

17. Cache and memory hierarchy: The basics

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

Instructor switch:

Professor Satish Narayanasamy

Taking over from Professor Bill Arthur

Covering caches and virtual memory (asynchronous lectures #17 - #25)

<https://powcoder.com>



Upcoming deadlines:

Add WeChat powcoder

HW4

due Nov 10th

Project 3

due Nov. 12th

Assignment Project Exam Help

Part 1: Memory Hierarchy and Caches: Introduction

Add WeChat powcoder

Memory seen in previous lectures

LC2K data-paths have these structures that hold data and instructions:

Register file (little array of words)

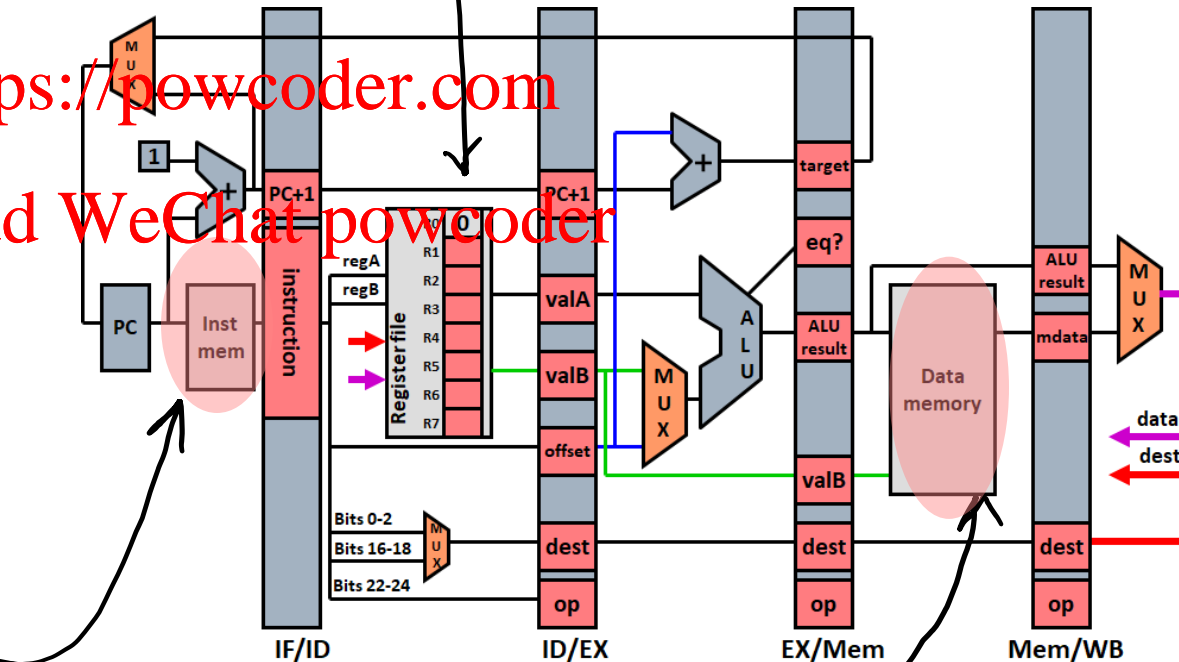
Memory (bigger array of words)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Upcoming lectures:
How to design “memory”
(highlighted parts)?



Memory System: Learning objective

LC2k program	can access	2^{18} bytes of memory
MIPS program	can access	2^{32} bytes of memory
ARM64 or x86-64 program	can access	2^{64} bytes of memory (18 billion billion bytes!)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Problem: No one memory technology is both fast and big to store all of program's data

Goal: Design a fast, big, and cheap memory system to store a program's data.

Memory System: Desirable Properties

Big memory

Fast memory

A load instruction would stall a data-path, if memory access takes longer than a cycle

Assignment Project Exam Help

<https://powcoder.com>

Cheap memory

Measured as cost per byte of data. A critical component that determines cost of computers.

Add WeChat powcoder

Volatile or not?

Does the data vanish or persist when power is turned off?

Memory Pyramid

Fast

Cost **Expensive**

Size **Small**

Cache
(SRAM)

<https://powcoder.com>

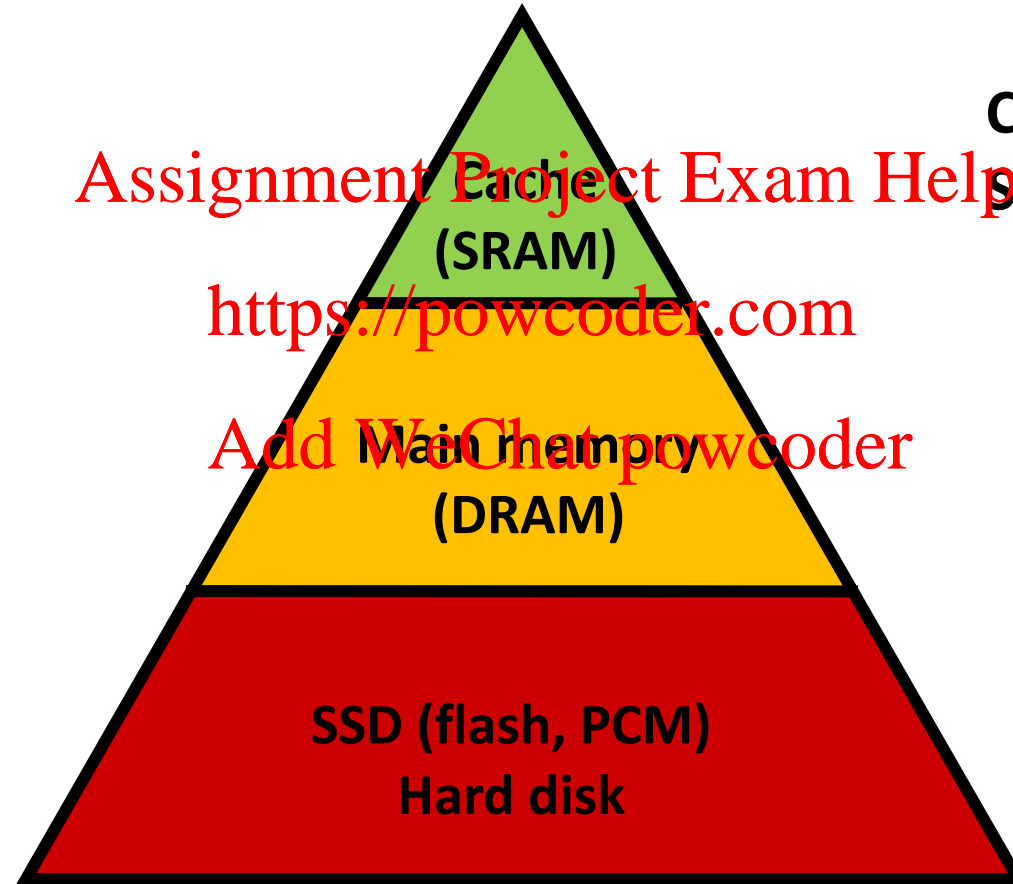
Main memory
(DRAM)

Slow

Cost **Cheap**

Size **Big**

SSD (flash, PCM)
Hard disk



SRAM (Static RAM)

Area: 6T: 6 transistors per bit (used on-chip within processor)

Fast: ~2ns access time, if size is small (few KBs)

Larger the size, longer the access time

Typical Size: Tens of KBs to a few MBs

Cost: Expensive

~\$5.0 per megabyte

\$0.13	for 2^{18} bytes of memory	(LC2K ISA)
\$20,000	for 2^{32} bytes of memory	(MIPS ISA)
\$88 trillion	for 2^{64} bytes of memory	(ARM64 ISA)

Volatile

DRAM (Dynamic RAM)

Area: A tiny capacitor and a transistor per bit

Slower: ~60ns access time (for few GBs of size)

Typical Size: Tens of GBs

Cost: Less expensive than SRAM

~\$0.004 per megabyte

\$0.00

for LC2K

\$16

for MIPS

\$70,000,000,000

for ARM64

Volatile

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Disks

Obnoxiously slow: 3,000,000ns access time

Typical size: tens of TBs

Cost: Cheaper than SSDs (flash, PCM)

\$0.000043 per megabyte

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

\$0.00 LC2

\$0.18 for MIPS

\$760,000,000 for ARM64

Non-volatile

Flash

Floating-gate transistors. SSDs have replaced hard disks in mobile phones and laptops.

Slower still: $\sim 250\text{ns}$ access time

Assignment Project Exam Help

Typical size: hundreds of GBs to a few TBs

<https://powcoder.com>

Cost: Less expensive than DRAM

Add WeChat powcoder

\$0.0012 per megabyte

\$0.00

for LC2

\$4.9

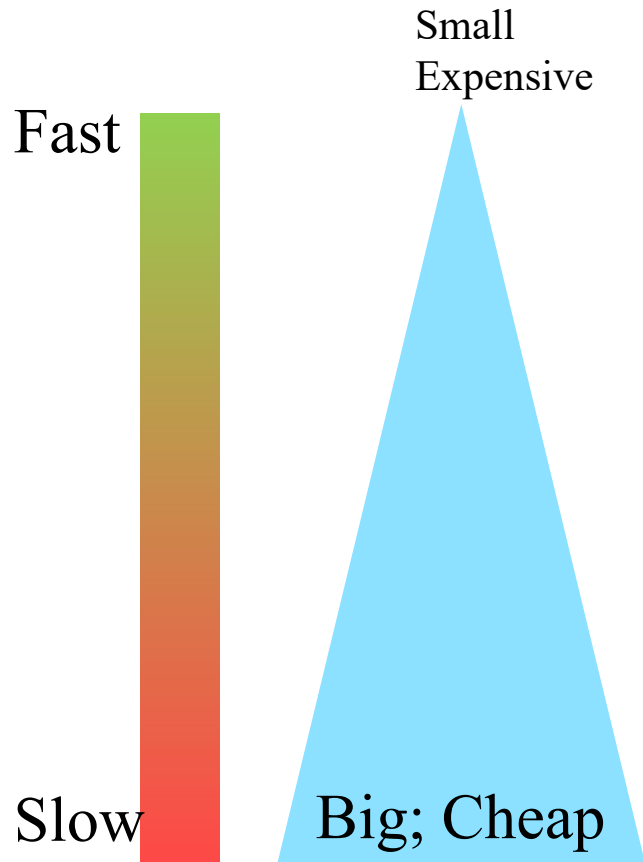
for MIPS

\$21,000,000,000

for ARM64

Non-volatile

Memory Technologies: Summary



Static RAM (SRAM)

Dynamic RAM (Dynamic memory)

Used inside processors

“Main” memory

volatile

Non-volatile

Solid state disks

Phase-change memory (PCM)

Flash

Emerging (e.g., [Intel Optane](#))

Common in mobile devices

Magnetic Disk (hard disk)

DVD

Magnetic tape

In research: [DNA storage](#)

Memory Hierarchy Goal

How to get best properties of different memory technologies?

A memory system that is as **fast** as SRAM, but as **big and cheap** as a hard-disk?

Assignment Project Exam Help

Fast:

Ideally run at processor's clock speed

1 ns access time

Add WeChat powcoder

Big and Cheap:

Sufficiently large to hold a program's data

Memory Hierarchy Analogy: Storing and retrieving a book

Option 1: Library stores all the books. Every time you switch to another book, return current book to library and get a new book.

Latency = few hours



Assignment Project Exam Help

Option 2: Borrow 20 frequently-used books and keep them at home book-shelf

Latency = few minutes (mostly, go to library once a week or so)

<https://powcoder.com>



Option 3: Keep 3 books in backpack

Latency = few seconds (mostly, go to book-shelf once a day or so)

Add WeChat powcoder



Memory hierarchy: Leveraging locality of reference

Fast

Cost **Expensive**

Size **Small**

Cache
(SRAM)

<https://powcoder.com>

Main memory
(DRAM)

Slow

Cost **Cheap**

Size **Big**

Disk
(magnetic or floating gate)

Temporarily move what you use here

For a program with good locality of reference, memory appears as fast as cache and as big as disk

Have a copy of everything here

A Realistic Memory Hierarchy

Cache

Few KBs to MBs of SRAM (within processor – on-chip cache)

Fast

Small, so cheap

Serves most loads and stores, provided program has good locality

Assignment Project Exam Help

<https://powcoder.com>

Main Memory

Tens of GBs of DRAM (outside processor – off-chip)

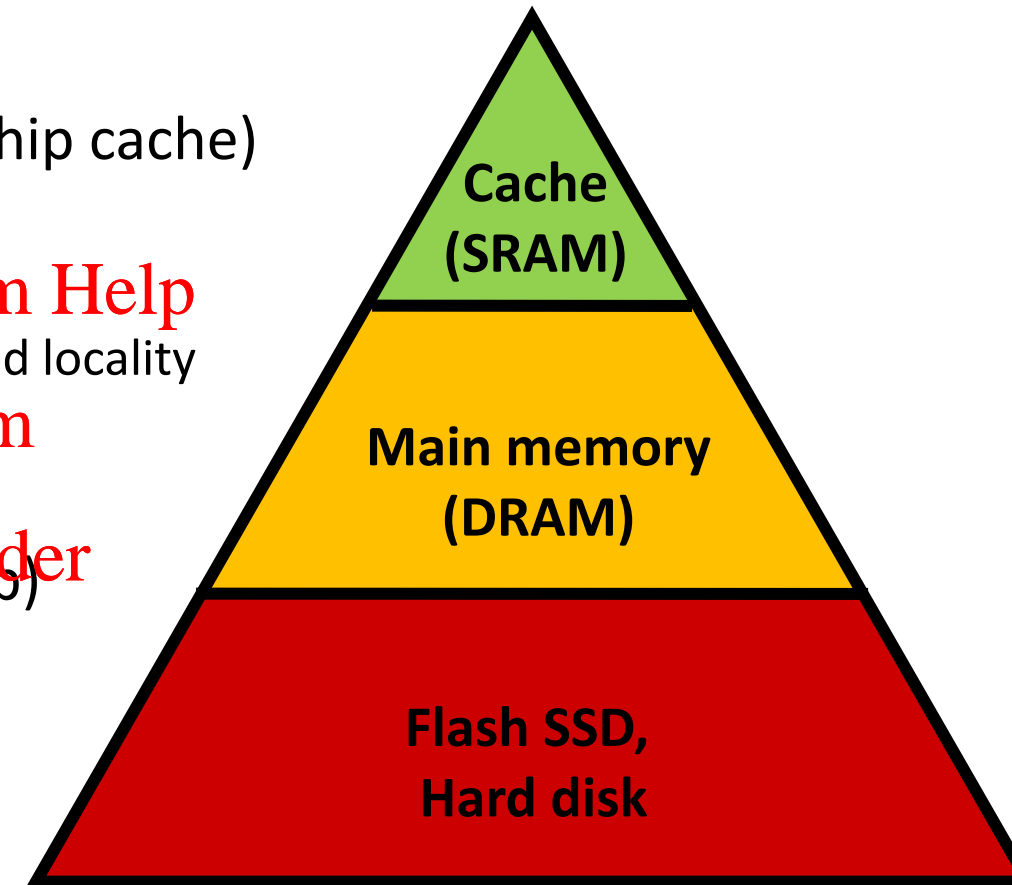
Cheaper than SRAM, faster than flash/disk

Add WeChat powcoder

“Swap space”

Few TBs OF flash and/or disk

Cheap, Big, Non-volatile.



No memory is enough for a 64-bit ISA (ARM64) program

Hard disk cost for storing all addresses accessible to a ARM64 program
\$760 million for 2^{64} bytes



Don't provision 2^{64} bytes of storage (even a hard disk is too expensive!)

Assignment Project Exam Help

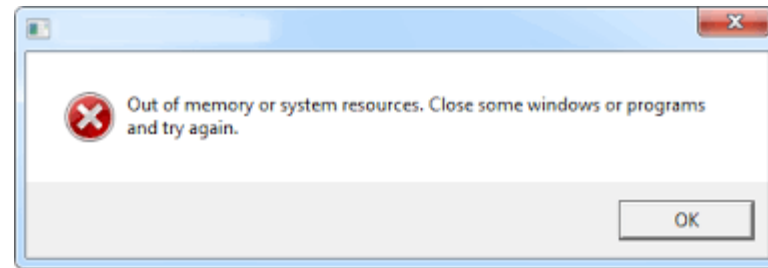
<https://powcoder.com>

Fake it. Use “virtual memory” to provide an illusion that ISA's entire address space is available.

Add WeChat powcoder

A few TB is enough for most desktop machines today, or a smartphone in a few years

Computer “crashes” if your program exceeds machine's available swap space on disk



ISA abstraction hides memory hierarchy from programmers

The architectural view of memory is

- What the machine language (or programmer) sees above ISA
- Just a big array

Assignment Project Exam Help

Breaking up the memory system into different pieces

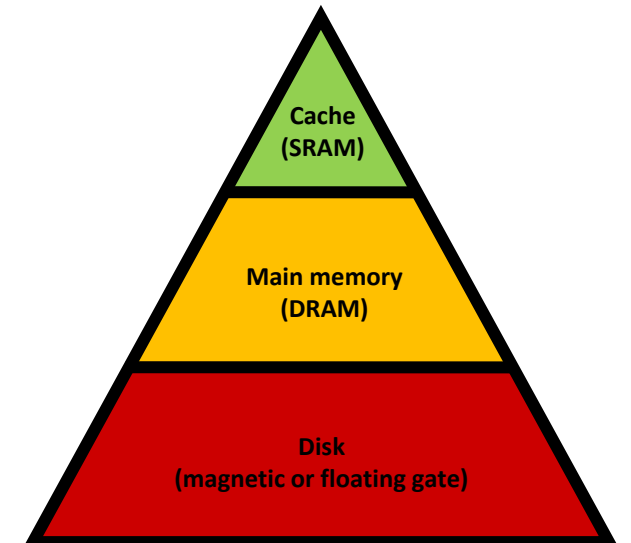
– cache (SRAM), main memory (DRAM) and disk –

is not architectural

- ISA does not expose these details to the programmer
- A new system implementation may break it up in a different way

Programmer
can load/store to
 2^{64} memory locations;
Can't see memory hierarchy

ARM-64 ISA



What is a cache?

Cache commonly refers to SRAM used on-chip within the processor.

However, even DRAM (main memory) is a “cache” in that it temporarily stores data fetched from hard disk

Assignment Project Exam Help

A cache is used to store data that is **most likely** to be **referenced** by a program

Try to maximize the number of references (loads/stores) that are serviced by the cache (avoid going slow, off-chip, main memory; or even worse, disk).

Thereby, minimize the average memory access time (AMAT) of a load/store

Cache: Importance

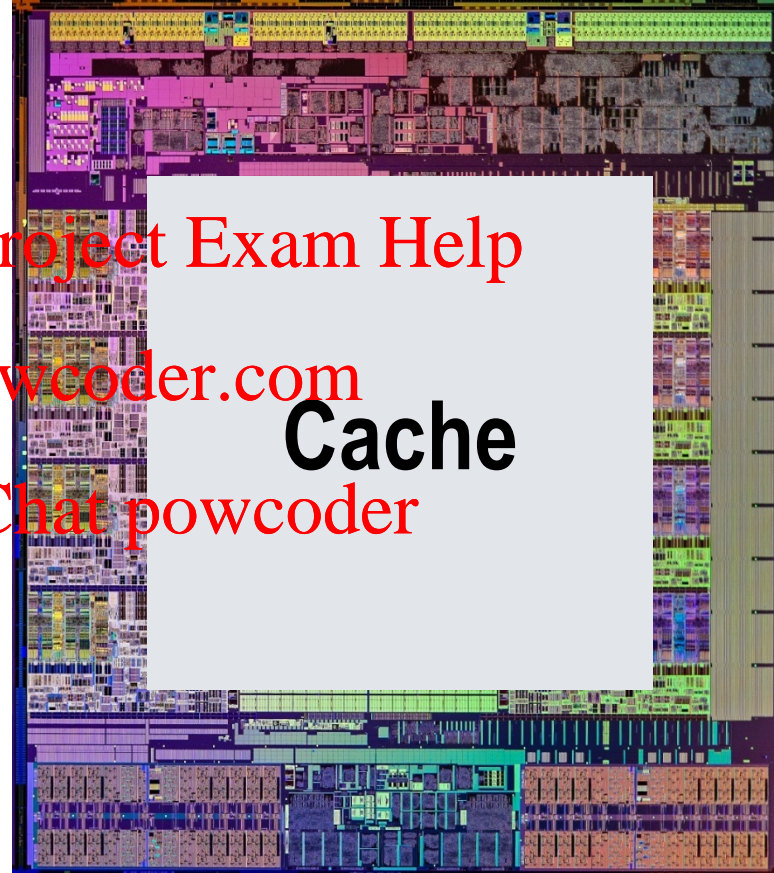
Caches consume
most of a processor's die area

Assignment Project Exam Help

<https://powcoder.com>

Cache

Add WeChat powcoder



Importance of Cache on Performance:

Cache Aware code can be several times faster than non-Aware code

```
#include<stdio.h>
#include<stdlib.h>

#define N 20000
int arrayInt[N][N];

int main(int argc, char **argv)
{
    int i, j;
    int count = 0;

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

    printf("Count :%d\n", count);
}
```

```
#include<stdio.h>
#include<stdlib.h>

#define N 20000
int arrayInt[N][N];

int main(int argc, char **argv)
{
    int i, j;
    int count = 0;

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

    printf("Count :%d\n", count);
}
```

Live demo:

See [L1_3_370_Course_Overview](#)

Video at minute 18:00

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Cache Design: This lecture

Basic Cache Architecture

Assignment Project Exam Help

How to select data to store in cache?

Principle of “Temporal locality”
<https://powcoder.com>

Add WeChat powcoder

Illustration

Performance metric: average memory access time

Assignment Project Exam Help

Part 2: Basic Cache Architecture

<https://powcoder.com>

Add WeChat powcoder

Basic Cache Design

Cache memory can copy data from any part of main memory. It has 2 parts:

- The **TAG (CAM)** holds the memory address
- The **BLOCK (SRAM)** holds the memory data

addr	data
addr	data

TAG BLOCK Add WeChat powcoder

<https://powcoder.com>

Accessing the cache: **compare reference address and tag**

- Match? Get the data from the cache block
- No Match? Get the data from main memory
- How to implement this functionality? **Solution: CAM for storing tags**

CAMs: content addressable memories

Instead of thinking of memory as an array of data indexed by a memory address

Think of memory as a set of data, that can search for a queried key

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Operations on CAMs

❑ **Search:** the primary way to access a CAM

- Send data to CAM memory
- Return “found” or “not found”: “hit” or “miss”
- If found, return location of where it was found or associated value

Assignment Project Exam Help

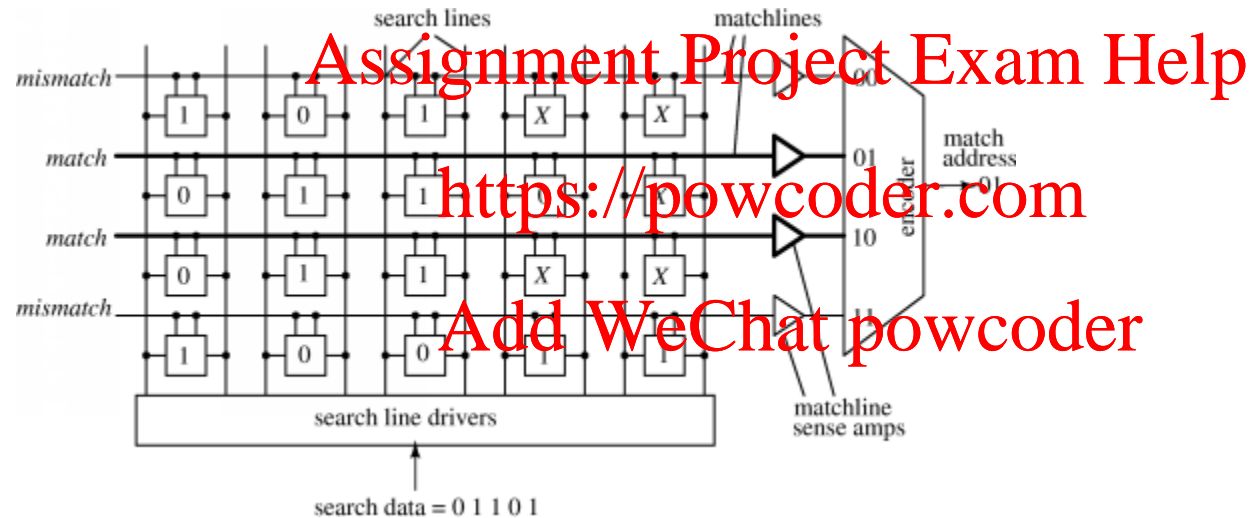
<https://powcoder.com>

Add WeChat powcoder

❑ **Write:**

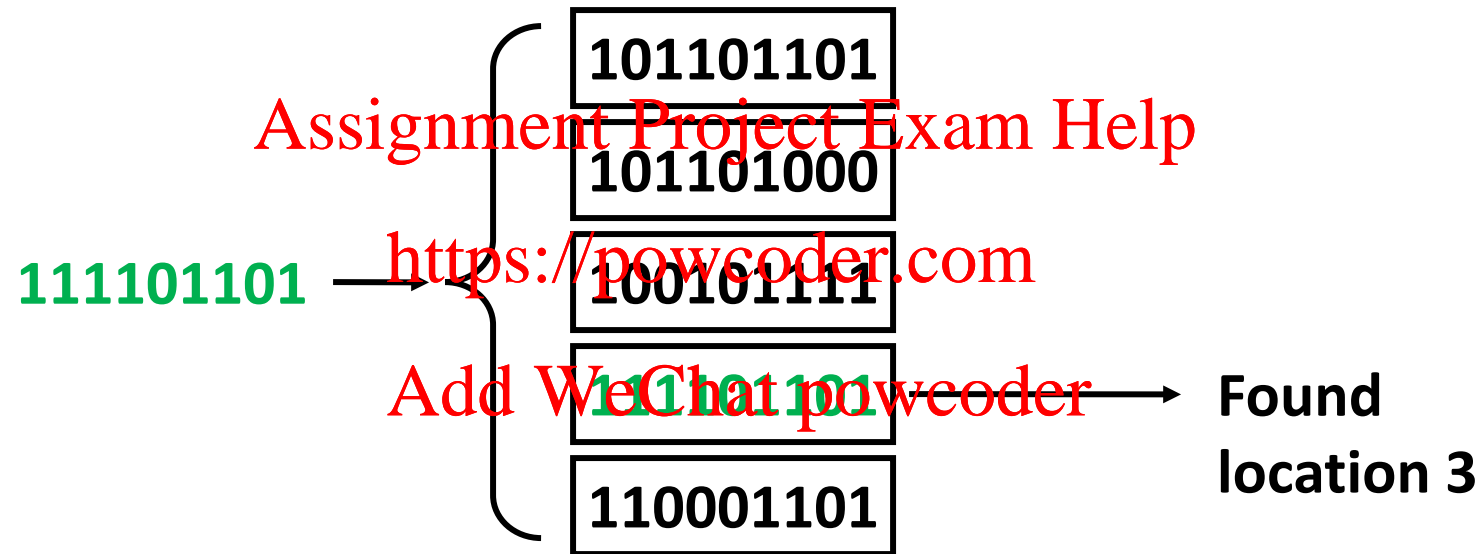
- Send data for CAM to remember
 - Where should it be stored if CAM is full?
 - Replacement policy
 - Replace oldest data in the CAM
 - Replace least recently searched data

CAM = content addressable memory



When used in caches, all tags are fully specified (no X – no don't cares)

CAM example



5 storage element CAM array of 9 bits each

Previous use of CAMs

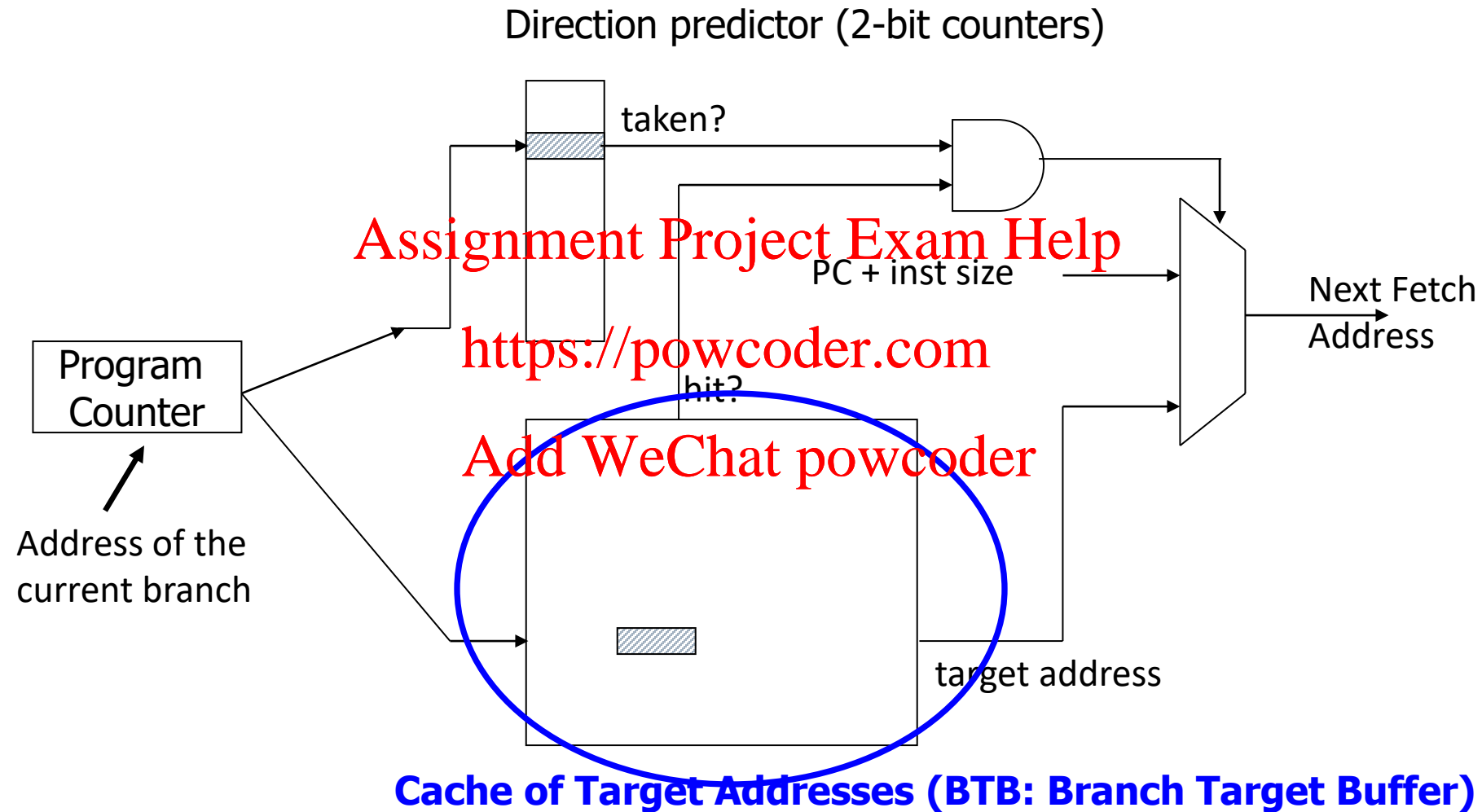
- ☐ You have seen a simple CAM used before. When?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Fetch Stage with Branch Prediction



Cache Organization

Cache memory can copy data from any part of main memory.

It has 2 parts:

- The **TAG (CAM)** holds the memory address
- The **BLOCK (SRAM)** holds the memory data

addr	data
addr	data

TAG

BLOCK

<https://powcoder.com>

Add WeChat powcoder

Cache Organization

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

addr	data
addr	data

TAG BLOCK

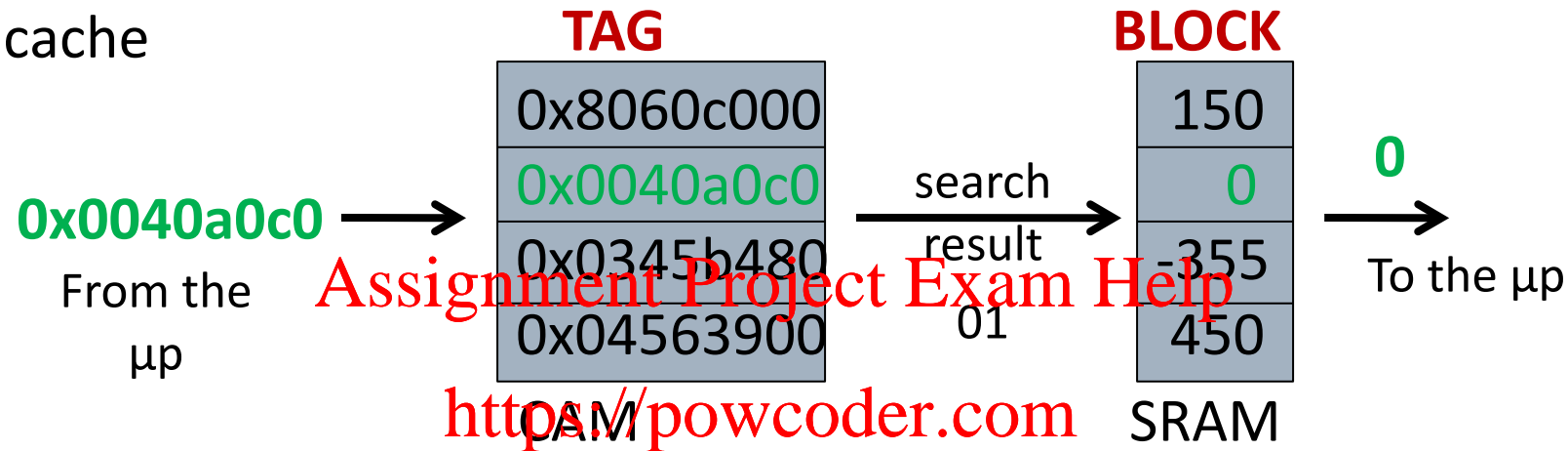
If there is a tag match, it is a cache **HIT**

If there is no tag match, it is a cache **MISS**

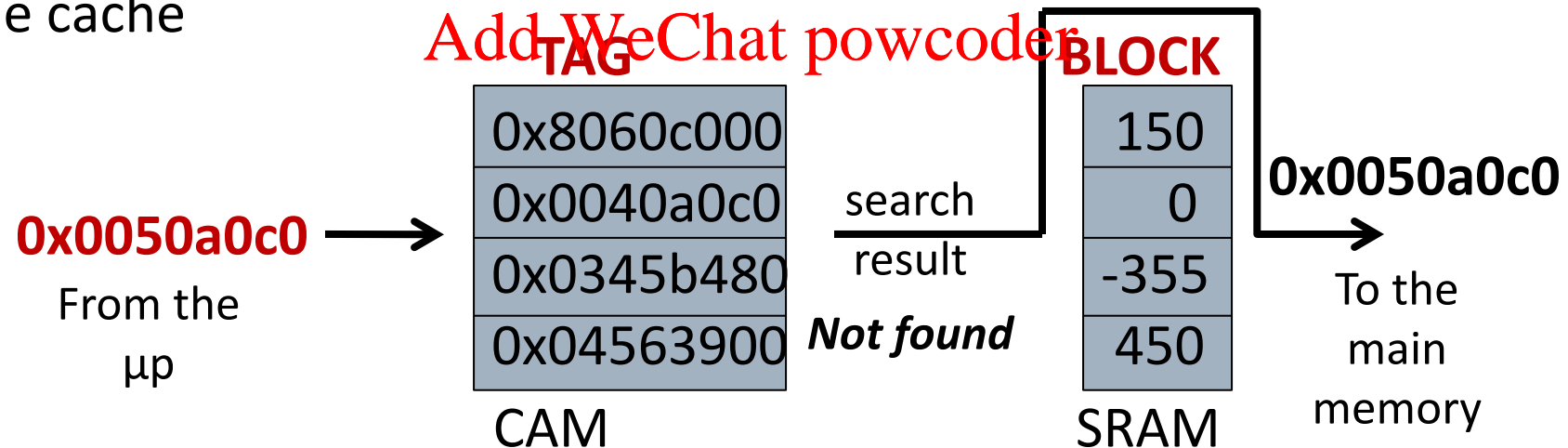
Goal: Cache data likely to be accessed in the future

Caches: the hardware view

A **hit** in the cache



A **miss** in the cache



Assignment Project Exam Help

<https://powcoder.com>
Part 3: Temporal locality

Add WeChat powcoder

Cache Operation

On a cache miss:

Fetch data from main memory and

Assignment Project Exam Help
Allocate a cache line and store fetched data in it
<https://powcoder.com>

Which cache line should be allocated? **Add WeChat powcoder**

If all cache lines are allocated, how to pick the victim for data replacement?

Something To Think About

Does an optimal replacement policy exist?

That is, given a choice of cache lines to replace, which one will result in the fewest total misses during program execution

Assignment Project Exam Help

<https://powcoder.com>

Why would we care?

Add WeChat powcoder

Picking the Most Likely Addresses

What is the probability of accessing a random memory location?

With no information, it is just as likely as any other address

Assignment Project Exam Help

But programs are not random

<https://powcoder.com>

They tend to use the same memory location over and over

Add WeChat powcoder

Temporal Locality

The principle of **temporal locality** in program references says that if you access a memory location (e.g., 0x1000) you will be more likely to re-access that location (e.g., 0x1000) than you will be to reference some other random location

Assignment Project Exam Help

Temporal locality says any miss location should be placed into the cache

It is the most recent reference location

Add WeChat powcoder

Temporal locality says that data in least recently referenced (or least recently used – **LRU**) cache line should be **evicted** to make room for the new line

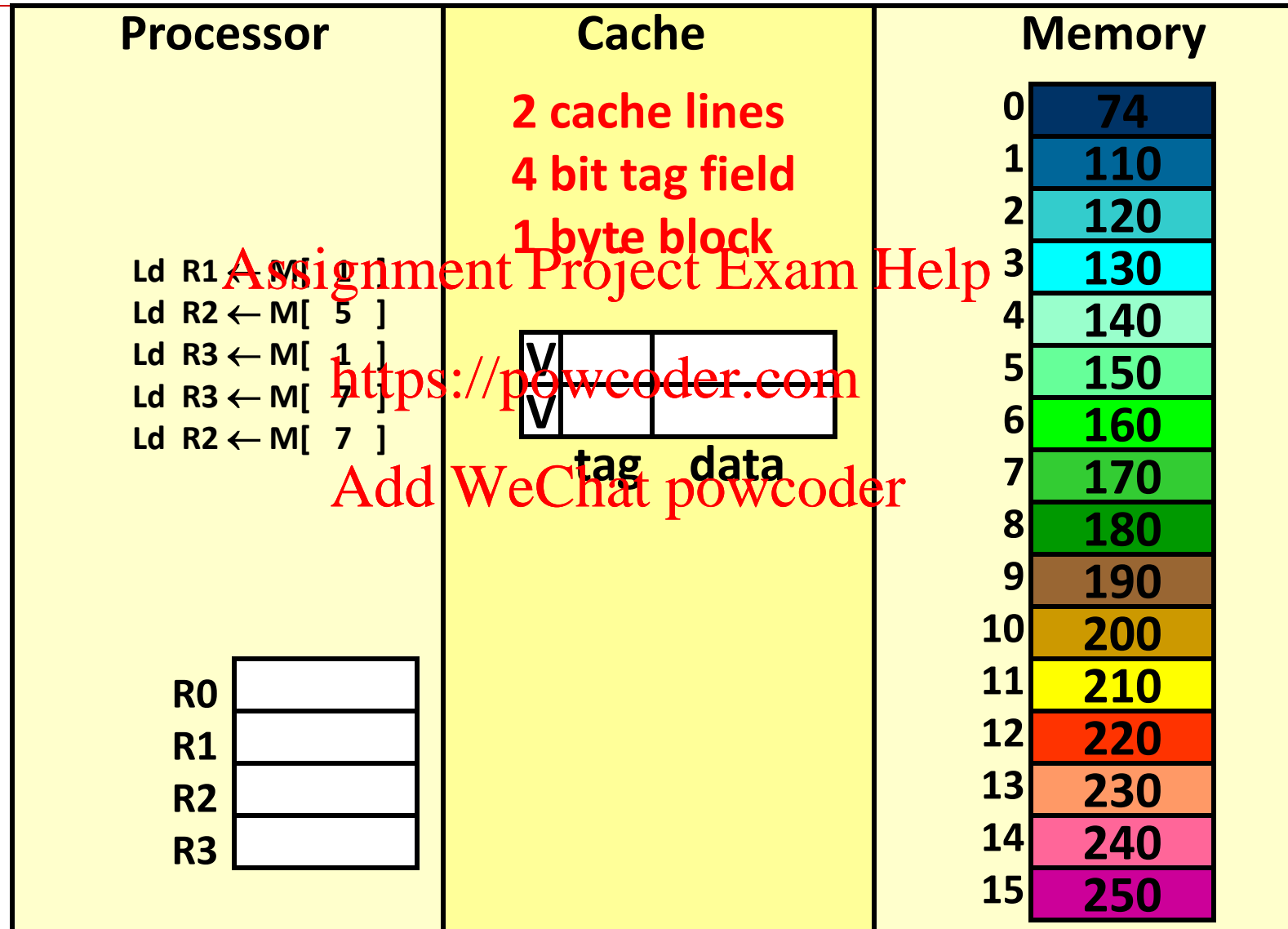
Because the re-access probability falls over time as a cache line isn't referenced, the LRU line is least likely to be re-referenced

Assignment Project Exam Help

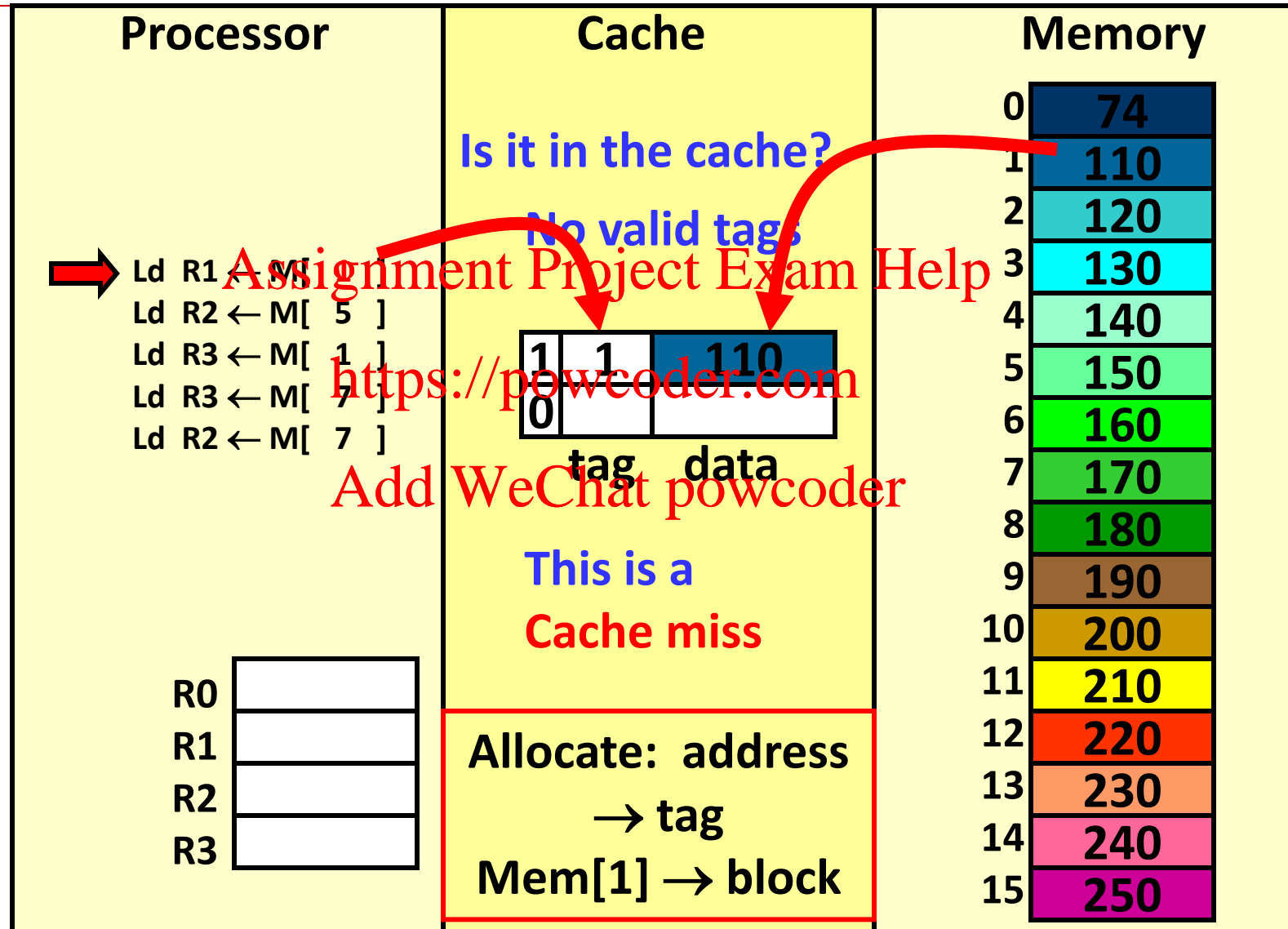
<https://powcoder.com>
Part 4: Cache Illustration

Add WeChat powcoder

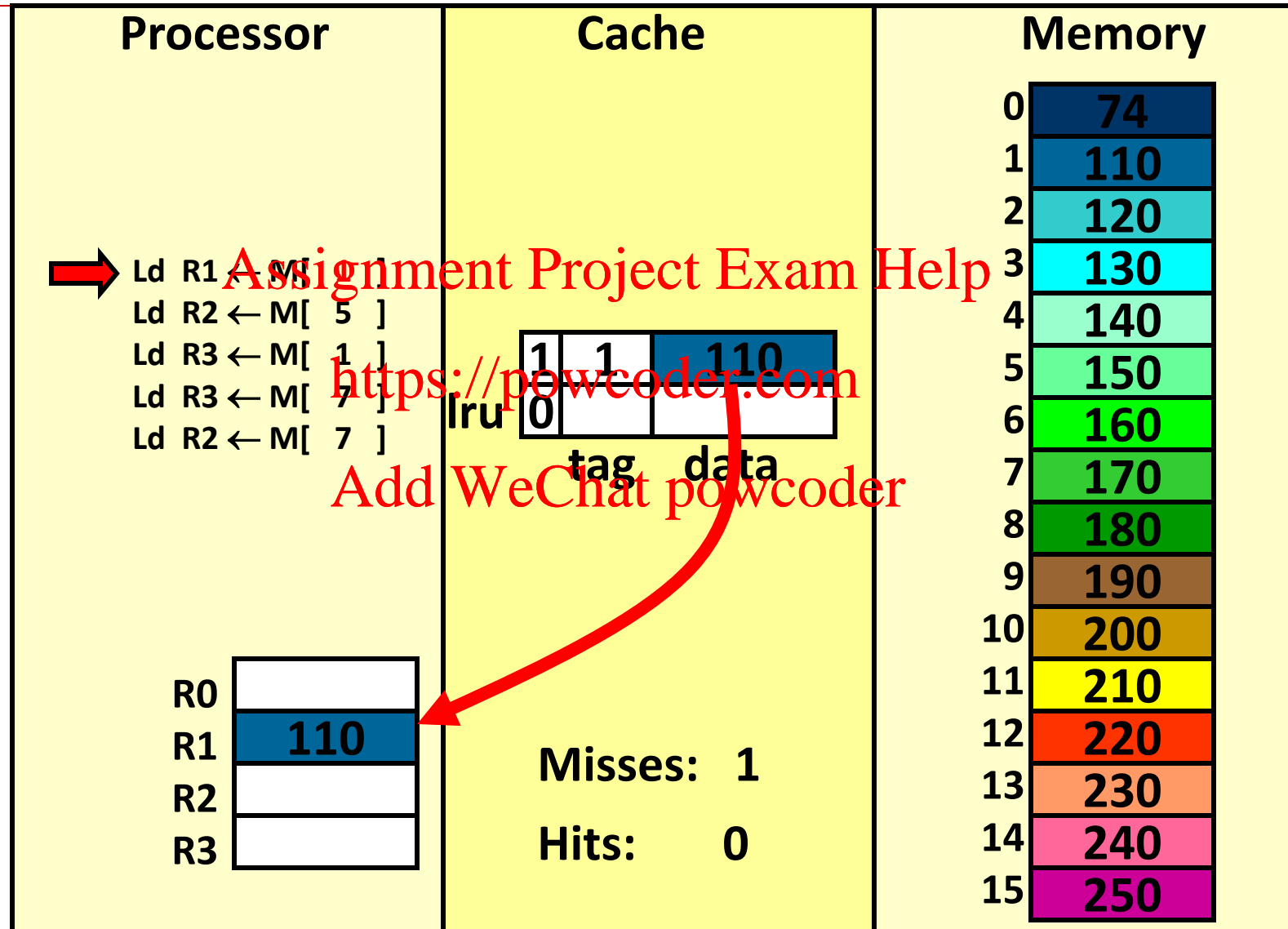
A Very Simple Memory System



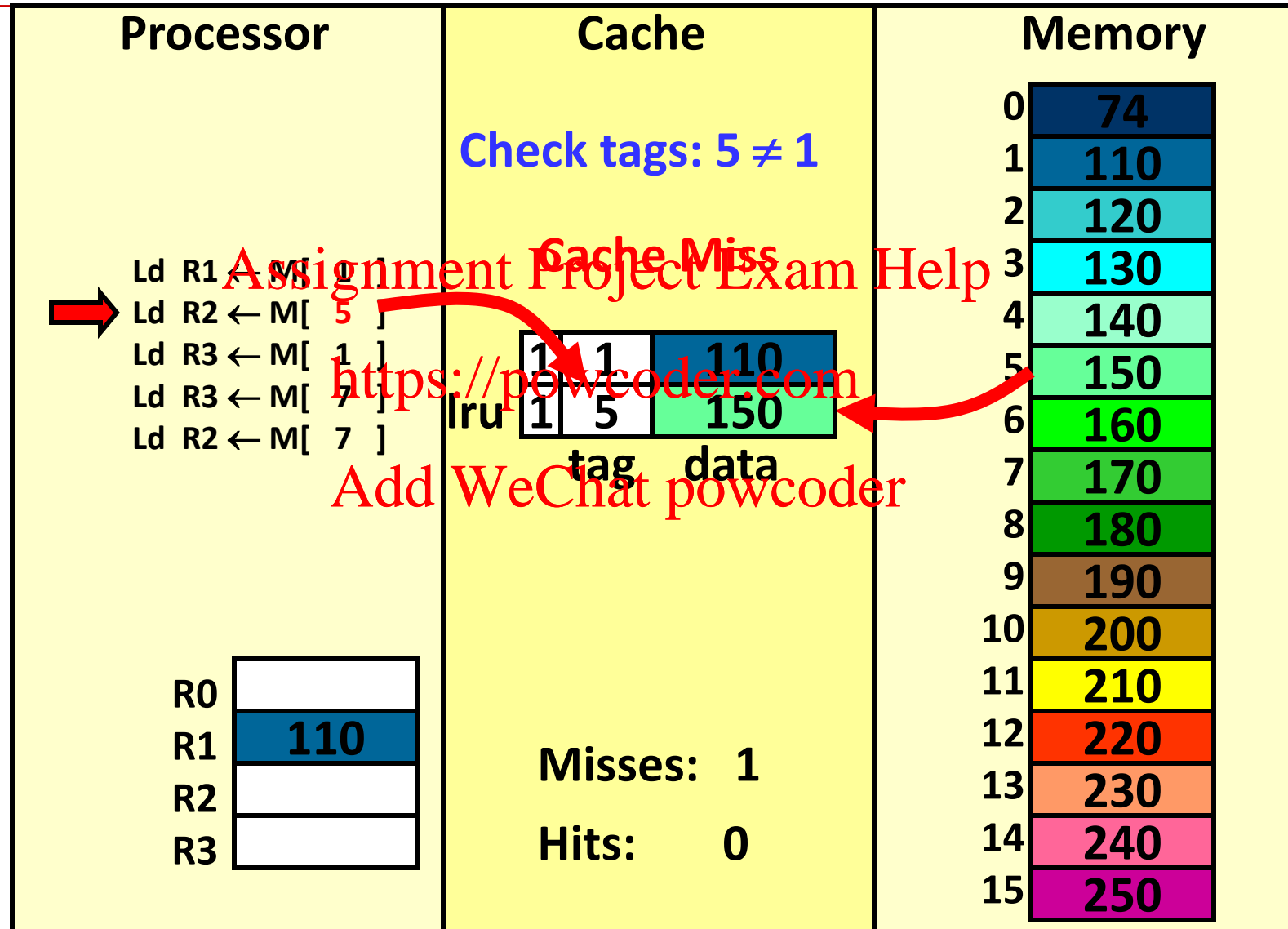
A Very Simple Memory System



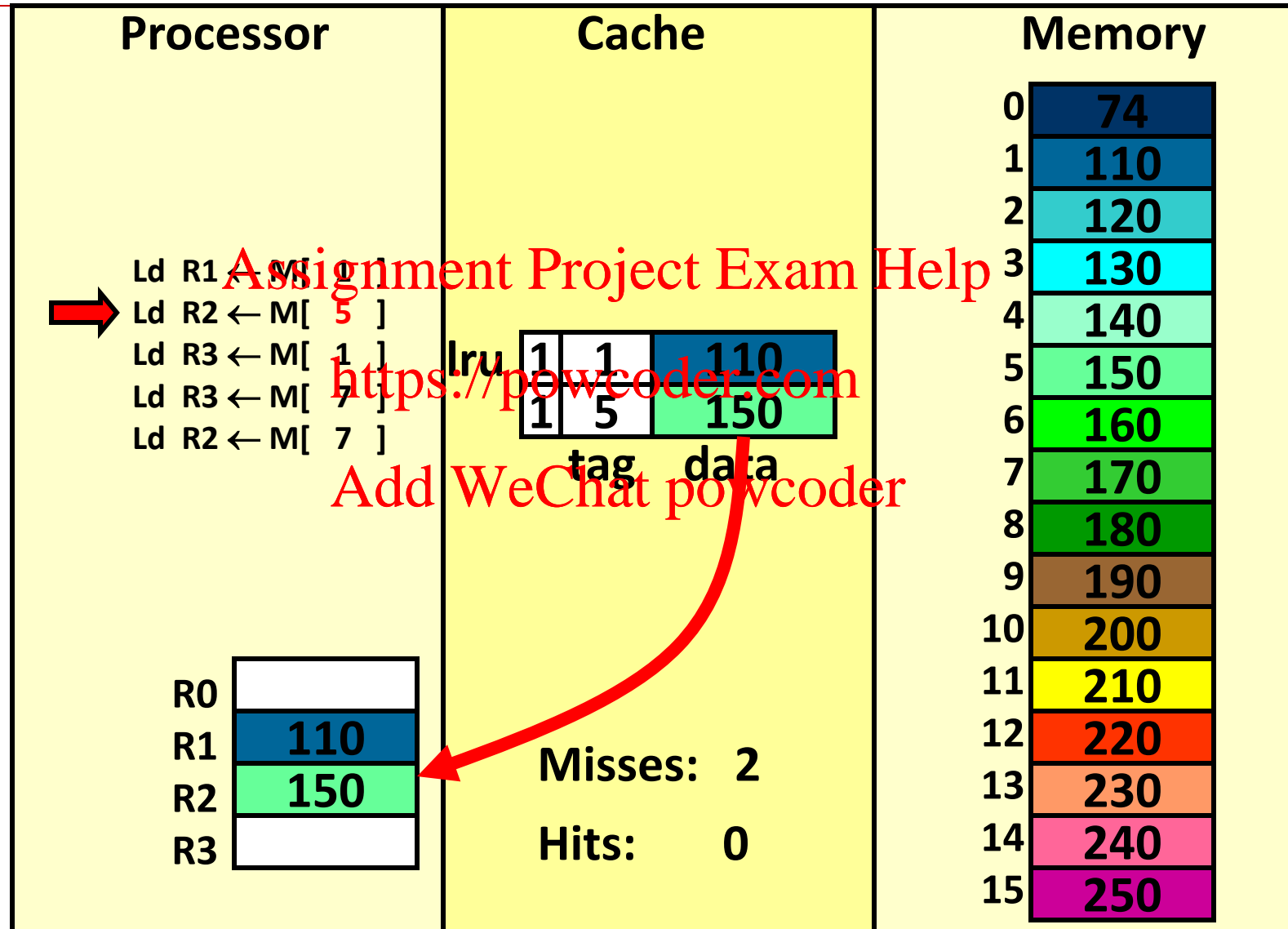
A Very Simple Memory System



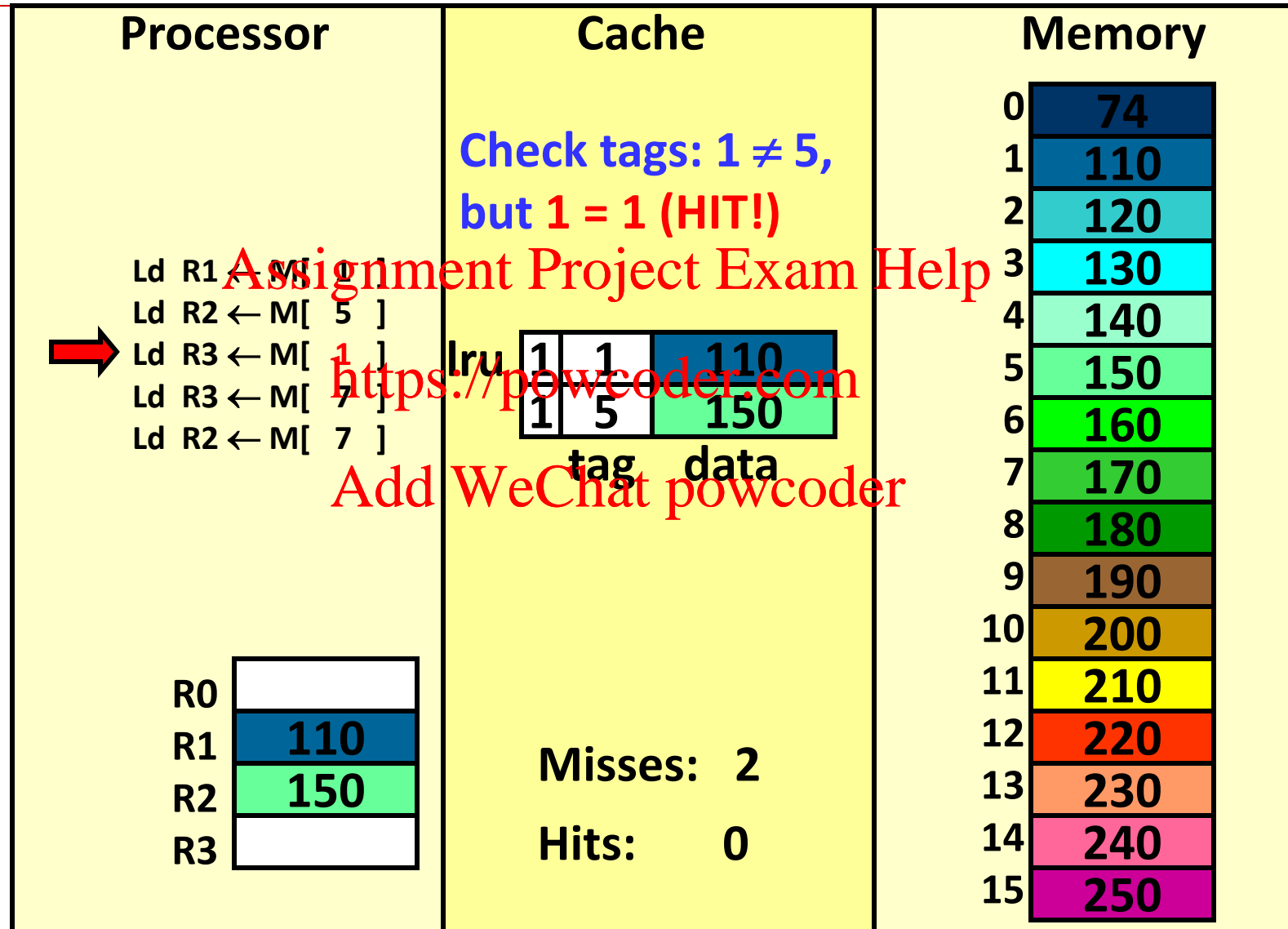
A Very Simple Memory System



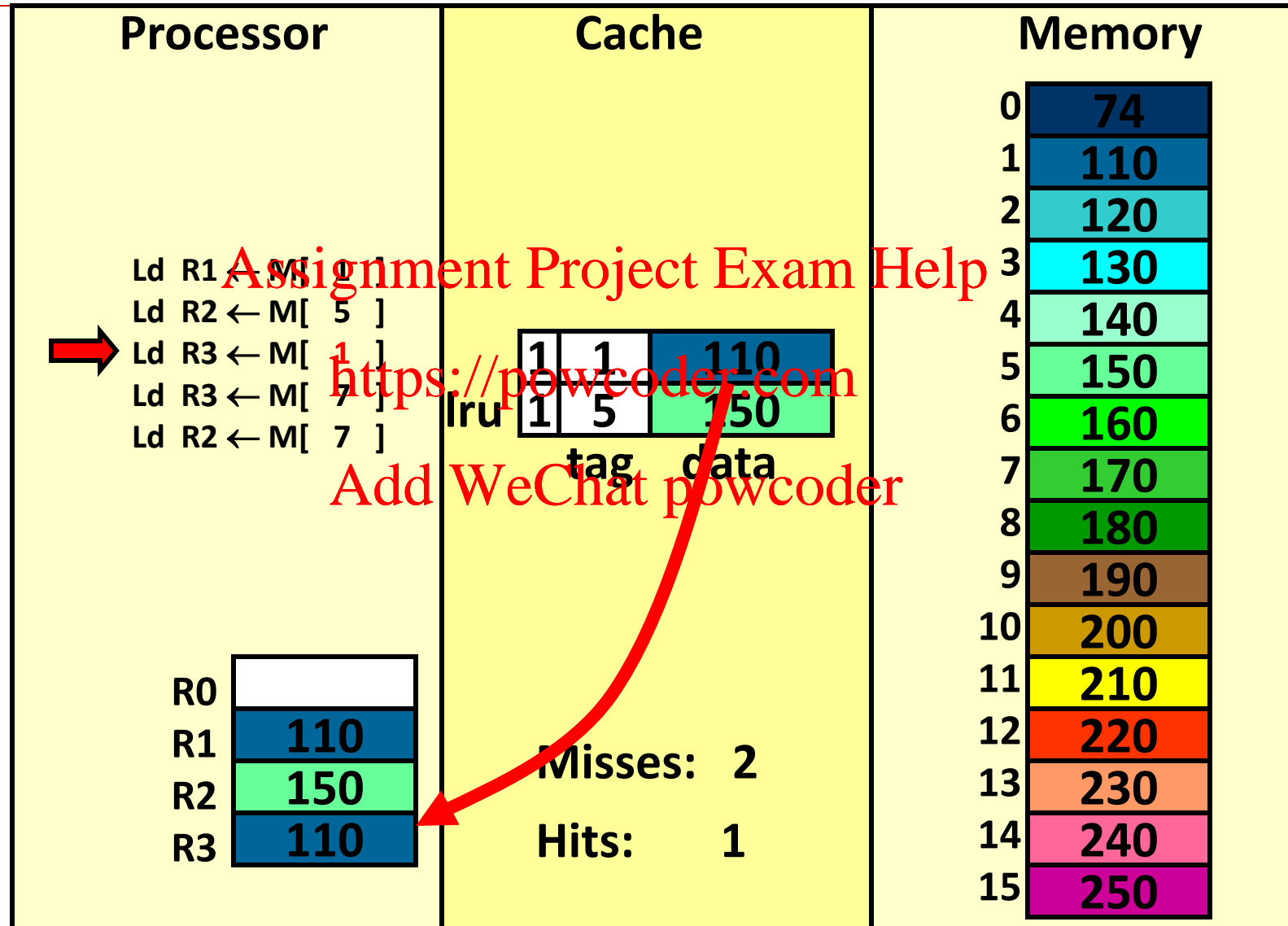
A Very Simple Memory System



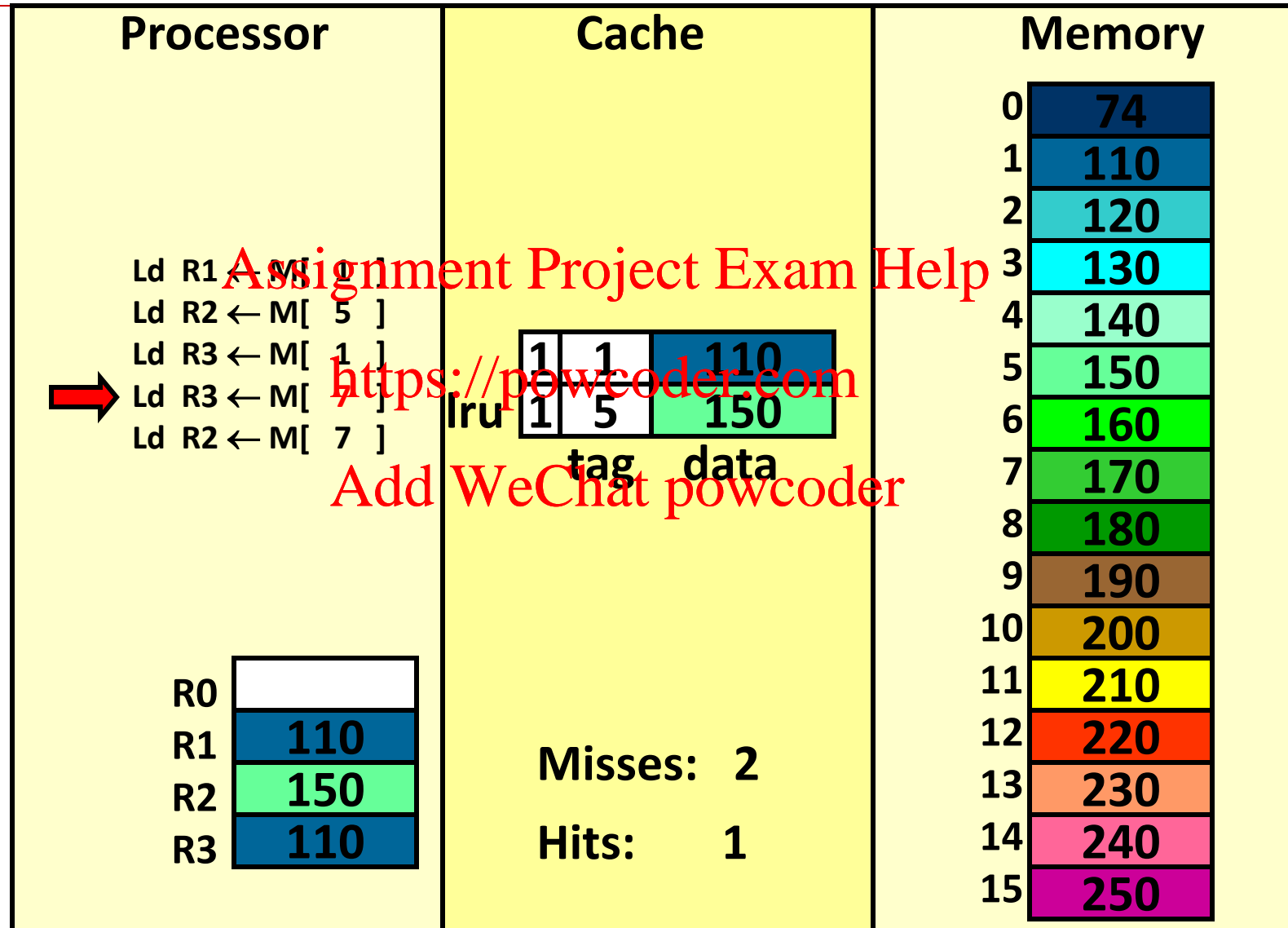
A Very Simple Memory System



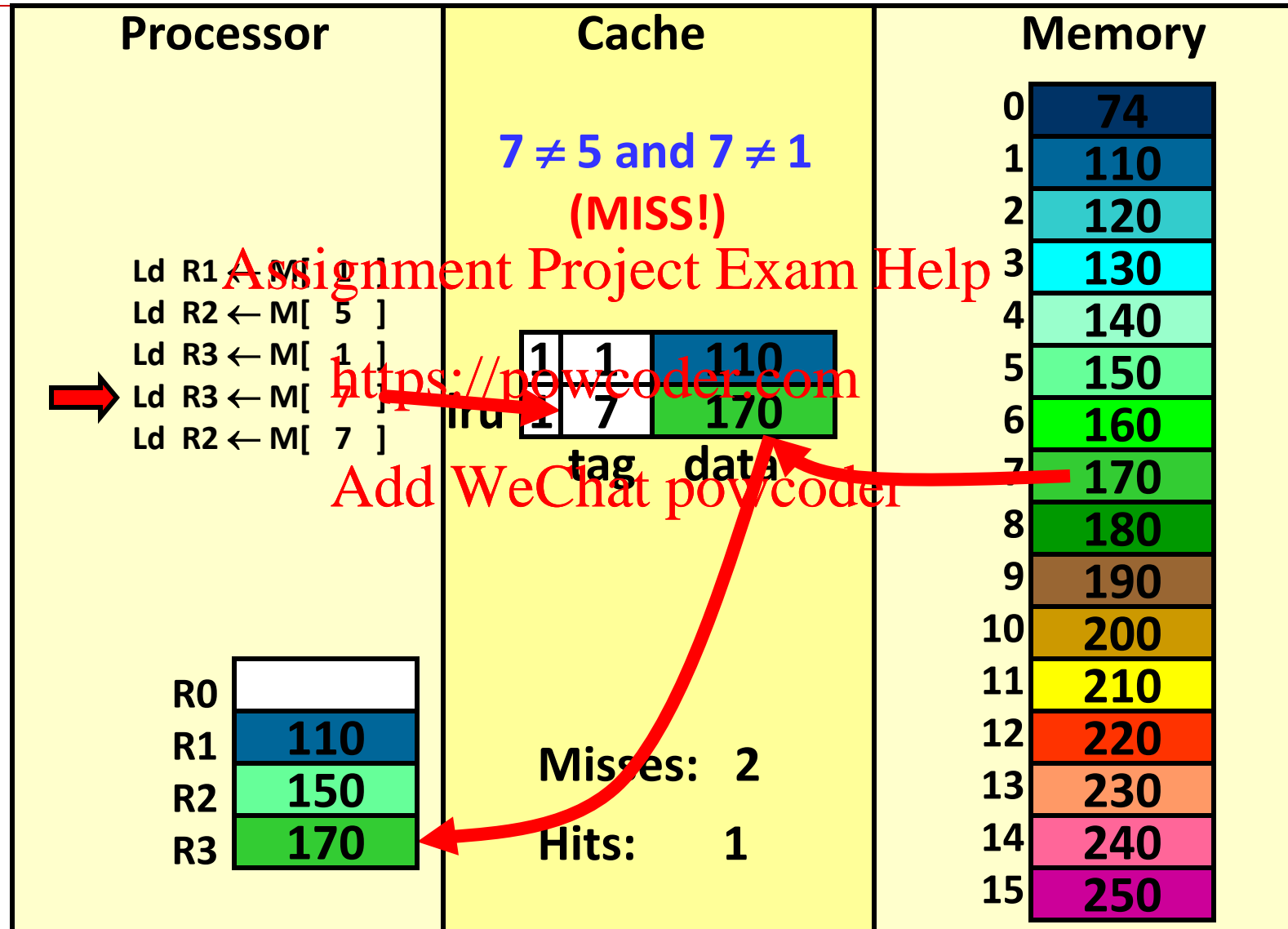
A Very Simple Memory System



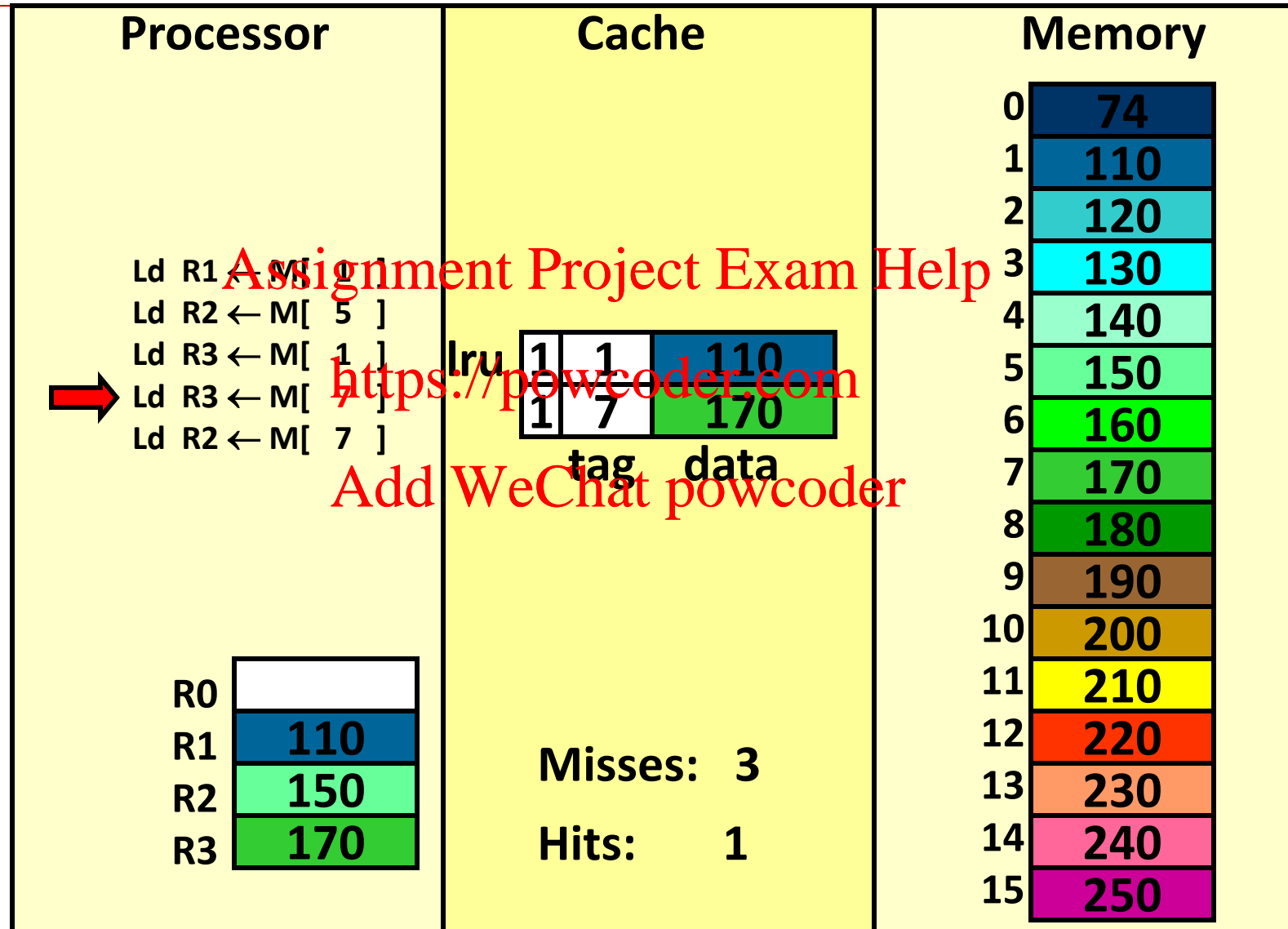
A Very Simple Memory System



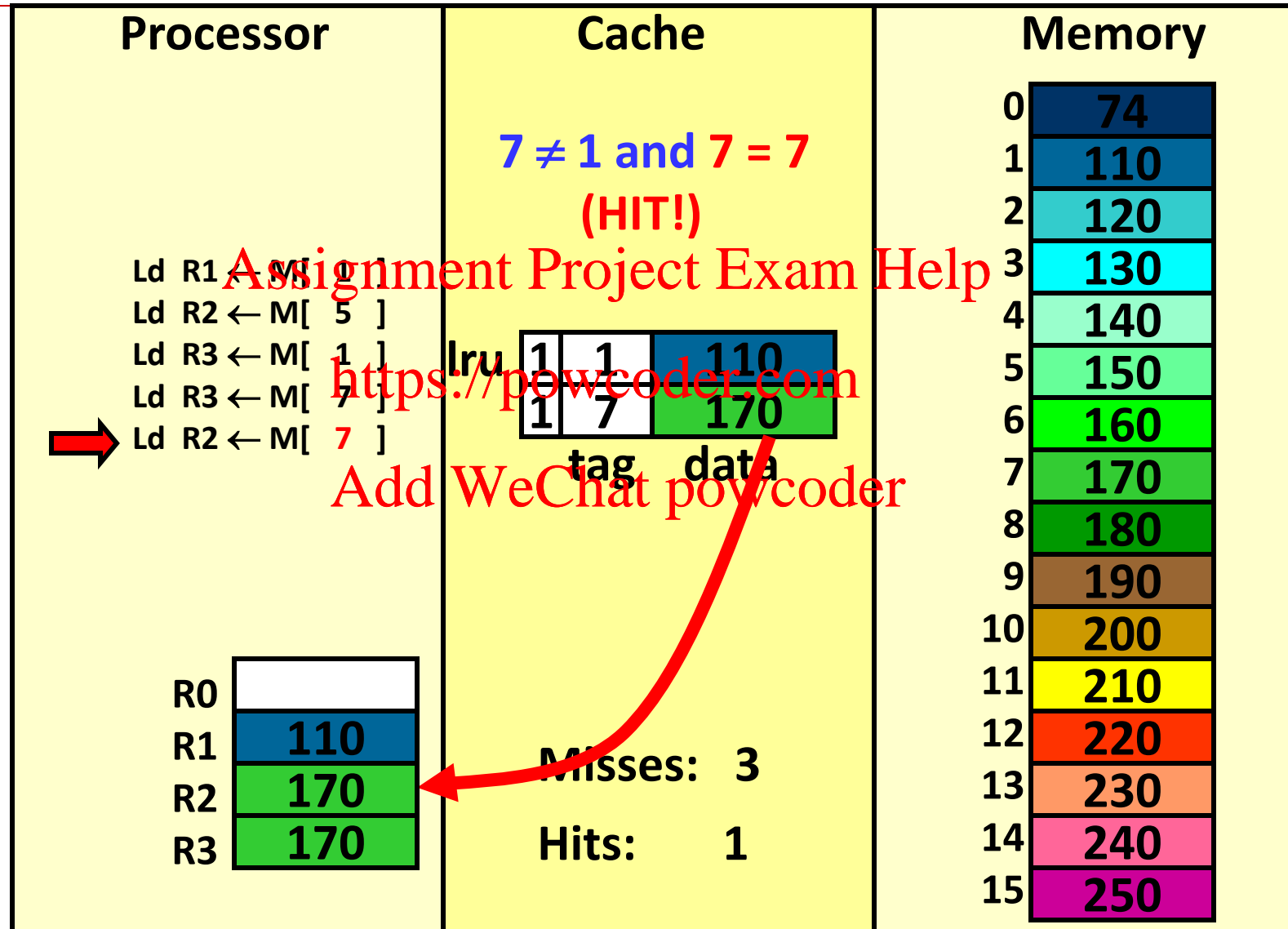
A Very Simple Memory System



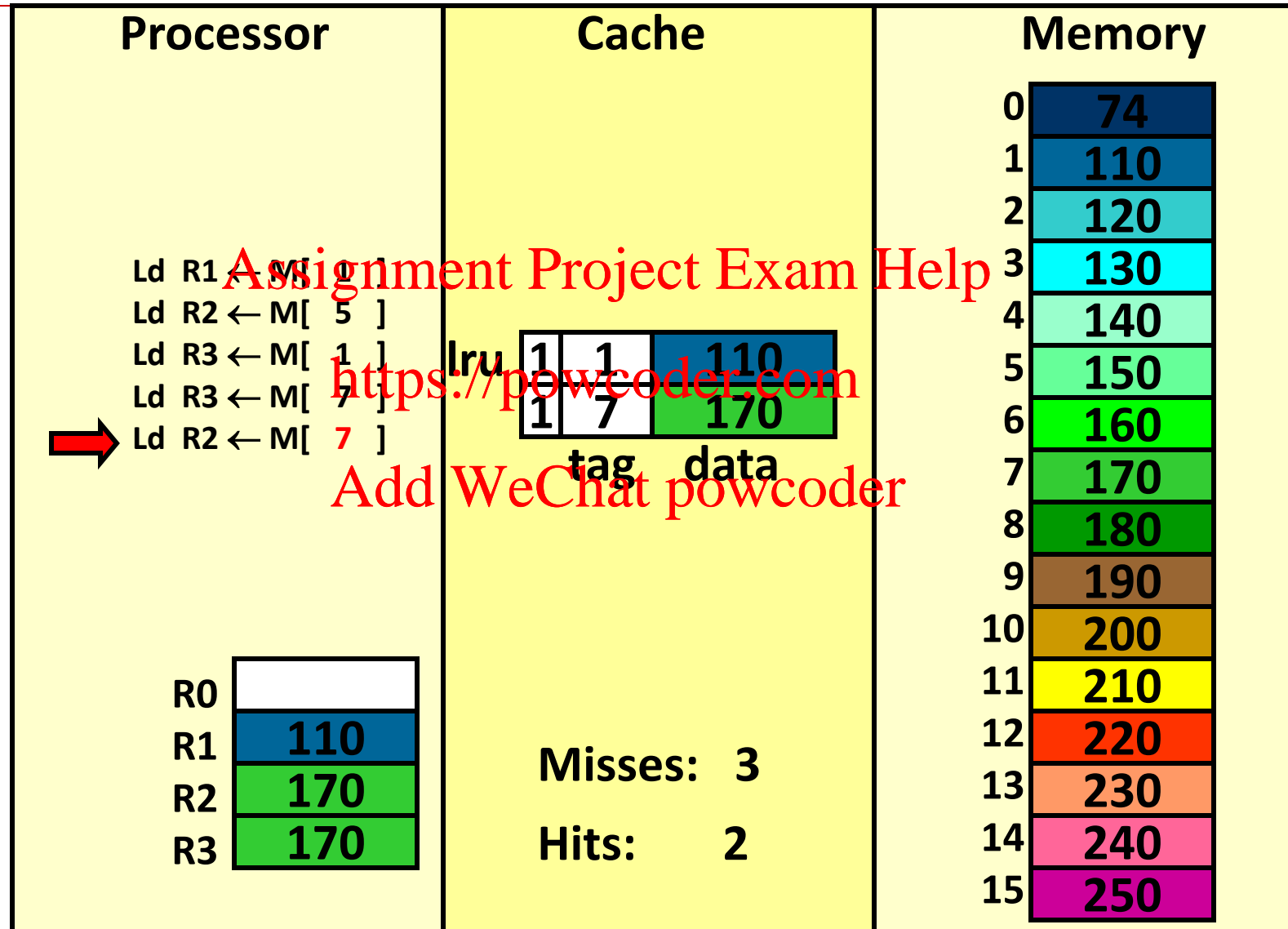
A Very Simple Memory System



A Very Simple Memory System



A Very Simple Memory System



Assignment Project Exam Help

<https://powcoder.com>
Part 5: Cache Performance and Area Overhead

Add WeChat powcoder

Calculating Average Memory Access Time (AMAT)

$$\text{AMAT} = \text{cache latency} \times \text{hit rate} + \text{memory latency} \times \text{miss rate}$$

Simple cache example: 3 misses, 2 hits

Assignment Project Exam Help

Assume following latencies:

Cache: 1 cycle

Memory: 15 cycles (assume it includes time to determine cache hit/miss)

<https://powcoder.com>

Add WeChat powcoder

AMAT for our example cache

$$= 1 \text{ cycle} \times (2/5) + 15 \times (3/5) = 9.4 \text{ cycles per reference}$$

AMAT: Example Problem

Assume the following latencies:

Cache	1	cycle
Main memory	100	cycles
Disk	10,000	cycles

Assignment Project Exam Help

Assume main memory latency (100 cycles) includes time to determine hit/miss in cache.

<https://powcoder.com>

Assume main memory is accessed on all cache misses, and that disk latency does **not** include time to determine hit/miss in cache.

Add WeChat powcoder

Assume a program with these characteristics:

- 100 memory references

- 90% of the cache accesses are hits

- 80% of the accesses to main memory are hits

What is the average memory access time (AMAT)?

$$0.9 * 1 + 0.1 * (100 + 0.2 * 10000) = 210.9$$

Reducing Average Memory Access Time

Assignment Project Exam Help

Reduce latency of cache, main memory, disk and/or
<https://powcoder.com>

Increase hit rate of cache and main memory
Add WeChat powder

Calculating Area Cost

How much does our example cache 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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

What is the area cost if a 32-bit address is used?

Next lecture: How can we reduce the area cost?

Have a small address.

Impractical, and caches are supposed to be micro-architectural

Assignment Project Exam Help
Solution: Cache bigger units of data larger than bytes

Each block has a single tag, and blocks can be whatever size we choose.

To Be Continued...

Add WeChat powcoder