

Module 6

Memory

Objectives:

- To master the concepts of hierarchical memory organization.
- To understand how each level of memory contributes to system performance, and how the performance is measured.
- To master the concepts behind cache memory and understands the basic concept of virtual memory.

Module 6

Memory

- 6.1 Introduction
- 6.2 Main Memory
- 6.3 Cache Memory
- 6.4 Virtual Memory
- 6.5 Summary

Module 6

Memory

- 6.1 Introduction
- 6.2 Main Memory
- 6.3 Cache Memory
- 6.4 Virtual Memory
- 6.6 Summary

- Overview
- Type of Memory
- The Memory Hierarchy

6 Overview

- Most computers are built using the *Von Neumann* model, which is centered on memory.
- The *programs* that perform the processing are stored in _____.
- The memory unit that communicates directly with the CPU is called _____.
- Devices that provide backup storage are called *auxiliary memory*.*

Only programs and data currently needed by the processor reside in main memory.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.233.
* Mano, Morris M. (1993). *Computer System Architecture* (3rd Edition), p.445.

6

- **Internal memory** is often equated with main memory, but there are other forms of internal memory:
 - The processor requires its own local memory → _____.
 - The control unit portion of the processor may also require its own internal memory → _____.
 - **Cache** is another form of internal memory.
- The complex subject of computer memory is made more manageable if we classify memory systems according to their key characteristics (see next table).

William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.121.

6

Table: Key Characteristics of Computer Memory Systems.

Location	Internal (e.g., processor registers, cache, main memory) External (e.g., optical disks, magnetic disks, tapes)	Performance	Access time Cycle time Transfer rate
Capacity	Number of words Number of bytes	Physical Type	Semiconductor Magnetic Optical Magneto-optical
Unit of Transfer	Word Block	Physical Characteristics	Volatile/nonvolatile Erasable/nonerasable
Access Method	Sequential Direct Random Associative	Organization	Memory modules

William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.121.

7

6

Type of Memory

Why are there so many different types of computer memory?



- The answer → new technologies continue to be introduced in an attempt to match the improvements in CPU design.
- Even though a large number of memory technologies exist, there are only two basic types of memory:

RAM (Random Access Memory)

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.233

8

6

(a) RAM (Random Access Memory)

- RAM is also the “main memory” or _____.
- Used to store programs and data that the computer needs when executing programs.



However, RAM is **volatile**, and **loses** this information once the power is turned off.

The organization of RAM is a key design issue → refers to the physical arrangement of bits to form words.*

- Two general types of chips used to build the bulk of RAM memory in today’s computers:
 - **SRAM (Static Random Access Memory)**.
 - _____.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.234
 * William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited. p.124

9

6

(b) ROM (Read-Only Memory)

- Stores critical information necessary to operate the system, such as the program necessary to boot the computer.



- ROM is not volatile and always retains its data.
- Maintain information when the power is shut off.

- Some different types of ROM:

□ **PROM (Programmable ROM)**

□ **EPROM (Erasable PROM)**

□ _____

□ **Flash memory** → essentially EEPROM with the added benefit .

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.234

10

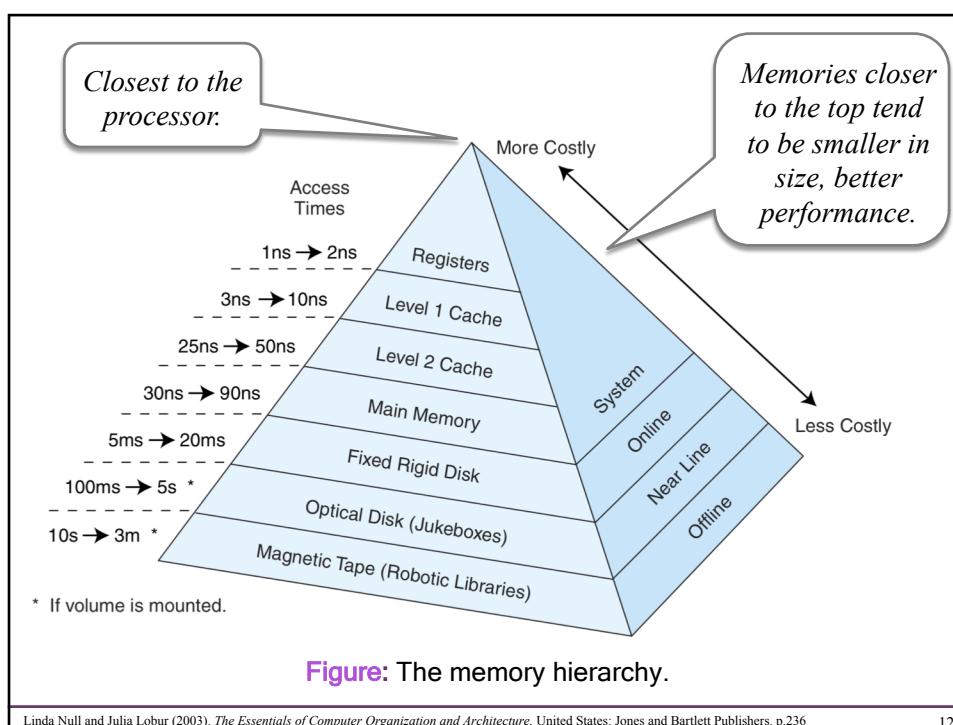
6

The Memory Hierarchy

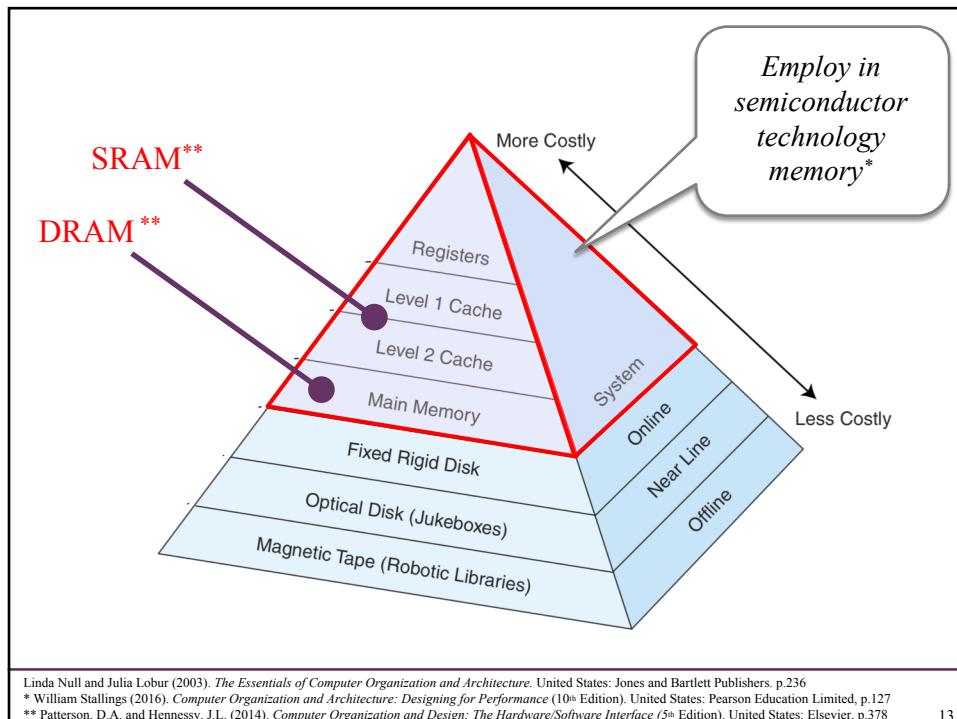
- Generally speaking, faster memory is _____ than slower memory.
- To provide the best performance at the lowest cost, memory is organized in a hierarchical fashion.
- Small, fast storage elements are kept in the CPU.
 - Larger, slower main memory is accessed through the data bus.
 - Larger, (almost) permanent storage in the form of disk and tape drives is still further from the CPU.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.235

11



12



13

6

- To access a particular piece of **data**, the CPU first sends a request to its nearest memory, usually _____.
 - If the data is not in cache, then **main memory** is queried.
 - If the data is not in main memory, then the request goes to **disk**.

- Once the data is located, then the data, and a number of its nearby data elements are **fetched** into _____.

The diagram shows a pyramid divided into horizontal layers. The top layer is labeled 'Registers'. Below it is 'Level 1 Cache', then 'Level 2 Cache'. The next layer down is 'Main Memory'. Following that is 'Fixed Rigid Disk', then 'Optical Disk (Jukeboxes)', and finally 'Magnetic Tape (Robotic Libraries)' at the base. A red arrow points from the label 'CPU' at the top left towards the 'Registers' layer at the top of the pyramid.

14

6

Table: The following terminology is used when referring to the memory hierarchy.

Terminology	Definition
<i>Hit</i>	The data is <u>found</u> at a given memory level.
<i>Miss</i>	The data is <u>not found</u> at a given memory level.
<i>Rate</i>	The percentage (%) of memory accesses <u>found</u> in a given level of memory.
<i>Rate</i>	The percentage (%) of memory accesses <u>not found</u> in a given level of memory → $(1 - \text{hit rate})$
<i>Hit Time</i>	The time required to access data at a given memory level.
<i>Miss Penalty</i>	The time required to process a <i>miss</i> , including the <u>time that it takes to replace</u> a block of memory plus the <u>time it takes to deliver</u> the data to the processor.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.235

15

6

Locality of Reference

- An entire blocks of data is copied after a *hit* because the principle of *locality* tells us that once a byte is accessed, it is likely that a nearby data element will be needed soon.
- There are three basic forms of *locality*:

- 1) *locality*—Recently accessed items tend to be accessed again in the near future.
- 2) *locality*—Accesses tend to be clustered in the address space.
- 3) *locality*—Instructions tend to be accessed sequentially.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.237

16

Module 6

Memory

6.1 Introduction

6.2 Main Memory

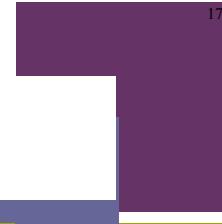
6.3 Cache Memory

6.4 Virtual Memory

6.6 Summary

- ❑ Overview
- ❑ Memory Organization
- ❑ Memory Capacity
- ❑ Memory Interleaving

17



RAM (Read-Access Memory)
ROM (Read-Only Memory)

6

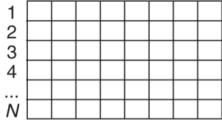
Overview

- The _____ is the central storage unit in a computer system.
 - The principle technology used for the **main memory** is based on semiconductor integrated circuits (RAM) either *static* or *dynamic* operating modes.
- Most of the **main memory** in a general-purpose computer is made up of RAM integrated circuit, but a portion of the memory may be constructed with ROM.

6 Memory Organization

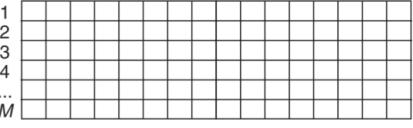
- We can envision memory as a matrix of bits.
- Each row, implemented by a _____, has a length typically equivalent to the word size of the machine.
- Each register (more commonly referred to as a *memory location*) has a unique address; memory addresses usually start at zero (0) and progress upward.

Address ←————— 8-bit —————→



N 8-Bit Memory Locations

Address ←————— 16-bit —————→



M 16-Bit Memory Locations

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.153

19

Memory Location

- Each part of memory has a separate memory location, which can be referred to using a _____. Number of bits for an address uniquely access to a memory location.

$$\text{No of bits} = \frac{\log (\text{memory capacity})}{\log 2}$$

$$\text{No of locations} = 2^{(\text{No.of bit in the address})}$$

Address :	Memory Words:
000:	
001:	
010:	
:	
$2^n - 1:$	

Where n is the bit of the address

20

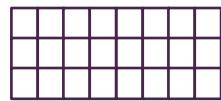
6

Memory Capacity

- Memory capacity usually measured in bits:

Total no. of memory locations (words) × size of memory word

- ### ■ **Examples:** 3 words (locations/rows)



(a)



(b)

$$\begin{aligned} \text{Size} &= 3 \text{ words} \times 8 \text{ bits} \\ &= 24 \text{ bits} \end{aligned}$$

$$\begin{aligned}Size &= \\&= \end{aligned}$$

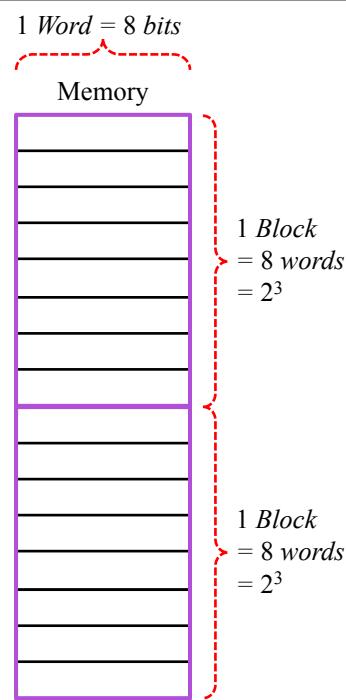
21

Example 1:

Main memory is divided into blocks.
The memory word is 8 bit and the size of a block is 8 words.

- (a) What is the **capacity** of the main memory, if the total number of blocks in the memory is 128?

 - (b) How many **blocks** in the main memory if the memory capacity is 32 Kbit?



22

6

Solution :

$$\begin{aligned}
 (a) \text{Memory Capacity} &= \text{Blocks} \times \text{Words} \times \text{Bits} \\
 &= 128 \times 8 \times 8 \\
 &= 2^7 \times 2^3 \times 8 \\
 &= 2^{10} \times 8 \\
 &= 1\text{Kbit} \times 8 = 8\text{Kbits}
 \end{aligned}$$

$$(b) \text{Memory Capacity} = \text{Blocks} \times \text{Words} \times \text{Bits}$$

$$32\text{Kbit} = \text{Blocks} \times 8 \times 8$$

23

6

Memory Interleaving

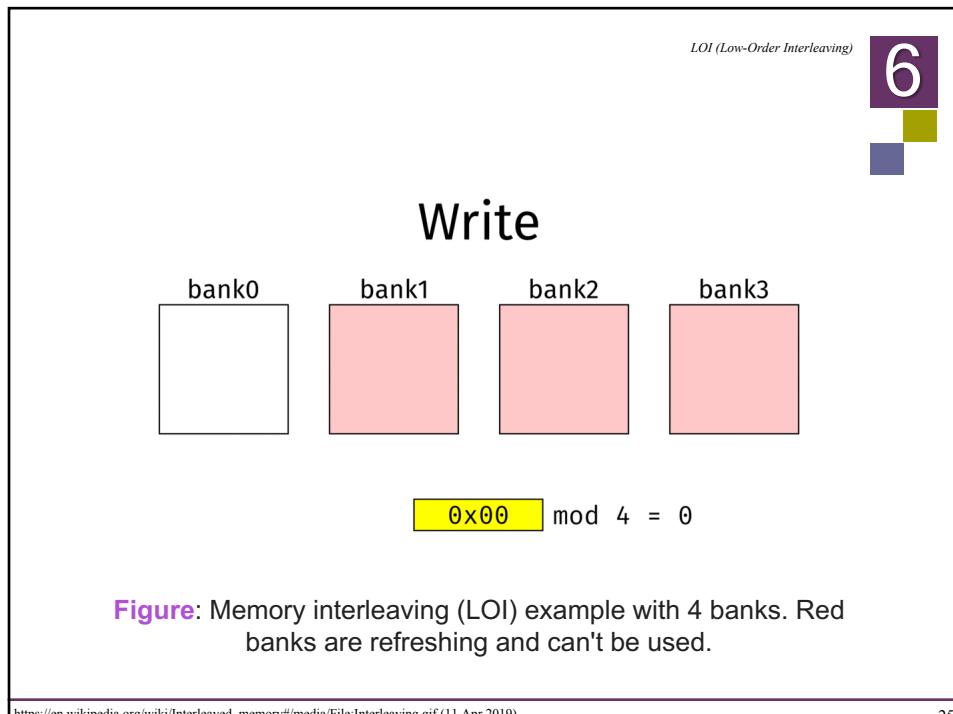
- A single shared memory module causes sequentialization of access.
- *Memory interleaving* → splits memory across multiple memory modules (or _____).

A number of memory chips can be grouped together to form a memory bank

Low-Order Interleaving (LOI)

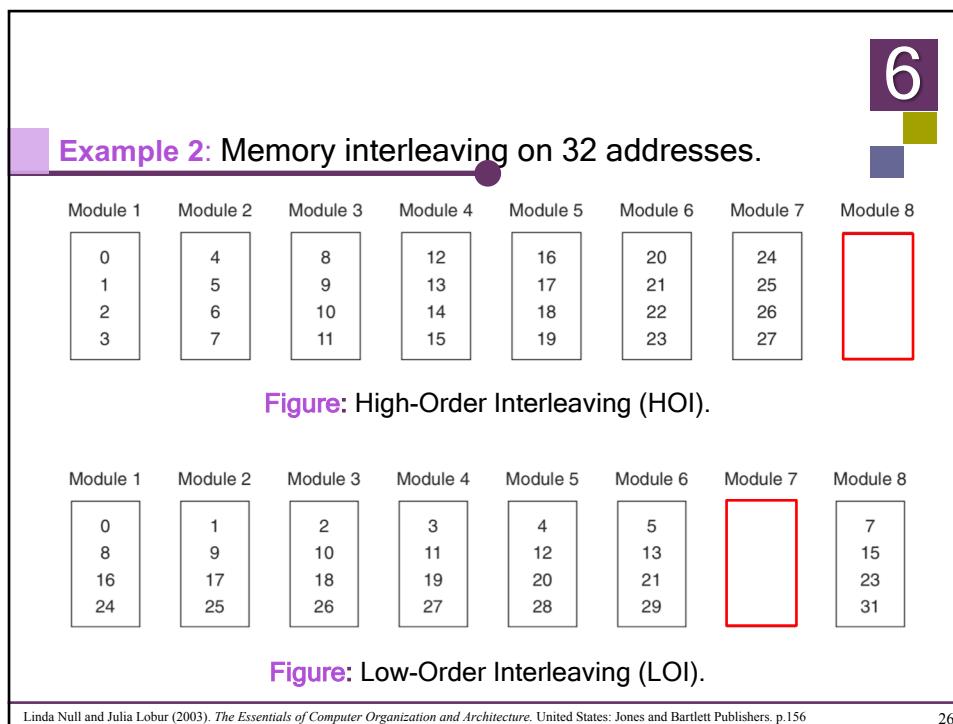
The low-order bits of the address are used to select the bank. e.g. → 00000101

The high-order bits of the address are used to select the bank. e.g. → 10100000



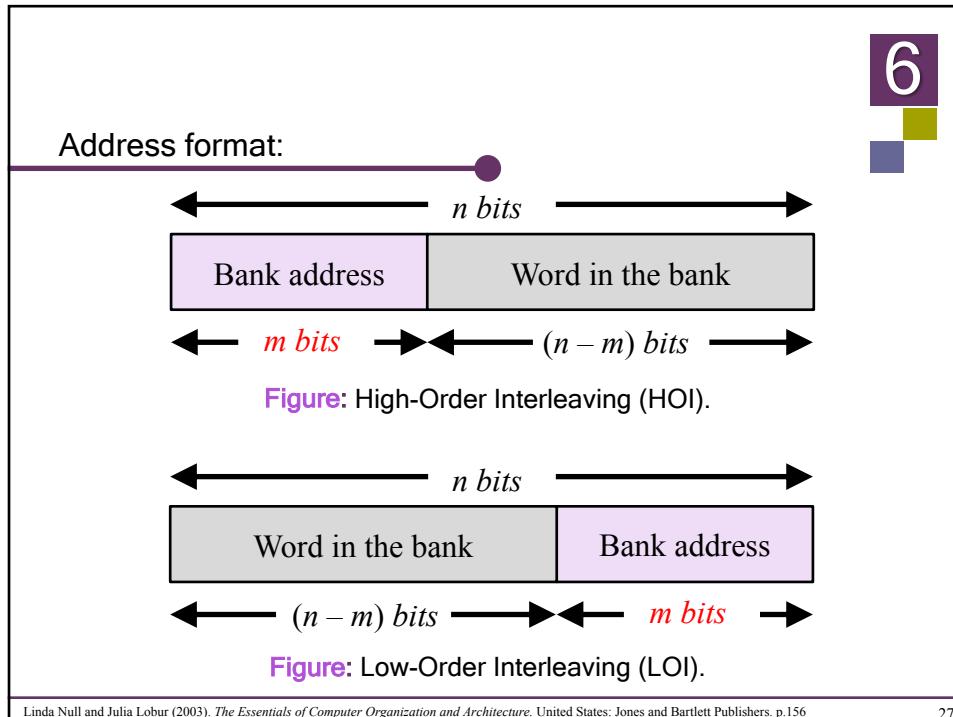
https://en.wikipedia.org/wiki/Interleaved_memory#/media/File:Interleaving.gif (11 Apr 2019)

25

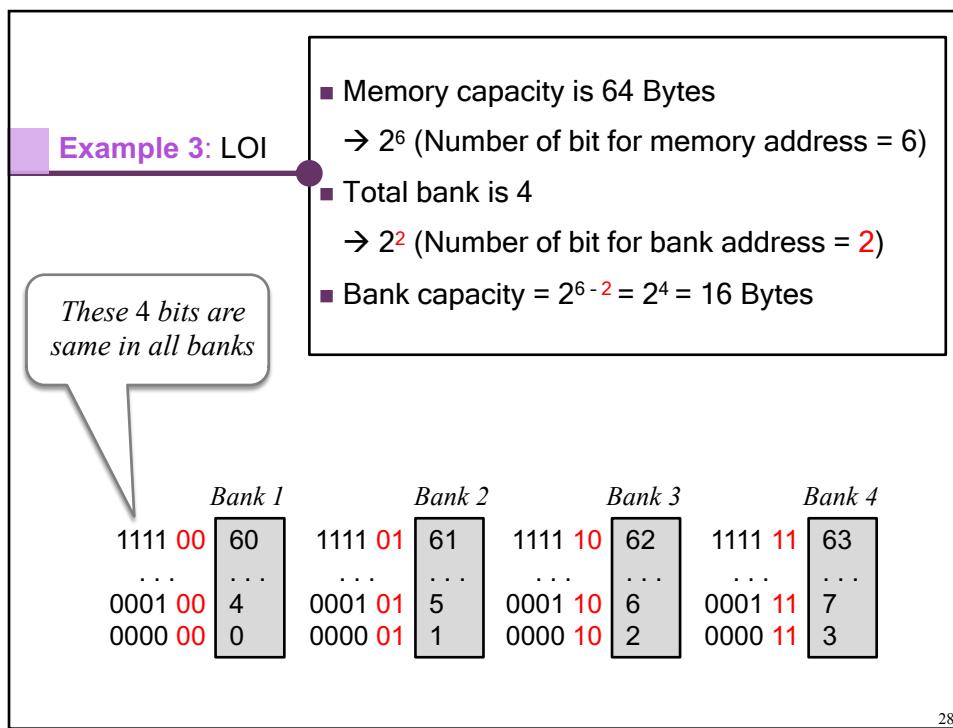


Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.156

26

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.156

27



28

Example 4: LOI

Given a memory address as 29Ch (10 bits) and there are 4 memory banks. Determine the memory bank address and the address of the word in the bank using LOI.

Solution:

- 4 bank $\rightarrow 2^2$ (Number of bit for bank address = 2)
- Bank capacity = $2^{10-2} = 2^8 = 256$ Bytes
- Memory address = 29Ch = 10 1001 1100

Bank 1		Bank 2		Bank 3		Bank 4	
1111 1111 00 ... 0000 0001 00 0000 0000 00	(3FCh) 1020 ... 4 0 (000h)	... 01 ... 01 1	1021 5 1 ... 01	... 10 ... 10 2 ... 10	1022 6 2 ... 10	... 11 ... 11 3 ... 11	1023 7 3 ... 11

29

LOI: Advantages and Disadvantages

6

It produces memory interference.

A failure of any single bank would be catastrophic to the whole system.

30

Example 5: HOI
(Same as Example 3)

These 4 bits are same in all banks

Bank 1		Bank 2		Bank 3		Bank 4	
00 1111	15	01 1111	31	10 1111	47	11 1111	63
...		
00 0001	1	01 0001	17	10 0001	33	11 0001	49
00 0000	0	01 0000	16	10 0000	32	11 0000	48

- Memory capacity is 64 Bytes
→ 2^6 (Number of bit for memory address = 6)
- Total bank is 4
→ 2^2 (Number of bit for bank address = 2)
- Bank capacity = $2^{6-2} = 2^4 = 16$ Bytes

31

Example 6: HOI

Given a main memory has 32 Mwords with the size of each word is 1 Byte, and there are 16 memory banks. Draw the modular memory address format if the system is implemented with HOI.

Solution:

- Main memory → 32 Mwords = 2^5 Mwords = $2^5 \times 2^{20}$ bits = 2^{25} (Number of bit for memory address = 25)
- 16 memory bank → 2^4 (Number of bit for bank address = 4)
- Word in the bank = $25 - 4 = 21$ bits

32

6

HOI: Advantages



- Easy memory extension by the addition of one or more memory modules to a maximum of $M - 1$.
- Provides better reliability, since a failed module affects only a localized area of the address space.
- This scheme would be used w/o conflict problems in multiprocessors if the modules are partitioned according to disjoint or non-interleaving processes (programs should be disjoint for its success).

33

6

HOI: Disadvantages



- Scheme will cause **memory conflicts** in case of pipelined, vector processors. The sequentially of instructions and data to be placed in the same module. Since memory cycle time is much greater than pipelined clock time, a previous memory request would not have completed its access before the arrival of the next request, thereby resulting in a **delay**.
- Process interacting and sharing instructions and data in multiprocessor system will **encounter considerable conflicts**.
- This technique is useful **only in one single user system / single user multitasking system**.

34

Module 6

Memory

6.1 Introduction

6.2 Main Memory

6.3 Cache Memory

6.4 Virtual Memory

6.6 Summary

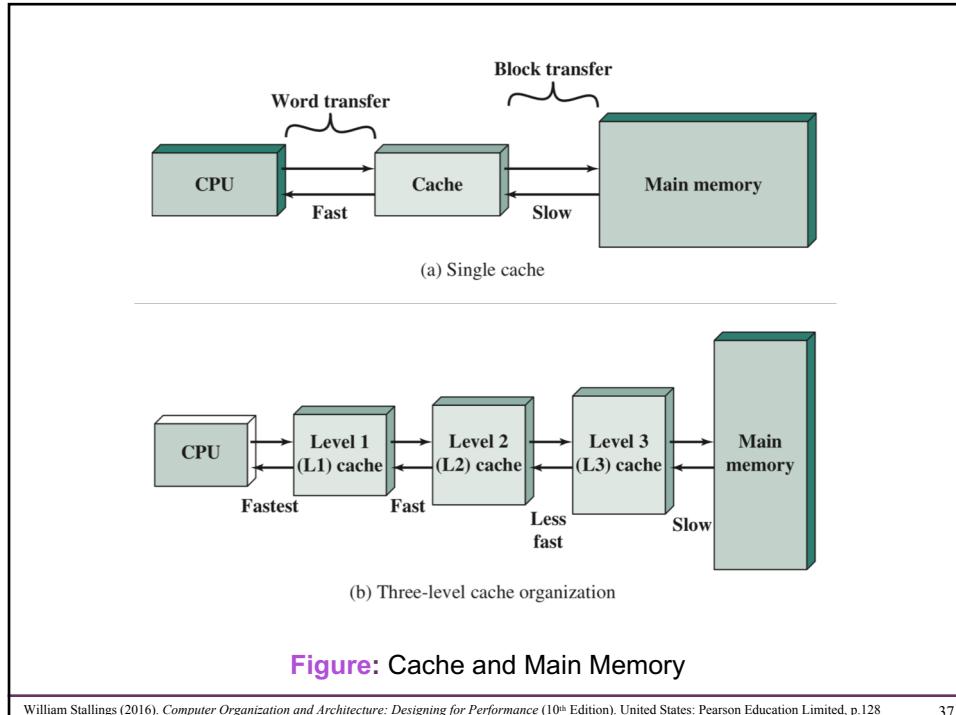
- ❑ Overview
- ❑ Cache Mapping Schemes
- ❑ Replacement Policy
- ❑ Cache Performances

35

6 Overview

- **Cache memory* is designed to combine:
 - ❑ the memory access time of expensive, high-speed memory with
 - ❑ the large memory size of less expensive, lower-speed memory.
- *The cache contains a copy of portions of main memory.
- Unlike main memory, which is accessed by _____, cache is typically accessed by content; hence, it is often called *content addressable memory*.

* William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.128



William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.128

37

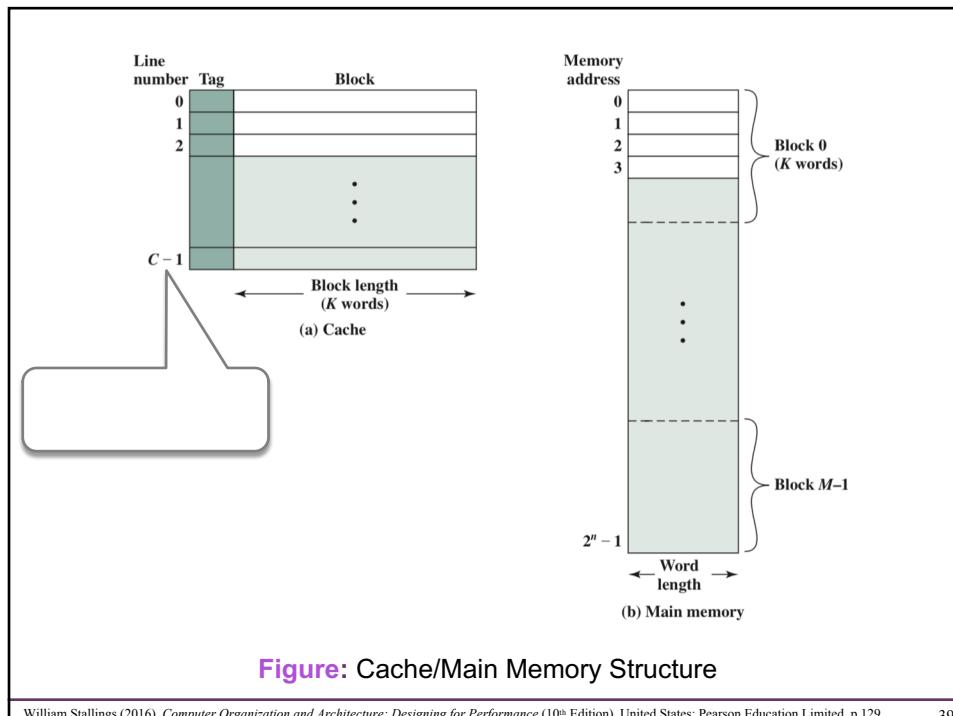
6

Cache Mapping Schemes

- The “content” that is addressed in *addressable cache memory* is a subset of the bits of a *main memory* address called a _____.
 - The fields into which a memory address is divided provide a many-to-one mapping between larger main memory and the smaller cache memory.
- Many blocks of main memory map to a single block of cache.
 - A _____ field in the cache block distinguishes one cached memory block from another.

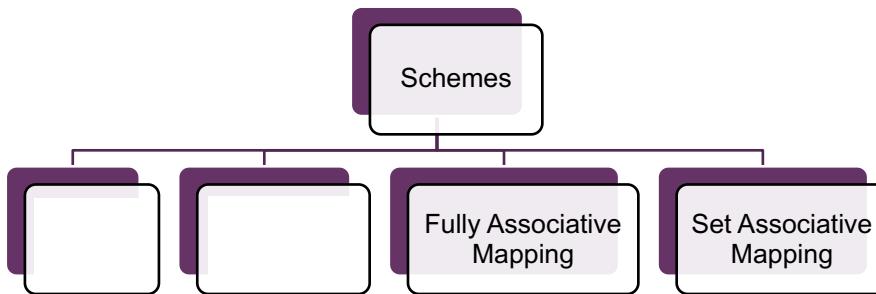
Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.239

38

**Figure:** Cache/Main Memory StructureWilliam Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.129

39

- Depending on the **mapping scheme**, cache may have two or three fields.
- These fields depend on the particular **mapping scheme** being used to determine where the data is placed when it is originally copied into cache.

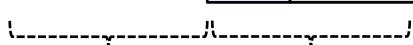
**Figure:** Cache mapping schemesLinda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.240

40

6

(a) Direct Mapping

- The simplest technique.
- *A cache structure in which each memory location is mapped to exactly one location in the cache.
- Address formats:

<i>Main memory address :</i>		<i>Cache Memory :</i>
		
		
<i>Cache Address</i>		<i>Cache Word</i>

*Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.384

41

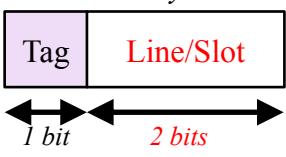
6

Example 7: Direct Mapping.

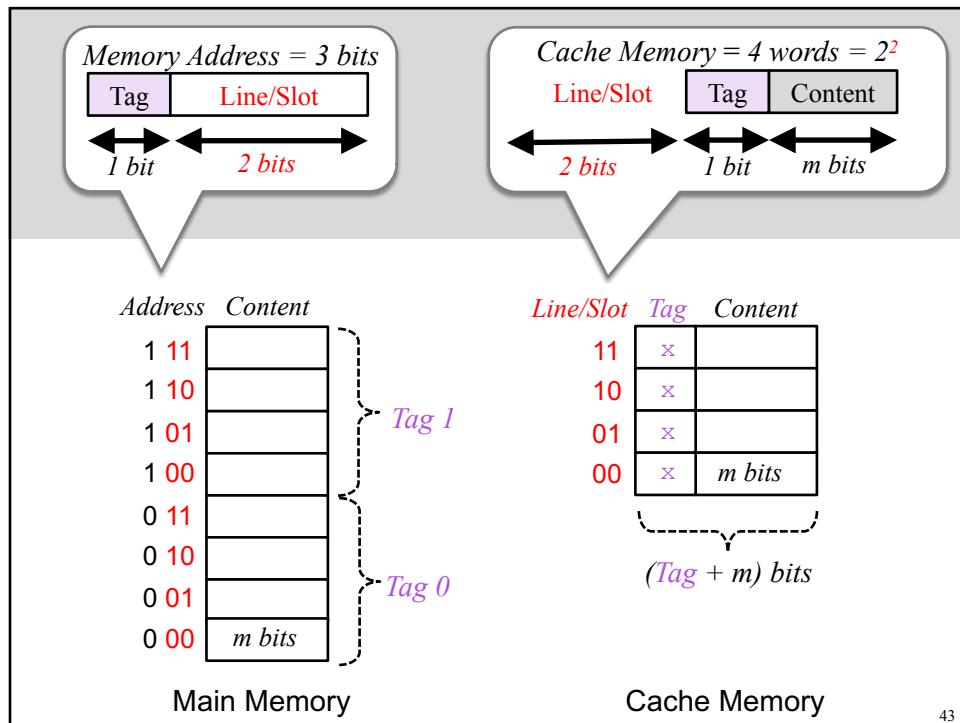
A main memory contains 8 words while the cache has only 4 words. Using direct address mapping, identify the fields of the main memory address.

Solution :

- Total memory words = $8 = 2^3$
→ Require 3 bits for main memory address.
- Total cache words = $4 = 2^2$
→ Require 2 bits for cache address (*Line/Slot*)
- Tag size = $3 - 2 = 1$ bit

<i>Main memory = 3 bits</i>		
-----------------------------	--	--

42



43

6

Example 8: Direct Mapping.

A main memory contains 32 words while the cache has only 8 words. Using direct address mapping, identify the fields of the main memory address.

Solution :

- Total memory words = $32 = 2^5$
→ Require 5 bits for main memory address.
- Total cache words = $8 = 2^3$
→ Require 3 bits for cache address (*Line/Slot*)
- Tag size = _____ bits

Main memory = 5 bits

Tag	Line/Slot
-----	-----------

2 bit 3 bits

44

Example 9: Direct Mapping.

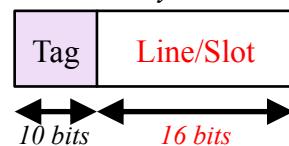
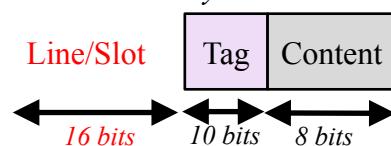
Size of cache memory is $64K$ word and the size of main memory is $64M \times 8$ bit word. Determine the word size of main memory, cache and the main memory address format.

Solution :

- Total words in main memory = $64M = 2^6 \times 2^{20} = 2^{26}$
→ Require 26 bits for main memory address.
- Total cache words = $64K = 2^6 \times 2^{10} = 2^{16}$ → 16 bits for *Line/Slot*
- Tag size = $26 - 16 = 10$ bits
- Size of main memory word = 8 bits
→ Size of cache word = Tag + (No. words in cache × size)
= _____ bits

45

Main memory = 26 bits

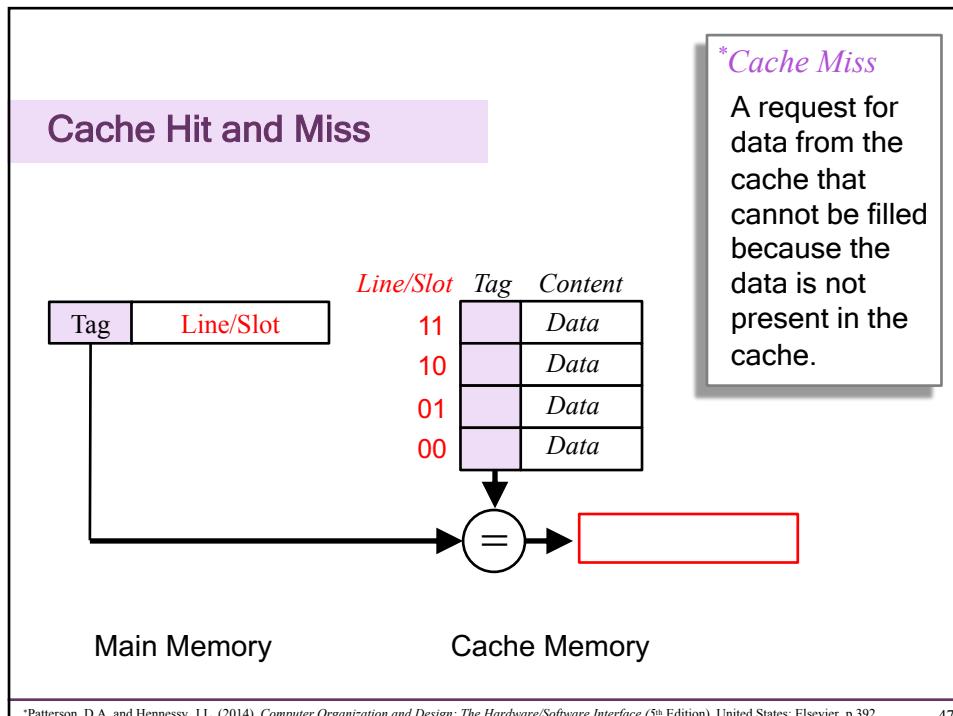
Cache Memory = $64K = 2^{16}$ 

Cache word

=

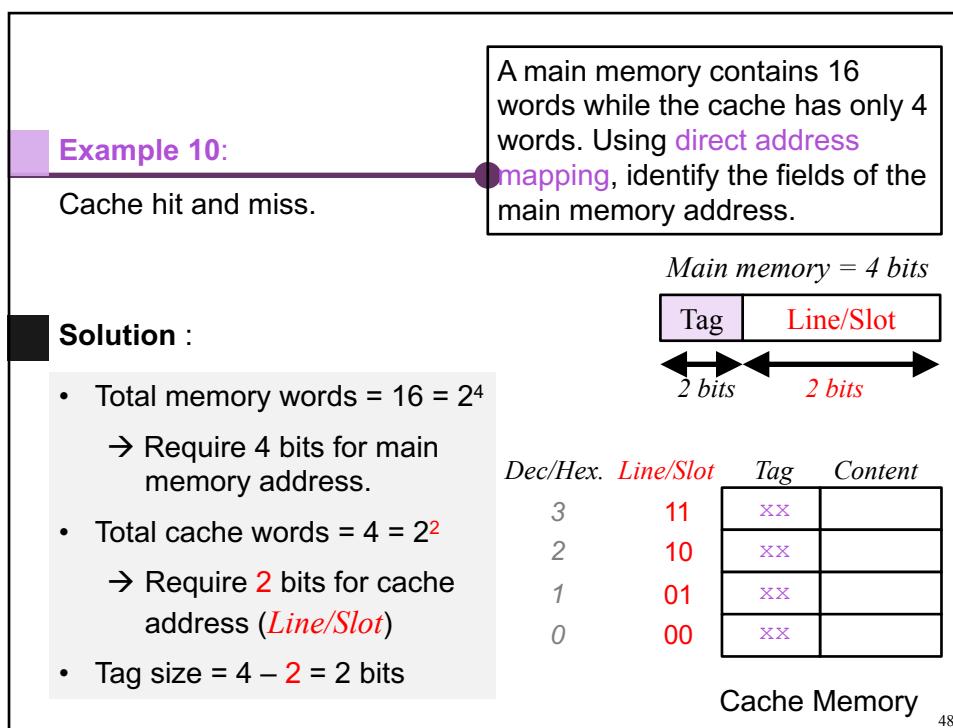
=

46



Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.392

47



48

Memory address:

→

Based on previous example, consider the following main memory with contents.

	Hex.	Address	Content (Hex)
F	11	11	11
E	11	10	10
D	11	01	01
C	11	00	A1
B	10	11	F5
A	10	10	F4
9	10	01	F3
8	10	00	F1
7	01	11	FF
6	01	10	91
5	01	01	10
4	01	00	19
3	00	11	13
2	00	10	02
1	00	01	01
0	00	00	EE

Dec/Hex. Line/Sot Tag Content

3	11	xx	xx
2	10	xx	xx
1	01	xx	xx
0	00	xx	xx

Cache Memory Main Memory

(File: Module 6 – CacheHitMiss.doc) – Example 1

49

Example 11: Cache hit and miss operation.
(Read Process)

A main memory contains 32 words while the cache has only 8 words. Using direct address mapping, identify the fields of the main memory address.

Solution :

- Total memory words = $32 = 2^5 \rightarrow 5$ bits
- Total cache words = $8 = 2^3 \rightarrow 3$ bits (Line/Sot)
- Tag size = $5 - 3 = 2$ bits

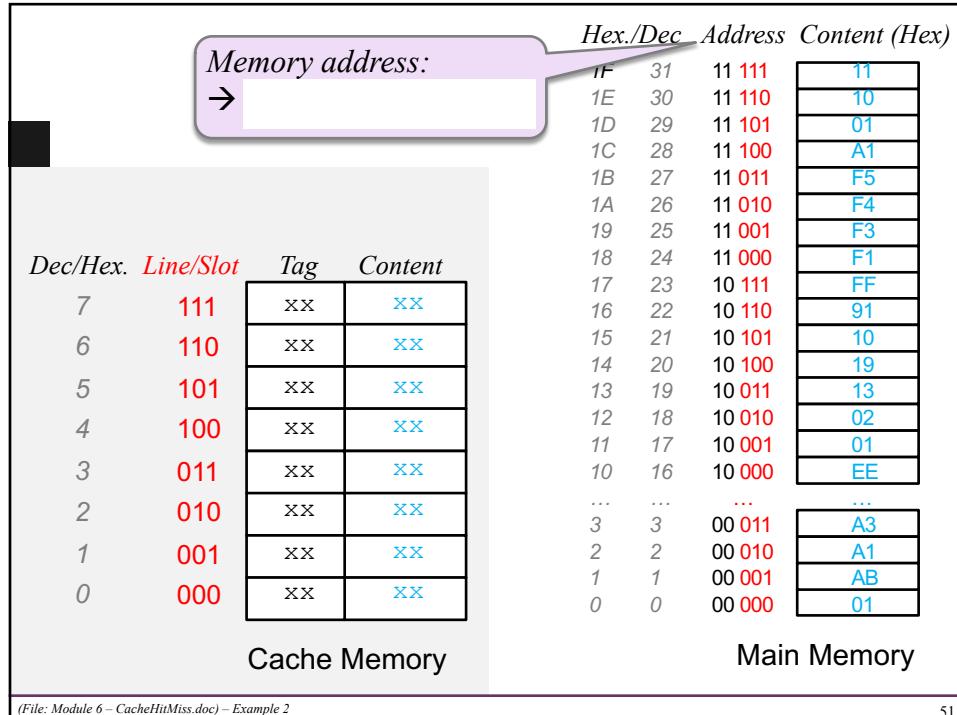
Main memory = 5 bits

Tag	Line/Sot
-----	----------

2 bits 3 bits

(File: Module 6 – CacheHitMiss.doc) – Example 2

50



51

Based on the example, show the contents of the cache as it responds to a series of request (decimal address):

22, 26, 22, 26, 16, 3, 16, 18

Generated Address			Format Address	
Address		Content	Tag	Line/S _{lot}
(Dec)	(Binary)	(Hex)		
22	10110	91		
26	11010	F4		
22	10110	91		
26	11010	F4		
16	10000	EE		
3	00011	A3		
16	10000	EE		
18	10010	02		

Main Memory

Complete the table
Main memory = 5 bits
Tag Line/S_{lot}
2 bits 3 bits

(File: Module 6 – CacheHitMiss.doc) – Example 2

52

Based on the example, show the contents of the cache as it responds to a series of request (decimal address):

22, 26, 22, 26, 16, 3, 16, 18

Format Address			Format Address		Read/Write operation cache				
Line/Slot	Tag	Content (Hex)	Tag	Line/Slot	Hit	Miss	Update cache	Read	Write
111	xx	xx	10	110					
110	10	91	11	010					
101	xx	xx	10	110					
100	xx	xx	11	010					
011	00	A3	10	000					
010	11 10	F4 02	00	011					
001	xx	xx	10	000					
000	10	EE	10	010					

Cache Memory
Main Memory

(File: Module 6 – CacheHitMiss.doc) – Example 2

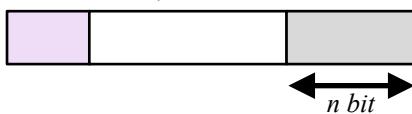
53

(b) Block Direct Mapping

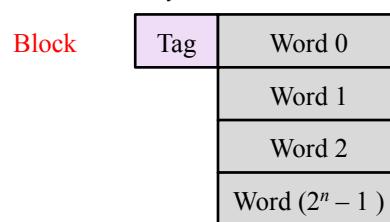
- The simplest technique of direct mapping can map each block of main memory into only one possible cache line.
- Address formats: (2^n words, e.g. $n = 2$)

$$\begin{aligned} &= 2^n - 1 \\ &= 2^2 - 1 \\ &= 4 - 1 = 3 \end{aligned}$$

Main memory address :



Cache Memory :

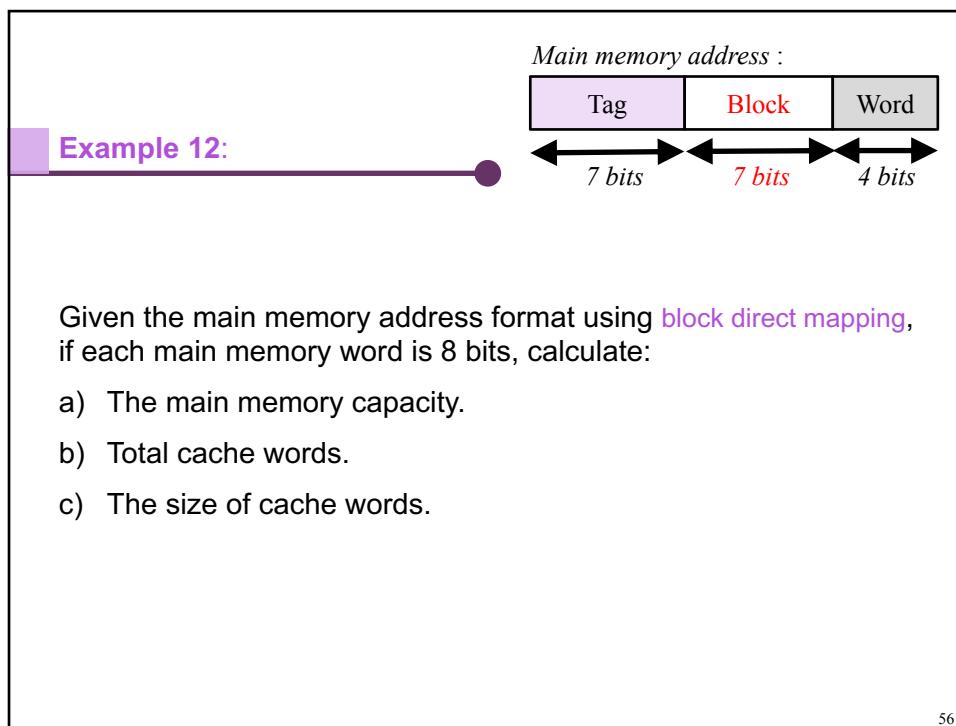
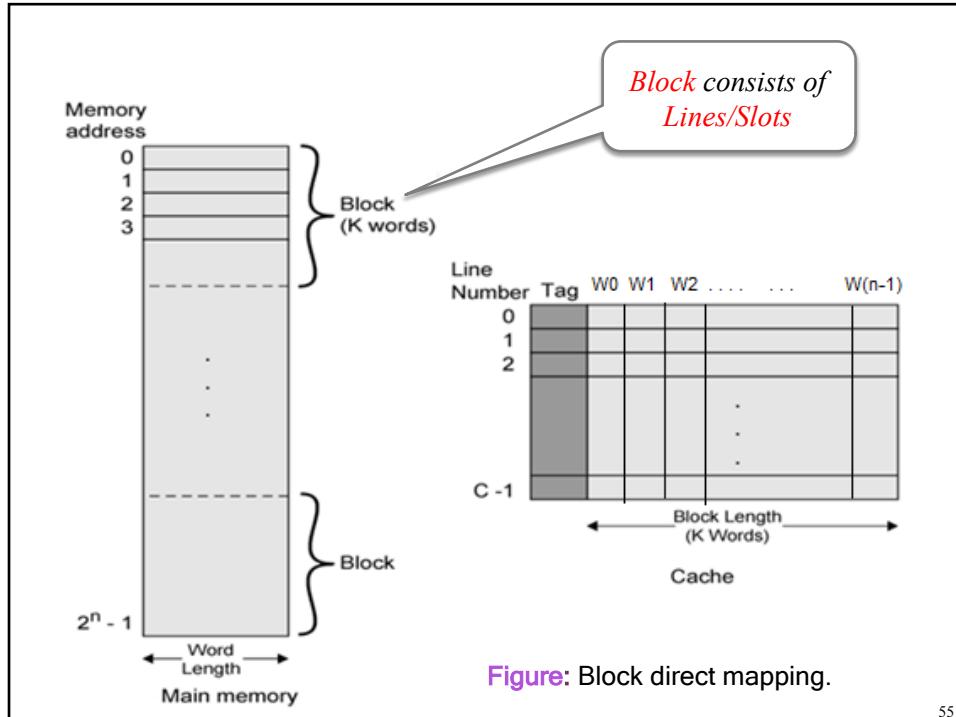


Cache Address

Cache Word

William Stallings (2016). Computer Organization and Architecture: Designing for Performance (10th Edition). United States: Pearson Education Limited, p.134

54



Solution :

Main memory address :

Tag	Block	Word
-----	-------	------

7 bits 7 bits 4 bits

a) Total main memory bits =

→ Total memory words = $2^{18} = 256\text{Kword}$

→ Total memory capacity = $256\text{Kword} \times 8\text{ bits} = 2\text{Mbit}$

b) Block = 7 bits

→ Total cache words = $2^7 \rightarrow 128$ words (*Block*)

c) Total words = $2^4 = 16$ words

→ Cache word size = Tag + (No. of words in cache × size)

=
=

57

Example 13:

Main memory address :

Tag	Block	Word
-----	-------	------

7 bits 4 bits 3 bits

Consider the main memory address.

Suppose a program generates the address 1AA.

→ In 14-bit binary, this number is: 00 0001 1010 1010

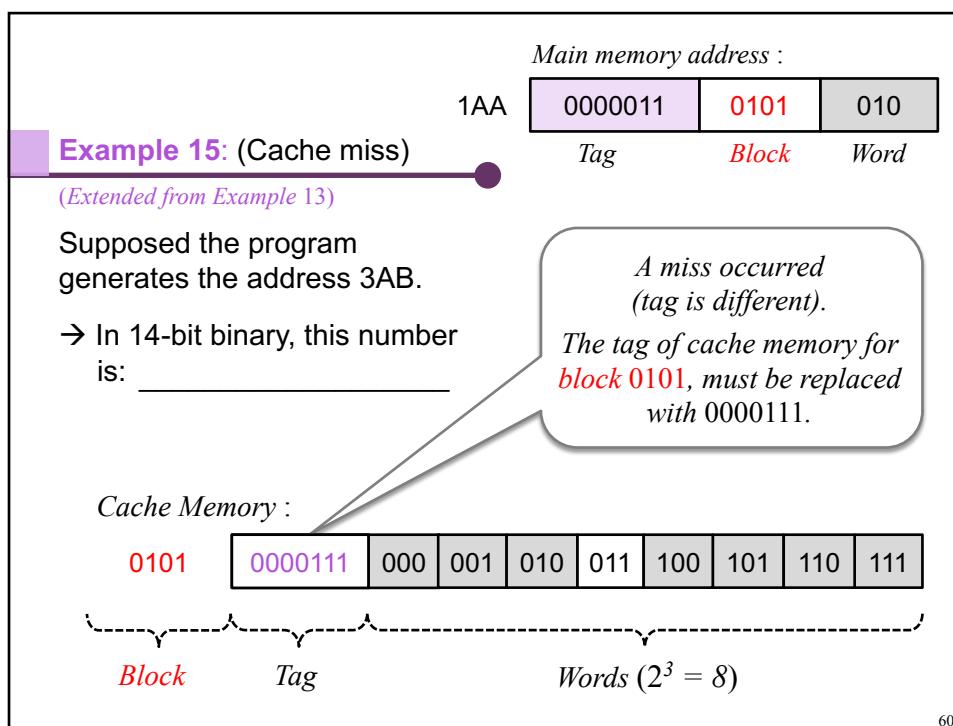
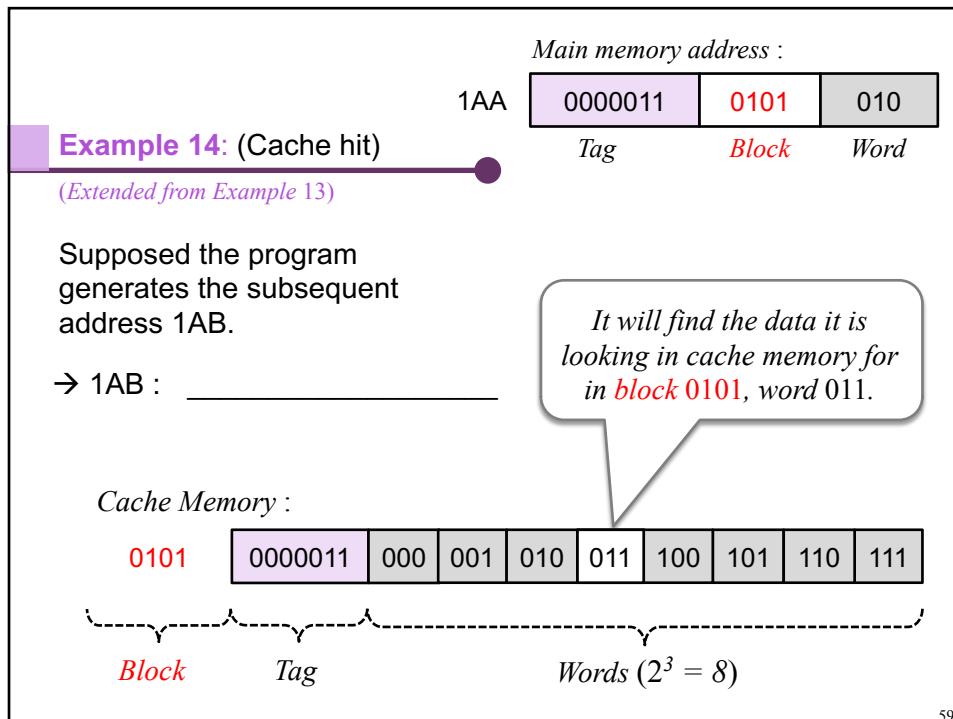
✓ The first 7 bits of the address go in the tag field.
✓ The next 4 bits go in the **block** field.
✓ The final 3 bits indicate the word within the block.

Padding with bit 0 at MSB to complete 14-bit

Main memory address :

--	--	--

58



Block Direct Mapping:

Advantages and Disadvantages



The technique is simple and inexpensive to implement.



- ❑ A fixed cache location for any given block. If a program accesses 2 line/slot/block that map to same line repeatedly, cache hit ratio will be low (Known as **thrashing**).
- ❑ Least effective in its utilization - that is, it may leave some cache lines unused because a given line/slot has fixed location.

*Other cache mapping schemes are designed to prevent this kind of **thrashing**.*

William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.138

61

6

62

m @ May 2023

31

6

(c) Fully Associative Mapping

- Overcomes the disadvantage of direct mapping by permitting each main memory block to be loaded into any location (line/slot/block) of the cache.
- A memory address is partitioned into only two fields:

Main memory address :

Tag	Word
-----	------

- With associative mapping, there is flexibility as to which block to replace when a new block is read into the cache.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.245
William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.138
Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.403

63

6

- When the cache is searched for a specific main memory block, the tag field of the main memory address is compared to all the valid tag fields in cache;

- ✓ if a match is found, _____.
- ✓ If there is no match, we have a _____ and the block must be transferred from main memory.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.245

64

Example 16:

Consider a memory configuration with 2^{14} words, a cache with 16 blocks, and blocks of 8 words.

- Each block = 8 words = $2^3 \rightarrow 3$ bits (*Word*)
- Tag size = $14 - 3 = 11$ bits

Main memory address = 14 bits		Cache Memory :	
Tag	Word	Tag	Word 0
		Word 1	
		Word 2	
		Word ($2^n - 1$)	

11 bits 3 bits

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.245 65

Example 17:

Main memory address :

Tag	Word

11 bits 3 bits

Supposed the program generates the memory address 666h.

→ In 14-bit binary, this number is: _____

It will find the data it is looking in cache memory for in tag 0CCh, word 110.

Cache Memory :

0CCh	000	001	010	011	100	101	110	111
------	-----	-----	-----	-----	-----	-----	------------	-----

Tag Words ($2^3 = 8$)

66

Example 18:

Consider a memory configuration with 2^{14} words, and blocks of 8 words. For data in cache as follows, what is the memory address (hexadecimal).

Cache Memory : (Data)

256h	000	001	010	011	100	101	110	111
------	-----	-----	-----	-----	-----	-----	-----	-----

Tag Words

→ In 14-bit binary, the number 256h is: 0010 0101 0110 **101**

→ The memory address = 001 0010 1011 **0101**
 $= \underline{\hspace{2cm}} \text{ h}$

67

Fully Associative Mapping:

Advantages and Disadvantages

6



Flexibility scheme in term of which block to be replaced when a new block is read into the cache.



Every line's tag is examined for a match in fully associative cache (associative mapping), cache searching gets **expensive**.



69

(d) Set Associative Mapping



$N \rightarrow$ number of
lines in each set

- In this scheme, instead of mapping anywhere in the entire cache, a memory reference can map only to the subset of cache slots.
- Set Associative Mapping is a compromise that exhibits the strengths of both the _____ and _____ approaches while reducing their disadvantages.

- A set-associative cache with N locations for a block is called an N -way set-associative cache.
- An N -way set-associative cache consists of a number of sets, each of which consists of N blocks.
- A block is directly mapped into a set, and then all the blocks in the set are searched for a match.

- In set-associative cache mapping, the main memory address is partitioned into three pieces:

Main memory address :



Cache memory : 1-way

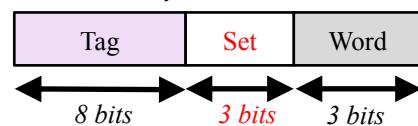


Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.246
 Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.403

71

Example 19: 2-way

Main memory address :



- Suppose we have a main memory of 2^{14} bytes.
- It is mapped to a 2-way set associative cache having 16 blocks where each block contains 8 words.
- Since this is a 2-way cache, each set consists of 2 blocks, and there are 8 sets.
- Thus, we need 3 bits for the set, 3 bits for the word, giving 8 leftover bits for the tag.

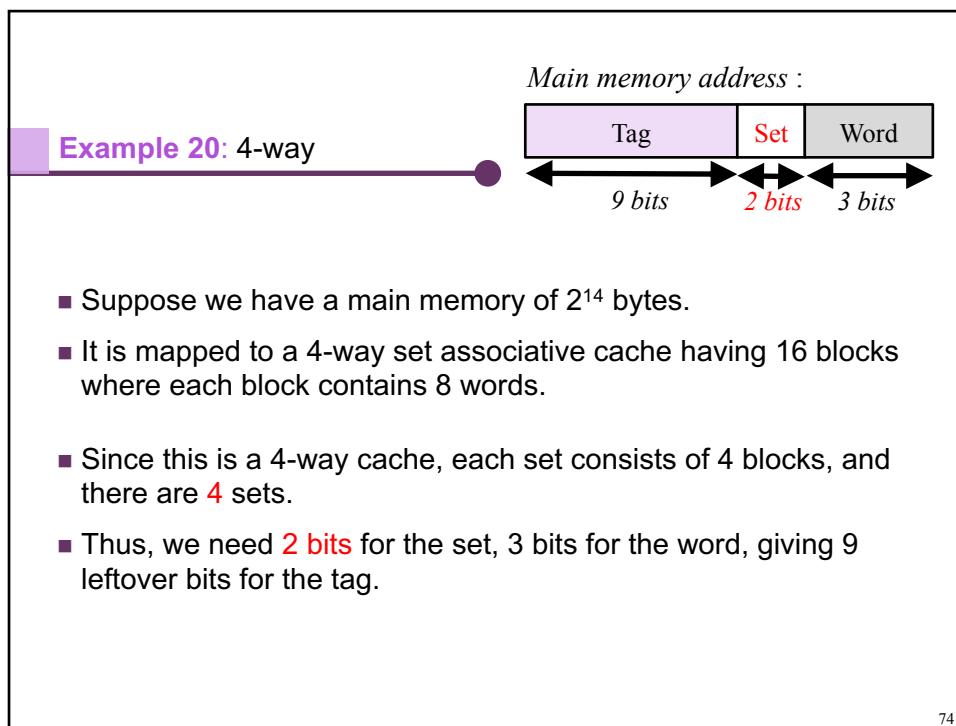
72

		<u>(Hex.) Memory Address:</u>		
		Tag	Set	Word
	3FFF	11111111111111		
	3FFE	11111111111110		
	3FFD	11111111111101		
	3FFC	11111111111100		
	3FFB	11111111111011		
	3FFA	11111111111010		
	3FF9	11111111111001		
	3FF8	11111111111000		
	3FF7	11111111110111		
	3FF6	11111111110110		
	3FF5	11111111110101		
	3FF4	11111111110100		
	3FF3	11111111110011		
	3FF2	11111111110010		
	3FF1	11111111110001		
	3FF0	11111111110000		

Cache memory (Hex.):

Set	Tag	Word	Tag	Word
7				
7				
7				
7				
6	FF	7	FF	6
6	FF	5	FF	4
6	FF	3	FF	2
6	FF	1	FF	0

73



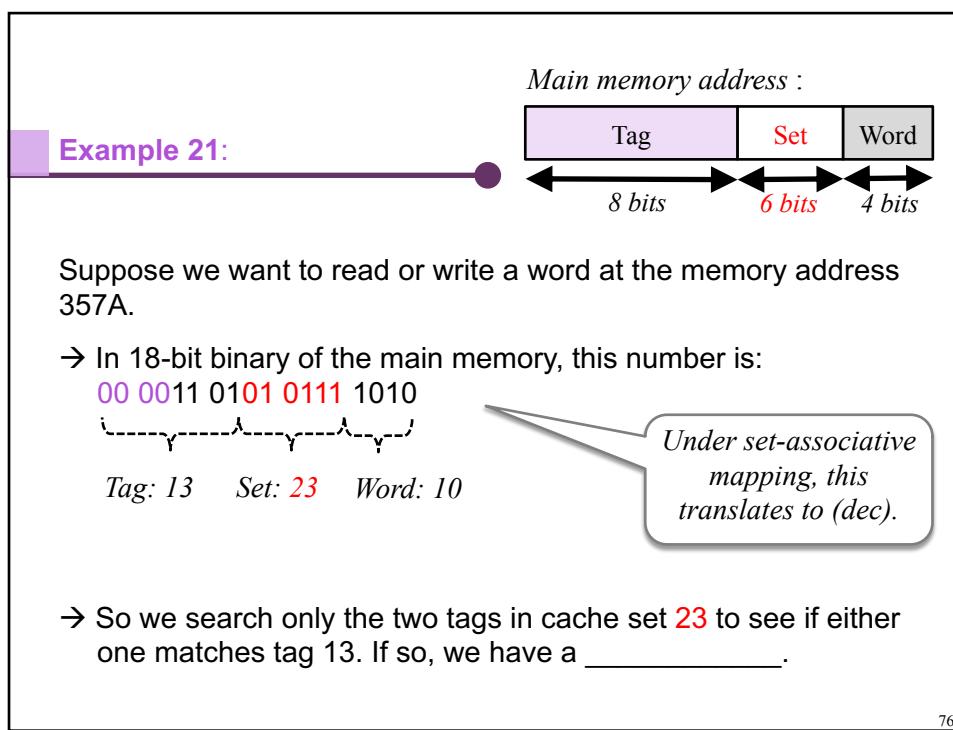
74

(Hex.)	<u>Memory Address:</u>	(Hex.)	<u>Memory Address:</u>
3FFF	11111111111111	3FF7	11111111110111
3FFE	11111111111110	3FF6	11111111110110
3FFD	11111111111101	3FF5	11111111110101
3FFC	11111111111100	3FF4	11111111110100
3FFB	11111111110111	3FF3	11111111110011
3FFA	11111111110101	3FF2	11111111110010
3FF9	11111111110001	3FF1	11111111110001
3FF8	11111111110000	3FF0	11111111110000

Cache memory (Hex.):

Set	Tag	Word	Tag	Word	Tag	Word	Tag	Word
3								
3								
2	FF	7	FF	6	FF	5	FF	4
2	FF	3	FF	2	FF	1	FF	0

75



76

Figure: An eight-block cache configured as direct mapped, two-way set associative, four-way set associative, and fully associative.

**One-way set associative
(direct mapped)**

Block	Tag	Data
0		000
1		001
2		010
3		011
4		100
5		101
6		110
7		111

Two-way set associative

Set	Tag	Data	Tag	Data
0		000		001
1		010		011
2		100		101
3		110		111

Four-way set associative

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0		000		001		010		011
1		100		101		110		111

Eight-way set associative (fully associative)

Tag	Data														
	000		001		010		011		100		101		110		111

Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.404

77

6

7

8

78

6

Replacement Policy

- *Once the cache has been filled, when a new block is brought into the cache, one of the existing blocks must be replaced.

How do we determine which block in cache should be replaced?



- The algorithm for determining replacement is called the *replacement policy*.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.247
* William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited. p.145

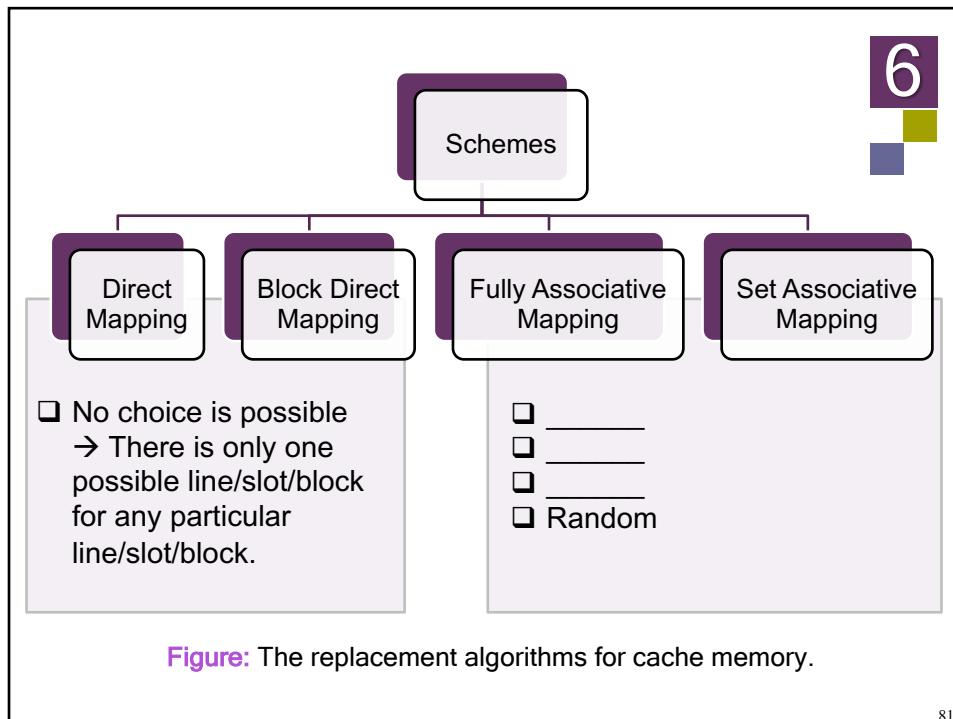
79

6

- A *replacement policy* is invoked when it becomes necessary to evict a line/slot/block from cache.
- *There are several popular replacement policies:
 - One that is not practical but that can be used as a benchmark by which to measure all others is the optimal algorithm.
- Although it is impossible to implement an optimal replacement algorithm, it is instructive to use it as a _____ for assessing the efficiency of any other scheme.

* Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.247
William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited. p.145

80

**Figure:** The replacement algorithms for cache memory.

81

Algorithms	Operational
FIFO <i>(First In First Out)</i>	Replace block that has been in cache longest.
LRU <i>(Least Recently Used)</i>	Keeps track of the last time that a block was accessed and evicts block that has been unused for the longest period of time.
LFU <i>(Least Frequently Used)</i>	Replace block which has had fewest hits.
Random	Picks a block at random and replaces it with a new block.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.248
William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited. p.145

82

6



Which replacement algorithm is the best?

- ✓ The algorithm selected often depends on how the system will be used.
- ✓ No single (practical) algorithm is best for all scenarios.
- ✓ For that reason, designers use algorithms that perform well under a wide variety of circumstances.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.248

83

6

Cache Write Policies

- In addition to determining which victim to select for replacement, designers must also decide what to do with so-called _____ of cache, or blocks that have been modified in cache.
- Dirty blocks must be written back to memory.
- A **write policy** determines how this will be done.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.249
*William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited. p.145

84

6

- *There are two basic write policies:

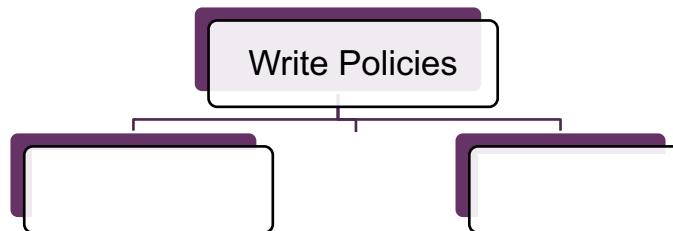


Figure: Cache write policy techniques

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.249
*William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.145
Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.393

85

6

Cache Write Policies:

(a) Write Through

- A write-through policy updates both the _____ and the _____ simultaneously on every write.



Every write to the cache requires a main memory access, essentially **slows** the system down to main memory speed.

- However, in real applications, the majority of accesses are reads so this slow-down is negligible.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.249

86

Cache Write Policies:

(b) Write Back

- A write-back policy (also called _____) only updates blocks in main memory when the cache block is selected as a victim and must be removed from cache.

 <ul style="list-style-type: none"> <input type="checkbox"/> Normally faster than <i>write-through</i> because time is not wasted writing information to memory on each write to cache. <input type="checkbox"/> Memory traffic is also reduced. 	 <ul style="list-style-type: none"> <input type="checkbox"/> Main memory & cache may not contain the same value at a given instant of time. <input type="checkbox"/> Data in cache may be lost, if a process terminates (crashes) before the write to main memory is done.
---	---

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.249

87

Cache Performances

6

- To improve the performance of cache, one must increase the _____ by using :

<ul style="list-style-type: none"> <input type="checkbox"/> a better mapping algorithm (up to roughly a 20% increase), <input type="checkbox"/> better strategies for write operations (potentially a 15% increase), <input type="checkbox"/> better replacement algorithms (up to a 10% increase), and <input type="checkbox"/> better coding practices (up to a 30% increase in hit ratio). 	
---	--

- Simply increasing the size of cache may improve the hit ratio by roughly 1– 4%, but is not guaranteed to do so.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.249

88

6

Average Memory Access Time (AMAT)

- To capture the fact that the time to access data for both **hits** and **misses** affects performance, designers sometime use AMAT as a way to examine alternative cache designs.

$$\text{AMAT} = \text{Time for a hit} + (\text{Miss rate} \times \text{Miss penalty})$$

- AMAT is the average time to access memory considering both _____ and _____ and the frequency of different accesses.

Patterson, D.A. and Hennessy, J.L. (2014). *Computer Organization and Design: The Hardware/Software Interface* (5th Edition). United States: Elsevier, p.402

89

Example 22:

Assume that 33% of the instructions in a program are data accesses. The cache hit ratio is 97% and the hit time is one cycle, but the miss penalty is 20 cycles.

$$\begin{aligned}\text{AMAT} &= \text{Time for a hit} + (\text{Miss rate} \times \text{Miss penalty}) \\ &= \\ &= \\ &= \end{aligned}$$

If the cache was perfect and never missed, the AMAT would be one cycle. But even with just a 3% miss rate, the AMAT here increases 1.6 times!

90

6

Effective Access Time (EAT)

- The performance of a **hierarchical memory** is measured by its EAT, or the average time per access.
- EAT is a weighted average that uses the hit ratio and the relative access times of the successive levels of the hierarchy.

$$EAT = (H \times Access_C) + ((1 - H) \times Access_{MM})$$

*H : Cache hit rate
Access_C : Access time for cache
Access_{MM} : Access time for main memory*

This formula can be extended to apply to three- or even four-level memories

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.248
 William Stallings (2016). *Computer Organization and Architecture: Designing for Performance* (10th Edition). United States: Pearson Education Limited, p.147

Example 23:

Suppose the cache access time is 10ns, main memory access time is 200ns, and the cache hit rate is 99%.

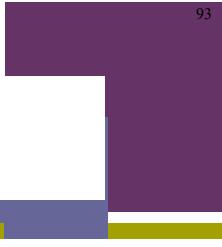
The average time for the processor to access an item in this two-level memory would then be:

$$\begin{aligned}
 EAT &= (H \times Access_C) + ((1 - H) \times Access_{MM}) \\
 &= \\
 &= \\
 &=
 \end{aligned}$$

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers, p.248

Module 6

Memory



6.1 Introduction
6.2 Main Memory
6.3 Cache Memory
6.4 Virtual Memory
6.6 Summary

□ Overview
□ A Real-World Example:
Pentium



■ We now know that caching allows a computer to access frequently used data from a smaller but faster cache memory. (Cache is found near the top of our memory hierarchy).

■ Another important concept inherent in the hierarchy is _____.

The purpose of virtual memory is to use the hard disk as an extension of RAM, thus increasing the available address space a process can use.



Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.250

- Virtual memory can be implemented with different techniques:

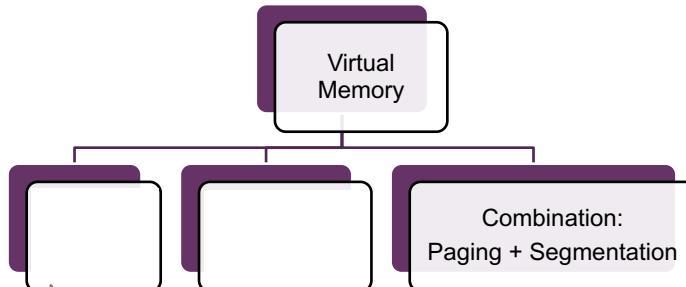


Figure: Virtual memory techniques

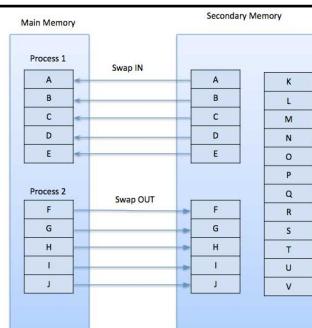
*The most popular
techniques !*

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.251

95

Paging

- The most common way to implement virtual memory is by using *paging*.
- _____ is divided into fixed-size blocks and programs are divided into the same size blocks.
- Typically, chunks of the program are brought into memory as needed.
- It is not necessary to store contiguous chunks of the program in contiguous chunks of main memory.
- Every *virtual address* for a program that generated by CPU must be translated into a _____.



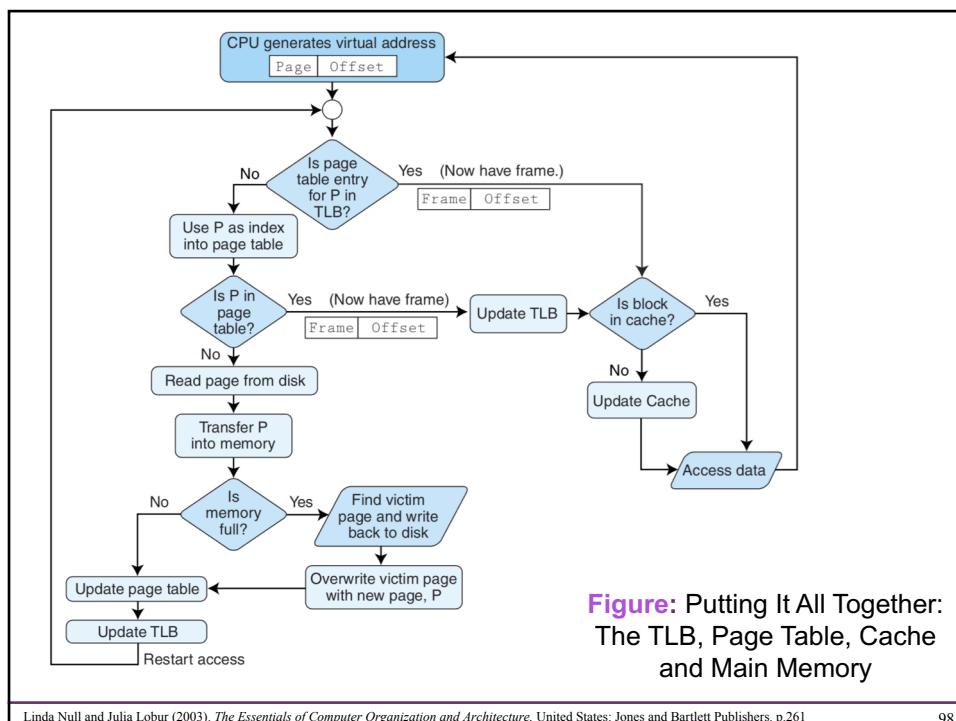
96

Table: Some frequently used terms for virtual memory implemented through paging.

Terms	Description
Virtual address	The logical or program address generated by CPU.
Physical address	The real address in physical memory.
Mapping	Mechanism by which virtual addresses are translated into physical memory.
.....
Pages	The equal-size chunks in virtual memory.
Paging	The process of copying a virtual page from disk to a page frame in main memory.
.....	Memory that becomes unusable.
.....
	An event when a requested page is not in main memory and must be copied into memory from disk.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.250-251

97



Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.261

98

6

Example 25:

Suppose a main memory access requires 200ns and that the page fault rate is 1% (99% of the time with the pages needed are in memory). Assume it costs about 10ms to access a page not in memory (this time of 10ms includes the time necessary to transfer the page into memory, update the page table, and access the data).

The effective access time for a memory access is now:

$$\begin{aligned} EAT &= 0.99(200\text{ns} + 200\text{ns}) + 0.01(10\text{ms}) \\ &= \\ &= \\ &= \end{aligned}$$

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.258

99

Example 26:

Based on previous example, even if 100% of the pages were in main memory, the effective access time would be:

$$\begin{aligned} EAT &= 100\%(200\text{ns} + 200\text{ns}) + 0 \\ &= \\ &= \end{aligned}$$

which is double the access time of memory.

Accessing the page table costs an additional memory access because the page table itself is stored in main memory.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.258

100

6

A Real-World Example: Pentium

- The Pentium architecture exhibits fairly characteristic traits of modern memory management.
 - 32-bit virtual addresses and 32-bit physical addresses.
 - Uses either 4KB or 4MB page sizes, when using paging in different combinations of paging and segmentation.
 - Two caches, L1 and L2, both utilizing a 32-byte block size.
 - Both L1 caches (*I-cache* and *D-cache*) utilize an LRU bit for dealing with block replacement.
 - Each L1 cache has a TLB (*Translation Lookaside Buffer*).
 - The L2 cache can be from 512KB up to 1MB.

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.263

101

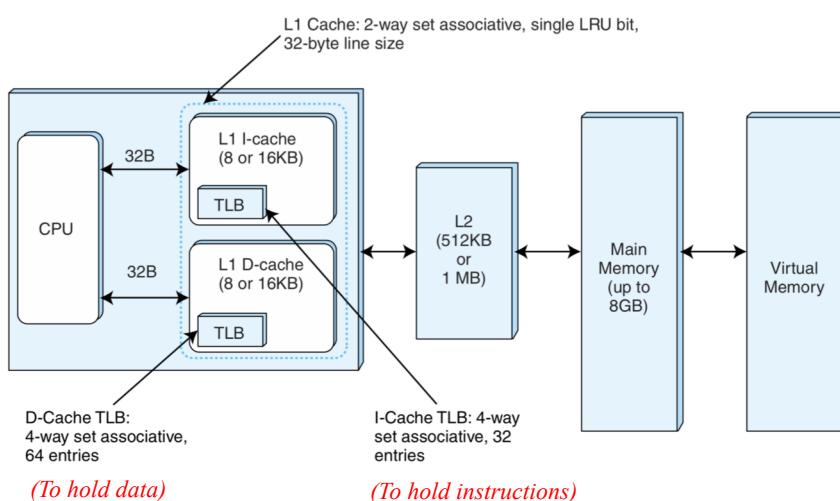


Figure: Pentium memory hierarchy

Linda Null and Julia Lobur (2003). *The Essentials of Computer Organization and Architecture*. United States: Jones and Bartlett Publishers. p.264

102

6.5 Summary

6

- Computer memory is organized in a hierarchy, with the smallest, fastest memory at the top and the largest, slowest memory at the bottom.
- Cache memory gives faster access to main memory, while virtual memory uses disk storage to give the illusion of having a large main memory.
- Cache maps blocks of main memory to blocks of cache memory. Virtual memory maps page frames to virtual pages.
- There are three general types of cache: Direct mapped, fully associative and set associative.

103

6

- With fully associative and set associative cache, as well as with virtual memory, replacement policies must be established.
- Replacement policies include LRU, FIFO, or LFU. These policies must also take into account what to do with dirty blocks.
- All virtual memory must deal with fragmentation, internal for paged memory, external for segmented memory.

104