



香港中文大學

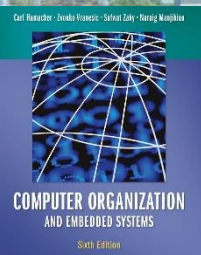
The Chinese University of Hong Kong

CSCI2510 Computer Organization

Lecture 07: Cache in Action

Ming-Chang YANG

mcyang@cse.cuhk.edu.hk



Reading: Chap. 8.6

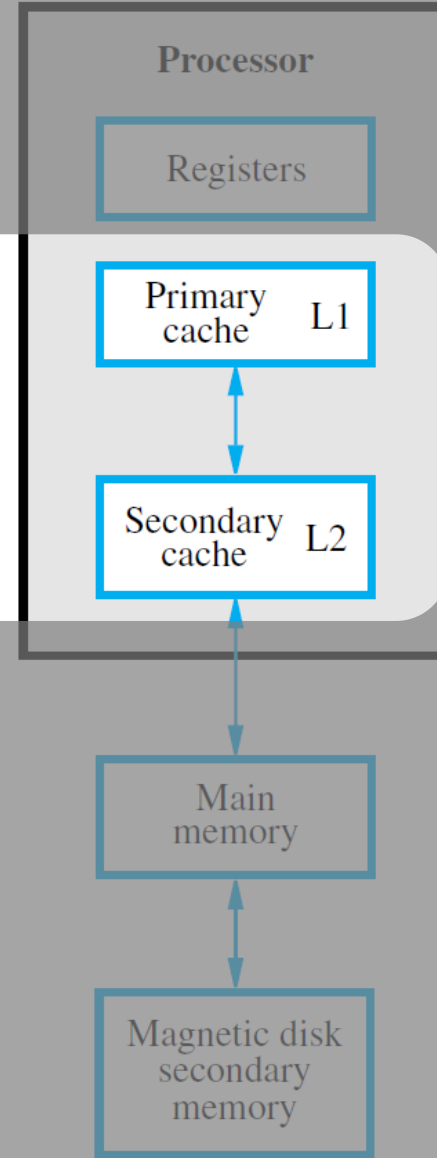
Recall: Memory Hierarchy



Processor

- Register: SRAM
- L1, L2 cache: SRAM
- Main memory: SDRAM
- Secondary storage: Hard disks or NVM

Increasing size
↓



↑
Increasing speed

↑
Increasing cost per bit

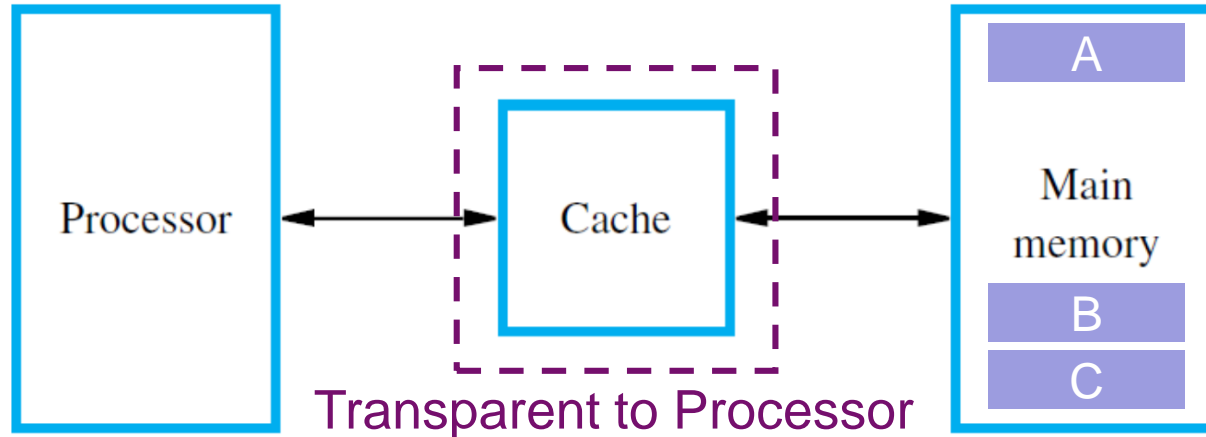


- Cache Basics
- Mapping Functions
 - Direct Mapping
 - Associative Mapping
 - Set Associative Mapping
- Replacement Algorithms
 - Least Recently Used (LRU) Replacement
 - Random Replacement
- Working Examples

Cache: Fast but Small



- The cache is a **small** but **very fast** memory.
 - Interposed between the processor and main memory.



- Its purpose is to make the main memory appear to the processor to be much faster than it actually is.
 - The processor does not need to know explicitly about the **existence of the cache**, but just feels faster!
- How to? Exploit the **locality of reference** to “properly” load some data from the main memory into the cache.

Locality of Reference



- **Temporal Locality** (locality in *time*)
 - If an item is referenced, it will tend to be **referenced again soon** (e.g. recent calls).
 - **Strategy:** When information item (instruction or data) is first needed, opportunistically bring it into cache (we hope it will be used soon).
- **Spatial Locality** (locality in *space*)
 - If an item is referenced, **neighboring items** whose addresses are close-by will tend to be **referenced** soon.
 - **Strategy:** Rather than a single word, fetch more data of adjacent addresses (unit: **cache block**) from main memory into cache.

Cache Usage



- **Cache Read (or Write) Hit/Miss**: The read (or write) operation **can/cannot** be performed on the cache.



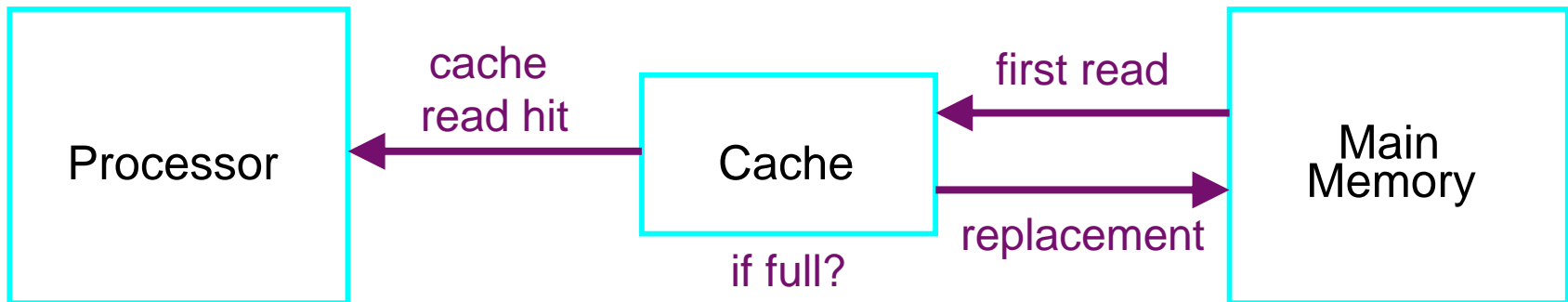
- **Cache Block / Line**: The unit composed of **multiple successive memory words** (size: cache block > word).
 - The contents of a cache block (of memory words) will be loaded into or unloaded from the cache at a time.
- **Mapping Functions**: Decide how cache is organized and how addresses are mapped to the main memory.
- **Replacement Algorithms**: Decide which item to be unloaded from cache when cache is full.

Read Operation in Cache



- **Read Operation:**

- Contents of a **cache block** are loaded from the memory into the cache for the **first read**.
- Subsequent accesses that can be (hopefully) performed on the cache, called a **cache read hit**.
- The number of cache entries is relatively small, we need to keep the most likely to-be-used data in cache.
- When an un-cached block is required (i.e., **cache read miss**), the **replacement algorithm** removes an old block and to create space for the new one if cache is full

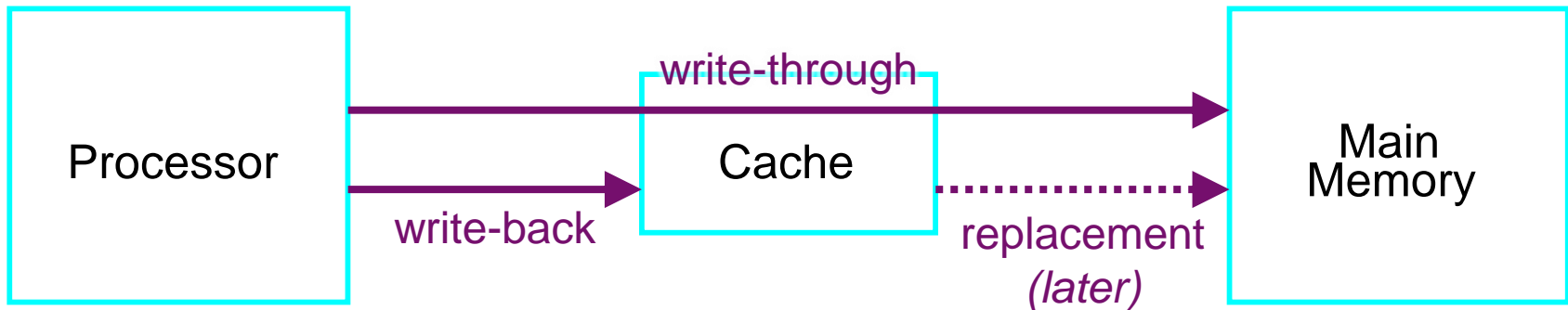


Write Operation in Cache



- **Write Operation:**

- **Scheme 1:** The contents of cache and main memory are updated at the same time (**write-through**).
- **Scheme 2:** Update cache only but mark the item as **dirty**. The corresponding contents in main memory will be updated later when cache block is unloaded (**write-back**).
 - **Dirty:** The data item needs to be written back to the main memory.



- Which scheme is simpler?
- Which one has better performance?

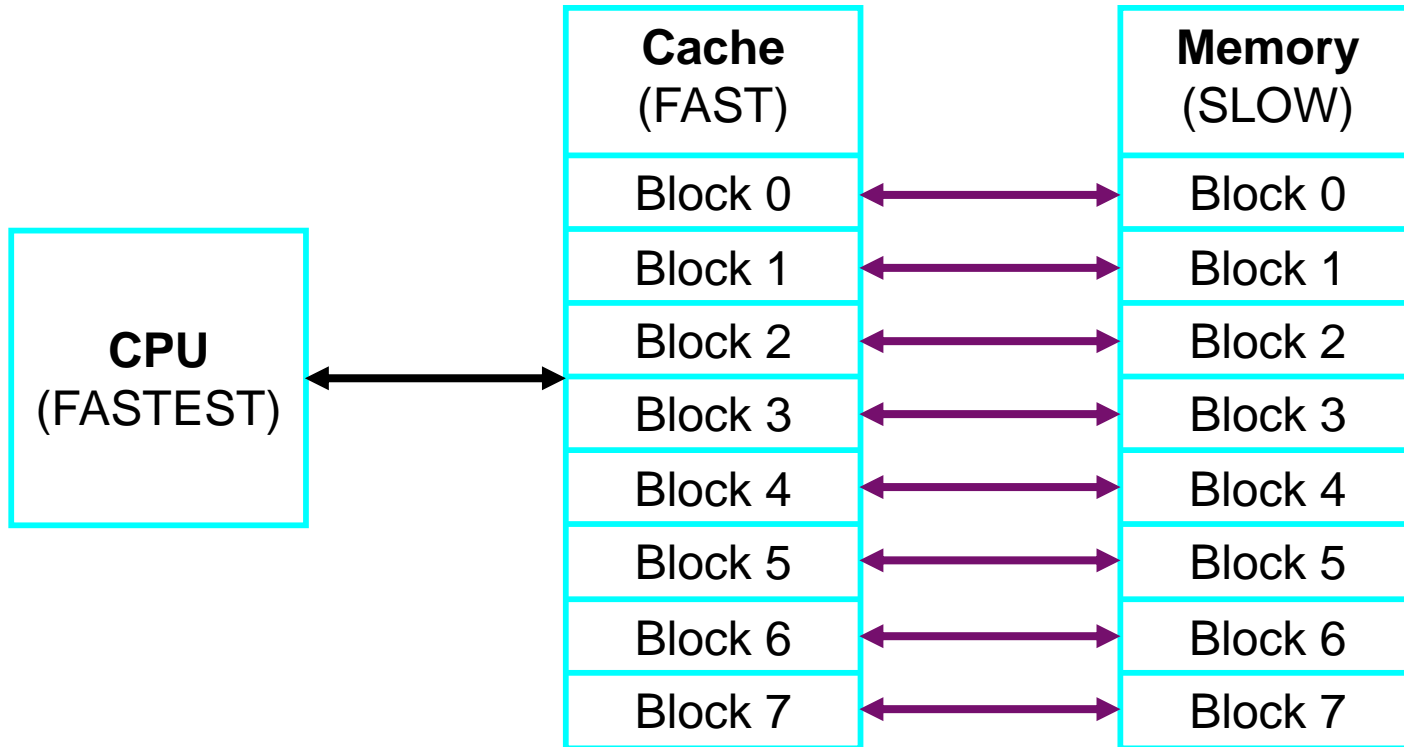


- Cache Basics
- Mapping Functions
 - Direct Mapping
 - Associative Mapping
 - Set Associative Mapping
- Replacement Algorithms
 - Least Recently Used (LRU) Replacement
 - Random Replacement
- Working Examples

Mapping Functions (1/3)



- **Cache-Memory** Mapping Function: A way to record which block of the main memory is now in cache.
- What if the case size == the main memory size?



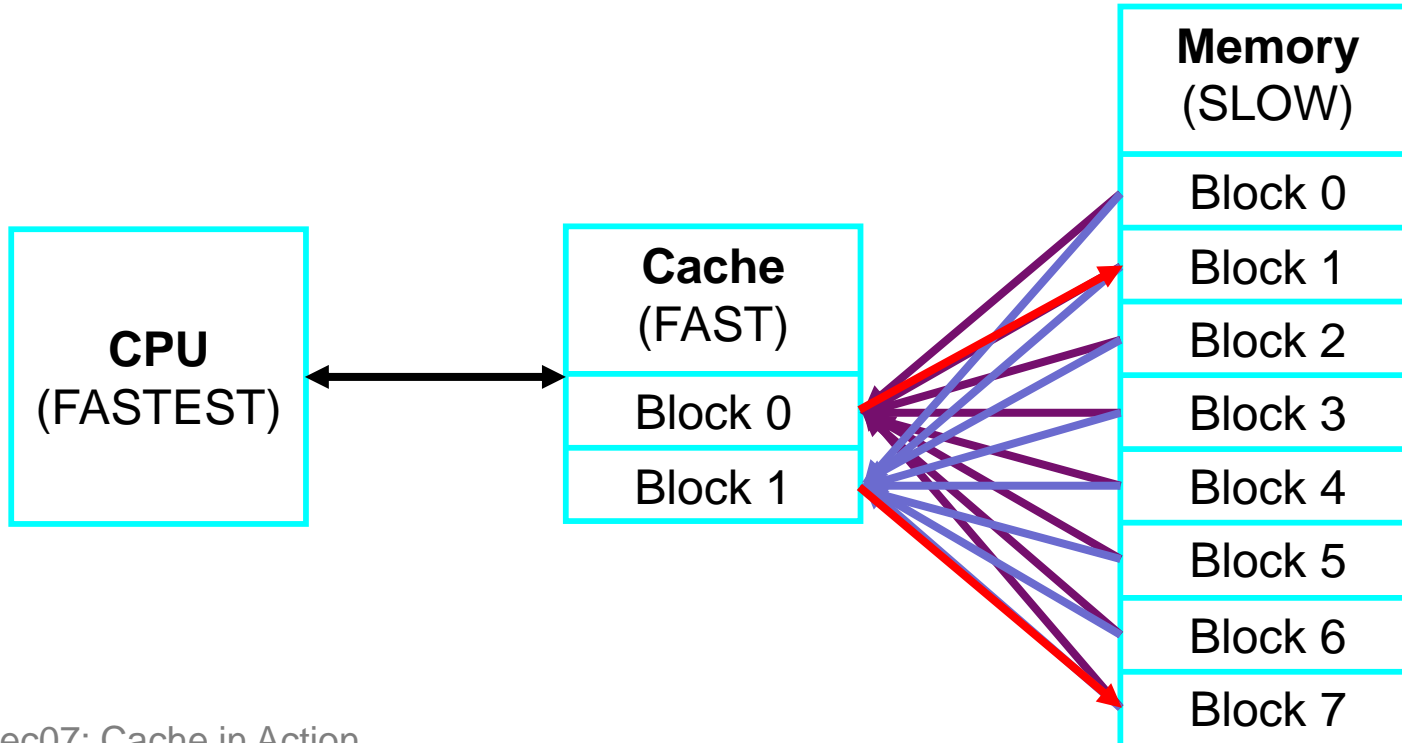
- Trivial! **One-to-one mapping** is enough!

Question: Do we still need the main memory?

Mapping Functions (2/3)



- **Reality:** The cache size is much smaller (\lll) than the main memory size.
- **Many-to-one mapping** is needed!
 - **Many** blocks in memory compete for **one** block in cache.
 - A block in cache can only represent **one** block in memory.



Mapping Functions (3/3)

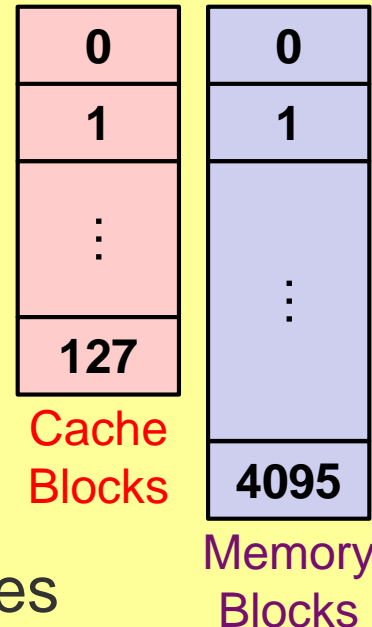


- **Design Considerations:**

- **Efficient:** Determine whether a block is in cache quickly.
- **Effective:** Make full use of cache to increase **cache hit ratio**.
 - **Cache Hit/Miss Ratio:** the probability of cache hits/misses.

- In the following discussion, we assume:

- **Synonym:** Cache Line = Cache Block = Block
- 1 Word = 1 Byte
- 1 Block = 16 Words = 16 Bytes
- **Cache Size:** 2K Bytes → **128 Cache Blocks**
- **Memory Size:** 16-bit Address → $2^{16} = 64\text{K Bytes}$
→ **4096 Memory Blocks**

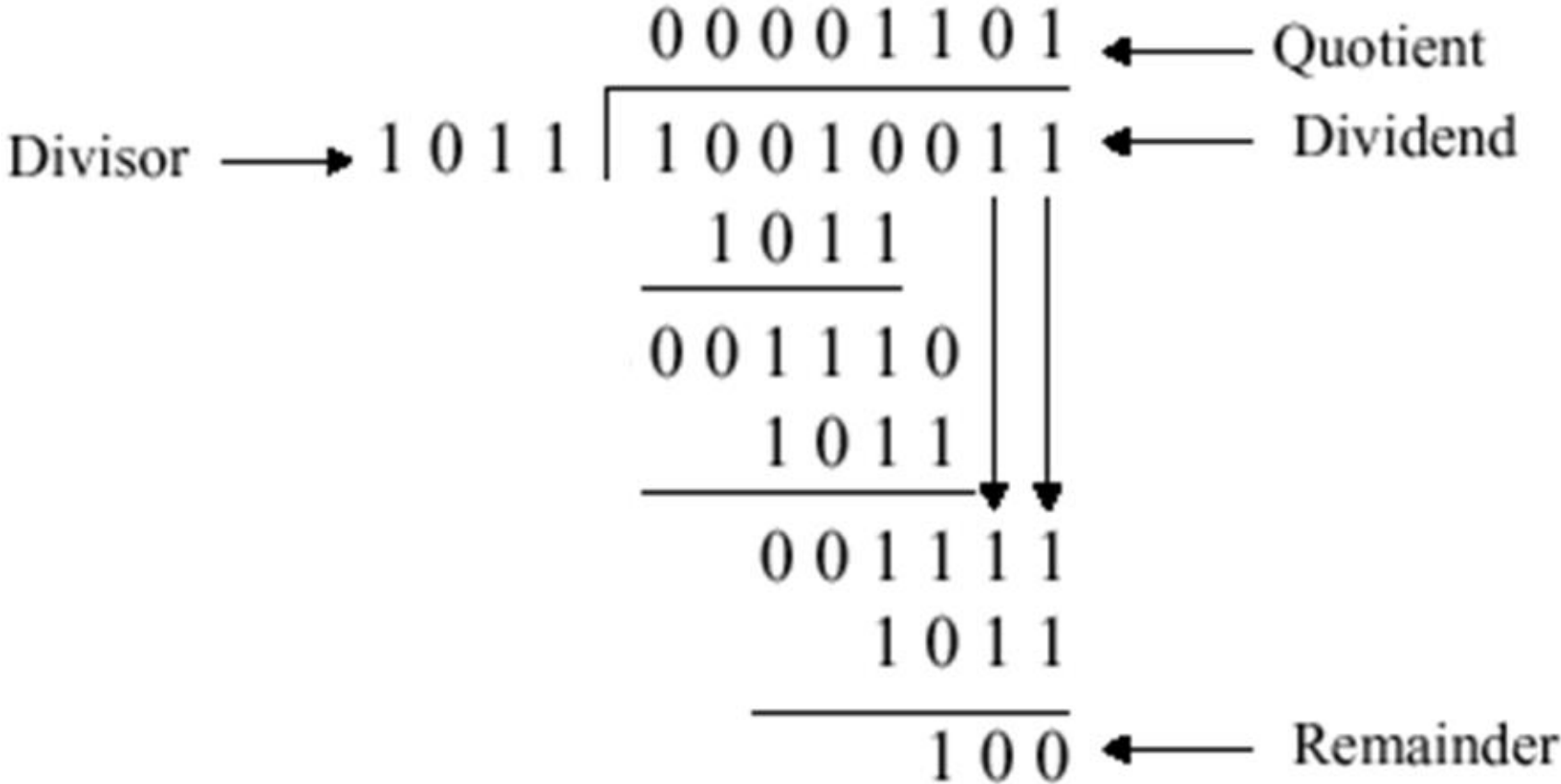


- **Memory Block (MB):** The block in the main memory.



Modulo (% , mod) Operator

- The **modulo (%)** operator is used to divide two numbers and get the **remainder**.
- Example:



Class Exercise 7.1

Student ID: _____ Date: _____

Name: _____

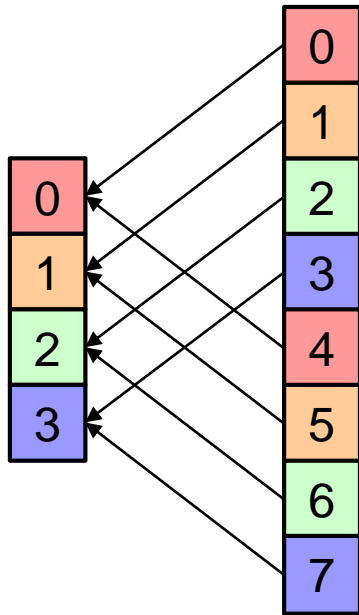
- Considering the previous example, what if the divisor equals to $(10)_2$, $(100)_2$, ..., $(10000000)_2$?

Direct Mapping (1/4)



Direct

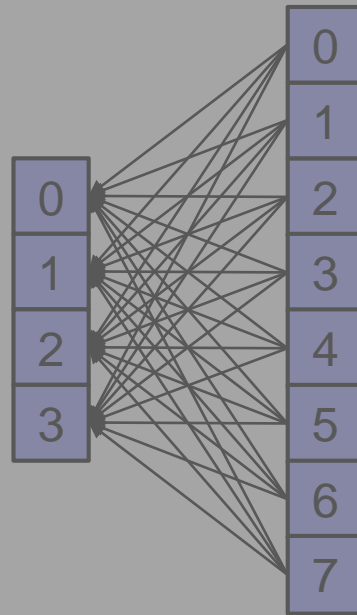
- A Memory Block is directly mapped (%) to a Cache Block.



Cache Blocks Memory Blocks

Associative

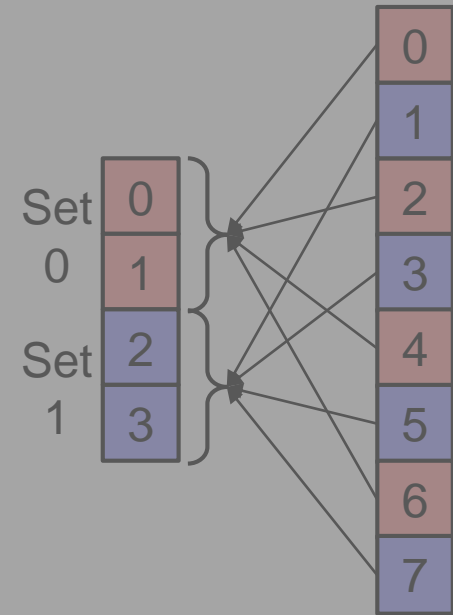
- A Memory Block can be mapped to any Cache Block. (First come first serve!)



Cache Blocks Memory Blocks

Set Associative

- A Memory Block is directly mapped (%) to a Cache Set. (In a set? Associative!)



Cache Blocks Memory Blocks

Direct Mapping (2/4)

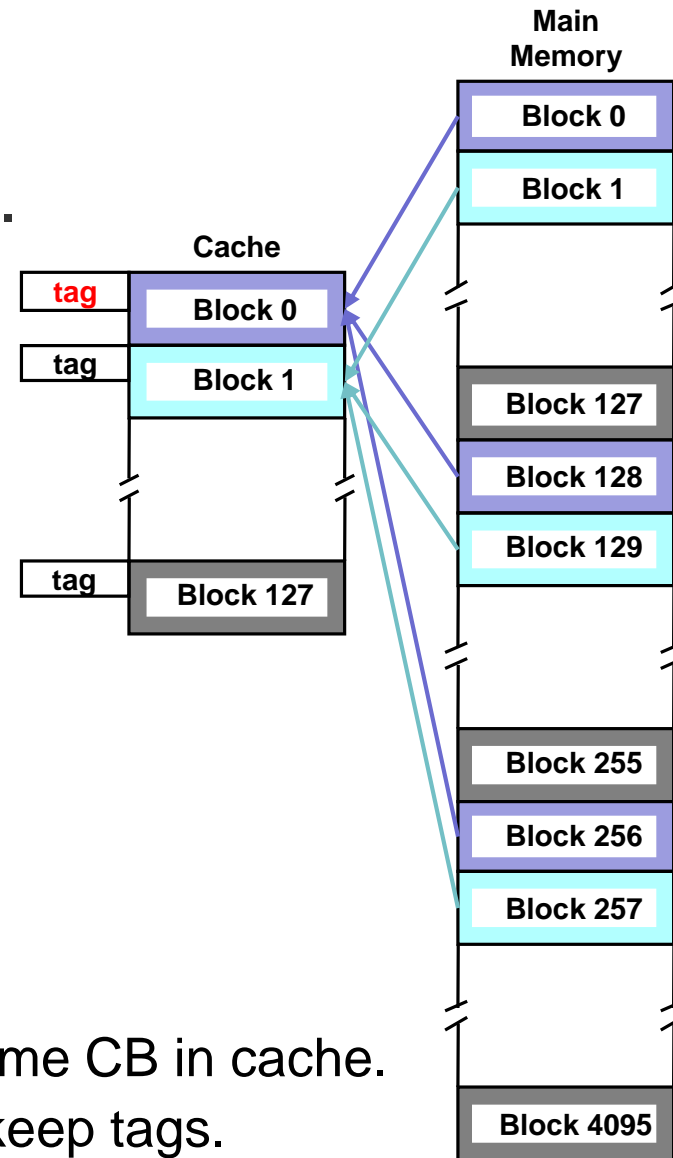


- **Direct Mapped Cache:**
Each Memory Block will be directly mapped to a Cache Block.

- **Direct Mapping Function:**

$$\text{MB } \#j \rightarrow \text{CB } \#(j \bmod 128)$$

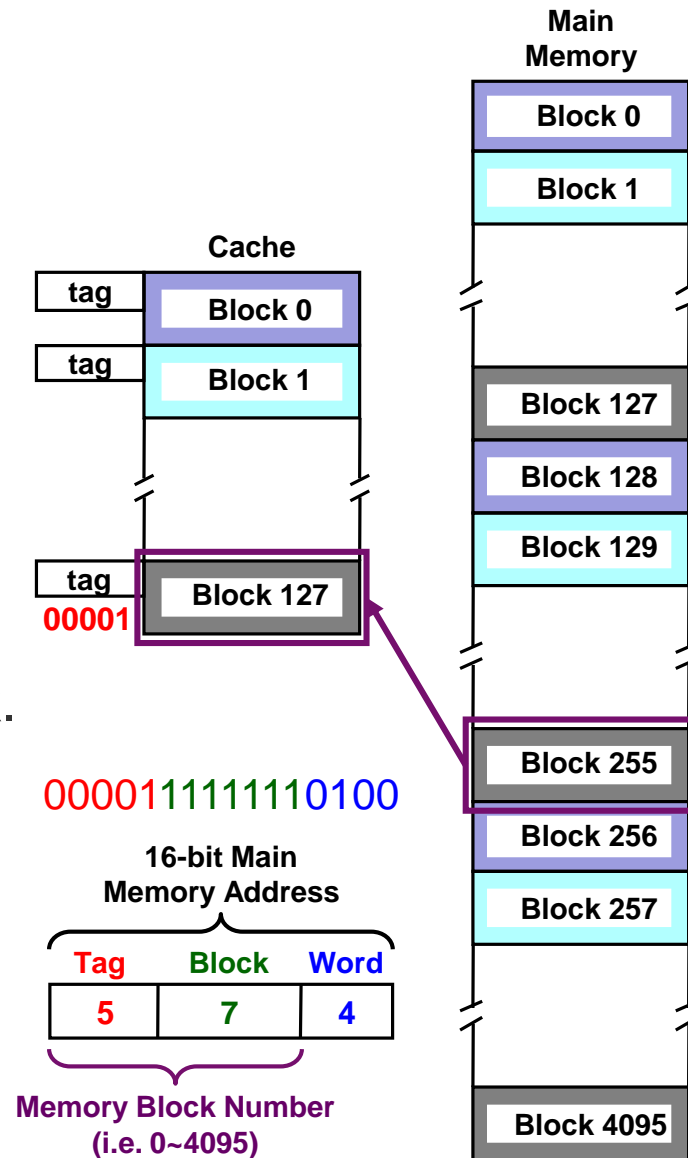
- **128**? There're 128 Cache Blocks.
- 32 MBs are mapped to 1 CB.
 - MBs **0, 128, 256, ..., 3968** \rightarrow CB **0**.
 - MBs **1, 129, 257, ..., 3969** \rightarrow CB **1**.
 - ...
 - MBs **127, 255, 383, ..., 4095** \rightarrow CB **127**.
- A **tag** is need for each CB.
 - Since many MBs will be mapped to a same CB in cache.
 - We need occupy some cache space to keep tags.



Direct Mapping (3/4)



- **Trick:** Interpret the 16-bit main memory address as three fields:
 - **Tag:** Keep track of which MB is placed in the corresponding CB.
 - **5** bits: $16 - (7 + 4) = 5$ bits.
 - **Block:** Determine the CB in cache.
 - **7** bits: There're $128 = 2^7$ cache blocks.
 - **Word:** Select one word in a block.
 - **4** bits: There're $16 = 2^4$ words in a block.
- Ex: CPU is looking for $(0FF4)_{16}$
 - MAR = $(0000111111110100)_2$
 - MB = $(000011111111)_2 = (255)_{10}$
 - CB = $(1111111)_2 = (127)_{10}$
 - Tag = $(00001)_2$

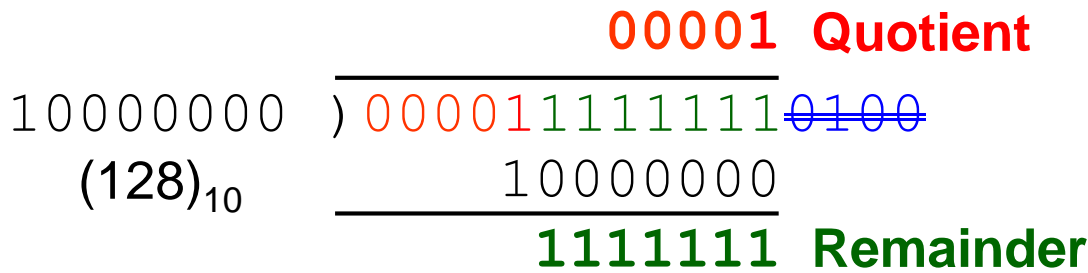


Direct Mapping (4/4)

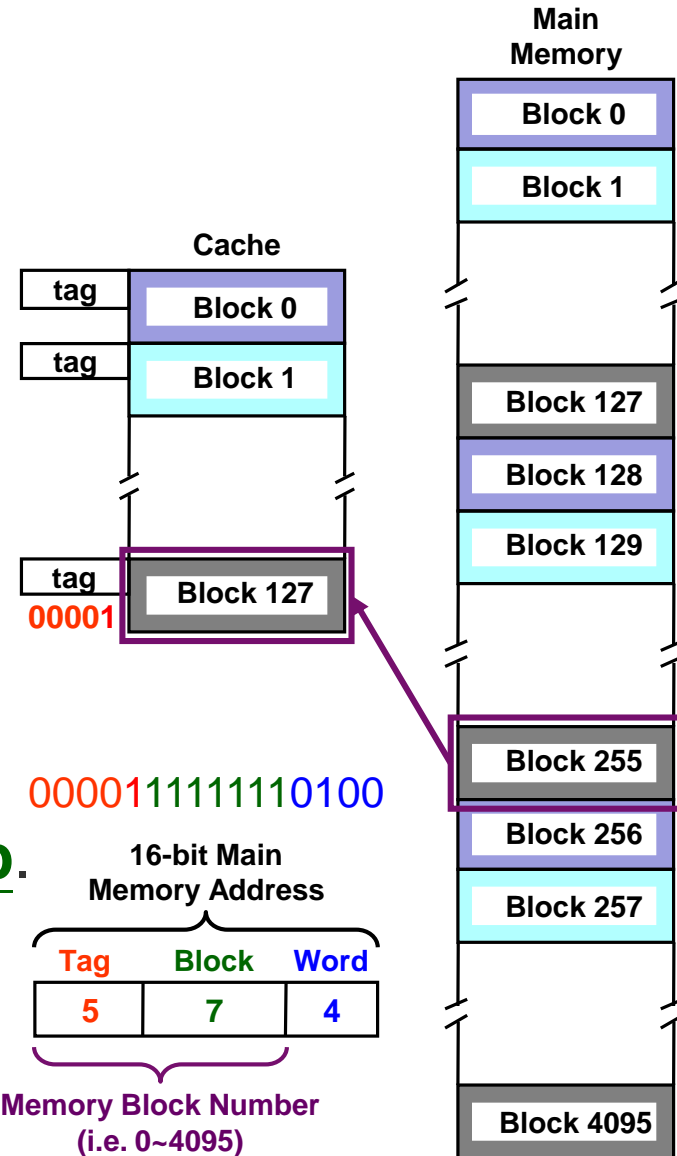


- Why the first 5 bits for **tag**? And why the middle 7 bits for **block**?

$$\text{MB } \#j \rightarrow \text{CB } \#(j \bmod 128)$$



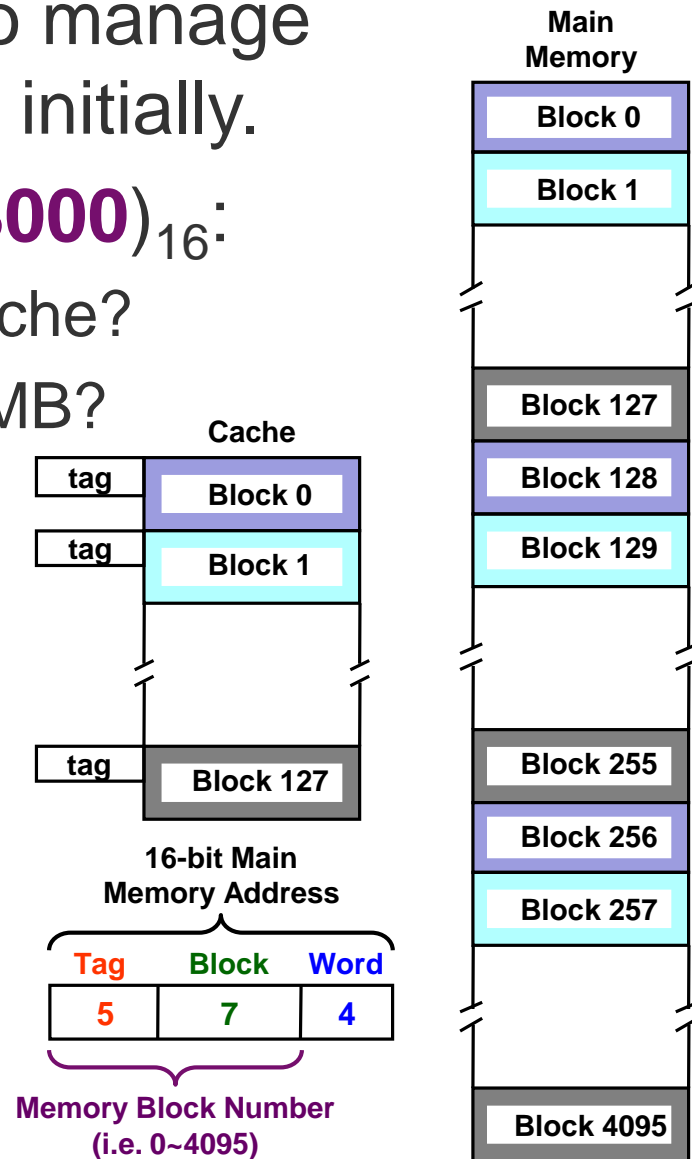
- Given a 16-bit address (**t**, **b**, **w**):
 - See if MB (**t**, **b**) is already in CB **b** by comparing **t** with the **tag** of CB **b**.
 - If not, replace CB **b** with MB (**t**, **b**) and update **tag** of CB **b** using **t**.
 - Finally access the word **w** in CB **b**.



Class Exercise 7.2



- Assume **direct mapping** is used to manage the cache, and all CBs are empty initially.
- Considering CPU is looking for $(8000)_{16}$:
 - Which MB will be loaded into the cache?
 - Which CB will be used to store the MB?
 - What is the new tag for the CB?

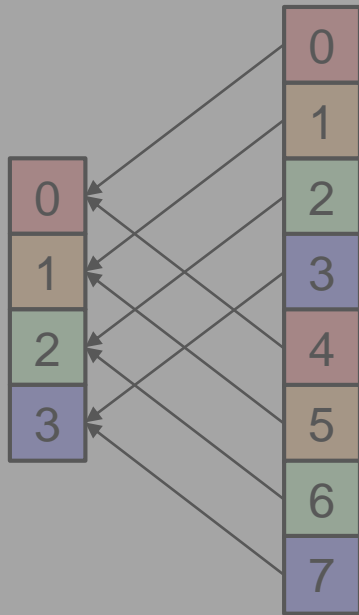


Associative Mapping (1/3)



Direct

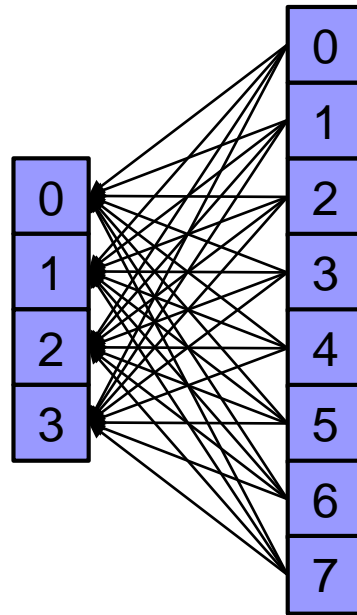
- A Memory Block is directly mapped (%) to a Cache Block.



Cache Blocks Memory Blocks

Associative

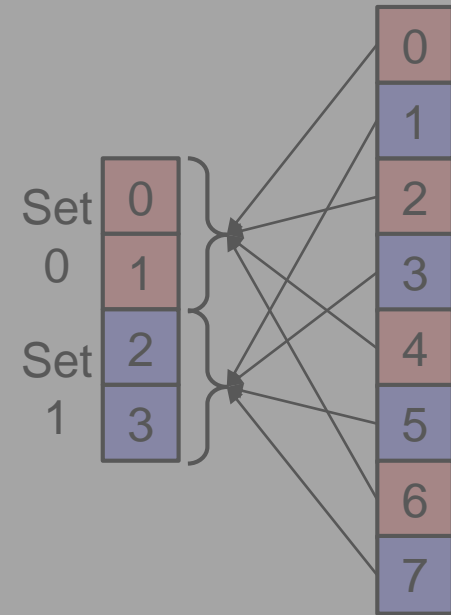
- A Memory Block can be mapped to any Cache Block.
(First come first serve!)



Cache Blocks Memory Blocks

Set Associative

- A Memory Block is directly mapped (%) to a Cache Set.
(In a set? Associative!)

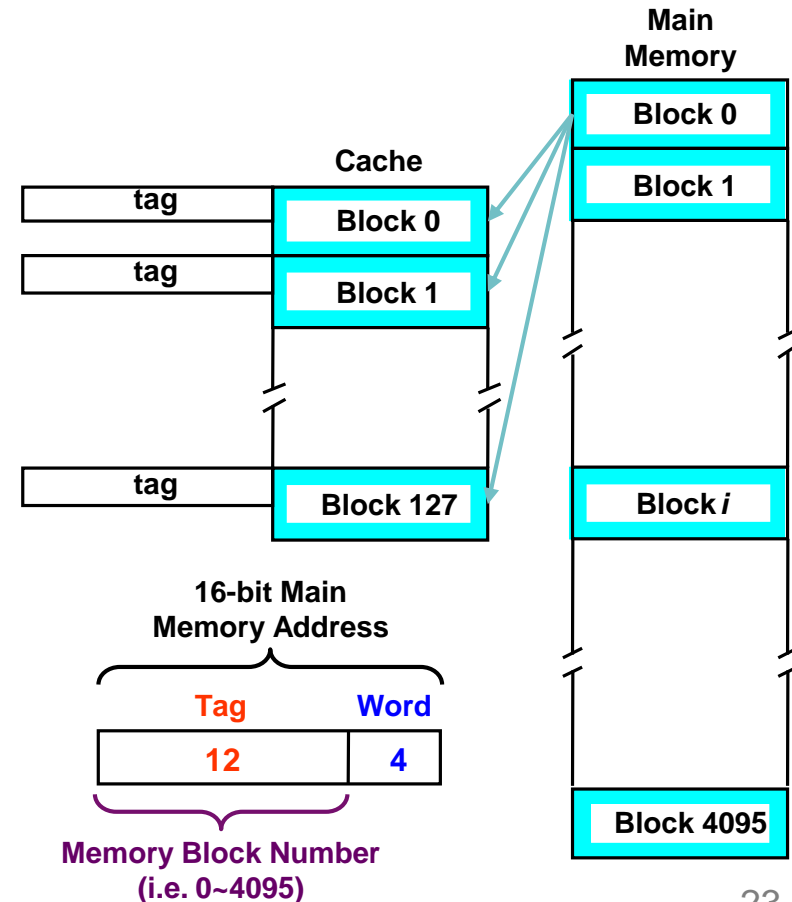


Cache Blocks Memory Blocks

Associative Mapping (2/3)



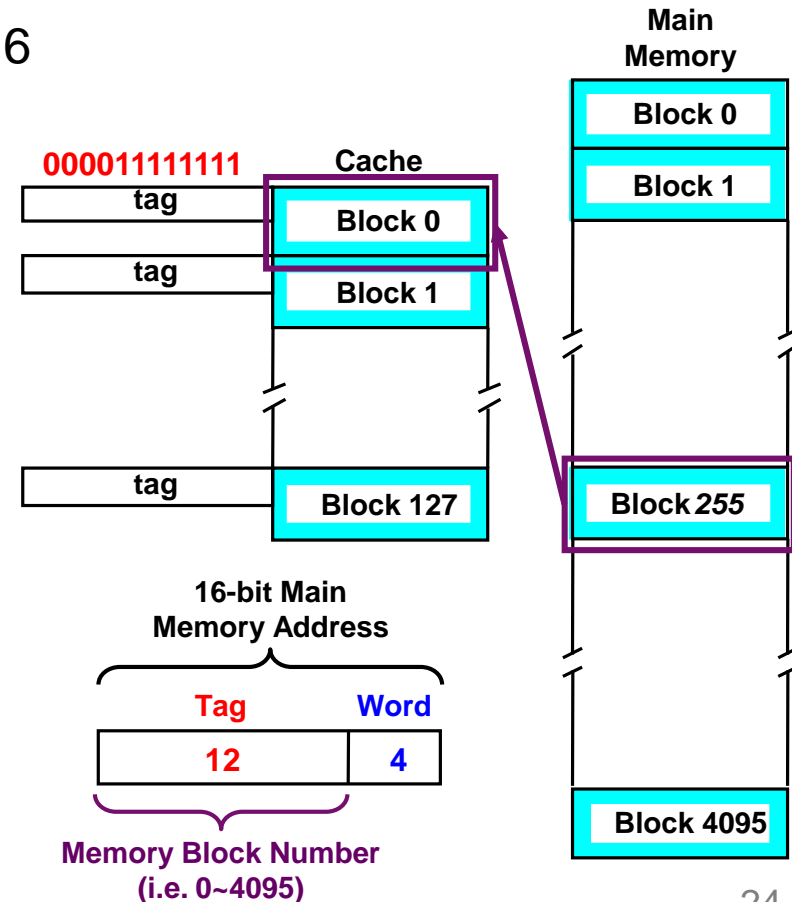
- **Direct Mapping:** A MB is restricted to a particular CB determined by mod operation.
- **Associative Mapping:**
Allow a MB to be mapped to any CB in the cache.
- **Trick:** Interpret the 16-bit main memory address as two fields:
 - **Tag:** The first **12** bits (i.e. the **MB number**) are all used to represent a MB.
 - **Word:** The last **4** bits for selecting a word in a block.



Associative Mapping (3/3)



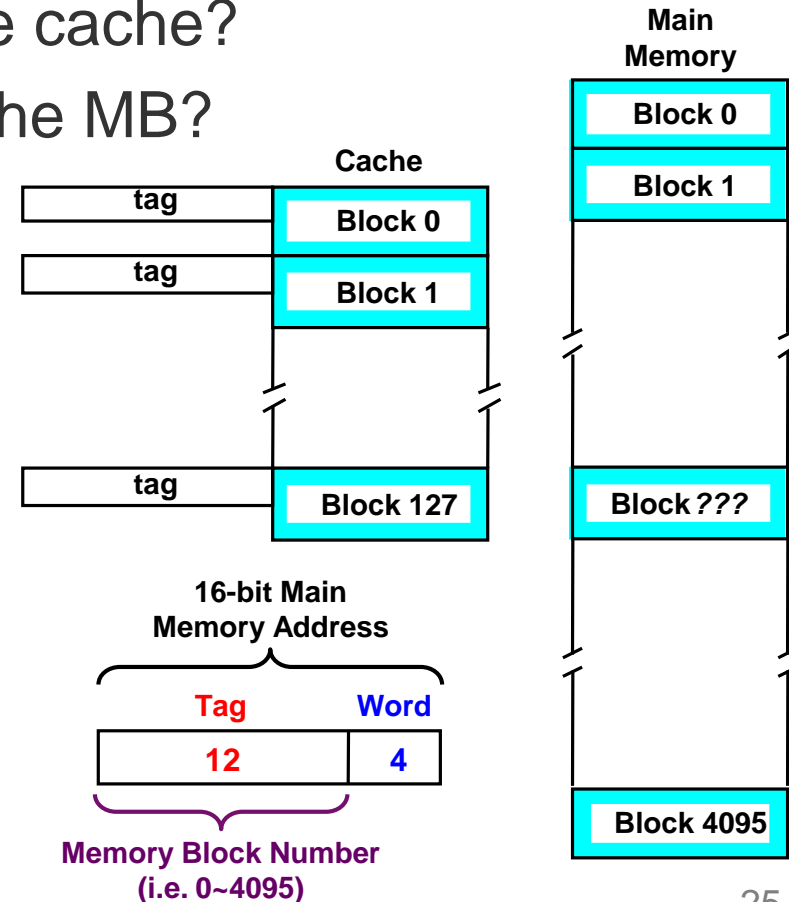
- How to determine the CB?
 - There's no pre-determined CB for any MB.
 - All CBs are used in the **first-come-first-serve (FCFS)** basis.
- Ex: CPU is looking for $(0FF4)_{16}$
 - Assume all CBs are empty.
 - $MAR = (0000111111110100)_2$
 - $MB = (000011111111)_2 = (255)_{10}$
 - $Tag = (000011111111)_2$
- Given a 16-bit address (**t**, **w**):
 - ALL tags of 128 CBs must be compared with **t** to see whether MB t is currently in the cache.
 - It can be done in parallel by HW.



Class Exercise 7.3



- Assume **associative mapping** is used to manage the cache, and all CBs are empty initially.
- Considering CPU is looking for $(8000)_{16}$:
 - Which MB will be loaded into the cache?
 - Which CB will be used to store the MB?
 - What is the new tag for the CB?

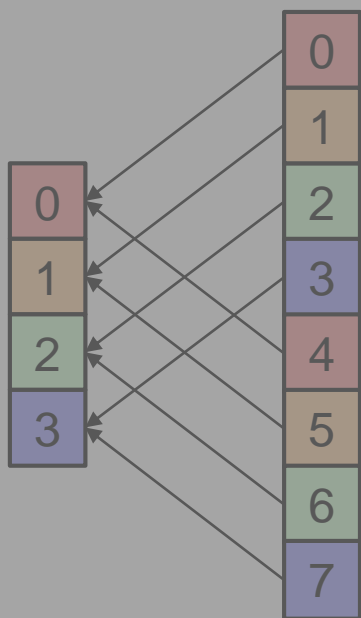


Set Associative Mapping (1/3)



Direct

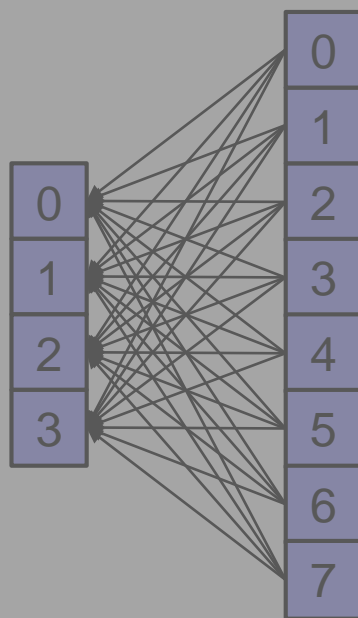
- A Memory Block is directly mapped (%) to a Cache Block.



Cache Blocks Memory Blocks

Associative

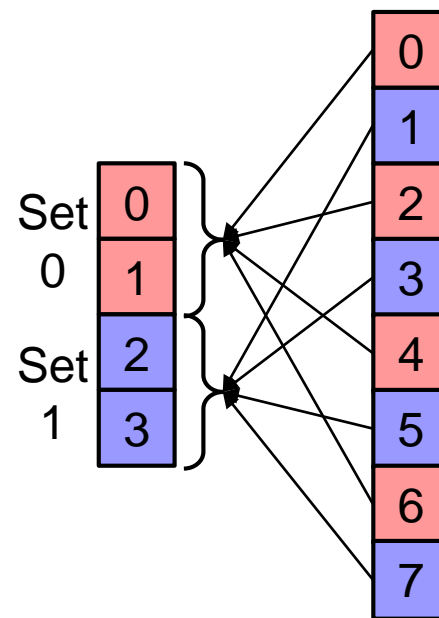
- A Memory Block can be mapped to any Cache Block. (First come first serve!)



Cache Blocks Memory Blocks

Set Associative

- A Memory Block is directly mapped (%) to a Cache Set. (In a set? Associative!)

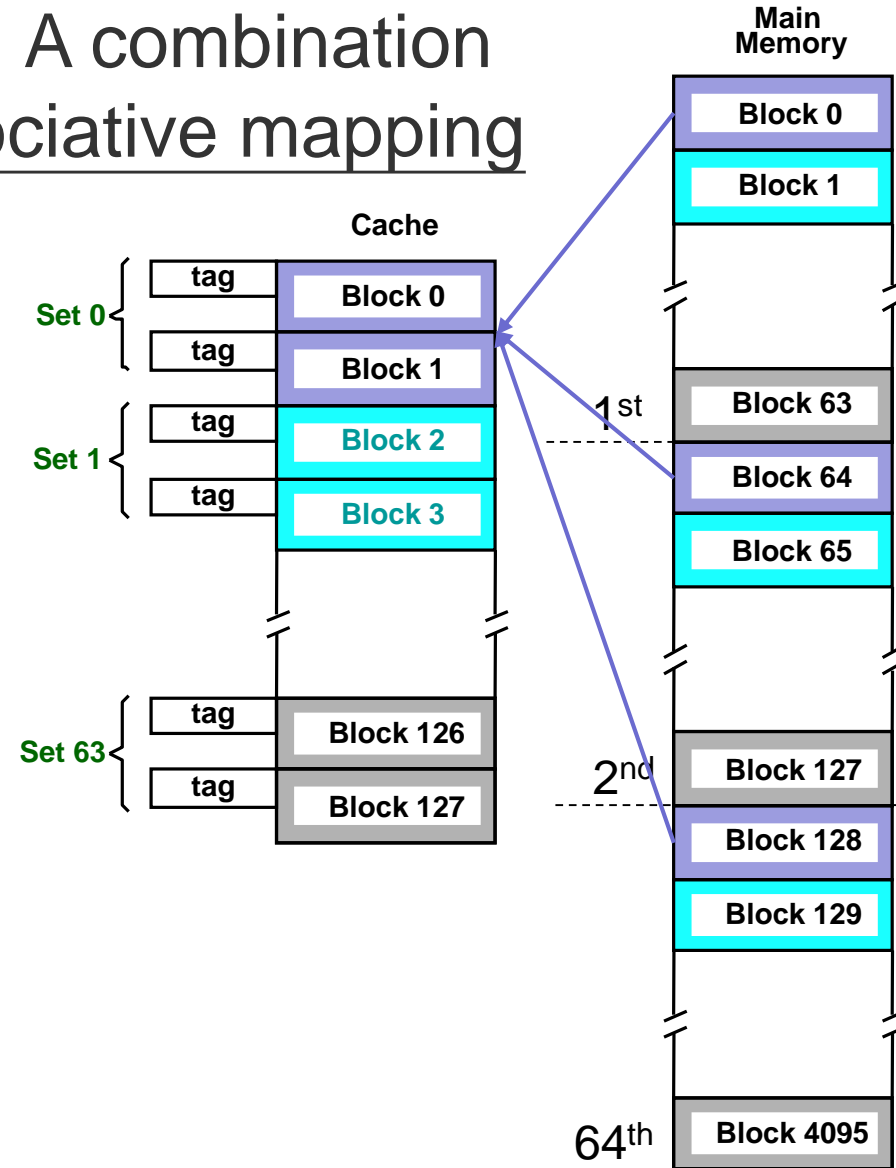


Cache Blocks Memory Blocks

Set Associative Mapping (2/3)



- **Set Associative Mapping:** A combination of direct mapping and associative mapping
 - **Direct:** First map a MB to a cache set (instead of a CB)
 - **Associative:** Then map to any CB in the cache set
- **K-way Set Associative:** A cache set is of *k* CBs.
 - Ex: **2**-way set associative
 - $128 \div 2 = 64$ (*sets*)
 - For MB #*j*, (*j mod 64*) derives the **Set** number.
 - E.g. MBs 0, 64, 128, ..., 4032
→ Cache Set #0.



Set Associative Mapping (3/3)



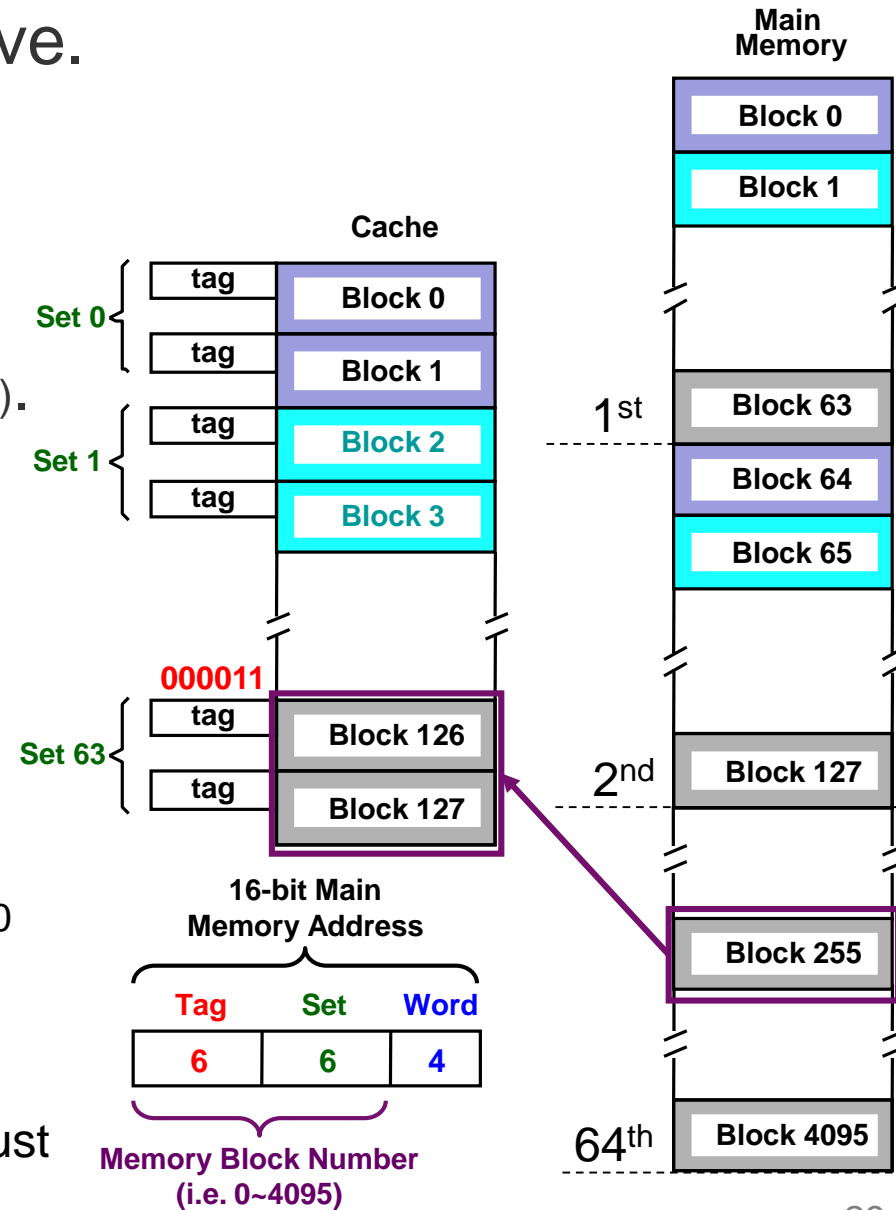
- Consider 2-way set associative.
- Trick:** Interpret the 16-bit address as three fields:

- Tag:** The first 6 bits (**quotient**).
- Set:** The middle 6 bits (**remainder**).
 - 6 bits: There're 2^6 cache sets.
- Word:** The last 4 bits.

Ex: CPU is looking for $(0FF4)_{16}$

- Assume all CBs are empty.
- MAR = $(0000111111110100)_2$
- MB = $(000011111111)_2 = (255)_{10}$
- Cache Set = $(111111)_2 = (63)_{10}$
- Tag = $(000011)_2$

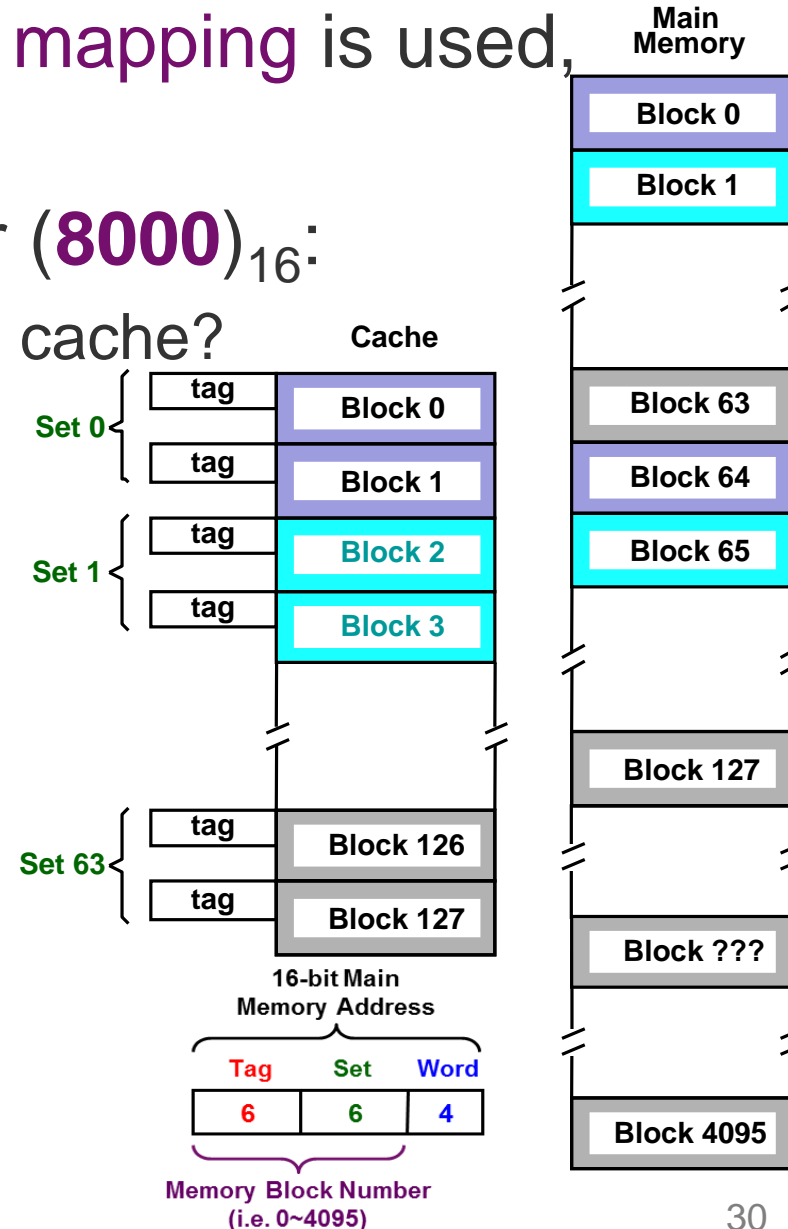
Note: **ALL tags** of CBs in a cache set must be compared (done in parallel by HW).



Class Exercise 7.4



- Assume **2-way set associative mapping** is used, and all CBs are empty initially.
- Considering CPU is looking for $(8000)_{16}$:
 - Which MB will be loaded into the cache?
 - Which CB will store the MB?
 - What is the new tag for the CB?

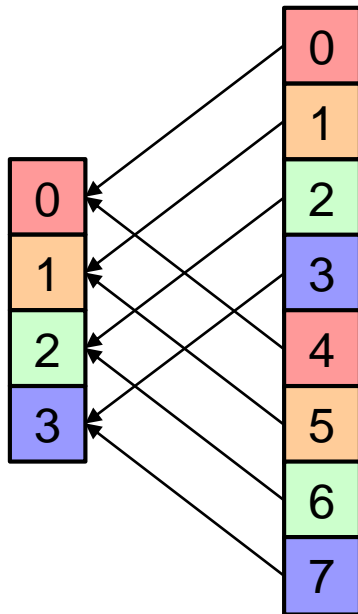


Summary of Mapping Functions (1/2)



Direct

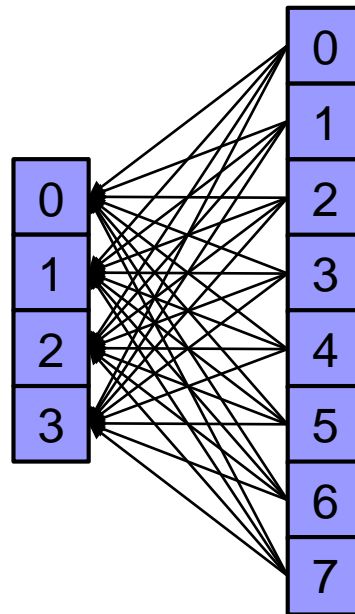
A Memory Block is directly mapped (%) to a Cache Block.



Cache Blocks Memory Blocks

Associative

A Memory Block can be mapped to any Cache Block.
(First come first serve!)

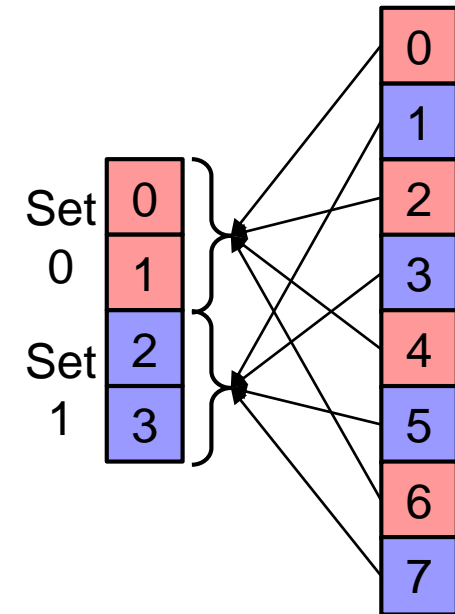


Cache Blocks Memory Blocks

Set Associative

A Memory Block is directly mapped (%) to a Cache Set.

In a Set? Associative!

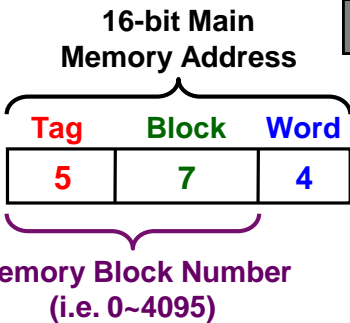
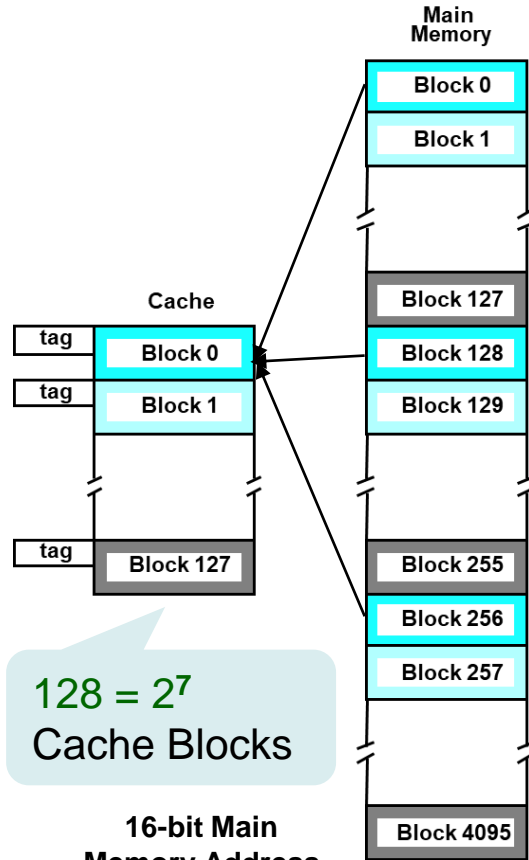


Cache Blocks Memory Blocks

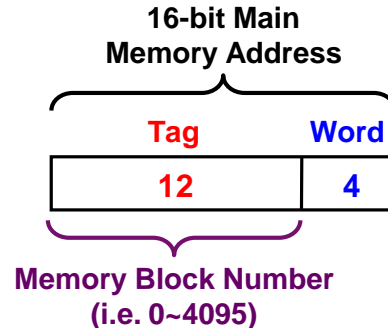
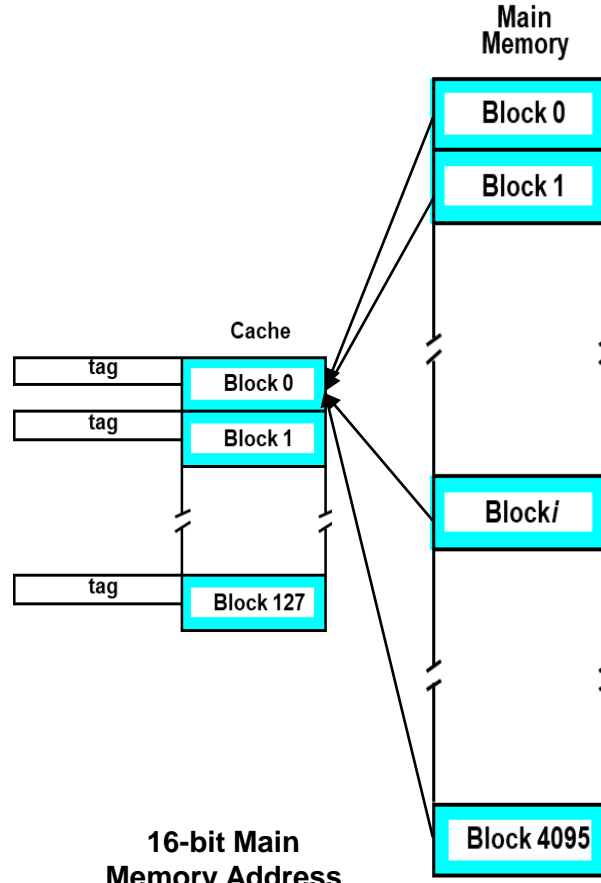
Summary of Mapping Functions (2/2)



Direct

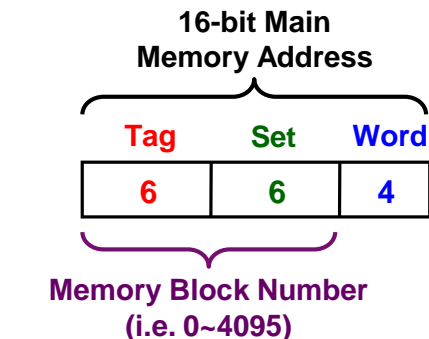
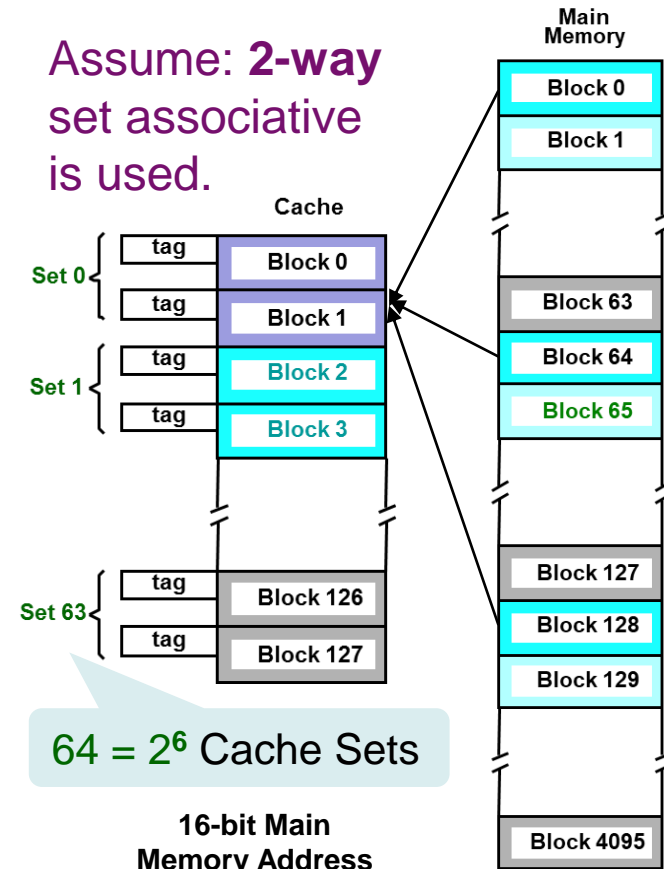


Associative



Set Associative

Assume: 2-way set associative is used.





- Cache Basics
- Mapping Functions
 - Direct Mapping
 - Associative Mapping
 - Set Associative Mapping
- Replacement Algorithms
 - Least Recently Used (LRU) Replacement
 - Random Replacement
- Working Examples

Replacement Algorithms

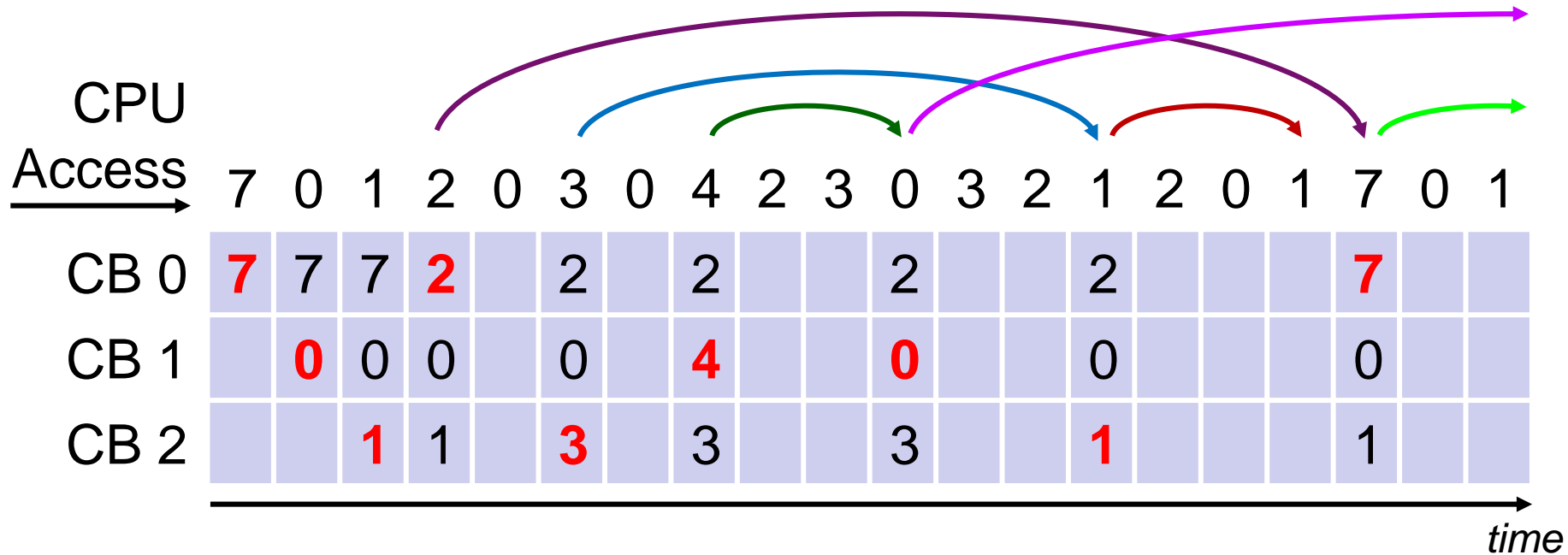


- **Replace:** **Write Back** (to old MB) & **Overwrite** (with new MB)
- **Direct Mapped Cache:**
 - The CB is pre-determined directly by the memory address.
 - The replacement strategy is trivial: Just replace the pre-determined CB with the new MB.
- **Associative and Set Associative Mapped Cache:**
 - Not trivial: Need to determine which block to replace.
 - **Optimal Replacement**: Always keep CBs, which will be used sooner, in the cache, if we can look into the future (**not practical!!!**).
 - **Least recently used (LRU)**: Replace the block that has gone the longest time without being accessed by looking back to the past.
 - Rationale: Based on temporal locality, CBs that have been referenced recently will be most likely to be referenced again soon.
 - **Random Replacement**: Replace a block randomly.
 - Easier to implement than LRU, and quite effective in practice.

Optimal Replacement Algorithm



- **Optimal Algorithm:** Replace the CB that will not be used for the longest period of time (in the future).
- Given an **associative mapped cache**, which is composed of 3 Cache Blocks (CBs 0~2).

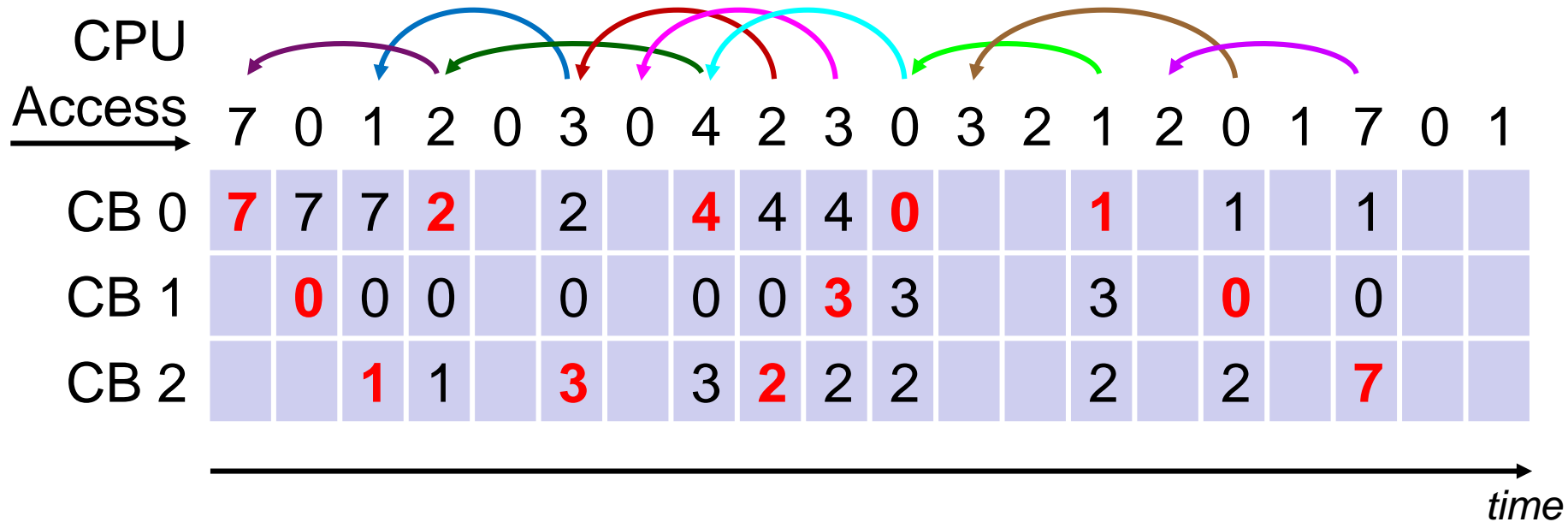


– The optimal algorithm causes **9** times of cache misses.

LRU Replacement Algorithm



- **LRU Algorithm:** Replace the CB that has not been used for the longest period of time (in the **past**).
- Given an **associative mapped cache**, which is composed of 3 Cache Blocks (CBs 0~2).

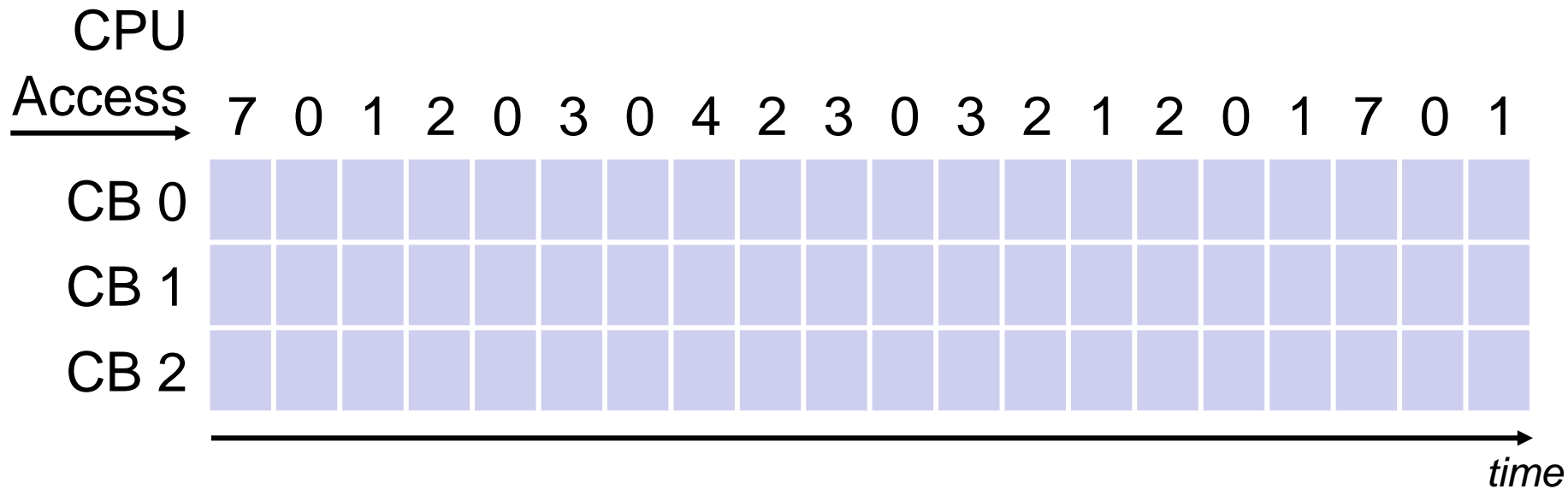


– The LRU algorithm causes **12** times of cache misses.

Class Exercise 7.5



- **First-In-First-Out Algorithm:** Replace the CB that has arrived for the longest period of time (in the **past**).
- Given an **associative mapped cache**, which is composed of 3 Cache Blocks (CBs 0~2).
- Please fill in the cache and state **cache misses**.





- Cache Basics
- Mapping Functions
 - Direct Mapping
 - Associative Mapping
 - Set Associative Mapping
- Replacement Algorithms
 - Least Recently Used (LRU) Replacement
 - Random Replacement
- Working Examples

Summary



- Cache Basics
- Mapping Functions
 - Direct Mapping
 - Associative Mapping
 - Set Associative Mapping
- Replacement Algorithms
 - Least Recently Used (LRU) Replacement
 - Random Replacement
- Working Examples