

Memory Hierarchy:

- * Memory Hierarchy is an enhancement to organize the memory such that it minimize the access time.
- * The Memory Hierarchy was developed based on a program behaviour known as locality of references.
- * The figure below clearly demonstrates the different levels of Memory hierarchy.

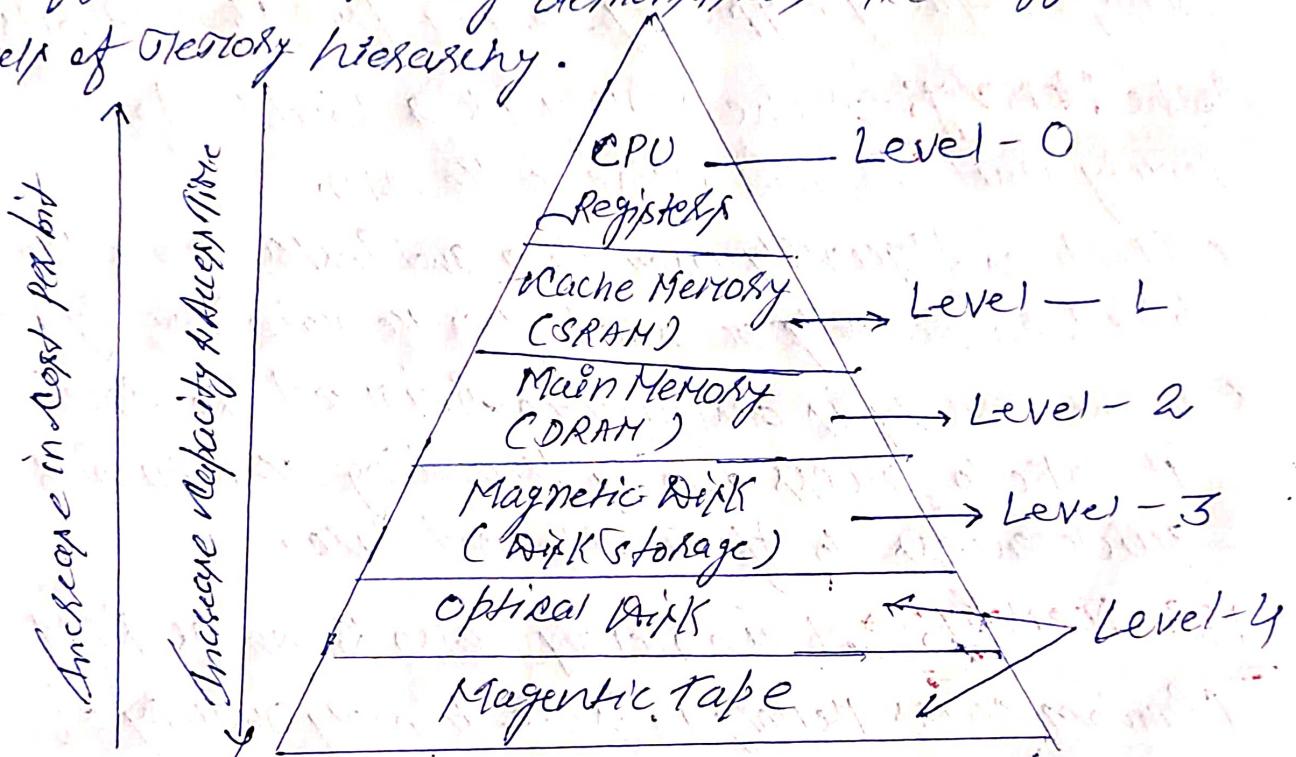


Fig: Memory Hierarchy Design

- * This memory hierarchy design is divided into 2 (Two) main types:-

1. External Memory or Secondary Memory:-

& Comprising of Magnetic disk, Optical disk, Magnetic tape i.e. peripheral storage devices which are accessible by the processor via A/I module.

2. Internal Memory or Primary Memory:-

& Comprising of Main memory, cache memory & CPU Registers. They are directly accessible by the processor.

* These are typically fast levels of Memory in a memory hierarchy:

Registers: Registers are small, high-speed Memory units located in the CPU.

* They are used to store the most frequently used data and instructions.

* Registers have the fastest access time and the smallest storage capacity, typically ranging from 16 to 64 bits.

Cache Memory: Cache Memory is a small, fast Memory unit located close to the CPU.

* It stores frequently used data and instructions that have been recently accessed from the Main Memory.

* Cache Memory is designed to minimize the time it takes to access data by providing the CPU with quick access to frequently used data.

Main Memory: → Main Memory, also known as RAM (Random Access Memory), is the primary memory of a computer system.

* It has a larger storage capacity than Cache Memory, but it is slower.

* Main Memory is used to store data and instructions that are currently in use by the CPU.

* Types of Main Memory:-

• Solid RAM: It stores the binary information in F-FP and information remaining valid until power is supplied.

* It has faster access time and is used in implementing Cache Memory.

- **Dynamical RAM:** → It stores the binary information as a charge on the capacitor. It requires self-refreshing circuitry to maintain the charge on the capacitor after few milliseconds. It contains more memory cells per unit area as compared to SRAM.

Secondary Storage:

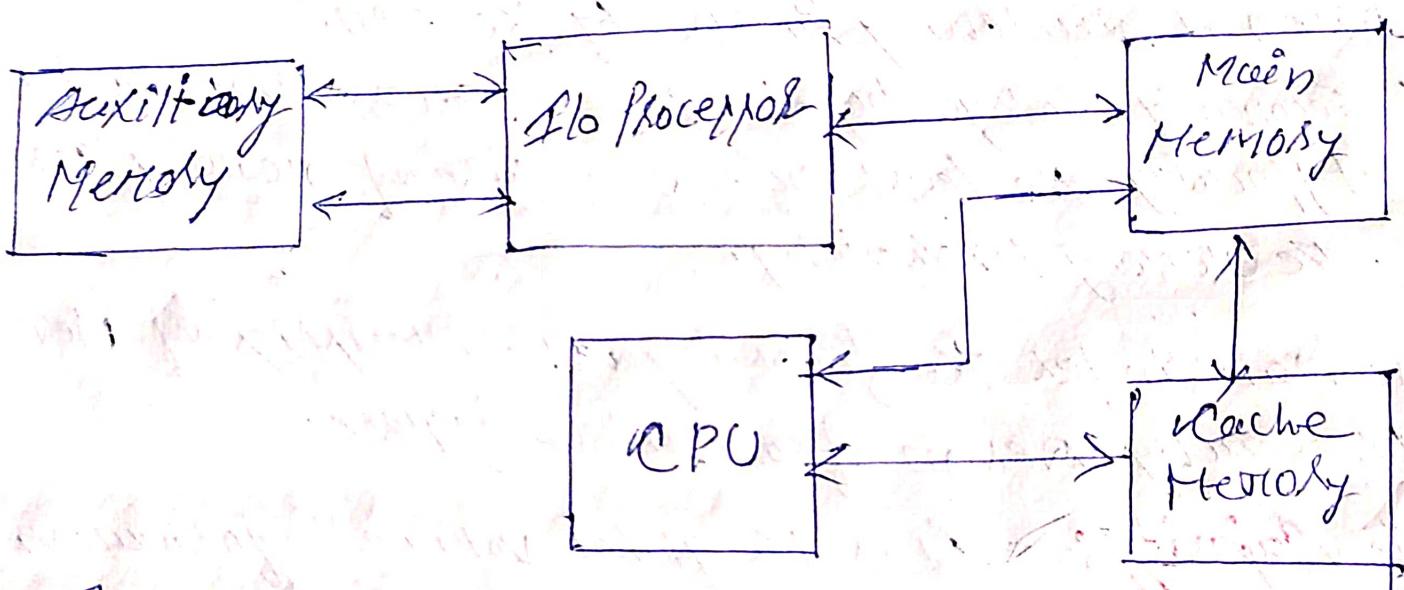
- * Secondary Storage, such as HDD, and SSD, is non-volatile memory unit that has a larger storage capacity than main memory.
- * It is used to store data and instructions that are not currently in use by the CPU.
- * Secondary storage has the slowest access time and is typically the least expensive type of memory in the memory hierarchy.

We can infer the following characteristics of Memory Hierarchy design from above figure.

1. **Capacity:** — It is the global volume of information the memory can store. As we move from top to bottom in the hierarchy, the capacity increases.
2. **Access Time:** — It is the time interval between the read/write request and the availability of the data. As we move from top to bottom in the hierarchy, the access time increases.
3. **Performance:** —

Faster when the computer system was designed without memory hierarchy design the speed gap increases between the CPU Register

- Segregated and Shared Memory due to large difference in access time.
- * This results in lower performance of the system and thus enhancement was required.
 - * This enhancement was made in the form of Memory Hierarchy design because of which the performance of the system increases. (Cache Memory speed)
 - * **Latency per bit:** As we move from bottom to top in the hierarchy, the latency per bit increases i.e. Internal Memory is faster than External Memory.

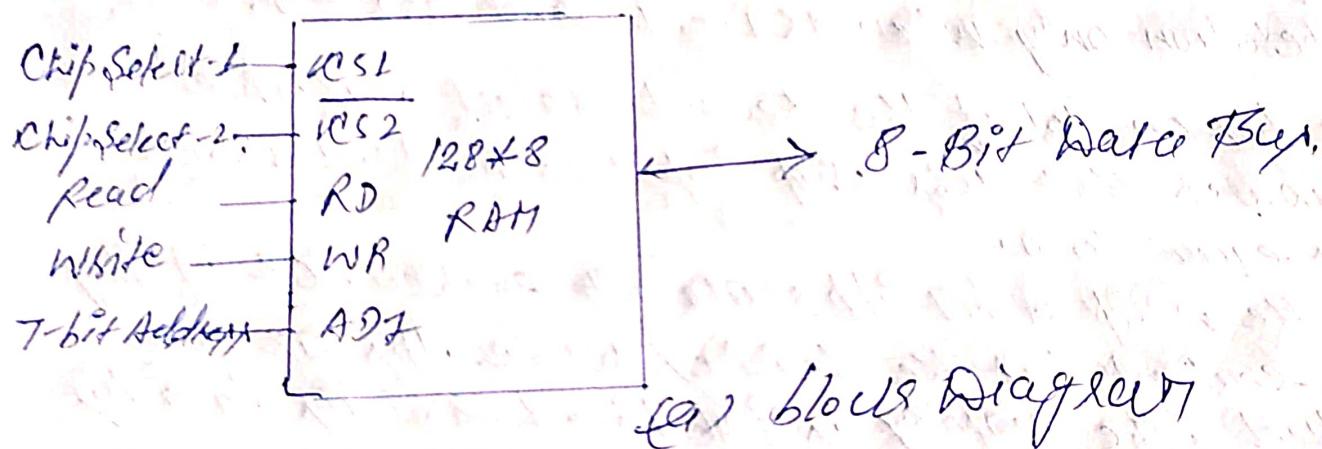


- * Flo Processor manages data transfer between Auxiliary Memory and Main Memory.

RAM and ROM chips:

A RAM chip is better-suited for communication with the CPU if it has one or more control pins that select the chip only when needed. A more common feature is a bidirectional data bus that follows the transfer of data either from memory to CPU during a read operation, or from CPU to memory during a write operation.

The block diagram of a RAM chip is shown in below fig: The capacity of the memory is 128 words of eight bits (one byte) per word. This requires a 7-bit address and an 8-bit bidirectional data bus.



(a) Block Diagram

CS1	CS2	RD	WR	Memory Function	State of Data Bus
0	0	X	X	Inhibit	High-Z Impedance
0	1	X	X	Inhibit	High-Z Impedance
1	0	0	0	Inhibit	High-Z Impedance
1	0	0	1	Write	Input data to RAM
1	0	1	X	Read	Output data from RAM
1	1	X	X	Inhibit	High-Z Impedance

(b) Function-table

Fig: Typical RAM Chip.

- * The read and write pins specify the memory operation and the two chip select (CS) control pins are for enabling the chip only when it is selected by the CS pin.
- * The availability of more than one control input to select the chip facilitates the decoding of the address lines when multiple chips are used in the microcomputer.
- * The read and write pins are sometimes combined into one line, labeled R/W. When this chip is selected, the two binary states in this line specify the two operations of read or write.
- * The function table listed in the previous page specifies the operation of the RAM chip. The unit is in operation only when $\overline{CS1} = 1$ and $\overline{CS2} = 0$.
- * The bar on top of the second select variable indicates that this chip input is enabled when it is equal to 0.
- * If the chip select pins are not enabled, OR if they are enabled but the read or write pins are not enabled, the memory is inhibited and its data bus is in a high impedance state.
- * When $\overline{CS1} = 1$ and $\overline{CS2} = 0$, the memory can be placed in a write or read mode.
- * When the WR pin is enabled, the memory stores a byte from the data bus onto a location specified by the address pins.
- * When RD pin is enabled, the content of the selected byte is placed onto the data bus.
- * The RD and WR signals control the memory operation as well as the bus buffer is associated with the bidirectional data bus.

- * A ROM chip is organized internally in a similar manner. However, since a ROM can only read, the data bus can only be an output stroke.
- * The block diagram of a ROM chip is shown in below diagram.
- * For a some-size chips, it is possible to have those bits of ROM plan of RAM, because the internal binary cells in ROM occupy less space than in RAM. For this reason, the diagram specifies a 512-byte ROM, while the RAM has only 128 bytes.

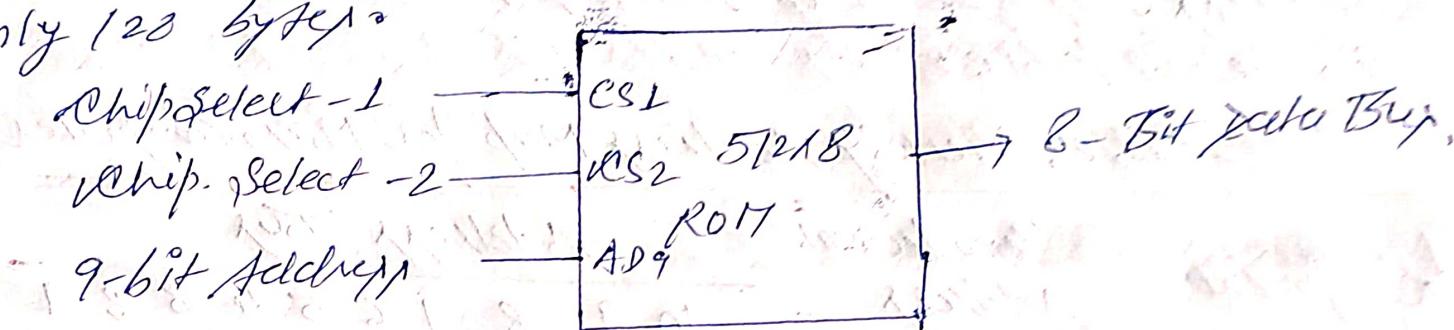


Fig: Typical ROM Chip.

- * The nine address lines in the ROM chip specify any one of the 512 bytes stored in it. The two chip select inputs must be $\overline{CS1} = 1$ and $\overline{CS2} = 0$ for the unit to operate.
- * Otherwise, the data bus is in a high-impedance state. There is no need for a read or write control because the unit can only read.
- * Then when the chip is enabled by the two select pins, the byte selected by the address lines appears on the data bus.

Memory Address Map:

- * The ~~decreasing~~ addressing of a memory can be established by means of a table that specifies the memory addresses assigned to each chip.
- * The table, called a memory address map, is a pictorial representation of assigned address space for each chip in the system.
- * To demonstrate with a particular example, suppose that a computer system needs 512 bytes of RAM and 512 bytes of ROM. The memory address map for this configuration is shown in below tables.

Table: Memory Address Map for Micro-computer

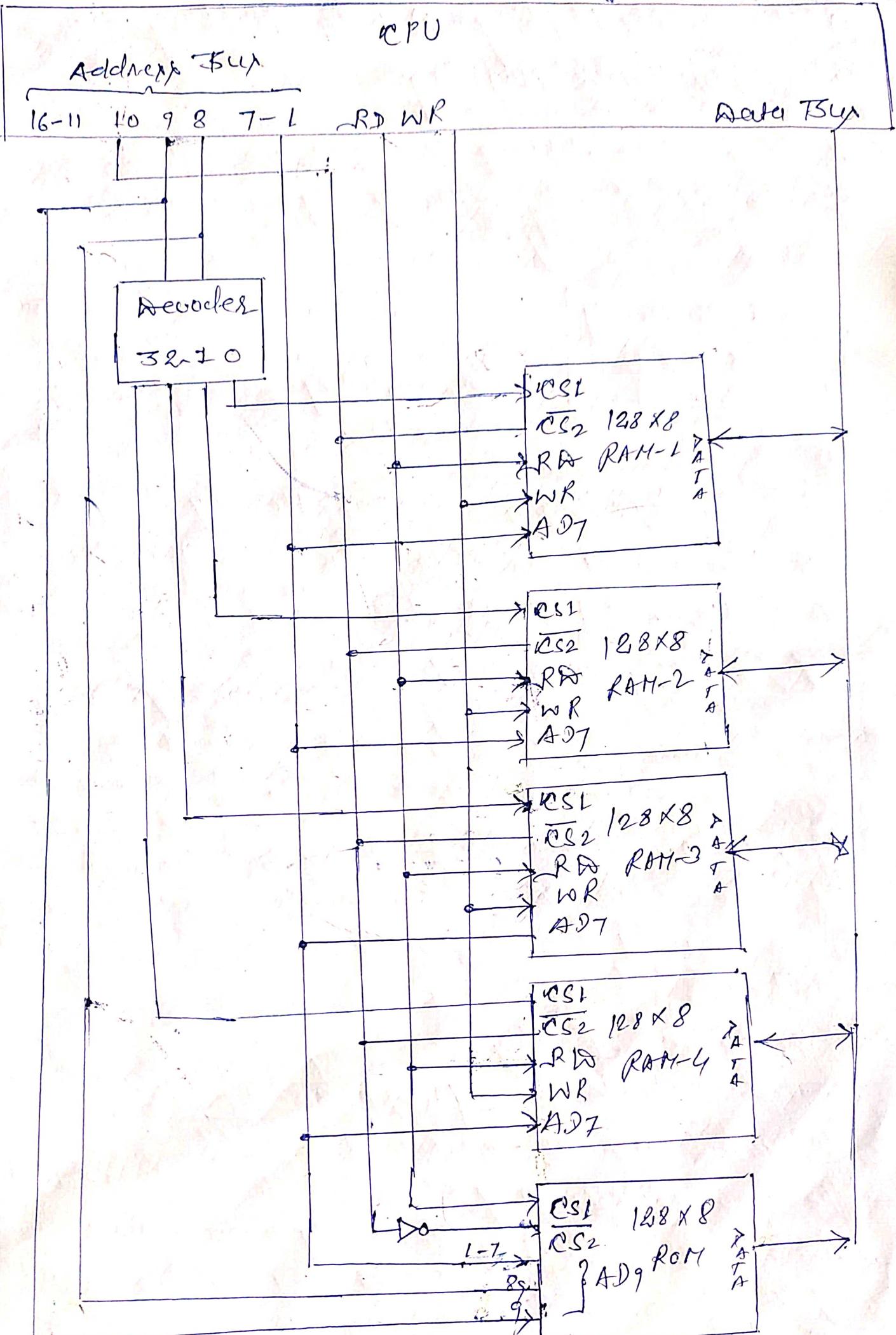
Component	Hexadecimal Address	Address Bus								
		10	9	8	7	6	5	4	3	2
RAM 1	0000-007F	00	0	X	X	X	X	X	X	X
RAM 2	0080-00FF	00	1	X	X	X	X	X	X	X
RAM 3	0100-017F	01	0	X	X	X	X	X	X	X
RAM 4	0180-01FF	01	1	X	X	X	X	X	X	X
ROM	0200-03FF	1	X	X	X	X	X	X	X	X

- * The component column specifies whether a RAM or ROM chip is used. The hexadecimal address column assigns a range of hexadecimal equivalent address for each chip. The address bus lines are listed in the third column.
- * Although there are 16 lines in the address bus, the table shows only 10 lines because the other

- * Bits are not used in this example which we assume to be zero.
- * The small x 's under the address bus lines designate those lines that must be connected to the address input in each chip.
- * The RAM chips have 128 bytes and needs seven address lines. The ROM chips has 512 bytes and needs 9 address lines.
- * The x 's are always assigned to the low-order bus lines: lines 1 through 7 for the RAM and lines 1 through 9 for the ROM.
- * It is now necessary to distinguish between four RAM chips by assigning to each different address.
- * For this particular example we chose bus lines 8 and 9 to represent four distinct binary combinations.
- * Note that any other pair of unused bus lines can be chosen for this purpose.
- * The table clearly shows that the nine low-order bus lines constitute a memory space for RAM equal to $2^9 = 512$ bytes.
- * The distinction between RAM and ROM address is done with another bus line. Here we chose line 10 for this purpose. When line 10 is 0, the CPU selects a RAM, and when this line is equal to 1, it selects the ROM.
- * The equivalent hexadecimal address for each chip is obtained from the information

- from the information about the address bus assignment.
- The address bus lines are subdivided into groups of four bits each so that each group can be represented with a hexadecimal digit.
- The first hexadecimal digit represents lines 13 to 16 and is always 0.
- The next hexadecimal digit represents lines 9 to 12, but lines 11 and 12 are always 0.
- The most ~~hexadecimal digit~~ significant digit.
- The range of hexadecimal address for each component is determined from the X'x associated with it. Here X'x represents a binary number that can range from an all 0's to an all 1's value.

Memory Connection to CPU:



Fog of Memory connection to the CPU

Cache Memory:

- * If the active portions of the program and data are placed in a fast small memory, the average memory access time can be reduced, thus reducing the total execution time of the program. Such a small fast memory is referred to as a cache memory.
- * It is placed between the CPU and Main Memory as illustrated in memory hierarchy.
- * The cache memory access time is less than the access time of Main Memory by a factor of 5 to 100.
- * The cache is the fastest component in the memory hierarchy and approaches the speed of CPU components.
- * The performance of cache memory is frequently measured in terms of a quantity called hit ratio.
- * When the CPU refers to memory and finds the word in cache, it is said to produce a hit. If the word is not found in cache, it is in Main Memory and it causes a miss.
- * The ratio of the number of hits divided by the total CPU references to memory (hits plus misses) is the hit ratio.
- * For ex: A computer with cache access time of 100 ns, a main memory access time of 1000 ns, and a hit ratio of 0.9 produces an avg. access time of 200 ns. Without a cache memory, whose access time is 1000 ns,
- * The transformation of data from Main Memory to cache memory is referred to as a mapping process.
- * Three types mapping procedure for cache memory.

- 1) Associative Mapping
- 2) Direct Mapping
- 3) Set-Associative Mapping

* For these three mapping procedure, we will take a specific example of a memory organization as shown in below fig.

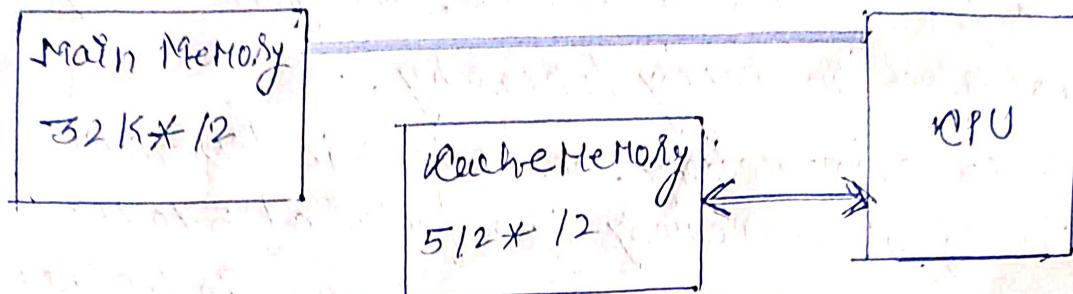


Fig: Example of Cache Memory.

* The Main Memory can store 32 K words of 12 bits each. The cache is capable of storing 512 of these words at any given time. For every word stored in cache, there is a duplicate copy in Main Memory. The CPU communicate with both memories. At first send a 15-bit address to cache. If there is a hit, the CPU accept the 12-bit data from cache. If there is a miss, the CPU reads the word from Main Memory and the word is then transferred to cache.

Associative Mapping

* The fastest and most flexible cache organization uses an associative memory.

* The associative memory stores both the address and content (data) of the memory word. This permits any location in cache to store any word from main memory.

- * The diagram shows three words properly stored in the cache. The address value of 15 bits of 1 is shown as a five-digit octal number, and the corresponding 12-bit word is shown as a four-digit octal number.
- * A CPU address of 15 bits is placed in the Address Register and the Associative Memory is searched for a matching address.
- * If the address is found, the corresponding 12-bit data is read and sent to the CPU.

CPU address (15 bits)

Associative Registers	
Address	Data
01000	3450
02777	67101
22345	1234

- * If no match occurs, the Main Memory is accessed for the word. The address-data pair is then transferred to the associative value memory.
- * If the cache is full, an address-data pair must be displaced to make room for a pair that is needed and not presently in the cache.
- * The decision for replacement fault is done on the basis of replacement algorithm that designs choices for the cache.

Direct Mapping

- * Associative Memories are expensive
- * Compare to RAMs because of the added logic associated with each cell.
- * The possibility of using a RAM for the cache is investigated in fig: below.

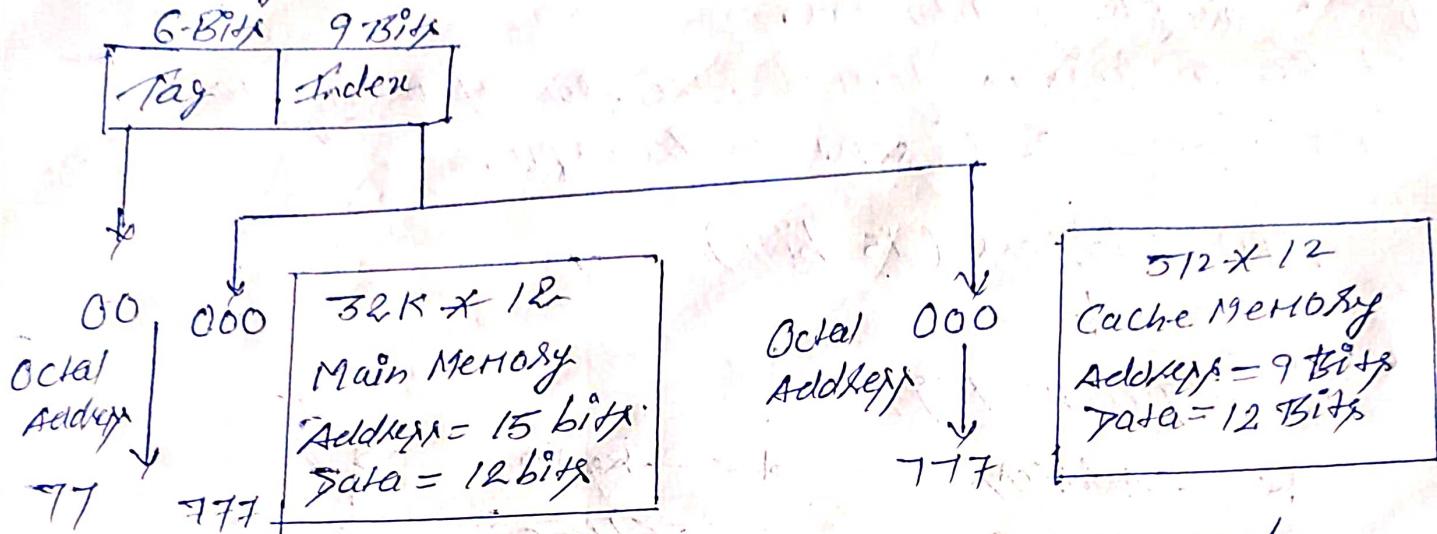


Fig: Addressing Relationship between Main and Cache Memory.

- * The 15-bit address is divided into two fields.
- * The most significant bit constitute the Index field and the remaining six bits form the Tag field.
- * The number of bits in the Index field is equal to the number of address required to access the Cache Memory.
- * In general case, there are 2^k words in Cache Memory and 2^m words in Main Memory.
- * The m -bit memory address is divided into two fields: k bits for the Index field and $m-k$ bits for the Tag field.
- * The direct mapping cache organisation uses the m -bit address to access the Main Memory and

- * The K-bit Index to select the cache.
- * The T-bit tag to identify the word in the cache.
- * The T-bit word in cache consists of the data word and its associated tag. Address forwarded because it is part of the data word and its associated tag.
- * When a new word is first brought into the cache, the tag bits are stored alongside the data bits.
- * When the CPU generates a memory, the Index field is used for the address to access the cache. The tag field of the CPU address is compared with the tag in the word read from the cache.
- * If there is no match, there is miss and the required word is read from Main Memory. It is then stored in the cache together with the new tag, replacing the previous value.

Main Memory Address	Main Memory
00000	1220
00777	2340
01000	3450
01777	4560
02000	5670
02777	6780

(a) Main Memory

Address	Index	Tag	Data
000	00	00	1220
777	02	6780	

(b) Cache Memory

Fig: Direct Mapping Cache Organization.

- The disadvantage of direct mapping is that the hit ratio can drop considerably if two or more words whose addresses have the same tag but different word addresses are accessed separately. However, this possibility is minimised by the fact that such words are relatively far apart in the address range (multiples of 512 locations in the example).
- The direct mapping is shown in ~~previous~~ fig.

Index Page data:

	6	6	5
Block	Tag	Block	Word
Block-0	000	01	3450
	007	02	6578
Block-1	010		
	017		
Block-2			
Block-3	770	02	
	777	02	6710

Fig: Direct Mapping Cache with block size of 8 words.

- The word at address zero is presently stored in the cache (Index=000, tag=00, data=1220) and RPU wants to access the word at address 002000.
- The Index address is 000, so it is used to access the cache. The two tags are compared. The cache tag is 00 but the address tag is 02, which does not produce a match. Therefore, the given memory is accessed and the data word 5670 is transferred to the RPU. The cache word at Index address 000 is then replaced with a tag of 02 and data 5670.

- * The direct-mapping example just described has a block size of one word. The same organization but using a block size of 8 words is shown in fig on the previous page.
- * The index field is now divided into two parts: the block field and the word field.
- * In a 512-word cache there are 64 blocks of 8 words each, since $64 \times 8 = 512$.
- * The block number is specified with a 6-bit field and the word within the block is specified with a 3-bit field.
- * The tag field stored within the cache is common to all eight words of the same block.
- * Every time a miss occurs, an entire block of eight word must be transferred from main memory to cache memory.
- * This requires enthalope, the hit ratio will most likely improve with a larger block size because of the sequential nature of computer program.

Set-Associative Mapping

- * Set-associative mapping is an improvement over the direct mapping organization in that each word of cache can store two 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 pairs in one word of cache is said to form a set.
- * Each index address refers to two data words and their associated tags.

Each tag requires six bits and each data word has 12 bits, so the word length is $2(6+12) = 36$ bits. An index address of nine bits can accommodate 512 words. Then the size of cache memory is 512×36 .

- * It can accommodate 1024 words of main memory since each word of cache contains two data words.
- * 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	3450	02	5670	
777	02	6710	00	2340	

Fig: Two-way Set-Associative Mapping Cache.

- * The octal numbers listed in above fig are with reference to the main memory contents.
- * When the CPU generates a memory request, the index value of the address is used to access the cache.
- * The tag field of the CPU address is then compared with both tags in the cache to determine if a match occurs.
- * The comparison logic is done by an associative search of the tags in the set similar to an associative memory search; thus name "Set-Associative".

- * The hit ratio will improve as the size increases because store word with the same index but different tags can reside in cache.
- * However, an increase in the set size increases the number of bits in words of cache and requires more complex comparison logic.
- * When a miss occurs in a set-associative cache and the set is full, it is necessary to replace one of the tag-data pairs with a new value. The most common replacement algorithm are, FIFO, LRU.

Writing into caches

- * When the CPU finds a word in cache during a read operation, the main memory is not involved in the transfers.
- * However, if the operation is a write, there are two ways that the system can proceed.
 - * The simplest and most commonly used procedure is to update Main Memory with every memory write operation with Cache Memory being updated with every memory write operation, with Cache Memory being updated in parallel if it contains the word at the specific address. This is called the write-through method.
 - * This method has the advantage that Main Memory always contains the same data as the cache. This characteristic is important in systems with direct memory access transfers. It ensures that the data residing in Main Memory are valid at all times so that an I/O device communicating through DMA would receive the most recent update data.

The Second procedure is called the write-back Method. In this method only the cache location is updated during a write operation.

* The location is then marked by a flag so that later when the word is removed from the cache it is copied into main memory.

* The reason for the write-back method is that during the time a word resides in the cache, it may be updated several times, however, as long as the word remains in the cache, it does not matter whether the copy in main memory is out of date, since requests from the word are filled from the cache.

* It is only when the word is displaced from the cache that an accurate copy needs to be written into main memory.

* Number of blocks written in a typical program ranges between 10 and 30 percent of the total reference of to memory.

Cache Initialization:

* The cache is initialized when the power is applied to the computer or when the main memory is loaded with a complete set of programs from auxiliary memory.

* After initialization the cache is considered to be empty, but in effect it contains some non valid data. It is customary to include with each word in cache a valid bit to indicate whether or not the word contains valid data.

* The cache is initialized by clearing all the valid bits to 0. The valid bit of a particular cache

word is set to 1 the first time this word is loaded from Main Memory and stays set unless the cache has to be initialized again.

- * The introduction of the valid bit ~~is~~ means that a word in cache is not replaced by another word unless the valid bit is set to 1 and a write of tags occurs.
- * If the valid bit happens to be 0, the new word automatically replaces the invalid data. Thus the Initialization condition has the effect of forcing writes from the cache until it fills with the valid data.