



The Chinese University of Hong Kong

# CSCI2510 Computer Organization Lecture 07: Cache in Action

### Ming-Chang YANG

mcyang@cse.cuhk.edu.hk

COMPUTER ORGANIZATION AND EMBEDDED SYSTEMS

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 the data are firstly needed,
    <u>opportunistically bring it into cache</u> (i.e., 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, fetching more data of adjacent addresses (unit: cache block) from main memory into cache at a time.
- Cache takes both types of locality into considerations.

### **Cache at a Glance**





- 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.
- Cache Read (or Write) Hit/Miss: The read (or write) operation can/cannot be performed on the cache.
- Cache Management:
  - 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) but the cache is already full, the replacement algorithm removes a cached block and to create space for the new one.



CSCI2510 Lec07: Cache in Action 2023-24 T1

### Write Operation in Cache



- Write Operation:
  - Write-Through Scheme: The contents of cache and main memory are updated at the same time.
  - Write-Back Scheme: Update cache only but mark the item as dirty. The corresponding contents in main memory will be updated later when cache block is unloaded.
    - **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 cache size equals 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.
- 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 at any given time.
     Unit: Memory



# Mapping Functions (3/3)

- **Design Considerations of Mapping Functions:** 
  - 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.



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

1 Block = 2<sup>3</sup> Words 1 Word = 2<sup>1</sup> Bytes

F	16	-b	it I	<b>Ne</b> Blo	em ck A — M	ory Addr	y A ress	<b>\d</b>	dre	ess	5 (	bi	na	ry)	)	. I	Byte Addr ( <i>decimal</i> )	- )	Men	nory B #0	lock	#	MBs ¢0~#409	)5
0	0	0	0	0	0	- <u>By</u>	0		0	0	0	0	0	0	0	- <b>→</b>	Byte #0	]↔	Word	Bvte	Bvte		0	
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	<b>~</b>	Byte #1		#0	#1	#0		1	1
0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	++	Byte #2	<b>~</b>	Word	Byte	Byte			1
0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	+	Byte #3	+	#1	#3	#2			
0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0		Byte #4		Word	Byte	Byte			
0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1		Byte #5		#2	#5	#4			
0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0		Byte #6		Word	Byte	Byte			
0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1		Byte #7		#3	#7	#6			
0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0		Byte #8		Word	Byte	Byte			
0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1		Byte #9		#4	#9	#8			
0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0		Byte #10		Word	Byte	Byte			
0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1		Byte #11	4	#5	#11	#10			
0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0		Byte #12		Word	Byte	Byte			
0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	1		Byte #13	4	#6	#13	#12			
0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	0		Byte #14	4	Word	Byte	Byte			
90	;Q;	590	$\rho_{\epsilon}$	¢.	0	R	he	iρ,	Qi	OR	202	21-:	22-	<b>T1</b>	1		Byte #15		#7	#15	#14		4095	ŀ

### Example: Memory Block #1

 $\mathbf{0}$ 

 $\mathbf{0}$ 



1 Block =  $2^3$  Words

### Example: Memory Block #4095 1 Block = 2<sup>3</sup> Words 1 Word = 2<sup>1</sup> Bytes

	16	-b	it l	<b>Me</b> Blo	ck A	Oľ \ddi Vorc	y A	<b>\d</b>	dre	es	s (	bi I	na	ry)	) 1	1	Byte Addr ( <i>decimal</i> )	-	Mem	nory B #0	lock	#	MBs 0~#4095
1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0		B#65520		Word	Bvte	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#65532	-	Word	Byte	Byte		
1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	+	B#65533	+	#32766	#65533	#65532		
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0		B#65534	┝	Word	Byte	Byte		
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	+	B#65535		#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
Name:	

Date:

Given the same dividend (10010011)<sub>2</sub> as the previous example, what will be the quotient and remainder if the divisor equals to (10)<sub>2</sub>, (100)<sub>2</sub>, ..., (10000000)<sub>2</sub>?

# **Direct Mapping (1/4)**



#### 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 <u>directly mapped</u> (%) to a Cache <u>Set</u>.

(In a set? Associative!)



# **Direct Mapping (2/4)**

- Direct Mapped Cache: Each Memory Block will be <u>directly mapped</u> 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 needed for each CB.
  - Many MBs will be mapped to a same CB in cache.
  - We need to use some cache space (cost!) 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 = 2<sup>7</sup> cache blocks.
  - Word: Select one word in a block.
    - 3 bits: There're 8 = 2<sup>3</sup> words in a block.
  - Byte: Select one byte in a word.
    - 1 bits: There're  $2 = 2^1$  bytes in a word.
- Ex: CPU is looking for (0FF4)<sub>16</sub>
  - MAR = (0000 1111 1111 0100)<sub>2</sub>
  - $MB = (0000 \ 1111 \ 1111)_2 = (255)_{10}$
  - $CB = (1111111)_2 = (127)_{10}$
  - $\text{Tag} = (00001)_2$

Cache

Block 0

Block 1

Block 127

0000111111110100

16-bit Main Memory Address

Block

7

**Memory Block Number** 

(i.e. 0~4095)

Word B

3

tag

tag

tag

Taq

5

0000

Main Memory

Block 0

Block 1

Block 127

Block 128

Block 129

Block 255

Block 256

Block 257

# **Direct Mapping (4/4)**

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

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

00001 Quotient 1000000 ) 000011111111 (128)<sub>10</sub> 10000000 1111111 Remainder

- Search a 16-bit address (t, b, w, b):
  - See if <u>MB (t, b)</u> is already in <u>CB b</u>
    by comparing t with the tag of <u>CB b</u>.
  - If not, replace <u>CB b</u> with <u>MB (t, b)</u> and update tag of <u>CB b</u> using t.
  - ③ Finally access the word w in <u>CB b</u>.





### **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?



1 Block =  $2^3$  Words

1 Word = 2<sup>1</sup> Bytes

Main Memory

# Associative Mapping (1/3)



#### **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 <u>directly mapped</u> (%) to a Cache <u>Set</u>.

(In a set? Associative!)



# 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)

- 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.
  - $-MAR = (0000 \ 1111 \ 1111 \ 0100)_2$
  - $-MB = (0000 \ 1111 \ 1111)_2 = (255)_{10}^{-1}$
  - $Tag = (0000 \ 1111 \ 1111)_2$
- Search 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.
    - 128 tag comparisons can be done in parallel by hardware (cost!).



1 Block =  $2^3$  Words

1 Word =  $2^1$  Bytes

Main

### **Class Exercise 7.3**

- 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?



1 Block =  $2^3$  Words

1 Word = 2<sup>1</sup> Bytes

Main

# Set Associative Mapping (1/3)



#### 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!)

 $\left( \right)$ 

2

3

Cache

Blocks

2

3

4

5

Memory

**Blocks** 

#### **Set Associative**

 A Memory Block is <u>directly mapped</u> (%) to a Cache Set.

(In a set? Associative!)







# 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.





Main

# Set Associative Mapping (3/3)

- 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 2<sup>6</sup> cache sets.
  - Word & Byte: The last 3 & 1 bits.

Ex: CPU is looking for (0FF4)<sub>16</sub>

- Assume all CBs are empty.
- MAR = (0000 1111 1111 0100)<sub>2</sub>
- MB = (0000 1111 1111)<sub>2</sub> = (255)<sub>10</sub>
- Cache Set =  $(111111)_2 = (63)_{10}$
- Tag = (000011)<sub>2</sub>

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



Main Memory

1 Block =  $2^3$  Words 1 Word =  $2^1$  Bytes

### **Class Exercise 7.4**

- Assume 2-way set associative mapping is used, and all CBs are empty initially.
- Considering CPU is looking for  $(8010)_{16}$ :
  - Which MB will be loaded into the cache?
  - Which CB will store the MB?
  - What is the new tag for the CB?



1 Block =  $2^3$  Words

1 Word = 2<sup>1</sup> Bytes

Main

Memory

# 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 <u>directly mapped (%)</u> to a <u>Cache Set</u>.









# 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).



time

– 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).



time

– 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.



### 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 <u>8 blocks</u>.
  - A block is of  $1 (= 2^{\circ})$  word.
  - A word is of <u>16 bits</u>.



- Consider a program:
  - 1) Computes the <u>sum</u> of the first column of an array using a forward loop.
  - 2) <u>Normalizes</u> the first column of an array by its mean (i.e. average) using a backward loop.
  - A[10][4] is <u>an array of</u> words located at the memory word addresses (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 organizing multi-dimensional arrays in main memory (which appears to programs as a single, continuous address space).
  - Row-Major: The consecutive elements of a row reside next to each other.

	0x0	0x1	0x2	0x3	0x4	0x5	0x6	0x7	0x8
memory	a11	a12	a13	a21	a22	a23	a31	a32	a33

Row-major order <sup>-</sup> a<sub>11</sub> a<sub>12</sub> a<sub>12</sub>



Column-major order

 Column-Major: The consecutive elements of a column reside next to each other.

	0x0	0x1	0x2	0x3	0x4	0x5	0x6	0x7	0x8
memory	a11	a21	a31	a12	a22	a32	a13	a23	a33

https://en.wikipedia.org/wiki/Row-\_and\_column-major\_order

# Cache Example (Cont'd)





CSCI2510 Lec07: Cache in Action 1

A word is of  $2^1$  bytes: There is **one** "byte" bit (X).<sub>47</sub>

# **Direct Mapping**



- The last 3-bits of address decide the CB.
  - Memory Block Num. % 8  $\rightarrow$  Cache Block Num.
- 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%

Program first column A[0][0]: (7A00) A[1][0]: (7A04) A[2][0]: (7A08) 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)

Tags not shown but are needed!

						Со	ntent	of Ca	ache	Bloc	ks af	ter Lo	oop 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 = 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	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																				
CSCI2	251	10 Leo	c07: (	Cache	e in Ad	ction	2023-	24 T1							Tag	s not	sho	vn bu	it are	need	led!

### **Class Exercise 7.6**

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



						Со	ntent	of Ca	ache	Bloc	ks aft	er Lo	oop 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[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	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																				
CSCI2	51	0 Leo	c07: (	Cache	e in A	ction 2	2023-	24 T1							Taq	s not	show	vn bu	it are	need	ded!



# **Associative Mapping**

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

						Со	ntent	of Ca	ache	Bloc	ks af	ter Lo	oop 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 = 9	i = 8	i = 7	i = 6	i = 5	i = 4	i = 3	i = 2	i = 1	i =
	0	A[0][0]	A[8][0]	A[0][																	
	1		A[1][0]	A[9][0]	A[1][0]	A[1][															
	2			A[2][0]	A[2][																
Cache Block Iumber	3				A[3][0]	A[3][															
	4					A[4][0]	A[4][														
	5						A[5][0]	A[5][													
	6							A[6][0]	A[6][												
	7								A[7][0]	A[7][											
																	_				

CSCI2510 Lec07: Cache in Action 2023-24 T1

#### Tags not shown but are needed!





### **Class Exercise 7.7**

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



	Co	ontent of Cache	Blocks after Loo	p Pass (i.e. Timelin	e)
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
<b>0 A[0][0] A[0][0] A</b> [0]	0][0] A[0][0] A[0][0]	A[0][0] A[0][0] A[0][0]	<mark>A[8][0]</mark> A[8][0]		
<b>1</b> A[1][0] A[ <sup>2</sup>	][0] A[1][0] A[1][0]	A[1][0] A[1][0] A[1][0]	A[1][0] <mark>A[9][0]</mark>		
2 A[2	2][0] A[2][0] A[2][0]	A[2][0] A[2][0] A[2][0]	A[2][0] A[2][0]		
Cache <sub>3</sub> Block	<mark>A[3][0]</mark> A[3][0]	A[3][0] A[3][0] A[3][0]	A[3][0] A[3][0]		
lumber <sup>4</sup>	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]		

CSCI2510 Lec07: Cache in Action 2023-24 T1

#### Tags not shown but are needed!

### 4-way Set Associative Mapping

• There are total  $8 \div 4 = 2$  Cache Sets.

− Memory Block Num. % 2 → Cache Set Num.

- The numbers of accessed MBs are all "even" (e.g. 7A00, 7A04) → Mapped to Cache Set #0.
- LRU replacement policy is used.
- When i = 9, 8, ..., 6: <u>4</u> cache hits in total.
- 4 out of the 8 cache positions are used (50% Util.).



#### CSCI2510 Lec07: Cache in Action 2023-24 T1

Tags not shown but are needed!



Program

first column

A[0][0]: (7A00)

A[1][0]: (7A04) A[2][0]: (7A08)

A[3][0]: (7A0C) A[4][0]: (7A10)

A[5][0]: (7A14) A[6][0]: (7A18)

A[7][0]: (7A1C)

### **Class Exercise 7.8**

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



						Со	ntent	of Ca	ache	Bloc	ks aft	ter Lo	oop 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[0][0]	A[0][0]	A[0][0]	A[4][0]	A[4][0]	A[4][0]	A[4][0]	A[8][0]	A[8][0]										
Sot 0 2	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]										
Selun	2	2    A[2][0]    A[2][0]    A[2][0]    A[2][0]    A[6][0]    <																			
CB#	3	3      A[3][0]      A[3][0]      A[3][0]      A[3][0]      A[3][0]      A[3][0]      A[7][0]      A[7][0]      A[7][0]      A[7][0]																			
00 #	4																				
Set 1 🖌	5																				
	6																				
	7																				
CSC	25	10 Le	c07: 0	Cache	e in Ao	ction 2	2023-	•24 T1							Taq	s not	shov	wn bu	it are	need	led!



### 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