

#### 香港中文大學

The Chinese University of Hong Kong

# CSCI2510 Computer Organization

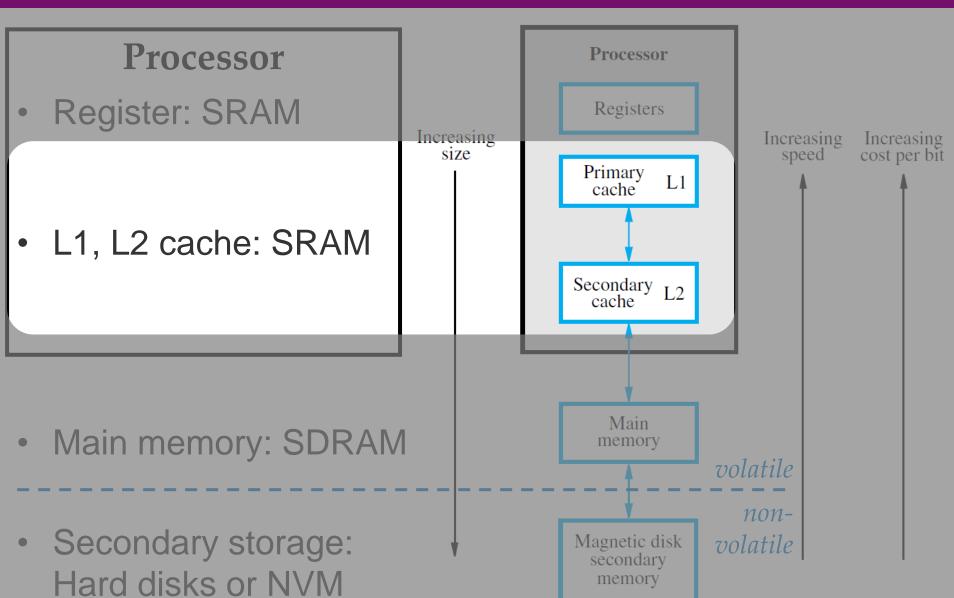
#### Lecture 07: Cache in Action

# Ming-Chang YANG mcyang@cse.cuhk.edu.hk

Reading: Chap. 8.6

### **Recall: Memory Hierarchy**





#### **Outline**

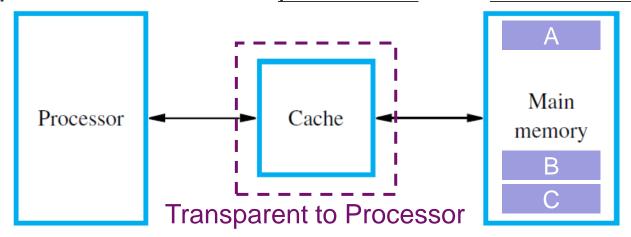


- Cache Basics
- Mapping Functions
  - Direct Mapping
  - Associative Mapping
  - Set Associative Mapping
- Replacement Algorithms
  - Optimal Replacement
  - Least Recently Used (LRU) Replacement
  - Random Replacement
- Working Examples

#### 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**

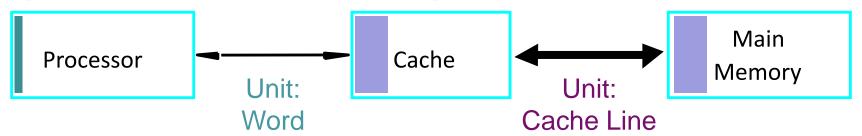


- Temporal Locality (locality in time)
  - If an item is referenced, it will tend to be referenced again soon (e.g. recent calls).
  - Strategy: When information item (instruction or data) is first needed, opportunistically bring it into cache (we hope it will be used soon).
- Spatial Locality (locality 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.



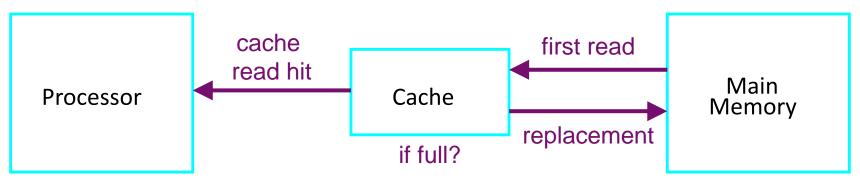
- Cache Block / Line: The unit composed of multiple successive memory words (size: cache block > word).
  - The contents of a cache block (of memory words) will be loaded into or unloaded from the cache at a time.
- Mapping Functions: Decide how cache is organized and how addresses are mapped to the 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 a cached block and to create space for the new one if cache is full

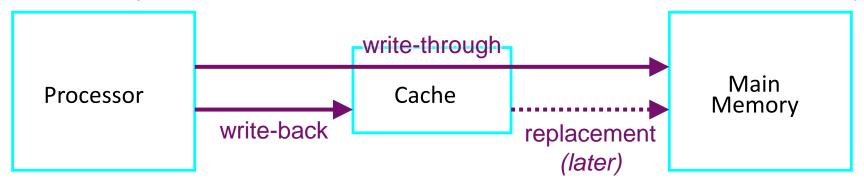


### **Write Operation in Cache**



#### 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?

#### **Outline**

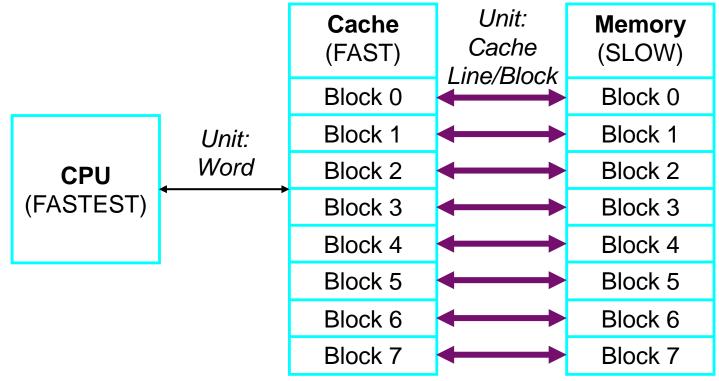


- Cache Basics
- Mapping Functions
  - Direct Mapping
  - Associative Mapping
  - Set Associative Mapping
- Replacement Algorithms
  - Optimal Replacement
  - Least Recently Used (LRU) Replacement
  - Random Replacement
- Working Examples

### **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?

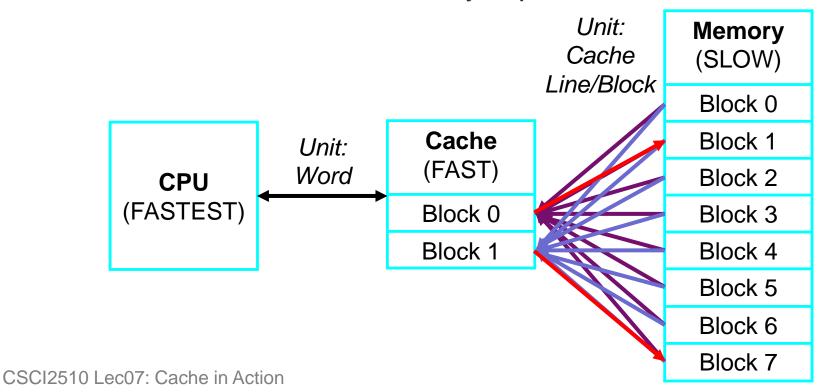


Trivial! One-to-one mapping is enough!

# **Mapping Functions (2/3)**



- Reality: The cache size is much smaller (<<<) than
  the main memory size.</li>
- Many-to-one mapping is needed!
  - Many blocks in memory compete for one block in cache.
  - One block in cache can only represent one block in memory.



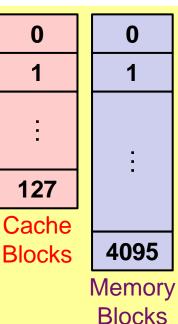
# **Mapping Functions (3/3)**



- 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
    - Note: A cache block is of successive memory words.
  - $1 \text{ Word} = 16 \text{ bits} = 2^1 \text{ Bytes}$
  - -1 Block = 8 Words =  $2^3$  Words
  - Cache Size: 2K Bytes → 128 Cache Blocks
    - Cache Block (CB): The block in the cache.
  - **Memory Size**: 16-bit Address  $\rightarrow$  2<sup>16</sup> = 64K Bytes

→ 4096 **Memory Blocks** 

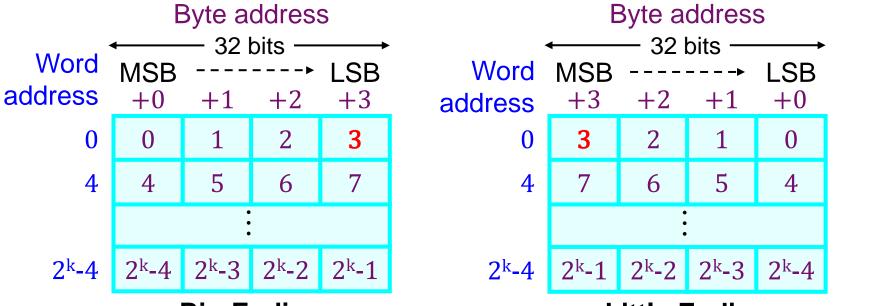
• Memory Block (MB): The block in the main memory.



#### Recall: Big-Endian and Little-Endian

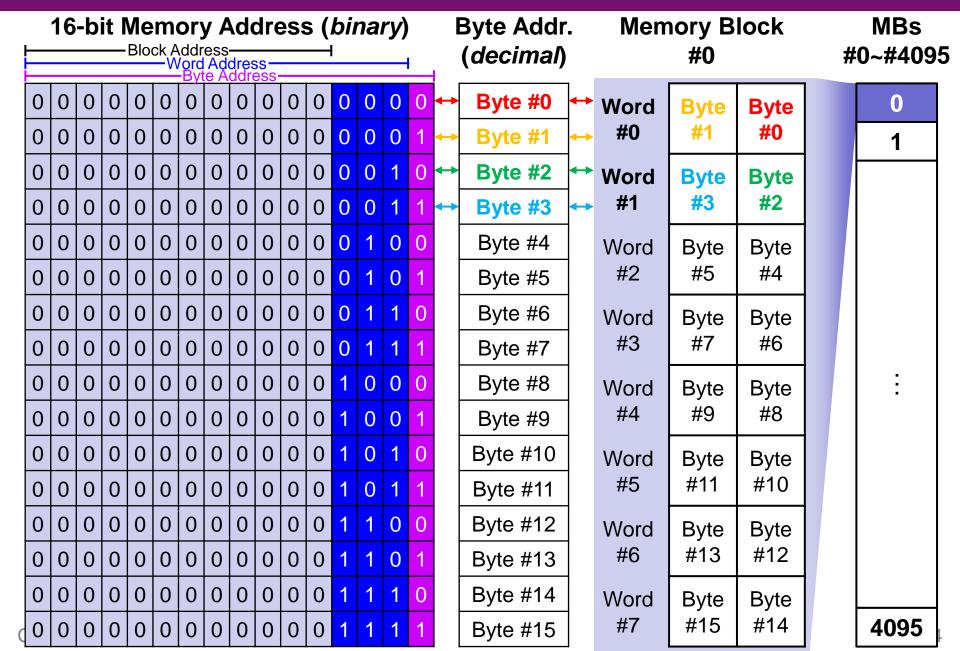


- Two ways to order byte addresses across a word:
  - Big-Endian: Bytes within a word are ordered left-to-right, and lower byte addresses are used for more significant bytes of a multi-byte data (e.g. Motorola).
  - Little-Endian: Bytes within a word are ordered right-to-left, and lower byte addresses are used for less significant bytes of a multi-byte data (e.g. Intel).



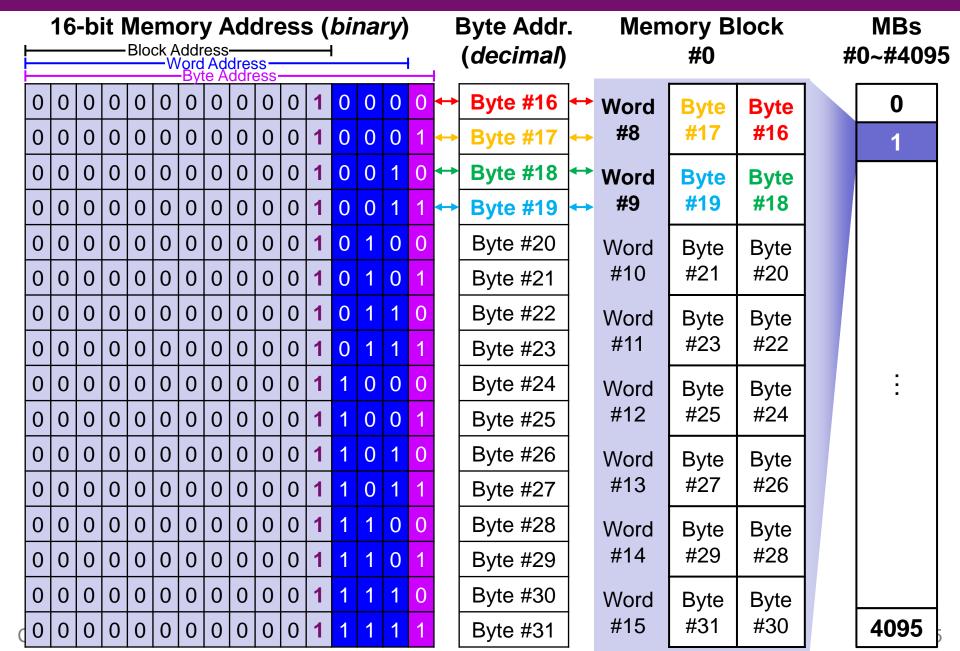
### **Example: Memory Block #0**





### **Example: Memory Block #1**





# **Example: Memory Block #4095**

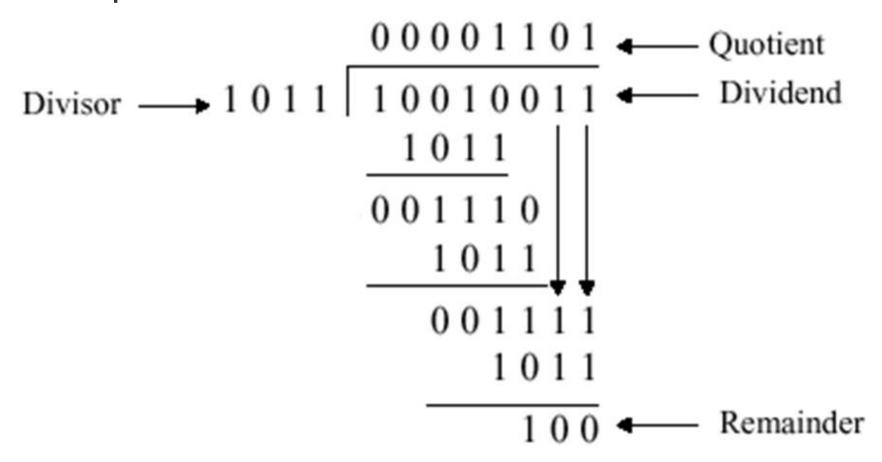


16-bit Memory Address (binary)  Block Address Word Address Byte Address							Ī				nory B #0	y Block 0		MBs #0~#4095										
1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0		B#65520		Word	Byte	Byte			0
1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	1		B#65521		#32760	#65525	#65520			1
1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	0		B#65522		Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1		B#65523		#32761	#65525	#65522			
1	1	1	1	1	1	1	1	1	1	1	1	0	1	0	0		B#65524		Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	0	1	0	1		B#65525		#32762	#65525	#65524			
1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0		B#65526		Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1		B#65527		#32763	#65527	#65526			
1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0		B#65528		Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	1		B#65529		#32764	#65529	#65528			
1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	0		B#65530		Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1		B#65531		#32765	#65531	#65530			
1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	<b>←→</b>	B#65532	<b>+</b>	Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	$\leftrightarrow$	B#65533	<b>+</b>	#32766	#65533	#65532			
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	<b>↔</b>	B#65534	<b>+</b>	Word	Byte	Byte			
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	$\leftrightarrow$	B#65535	<b>*</b>	#32767	#65535	#65534			4095

### Prior Knowledge: Modulo Operator



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



#### Class Exercise 7.1

Student ID:	Date:
Name:	

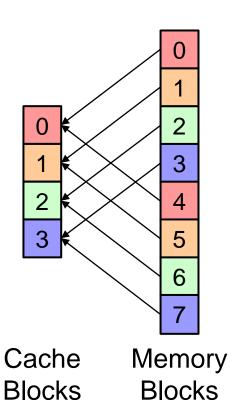
• Considering the previous example, what if the divisor equals to  $(10)_2$ ,  $(100)_2$ , ...,  $(10000000)_2$ ?

# Direct Mapping (1/4)



#### **Direct**

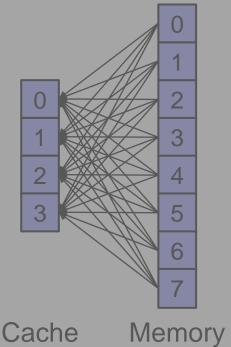
•A Memory Block is directly mapped (%) to a Cache Block.



#### **Associative**

•A Memory Block can be mapped to any Cache Block.

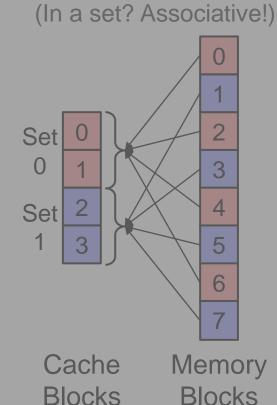
(First come first serve!)



Cache Memory Blocks

#### **Set Associative**

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



# Direct Mapping (2/4)



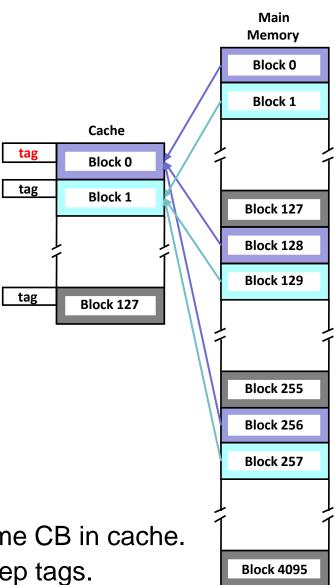
Direct Mapped Cache:

 Each Memory Block will be
 directly mapped to a Cache Block.

#### Direct Mapping Function:

 $MB \# j \rightarrow CB \# (j \mod 128)$ 

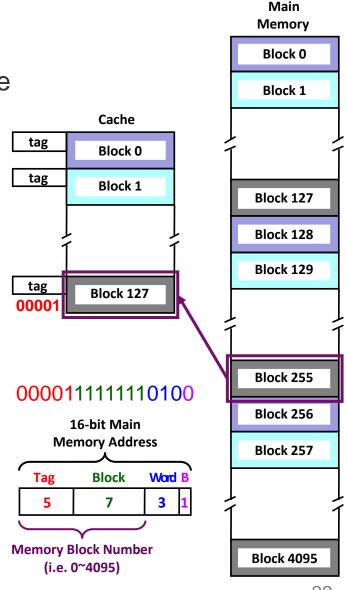
- 128? There're 128 Cache Blocks.
- 32 MBs are mapped to 1 CB.
  - MBs 0, 128, 256, ..., 3968 → CB 0.
  - MBs 1, 129, 257, ..., 3969 → CB 1.
  - ...
  - MBs **127**, **255**, **383**, ..., **4095** → CB **127**.
- A tag is need for each CB.
  - Since many MBs will be mapped to a same CB in cache.
  -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



# Direct Mapping (3/4)



- Trick: Interpret the 16-bit main memory address as follows:
  - 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 = 27 cache blocks.
  - Word: Select one word in a block.
    - 3 bits: There're 8 = 23 words in a block.
  - Byte: Select one byte in a word.
    - 1 bits: There're 2 = 21 bytes in a word.
- Ex: CPU is looking for (0FF4)<sub>16</sub>
  - $MAR = (0000 1111 1111 0100)_2$
  - MB =  $(0000 \ 1111 \ 1111)_2 = (255)_{10}$
  - CB =  $(11111111)_2$  =  $(127)_{10}$
  - $Tag = (00001)_2$



# Direct Mapping (4/4)



Main Memory

Block 0

Block 1

**Block 127** 

Block 128

Block 129

Block 255

Block 256

Block 257

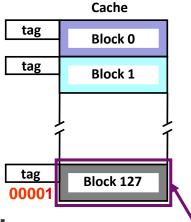
**Block 4095** 

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

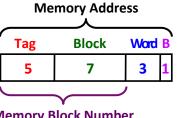
 $MB \#j \rightarrow CB \#(j \mod 128)$ 

1111111 Remainder

- Given a 16-bit address (t, b, w, b):
  - See if MB (t, b) is already in CB b
     by comparing t with the tag of CB b.
  - ② 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.





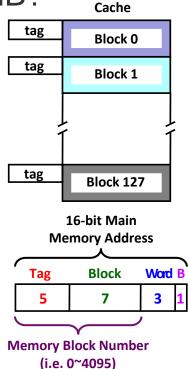


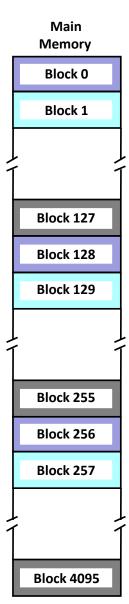
Memory Block Number (i.e. 0~4095)

#### Class Exercise 7.2



- Assume direct mapping is used to manage the cache, and all CBs are empty initially.
- Considering CPU is looking for (8010)<sub>16</sub>:
  - 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?



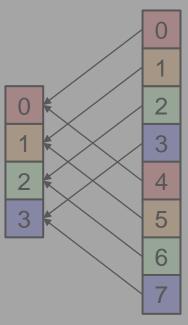


### **Associative Mapping (1/3)**



#### **Direct**

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

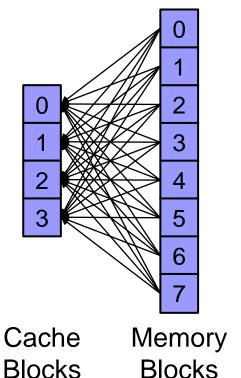


Cache Memory Blocks

#### **Associative**

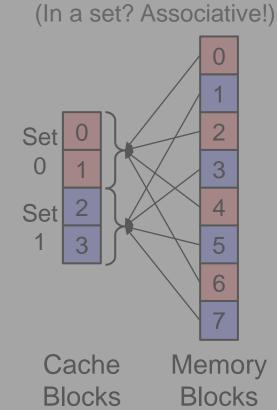
•A Memory Block can be mapped to any Cache Block.

(First come first serve!)



#### **Set Associative**

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



# Associative Mapping (2/3)



 Direct Mapping: A MB is restricted to a particular CB determined by mod operation.

#### Associative Mapping:

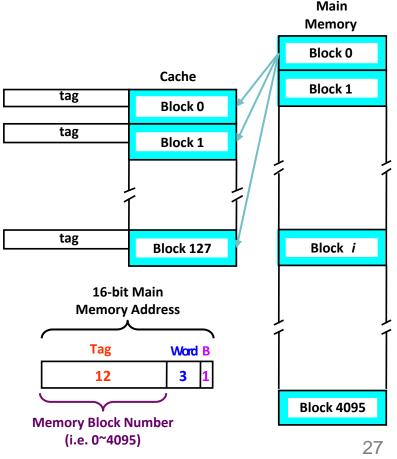
Allow a MB to be mapped to any CB in the cache.

 Trick: Interpret the 16-bit main memory address as follows:

 Tag: The first 12 bits (i.e. the MB) number) are all used to

represent a MB.

– Word & Byte: The last 3 & 1 bits for selecting a word & byte in a block.

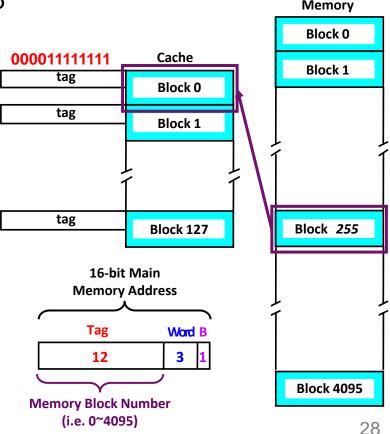


# **Associative Mapping (3/3)**



Main

- 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)<sub>16</sub>
  - Assume all CBs are empty.
  - $-MAR = (0000 1111 1111 0100)_2$
  - $-MB = (0000 1111 1111)_2 = (255)_{10}^{11}$
  - $Tag = (0000 1111 1111)_2$
- Given a 16-bit addr. (t, w, b):
  - 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



Main Memory

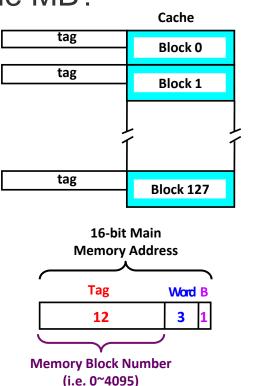
Block 0

Block 1

Block ???

**Block 4095** 

- Assume associative mapping is used to manage the cache, and all CBs are empty initially.
- Considering CPU is looking for (8010)<sub>16</sub>:
  - 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?

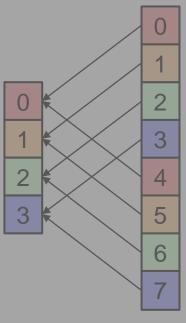


# **Set Associative Mapping (1/3)**



#### **Direct**

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

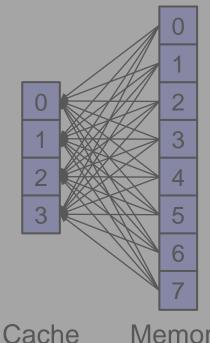


Cache Memory Blocks

#### **Associative**

•A Memory Block can be mapped to any Cache Block.

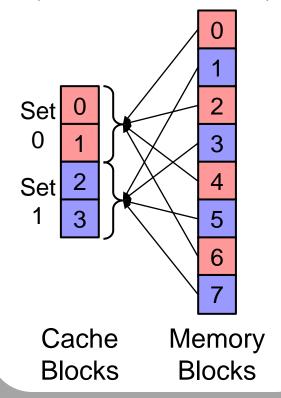
(First come first serve!)



Cache Memory Blocks

#### **Set Associative**

 A Memory Block is directly mapped (%) to a Cache <u>Set</u>. (In a set? Associative!)



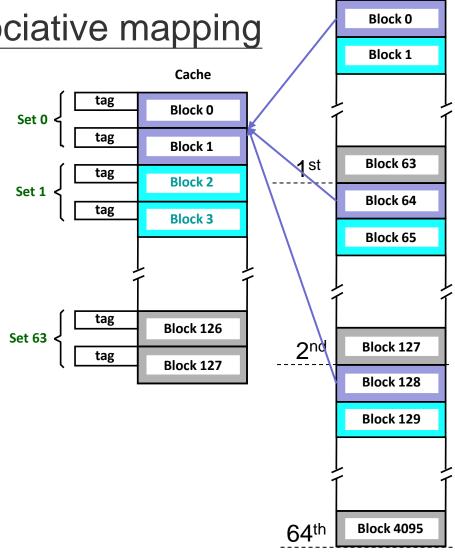
# **Set Associative Mapping (2/3)**



Main

Memory

- 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)**



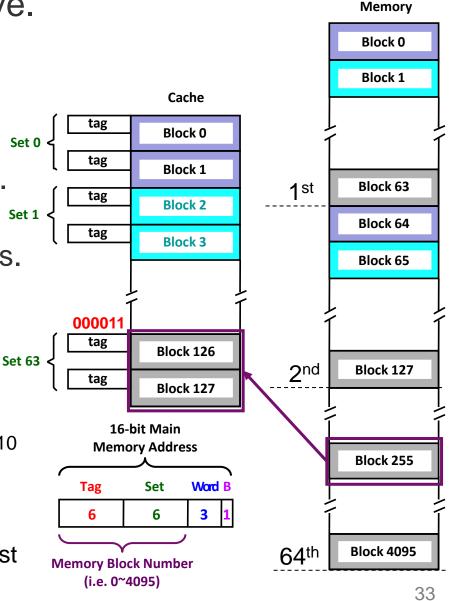
Main

- Consider 2-way set associative.
- Trick: Interpret the 16-bit address as follows:
  - Tag: The first 6 bits (quotient).
  - Set: The middle 6 bits (remainder).
    - 6 bits: There're 26 cache sets.
  - Word & Byte: The last 3 & 1 bits.

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

- Assume all CBs are empty.
- $MAR = (0000 1111 1111 0100)_2$
- $MB = (0000 1111 11111)_2 = (255)_{10}$
- Cache Set =  $(1111111)_2$  =  $(63)_{10}$
- $Tag = (000011)_2$

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



#### Class Exercise 7.4



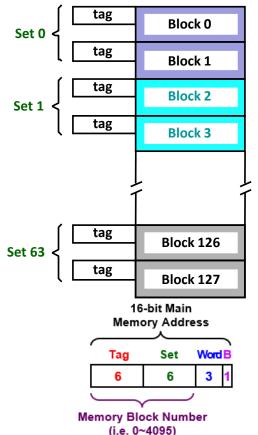
Main

 Assume 2-way set associative mapping is used, and all CBs are empty initially.

Considering CPU is looking for (8010)<sub>16</sub>:

— Which MB will be loaded into the cache?

- Which CB will store the MB?
- What is the new tag for the CB?



Memory Block 0 Block 1 Block 63 Block 64 Block 65 Block 127 Block ??? **Block 4095** 

# **Summary of Mapping Functions (1/2)**



#### **Direct**

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

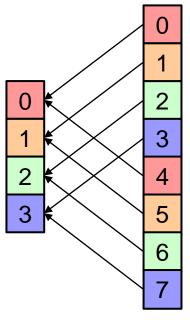
#### **Associative**

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

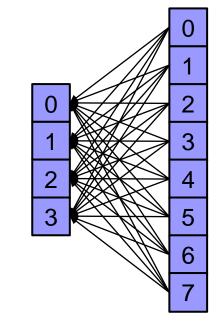
#### **Set Associative**

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

In a **Set**? **Associative**!

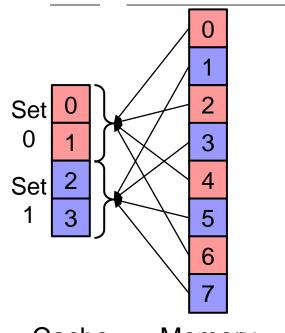


Cache Memory Blocks



Cache

Memory Blocks

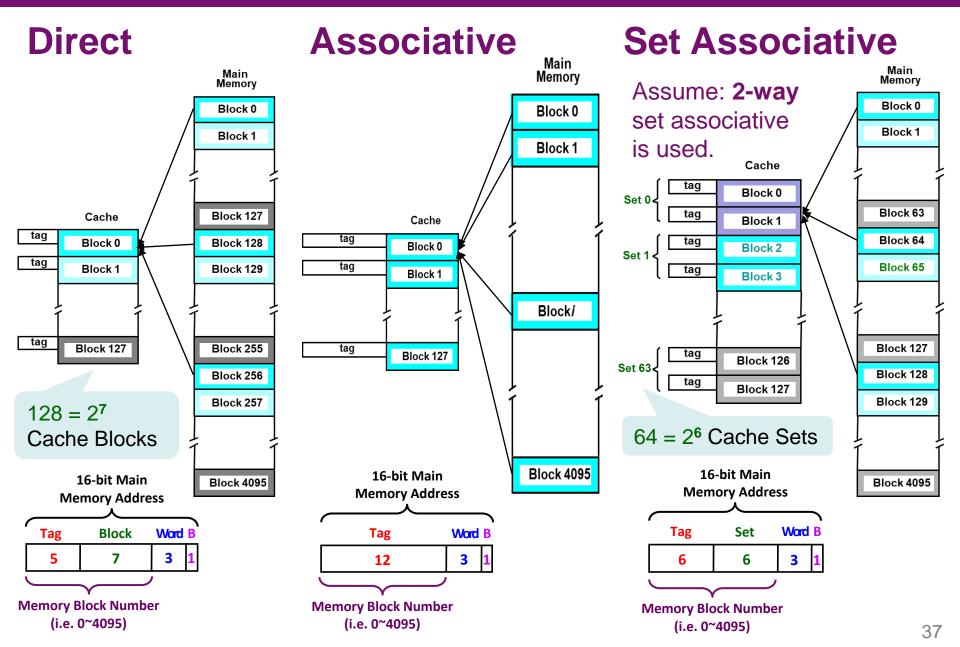


Cache

Memory Blocks

# Summary of Mapping Functions (2/2)





#### **Outline**



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

### Replacement Algorithms

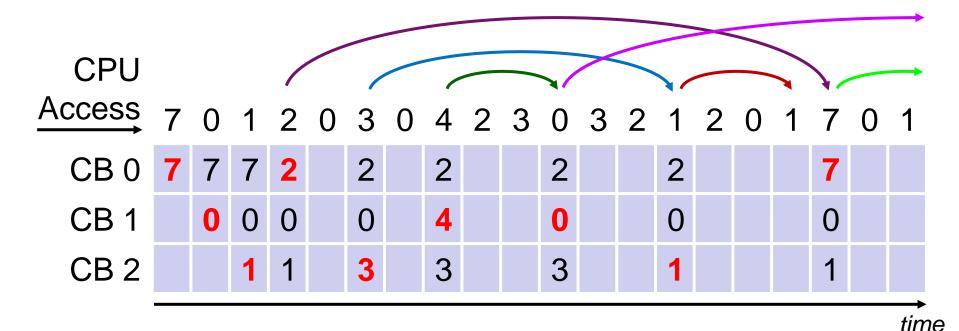


- Replace: Write Back (to old MB) & Overwrite (with new MB)
- Direct Mapped Cache:
  - The CB is pre-determined directly by the memory address.
  - The replacement strategy is trivial: Just replace the predetermined CB with the new MB.
- Associative and Set Associative Mapped Cache:
  - Not trivial: Need to determine which block to replace.
    - Optimal Replacement: Always keep CBs, which will <u>be used</u> sooner, in the cache, if we can <u>look into the future</u> (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 <u>temporal locality</u>, 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).

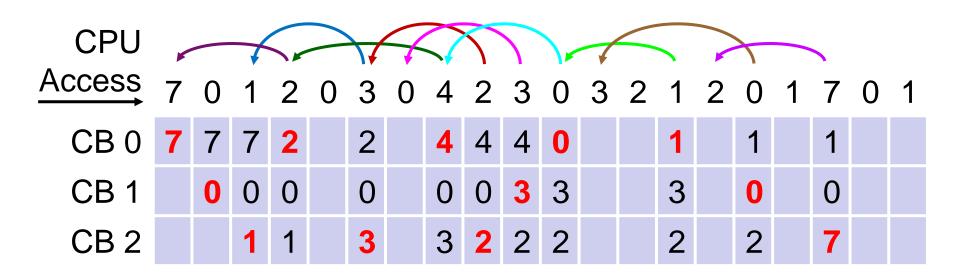


- The optimal algorithm causes 9 times of cache misses.

# 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).

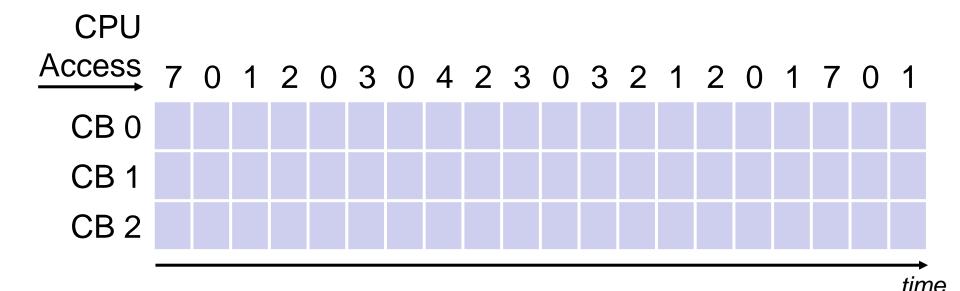


- The LRU algorithm causes 12 times of cache misses.

time



- 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.



#### **Outline**



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

## Cache Example



- Cache Configuration:
  - Cache has 8 blocks.
  - A block contains one word.
  - A word is of 16 bits.

```
short A[10][4];
int sum = 0;
int j, i;
double mean;
// 1) forward loop
for (j = 0; j \le 9; j++)
  sum += A[j][0];
mean = sum / 10.0;
// 2) backward loop
for (i = 9; i >= 0; i--)
  A[i][0] = A[i][0] / mean;
```

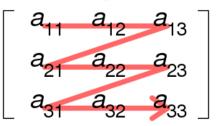
- Consider a program:
  - 1) Computes the <u>sum</u> of the first column of an array using a forward loop.
  - Normalizes the first column of an array by its mean (i.e. average) using a backward loop.
  - A[10][4] is an array of words located at memory (7A00)<sub>16</sub>~(7A27)<sub>16</sub> in row-major order.

### Row-Major vs. Column-Major Order

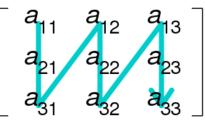


- Row-major order and column-major order are methods for storing multidimensional arrays in memory.
  - Row-Major: The consecutive elements of a row reside next to each other.
  - Column-Major: The consecutive elements of a column reside next to each other.
- For example,

Row-major order



Column-major order



Values as stored in Memory: 1

Column major:  $\begin{pmatrix} 1 & 5 & 9 & 13 \\ 2 & 6 & 10 & 14 \\ 3 & 7 & 11 & 15 \end{pmatrix}$ 

Row major:  $\begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 10 & 14 & 15 & 10 \end{pmatrix}$ 

## Cache Example (Cont'd)



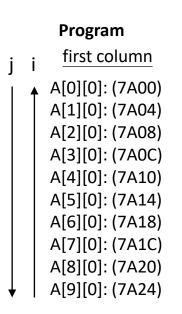
<b>A[10][4];</b> at <b>(7A00)</b> <sub>16</sub> ~ <b>(7A27)</b> <sub>16</sub> in <b>row-major order</b> . A[0][0] A[0][1] A[0][2] A[0][3] A[1][0] A[1][1] A[1][2] A[1][3] A[2][0] A[2][1] A[2][2] A[2][3] A[3][0] A[3][1] A[3][2] A[3][3] A[4][0] A[4][1] A[4][2] A[4][3]	Hex. (7A00) <sub>16</sub> (7A01) <sub>16</sub> (7A02) <sub>16</sub>	Memory Word Address (15-bit)  Binary  ( 1 1 1 1 0 1 0 0 0 0 0 0 0 0 0 0) <sub>2</sub> ( 1 1 1 1 0 1 0 0 0 0 0 0 0 0 1) <sub>2</sub> ( 1 1 1 1 0 1 0 0 0 0 0 0 0 1 1) <sub>2</sub> ( 1 1 1 1 0 1 0 0 0 0 0 0 0 1 1) <sub>2</sub>	Memory Contents (40 array elements)  A[0][0]  A[0][1]  A[0][2]  A[0][3]
A[5][0] A[5][1] A[5][2] A[5][3]	(7A03) <sub>16</sub> (7A04) <sub>16</sub>	$(1111010000011)_{2}$	A[1][0]
Program	:	<b>:</b>	
j i <u>first column</u> A[0][0]: (7A00)  A[1][0]: (7A04)  A[2][0]: (7A08)  A[3][0]: (7A0C)	(7A24) <sub>16</sub> (7A25) <sub>16</sub> (7A26) <sub>16</sub> (7A27) <sub>16</sub>	( 1 1 1 1 0 1 0 0 0 1 0 0 1 0 0) <sub>2</sub> ( 1 1 1 1 0 1 0 0 0 1 0 0 1 0 1) <sub>2</sub> ( 1 1 1 1 0 1 0 0 0 1 0 0 1 1 0) <sub>2</sub> ( 1 1 1 1 0 1 0 0 0 1 0 0 1 1 1) <sub>2</sub>	A[9][0] A[9][1] A[9][2] A[9][3]
A[4][0]: (7A10) A[5][0]: (7A14) A[6][0]: (7A18) A[7][0]: (7A1C) A[8][0]: (7A20) A[9][0]: (7A24)	Tag: 12 bits Tag: 14 bits	$8 = 2^3$ blocks in cache $\rightarrow$ 3 bits encodes cache by	
,	Tag: 15 bits	Tag for Associative	

• Why there's no "word" bit? One block contains one word (2°).

### **Direct Mapping**



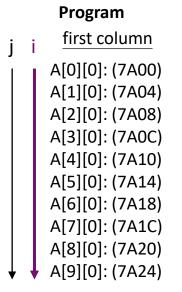
- The last 3-bits of address decide the CB.
  - Address % 8 → Cache Block Number
- No replacement algorithm is needed.
- When i = 9 and i = 8: 2 cache hits in total.
- Only 2 out of the 8 cache positions are used.
  - Very poor cache utilization: 25%



						Coi	ntent	of C	ache	Bloc	ks aft	ter Lo	оор Р	ass (	i.e. T	ïmeli	ne)				
		j = 0	j = 1	j = 2	j = 3	j = 4	j = 5	j = 6	j = 7	j = 8	j = 9	i = 9	i = 8	i = 7	i = 6	i = 5	i = 4	i = 3	i = 2	i = 1	i = 0
	0	A[0][0]	A[0][0]	A[2][0]	A[2][0]	A[4][0]	A[4][0]	A[6][0]	A[6][0]	A[8][0]	A[8][0]	A[8][0]	A[8][0]	A[8][0]	A[6][0]	A[6][0]	A[4][0]	A[4][0]	A[2][0]	A[2][0]	A[0][0]
	1																				
	2																				
Cache Block	3																				
Number	4		A[1][0]	A[1][0]	A[3][0]	A[3][0]	A[5][0]	A[5][0]	A[7][0]	A[7][0]	A[9][0]	A[9][0]	A[9][0]	A[7][0]	A[7][0]	A[5][0]	A[5][0]	A[3][0]	A[3][0]	A[1][0]	A[1][0]
	5																				
	6																				
	7																				



- Assume direct mapped cache is used.
- What if the *i* loop is a forward loop?

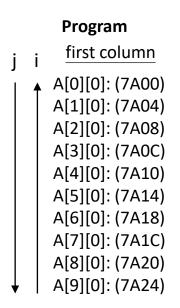


						Co	ntent	of C	ache	Bloc	ks aft	er Lo	ор Р	ass (	i.e. T	imeli	ne)				
		j = 0	j = 1	j = 2	j = 3	j = 4	j = 5	j = 6	j = 7	j = 8	j = 9	i = 0	i = 1	i = 2	i = 3	i = 4	i = 5	i = 6	i = 7	i = 8	i = 9
	0	A[0][0]	A[0][0]	A[2][0]	A[2][0]	A[4][0]	A[4][0]	A[6][0]	A[6][0]	A[8][0]	A[8][0]										
	1																				
	2																				
Cache Block	3																				
Number	4		A[1][0]	A[1][0]	A[3][0]	A[3][0]	A[5][0]	A[5][0]	A[7][0]	A[7][0]	A[9][0]										
	5																				
	6																				
	7																				

## **Associative Mapping**



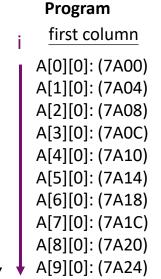
- All CBs are used in the FCFS basis.
- LRU replacement policy is used.
- When i = 9, 8, ..., 2: 8 cache hits in total.
- 8 out of the 8 cache positions are used.
  - Optimal cache utilization: 100%



						Co	ntent	of Ca	ache	Bloc	ks af	ter Lo	ор Р	ass (	i.e. T	imeli	ne)				
		j = 0	j = 1	j = 2	j = 3	j = 4	j = 5	j = 6	j = 7	j = 8	j = 9	i = 9	i = 8	i = 7	i = 6	i = 5	i = 4	i = 3	i = 2	i = 1	i = 0
	0	A[0][0]	A[8][0]	A[0][0]																	
	1		A[1][0]	A[9][0]	A[1][0]	A[1][0]															
	2			A[2][0]																	
Cache Block	3				A[3][0]																
Number	4					A[4][0]															
	5						A[5][0]														
	6							A[6][0]													
	7								A[7][0]												



- Assume associative mapped cache is used.
- What if the *i* loop is a forward loop?



						Co	ntent	of C	ache	Bloc	ks aft	er Lo	op P	ass (	i.e. T	imeli	ne)				
		j = 0	j = 1	j = 2	j = 3	j = 4	j = 5	j = 6	j = 7	j = 8	j = 9	i = 0	i = 1	i = 2	i = 3	i = 4	i = 5	i = 6	i = 7	i = 8	i = 9
	0	A[0][0]	A[8][0]	A[8][0]																	
	1		A[1][0]	A[9][0]																	
	2			A[2][0]																	
Cache Block	3				A[3][0]																
Number	4					A[4][0]	A[4][0]	A[4][0]	A[4][0]	A[4][0]	A[4][0]										
	5						A[5][0]	A[5][0]	A[5][0]	A[5][0]	A[5][0]										
	6							A[6][0]	A[6][0]	A[6][0]	A[6][0]										
	7								A[7][0]	A[7][0]	A[7][0]										

# 4-way Set Associative Mapping



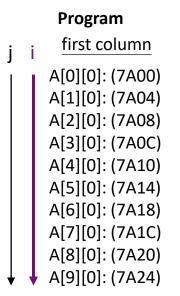
- There are total 8 ÷ 4 = 2 Cache Sets.
  - Address % 2 → Cache Set Number
- All accessed addresses are "even" (e.g. 7A00,
   7A04) → They will all be mapped to Cache Set 0.
- LRU replacement policy is used.
- When i = 9, 8, ..., 6: 4 cache hits in total.
- 4 out of the 8 cache positions are used (50% Util.).

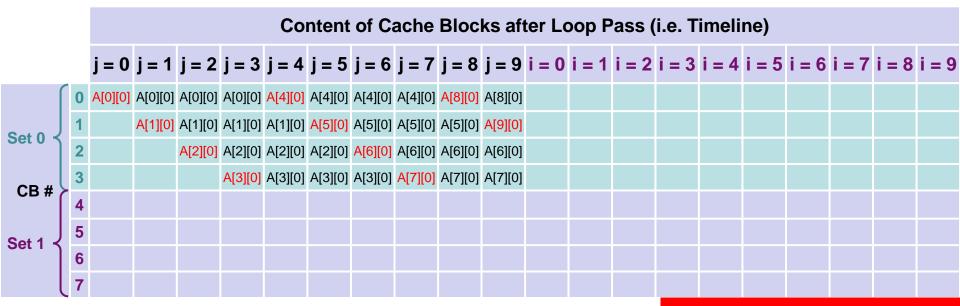
i first column
i A[0][0]: (7A00)
A[1][0]: (7A04)
A[2][0]: (7A04)
A[2][0]: (7A06)
A[3][0]: (7A0C)
A[4][0]: (7A10)
A[5][0]: (7A14)
A[6][0]: (7A18)
A[7][0]: (7A1C)
A[8][0]: (7A20)
A[9][0]: (7A24)

						Coi	ntent	of Ca	ache	Bloc	ks aft	er Lo	op P	ass (	i.e. T	imelii	ne)				
		j = 0	j = 1	j = 2	j = 3	j = 4	j = 5	j = 6	j = 7	j = 8	j = 9	i = 9	i = 8	i = 7	i = 6	i = 5	i = 4	i = 3	i = 2	i = 1	i = 0
	0	A[0][0]	A[0][0]	A[0][0]	A[0][0]	A[4][0]	A[4][0]	A[4][0]	A[4][0]	A[8][0]	A[4][0]	A[4][0]	A[4][0]	A[4][0]	A[0][0]						
Set 0	1		A[1][0]	A[1][0]	A[1][0]	A[1][0]	A[5][0]	A[5][0]	A[5][0]	A[5][0]	A[9][0]	A[9][0]	A[9][0]	A[9][0]	A[9][0]	A[5][0]	A[5][0]	A[5][0]	A[5][0]	A[1][0]	A[1][0]
	2			A[2][0]	A[2][0]	A[2][0]	A[2][0]	A[6][0]	A[2][0]	A[2][0]	A[2][0]										
CB#	3				A[3][0]	A[3][0]	A[3][0]	A[3][0]	A[7][0]	A[3][0]	A[3][0]	A[3][0]	A[3][0]								
OD #	4																				
Set 1 <	5																				
Oct 1	6																				
	7																				



- Assume 4-way set associative mapped cache is used.
- What if the *i* loop is a forward loop?





# **Summary**



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