

EECS 370 - Lecture 18

Cache Blocks, Writeback, and Associativity



Announcements

- P3 published
 - Checkpoint due **today**
 - 5% of project – have pipeline working without data hazards or branches
- Midterm scores published **soon**
- HW 4 due Monday

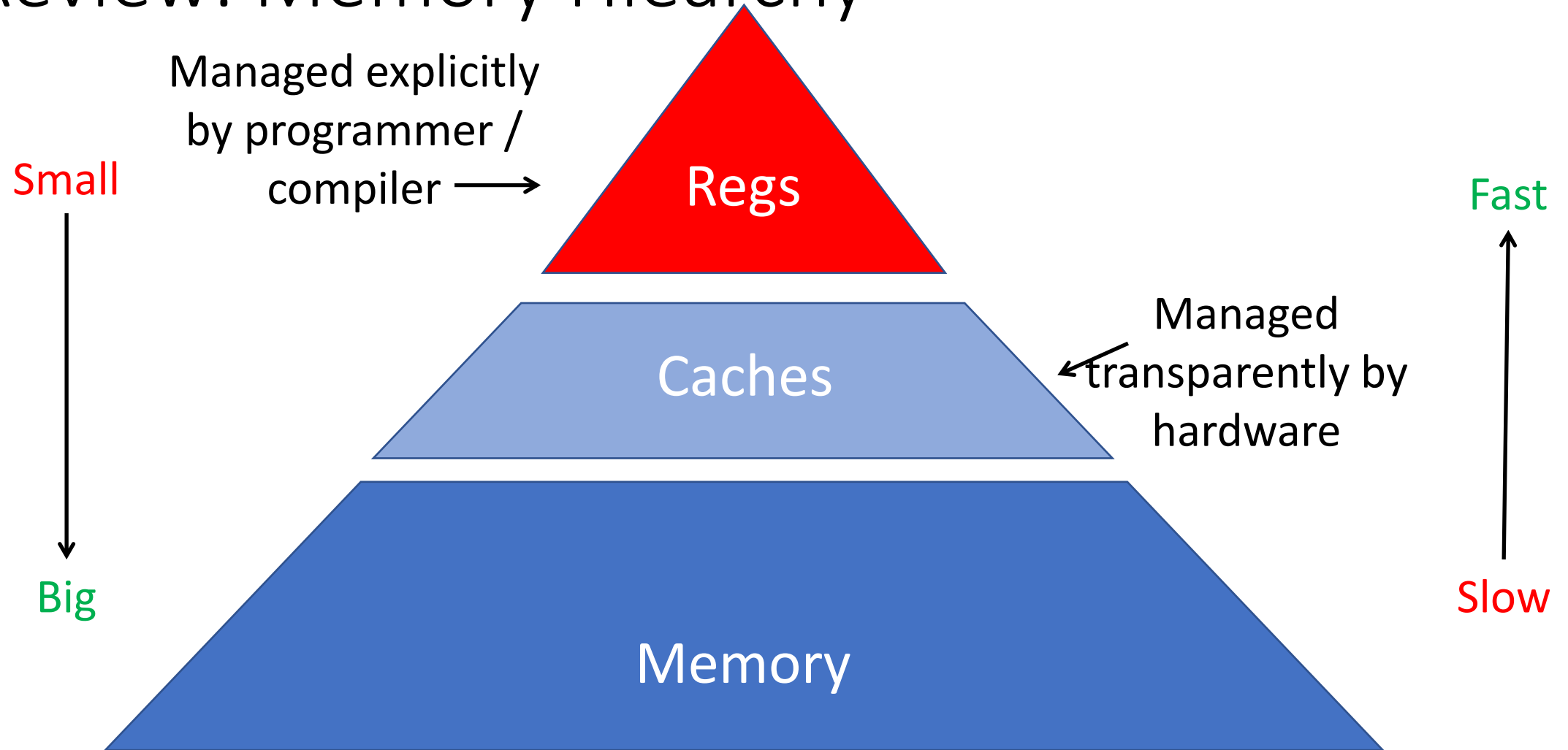
Resources

- Want extra examples with pipelining? Try playing with the "Pipeline Simulator" under "Resources" on the website
 - <https://vhosts.eecs.umich.edu/370simulators/pipeline/simulator.html>
 - Several pre-written programs you can step through to understand what's going on
 - Note that the project pipeline is slightly different

Review: Memory Hierarchy

- Problem: Large memory is much slower than processor
 - We'd have to wait ~100s of cycles for every load
- Observation:
 1. Small memory can be made as fast as processor
 2. We only need a small amount of memory at a time
- Idea: let's include a small amount of memory – a **cache** – that holds data hardware thinks is likely to be needed in the near future
 - Check the cache before going to main memory

Review: Memory Hierarchy



Increasing Block Size

Case 1:

Block size: 1 bytes

1	0	74
1	6	160

V tag data (block)

How many bits needed per tag?

$$= \log_2(\text{number of blocks in memory}) = \log_2(16)$$

$$= 4 \text{ bits}$$

$$\text{Overhead} = (4+1) / 8 = 62.5\%$$

Case 2:

Block size: 2 bytes

1	0	74	110
1	3	160	170

V tag data (block)

How many bits needed per tag?

$$= \log_2(\text{number of blocks in memory}) = \log_2(8)$$

$$= 3 \text{ bits}$$

$$\text{Overhead} = (3+1) / 16 = 25\%$$

Memory	Tag (case 1)	Tag (case 2)
0	0	0
1	1	0
2	2	1
3	3	1
4	4	2
5	5	2
6	6	3
7	7	3
8	8	4
9	9	4
10	10	5
11	11	5
12	12	6
13	13	6
14	14	7
15	15	7

Poll: What will the overhead of this cache be?

Figuring out the tag

Poll: What's the pattern?

- If block size is N, what's the pattern for figuring out the tag from the address?

$$\text{tag} = \left\lfloor \frac{\text{addr}}{\text{block size}} \right\rfloor$$

- If block size is power of 2, then this is just everything except the $\log_2(\text{block size})$ bits of the address in binary!

- E.g.

$$0d11 = 0b1011$$

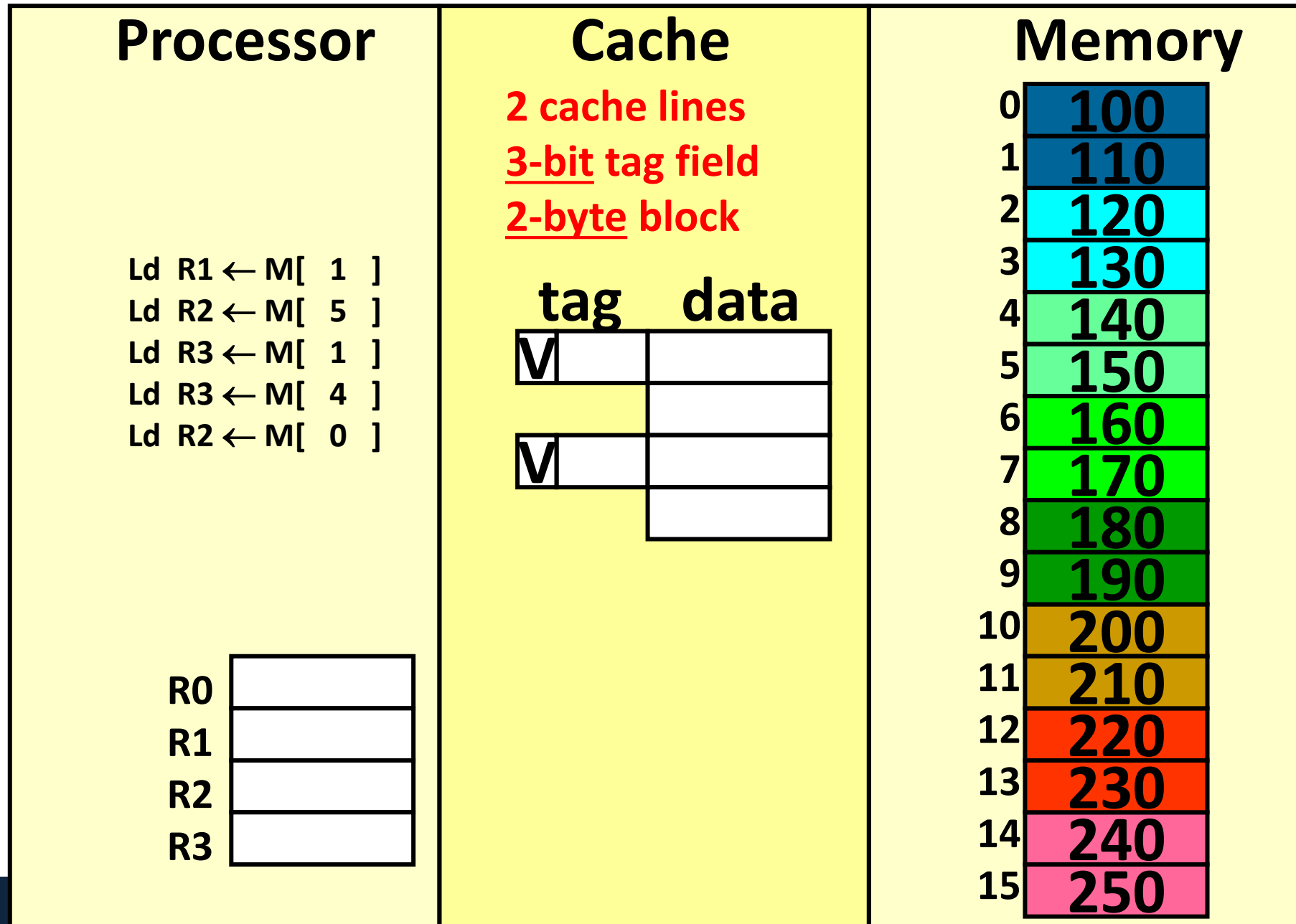
$$\text{Tag} = 0b101 = 0d5$$

$$\text{Block Offset} = 1$$

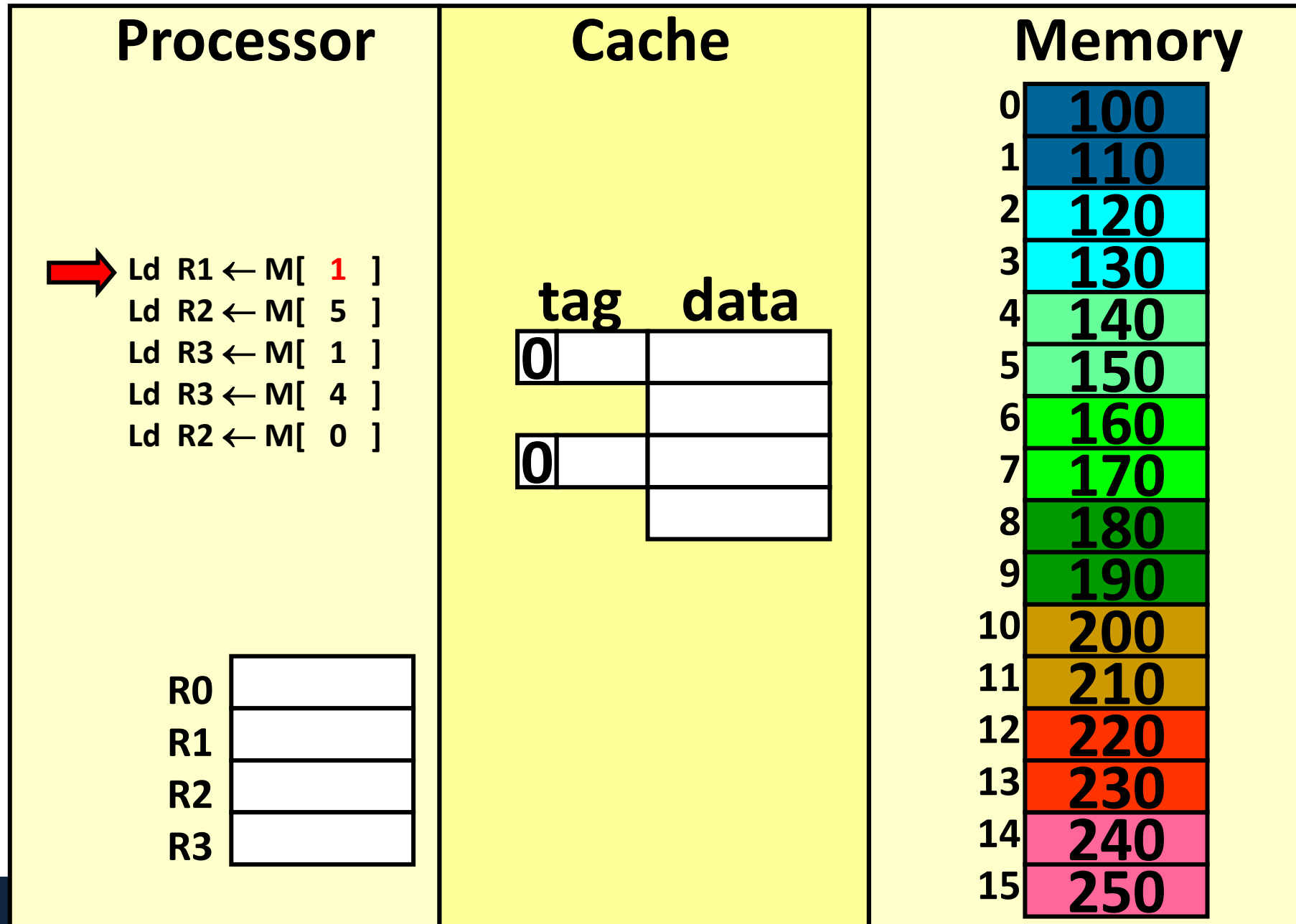
- Remaining bits (block offset) tells us how far into the block the data is

Memory		Tag (case 2)
0	74	0
1	110	0
2	120	1
3	130	1
4	140	2
5	150	2
6	160	3
7	170	3
8	180	4
9	190	4
10	200	5
11	210	5
12	220	6
13	230	6
14	240	7
15	250	7

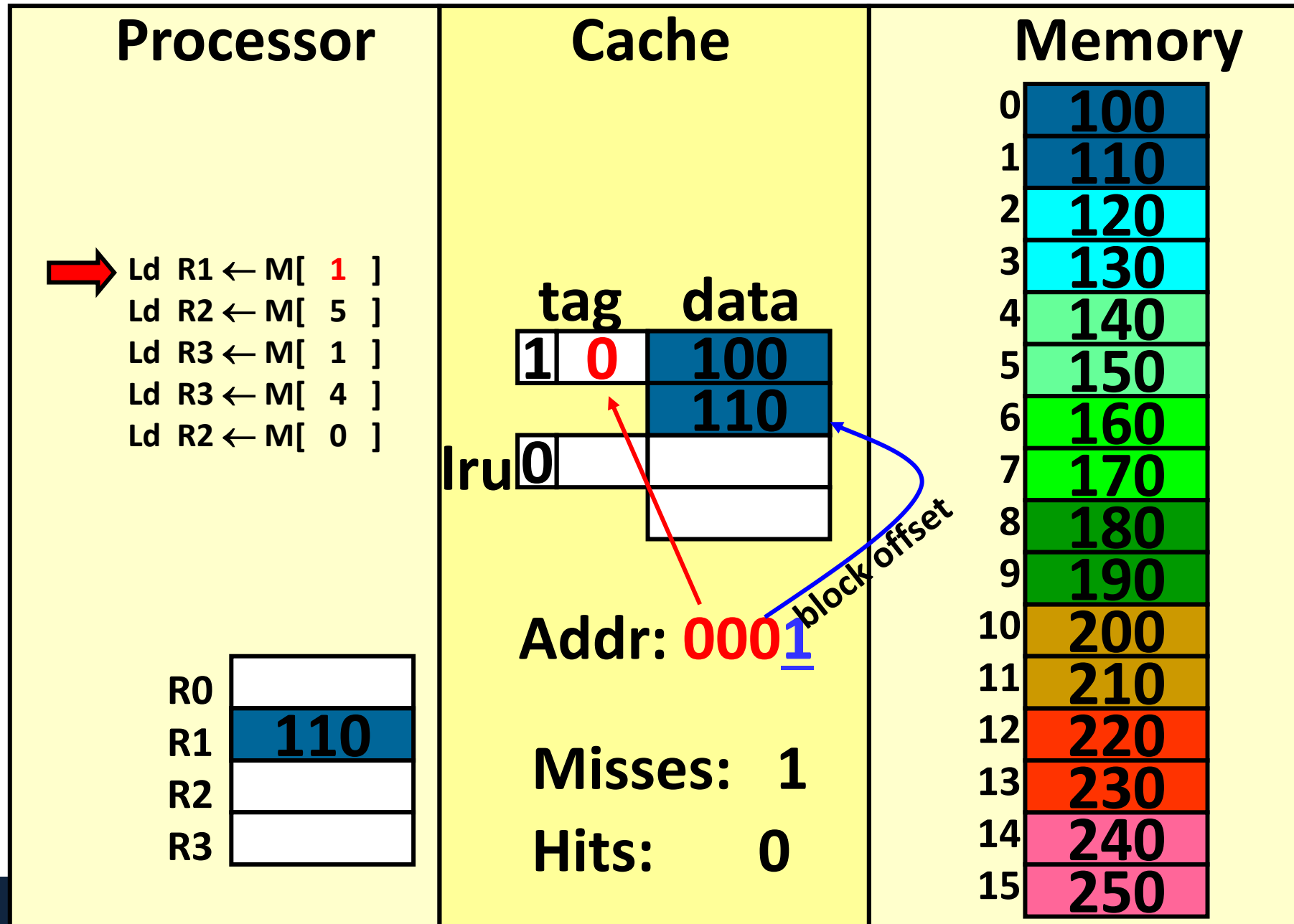
Block size for caches



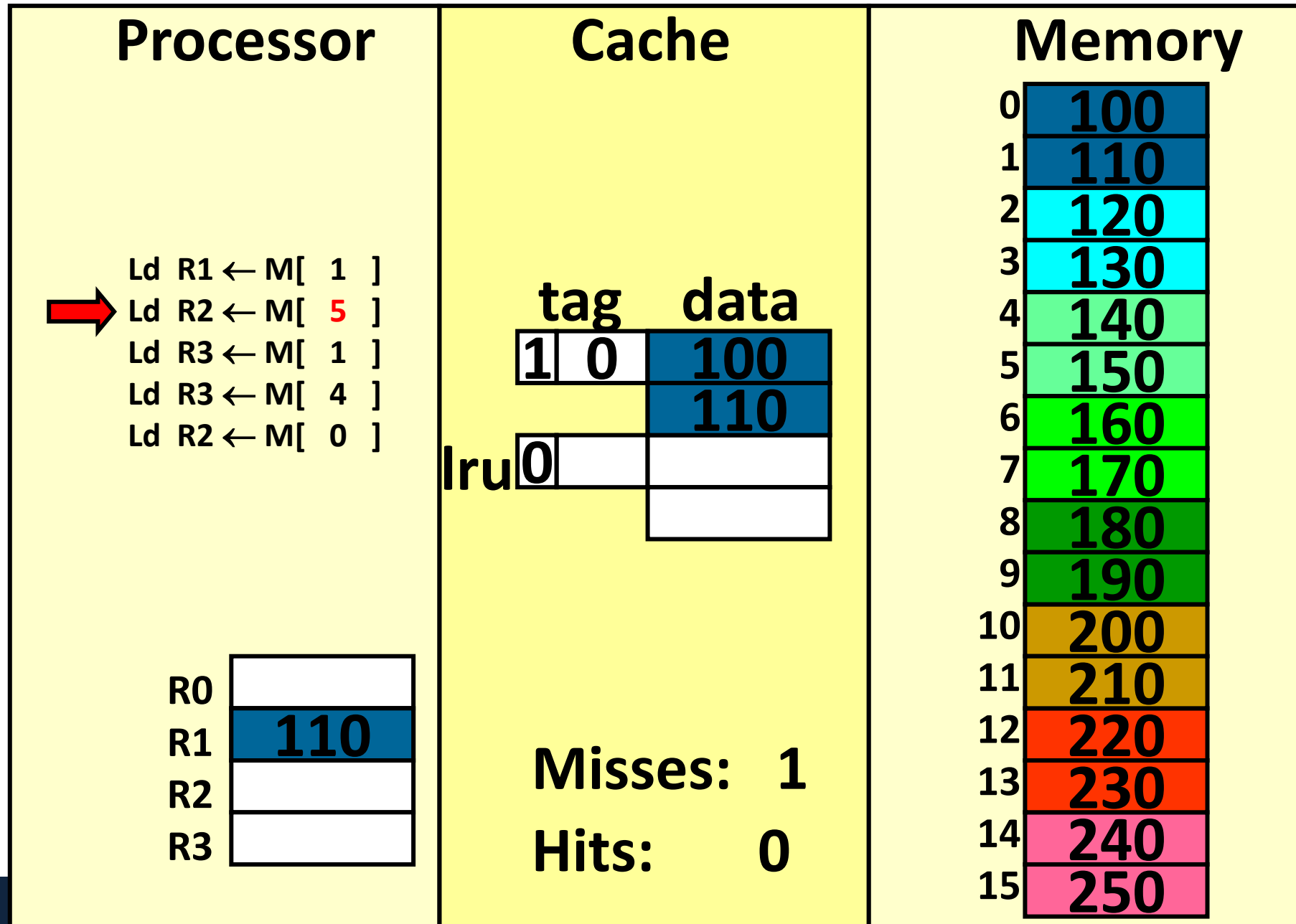
Block size for caches



Block size for caches

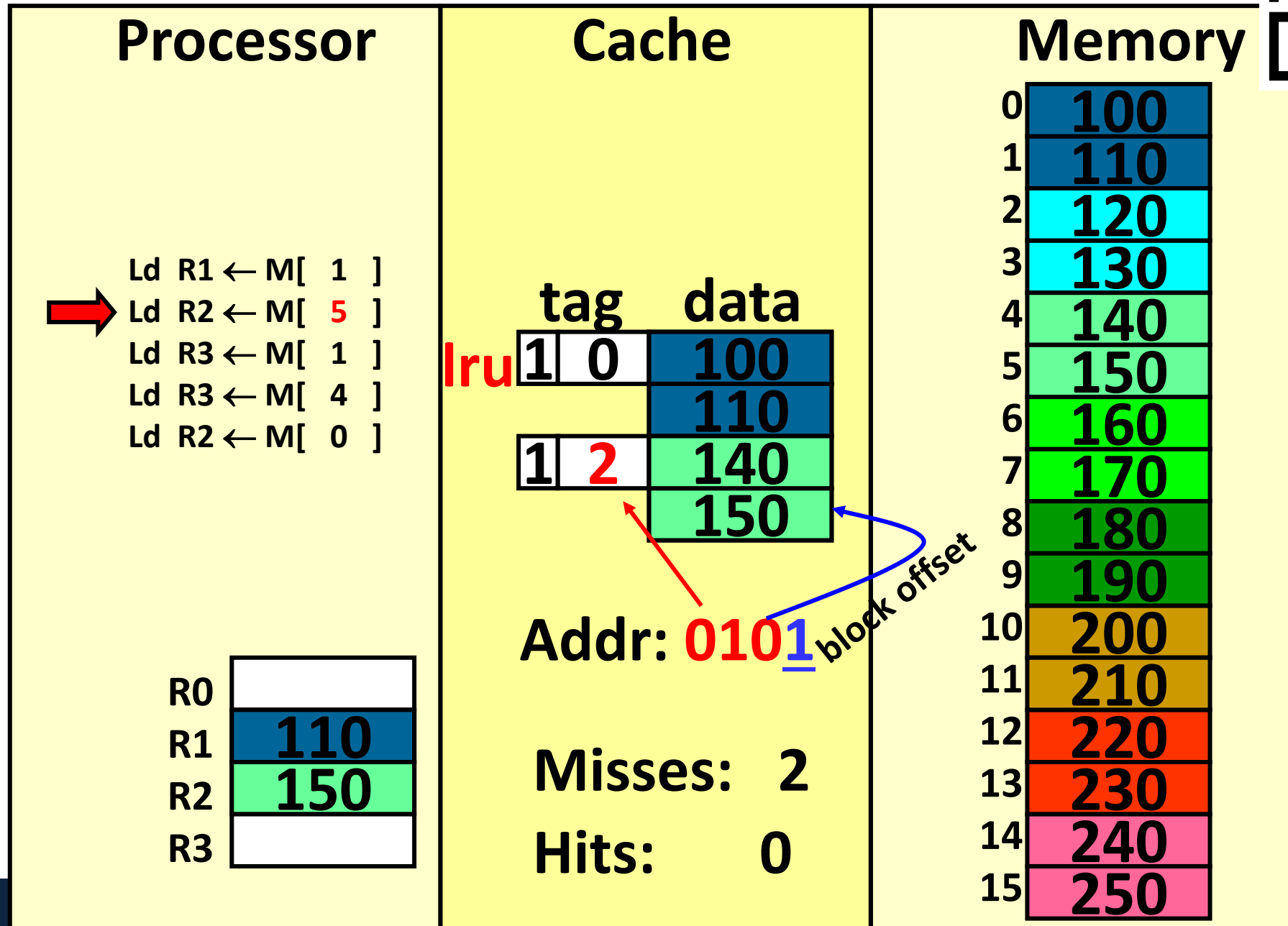
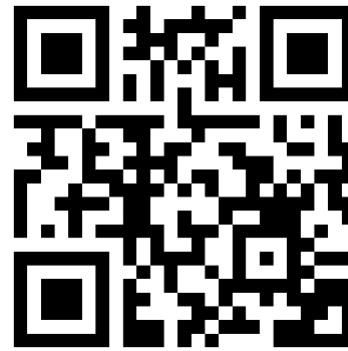


Block size for caches

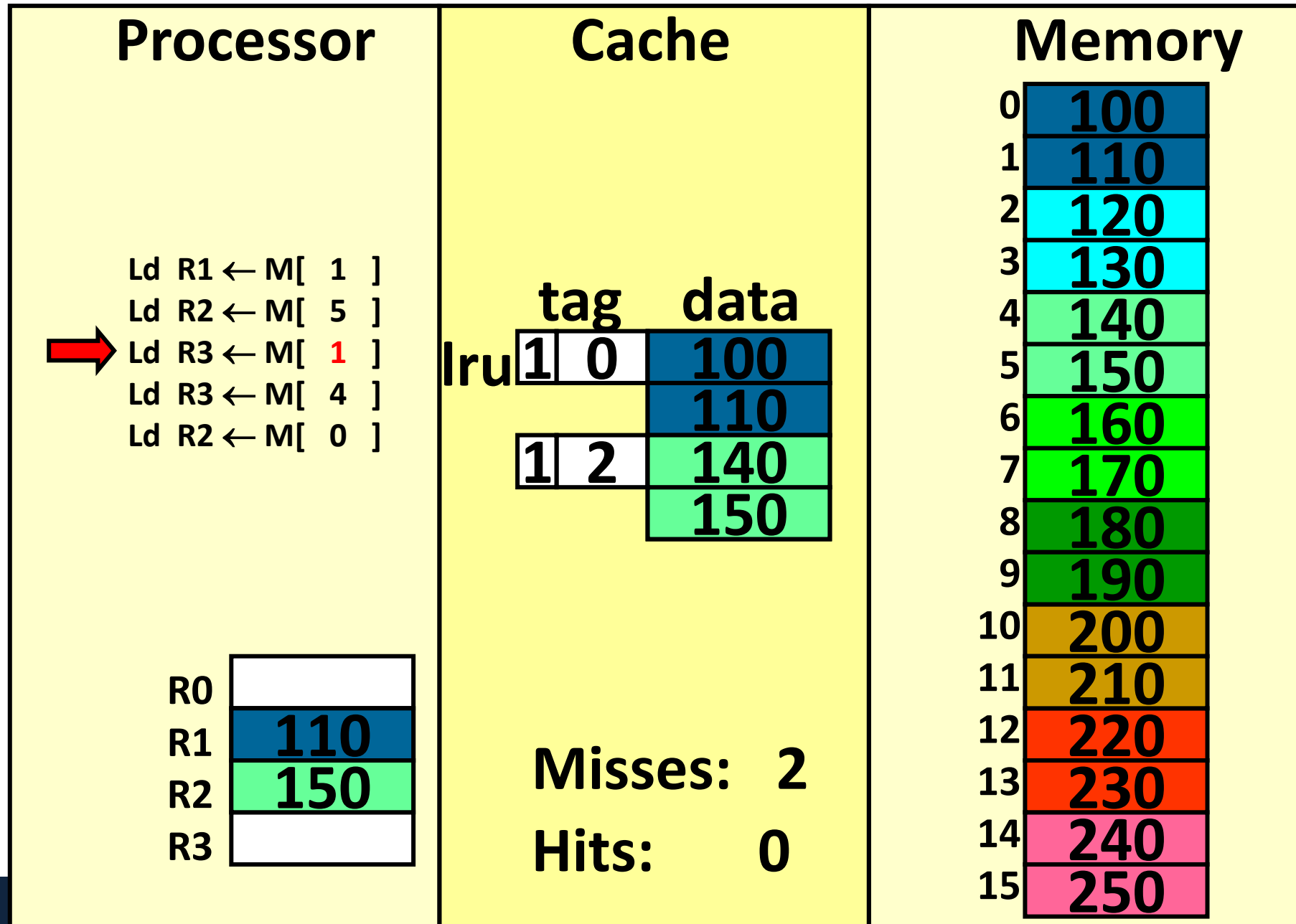


Block size for caches

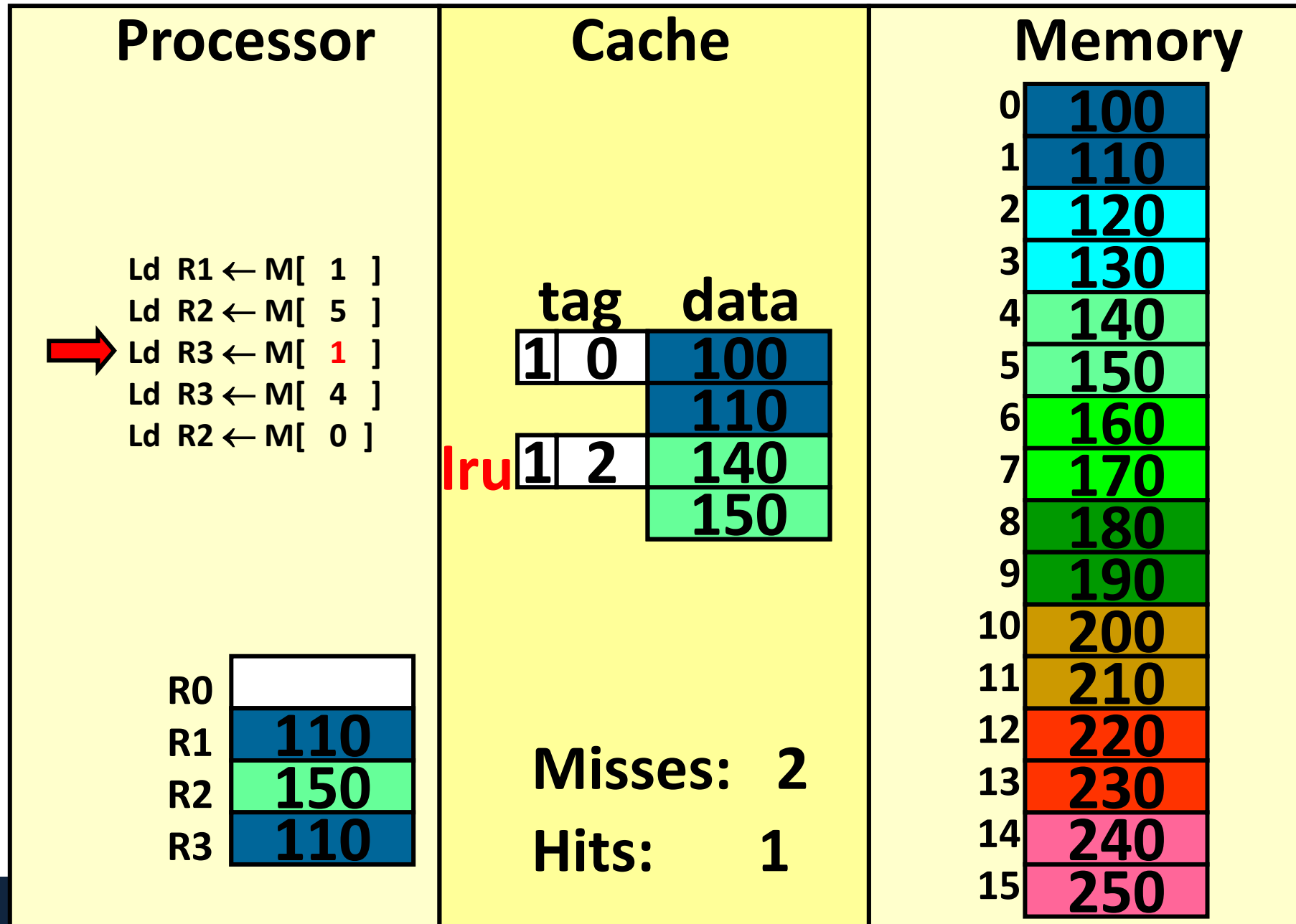
Poll: Complete the last 3 instructions yourself



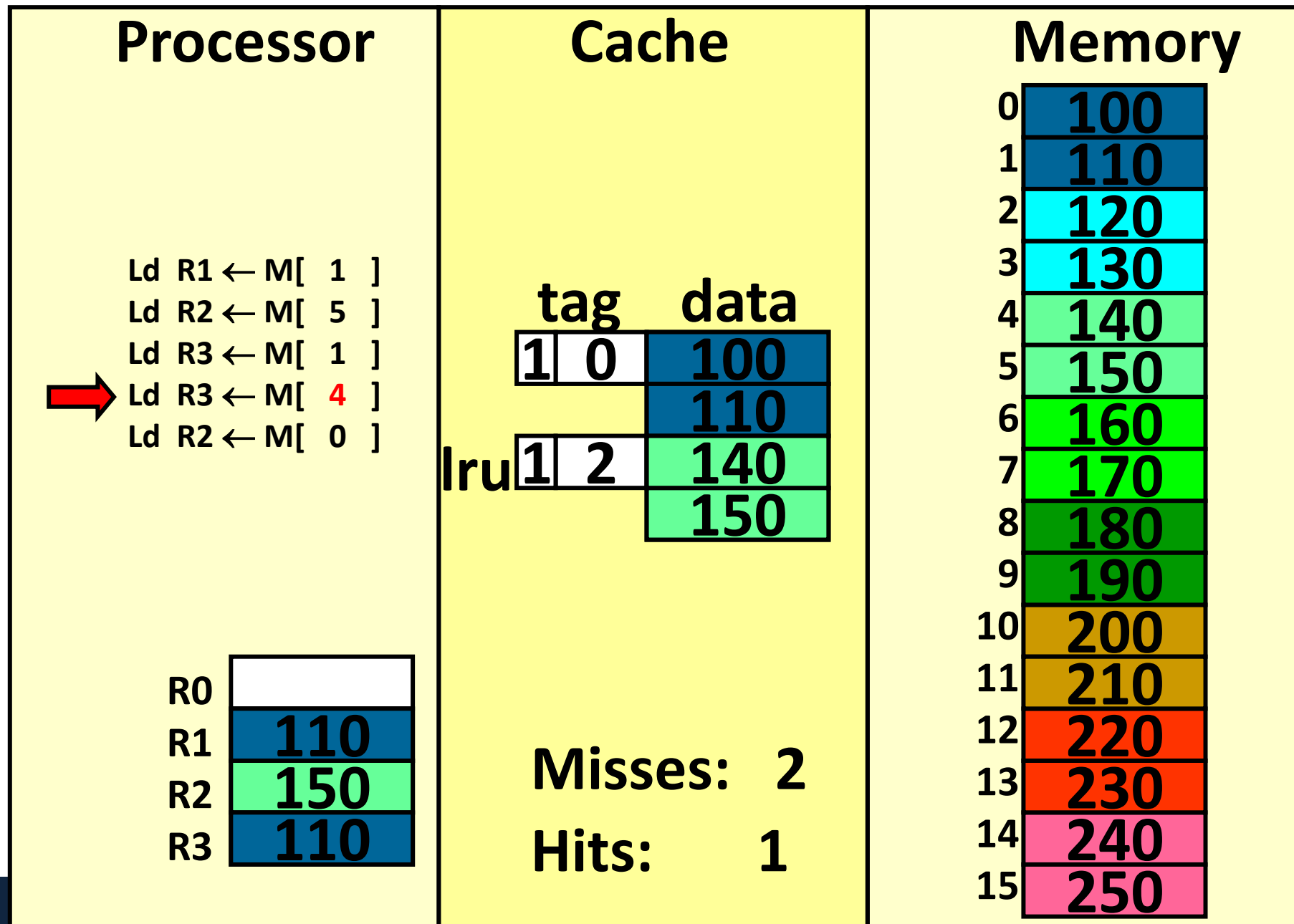
Block size for caches



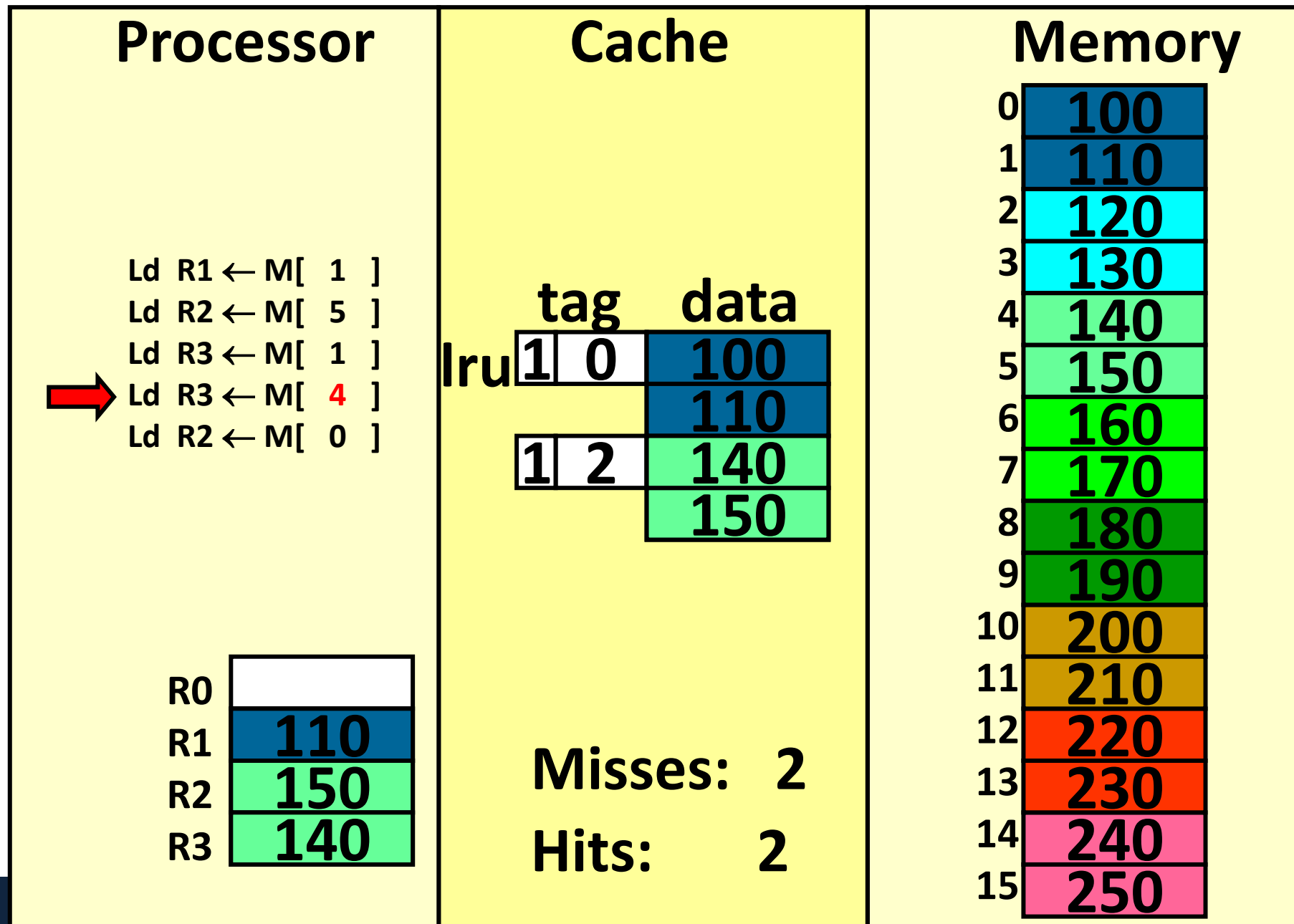
Block size for caches



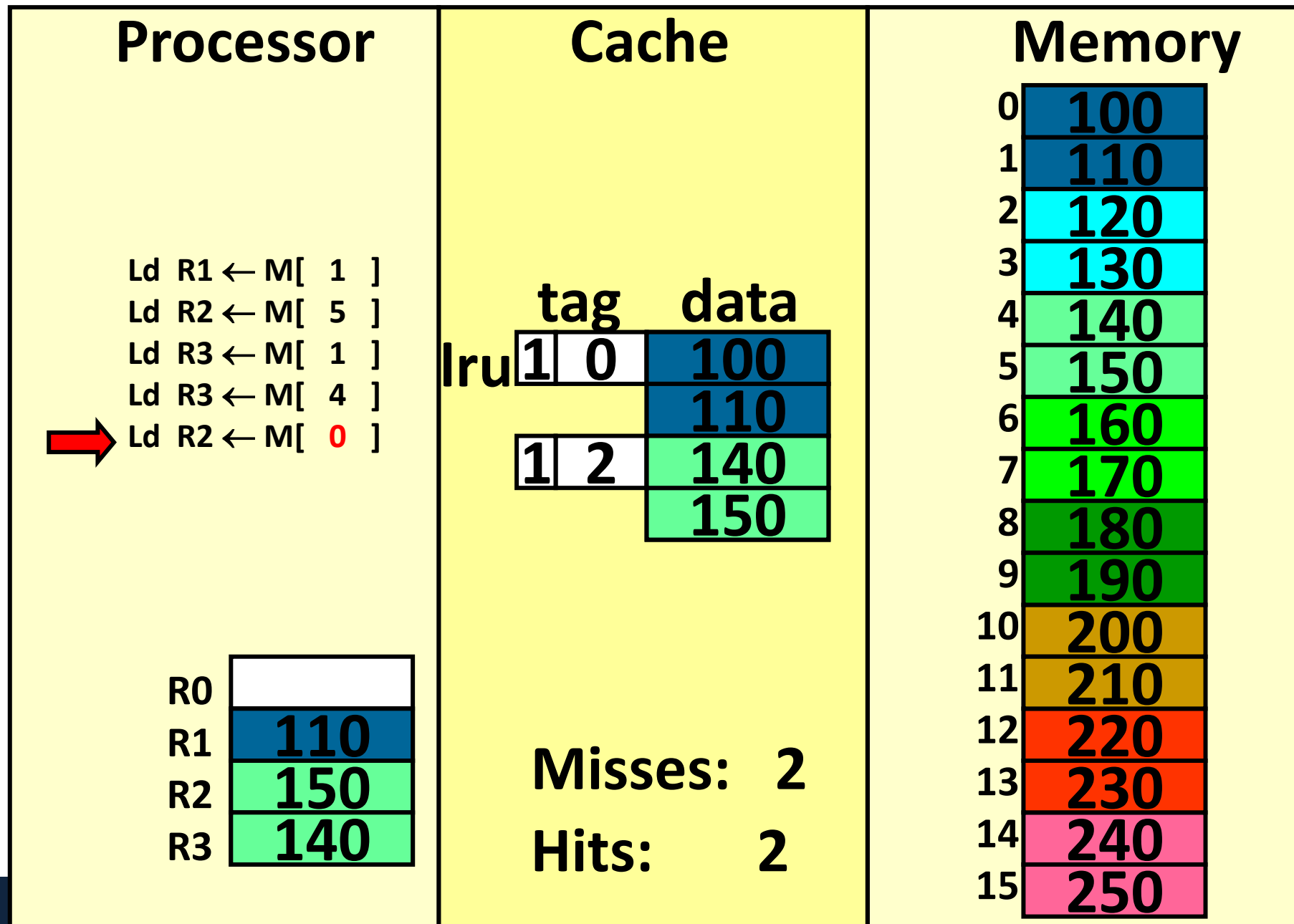
Block size for caches



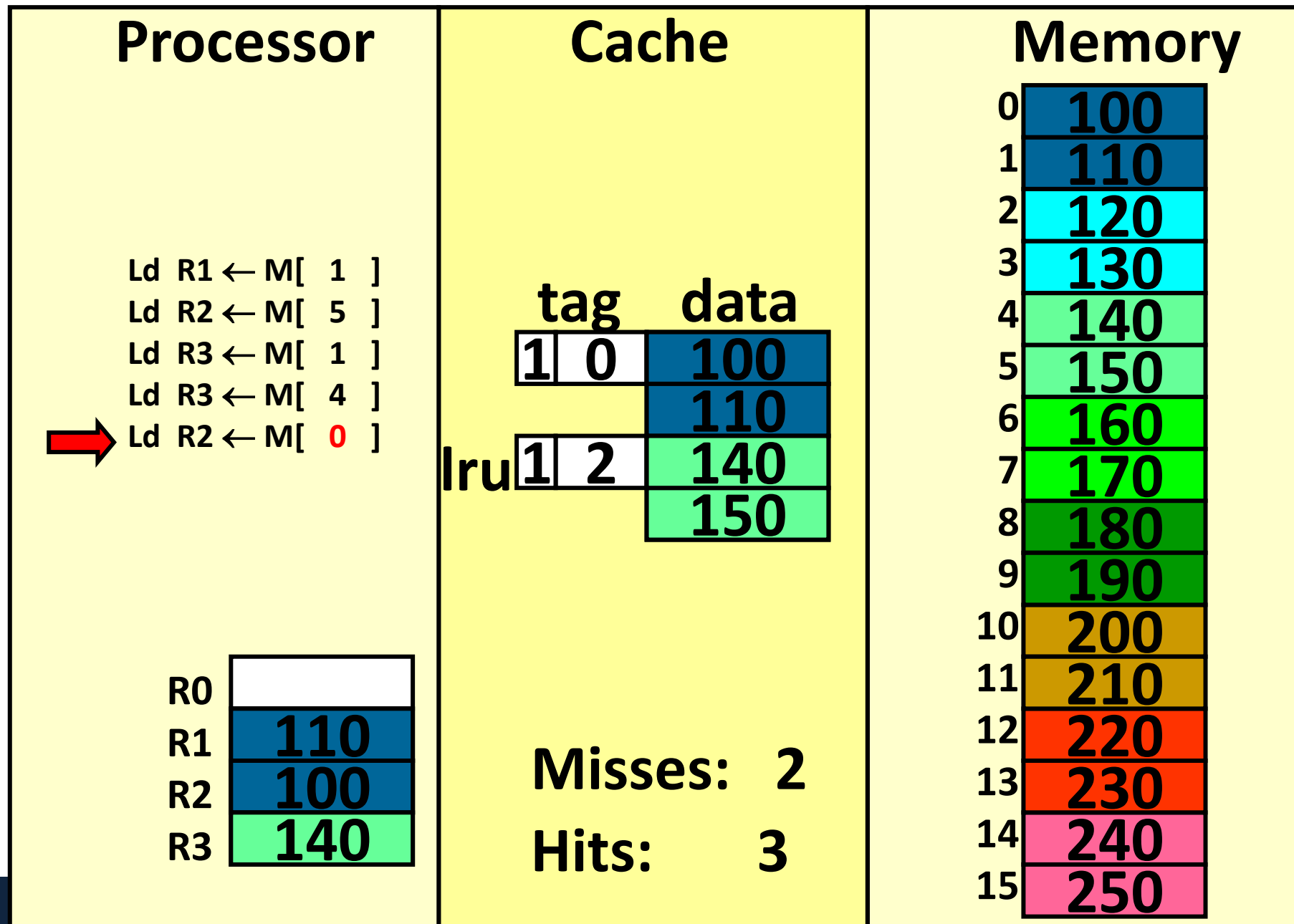
Block size for caches



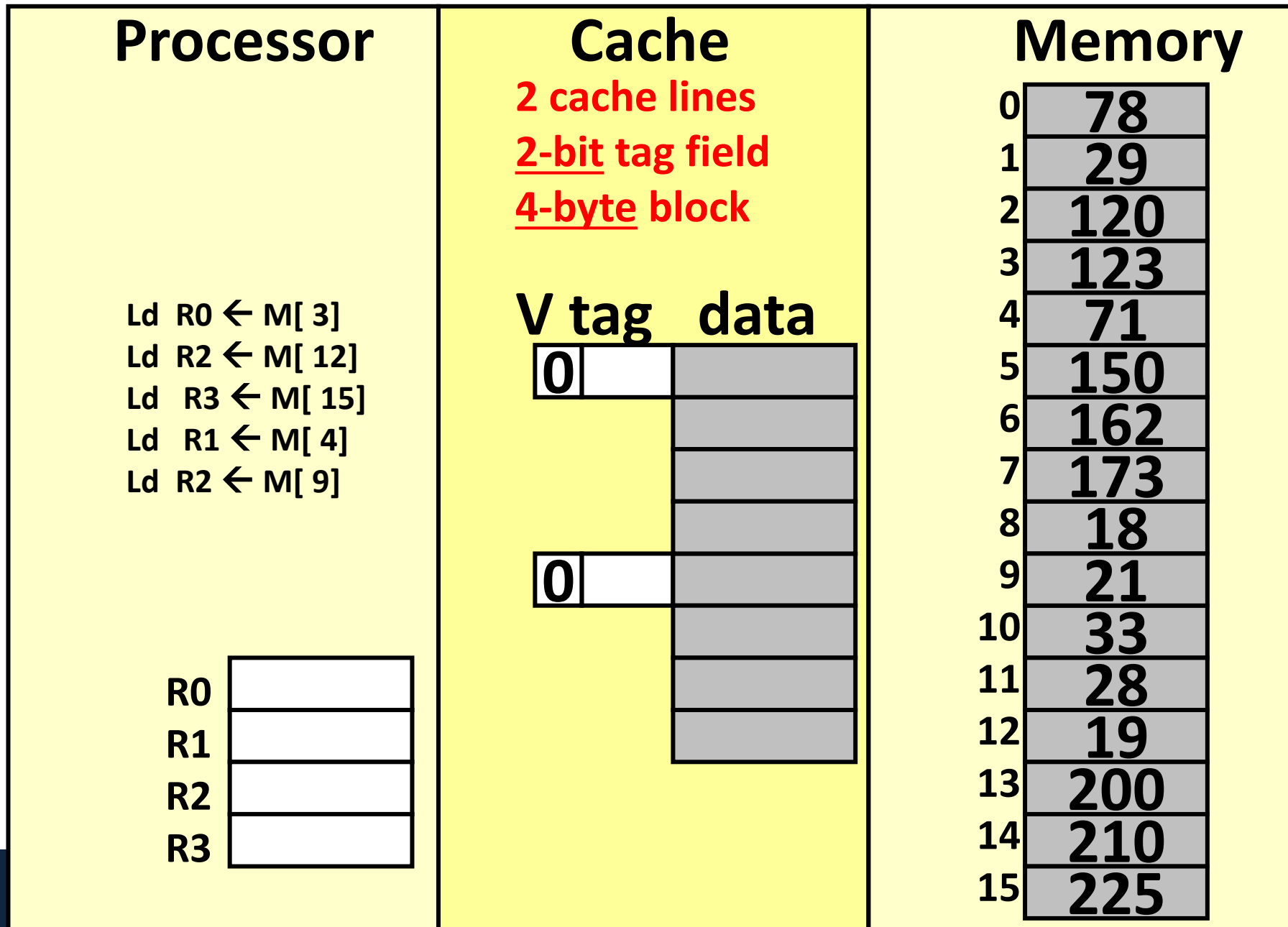
Block size for caches



Block size for caches



Extra Practice Problem



Solution to Practice Problem

Ld R0 \leftarrow M[3]

Ld R2 \leftarrow M[12]

Ld R3 \leftarrow M[15]

Ld R1 \leftarrow M[4]

Ld R2 \leftarrow M[9]

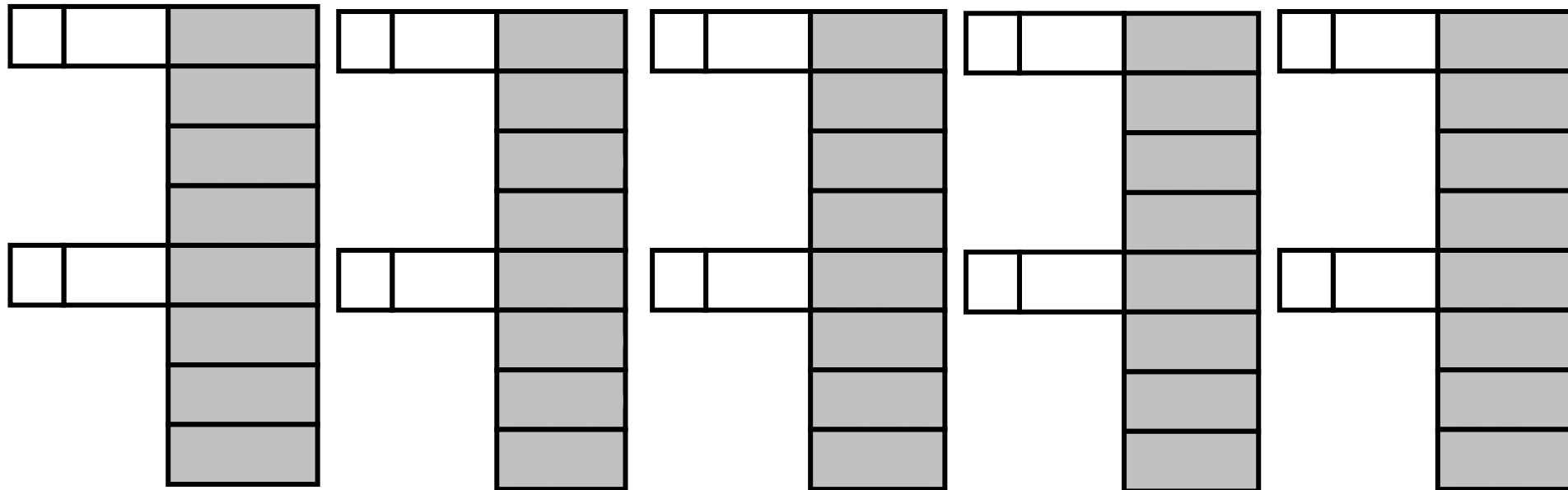
V tag data

V tag data

V tag data

V tag data

V tag data



Solution to Practice Problem

Ld R0 ← M[3]

Ld R2 ← M[12]

Ld R3 ← M[15]

Ld R1 ← M[4]

Ld R2 ← M[9]

V	tag	data	V	tag	data	V	tag	data	V	tag	data	V	tag	data
1	0	78	1	0	78	1	0	78	1	1	71	1	1	71
		29	lru		29	lru		29			150	lru		150
		120			120			120			162			162
		123			123			123			173			173
0			1	3	19	1	3	19	1	3	19	1	2	18
lru					200			200	lru		200			21
					210			210			210			33
					225			225			225			28
miss			miss			hit			miss			miss		

Poll: Answer the below questions

Extra Class Problem

*We'll see later that this is called a "fully-associative cache"

- Given a cache that works as we've described* with the following configuration: total size is 8 bytes, block size is 2 bytes, LRU replacement. The memory address size is 16 bits and is byte addressable.
 - How many bits are for each tag? How many blocks in the cache?
 $\text{Tag} = 16 - 1 \text{ (block offset)} = 15 \text{ bits};$
 $2 \text{ byte blocks, } 8 \text{ bytes total} = 4 \text{ blocks}.$
 - For the following reference stream, indicate whether each reference is a hit or miss: 0, 1, 3, 5, 12, 1, 2, 9, 4
 $M, H, M, M, M, H, H, M, M$
 - What is the hit rate?
 $3/9 = 33 \%$
 - How many bits are needed for storage overhead for each block?
 $\text{Overhead} = 15 \text{ (Tag)} + 1 \text{ (V)} + 2 \text{ (LRU)} = 18 \text{ bits}$

"Way" explanation

- A "way" in a cache is a particular location that a piece of memory **could** exist in
 - E.g. our last cache example was a "two-way" cache
 - Is it the same as the number of entries in the cache?
 - In this example, yes
 - But generally no. Hold on until we talk about set-associate caches

V		
V		
	tag	data

LRU Implementation with Counters

◆ 2 ways

- ❖ 1 bit per set to mark latest way accessed in set
- ❖ Evict way not pointed by bit

◆ k-way set associative LRU

- ❖ Requires full ordering of way accesses
- ❖ Hold a $\log_2 k$ bit counter per line
- ❖ When a way i is accessed
 - $X = \text{Counter}[i]$
 - $\text{Counter}[i] = k-1$
 - for ($j = 0$ to $k-1$)
 - if ($(j \neq i)$ AND ($\text{Counter}[j] > X$)) $\text{Counter}[j]--$;
- ❖ When replacement is needed
 - evict way with counter = 0
- ❖ Expensive for even small k 's

Initial State				
Way	0	1	2	3
Count	0	1	2	3

Access way 2				
Way	0	1	2	3
Count	0	1	3	2

Access way 0				
Way	0	1	2	3
Count	3	1	2	1

Spatial Locality

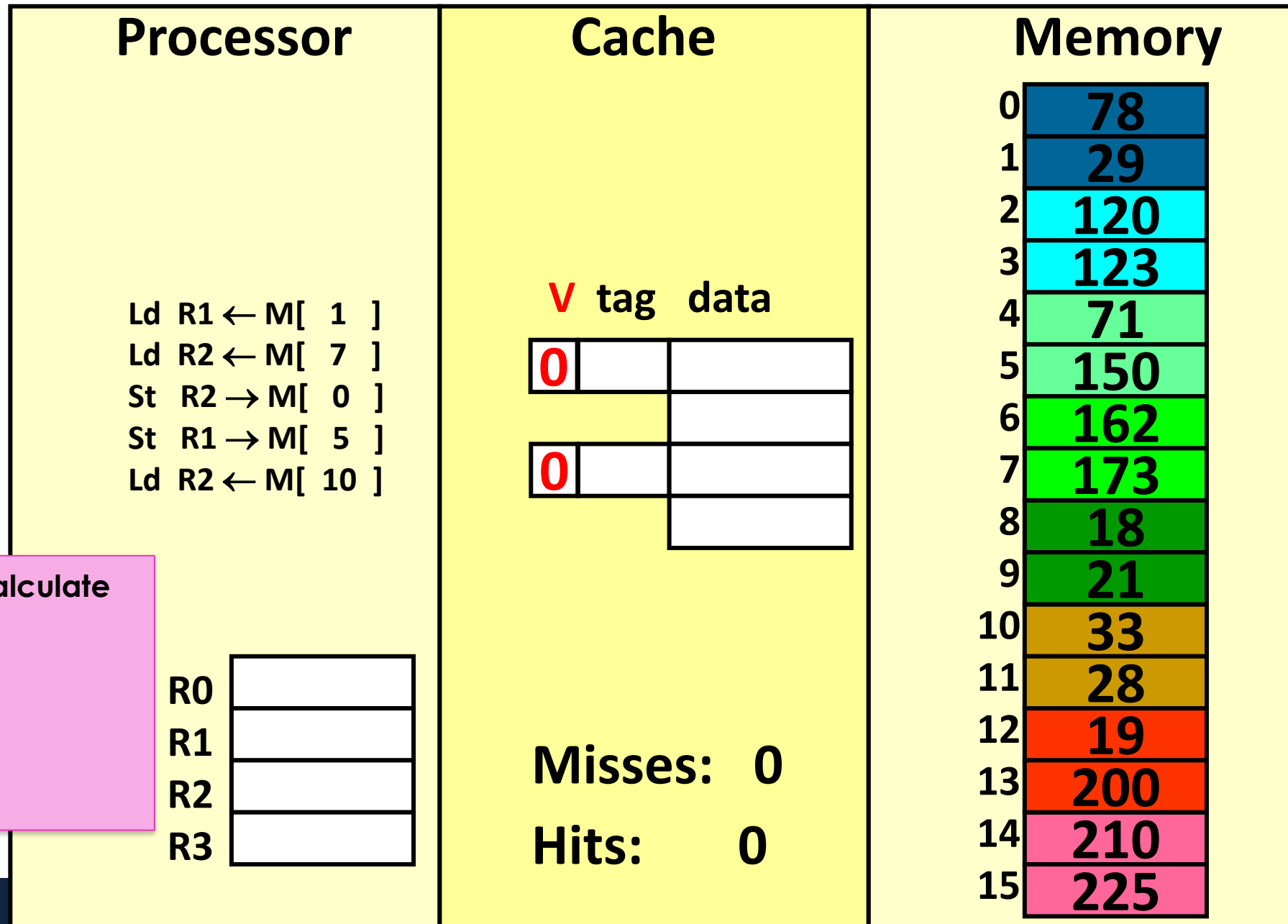
- Notice that when we accessed address 1, we also brought in address 0
 - This turned out to be a good thing, since we later referenced address 0 and found it in the cache
- This is taking advantage of **spatial locality**:
 - If we access a memory location (e.g. 1000), we are more likely to access a location near it (e.g. 1001) than some random location
 - Arrays and structs are a big reason for this

```
for(i=0; i< N; i++)  
    for(j = 0; j < N; j++ )  
    {  
        count++;  
        arrayInt[i][j] = 10;  
    }
```

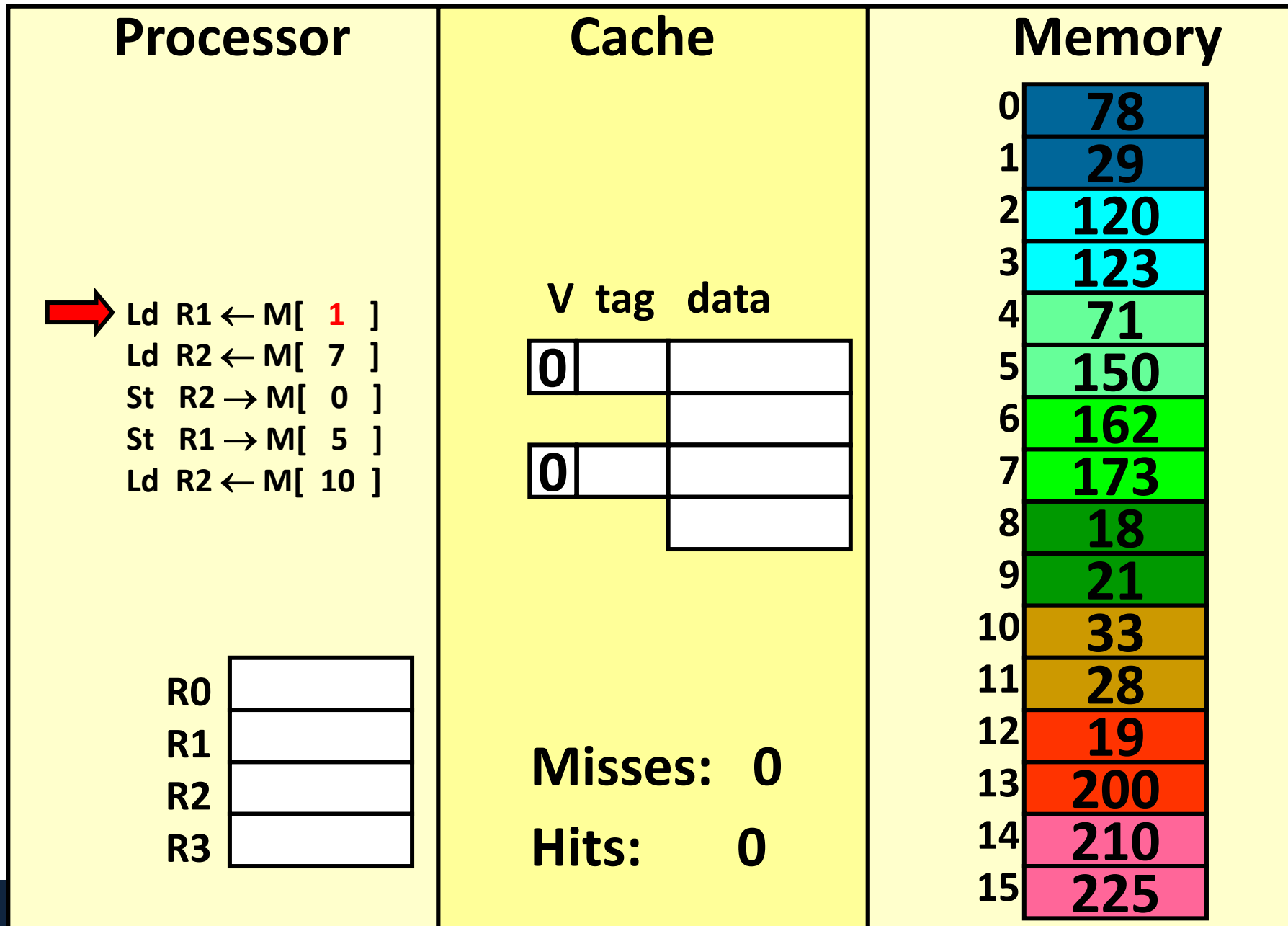
What about stores?

- 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?
(no allocate-on-write policy)

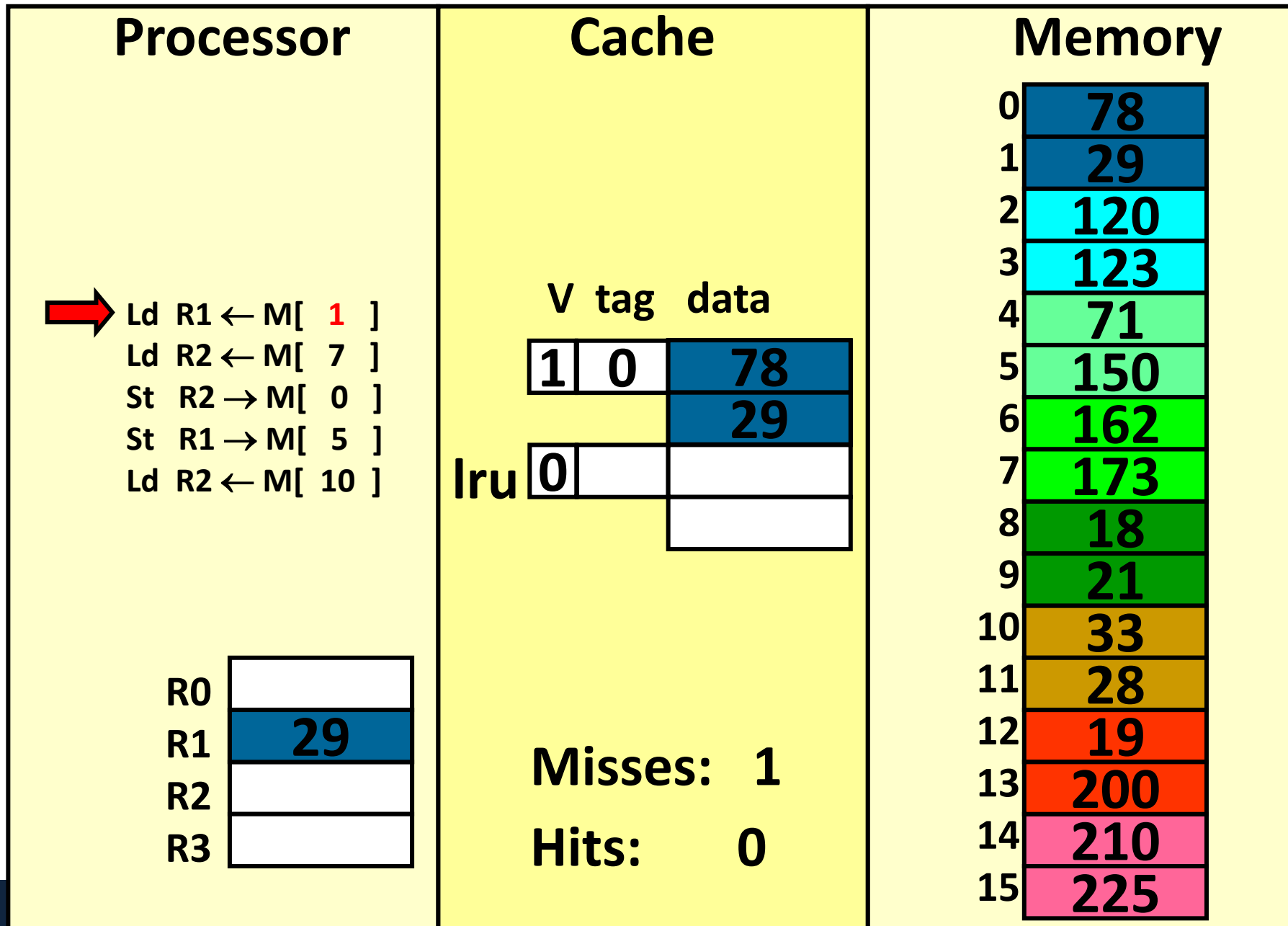
Handling stores (write-through, allocate on write)



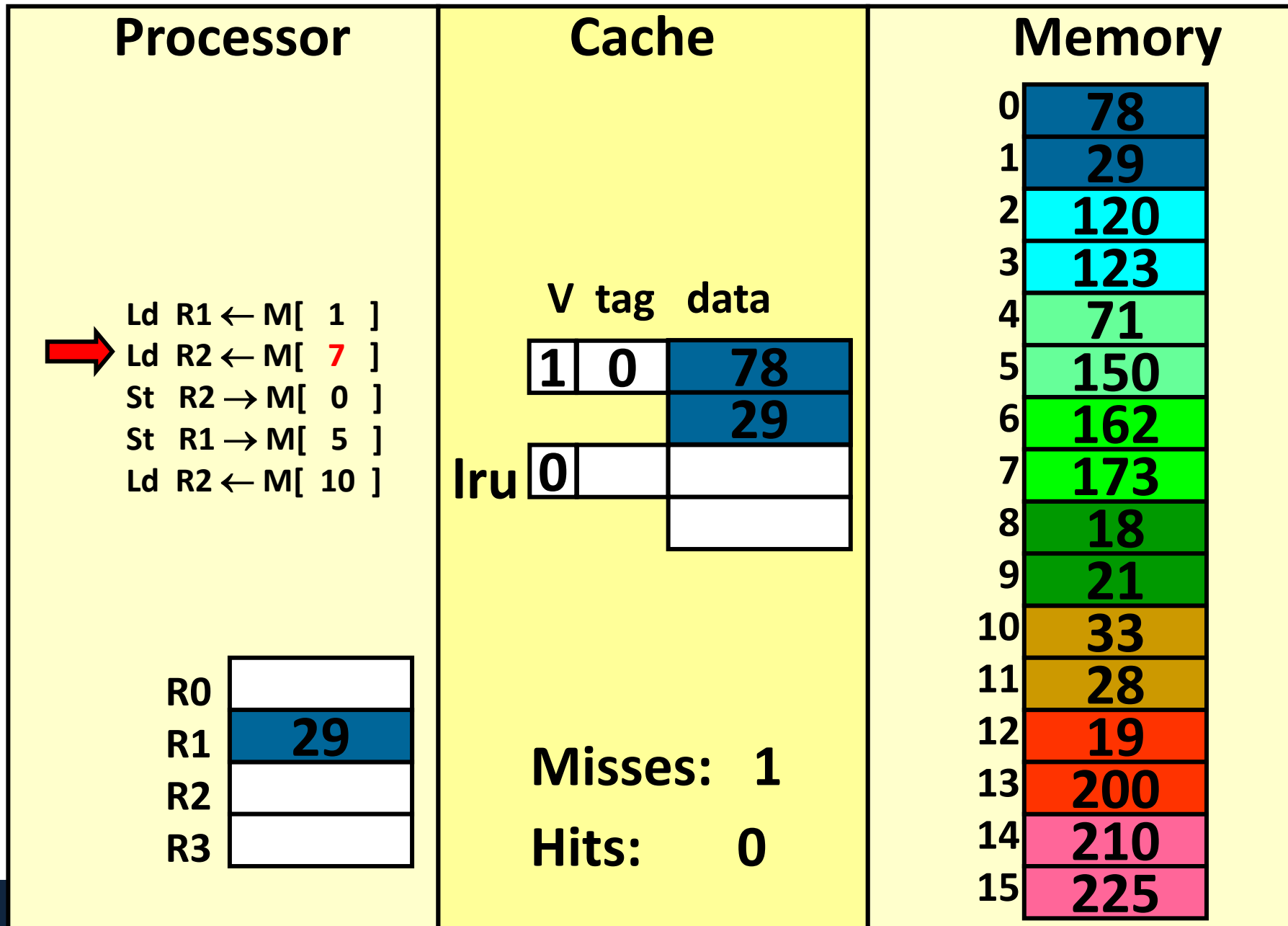
write-through, allocate on write (REF 1)



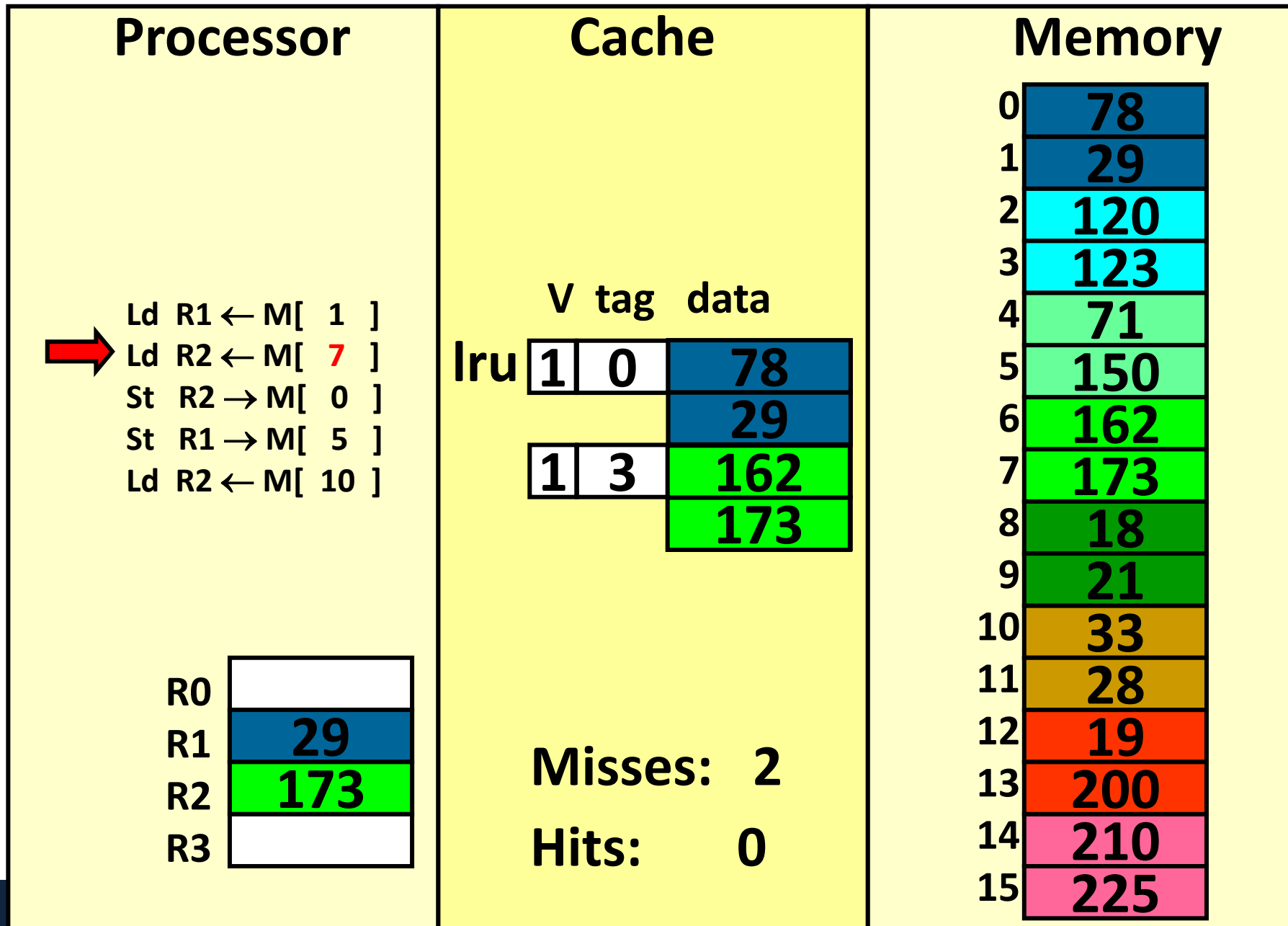
write-through, allocate on write (REF 1)



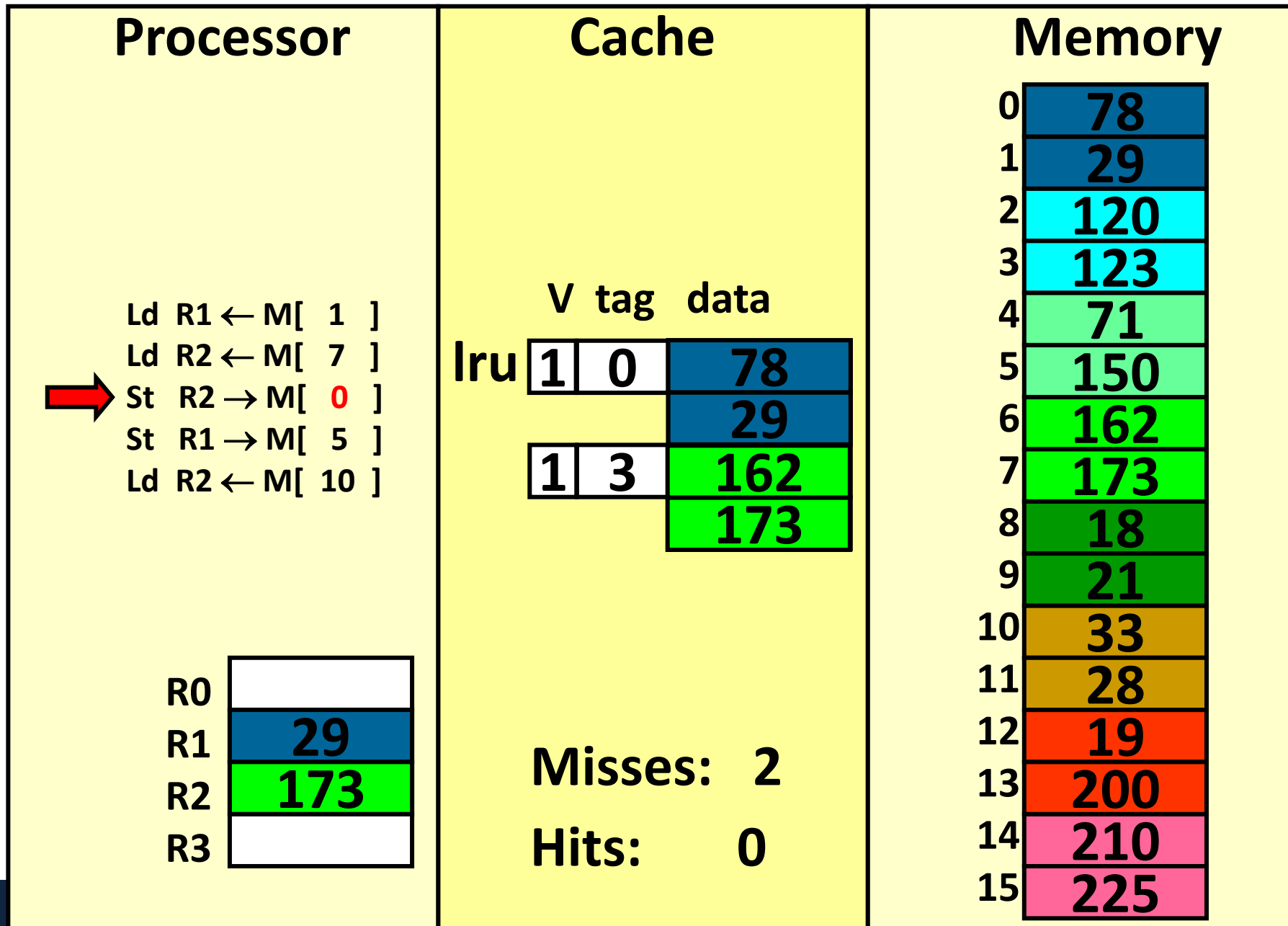
write-through, allocate on write (REF 2)



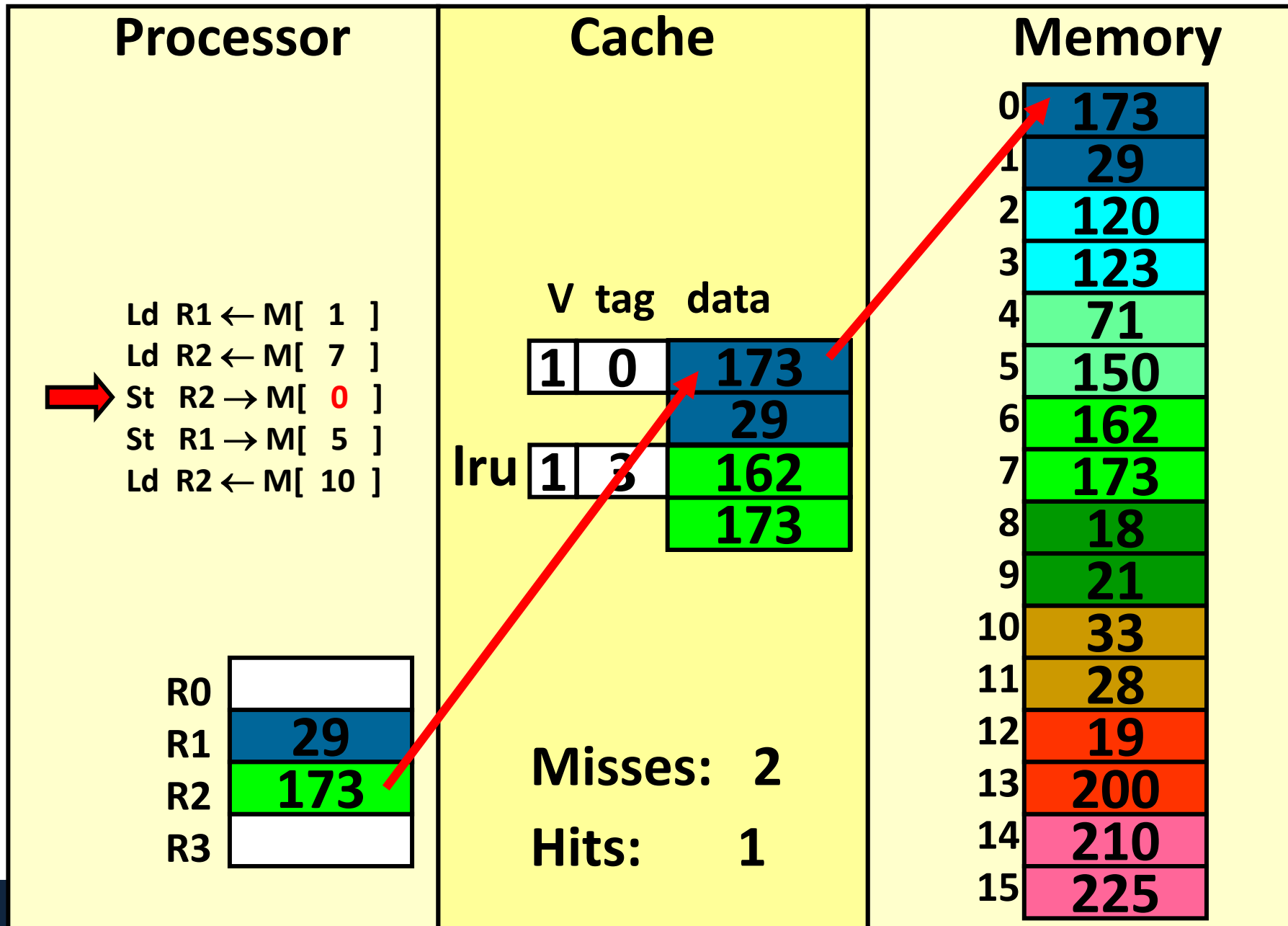
write-through, allocate on write (REF 2)



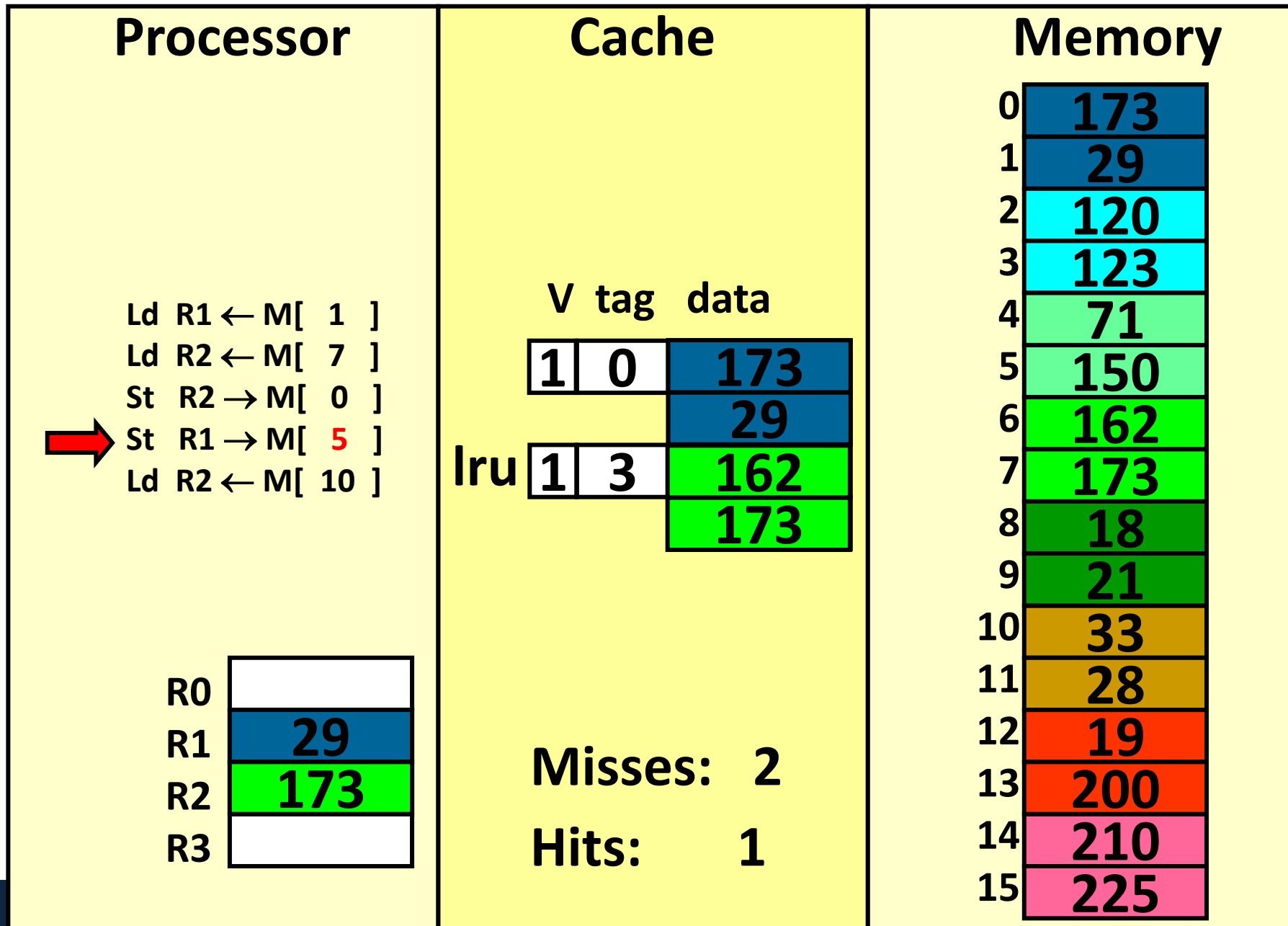
write-through, allocate on write (REF 3)



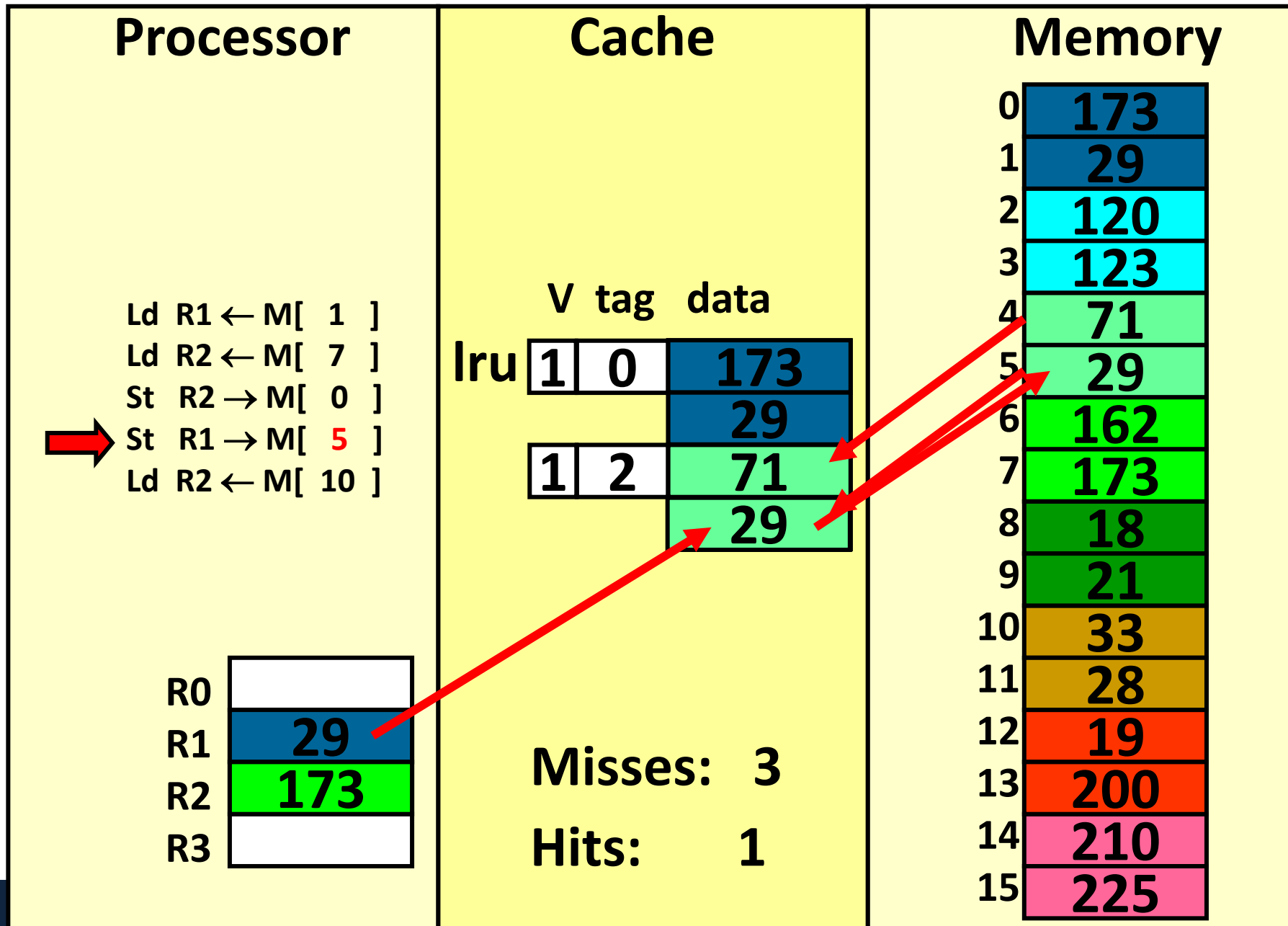
write-through, allocate on write (REF 3)



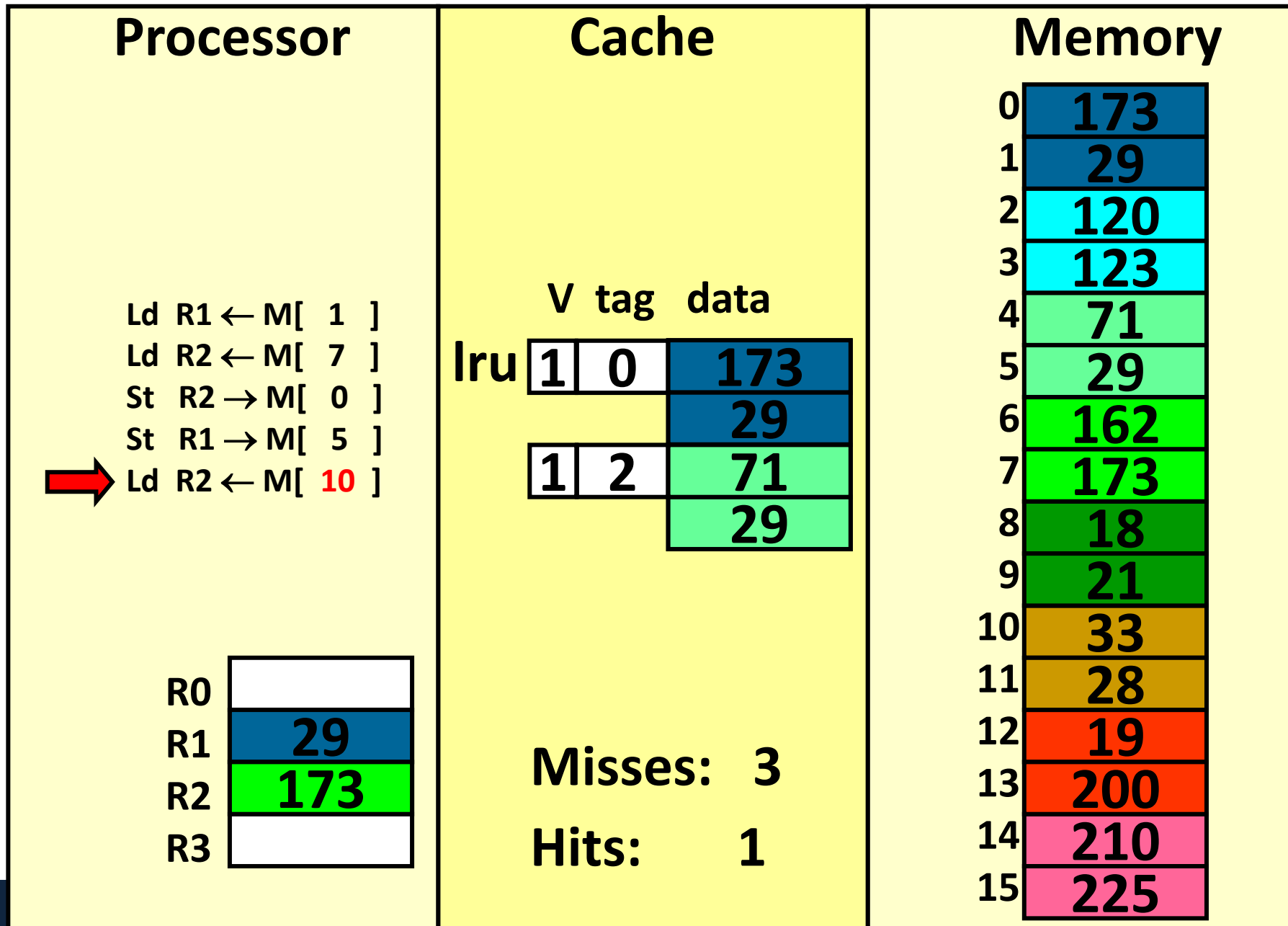
write-through, allocate on write (REF 4)



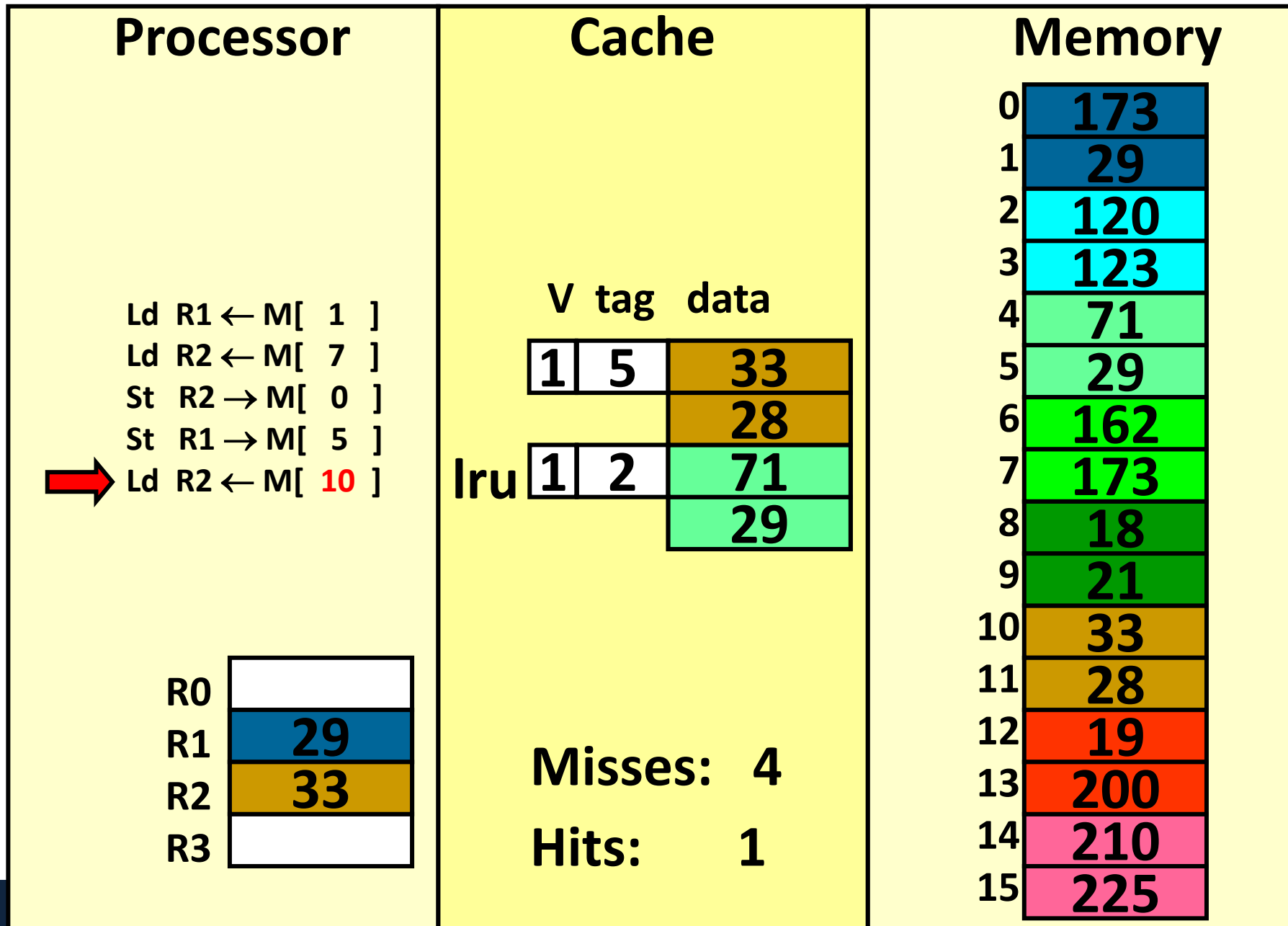
write-through, allocate on write (REF 4)



write-through, allocate on write (REF 6)



write-through, allocate on write (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 $< 20\%$
 - Can we take advantage of that?
 - Multi-core processors have limited bandwidth between caches and memory
 - Extra stores also cost power

Next time

- Write-back caches
- Direct-mapped vs associative caches.
- Lingering questions / feedback? I'll include an anonymous form at the end of every lecture: <https://bit.ly/3oXr4Ah>

