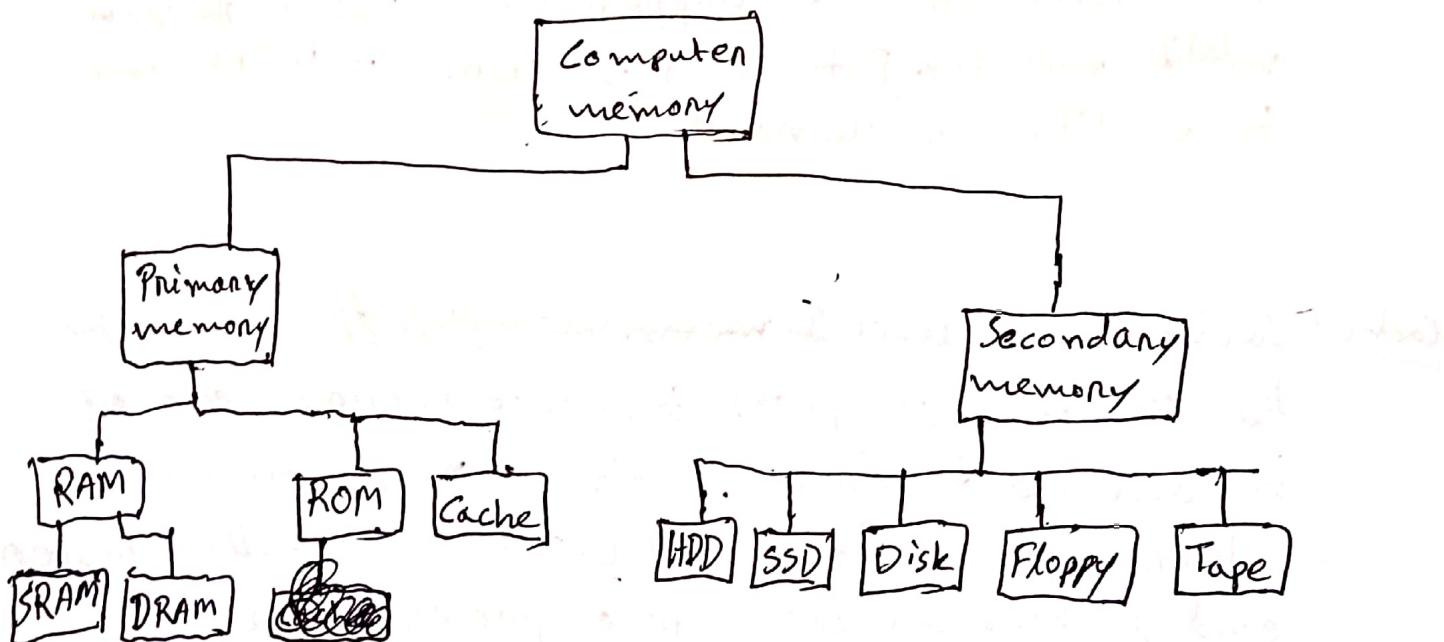


Name: Rashik Rahman
ID: 17201012

1

Answer to the Q. No. 1

A memory is just like a human brain. In computer science memory is considered as the a storage space where data to be processed and instructions required for processing are stored. Memory types are given below.



RAM: RAM stands for Random Access Memory. As the name suggests it can be accessed in any random order. Read, write operation can be done on it and it is a volatile memory. Data is fetched from secondary memory and is stored in RAM for processing. After process result is given to RAM and RAM saves it to secondary memory. There are two types of RAM, SRAM and DRAM.

(2)

DRAM: DRAM stands for Dynamic RAM. DRAM is a type of ~~memory~~ random access ~~memory~~ semiconductor memory that stores each bit of data in a memory cell consisting of a tiny capacitor and a transistor.

SRAM: SRAM stands for Static RAM. It is a type of memory that uses latching circuitry to store each bit. ~~It is~~ As it is a type of RAM so it is also a volatile memory.

ROM: ROM stands for Read Only Memory. That is only read operation can happen on it. It is ~~the~~ non volatile memory. Data that is stored in ROM can't be modified or changed.

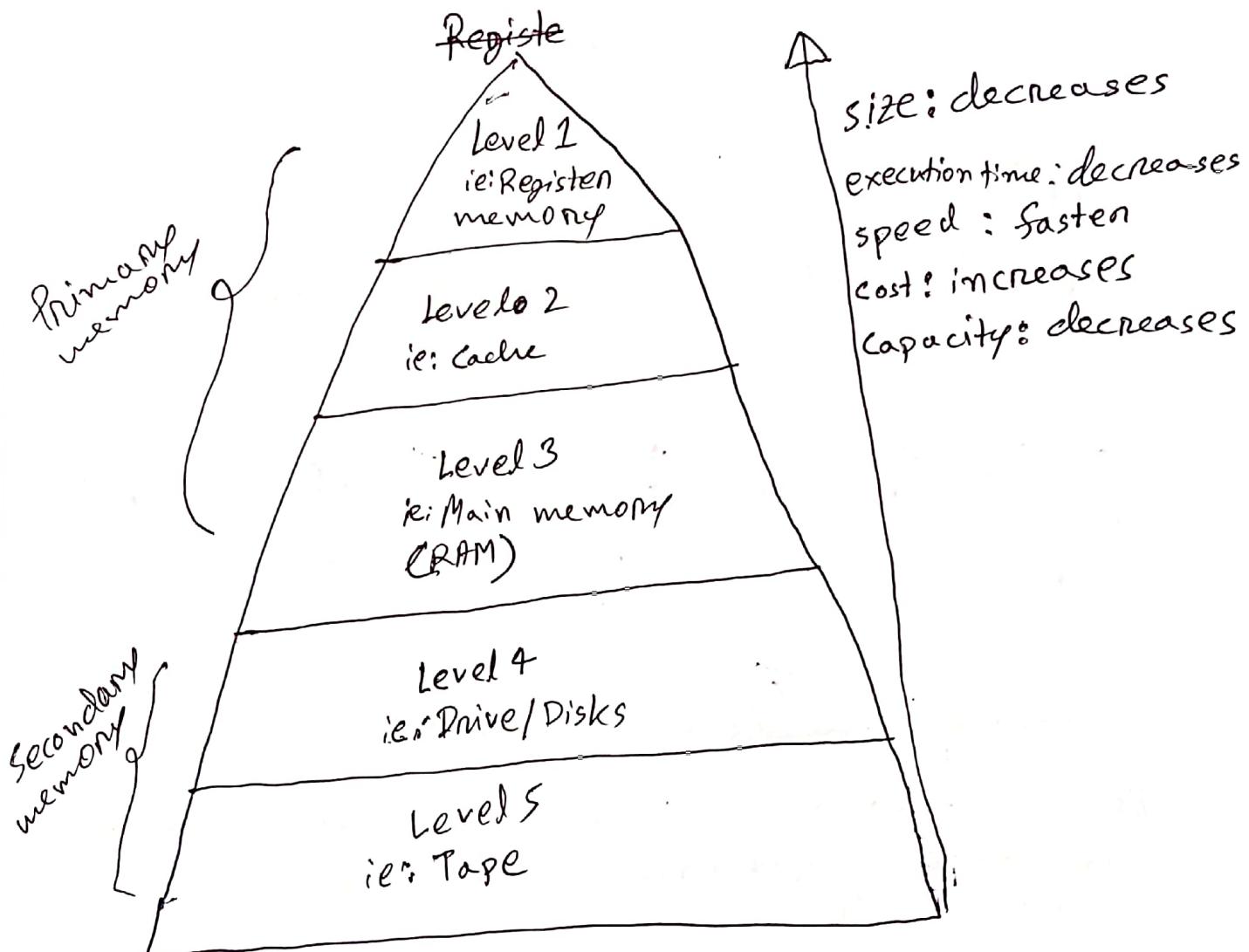
Cache: Cache is a level ~~memory~~ memory. It is used by the CPU of a computer to reduce average cost of to access data from main memory or in other words from RAM. Cache is smaller, faster and located closer to a processor core.

HDD: Hard Disk Drive

HDD/SSD/Disk/Tape: Hard Disk Drive, Solid State Drive, compact Disk, floppy & tape these all are secondary memory that is used to store data permanently. Data can also be changed or modify or delete in these memories.

Answer to the Q. No. 2

Memory hierarchy drawn below:



In computer architecture the memory hierarchy separates computer storage into a hierarchy, based on response time. Since the response time, complexity and capacity are related the levels may also be distinguished by their performance and controlling technologies. From the fig we can see that

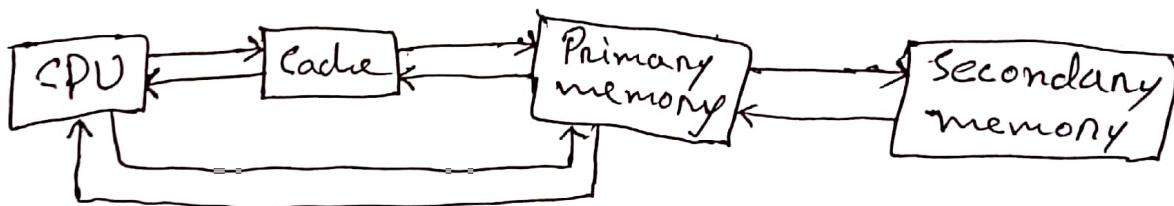
(4)

When we go to upper level from lower level the size or form factor decreases, execution time decrease thus increases speed, and capacity also decreases.

Answer to the Q. NO. 3

~~ERD~~

Steps of cache memory are given below.

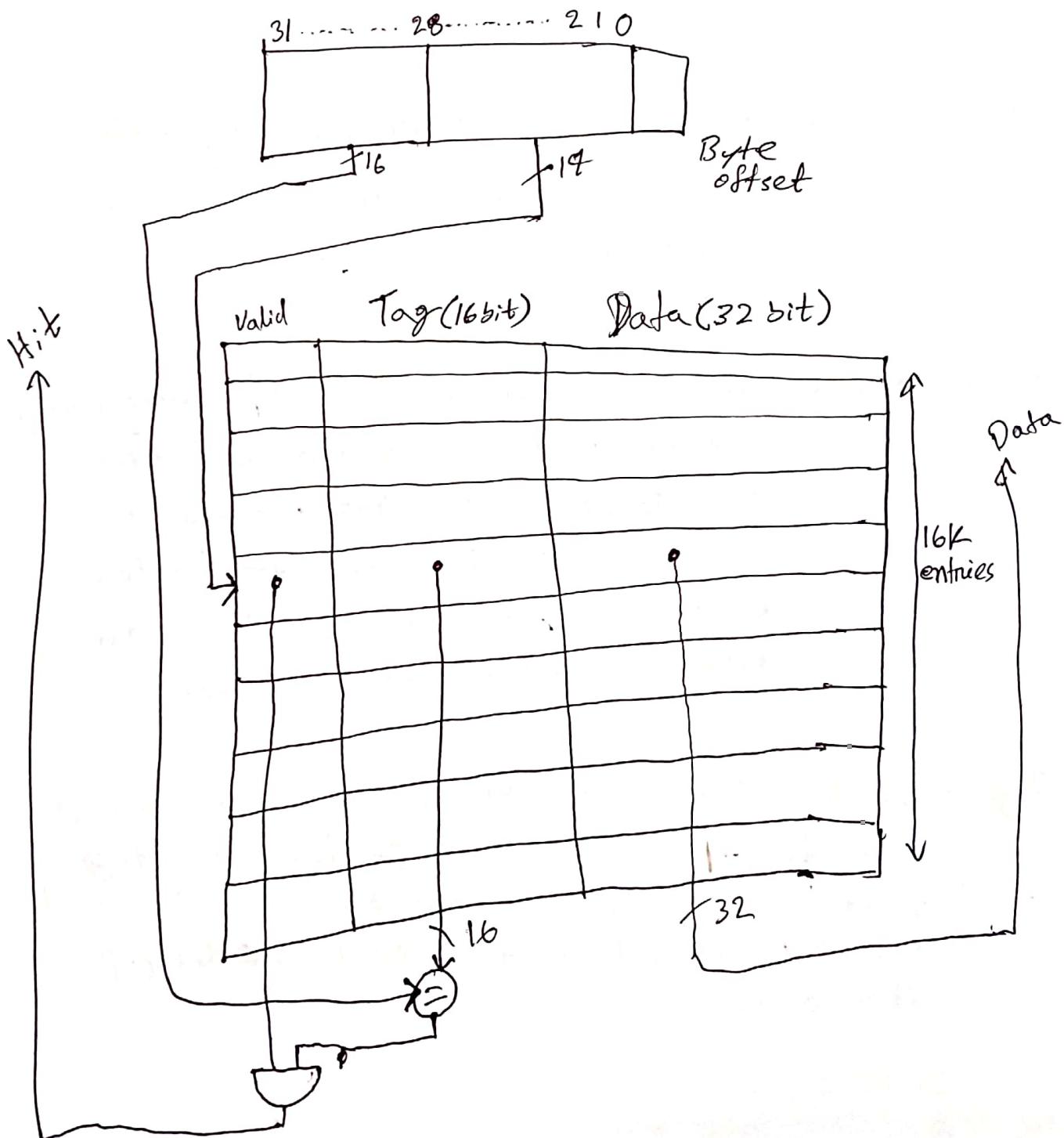


- A request is made by the CPU to cache.
- If the data is found in cache then it is retrieved & returned to the CPU. This is called hit.
- If the data isn't in cache, then it is retrieved returned from lower level memory like main memory or secondary memory. This is called miss. Data flows from lower to upper level by ~~two~~ moving up ~~two~~ between two consecutive memory level at a time.

(5)

Answer to the Q. No. 4

Hardware logic of cache hit on miss drawn below:



Answer to the Q.No. 5

Definition of the following terms:

- **Cache:** It is a hardware cache used by the central processing unit of a computer to reduce the average cost to access data from the main memory.
- **Memory hierarchy:** The memory hierarchy separates computer storage into a hierarchy based on response time.
- **Direct map cache:** A cache where the cache location for a given address is determined from the middle address bits. If the cache line is 2^n then the bottom n address bits correspond to an offset within a cache entry.
- **Tag :** The tag field of the CPU address is compared with the tag of the line. If the two tag match then hit occurs and the desired word is found in the cache. Otherwise it's a miss.

- Hit time: Hit time is the time it takes to access a memory location in the cache.
- Miss rate: Miss rate is the ratio of total miss and total request.
- Block : A memory block is a ~~cont~~ contiguous chunk of memory.
- Hit: When CPU request for data. If the data is in cache then it is hit.
- Miss: When CPU requests for data. If the isn't in cache and has to be retrieved from lower level memory then it is a miss.

Answer to the Q. No. 6

Given,

Cache size = 32 words

Main memory size = 512 words

Block size = 4 words.

Thus

$$\therefore \text{Memory blocks} = \frac{512}{4} = 128 \text{ blocks}$$

$$\therefore \text{Cache lines} = \frac{32}{4} = 8 \text{ blocks.}$$

$$\begin{aligned}\therefore \text{Total bits required for physical address} &= \log_2 (\text{no. of memory size}) \\ &= \log_2 (512) \\ &= 9 \text{ bits.}\end{aligned}$$

We know that physical address is divided into 3 sections, index, tag and ~~block~~ block offset.

$$\begin{aligned}\therefore \text{Index} &= \log_2 (\text{no. of cache lines}) \\ &= \log_2 (8) \\ &= 3\end{aligned}$$

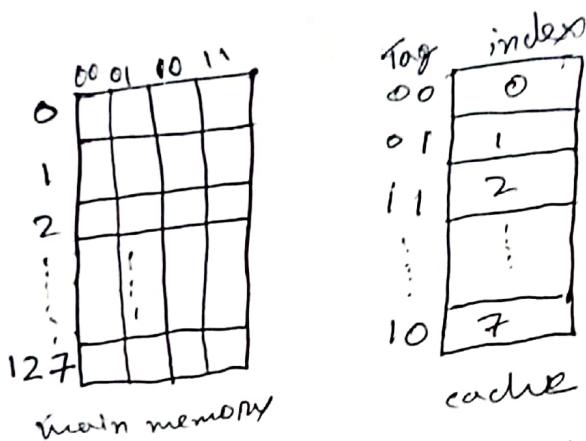
$$\begin{aligned}\text{Offset} &= \log_2 (\text{Block size}) \\ &= \log_2 (4) \\ &= 2\end{aligned}$$

$$\therefore \text{Tag bit} = 9 - 3 - 2 = 4$$

So, $\xleftarrow{\quad} 9 \xrightarrow{\quad}$

Tag	Index	Offset
4bit	3bit	2bit

Direct mapping cache config:



Dot In case of data miss data will flow from memory to cache. So to know if it is hit or miss we need to know cache config. If the requested tag's index matches with the index bits of the physical address then its a hit otherwise its a miss. For miss we need to map data from memory to cache. Let assume we move block 12 to cache. For direct method we mod 12 with total cache block so we do $12 \% 8 = 4$. So block 12 of memory will goto block 4 of cache.

Answer to the Q. No. 7

Four questions to design cache are given below:

- Where can a block be placed in upper level? (Block placement → mapping function)
- How is a block found if it is in the upper level? (Block identification)
- Which block should be replaced on a miss? (Block replacement)
- What happens on a write? (Write strategy).

Answer to the Q. No. 8

Sequence: 0, 8, 0, 6, 6, 8, 7, 11

① Direct method:

Requested memory address	Hit/ Miss	Blocks			
		0	1	2	3
0	miss	mem[0]			
8	miss	mem[8]			
0	miss	mem[0]			
6	miss			mem[6]	
6	hit			mem[6]	
8	miss	mem[8]			
7	miss				mem[7]
11	miss				mem[11]

ii) 2 way set associative:

(12)

$$\text{Set} = \frac{4}{2} = 2$$

Now we mod block address by 2.

requested memory address	hit/miss	Blocks			
		set 0		set 1	
0	miss	mem[0]		2	3
8	miss		mem[8]		
0	hit	mem[0]			
6	miss		mem[6]		
6	hit		mem[6]		
8	miss	mem[8]			
7	miss			mem[7]	
11	miss				mem[11]

iii) Full associative:

~~$$\text{Set} = \frac{4}{4} = 1; \text{ we mod block address by } 1$$~~

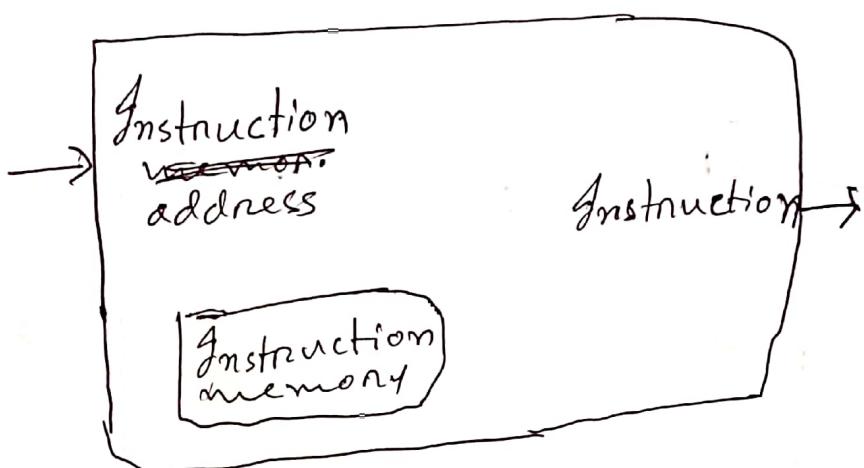
requested memory address	hit/miss	Blocks			
		0	1	2	3
0	miss	mem[0]			
8	miss		mem[8]		
0	hit	mem[0]			
6	miss			mem[6]	
6	hit			mem[6]	
8	hit		mem[8]		
7	miss				mem[7]
11	miss	mem[11]			

Answer to the Q. No. 9

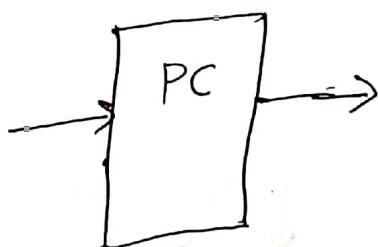
of processor

All the functional units are drawn below:

(i) Instruction memory:

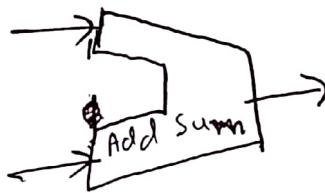


(ii) Program counter:

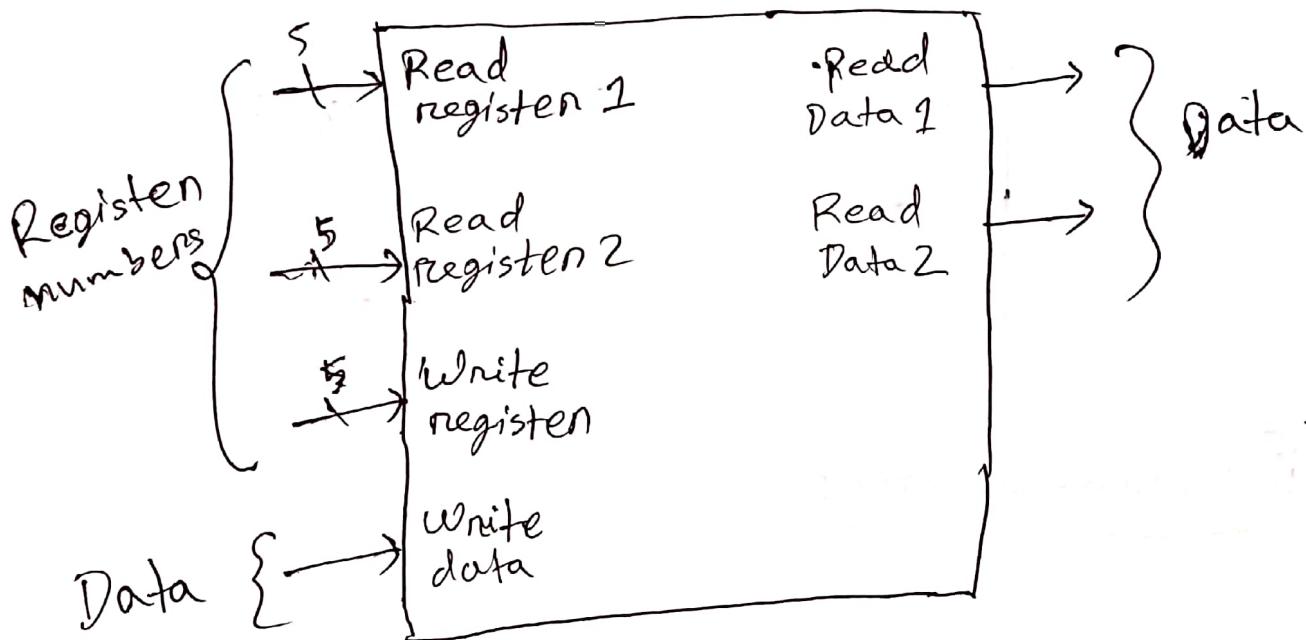


~~iii) Address~~

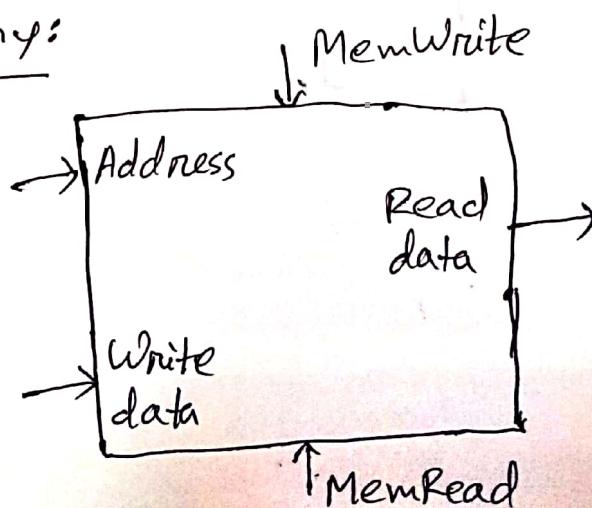
iii) Address:



iv) Register files:



v) Data memory:

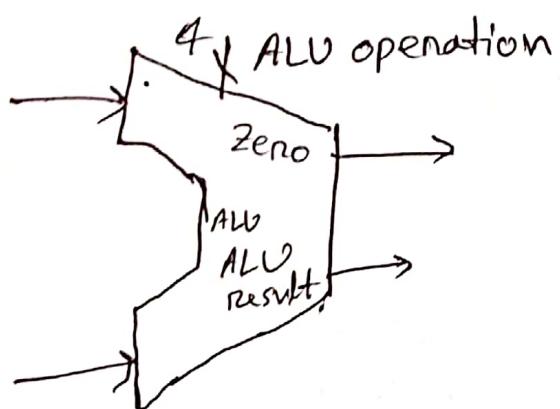


15

(vi) Sign extend:



(vii) ALU



(16)

Answer to the Q. No. 10

$$R[i] = R[i+2] - Y$$

here,
i = 12

so,

$$R[12] = R[14] - Y$$

Mips instructions,

~~i~~ ~~for~~

Considering,

$$R \rightarrow \$S_0(16)$$

$$Y \rightarrow \$S_1(17)$$

$$\text{temp} \rightarrow \$t_0(8)$$

(i)

lw \$t_0, 56(\$S₀)

(ii) ~~sub~~ sub \$t_0, \$t_0, \$S₁

(iii) sw \$t_0, 48(\$S₀)

~~i For the 1st~~

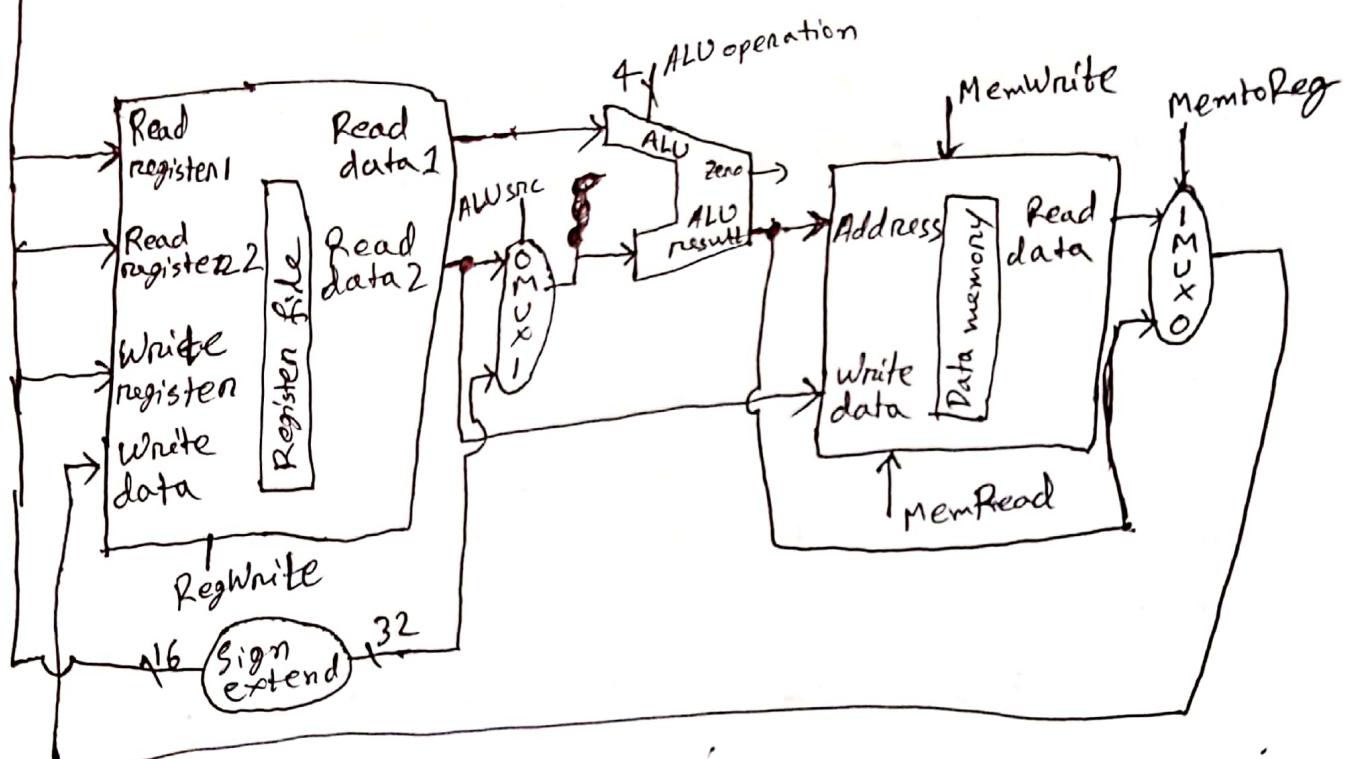
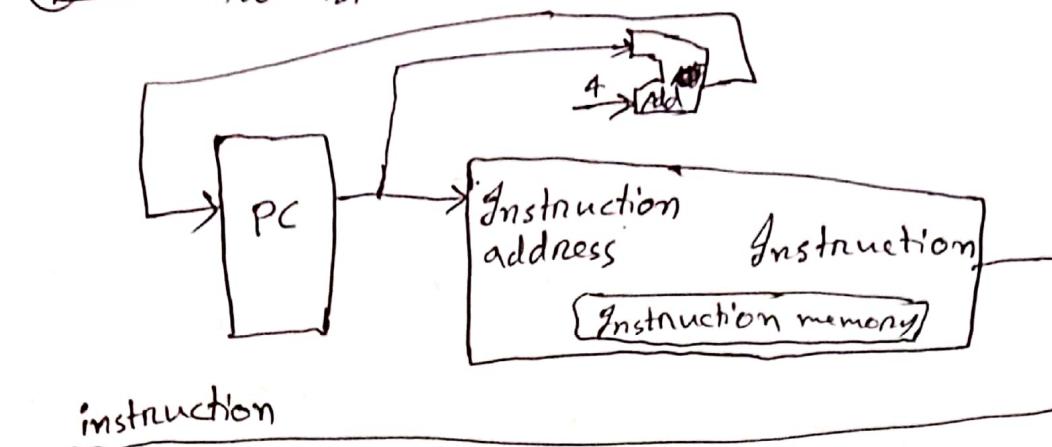


fig: Datapath organization.

Instruction address is placed on the instruction memory from program counter. As each instruction is of 4 words, so every time PC passes an address to instruction memory it increments by 4 by using an adder. Instruction memory retrieves and passes the instruction, corresponding to that instruction address.

(i) For the 1st sub instruction

As it is an I-type instruction so the number of \$S0 that is 16 would be passed to read register 1 and the content of \$S0 that is the base address would come out of read data 1. The ALUsrc MUX is SET to indicate it's an I-type instruction. ~~offset~~ 16 bit offset would be passed to sign extend and will be converted to 32 bit offset. Now ~~offset~~ offset and base address would be passed to ALU and we'll get the physical address from ALU result. This physical address is passed to the address port of data memory. Now MemRead and MemtoReg are SET so data memory will fetch the data corresponding to that physical address and pass it to write data in register file. RegWrite is SET and the number of \$to that is 8 is passed to write register. This is how data is fetched from memory and stored in register.

(ii) For the 2nd sub instruction:

The 2nd instruction is R-type so the ALUsrc is CLEAR. Now we pass the number of ^{the} register \$to that is 8 to Read register 1 and the number of the register \$S1 that is 17 to Read register 2. The content of these two registers are passed to ALU then sub operation happen. The result is passed to MUX. MemtoReg is CLR so the result is passed to write data and ~~the~~ RegWrite is SET so the number of \$to is passed to write register. Thus the result is stored in the register \$to.

(iii) For the 3rd instruction (sw);

This sw instruction is an I-type instruction so ALUsnc is SFT. The number of \$S0 that is 16 is passed to Read register 2 and the 16 bit offset is converted to 32 bit offset via sign extend. Then the content of \$S0 that is the base address and the offset is passed to ALU. ALU returns the physical address and passes it to Data memory. The number of \$T0 that is 8 is passed to read register 2 and the content of \$T0 that is the result comes out of Read data 2 and is passed to write ~~data~~ data section of Data memory. MemWrite & MemtoReg are SFT so the result is stored at that ~~give~~ calculated physical address. This is how save instruction works.