

17. Cache organization: The basics

EECS 370 – Introduction to Computer Organization – Winter 2023

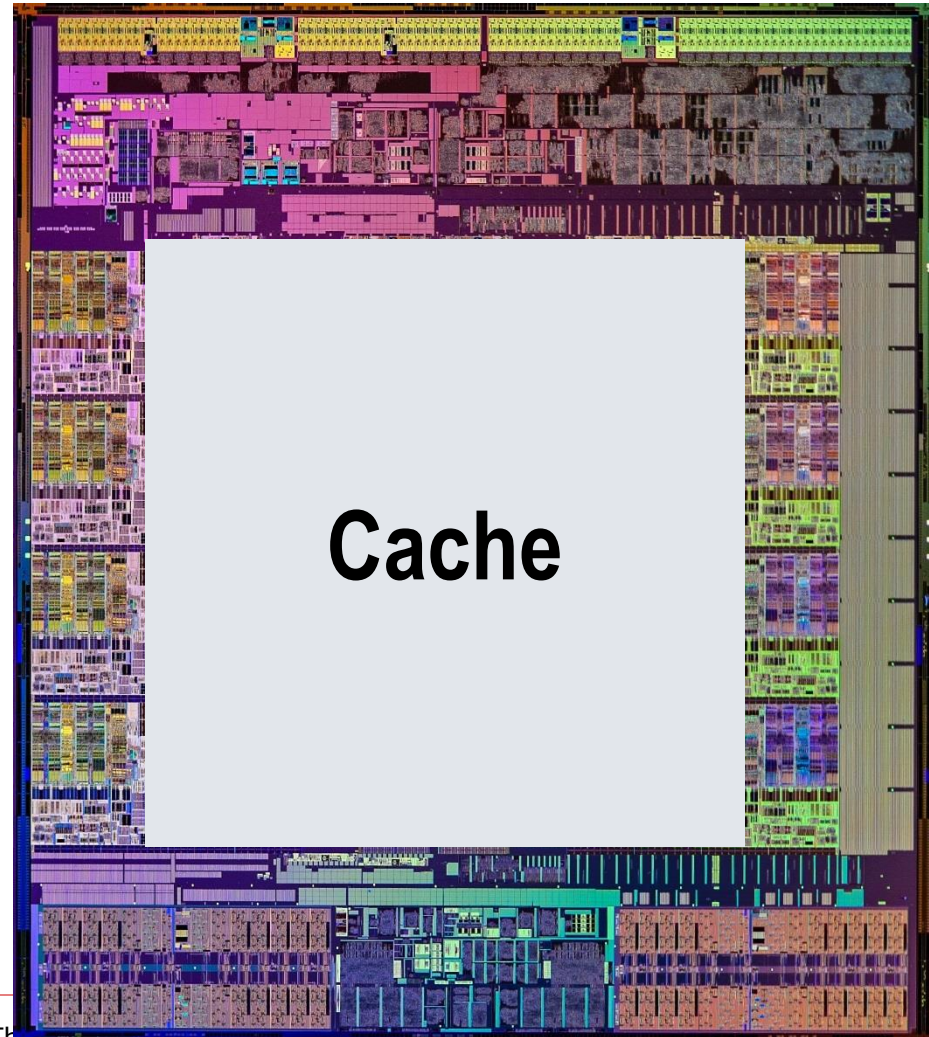
**EECS Department
University of Michigan in Ann Arbor, USA**

Checking in and stuff upcoming

- ❑ Thursday 3/16: Project 3 milestone
- ❑ Monday 3/20: Homework 4 due
- ❑ Thursday 3/23: Project 3 due

Today's Lecture: Cache Organization Basics

Caches consume most of a processor's die area



Cache Aware vs 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);
}
```

Memory

- ❑ So far, we have discussed two structures that hold data:
 - Register file (little array of words)
 - Memory (bigger array of words)

- ❑ We have discussed several methods of implementing storage devices:
 - Static memory (made with logic gates)
 - Dynamic memory (transistor and capacitor)
 - ROM, and other ROM-like storage, e.g., flash (floating gate transistors)

Memory Hierarchy

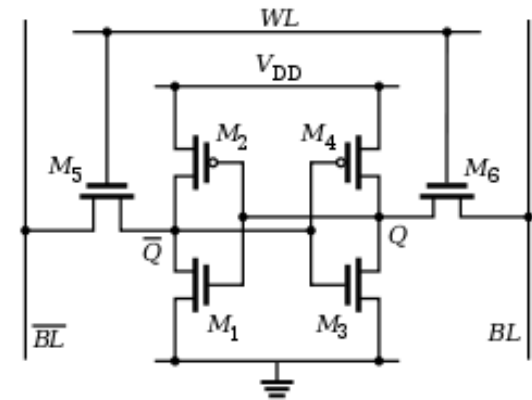
- ❑ We want a lot of memory
 - LC2 can handle 2^{18} bytes of memory
 - MIPS can handle 2^{32} bytes of memory
 - Athlon-64 or EM64T can handle 2^{64} bytes of memory

- ❑ What are our choices?
 - SRAM, DRAM, ROM, disk, tape, DVD?

Option 1: SRAM

❑ Fast: ~2ns access time or faster

- Decoders are big
- Array is big
 - Why?



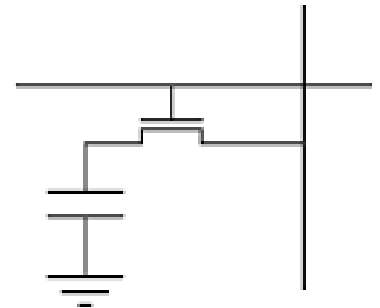
❑ Expensive, high area requirement

- SRAM: ~\$5.0 per megabyte

Option 2: DRAM

- ❑ Slower: ~50ns access time
 - Must stall for dozens of cycles on each memory load

- ❑ Less expensive than SRAM.
 - DRAM costs ~\$0.004 per megabyte



Option 3: Flash

- ❑ Slow: access time varies wildly (ns of ns to μ s)
 - Must stall for dozens of cycles on each memory load
- ❑ Less expensive than DRAM
 - Flash costs ~\$0.0001 per megabyte for large things.
- ❑ Non-volatile

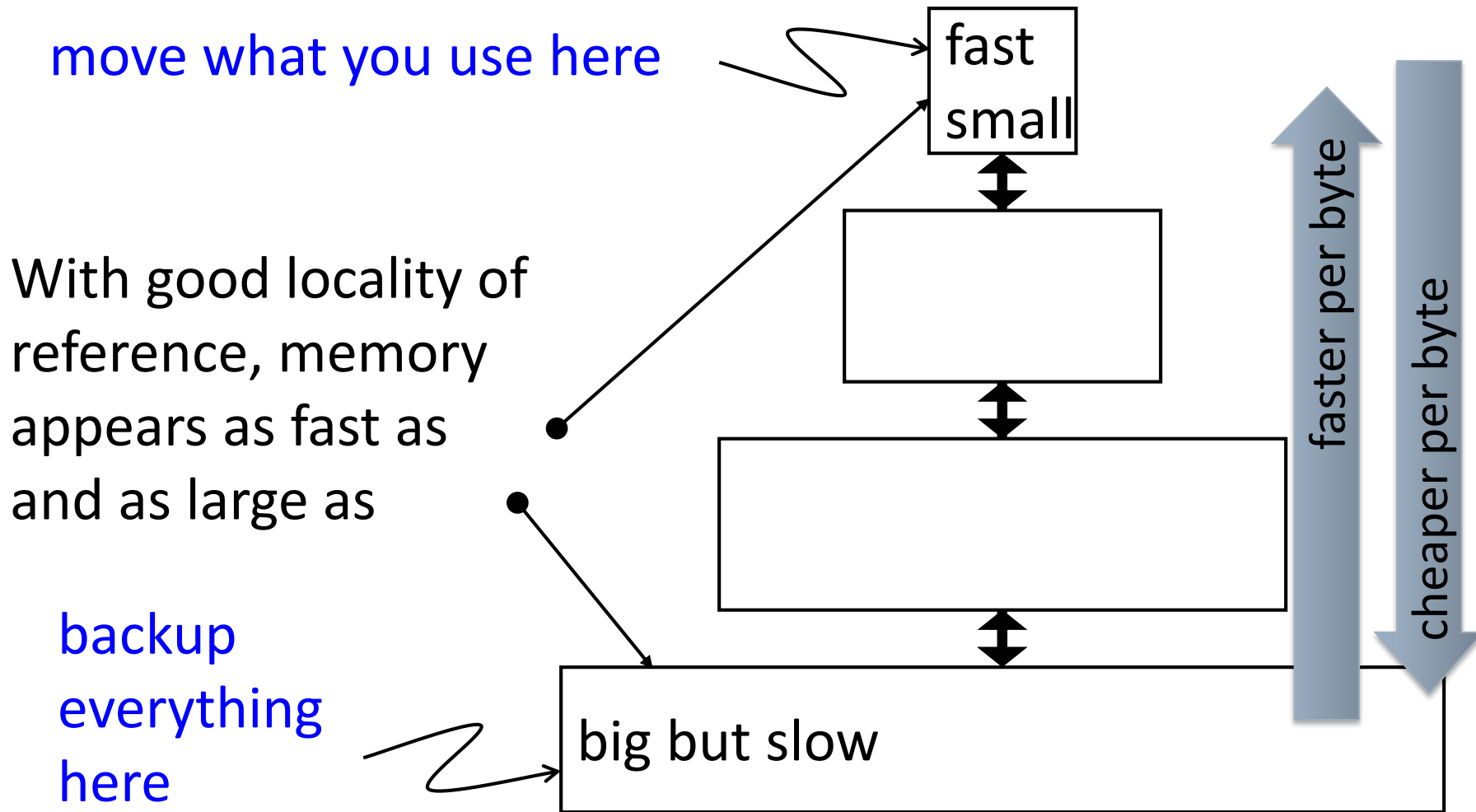
Option 4: Disks

- ❑ Obnoxiously slow: 3,000,000ns access time
 - We could have stopped with the Intel 4004
- ❑ Cheap
 - Disk storage costs ~\$0.00002 per megabyte
- ❑ Non-volatile

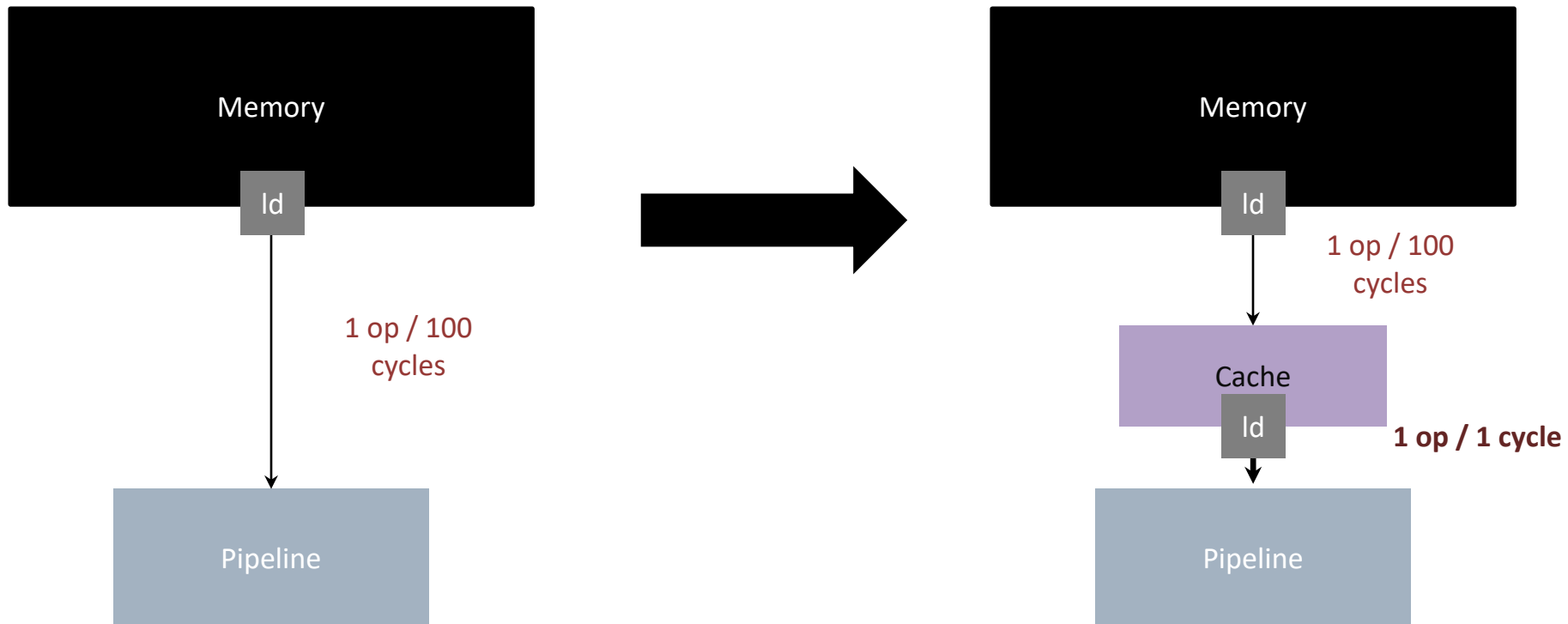
Memory Hierarchy Goals

- ❑ **Fast:** Ideally run at processor clock speed
 - 1 ns access
- ❑ **Cheap:** Ideally free
 - Not more expensive than rest of system
- ❑ Options DRAM, flash, disks are too slow
- ❑ Option SRAM is too expensive
- ❑ How to get best properties of multiple memory technologies?

Memory Hierarchy Goals



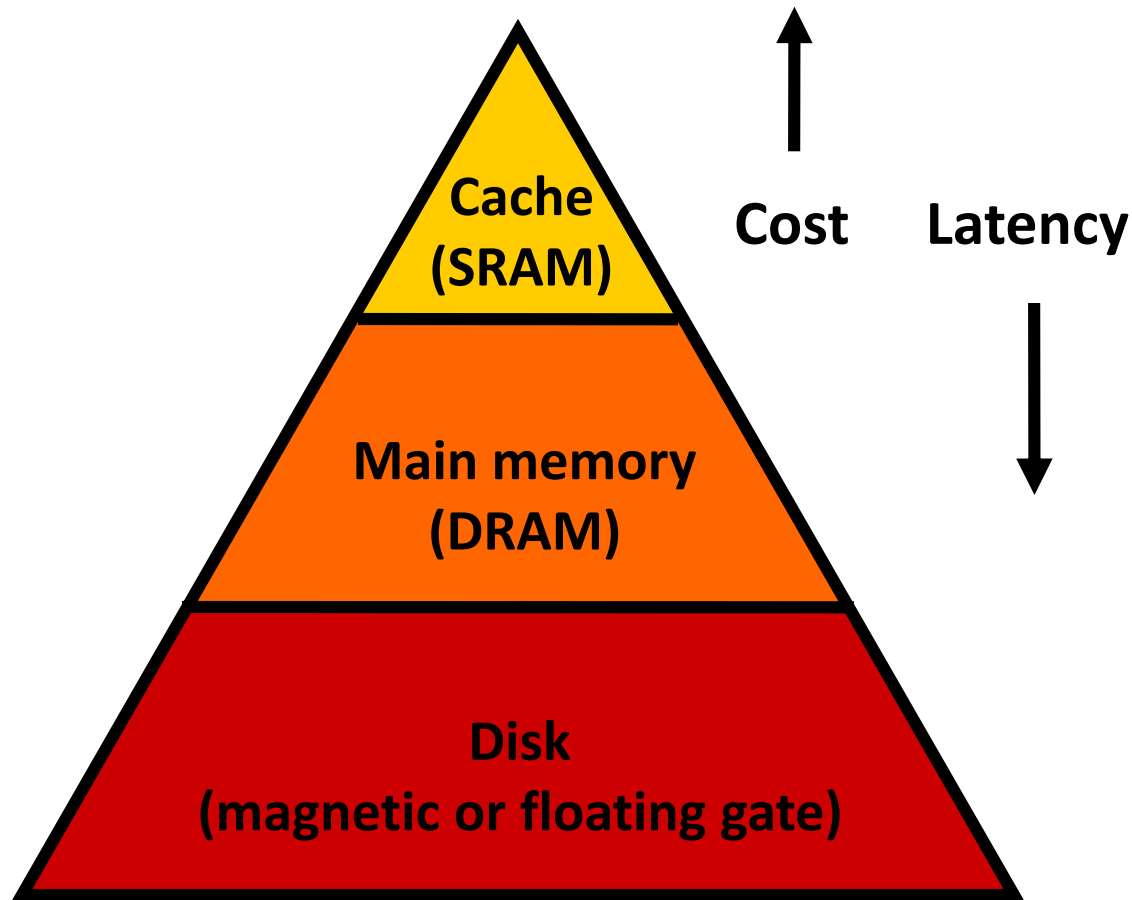
Caches - Overview



Memory Hierarchy

- ❑ Use a small array of SRAM
 - Small so fast and cheap
 - For the cache (hopefully covers most loads and stores)
- ❑ Use a larger amount of DRAM
 - Cheaper than SRAM, faster than flash/disk
 - For the main memory
- ❑ Use a lot of flash and/or disk
 - Non-volatile. Cheap. Big
 - For Virtual memory
- ❑ Don't try to buy 2^{64} bytes of anything
 - Use “virtual memory” to make it look like the entire address range is available
 - A few TB is enough for most desktop machines today, or a smartphone in a few years

Memory hierarchy



Definitions

- ❑ The architectural view of memory is
 - What the machine language (or programmer) sees
 - Just a big array
- ❑ Breaking up the memory system into different pieces – cache, main memory (made up of DRAM) and Disk – is not architectural
 - The machine language doesn't know about it
 - A new implementation may not break it up in the same way

Function of the Cache

- ❑ The cache will hold the data that we think is **most likely** to be **referenced**
 - Because we want to maximize the number of references that are serviced by the cache to minimize the average memory access latency
 - How do we decide what the most likely accessed memory locations are?

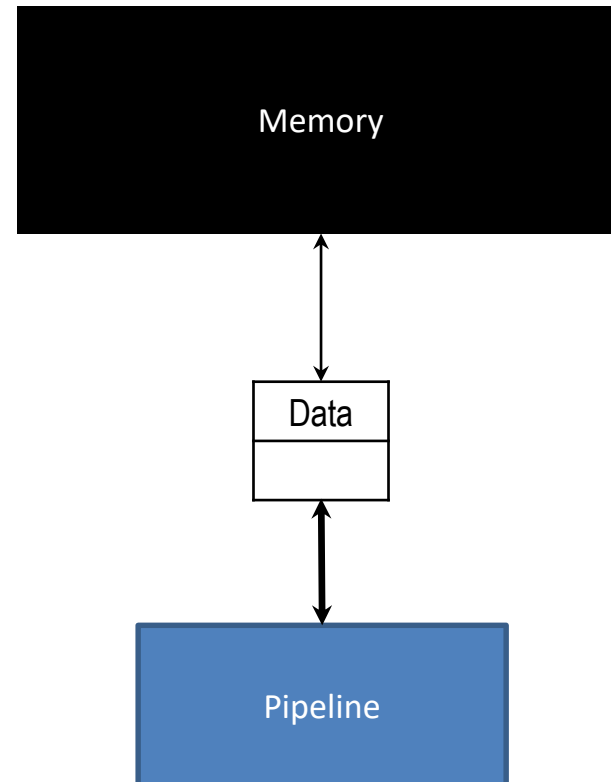
Cache Analogy

❑ Studying books in library

- **Option 1: Every time you switch to another book, return current book to shelf and get new book from shelf**
 - Latency = 5 minutes
- **Option 2: Keep 10 commonly-used books on shelf above desk**
 - Latency = 1 minute
- **Option 3: Keep three books open to appropriate locations on desk**
 - Latency = 10 seconds

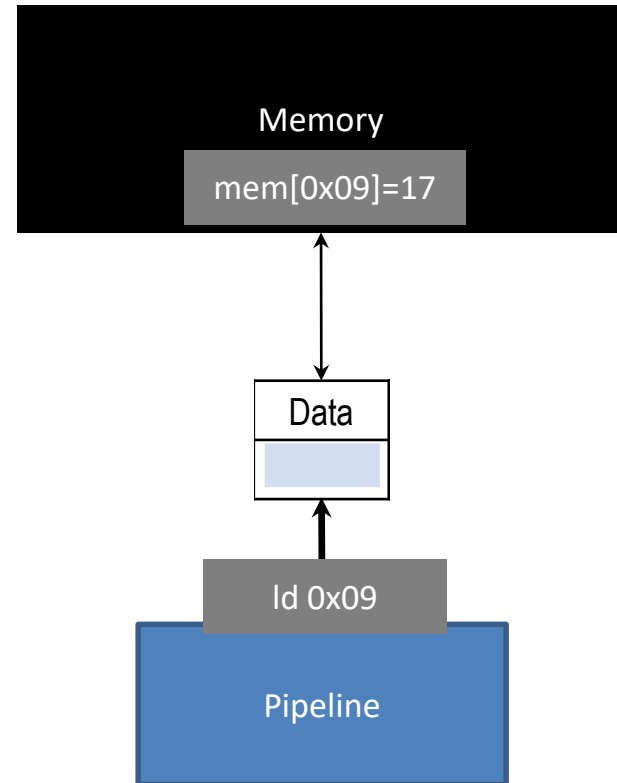
The Simplest Cache

- ❑ Only worry about load instructions for now
 - Word-addressable address space
 - Consists of a single, word-size storage location to remember last loaded value



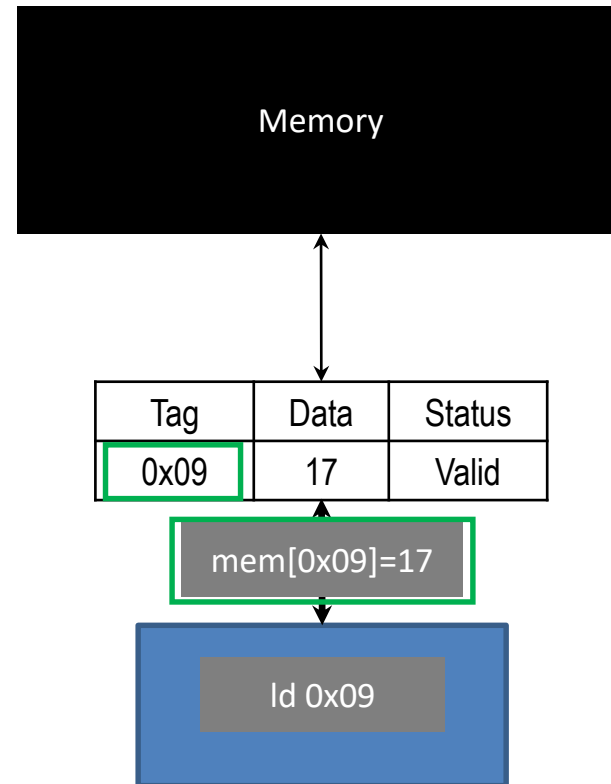
The Simplest Cache

- ❑ Whenever memory returns data, store it in the cache
- ❑ We'll also need to know what address this data corresponds to
 - Store that as “tag”
- ❑ Also include a “valid” status bit



The Simplest Cache

- ❑ Next memory access, first check if the tag matches address
 - Yes? Return cache data
 - No? Go to memory as before



Definitions

- ❑ **Hit:** when data for a memory access is found in the cache
- ❑ **Miss:** when data for a memory access is not found in the cache
- ❑ **Hit/Miss rate:** percentage of memory accesses that hit/miss in the cache

Slightly Less Simple:

More than one thing stored

- ❑ 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**

- ❑ Now **compare reference address and tag** for all cache lines
 - Match? Get the data from the cache block
 - No Match? Get the data from main memory

CAMs: content addressable memories

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

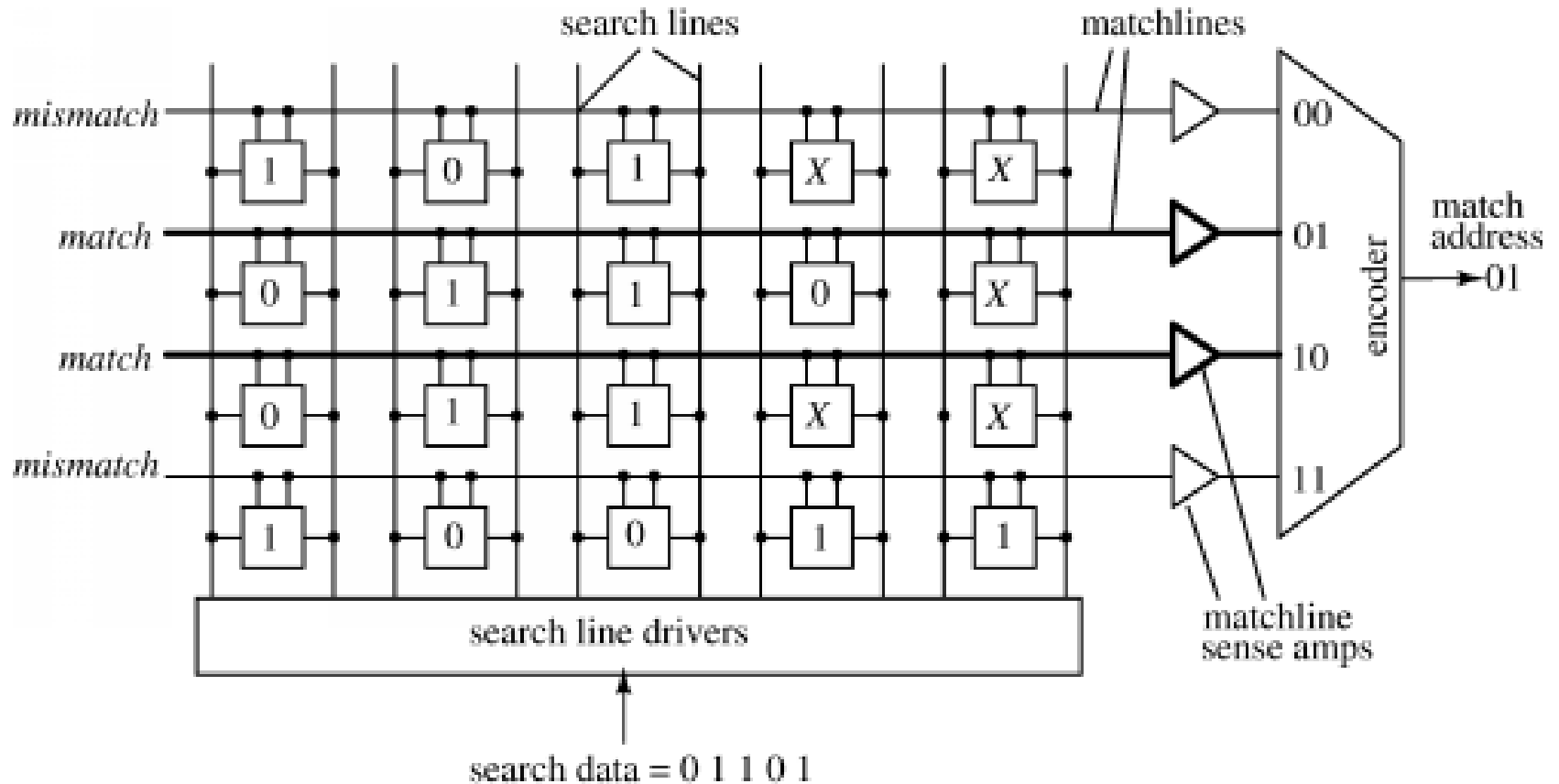
- ❑ Think of memory as a set of data matching a query
 - Instead of an address, we send a key to the memory, asking whether the key exists and, if so, what value it is associated with
 - Memory answers: yes/no (hit/miss for caches) and gives associated value (if there is one)

Operations on CAMs

- ❑ **Search:** the primary way to access a CAM
 - Send data to CAM memory
 - Return “found” or “not found”
 - If found, return location of where it was found or associated value

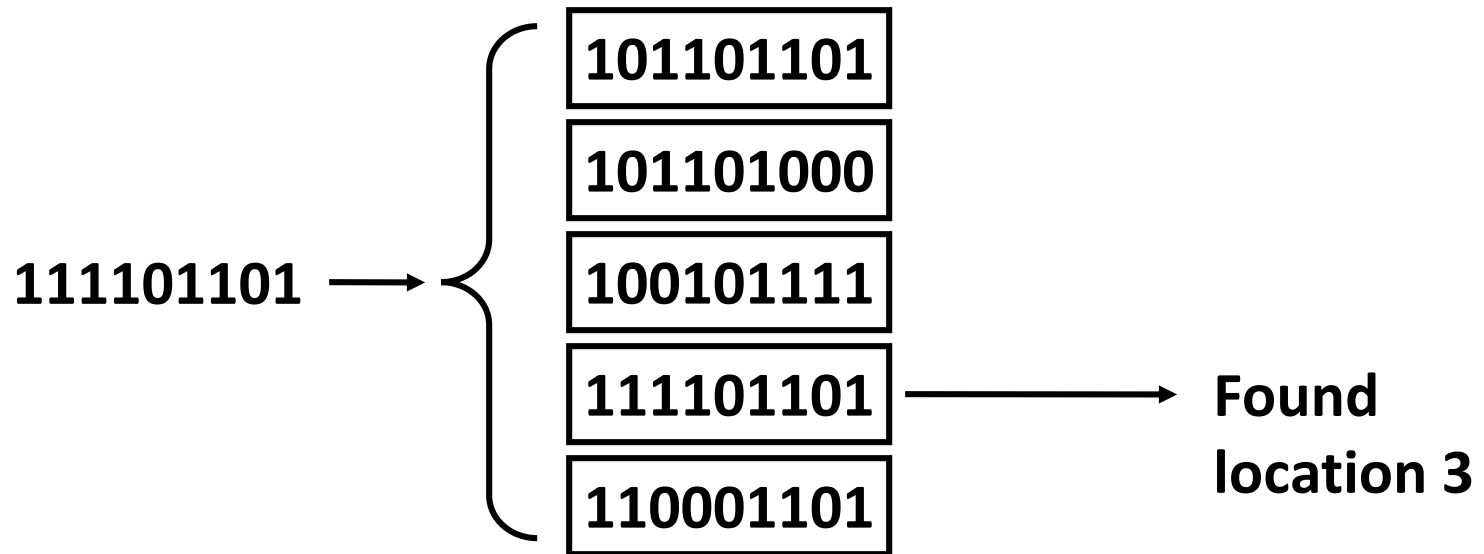
- ❑ **Write:**
 - Send data for CAM to remember
 - Where should it be stored if CAM is full?
 - Replacement policy
 - Replace oldest data in the CAM
 - Replace least recently searched data

CAM = content addressable memory



When used in caches, all tags are fully specified (no X)

CAM example



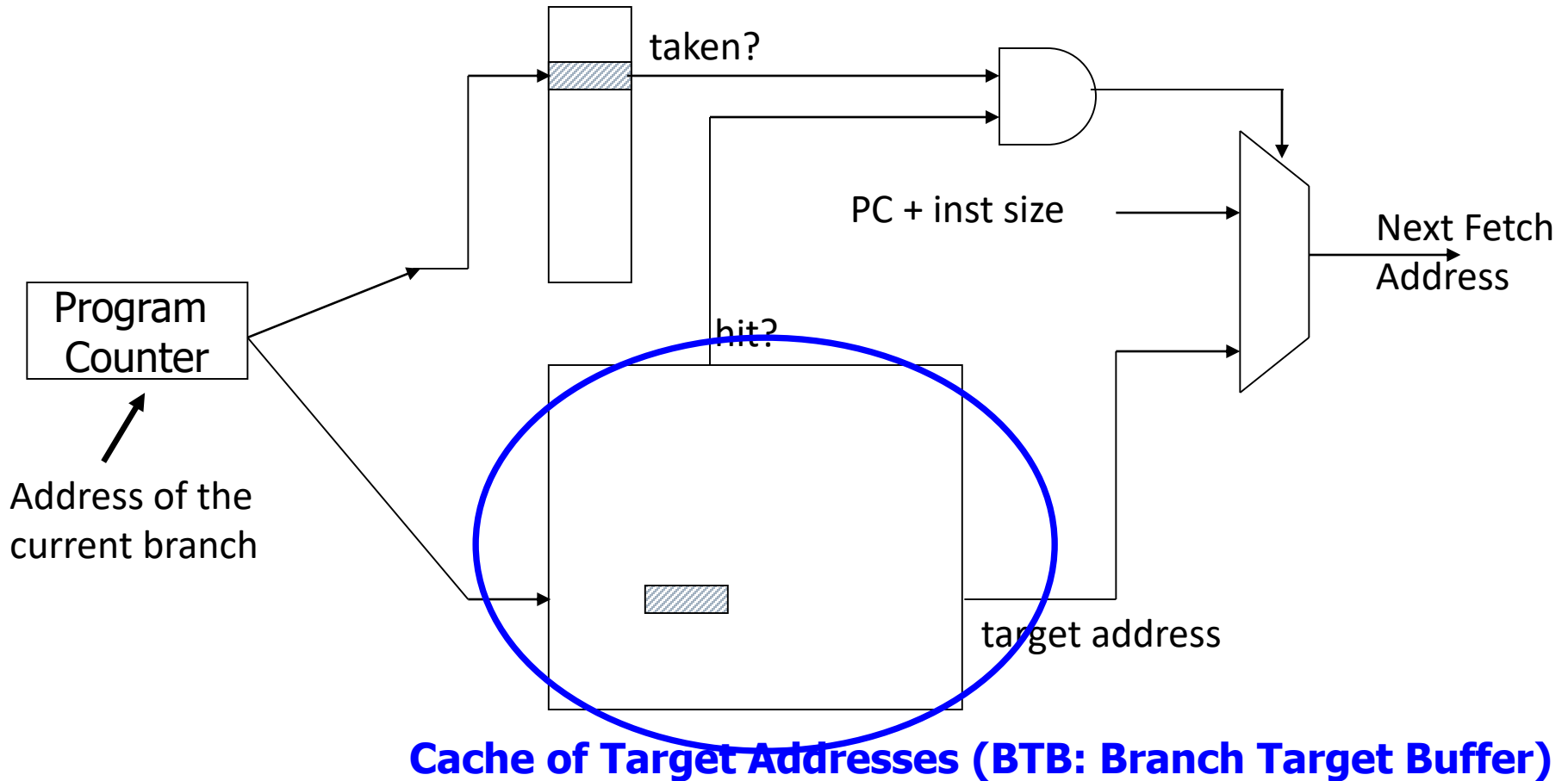
5 storage element CAM array of 9 bits each

Previous use of CAMs

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

Fetch Stage with Branch Prediction

Direction predictor (2-bit counters)



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

Cache Organization

addr	data
addr	data

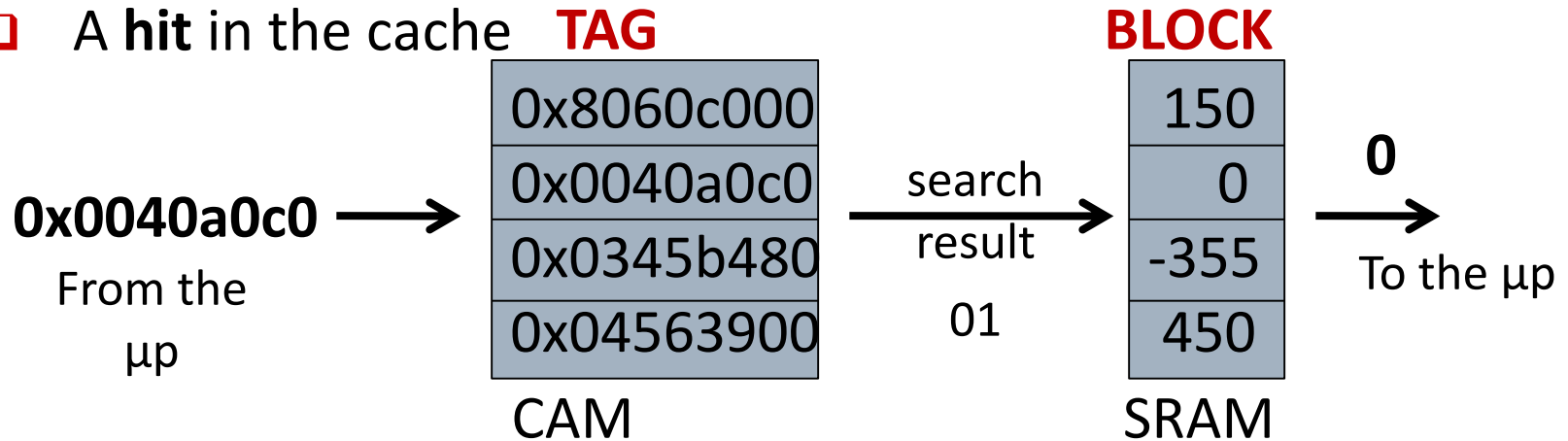
TAG BLOCK

- ❑ 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**

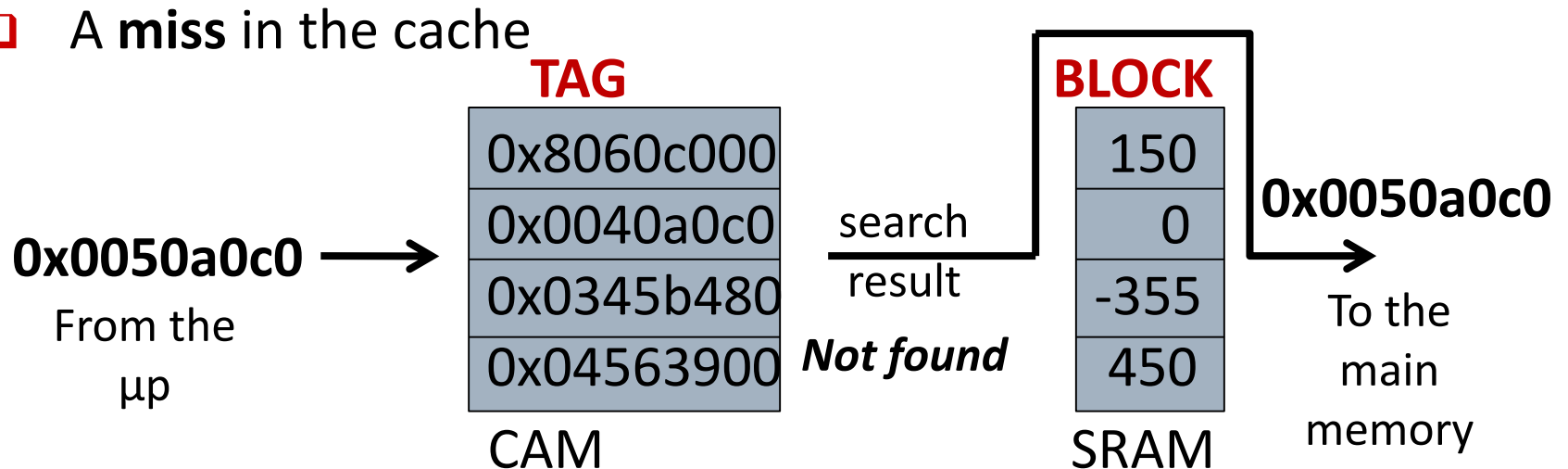
Our goal is to keep the data we think will be accessed in the near future in the cache

Caches: the hardware view

❑ A **hit** in the cache



❑ A **miss** in the cache



Example Problem

Given the following:

Cache has 1 cycle access time

Main memory has 100 cycle access time

Disk has 10,000 cycles access time

What is the average access time for 100 memory references if 90% of the cache accesses are hits and 80% of the accesses to main memory are hits? **Assume main memory access time (100 cycle) includes tag array access to determine hit/miss.**

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

Cache Operation

- ❑ Every cache **miss** will get the data from memory and **ALLOCATE** a cache line to put the data in
 - Just like any CAM write
- ❑ Which line should be allocated?
 - Random? OK, but hard to grade test questions
 - Better than random? How?

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
 - Hint: a crystal ball will come in handy in solving this problem...
- ❑ Why would we care?

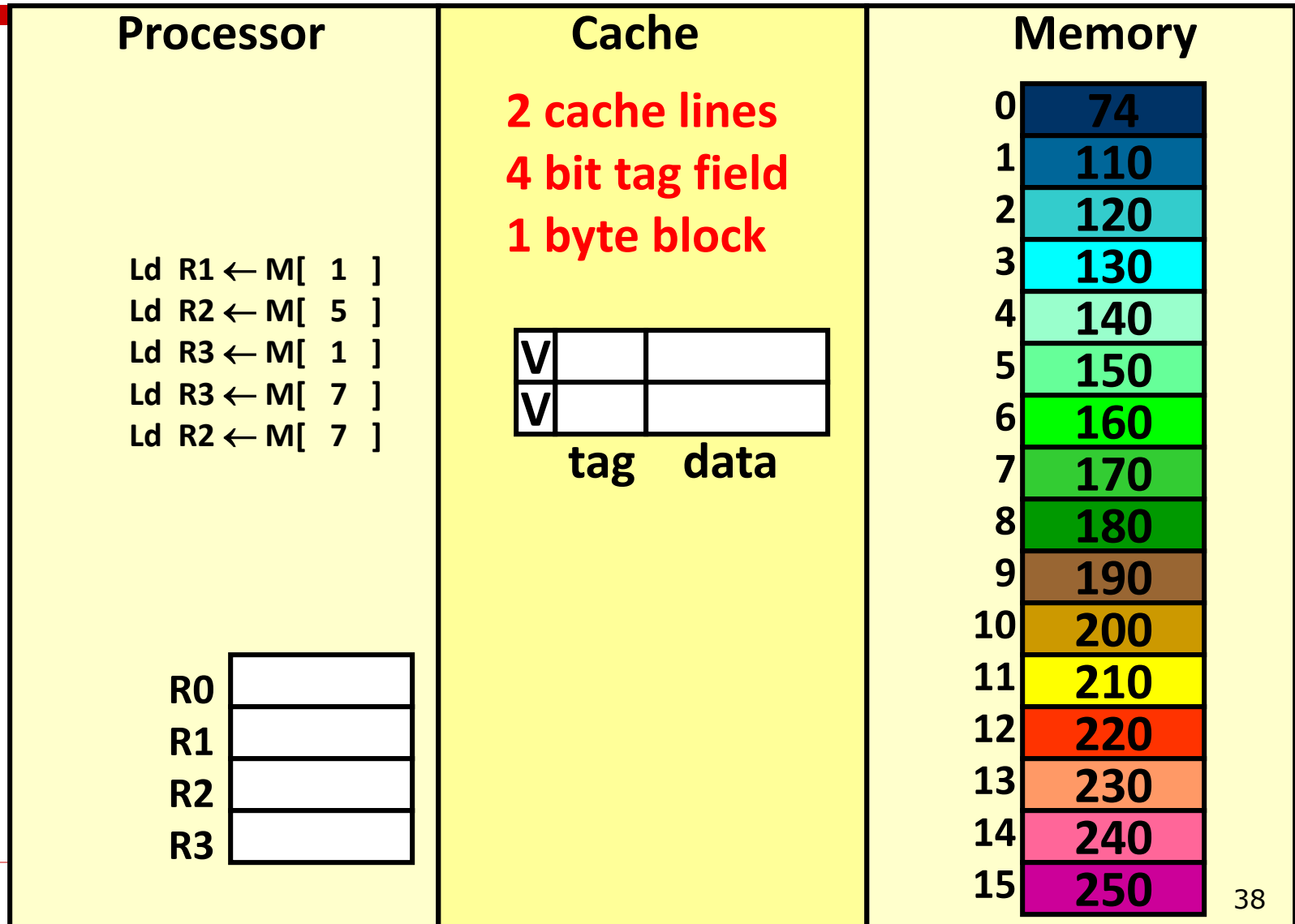
Picking the Most Likely Addresses

- ❑ What is the probability of accessing a random memory location?
 - With no information, it is just as likely as any other address
- ❑ But programs are not random
 - They tend to use the same memory locations over and over
 - We can use this to pick the most referenced locations to put into the cache

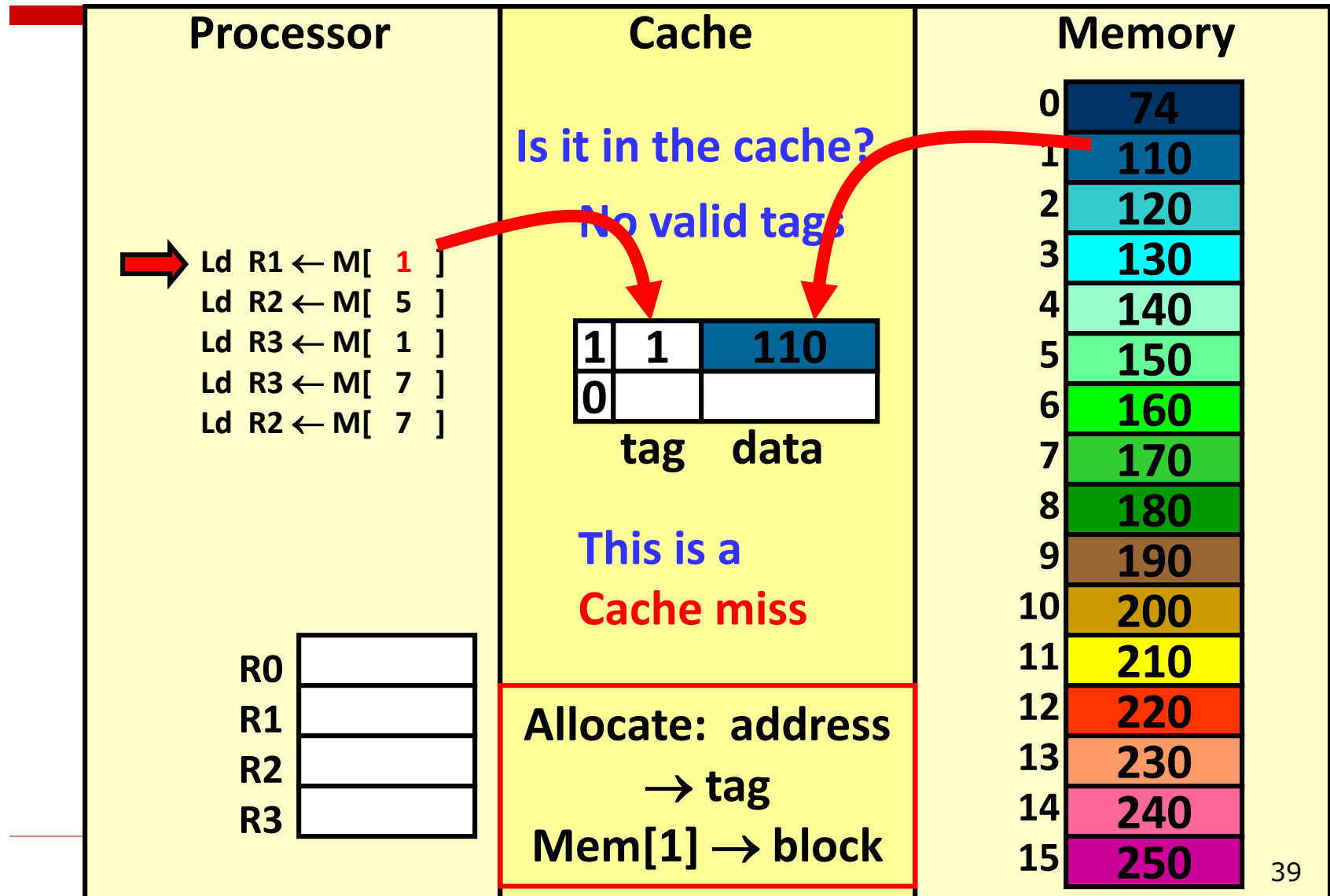
Temporal Locality

- ❑ The principle of **temporal locality** in program references says that if you access a memory location (e.g., 1000) you will be more likely to re-access that location than you will be to reference some other random location
- ❑ Temporal locality says any miss data should be placed into the cache
 - It is the most recent reference location
- ❑ Temporal locality says that the 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

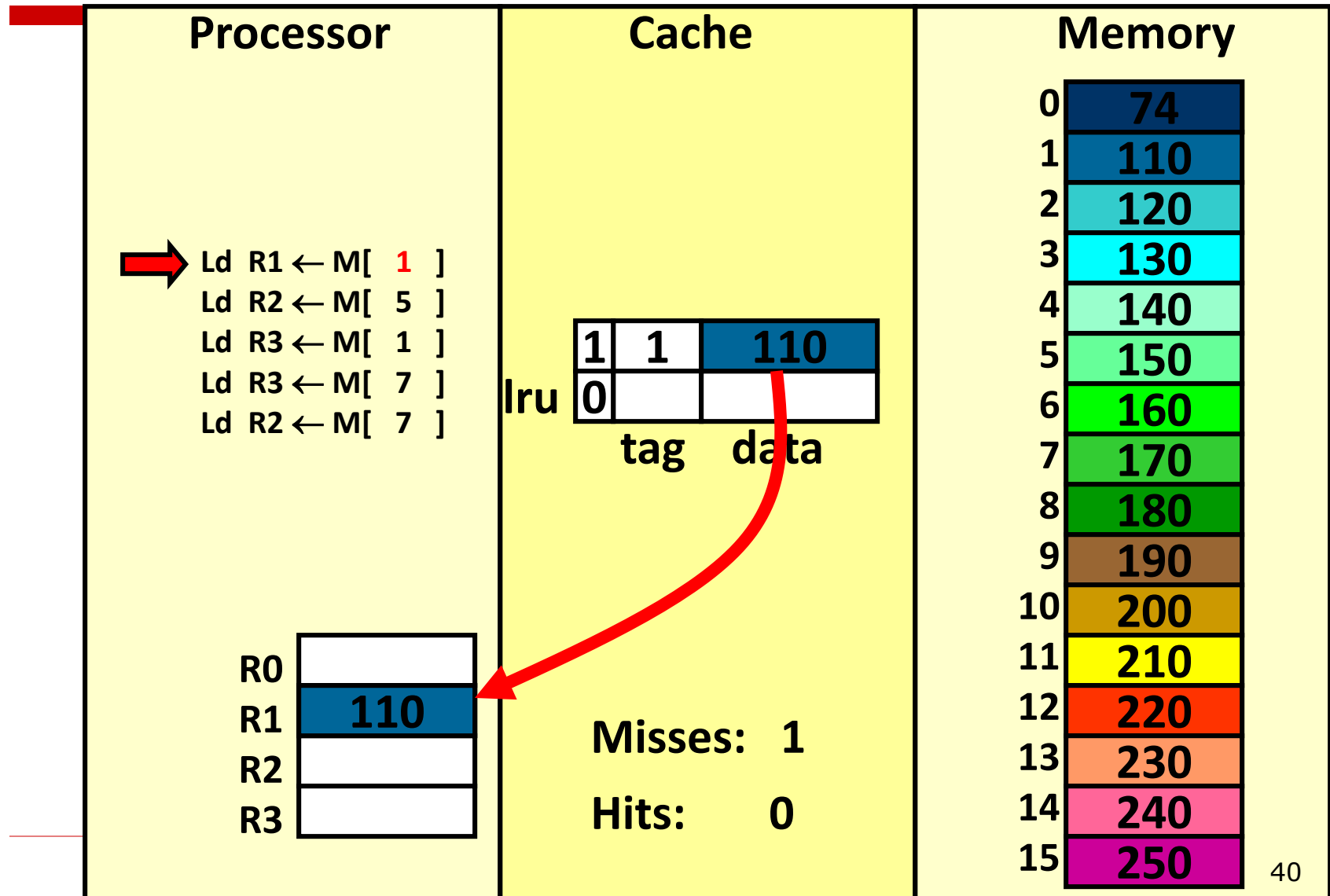
A Very Simple Memory System



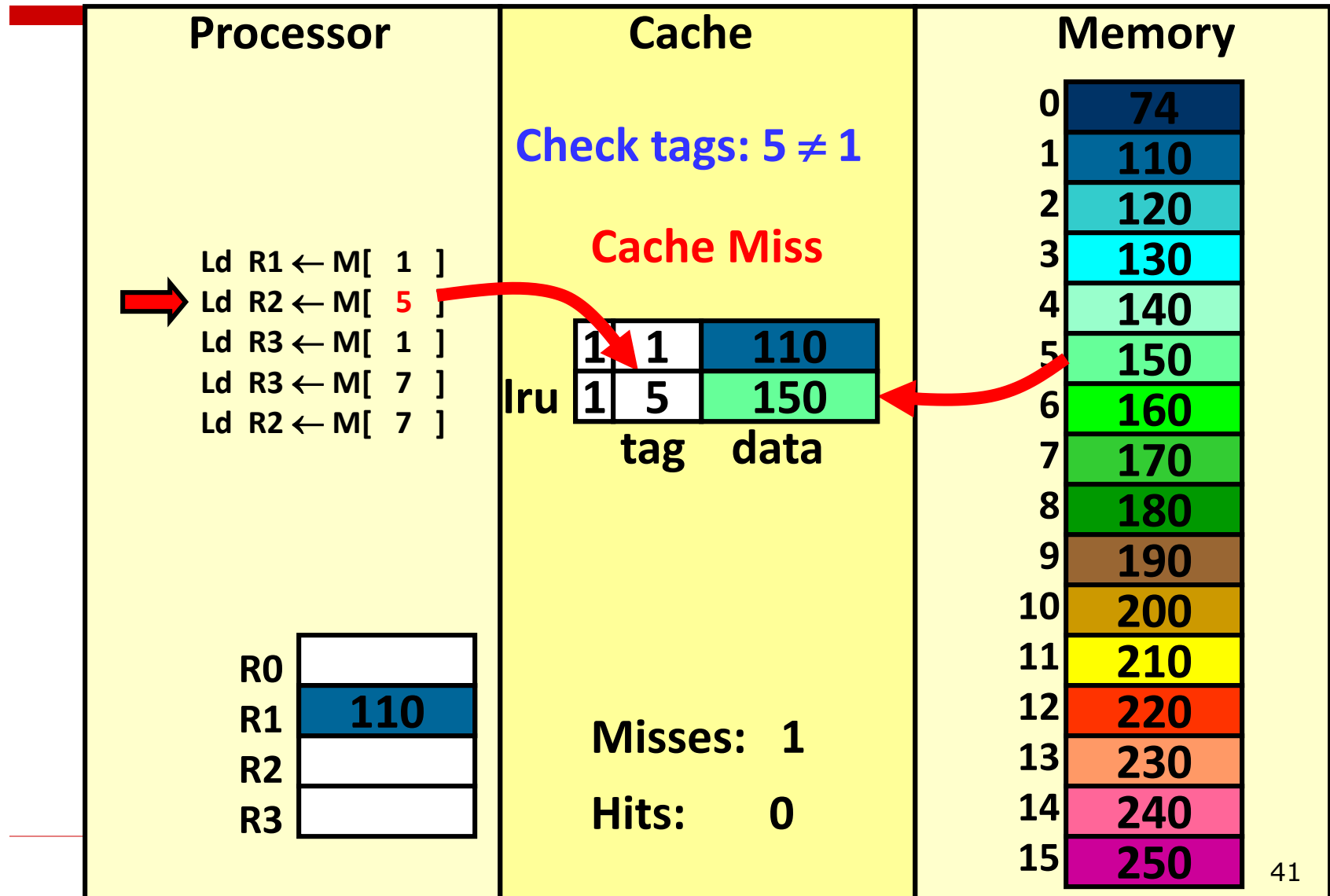
A Very Simple Memory System



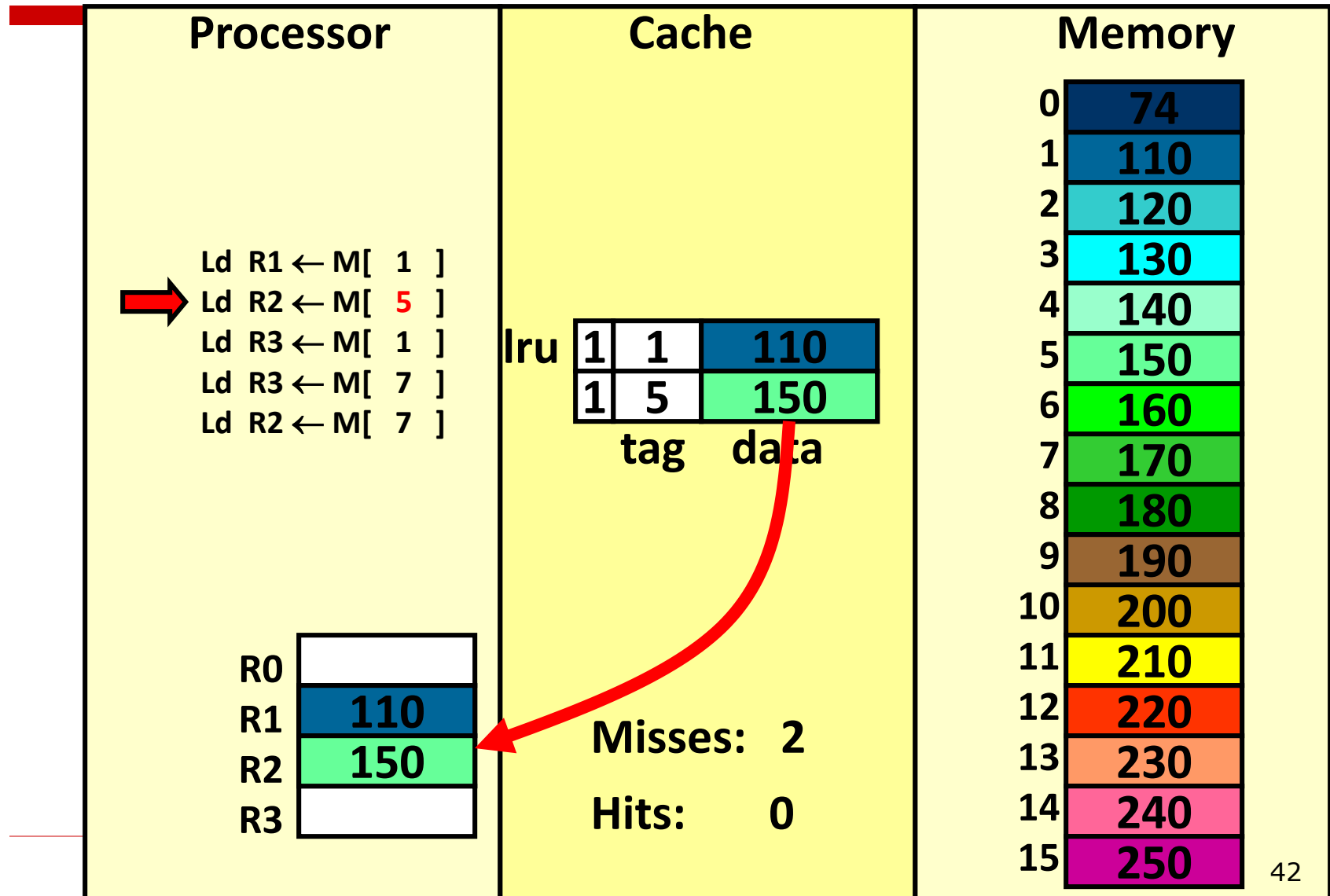
A Very Simple Memory System



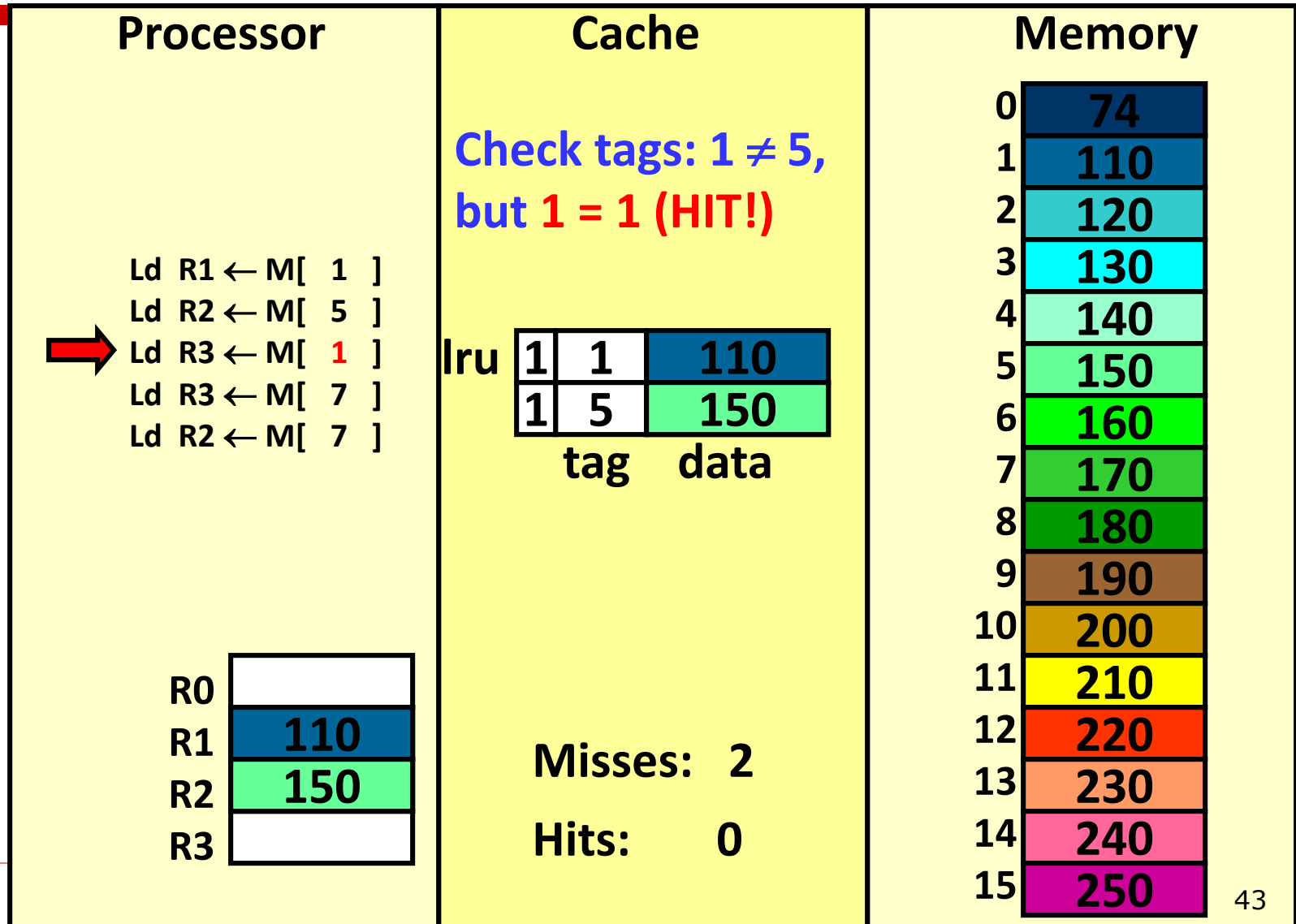
A Very Simple Memory System



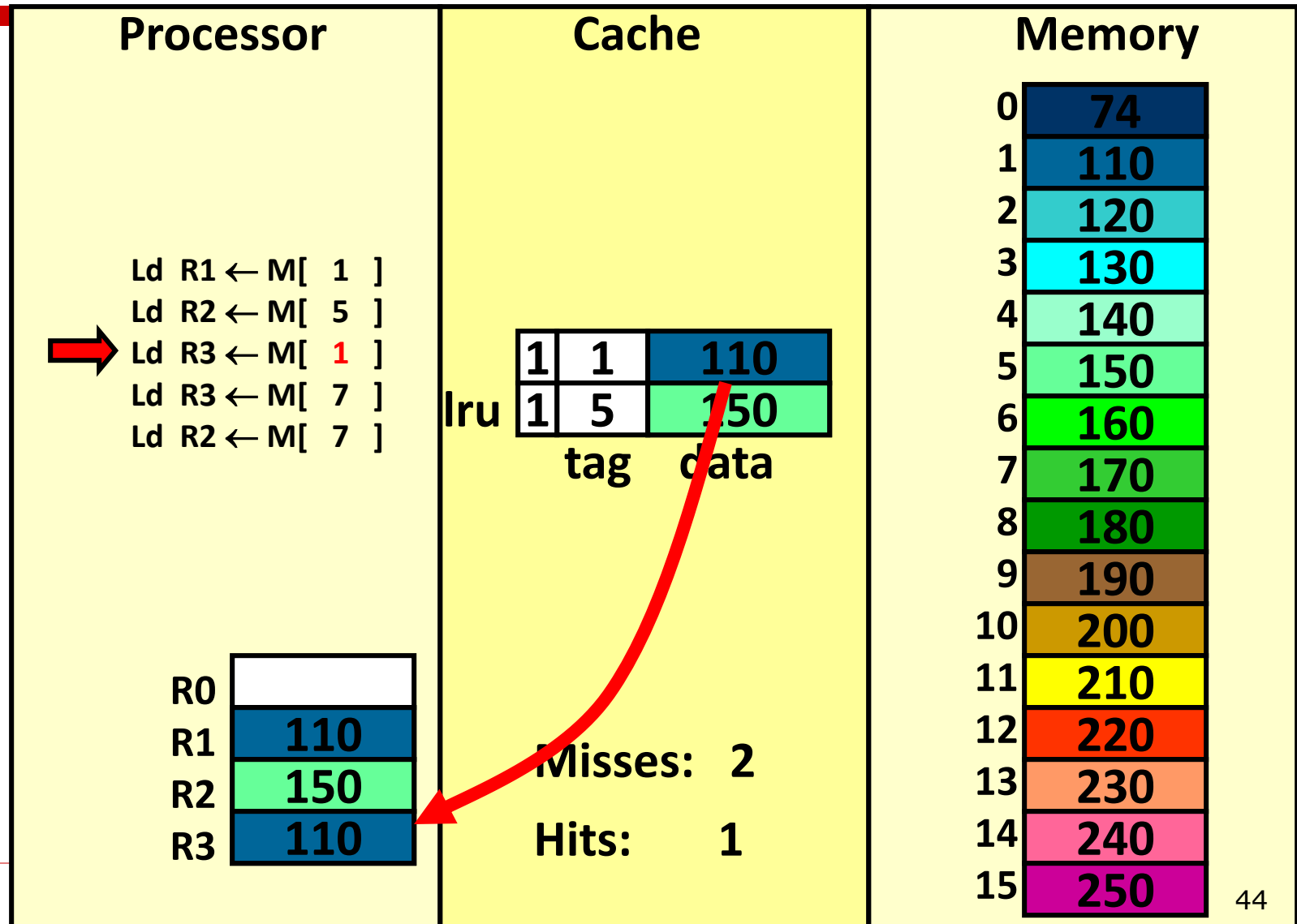
A Very Simple Memory System



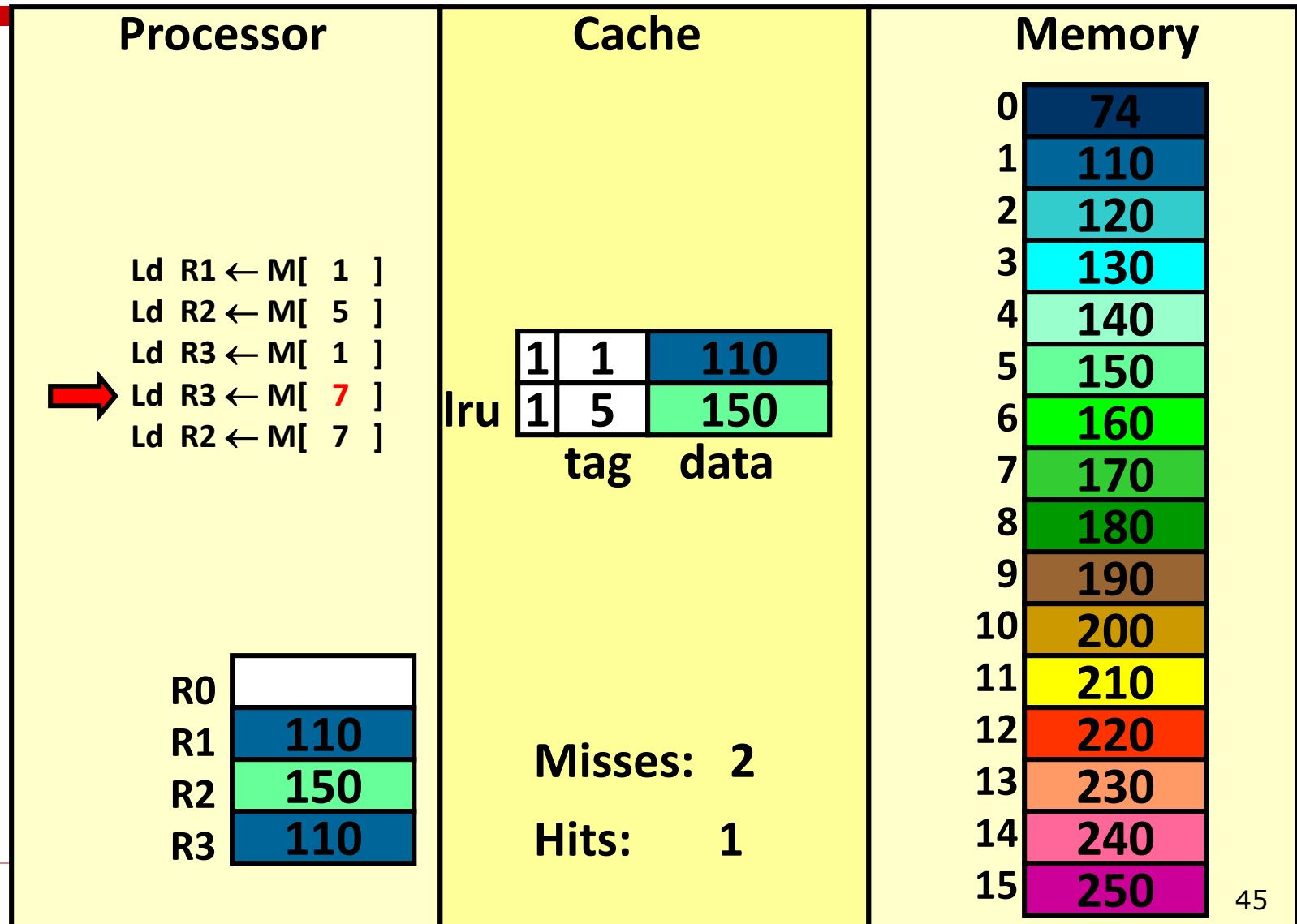
A Very Simple Memory System



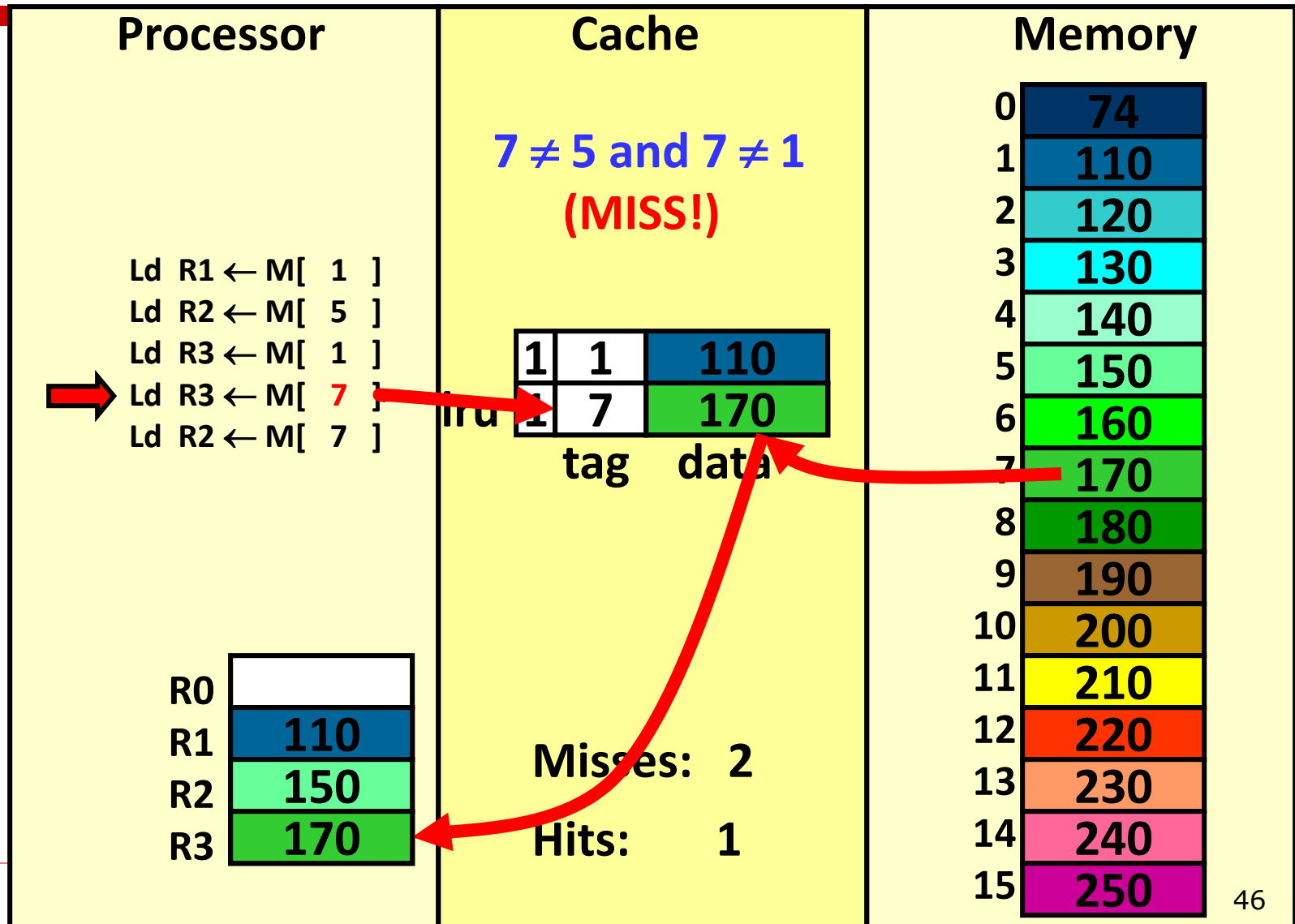
A Very Simple Memory System



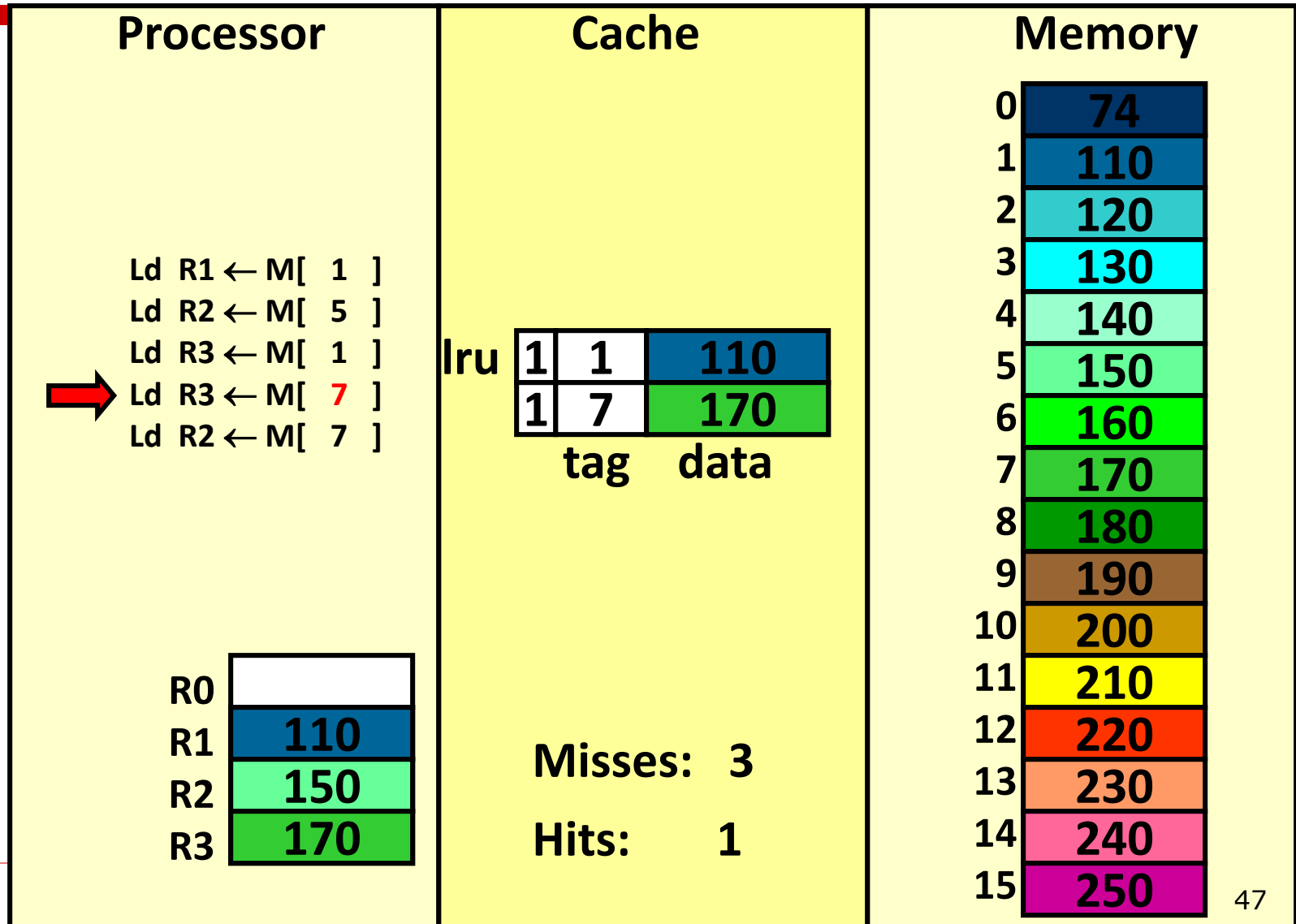
A Very Simple Memory System



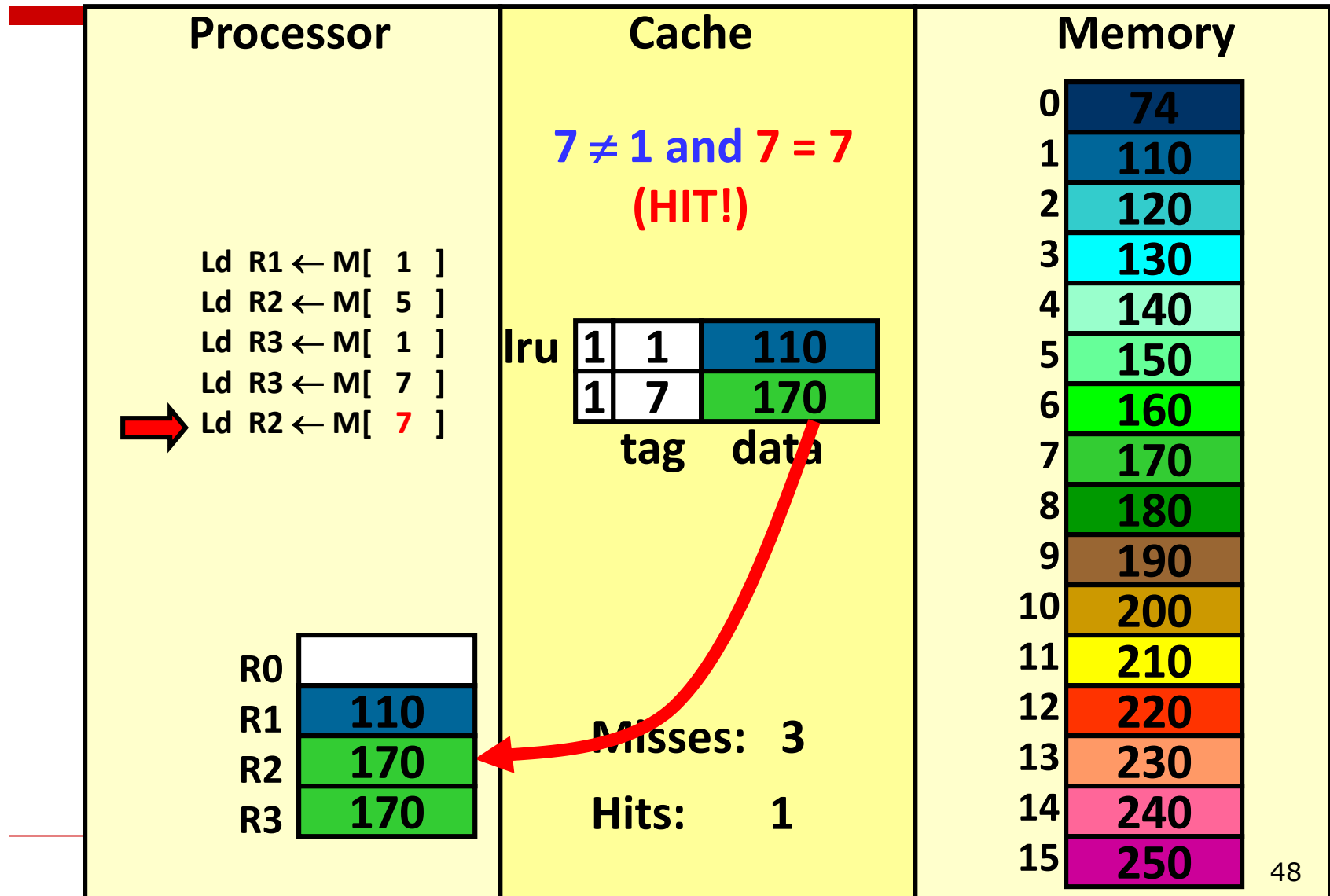
A Very Simple Memory System



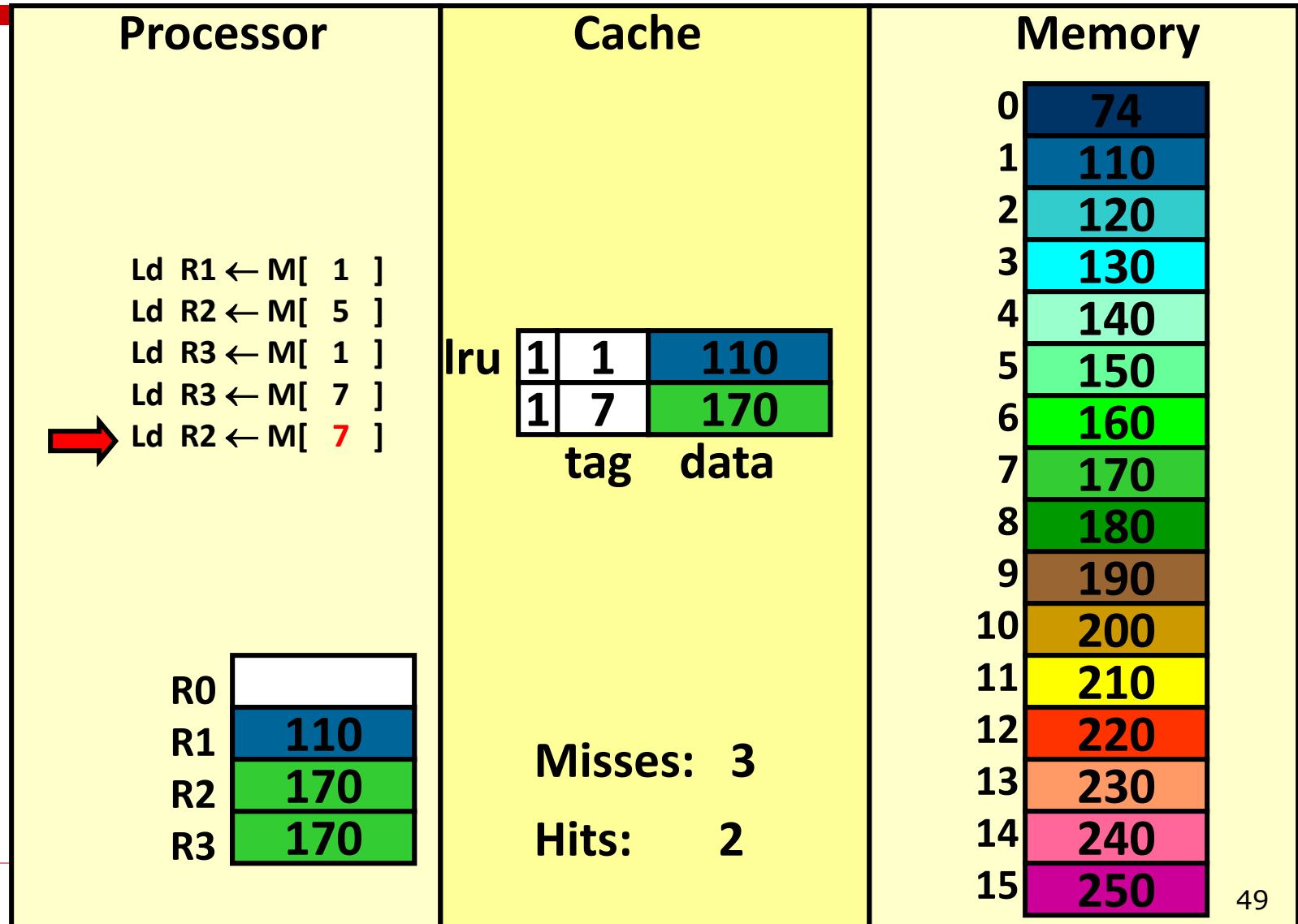
A Very Simple Memory System



A Very Simple Memory System



A Very Simple Memory System



Calculating Average Access Latency

- ❑ Avg latency
 - = cache latency \times hit rate + memory latency \times miss rate

- ❑ Avg latency for our example cache
 - = 1 cycle \times (2/5) + 15 \times (3/5)
 - = 9.4 cycles per reference

- ❑ To improve average latency:
 - Improve memory access latency, or
 - Improve cache access latency, or
 - Improve cache hit rate

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

How can we reduce the overhead?

- ❑ 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.
- ❑ To Be Continued...