

# 16. Cache organization: block size, writes

---

EECS 370 – Introduction to Computer Organization - Winter 2016

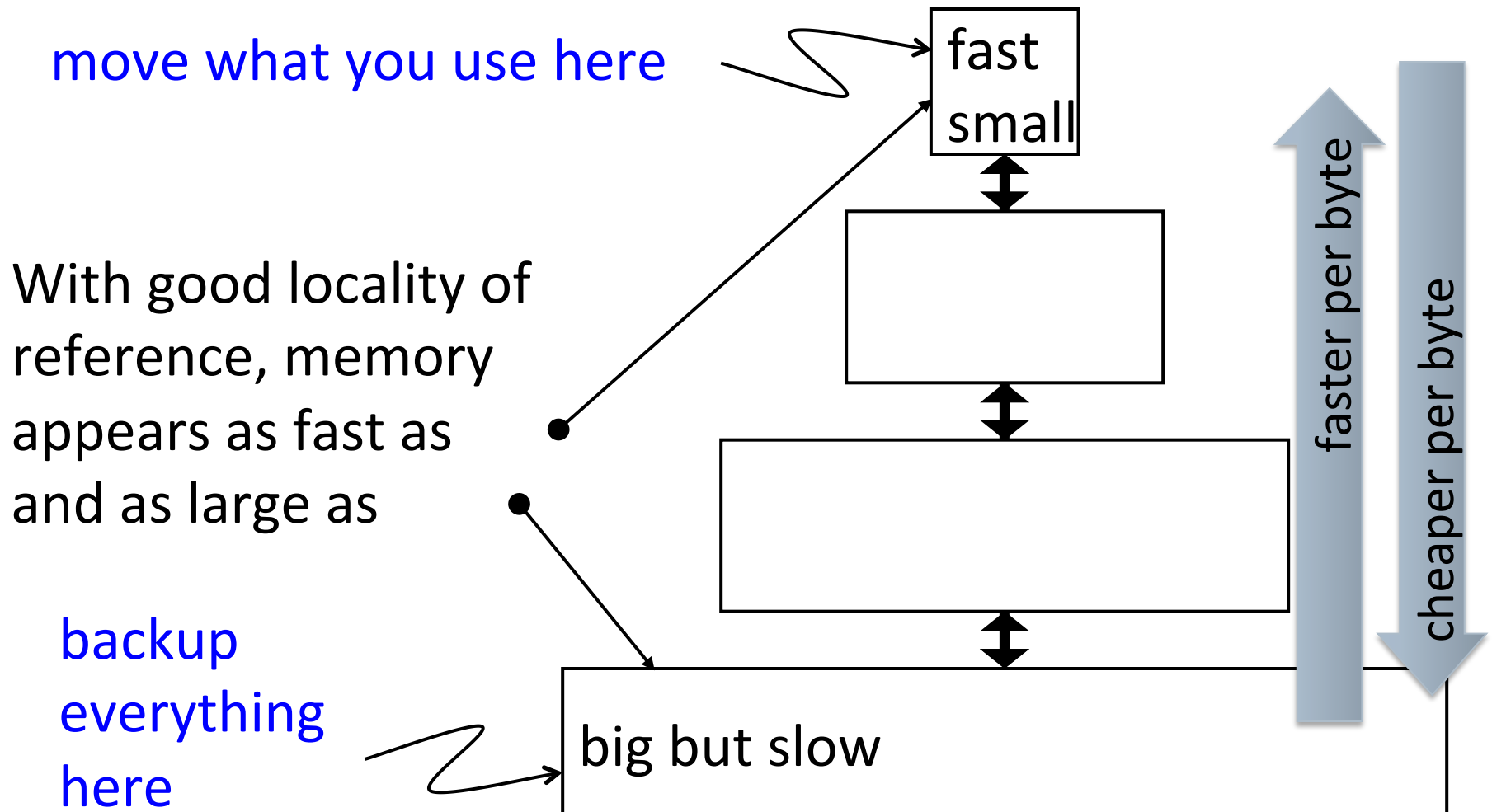
**Profs. Valeria Bertacco & Reetu Das**

EECS Department  
University of Michigan in Ann Arbor, USA

© Bertacco-Das, 2016

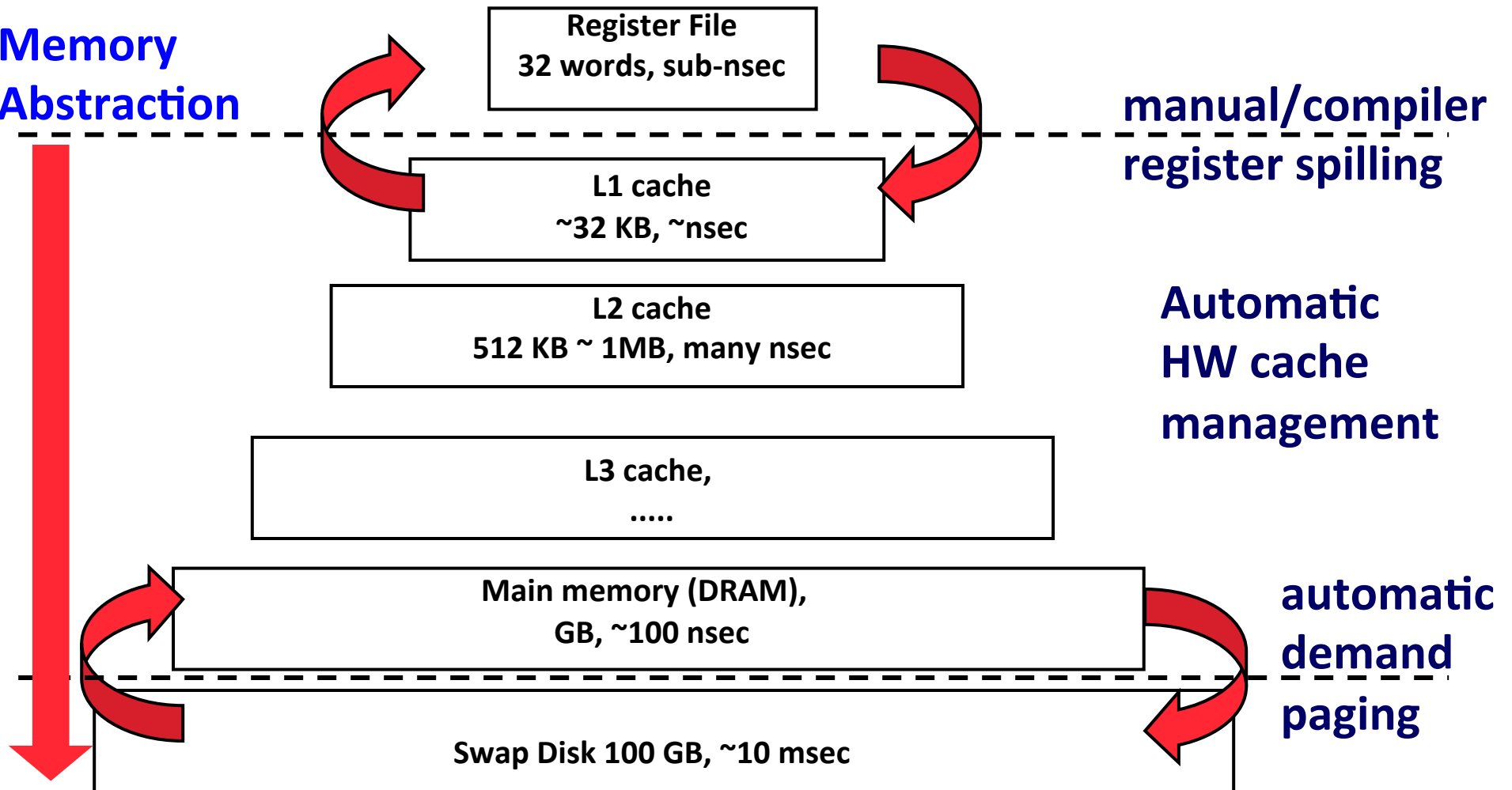
The material in this presentation cannot be  
copied in any form without our written permission

## Review: Memory Hierarchy Goals

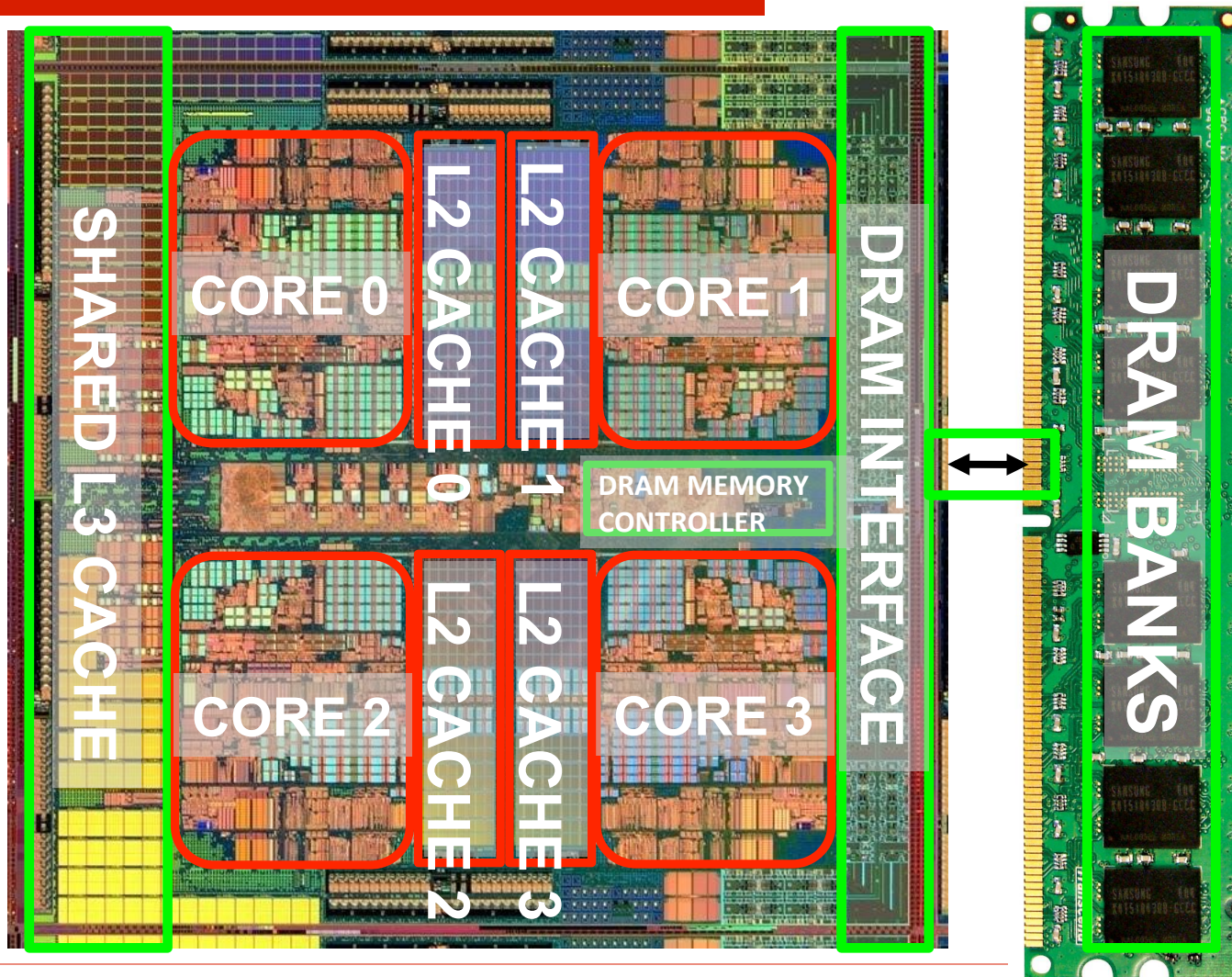


# Memory in a Modern System

**Memory  
Abstraction**



# Memory in a Modern System



## Review: Cache Organization

---

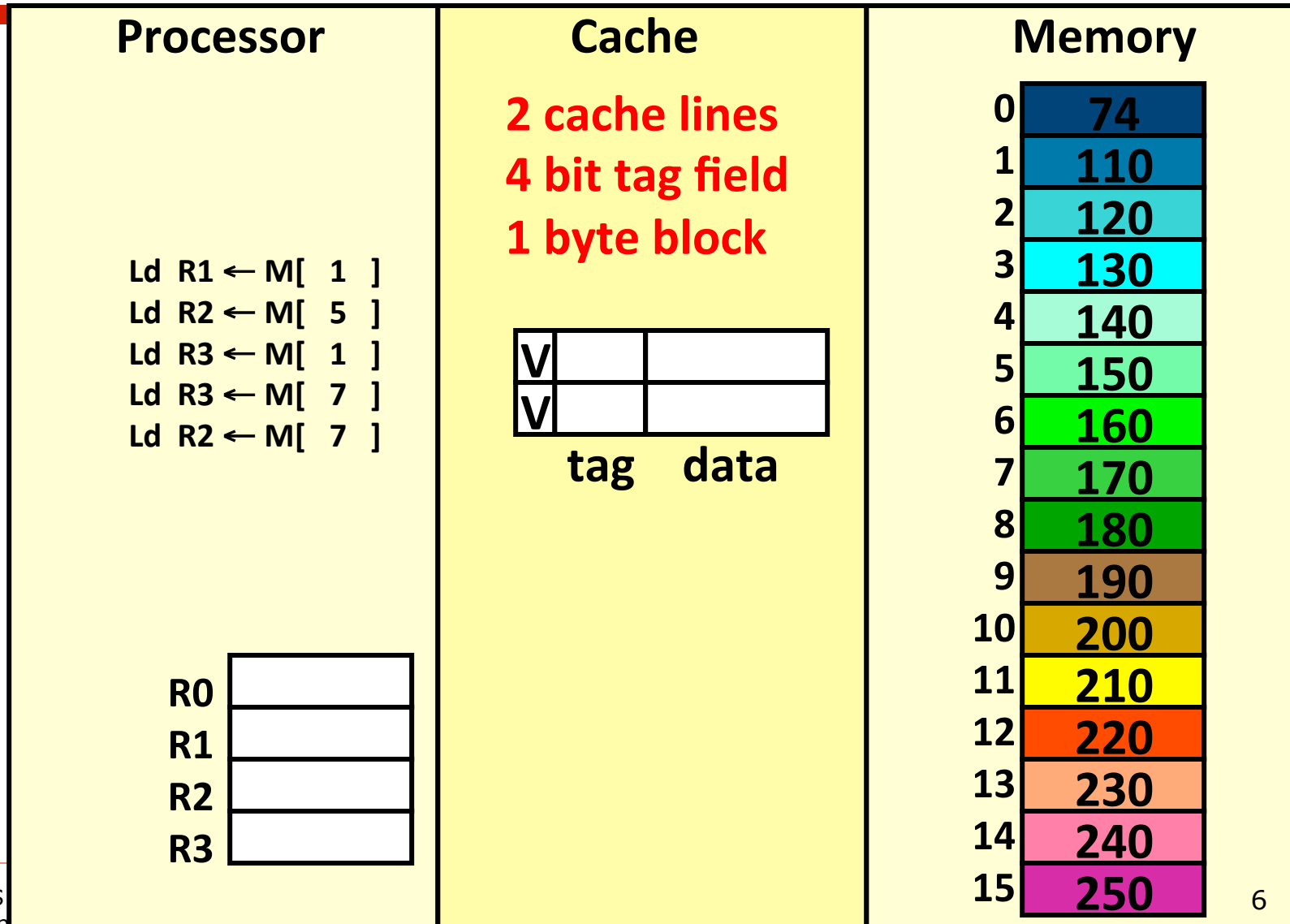
addr	data
addr	data

**TAG    BLOCK**

- ❑ A cache memory consists of multiple tag/block pairs (called **cache lines**)
  - Searches can be done in parallel (within reason)
  - At most one tag will match
- ❑ If there is a tag match, it is a cache **HIT**
- ❑ If there is no tag match, it is a cache **MISS**

Our goal is to keep the data we think will be accessed in the near future in the cache

# Recap: A Very Simple Cache



# Calculating Cost

---

- ❑ How much does our example cache from last lecture cost (in bits)?
  - Calculate storage requirements
    - 2 bytes of SRAM
  - Calculate overhead to support access (tags)
    - 2 4-bit tags
    - The cost of the tags is often forgotten for caches, but this cost drives the design of real caches
    - 2 valid bits
- ❑ What is the cost if a 32 bit address is used?
  - **2 32-bit tags**
  - 2 valid bits
  - 2 bytes of SRAM

# How can we reduce the overhead for each cache line?

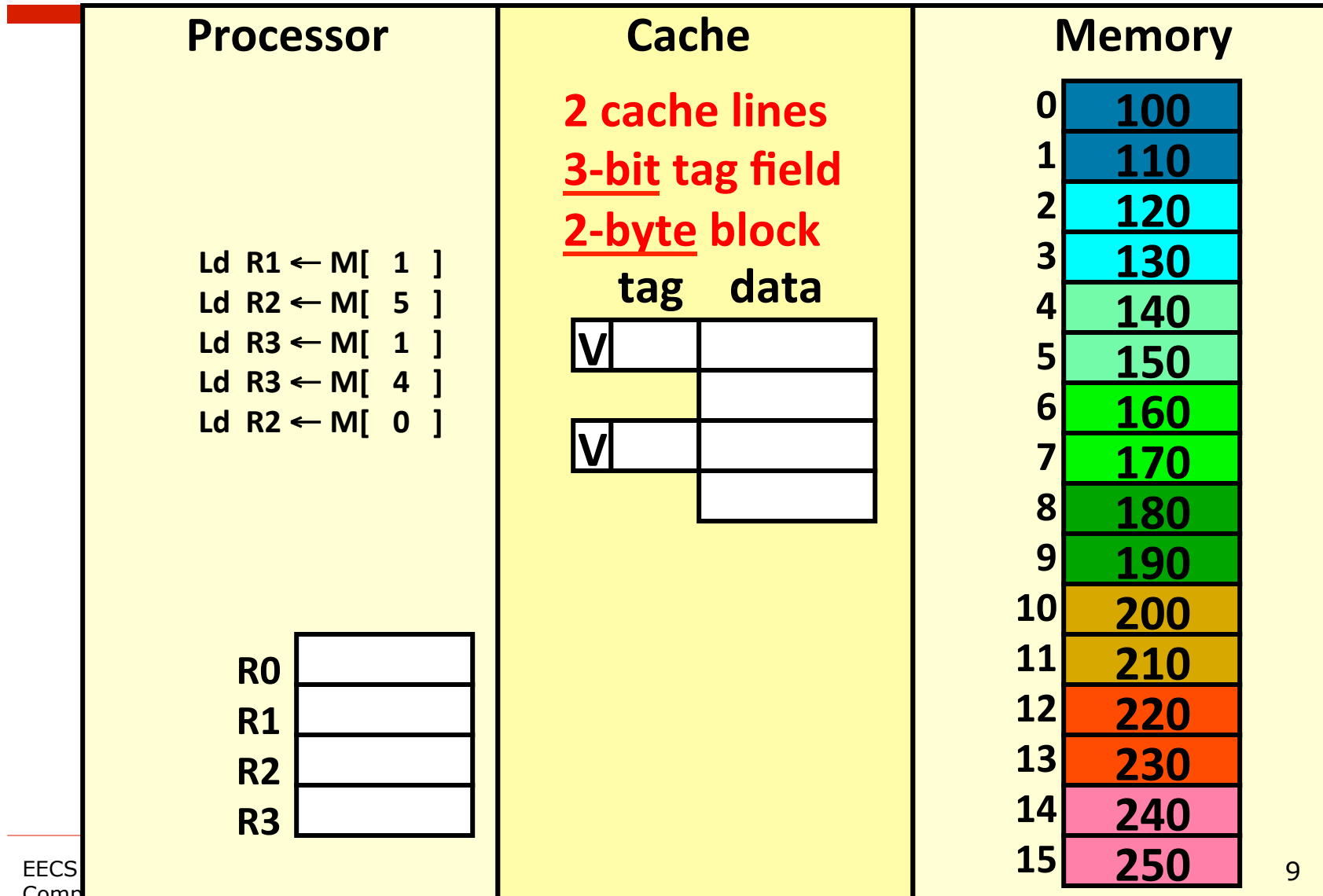
---

lru	1	1	110
	1	7	170
	tag		data

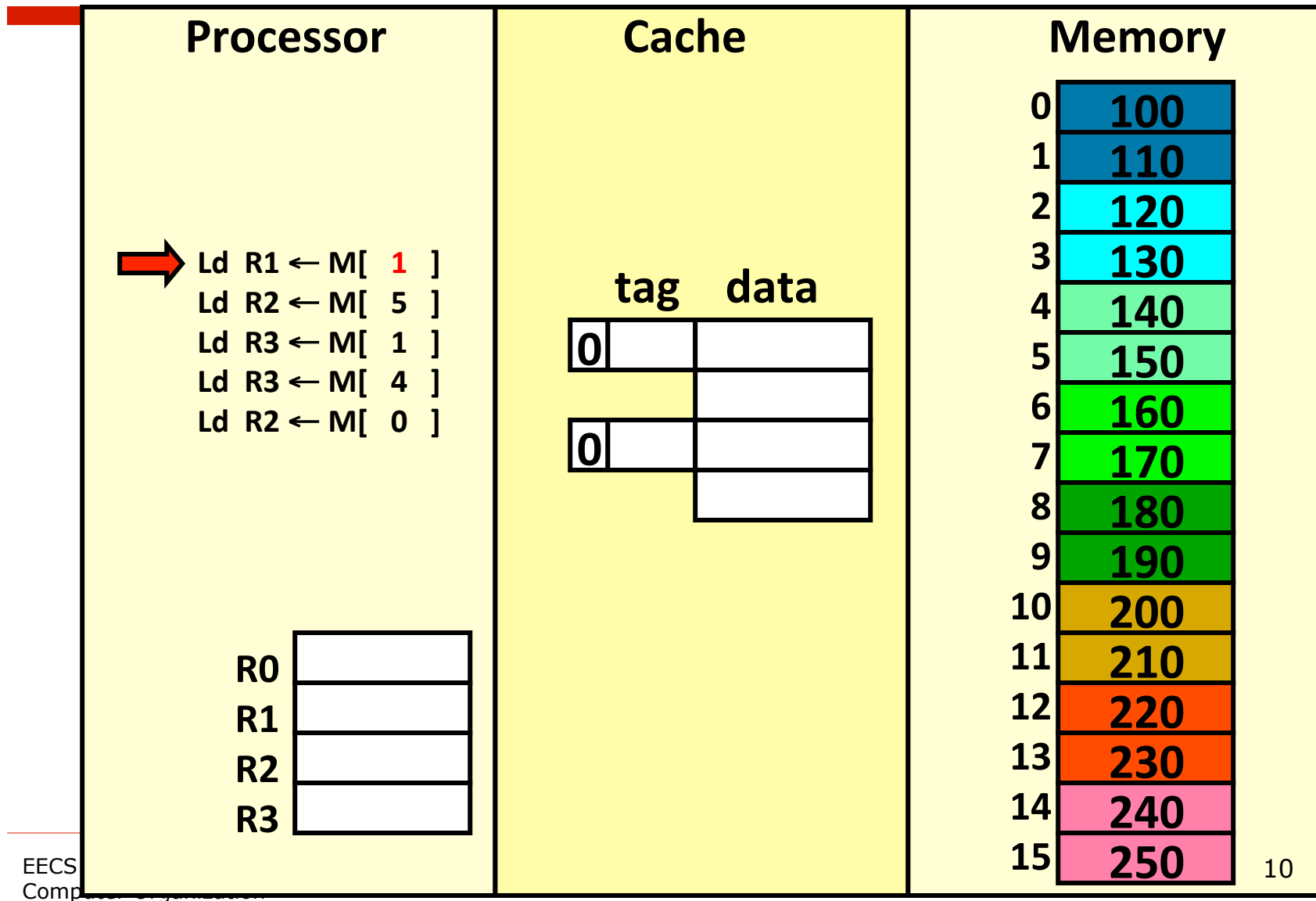
- ❑ Have a small address.
  - Impractical, and caches are supposed to be micro-architectural
- ❑ Cache bigger units than bytes
  - Each block has a single tag, and blocks can be whatever size we choose.



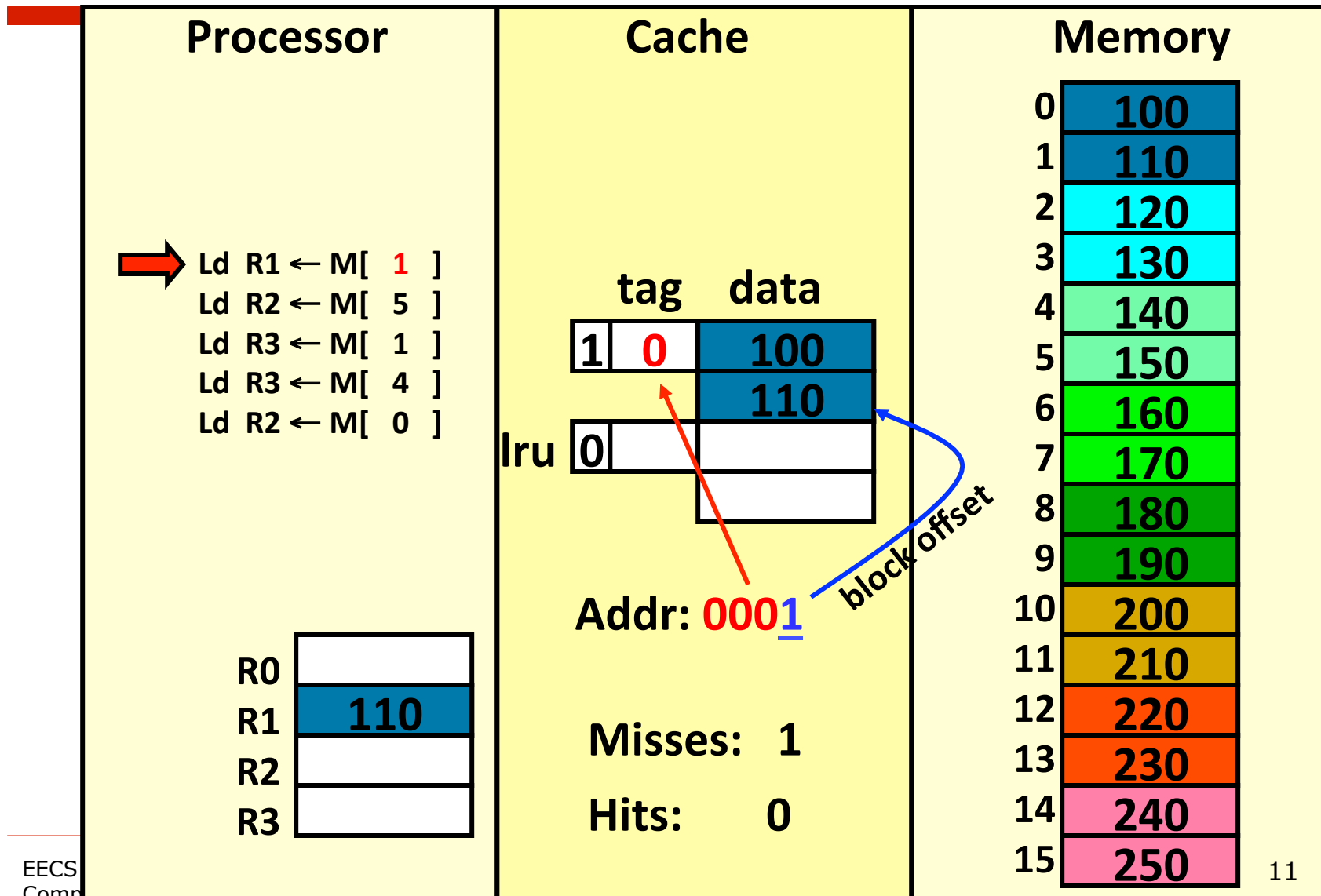
# Block size for caches



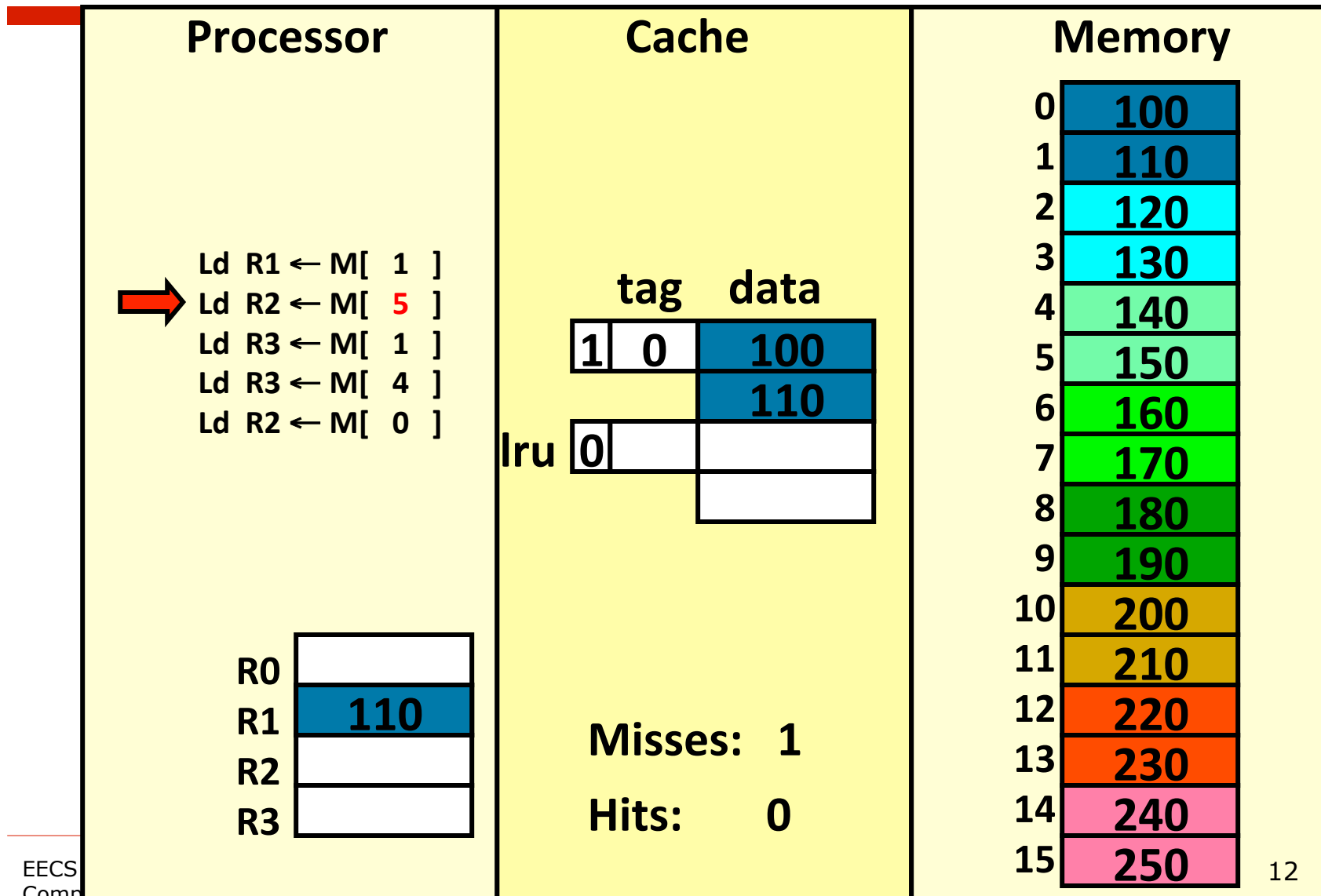
# Block size for caches



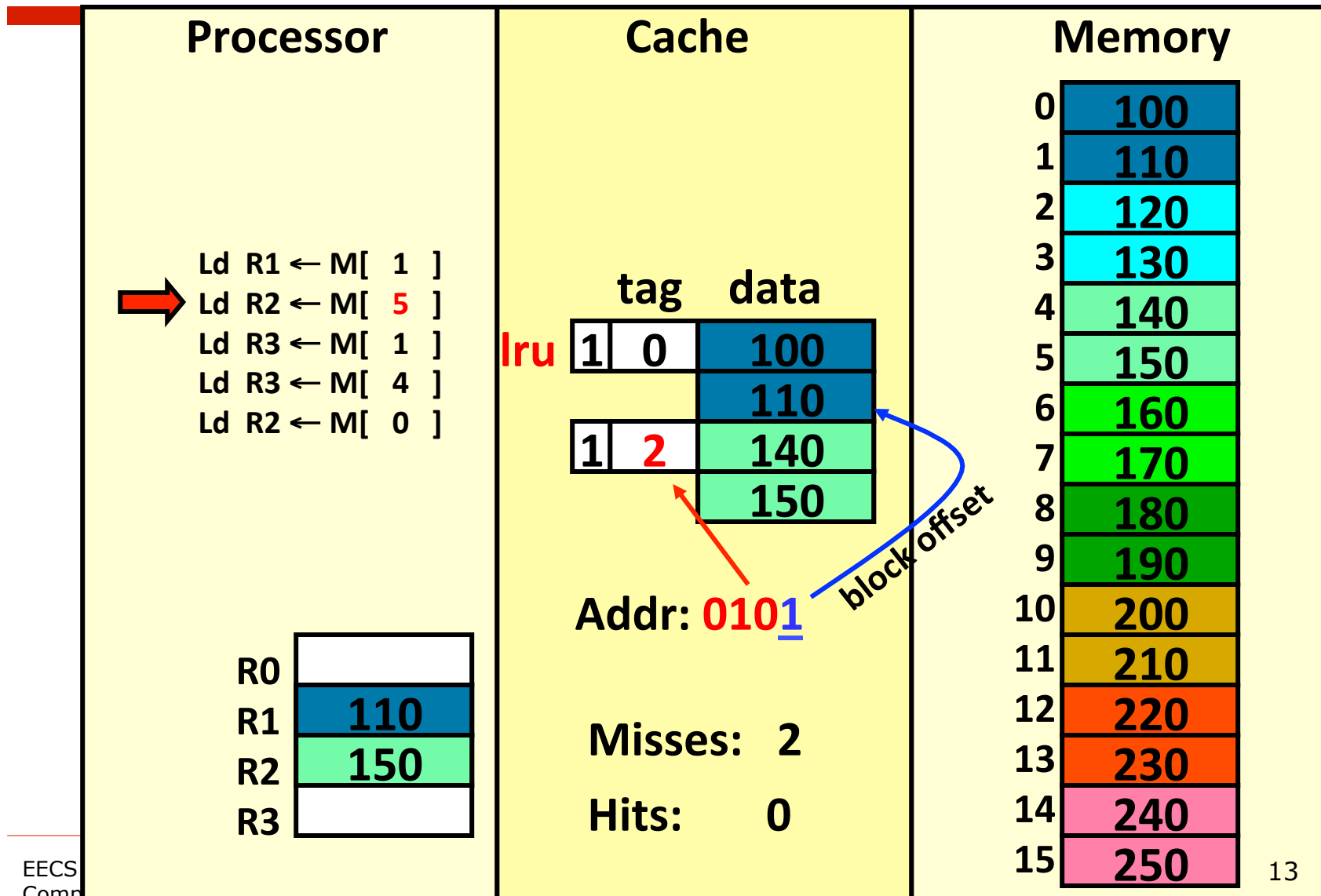
# Block size for caches



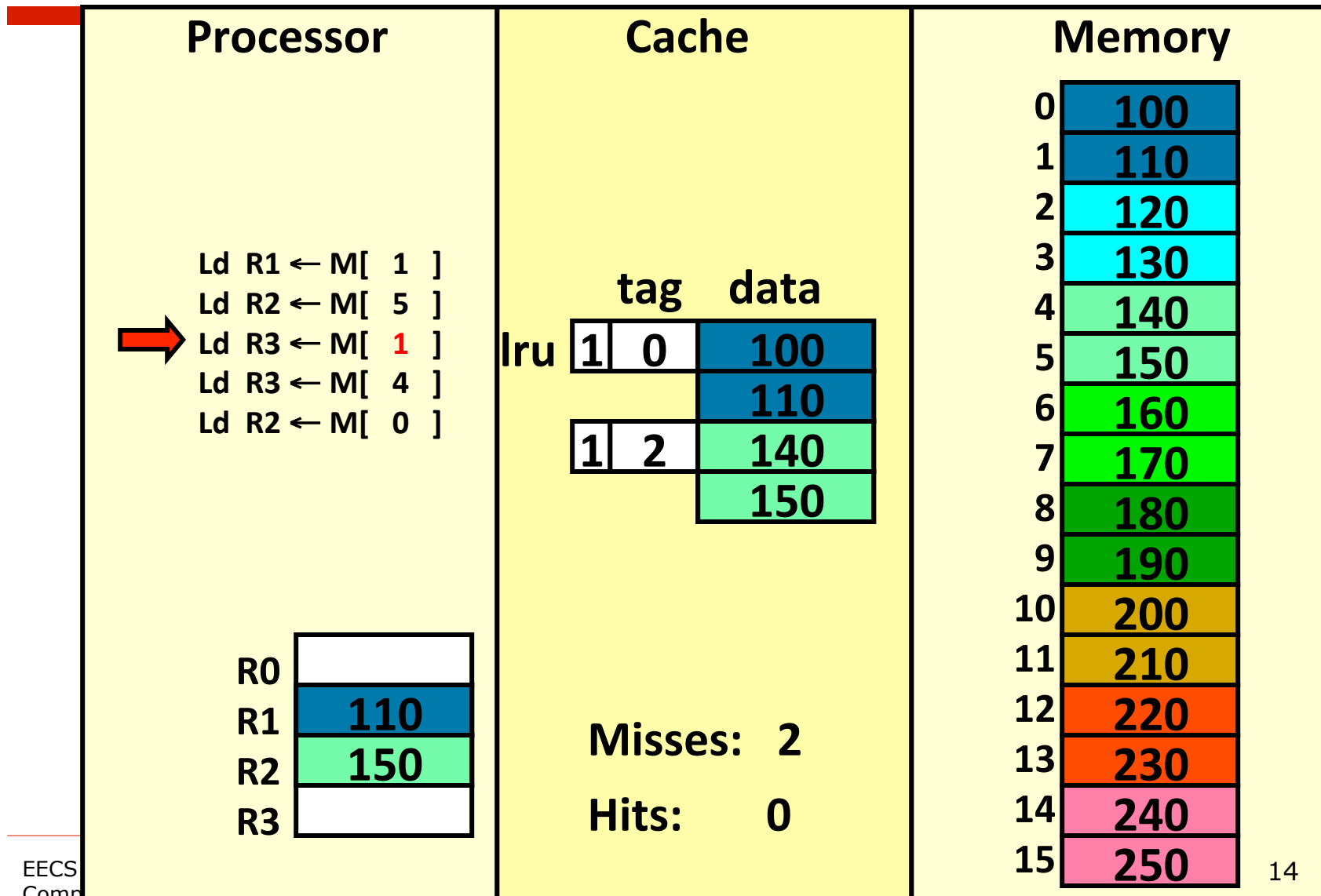
# Block size for caches



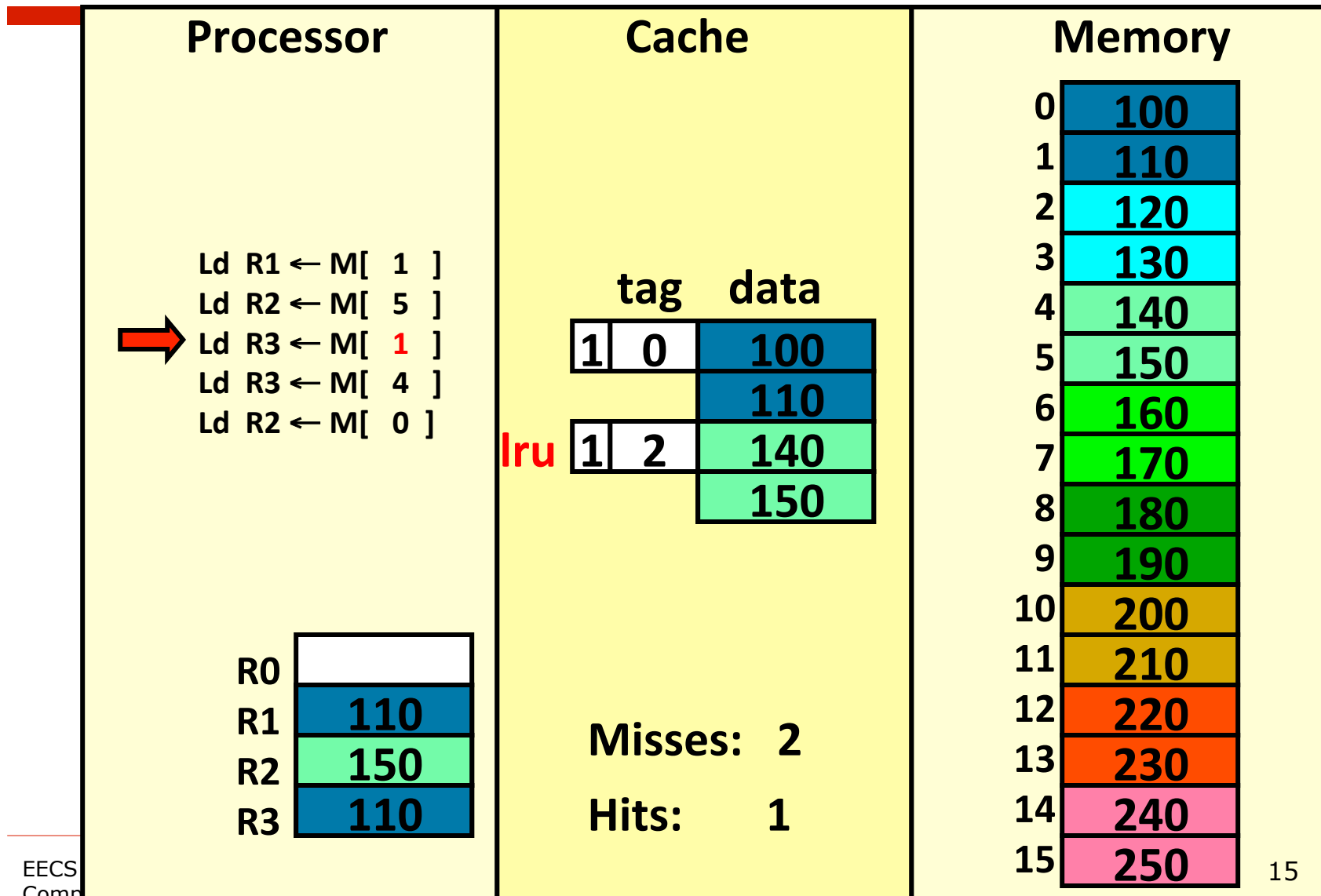
# Block size for caches



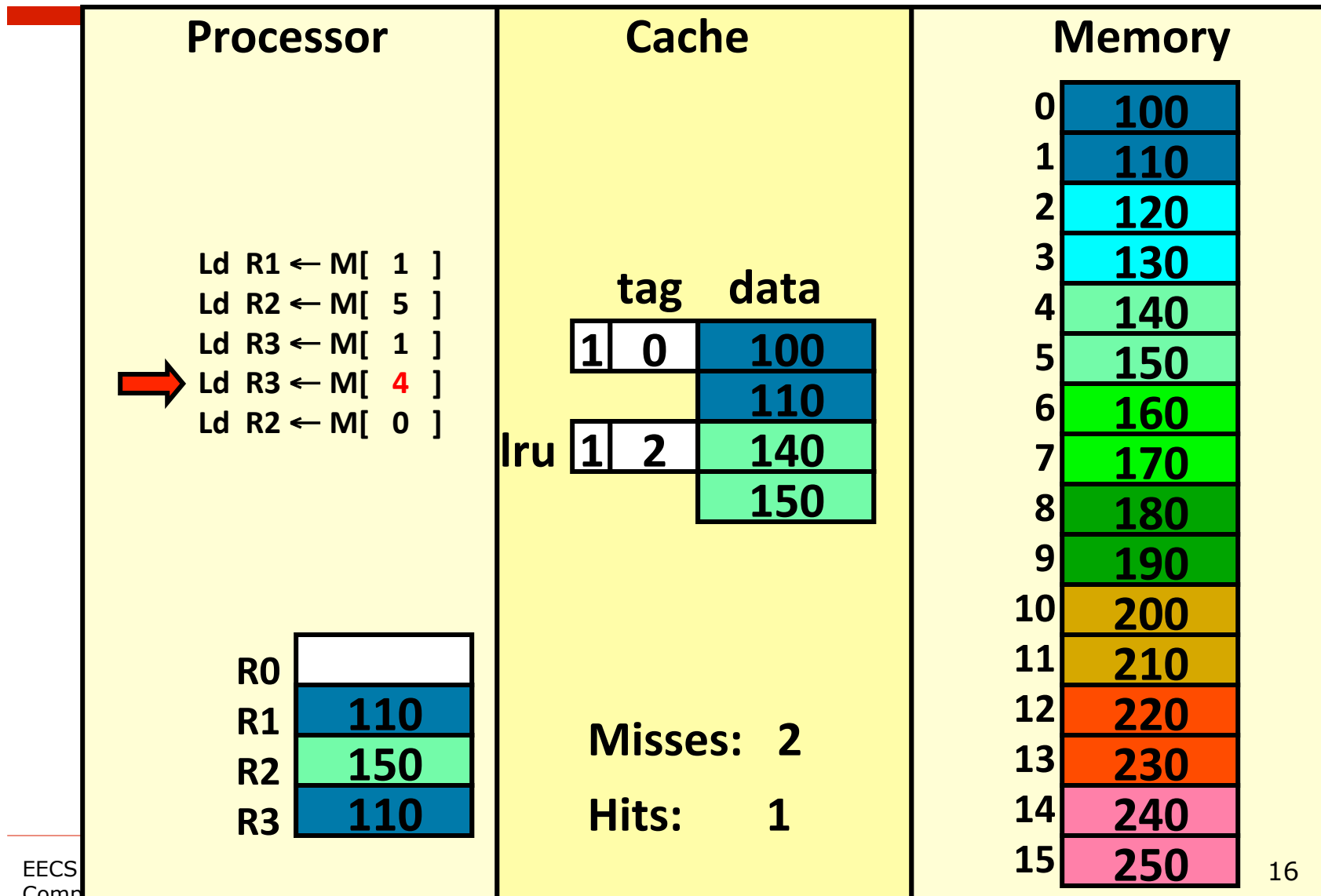
# Block size for caches



# Block size for caches

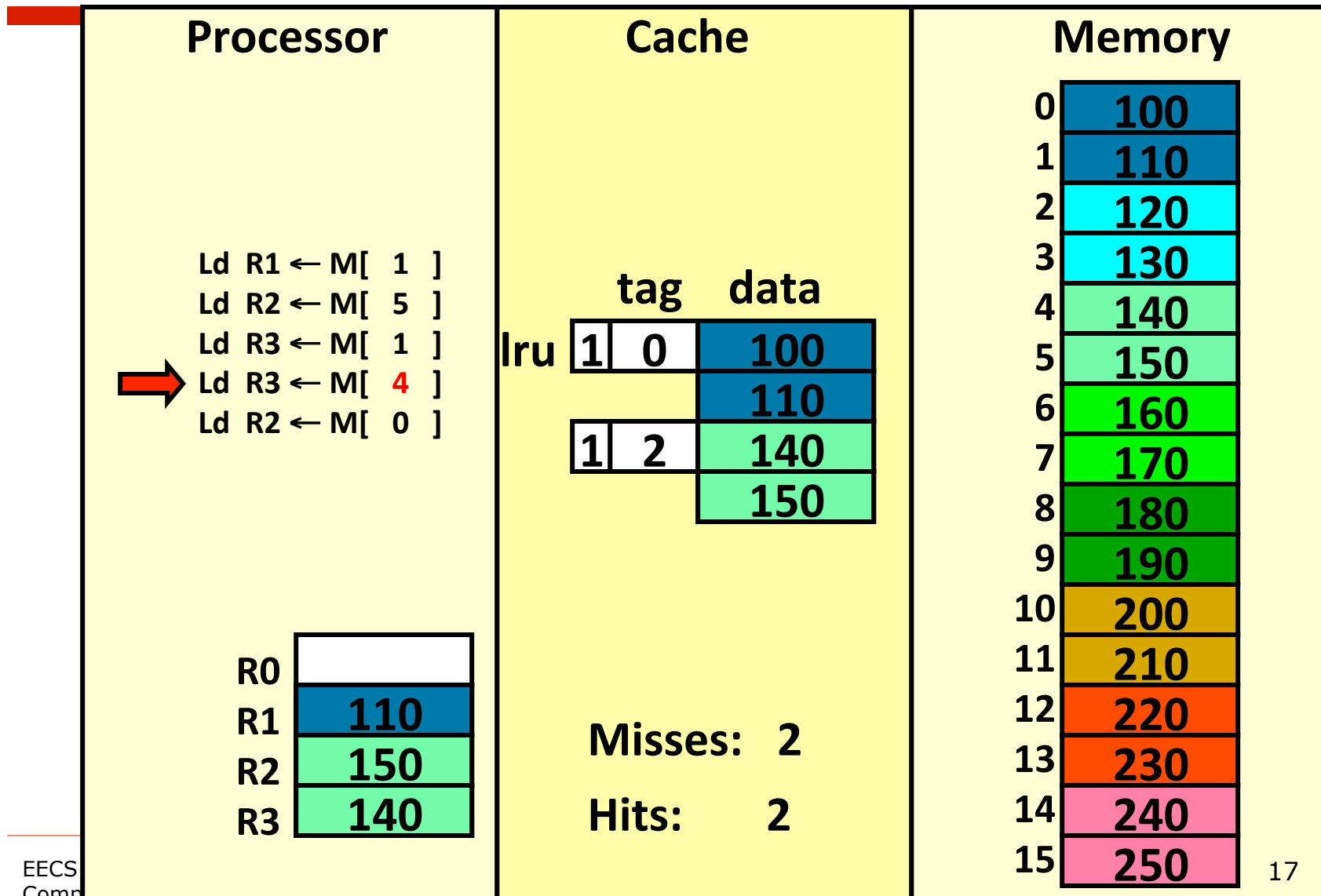


# Block size for caches

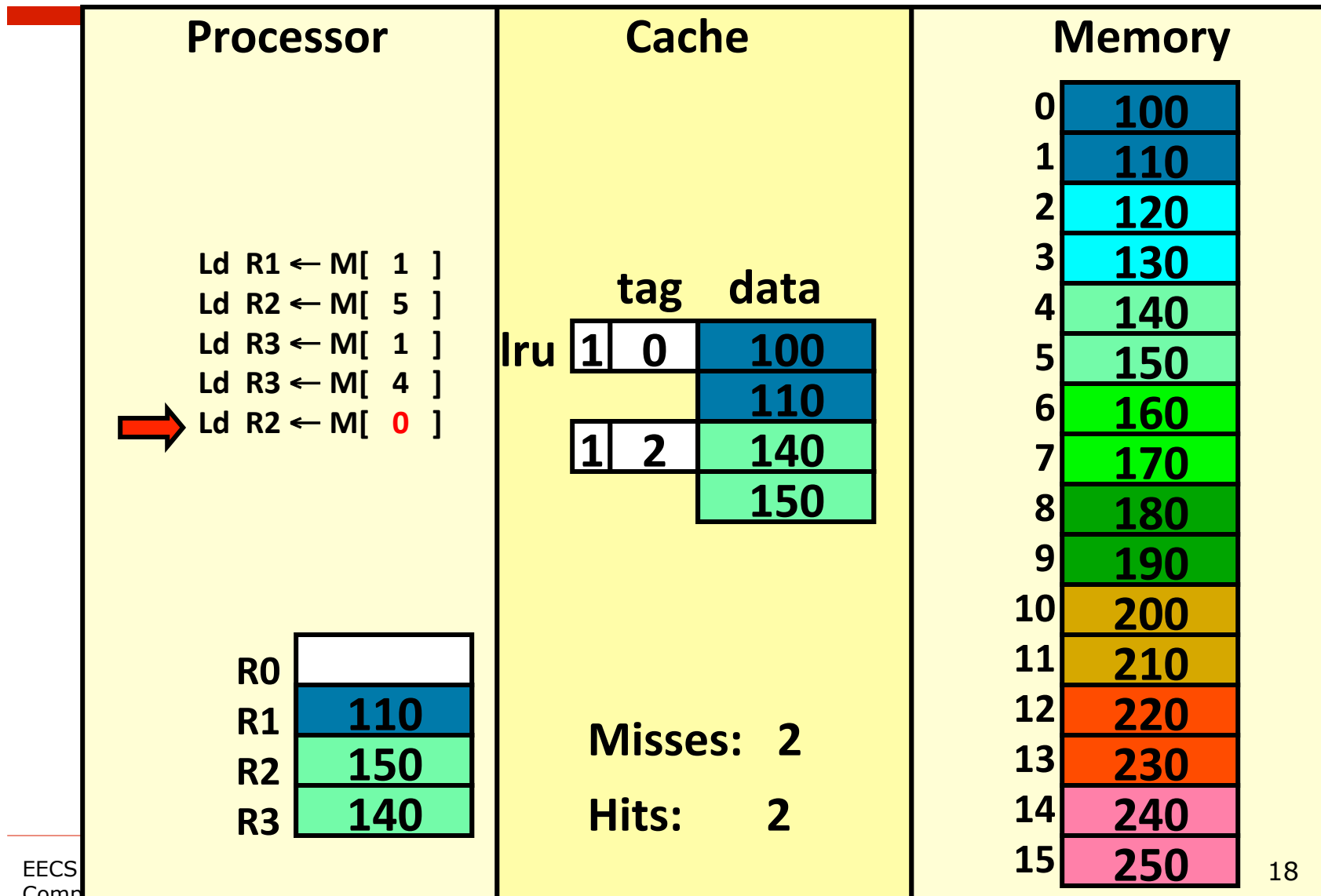




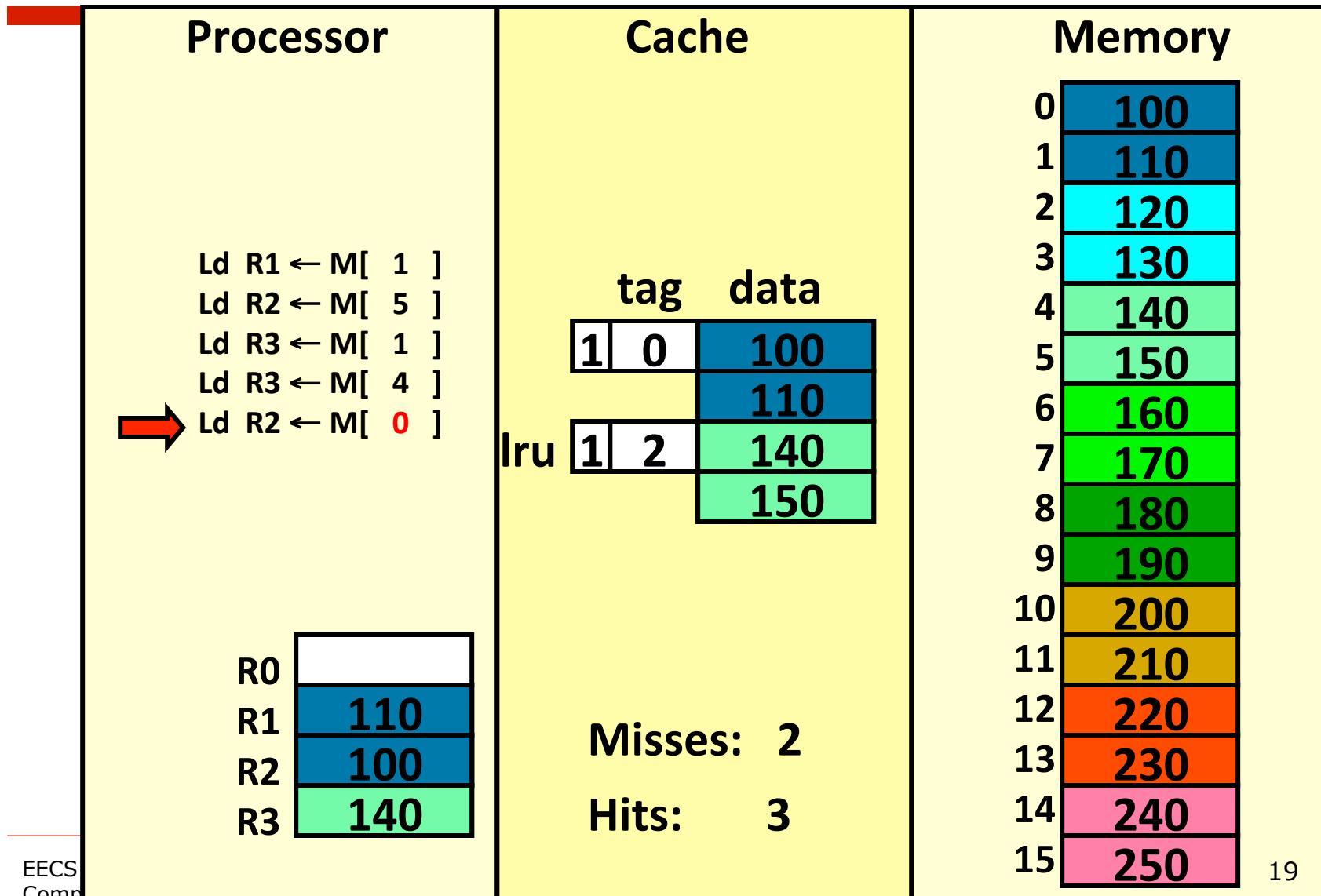
# Block size for caches



# Block size for caches



# Block size for caches



# Spatial Locality

---

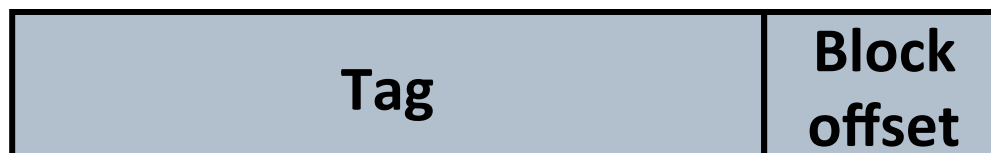
- ❑ Notice that when we accessed address 1, we also brought in address 0.
  - This turned out to be a good thing since we later referenced address 0 and found it in the cache.
- ❑ **Spatial locality** in a program says that if we reference a memory location (e.g., 1000), we are more likely to reference a location near it (e.g. 1001) than some random location.

# Basic Cache organization

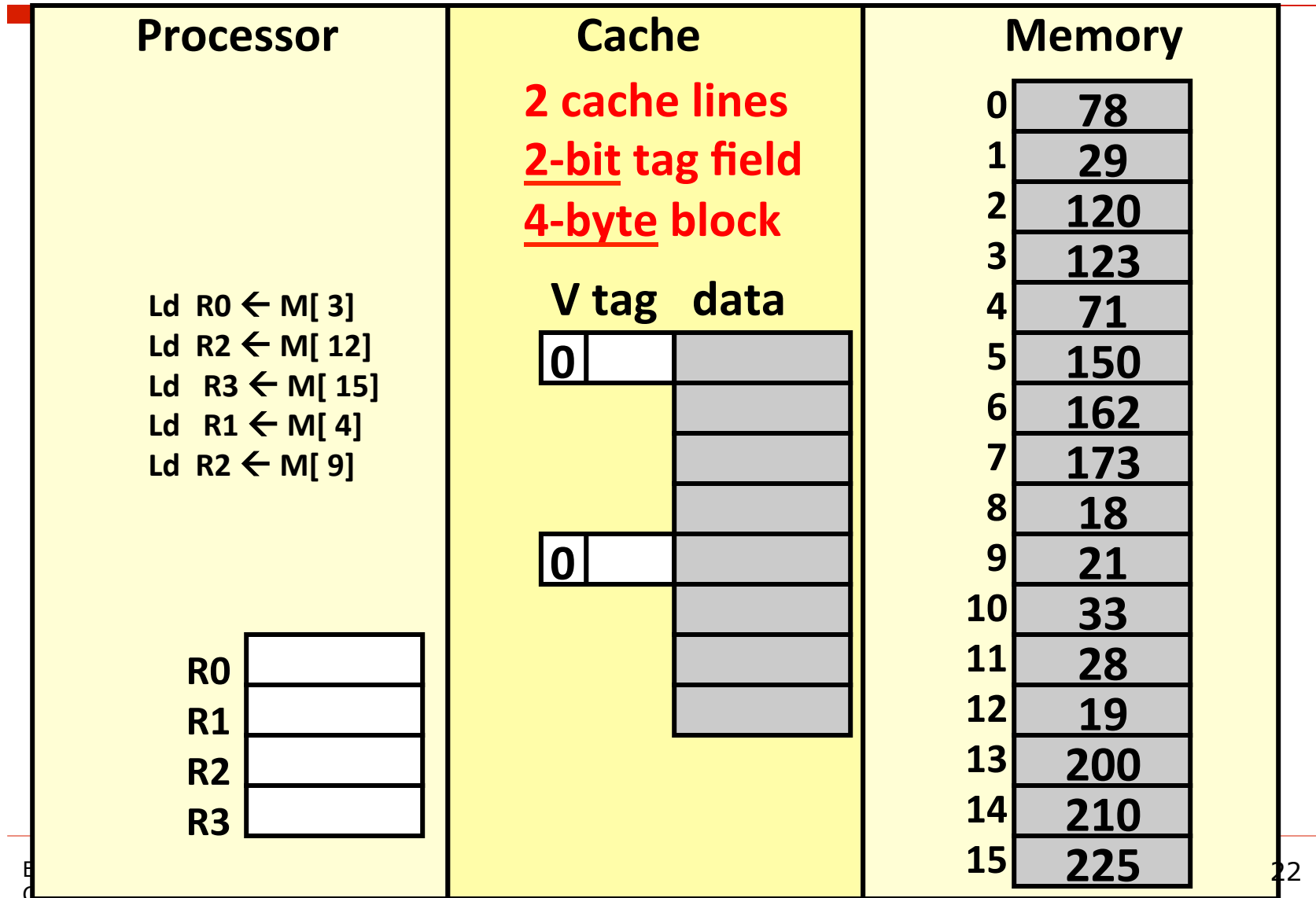
---

## ❑ Decide on the block size

- How? Simulate lots of different block sizes and see which one gives the best performance
- Most systems use a block size between 32 bytes and 128 bytes
- Longer sizes reduce the overhead by:
  - Reducing the number of CAM entries
  - Reducing the size of each CAM entry



# Practice Problem– Compute the register and cache state after executing the instruction sequence



# Solution to Practice Problem

Ld R0  $\leftarrow$  M[ 3]

V tag data

1	0	78
		29
		120
		123
0		

lru

miss

Ld R2  $\leftarrow$  M[ 12]

V tag data

1	0	78
		29
		120
		123
1	3	19
		200
		210
		225

lru

miss

Ld R3  $\leftarrow$  M[ 15]

V tag data

1	0	78
		29
		120
		123
1	3	19
		200
		210
		225

lru

hit

Ld R1  $\leftarrow$  M[ 4]

V tag data

1	1	71
		150
		162
		173
1	3	19
		200
		210
		225

lru

miss

Ld R2  $\leftarrow$  M[ 9]

V tag data

1	1	71
		150
		162
		173
1	2	18
		21
		33
		28

lru

miss

# Class Problem 1

---

- ❑ Given a cache with the following configuration: total size is 8 bytes, block size is 2 bytes, *fully associative*, LRU replacement. The memory address size is 16 bits and is byte addressable.

1. How many bits are for each tag? How many blocks in the cache?

$$\text{Tag} = 16 - 1 (\text{block offset}) = 15 \text{ bits}$$

2. For the following reference stream, indicate whether each reference is a hit or miss: 0, 1, 3, 5, 12, 1, 2, 9, 4

M, H, M, M, M, H, H, M, M

3. What is the hit rate?  $3/9 = 33\%$
4. How many bits are needed for storage overhead?

$$\text{Overhead} = 15 (\text{Tag}) + 1 (\text{V}) + 2 (\text{LRU}) = 18 \text{ bits}$$

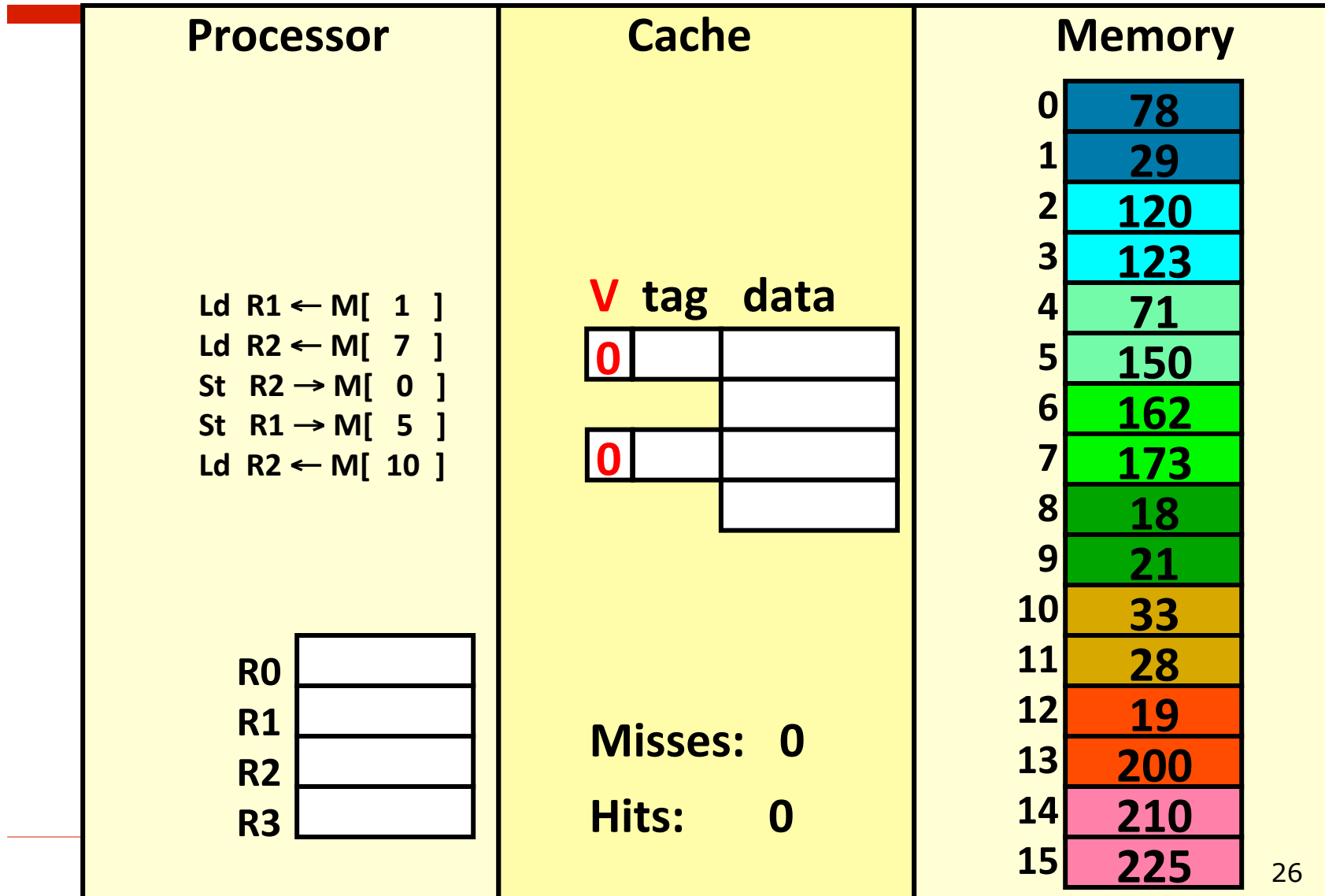


# What about stores?

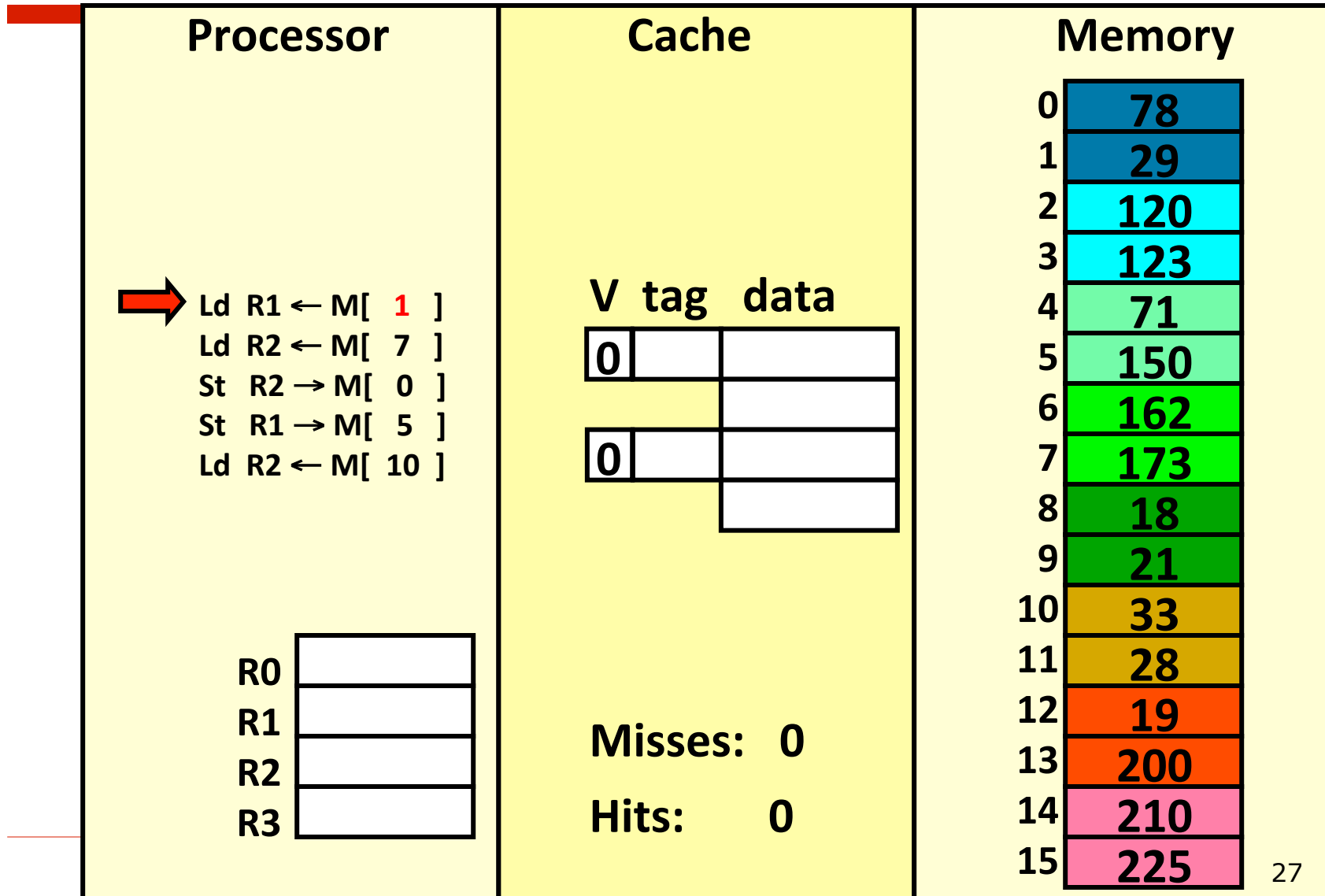
---

- ❑ Where should you write the result of a store?
  - If that memory location is in the cache?
    - Send it to the cache.
    - Should we also send it to memory?  
(write-through policy)
  - If it is not in the cache?
    - Allocate the line (put it in the cache)?  
(allocate-on-write policy)
    - Write it directly to memory without allocation?

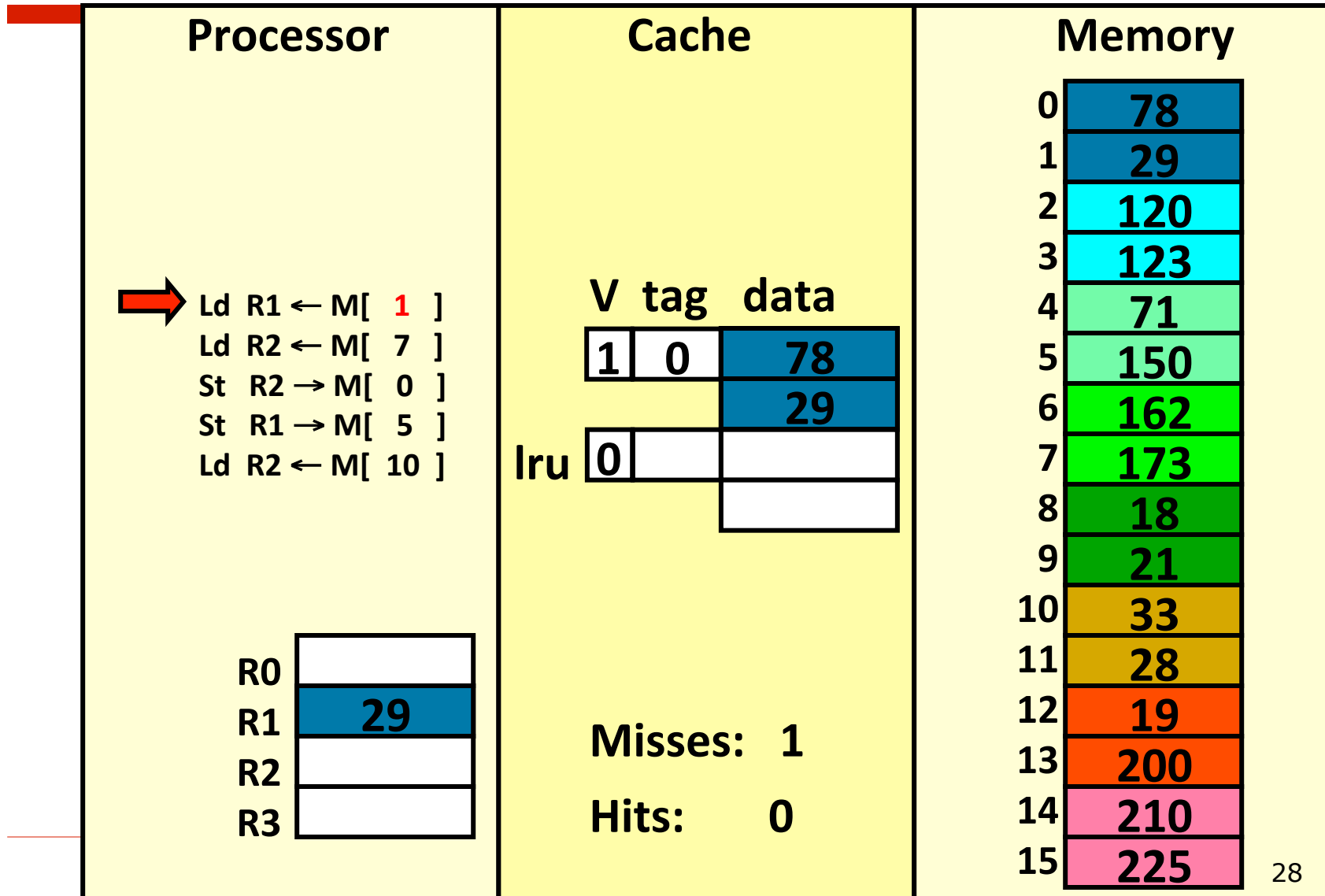
# Handling stores (write-through)



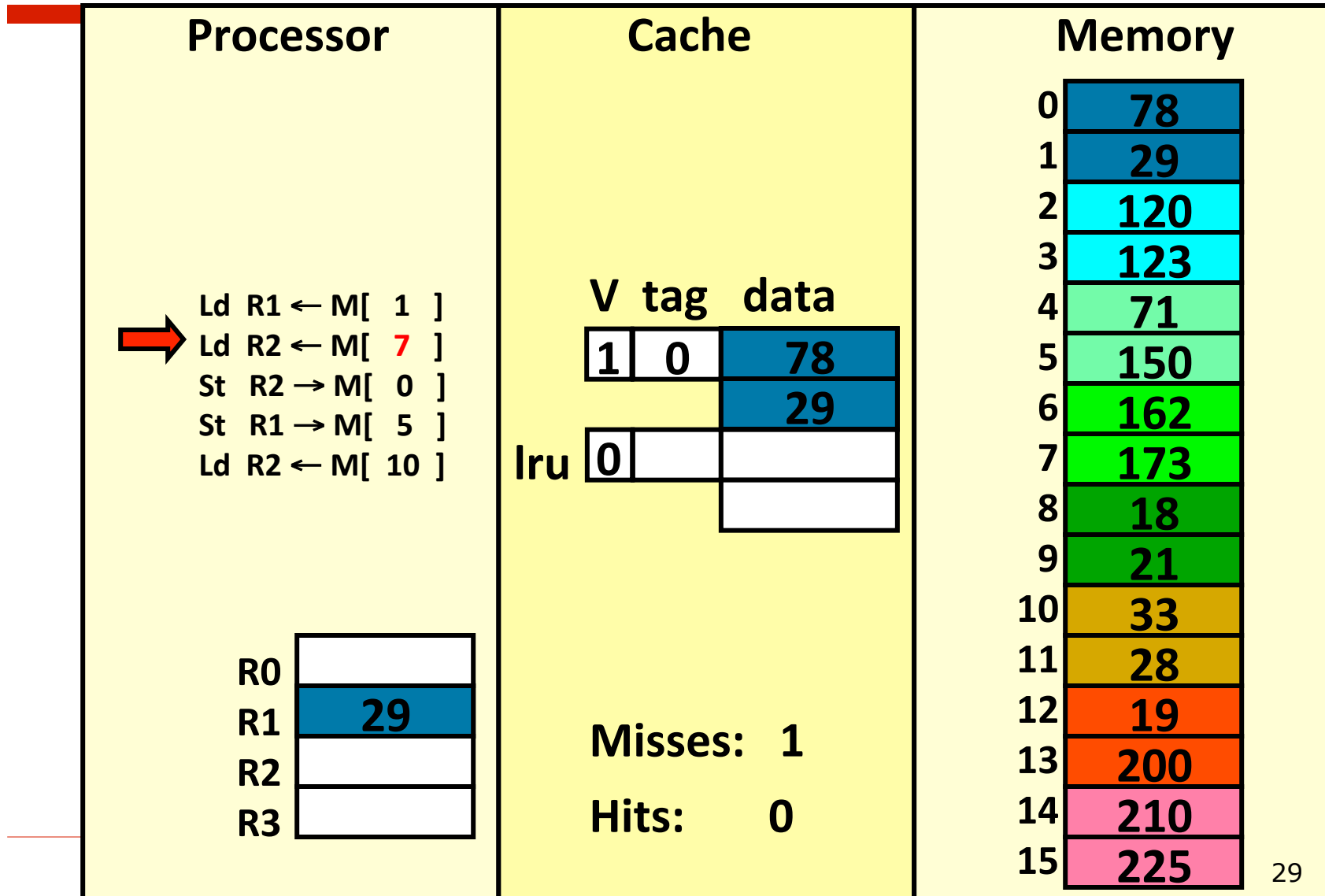
# write-through (REF 1)



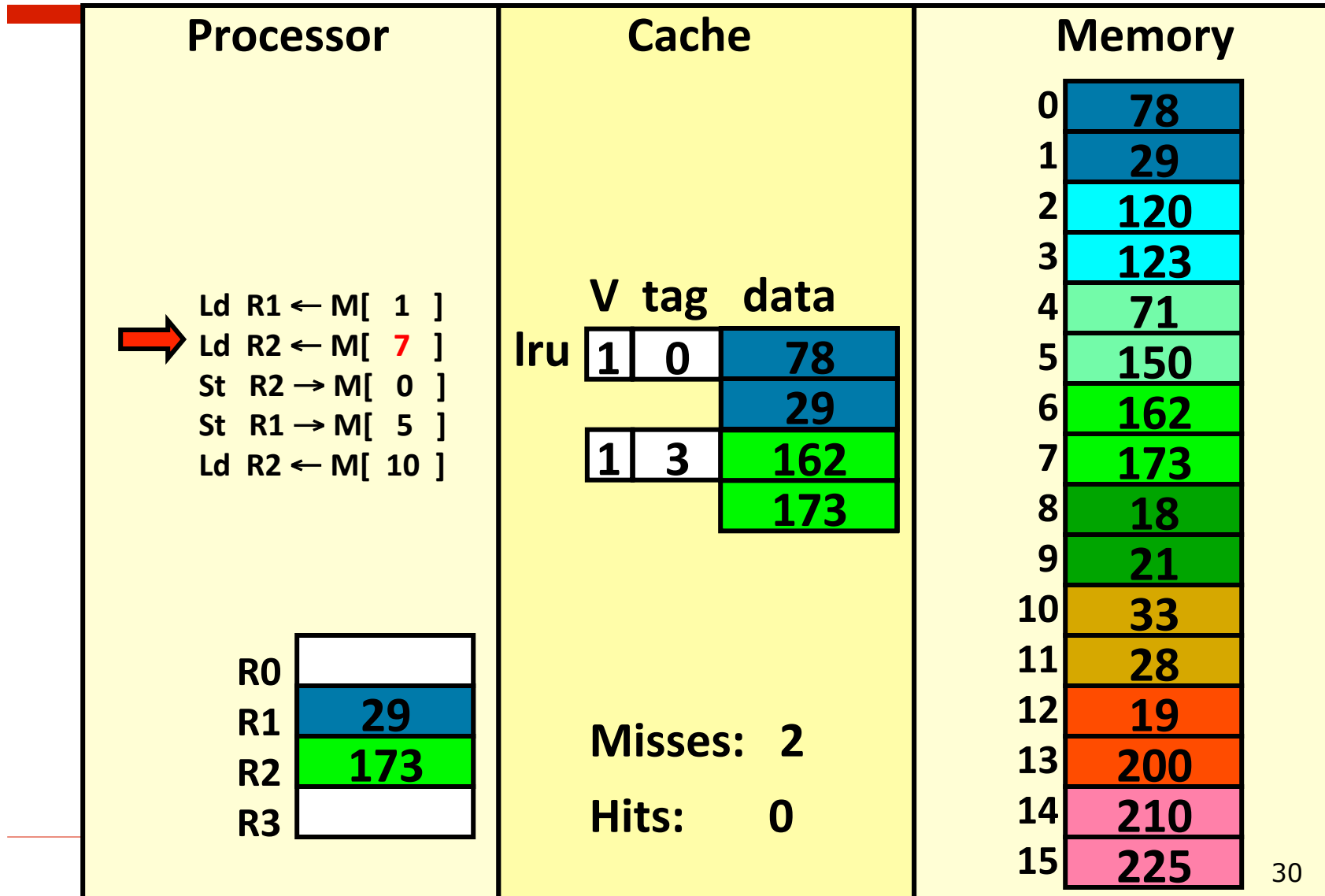
# write-through (REF 1)



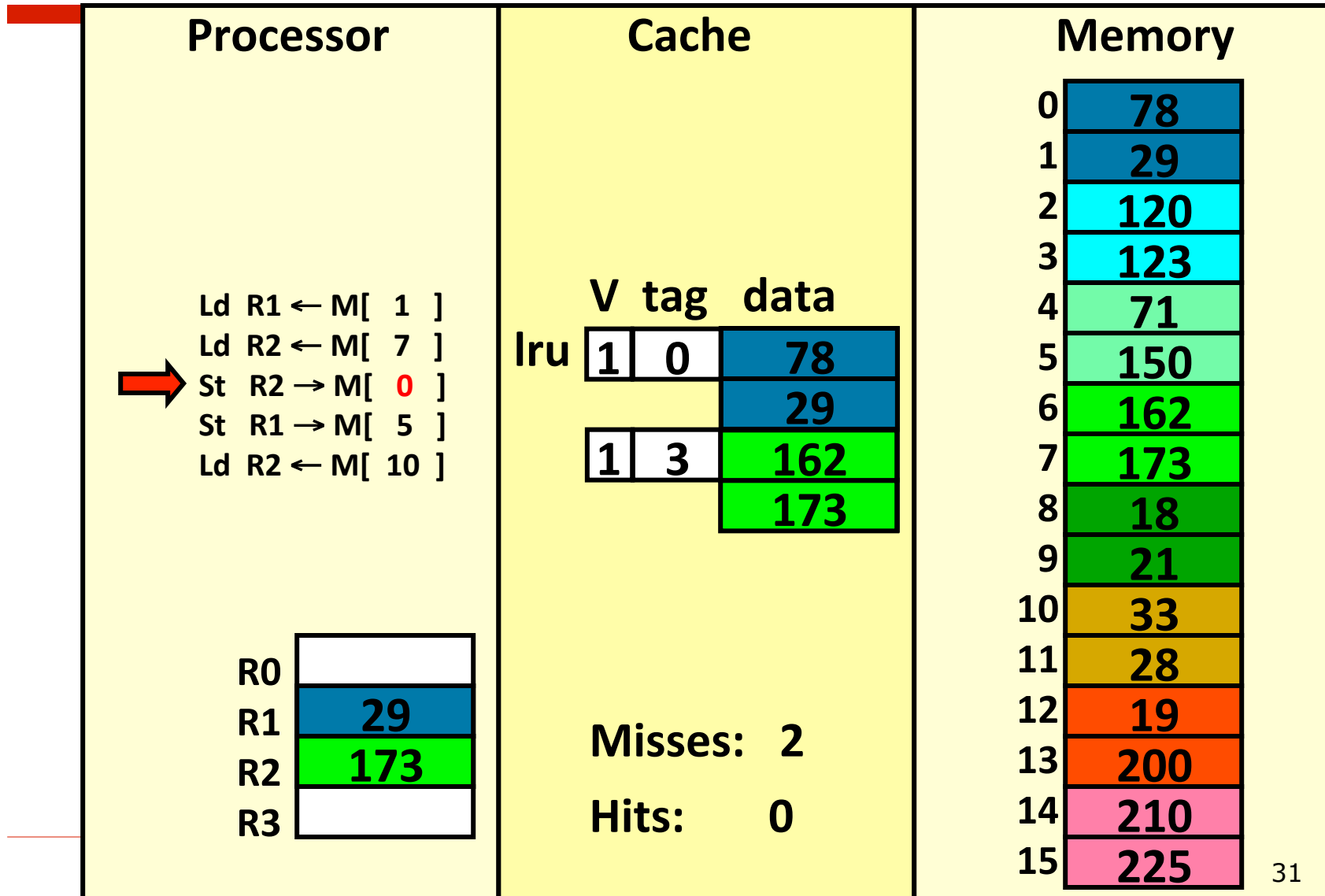
# write-through (REF 2)



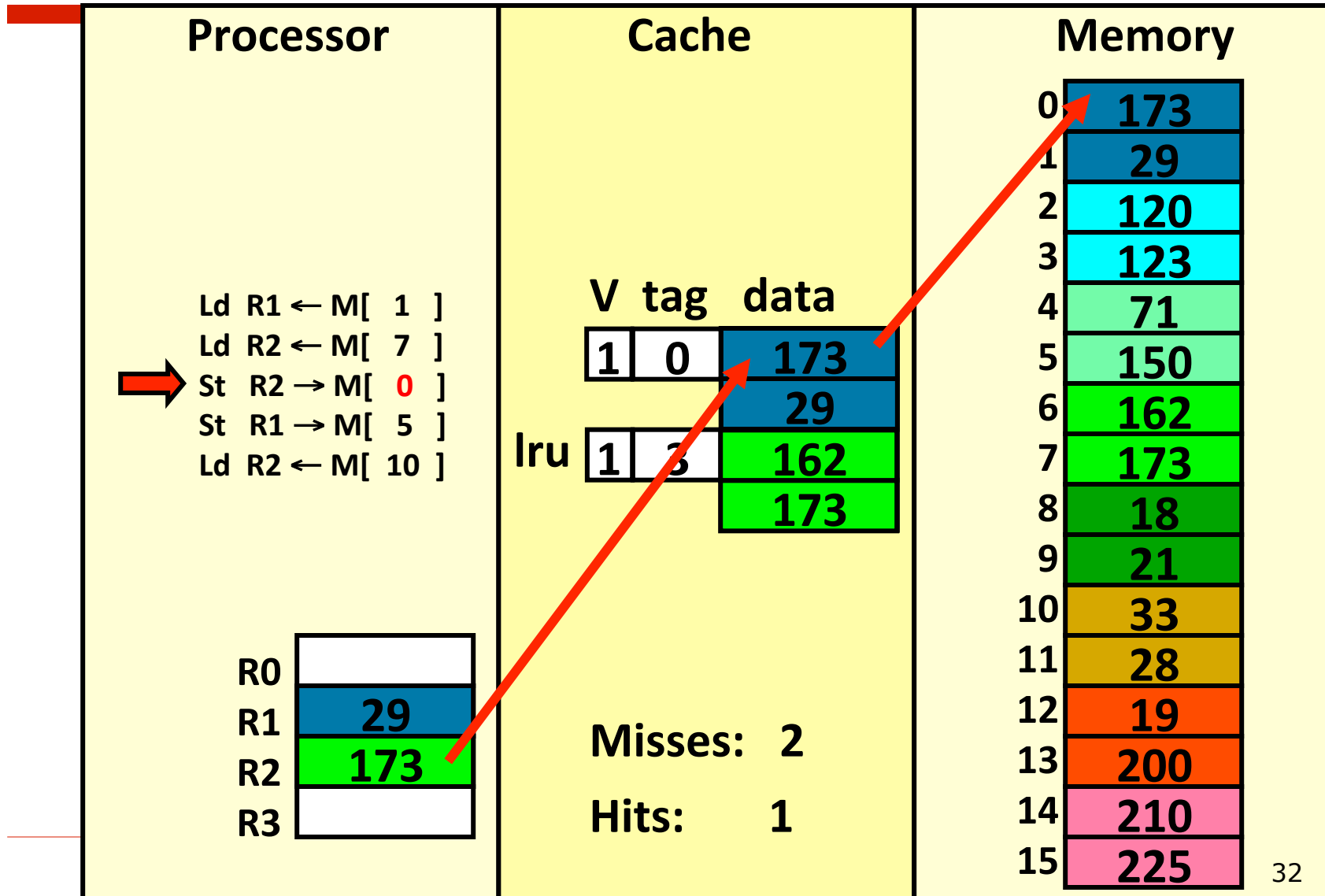
# write-through (REF 2)



# write-through (REF 3)

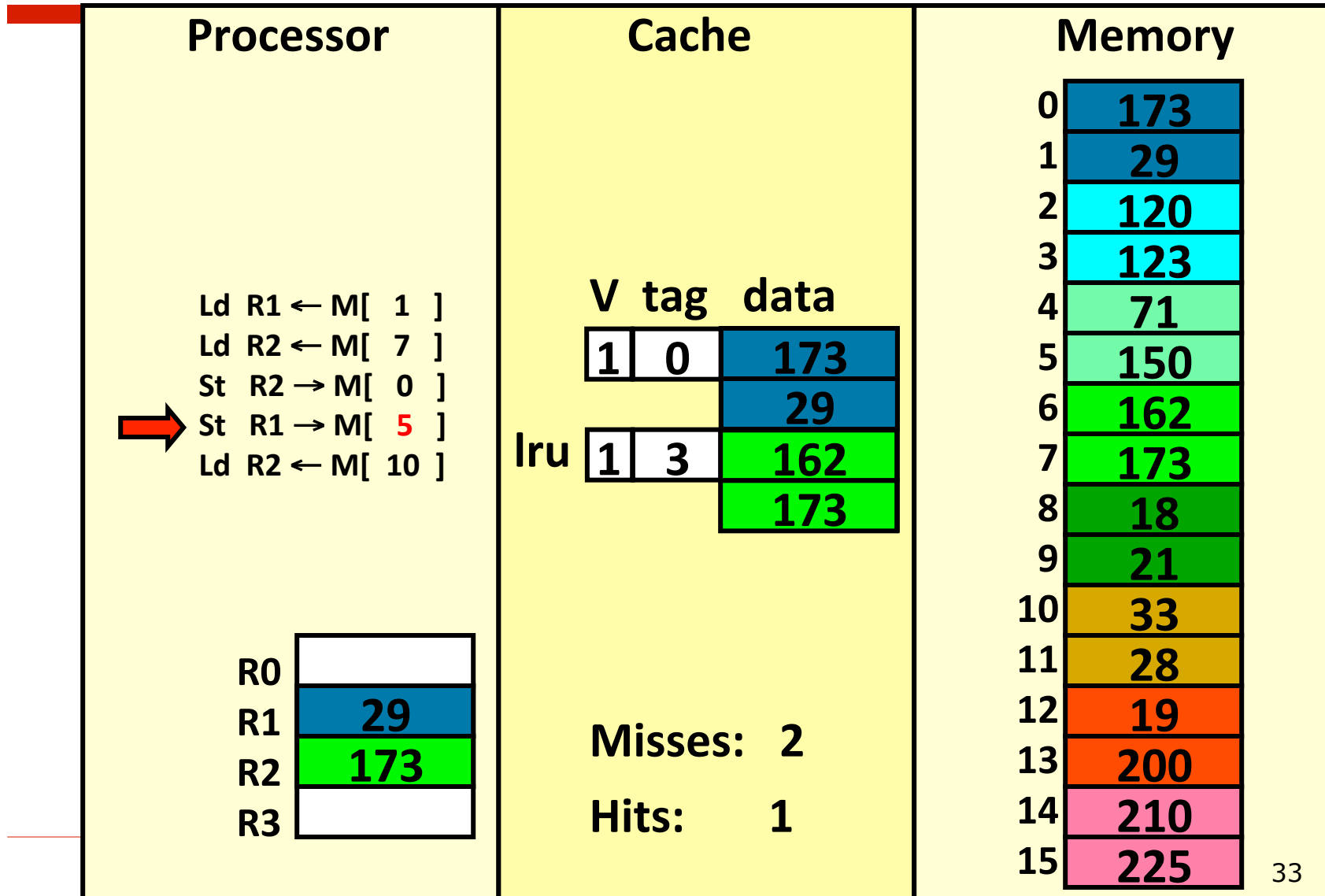


# write-through (REF 3)

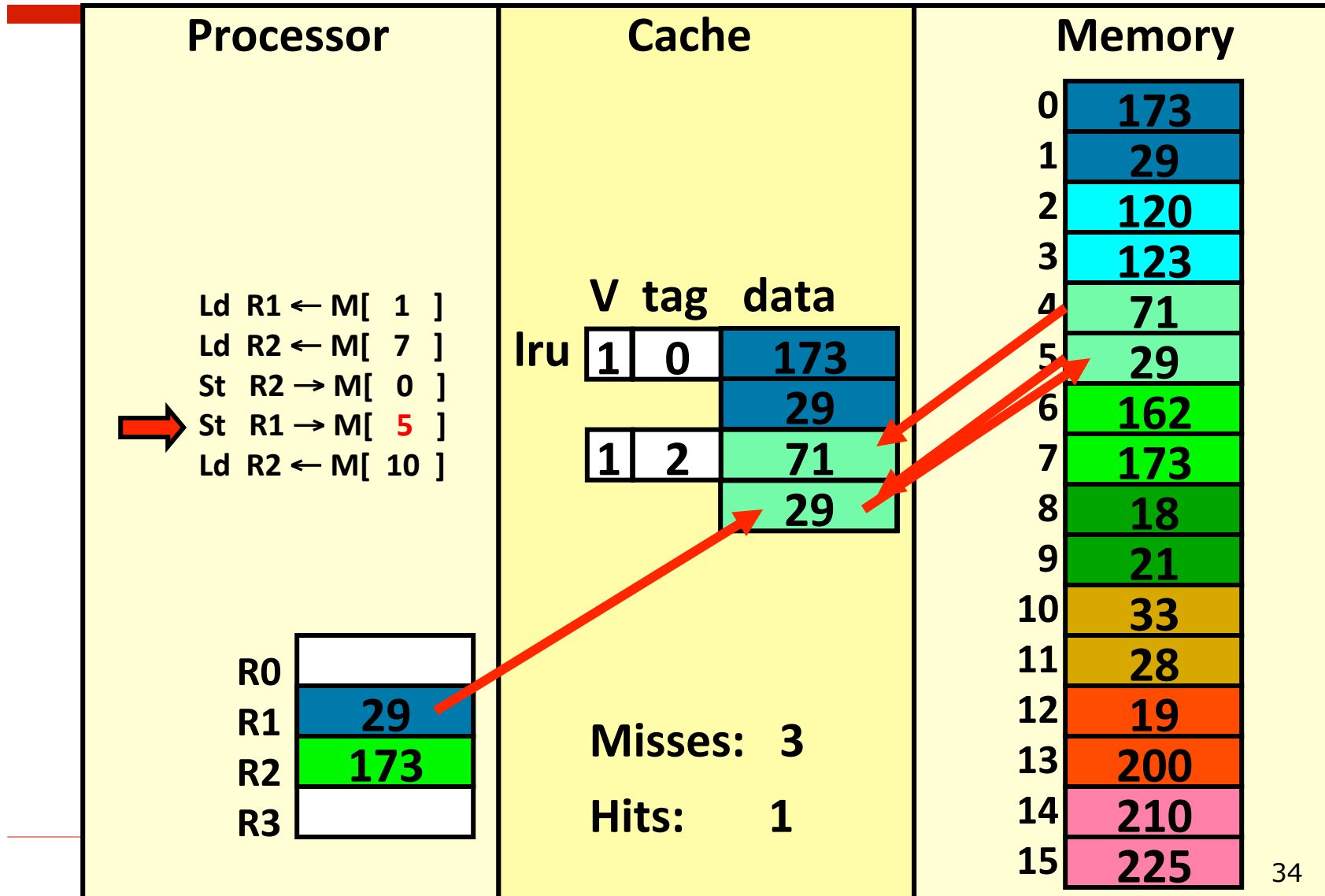




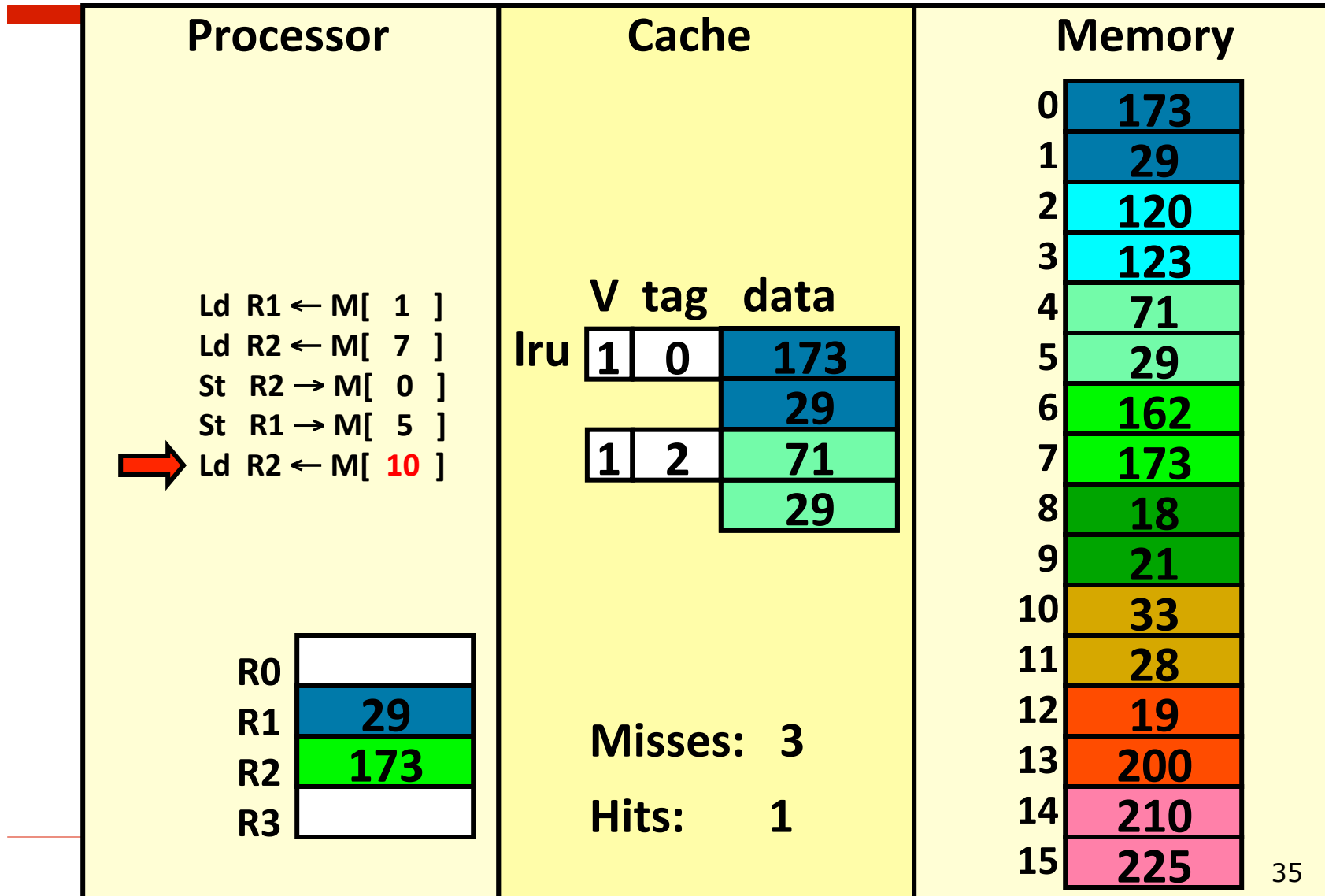
# write-through (REF 4)



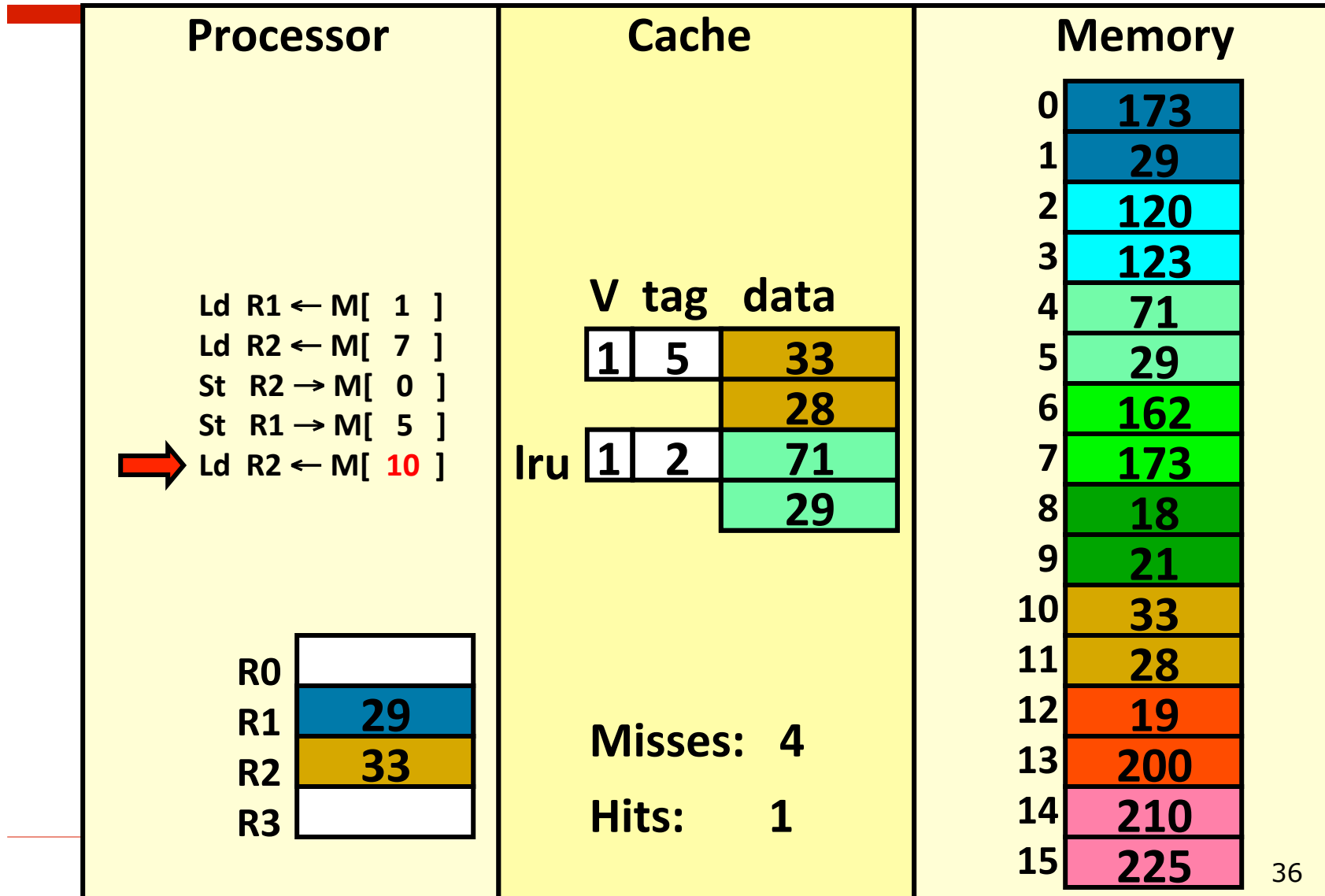
# write-through (REF 4)



# write-through (REF 6)



# write-through (REF 6)



# How many memory references?

---

- ❑ Each miss reads a block
  - 2 bytes in this cache
- ❑ Each store writes a byte
- ❑ Total reads: 8 bytes
- ❑ Total writes: 2 bytes

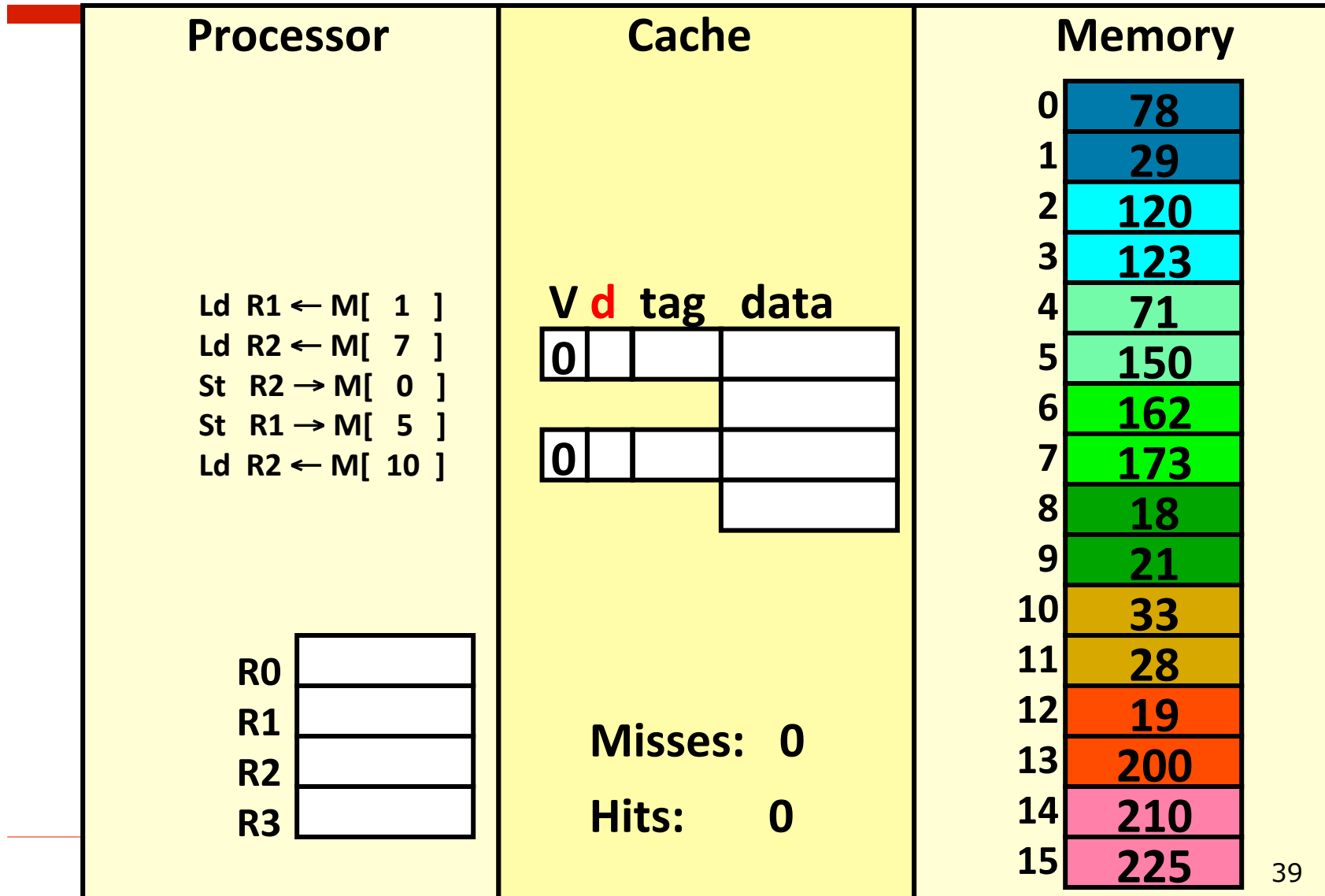
but caches generally miss < 20%

# Write-through vs. write-back

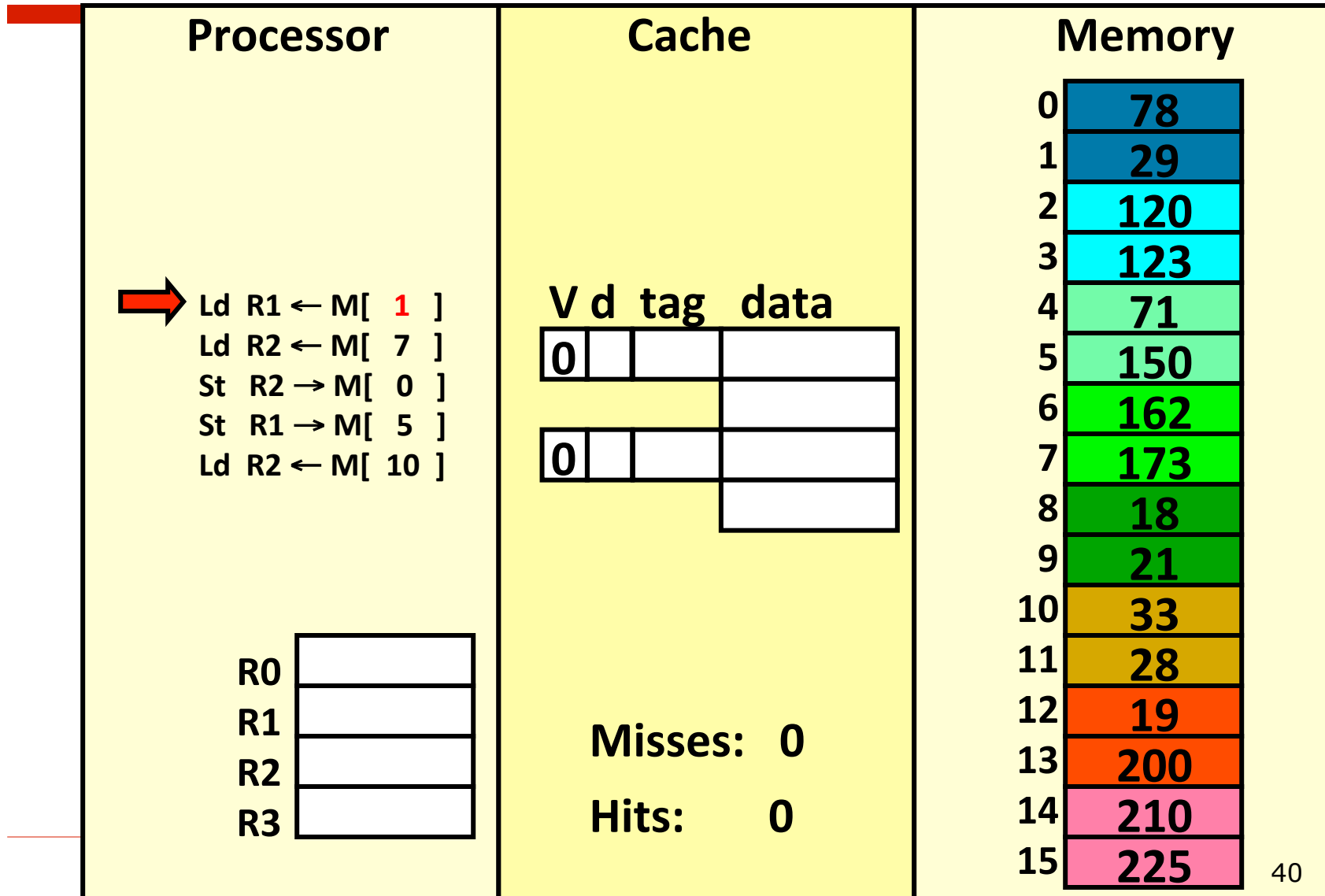
---

- ❑ Can we also design the cache to **NOT** write all stores to memory immediately?
  - We can keep the most recent copy in the cache and update the memory **only when** that data is evicted from the cache (a **write-back** policy).
  - Do we need to write-back all evicted lines?
    - No, only blocks that have been stored into
    - Keep a “**dirty bit**”, reset when the line is allocated, set when the block is stored into. If a block is “dirty” when evicted, write its data back into memory.

# Handling stores (write-back)

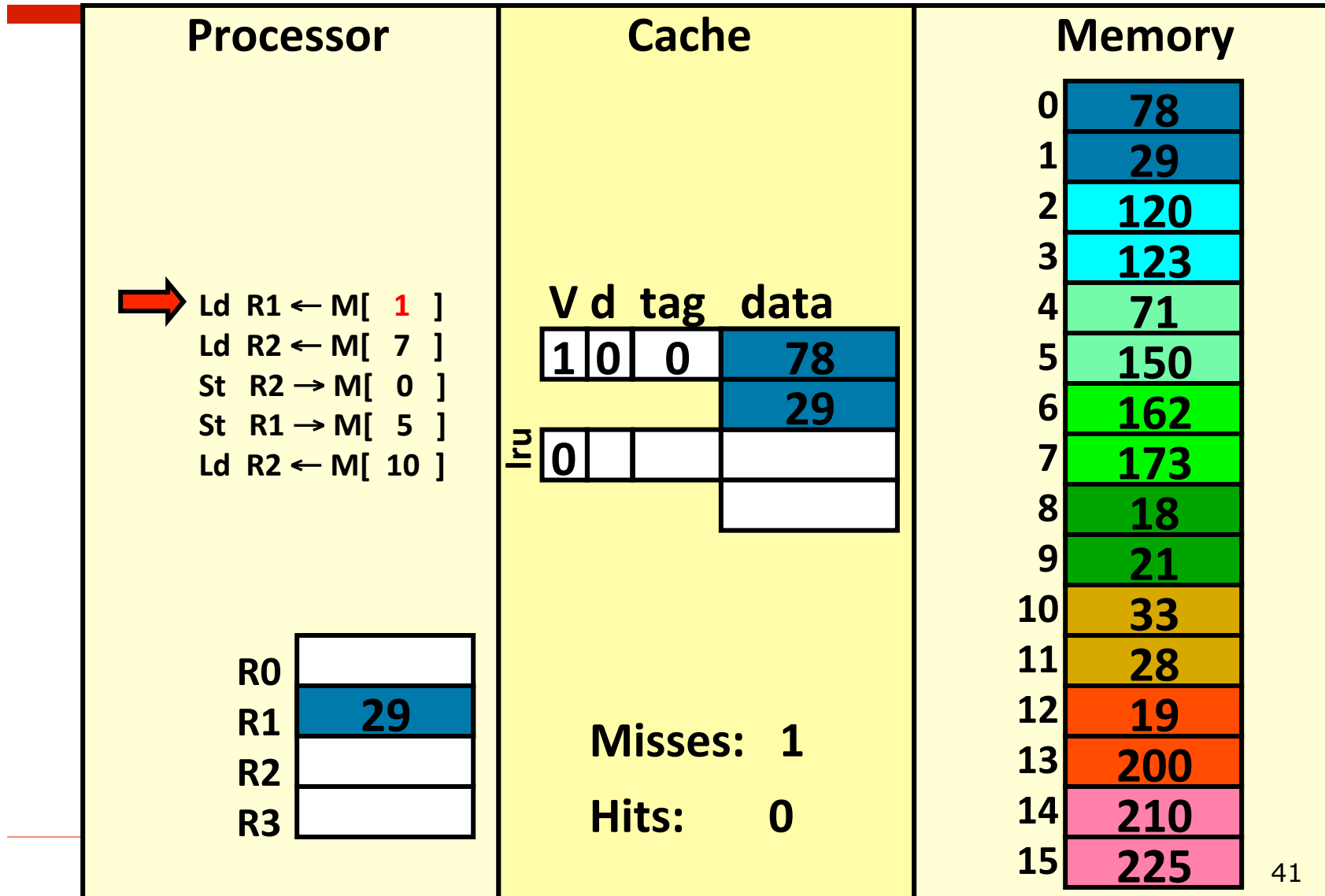


# write-back (REF 1)

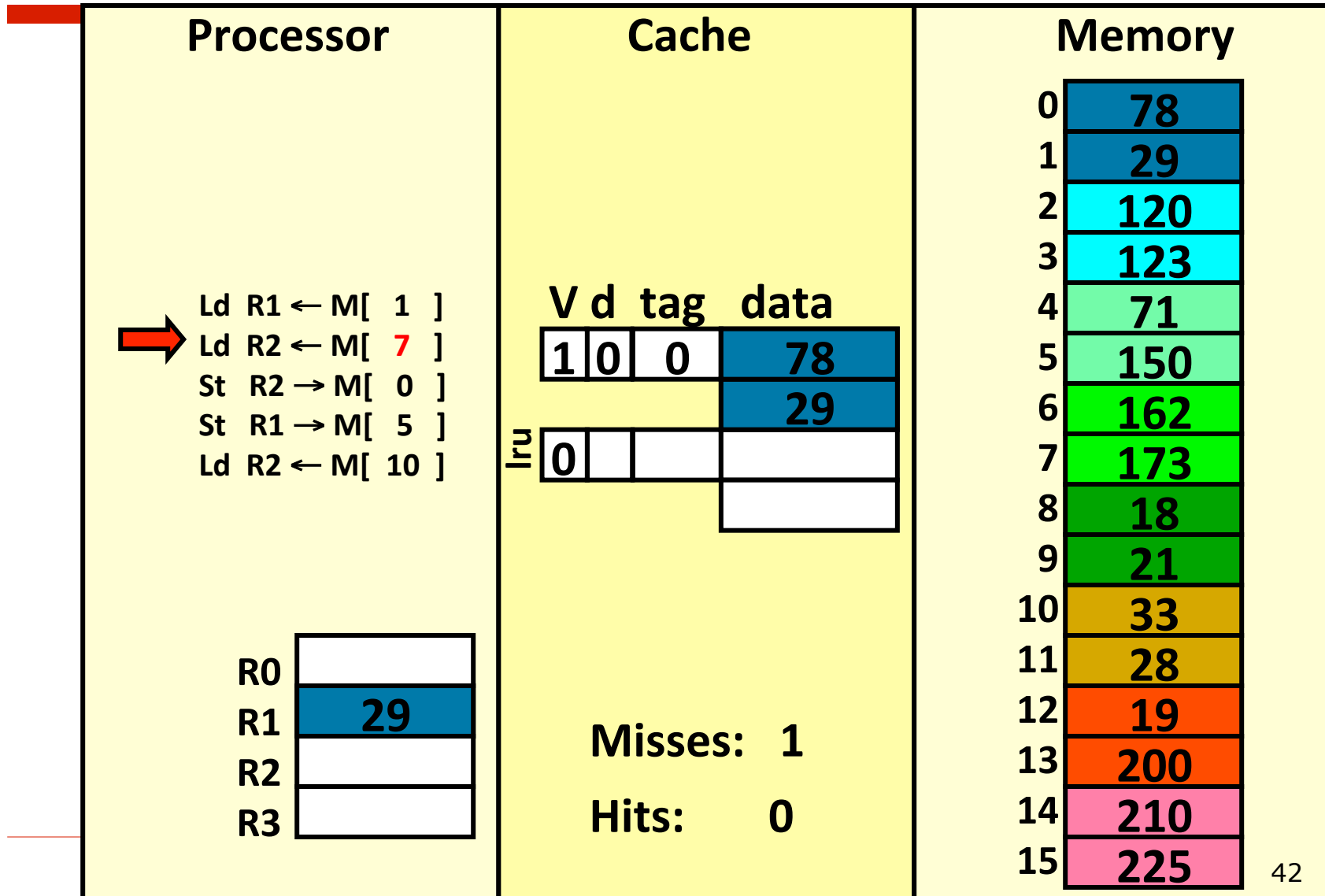




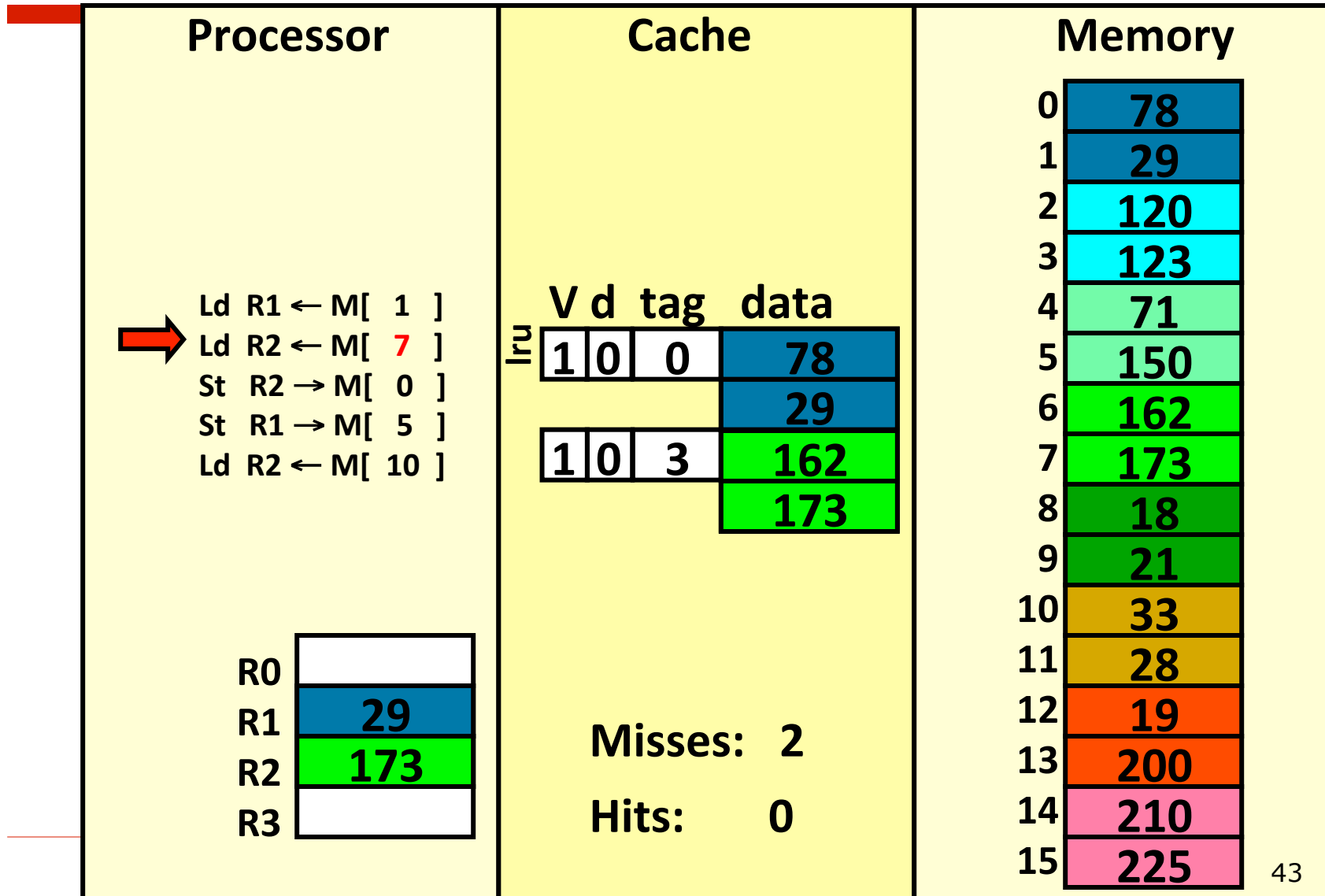
# write-back (REF 1)



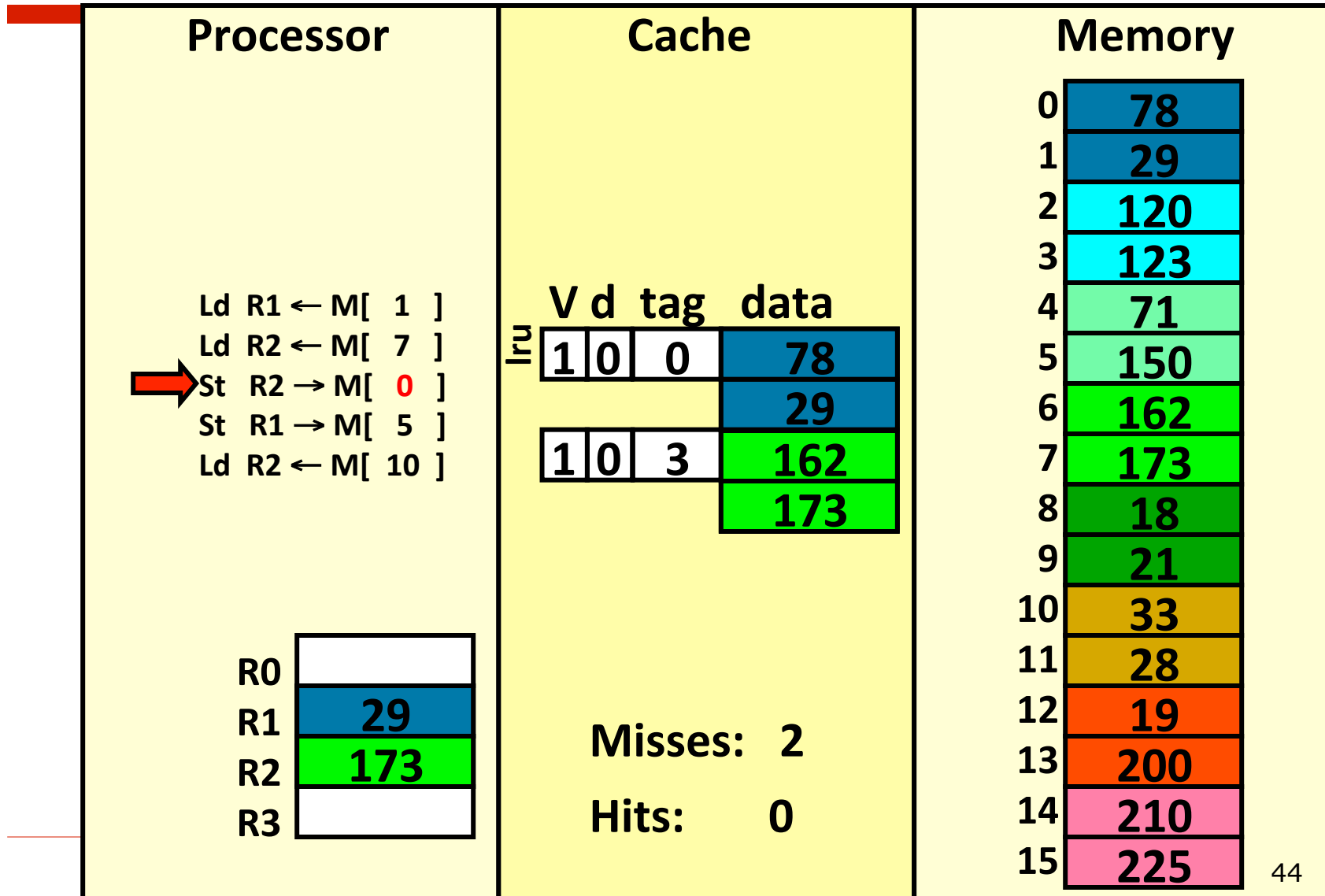
# write-back (REF 2)



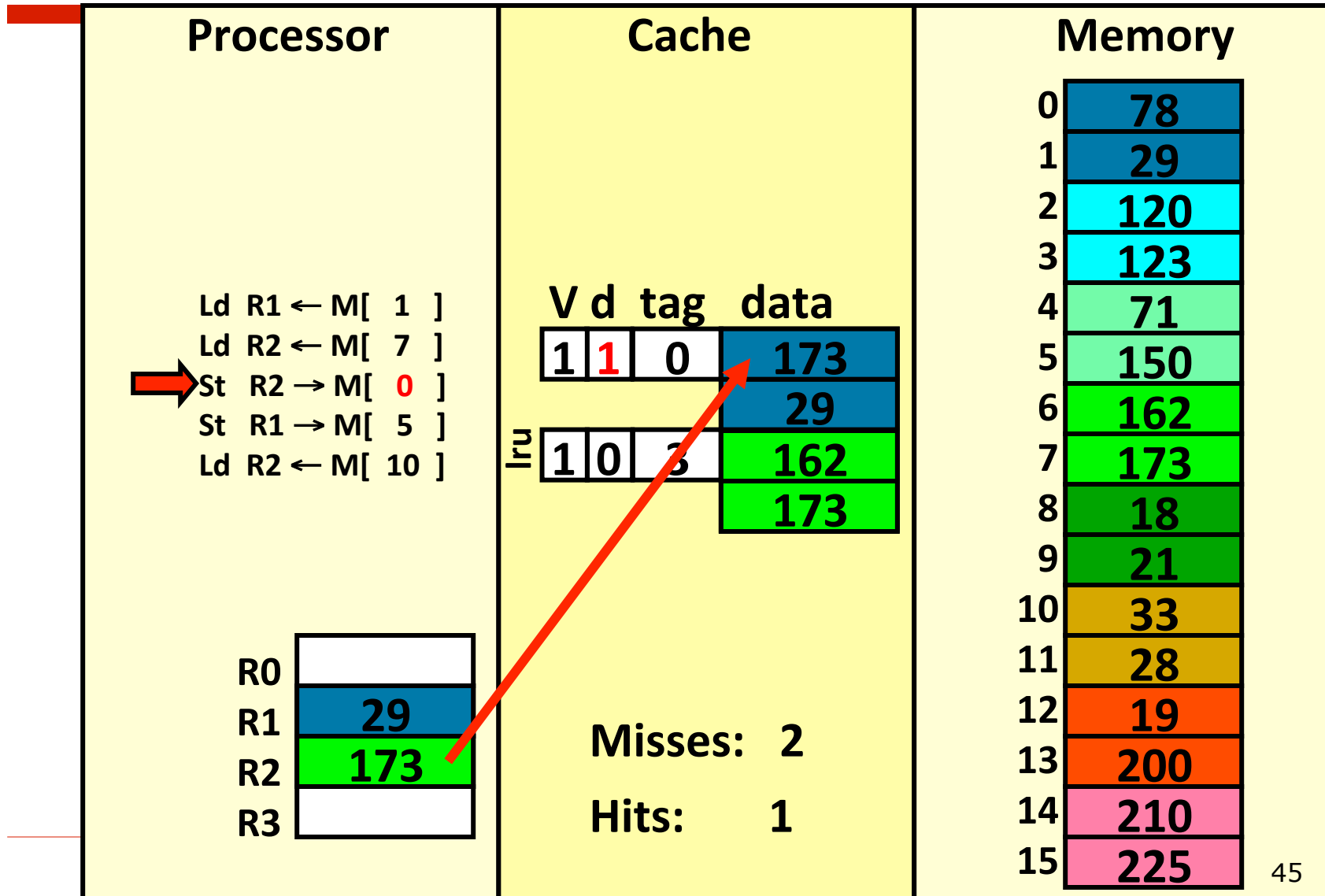
# write-back (REF 2)



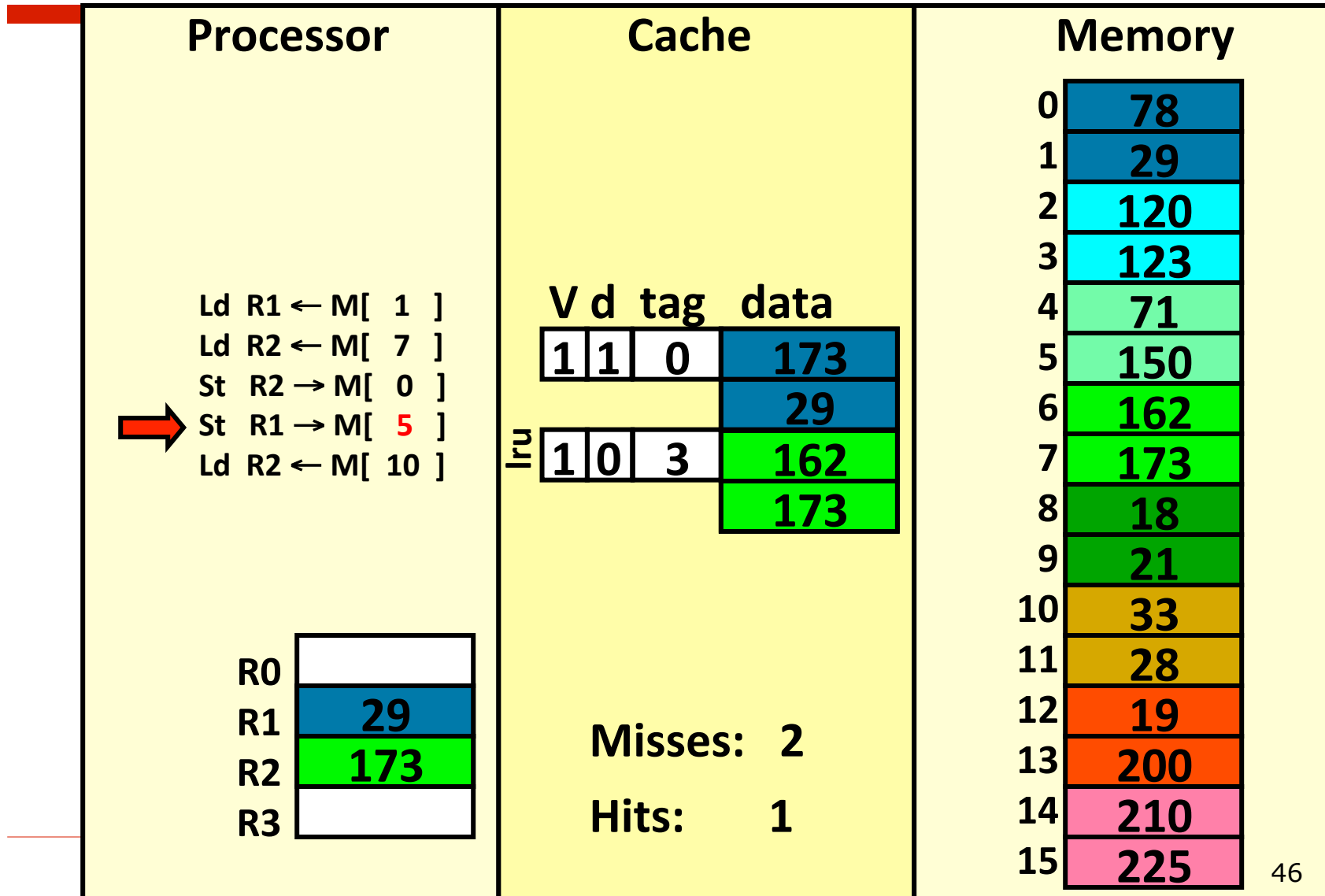
# write-back (REF 3)



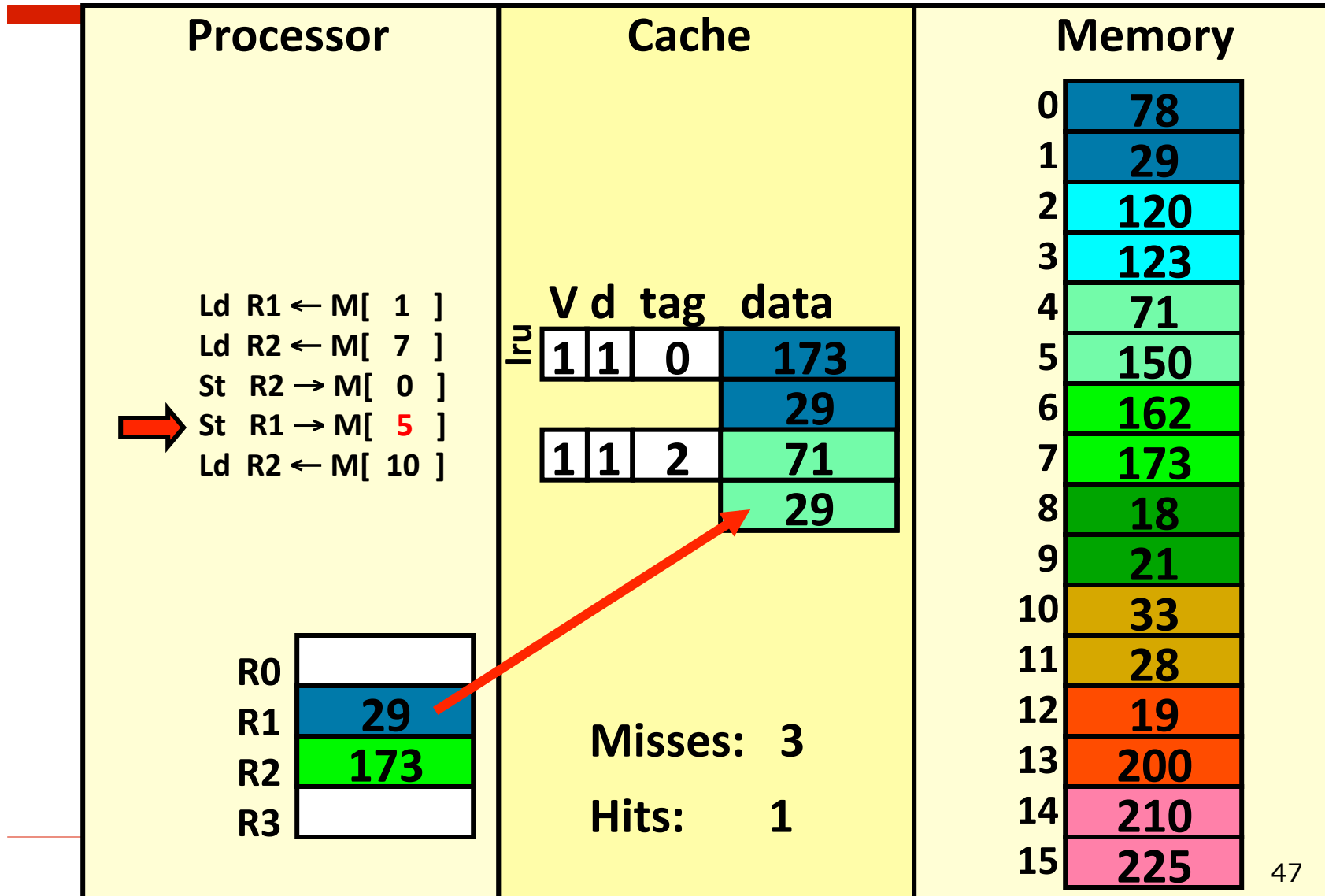
# write-back (REF 3)



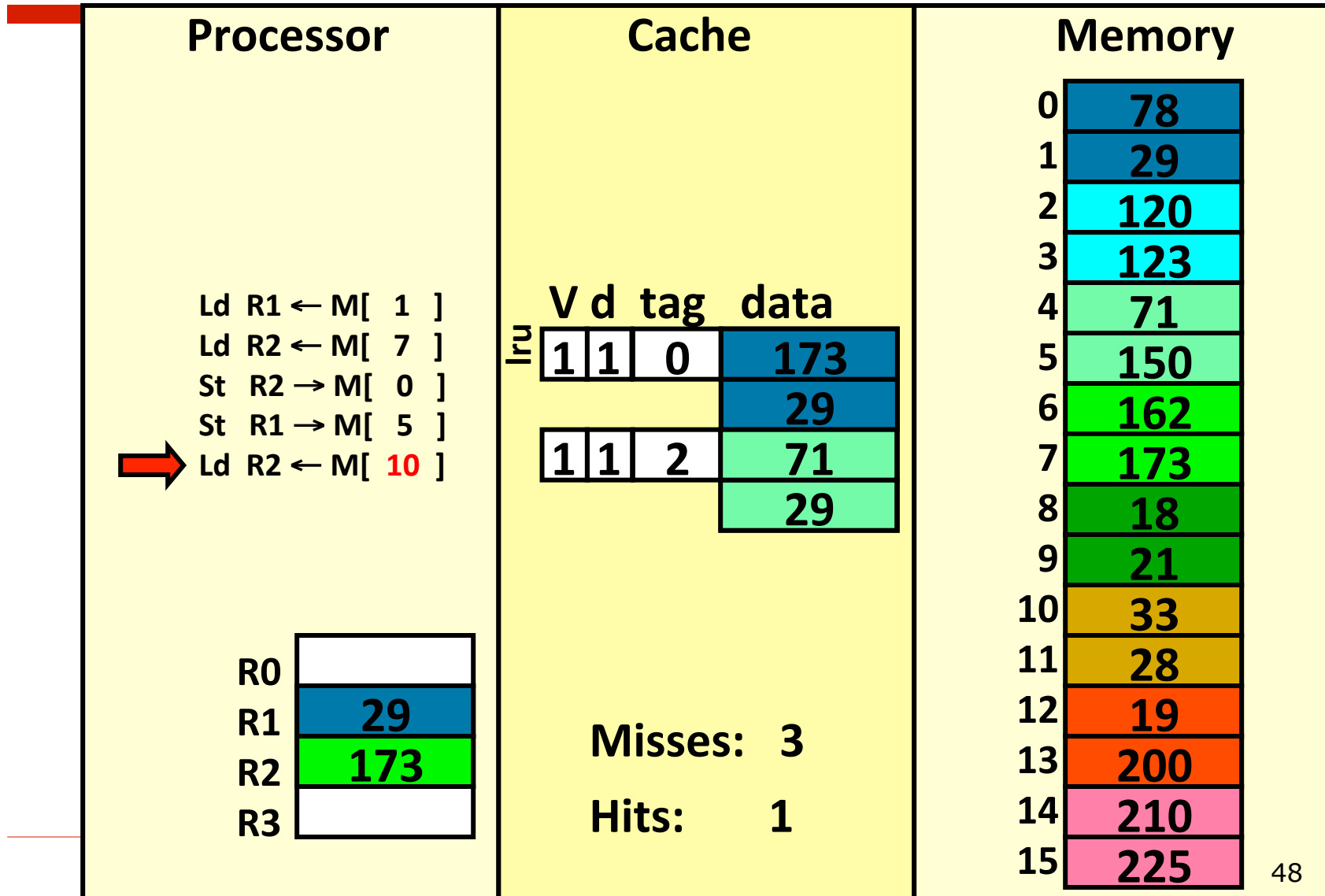
# write-back (REF 4)



# write-back (REF 4)

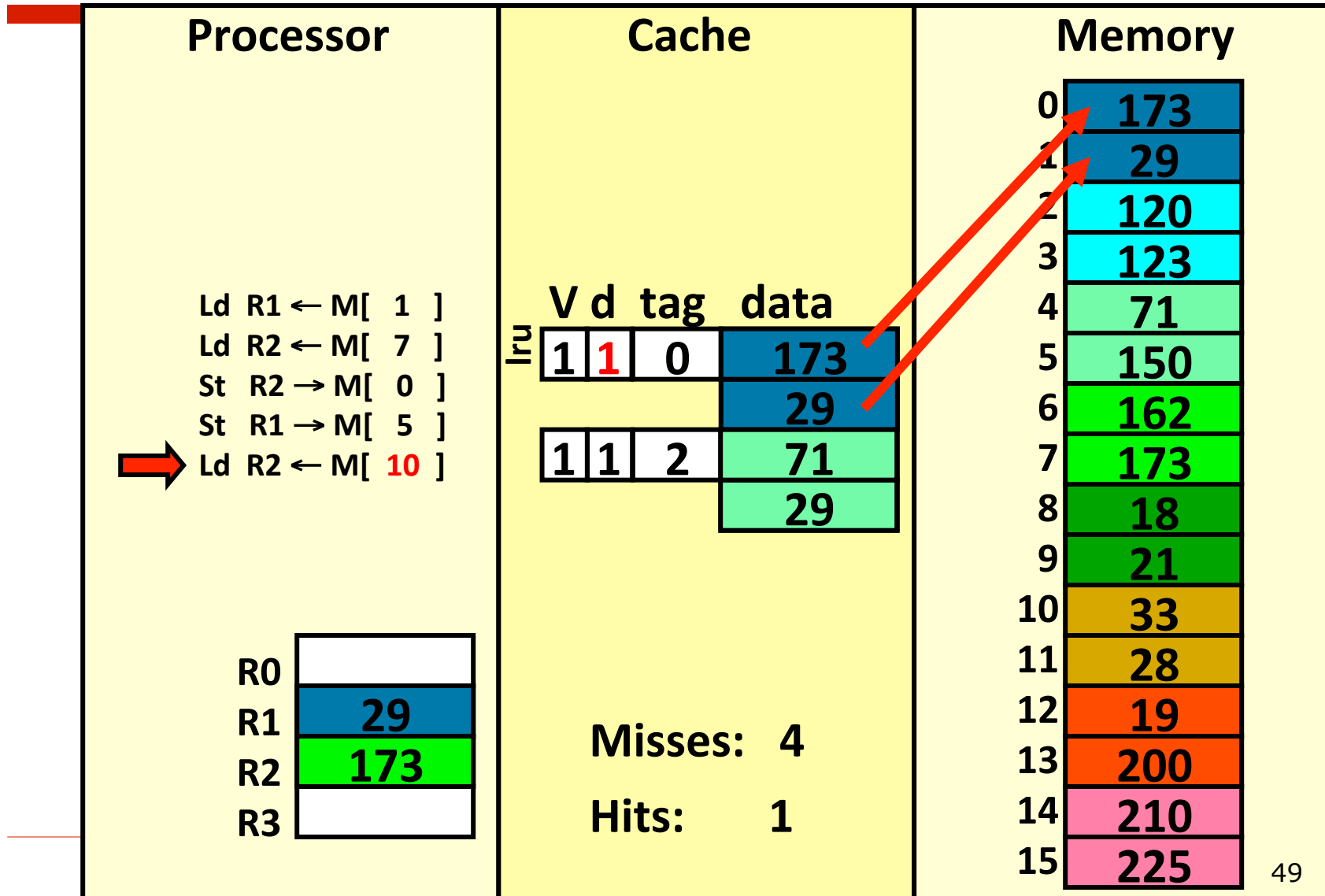


# write-back (REF 5)

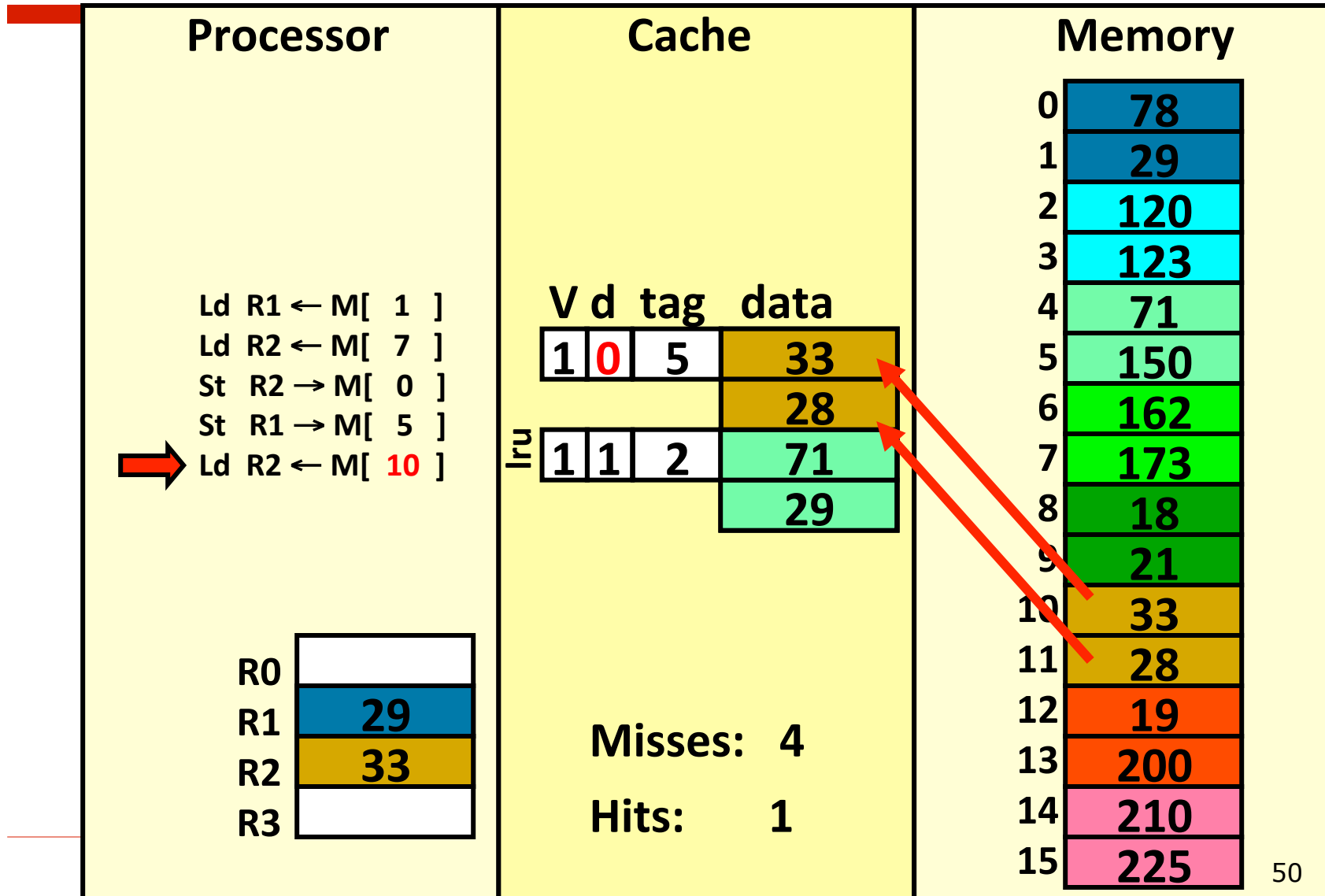




# write-back (REF 5)



# write-back (REF 5)



# How many memory references?

---

- ❑ Each miss reads a block
  - 2 bytes in this cache
- ❑ Each evicted dirty cache line writes a block
- ❑ Total reads: 8 bytes
- ❑ Total writes: 4 bytes (after final eviction)

For this example, would you choose write-back or write-through?

# Class Problem 2

---

- ❑ Consider the following cache:

32-bit memory addresses, byte addressable, 64KB cache

64B cache block size, write-allocate, write-back, *fully associative*

This cache will need 512 kilobits for the data area (64 kilobytes times 8 bits per byte). Note that in this context, 1 kilobyte = 1024 bytes (NOT 1000 bytes!)

Besides the actual cached data, this cache will need other storage. Consider tags, valid bits, dirty bits, bits to keep track of LRU, and anything else that you think is necessary.

- ❑ How many additional bits (not counting the data) will be needed to implement this cache ?

# Class Problem 2

---

- ❑ Consider the following cache:  
32-bit memory addresses, byte addressable, 64KB cache  
64B cache block size, write-allocate, write-back, *fully associative*  
This cache will need 512 kilobits for the data area (64 kilobytes times 8 bits per byte). Note that in this context, 1 kilobyte = 1024 bytes (NOT 1000 bytes!)  
Besides the actual cached data, this cache will need other storage. Consider tags, valid bits, dirty bits, bits to keep track of LRU, and anything else that you think is necessary.
- ❑ How many additional bits (not counting the data) will be needed to implement this cache ?

Tag = 32 (Address) – 6 (block offset) = 26 bits

#blocks = 64K / 64 = 1024 → LRU bits = 10

Overhead per block = 26(Tag) + 1(V) + 1(D) + 10 (LRU) = 38 bits

## Class Problem 3

---

- ❑ Suppose that accessing a cache takes 10ns while accessing main memory in case of cache-miss takes 100ns. What is the average memory access time if the cache hit rate is 97%?

$$AMAT = 10 + (1 - 0.97) * 100 = 13 \text{ ns}$$

- ❑ To improve performance, the cache size is increased. It is determined that this will increase the hit rate by 1%, but it will also increase the time for accessing the cache by 2ns. Will this improve the overall average memory access time?

$$AMAT = 12 + (1 - 0.98) * 100 = 14 \text{ ns}$$