# Lab 6 - Cache

CSED 311 Computer Architecture Lab

Sungjun Cho

2021 - 05 - 11

allencho1222@postech.ac.kr
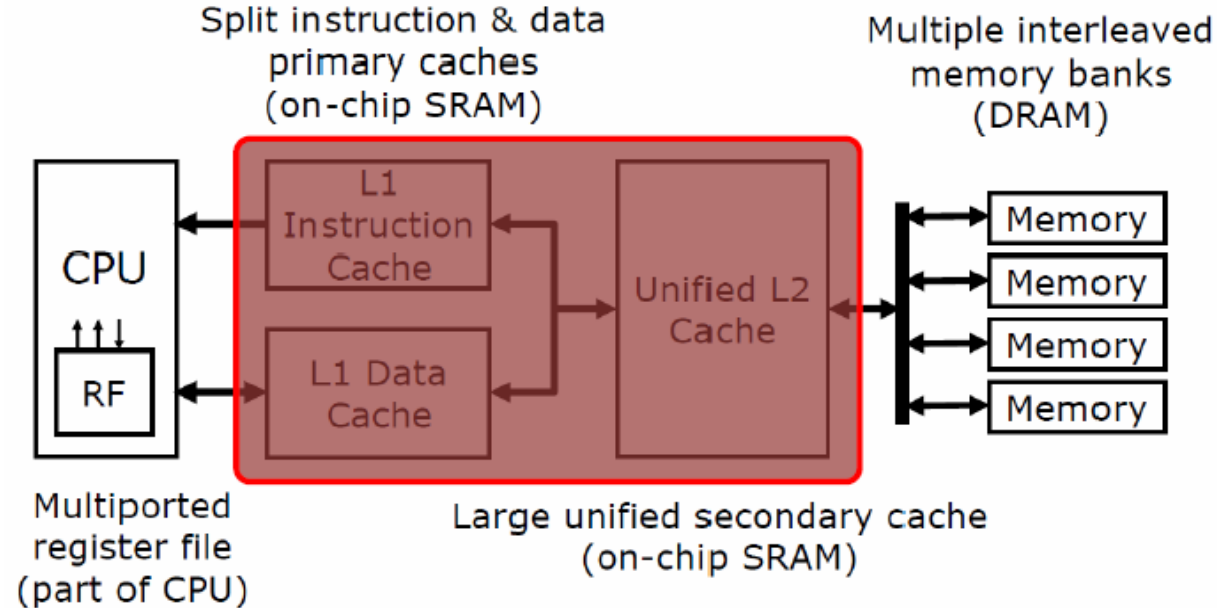csed311-ta@postech.ac.kr

# Objectives

- Understand how cache works

- <span style="color:red">Implement a set-associative cache on your pipelined CPU</span>

- Evaluate the speed-up achieved by using cache
  - Hit ratio
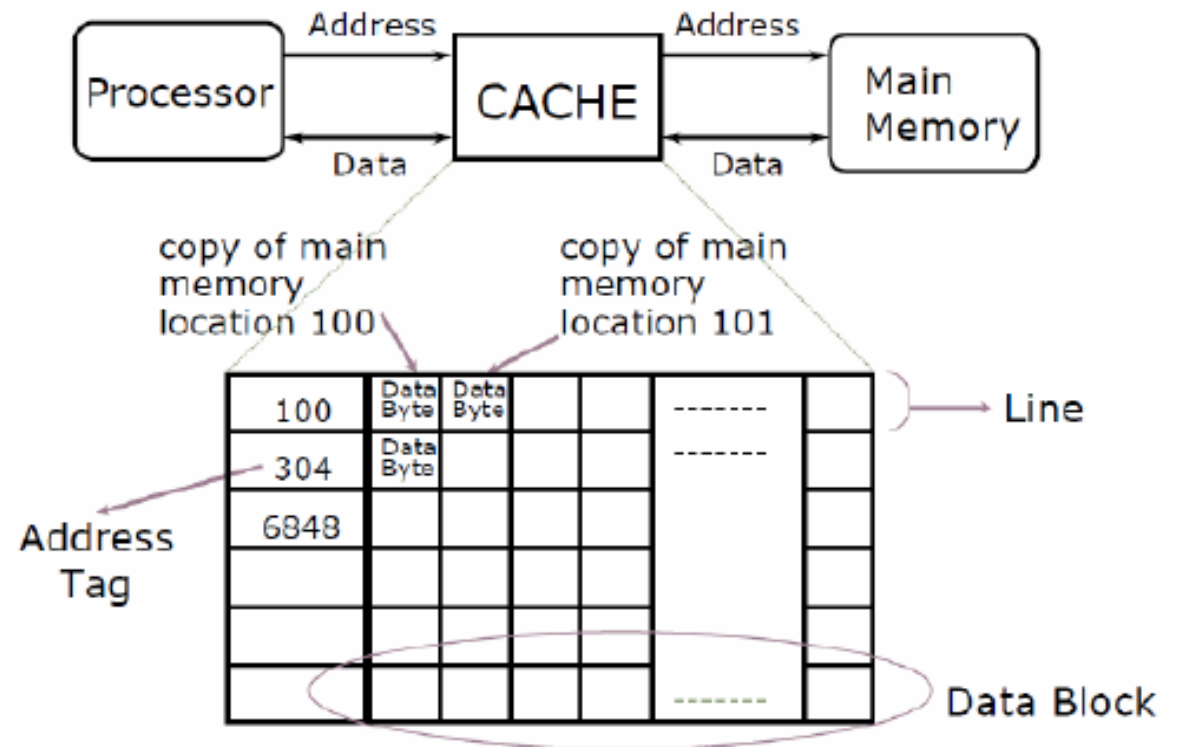  - Corresponding speed-up (vs. no-cache CPU)

# Cache

- Mitigating the gap between CPU and memory
  - Memory access: few hundred cycles
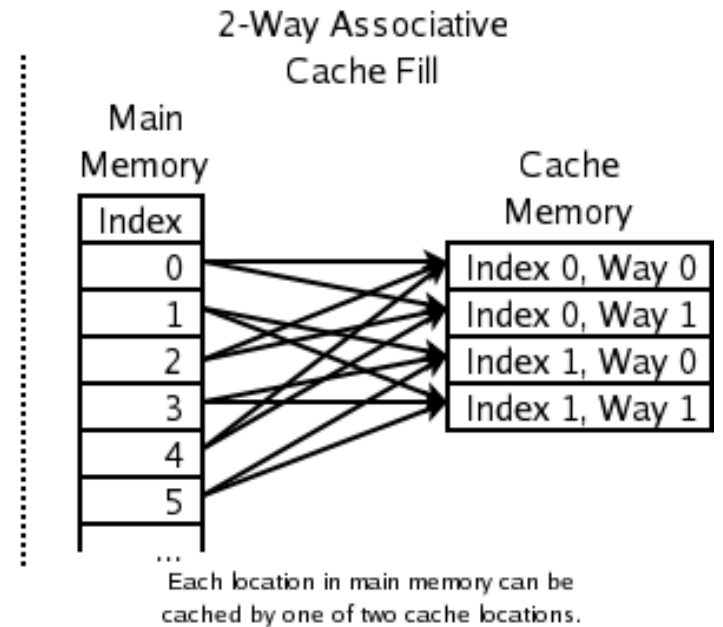  - Cache access: few cycles

- Why does it work?
  - Locality!



Split instruction & data primary caches (on-chip SRAM)

Multiple interleaved memory banks (DRAM)

CPU

RF

L1 Instruction Cache

L1 Data Cache

Unified L2 Cache

Memory
Memory
Memory
Memory

Multiported register file (part of CPU)

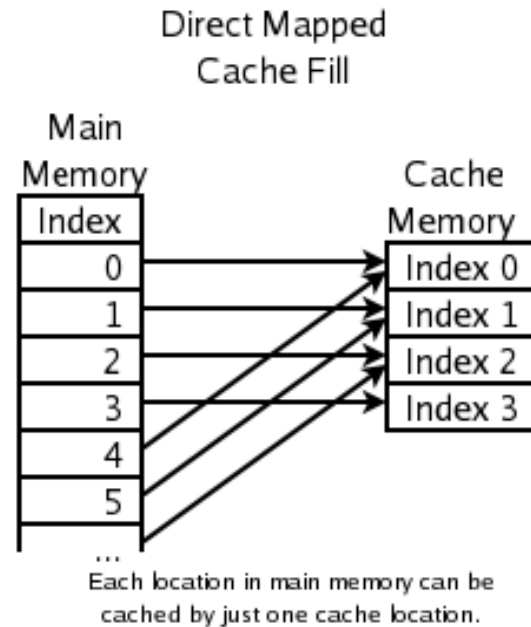Large unified secondary cache (on-chip SRAM)

# Cache: internal structure

- **Tag**
  - Detect address conflicts

- **Data**
  - Fetch by line: exploiting spatial locality

# Cache: associativity

■ **Associativity**
  • Reducing address conflicts

■ **Direct mapped, n-way, fully associative**
  • Tradeoffs exist



Direct Mapped Cache Fill — Each location in main memory can be cached by just one cache location.

2-Way Associative Cache Fill — Each location in main memory can be cached by one of two cache locations.

# Cache: other design choices

- Replacement policy
  - Random, LRU, FIFO, …,
  - Each one has strengths & drawbacks

- Write policy
  - Write-through, writeback, write-no-allocate
  - Related to coherency management

# Lab Assignment 06 (1/5)

- Design and implement your own cache with the following requirements:
  - 2-way set associative, single-level cache
  - Capacity: 32 words / Line size: 4 words
  - If hit, return data in the following cycle
  - It should be a part of CPU (not TB)

- You have to choose the following design choices:
  - Replacement policy, write policy
  - Unified (32 words == 8 Lines) or separate I/D (4 lines ea.) cache
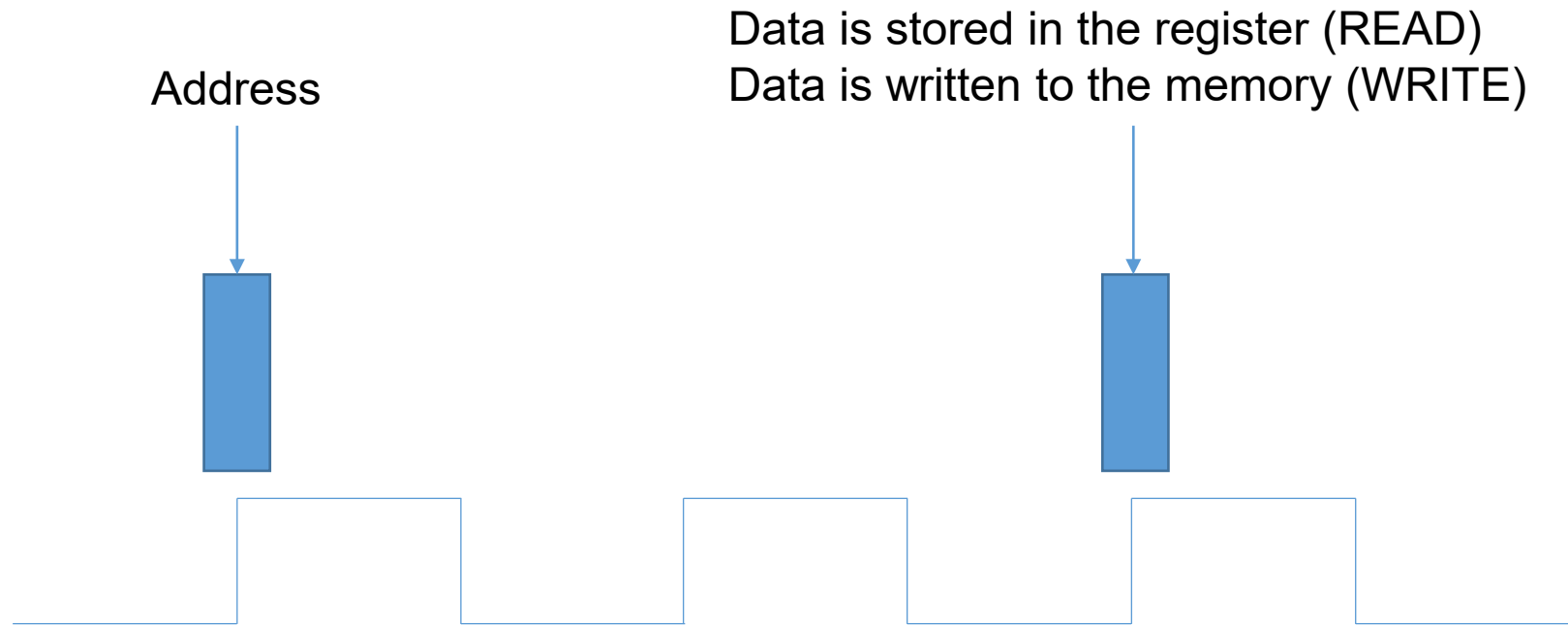
# Lab Assignment 06 (2/5)

- You may implement direct mapped cache with reduced score of 80%
  - Direct-mapped, single-level cache
  - Capacity: 16 words / Line size: 4 words
  - If hit, return data in the following cycle
  - It should be a part of CPU (not TB)

# Lab Assignment 06 (3/5)

- **Memory access latency – Previous pipelined CPU**
  - One memory access fetches one word with one cycle

- **New memory access latency**
  - <u>You're required to modify memory.v</u>
  - Cache hit takes <span style="color:red">one cycle</span>
  - One memory access should fetch <span style="color:red">four words into cache (= one cache line) and take six cycles</span>
  - For your baseline CPU (CPU without cache), one memory access should fetch <span style="color:red">one word</span> and take <span style="color:red">two cycles</span>

  - Then, you need to implement two different memory models
    - ✓ Memory for CPU with cache: return four words in six cycles
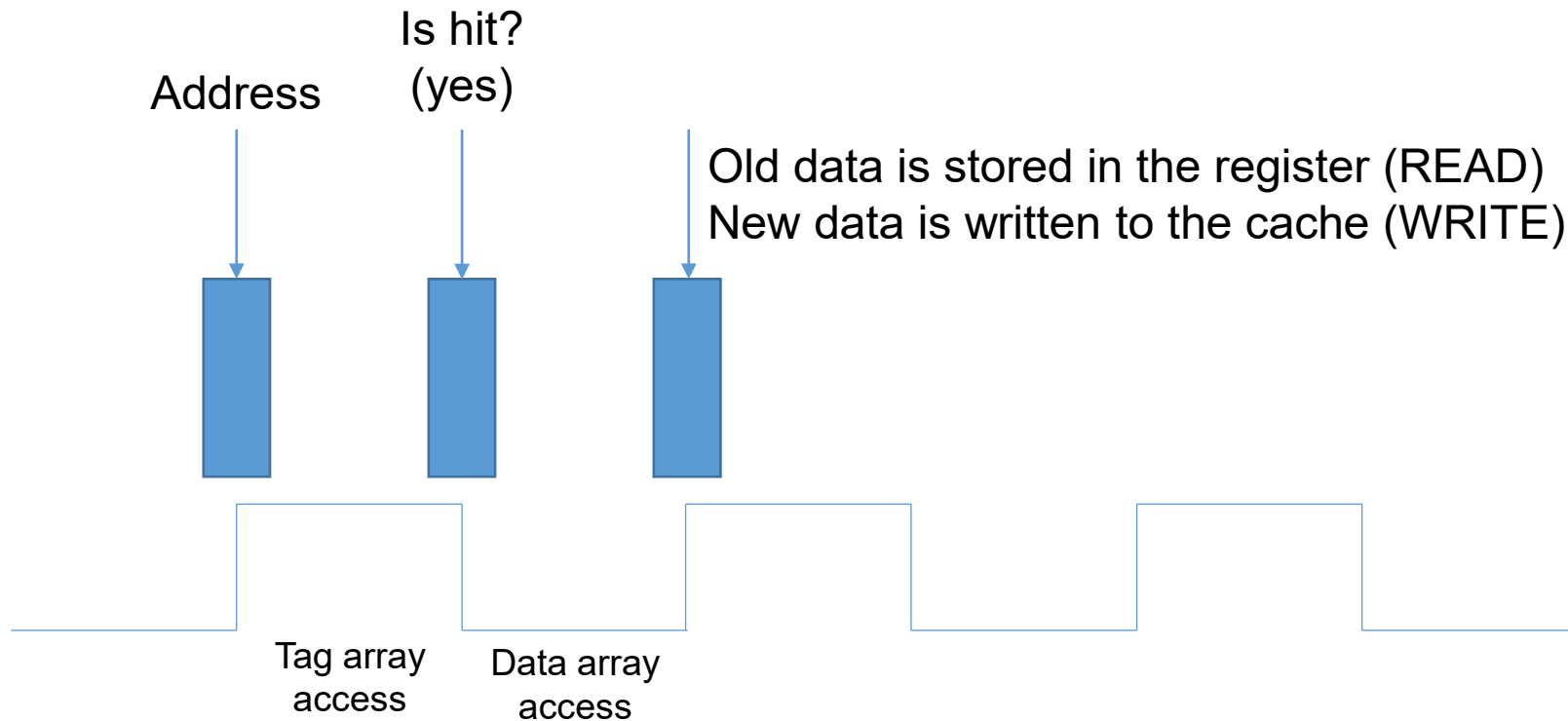    - ✓ Memory for CPU without cache: return one word in two cycles
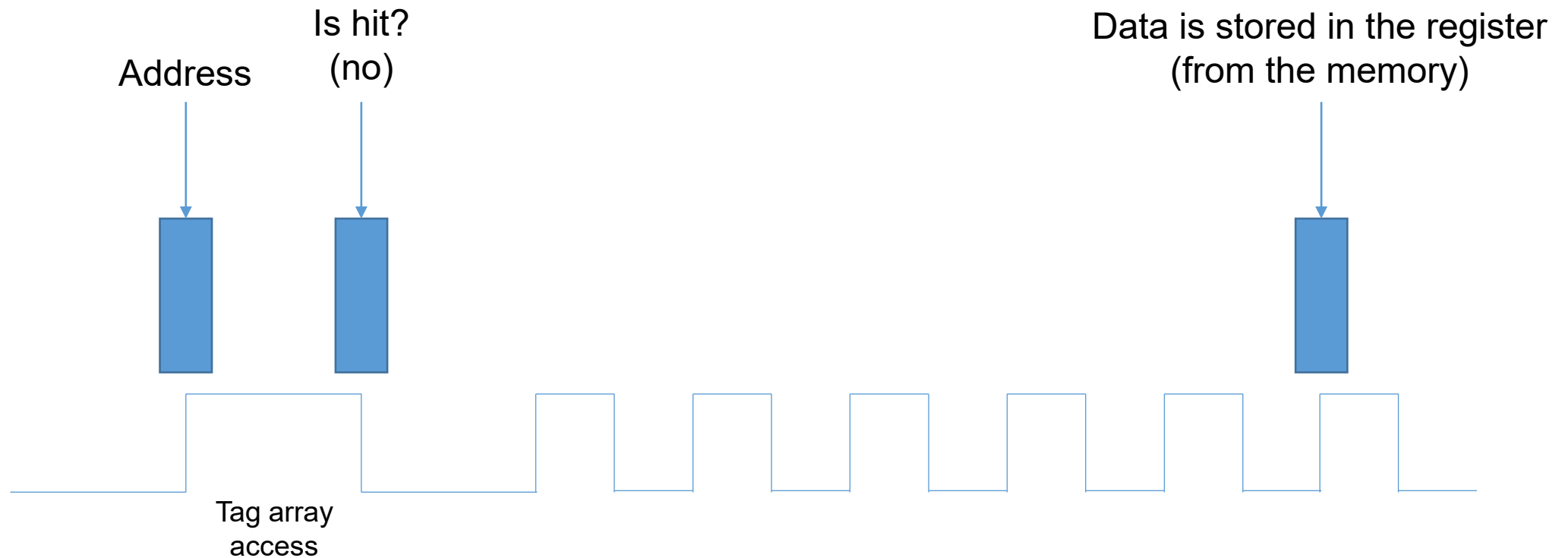
# Waveform (Baseline CPU)

- Cache does not exist

Address

Data is stored in the register (READ)
Data is written to the memory (WRITE)

# Waveform (Cache hit)

- Cache miss/hit and data can be checked in a single cycle

Address

Is hit?
(yes)

Old data is stored in the register (READ)
New data is written to the cache (WRITE)

Tag array
access

Data array
access

# Waveform (Cache miss)

Address

Is hit?
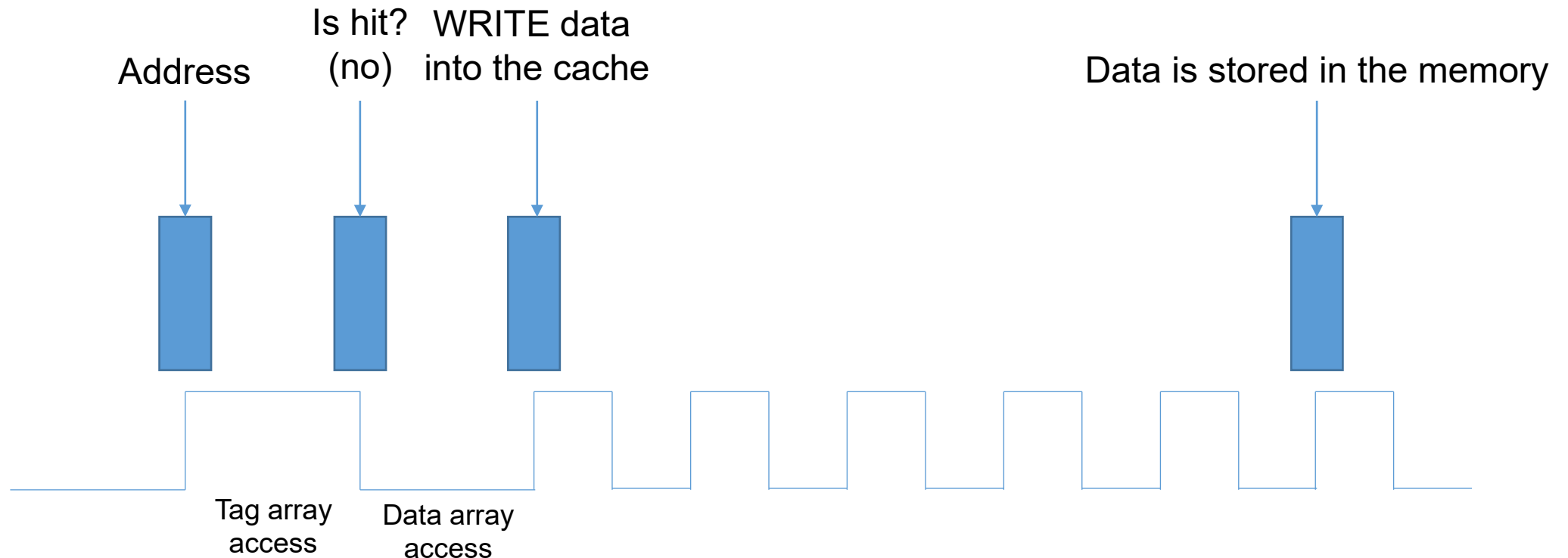(no)

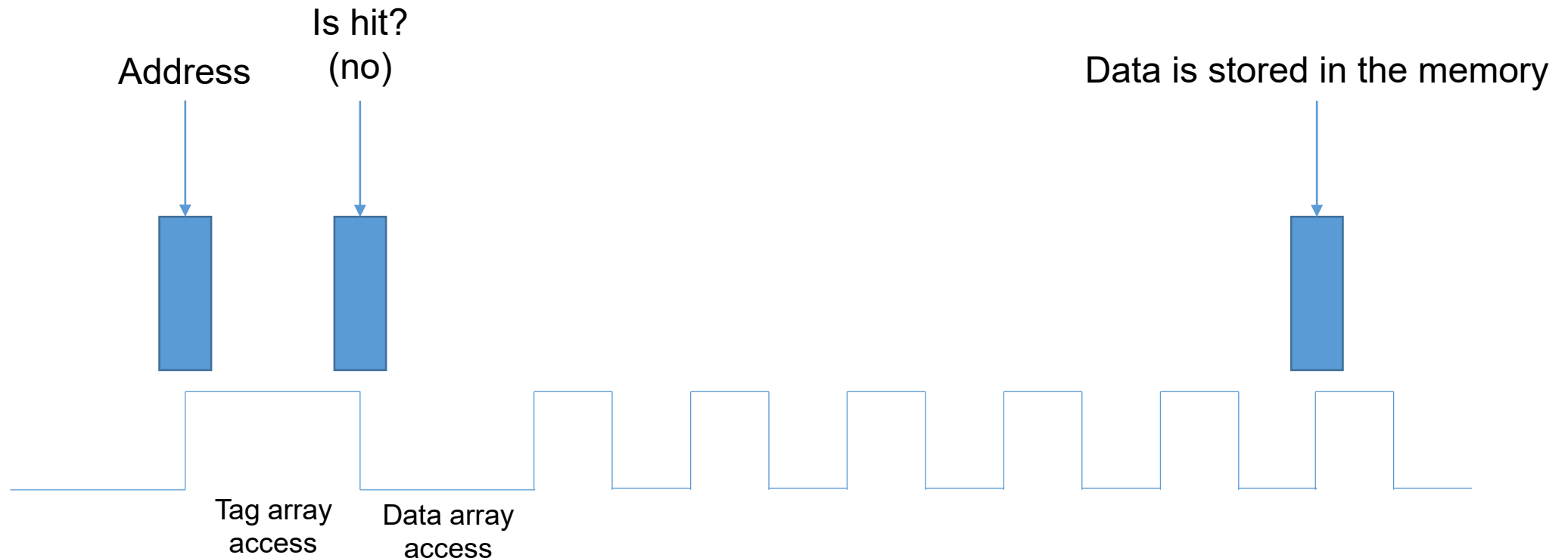Data is stored in the register
(from the memory)

Tag array
access

# Waveform (WRITE-THROUGH)

- Usually comes with write-no-allocate
  - Cache hit (WRITE request)
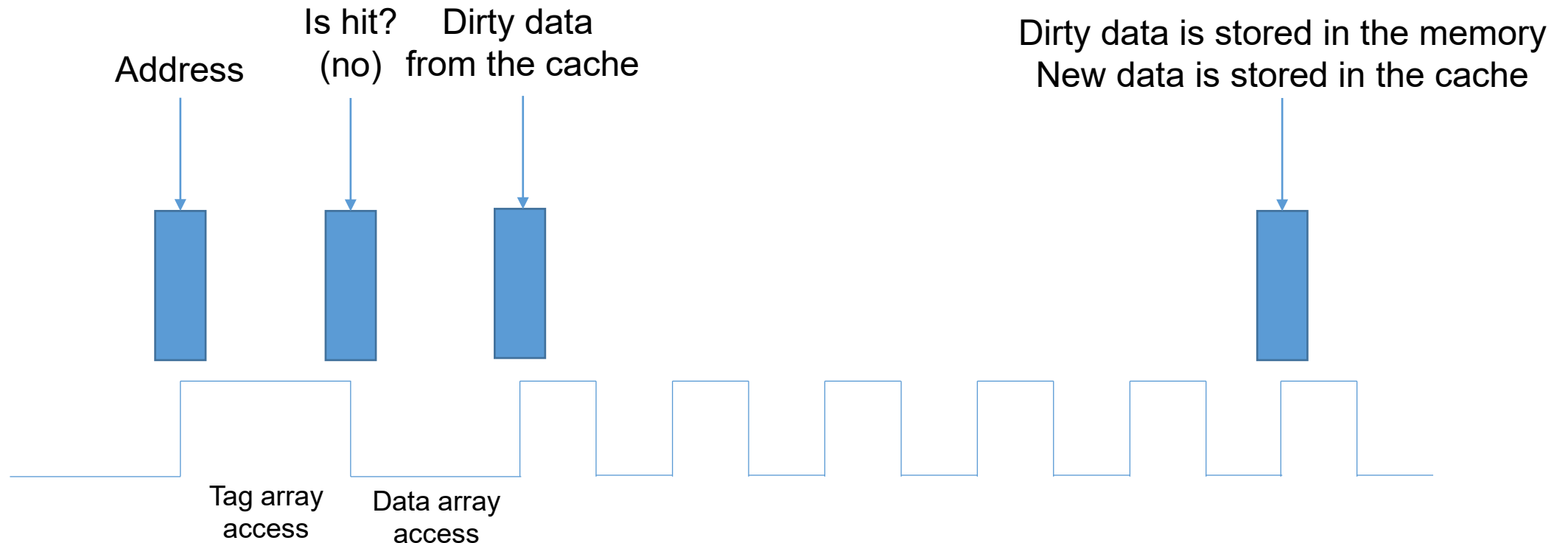    - ✓ WRITE data into the cache + WRITE data into the memory



Address

Is hit? (no)

WRITE data into the cache

Data is stored in the memory

Tag array access

Data array access

# Waveform (WRITE-THROUGH)

- Usually comes with write-no-allocate
  - Cache miss (WRITE request)
    - ✓ WRITE data into the memory



Address

Is hit?
(no)

Data is stored in the memory

Tag array
access

Data array
access

# Waveform (WRITE-BACK)

- Usually comes with write-allocate
  - Cache miss (WRITE request)
    - ✓ READ data from the memory + WRITE dirty data into the memory

# Lab Assignment 06 (4/5)

- **Memory requirements**
  - You can use either a 2-port RAM or a single-port RAM
  - Latencies of RAM
    - ✓ Should be serialized



⇒ Different ports can handle independent requests!

# Lab Assignment 06 (5/5)

- For the report, the following contents should be included
  - Describe your design choice
  - Calculate the hit (or miss) ratio
    - ✓ Hit ratio = (# of hits) / (# of memory accesses)
  - Compare the performance
    - ✓ CPU without cache vs. CPU with cache
    - ✓ You should use the new latencies!

# Announcement

- A skeleton code from this lab is not provided.
  You should implement cache *on your previous implementation of pipelined CPU*.


- You *cannot begin* this assignment unless you have finished previous work.