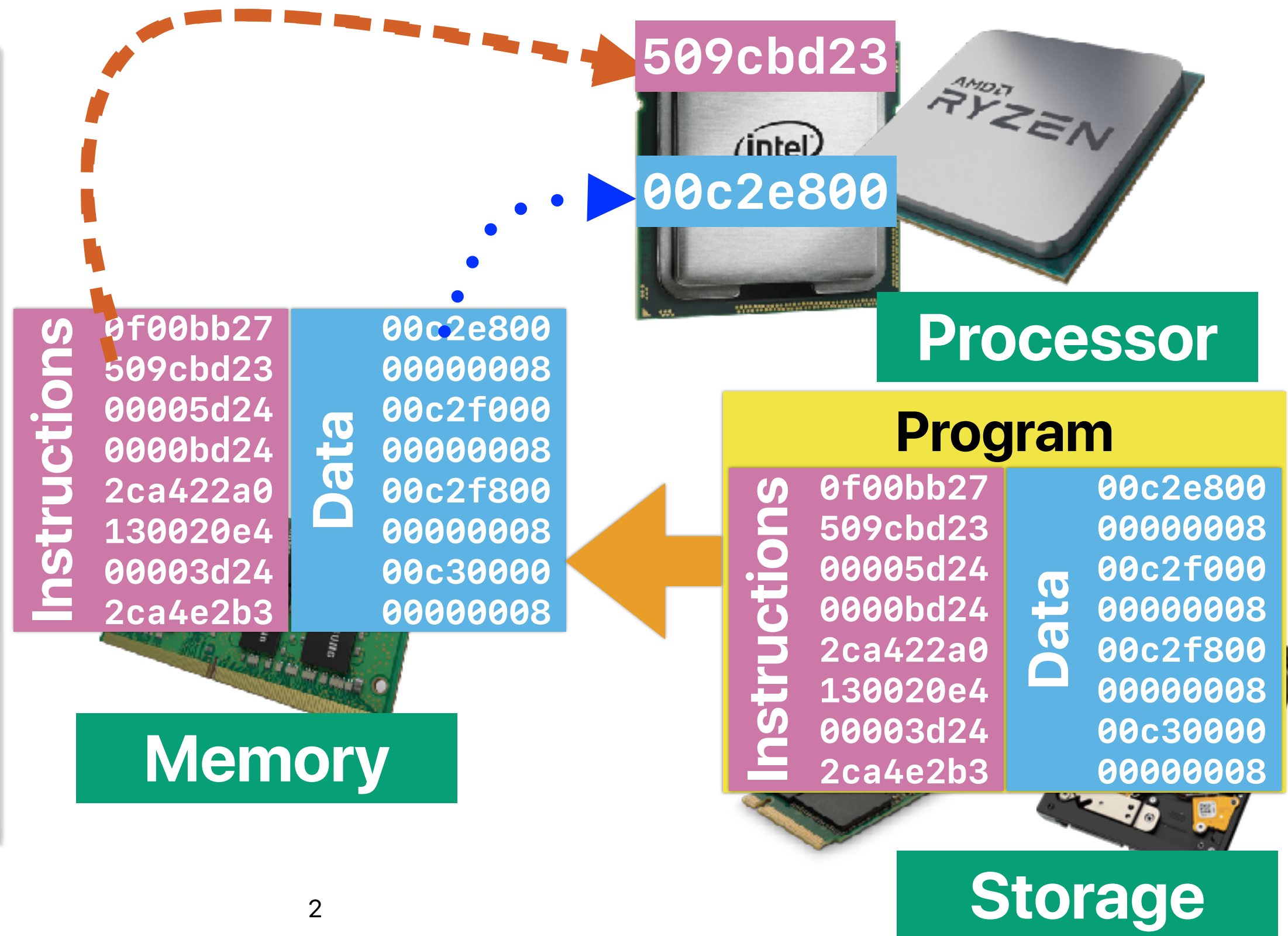


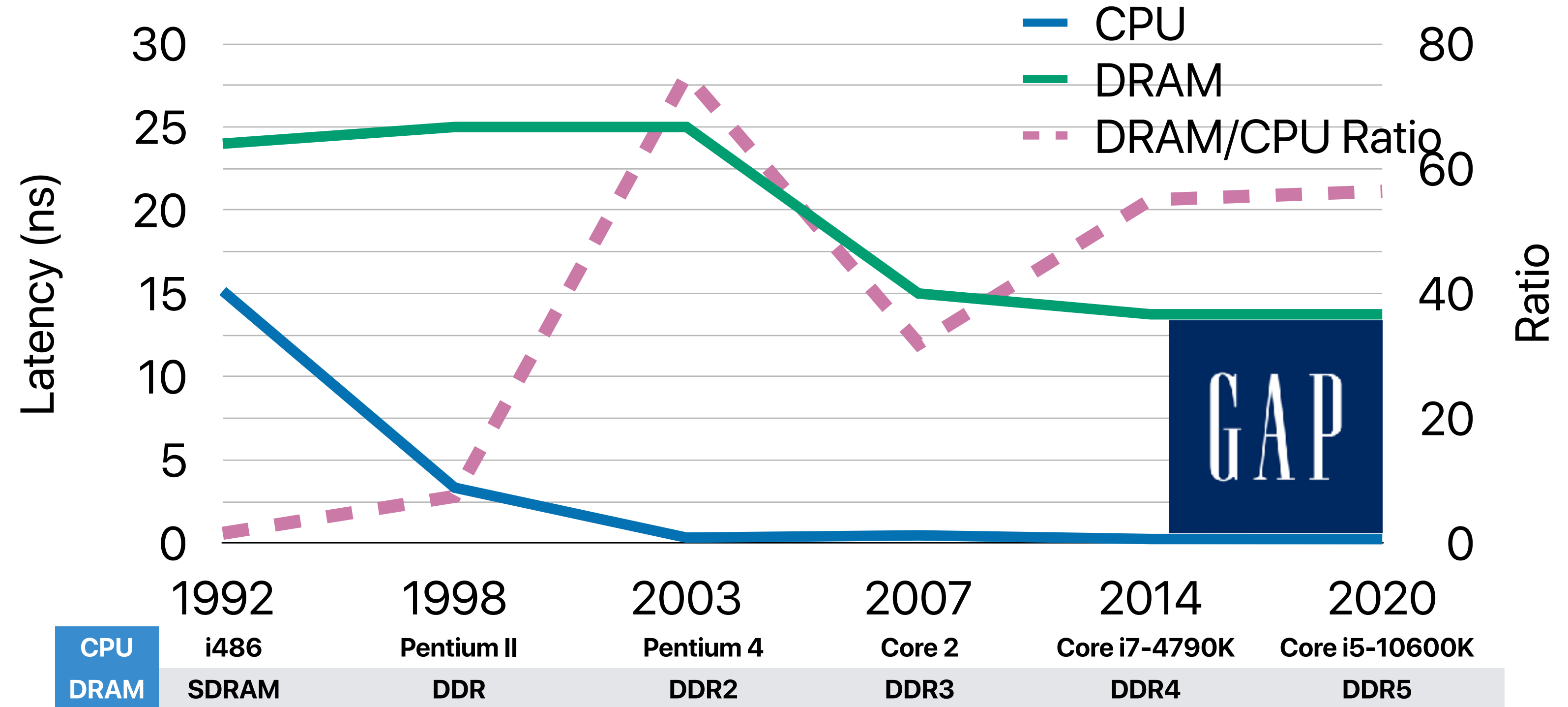
# Memory Hierarchy: Basics

Hung-Wei Tseng

# Recap: von Neumann Architecture



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

- Corollary 1 — each optimization has an upper bound  $Speedup_{max}(f, \infty) = \frac{1}{(1-f)}$

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

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

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

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

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

- 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

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

- Corollary 6: Don't hurt the non-common case too much

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

# 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!



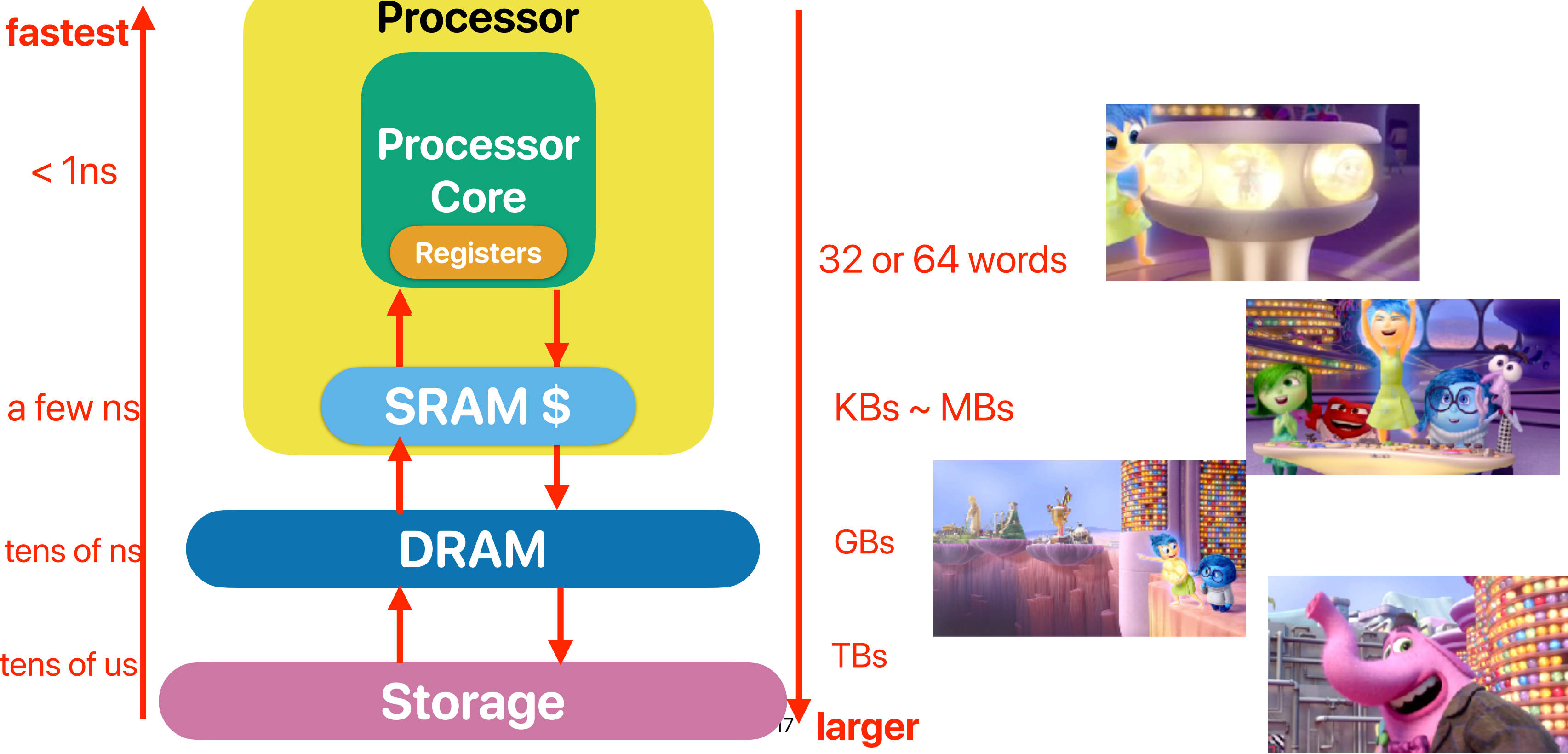
# Alternatives?

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	

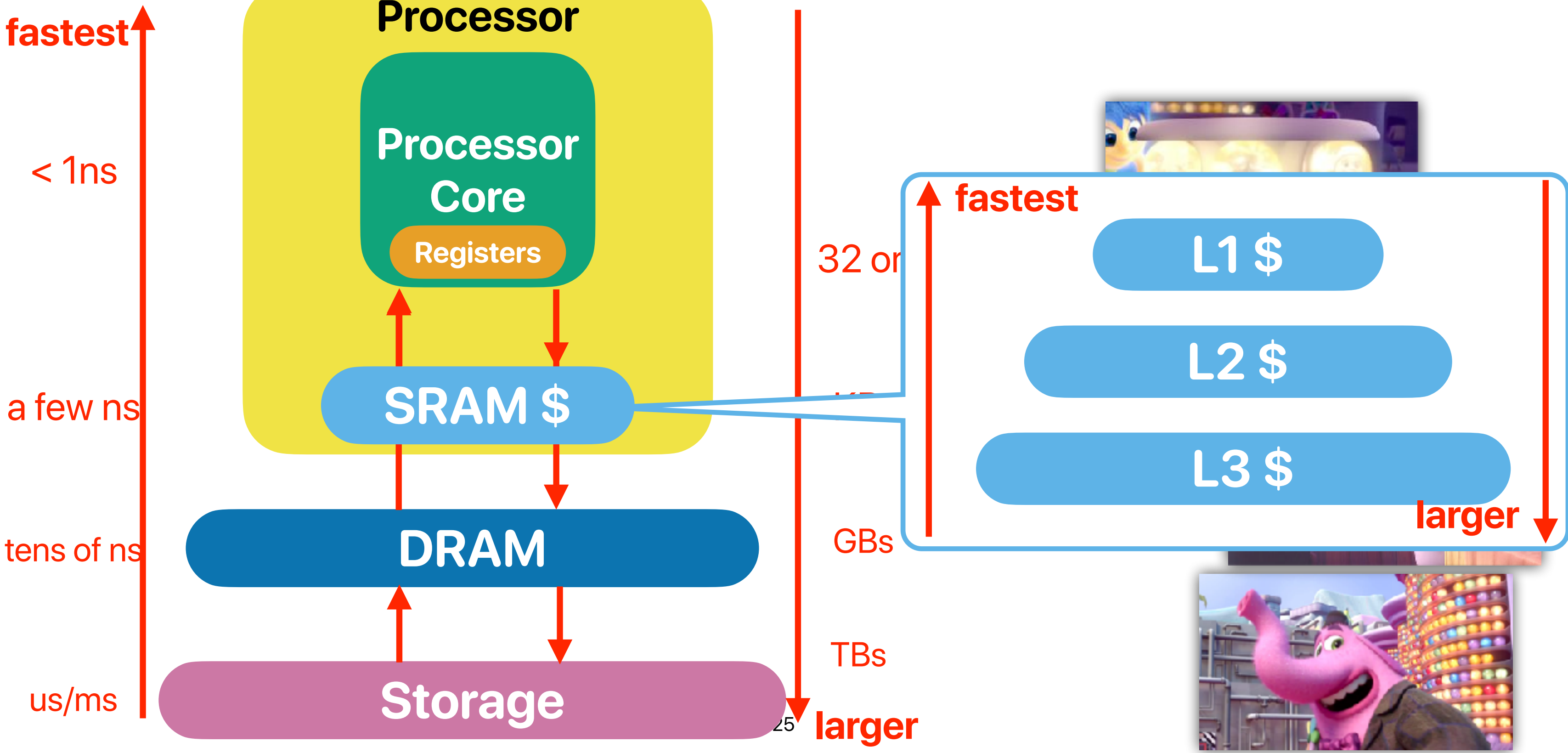
Cache

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	1.2 V
Specification	
AMD Ryzen 7 7700X 8-Core Processor	
Ext. Family	7
Ext. Model	100
Revision	RPL
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	14
Ext. Model	07
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	

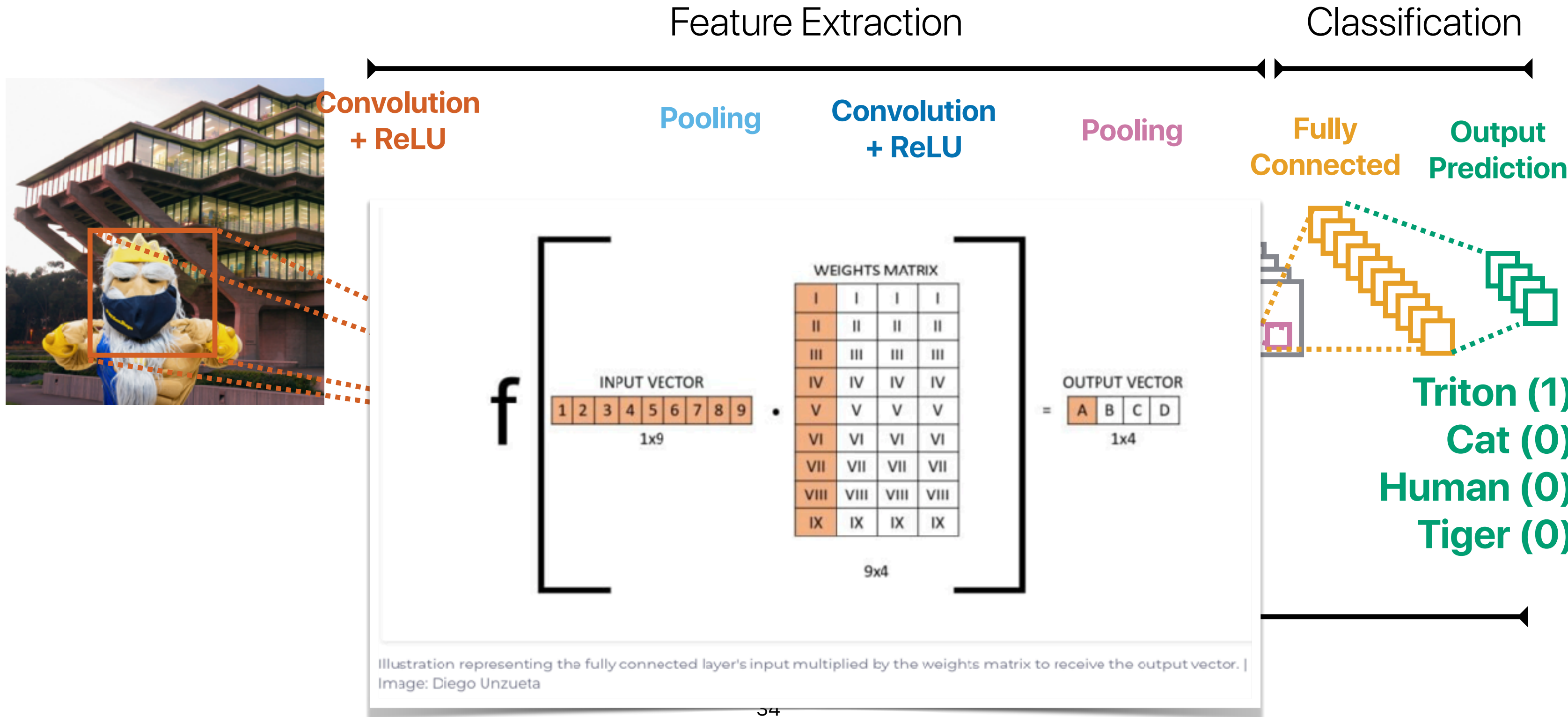
# Take-aways: inside out our memory hierarchy

- 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**



# The Machine Learning Inference Pipeline



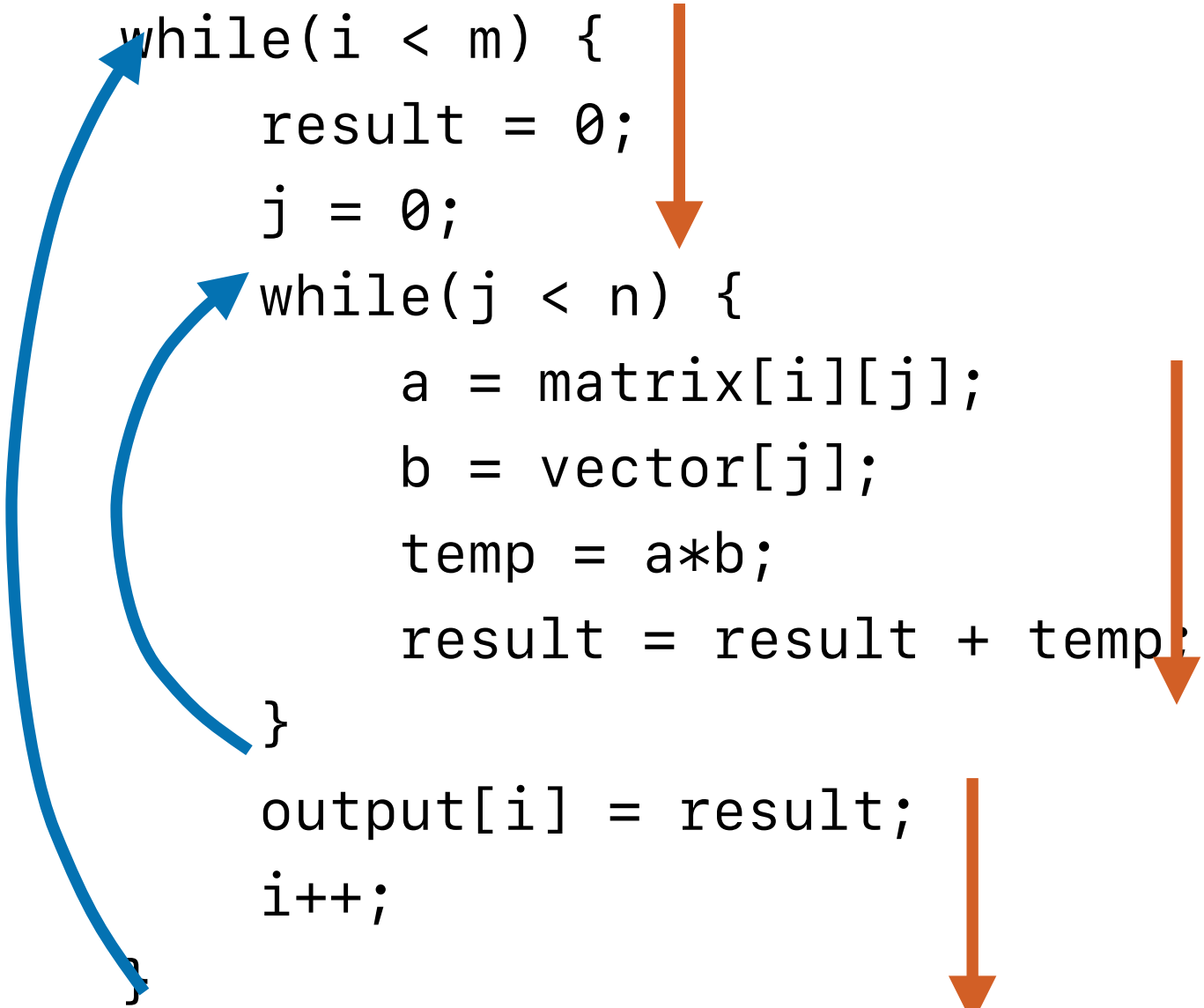
# Code also has locality

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 the next instruction

**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!
  - Spatial locality
  - Temporal locality

# Designing a hardware to exploit locality

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

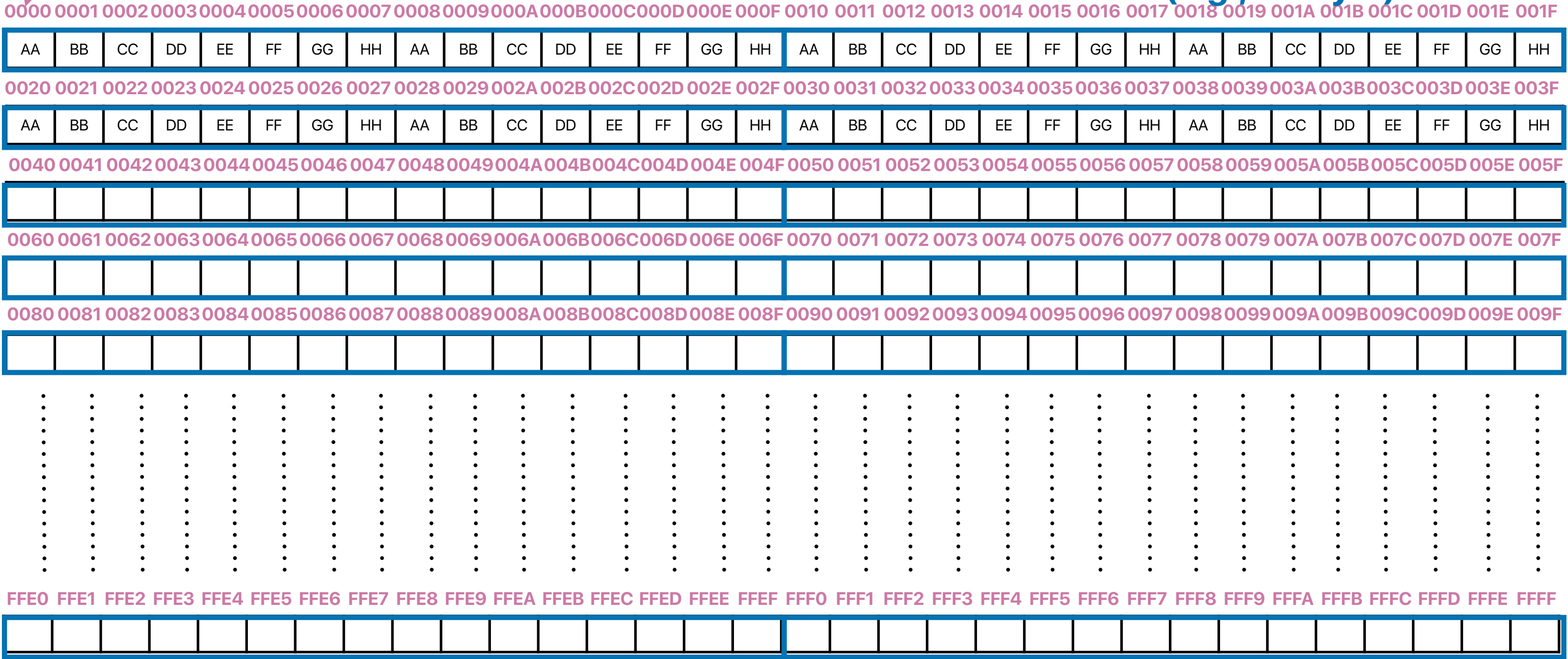
**We need to "cache consecutive memory locations" every time — the cache should store a "block" of code/data**

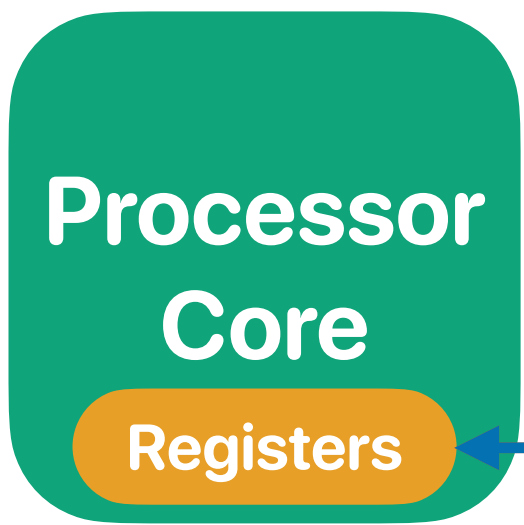
- Temporal locality — application revisit the same thing again and again
  - Code — loops, frequently invoked functions
    - Typically tens of static instructions — at most several KBs
  - Data — program can read/write the same data many times (e.g., vectors in matrix-vector product)



# Block and the memory space

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



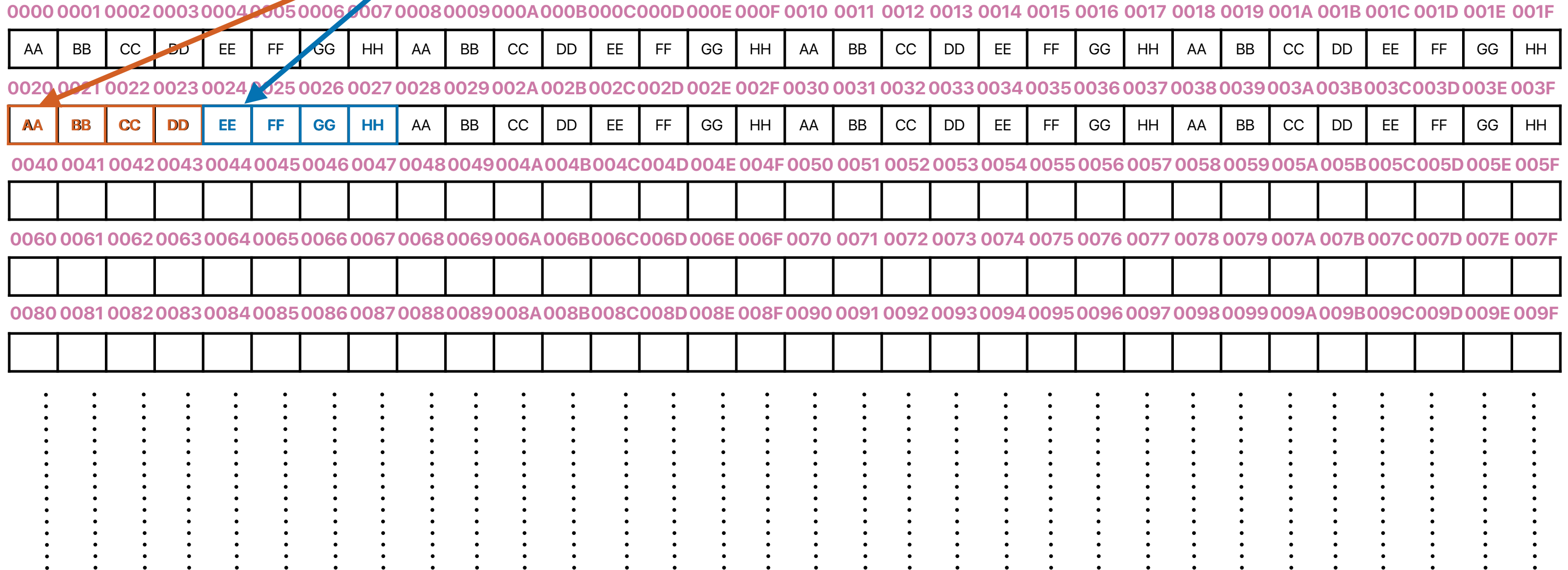


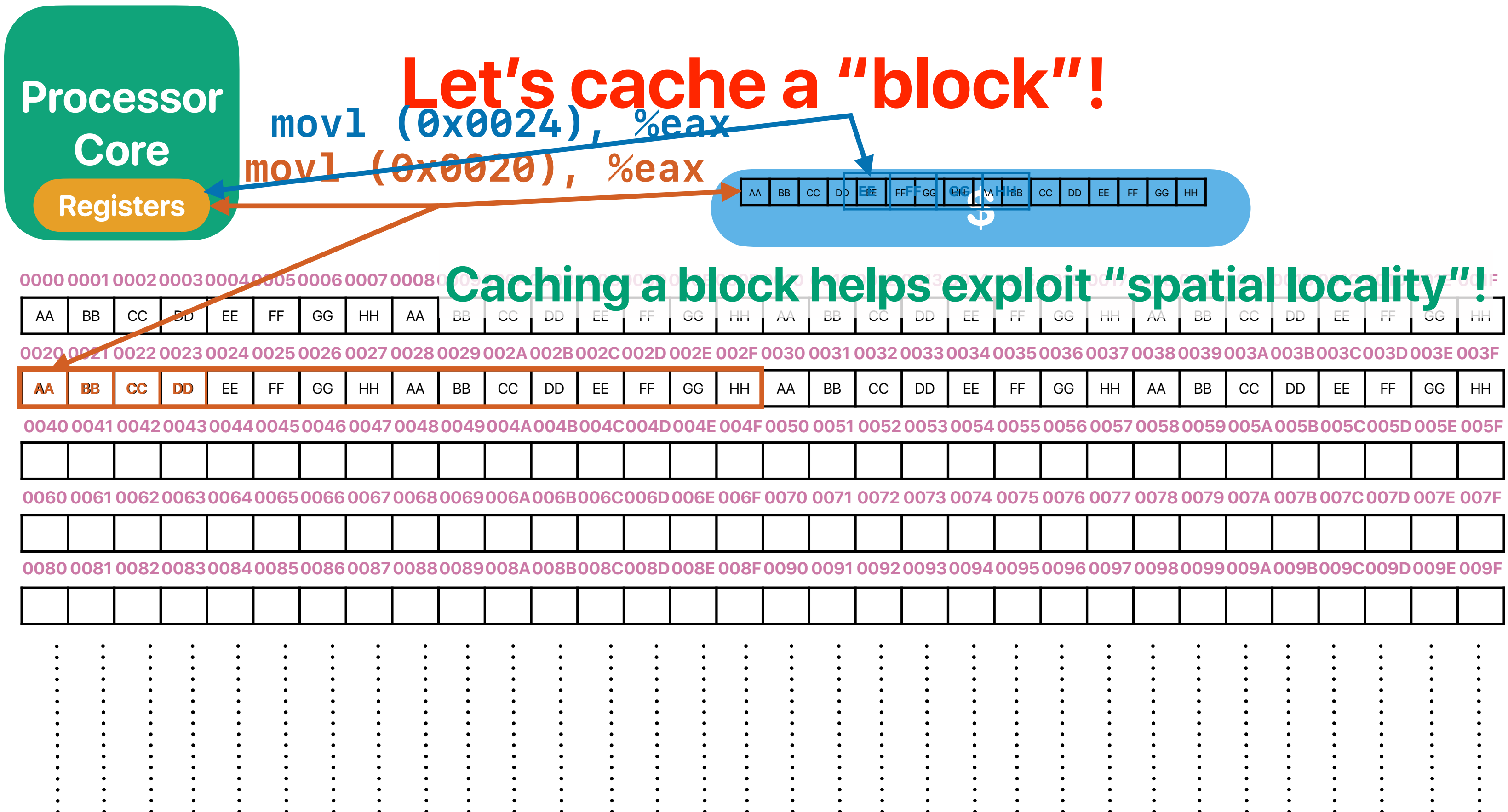
# When there is a "movl"

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**

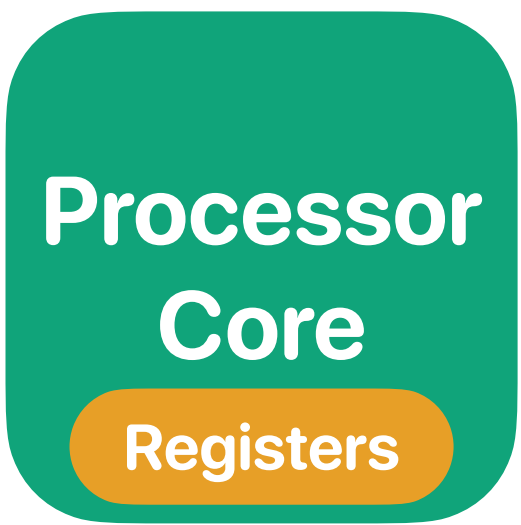
- Temporal locality — application revisit the same thing again and again

**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)





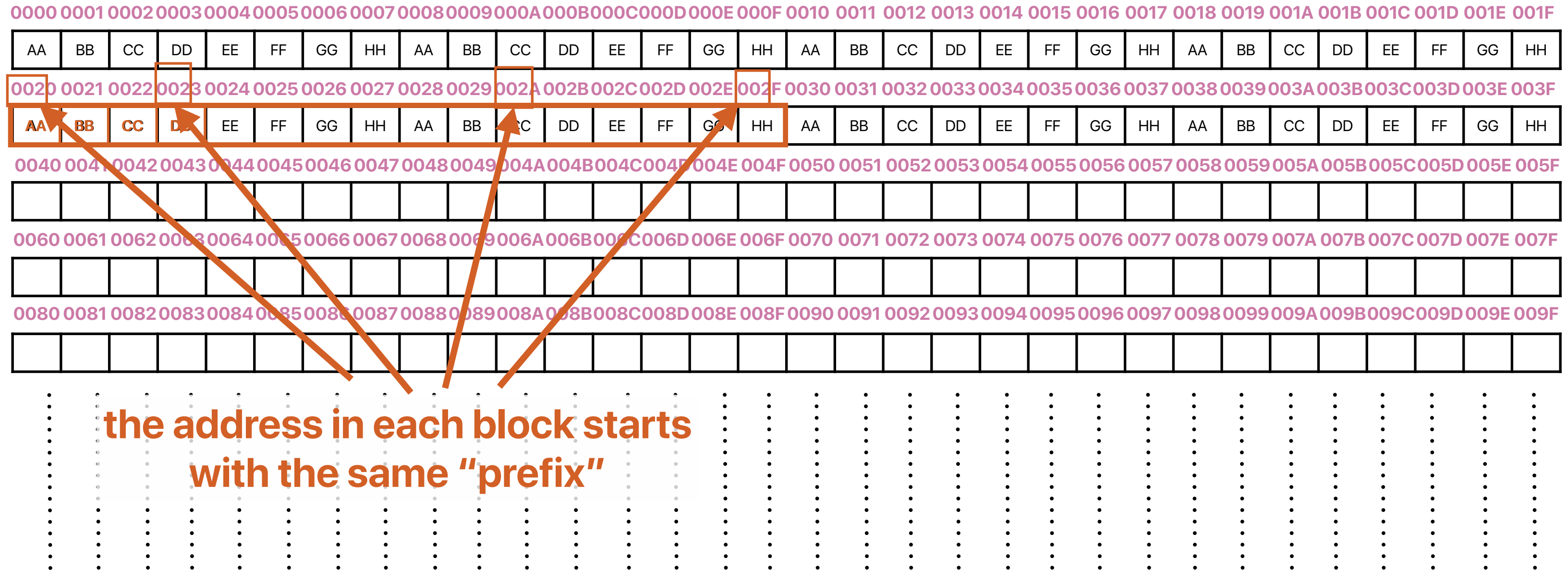
# How to tell who is there?

?															
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:															

# Registers

# Let's cache a "block"!

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

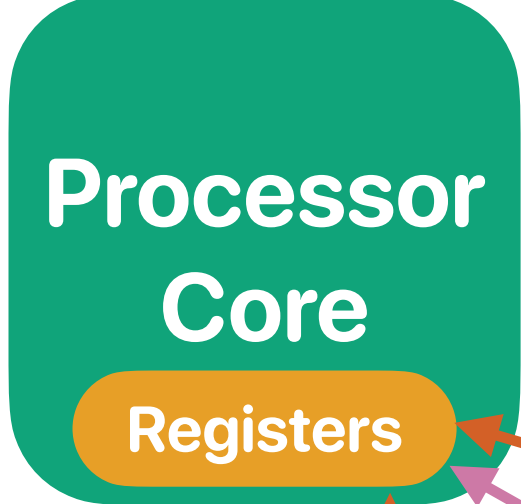
~~movl (0x0020), %eax~~

# Registers

0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
0x0000	0x0001	0x0002	0x0003	0x0004	0x0005	0x0006	0x0007	0x0008	0x0009	0x000A	0x000B	0x000C	0x000D	0x000E	0x000F
01	2	3	4	5	6	7	8	9	A	B	C	D	E	F	



## 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

	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!
ry	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:

# 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
```



# Take-aways: inside out our memory hierarchy

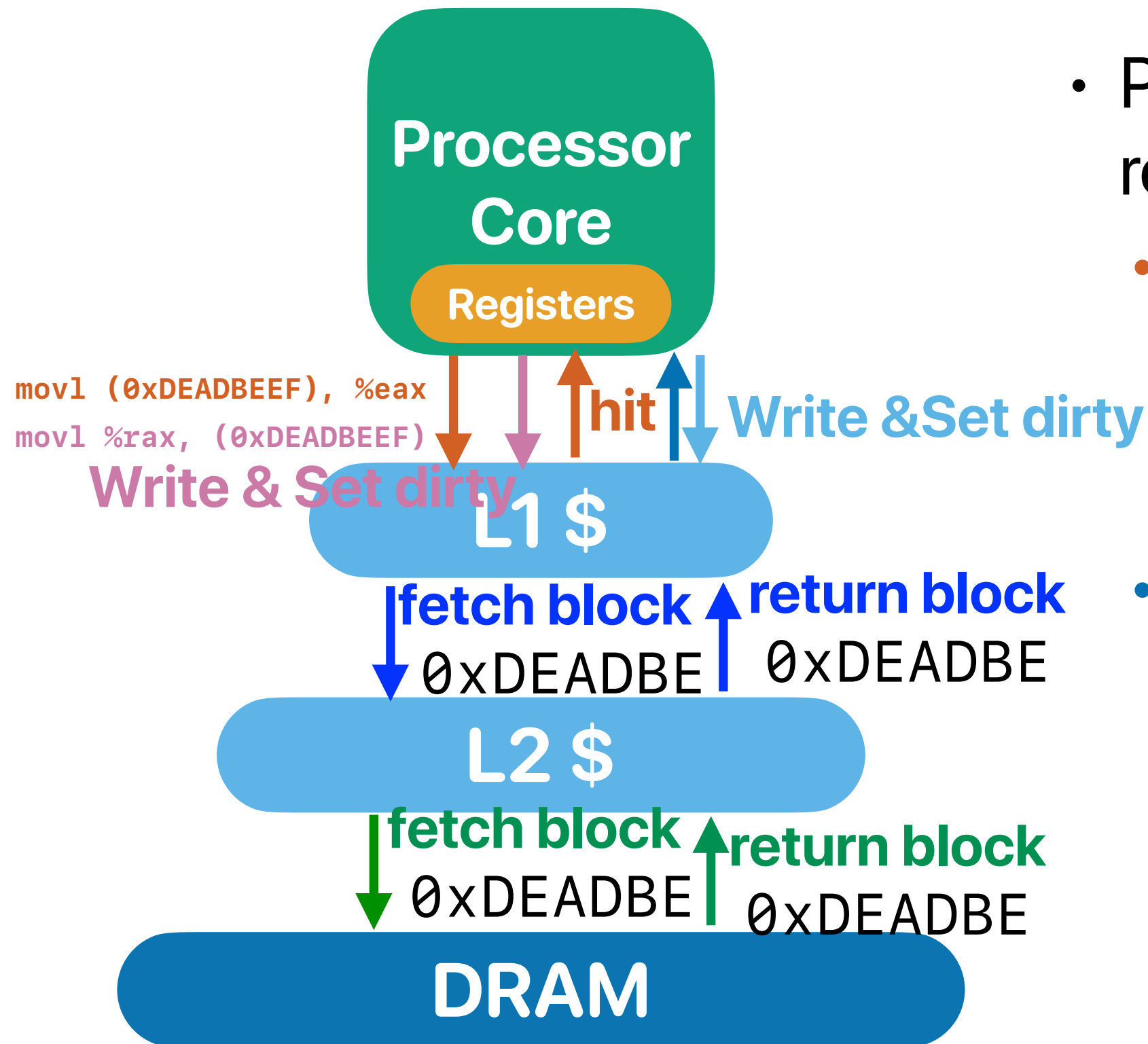
- Memory access time is the most critical performance problem
  - One memory operation is as expensive as 50 arithmetic operations
  - Processor has to fetch instructions from memory
  - We have an average of 33% of data memory access instructions!
- Hierarchical caching with small amount of SRAMs will work if we can efficiently capture data and instructions
- Caching is possible! Most of time, we only work on a small amount of data!
  - Spatial locality
  - Temporal locality
- Basic cache structures —
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Tags to distinguish cached blocks

# Take-aways: designing caches

- Basic cache structures —
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Tags to distinguish cached blocks

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

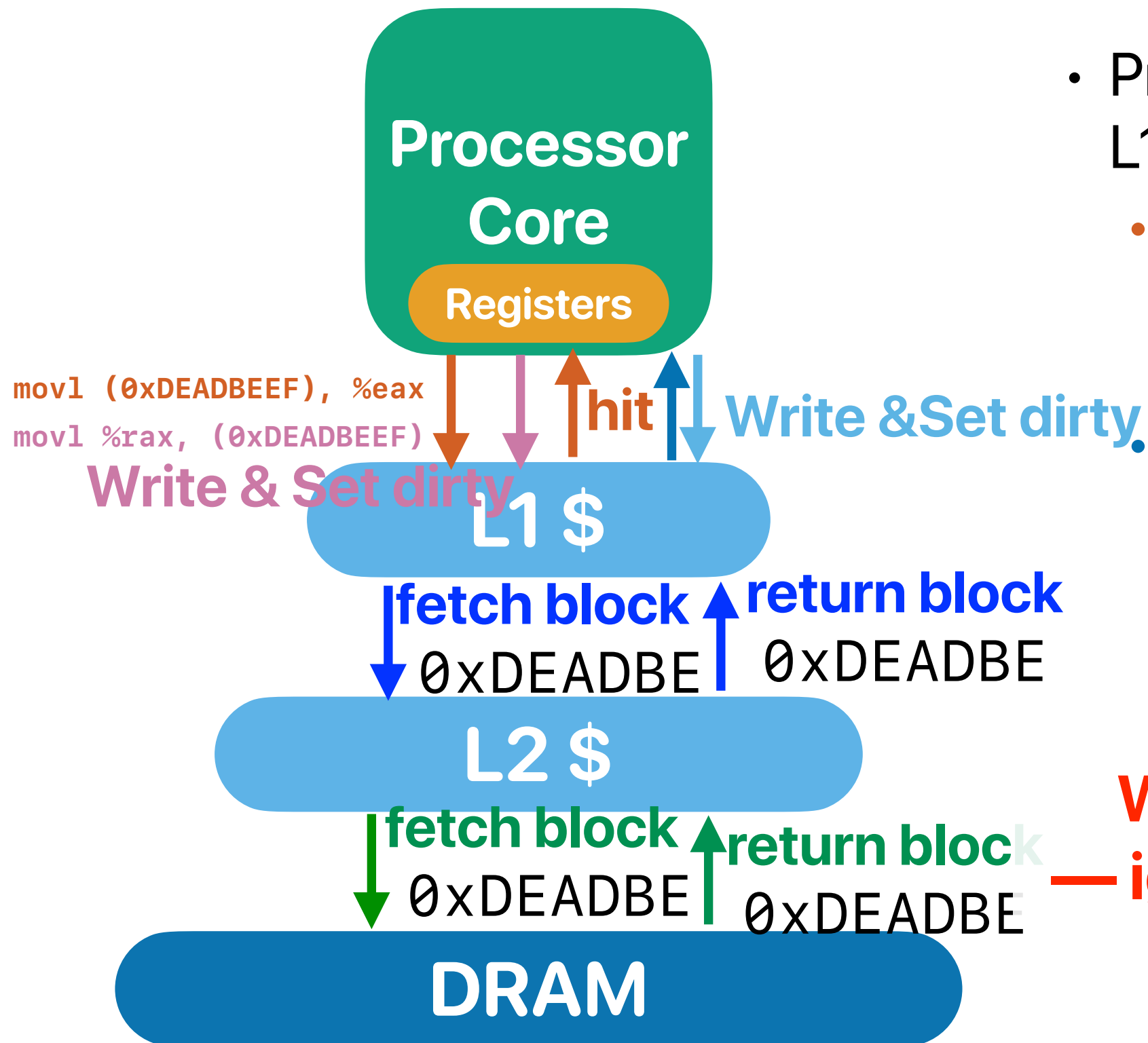
# Processor/cache interaction



- Processor sends memory access request to L1-\$
  - **if hit & it's a read**
    - **Read: return data**
    - **Write: Update "ONLY" in L1 and set DIRTY** **Why don't we write to L2?**  
**— Too slow**
  - **if miss**
    - Fetch the requesting block from lower-level memory hierarchy and place in the cache
    - **Present the write "ONLY" in L1 and set DIRTY** **What if we run out of \$**

# What if we run out of \$ blocks?

# Processor/cache interaction

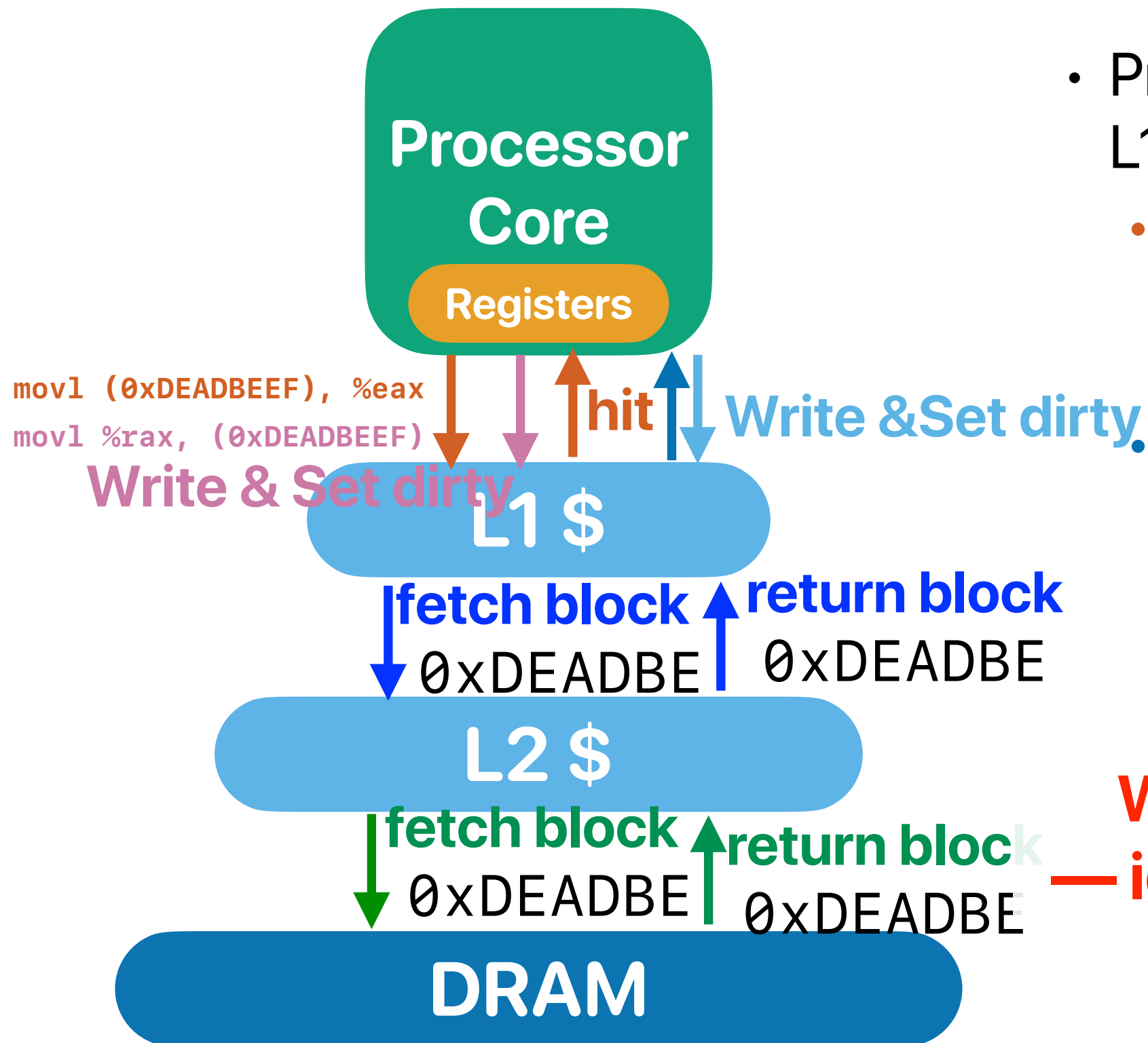


- Processor sends memory access request to L1-\$
  - **if hit & it's a read**
    - Read: return data
    - Write: Update "ONLY" in L1 and set DIRTY
  - **if miss**
    - If there an empty block — place the data there
    - If NOT (most frequent case) — select a **victim block**
      - Least Recently Used (LRU) policy

**What if the victim block is modified?**  
**— ignoring the update is not acceptable!**

- Present the write "ONLY" in L1 and set DIRTY

# Processor/cache interaction

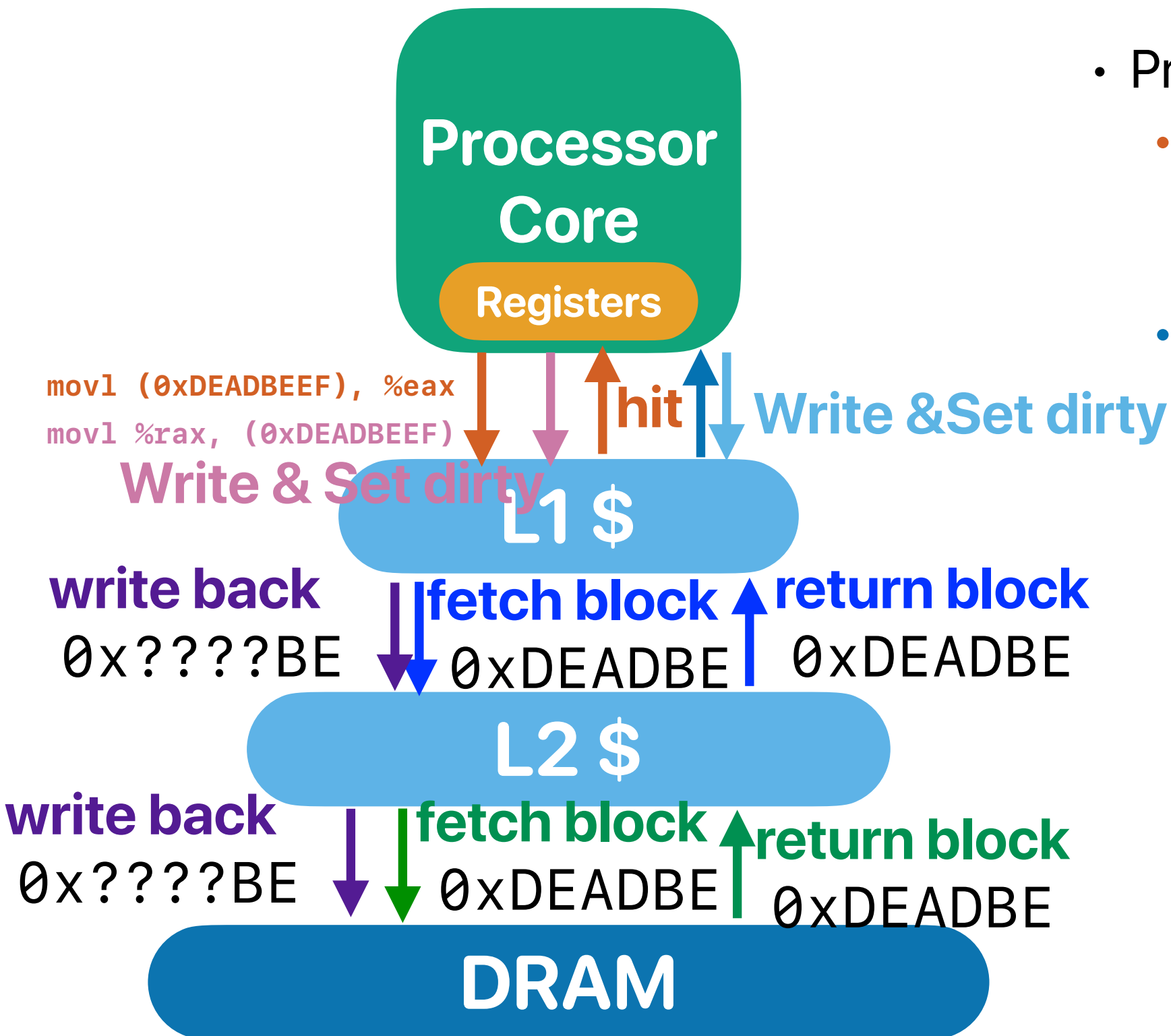


- Processor sends memory access request to L1-\$
  - if hit & it's a read
    - Read: return data
    - Write: Update "ONLY" in L1 and set DIRTY
  - if miss
    - If there an empty block — place the data there
    - If NOT (most frequent case) — select a **victim block**
      - Least Recently Used (LRU) policy

**What if the victim block is modified?**  
**— ignoring the update is not acceptable!**

- Present the write "ONLY" in L1 and set DIRTY

# Processor/cache interaction



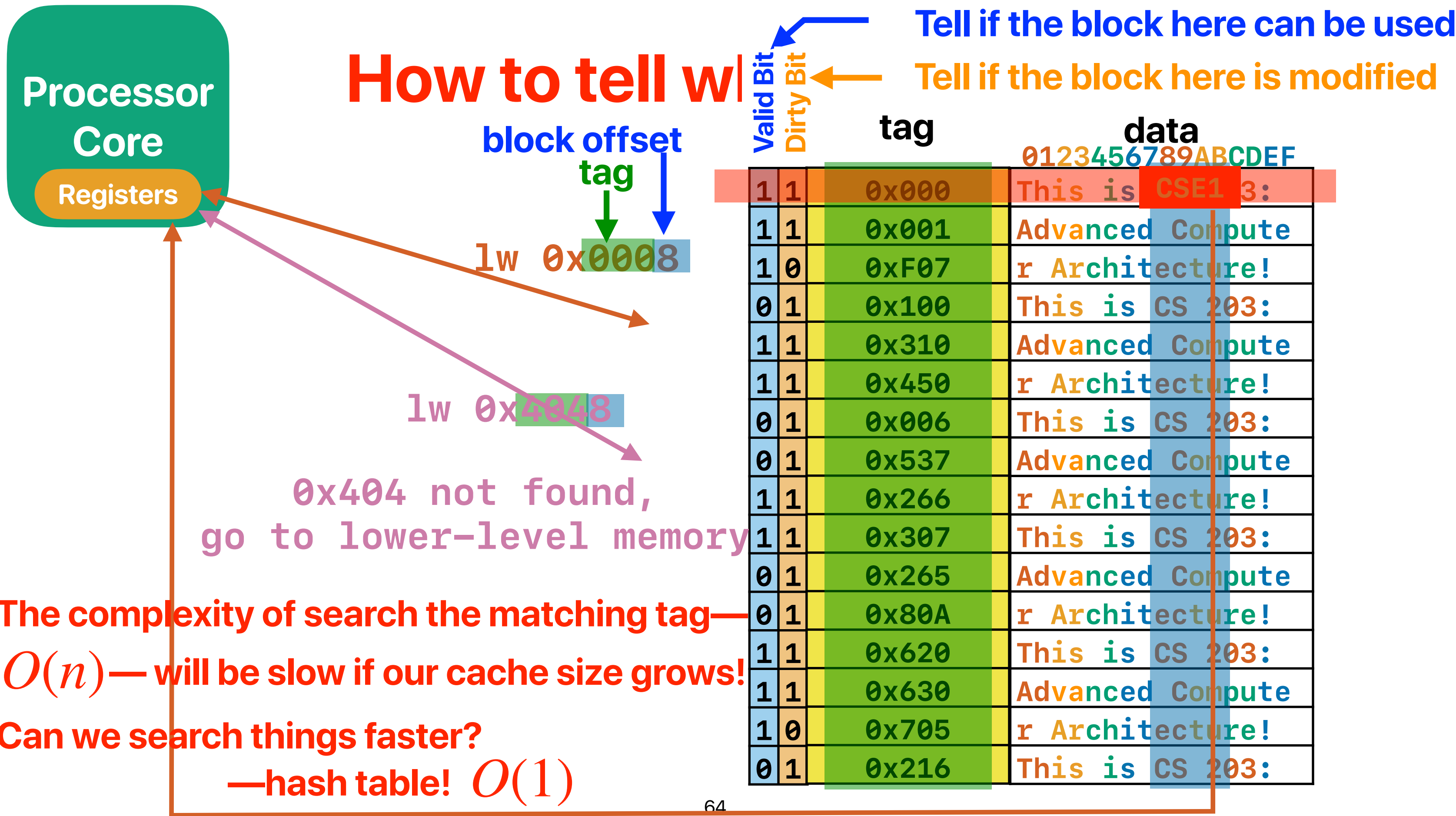
- Processor sends memory access request to L1-\$
  - **if hit & it's a read**
    - **Read:** return data
    - **Write:** Update "ONLY" in L1 and set DIRTY
  - **if miss**
    - If there an empty block — place the data there
    - If NOT (most frequent case) — select a **victim block**
      - Least Recently Used (LRU) policy
    - If the victim block is "dirty" & "valid"
      - **Write back** the block to lower-level memory hierarchy
      - If write-back or fetching causes any miss, repeat the same process
    - Fetch the requesting block from lower-level memory hierarchy and place in the cache
    - **Present the write "ONLY" in L1 and set DIRTY**

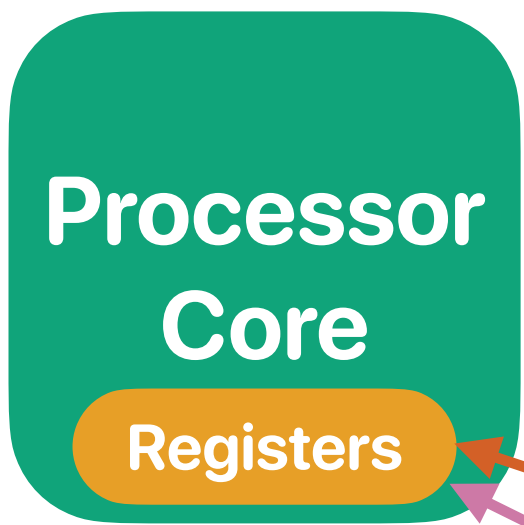


# Take-aways: designing caches

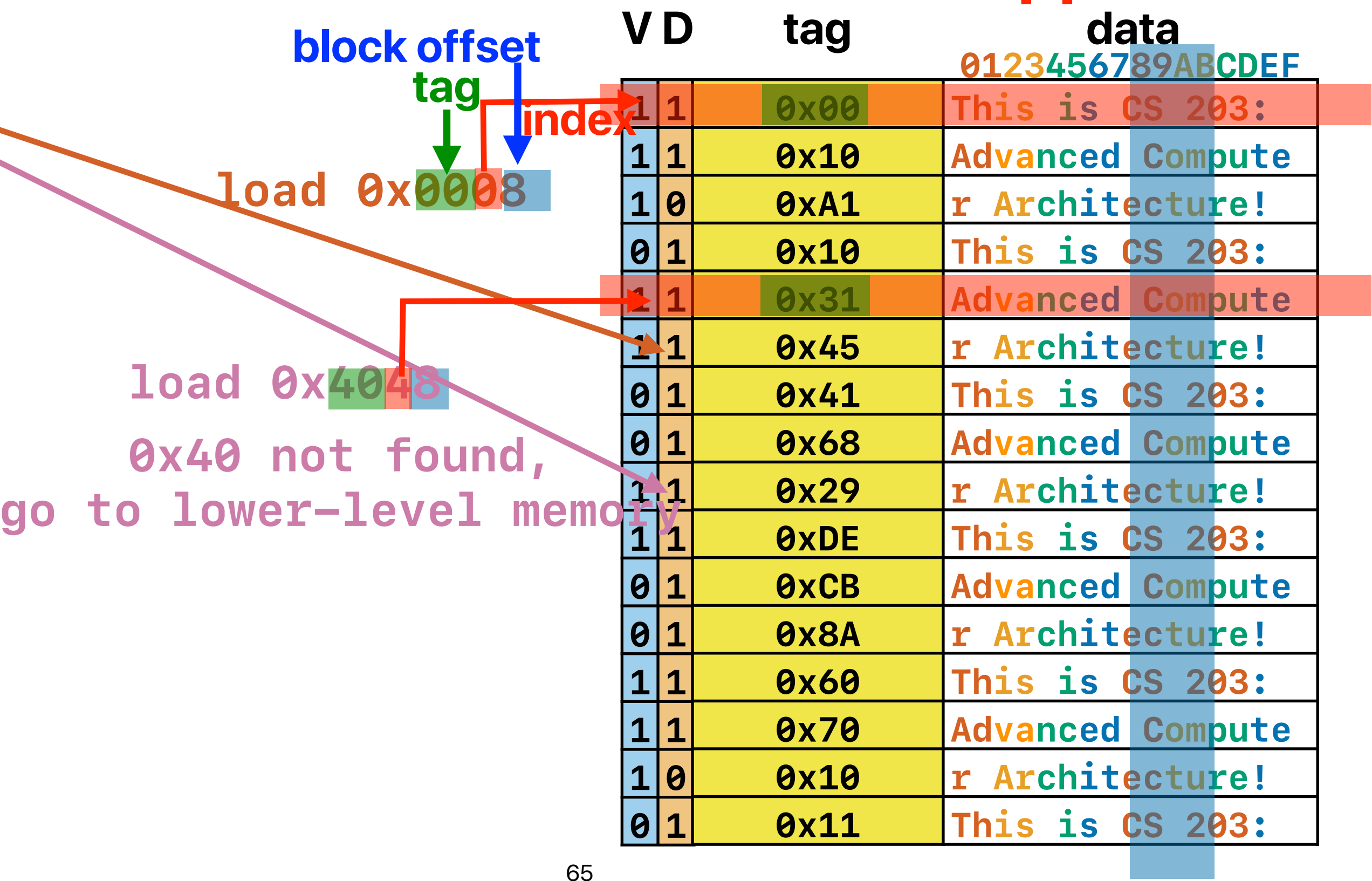
- Basic cache structures —
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Tags to distinguish cached blocks
- Hierarchical caching — data must be presented on the top level (L1) before the processor can use

# How to tell w



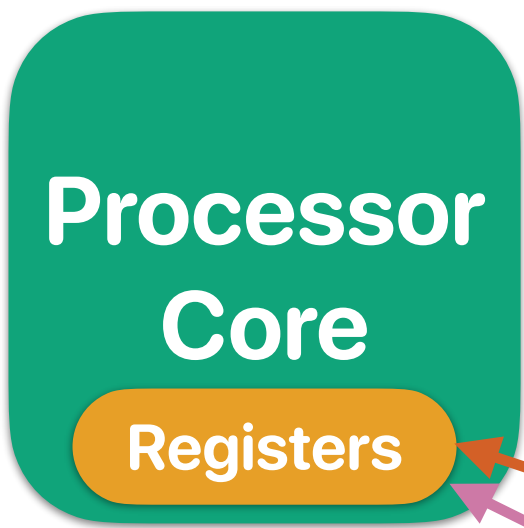


# Hash-like structure — direct-mapped cache

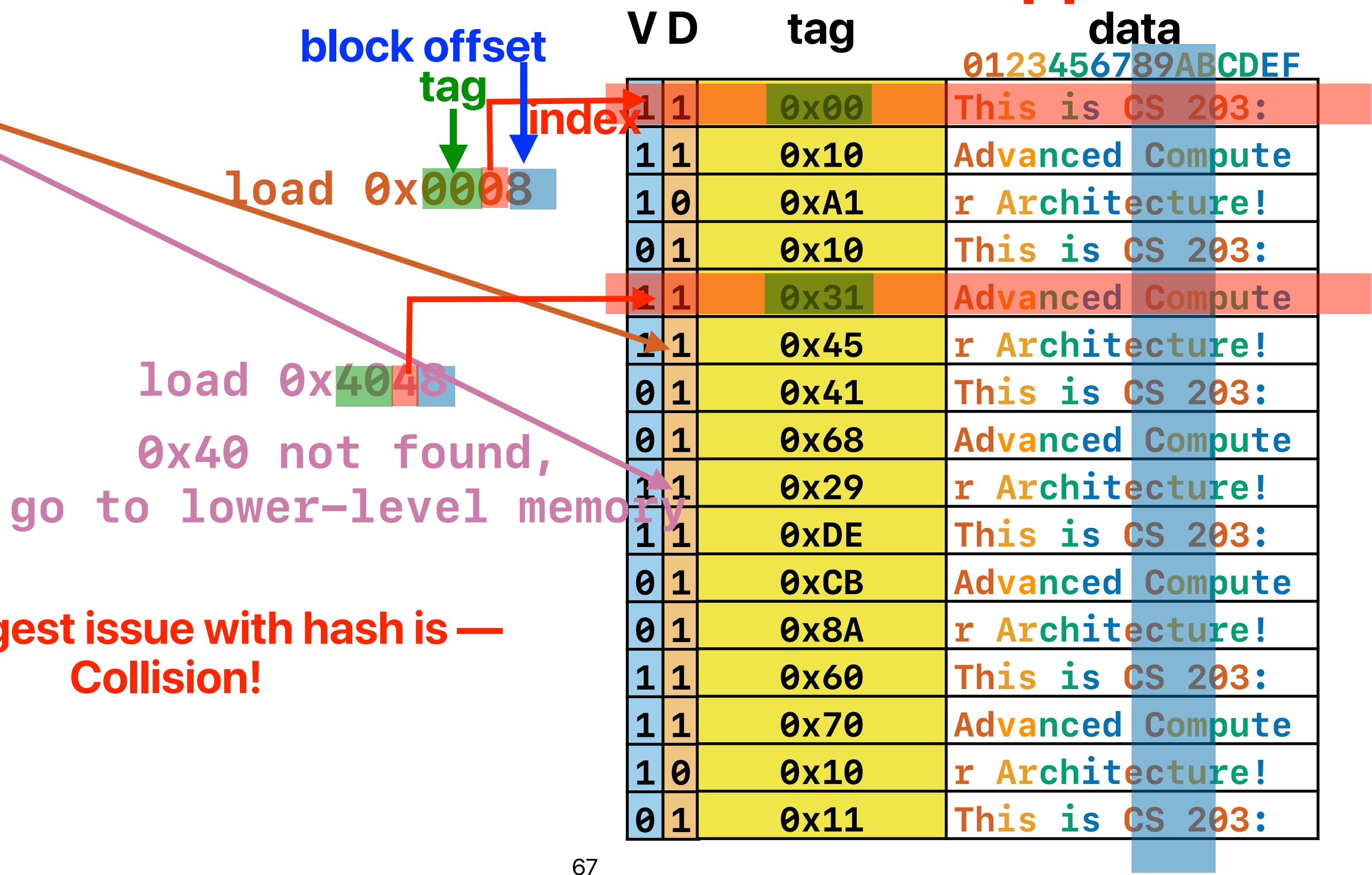


# Take-aways: designing caches

- Basic cache structures —
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Tags to distinguish cached blocks
- Hierarchical caching — data must be presented on the top level (L1) before the processor can use
- Optimizing cache structures
  - Hash block into "sets" to reduce the search time



# Hash-like structure — direct-mapped cache

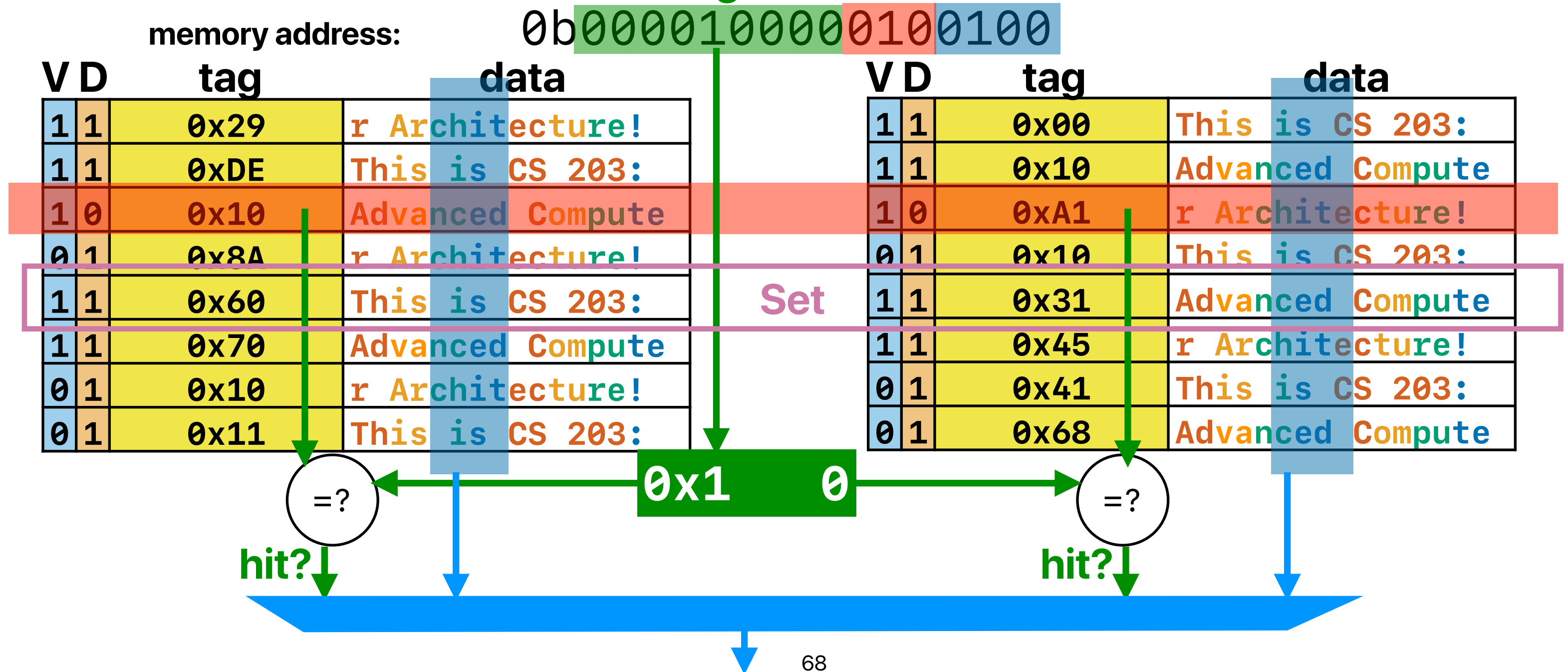


The biggest issue with hash is — Collision!

# Way-associative cache

memory address: 0x0 8 2 4

tag set index block offset



# 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
```



# Take-aways: designing caches

- Basic cache structures —
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Tags to distinguish cached blocks
- Hierarchical caching — data must be presented on the top level (L1) before the processor can use
- Optimizing cache structures
  - Hash block into "sets" to reduce the search time
  - Set-associativity to reduce the "collision" problem

# Way-associative cache

memory address:  $0x0$  8 2 4

set block

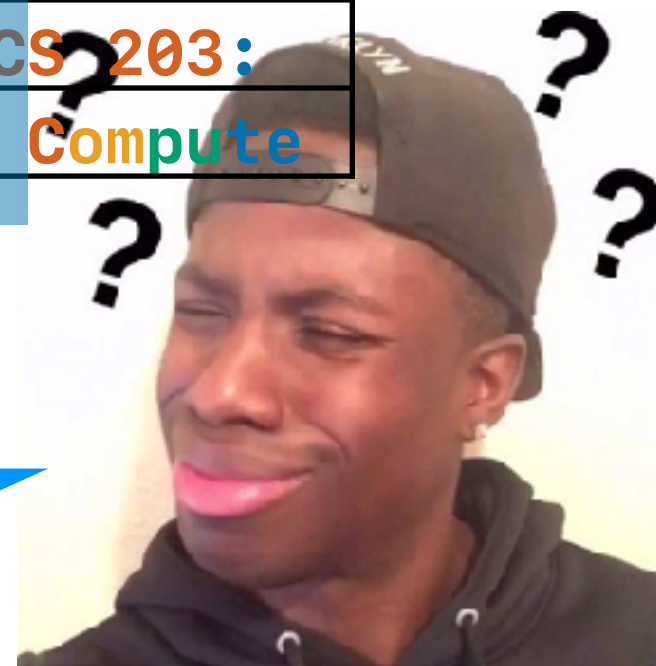
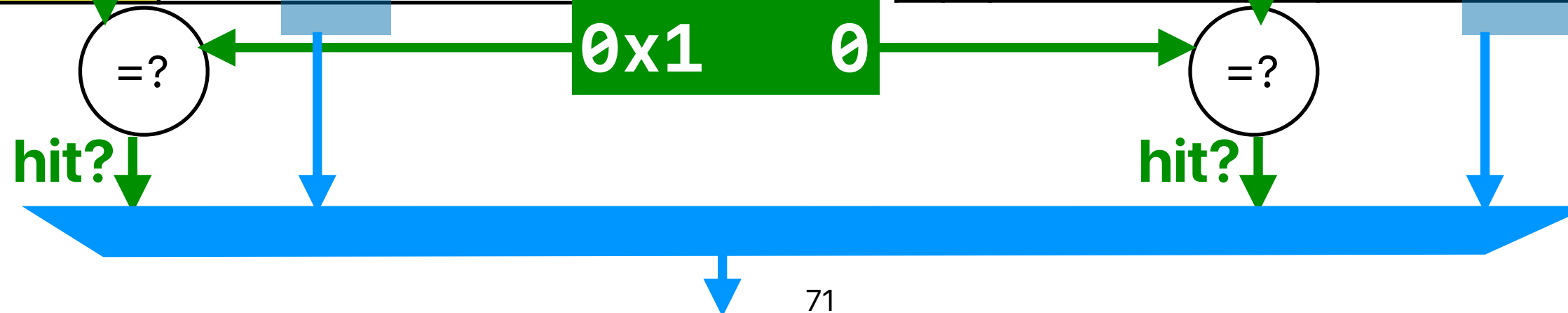
tag index offset

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



# **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

- Basic cache structures —
  - Caching in granularity of a block to capture spatial locality
  - Caching multiple blocks to keep frequently used data — temporal locality
  - Tags to distinguish cached blocks
- Hierarchical caching — data must be presented on the top level (L1) before the processor can use
- Optimizing cache structures
  - Hash block into "sets" to reduce the search time
  - Set-associativity to reduce the "collision" problem
- $C = A B S$ 
  - C: capacity
  - A: Associativity
  - S: Number of sets
  - $\lg(S)$ : Number of bits in set index
  - $\lg(B)$ : Number of bits in block offset

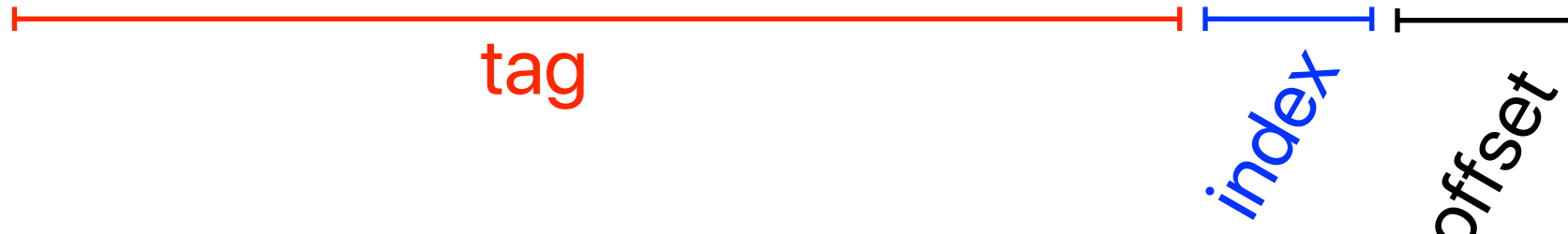
# **Simulate the cache!**



# 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:  $0x \quad 8 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 8 \quad 0$   
=  $0b \quad 1000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 1000 \quad 0000$   


# 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

```
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;
}
```

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	

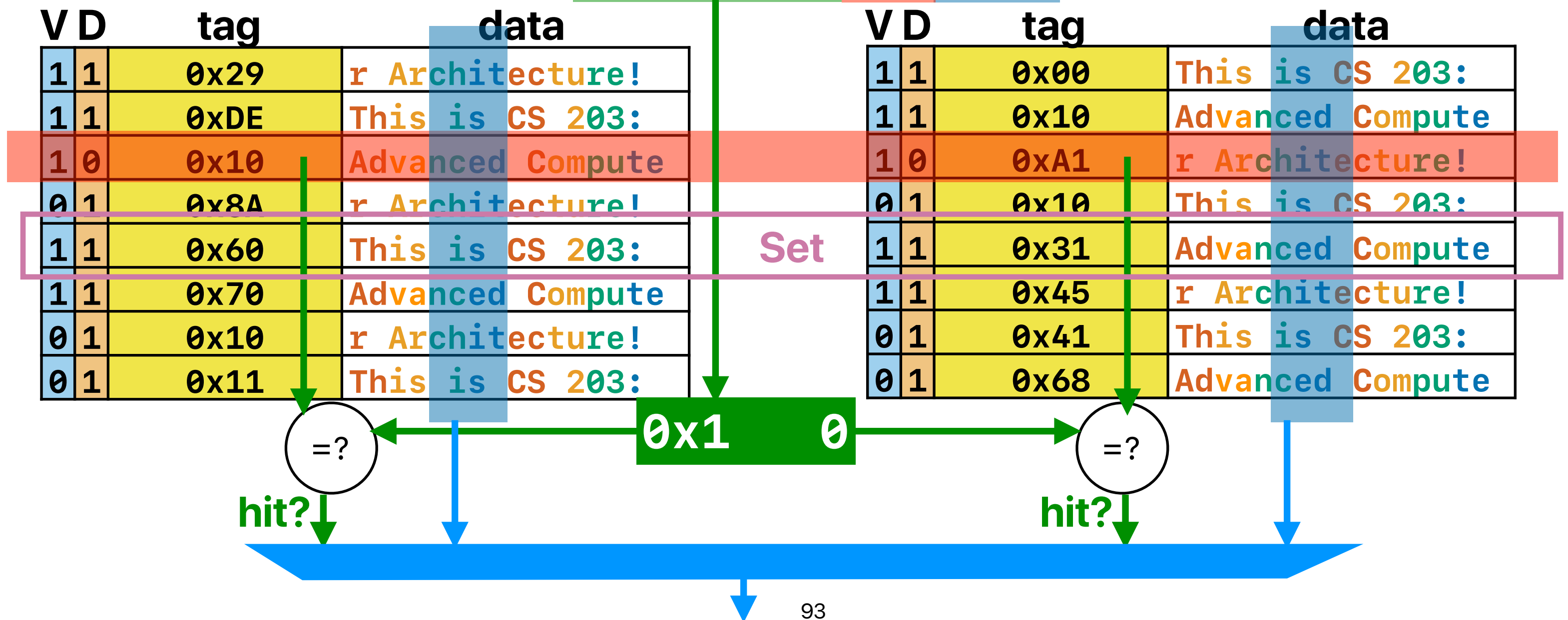
# Way-associative cache

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(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;  
}

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

# **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