

MEMORY ORGANISATION

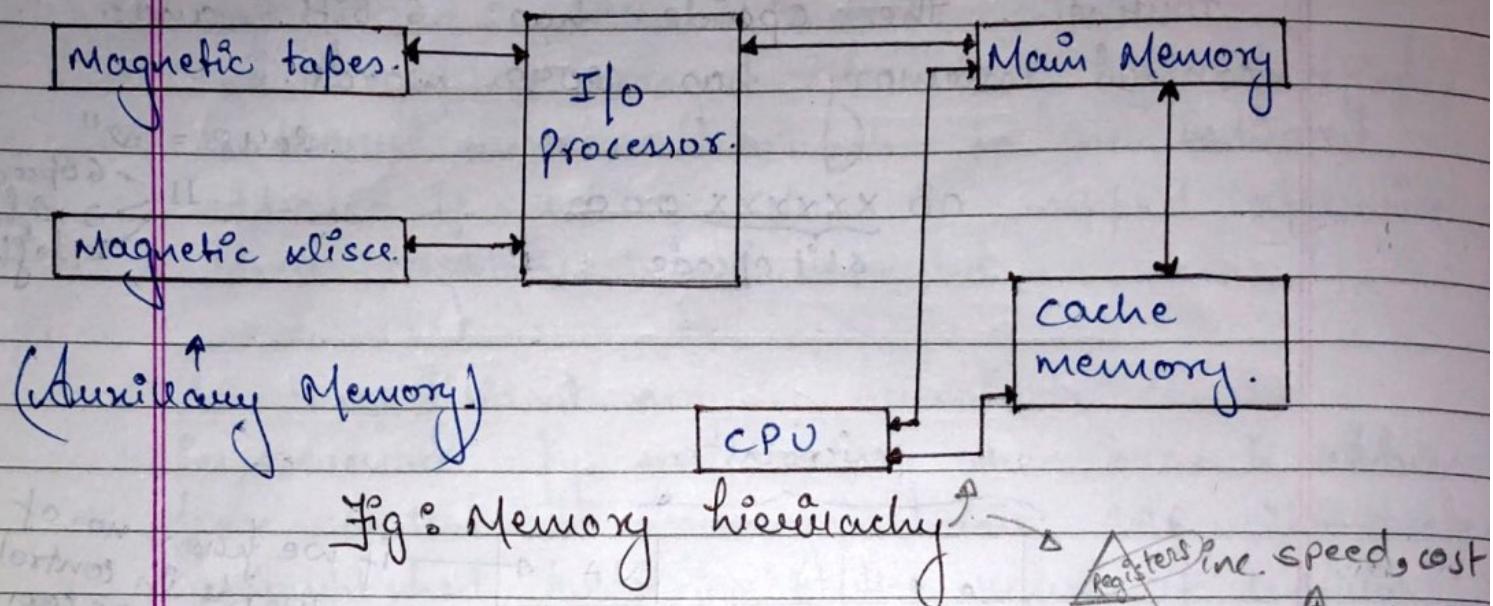
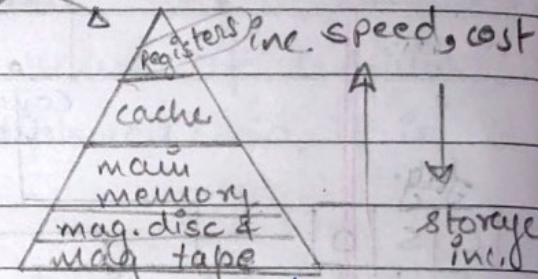


Fig: Memory hierarchy

→ The main goal of memory hierarchy is to obtain highest possible average access speed by minimizing the total cost of entire memory system.



Main Memory.

It is relatively large & fast memory used to store programs & data during the computer operation.

The principle technology used for main memory is based on semiconductor integrated circuits.

Integrated RAM chips are available in two possible modes - static & Dynamic.

→ The static RAM consists of internal flip flops that stores the binary information. The stored information remains valid as long as power is applied to the unit.

→ The Dynamic Ram stores the binary information

$$\begin{aligned}
 & \text{1024} \times 8 \\
 & 512 \text{ bytes RAM} \rightarrow 128 \times 8 = \frac{512 \times 8}{128 \times 8} = 4 \text{ RAM chip} \\
 & 512 \text{ bytes ROM} \rightarrow 512 \times 8 = 1 \text{ ROM chip.}
 \end{aligned}$$

Date: _____
Page: _____

in the form of electric charges that are applied to the capacitors. The stored charge on capacitors tends to discharge with time & the capacitors must be periodically recharged by refreshing the dynamic memory.

Volatile. RAM is used for storing bulk of programs & data that are subject to change.

ROM is used for storing programs that are permanently resident in computer.

Non-volatile. ROM portion of main memory is needed for storing an initial program called a bootstrap loader.

The bootstrap loader is a program whose function is to start the computer software operating when power is turned ON.

RAM chips & ROM chips.

1024x8 memory constructed with 128x8 RAM chips & 512x8 ROM chips.

[space taken by 128 RAM = space taken by 512 ROM]
 we can take 128 ROM chip also
 but to prove above equality we're designing with 512 ROM.

17/10/18

RAM chip

A RAM chip has one or more control inputs that select the chip only when needed.

RAM chip select only when
 $\overline{CS1} = 1$
 $\overline{CS2} = 0$

Date: _____
Page: _____

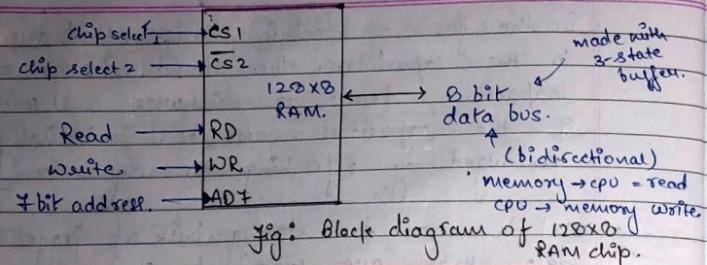


Fig: Block diagram of 128x8 RAM chip.

CS1	CS2	RD	WR	Memory function	state of data bus
0	0	x	x	Inhibit (high)	high impedance state
0	1.	x	x	Inhibit	high impedance state.
1	0	0	0	Inhibit	high impedance state.
1	0	0	1	Write	Input data to RAM.
1	0	1	x	Read.	Output data from RAM
1	1	x	x	Inhibit	high impedance state

Functional Table

Rom chip.

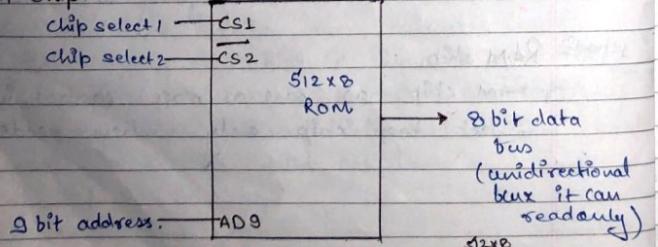


Fig: Block diagram of ROM chip.

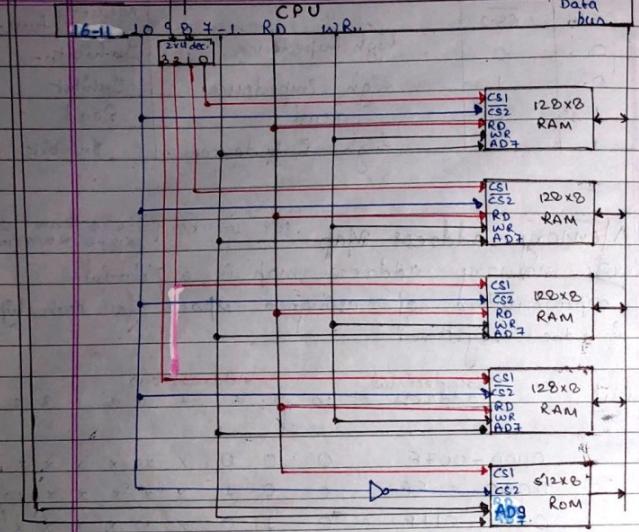
* Internal cells of Rom chip take less space compared to Ram chip.

Date: _____
Page: _____

CS1	CS2		Memory func.
0	0	high impedance	Inhibit
0	1	high impedance	Inhibit
1	0	Read	• Read
1	1	high impedance	Inhibit

Memory address Map. for 1024×8 [128×8 RAM 512×8 ROM 512×8]
A memory address map is a pictorial representation of assigned space for each chip in the system.

Memory connection to CPU



- Ques.
- How many 128x8 RAM chips are needed to provide a memory capacity of 2048 bytes.
 - How many lines of address bus must be used to access 2048 bytes of memory.
 - How many of these lines will be common to all chips.
 - How many lines should be decoded for chip select & specify the size of decoder.

a) $2048 / 128 = 16$

b) $\frac{14}{common} + 4 \text{ select for } 16 \text{ RAM} = 7 + 4 = 11$

c) 4 \leftarrow 4x16 decoder.

Ques. A computer uses RAM chips of 1024x8 capacity.

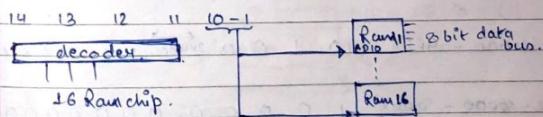
- How many chips are needed & how should these address lines be connected to provide a memory capacity of 1024 bytes. (common 3 decoder to 1024)
- How many chips are needed to provide a memory capacity of 16K bytes. Explain in words how the chips are to be connected to the address bus.

g) Ram chip = $\frac{1024 \times 8}{1024 \times 1} = 8$ RAM chip.

8 chips are needed with add. line connected in parallel.

b) $\frac{16K \times 8}{1024 \times 1} = \frac{2^4 \times 2^{10} \times 8}{2^{10}} = 128$ RAM chip.

$16K = 2^{14}$



Set 16, each 8 parallelly connected.

128 chips are needed use 14 address lines, 10 lines specify chip address, 4 lines are decoded into 16 chip select lines.

Ques. A computer employs RAM chips of 64×8 & ROM chips of 1024×8 . The comp. sys. needs 8K bytes of RAM, 4K bytes of ROM and 4 interface units. each with 16 registers. A memory mapped I/O configuration is used. The Highest order bits of address bus are assigned 00 for RAM, 01 for ROM & 10 for Interface. How many RAM & ROM chips are needed.

- Draw a memory address map for system.
- Give address range in hexadecimal for RAM, ROM & Interface.

$$a) \text{ RAM chips} = \frac{8 \times 1024 \times 8}{256 \times 8} = 8$$

$$\text{ROM chip} = \frac{4 \times 1024 \times 8}{1024 \times 8} = 4$$

$$\text{Interface reg.} = 4 \times 4 = 16$$

ROM 4000 - 4FFF ROM 100 2x1 dec. x #3 x x x x x x x x

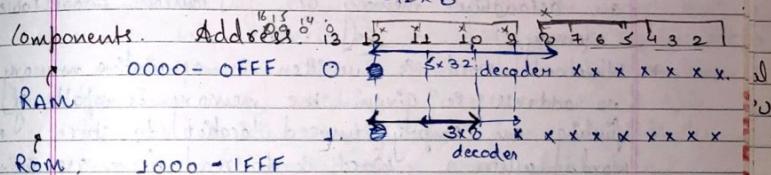
Interface: 8000 - 800F I ^{Interval} 0 0 0 0 0 0 0 0 0 0 x x x x

RAM₁ 0000-00FF ROM1 → 4000-43FF
RAM₂ 0100-01FF ROM2 → 4400-47FF

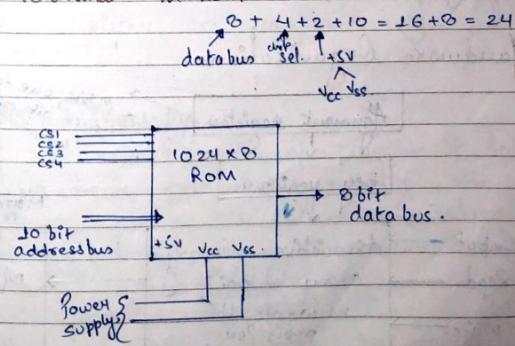
Ques. Draw memory address map for 4096 byte of RAM & 4096 byte of ROM.
Size of 1 RAM chip 128×8 , Ram chip 512×8

$$\text{No. of RAM} = \frac{4096 \times 8}{128 \times 8} = 32$$

$$\text{ROM chip} = \frac{4096 \times 8}{512 \times 8} = 8$$



Ques. A ROM chip of 1024×8 bits has 4 select inputs and operates from a 5V power supply. How many pins are needed for the IC package. Draw a block diagram & label its input and output terminal in ROM.

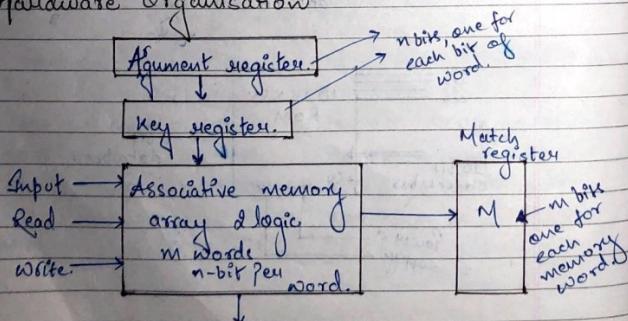


Temporal Locality - means current data or instr that is being fetched may be needed soon. So we should store that data in cache memory to avoid again searching in main memory to retrieve data and saving time.

Associative Memory (Content Addressable Memory) CAM

- The time required to find an item stored in memory can be reduced if stored data can be identified for access by the content of data itself rather than by an address.
- A memory unit accessed by content is called an associative memory or content addressable memory.
- When a word is written in associative memory no address is given. The memory is capable of finding an empty unused location(s) to store the word. When a word is to read from associative memory, the content of the word or part of the word is specified.
- An associative memory is more expensive than a RAM because each cell must have storage capability as well as logic circuits for matching its content with an external argument.

Hardware Organisation



Spatial Locality - instr or data near to the current memory location that is being fetched, may be needed soon (in near future)

Date: _____
Page: _____
15

Sum = 0
for ($i=0$; $i < \text{arr.length}$; $i++$)
sum += arr[i];
return sum;

A: 101 111100
K: 111 000000
word 1. 100 111100 no match spatial locality.
word 2. 101 000001 match.

Each word in memory is compared in parallel with the content of argument register. The words that match will be the k bits of argument register set a corresponding bit in the match register. After the matching process, those bits in match register that have been set indicate that these corresponding words have been matched.

The key register provides a mask for choosing a particular field for key in argument word. The entire argument is compared with each memory word. If key register contains all 1's, otherwise only those bits in argument that have 1's in corresponding position of key register are compared.

CACHE MEMORY

If the active portion of program & data are placed in fast small memory, the average memory access time can be reduced. Such a fast small memory is known as cache memory.

Basic operation of cache

When the CPU needs to access memory, the cache is examined, if the word is

found in cache, it is read from the cache memory. If the word addressed by CPU is not found in cache, main memory is accessed to read the word.

A block of word is then transferred from main memory to cache memory.

$$\text{Hit ratio} = \frac{\text{no. of hits}}{\text{total CPU references.}}$$

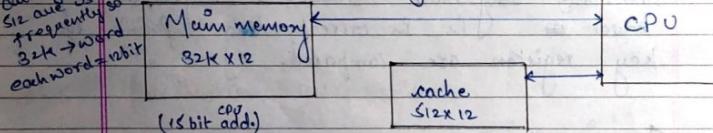
↑
Performance
of cache
memory.

hit ratio
more,
more
performance.

26/10/10

The transformation of data from main memory to cache memory is known as **Mapping**. There are 3 types of mapping -
 1) Associative mapping.
 2) Direct mapping.
 3) Set associative mapping.

Ex
 Out of 32k word
 512 are used
 frequently so in cache
 32k → word
 each word = 12bit



Associative mapping

The fastest & most flexible cache organisation uses an associative memory.

The associative memory stores both the address and content of the words.

1st search in cache, if not found search in main memory provide to CPU or well as store in cache (provided by CPU to search)

CPU address (15 bits) by CPU to search

Argument register

If cache is full we apply replacement method (FIFO, Round Robin)

Address (Octal)	Data (Octal)
01000	3450
02777	6710
22315	1234
03777	7453

Associative Mapping cache

2) Direct mapping.

Associative memory are expensive compared to RAM because of added logic associated with each cell.

(15 bit) CPU address

* In general if we have, n^k words in cache.

Tag. index.
 $(15 - k = 6 \text{ bits})$
 (add. bit of cache)
 is 9-bit

n^k word in main memory
 Index = K bits

Tag = $n - k$ bits.

6 bits shift	
Tag	Index
000000	(6 bits of LSB = Index bit)
111111	
777777	(Octal)

Address	Data	Index address	Tag	Cache	Tag	Data	Index address	CPU
00000	1220	000	00	1220	02777	1220	00000	1st
00777	2340		01	3450				
01000	3450							
01777	4560	***	02	6710				
02000	5670	***						
02777	6710							

CPU tries to check if data is available in cache.
it again gives address 00000
CPU breaks into 2 parts

00 000
Tag. Index → Search in cache,
and matches with tag of cache
(memory & its own tag).
if both same then hit.

If CPU gives 01000 (read from cache)
tag index else find in main
CPU tag = 01 cache tag = 00
at index 000

Hence, miss condition, so finds it
in main memory, tag and data corresponding to
01000 so replaces tag and data of cache &
overwrites that read from main.

(In such practice hit ratio falls it is a disadvantage.)

Date: 26/10/18
Page: 154

CPU 1st
tag, Index
00000

02777
tag, Index

Disadvantage of direct mapping is that the hit ratio can drop if two or more words whose addresses have same index but different tags are accessed repeatedly.

If 1 block = 8 words.

$$\frac{512}{8} = 64 \text{ blocks.}$$

(6) Tag Index (8)

is bit CPU add.

block (6) word.
(Octal (3))
no. of blocks = 2^6
each block has 2^3 words.

* More block size, more hit ratio improvement

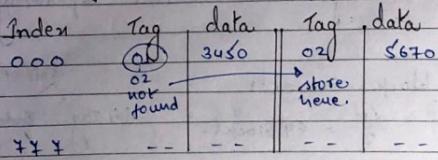
Octal Index	Tag	Data
01	4370	
00	5610	
01	02	+
01	02	-
02	-	
03	-	
03	-	

Tag within each block
should be same.

3) Set associative mapping.

Set associative mapping is an improvement over the direct mapping where each word of cache can store one or more words of memory under the same index address.

Each data word is stored together with its tag and the number of tag data items in one word of cache is said to form a set



Two way set associative.

$$\text{Index} = 512 \quad \text{Tag} = 6 + \text{data} = 4 \times 2 = 12 = 10_2 \quad \text{word} = \text{tag} + (\text{data})_2$$

Two word = $10_2 \times 2$

$$\text{Size of cache memory} = 512 \times 36$$

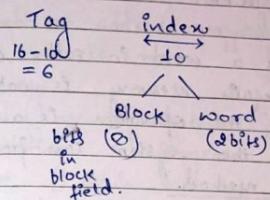
Ques. A digital computer has a memory unit of $64K \times 64$ and a cache memory of $1K$ words. The cache uses direct mapping with a block size of 4 words.

a) How many bits are there in tag, index, block and word field of address format?

b) How many bits are there in each words of cache, and how are they divided into functions. include a valid bit.

c) How many blocks can cache accommodate.
when we power on, cache memory initialize, and it is not empty it has non-valid data.
when in first tries (and we bring data from main) we set valid bit = 1 and if it is equal to 1 then we don't apply any replacement.

$$2^{16} \times 16$$



$$2^8 = 256$$

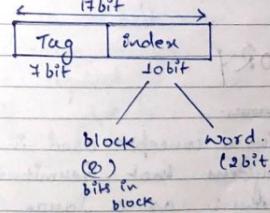
($256 \times 4 = 1024$ in cache) *Vanity.*

$$\text{word} = \text{tag} + \text{data}$$

$$\boxed{V \mid \text{tag} \& \text{data}} \quad 1 + 6 + 16 = 23 \text{ bits.}$$

Ans. A 2-way associative cache memory uses block of 4 words. The cache can accommodate a total of 1024 words from main memory. The main memory size is $128K \times 32$. Formulate all the information required to construct the cache memory what is size of cache memory.

$$\frac{2048}{2} = 1024 = 2^{10} = 10 \text{ bit (index) address.}$$



$$128K \times 32$$

$$2^{10} \times 32$$

$$1024 \times 32$$

$$\text{no. of block} = 2^8 = 256$$

$$\text{data} = \frac{32}{4} = 32 \times 2 = 78$$

$$1024 \times 78$$

Writing into cache

When a CPU finds a word in cache during read operation the main memory is not involved in the transfer. If the operation is a write there are two ways the system can proceed.

1) Write through method.

The simplest and most commonly used procedure is to update main memory with every memory write operation, with cache memory being updated in parallel. If it contains the word at specified address. This is called write through method.

This method has the advantage that main memory always contains the same data as the cache.

2) Write back method.

In this method, only the cache location is updated during a write operation. The location is then mapped by flag. So that, later when the word is removed from cache it is copied into main memory.

VIRTUAL MEMORY

~~is a memory manager.~~

Virtual memory is a concept used in some large computer systems that permit the user to construct programs as a large memory space is available equal to the totality of auxiliary memory.

Date: 30/10/18
Page:

add. space = set of virtual add.
mem. space = set of physical add.

Date:
Page:

Each address that is referenced by the CPU goes through an address mapping from virtual address to a physical address in main memory.

Virtual memory system provides a mechanism for translating program generated addresses into correct main memory locations.

The translation or mapping is handled automatically by hardware by means of a mapping table.

Address space & Memory space

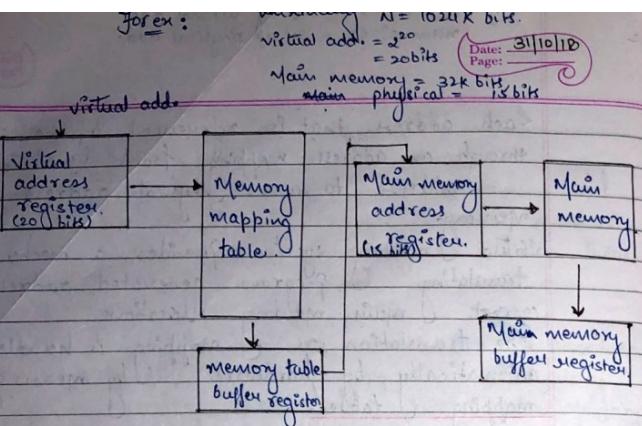
An address used by programmers is virtual address and set of such addresses is address space.

An address in main memory is called location or physical address and set of such location is called memory space.

The mapping table may be stored in separate memory or in main memory.

In first case, an additional memory unit is required as well as one extra memory access time. In second case, the table takes space from main memory and two accesses to memory are required. with the program running at half speed.

For ex:



Address mapping using pages

1 page size = 1 block size = 1K words.

* memory space
↓
blocks
↓
address space
↓
Pages.

$$\text{no. of pages} = \frac{1024 \text{ K}}{1 \text{ K}} = 1024$$

$$\text{no. of blocks} = \frac{32 \text{ K}}{1 \text{ K}} = 32$$

Page → map → block.
virtual add. ← line number
 Page number

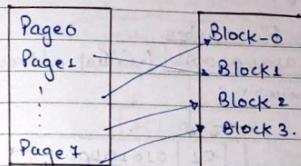
→ The table implementation of address mapping is simplified if the information in the address space and memory space are each divided into groups of fixed size.

→ The physical memory is broken down into groups of equal size called blocks, and the term page refers to groups of address space of same size.

If address space = 8 K, virtual add = 13 bit
memory space = 4 K Physical add = 12 bit
1 page size = 1 block size = 1K word.

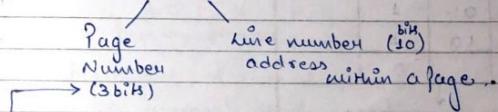
$$\text{no. of pages} = \frac{8 \text{ K}}{1 \text{ K}} = 8$$

$$\text{no. of block} = \frac{4 \text{ K}}{1 \text{ K}} = 4$$



Max. 4 page can reside in main memory.

virtual address (13 bit)



If 1 page has 1024 words
Line address = 10 bit.

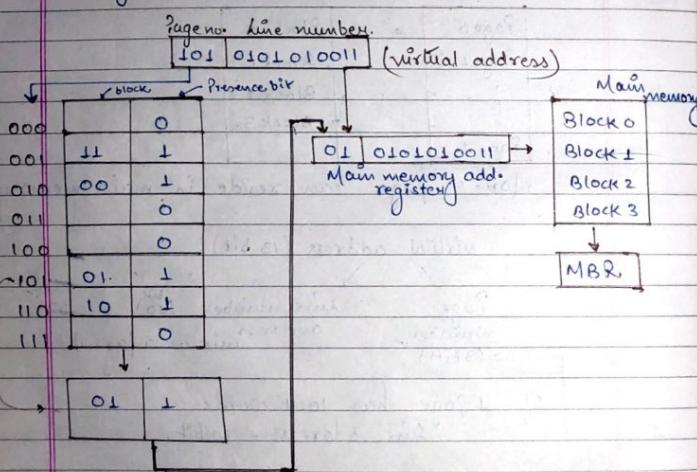
implies no. of pages = 2^3

when Page → block, corresponding line number doesn't change.

Block number = 2 bit (since we've 4 blocks)

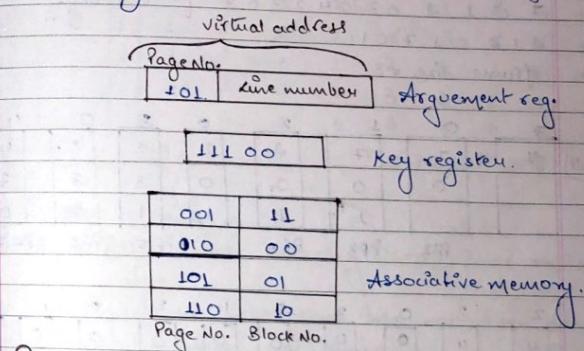
Date: 31/10/18
Page:

The mapping from address space to memory space is facilitated if each virtual address is considered to be represented by a number of page numbers address and a line within the page. In a computer with 2^P words per page, P bits are used to specify line address and the remaining high order bits of virtual address specify the page number.



Date: 1/11/18
Page:

associative memory with each word in memory containing a page number together with its corresponding block number. The page field in each word is compared with page number in virtual address, if a match occurs the word is read from memory and its corresponding block number is extracted.



Page replacement algorithms.

When a program starts execution one or more pages are transferred into main memory and the page table is said to indicate their position. The program is executed from main memory until it attempts to reference a page that is still in auxiliary memory. This condition is called **page fault**. When a page fault occurs in virtual memory system it signifies that the page referenced by the CPU is not in main memory. A new page is then transferred from auxiliary memory to main memory. If main memory is full, it is necessary to

Associative memory Page table. (no storage wastage)
A more efficient way to organise the page table is to construct it with a number of words equal to the number of blocks in main memory. In this way, size of memory is reduced and each location is fully utilized. This method can be implemented by means of an

Remove a page from memory block to make room for new page. The policy for choosing pages to remove is determined from the replacement algorithm.

1) FIFO page replacement.

For ex. reference string =

7 0 1 2 0 3 0 4 2 3 0 3

frame size = 3

Given -
frame size = Max. Page
in main memory.

Reference string.

(PF _i) Page fault	0	1	2	0	3	0	4
	PF ₂	PF ₃	PF ₄	PF ₅	PF ₆	PF ₇	
4	7	7	9	2	2	2	4
	0	0	0	0	3	3	3
	1	1	1	1	0	0	0

2	3	0	3	2	1	2	0	1	7
4	4	9	0	0	0	0	0	0	4
2	2	2	2	2	1	1	1	1	1
0	3	3	3	3	3	2	?	2	2

PF₀ PF₉ PF₁₀ PF₁₁ PF₁₂ PF₁₃

PF₁₄ PF₁₅

Total Page fault = 15

Blady's anomaly.

Reference String.

1 2 3 4 1 2 5 1 2 3 4 5

frame size = 3

Page fault = 9

frame size = 4

Page fault = 10

Date: _____
Page: _____

Date: 2/11/18
Page: _____

1	2	3	4	1	2	5	1	2	3	4	5	2nd
	1	1	1	4	4	4	5	5	3	3	5	e
	2	2	2	2	1	1	1	1	3	3	3	
	3	3	3	3	3	2	2	2	2	4	4	
	3	3	3	3	3	2	2	2	2	4	4	
	4	4	4	4	4	4	4	4	4	4	4	

1	2	3	4	5	2	5	1	2	3	4	5
	1	1	1	1	5	5	5	5	4	4	4
	2	2	2	2	2	2	2	2	1	1	1
	3	3	3	3	3	3	3	3	3	3	3
	3	3	3	3	3	3	3	3	3	3	3
	4	4	4	4	4	4	4	4	4	4	4

2) LRU (Least Recently Used) page replacement algo.

Replace the page that has not been used for longest period of time.

String 8:-

4 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

frame size = 3

4	7	9	2	2	4	4	4	0	0	0	0	1
	0	0	0	0	0	0	0	3	3	3	3	0
	1	1	1	3	3	3	2	2	2	2	2	1
	1	1	1	3	3	3	2	2	2	2	2	2
	1	1	1	3	3	3	2	2	2	2	2	2
	1	1	1	3	3	3	2	2	2	2	2	2

Page fault = 12.

Ques. 3) Optimal page replacement.

Replace the page that will not be used for the longest period of time.

4	0	1	2	0	3	0	4	2	3	0	3	2	1	2
4	4	7	2	2	2	2	2	2	2	2	2	2	2	2
0	0	0	0	0	4	0	0	0	0	0	0	0	0	0
1	1	1	3	1	3	1	3	1	3	1	3	1	3	1
0	1	7	0	1	1	1	1	1	1	1	1	1	1	1

Page fault = 9.

Ques. An address space is specified by 24 bits. and the corresponding memory space by 16 bits.

- How many bits are there in address space?
- How many bits are there in memory space?
- If a page consists of 2K words then how many pages & blocks are there in the system?

a) 2^{24} words

b) 2^{16} words

c) No. of pages = $\frac{2^{24}}{2^{11}} = 2^{13}$

No. of blocks = $\frac{2^{16}}{2^{11}} = 2^5$

Ques.

If virtual memory has a page size of 1K words. There are 8 pages and 4 blocks. The associative memory page table contains the following entries.

Page	block
0	3
1	1
4	2
6	0

Make a list of all virtual addresses in decimal that will cause a page fault if used by CPU. Page no. 2, 3, 5, 7 will cause a page fault.

Pg No.	Add.	Address that will cause Page fault	0 → 0 - 1023
2	2K	2048 - 3071	1 → 1024 → 2047
3	3K	3072 - 4095	2 → 2048 → 3071
5	5K	5120 - 6143	3 → 3072 - 4095
7	7K	7168 - 8191	4 → 4096 - 5119
			5 → 5120 - 6143
			6 → 6144 - 7167
			7 → 7168 - 8191

Ques. Consider a system with 2 level cache. The access time of L1 cache, L2 cache & main memory are 1ns, 10ns, 50ns. The hit ratio of L1 and L2 cache are 0.8 and 0.9 respectively. What is the average access time.

$$t_1 = 1\text{ ns}, t_2 = 10\text{ ns}, t_m = 50\text{ ns}; \text{L1 cache hit} = 0.8, \text{L2 cache hit} = 0.9$$

Hierarchical access

$$T_{avg} = h_{l1} t_{l1} + (1-h_{l1}) h_{l2} (t_{l1} + t_{l2}) + (1-h_{l1})(1-h_{l2}) (t_{l1} + t_{l2} + t_m)$$

$$= 0.8 + 0.2 \times 0.9 \times 11 + 0.2 \times 0.1 \times 511$$

$$= 13\text{ ns}$$

Simultaneous access

$$T_{avg} = h_{l1} t_{l1} + (1-h_{l1}) h_{l2} (t_{l1} + t_{l2}) + (1-h_{l1})(1-h_{l2}) t_m$$

$$= 0.8 + 0.2 \times 0.9 \times 10 + 0.2 \times 0.1 \times 500 = 12.6\text{ ns}$$

Date: _____
Page: _____

what is the average access time of the system
ignoring the search within the cache?

Consider the design of a level hierarchy with the
following specifications for the memory characteristic.

Memory Level	Access Time	Capacity	Cost /KB
Cache	$t_c = 25\text{ns}$	$S_c = 512\text{ KB}$	$C_c = \$125$
Main memory	$t_m = -$	$S_m = 32\text{ MB}$	$C_m = \$0.2$
Disk array	$t_d = 4\text{ms}$	$S_d = -$	$C_d = \$0.0002$

The goal is to achieve an average memory access time, $t_a = 10.04\text{usec}$. with a cache hit ratio $H_c = 0.98$ and $H_m = 0.9$ in main memory.
Also the total of the memory hierarchy is upper bounded by \$15,000. Find out the unknown quantities.

$$t_a = H_c t_c + (1-H_c) H_m t_m + (1-H_c)(1-H_m) t_d \leq 10.04 \times 10^{-6}$$

$$= 0.98 \times 25 \times 10^{-9} + (1-0.98) 0.9 t_m + (1-0.98)(1-0.9) \times 4 \times 10^{-3}$$

$$t_m = 1.11 \times 10^{-4} \text{ sec}$$

$$= 1.11 \text{ usec.}$$

$$512 \times 1.25 + 32 \times 1024 \times 0.2 + S_d \times 0.0002 \leq 15,000$$

$$S_d = 39032000 \text{ KB} = 37.22 \text{ GB}$$

Ques. A virtual memory system has an address space of 8K words. Memory space of 4K words and page & block size of 1K words. The following page reference changes occur during a

Date: _____
Page: _____

given time interval.

4 2 0 1 2 6 1 4 0 1 0 2 3 5 7 find the
page fault if the replacement algo. is
i) FIFO ii) LRU iii) optimal.

i) FIFO

Page size = 4

4	4	4	4	4	6	6	6	6	6	6
2	2	2	2	2	2	2	4	4	4	4
0	0	0	0	0	0	0	0	0	0	0
(1)	(2)	(3)	1	1	1	1	1	1	1	1
(4)	(5)	(6)								

6 6 5 5

4 4 4 7

2 2 2 2

1 3 3 3

(7) (8) (9) (10)

to page faults.

ii) LRU

4	4	4	4	4	4	4	4	4	4	4
2	2	2	2	2	2	2	2	2	2	2
0	0	0	6	6	6	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	3

4 4

2 2

5 5

3 3

optimal.

4	4	4	4	4	4	4	4	4	4	4
2	2	2	2	6	6	6	6	2	2	2
0	0	0	0	0	0	0	0	3	5	5
1	1	1	1	1	1	1	1	1	1	1

9 Page fault.