

19. Caches: Direct Mapped

Assignment Project Exam Help

EECS 370 – Introduction to Computer Organization – Fall 2020

<https://powcoder.com>

Satish Narayanasamy
Add WeChat powcoder

EECS Department
University of Michigan in Ann Arbor, USA

© Narayanasamy 2020

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

Announcements

Upcoming deadlines:

HW4

due Nov 10th

Project 3

due Nov. 12th

Assignment Project Exam Help

<https://powcoder.com>

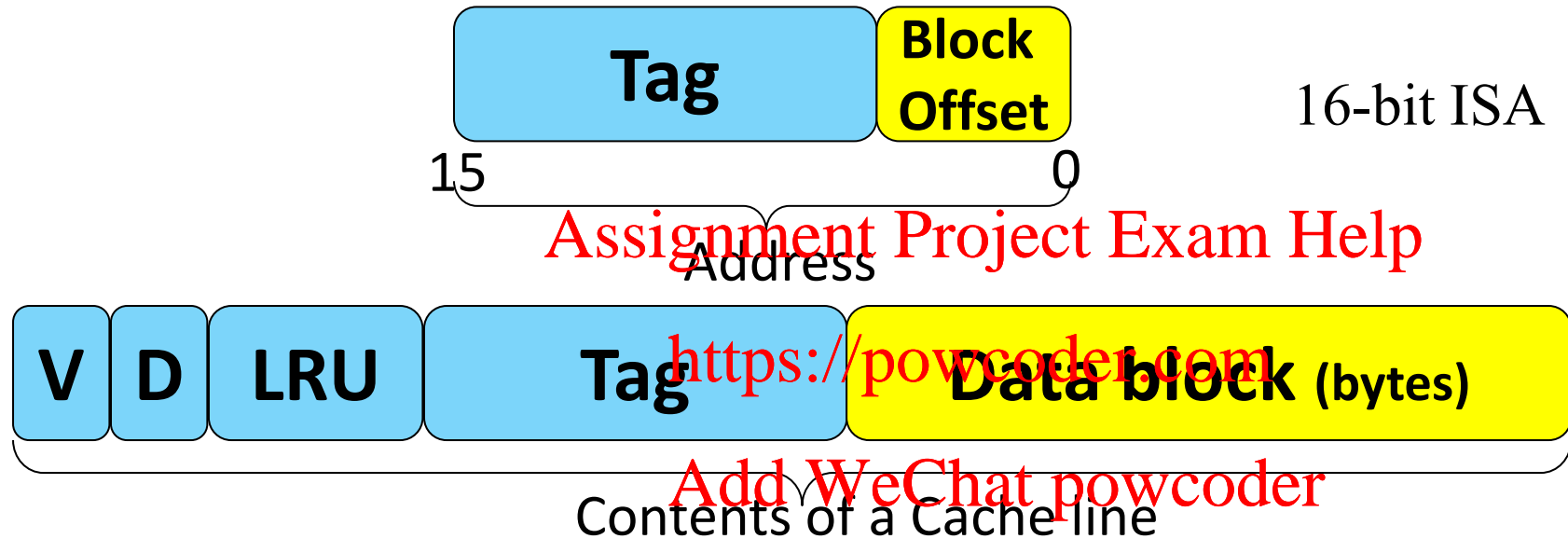
Add WeChat powcoder

Assignment Project Exam Help

Recap: Cache Blocks and Write policy

Add WeChat powcoder

Review: Cache Organization



Cache blocks:

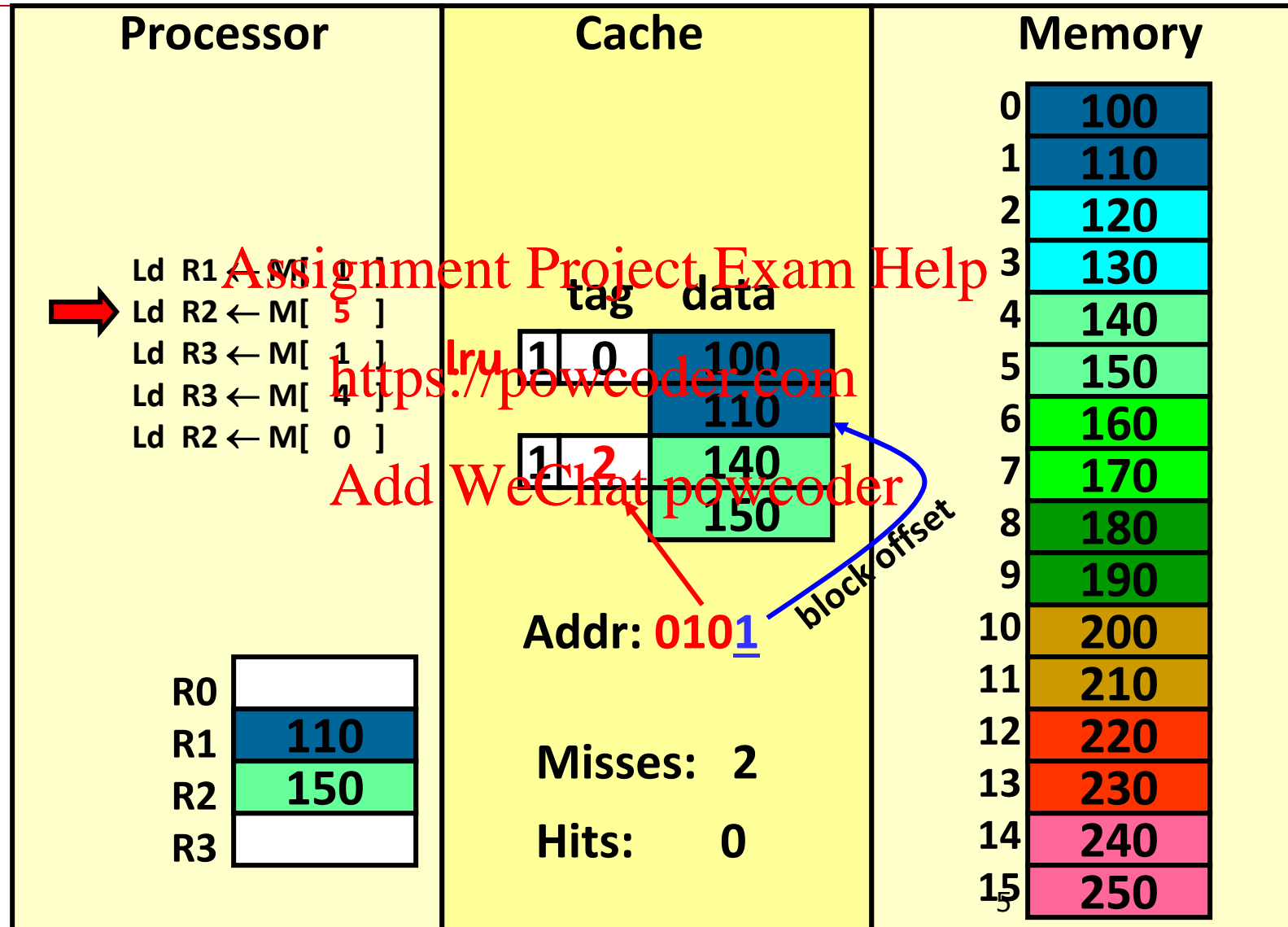
- Captures spatial locality (increase cache hit rate)

- Reduces tag overhead (number and size of tags)

Need not store block offset in the cache line

Determine byte to be read/written from the address directly

Review: How to find tag from address?



Review: Writes

Write-allocate vs. no-write-allocate caches

Policy that decides what to do with a cache-miss on a store instruction.

Assignment Project Exam Help

Write-allocate: First bring data from memory into the cache, then write

<https://powcoder.com>

Add WeChat powcoder

No-write-allocate: do not bring data in the cache, just write directly to the memory, not to the cache

Review: Writes

Write-through vs. write-back caches

Policy that decides when to **write to cache vs. memory vs. both**

Assignment Project Exam Help

Write-through: write to both cache and memory

<https://powcoder.com>

Write-back: write only to cache, keep track of dirty cache line, write to memory when dirty cache line is evicted

Add WeChat powcoder

Review: Writes

Store w No Allocate	Write-Back	Write-Through
Hit?	Write Cache	Write to Cache + Memory
Miss?	Write to Memory	Write to Memory
Replace block?	If evicted block is dirty, write to Memory	Do Nothing

<https://powcoder.com>

Store w Allocate	Write-Back	Write-Through
Hit?	Write Cache	Write to Cache + Memory
Miss?	Read from Memory to Cache, Allocate to LRU block Write to Cache	Read from Memory to Cache, Allocate to LRU block Write to Cache + Memory
Replace block?	If evicted block is dirty, write to Memory	Do Nothing

Assignment Project Exam Help

Direct Mapped Caches

<https://powcoder.com>

Add WeChat powcoder

Fully-associative caches

A block can go
to **any** location

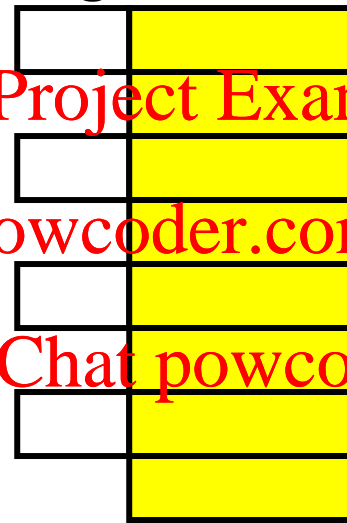
Address:



3 bits

1 bit

tag data



Memory

Tag

Block_offset

0	100	00 0	0
1	110	00 0	1
2	120	00 1	0
3	130	00 1	1
4	140	01 0	0
5	150	01 0	1
6	160	01 1	0
7	170	01 1	1
8	180	10 0	0
9	190	10 0	1
10	200	10 1	0
11	210	10 1	1
12	220	11 0	0
13	230	11 0	1
14	240	11 1	0
15	250	11 1	1

Fully-associative caches

We designed a fully-associative cache

- A memory location can be copied to any cache line.
- We **check every cache tag** to determine whether the data is in the cache.

Assignment Project Exam Help

This approach can be too slow sometimes

- Parallel tag searches are expensive and can be slow. Why?

Add WeChat powcoder

Direct mapped caches

We can redesign the cache to eliminate the requirement for parallel tag lookups

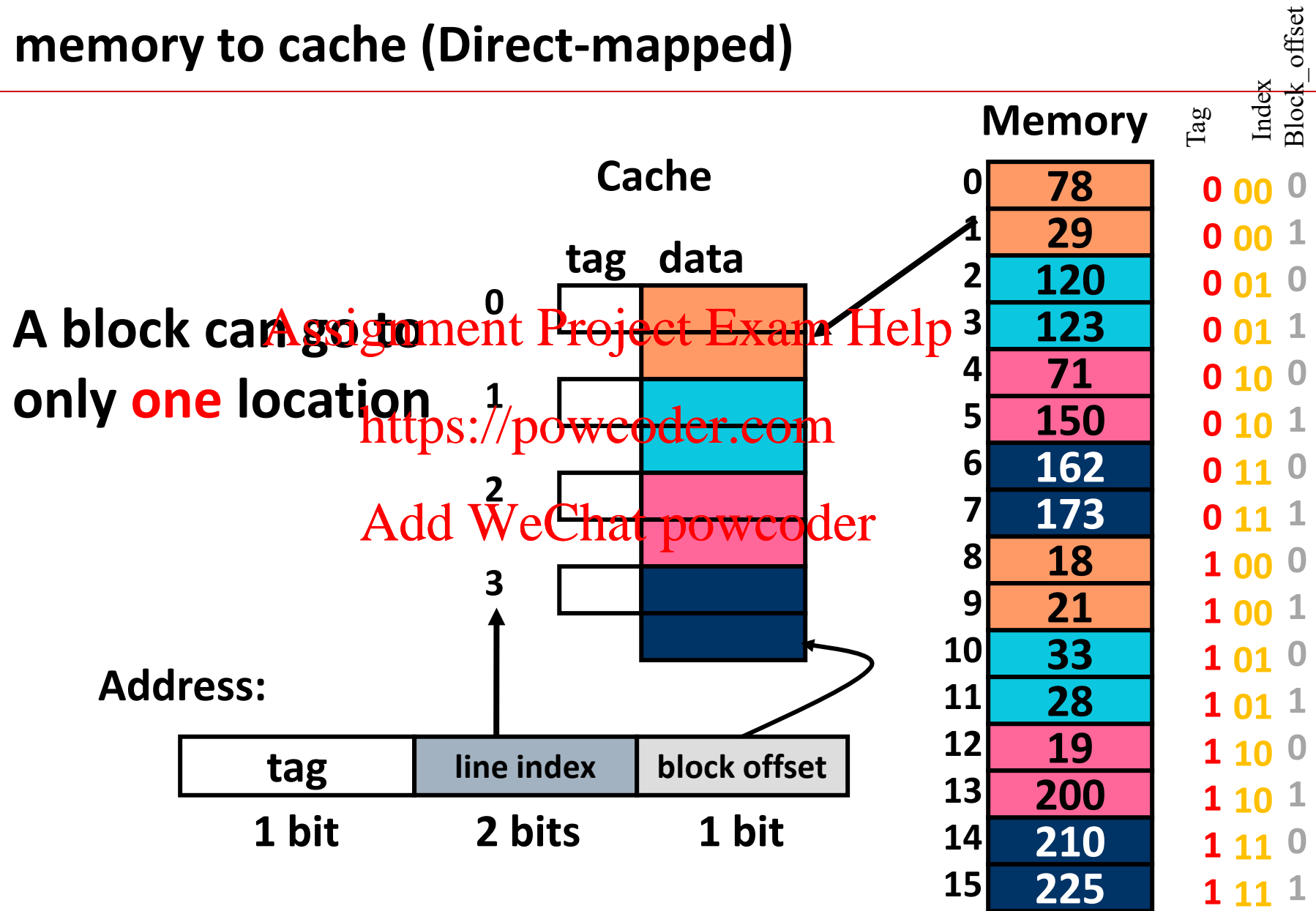
- Direct mapped caches partition memory into as many regions as there are cache lines
- Each memory region maps to a single cache line in which data can be placed
- You then only need to check a single tag – the one associated with the region the reference is located in

Assignment Project Exam Help

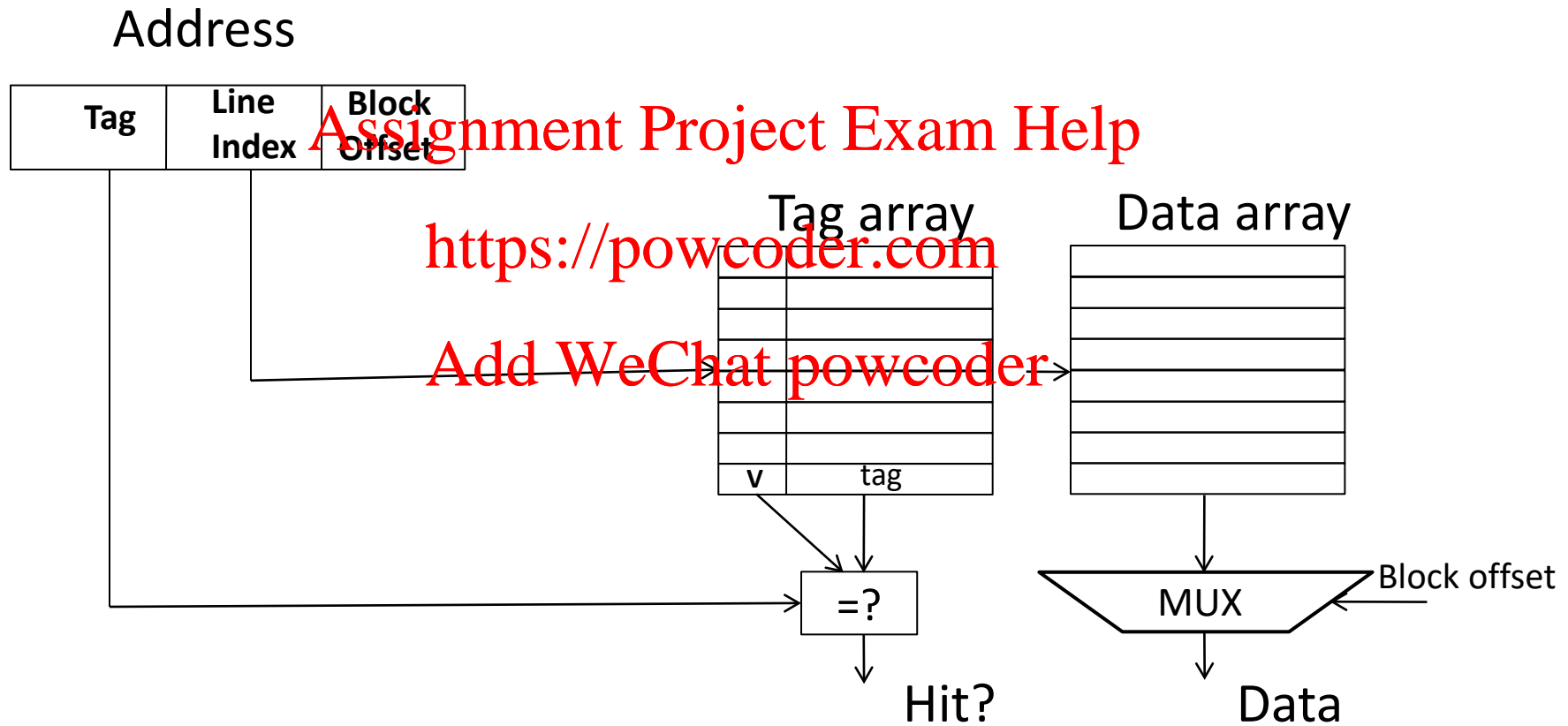
<https://powcoder.com>

Add WeChat powcoder

Mapping memory to cache (Direct-mapped)



Direct-mapped cache: Placement & Access



Direct mapped caches

Two blocks in memory that map to the same cache index cannot be present in the cache at the same time (conflict)

One index → one entry

Assignment Project Exam Help

Can lead to 0% hit rate if more than one block accessed in an interleaved manner map to the same index

<https://powcoder.com>

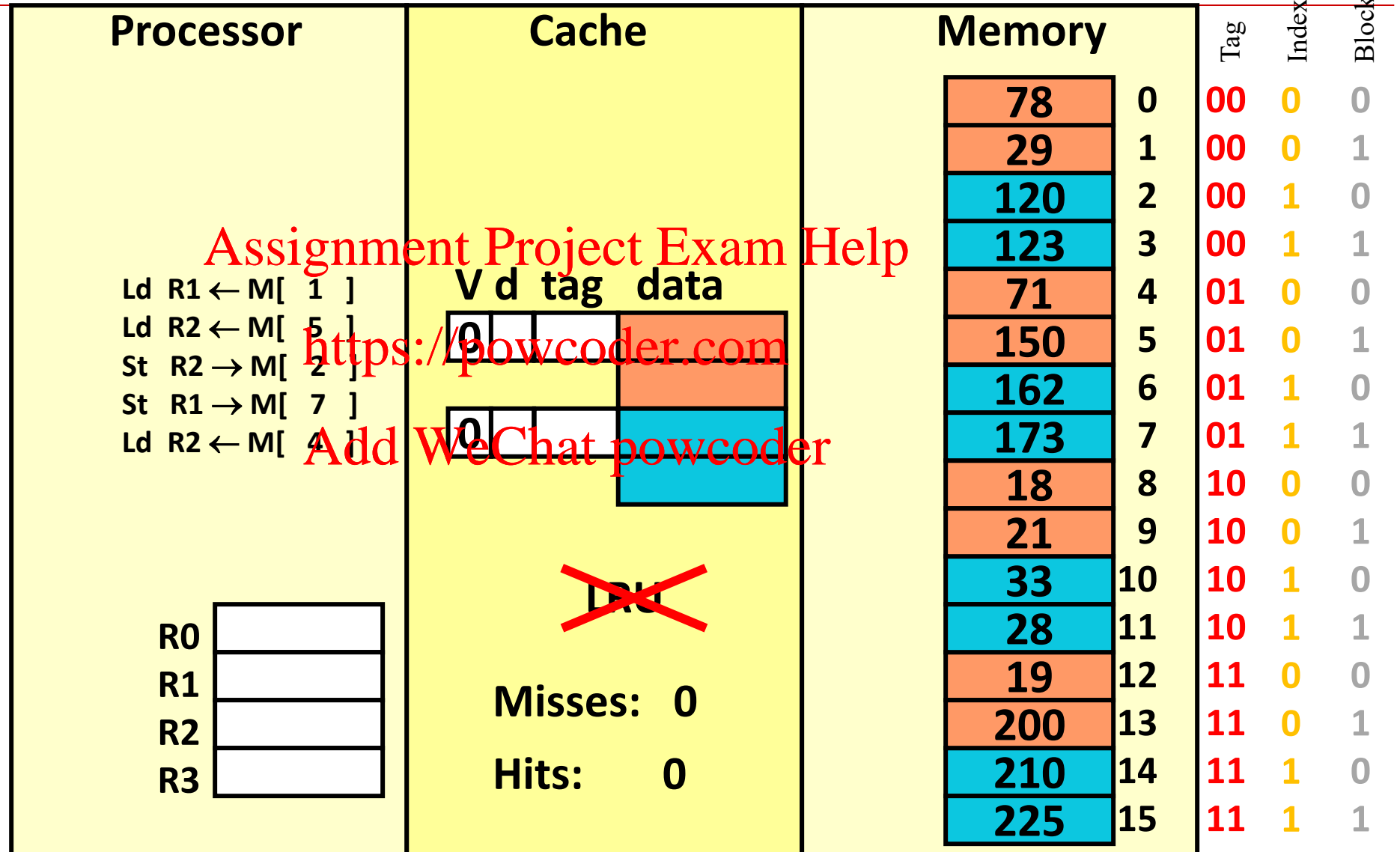
Assume addresses A and B have the same index bits but different tag bits

A, B, A, B, A, B, A, B, ...

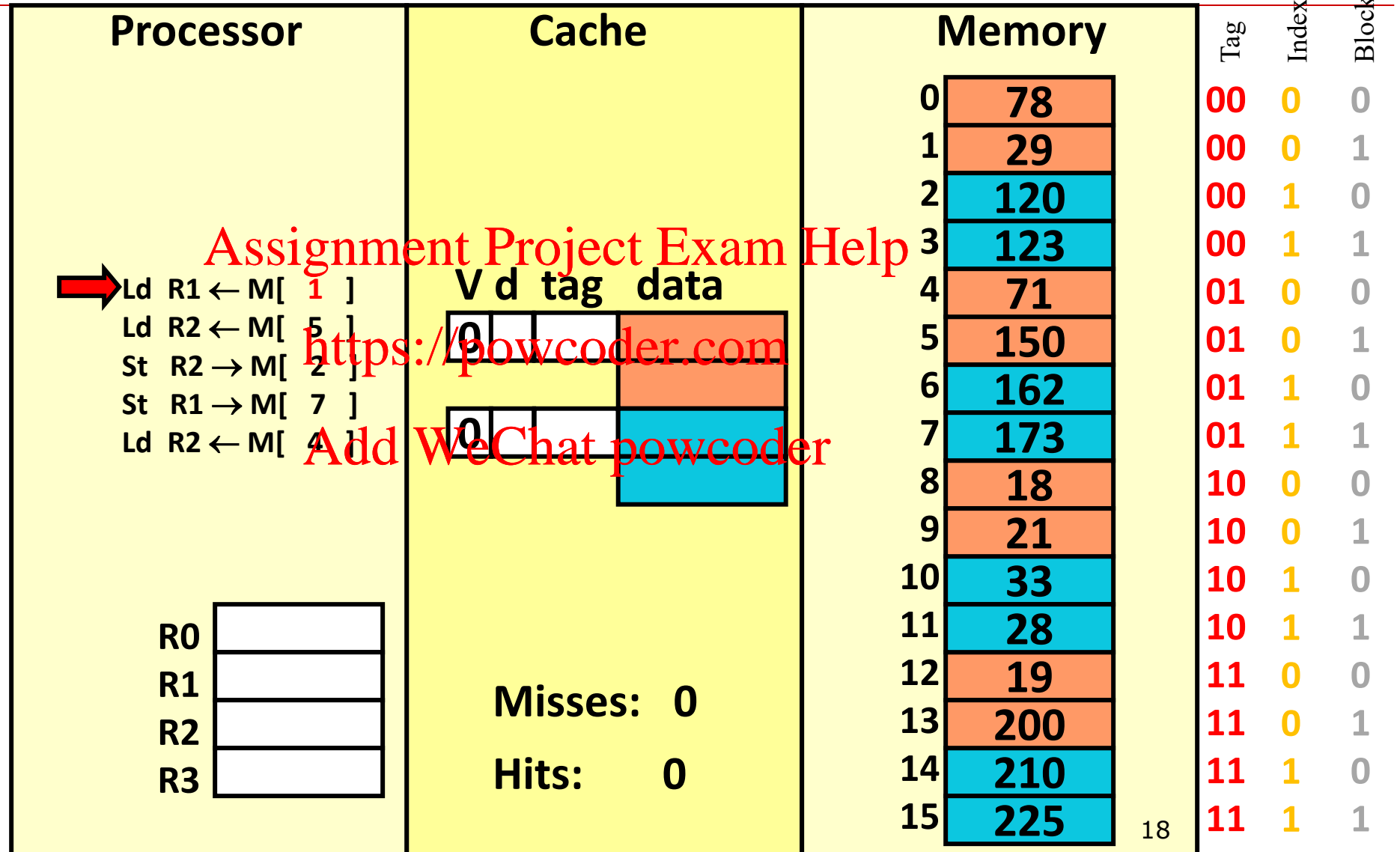
Add WeChat powcoder

All accesses are conflict misses

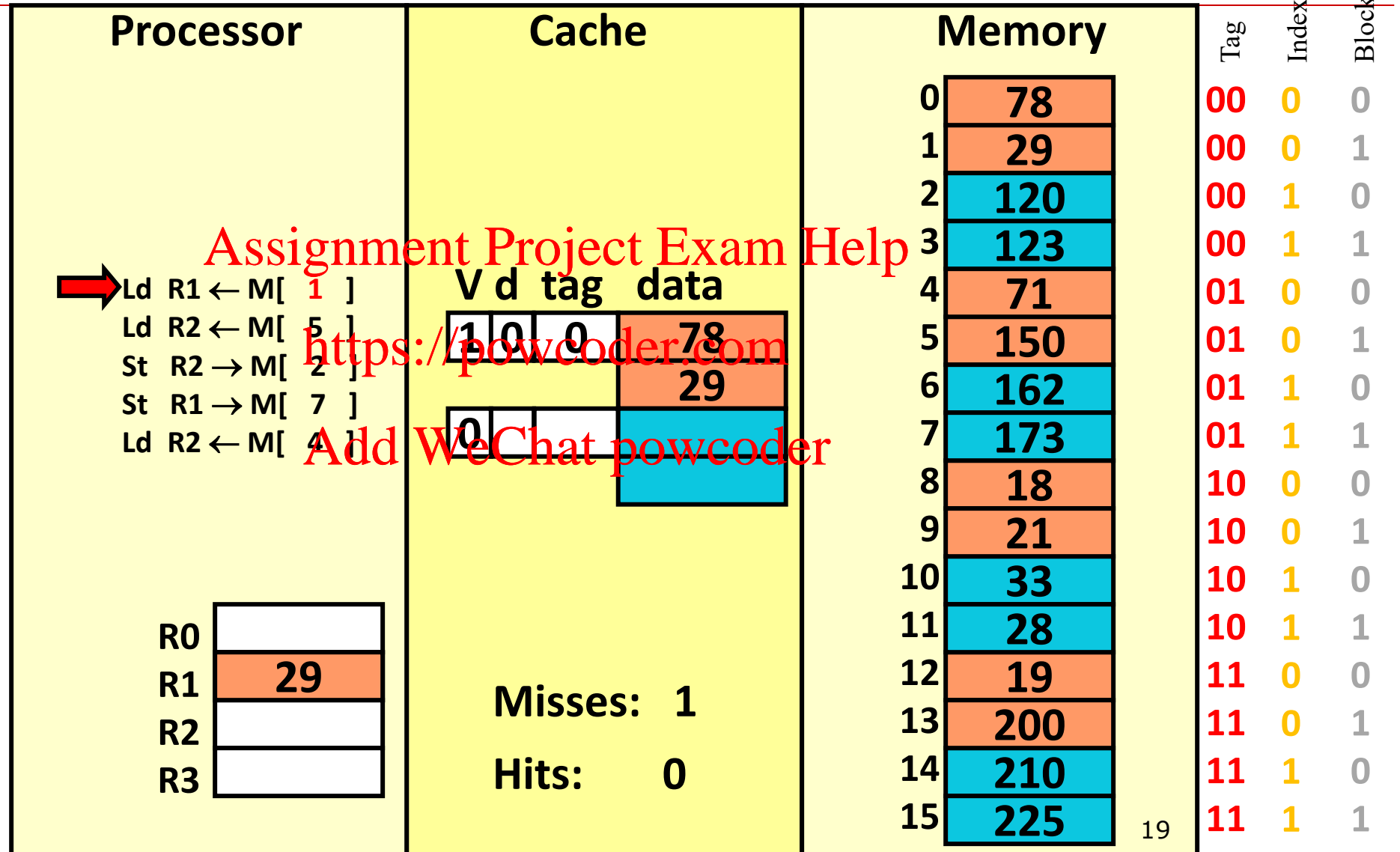
Direct-mapped cache



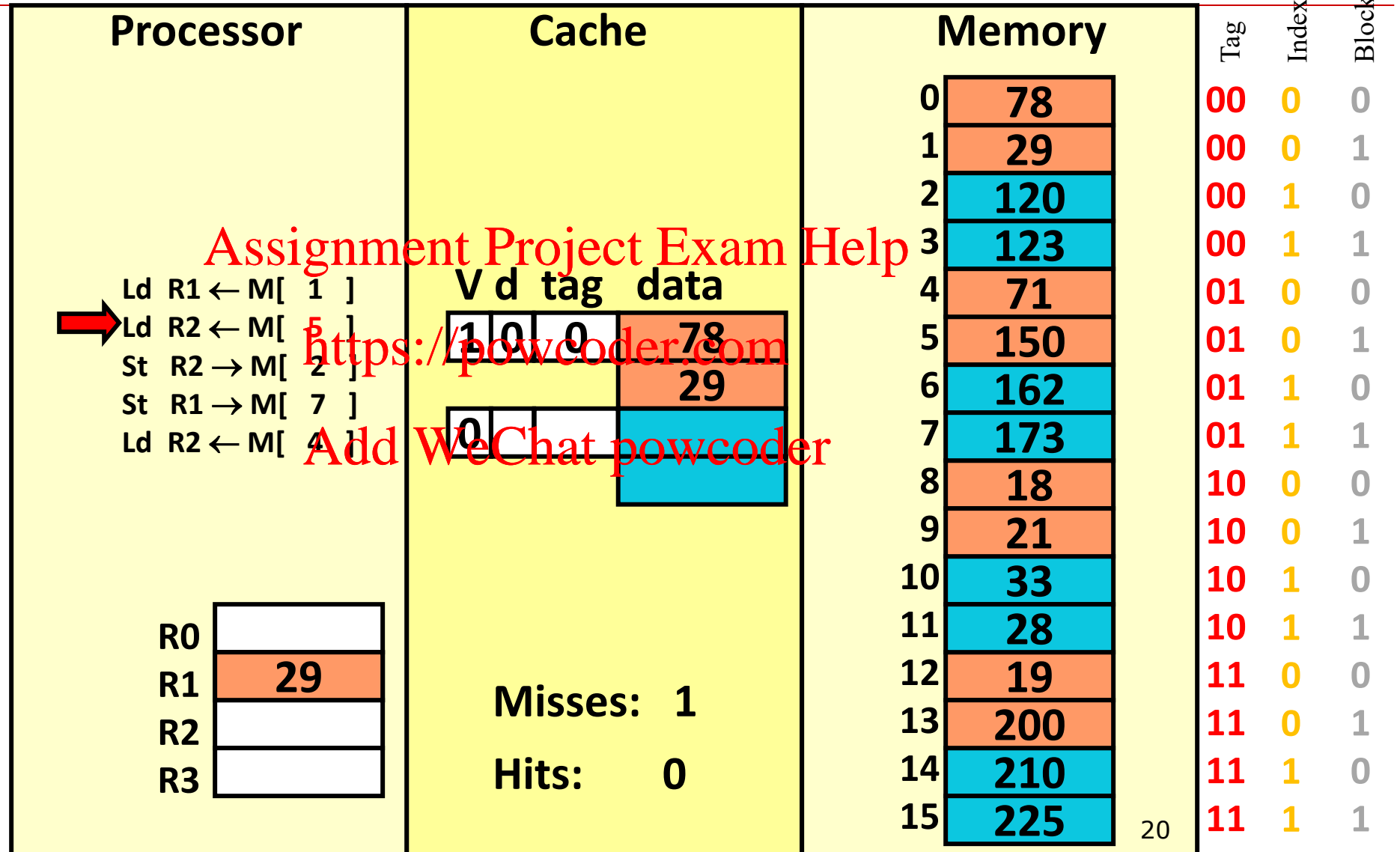
Direct-mapped (REF 1)



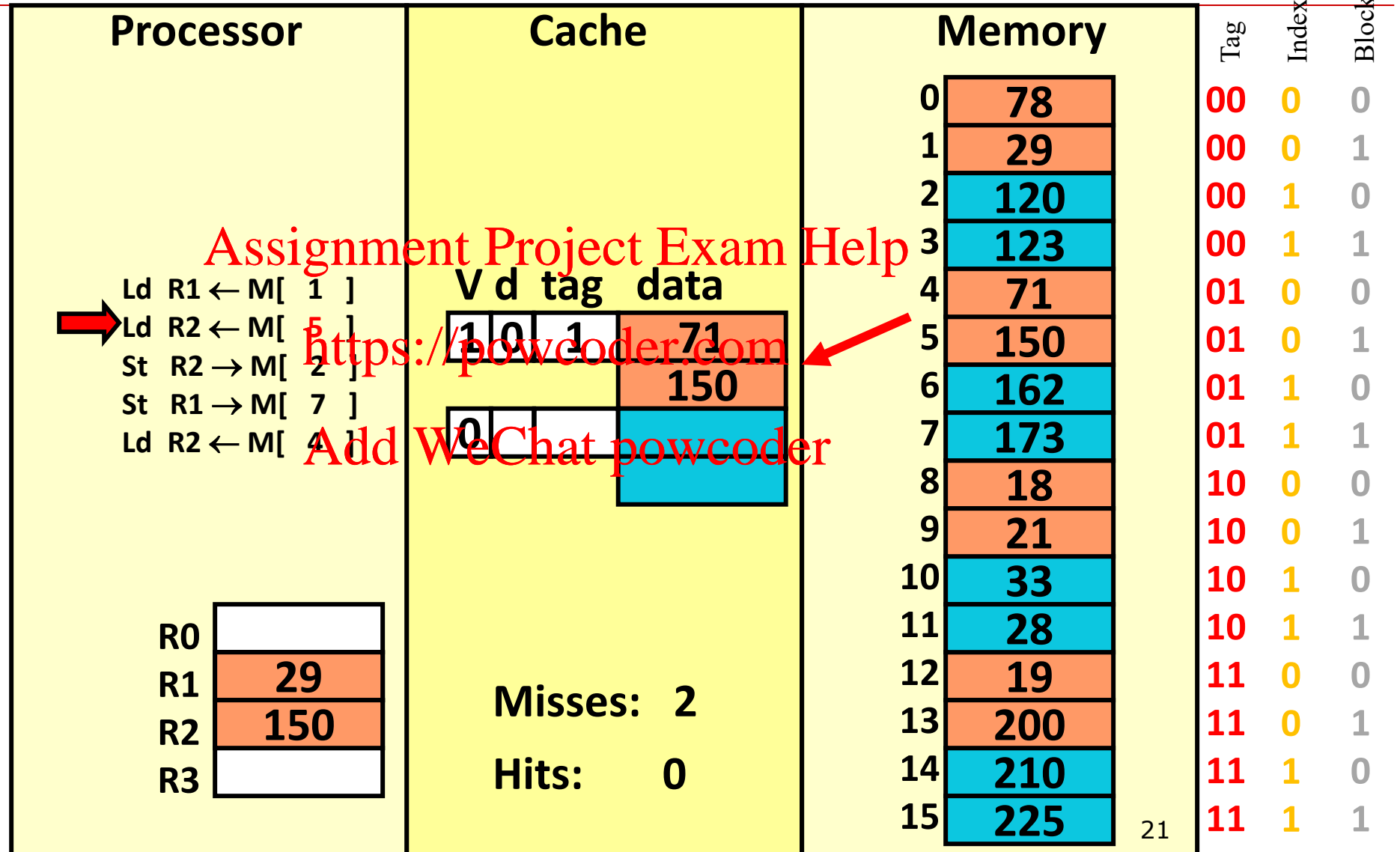
Direct-mapped (REF 1)



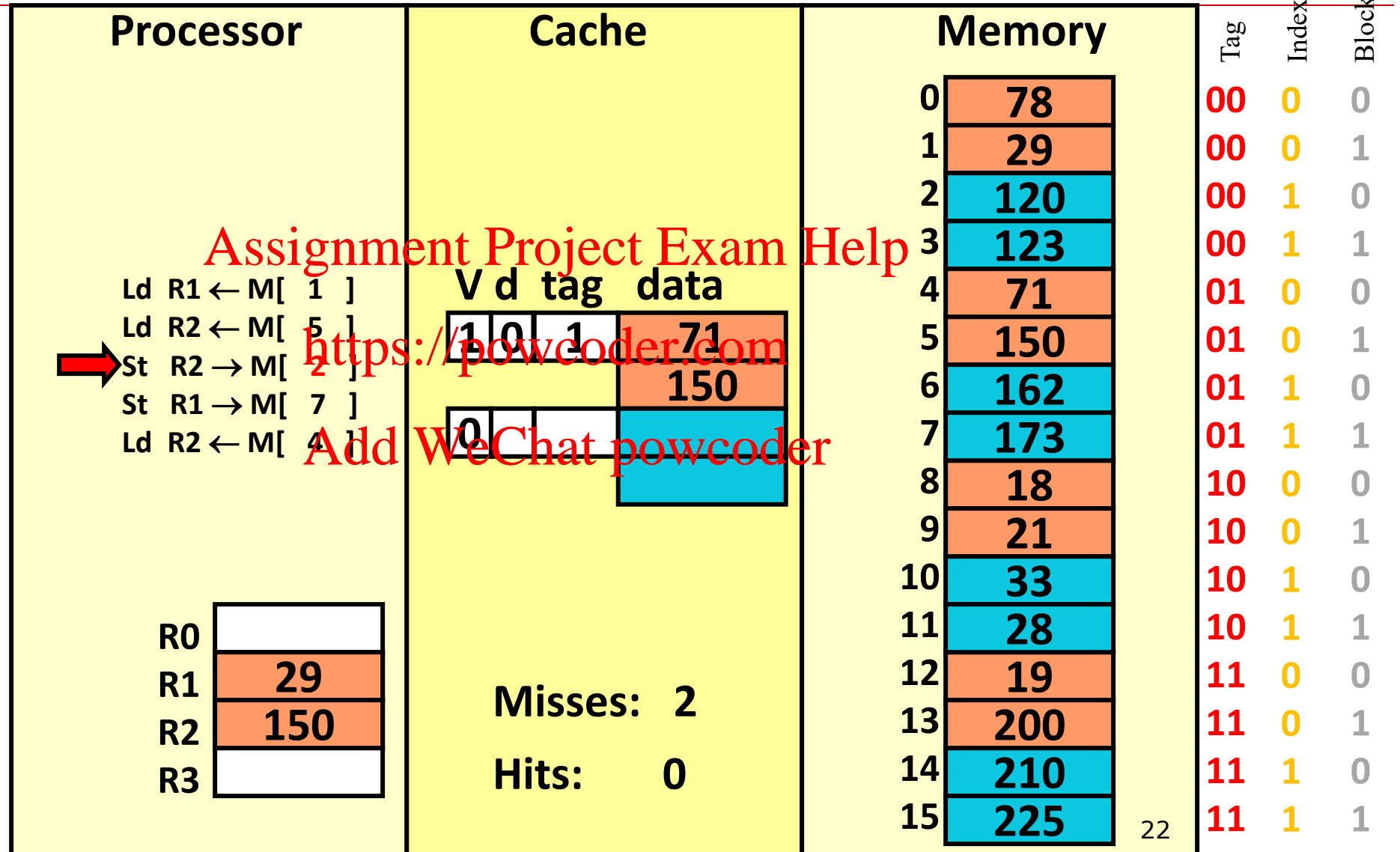
Direct-mapped (REF 2)



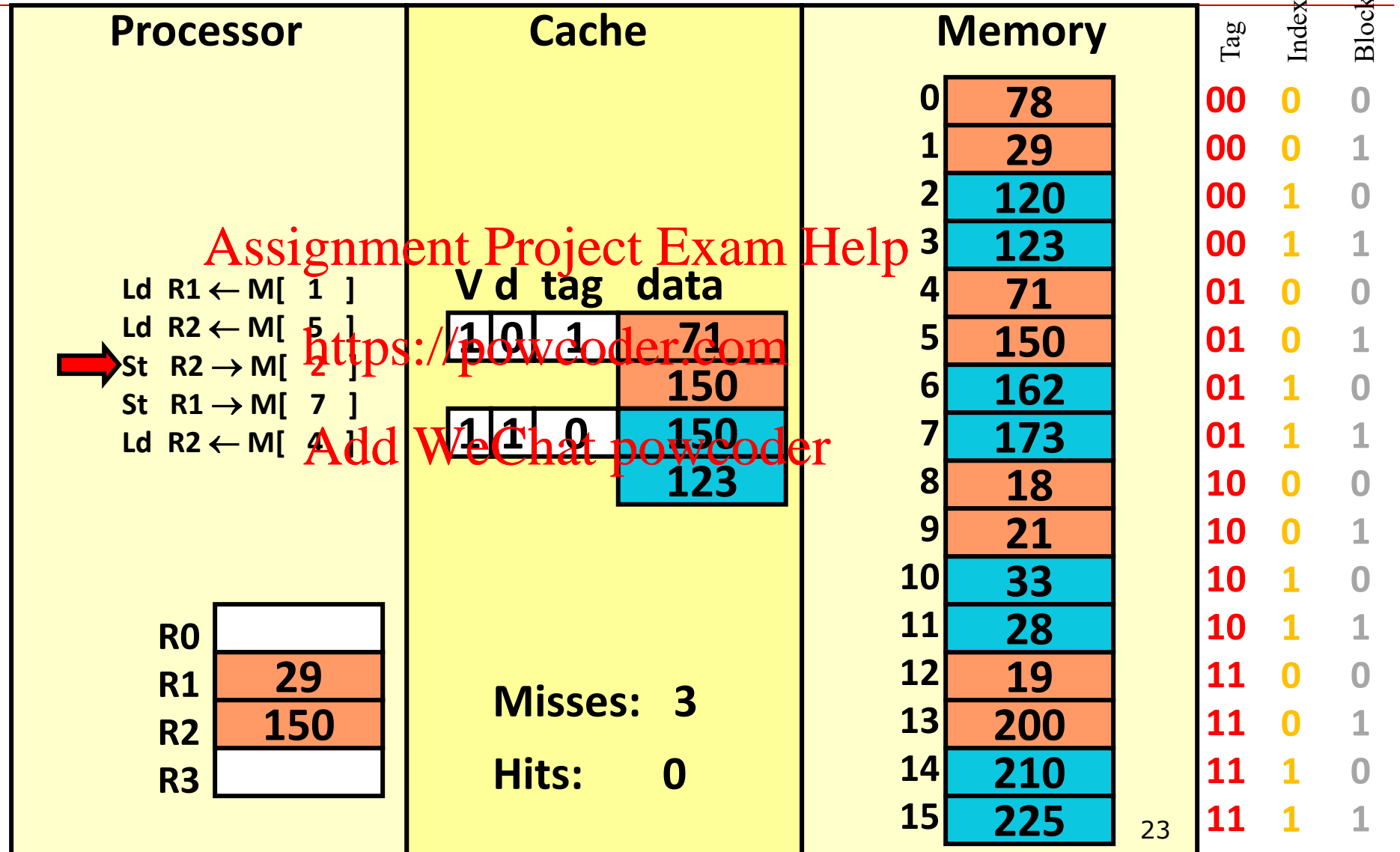
Direct-mapped (REF 2)



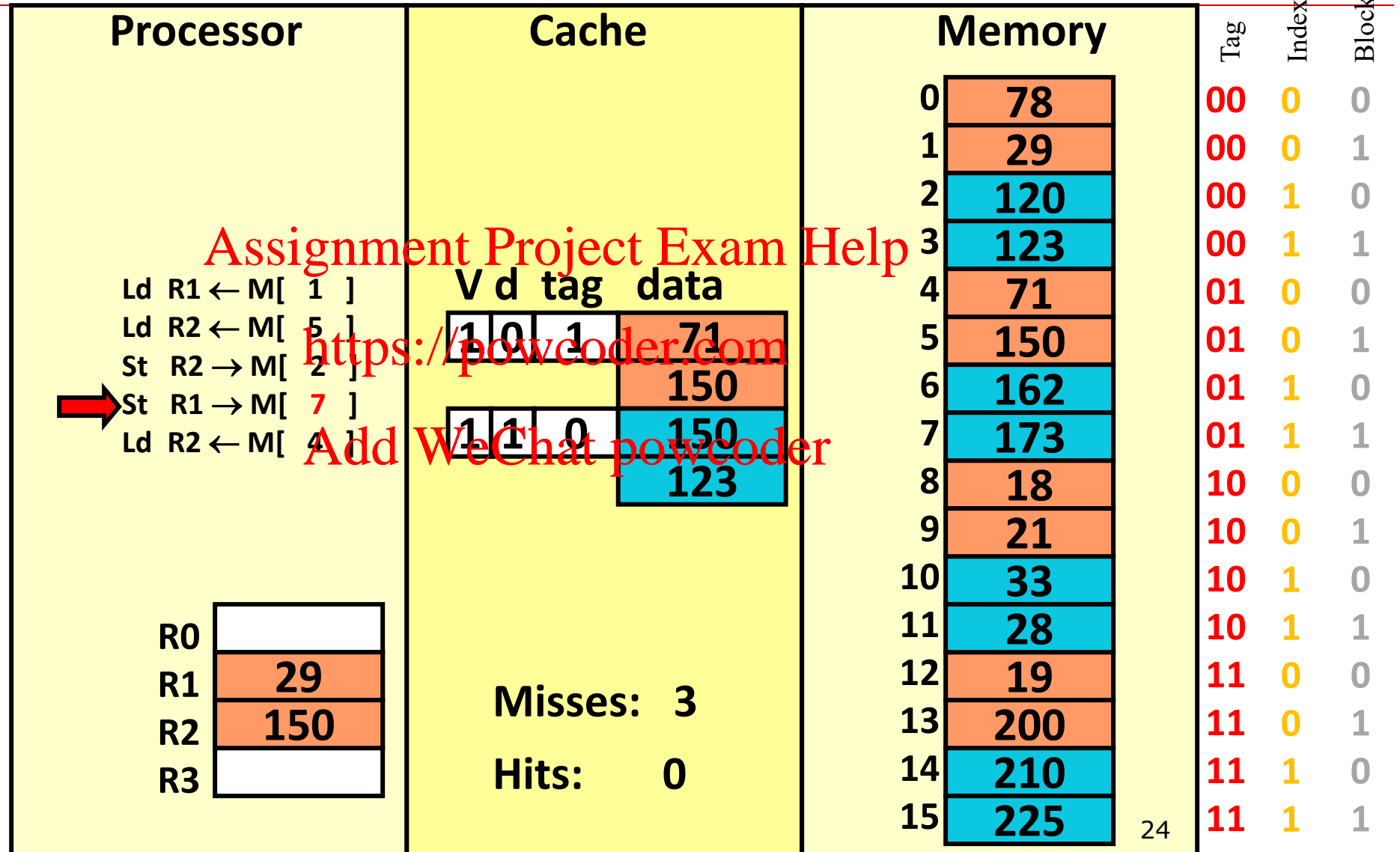
Direct-mapped (REF 3)



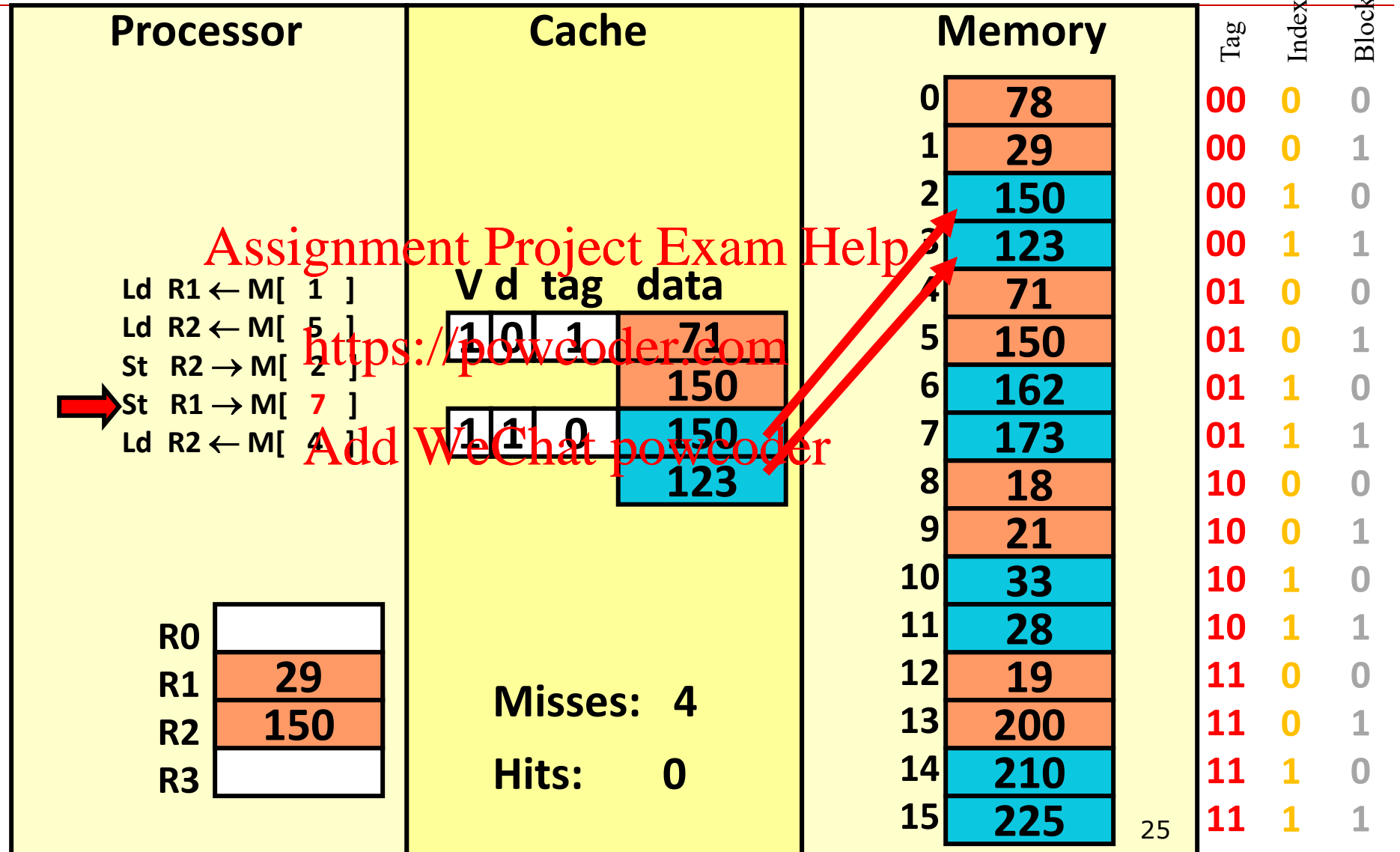
Direct-mapped (REF 3)



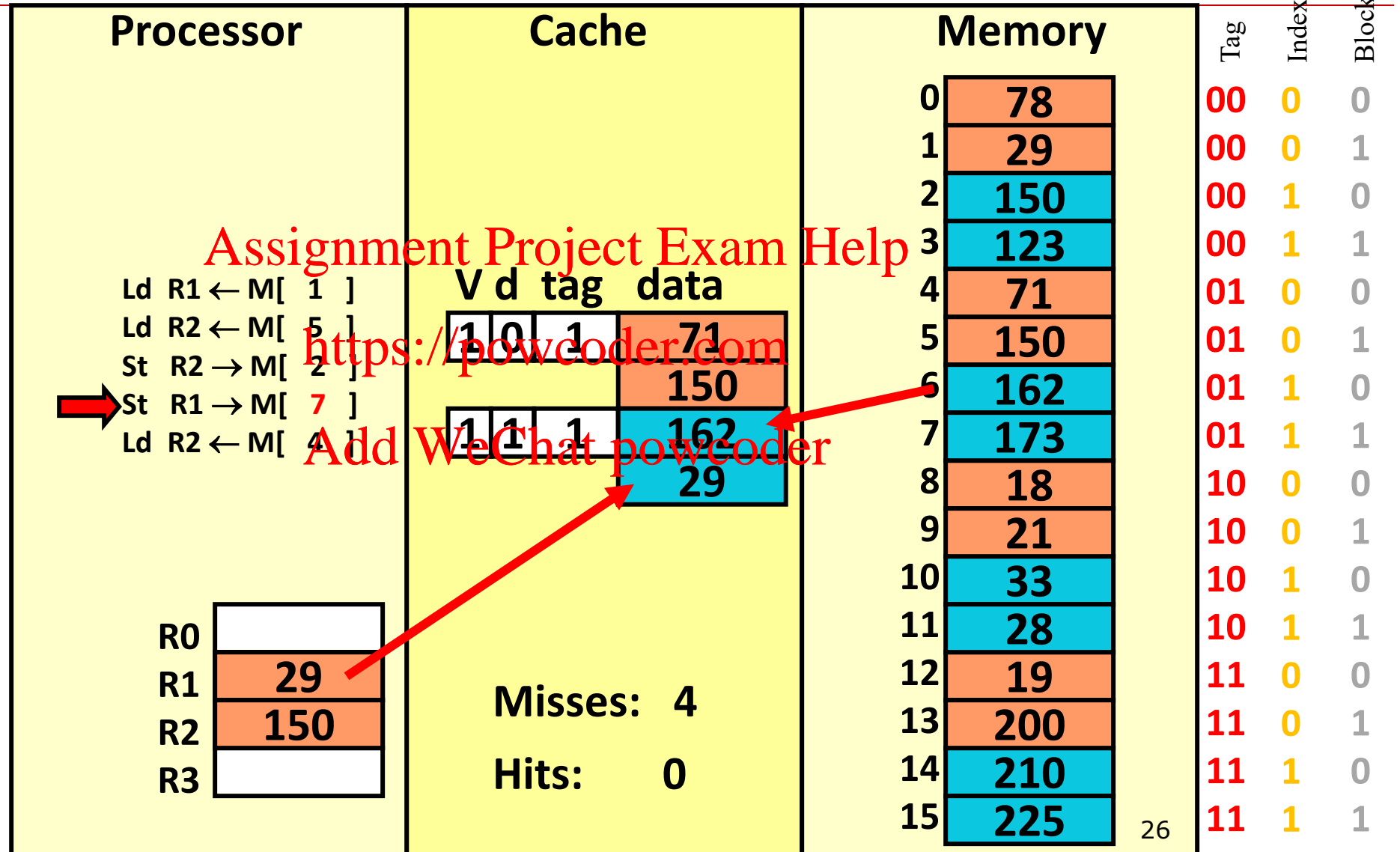
Direct-mapped (REF 4)



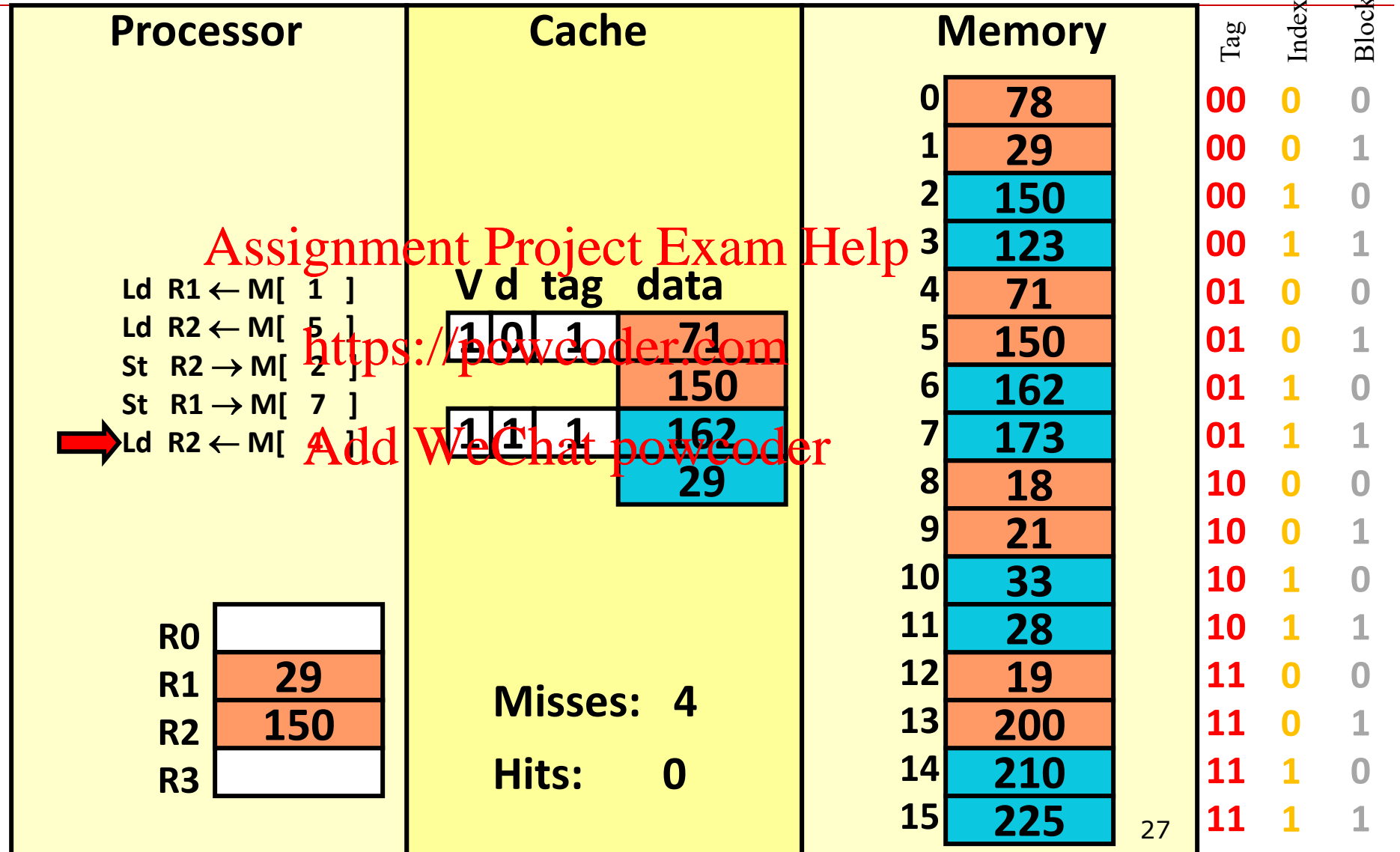
Direct-mapped (REF 4)



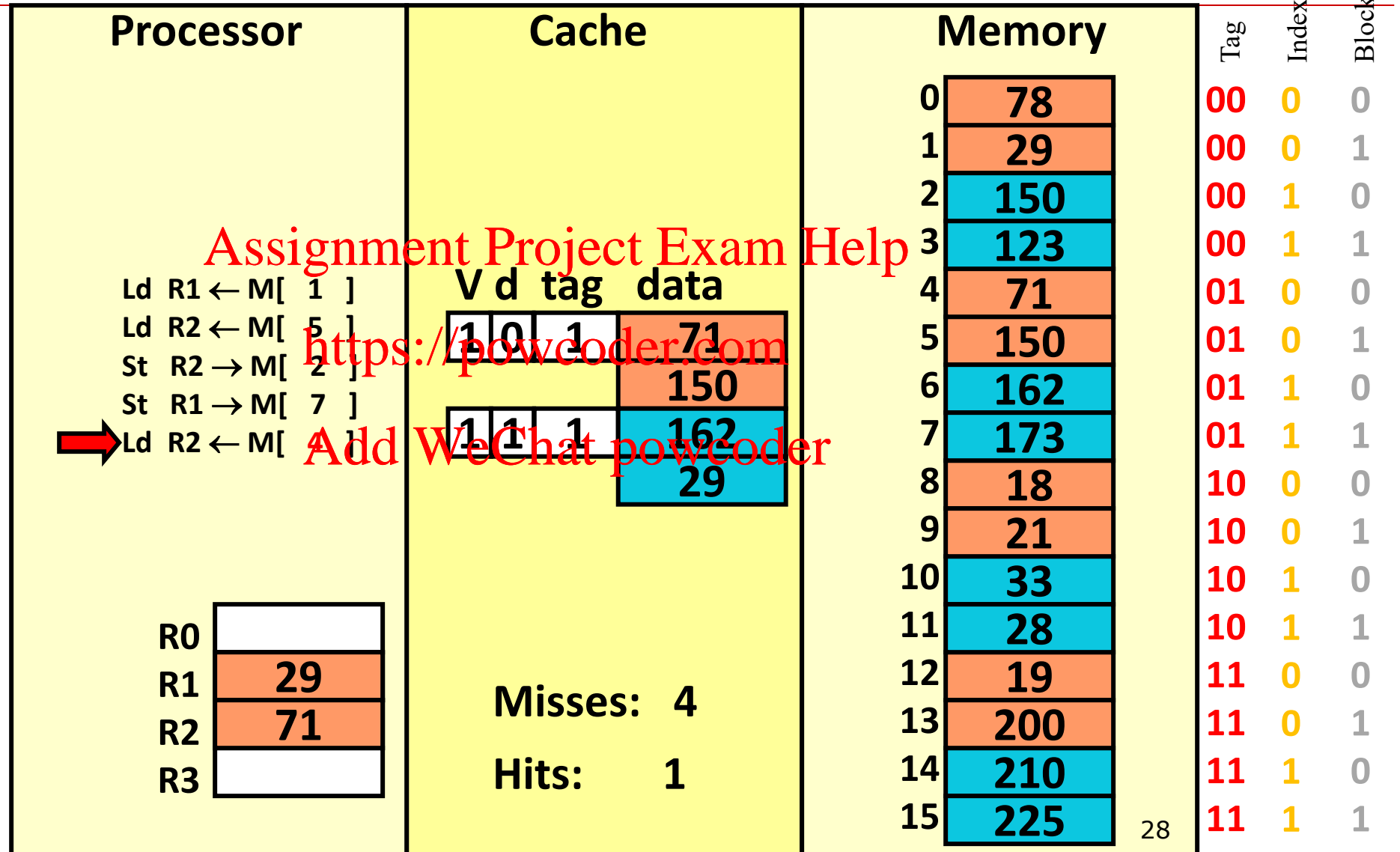
Direct-mapped (REF 4)



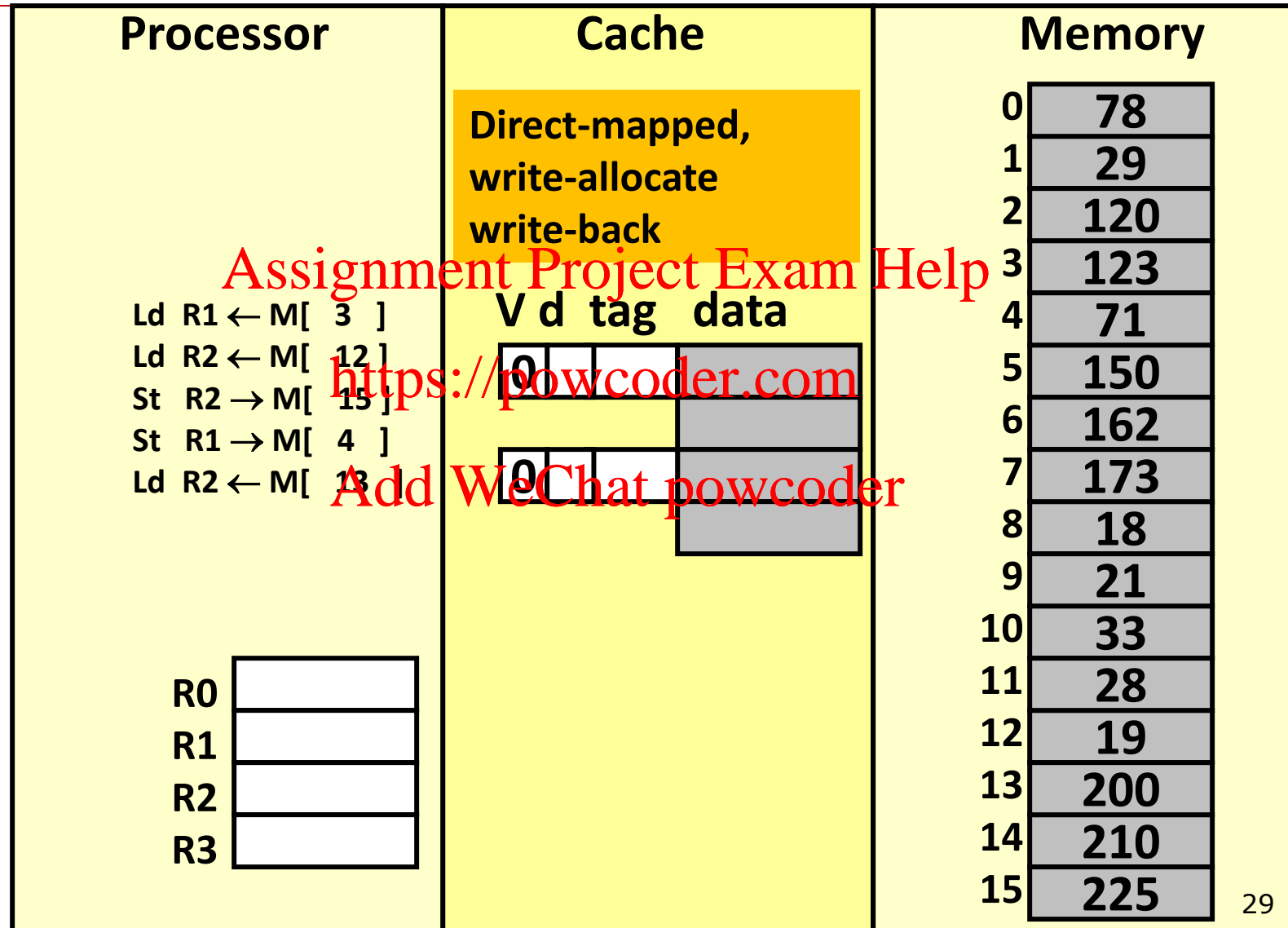
Direct-mapped (REF 5)



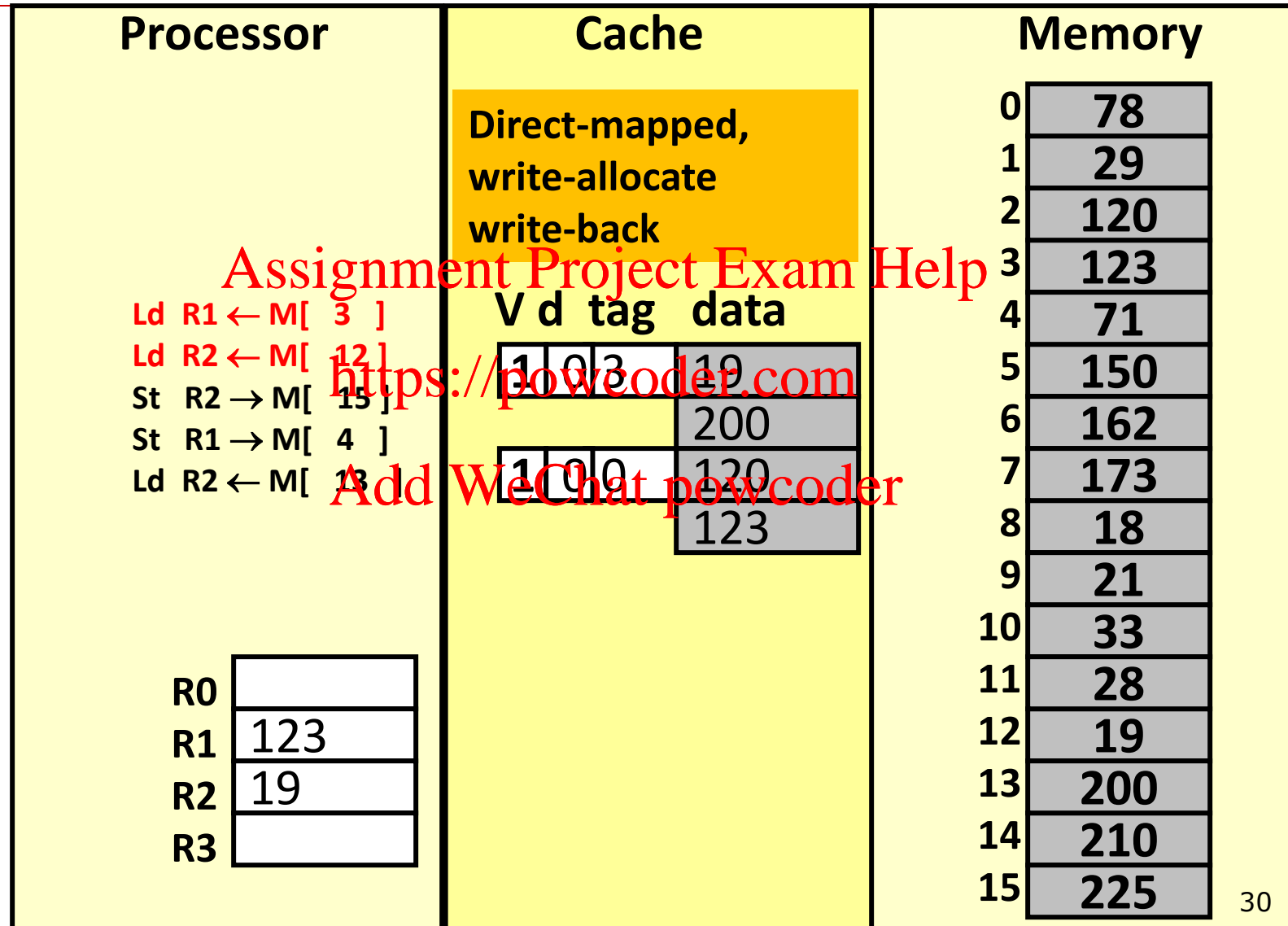
Direct-mapped (REF 5)



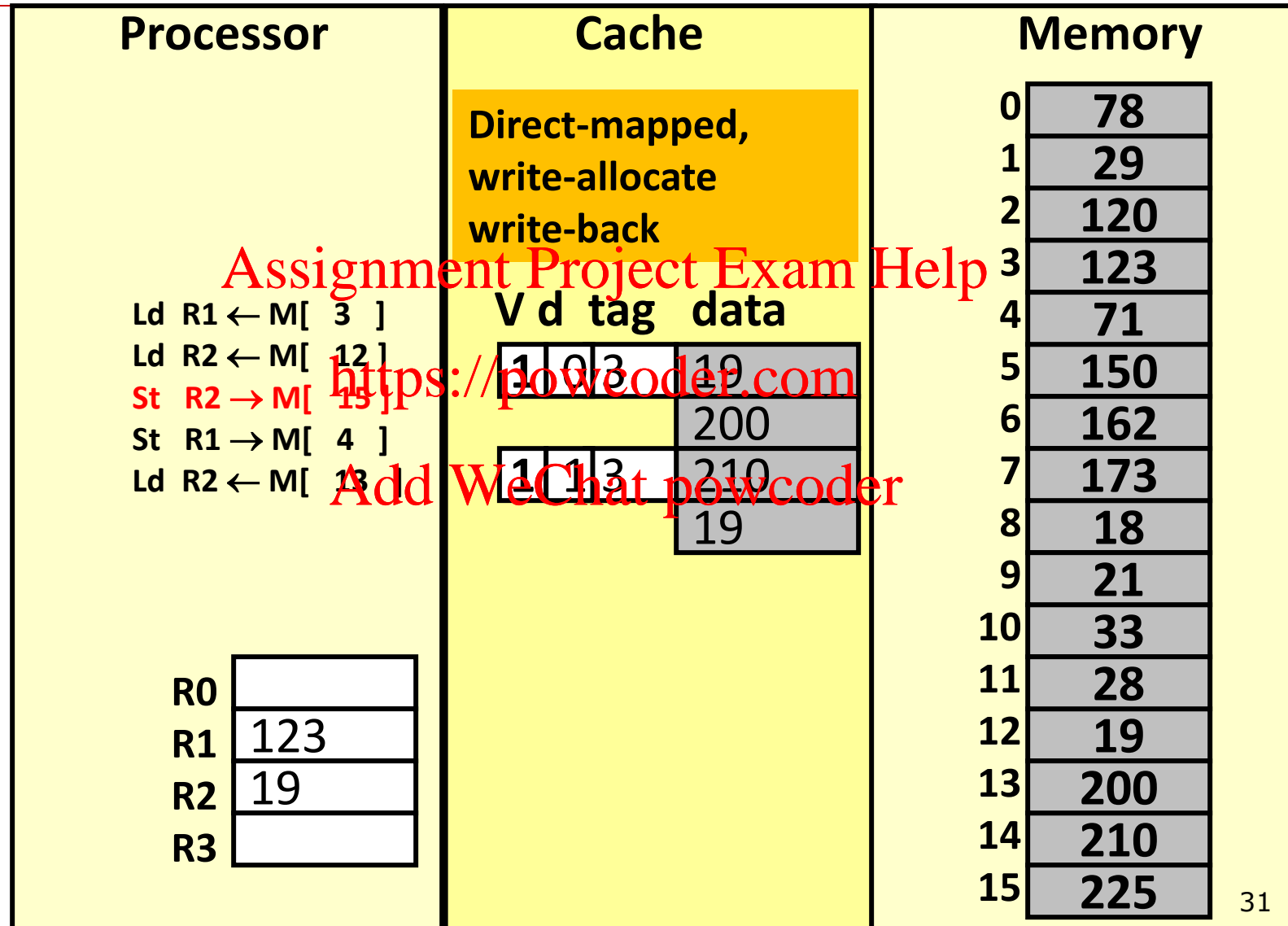
Class Problem – What is the state of the cache after executing the following instruction sequence?



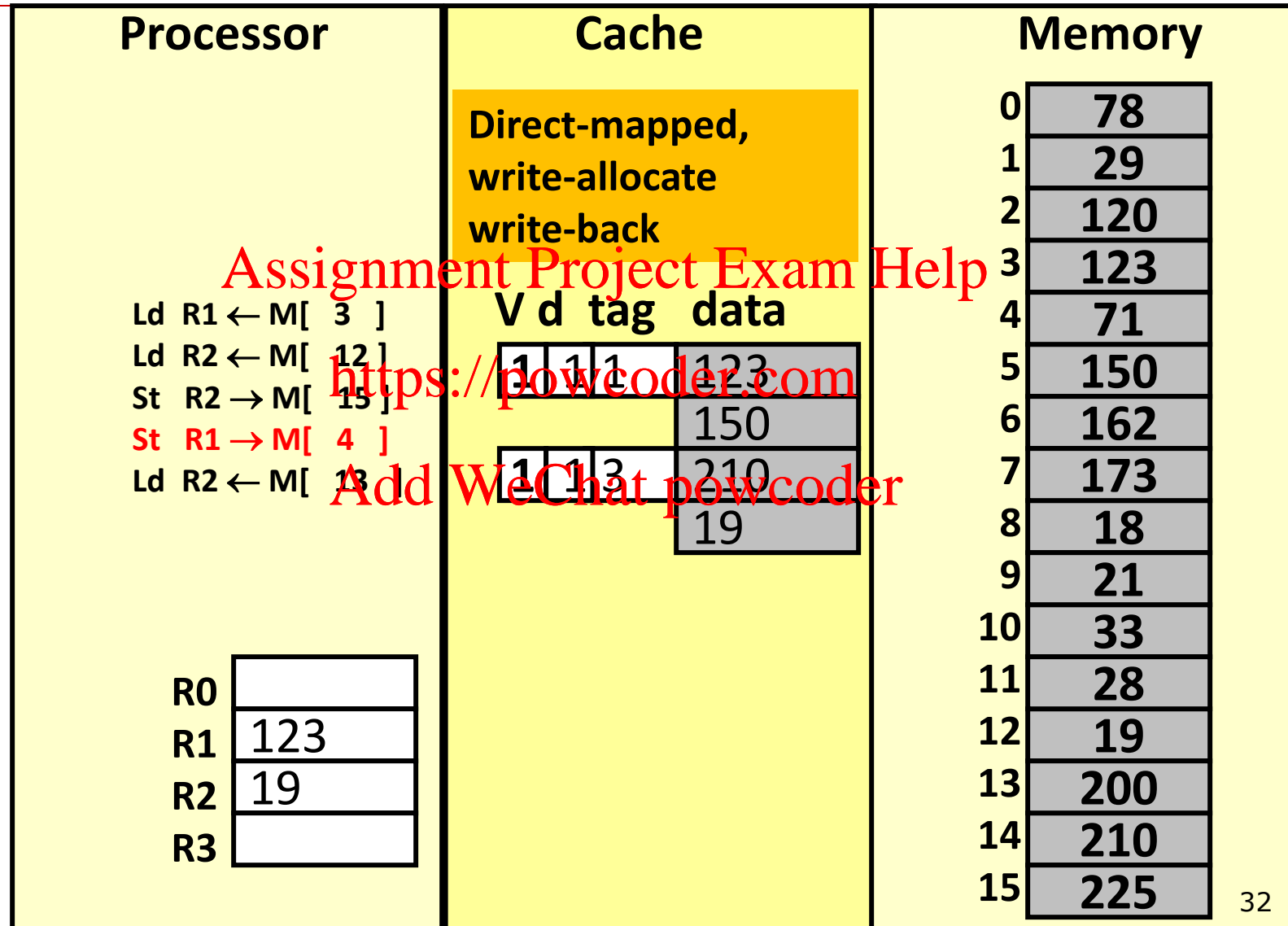
Class Problem – What is the state of the cache after executing the following instruction sequence?



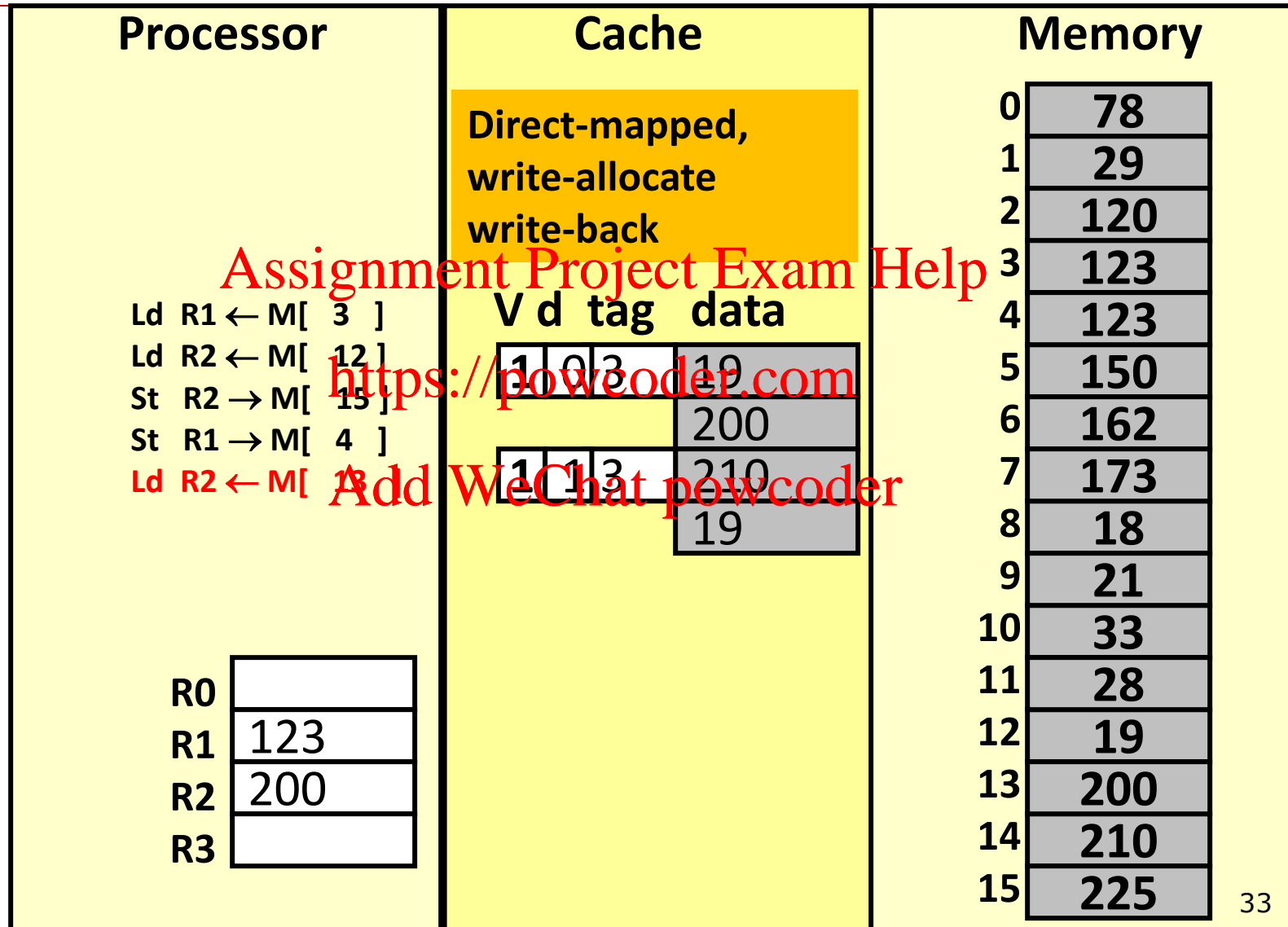
Class Problem – What is the state of the cache after executing the following instruction sequence?



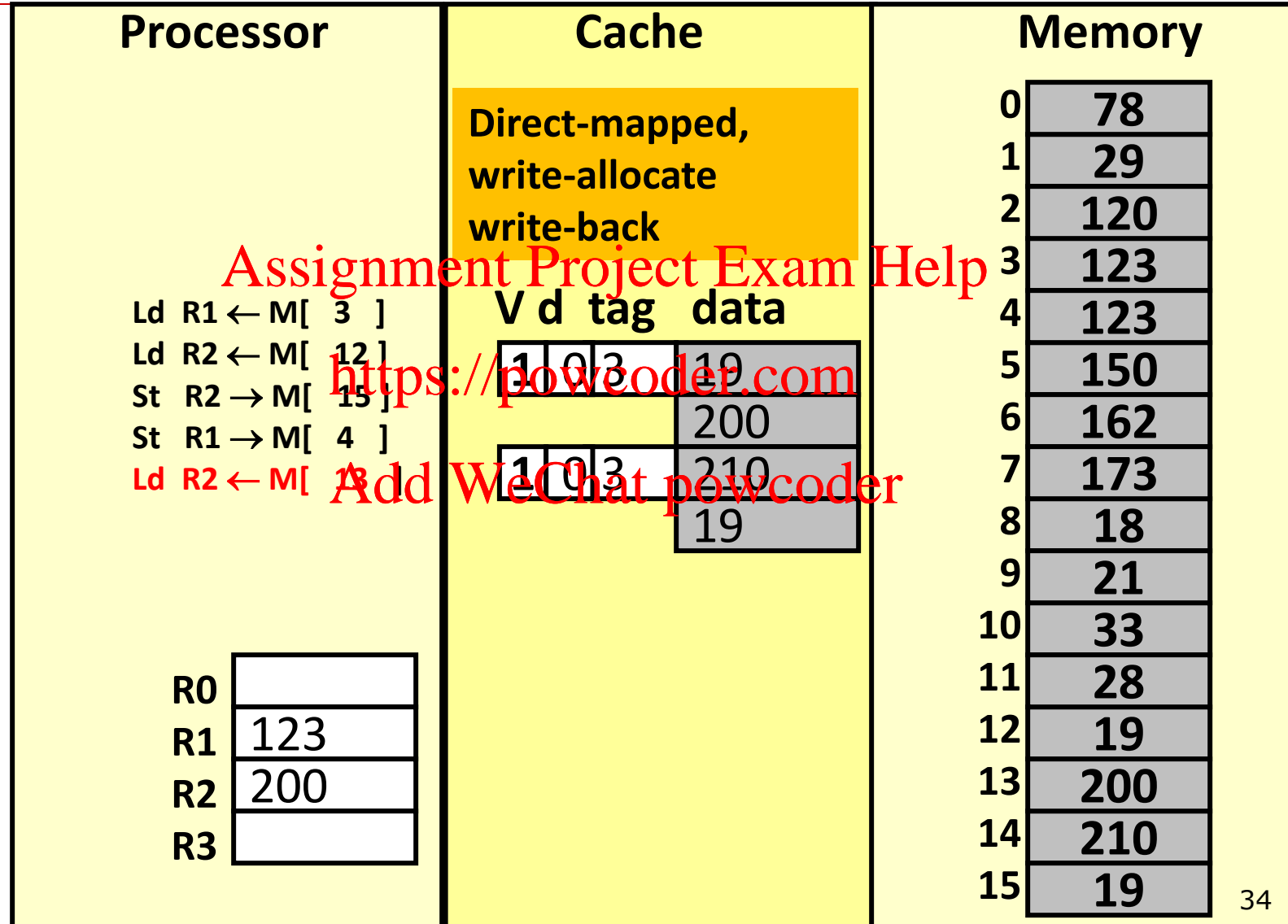
Class Problem – What is the state of the cache after executing the following instruction sequence?



Class Problem – What is the state of the cache after executing the following instruction sequence?



Class Problem – What is the state of the cache after executing the following instruction sequence?



Class Problem

How many tag bits are required for:

32-bit address, byte addressed, direct-mapped 32k cache, 128 byte block size, write-back

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

What are the overheads of this cache?

Class Problem

How many tag bits are required for:

32-bit address, byte addressed, direct-mapped 32k cache, 128 byte block size, write-back

Bytes in block = 128 \Rightarrow Block offset size = 7 bits (*byte addressable*)

Lines = 32k / 128 = 256 \Rightarrow Line index = 8 bits

Tag bits = 32 - 7 - 8 = 17 bits

Add WeChat powcoder

What is the overhead of this cache?

17 bits (Tag) + 1 bit (Valid) + 1 bit (Dirty) = 19 bits / line

19 bits / line * 256 lines = 4864 bits

4864 bits / 32KB = 1.9% overhead

What about cache for instructions?

Instructions should be cached as well

We have two choices:

1. Treat instruction fetches as normal data and allocate cache lines when fetched
2. Create a second cache (called the instruction cache or ICache) which caches instructions only

<https://powcoder.com>

How do you know which cache to use?

Add WeChat powcoder
What are advantages of a separate ICache?

Integrating Caches into a Pipeline

How are caches integrated into a pipelined implementation?

Replace instruction memory with Icache

Replace data memory with Dcache

Assignment Project Exam Help

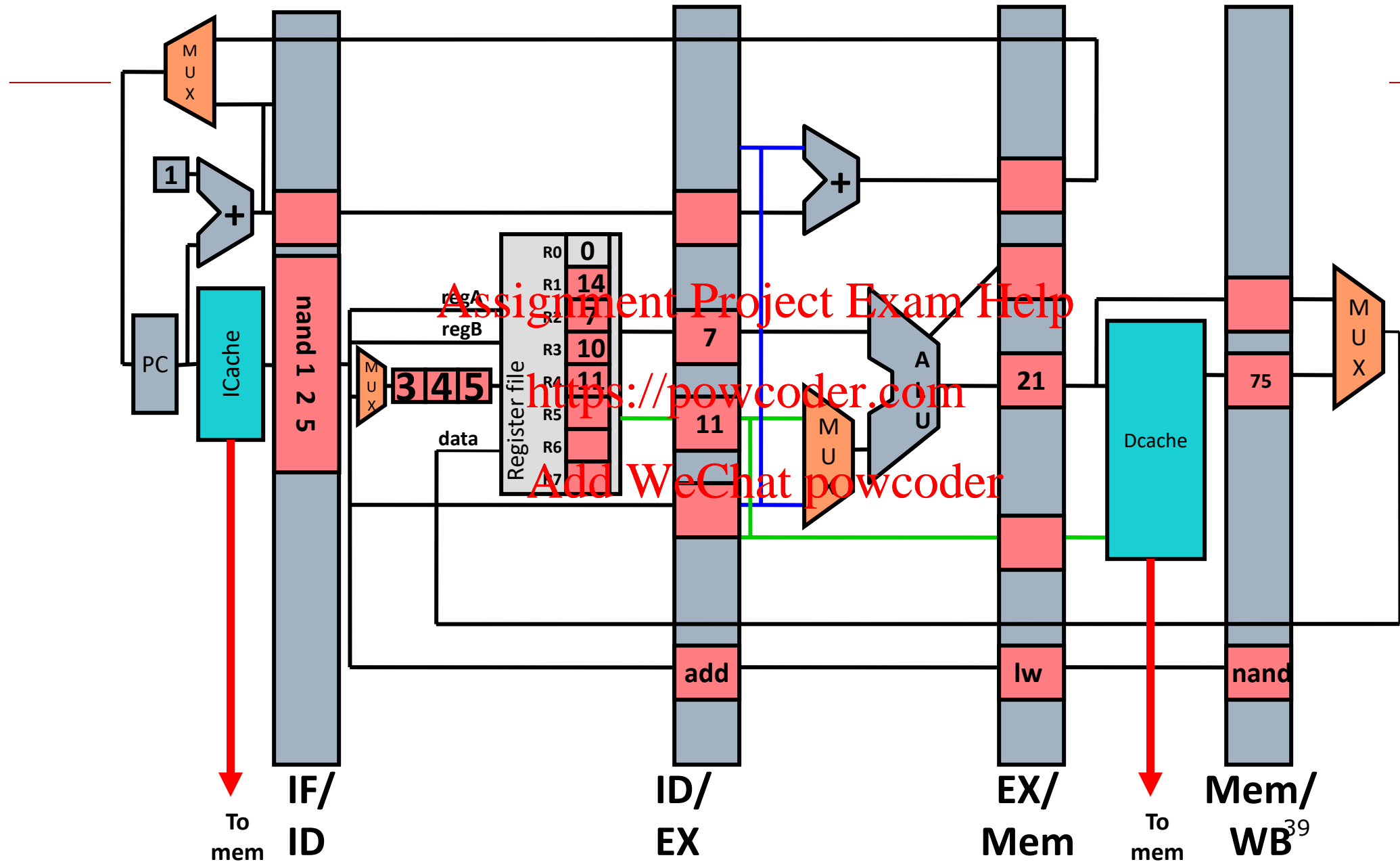
Issues:

Memory accesses now have variable latency

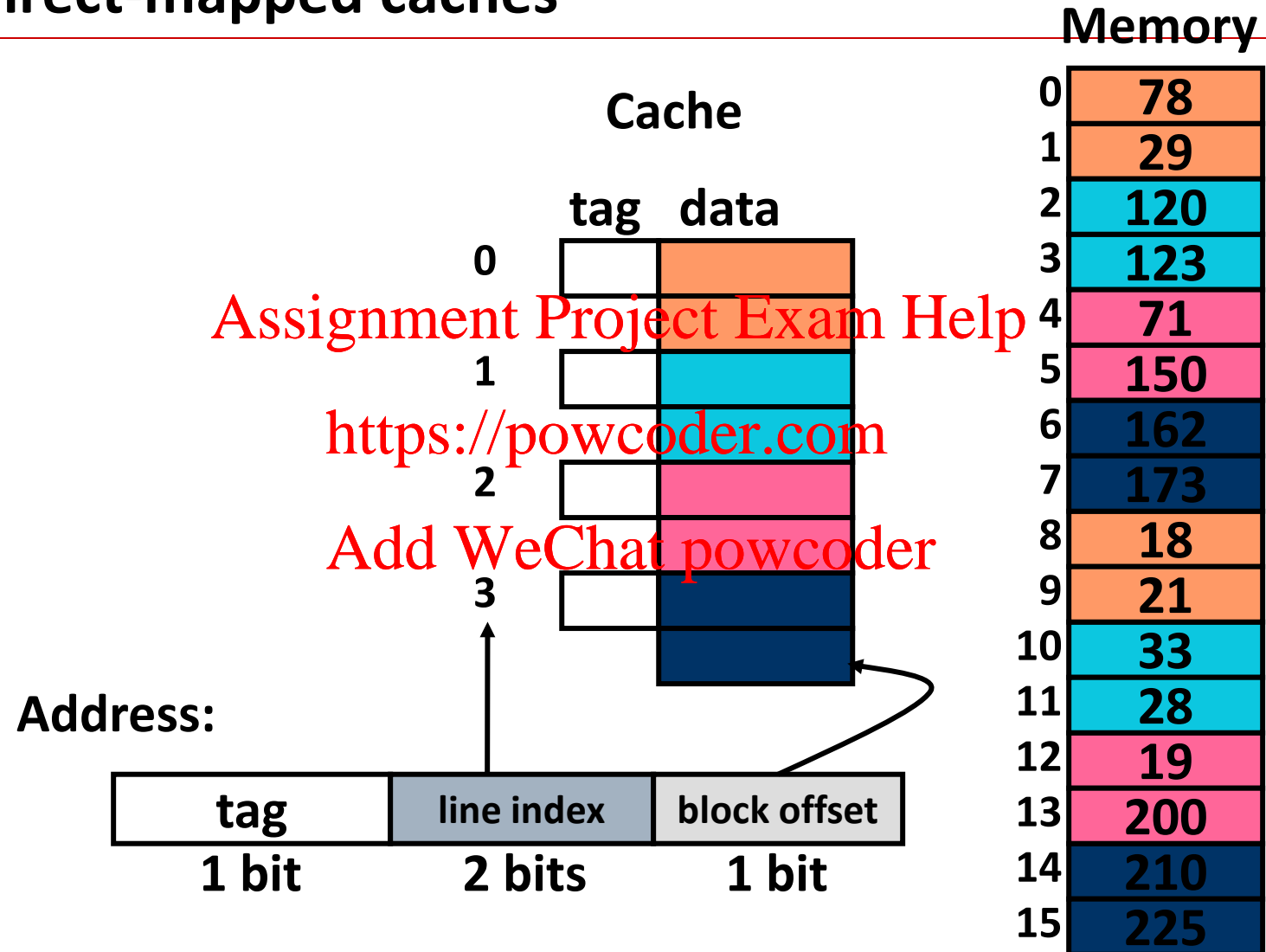
Both caches may miss at the same time

Add WeChat powcoder

LC2K Pipeline with Caches



Summary: Direct-mapped caches



Next lecture: Get the advantage of both...

Set **associative** caches:

- Partition memory into regions

 - like direct mapped but fewer partitions

- Associate a region to a set of cache lines

 - Check tags for all lines in a set to determine a HIT

- Treat each line in a set like a small fully associative cache

- LRU (or LRU-like) policy generally used

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Set-associative cache

