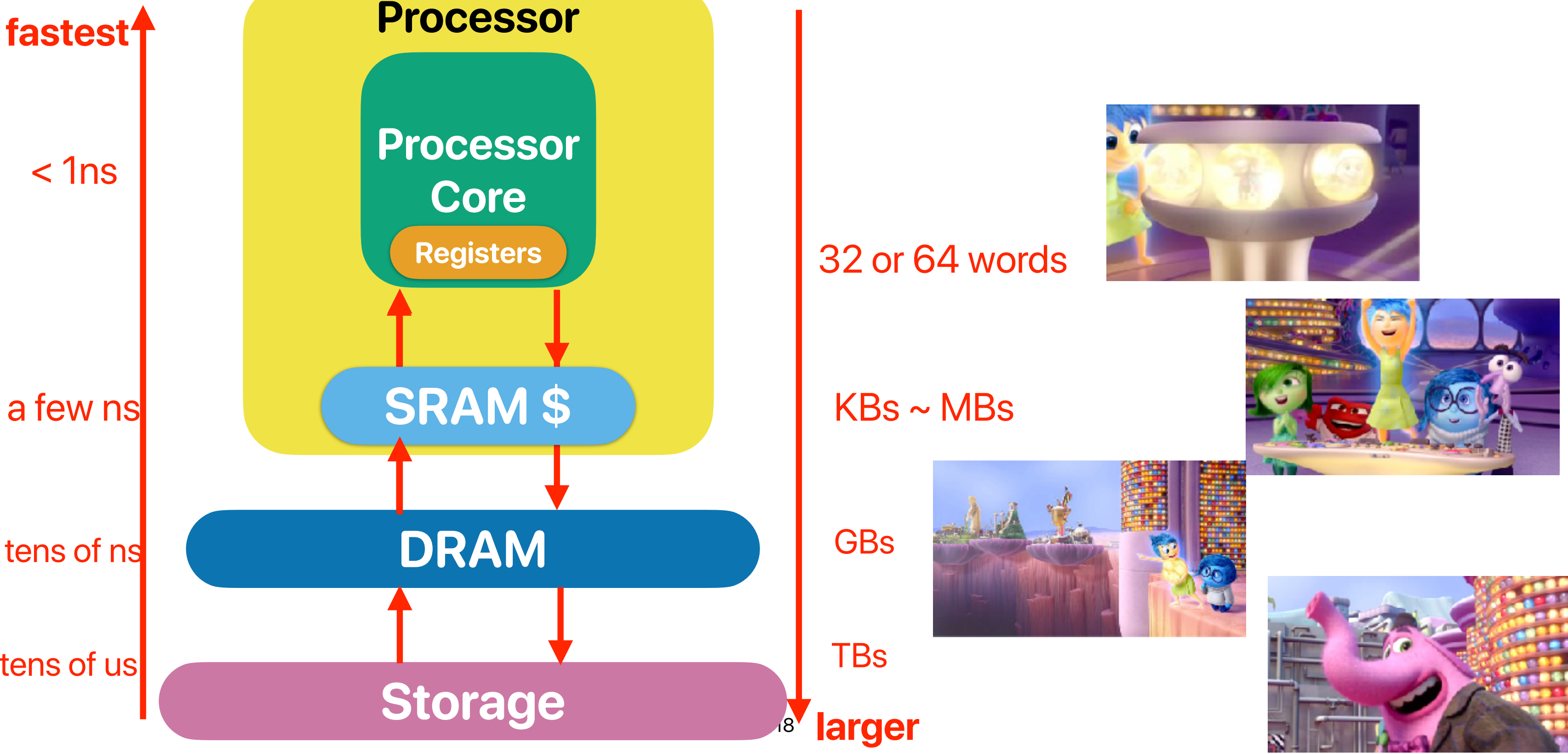
| Memory technology          | Typical access time     | \$ per GiB in 2012 |
|----------------------------|-------------------------|--------------------|
| SRAM semiconductor memory  | 0.5–2.5 ns              | \$500–\$1000       |
| DRAM semiconductor memory  | 50–70 ns                | \$10–\$20          |
| 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 \$\$\$**



# Memory Hierarchy




# L1? L2? L3?

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    | Socket AM5 (LGA1718) |              |         |
| Technology | 5 nm                 | Core Voltage | 1.288 V |



Specification

|                                    |    |            |        |
|------------------------------------|----|------------|--------|
| AMD Ryzen 7 7700X 8-Core Processor |    |            |        |
| Family                             | F  | Model      | 1      |
| Ext. Family                        | 19 | Ext. Model | 61     |
| Stepping                           | 2  | Revision   | RPL-B2 |

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    | Socket 1700 LGA      |              |         |
| Technology | 10 nm                | Core Voltage | 1.412 V |



Specification

|                             |   |            |    |
|-----------------------------|---|------------|----|
| Intel(R) Core(TM) i7-14700K |   |            |    |
| Family                      | 6 | Model      | 7  |
| Ext. Family                 | 6 | Ext. Model | B7 |
| Stepping                    | 1 | 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

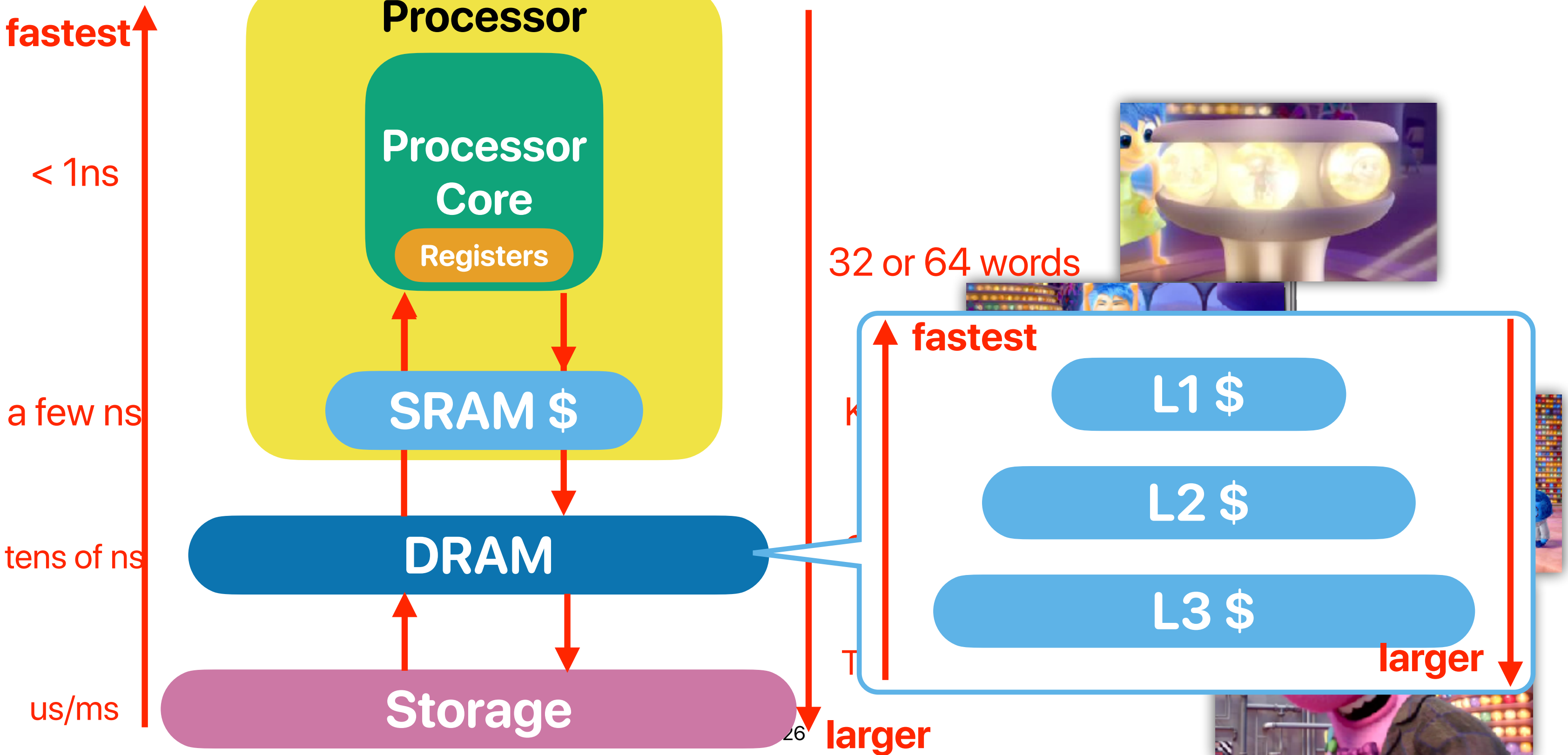
|          |                        |
|----------|------------------------|
| L1 Data  | 8 x 48 KB + 12 x 32 KB |
| L1 Inst. | 8 x 32 KB + 12 x 64 KB |
| Level 2  | 8 x 2 MB + 3 x 4 MB    |
| Level 3  | 33 MBytes              |

Selection Socket #1 Cores 8 + 12 Threads 28

CPU-Z Ver. 2.08.0.x64 Tools Validate Close



# Memory Hierarchy



# 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  | 1.25 V              |
| Specification   |                     |
| AMD Ryzen 7 7700X 8-Core Processor  |                     |
| Ext. Family   | 19                  |
| Ext. Model  | 1                   |
| Revision  | RPL12               |
| 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  | 1.12 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  |                        |

# Take-aways: inside out our memory hierarchy

- Memory access time is the most critical performance problem
  - One memory operation is as expensive as 50 arithmetic operations
  - Processor has to fetch instructions from memory
  - We have an average of 33% of data memory access instructions!
- Hierarchical caching with small amount of SRAMs will work if we can efficiently capture data and instructions

# **The predictability of your code**

# Code also has locality

keep going to the  
next instruction —  
**spatial locality**

```
for(uint64_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint64_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)



# Take-aways: inside out our memory hierarchy

- Memory access time is the most critical performance problem
  - One memory operation is as expensive as 50 arithmetic operations
  - Processor has to fetch instructions from memory
  - We have an average of 33% of data memory access instructions!
- Hierarchical caching with small amount of SRAMs will work if we can efficiently capture data and instructions
- Caching is possible! Most of time, we only work on a small amount of data!

# Designing a hardware to exploit locality

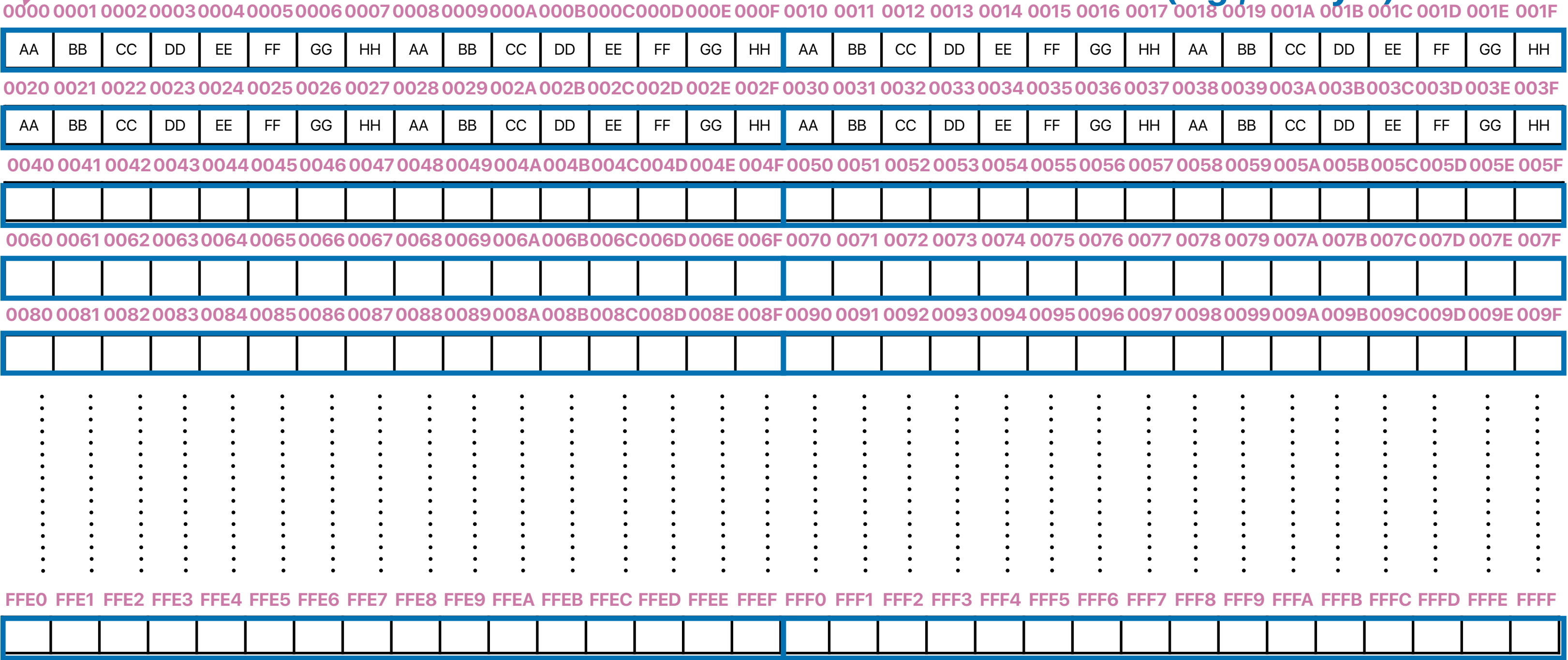
- Spatial locality — application tends to visit nearby stuffs in the memory

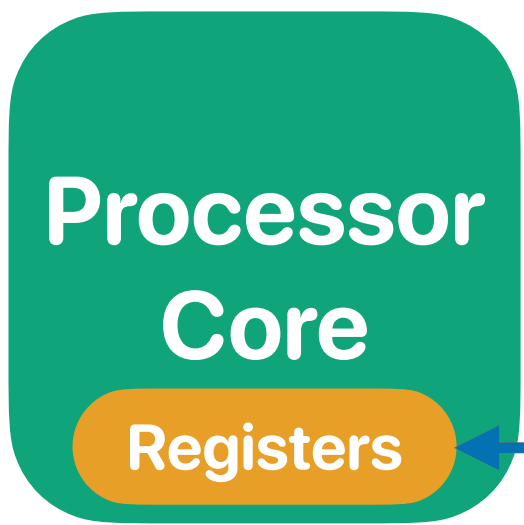
**We need to "cache consecutive memory locations" every time — the cache should store a "block" of code/data**

- 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)

# Block and the memory space

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



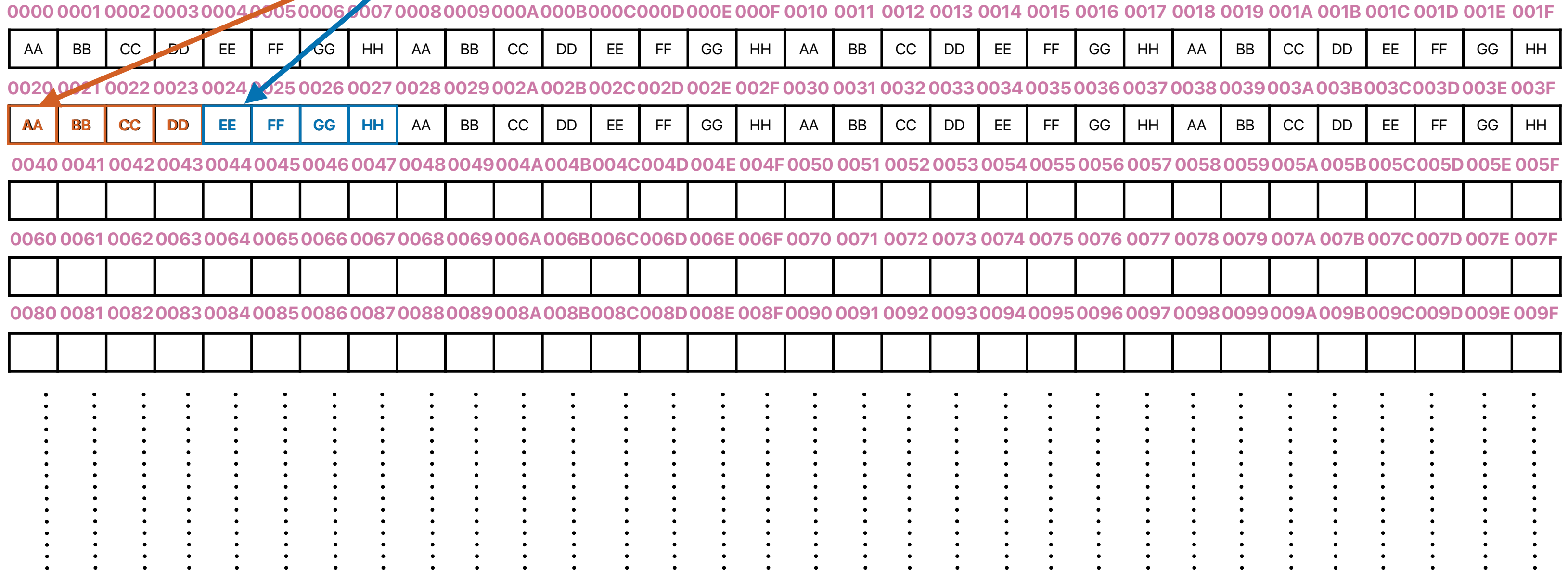


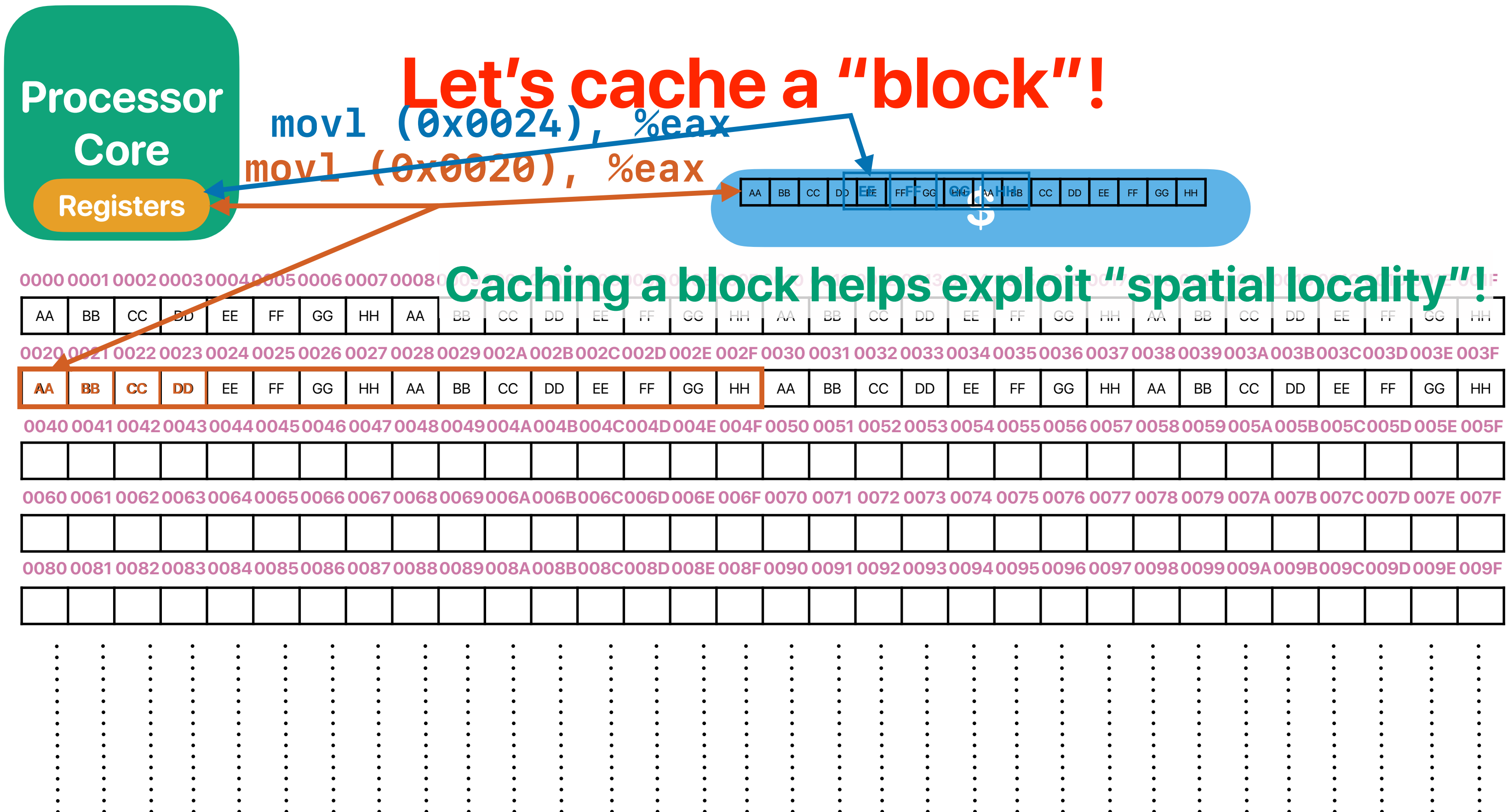
# When there is a "movl"

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(uint64_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint64_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

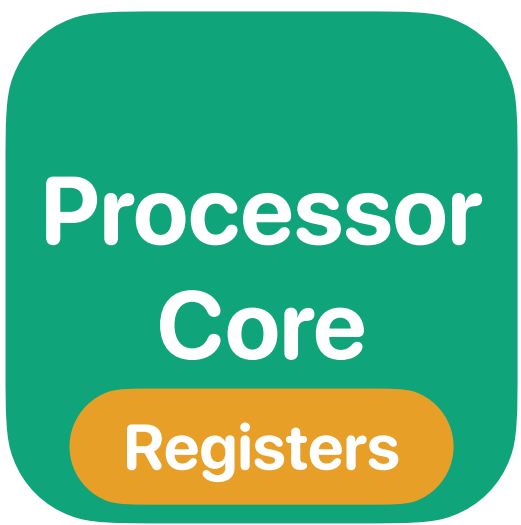
**We need to “cache consecutive memory locations” every time — the cache should store a “block” of code/data**

- 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)



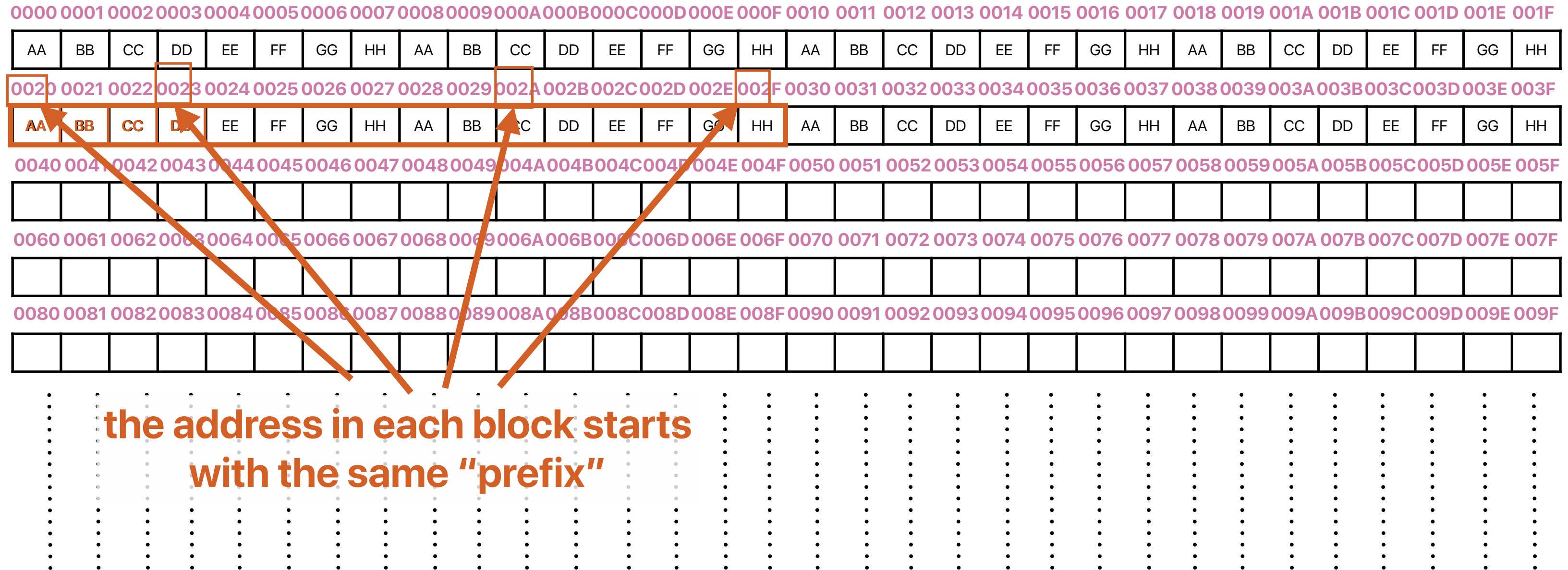
# How to tell who is there?

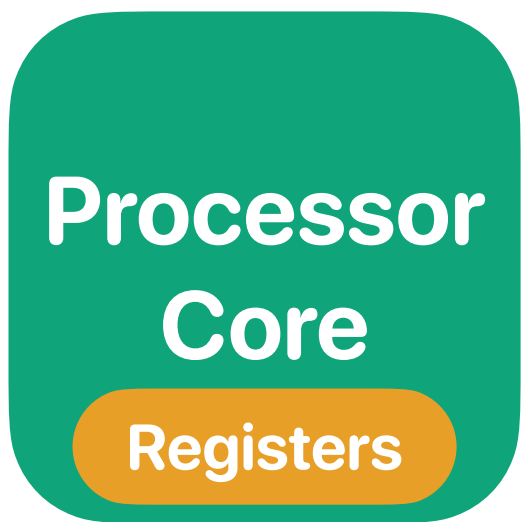
|   |                                 |
|---|---------------------------------|
| ? | 0123456789ABCDEF                |
|   | This is CS 203:                 |
|   | Advanced Computer Architecture! |
|   | This is CS 203:                 |
|   | Advanced Computer Architecture! |
|   | This is CS 203:                 |
|   | Advanced Computer Architecture! |
|   | This is CS 203:                 |
|   | Advanced Computer Architecture! |
|   | This is CS 203:                 |
|   | Advanced Computer Architecture! |
|   | This is CS 203:                 |
|   | Advanced Computer Architecture! |

# Registers

# Let's cache a "block"!

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

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



# How to tell who is there?

00000000000000000000000000000000  
0x00000000000000000000000000000000  
0123456789ABCDEF

the common address  
prefix in each block



| tag array |                  |
|-----------|------------------|
| 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:  |



# 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

|  |   |   |       |                  |
|--|---|---|-------|------------------|
|  | 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:  |

# 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

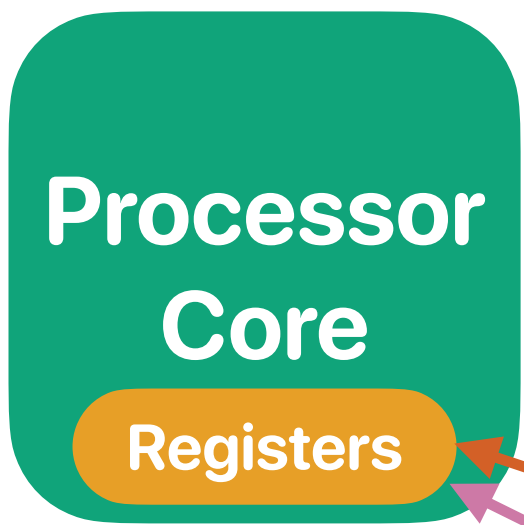
Dirty Bit

tag

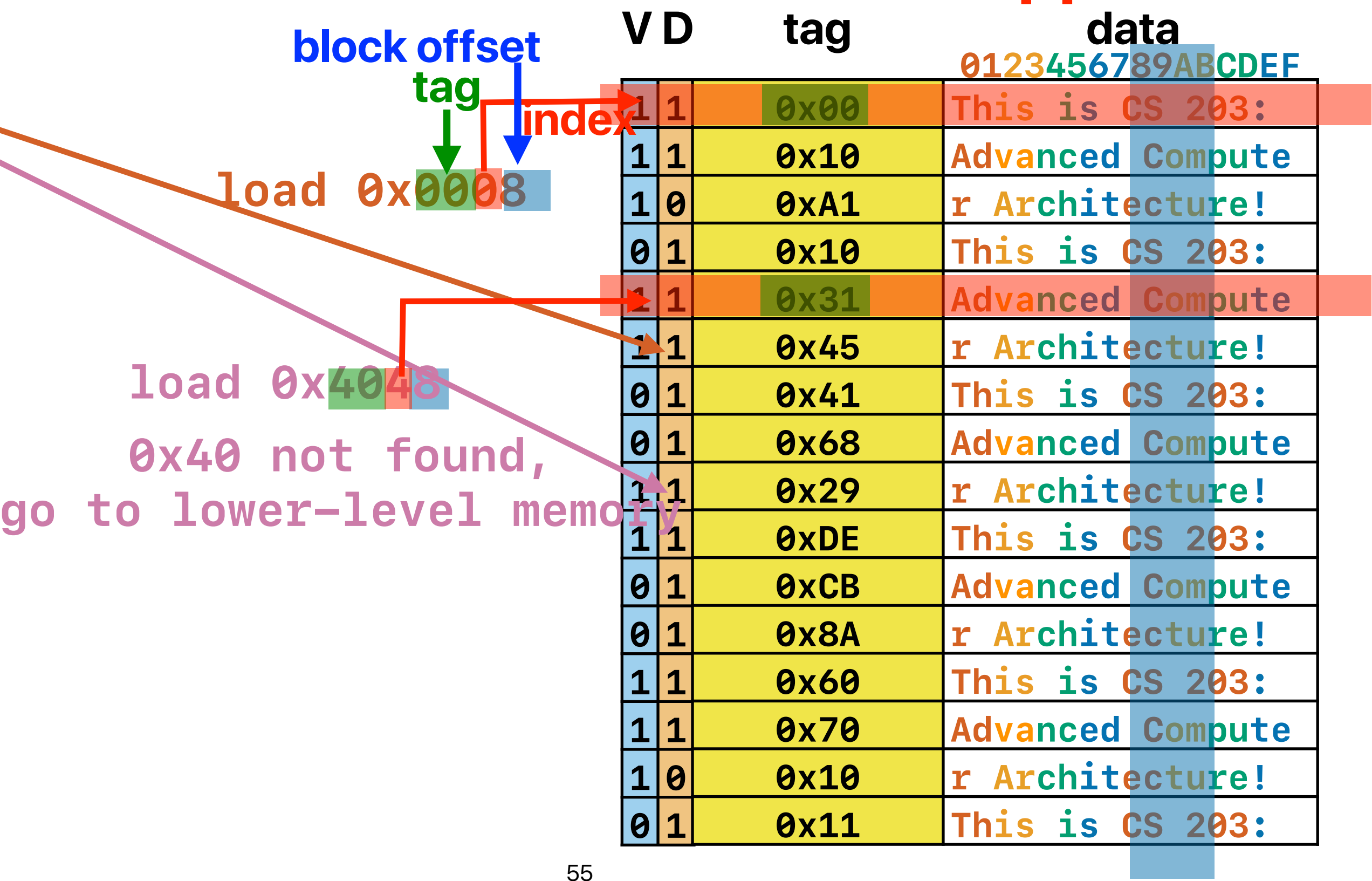
data

|   |   |  |       | 0123456789ABCDEF |                  |  |  |  |               |  |  |  |               |  |  |  |  |  |  |
|---|---|--|-------|------------------|------------------|--|--|--|---------------|--|--|--|---------------|--|--|--|--|--|--|
| 1 | 1 |  | 0x000 |                  | This is          |  |  |  | CSE13:        |  |  |  |               |  |  |  |  |  |  |
| 1 | 1 |  | 0x001 |                  | Advanced Compute |  |  |  | r             |  |  |  | Architecture! |  |  |  |  |  |  |
| 1 | 0 |  | 0xF07 |                  | r                |  |  |  | Architecture! |  |  |  |               |  |  |  |  |  |  |
| 0 | 1 |  | 0x100 |                  | This is          |  |  |  | CS 203:       |  |  |  |               |  |  |  |  |  |  |
| 1 | 1 |  | 0x310 |                  | Advanced Compute |  |  |  | r             |  |  |  | Architecture! |  |  |  |  |  |  |
| 1 | 1 |  | 0x450 |                  | r                |  |  |  | Architecture! |  |  |  |               |  |  |  |  |  |  |
| 0 | 1 |  | 0x006 |                  | This is          |  |  |  | CS 203:       |  |  |  |               |  |  |  |  |  |  |
| 0 | 1 |  | 0x537 |                  | Advanced Compute |  |  |  | r             |  |  |  | Architecture! |  |  |  |  |  |  |
| 1 | 1 |  | 0x266 |                  | r                |  |  |  | Architecture! |  |  |  |               |  |  |  |  |  |  |
| 1 | 1 |  | 0x307 |                  | This is          |  |  |  | CS 203:       |  |  |  |               |  |  |  |  |  |  |
| 0 | 1 |  | 0x265 |                  | Advanced Compute |  |  |  | r             |  |  |  | Architecture! |  |  |  |  |  |  |
| 0 | 1 |  | 0x80A |                  | r                |  |  |  | Architecture! |  |  |  |               |  |  |  |  |  |  |
| 1 | 1 |  | 0x620 |                  | This is          |  |  |  | CS 203:       |  |  |  |               |  |  |  |  |  |  |
| 1 | 1 |  | 0x630 |                  | Advanced Compute |  |  |  | r             |  |  |  | Architecture! |  |  |  |  |  |  |
| 1 | 0 |  | 0x705 |                  | r                |  |  |  | Architecture! |  |  |  |               |  |  |  |  |  |  |
| 0 | 1 |  | 0x216 |                  | This is          |  |  |  | CS 203:       |  |  |  |               |  |  |  |  |  |  |





# Hash-like structure — direct-mapped cache





# Blocksize == Linesize

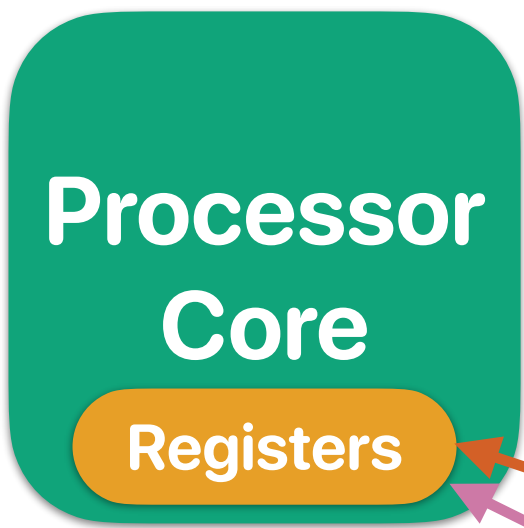
```
[5]: # Your CS203 Cluster
! cs203 demo "lscpu | grep 'Model name'; getconf -a | grep CACHE"

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
```

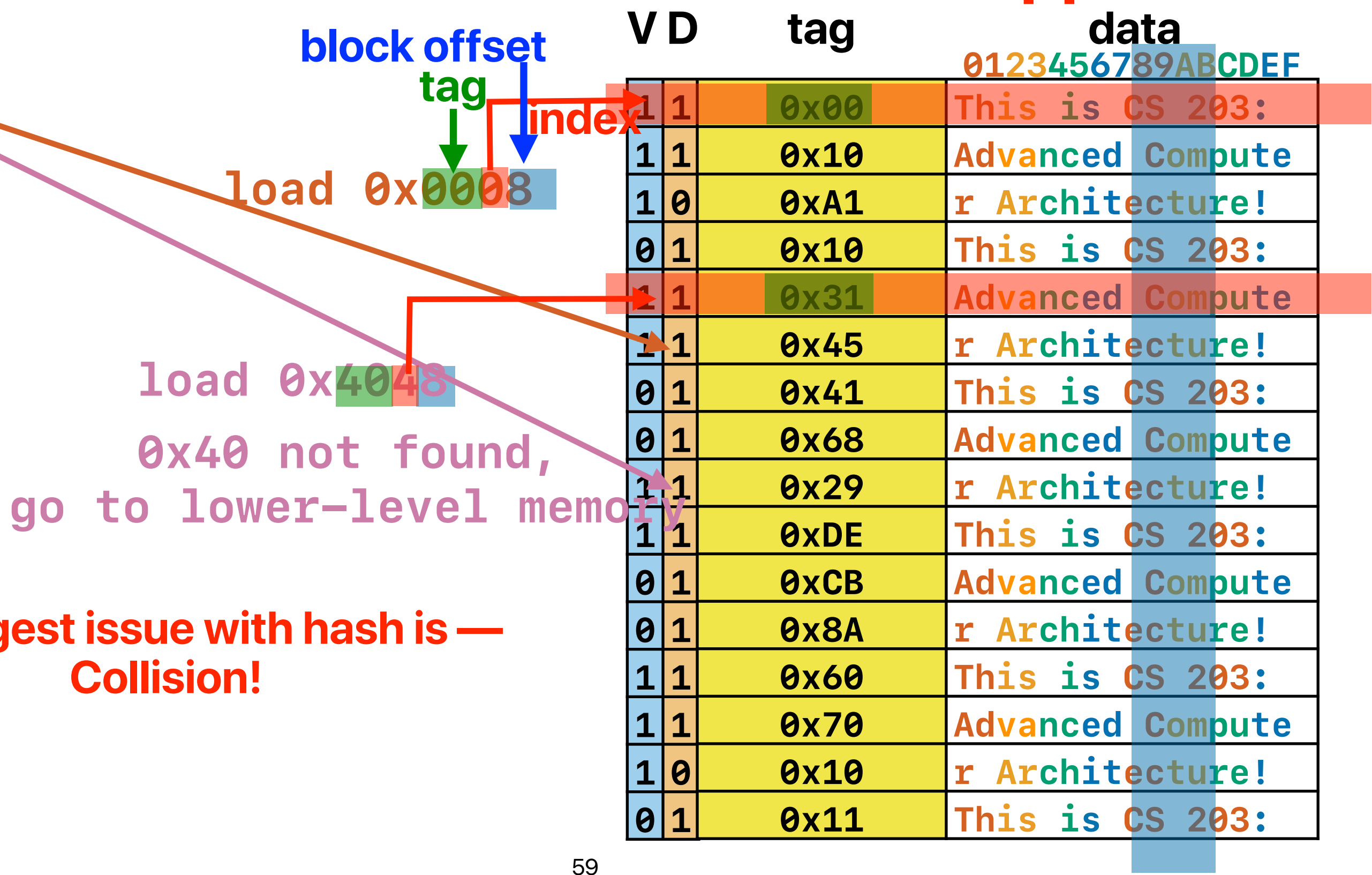
# What is Associativity?

```
[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
```



# Hash-like structure — direct-mapped cache



The biggest issue with hash is — Collision!

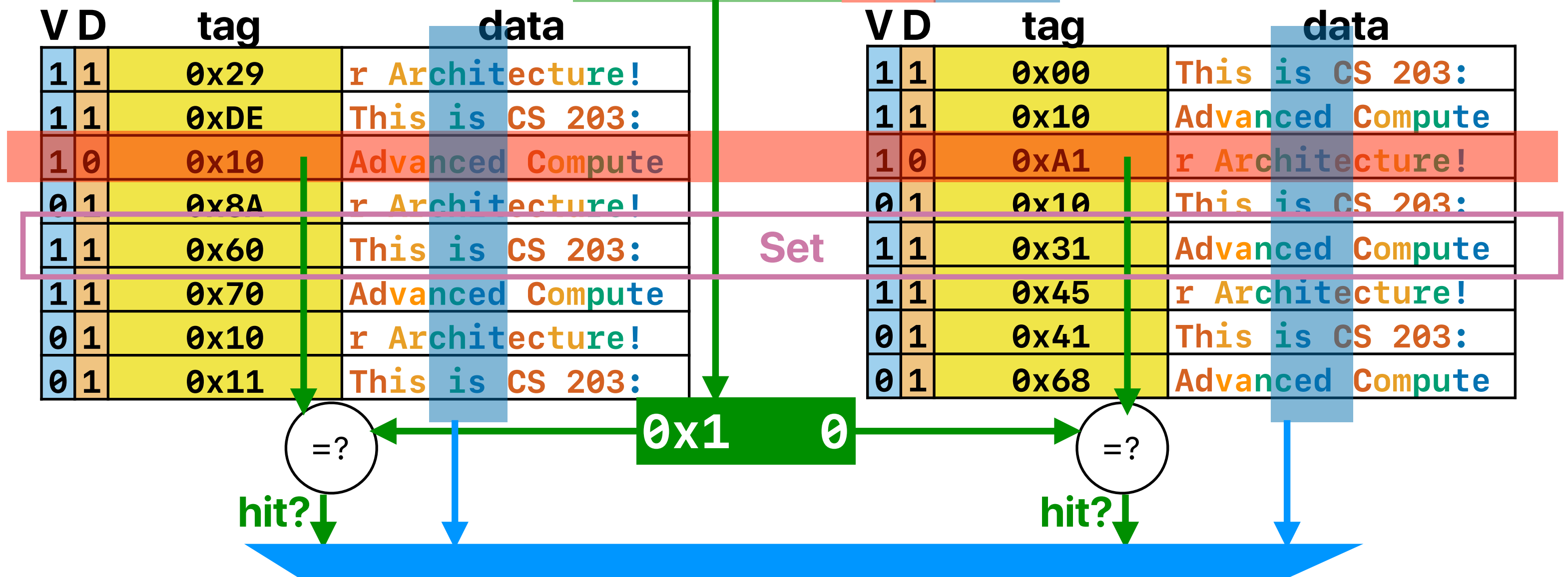
# Way-associative cache

memory address:      0x0      8      2      4

set      block

tag      index      offset

memory address:      0b00001000000100100





# Take-aways: designing caches

- Cache structures — blocks and sets
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Set associative to avoid “collisions”

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

# Processor/cache interaction



- Processor sends memory access request to L1-\$

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

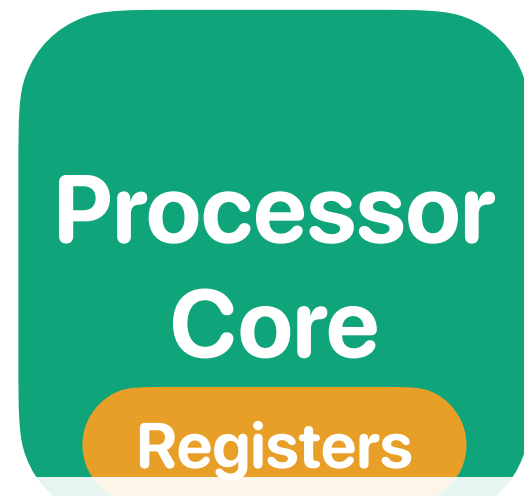
- **if miss**

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

**What if we run out of \$ blocks?**



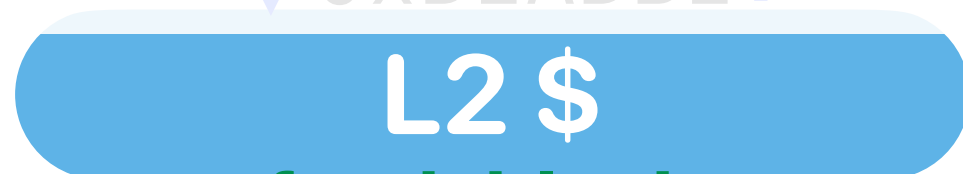
# Considering we have limited space in \$



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

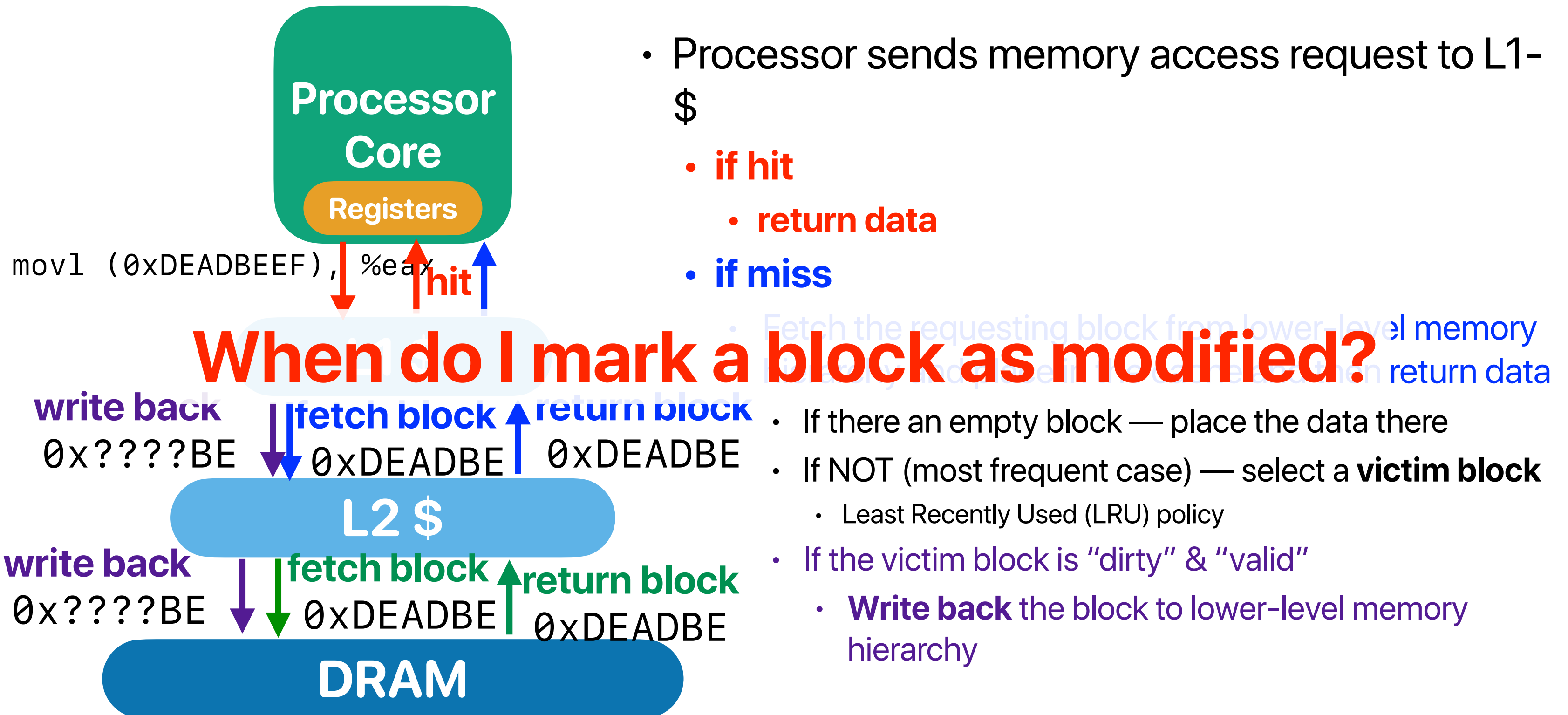
**What if the victim block is modified?**

**— ignoring the update is not acceptable!**

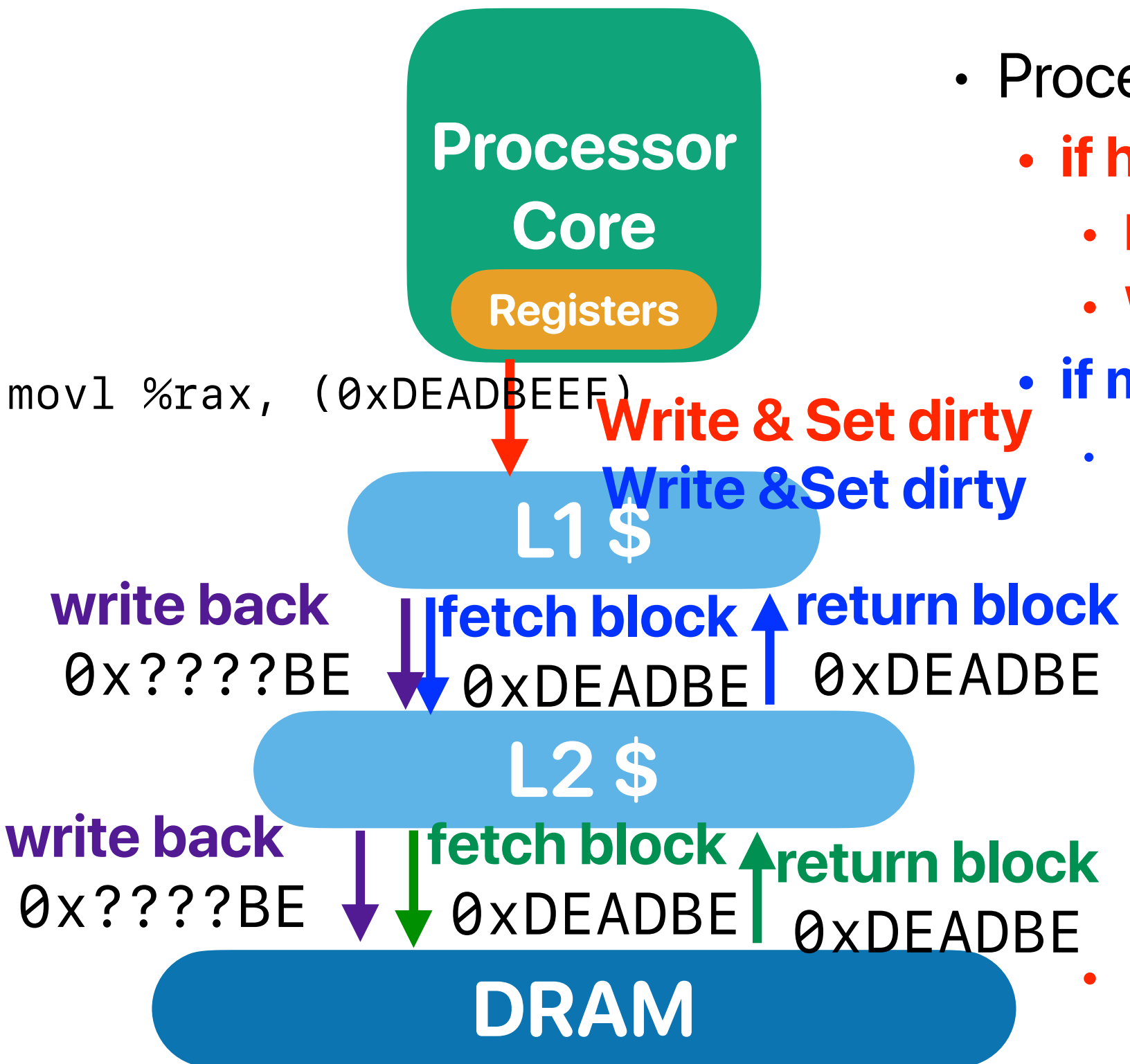


- If there an empty block — place the data there
- If NOT (most frequent case) — select a **victim block**
  - Least Recently Used (LRU) policy

# 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**

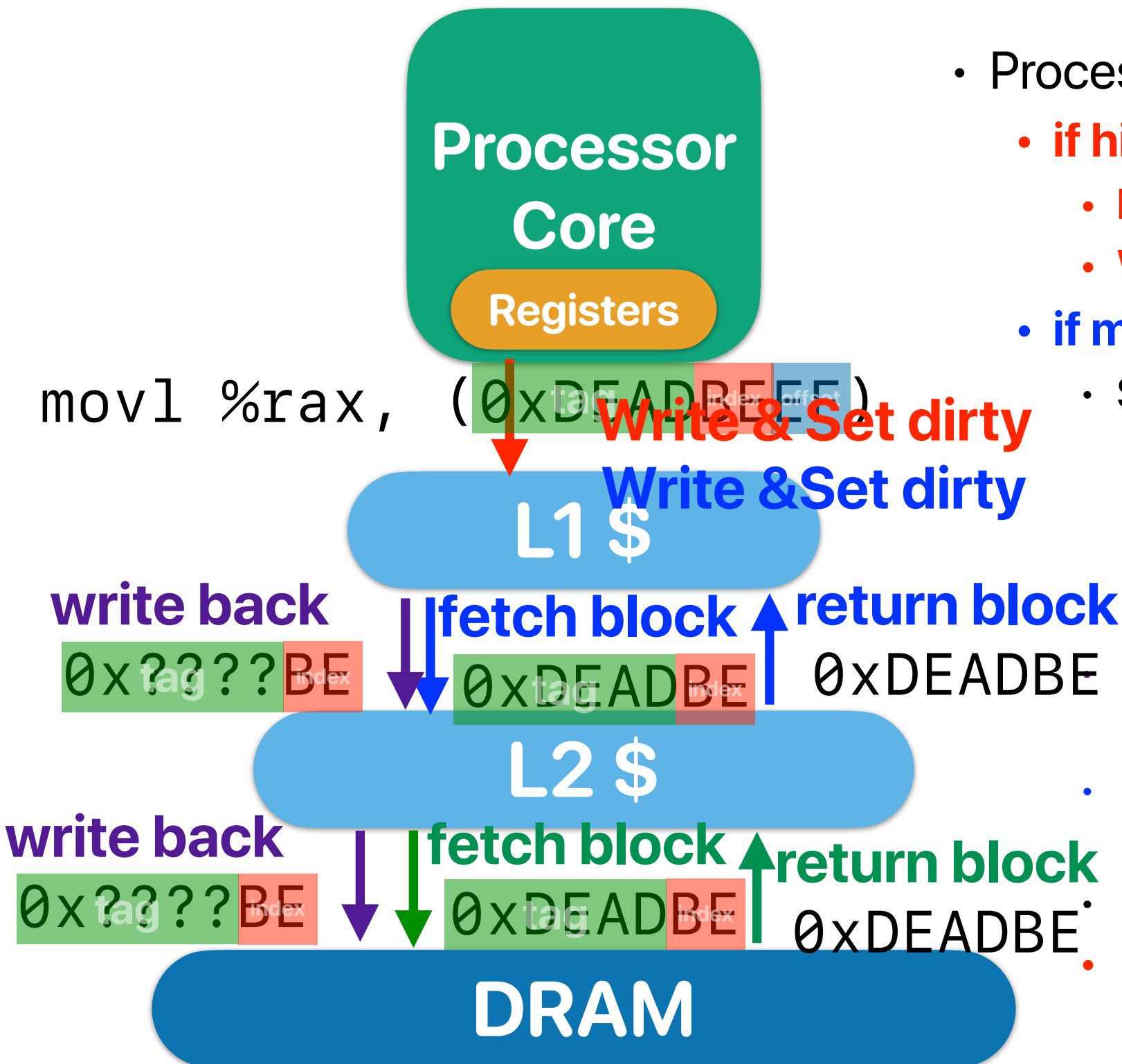
- Least Recently Used (LRU) policy

- If the victim block is "dirty" & "valid"

- **Write back** the block to lower-level memory hierarchy

- **Write: Update "ONLY" in L1 and set DIRTY**

# 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**

# Way-associative cache

memory address:  $0x0$  8 2 4

set block

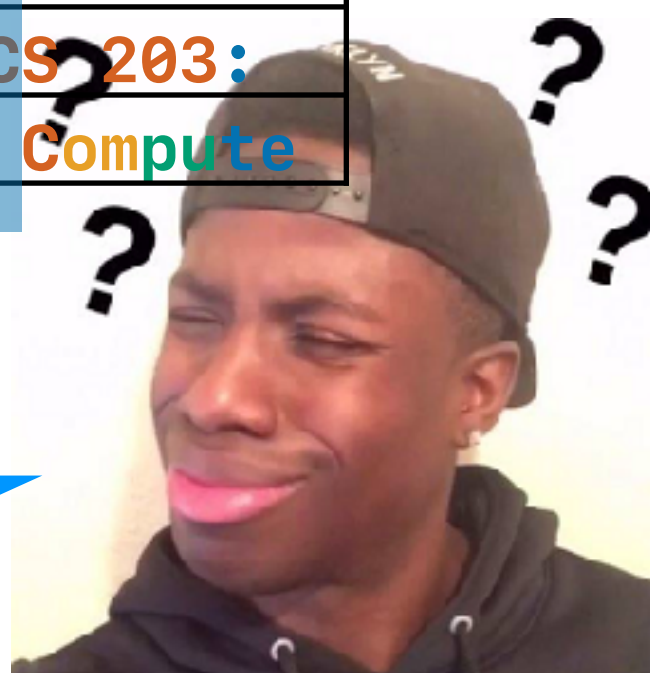
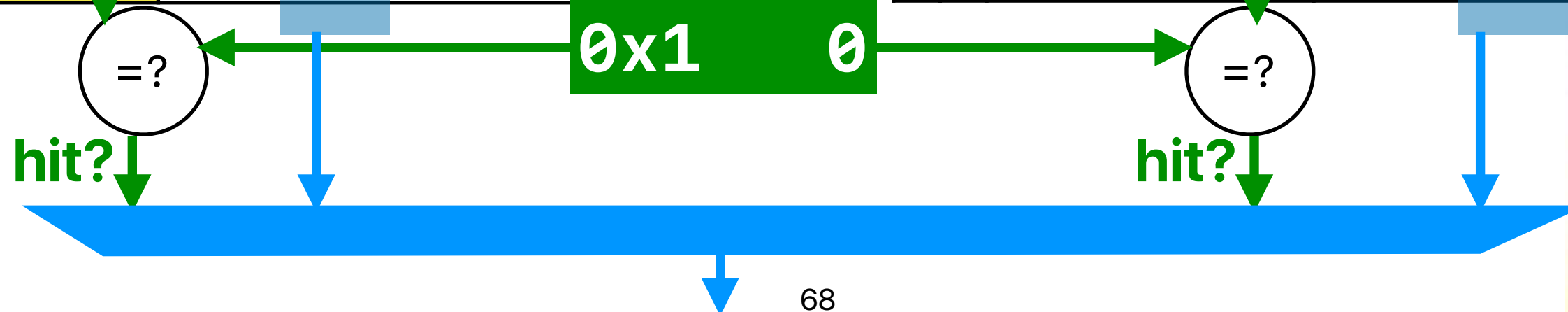
tag index offset

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





# Take-aways: designing caches

- Cache structures — blocks and sets
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Set associative to avoid “collisions”
- Data must present in cache before processor can use

# 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 

- 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)
- $(\text{address} / \text{block\_size}) \% \mathbf{S} = \text{set index}$

# Take-aways: designing caches

- Cache structures — blocks and sets
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Set associative to avoid “collisions”
- Data must present in cache before processor can use
- $C = A B S$ 
  - C: capacity
  - A: Associativity
  - S: Number of sets
  - $\lg(S)$ : Number of bits in set index
  - $\lg(B)$ : Number of bits in block offset

# **Simulate the 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(uint64_t i = 0; i < m; i++) {  
    result = 0;  
    for(uint64_t j = 0; j < n; j++) {  
        result += matrix[i][j]*vector[j];  
    }  
    output[i] = result;  
}
```



# Matrix vector revisited

|                                   | tag index     |                | tag index   |  |
|-----------------------------------|---------------|----------------|---|--|
|                                   | tag           |                | index   |  |
|                                   | Address (Hex) |                | Address (Binary)                                    |  |
| for(uint64_t i = 0; i < m; i++) { | &a[0][0]      | 0x558FE0A1D330 | 0b10101011000111111100000101000011101001100110000   |  |
| result = 0;                       | &b[0]         | 0x558FE0A1DC30 | 0b10101011000111111100000101000011101110000110000   |  |
| for(uint64_t j = 0; j < n; j++) { | &a[0][1]      | 0x558FE0A1D338 | 0b10101011000111111100000101000011101001100111000   |  |
| result += matrix[i][j]*vector[j]; | &b[1]         | 0x558FE0A1DC38 | 0b101010110001111111000001010000111011100000111000  |  |
| }                                 | &a[0][2]      | 0x558FE0A1D340 | 0b10101011000111111100000101000011101001101000000   |  |
| output[i] = result;               | &b[2]         | 0x558FE0A1DC40 | 0b1010101100011111110000010100001110111000010000000 |  |
| }                                 | &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 |      |

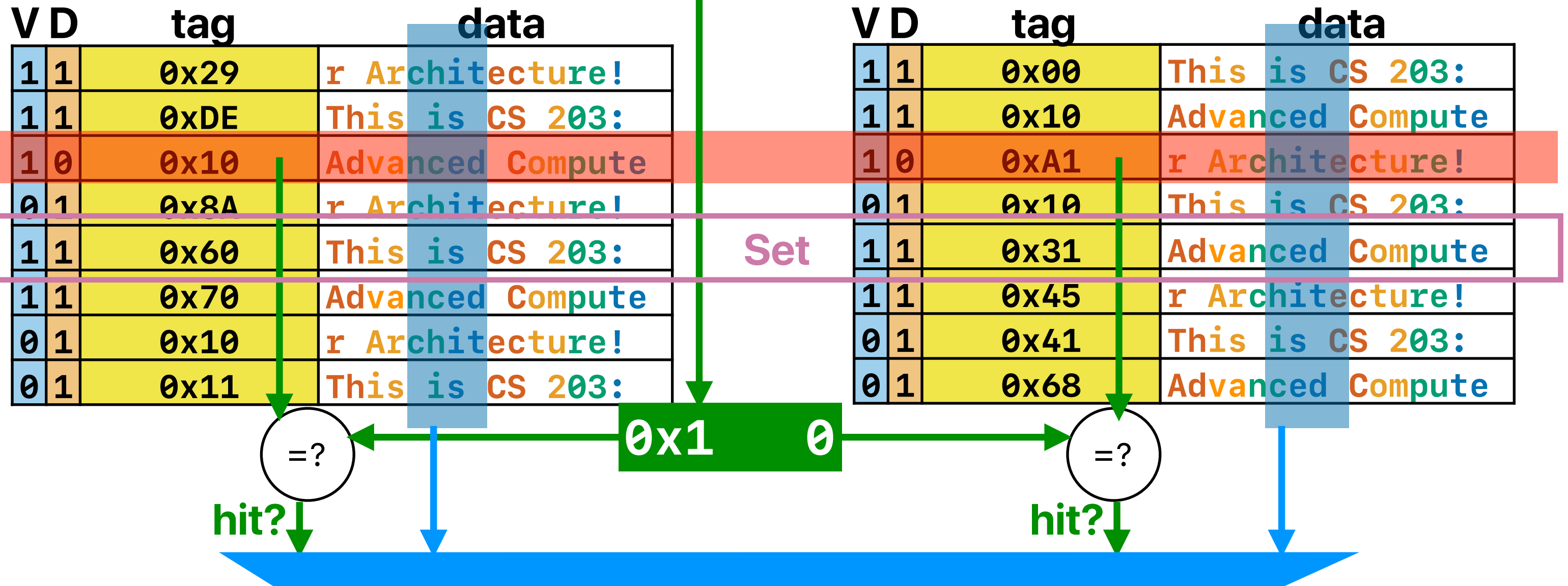
# Way-associative cache

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:  $0x$       8      0      0      0      0      0      8      0

= 0b1000 0000 0000 0000 0000 0000 0000 1000 0000

The diagram shows a 64-bit address split into three fields: a 56-bit tag (red), a 3-bit index (blue), and a 4-bit offset (black). The tag field is the largest, followed by the index, and then the offset.

tag      index      offset

# Matrix vector revisited

|                                   | tag index     |                | tag index   |  |
|-----------------------------------|---------------|----------------|---|--|
|                                   | tag           |                | index   |  |
|                                   | Address (Hex) |                | Address (Binary)                                  |  |
| for(uint64_t i = 0; i < m; i++) { | &a[0][0]      | 0x558FE0A1D330 | 0b10101011000111111100000101000011101001100110000 |  |
| result = 0;                       | &b[0]         | 0x558FE0A1DC30 | 0b10101011000111111100000101000011101110000110000 |  |
| for(uint64_t j = 0; j < n; j++) { | &a[0][1]      | 0x558FE0A1D338 | 0b10101011000111111100000101000011101001100111000 |  |
| result += matrix[i][j]*vector[j]; | &b[1]         | 0x558FE0A1DC38 | 0b10101011000111111100000101000011101110000111000 |  |
| }                                 | &a[0][2]      | 0x558FE0A1D340 | 0b10101011000111111100000101000011101001101000000 |  |
| output[i] = result;               | &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  |



# **Taxonomy/reasons of cache misses**

# 3Cs of misses

- Compulsory miss
  - Cold start miss. First-time access to a block
- Capacity miss
  - The working set size of an application is bigger than cache size
- Conflict miss
  - Required data replaced by block(s) mapping to the same set
  - Similar collision in hash

**How can programmer improve  
memory performance?**

# Data structures



# Column-store or row-store

- Considering your the most frequently used queries in your database system are similar to

`SELECT AVG(assignment_1) FROM table`

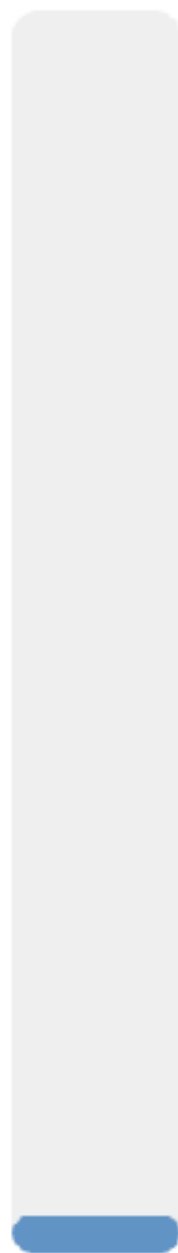
Which of the following would be a data structure that better implements the table supporting this type of queries?

| Array of objects   | object of arrays  |
|--|---|
| <pre>struct grades {     int id;     double *homework;     double average; }; table = (struct grades *) \ malloc(num_of_students*sizeof(struct</pre> | <pre>struct grades {     int *id;     double **homework;     double *average; }; table =(struct grades *)malloc(sizeof(struct grades));</pre> |

- A. Array of objects
- B. Object of arrays

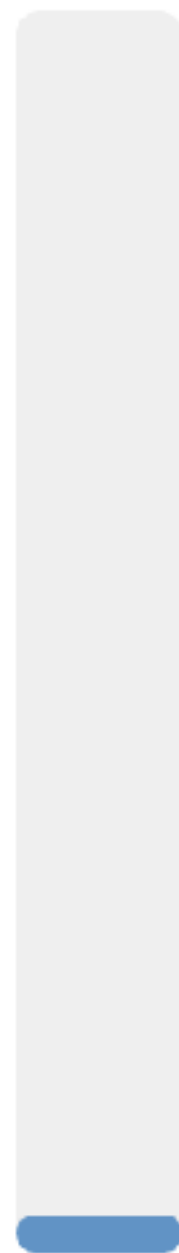


0



A

0



B

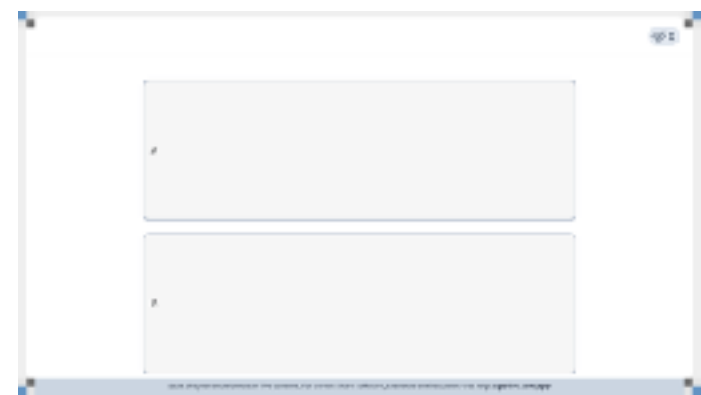


# Column-store or row-store

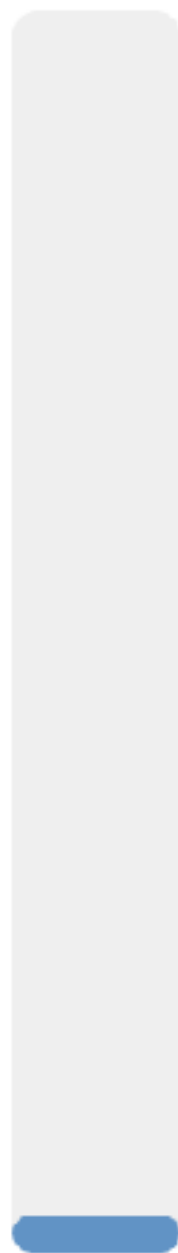
- Considering your the most frequently used queries in your database system are similar to  
`SELECT AVG(assignment_1) FROM table`  
Which of the following would be a data structure that better implements the table supporting this type of queries?

| Array of objects   | object of arrays  |
|--|---|
| <pre>struct grades {     int id;     double *homework;     double average; }; table = (struct grades *) \ malloc(num_of_students*sizeof(struct</pre> | <pre>struct grades {     int *id;     double **homework;     double *average; }; table =(struct grades *)malloc(sizeof(struct grades));</pre> |

- A. Array of objects
- B. Object of arrays

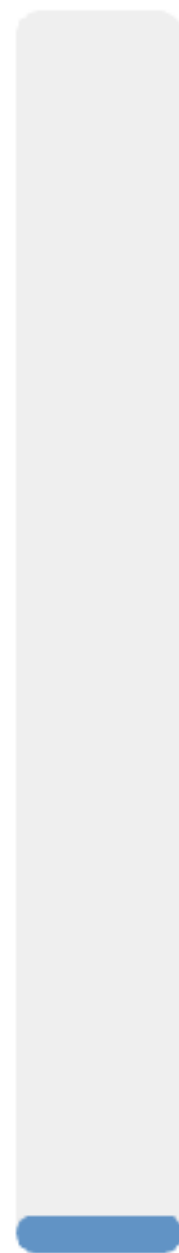


0



A

0



B



# Column-store or row-store

- Considering your the most frequently used queries in your database system are similar to

`SELECT AVG(assignment_1) FROM table`

Which of the following would be a data structure that better implements the table supporting this type of queries?

| Array of objects   | object of arrays  |
|--|---|
| <pre>struct grades {     int id;     double *homework;     double average; }; table = (struct grades *) \ malloc(num_of_students*sizeof(struct</pre> | <pre>struct grades {     int *id;     double **homework;     double *average; }; table =(struct grades *)malloc(sizeof(struct grades));</pre> |

A. Array of objects **What if we want to calculate average scores for each student?**

B. Object of arrays



# Array of structures or structure of arrays

| Array of objects         |  |           | object of arrays |  |           |         |    |    |    |          |          |          |         |         |         |
|--------------------------|--|-----------|------------------|--|-----------|---------|----|----|----|----------|----------|----------|---------|---------|---------|
|                          | <pre>struct grades {     int id;     double *homework;     double average; };</pre>  |           |                  | <pre>struct grades {     int *id;     double **homework;     double *average; };</pre>   |           |         |    |    |    |          |          |          |         |         |         |
|                          | ID   | *homework | average          | ID   | *homework | average |    |    |    |          |          |          |         |         |         |
| average of each homework | <pre>for(i=0;i&lt;homework_items; i++) {     gradesheet[total_number_students].homework[i] = 0.0;     for(j=0;j&lt;total_number_students;j++)     gradesheet[total_number_students].homework[i]     +=gradesheet[j].homework[i];     gradesheet[total_number_students].homework[i] /=     (double)total_number_students; }</pre> |           |                  | <table><tr><td>ID</td><td>ID</td><td>ID</td></tr><tr><td>homework</td><td>homework</td><td>homework</td></tr><tr><td>average</td><td>average</td><td>average</td></tr></table> <pre>for(i = 0;i &lt; homework_items; i++) {     gradesheet.homework[i][total_number_students] = 0.0;     for(j = 0; j &lt;total_number_students;j++)     {         gradesheet.homework[i][total_number_students] +=         gradesheet.homework[i][j];     }     gradesheet.homework[i][total_number_students] /=     total_number_students; }</pre> |           |         | ID | ID | ID | homework | homework | homework | average | average | average |
|                          | ID   | ID        | ID               |  |           |         |    |    |    |          |          |          |         |         |         |
| homework                 | homework   | homework  |                  |  |           |         |    |    |    |          |          |          |         |         |         |
| average                  | average  | average   |                  |  |           |         |    |    |    |          |          |          |         |         |         |

# Column-store or row-store

- If you're designing an in-memory database system, will you be using

| RowId | Empld | Lastname | Firstname | Salary |
|-------|-------|----------|-----------|--------|
| 1     | 10    | Smith    | Joe       | 40000  |
| 2     | 12    | Jones    | Mary      | 50000  |
| 3     | 11    | Johnson  | Cathy     | 44000  |
| 4     | 22    | Jones    | Bob       | 55000  |

- column-store — stores data tables column by column

10:001, 12:002, 11:003, 22:004;

Smith:001, Jones:002, Johnson:003, Jones:004;

Joe:001, Mary:002, Cathy:003, Bob:004;

40000:001, 50000:002, 44000:003, 55000:004;

if the most frequently used query looks like —  
**select Lastname, Firstname from table**

- row-store — stores data tables row by row

001:10, Smith, Joe, 40000;

002:12, Jones, Mary, 50000;

003:11, Johnson, Cathy, 44000;

004:22, Jones, Bob, 55000;

# Takeaways: Software Optimizations

- Data layout — capacity miss, conflict miss, compulsory miss

# Loop interchange/fission/fusion

# Demo — programmer & performance

A

```
for(i = 0; i < ARRAY_SIZE; i++)
{
    for(j = 0; j < ARRAY_SIZE; j++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

B

```
for(j = 0; j < ARRAY_SIZE; j++)
{
    for(i = 0; i < ARRAY_SIZE; i++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

$O(n^2)$

Complexity

$O(n^2)$

Same

Instruction Count?

Same

Same

Clock Rate

Same

Better

CPI

Worse

# Loop optimizations

## Loop interchange

A

```
for(i = 0; i < ARRAY_SIZE; i++)  
{  
    for(j = 0; j < ARRAY_SIZE; j++)  
    {  
        c[i][j] = a[i][j]+b[i][j];  
    }  
}
```



B

```
for(j = 0; j < ARRAY_SIZE; j++)  
{  
    for(i = 0; i < ARRAY_SIZE; i++)  
    {  
        c[i][j] = a[i][j]+b[i][j];  
    }  
}
```



# Takeaways: Software Optimizations

- Data layout — capacity miss, conflict miss, compulsory miss
- Loop interchange — conflict/capacity miss

# Loop optimizations

## Loop interchange

A

```
for(i = 0; i < ARRAY_SIZE; i++)
{
    for(j = 0; j < ARRAY_SIZE; j++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

B

```
for(j = 0; j < ARRAY_SIZE; j++)
{
    for(i = 0; i < ARRAY_SIZE; i++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```



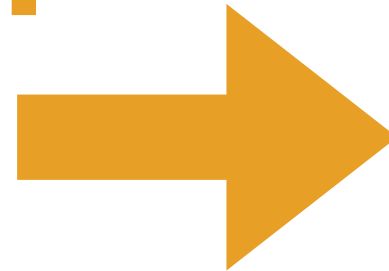
B

```
double a[8192], b[8192], c[8192], \
       d[8192], e[8192];
for(i = 0; i < 512; i++) {
    e[i] = (a[i] * b[i] + c[i])/d[i];
}
```

## Loop fission

A

```
double a[8192], b[8192], c[8192], \
       d[8192], e[8192];
for(i = 0; i < 512; i++)
    e[i] = a[i] * b[i] + c[i];
for(i = 0; i < 512; i++)
    e[i] /= d[i];
```



# Takeaways: Software Optimizations

- Data layout — capacity miss, conflict miss, compulsory miss
- Loop interchange — conflict/capacity miss
- Loop fission — conflict miss — when \$ has limited way associativity

# What if we change the processor?

- If we have an intel processor with a 32KB, 8-way, 64B-blocked L1 cache, which version of code performs better?
  - A. Version A, because the code incurs fewer cache misses
  - B. Version B, because the code incurs fewer cache misses
  - C. Version A, because the code incurs fewer memory references
  - D. Version B, because the code incurs fewer memory references**
  - E. They are about the same

A

```
double a[8192], b[8192], c[8192], \
      d[8192], e[8192];
for(i = 0; i < 512; i++)
    e[i] = a[i] * b[i] + c[i];
for(i = 0; i < 512; i++)
    e[i] /= d[i];
```

B

```
double a[8192], b[8192], c[8192], \
      d[8192], e[8192];
for(i = 0; i < 512; i++) {
    e[i] = (a[i] * b[i] + c[i])/d[i];
}
```

# Loop optimizations

## Loop interchange

A

```
for(i = 0; i < ARRAY_SIZE; i++)
{
    for(j = 0; j < ARRAY_SIZE; j++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

B

```
for(j = 0; j < ARRAY_SIZE; j++)
{
    for(i = 0; i < ARRAY_SIZE; i++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```



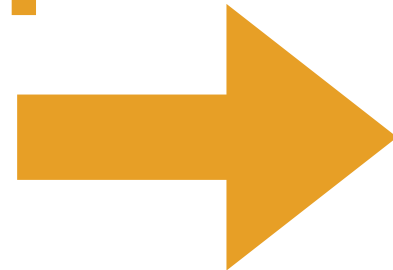
B

```
double a[8192], b[8192], c[8192], \
       d[8192], e[8192];
for(i = 0; i < 512; i++) {
    e[i] = (a[i] * b[i] + c[i])/d[i];
}
```

## Loop fission

A

```
double a[8192], b[8192], c[8192], \
       d[8192], e[8192];
for(i = 0; i < 512; i++)
    e[i] = a[i] * b[i] + c[i];
for(i = 0; i < 512; i++)
    e[i] /= d[i];
```



A

```
double a[8192], b[8192], c[8192], \
       d[8192], e[8192];
for(i = 0; i < 512; i++)
    e[i] = a[i] * b[i] + c[i];
for(i = 0; i < 512; i++)
    e[i] /= d[i];
```

## Loop fusion

B

```
double a[8192], b[8192], c[8192], \
       d[8192], e[8192];
for(i = 0; i < 512; i++) {
    e[i] = (a[i] * b[i] + c[i])/d[i];
}
```



# Takeaways: Software Optimizations

- Data layout — capacity miss, conflict miss, compulsory miss
- Loop interchange — conflict/capacity miss
- Loop fission — conflict miss — when \$ has limited way associativity
- Loop fusion — capacity miss — when \$ has enough way associativity

# Tiling

# Case study: Matrix Multiplication

```
for(i = 0; i < ARRAY_SIZE; i++) {  
    for(j = 0; j < ARRAY_SIZE; j++) {  
        for(k = 0; k < ARRAY_SIZE; k++) {  
            c[i][j] += a[i][k]*b[k][j];  
        }  
    }  
}
```

**Algorithm class tells you it's  $O(n^3)$**

**If  $n=1024$ , it takes about 1 sec**

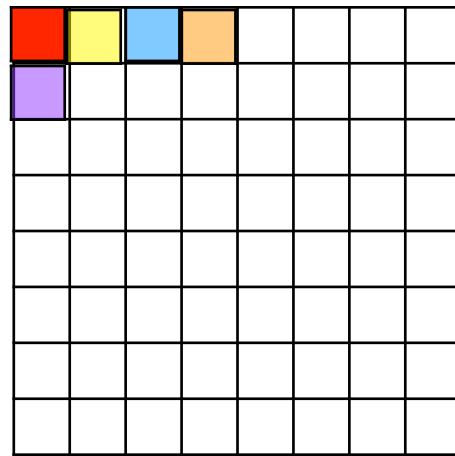
**How long is it take when  $n=2048$ ?**



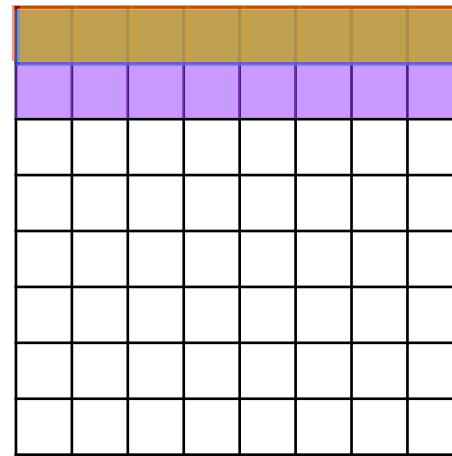
# Matrix Multiplication

```
for(i = 0; i < ARRAY_SIZE; i++) {  
    for(j = 0; j < ARRAY_SIZE; j++) {  
        for(k = 0; k < ARRAY_SIZE; k++) {  
            c[i][j] += a[i][k]*b[k][j];  
        }  
    }  
}
```

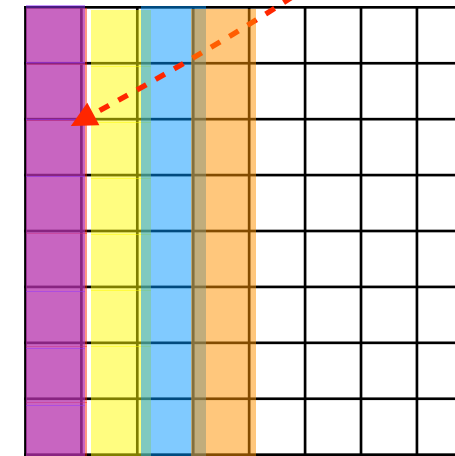
Very likely a miss if  
array is large



c



a



b

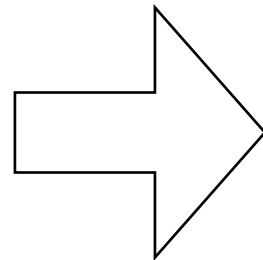
- If each dimension of your matrix is 2048
  - Each row takes  $2048 \times 8$  bytes = 16KB
  - The L1 \$ of intel Core i7 is 48 KB, 12-way, 64-byte blocked
  - You can only hold at most 3 rows/columns of each matrix!
  - You need the same row when j increase!

# Tiling algorithm for matrix multiplication

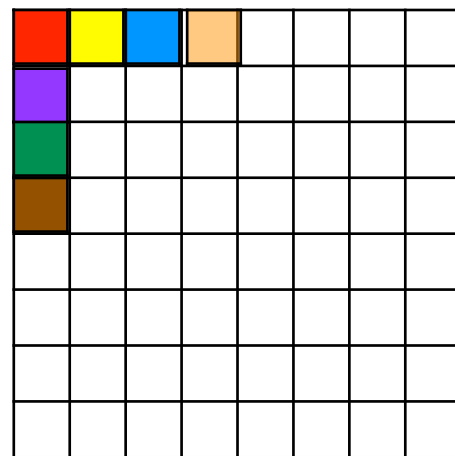
- Discover the cache miss rate
  - `valgrind --tool=cachegrind cmd`
    - cachegrind is a tool profiling the cache performance
- Performance counter
  - Intel® Performance Counter Monitor <http://www.intel.com/software/pcm/>

# Tiling algorithm for matrix multiplication

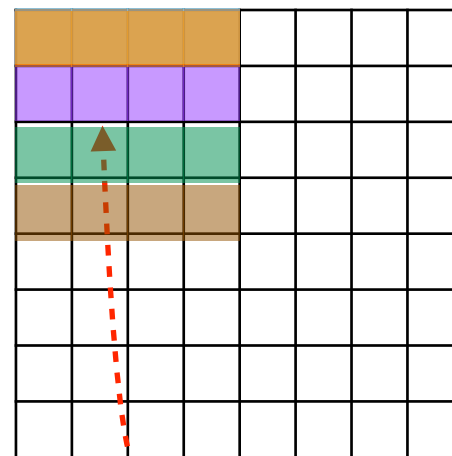
```
for(i = 0; i < ARRAY_SIZE; i++) {
  for(j = 0; j < ARRAY_SIZE; j++) {
    for(k = 0; k < ARRAY_SIZE; k++) {
      c[i][j] += a[i][k]*b[k][j];
    }
  }
}
```



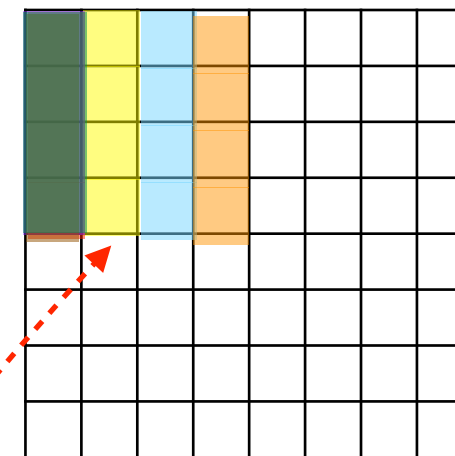
```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {
  for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {
    for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {
      for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)
        for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)
          for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)
            c[ii][jj] += a[ii][kk]*b[kk][jj];
    }
  }
}
```



c



a

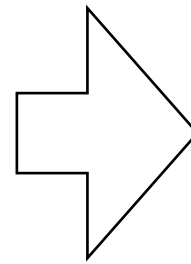


b

**You only need to hold these  
sub-matrices in your cache**

# Matrix Transpose

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```



```
// Transpose matrix b into b_t  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        b_t[i][j] += b[j][i];  
    }  
}  
  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        // Compute on b_t  
                        c[ii][jj] += a[ii][kk]*b_t[jj][kk];  
        }  
    }  
}
```

# Takeaways: Software Optimizations

- Data layout — capacity miss, conflict miss, compulsory miss
- Loop interchange — conflict/capacity miss
- Loop fission — conflict miss — when \$ has limited way associativity
- Loop fusion — capacity miss — when \$ has enough way associativity
- Blocking/tiling — capacity miss, conflict miss

# Computer Science & Engineering

142

つづく

