

Memory System Basics

Introduction—

Memory system is an important part of the computer. Memory system is where the program is stored and retrieved by the CPU for execution. In memory partial and final results are stored after computation.

Memory classifications—

With respect to the way of data access we can classify memories as:

- Random Access Memory (RAM)
- Sequentially accessible Memory (SAM)
- Direct Access Memory (DAM)
- Contents Addressable Memory (CAM)

Random Access Memory → In random access memory the access time to any piece of data is independent to the physical location of data. Access time is constant.

(Access time is the time interval between the instant of data read/write request, and the instant at which the delivery of data is completed or its storage is started.)

Random Access memory is of two types

1. Read write memories (RAM)
2. Read only memories (ROM)

Read-write memory (RAM) is further divided in three parts

- * Static RAM (SRAM)
- * Dynamic RAM (DRAM)
- * Non-Volatile RAM (NVRAM)

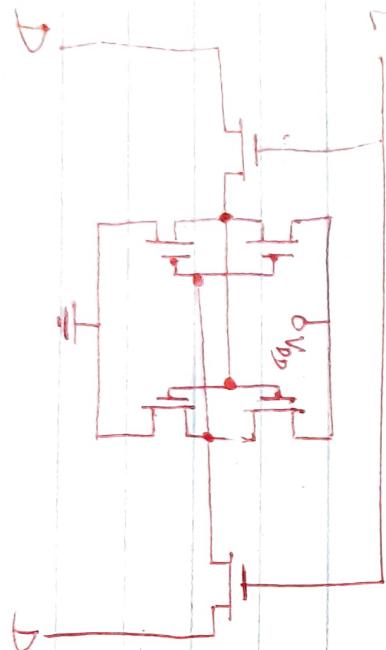
Static Random Access Memory (SRAM)

In case of static random access memory (SRAM), the data is stored essentially in flip-flop circuits.

DRAM

Each one-bit memory cell uses a capacitor for data storage. Since capacitors leak, there is a need to refresh.

SEL



D

Schematics of one-bit cell for static RAM

SEL

INPUT



C

Memory cell for DRAM

DRAM required to be refreshed every 2 or 3 ms, to keep the data or the charge in the capacitance.

Pros & Cons

Static RAM

1. faster and less power hungry than DRAMs.
2. More expensive (6 transistors/cell)

DRAM

- 1. less expensive (only one transistor per cell)
- 2. slower than SRAM

Programmable read only memories (PROM)

PROM are

programmed during manufacturing process. The contents of each memory cell is locked by a fuse or antifuse (diodes). PROMs are used for permanent data storage.

Erasable read only memory (EPROM)

Can

Erase with ultraviolet light (about 20 minutes). Memory cells are built with floating gate transistors. Data can be stored in EPROMs for about 10 years.

Electrically erasable read only memory (EEPROM)

Erasing does not require ultraviolet light. But higher voltage and can be applied not to the whole circuit but to single memory cell separately.

Lecture-28

Non-volatile access memory (NVRAM) :-

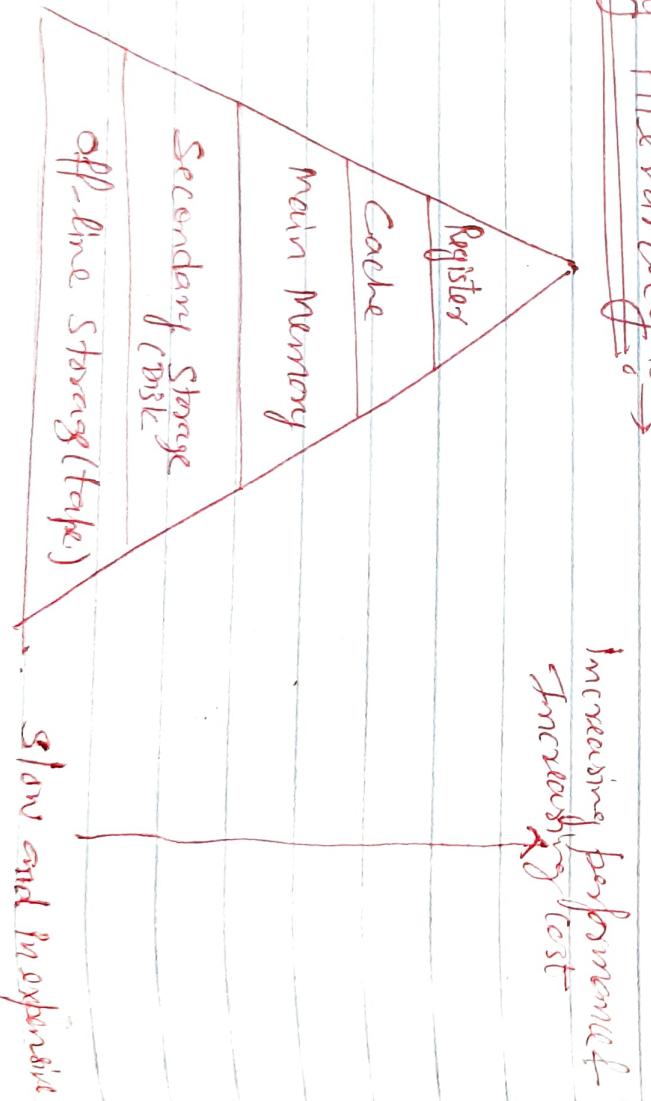
With this

Name we refer generally to any memory which does not lose information when power is turned off. Except from ROMs such as BIOS memory (Basic Input Output System).

Flash Memory - (FLASH) -

By this name the cheaper variant of EEPROM is described. In case of FLASH memory not separate bytes but blocks of bytes are being erased at the same time. It makes the construction of such memories cheaper in comparison to regular EEPROMs.

Memory Hierarchy →



Memory type Access time Cost / mB Typical amount used Typical cost

Registers

0.5 ns

High

2 kB

—

Cache

5-20 ns

\$1.80

2 MB

\$160

Main memory

40-80 ns

\$0.40

512 MB

\$205

Disk memory

5 ms

\$0.005

40 GB

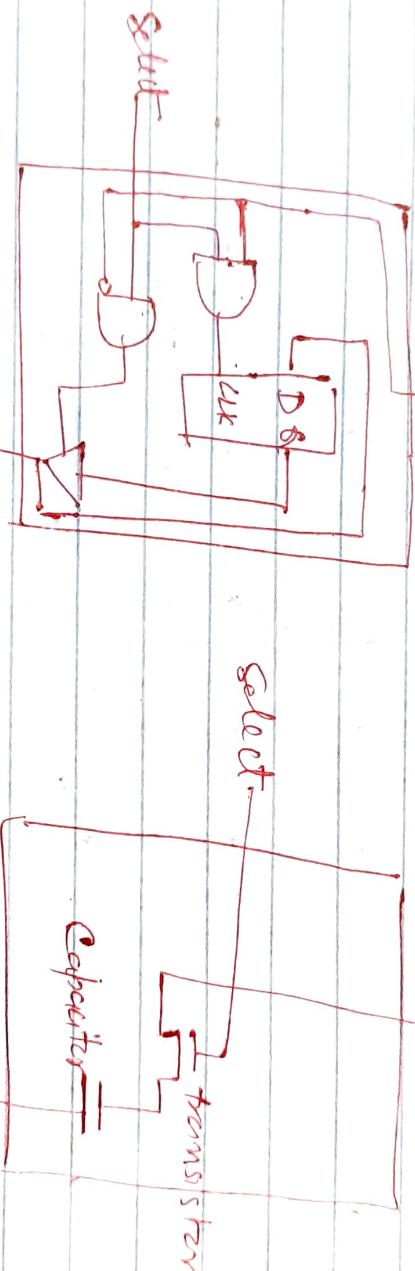
\$200

→

→

Read

Read

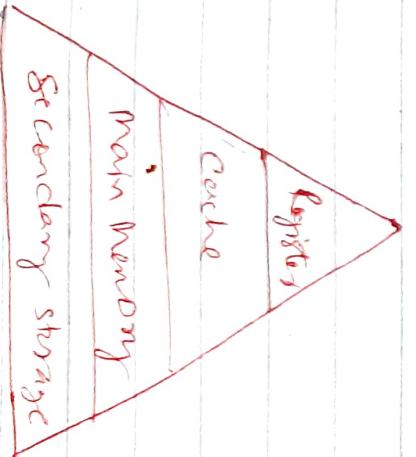


→ Data In/Out → Data In/Out

DRAM

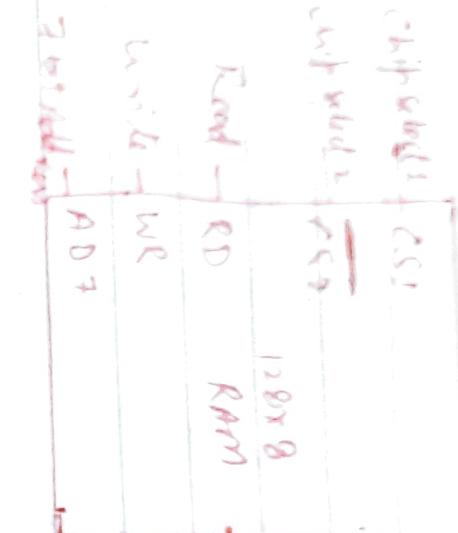
→ RAM

Optimization of memory Hierarchy →



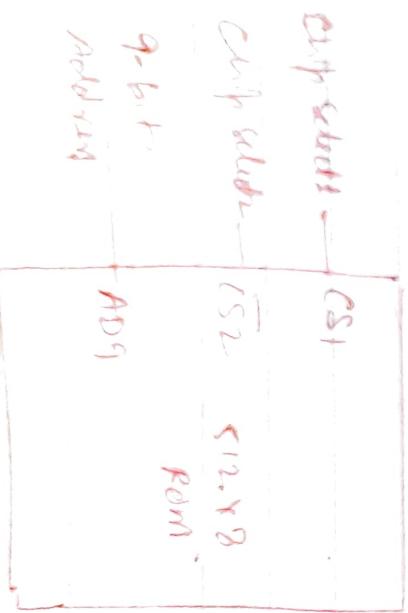
Lecture - 29

RAM Chip

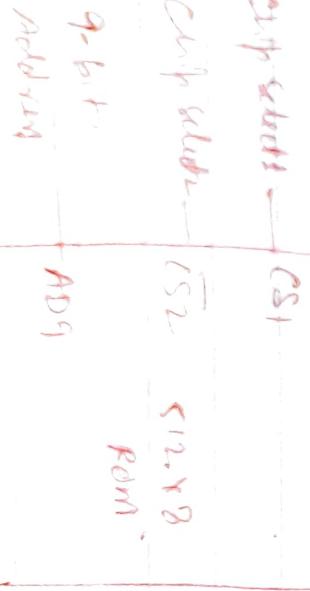


CS ₁	CS ₂	RD	WR	Memory function	State of Data bus
0	0	X	X	Inhibit	High-impedance
0	1	X	X	Inhibit	High-impedance
1	0	0	0	"	"
1	0	0	1	Write	Input data in RAM
1	1	1	X	Read	Output data from RAM
1	1	X	X	Inhibit	High-impedance

RAM Chip



8-bit Data bus



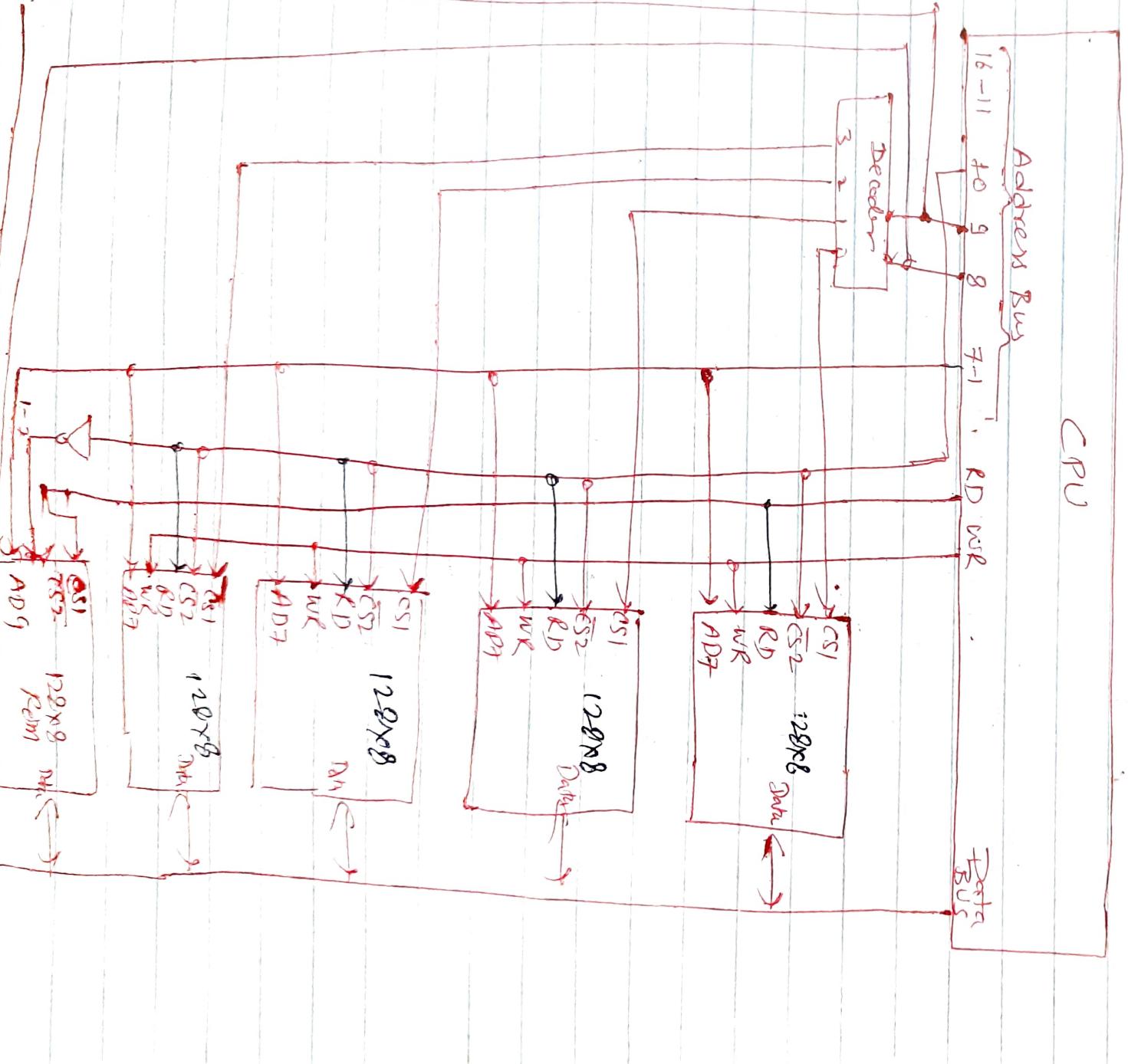
8-bit Data bus

Memory Address Map for Microcomputer

Component Hexadecimal Address Bus Address Bus

	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0
RAM1	0000 - 00FF	0 0 0 X X X	X X X X				
RAM2	0000 - 01FF	0 0 1 X X X	X X X X				
RAM3	0100 - 01FF	0 1 0 X X X	X X X X				
RAM4	0100 - 01FF	0 1 1 X X X	X X X X				
ROM	0200 - 03FF	1 X X X X X	X X X X				

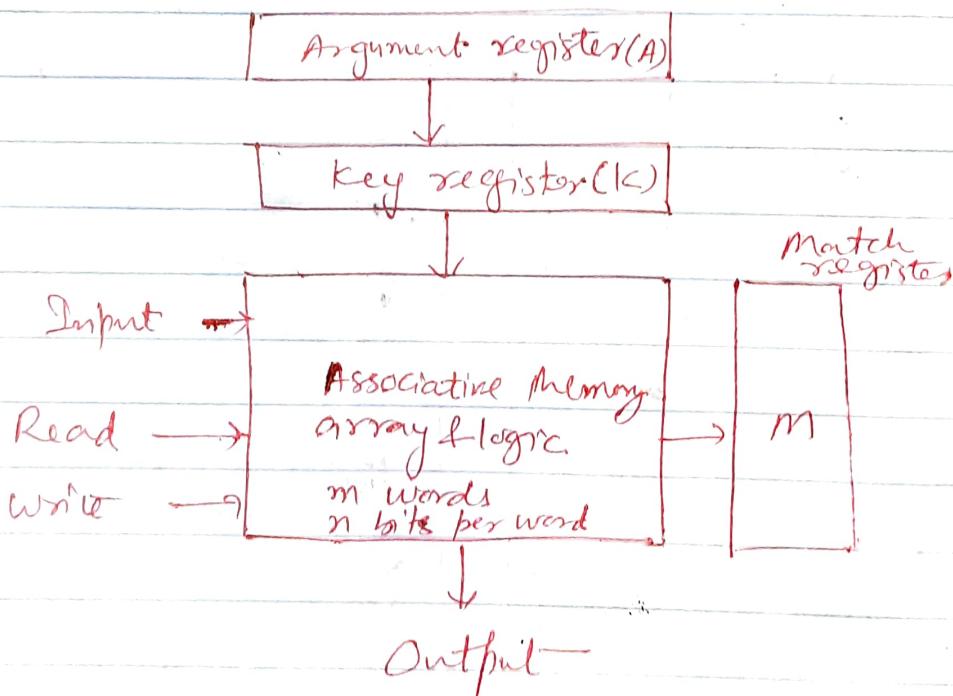
CPU



Lecture-30

Associative Memory :-

The time required to find an item stored in memory can be reduced considerably if stored data can be identified for access by the content of the data itself rather than by an address. A memory unit accessed by content is called an associative memory or content addressable memory (CAM).

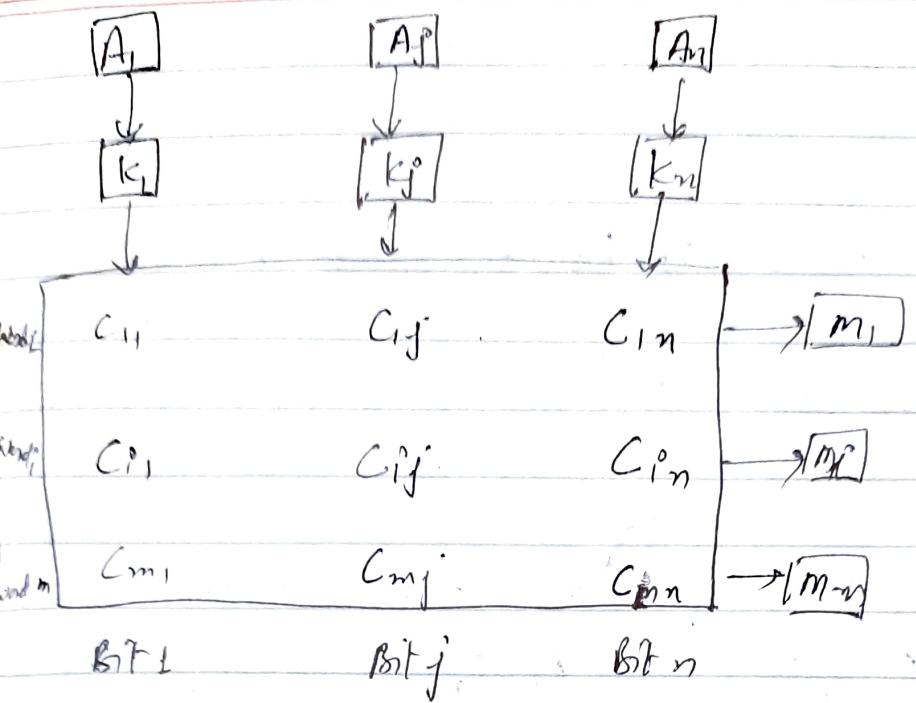


A 101 111100

K 111 000,000

Word 1 100 111100 no match

Word 2 101 000001 match

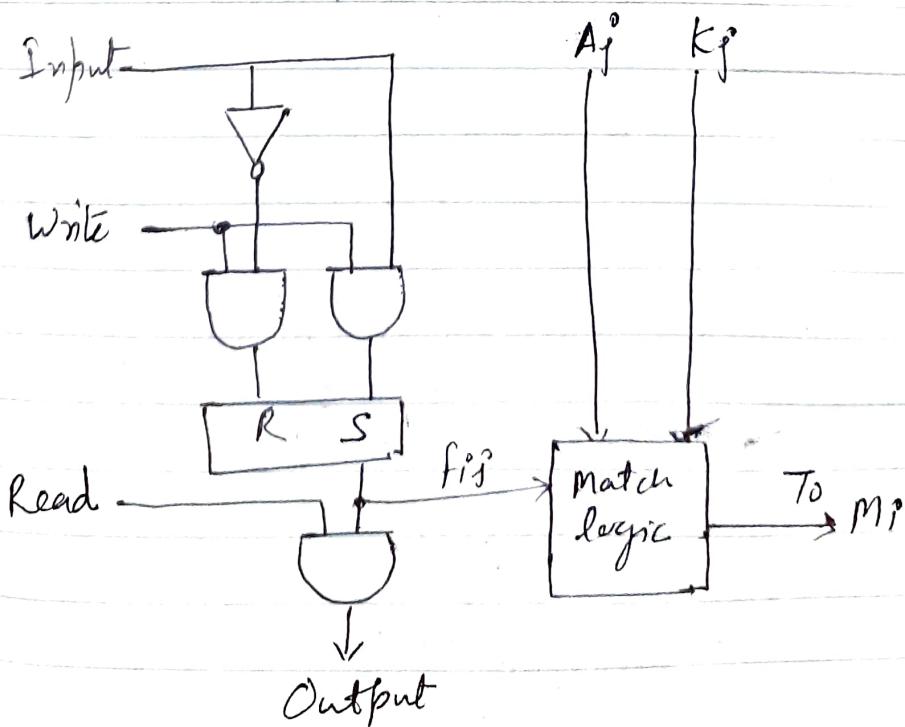


Match logic

$$x_f = A_j f_{ij} + A_j' f_{ij}'$$

$$M_i = x_1 x_2 x_3 \dots x_n$$

One cell of associative memory

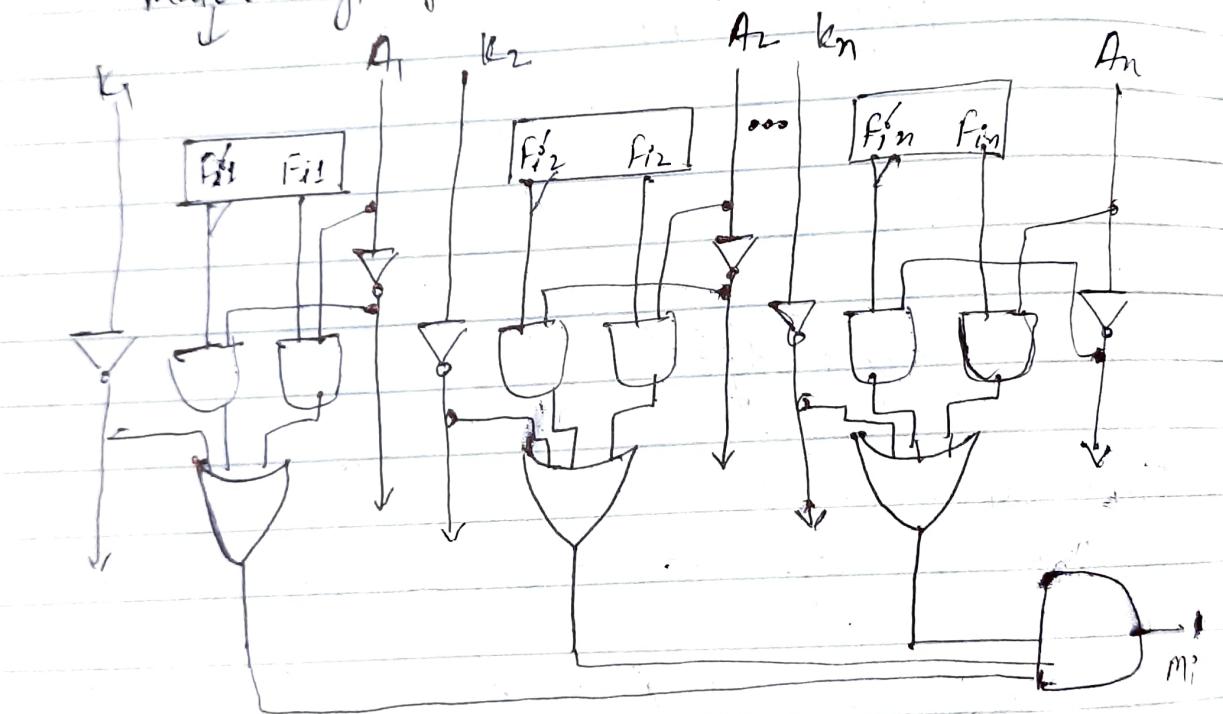


$$x_j + k'_j = \begin{cases} x_j & \text{if } k_j \neq 1 \\ 1 & \text{if } k_j = 0 \end{cases}$$

$$M_1 = (x_1 + k'_1)(x_2 + k'_2)(x_3 + k'_3) \dots (x_n + k'_n)$$

$$m_p = \prod_{j=1}^r (A_j f_{ij} + A_j' f_{ij}' + k'_j)$$

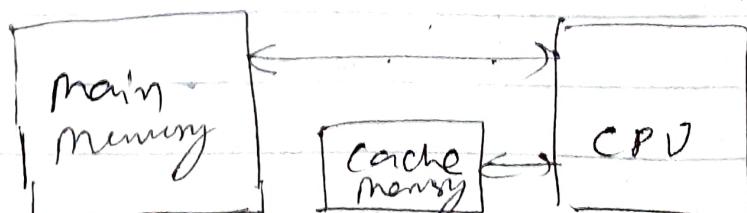
match logic for one word of associative memory



Cache Memory :-

locality of Reference :- Program access a relatively small portion of the address space at any instant of time.

Hit Ratios and Effective Access Times



Temporal Locality (Locality in Time) :- If an item is referenced, it will tend to be referenced again soon.

Spatial Locality (Locality in Space) :- If an item is referenced, items whose addresses are close by tend to be referenced soon.

$$\text{Hit ratio} = \frac{\text{No. of times referenced words are in cache}}{\text{Total number of memory accesses}}$$

$$\text{Eff. access time} = \frac{(\# \text{hits})(\text{Time per hit}) + (\#\text{misses})(\text{Time per miss})}{\text{Total number of memory access}}$$

Hit ratios and effective access time for multi-level cache →

$$H_1 = \frac{\text{No. times accessed word is in on-chip cache}}{\text{Total number of memory accesses}}$$

$$H_2 = \frac{\text{No. times accessed word is in off-chip cache}}{\text{No. times accessed word is not in on-chip cache}}$$

~~TOP~~

$$T_{\text{EFF}} = \frac{(\text{No. on-chip cache hits})(\text{on-chip cache hit time}) + (\text{No. off-chip cache hits})(\text{off-chip cache hit time}) + (\text{No. off-chip cache misses})(\text{off-chip cache miss time})}{\text{Total number of memory accesses}}$$

Cache Mapping →

① Associative mapping →

In case of associative mapping ~~the~~ Address & Data of main memory are in the Cache Memory.

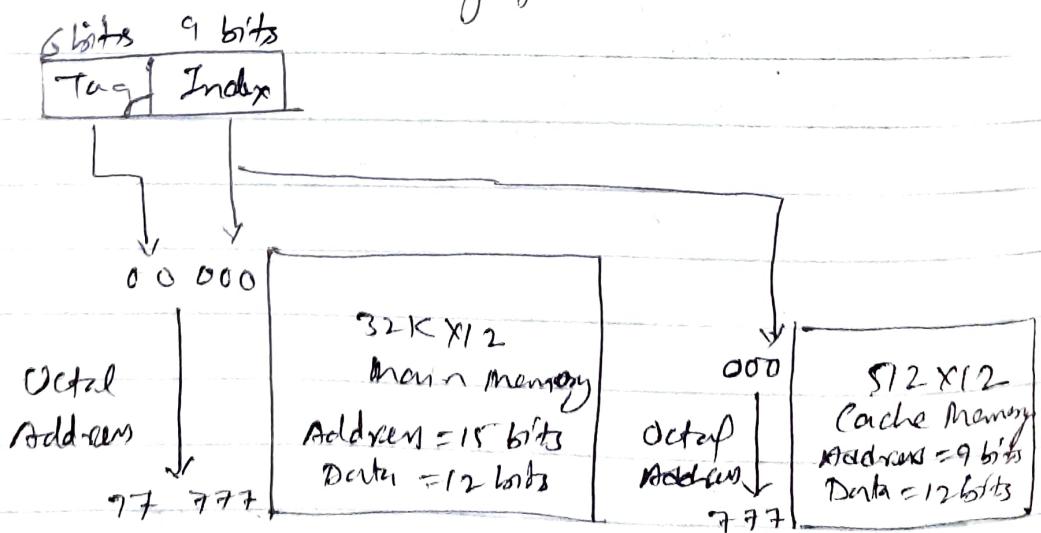
CPU Address (15 bits)

22
Argument register

Address	Data
01000	3450
02777	6710
22345	1234

Direct Mapping

In the general case, there are 2^k words in cache memory and 2^n words in main memory. The n -bit memory address is divided into two fields; k bits for the index field and $n-k$ bits for the tag field.



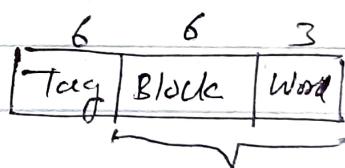
Memory address	Memory data
00000	1220
00777	2340
01000	3450
01777	4570
02000	5670
02777	6710

Index address	Tag	Data
000	00	1220
777	02	6710

(b) Cache memory

as main memory

Index	Tag	Data
000	01	3450
007	01	6570
010		
017		
770	02	
777	02	6710



Direct mapping. Cache with block size of 8 words

Lecture-31

Set-Associative Mapping

In general, a set associative cache of set size k will accommodate k words of main memory in each word of cache.

Index	Tag	Data	Tag	Data
000	01	3452	02	5678
777	02	6710	00	2340

Two-way set associative mapping cache.

Write into Cache

- ① write through— parallelly update cache & main memory
- ② write back → Only the cache locator is updated during write back operation. The locator is then marked by a flag so that later when the word is removed from the cache it is copied into main memory.

Cache Initialization?

Initially cache has some invalid data. We have a valid bit field for it. When the value of valid

~~bits~~

bit is 1 then its data is valid. Initially we make all the ^{valid} bits zero.

Block Replacement Policies :-

Random Replacement

FIFO (First In first Out)

LRU (Least Recently Used)

Random Replacement →

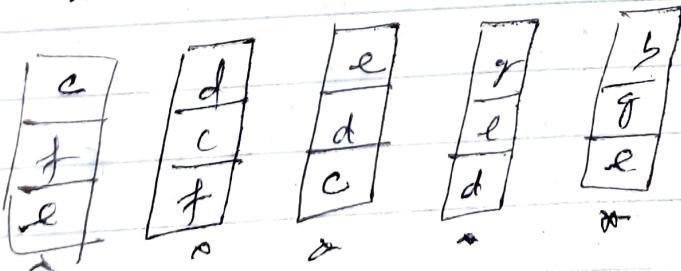
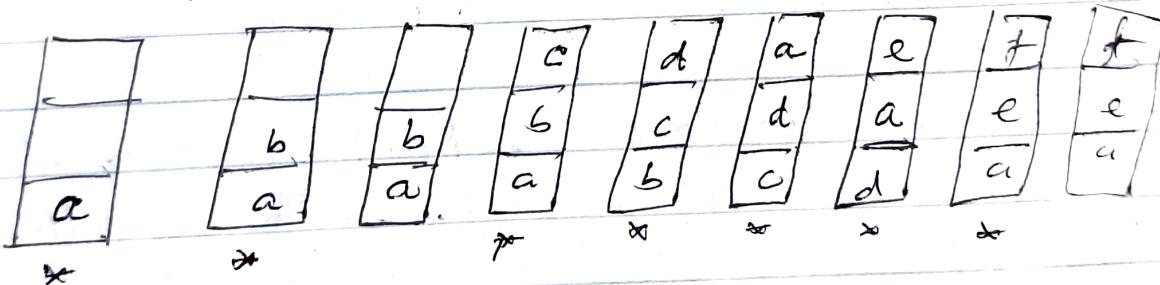
Randomly any block

can changed.

FIFO

a, b, a, c, d, a, e, f, a, c, d, e, g, b

Block size for 3.

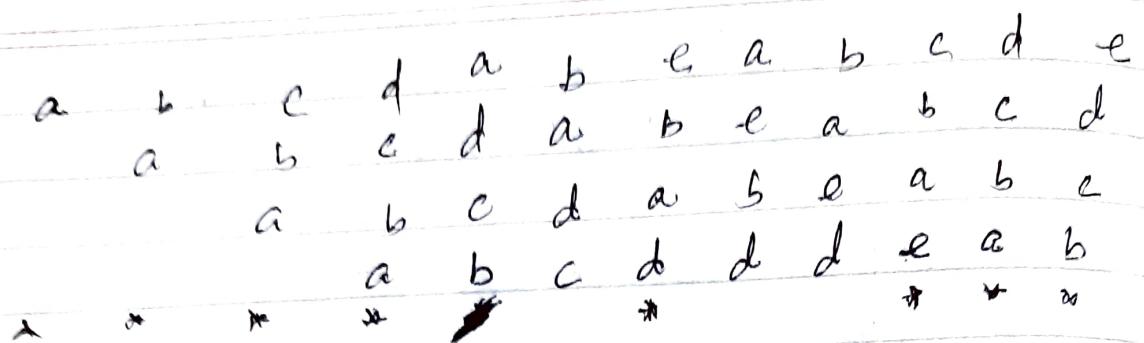


LRU :-

a b c d a b e a b c d e

Block size for 4.

Lecture-32



Q1. What do you mean by write-back method?

(2012-13) 2 marks

Q2. Explain the various types of mapping procedures used by cache memory

(2012-13) -10 marks

Q3. Explain the need of memory hierarchy.

What is the main reason for not having a large main memory for storing the totality of information in a computer system?

(2012-13) 10 marks

Q4. A computer employs RAM chips of 256×8 and ROM chips of 1024×8 . The computer system needs 2k bytes of RAM, 4k bytes of ROM and 4 interface unit each with 4 registers. A memory mapped configuration is used. The 2-highest order bits of address bus are assigned '00' for RAM '0' for ROM and '10' for interface registers.

- How many RAM & ROM chips are needed?
- Draw a memory address map for the system.
- Give the address range in hexadecimal for

t_{eff} = effective access time.

$$t_{eff} = t_{cache} + (1-h)t_{main}$$

e.g. for a hit ratio of 90% (0.9), a cache access time of 10 ns, and a main memory access time of 60 ns, we have

$$t_{eff} = 10 + (1 - 0.9)60 = 16 \text{ ns.}$$

RAM, ROM and Interface.

Write-through \rightarrow

Parallel updated Cache & Main

memory.

Write-back:-

After removing from the cache it may be updated in main memory.

Virtual Memory :-

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

CPU generates virtual address & Main memory has physical address

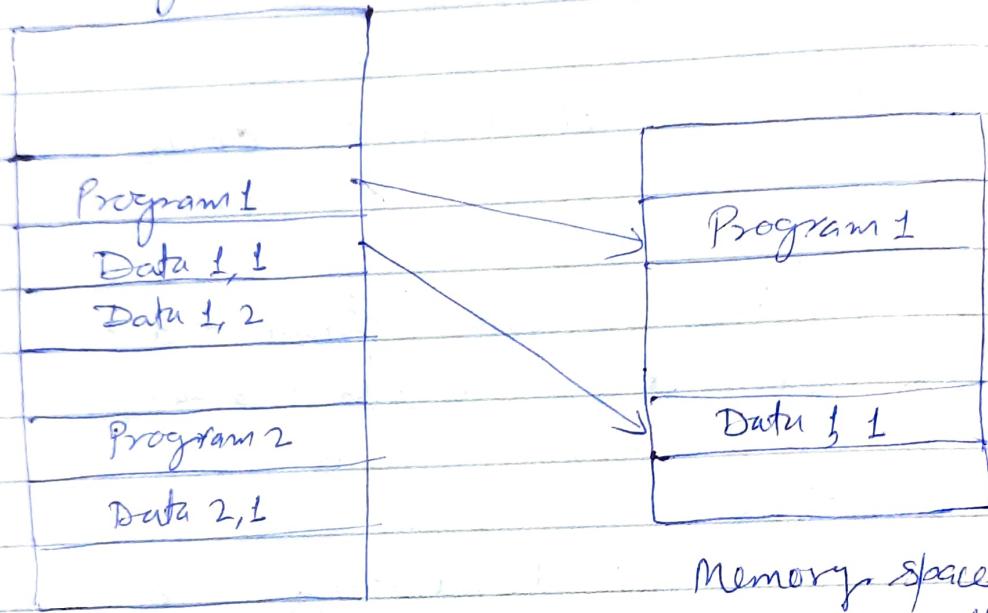
Address Space

* memory space

e.g. \rightarrow main memory capacity of 32 k words ($k = 1024$). Fifteen bits are needed to specify a physical address in memory since $32k = 2^{15}$. Suppose that the computer has available auxiliary memory for storing $2^{20} = 1024$ words. Thus auxiliary memory has a capacity for storing information equivalent to the capacity of 32 main memories. Denoting

the address space by N and the memory space by M , we then have for this example $N=1024$ and $M=32K$.

Auxiliary memory



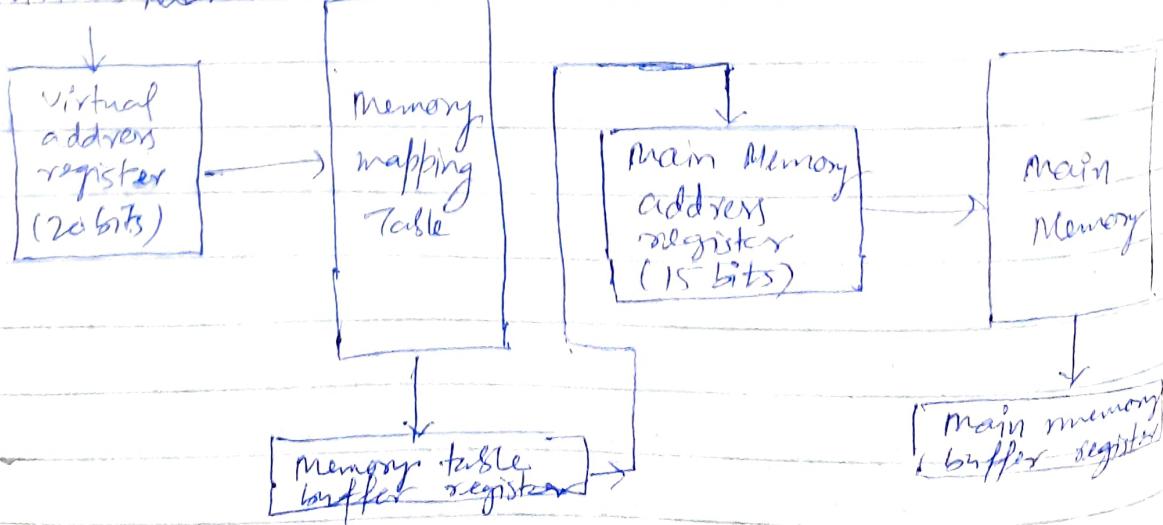
Address space

$$N = 1024k = 2^{20}$$

$$\text{Memory space} \\ M = 32k = 2^{15}$$

Memory table for mapping a virtual address

Virtual Address



For mapping

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 programs running at half speed.

A third alternative is to use an associative memory. (We will explain)

Address Mapping Using Pages →

Physical memory is broken down into groups of equal size called blocks which may range from 64 to 4096 words each. & ^(page, page) virtual memory & (programs) are divided into equal size pages.

e.g. Consider a computer with an address space of 8K and a memory space of 4K. If we split each into groups of 1K words we obtain eight pages and four blocks.

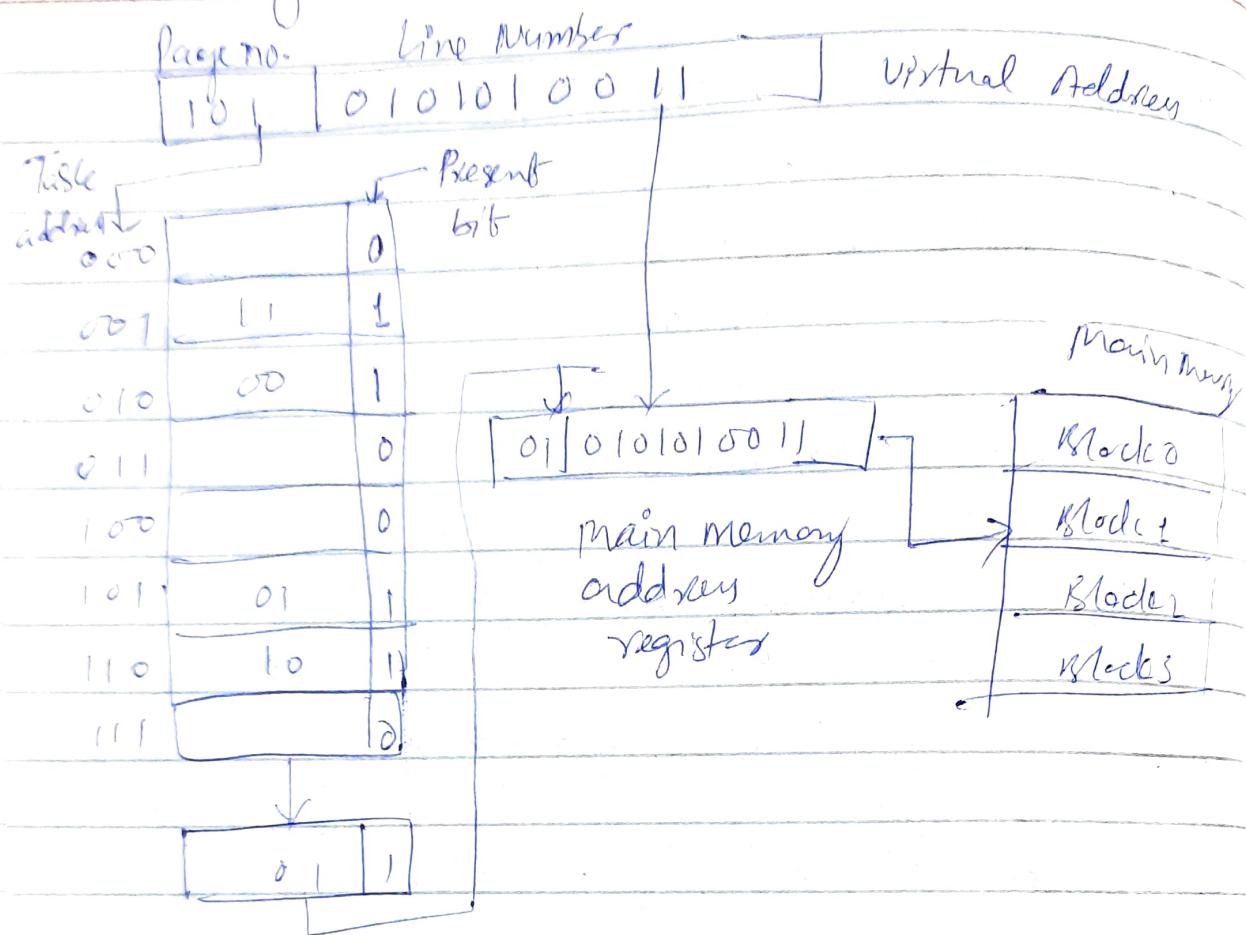
Page 0
Page 1
Page 2
Page 3
Page 4
Page 5
Page 6
Page 7

Address space
 $N = 2^8 K = 2^{13}$

Block 0
Block 1
Block 2
Block 3

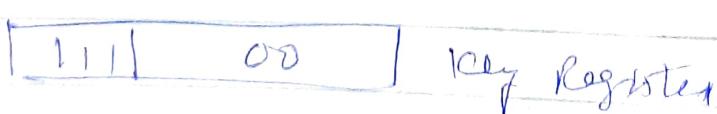
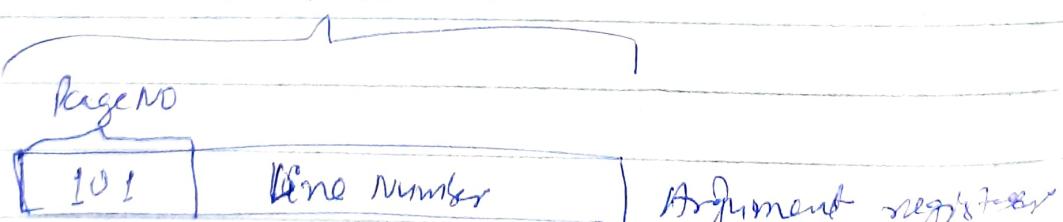
Memory space
 $m = 4 K = 2^{12}$

Memory tables in page system

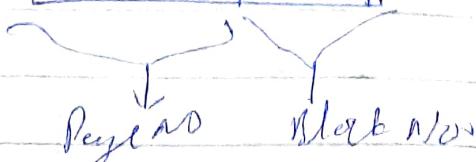


Associative Memory Page Table

virtual Address



001	11	
010	00	
101	01	Associative memory
110	10	



Lecture-33

Page Replacement

FIFO, LRU, Optimal

Memory Management Hardware

MMU is combination of software & hardware.

→ The basic components of a memory management unit are:

1. A facility for dynamic storage relocation that maps logical memory references into physical memory address.
2. A provision for sharing common programs stored in memory by different users.
3. Protection of information against unauthorized access between users and preventing users from changing operating system functions.

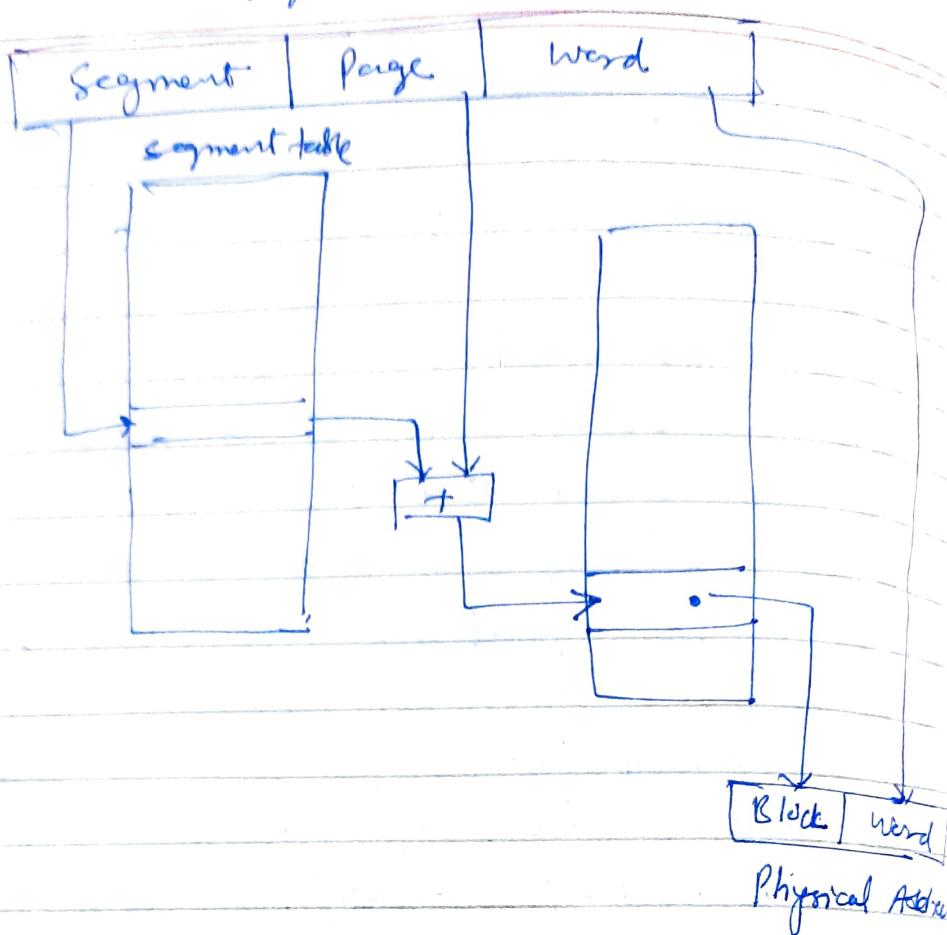
Segment →

To divide programs and data into logical parts called segments. A segment is a set of logically related instructions or data elements associated with a given name.

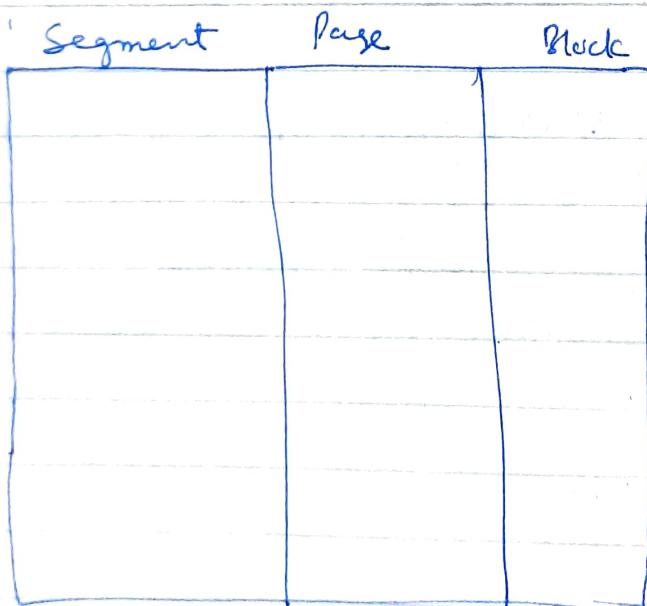
Logical Address

Segmented page Mapping →

Logical Address



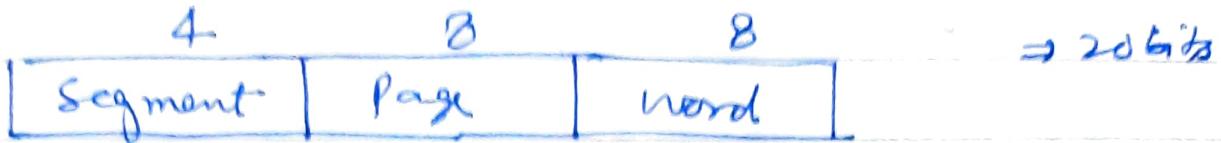
Argument registers



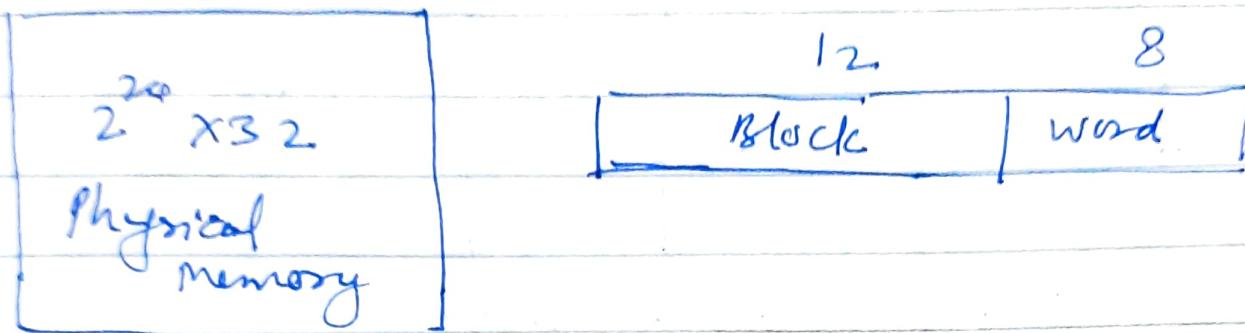
(b) Associative memory translation of off-aside buffer (TLB)

Lecture - 34

Numerical Example :-



- Q) Logical address format : 16 segments of 256 pages each, each page has 256 words



- ① Physical address format : 4096 blocks ~~out of~~ of ~~of~~ 256 words each, each word has 32 bits.

Hexadecimal Address	Page Numbers
60000	Page 0
60100	Page 1
60200	Page 2
60300	Page 3
60400	Page 4
604FF	Page 5

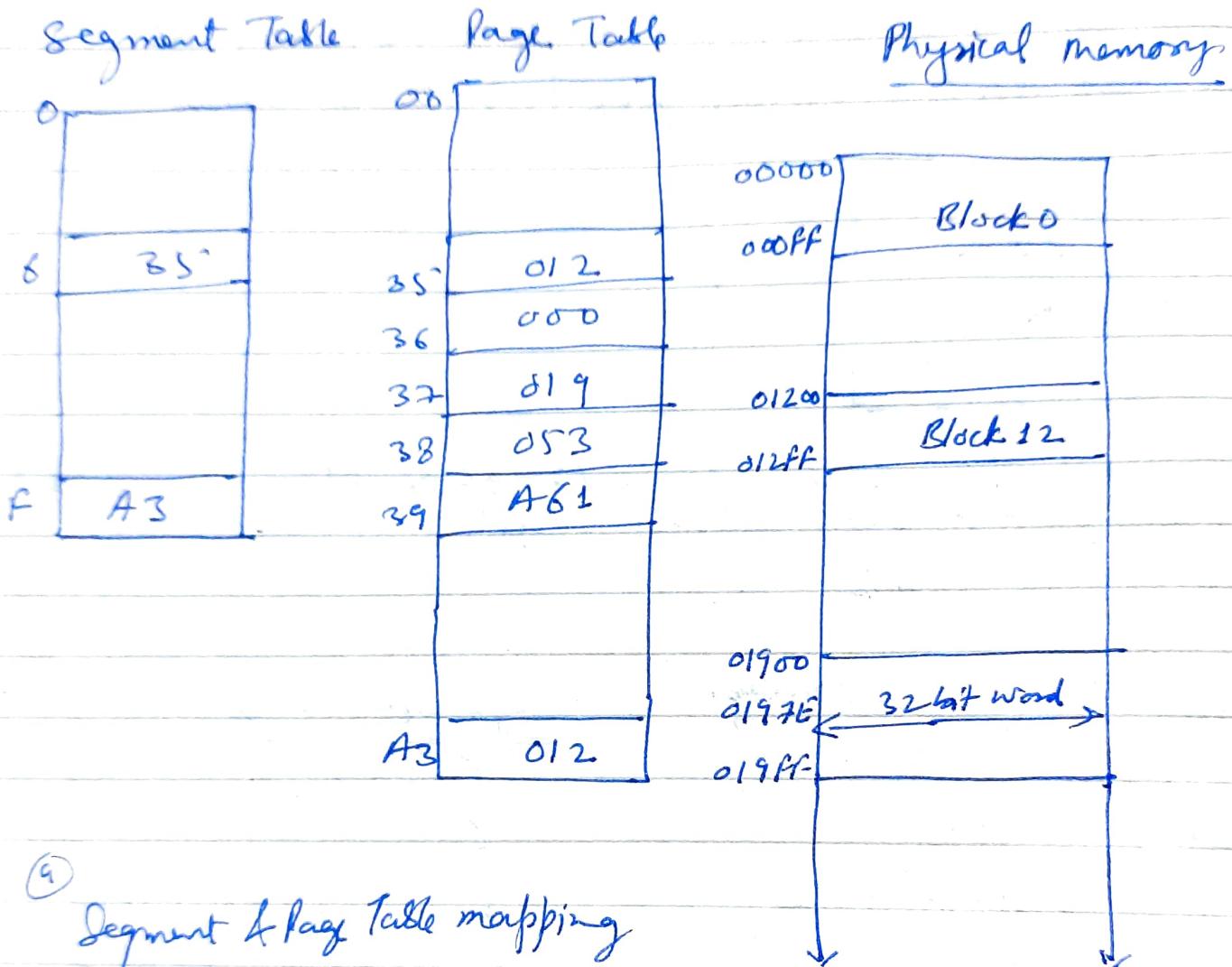
Segment	Page	Block
6	00	012
6	01	000
6	02	019
6	03	053
6	04	A61

- ② Logical address Assignment

(i) Segment-page versus memory block assignment

Logical Address (in hexadecimal)

6	02	7E
---	----	----



(a) Segment & Page Table mapping

Segment	Page	Blocks
6	02	019
6	04	A61

(b) Associative memory (TLB)

- Q1. Q) How many 128×8 RAM chips are needed to provide a memory capacity of 2048 bytes?
B) How many lines of the address bus must be used to access 2048 bytes of memory? ~~How many of these lines will be common to all chips?~~
C) How many lines must be decoded for chip select? Specify the size of the decoder.