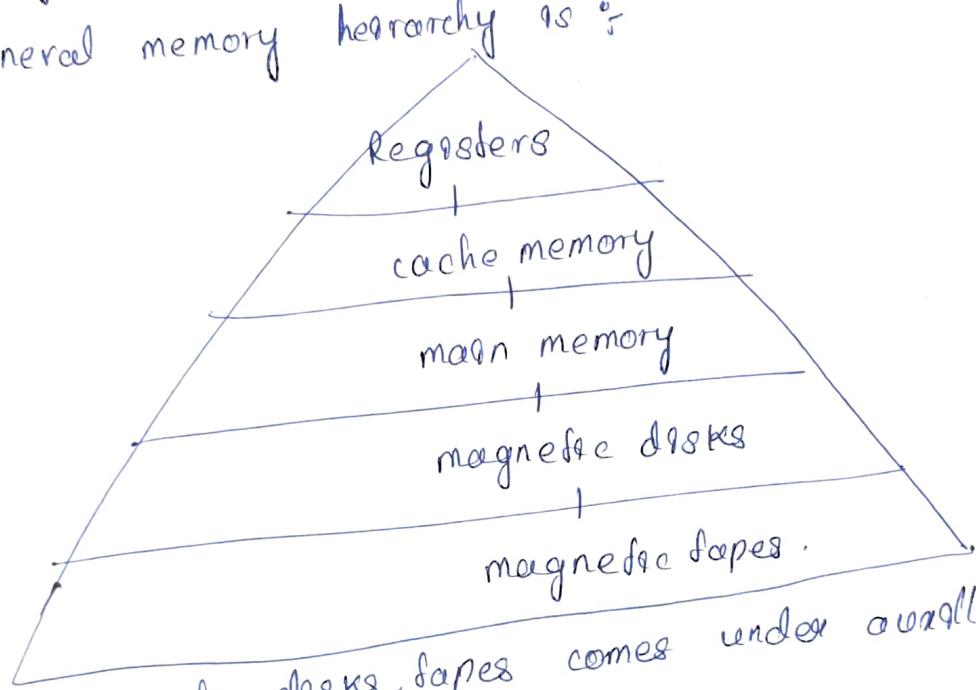


- ④ memory hierarchy of memory devices are used to store the data, instructions, information.
- ⑤ memory devices can perform only read and write operations.
- ⑥ different types of memory devices are of main (or) primary memory, secondary (or) storage devices.
- ⑦ general memory hierarchy is :-



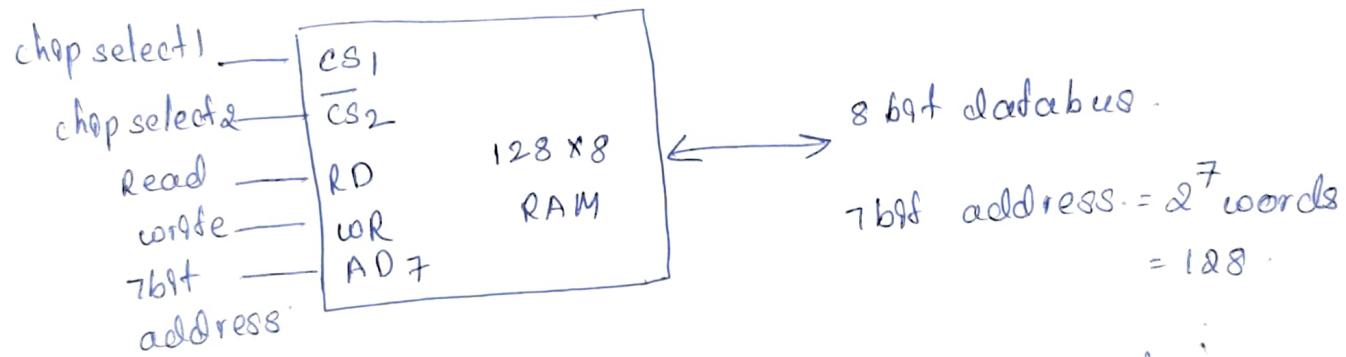
- ⑧ where magnetic disks, tapes comes under availability (or) secondary storage devices.
- ⑨ In the above hierarchy, speed is decreased when we move from top to bottom.
- ⑩ storage capacity is increased from top to bottom.
- ⑪ cost is decreased from top to bottom.
- memory access methods of following memory locations -
- information from memory locations -
- in this each memory location has unique address, using this address any memory location can be reached in some amount of time in any order.
- are methods to access

- b) sequential access of memory is accessed sequentially.
- c) direct access of in this method, information is stored on tracks, with each track having separate read/write head or.

## (ii) main memory of it communicates directly with CPU.

- it is central storage unit of computer.
- it is large and fast memory used to store data during computer operations.
- main memory is made up of RAM and ROM.

## RAM (Random access memory)



- it is random access memory i.e. directly searches for element.
- it is volatile memory, used for both read & write operations.

○ types are - static, dynamic.  
static RAM → data is stored on constant power flow, because of doesn't need to be refreshed stored.

transistors and requires a continuous power, SRAM to remember data being

- binary data present here will be stored until power is off.
- the word storage and decodes that memory contains data

- as long as power is supplied.
- ① data is lost when power goes down due to volatile nature (2)  
② SRAM chip uses matrix of 6 transistors and no capacitors  
③ transistors do not require power to prevent leakage, so SRAM not be refreshed on regular basis.
- ④ SRAM uses more chips than DRAM for some amount of storage space, making manufacturing cost high.
- ⑤ characteristics of static RAM →  
• long life • no need to refresh • faster.  
• used as cache. Large size. expensive.  
advantages → low power consumption, faster access.  
disadvantages → less memory capacity, high cost.
- DRAM (Dynamic Random access memory) is used to store data. It is continuously refreshed in order to maintain the data. It is used in many systems as it is cheap & small.
- ⑥ it is used in DRAM →  
⑦ characteristics of DRAM →  
• short data lifetime.  
• needs to be refreshed continuously.  
• slower compared to SRAM.  
• used as RAM, smaller size, less expensive.  
advantage → low cost, more memory capacity.  
disadvantage → slow access, high power consumption.
- SRAM
- ⑧ Transistors are used to store information in SRAM.  
⑨ capacitors are not used hence no refreshing is required.
- ⑩ faster compared to DRAM  
⑪ expensive  
⑫ low density devices  
⑬ used on cache
- DRAM
- ⑭ capacitors are used to store data in DRAM.  
⑮ to store information for a longer time, contends of capacitor needs to be refreshed periodically.  
⑯ slow  
⑰ cheap  
⑱ high density.  
⑲ used on main memory.

ROM (read only memory) is used only to read not to write.

- it is non volatile memory.
  - it generally stores such instructions that are required to start a computer.
  - ROM chips are also used in electronic items like washing machines, ovens etc.
  - data is recorded at the time of manufacturing.
- Types of ROM are PROM, EPROM, EEPROM.

PROM (programmable read only memory) is it can be modified only once by user.

- user buys a blank PROM and enters required data using PROM program.
  - it can be programmed only once and is not erasable.
- EPROM (erasable and programmable Read only memory) is it can be erased by using UV (ultra violet) light for duration upto 40 minutes.

we can erase and write new content by multiple times.

- written instructions can be erased and changed with help of UV light.
- EEPROM (electrically erasable and programmable Read only memory) is

data can be erased electrically.

- it can be erased and reprogrammed about ten thousand times.

Advantages of ROM are non volatile in nature, cheaper than RAM, easy to read, reliable than RAM, static & no need of refreshment.

Auxiliary memory: most common auxiliary memory devices used in computer system are magnetic disks and tapes. ③

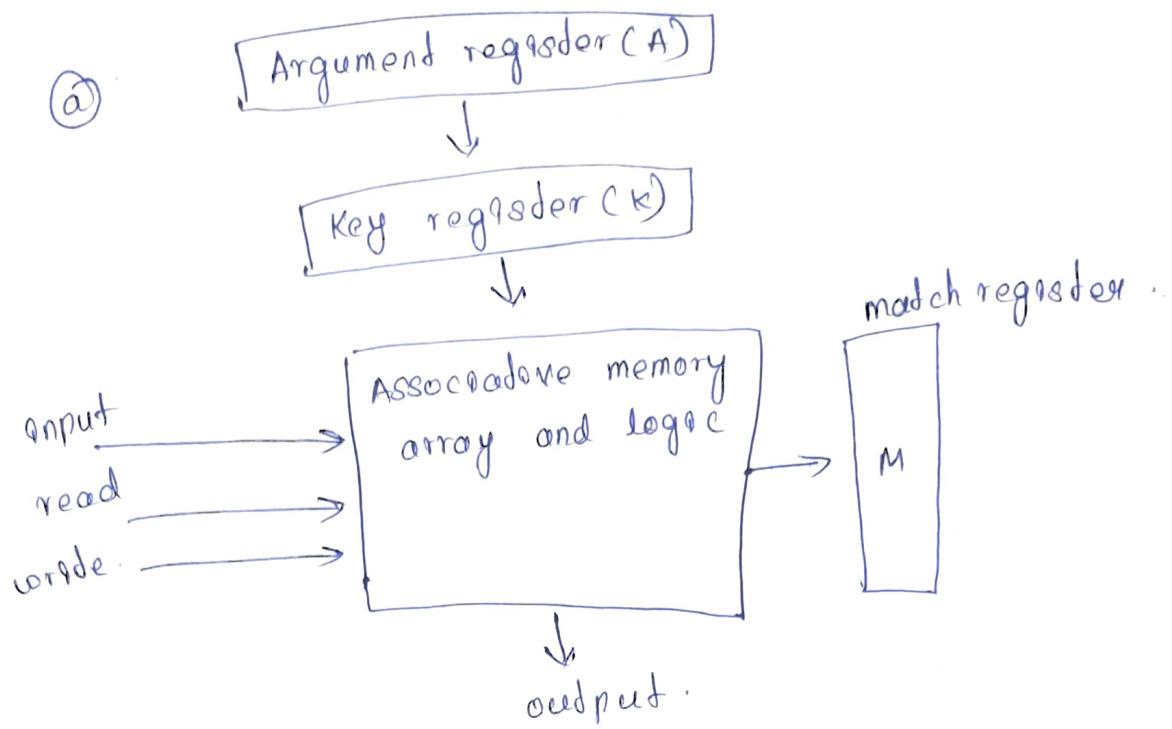
- ⓐ magnetic disks & explain about hard disks.
- ⓑ magnetic tape & data is read/write sequentially.
- ⓒ only one side of ribbon is used for storing data.
- ⓓ it is a sequential memory which contains thin plastic ribbon to store data and coated by magnetic oxide.
- ⓔ data read/write slower because of sequential access.
- ⓕ it is highly reliable which requires magnetic tape drive writing and reading data.

advantages: low cost, provides backup or archival storage, reusable memory.

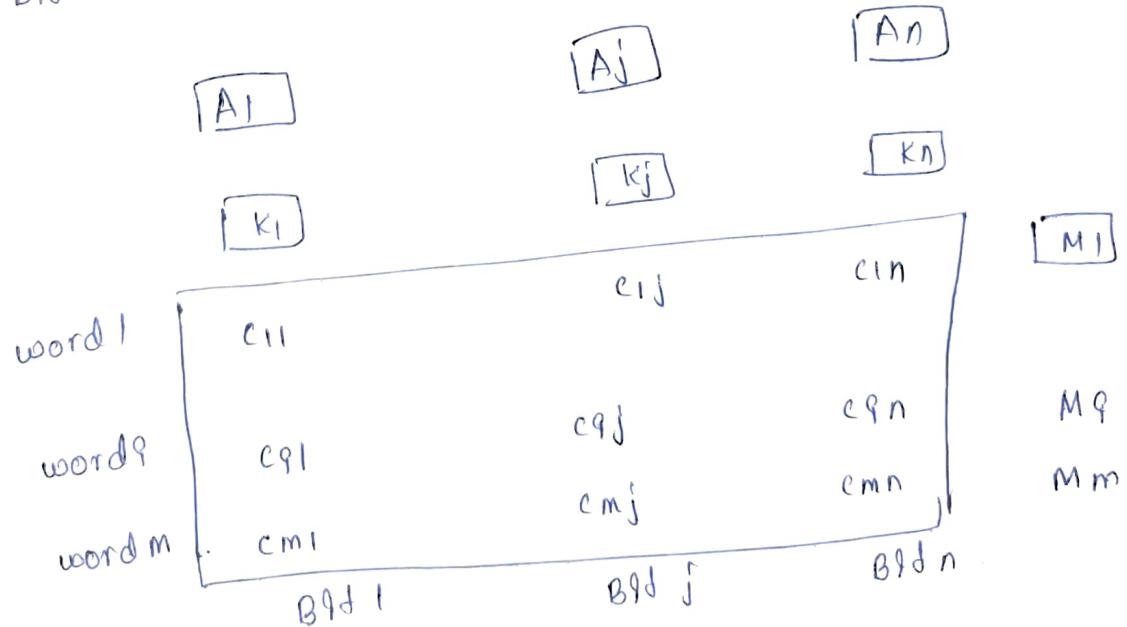
disadvantages: used for large files, sequential access, data stored cannot be easily updated (or) modified. i.e. difficult to make updates on data.

Associative memory: it can be considered as memory unit whose data stored can be selected for access by content of data rather than address.

- ⓐ it is referred as content Addressable memory (CAM).
- ⓑ when write operation is performed, no address or memory location is given to word.
- ⓒ the word itself is capable of sending empty location to store word.
- ⓓ In read operation, content of word is specified.
- ⓔ the word matching with specified content are selected for reading.

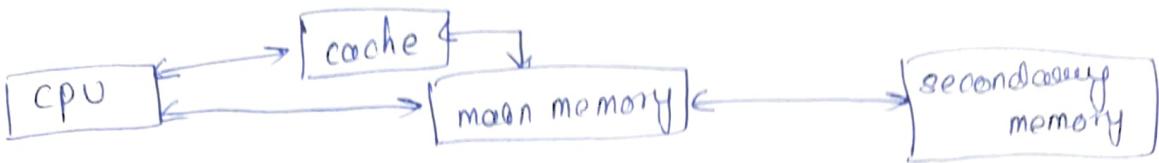


- ① Argument register (A) contains part of word.
  - ② Key register contains information, if  $bfd = 1$ , will be compared Argument value and word value in memory, if  $bfd = 0$ , comparison neglected.



## Cache memory

(4)



- ① Cache is the fastest memory device.
- ② Whenever processor requires any data, first it will search for data in cache, if not found goes to main memory, if found, a copy will be sent to CPU and a copy to cache.
- ③ If data is not in main, then system searches for data in secondary, if found copy will be sent to main memory.
- ④ It is very costly so we take less storage capacity.
- ⑤ The data (or) contents of main memory that are used again and again by CPU are stored in cache, so that it can easily accessed in short time.

Hit ratio of the performance of cache memory is calculated in terms of quantity called Hit ratio.

Hit ratio  $\Rightarrow$  whenever CPU refers to memory and if it finds that data in cache is present then it is called hit.

- ⑥ If data not found, it is known as miss.
- ⑦ If data not found, it is known as miss.

Hit ratio =  $\frac{\text{hit}}{\text{hit} + \text{miss}}$ .

- ⑧ Cache performance depends on hit ratio, we can improve cache blocks size, cache performance using higher associativity, reduce miss ratio, and reduce time to hit on cache.

Cache mapping of three different types of mappings used.

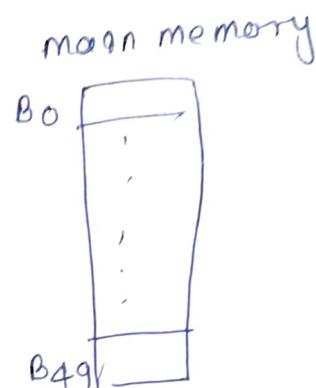
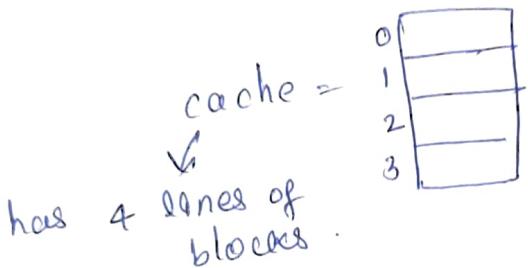
1. Direct mapping
2. Associative mapping
3. Set associative mapping.

① Direct mapping is as simple, less expensive to implement.

② certain blocks of main memory would be able to map cache only up to a certain line of cache.

$$\text{cache line no} = \frac{(\text{Address of main memory block}) \text{ mod} (\text{total no. of lines in cache})}{}$$

e.g. Let



③ If  $B_0$  has to be placed in cache, then

$$\text{formula} =$$

$$\begin{aligned} \text{block no} &= B_0 \text{ mod } 4 \\ &= 0 \text{ mod } 4 = 0 \end{aligned}$$

$$(C = K \text{ mod } n)$$

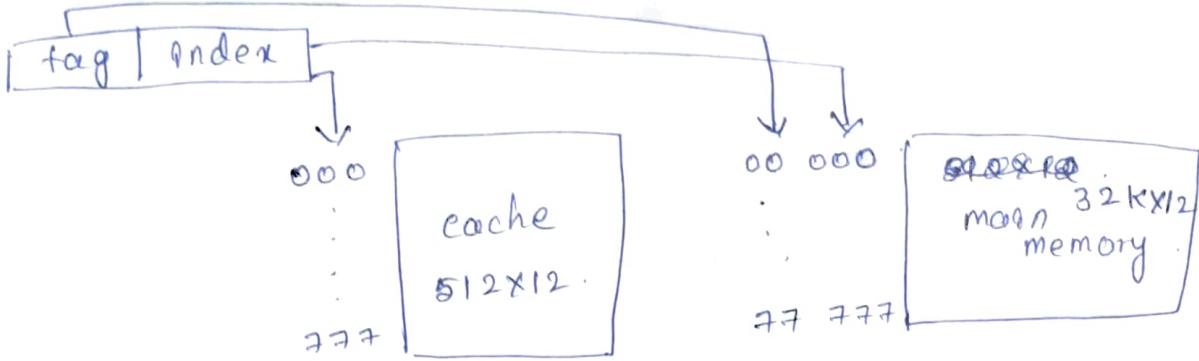
$\therefore B_0$  has to be placed at 0.

If  $B_3$  has to be placed, then

$$\begin{aligned} \text{block no} &= B_3 \text{ mod } 4 \\ &= 3 \text{ mod } 4 \end{aligned}$$

$B_3$  has to be placed only at 3rd block of cache.

- Q address is divided into 2 parts i.e. tag and index. (5)



Associative mapping of the main memory block is capable of mapping to any given line of cache, which is free at that particular moment.

- particular moment.

  - ① it helps us make a fully associative mapping comparatively more flexible than direct mapping.
  - ② searching for content done with content only.
  - take associative memory.

Set associative mapping is of the nature of both direct and associative.

- both direct and associative  
cache is divided into number of sets.

Eg: 2 way sed association - - -

- ① if any memory block has to be mapped with cache, then sets & selected as direct mapping i.e.  

$$= (\text{Block of memory}) \bmod (\text{Set numbers})$$

if B<sub>4</sub> has to be placed then  
 sed no = 4 \* 1.2 = ①

- ① Block B<sub>4</sub> has to be placed at S<sub>0</sub>, but on S<sub>0</sub>, it can be placed anywhere (on any block), like associative.

## Cache Memory

### Types of cache memory:

L<sub>1</sub>: is first level of cache, called as Level 1 cache (or) L<sub>1</sub> cache; here small amount of memory is present inside the CPU itself.

- ① As the memory is present in CPU, it can work at same speed as CPU.

② The size of memory ranges from 2KB to 64KB.

③ The size of memory ranges from 256KB to 1MB.

L<sub>2</sub>: This level of cache may be inside or outside CPU.

④ In case of outside of CPU, it is connected with CPU with very high speed bus.

⑤ All the cores of CPU can have their own separate level L<sub>2</sub> cache among themselves.

⑥ They can share one L<sub>2</sub> cache among them. The size of this cache is in range of 256KB to 1MB.

⑦ The memory size of this cache is 512KB.

⑧ These are slower than L<sub>1</sub> cache.

L<sub>3</sub>: This is not present in all processors.

⑨ This is used to enhance performance of Level and Level 2 cache.

⑩ It is located outside of CPU and is shared by all cores of CPU.

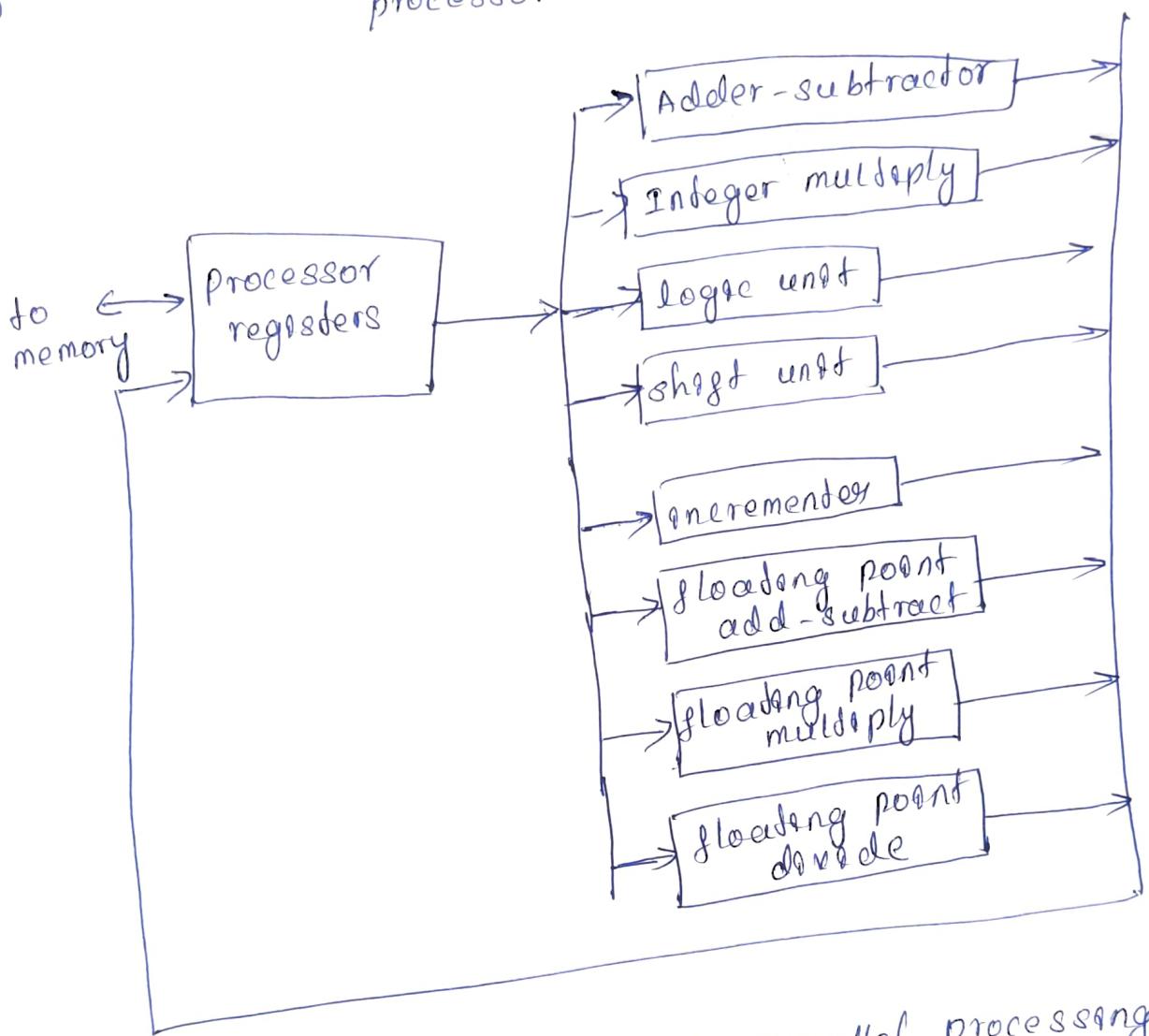
1MB to 8MB.

⑪ Memory size ranges from L<sub>1</sub> and L<sub>2</sub> cache, but faster than RAM.

⑫ Slower than L<sub>1</sub> and L<sub>2</sub> cache.

Parallel processing is a term used to denote large class of technique used to provide simultaneous data processing tasks to increase speed of computer.

- ① It can be achieved by using multiple functional units.
- ② processor with multiple functional units.



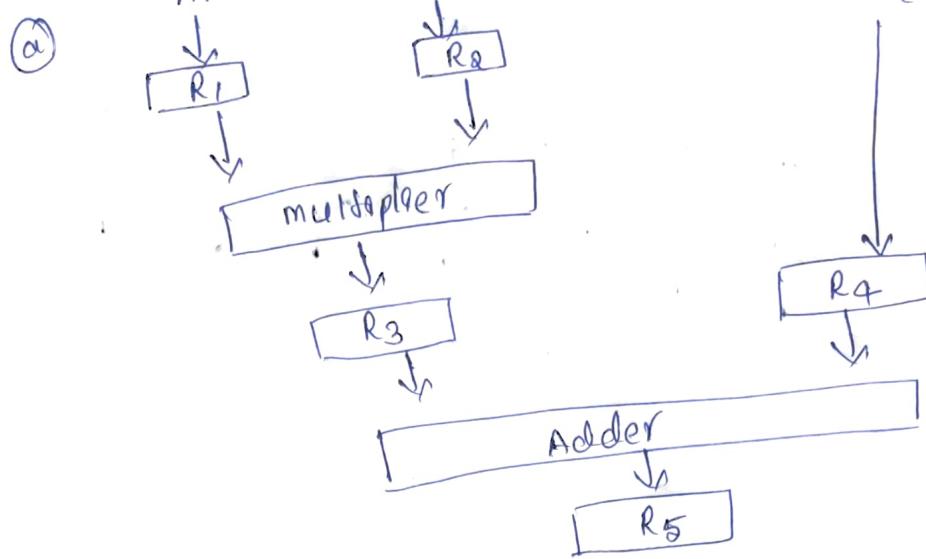
- ① There are different ways that parallel processing can be classified.
- ② One classification introduced by M. J. Flynn.
- ③ Flynn's classification divides computers into 4 major groups as follows:
- |  |                                     |
|--|-------------------------------------|
| single instruction stream, single data stream (SISD) | stream, multiple data stream (SIMD) |
| single instruction " " , single " (MISD)             | " , multiple " " (MIMD)             |
| multiple " "   | " "                                 |
| multiple " "   | " "                                 |

- ① instruction stream  $\rightarrow$  sequence of instructions read from memory.
- ② data stream  $\rightarrow$  operation performed on data in processor.
- ③ SIMD  $\rightarrow$  represents a system with control unit, processor and memory. Instructions are executed sequentially.
- ④ here parallel processing can be achieved with multiple functional units (or) by pipeline processing.
- ⑤ SIMD  $\rightarrow$  contains many processing units under common control unit; at a time only one instruction taken but different operations on different data.
- ⑥ MIMD  $\rightarrow$  only theoretical interest no practical system has concerned with this concept.
- ⑦ MIMD  $\rightarrow$  capable systems capable of processing many programs at a time. Most multiprocessor and multi-computer systems comes under this concept.
- ⑧ pipelining  $\rightarrow$  parallel processing can be done with help of pipelining also.
- ⑨ it is a technique of dividing sequential process into suboperations.
- ⑩ each subprocess executed in a special dedicated segment that operates concurrently with all other segments.
- ⑪ it is a collection of segments, where each segment performs fixed tasks.

⑫ e.g.:  $A^q * B^q + C^q$   
is performed with operations:

for  $q = 1, 2, 3, \dots, 7$ .

$$\begin{aligned}
 R_1 &\leftarrow A^q, R_2 \leftarrow B^q \\
 R_3 &\leftarrow R_1 * R_2, R_4 \leftarrow C^q \\
 R_5 &\leftarrow R_3 + R_4
 \end{aligned}$$



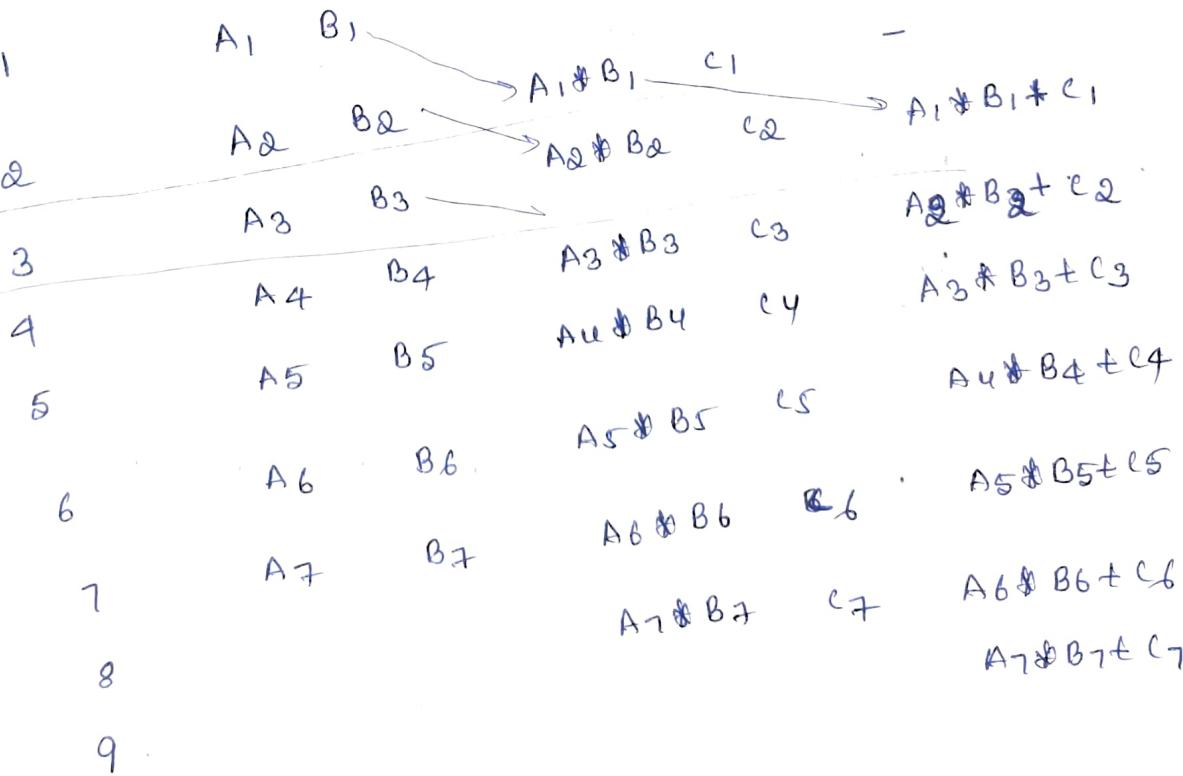
(b) content of Registers in pipeline eg:

clock  
pulse  
number

segment 1  
R<sub>1</sub>    R<sub>2</sub>

segment 2  
R<sub>3</sub>    R<sub>4</sub>

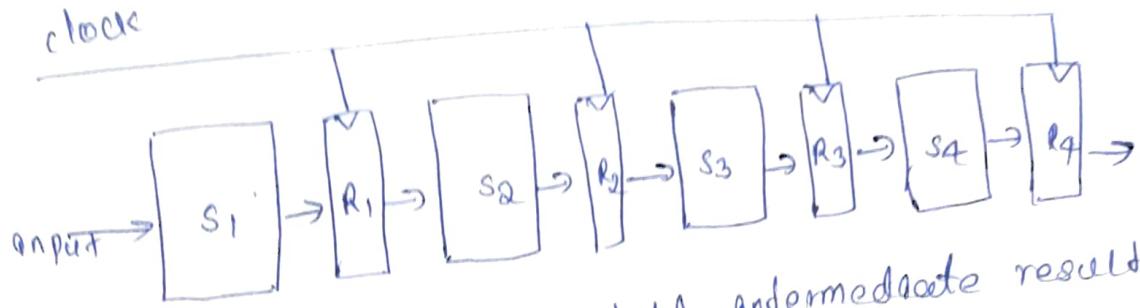
segment 3  
R<sub>5</sub>



9 clock cycles.

## 4 segment pipeline

①



$S = \text{segments}$ ,  $R = \text{registers}$  so hold intermediate result.  
Space-time diagram of pipeline can be seen with space-time

② The behavior of pipeline

③ diagram

④ space space diagram of a 4-segment pipeline as (horizontal axis goes vertical axis segment number.)

$T = \text{task}$ .

space space diagram for pipeline.

⑤

	1	2	3	4	5	6	7	8	9	clock cycles
segment: $S_1$	$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$				
$S_2$		$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$			
$S_3$			$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$		
$S_4$				$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	

⑥ 6 different tasks  $T_1$  to  $T_6$ .

⑦ segments  $\Rightarrow S_1, S_2, S_3, S_4$ .

⑧ clock cycles  $\Rightarrow 9$ .

- there are 2 areas of computer design where pipeline is applicable.

Arithmetic pipeline,  
Instruction pipeline.

- Arithmetic pipeline is usually found on very high speed computers.
- They are used to implement floating point operations, multiplication of fixed point numbers.

○

Eg°:

$$x = A \times 2^a$$

$$y = B \times 2^b$$

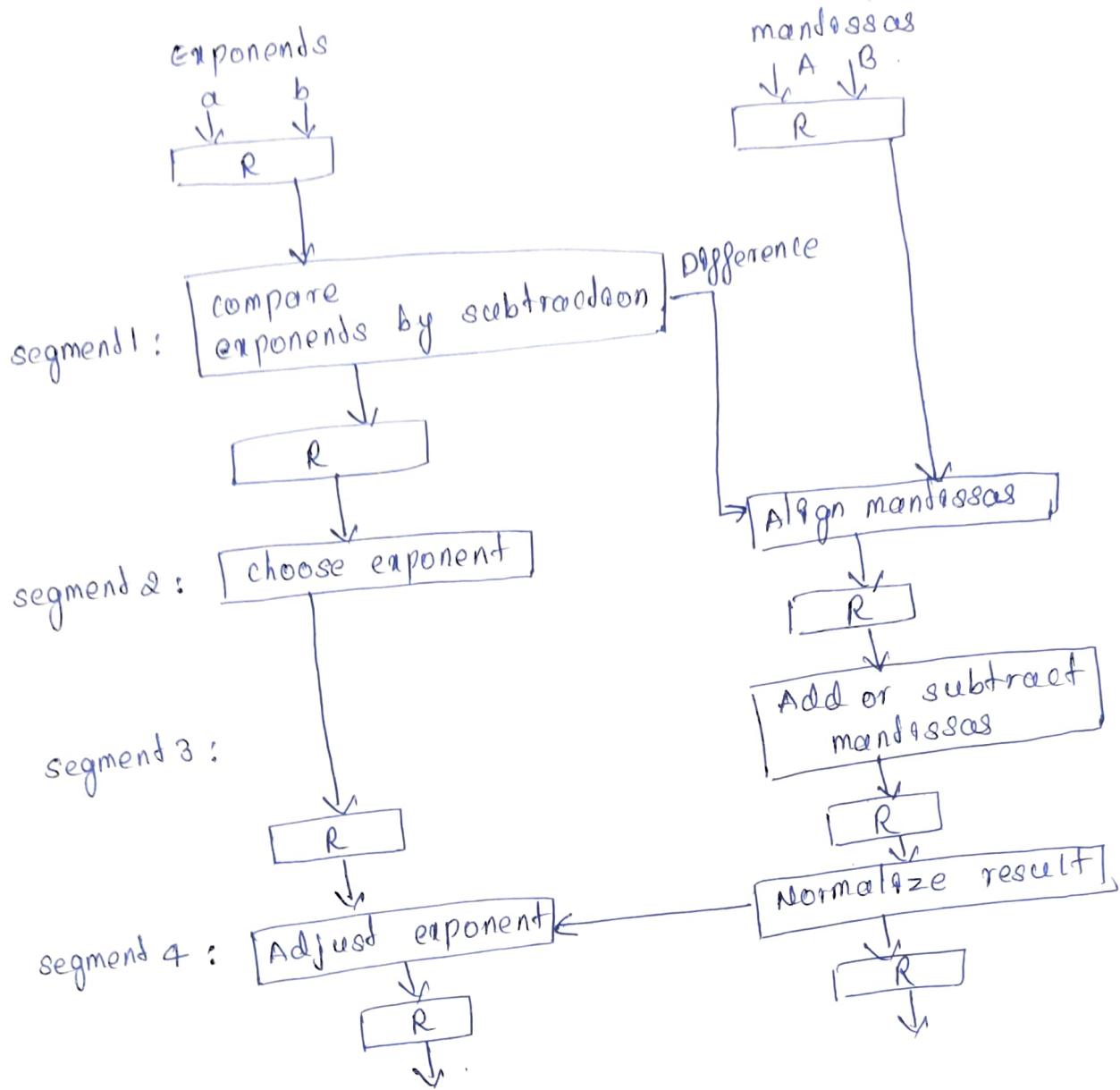
$A, B$  = mantissas.  
 $a, b$  = exponents.

- floating point addition and subtraction can be performed on 4 segments.

1. compare exponents
2. Align mantissas
3. Add or subtract mantissas
4. normalize result.

Ex:

a) pipeline for floating point addition and subtraction.



instruction pipeline or pipeline processing can occur not only in data stream but on the instruction stream as well.

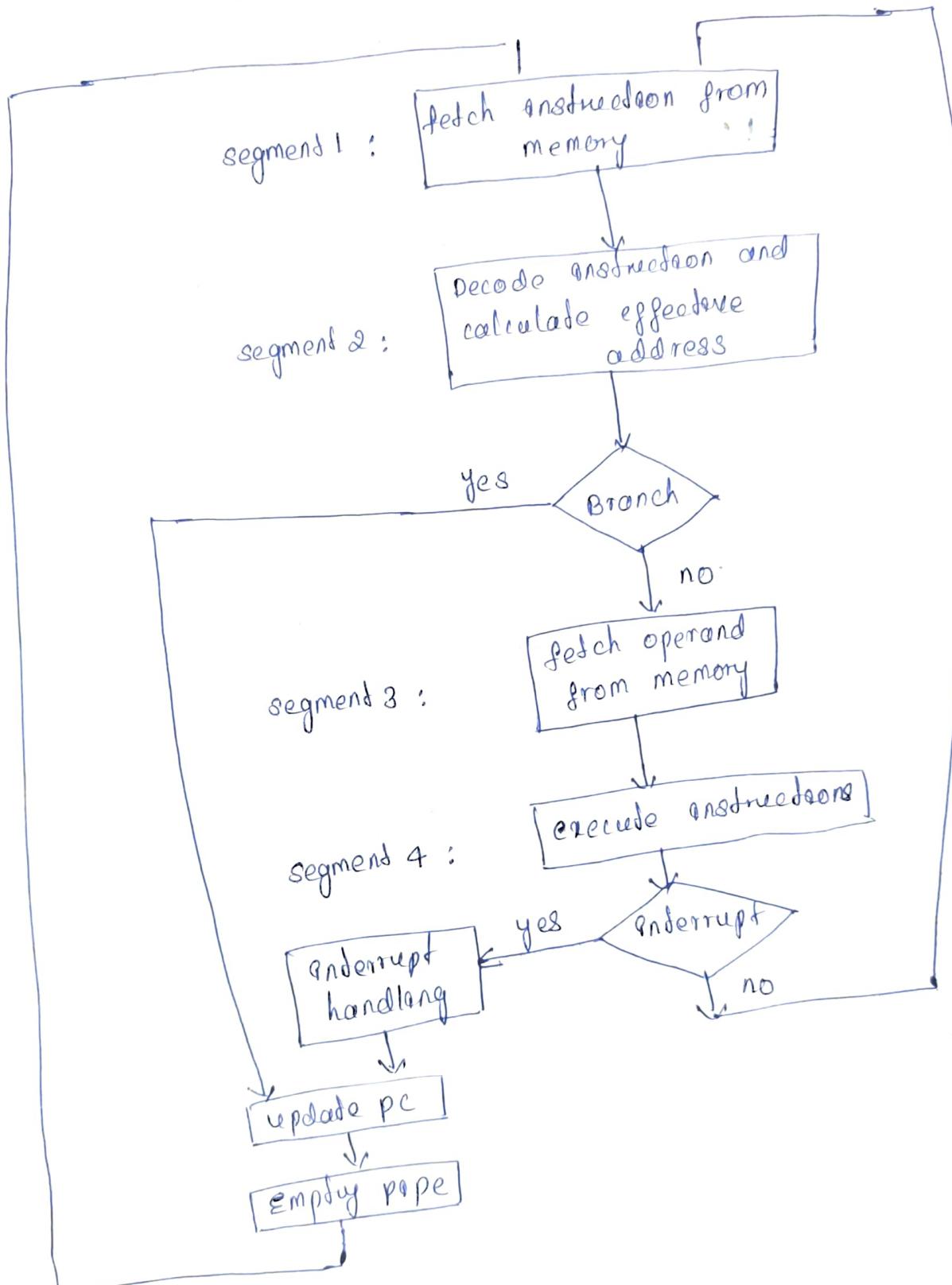
① 4 segment instruction pipeline :-

1. fetch an instruction (RI)
2. decodes the instruction and calculating effective address (DA).
3. fetch operand (FO)
4. executing instruction (EX).

@

## 4 segment cpu pipeline

⑨



## (b) Among of instruction pipeline.

Step:	1	2	3	4	5	6	7	8	9	10	11	12	13
Instruction:	1 FI DA FO EX												
2		FI DA FO EX											
(branch) 3			FI DA FO EX										
4				FI - - FI DA FO EX									
5					- - - FI DA FO EX								
6						FI DA FO EX							
7							FI DA FO EX						

① At step 3 is a branch instruction, as soon as this instruction is decoded, on segment DA in step 4, transfer from FI to DA of other instruction is halted until branch instruction is executed in step 6.

② If branch is taken, new instruction is fetched in step 7.

③ If branch not taken, instruction fetched previously in step 4 can be used.

④ In general there are 3 major difficulties cause instruction pipeline to deviate from normal operation.

⑤ Resource conflicts → caused by access to memory by 2 segments at same time, can be solved by using separate instruction and data memories.

⑥ Data dependency → arises when instruction depends on result of previous instruction, but result is not available.

⑦ Branch difficulties → arise from branch and other instructions that change value of PC.