

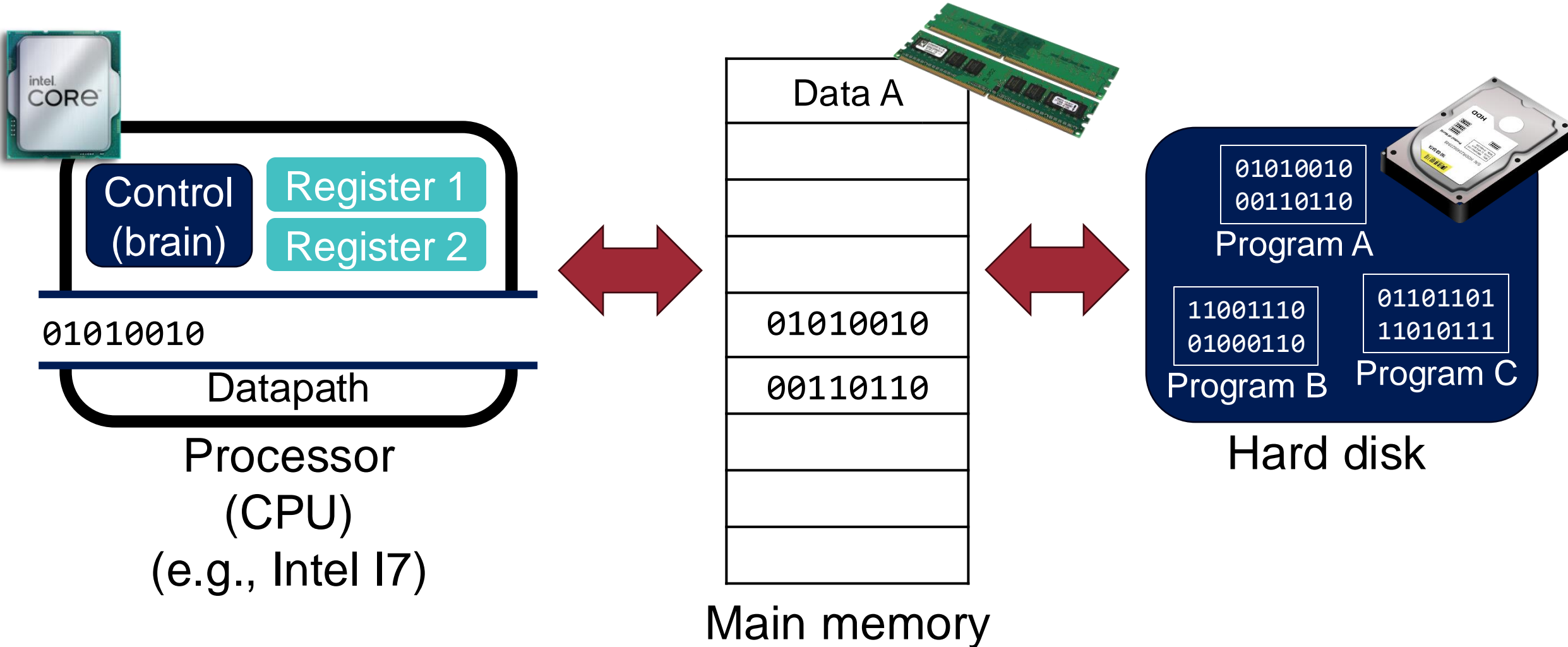
CSE261: Computer Architecture

16. Memory Hierarchy (2): Caches #1

Seongil Wi

Recap: Memory Hierarchy

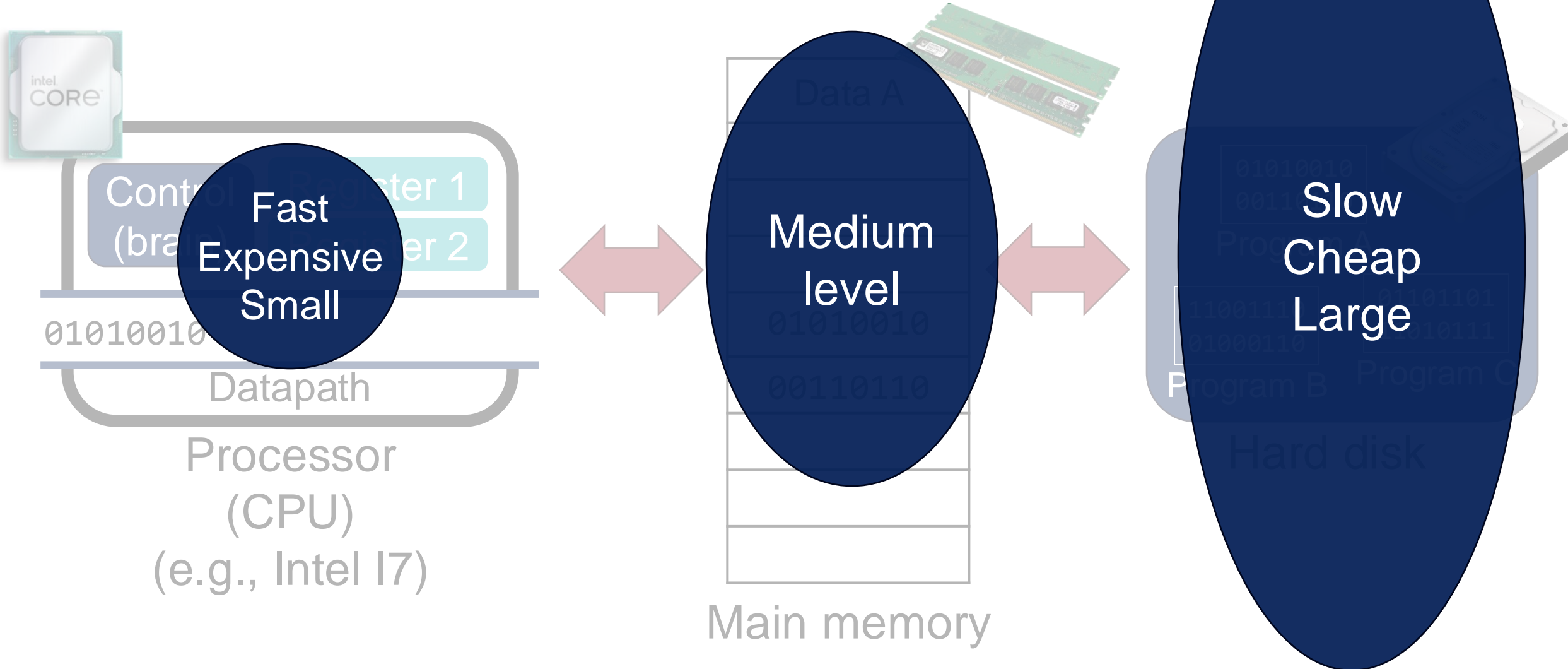
A structure that uses multiple levels of memories



Recap: Memory Hierarchy

3

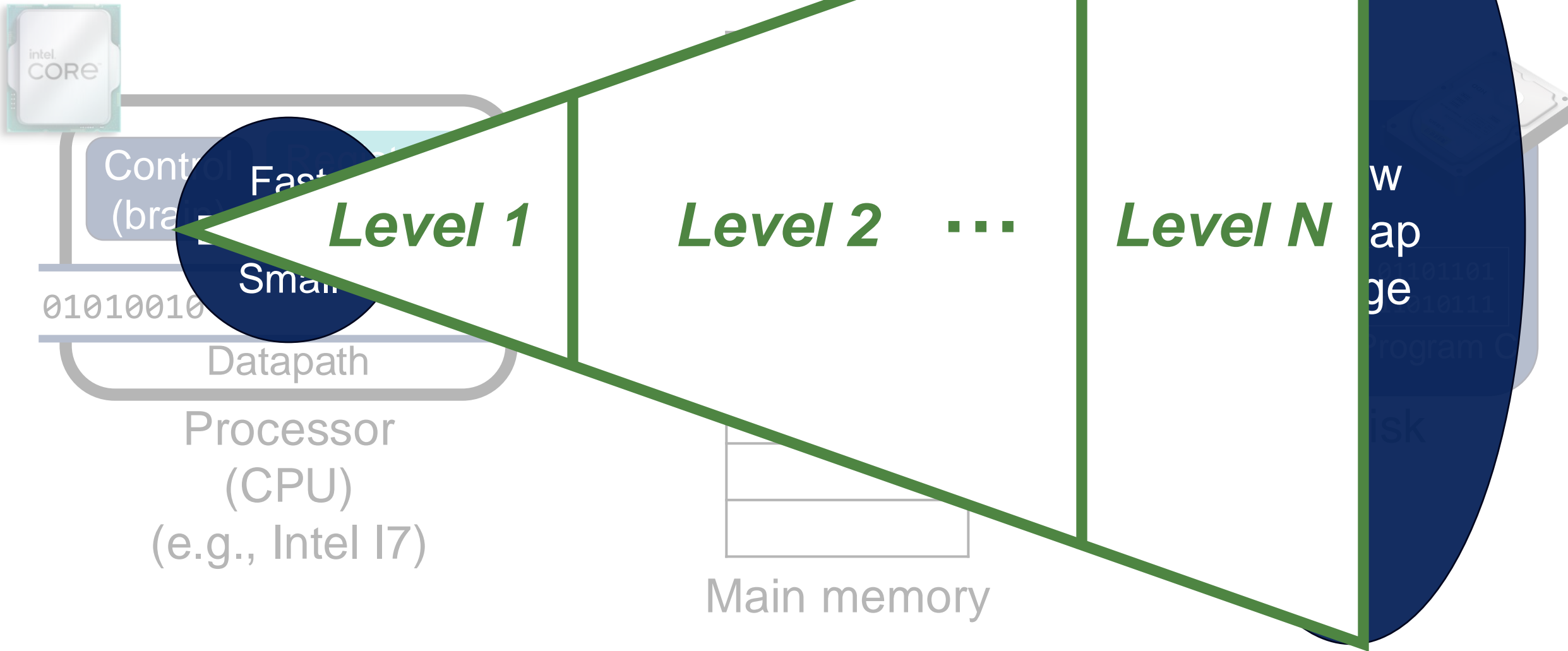
A structure that uses multiple levels of memories



Recap: Memory Hierarchy

4

A structure that uses multiple levels of mem



Recap: Locality (지역성)

The tendency to access the same set of memory locations repetitively over a short period of time

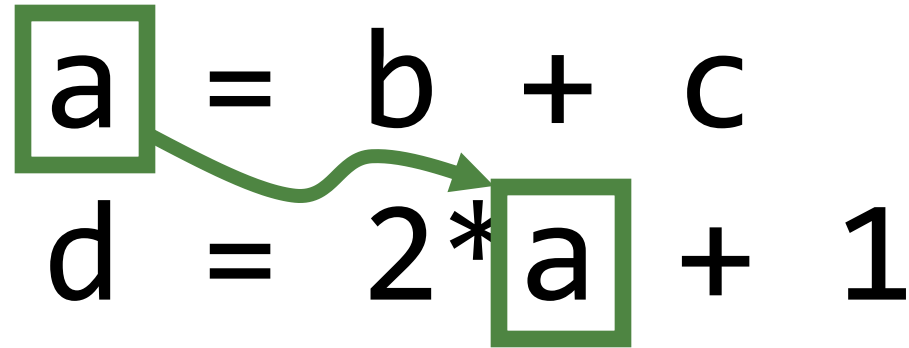
1. Temporal locality (locality in time)

- If an item is referenced, the same item will tend to be referenced again soon

2. Spatial locality (locality in space)

- If an item is referenced, nearby items will tend to be referenced soon

Recap: Temporal Locality Example


$$\begin{array}{l} \boxed{a} = b + c \\ d = 2 * \boxed{a} + 1 \end{array}$$

Temporal locality: items accessed recently are likely to be accessed again soon!

Recap: Spatial Locality Example

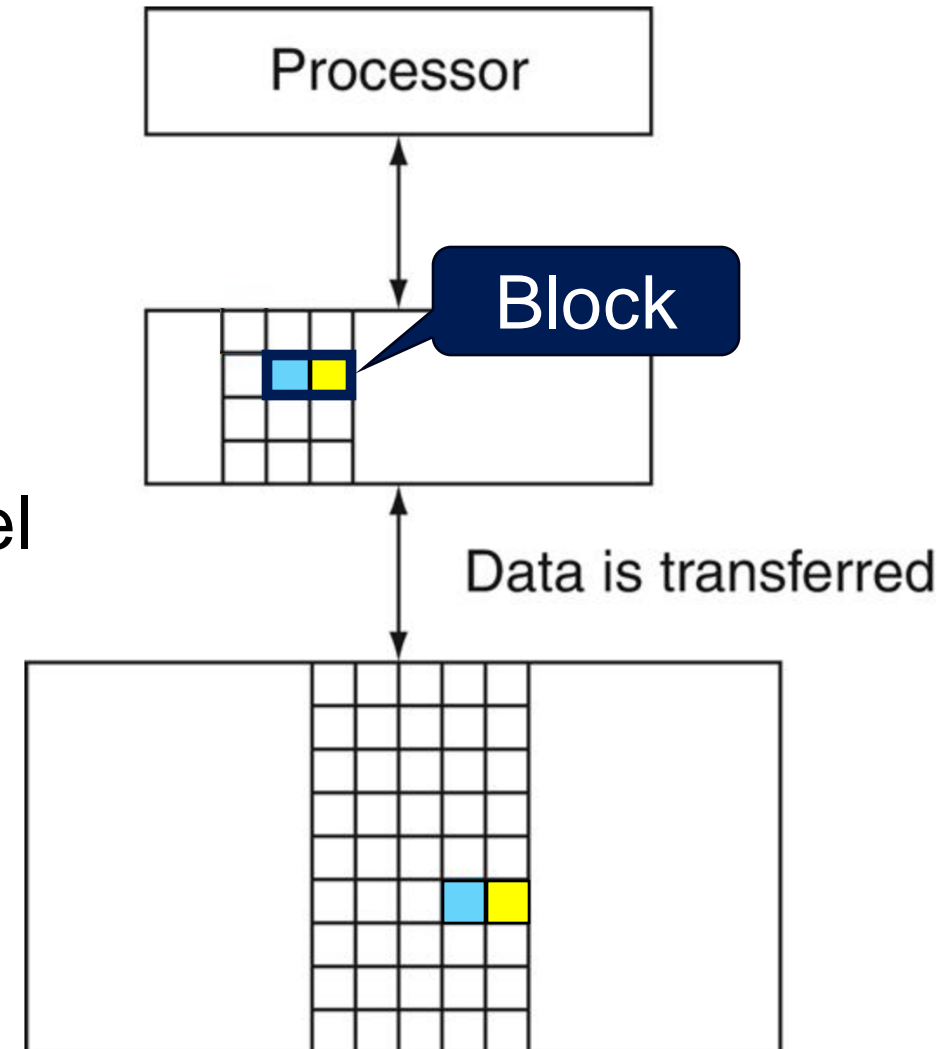
7

```
for (i=0; i<10; i++)  
    sum = sum + a[i]
```

Spatial locality: Items near those accessed recently are likely to be accessed soon

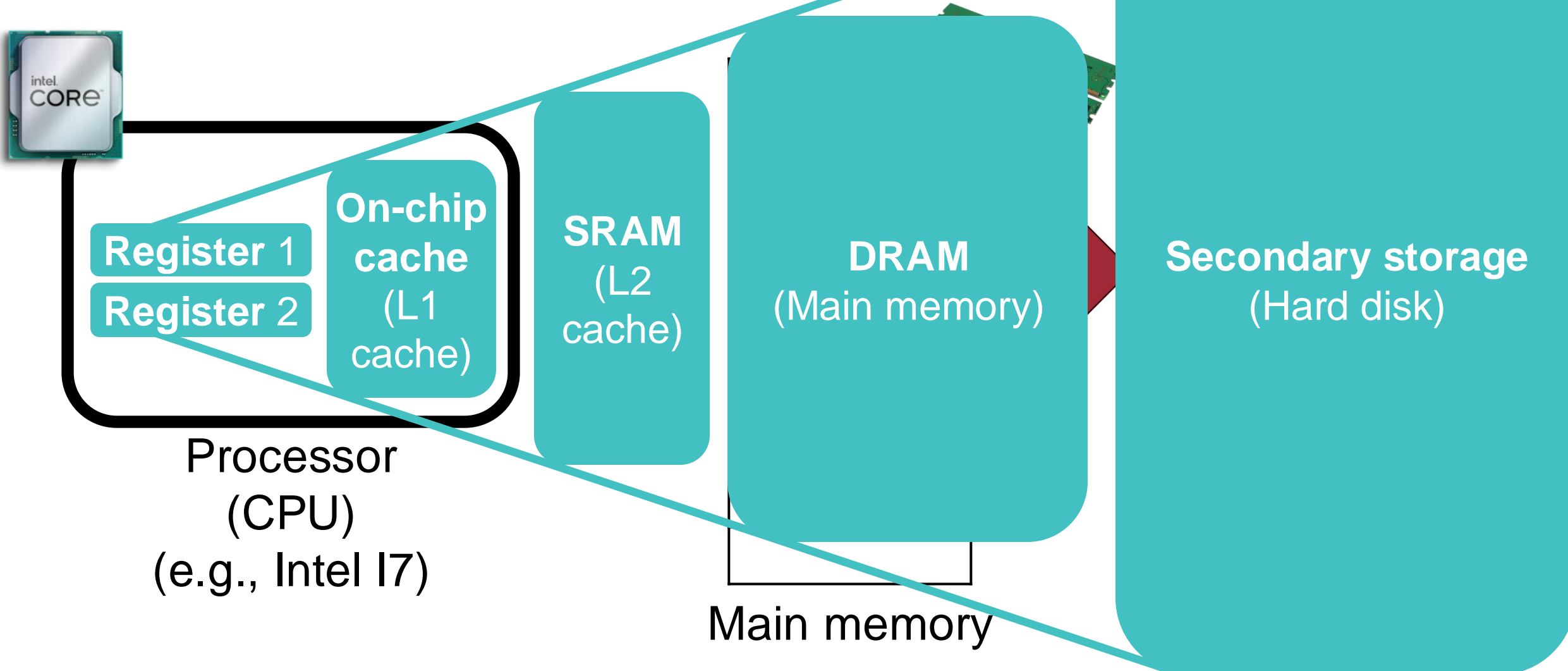
Recap: Terms

- Block (a.k.a., line): unit of copying
 - Several words in cache memory
- **Hit**: data requested is in the upper level
 - Hit ratio: $\text{hits}/\text{accesses}$
- **Miss**: data requested is not in the upper level
 - Block copied from lower level
 - Miss penalty: time taken to resolve miss
 - Miss ratio: $\text{misses}/\text{accesses}$
 $= 1 - \text{hit ratio}$



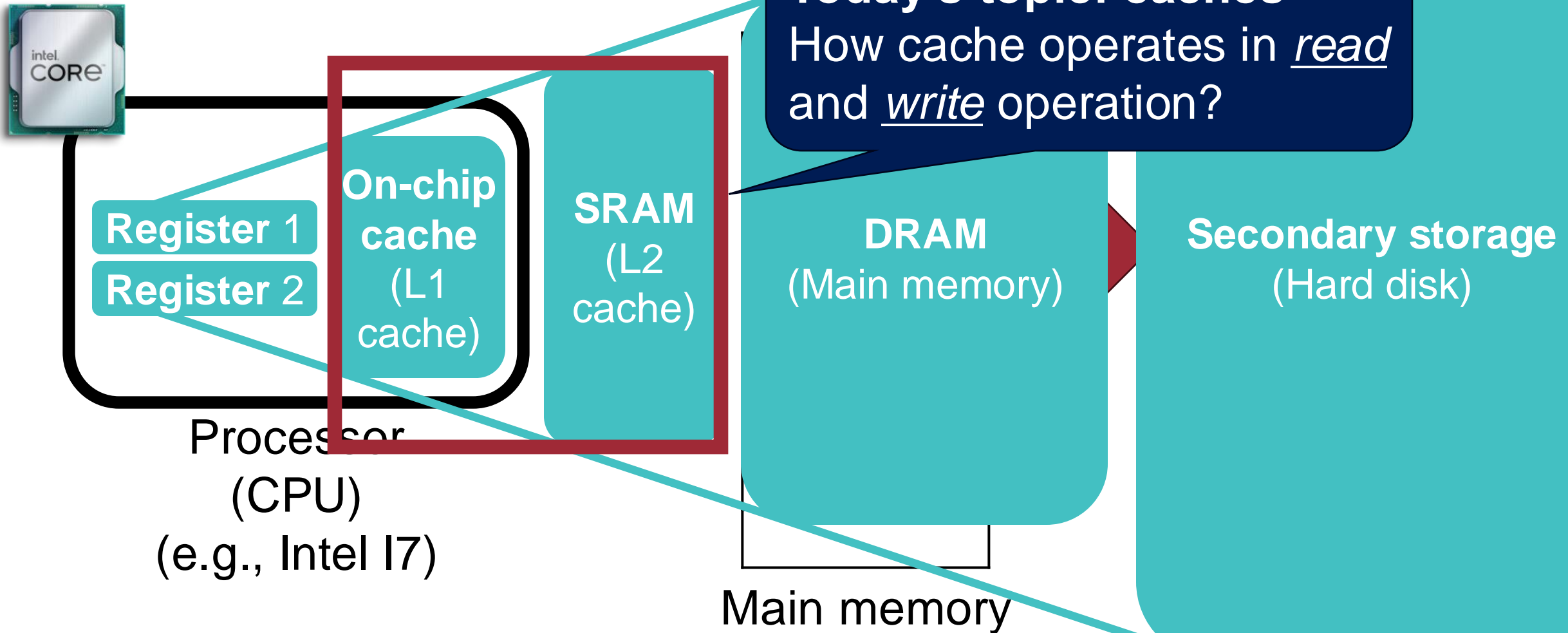
Today's Topic: Caches

9



Today's Topic: Caches

10

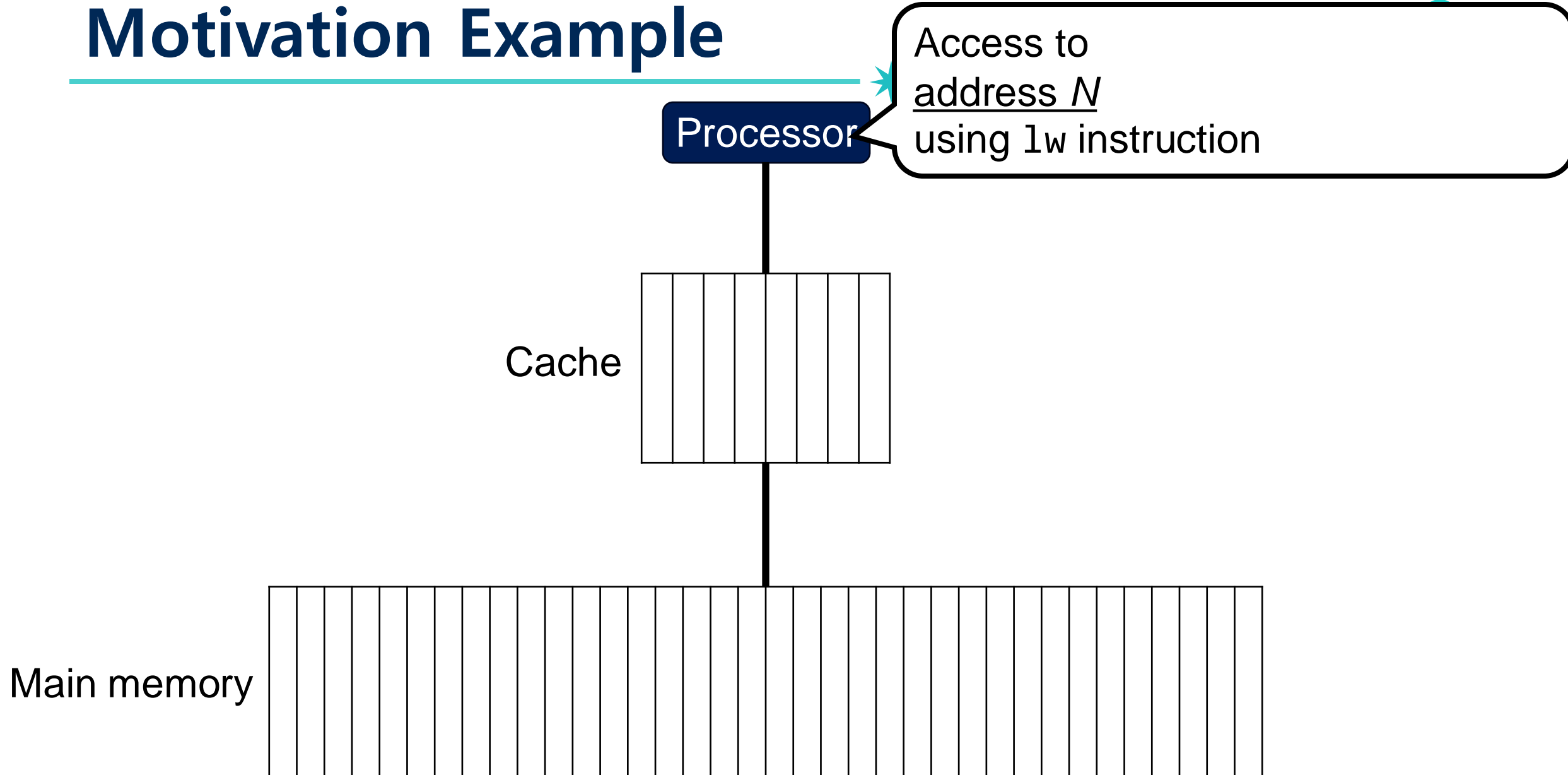


Caches

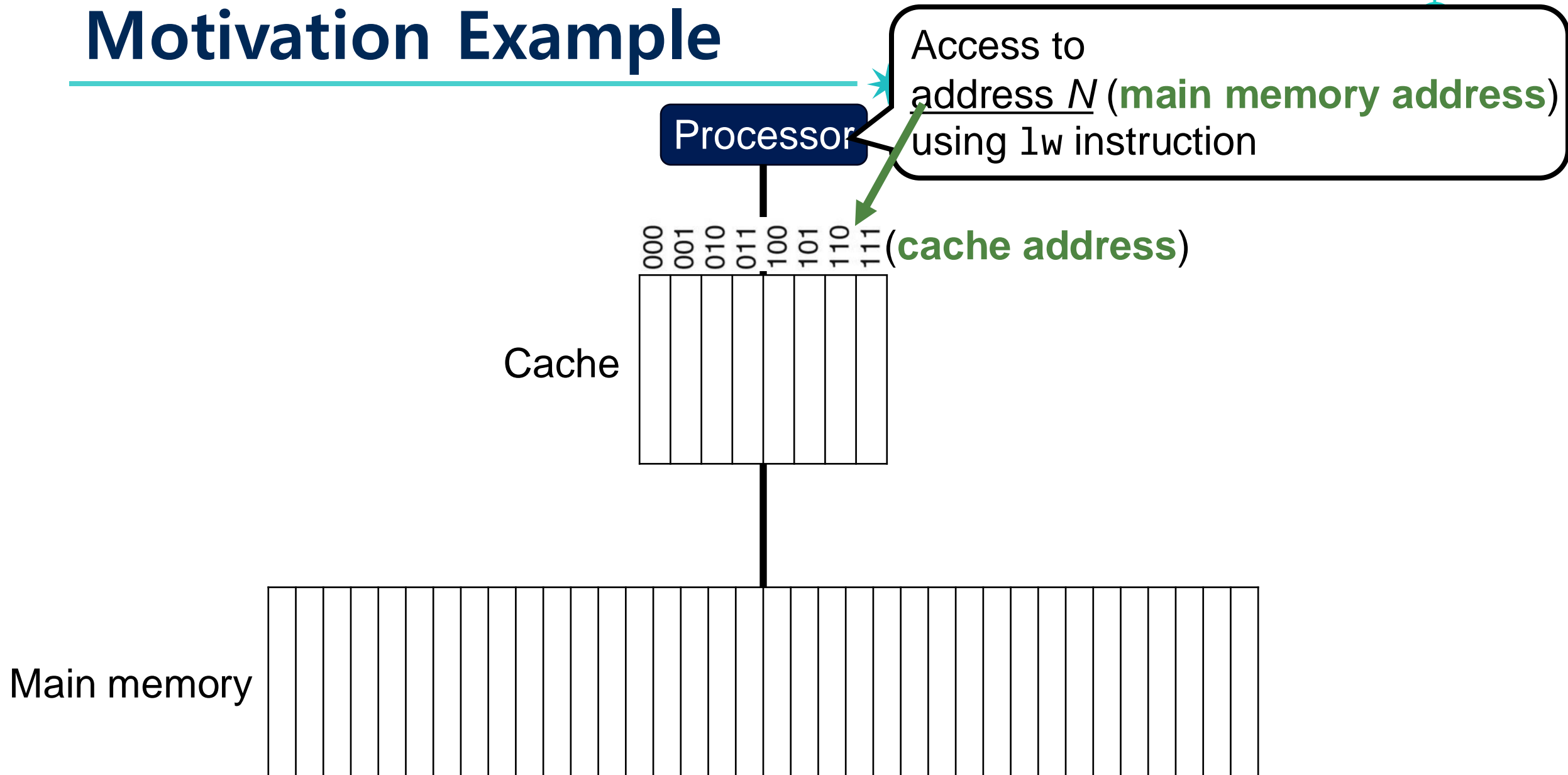


The level of the memory hierarchy closest to the CPU

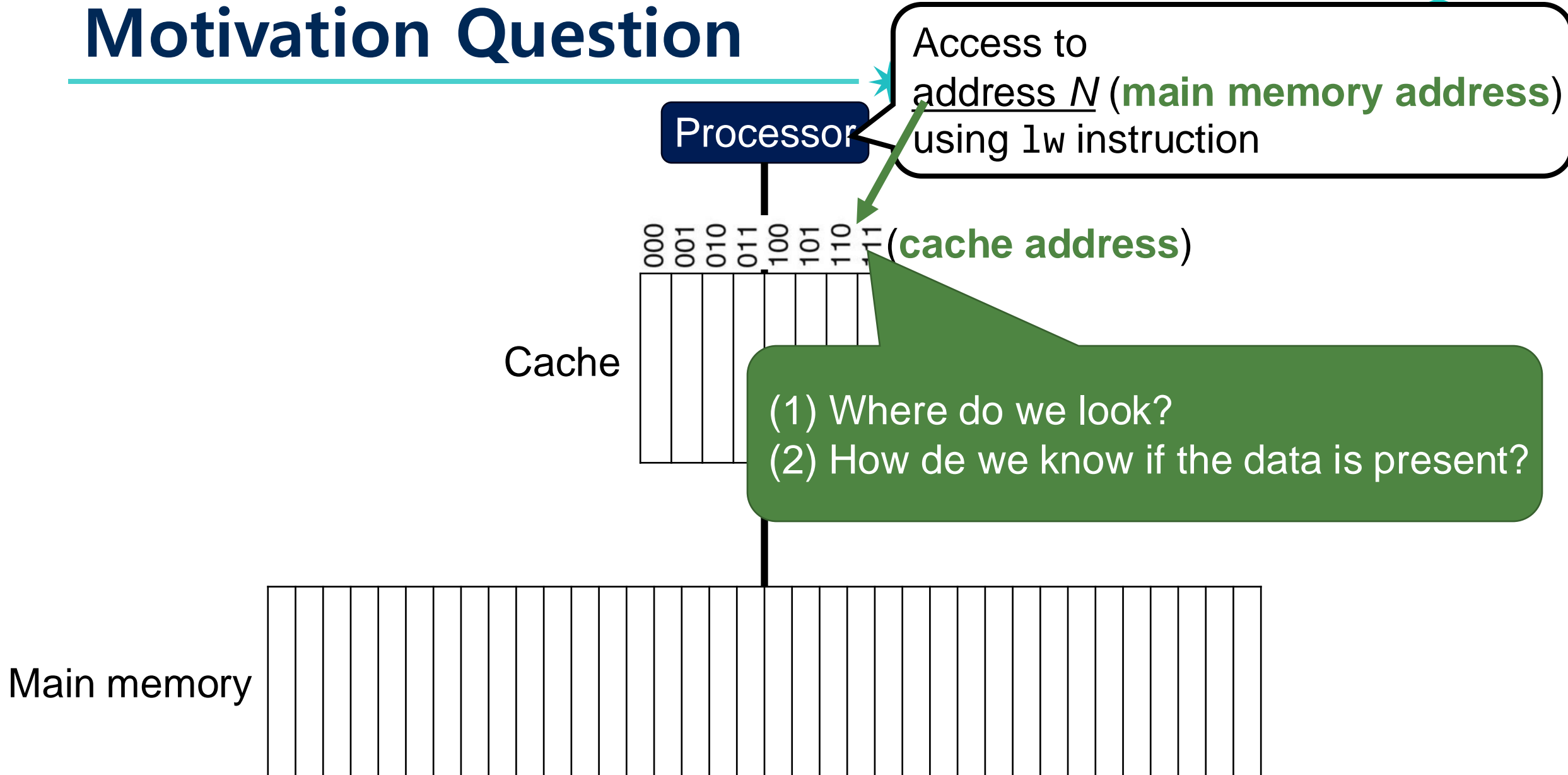
Motivation Example



Motivation Example



Motivation Question



Direct Mapped Cache

Direct Mapped Cache

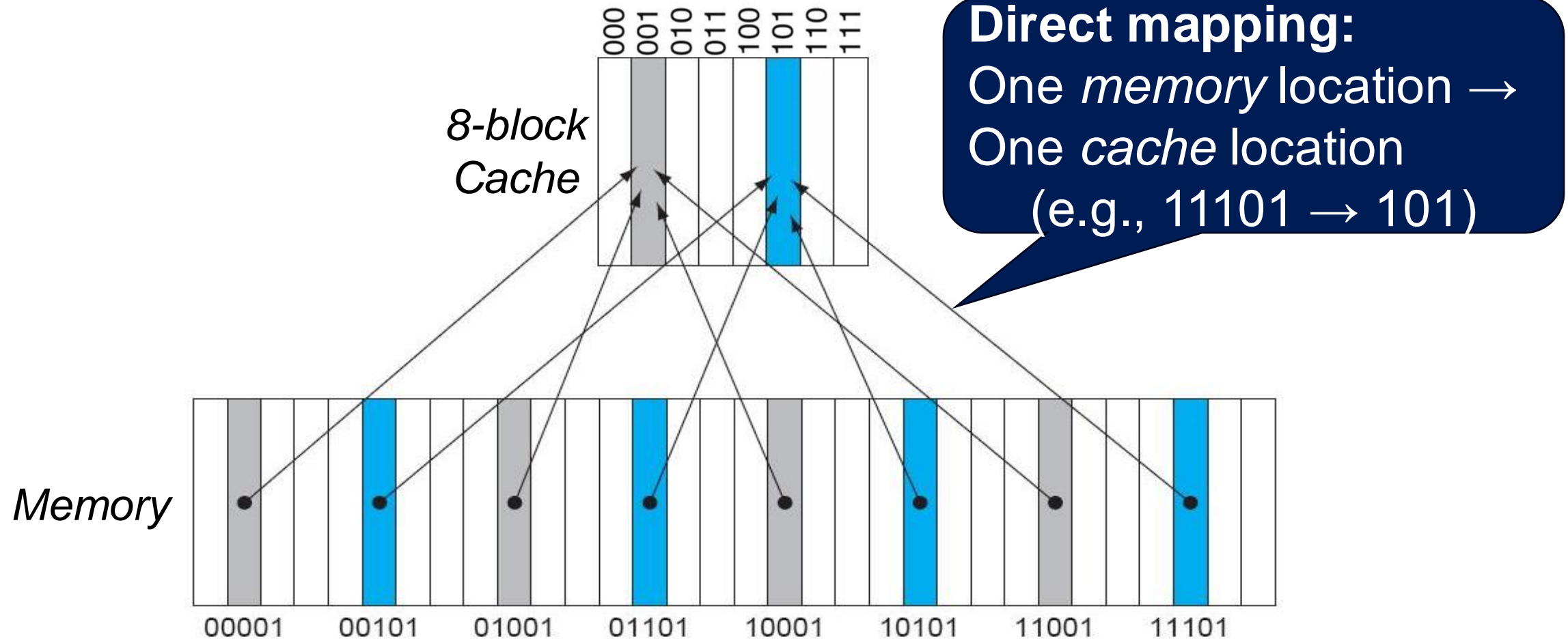


Each *memory* location is mapped to exactly one location in the *cache*

Direct Mapped Cache



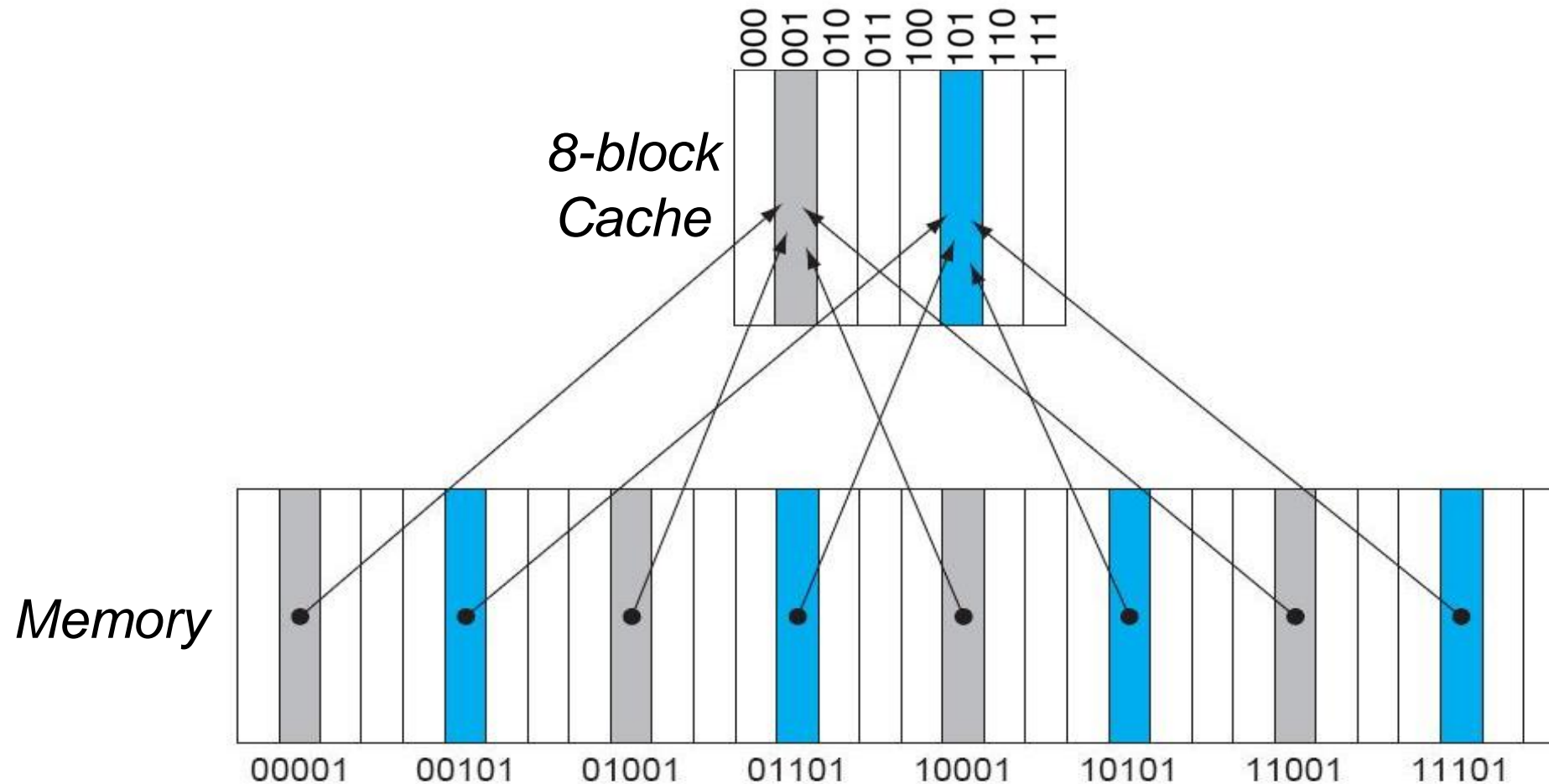
Each *memory* location is mapped to exactly one location in the *cache*



Mapping Algorithm



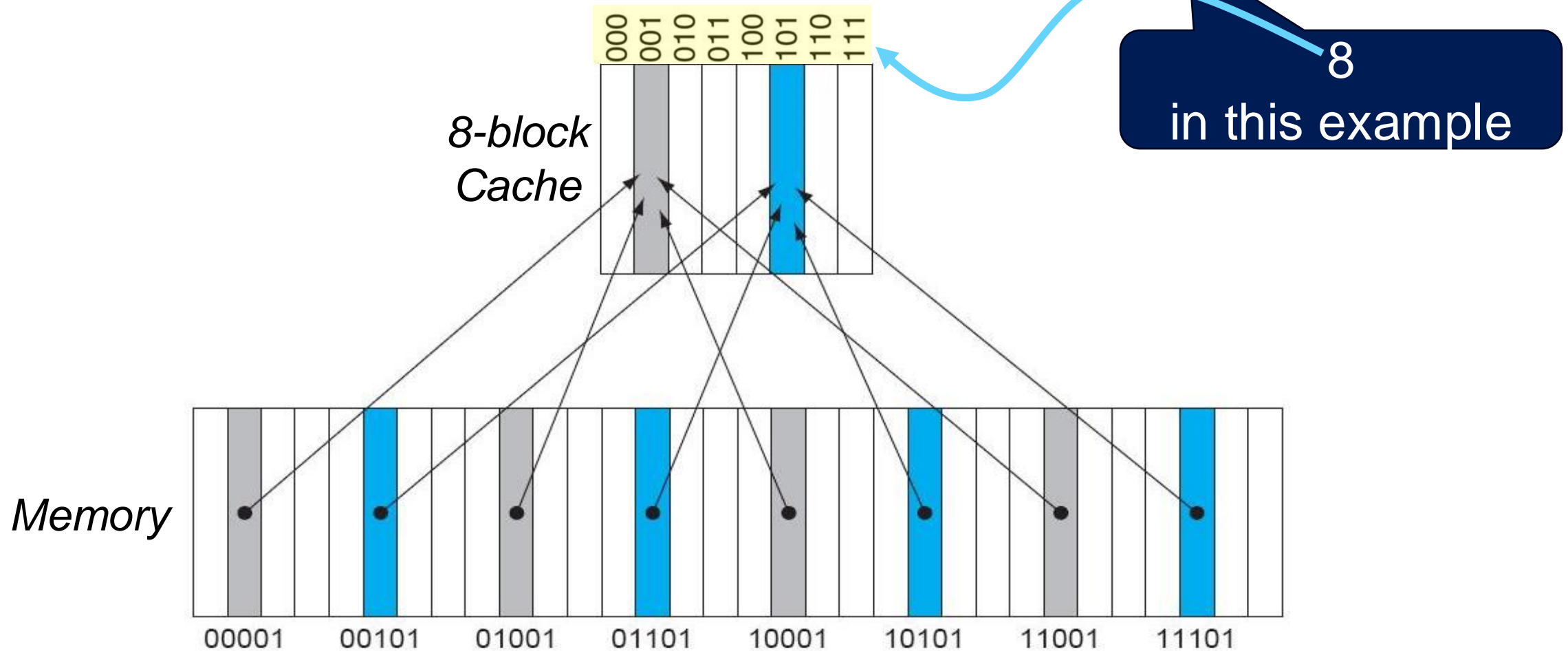
Location in the cache = (**Block address**) modulo (**# Blocks in cache**)



Mapping Algorithm



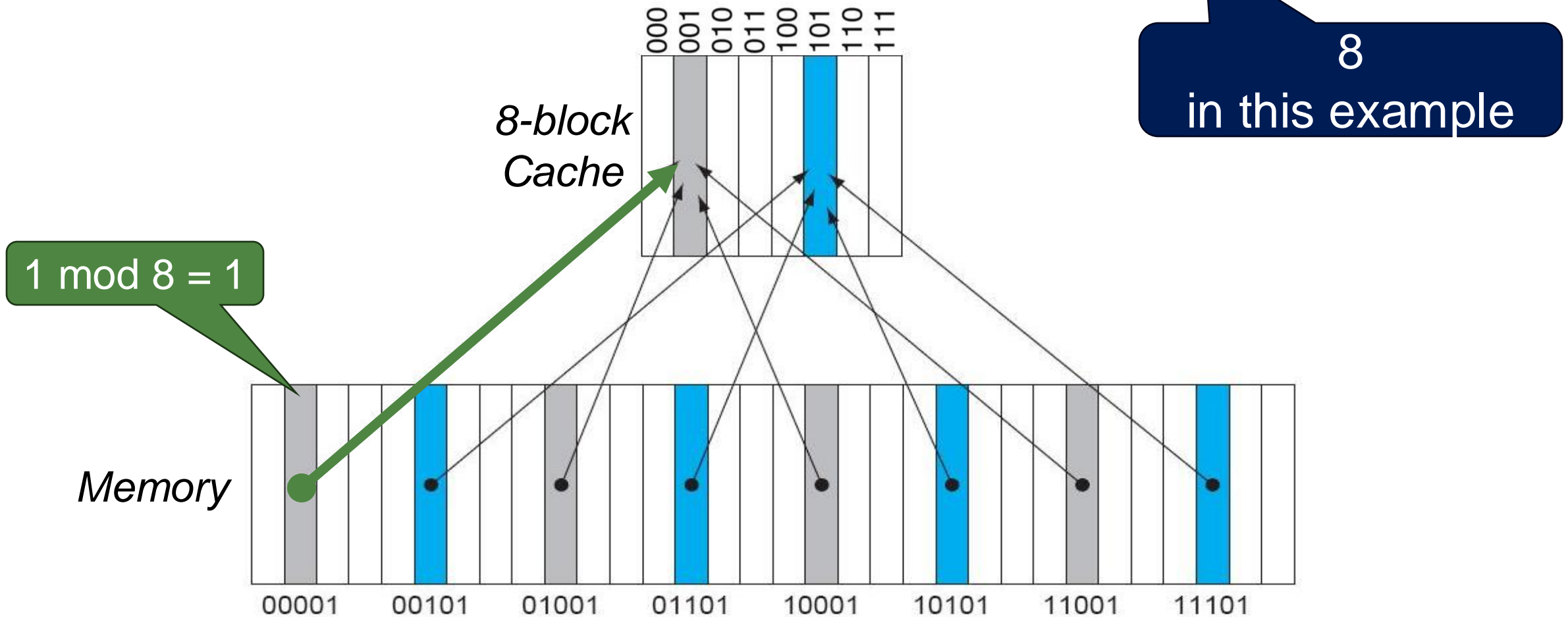
Location in the cache = (**Block address**) modulo (**# Blocks in cache**)



Mapping Algorithm



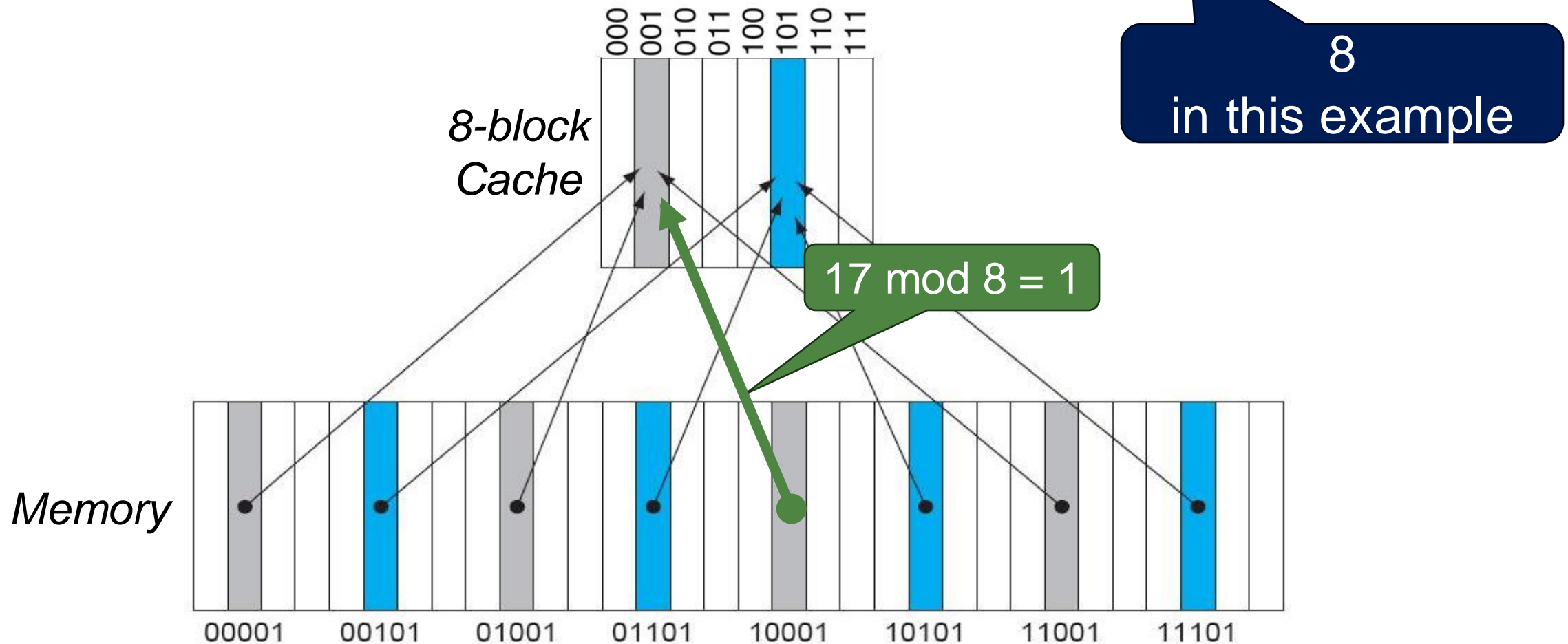
Location in the cache = (**Block address**) modulo (**# Blocks in cache**)



Mapping Algorithm

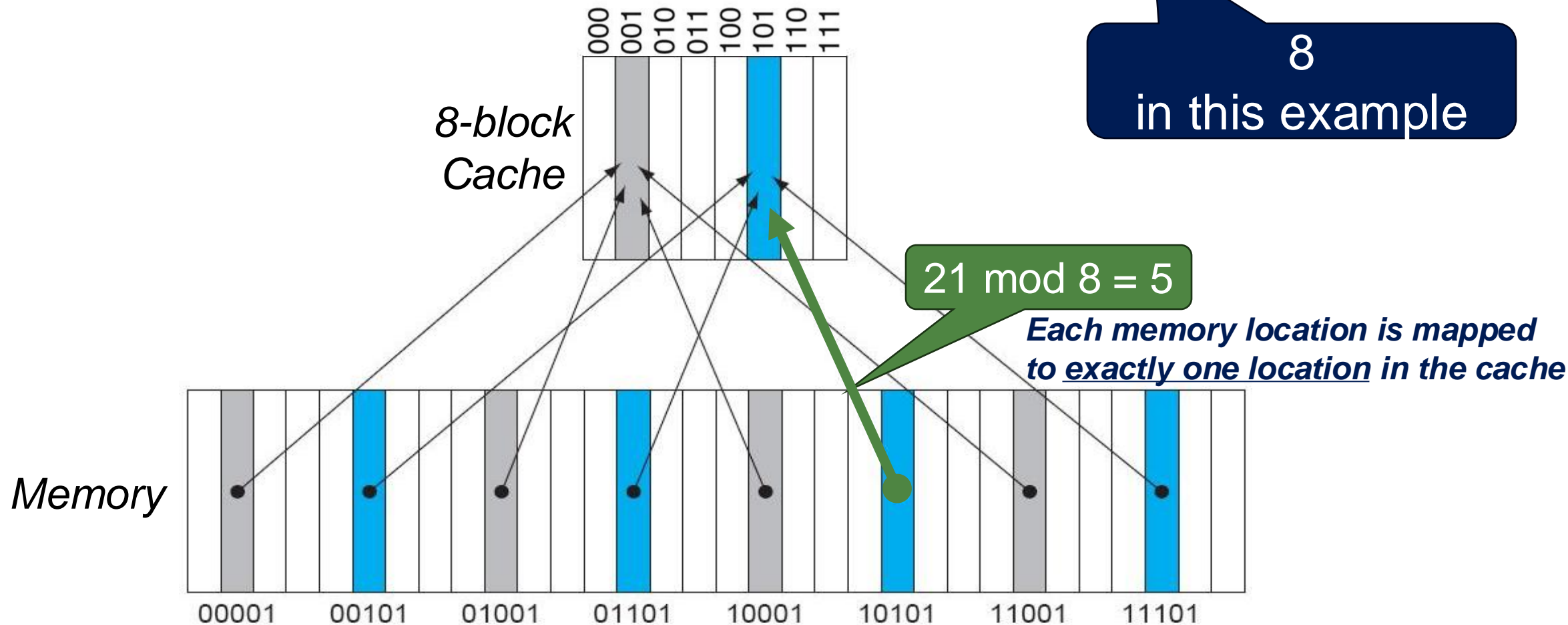


Location in the cache = (**Block address**) modulo (**# Blocks in cache**)



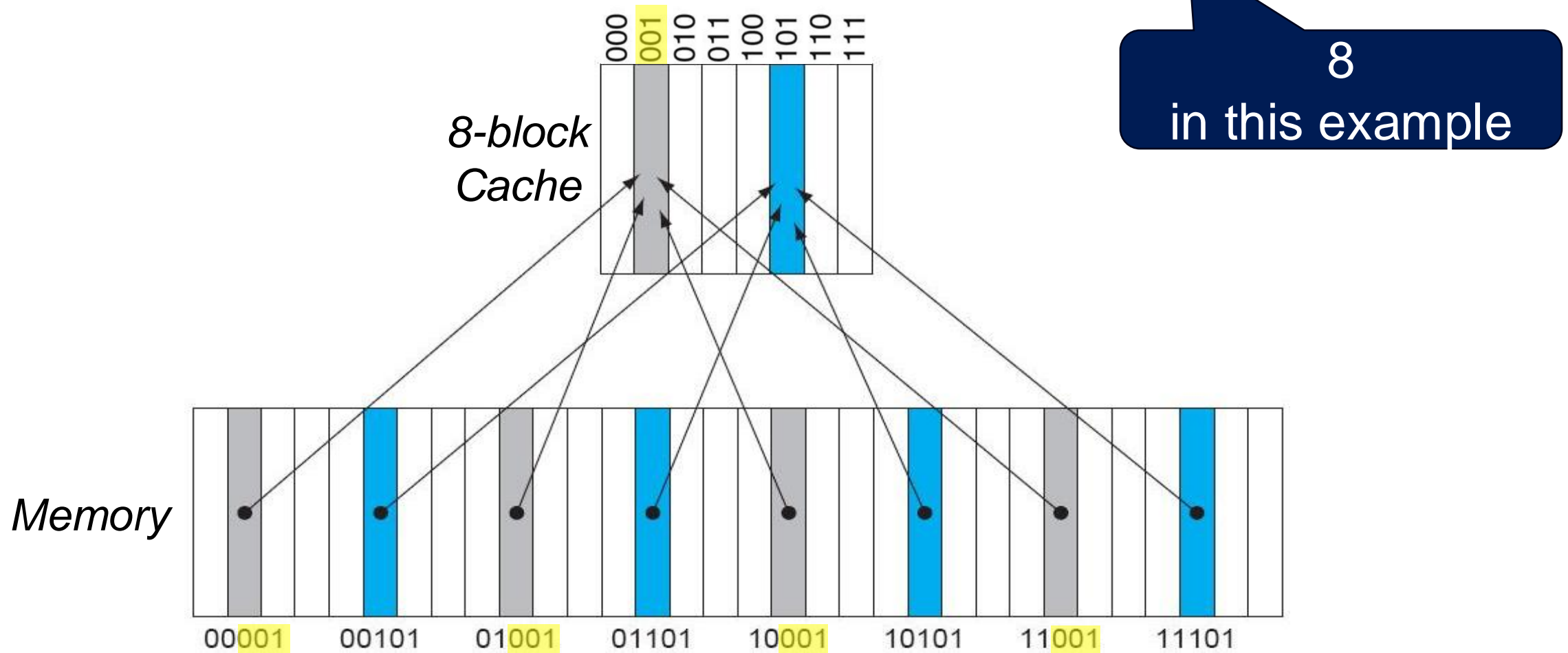
Mapping Algorithm

Location in the cache = (**Block address**) modulo (**# Blocks in cache**)



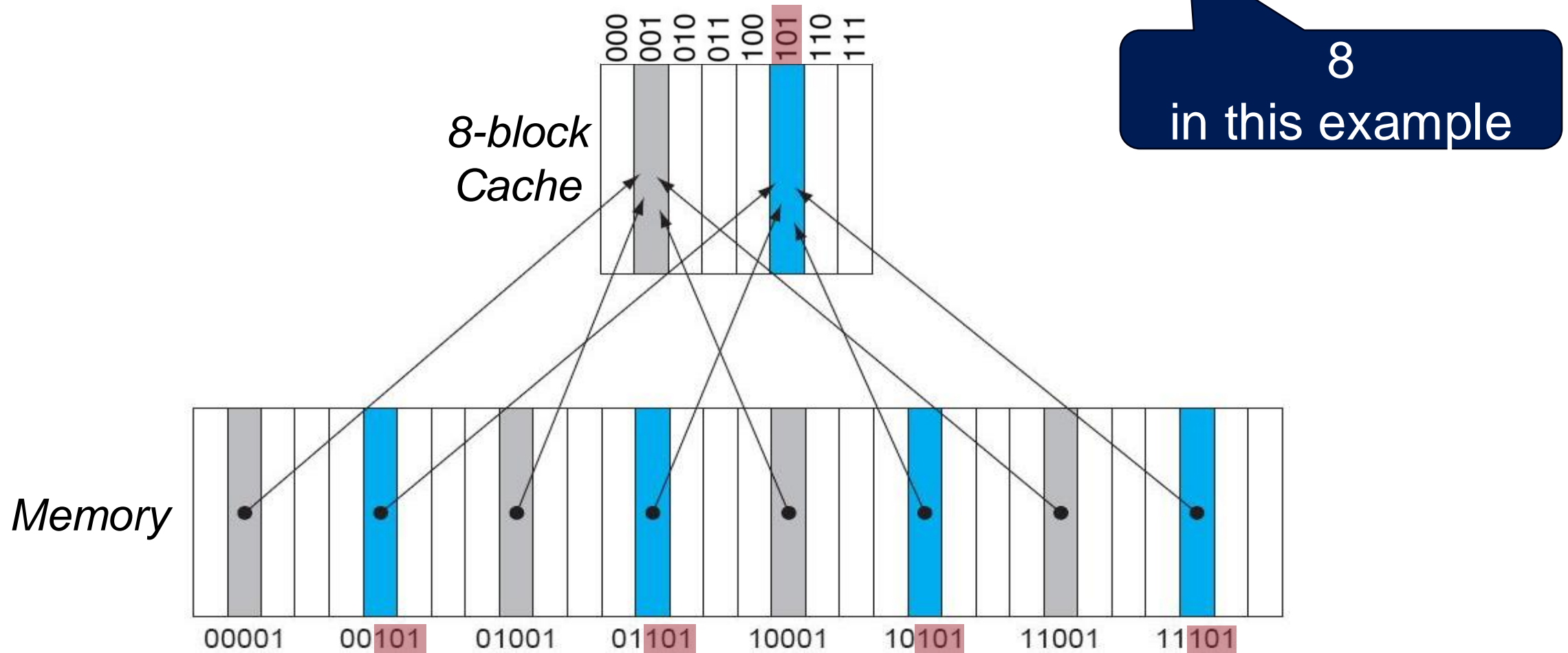
Mapping Algorithm: Binary View

Location in the cache = (**Block address**) modulo (**# Blocks in cache**)

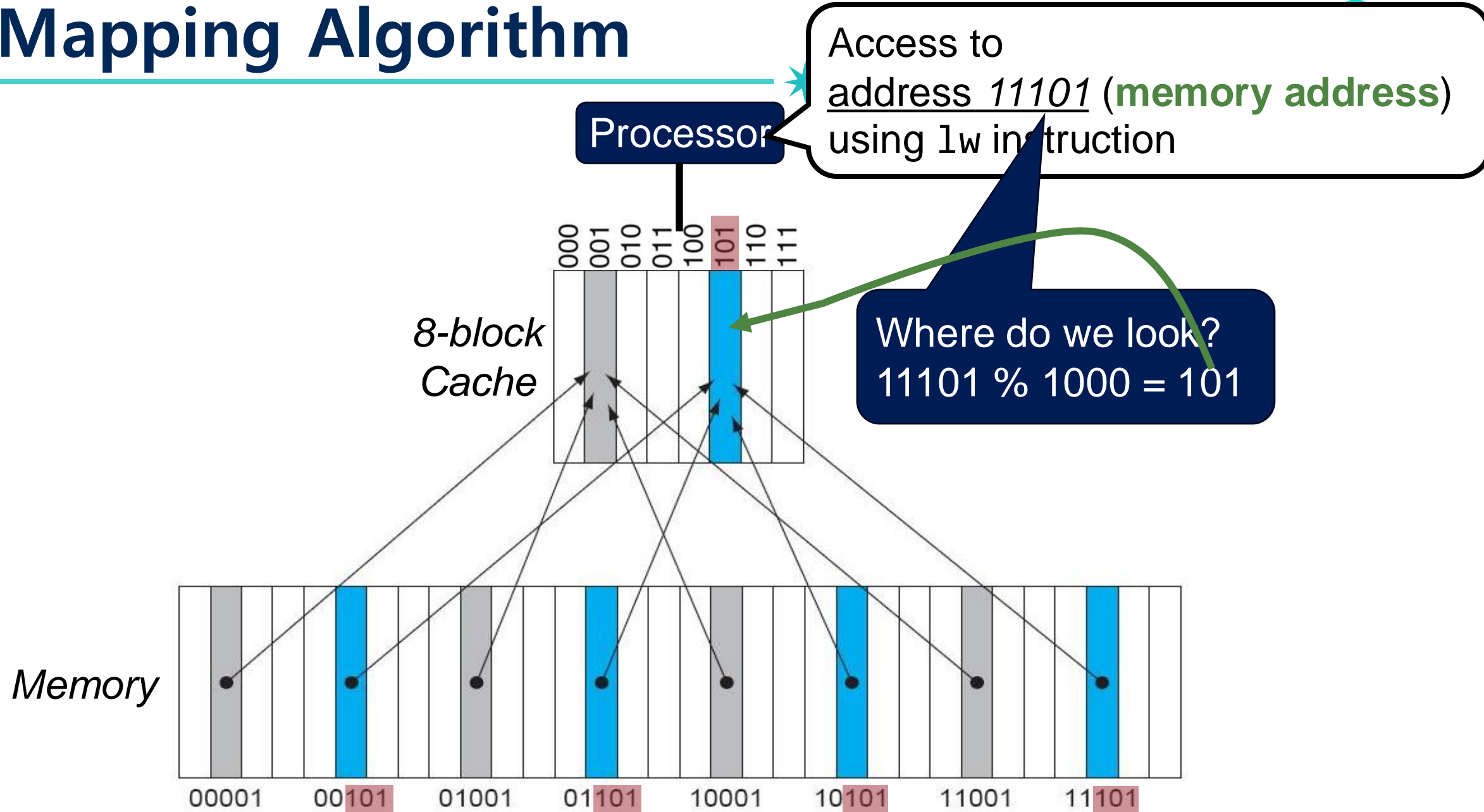


Mapping Algorithm: Binary View

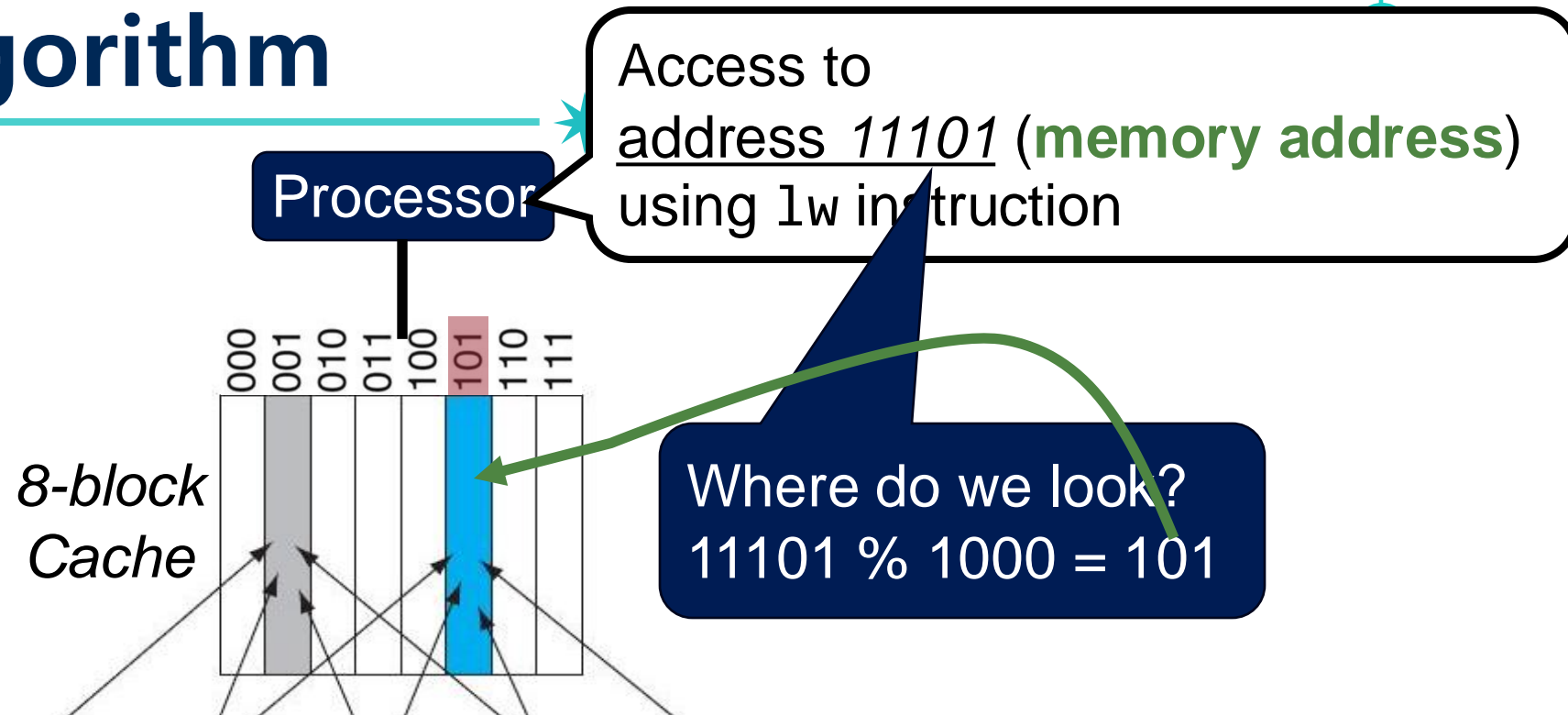
Location in the cache = (**Block address**) modulo (**# Blocks in cache**)



Mapping Algorithm



Mapping Algorithm



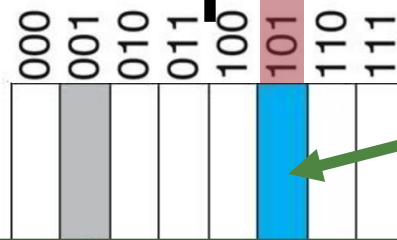
The value might not be from 11101
but from another location (e.g., 00101).
How can you confirm that it's the value from 11101?

Mapping Algorithm

Processor

Access to
address 11101 (**memory address**)
using lw instruction

8-block



Where do we look?
 $11101 \% 1000 = 101$

Idea: store the *high-order bits* in the
cache to specify the exact location

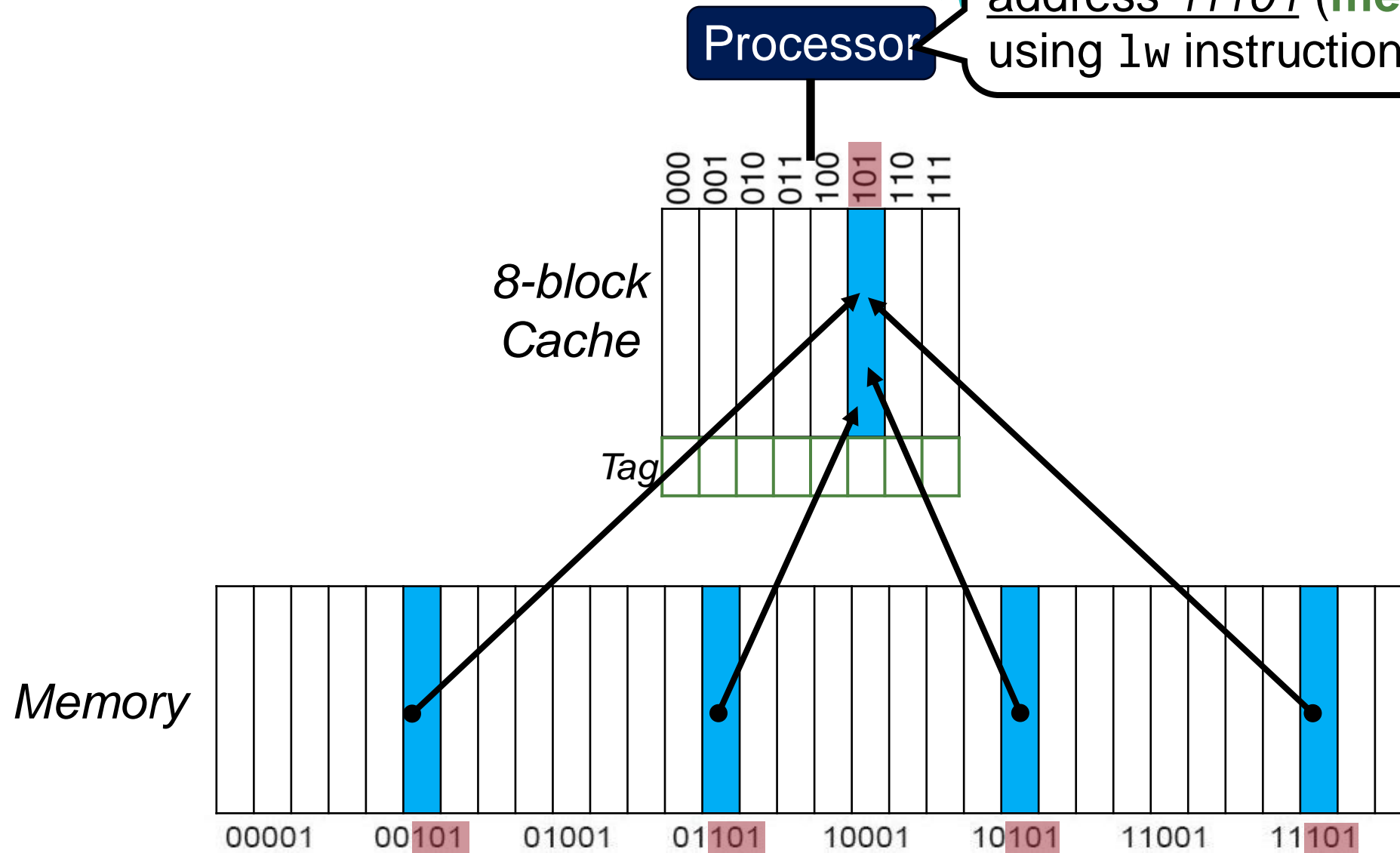


The value might not be from **11101**
but from another location (e.g., **00101**).

How can you confirm that it's the value from 11101?

Additional Field: Tag Bits

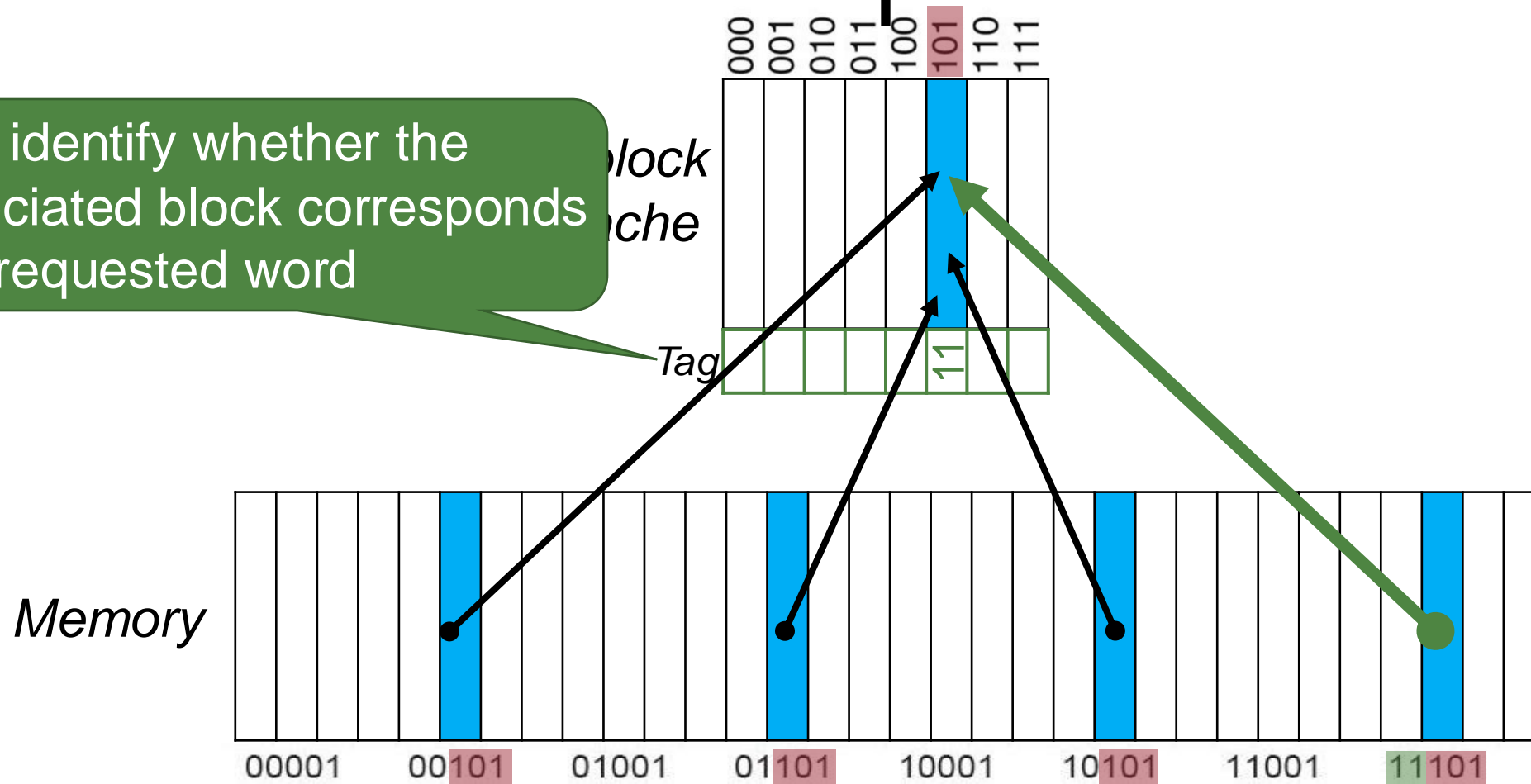
Access to
address 11101 (**memory address**)
using lw instruction



Additional Field: Tag Bits

Access to
address 11101 (**memory address**)
using lw instruction

Tag: identify whether the associated block corresponds to a requested word



Additional Field: Tag Bits

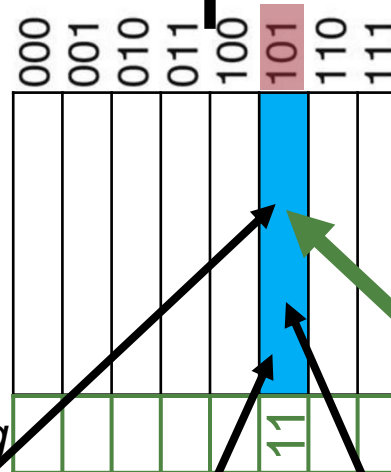
Access to
address 11101 (**memory address**)
using lw instruction

Processor

Tag: identify whether the
associated block corresponds
to a requested word

Block
Cache

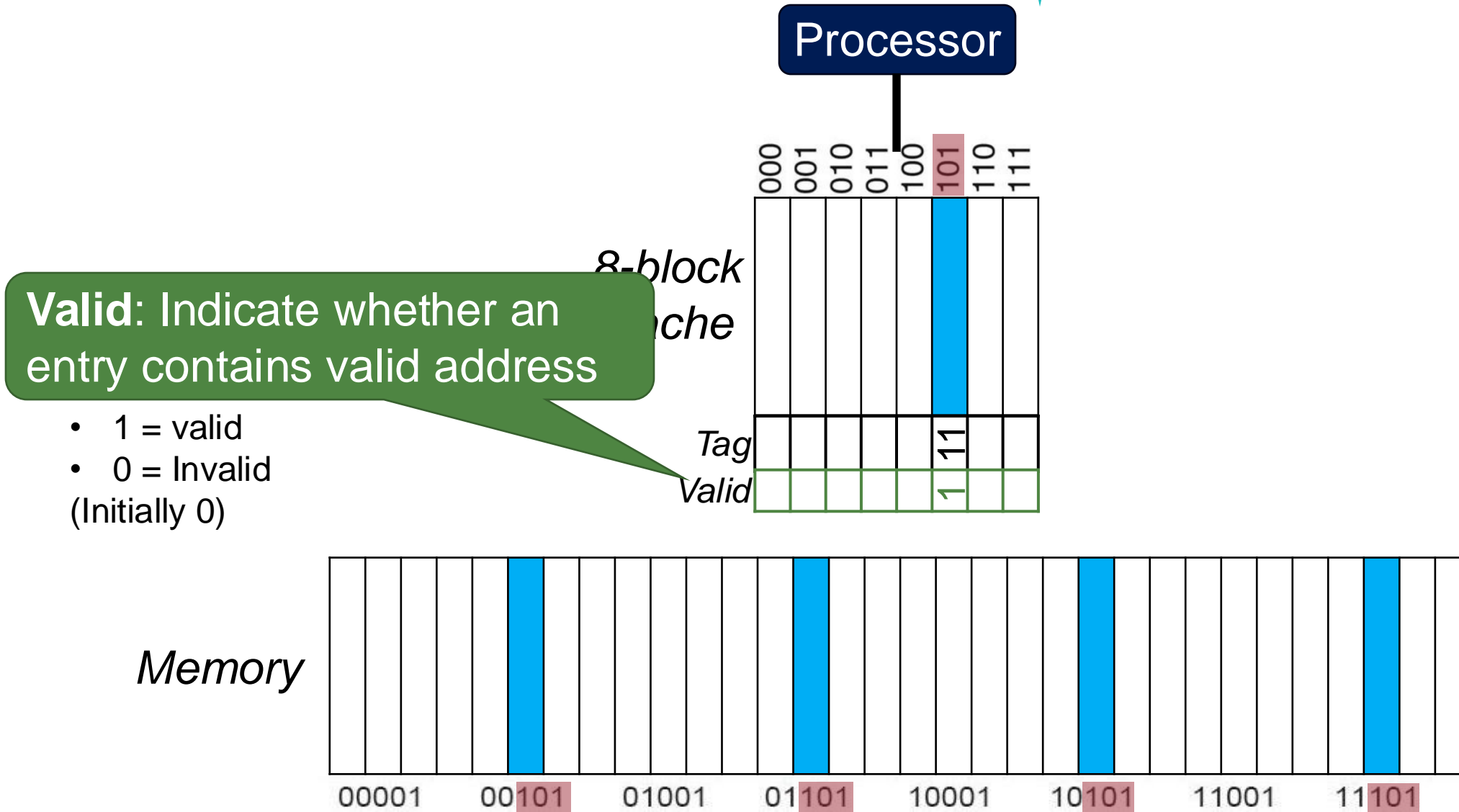
Tag



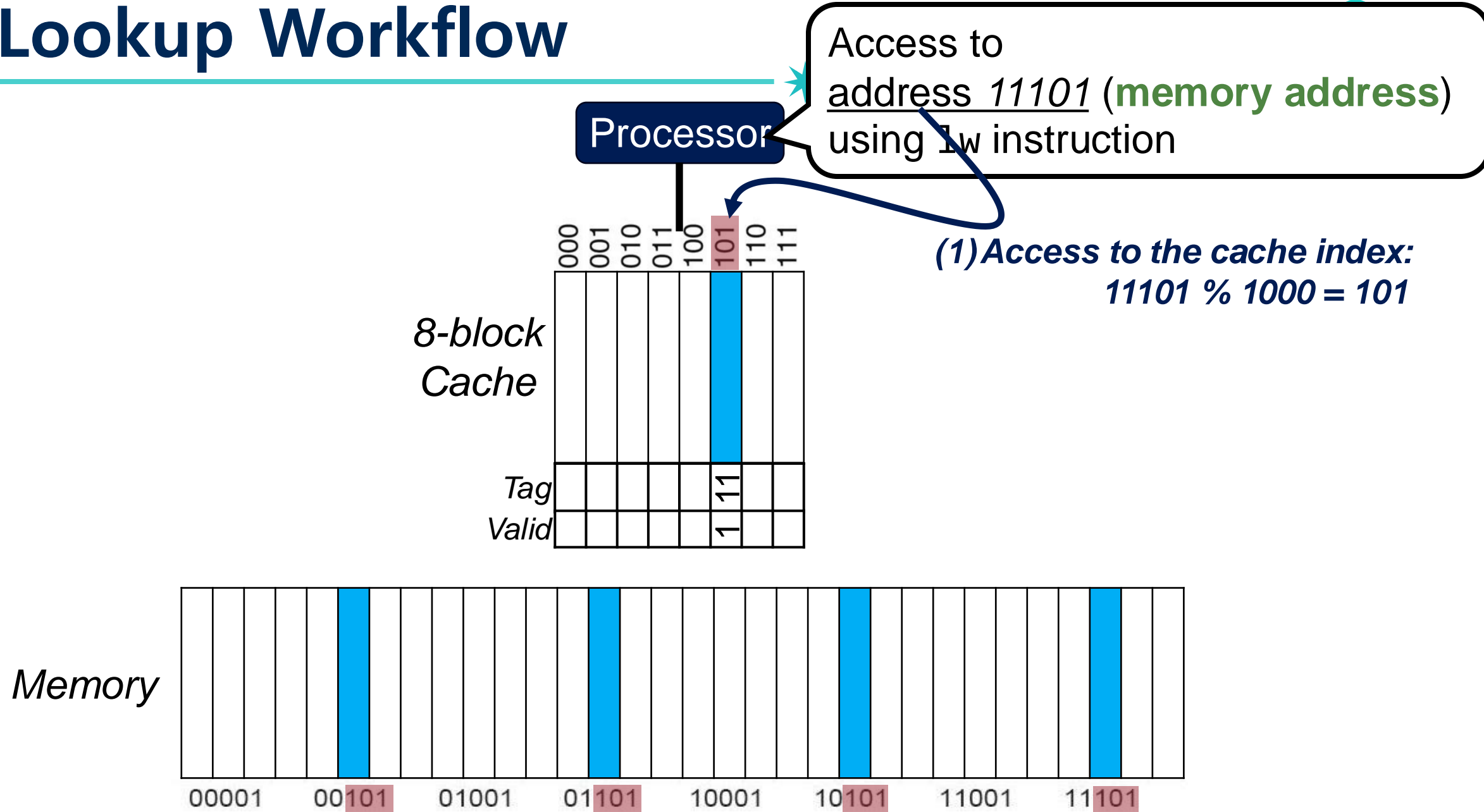
How do you know if
the data being accessed is *valid or not*?

(When a processor starts up, the cache does not have valid data)

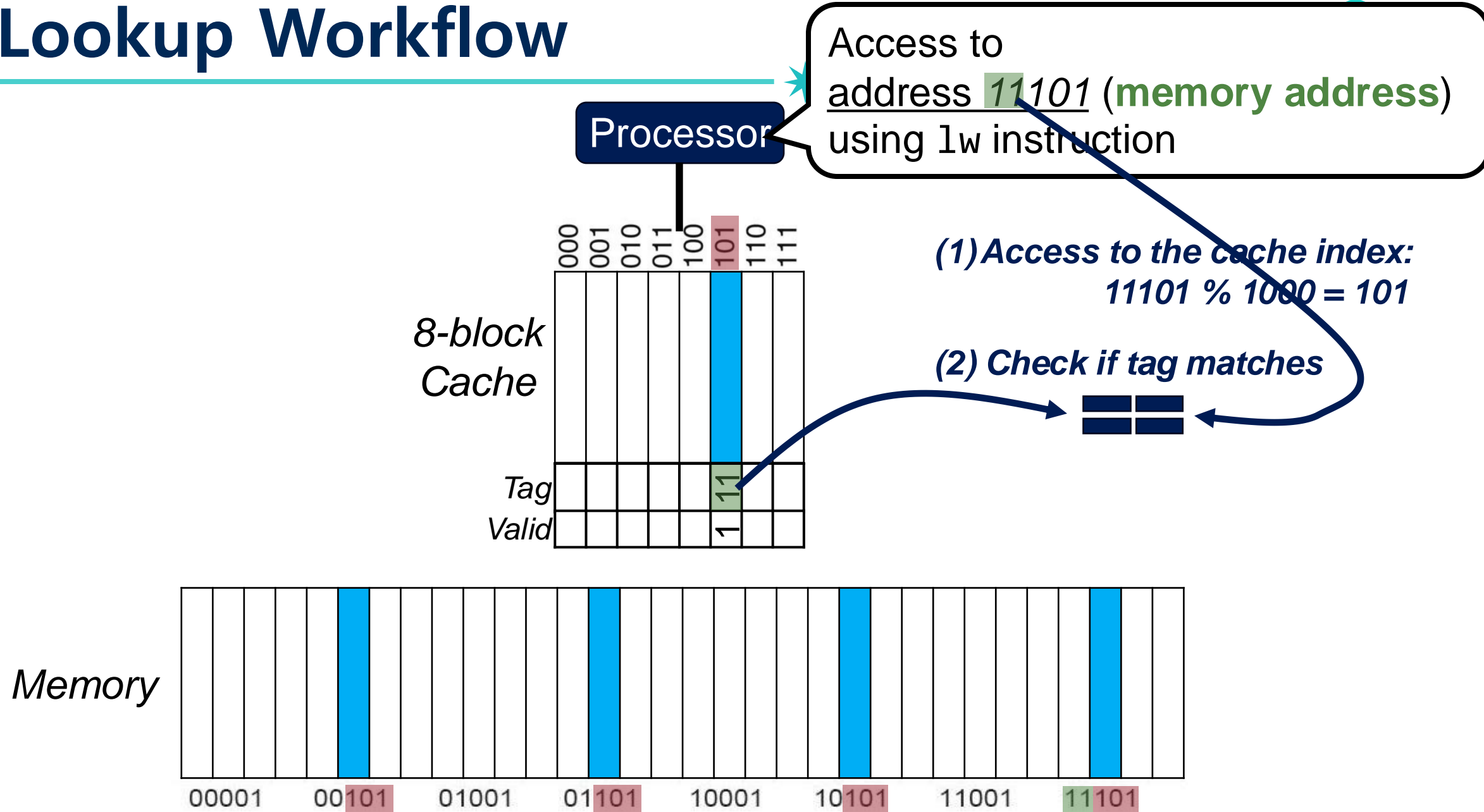
Additional Field: Valid Bit



Lookup Workflow



Lookup Workflow



Lookup Workflow

Processor

Access to
address 11101 (**memory address**)
using lw instruction

8-block Cache

[illegible]

(1) Access to the cache index:
 $11101 \% 1000 = 101$

(2) Check if tag matches



(3) Check if valid bit sets

Memory

00001 00101 01001 01101 10001 10101 11001 11101

Lookup Workflow

(4) If it's a hit: access the **cache data** (no need to access the memory)

Processor

Access to address **11101** (**memory address**) using lw instruction

8 block
Cache

	000	001	010	011	100	101	110	111
Tag						11		
Valid					1			

(1) Access to the cache index:
 $11101 \% 1000 = 101$

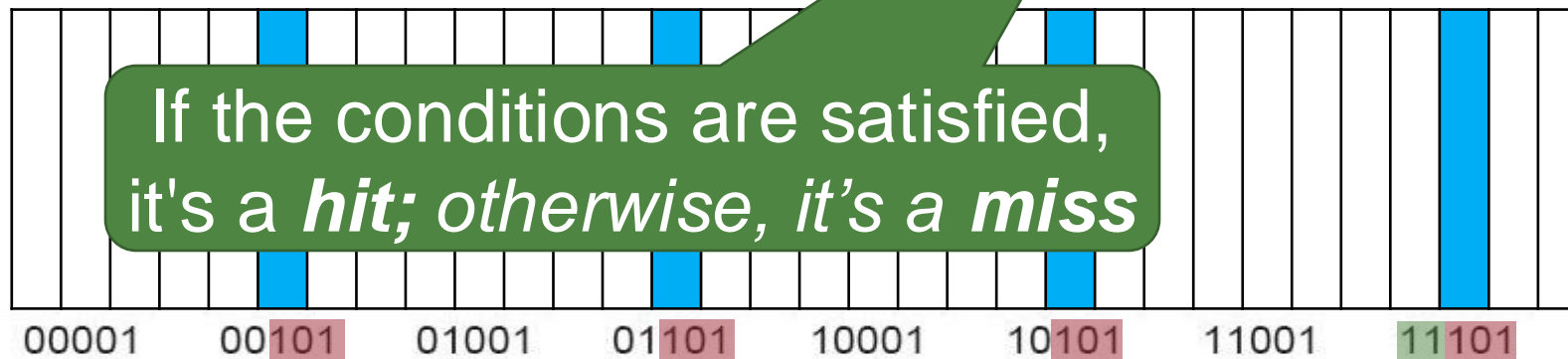
(2) Check if tag matches

==

(3) Check if valid bit sets

Memory

If the conditions are satisfied, it's a **hit**; otherwise, it's a **miss**



Direct Mapped Cache Example



The initial state of the cache after power-on

Memory access sequence:

(1) 10110

(2) 11010

(3) 10000

(4) 00011

(5) 10010

(6) 10110

Index	Valid	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

Direct Mapped Cache Example

a. The initial state of the cache after power-on

Memory access sequence:

Miss (1) 10110

(2) 11010

(3) 10000

(4) 00011

(5) 10010

(6) 10110

Index	Valid	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

Direct Mapped Cache Example

Memory access sequence:

Miss (1) 10110

(2) 11010

(3) 10000

(4) 00011

(5) 10010

(6) 10110

b. After handling a miss of address 10110

Index	Valid	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

Direct Mapped Cache Example

b. After handling a miss of address 10110

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

(3) 10000

(4) 00011

(5) 10010

(6) 10110

Index	Valid	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

Direct Mapped Cache Example

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

(3) 10000

(4) 00011

(5) 10010

(6) 10110

c. After handling a miss of address 11010

Index	Valid	Tag	Data
000	N		
001	N		
010	Y	11 _{two}	Memory(11010 _{two})
011	N		
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

Direct Mapped Cache Example

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

Miss (3) 10000

(4) 00011

(5) 10010

(6) 10110

d. After handling a miss of address 10000

Index	Valid	Tag	Data
000	Y	10_{two}	Memory(10000_{two})
001	N		
010	Y	11_{two}	Memory(11010_{two})
011	N		
100	N		
101	N		
110	Y	10_{two}	Memory(10110_{two})
111	N		

Direct Mapped Cache Example



e. After handling a miss of address 00011

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

Miss (3) 10000

Miss (4) 00011

(5) 10010

(6) 10110

Index	Valid	Tag	Data
000	Y	10 _{two}	Memory(10000 _{two})
001	N		
010	Y	11 _{two}	Memory(11010 _{two})
011	Y	00 _{two}	Memory(00011 _{two})
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

Direct Mapped Cache Example

e. After handling a miss of address 00011

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

Miss (3) 10000

Miss (4) 00011

Miss (5) 10010

(6) 10110

Index	Valid	Tag	Data
000	Y	10 _{two}	Memory(10000 _{two})
001	N		
010	Y	10 ≠ 11 _{two}	Memory(11010 _{two})
011	Y	00 _{two}	Memory(00011 _{two})
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

Direct Mapped Cache Example

f. After handling a miss of address 10010

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

Miss (3) 10000

Miss (4) 00011

Miss (5) 10010

(6) 10110

Index	Valid	Tag	Data
000	Y	10 _{two}	Memory(10000 _{two})
001	N		
010	Y	10 _{two}	Memory(10010 _{two})
011	Y	00 _{two}	Memory(00011 _{two})
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

Direct Mapped Cache Example

Memory access sequence:

Miss (1) 10110

Miss (2) 11010

Miss (3) 10000

Miss (4) 00011

Miss (5) 10010

Hit (6) 10110

Index	Valid	Tag	Data
000	Y	10 _{two}	Memory(10000 _{two})
001	N		
010	Y	10 _{two}	Memory(10010 _{two})
011	Y	00 _{two}	Memory(00011 _{two})
100	N		
101	N		
110	Y	10 _{two}	Memory(10110 _{two})
111	N		

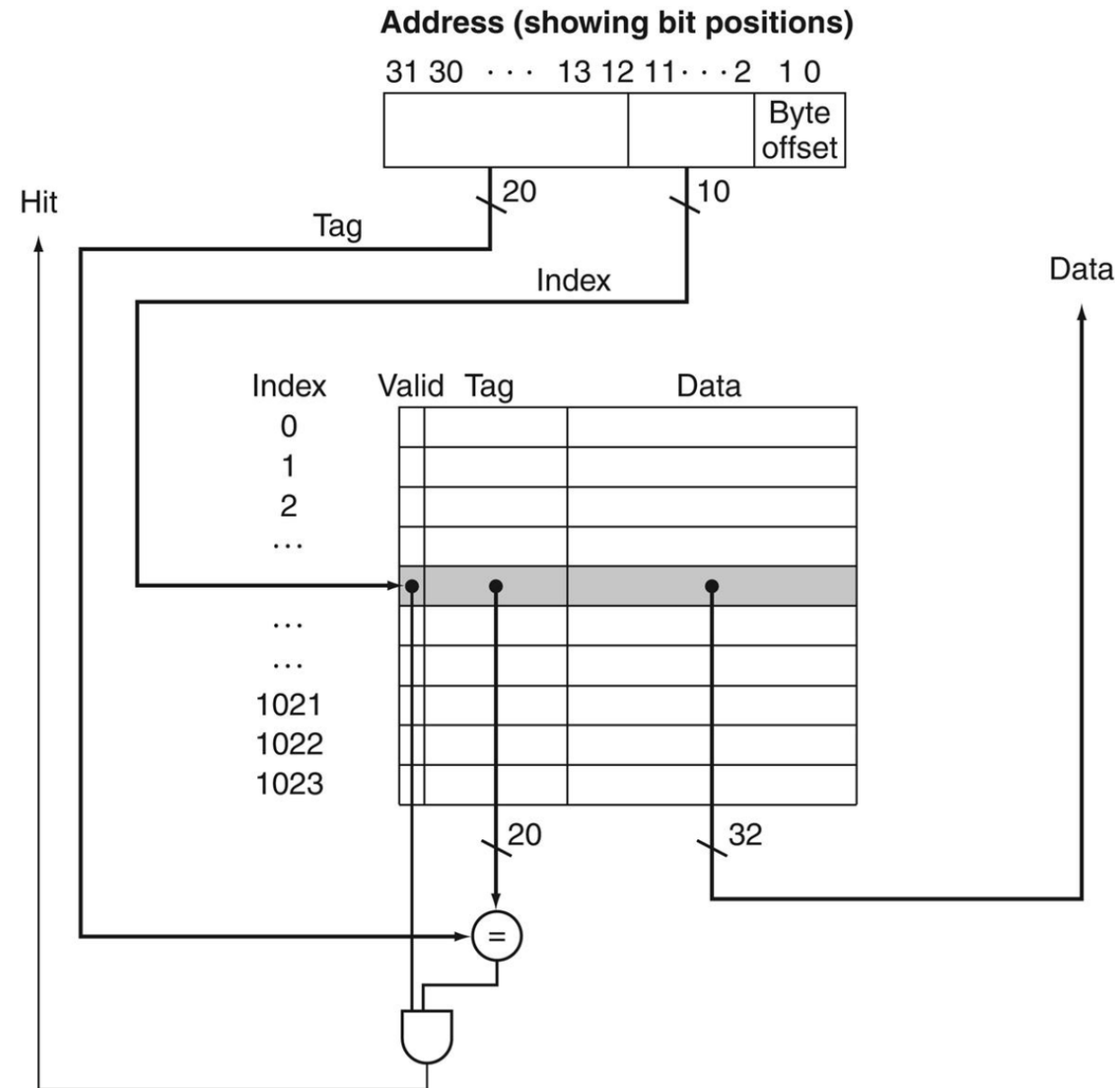
Direct Mapped Cache Structure

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



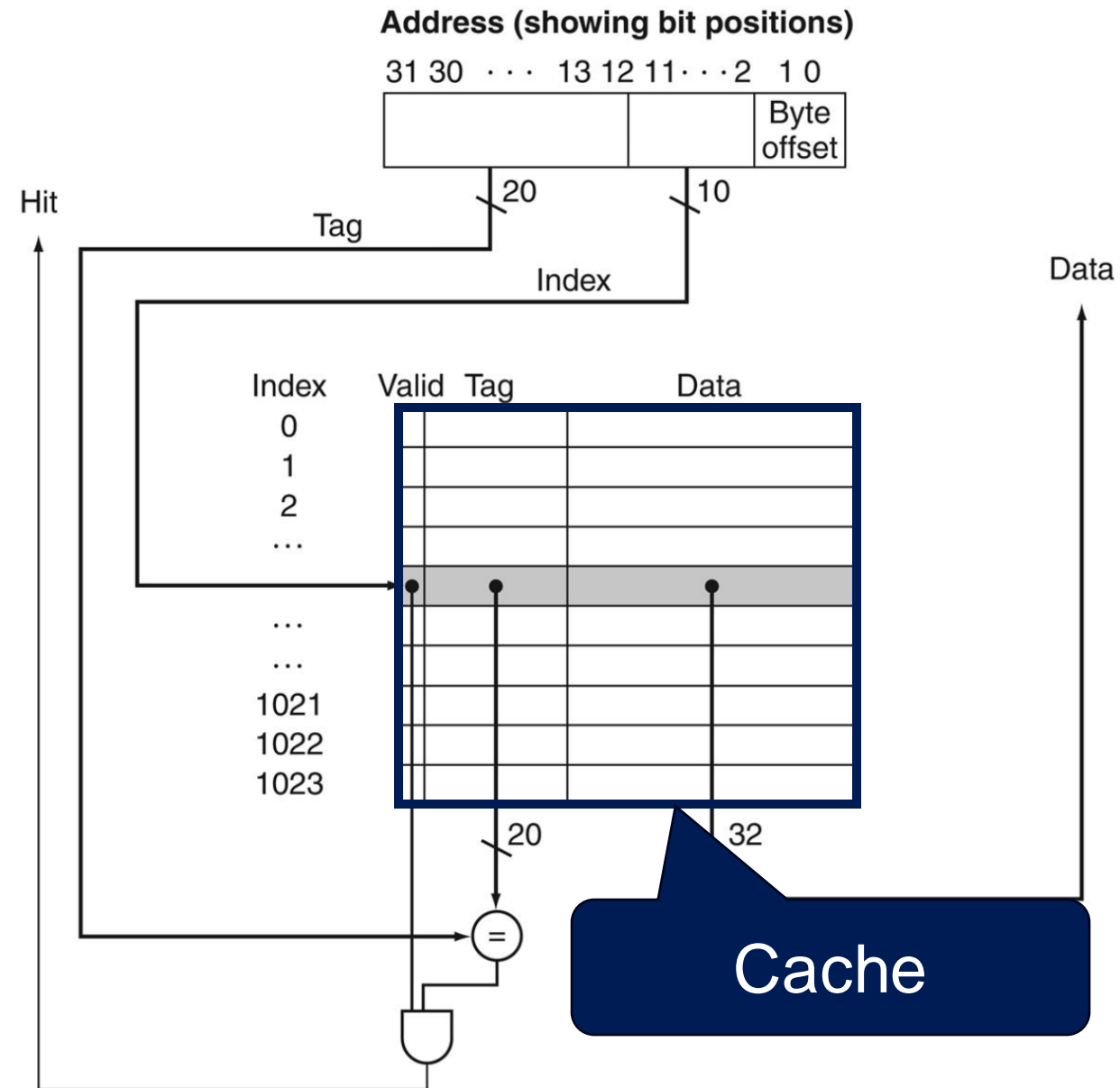
Direct Mapped Cache Structure

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



Direct Mapped Cache Structure

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?

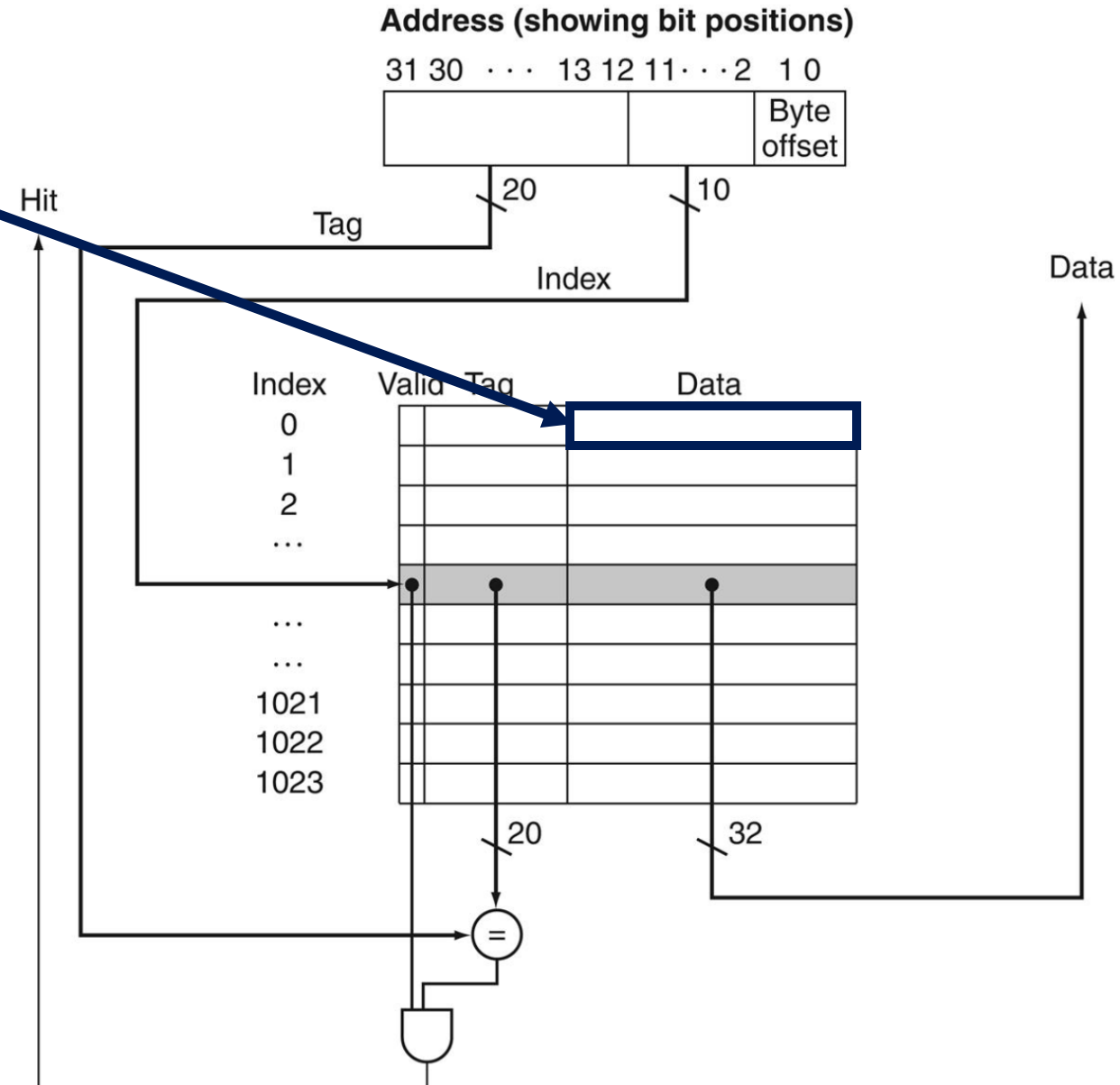


Direct Mapped Cache Structure

- **Assumption:**

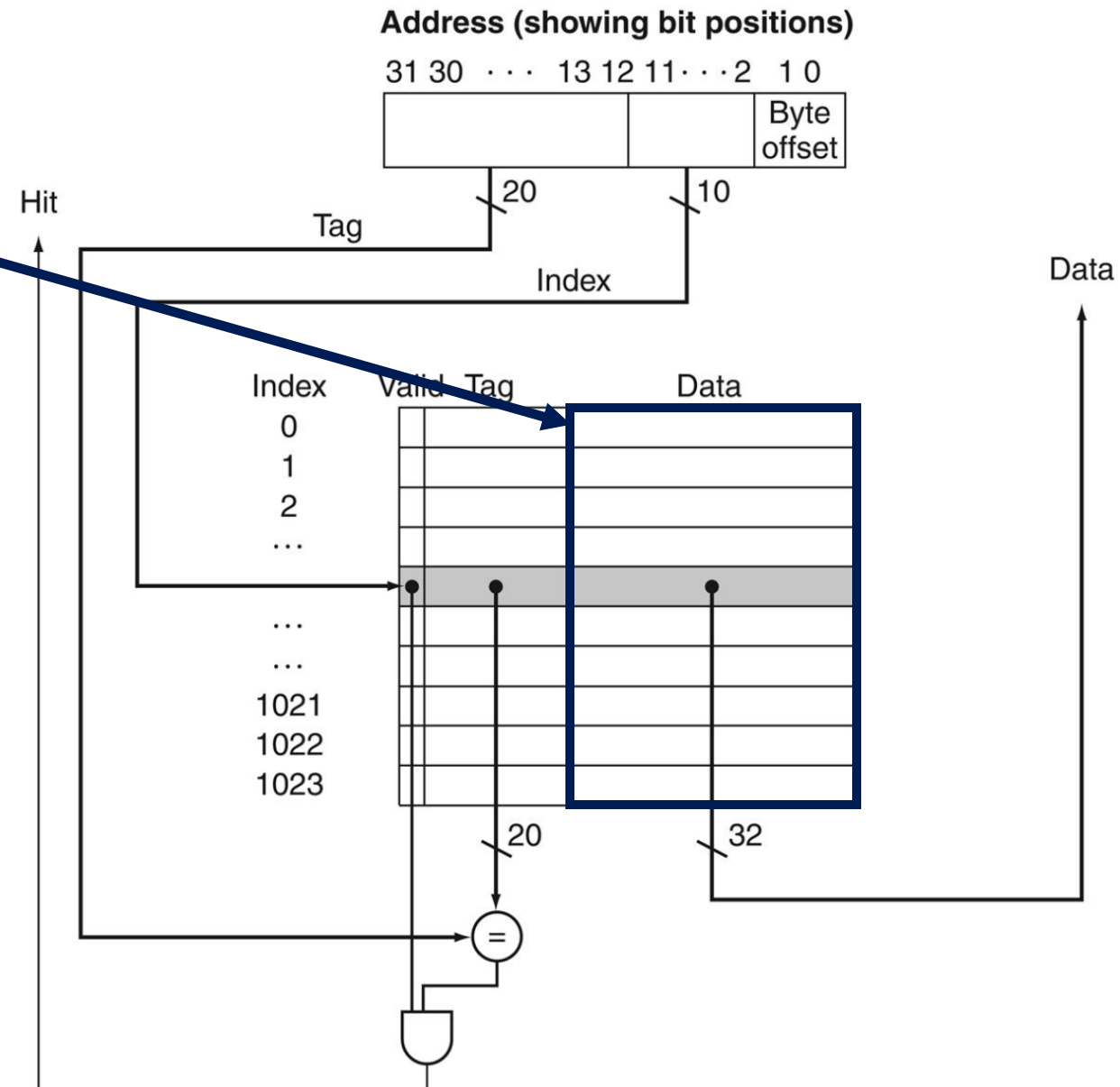
- Each block is 1 word (4 bytes)
- Cache holds 1024 words as data

- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



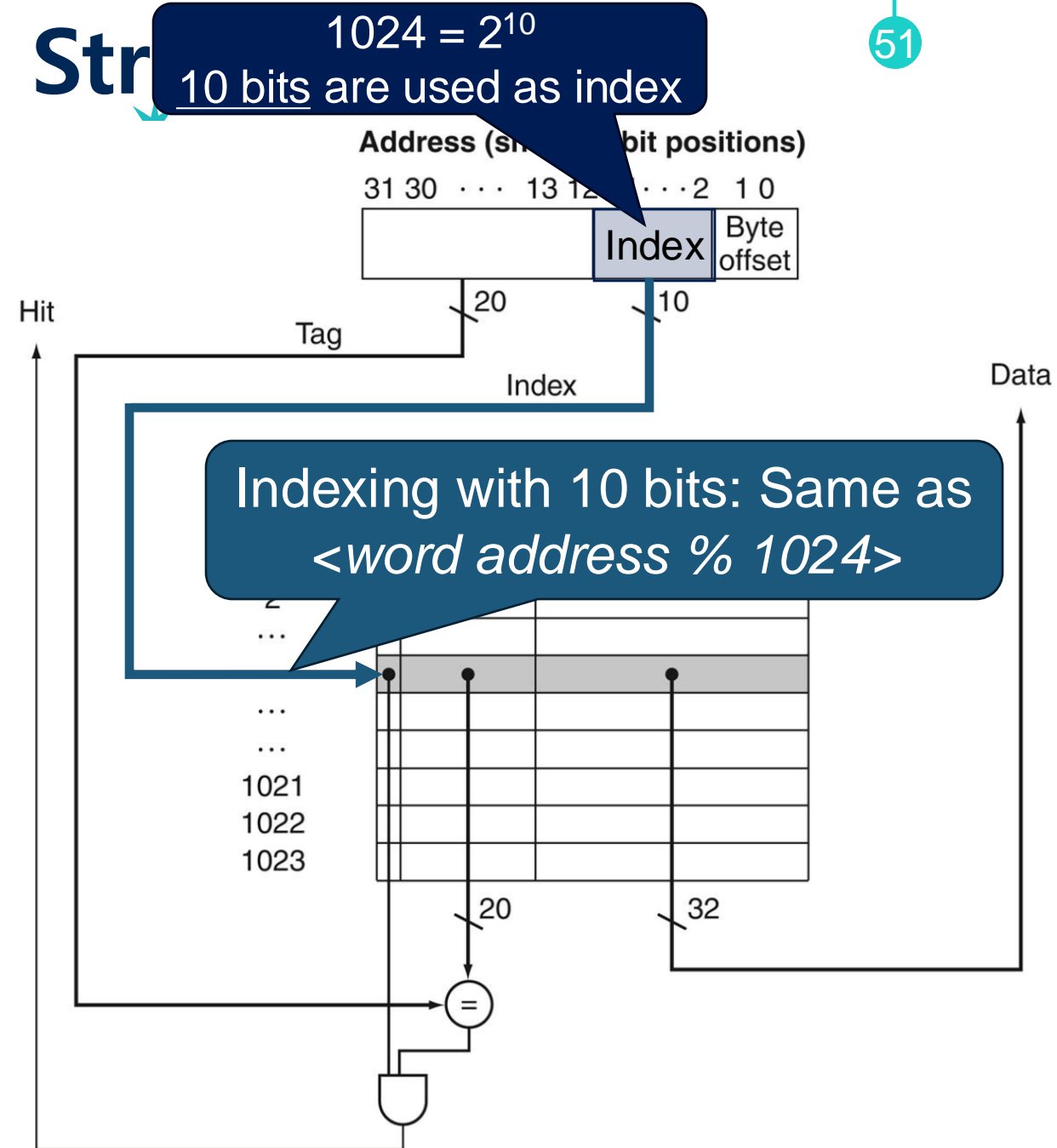
Direct Mapped Cache Structure

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



Direct Mapped Cache Str

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



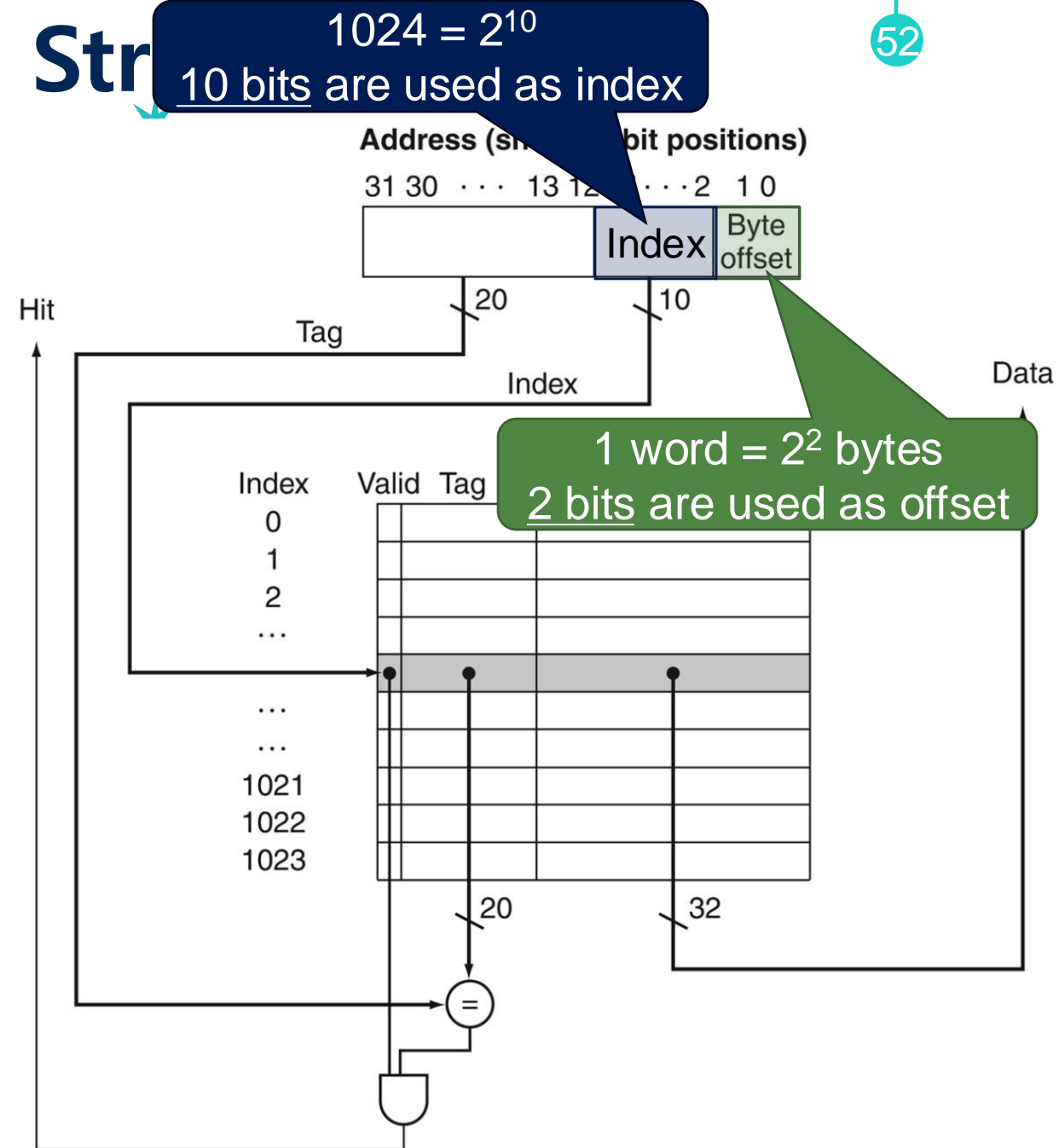
Direct Mapped Cache Structure

52

- **Assumption:**

- Each block is 1 word (4 bytes)
- Cache holds 1024 words as data

- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



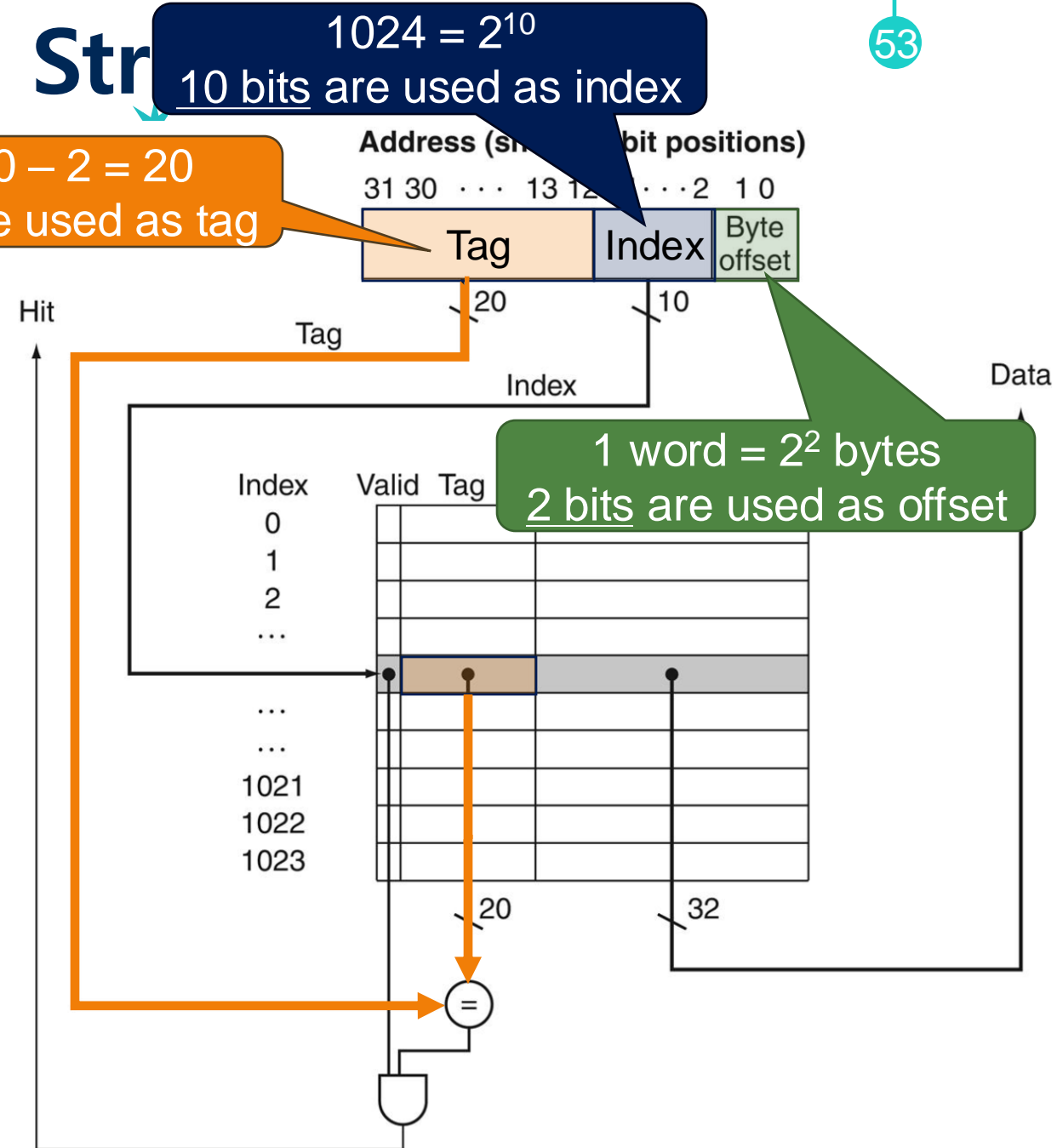
Direct Mapped Cache Structure

53

- **Assumption:**

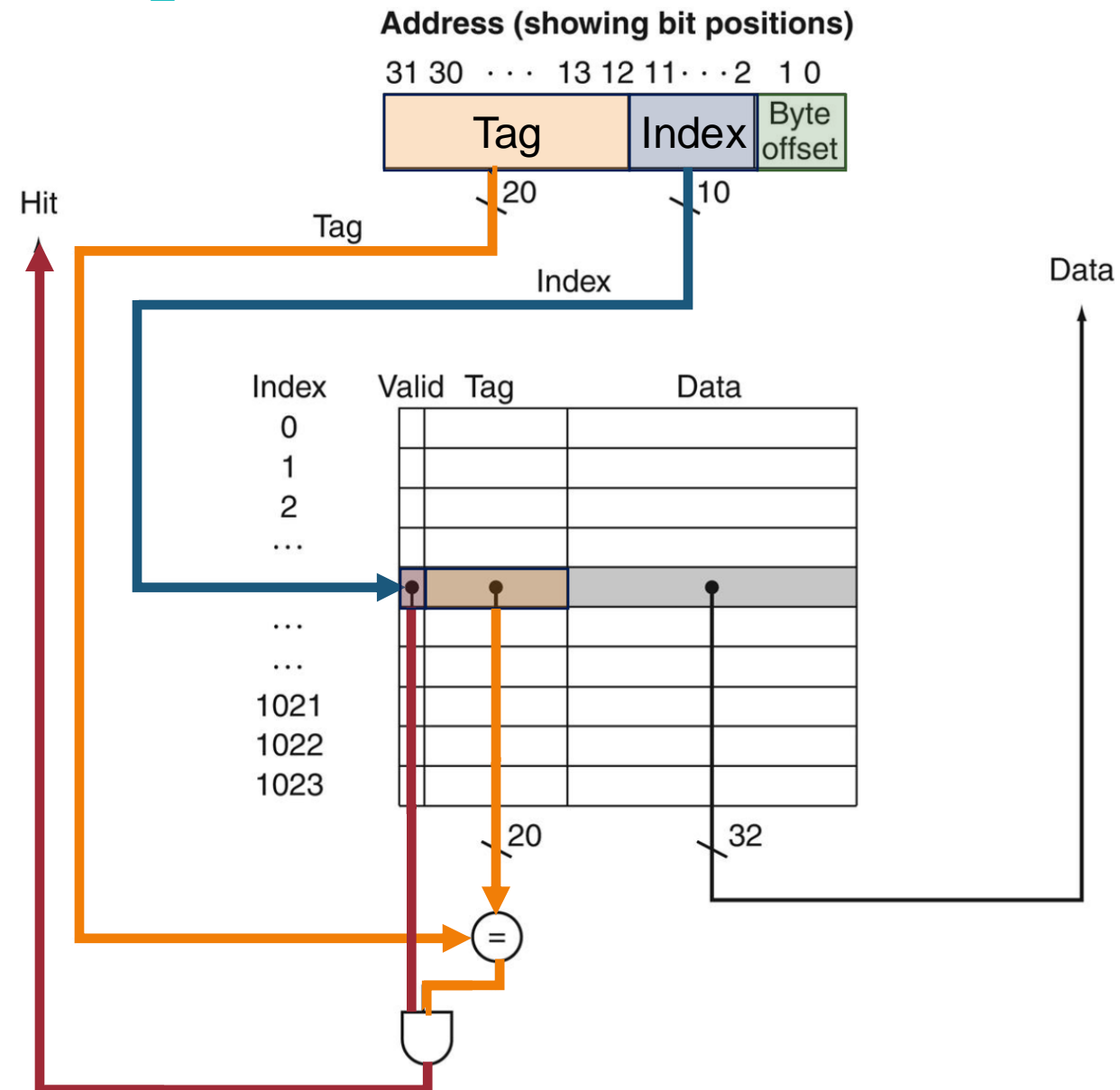
- Each block is 1 word (4 bytes)
- Cache holds 1024 words as data

- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



Direct Mapped Cache Structure

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?



Recap: Hit

(4) If it's a hit: access the **cache data** (no need to access the memory)

Processor

Access to address **11101** (**memory address**) using lw instruction

8 block
Cache

	000	001	010	011	100	101	110	111
Tag						11		
Valid					1			

(1) Access to the cache index:
 $11101 \% 1000 = 101$

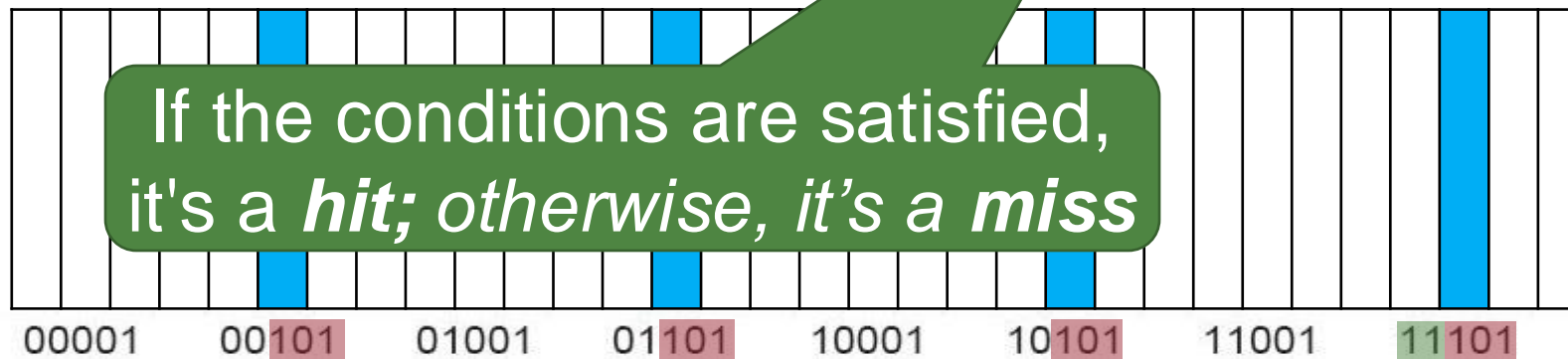
(2) Check if tag matches

==

(3) Check if valid bit sets

Memory

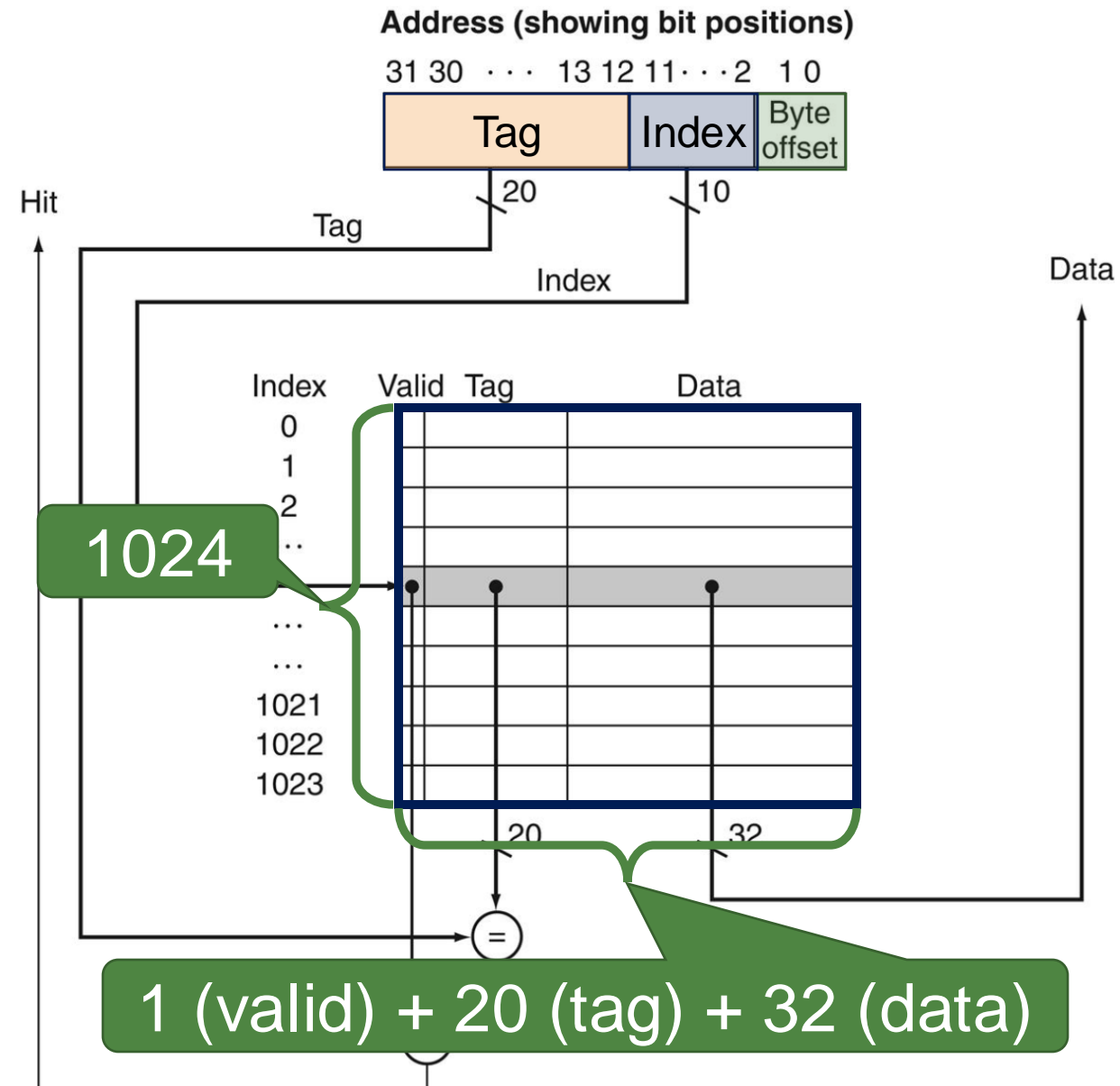
If the conditions are satisfied, it's a **hit**; otherwise, it's a **miss**



Direct Mapped Cache Structure

- **Assumption:**
 - Each block is 1 word (4 bytes)
 - Cache holds 1024 words as data
- How should the memory address be structured?
- What is the tag field size?
- What is the total number of bits in cache?

1024 x (1+20+32) bits



Handling Cache Hits & Misses

Hits vs. Misses



- Read hits
- Read misses
- Write hits
- Write misses

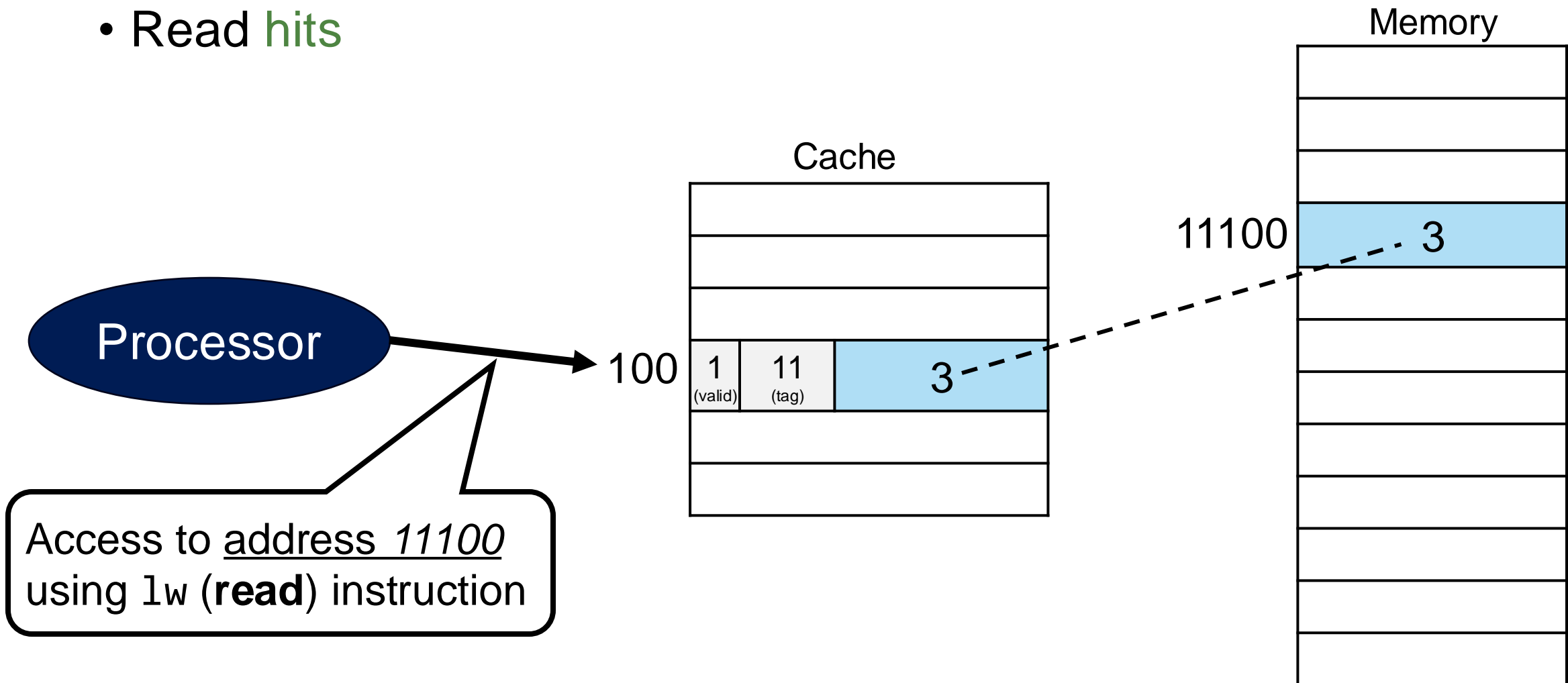
Hits vs. Misses



- Read hits
- Read misses
- Write hits
- Write misses

Hits vs. Misses: Read Hits

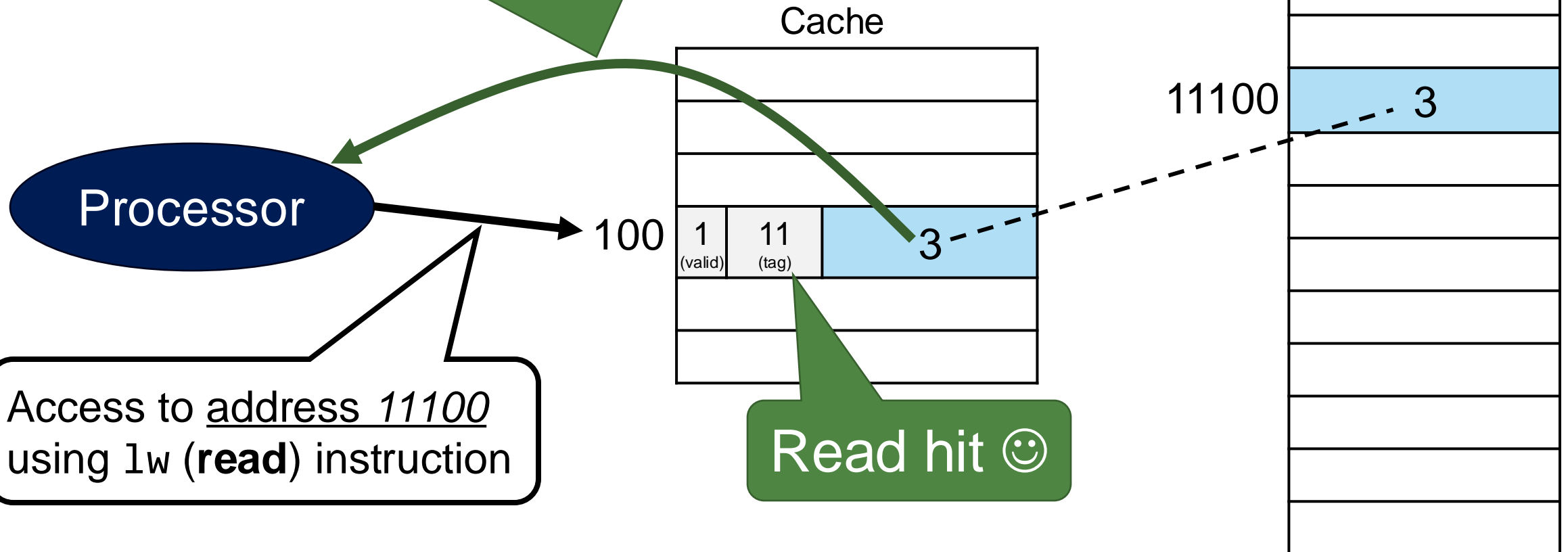
- Read hits



Hits vs. Misses: Read Hits

- Read hits

No penalty!
This is what we want!



Summary: Hits vs. Misses

- Read hits
 - This is what we want!
- Read misses
- Write hits
- Write misses

Summary: Hits vs. Misses

- Read hits
 - This is what we want!
- Read misses
- Write hits
- Write misses

Hits vs. Misses: Read Misses

- Read **misses**

Processor

Access to address 11100
using lw (**read**) instruction

100

Cache

1 (valid)	10 (tag)	1

Read miss ☹️

Memory

11100	3
10100	1



Read Misses



(1) Stall the CPU

Stall

100

Access to address 11100
using lw (**read**) instruction

Cache

1 (valid)	10 (tag)	1

Read miss ☹️

Memory

11100	3
10100	1

Hits vs. Misses: Read Misses

- Read misses

(1) Stall the CPU

Stall

Access to address 11100
using lw (**read**) instruction

100

1 (valid)	11 (tag)	3

(2) Fetch block
from memory

11100

3

10100

1

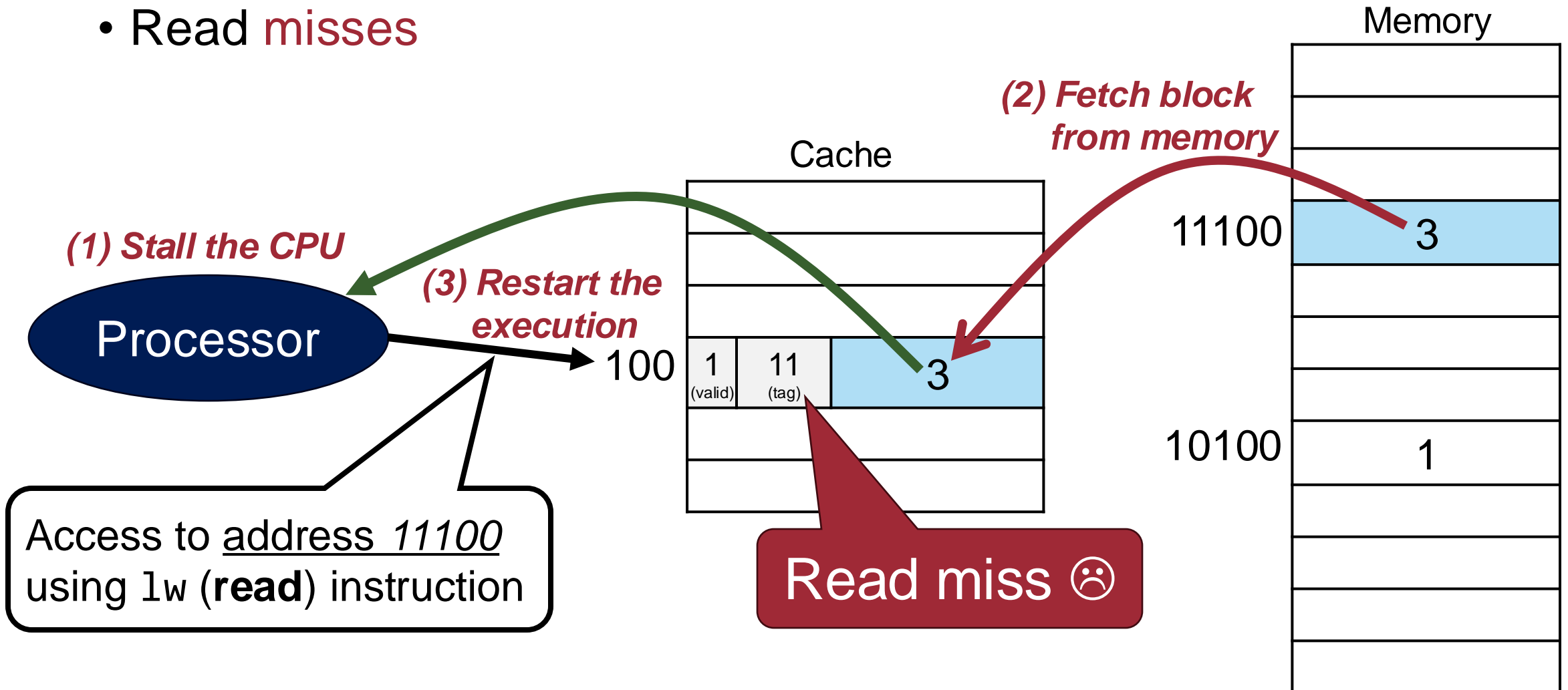
Read miss ☹️

Memory

3
1

Hits vs. Misses: Read Misses

- Read **misses**



Hits vs. Misses: Read Misses



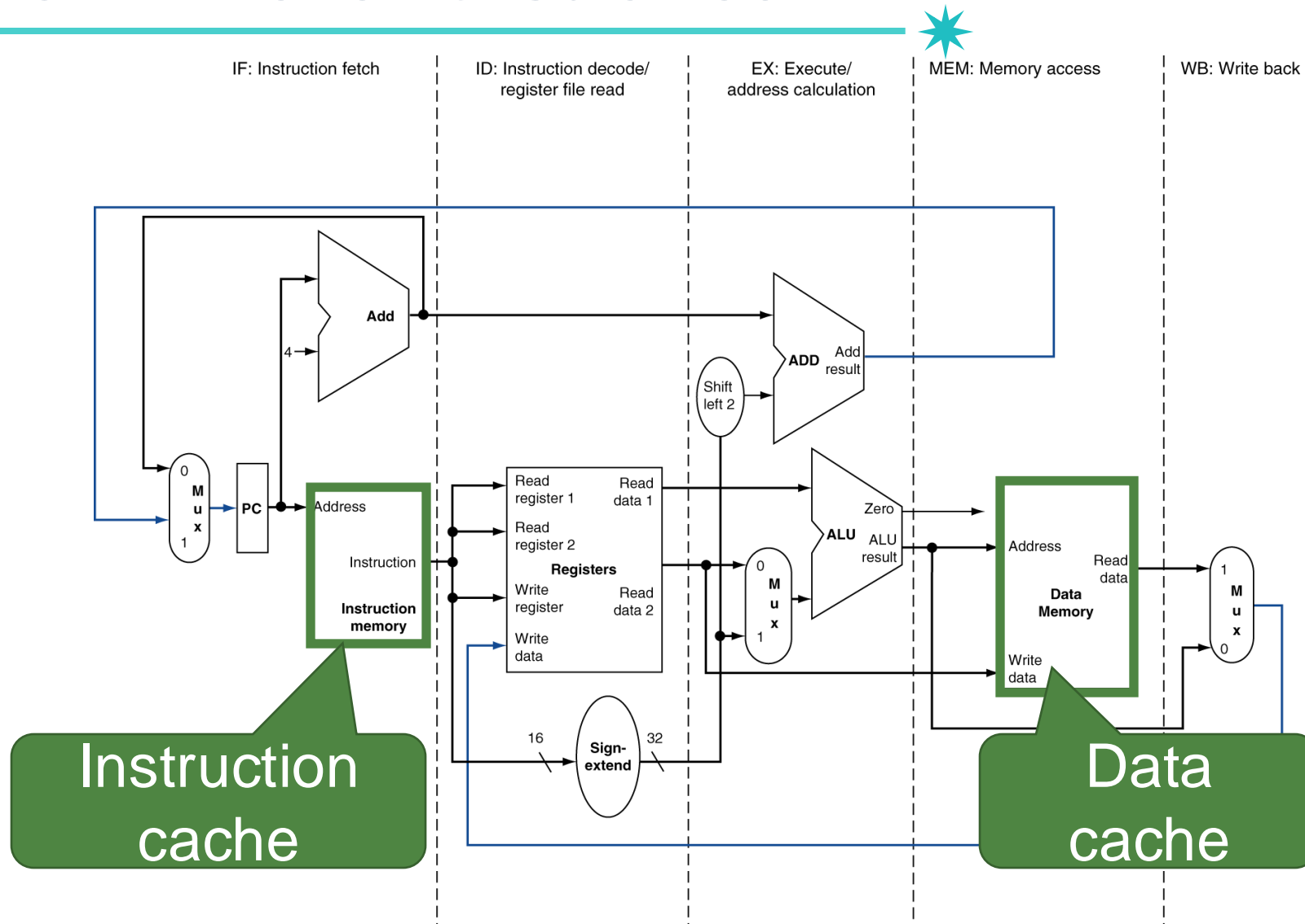
- Read misses

(1) **CPU stalls**: freezing the contents of all the registers while waiting for memory

(2) **Fetch block** from memory

(3) **Restart** the execution

Two Different Caches



Hits vs. Misses: Read Misses



- Read **misses**

(0) (For the instruction cache) Set PC-4

(1) **CPU stalls**: freezing the contents of all the registers while waiting for memory

(2) **Fetch block** from memory

(3) **Restart** the execution (= Re-fetch instruction)

Summary: Hits vs. Misses

- Read hits
 - This is what we want!
- Read misses
 - Stall the CPU, fetch block from memory, restart
- Write hits
- Write misses

Summary: Hits vs. Misses

- Read hits
 - This is what we want!

- Read misses
 - Stall the CPU, fetch block from memory, restart

*Cases to consider
for the instruction cache*

- Write hits

- Write misses

*Cases to consider
for the data cache*

Summary: Hits vs. Misses

- Read hits
 - This is what we want!

- Read misses
 - Stall the CPU, fetch block from memory, restart

*Cases to consider
for the instruction cache*

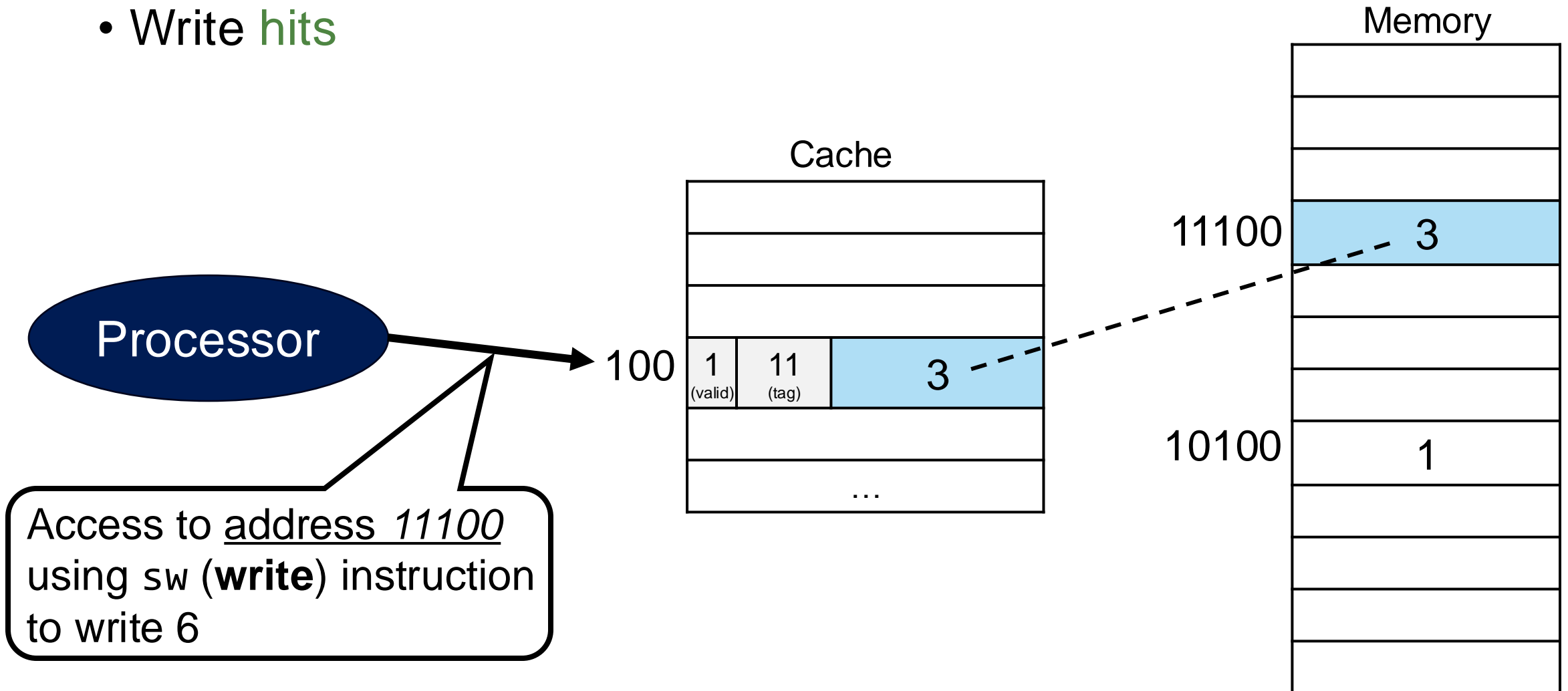
- Write hits

- Write misses

*Cases to consider
for the data cache*

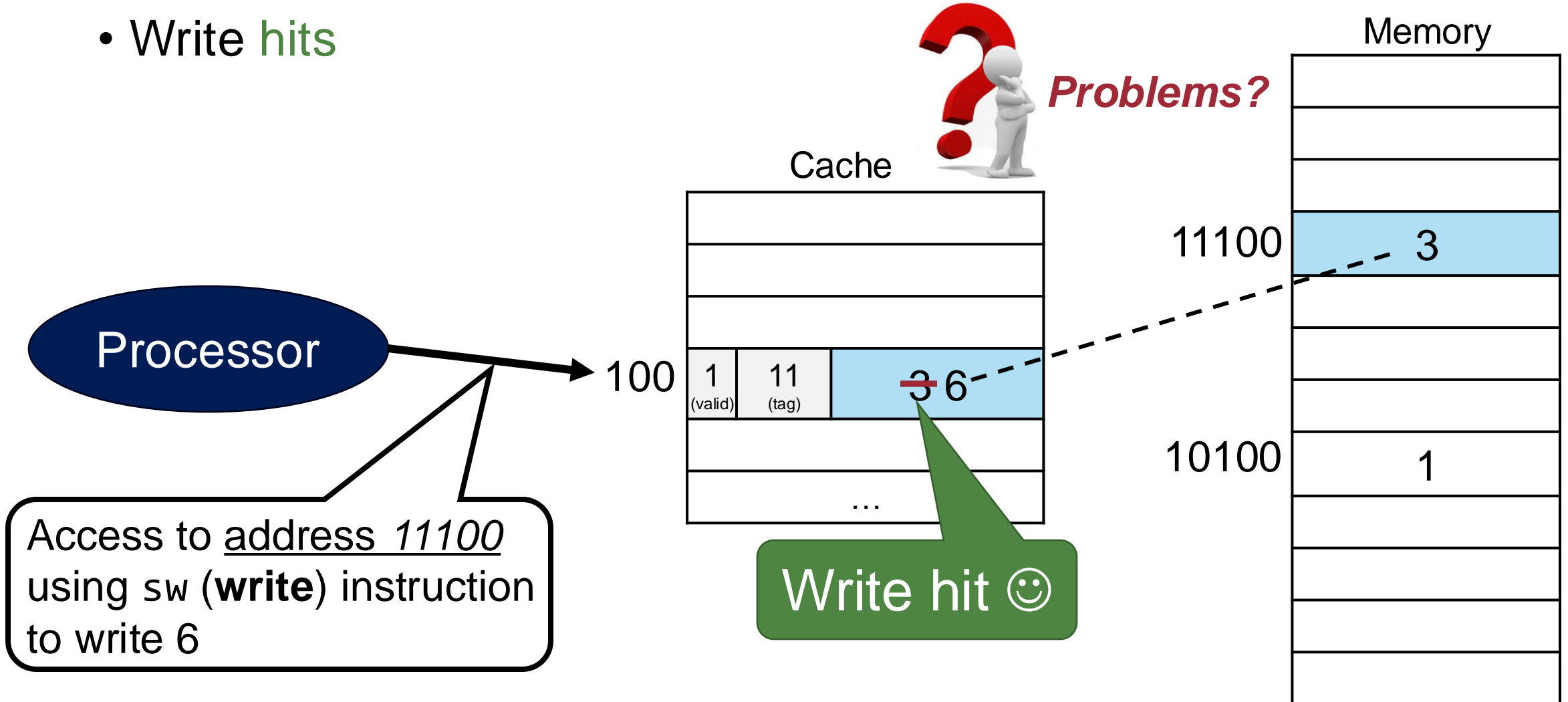
Hits vs. Misses: Write Hits

- Write hits



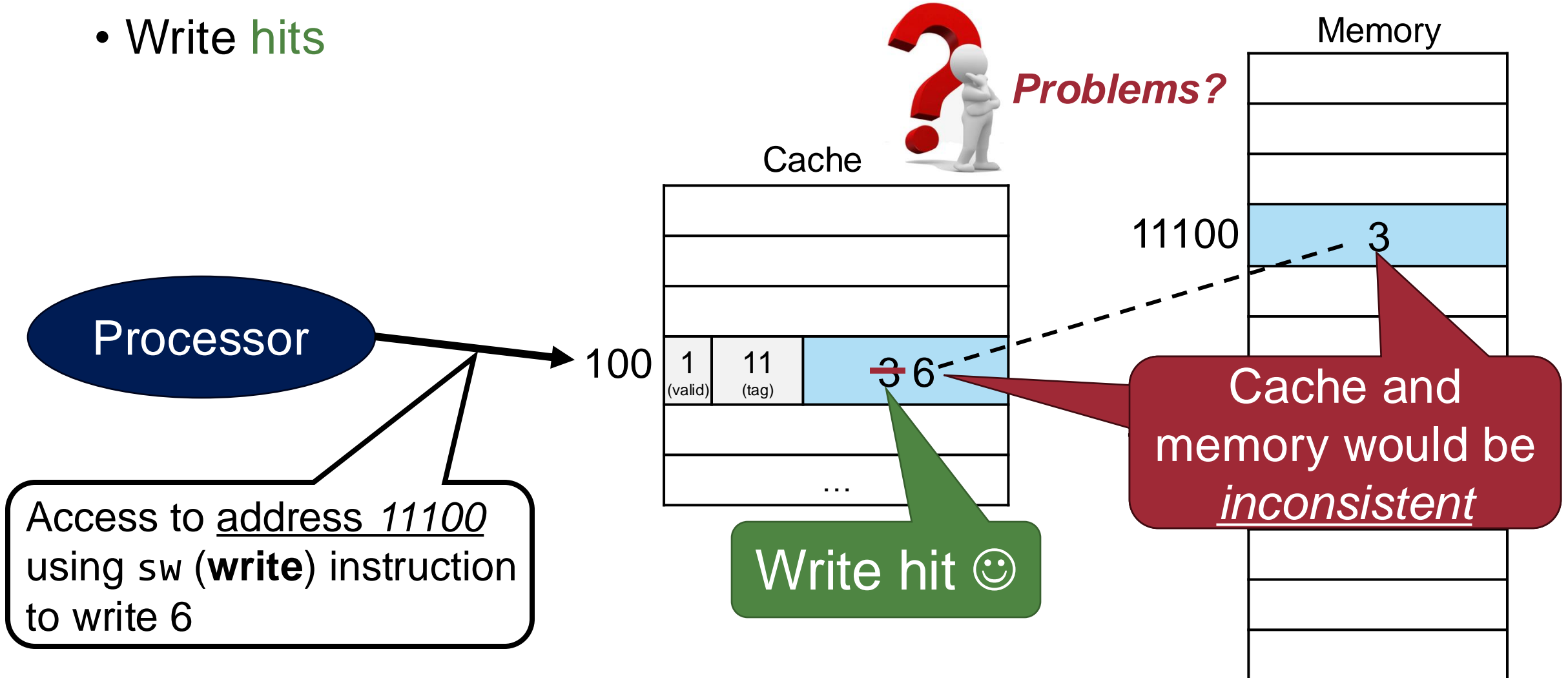
Hits vs. Misses: Write Hits

- Write hits



Hits vs. Misses: Write Hits

- Write hits



Write Policies



- Write **hits**
 - Cache and memory would be inconsistent! What can we do?

(1) Write through

(2) Write back

Write Policies



- Write **hits**
 - Cache and memory would be inconsistent! What can we do?
- (1) **Write through:** On each write hit, the information is written to both in the cache and in the memory
- (2) **Write back**

Write Through



Processor

Access to address 11100
using sw (**write**) instruction
to write 6

100

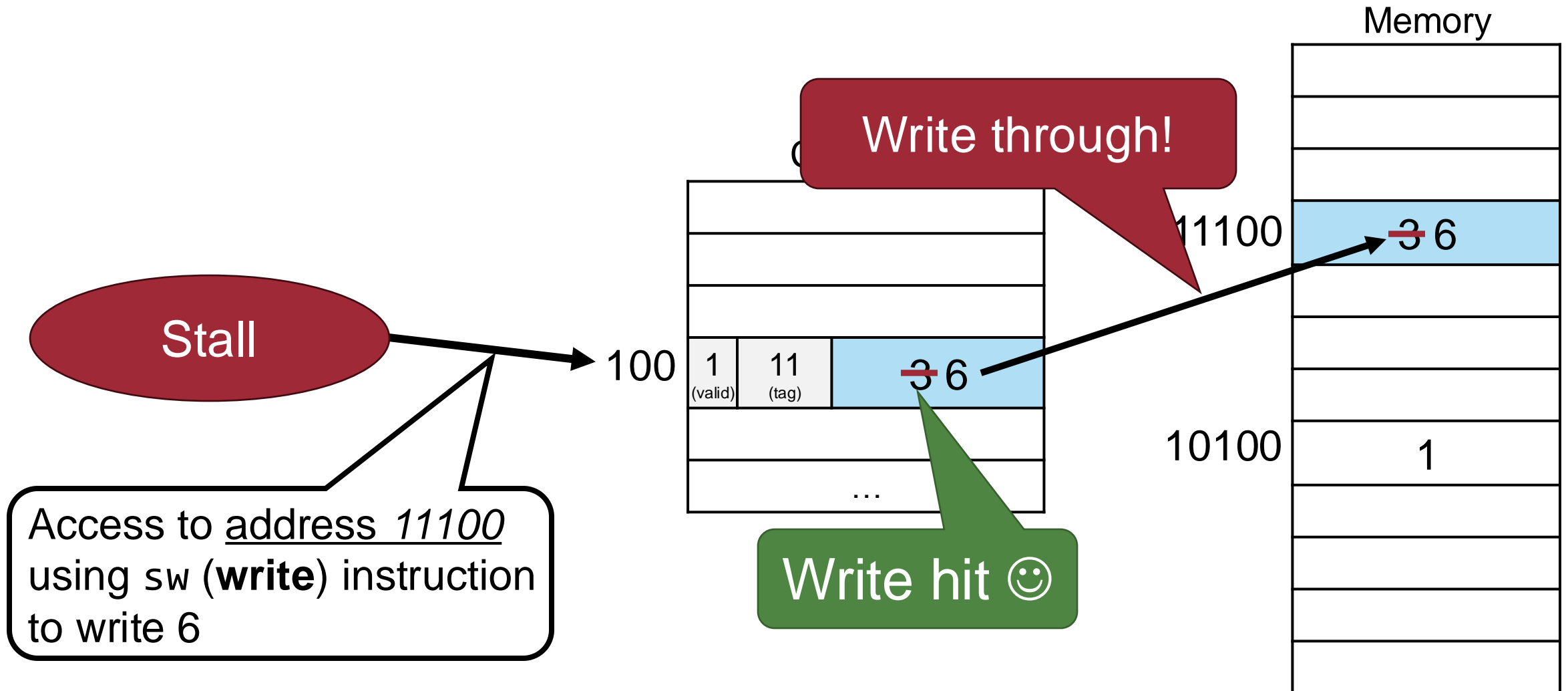
1 (valid)	11 (tag)	3 6
		...

Write hit 😊

Memory

11100	3
10100	1

Write Through



Write Through

82

Processor

Access to address 11100
using sw (**write**) instruction
to write 9

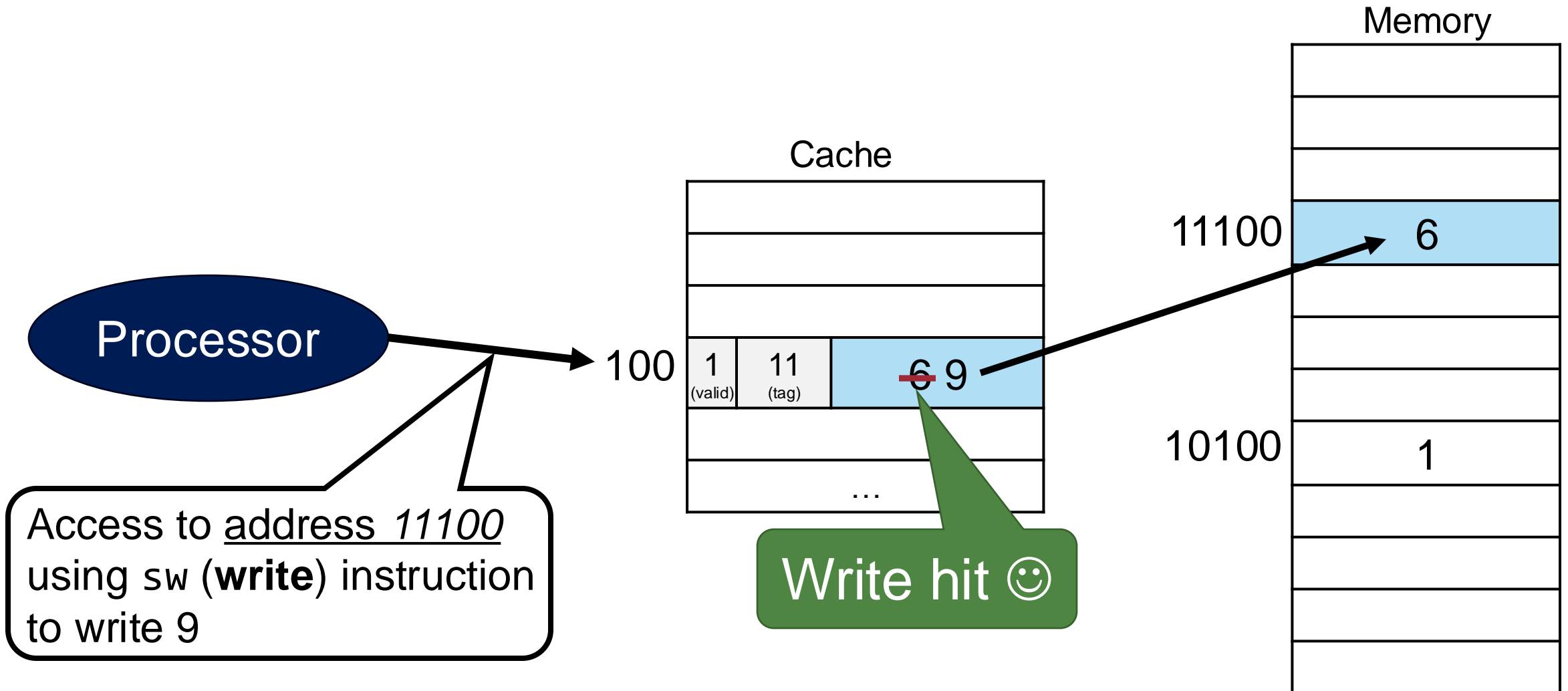
100

1 (valid)	11 (tag)	3 6
		...

11100	3 6
10100	1

Write Through

83



Write Through

Performance degradation:
whenever data cache is updated,
there should a memory write

Stall

Access to address 11100
using sw (**write**) instruction
to write 9

100

1 (valid)	11 (tag)	6 9
		...

Write hit 😊

11100

Memory

10100

1

~~6~~ 9

Write Through



Processor

Access to address **10100**
using lw (**read**) instruction

100

Cache

1 (valid)	11 (tag)	6 9
		...

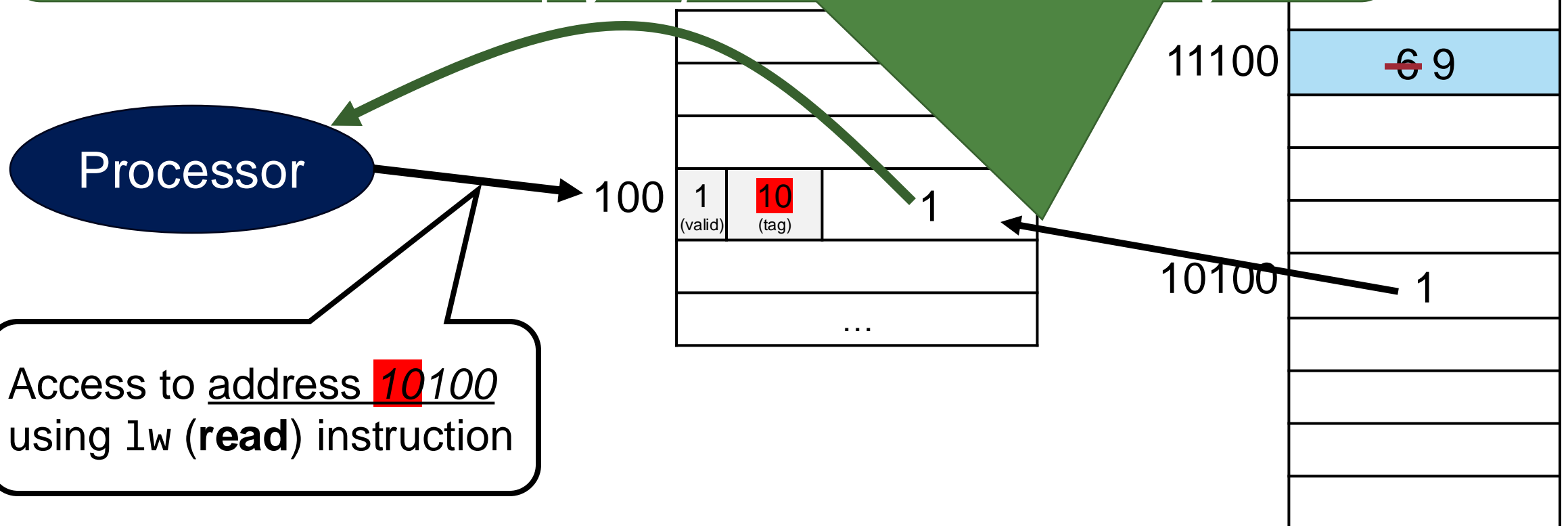
Read miss ☹️

Memory

11100	6 9
10100	1

Write Through

Fast processing in case of miss:
The value at address 11100 is already synchronized, so we can simply replace the cache entry



Write Through

87

Processor

Access to address **10**100
using sw (**write**) instruction
to write 7

100

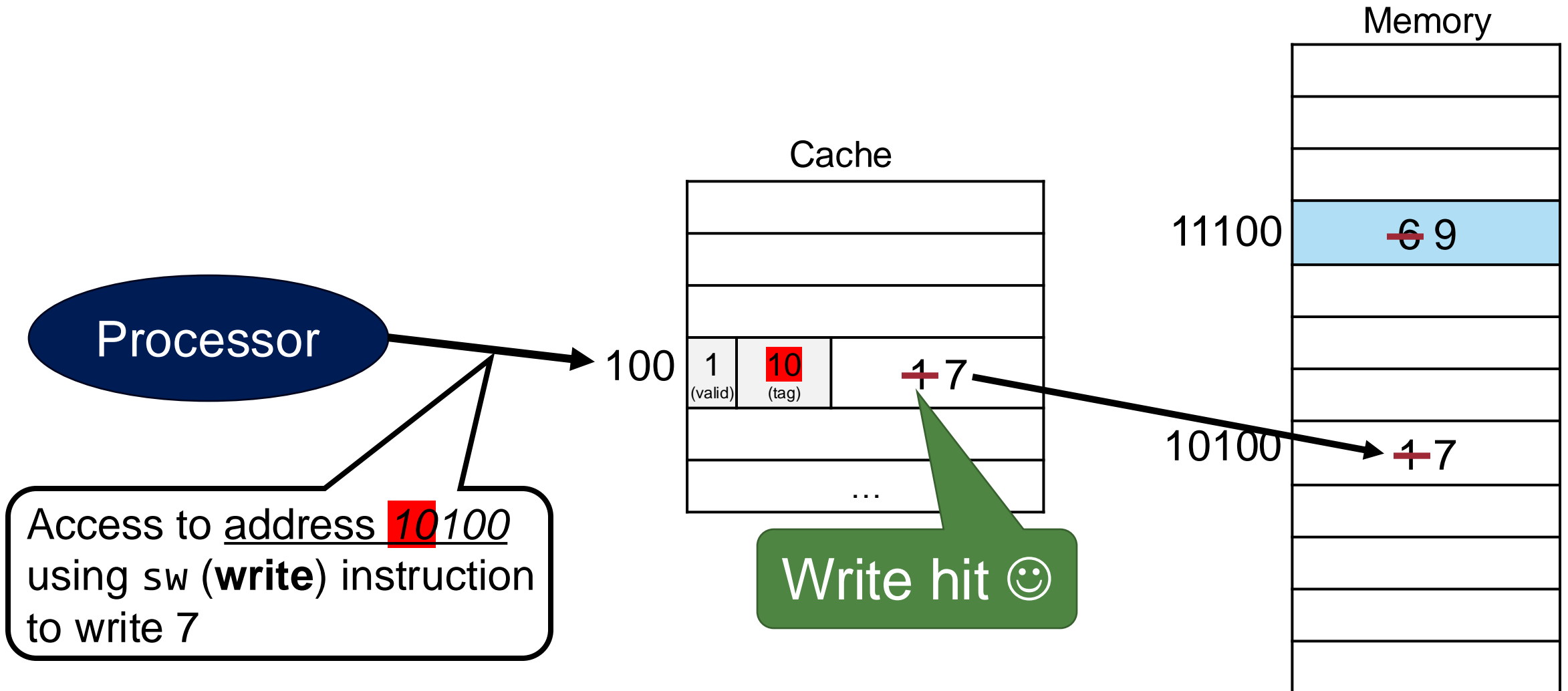
Cache

1 (valid)	10 (tag)	1
		...

Memory

11100	6 9
10100	1

Write Through



Main Disadvantage: Performance Degradation ⁹⁰

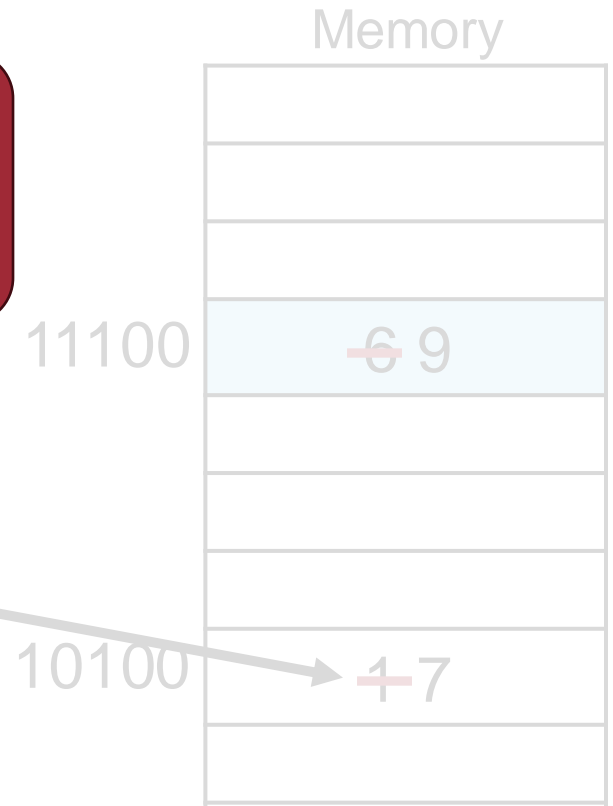
Performance degradation:
whenever data cache is updated,
there should a memory write

Processor



100 1 10 17
(valid) (tag)
Solution? → Write Buffer

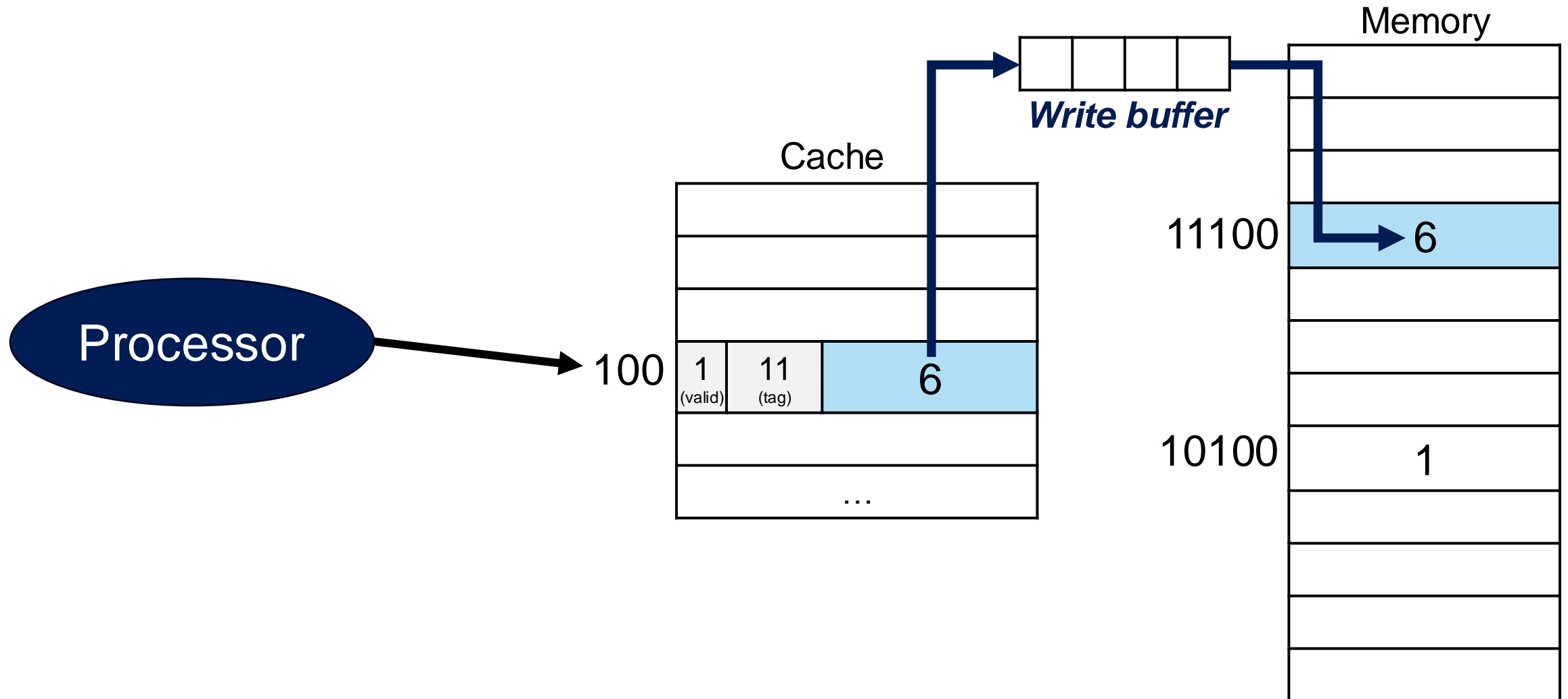
Access to address 10100



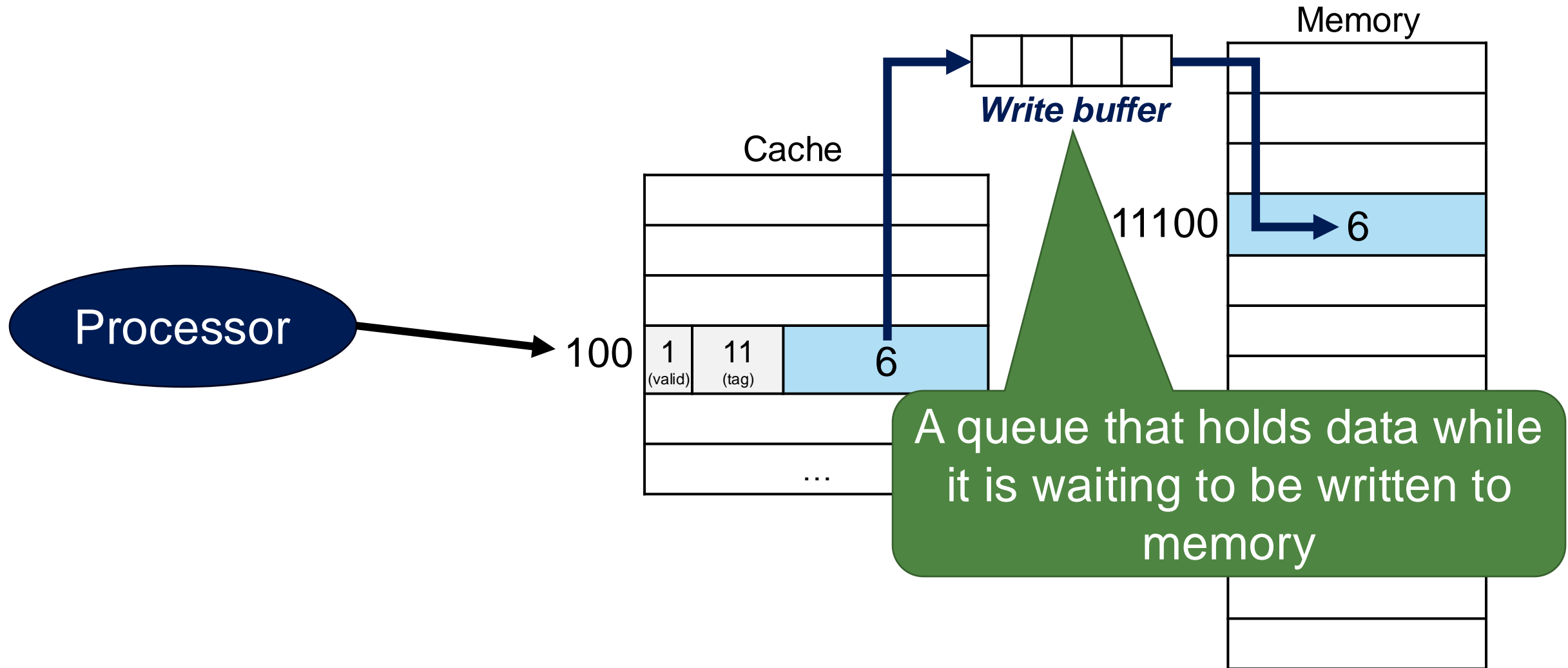
If base CPI = 1, 10% of instructions are stores,
write to memory takes 100 extra cycles

$$\text{CPI} = 1 + 0.1 \times 100 = 11$$

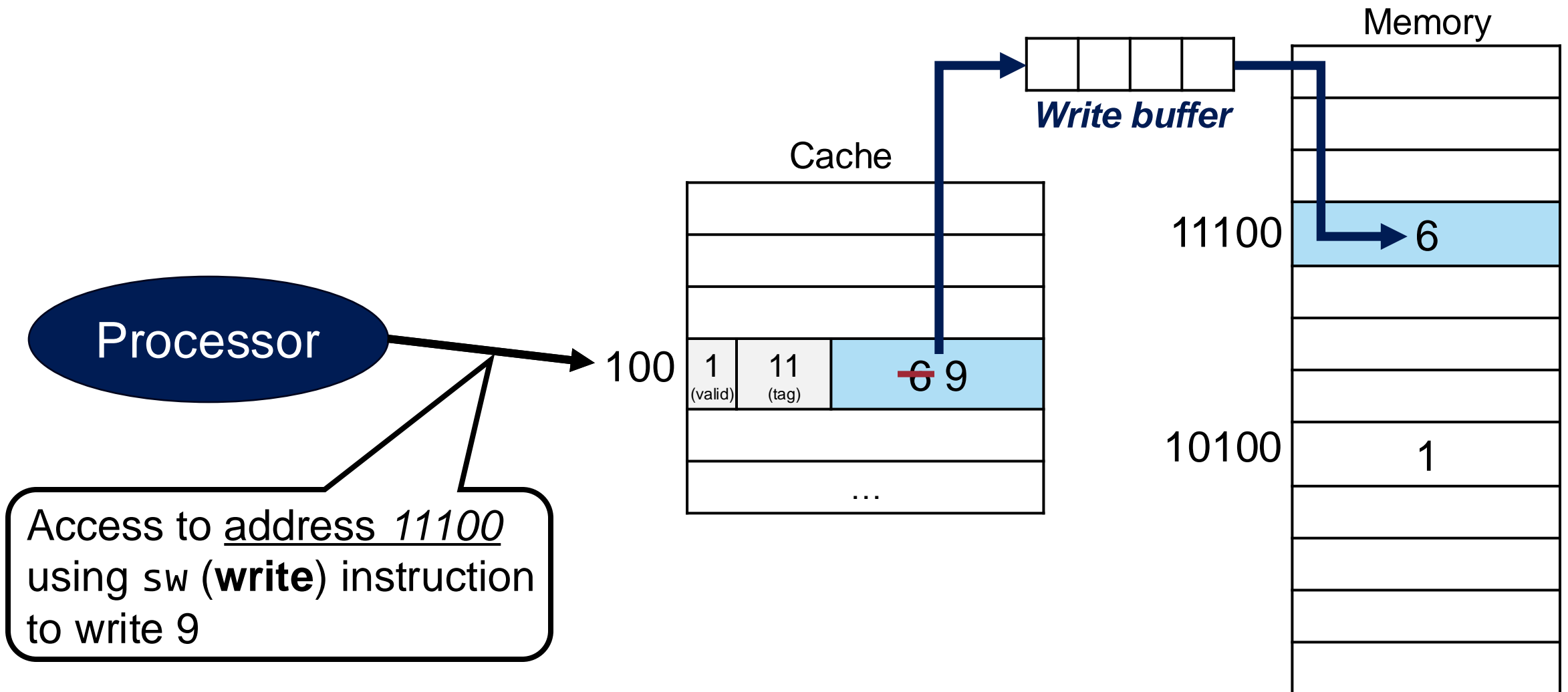
Simple Idea: Write Buffer + Write Through



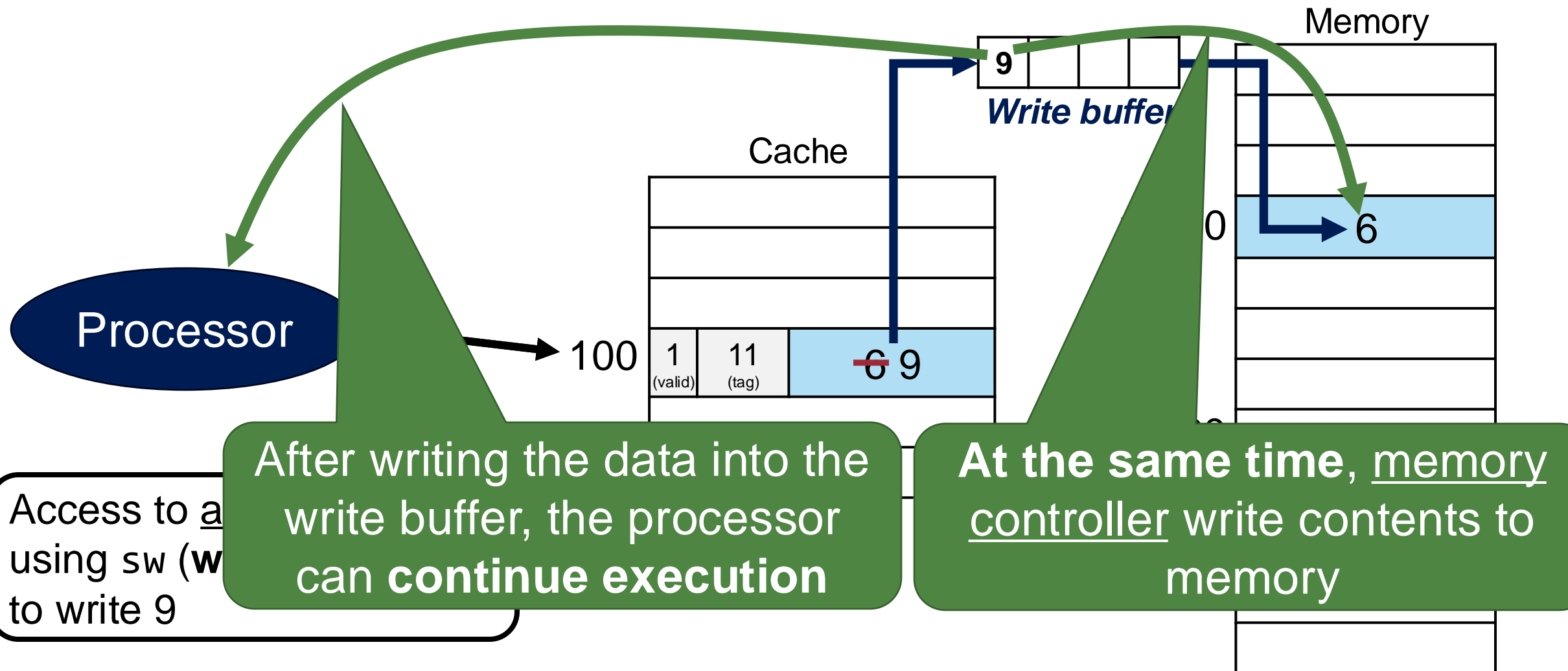
Simple Idea: Write Buffer + Write Through



Simple Idea: Write Buffer + Write Through



Simple Idea: Write Buffer + Write Through



Write Policies



- Write **hits**
 - Cache and memory would be inconsistent! What can we do?
- (1) **Write through:** On each write hit, the information is written to both in the cache and in the memory
 - **Pros:** faster processing in case of 'miss'
 - **Cons:** make writes take longer (whenever data cache is updated, there should a memory write)
 - ✓ Solution: write buffer
- (2) **Write back:** Next lecture!

Question?