

Seongil Wi



Notification: Final Exam

2

- Date: Dec. 19 (Thursday)
- Class Time (1h 15m), Closed book

T/F problems + Computation problems + Descriptive problems

Scope: All the material learned after the midterm

Notification: HW3



- Implementing a cache simulator
- Due: Dec 15, 11:59 PM

 Since the final exam period is approaching, it is recommended to complete this assignment <u>as soon as possible</u>

Handling Cache Hits & Misses

Hits vs. Misses

*

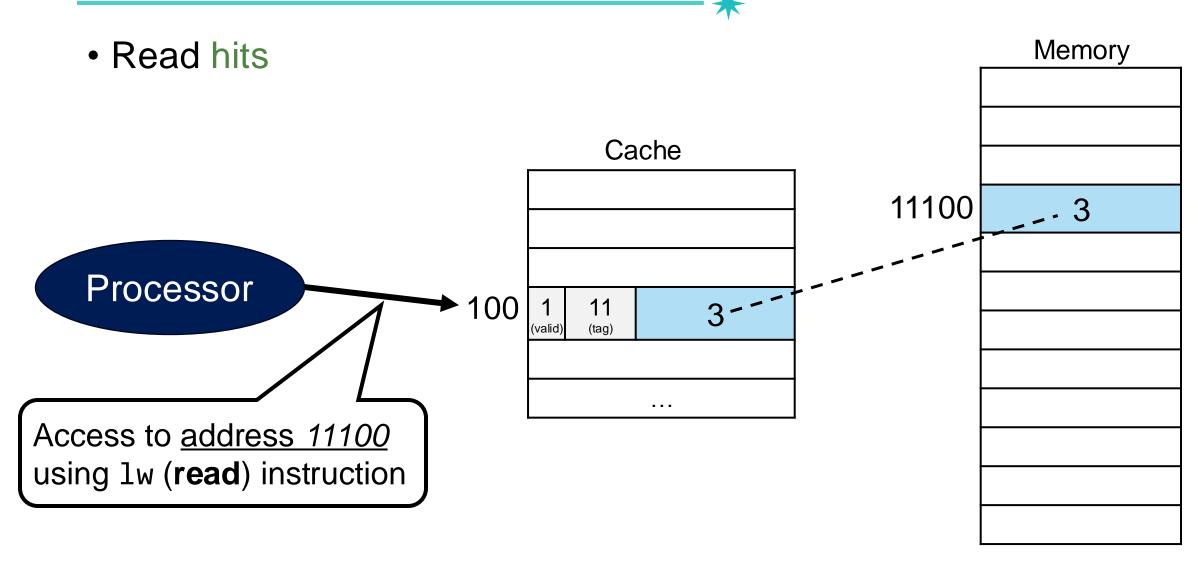
Read hits

Read misses

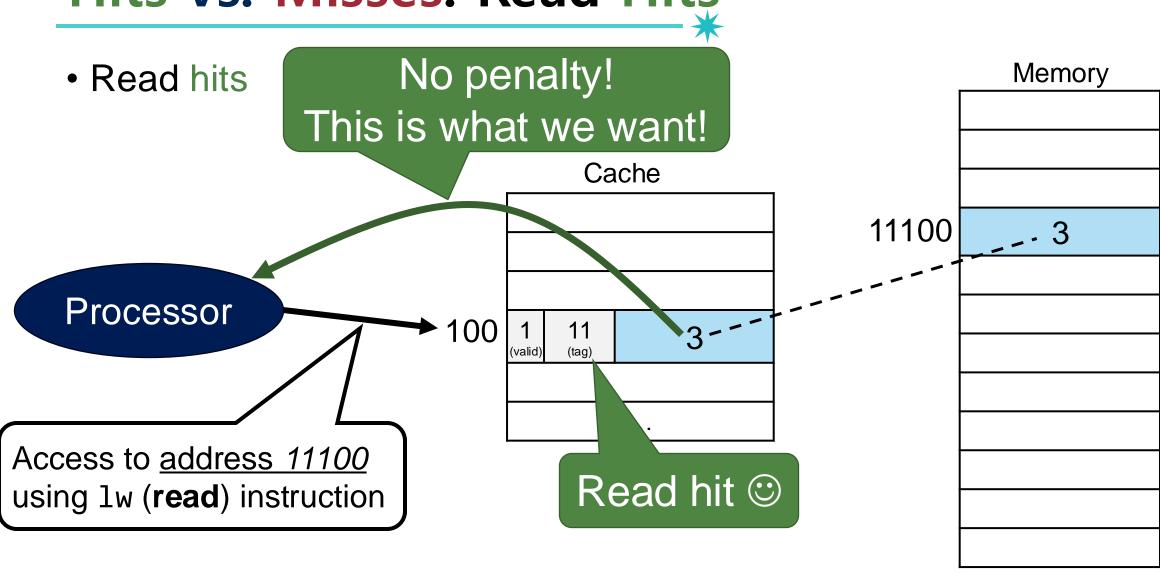
Write hits

Write misses

Hits vs. Misses: Read Hits



Hits vs. Misses: Read Hits



Summary: Hits vs. Misses

- Read hits
 - This is what we want!

Read misses

Write hits

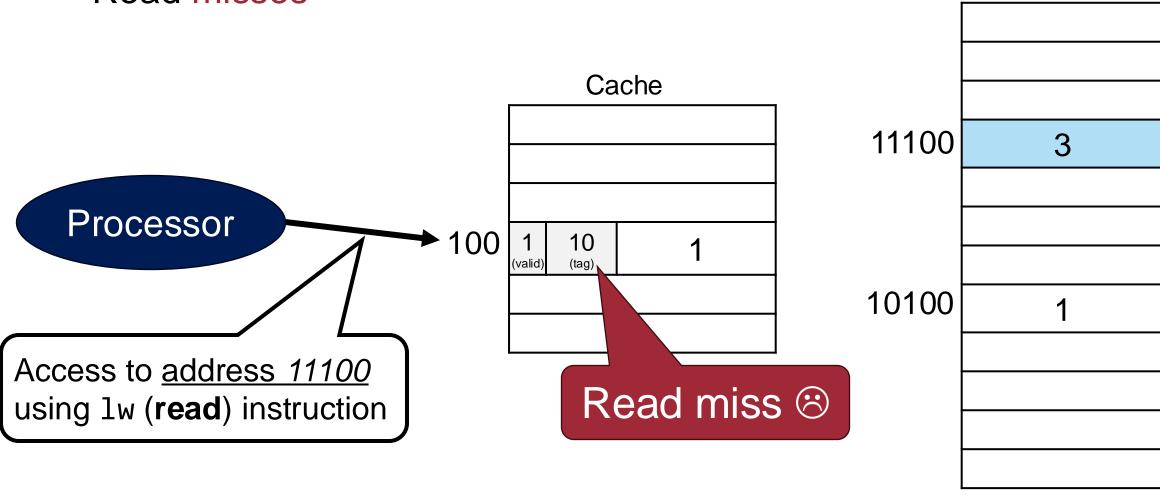
Write misses

Hits vs. Misses: Read Misses



Memory

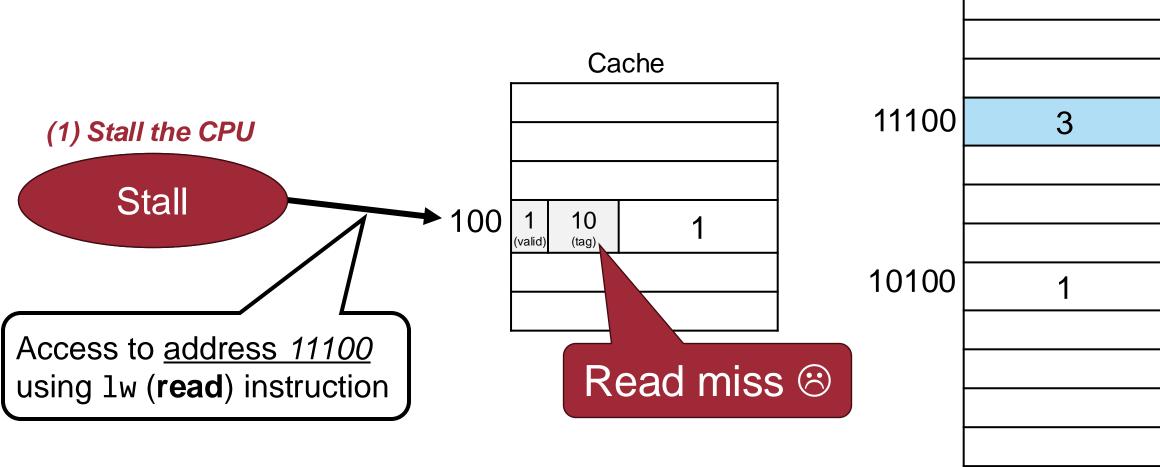
Read misses



Hits vs. Misses: Read Misses

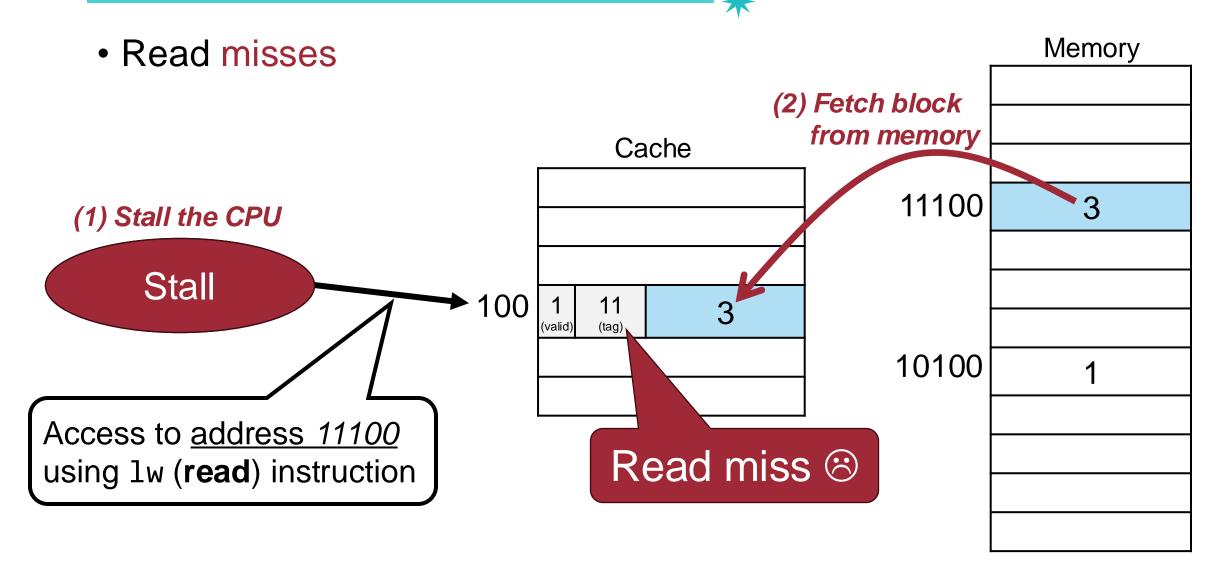


Read misses



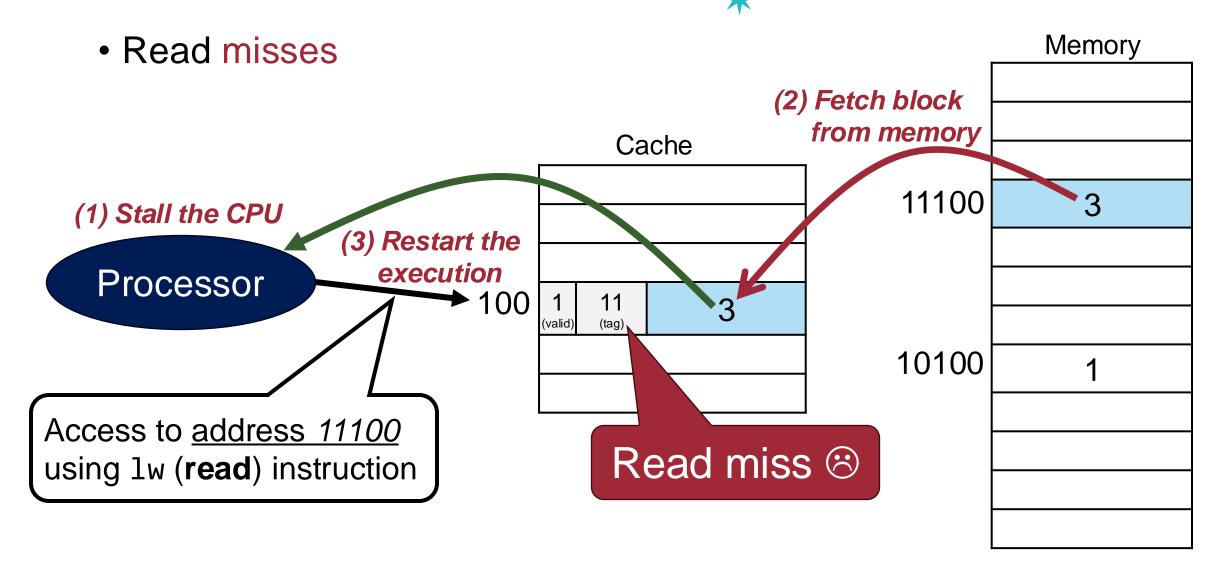
Memory

Hits vs. Misses: Read Misses



Hits vs. Misses: Read Misses





Summary: Hits vs. Misses

- Read hits
 - This is what we want!

Read misses

- Stall the CPU, fetch block from memory, restart

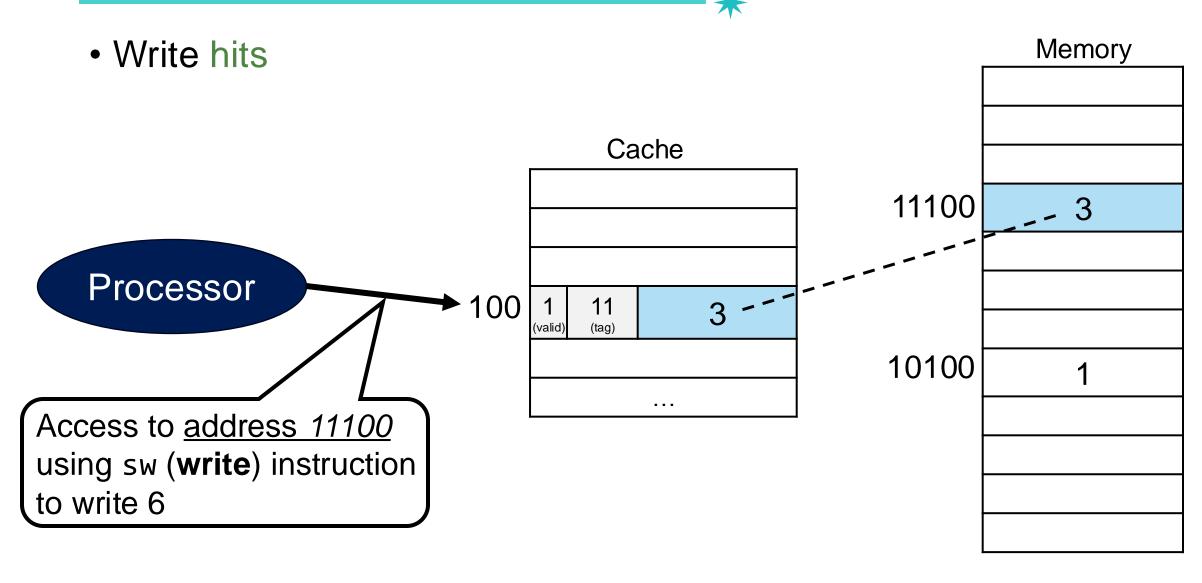
Write hits

Write misses

Cases to consider for the <u>instruction cache</u>

Cases to consider for the data cache

Hits vs. Misses: Write Hits



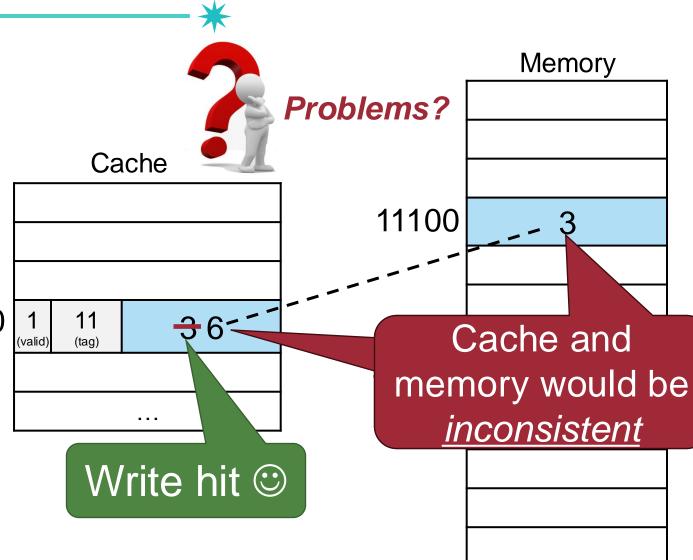
Hits vs. Misses: Write Hits

 Write hits Memory **Problems?** Cache 11100 Processor 100 \$6--11 10100 Access to address 11100 using sw (write) instruction Write hit © to write 6

Hits vs. Misses: Write Hits

Write hits

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



Write Policies





- Write hits
 - Cache and memory would be inconsistent! What can we do?
 - (1) Write through

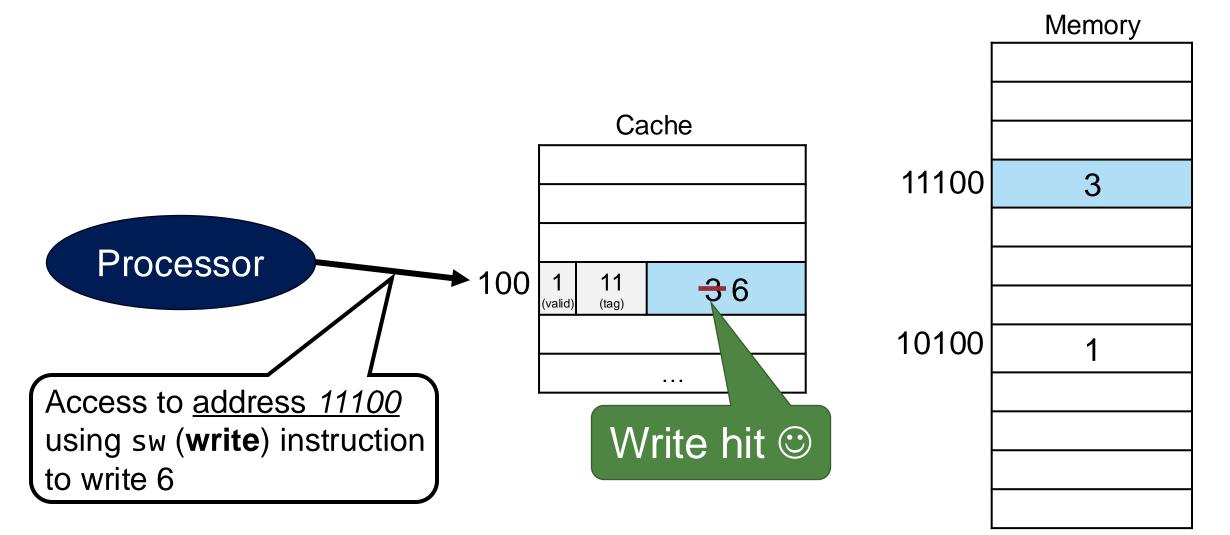
(2) Write back

Write Policies



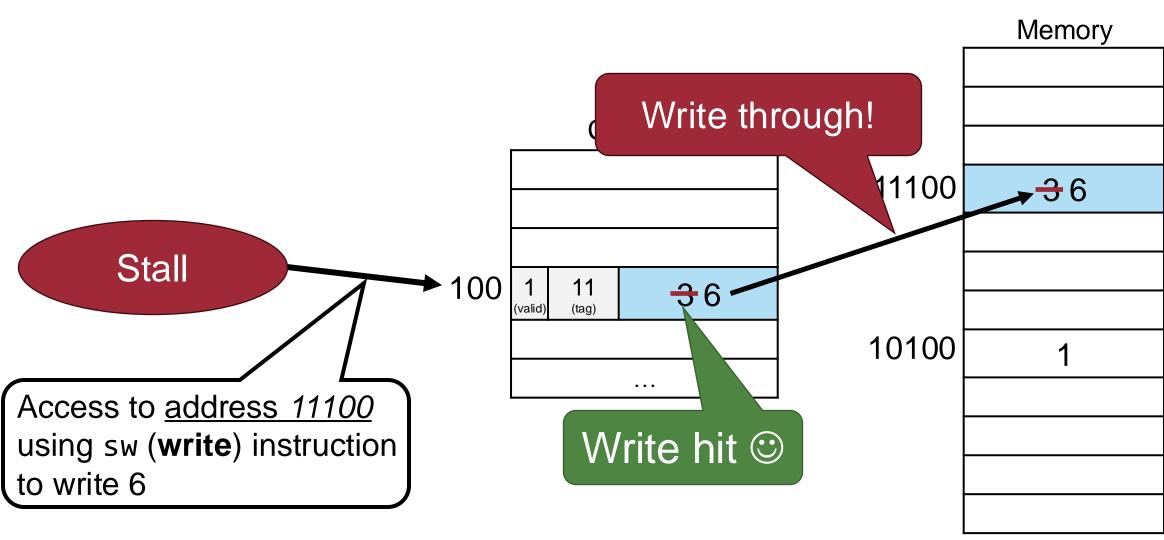
- Write hits
 - Cache and memory would be inconsistent! What can we do?
 - (1) Write through: On each write hit, the information is written to both in the cache and in the memory

(2) Write back

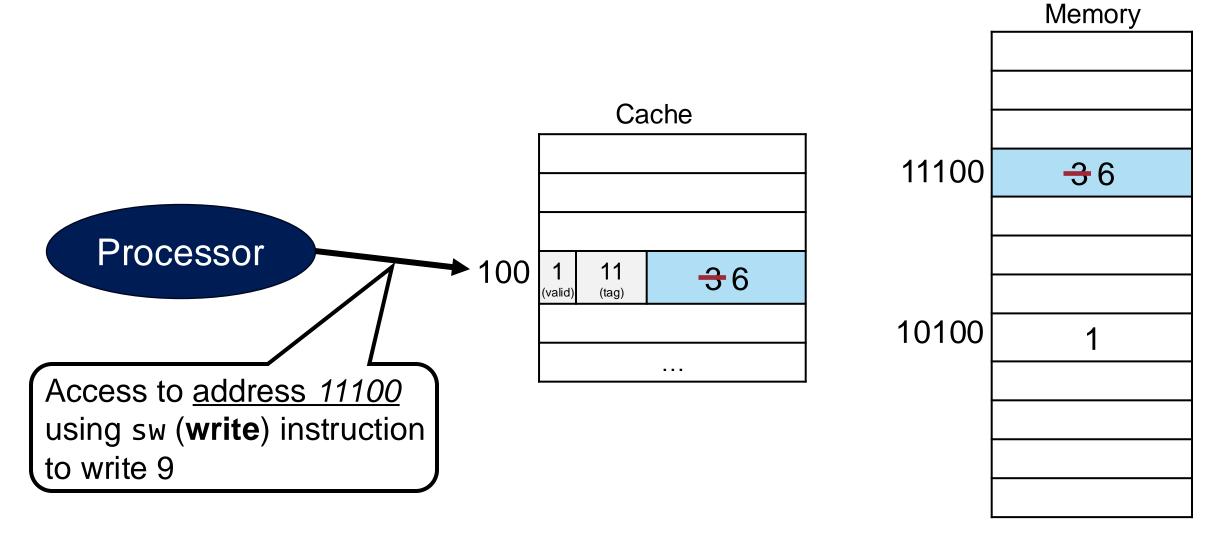




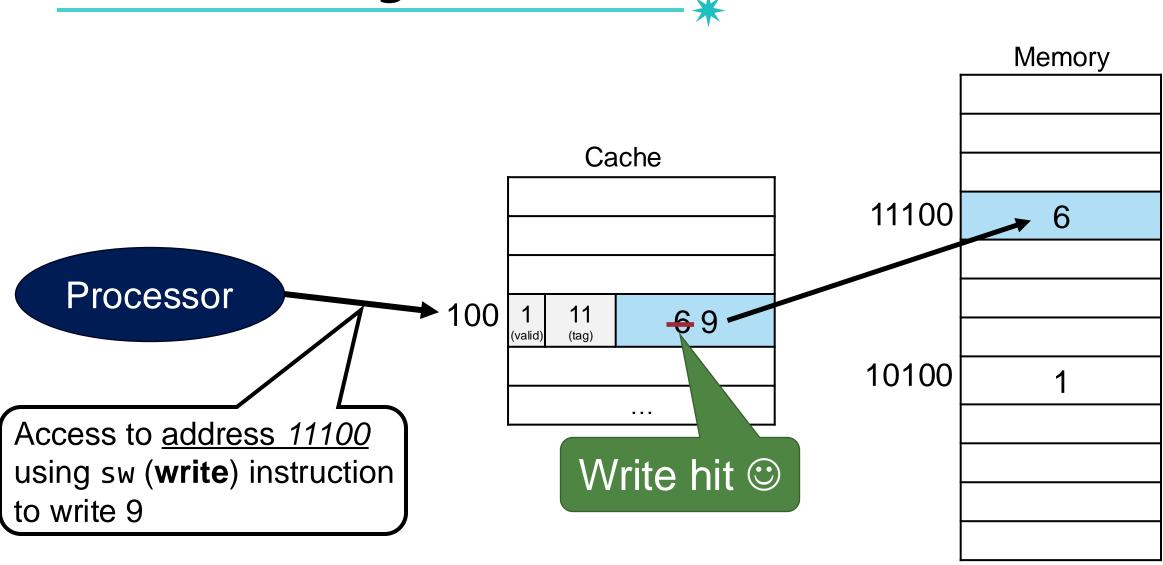






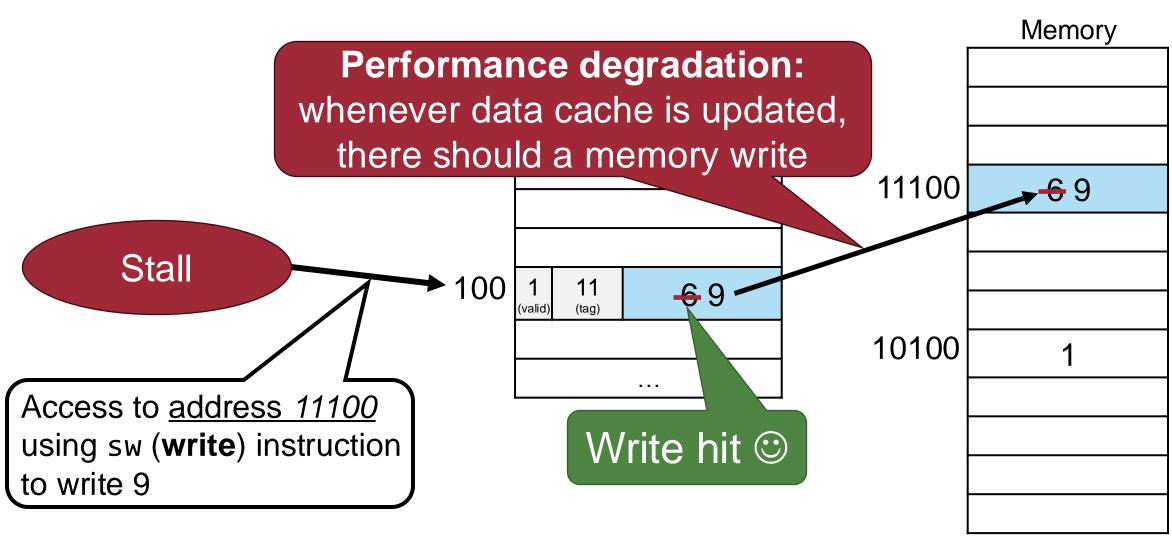




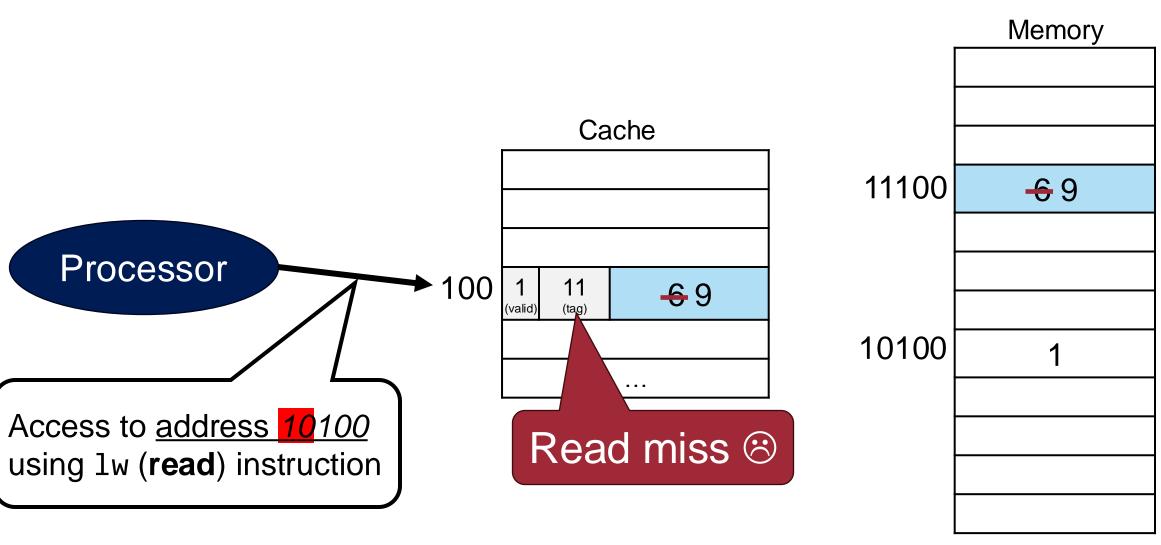


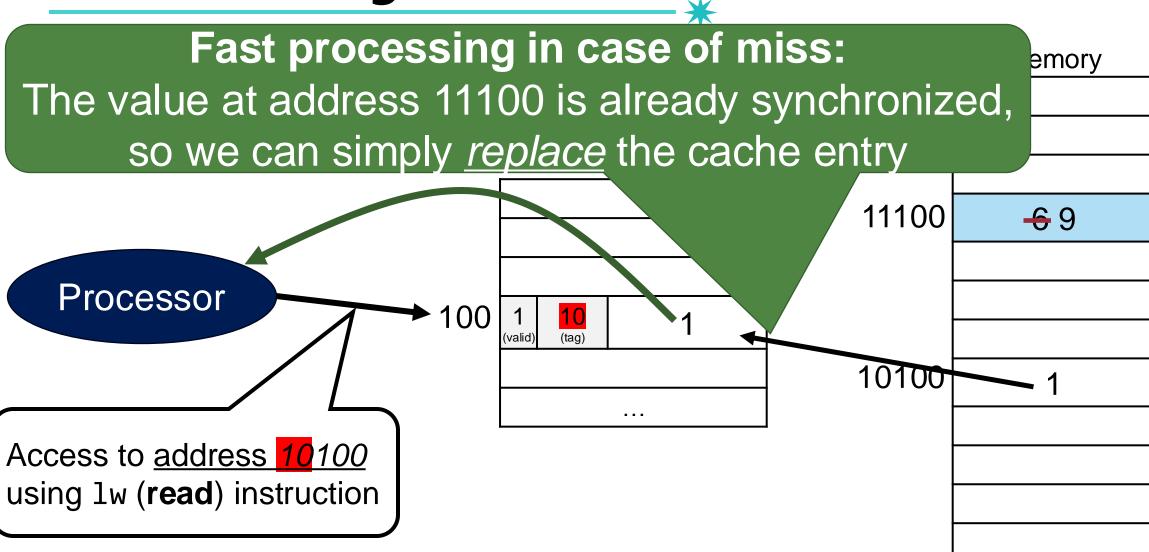




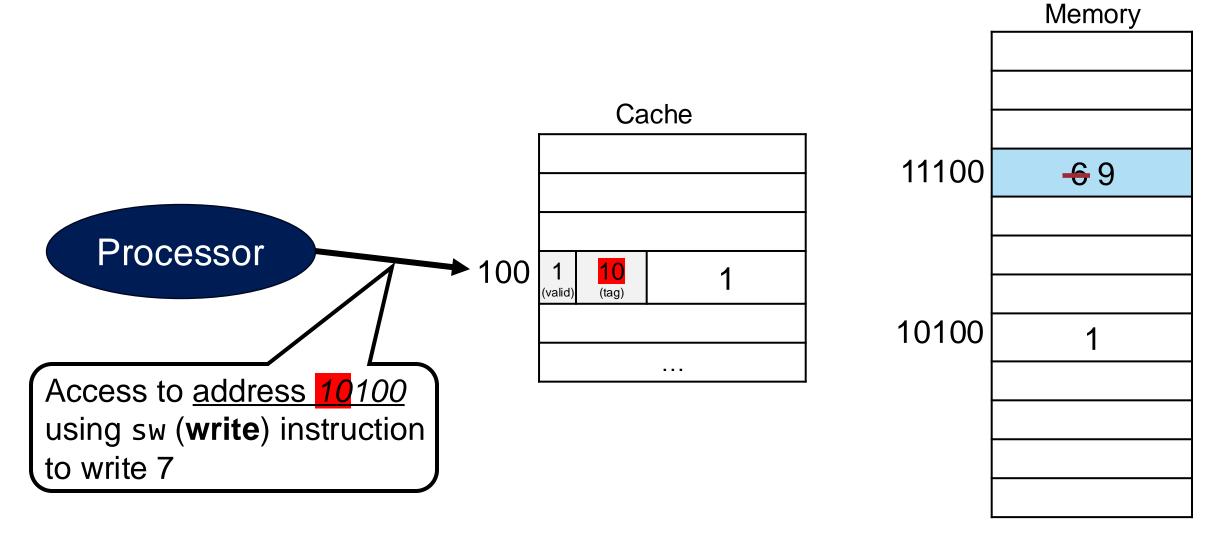




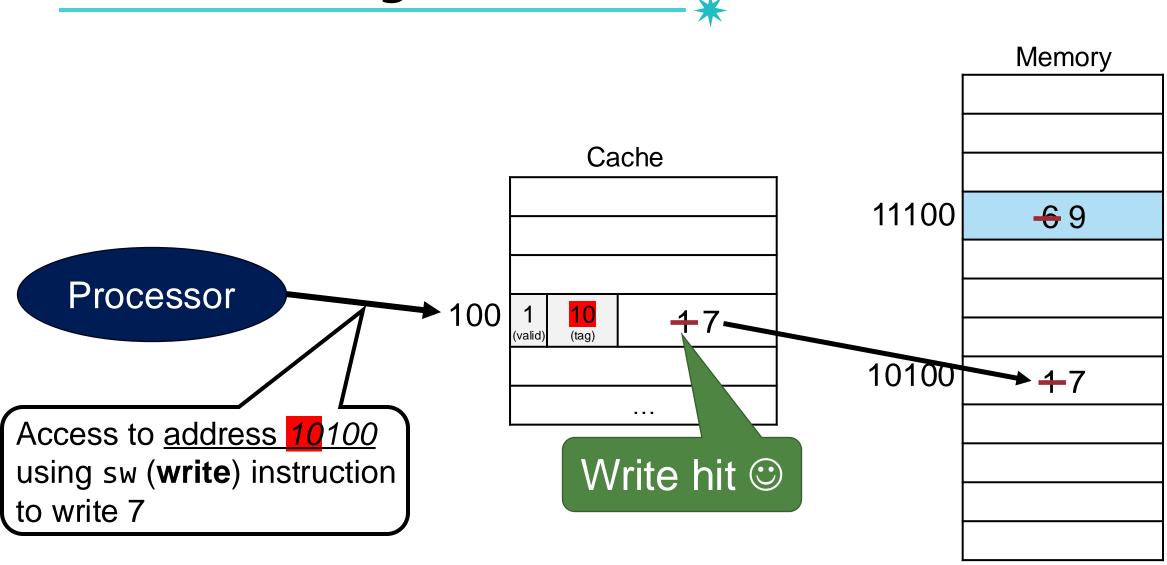




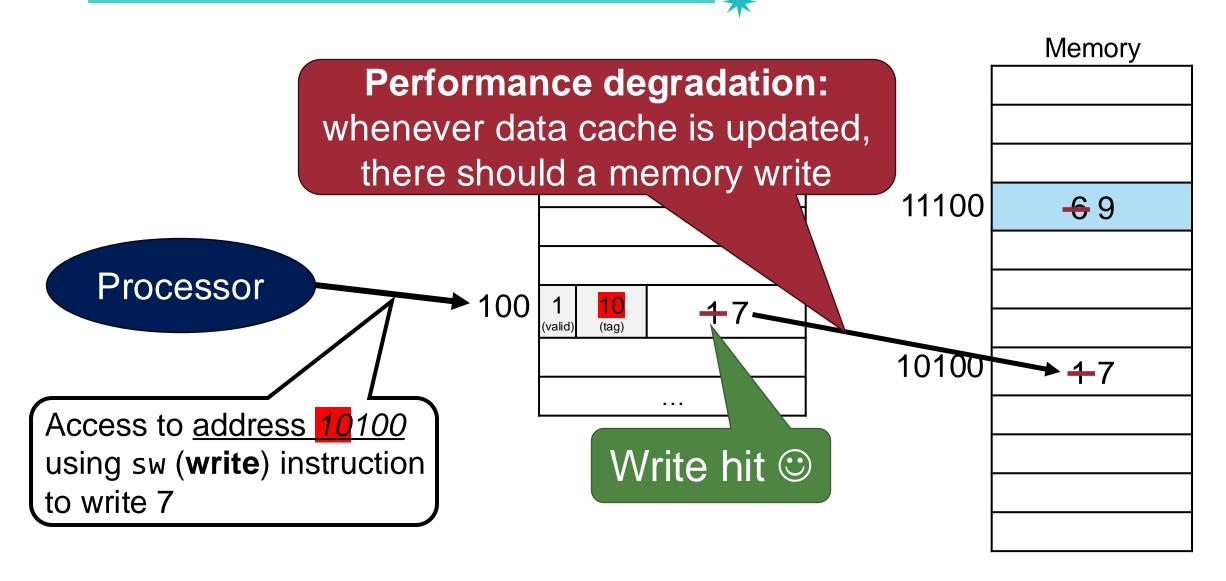




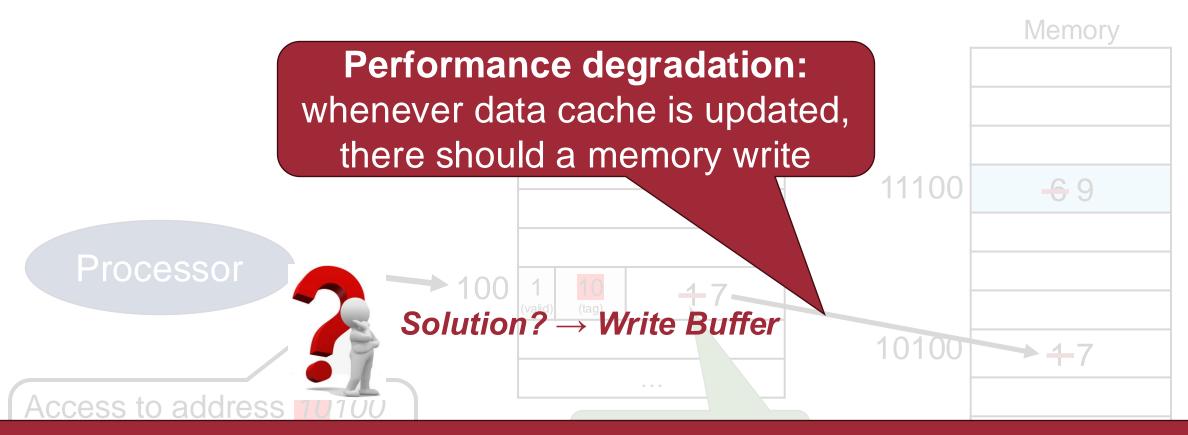




Main Disadvantage: Performance Degradation 28

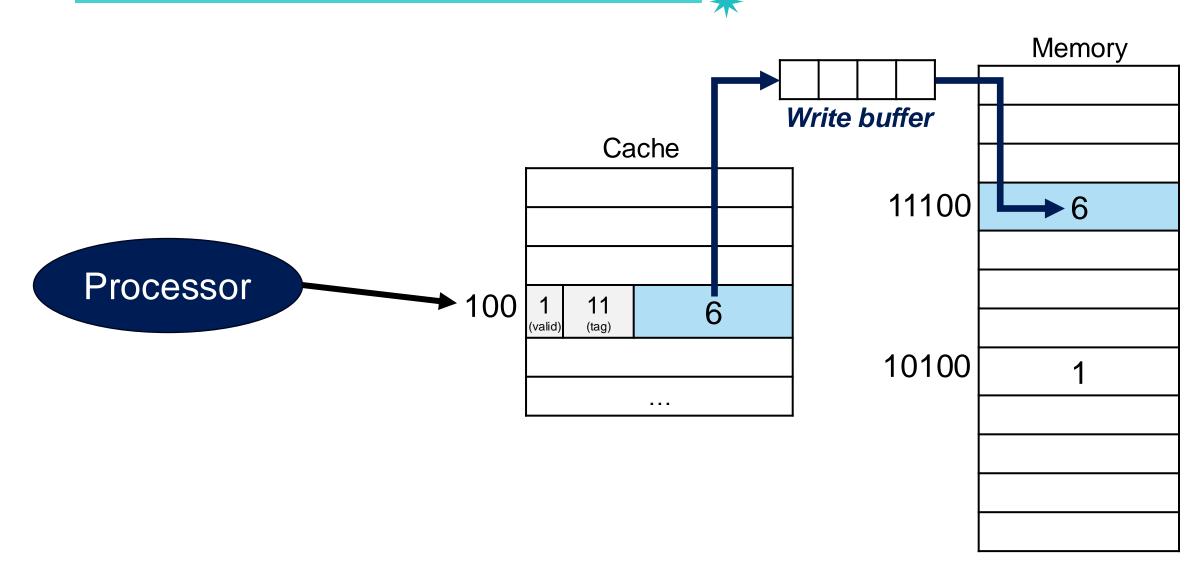


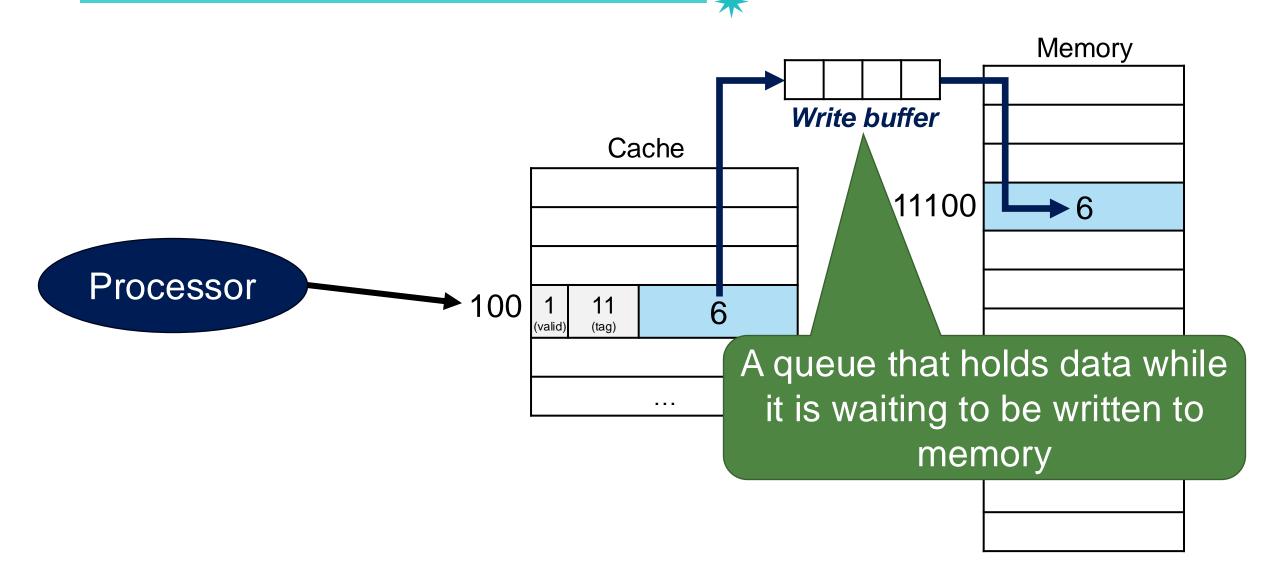
Main Disadvantage: Performance Degradation 29

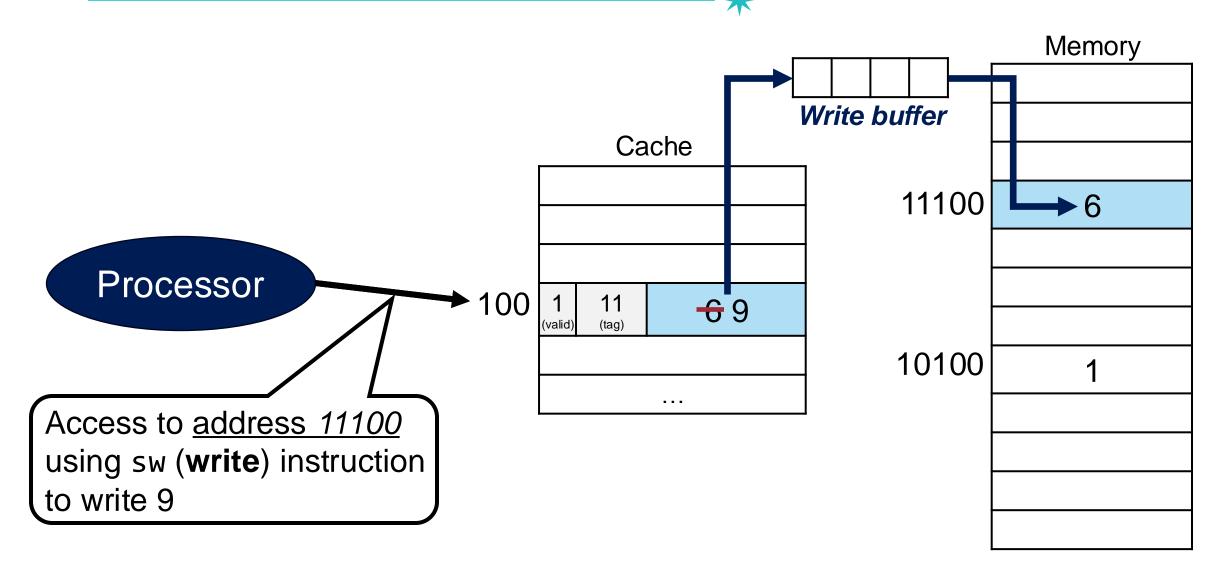


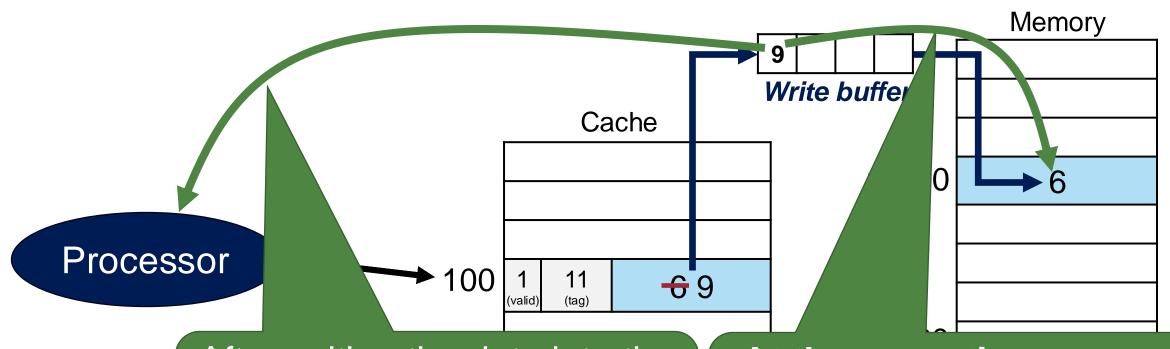
If base CPI = 1, 10% of instructions are stores, write to memory takes 100 extra cycles

CPI = 1 + 0.1 x 100 = 11









Access to a using sw (w to write 9

After writing the data into the write buffer, the processor can **continue execution**

At the same time, memory controller write contents to memory

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

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: On each write hit, update only the block in cache. The modified cache block is written to memory only when it is replaced
 - Keep track of whether each block is dirty (additional bit).

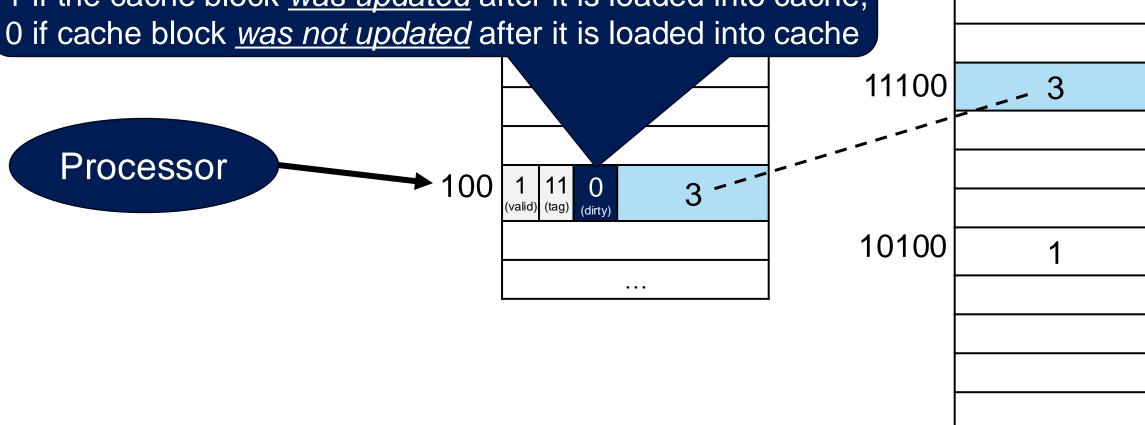
Memory

Write Back – Dirty Bit



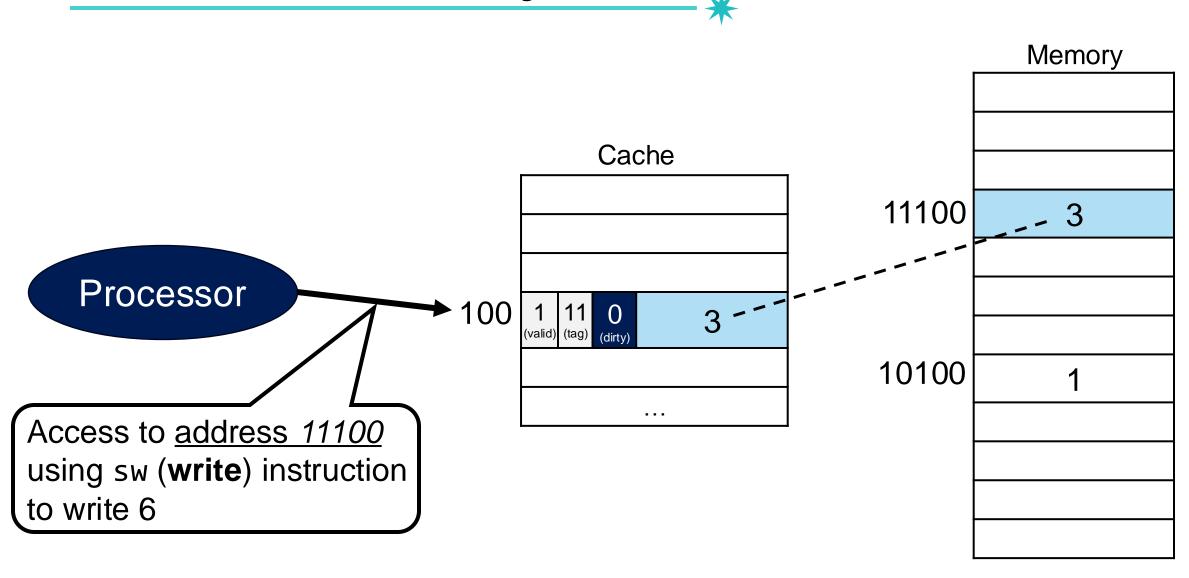


1 if the cache block was updated after it is loaded into cache,

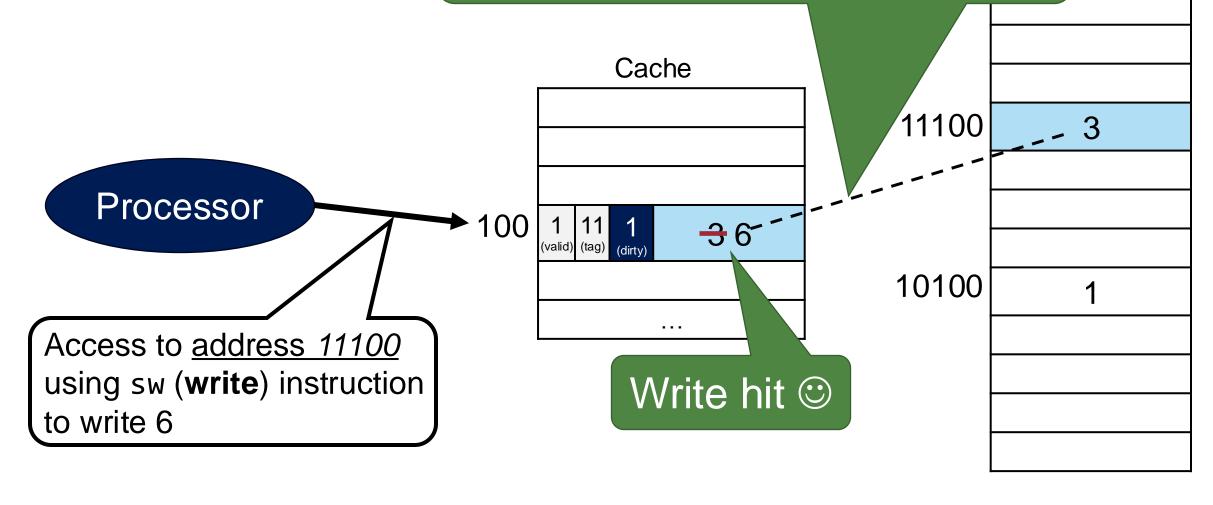


37

Write Back - Dirty Bit

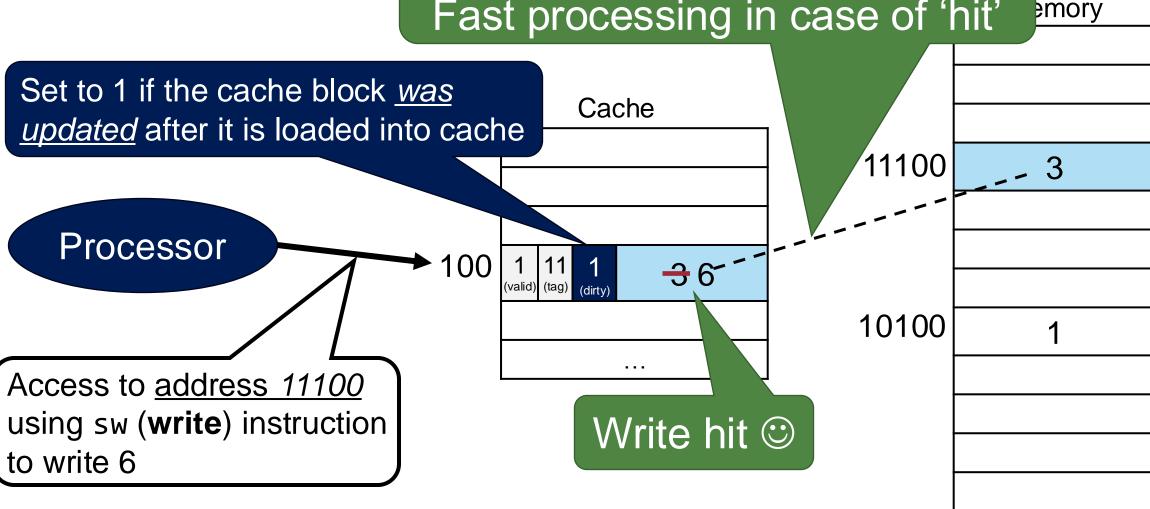


No writes to memory →
Fast processing in case of 'hit'

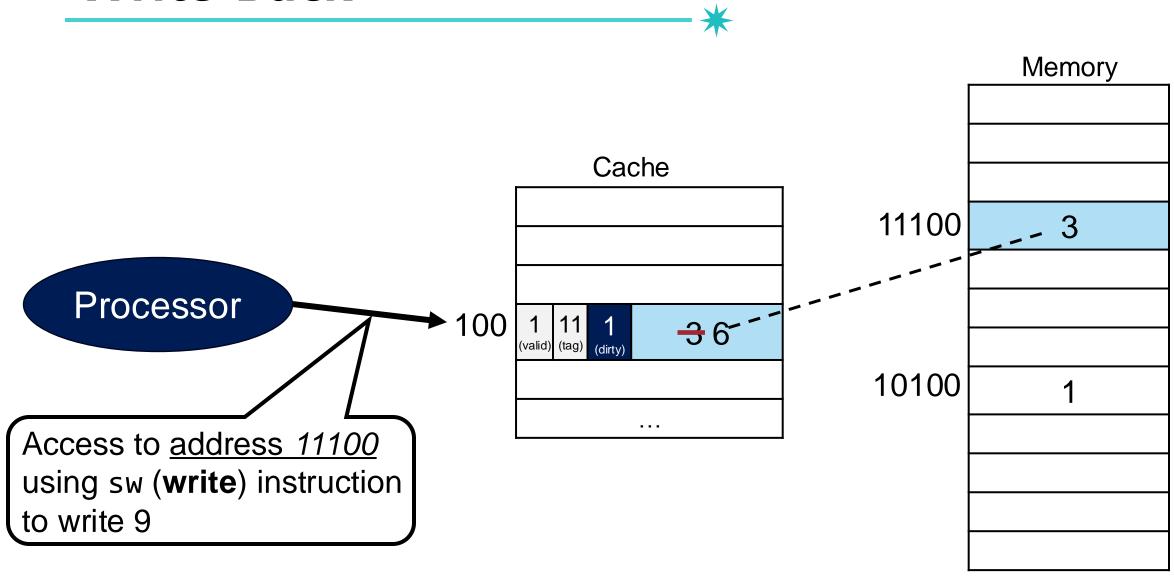


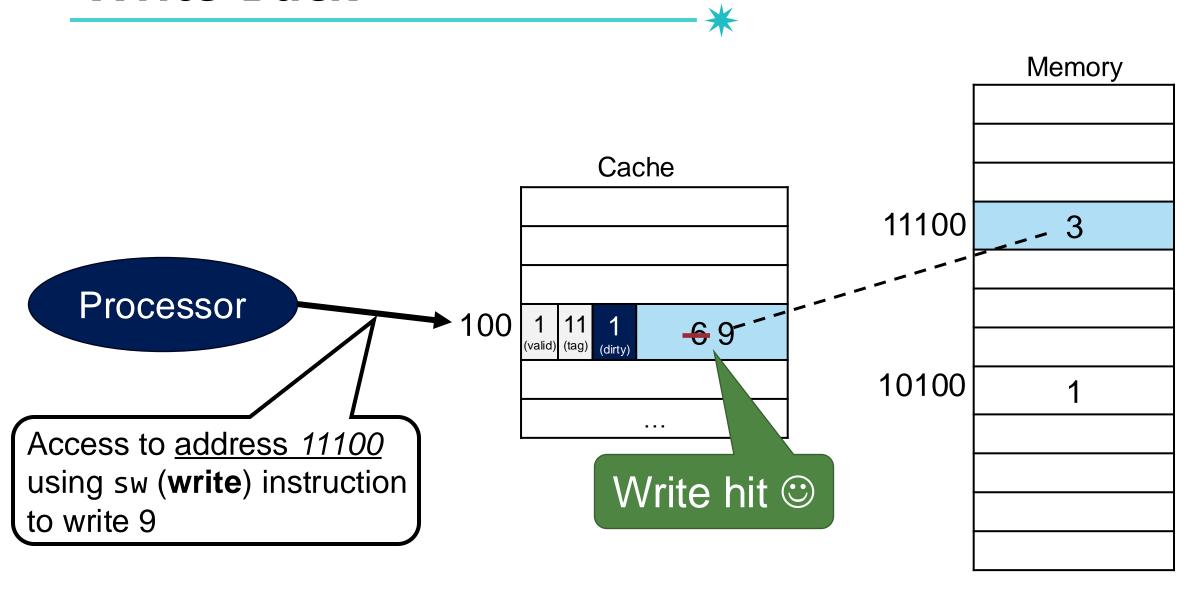


No writes to memory → Fast processing in case of 'hit'

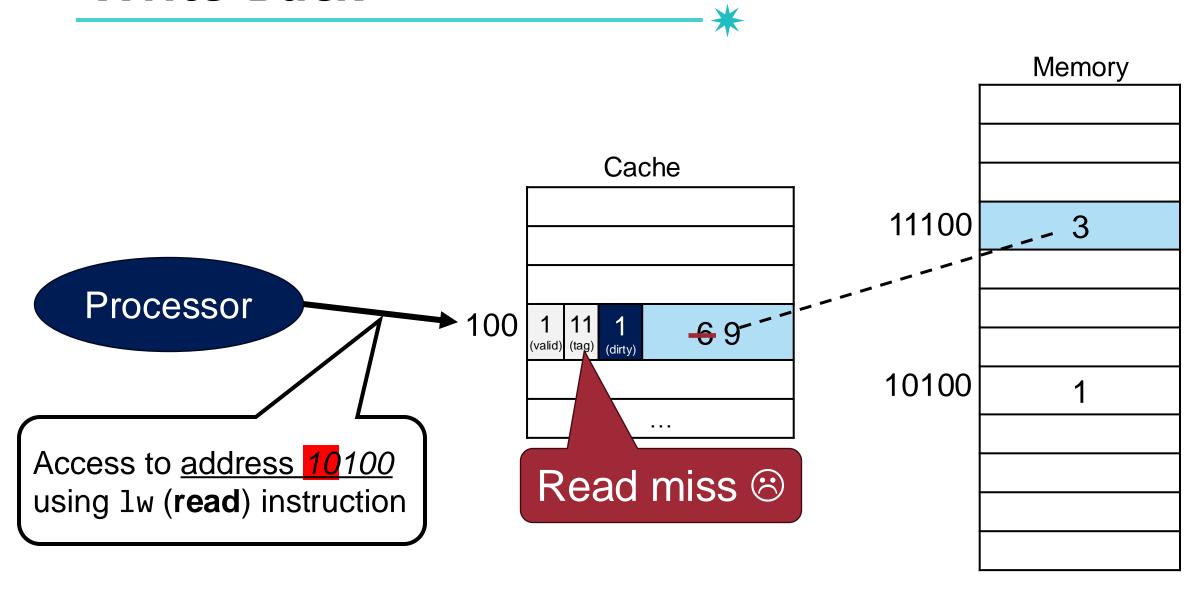




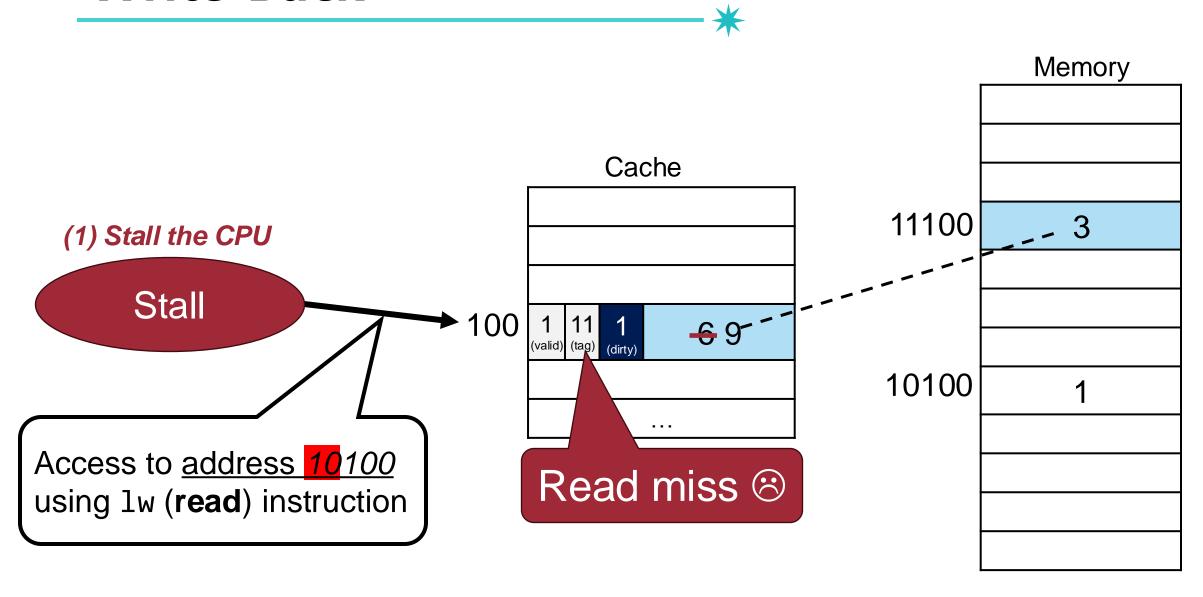




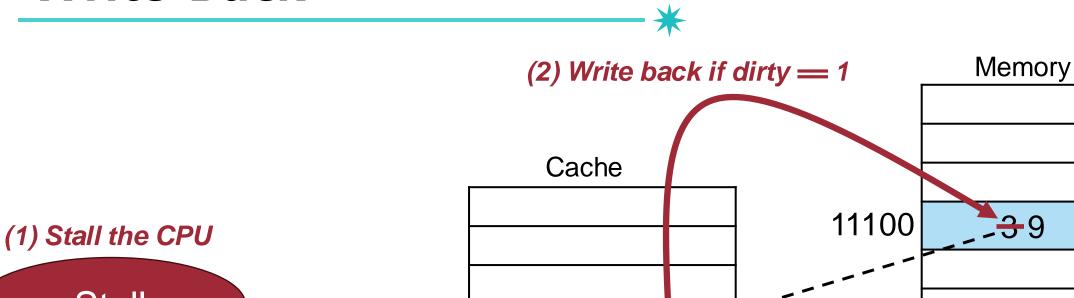




43







69

10100

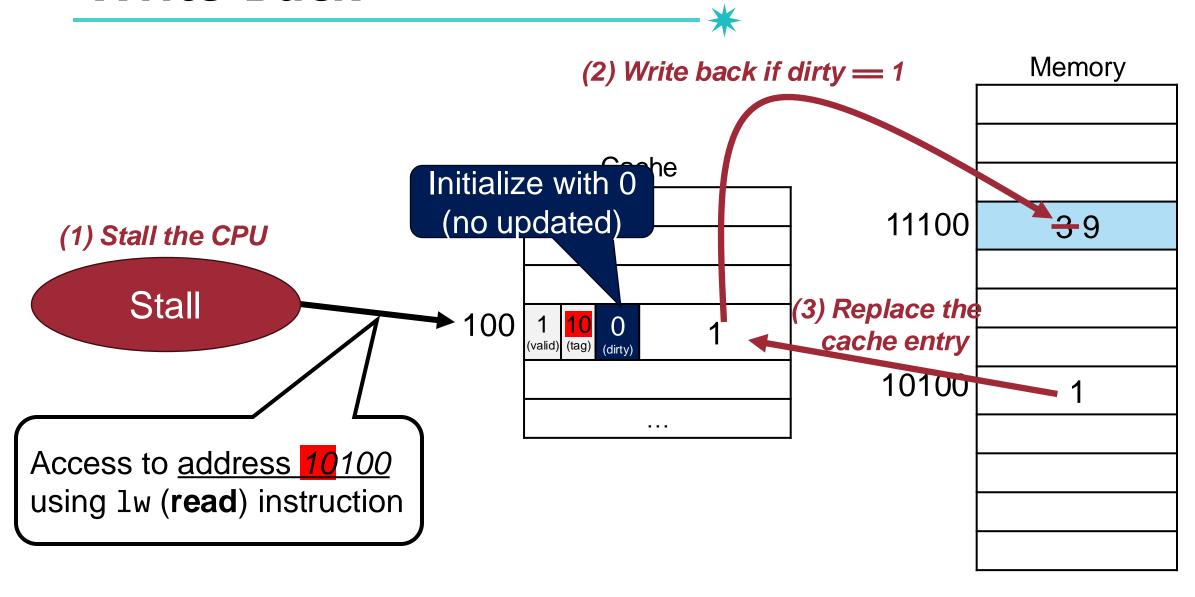
Access to <u>address</u> <u>10</u>100 using lw (**read**) instruction

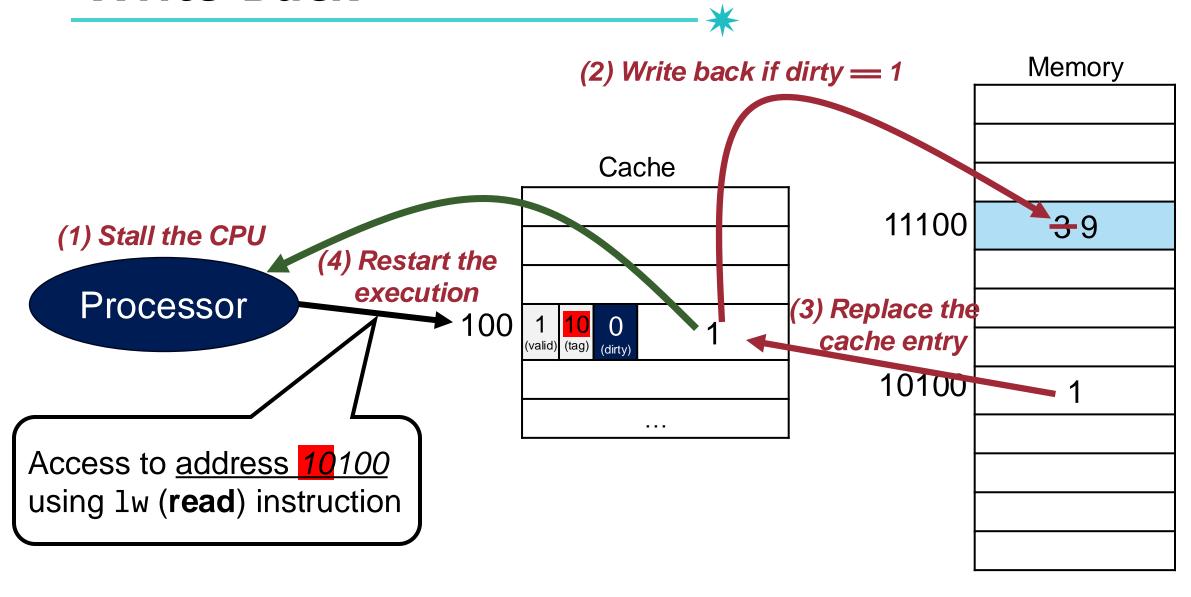
Stall

Read miss 🕾

100

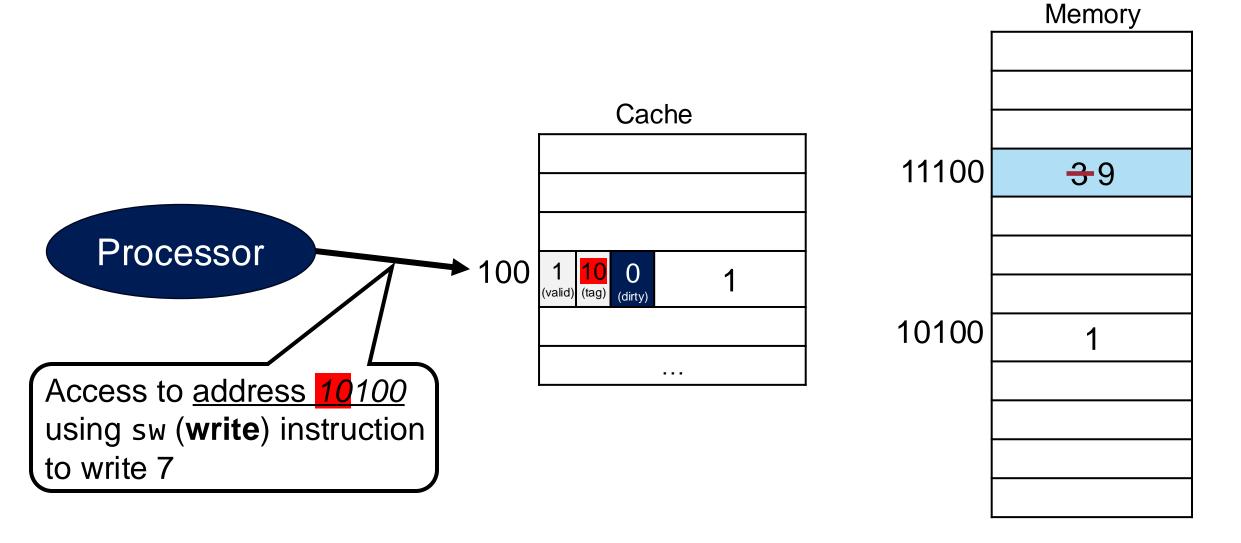
(valid) (tag) (dirty)



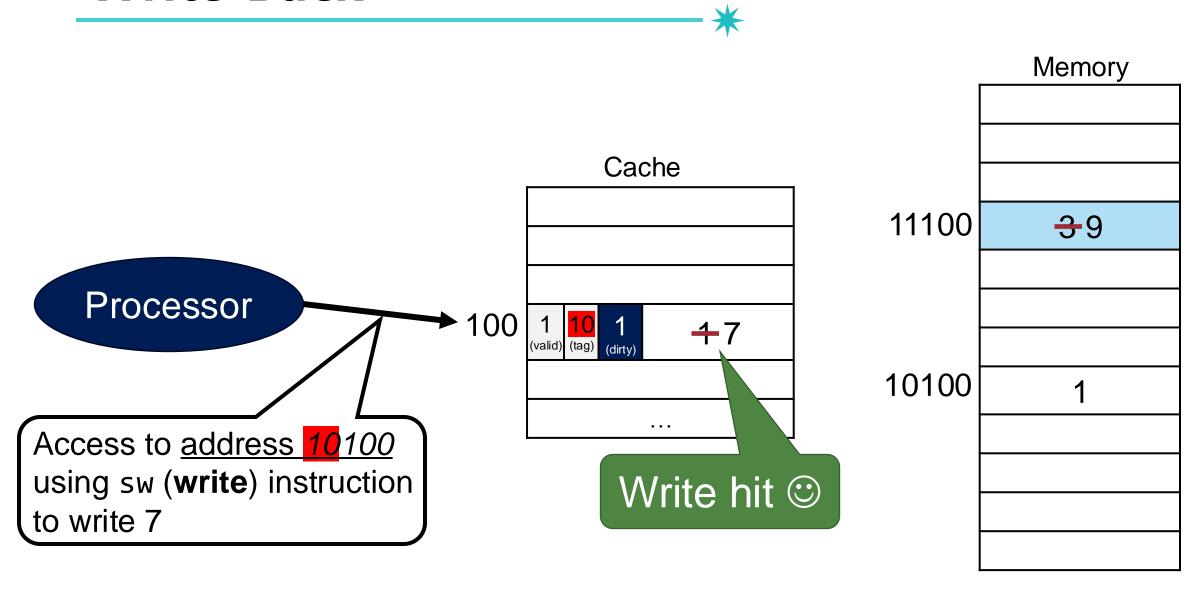








48



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: On each write hit, update only the block in cache. The modified cache block is written to memory only when it is replaced
 - Keep track of whether each block is dirty (additional bit).
 - Pros: faster processing in case of 'hit' (No repeated writes to memory)
 - Cons: slower in case of 'miss'

Summary: Hits vs. Misses

- Read hits
 - This is what we want!
- Read misses
 - Stall the CPU, fetch block from memory, restart
- Write hits
 - (option #1) Write the data both in cache and memory (write through)

Cases to consider

for the instruction cache

- (option #2) Write the data only into the cache (write back)

Write misses

Cases to consider for the data cache

5

Write Misses



- Write misses
 - Write allocate (common strategy)

- No write allocation:

Write Misses



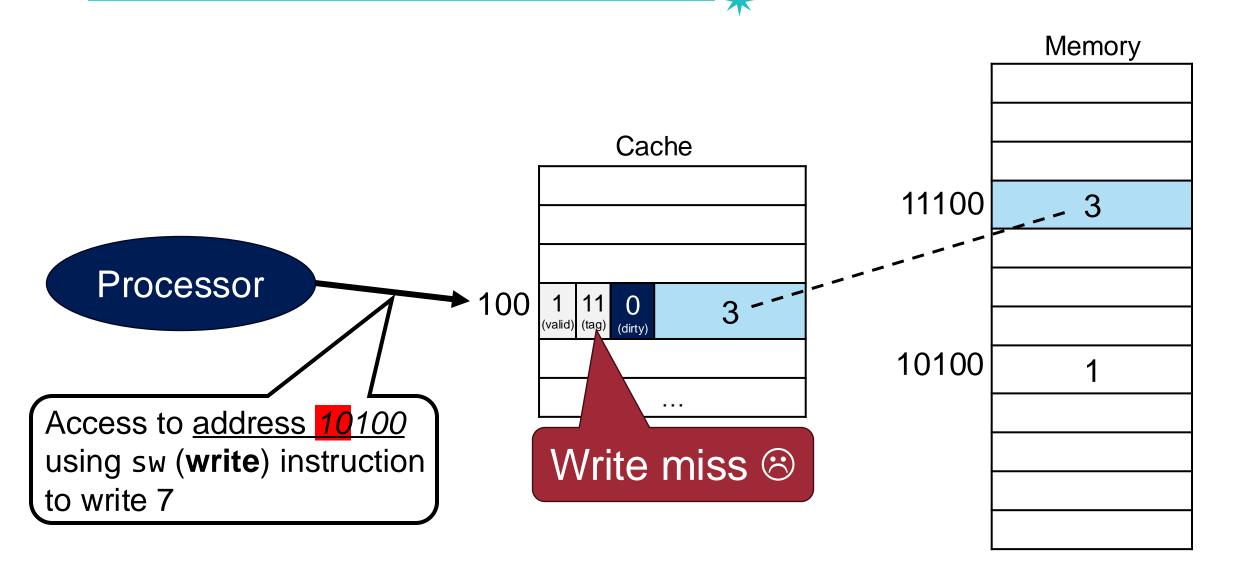


- Write misses
 - Write allocate (common strategy): allocate block on miss in the cache
 - The block is fetched from memory and then the appropriate portion of the block is overwritten (Similar with read miss)

- No write allocation:

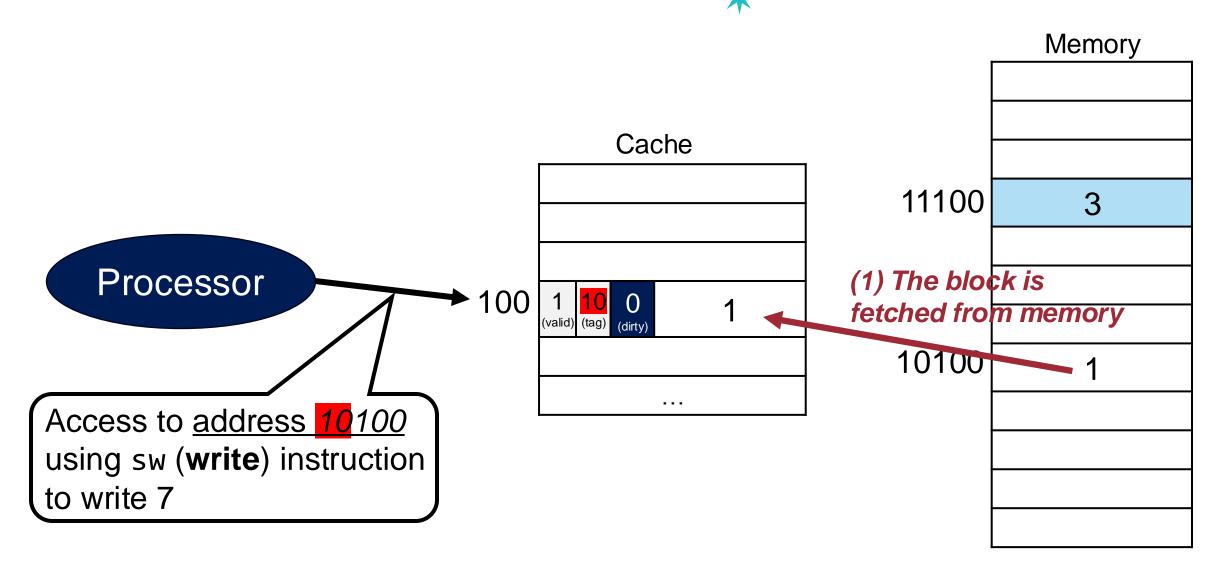
Write Misses: Write Allocate



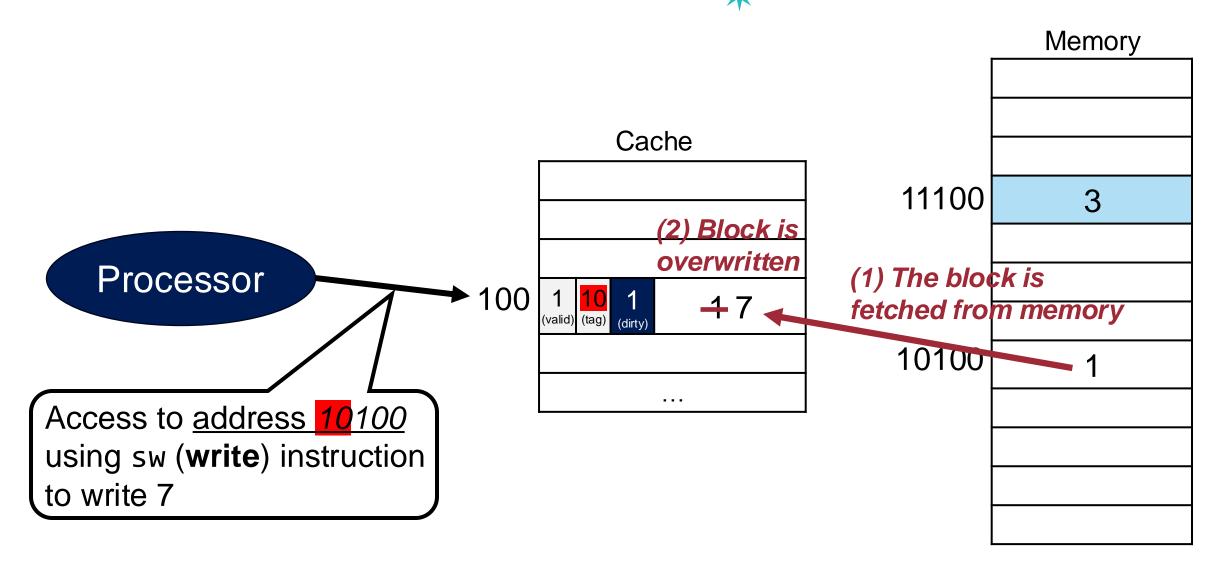


54

Write Misses: Write Allocate



Write Misses: Write Allocate



Write Misses



- Write misses
 - Write allocate (common strategy): allocate block on miss in the cache
 - The block is fetched from memory and then the appropriate portion of the block is overwritten (Similar with read miss)

 No write allocation: update the portion of the block in memory in direct manner but not put it in the cache

Summary: Cache Hits vs. Misses

- Read hits
 - This is what we want!
- Read misses
 - Stall the CPU, fetch block from memory, restart
- Write hits
 - (option #1) Write the data both in cache and memory (write through)
 - (option #2) Write the data only into the cache (write back)
- Write misses
 - (option #1) Write allocate
 - (option #2) Write no allocate

Cases to consider for the instruction cache

> Cases to consider for the data cache

Improving Cache Performance

How to Improve the Cache Performance? 59

Adjust the block size

Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches

How to Improve the Cache Performance? 60

Adjust the block size

Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches

6

Data

Problem with a One-word Block?

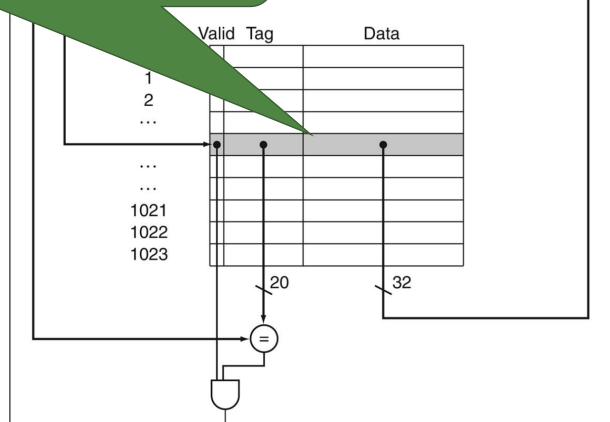
Address (showing bit positions) 31 30 · · · 13 12 11 · · · 2 1 0

We can exploit temporal locality:

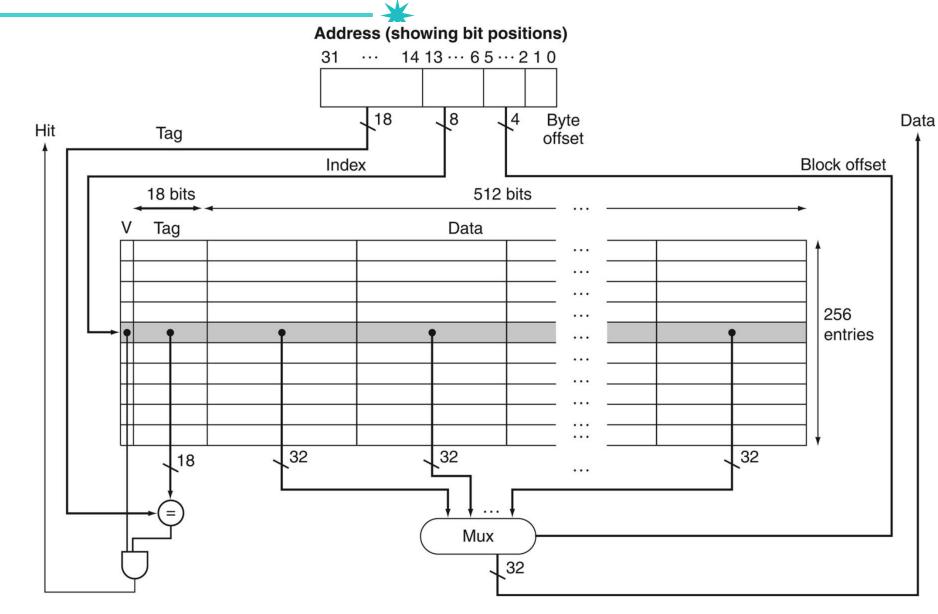
If the same address is re-referenced,
a cache hit occurs



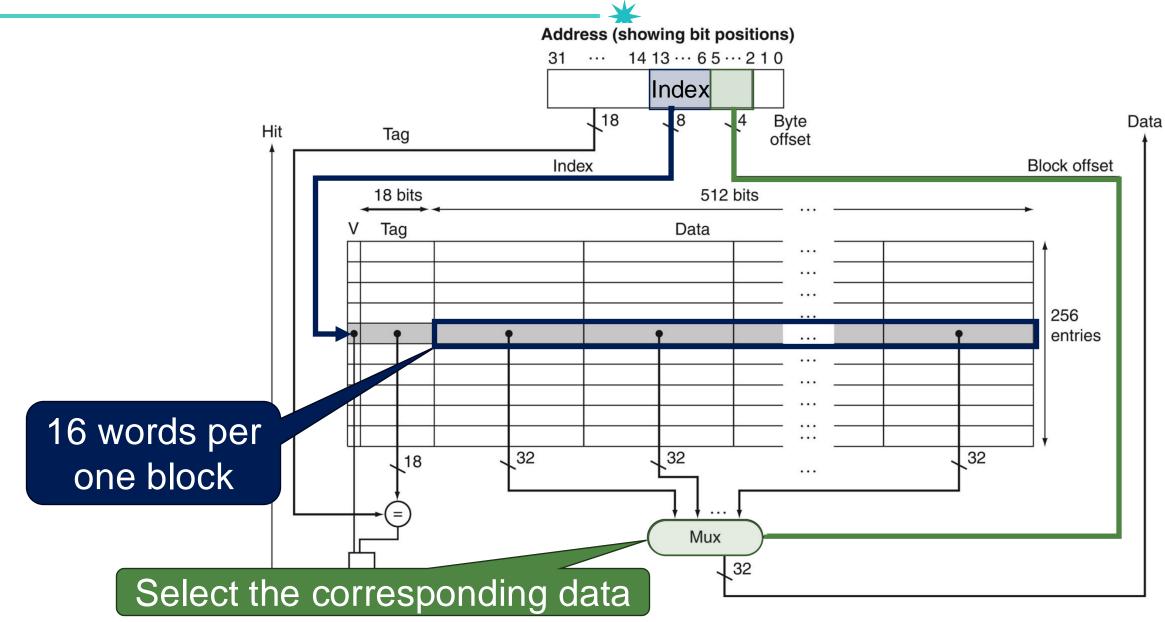
How about spatial locality?



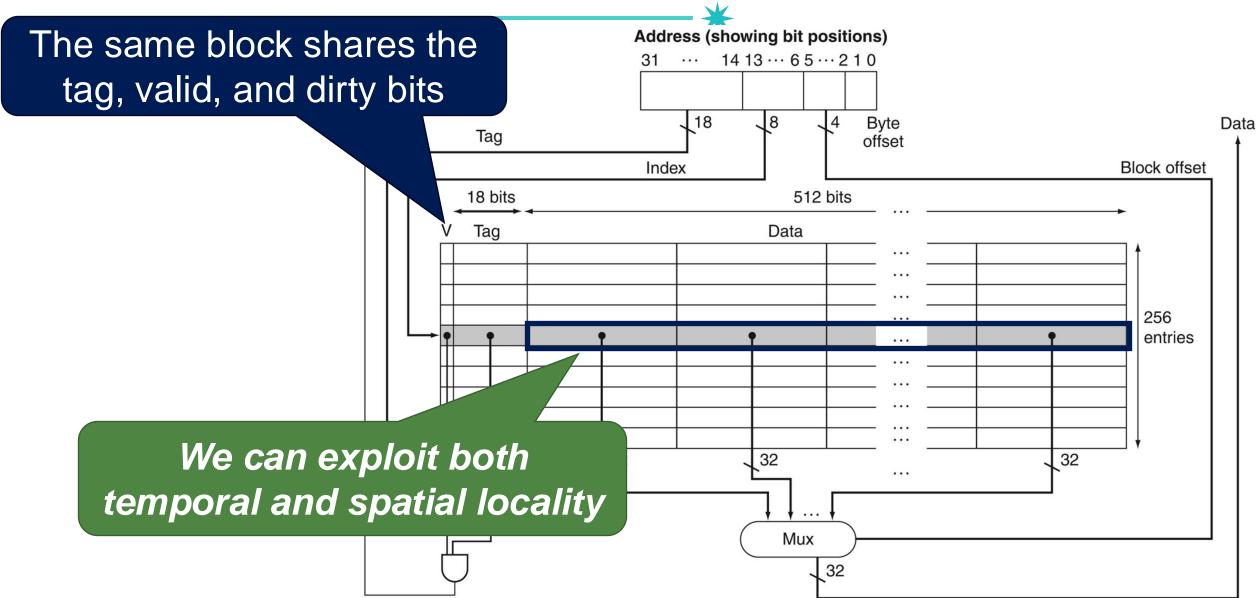




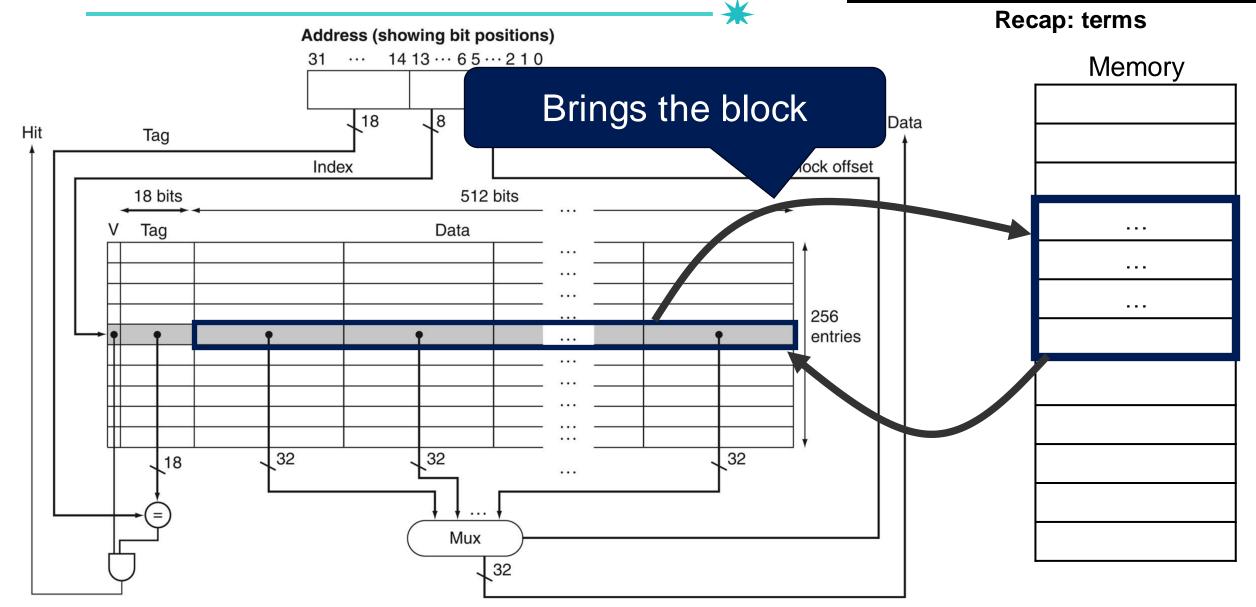




64



- Block (a.k.a., line): unit of copying
 - Several words in cache memory



Increasing the Block Size

- Pros: tends to decrease miss rate
 - Due to the spatial locality

- **Hit**: data requested is in the upper level
 - Hit ratio: hits/accesses
- Miss: data requested is not in the upper level
 - Block copied from lower level
 - Miss penalty: time taken to resolve miss
 - Miss ratio: misses/accesses

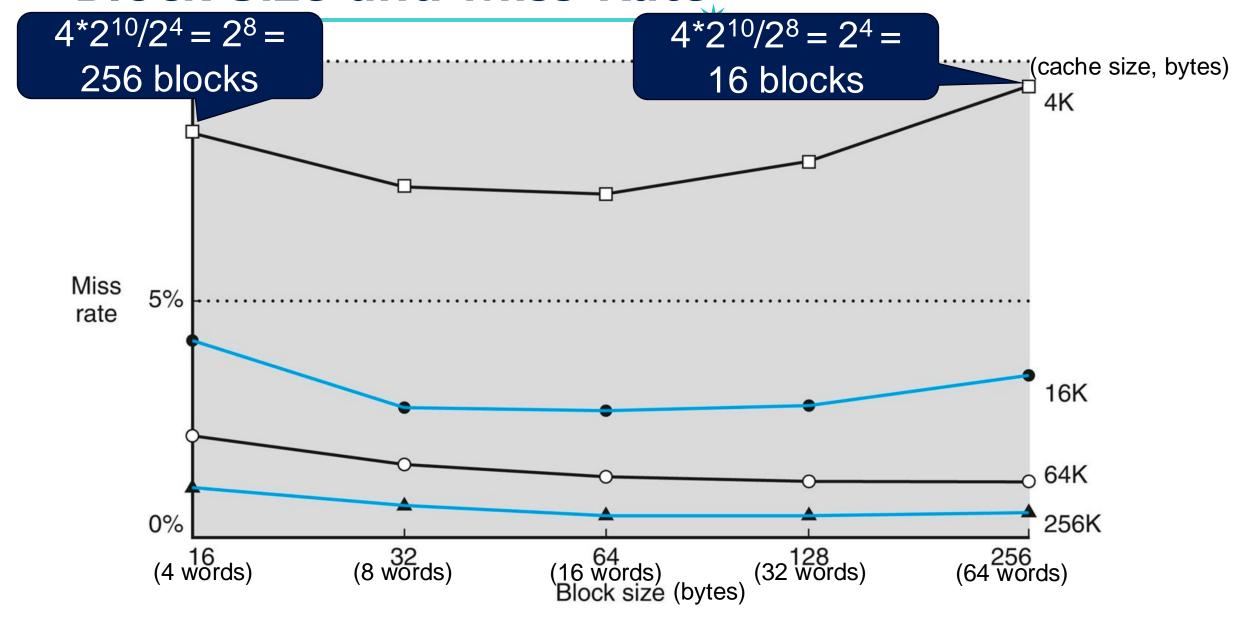
= 1 - hit ratio

Recap: terms

Cons

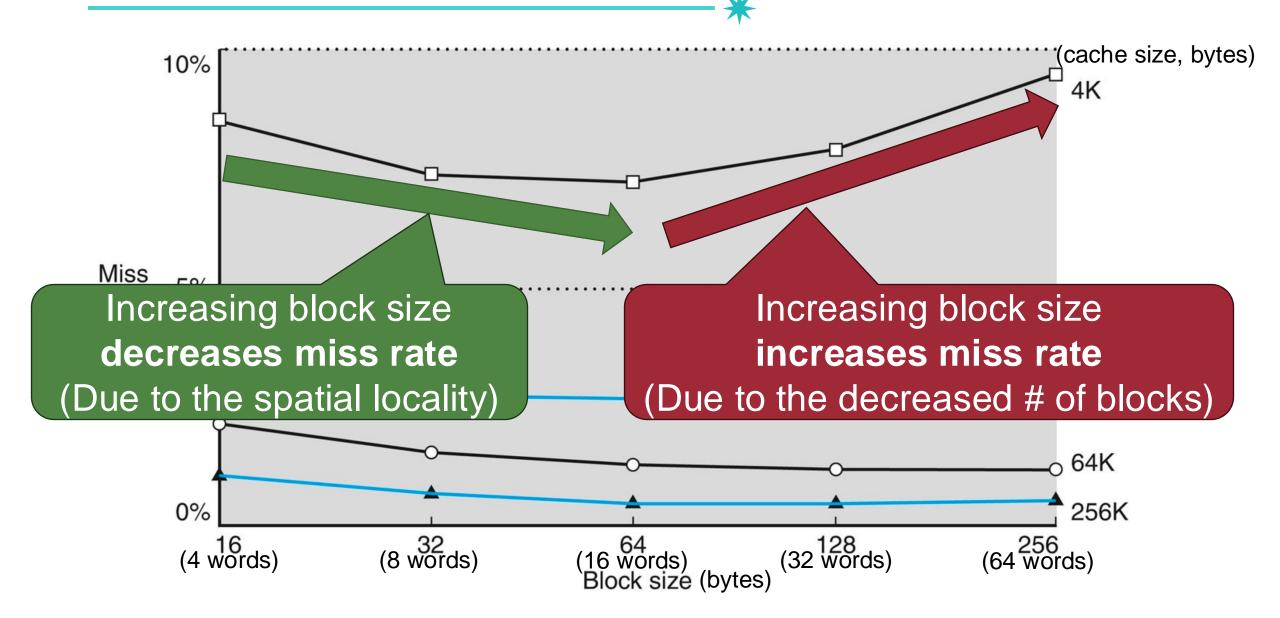
- However, if the block size is significantly increased (while keeping the total cache size fixed), the number of indexes (cache lines) decreases, which can actually lead to an increase in the miss rate
- Also, the miss penalty increases because transferring such a large block takes more time

Block Size and Miss Rate



Block Size and Miss Rate





Cache Performance



Example:

- -I-cache miss rate = 2%
- -D-cache miss rate = 4%
- -Miss penalty = 100 cycles
- -Base CPI (ideal cache) = 2
- -Load & stores are 36% of instructions

· Miss cycles (miss penalty) per instruction

- -1-cache: $0.02 \times 100 = 2$
- -D-cache: $0.36 \times 0.04 \times 100 = 1.44$
- Actual CPI = 2 + 2 + 1.44 = 5.44

How to Improve the Cache Performance?

Adjust the block size

Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches

How to Improve the Cache Performance?

Adjust the block size

Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches





Cache placement policy: where can a block be placed?

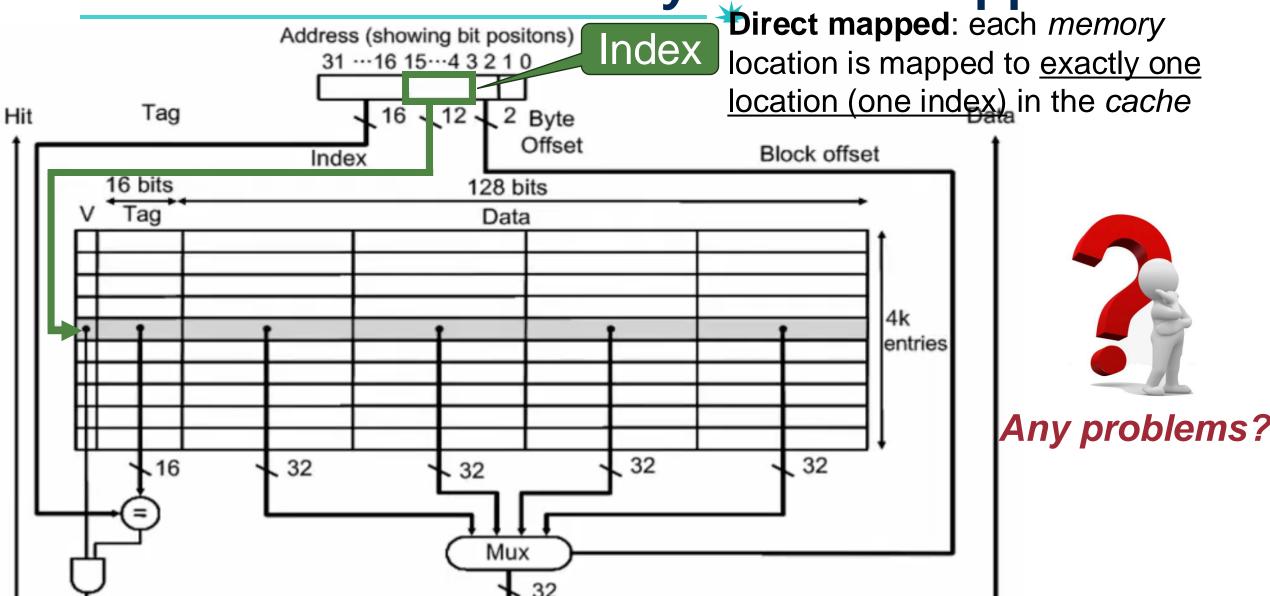
- Direct mapped
- -Fully associative
- -Set associative

Cache Placement Policy: Direct Mapped

73

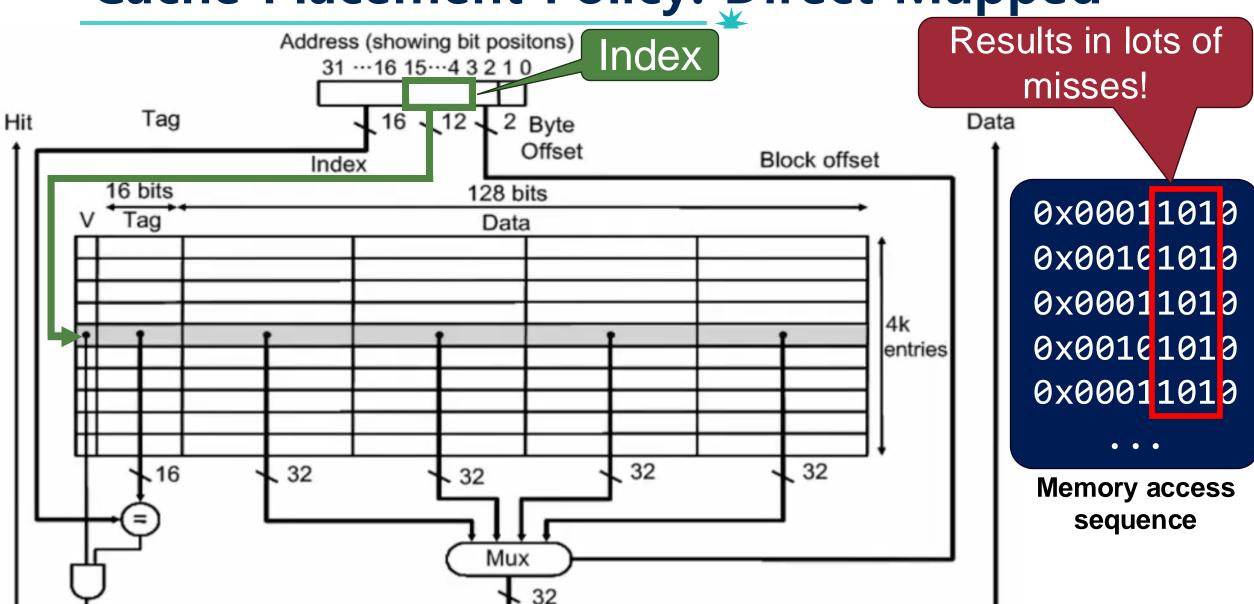
Address is modulo the number of blocks in the cache

Cache Placement Policy: Direct Mapped



Cache Placement Policy: Direct Mapped

75



76

Direct Mapped: Summary

Address is modulo the number of blocks in the cache

- **Disadvantage**: imagine two frequently-accessed variables. What if their addresses have the same index bits?
 - Such addresses are in conflict in the cache
 - Can't hold both in the cache at once
 - Result in lots of misses (bad!)

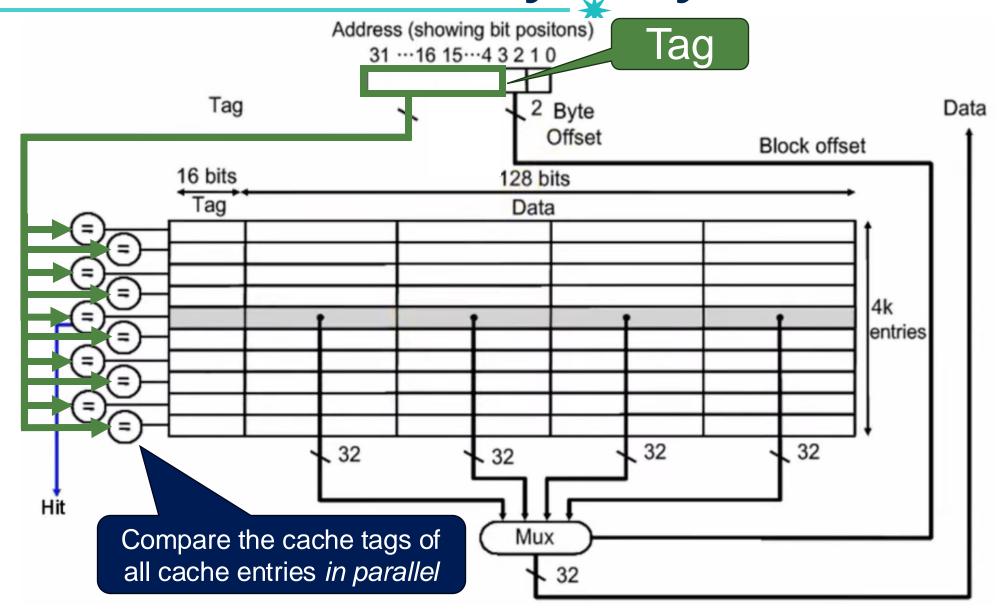
Cache Placement Policy: Fully Associative

- No cache index
- Block can be placed in any location in the cache
- Requires all entries to be searched at once (compare <u>the cache</u> tags of all cache entries in parallel)

Tag Block offset Byte offset

The index field does not exist

Cache Placement Policy: Fully Associative ®



Fully Associative: Summary

79

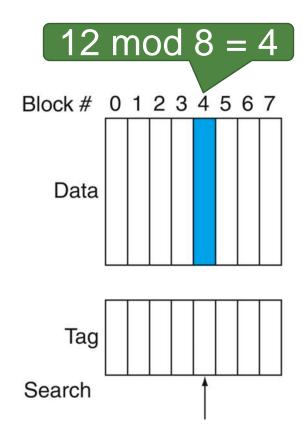
- No cache index
- Block can be placed in any location in the cache
- Requires all entries to be searched at once (compare the cache tags of all cache entries in parallel)

Tag Block offset Byte offset

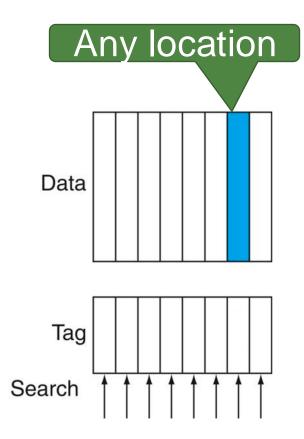
The index field does not exist

Disadvantage: Comparator per entry (expensive hardware)

Associativity – Access to Address 12



Direct mapped

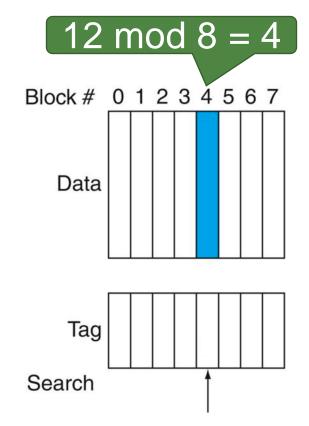


Fully associative

Increasing associativity (shrink index & expand tag)
Reducing miss rate by more flexible placement

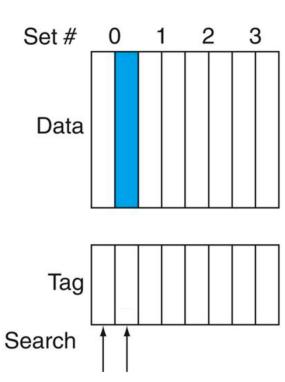
Introduction to Set-associative Cache



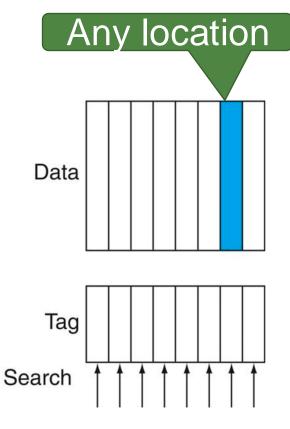


Direct mapped

Miss rate ↑



Set-associative



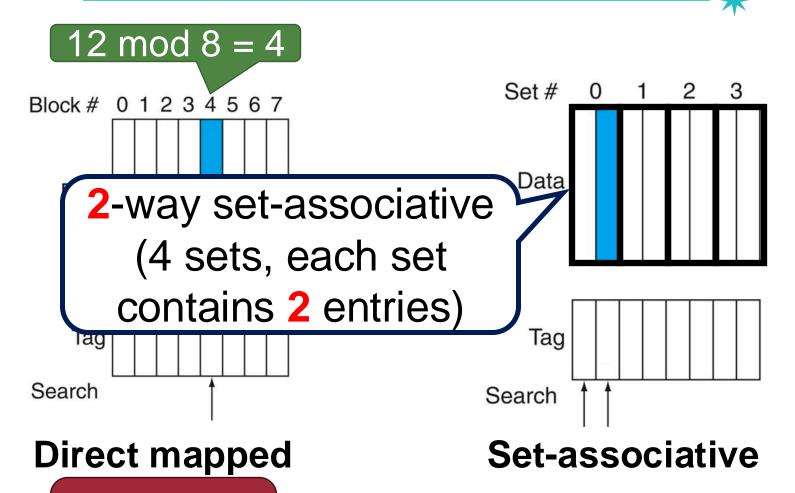
Fully associative

easing associativity (shrink index & expand tag Reducing miss rate by more flexible placement

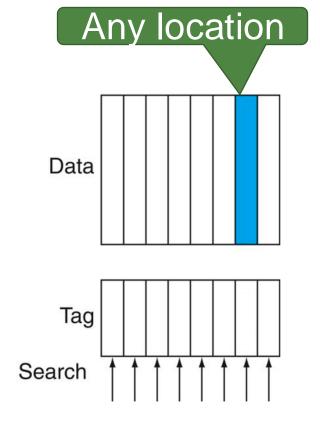
HVV complexity_1

Introduction to Set-associative Cache





Miss rate ↑



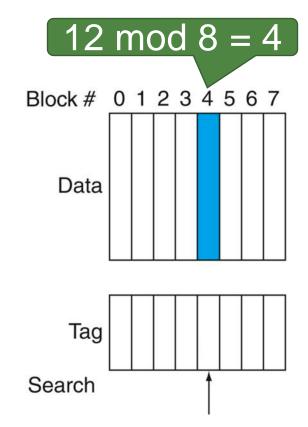
Fully associative

easing associativity (shrink index & expand tag Reducing miss rate by more flexible placement

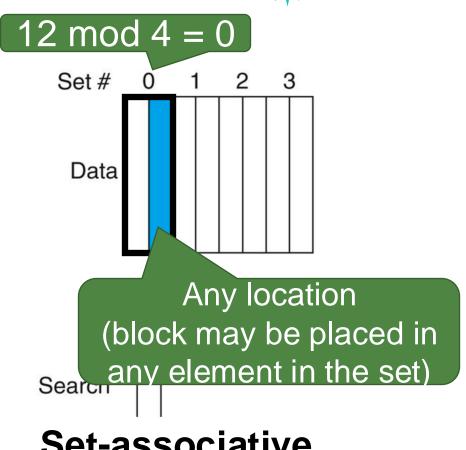
HVV complexity ↑

Introduction to Set-associative Cache

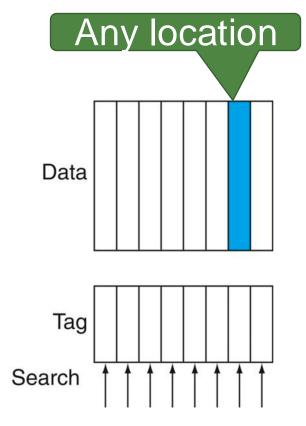








Set-associative



Fully associative

Miss rate ↑

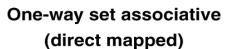
easing associativity (shrink index & expand tag Reducing miss rate by more flexible placement

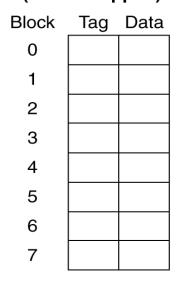
complexity_

Cache Placement Policy: N-way Set Associative

- Allow a given block to go in a target set
 - Position of memory block = (block #) mod (# of sets in the cache)
- Each set contains N entries
- Search all entries in a given set at once
 - N comparators (less expensive)

Spectrum of Associative





Set

Two-way set associative

Set	Tag	Data	Tag	Data
0				
1				
2				
3				

4 way

2 sets

Four-way set associative

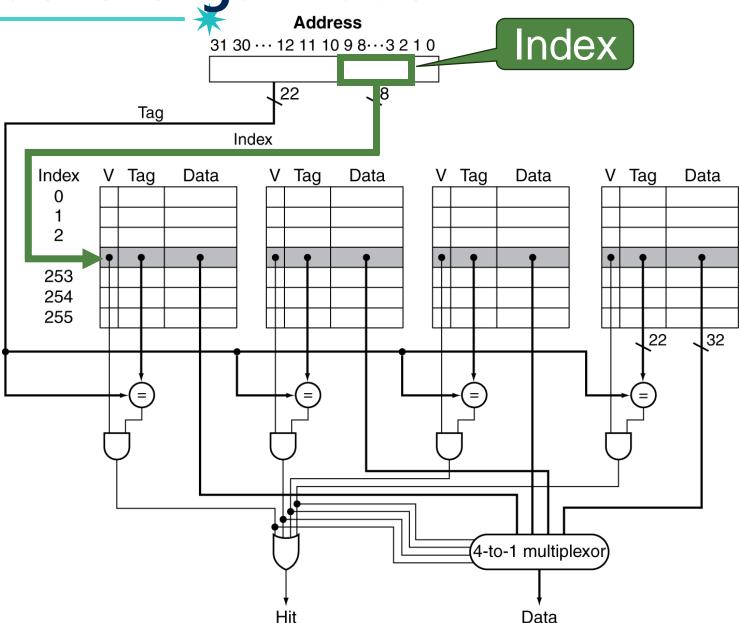
	Tag	Data	Tag	Data	Tag	Data	Tag	Dana
Ī								

Eight-way set associative (fully associative)

_Ta	ag	Data	Tag	Data												

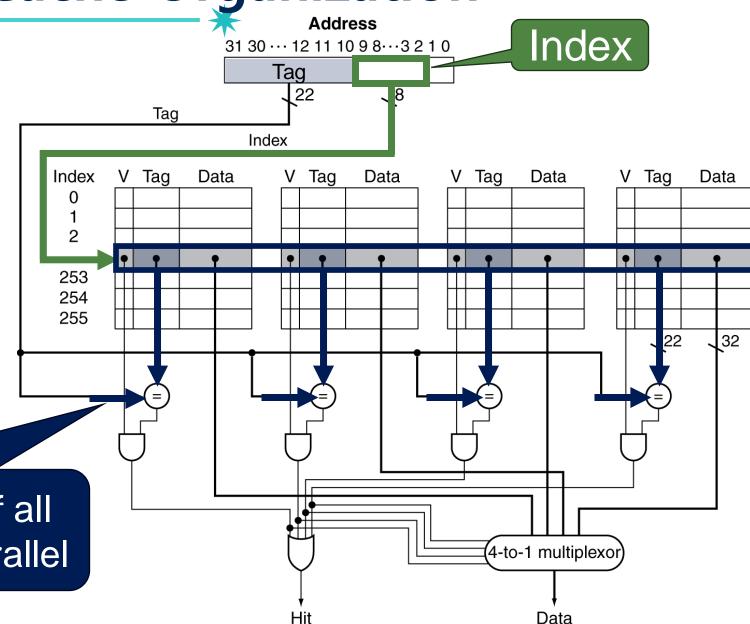
Set Associative Cache Organization

- 1024 blocks
- 4-way set associativity
- 1 block = 1 word



Set Associative Cache Organization

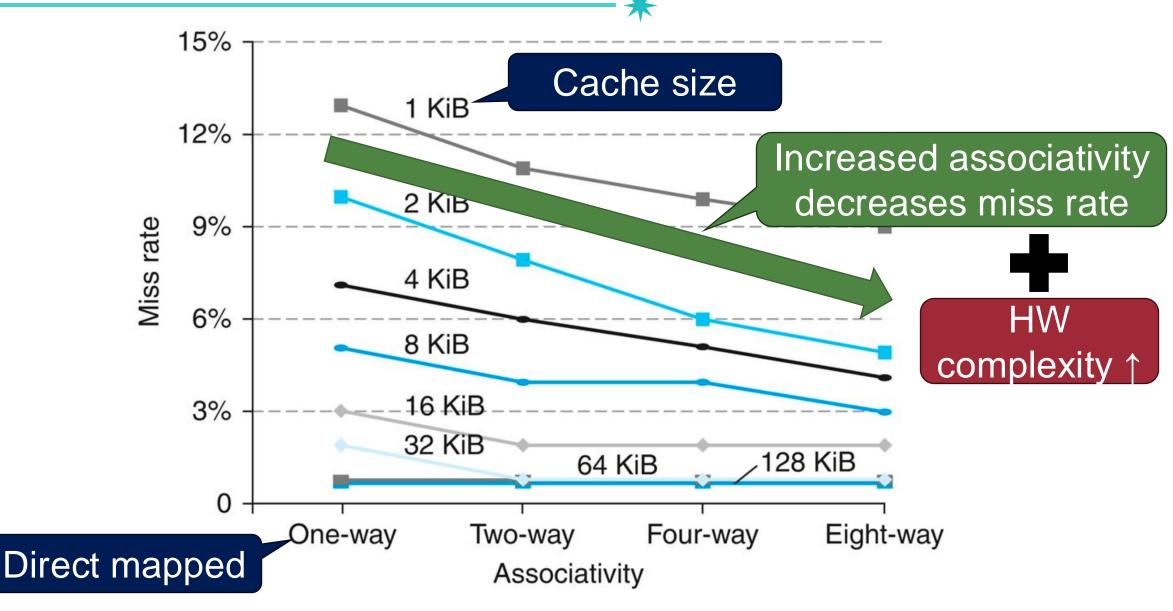
- 1024 blocks
- 4-way set associativity
- 1 block = 1 word



Compare the cache tags of all entries in a given set in parallel









Summary: Cache Placement Policy

Cache placement policy: where can a block be placed?

- Direct mapped
- -Fully associative
- -Set associative

How to Improve the Cache Performance? 91

Adjust the block size

Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches

How to Improve the Cache Performance?

Adjust the block size

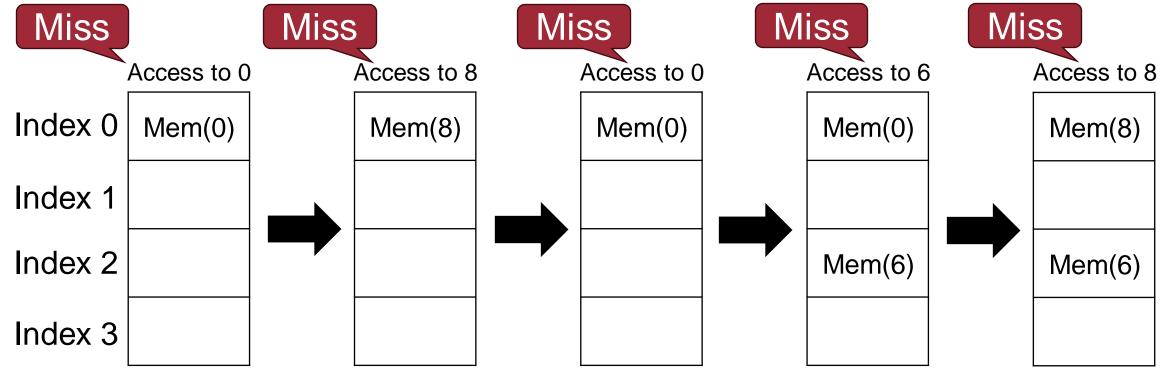
Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches

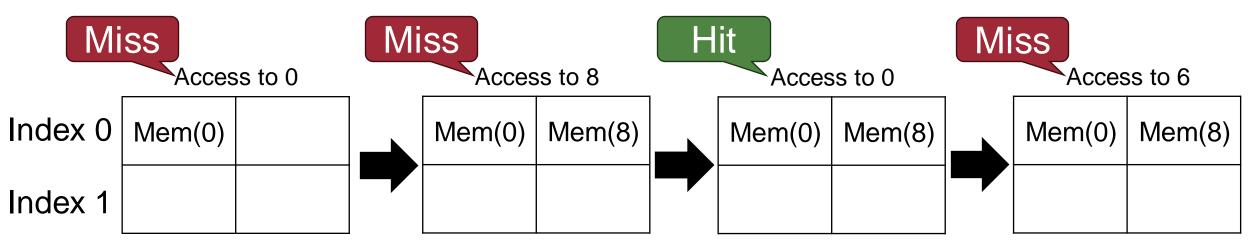
- Block access sequence: 0, 8, 0, 6, 8
- 4 Blocks
- Block size: 4 bytes (1 word)

- Block access sequence: 0, 8, 0, 6, 8
- 4 Blocks
- Block size: 4 bytes (1 word)



Direct mapped

- Block access sequence: 0, 8, 0, 6, 8
- 4 Blocks
- Block size: 4 bytes (1 word)

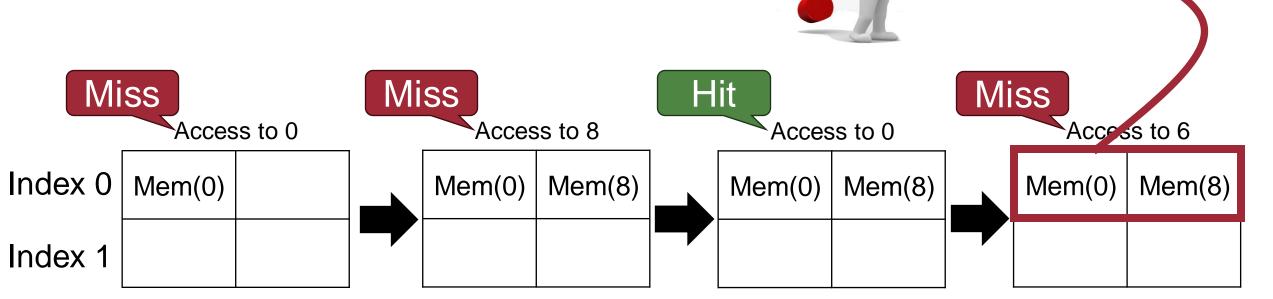


2-way set associative

Which block should be

replaced?

- Block access sequence: 0, 8, 0, 6, 8
- 4 Blocks
- Block size: 4 bytes (1 word)



2-way set associative

Which Block Should be Replaced on a Cache Miss?

- Direct mapped: no choice
- Set associative & fully associative: follow the replacement policy

98

Replacement Policy



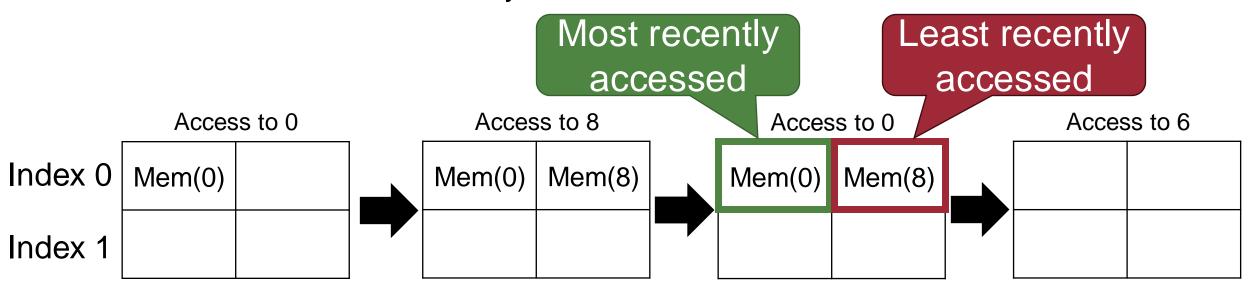
 Least Recently Used (LRU): replace the one NOT used (accessed) for the <u>longest time</u>

99

Replacement Policy



- Least Recently Used (LRU): replace the one NOT used (accessed) for the longest time
 - Temporal locality of access is considered
 - Need a reference history information

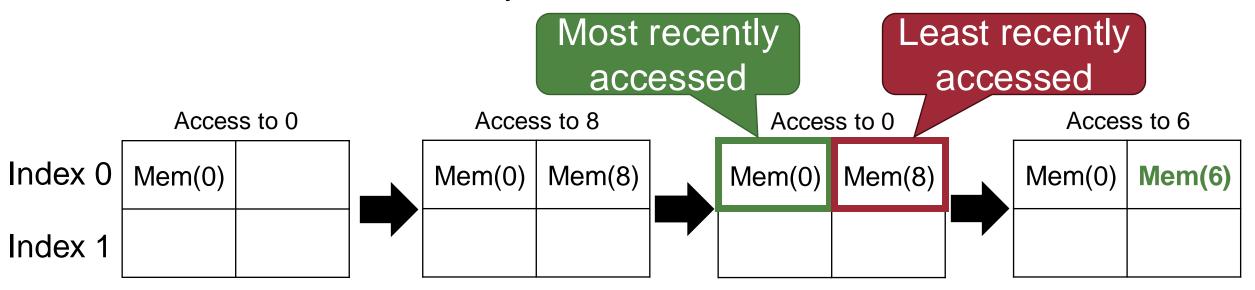


2-way set associative

Replacement Policy

10

- Least Recently Used (LRU): replace the one NOT used (accessed) for the <u>longest time</u>
 - Temporal locality of access is considered
 - Need a reference history information



2-way set associative

Replacement Policy





- Least Recently Used (LRU): replace the one NOT used (accessed) for the <u>longest time</u>
 - Temporal locality of access is considered
 - Need a reference history information

Random

- Gives approximately the same performance as LRU for high associativity
- + First In First Out (FIFO): replace the block with long-lived in the cache
- + Least Frequently Used (LFU): replace the block with fewest reference

How to Improve the Cache Performance?

Adjust the block size

Adjust the cache placement policy

Adjust the cache replacement policy

Multilevel caches

How to Improve the Cache Performance?

Adjust the block size

Adjust the cache placement policy

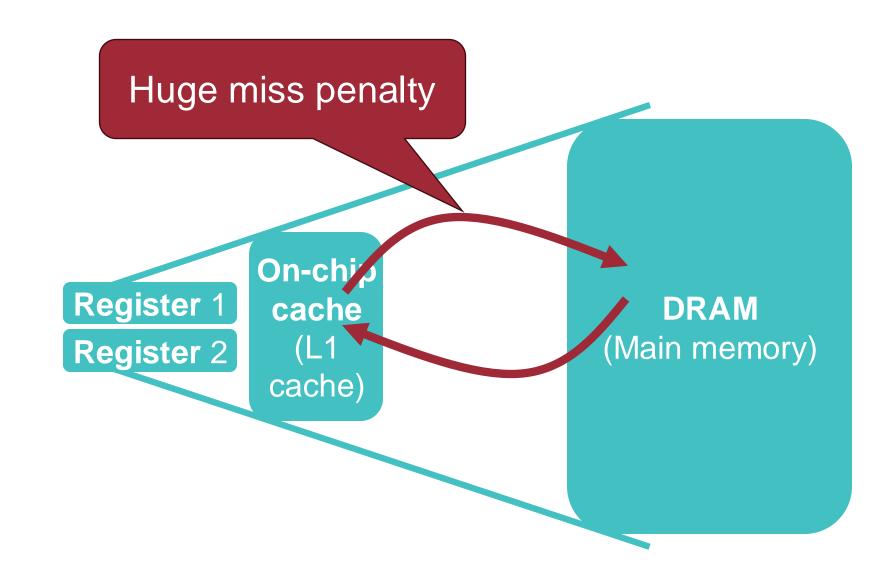
Adjust the cache replacement policy

Multilevel caches

Multilevel Caches



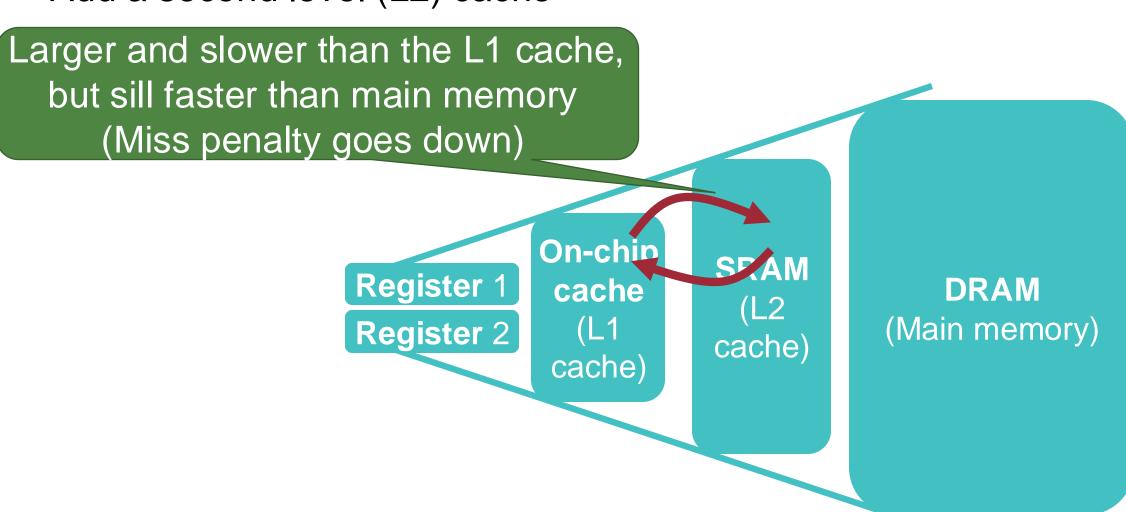




10

Multilevel Caches

Add a second level (L2) cache



Summary: How to Improve the Cache Performance?

Adjust the block size

 When the cache block size is too small or too large, cache shows poor performance

Adjust the cache placement policy

- Direct mapped, fully associative, set associative
- A fully associative cache theoretically shows best performance, but it is complex

Adjust the cache replacement policy

- Generally, LRU usually shows best performance

Multilevel caches

- Nowadays, there are more than one level of cache

Question?