

Memory unit is an essential component of computer as it is used for storing programs and data. The memory unit that communicates directly with the CPU is called main memory. A memory device that provides backup storage is called auxiliary memory. E.g. magnetic disks and tapes. Auxiliary memory stores large data files and other information. Programs and data that are currently needed by processor reside in main memory.

Memory hierarchy

The total memory capacity of a computer can be visualized as being hierarchy of components. Memory hierarchy consists of all storage device in a computer system from slow but high capacity auxiliary memory to fast main memory, to a smaller and faster cache memory.

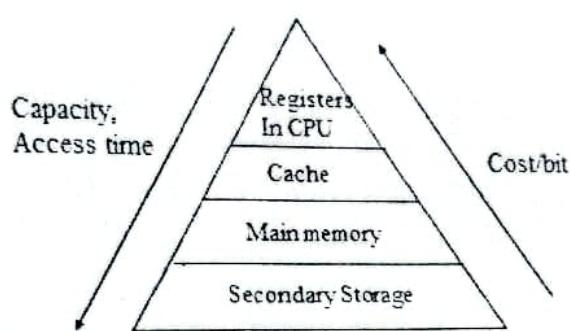


Fig. the memory hierarchy

Main memory

Main memory is the central storage unit in a system. It is a large and fast memory. Main memory is based on semiconductor integrated circuits. RAM and ROM chips are semiconductor ICs. RAM is available in two types:

1. Static RAM
2. Dynamic RAM

Static RAM consists of flip-flop and stores binary information. The information remains valid as long as power is applied to the unit.

Dynamic RAM stores information in the form of electric charge applied for capacitors. The stored charge on capacitors discharges with time thus needs to be charged by refreshing dynamic memory. DRAM offers reduced power consumption and larger storage capacity in a single memory chip.

ROM is used for storing programs that permanently resides in computer. ROM of main memory is needed for storing an initial program called bootstrap loader. Bootstrap loader is the program whose function is to start computer software operating when power is turned on. Contents of ROM remain unchanged when power is turned off. Types of ROM:

1. PROM
2. EEPROM
3. EEPROM
4. Flash EEPROM

RAM and ROM chips

A RAM chip is suited for communication with CPU if it has one or more control inputs that selects the chip only when needed. Another feature is bidirectional bus i.e. data transfer takes place either from memory to CPU during read operation and from CPU to memory during write operation. Bidirectional bus is constructed with three state buffer: logic 1 and logic 0 signals or a high impedance state (behaves like open circuit i.e. output does not carry signal with no logic).

Capacity of memory is 128 words of eight bits(one byte) per word. This requires a 7 bit address an 8 bidirectional data bus. Read and write specify memory operation. Chip select are for enabling chip only when it is selected by processor. When chip is selected the two binary states in the line specify two operations of read or write.

CS1	CS2	RD	WR	Memory fxn	State of data bus
0	0	X	X	Inhibit	High Impedance
0	1	X	X	Inhibit	High Impedance
1	0	X	X	Inhibit	High 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 Impedance

Fig: Function Table

ROM chip is organized in a similar manner. Since a ROM can only read the data bus is in output mode. For same chip more bits of ROM is possible than of RAM because internal binary cells in ROM occupy less space than in RAM.

The nine address lines to the ROM chip specify one of 512 bytes stored in it. CS=1 and (CS2)=0 selects unit to operate. Otherwise it is in high impedance state. No need for read or write control as it can only read.

Auxiliary Memory

Devices provides backup storage are called auxiliary memory. Most common auxiliary memory used in systems are magnetic disks and magnetic tapes. The important characteristics of auxiliary memory devices are access mode, access time, transfer rate, capacity and cost.

The average time required to reach a storage location in memory and obtain its time is called access time. In electromechanical devices with moving parts, the access time consists of seek time required to position read-write head to a location and transfer time required to transfer data to or from the device. Seek time is usually larger than transfer time. Thus data is organized in records or blocks. Reading or writing is done in entire block. The transfer rate is the number of characters or words that the device can transfer per second, after it has been positioned at the beginning of record.

Magnetic Disks

A magnetic disk is a circular plate constructed of both metal or plastic coated with magnetic material. Both sides of disks are used and many disks can be stacked on one spindle with read/write heads available on each surface. Bits are stored in tracks. Tracks are divided into sections called sectors. The minimum quantity of information that can be transferred is sector.

Some units may use single read/write head for each disk surface. In this type track address bits are used by a mechanical assembly to move head into specified track position before reading or writing.

Permanent timing tracks are used in disks to synchronize the bits and recognize the sectors. Here address bits specify disk number, disk surface, the sector number and track within sector. Read/write head must be positioned on specified track for information transfer.

Track nearer to center is smaller than track nearer to circumference. Thus track recording may vary. To make all records in a sector of equal length, disk use variable recording density with higher density on tracks near center than on circumference.

Magnetic Tapes

A magnetic tape transport consists of the electrical, mechanical and electronic components to provide the parts the control mechanism for a magnetic tape. It is a strip of plastic coated with magnetic recording. Bits are recorded as magnetic spots on the tape along tracks. Seven or nine bits data recorded simultaneously with a parity bit. Magnetic tape can be started, stopped can be

move forward or reverse or can be rewound. For this purpose information is stored in blocks. Gaps of unrecorded tape are inserted between records where tape can be stopped. Each record has identification bit at the beginning and at the end. Records may be fixed or variable length.

Associative Memory

It is a memory unit accessed in parallel by the content of the data itself rather than by an address. Hence, it is also called *Content Addressable Memory (CAM)*. When a word is written in an associative memory, no address is given but the memory is capable of finding an empty unused location to store the word. When the word is to be read from an associative memory, the content of the word, or part of word is specified. The memory locates all words which match the specified content and marks them for reading. It is suited for parallel searching and is more expensive than sequential memory. Associative memory is used in applications where the search time is very critical and must be very fast.

Hardware Organization

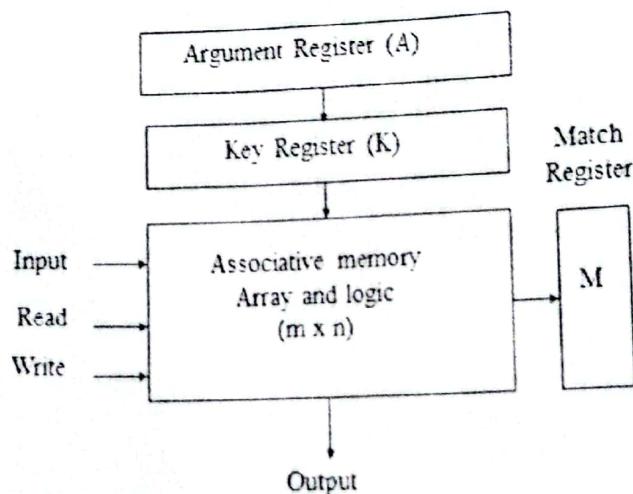


Fig. Associative memory

Associative memory consists of a memory array and logic for m words with n bits per word. The argument register A and key register K each have n bits, one for each bit of a word. The match register M has m bits, one for each memory word. Each word in memory is compared in parallel with the content of argument register. The words that match bits of argument register set a corresponding bit in two words whose corresponding bits in the match register have been set.

The key register provides a mask for choosing a particular field or key in argument word. The argument is compared with each memory word if key register contains all 1's. Otherwise only those bits that have 1's in their corresponding position of key register are compared. Thus key provides a mask or identifying piece of information.

Example of operation (from the figure below):

Suppose –

Data Register = 101 111100

Mask Register = 111 000000

Word 1 = 100 111100 => no match

Word 2 = 101 000001 => match

The relation between memory array and external registers in an associative memory is shown in figure. The cells are marked by C with two subscripts. First gives the word number and second specify the bit position in the word. Thus cell C_{ij} is the cell for bit j in word i. A_j in the argument register is compared with all bits in column j of array provided that K_j=1. This is done for all columns. If a match occurs between argument and bits in words, the corresponding M_i is set to 1. If one or more unmasked bits of argument and the word do not match, M_i is cleared to 0.

Note: Refer to your class note for match logic and read write operations

Cache Memory

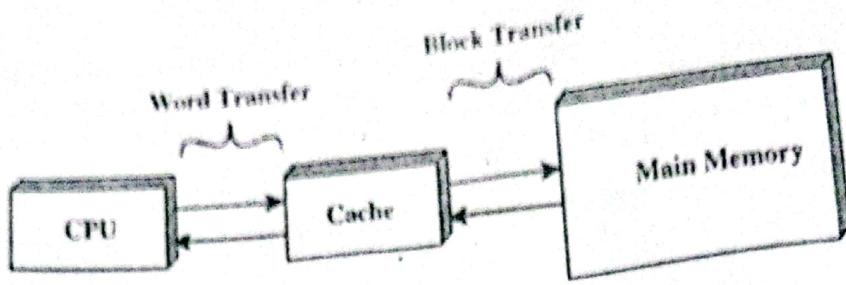
Cache memory is the fast small memory where active portion of programs and data are stored, so that the average access time is reduced. Thus it reduces the total execution time of program. Frequently used data are stored in cache memory. The cache memory access is 5 to 10 times faster than main memory. Average access time of memory can be approached to access time of cache by placing frequently used instruction into cache.

The basic operation of cache is as follows. When CPU needs to access memory, the cache is examined. If the word is found in cache, it is read from there. If the word is not found in cache, main memory is accessed to read the word. A block of word is then transferred to cache memory from main memory.

Performance of cache memory is measured in terms of a quantity called hit ratio. When the CPU refers to memory and finds it in cache, it is said to produce a hit. If the word is not found in cache, it is said to be a miss. The ratio of number of hits divided by total references to memory is the hit ratio. If the hit ratio is high, so that the most of CPU access is in cache than main memory, average access time is closer to access time of cache.

Cache – Main Memory interface

Assume an access to main memory causes a block of K words to be transferred to the cache. The block transferred from main memory is stored in the cache as a single unit called a *slot, line, or page*. Once copied to the cache, individual words within a line can be accessed by the CPU. Because of the high speeds involved with the cache, management of the data transfer and storage in the cache is done in hardware – the O/S does not know about the cache. If there are 2^n words of memory, then there will be $M = 2^n/K$ blocks in the memory.



M will be much greater than the number of lines, C , in the cache. Every line of data in the cache must be tagged in some way to identify what main memory block it is. The line of data and its tag are stored in the cache.

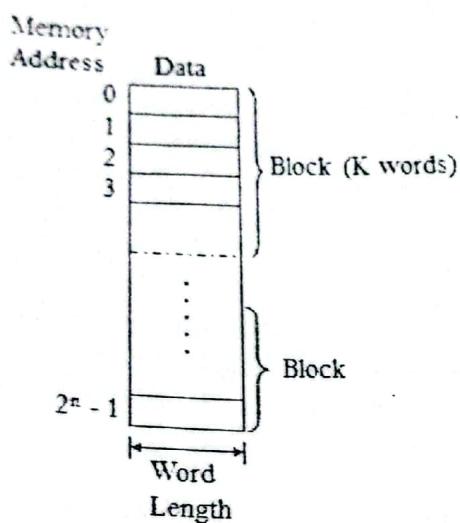


Fig. Main Memory Structure

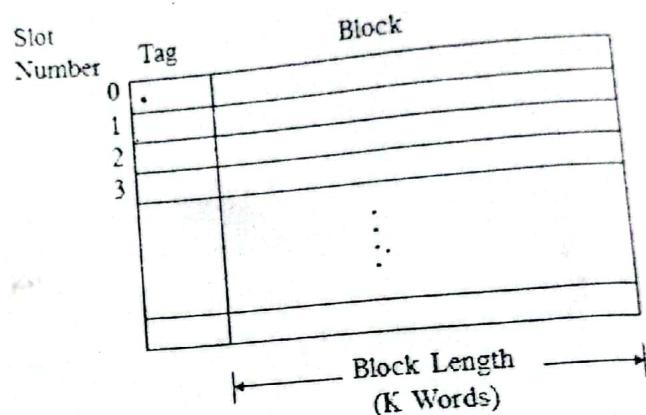


Fig. Cache Memory Structure

The transformation of data from main memory to cache memory is referred as mapping process.
Three types of mapping procedure:

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

Direct Mapping

Each main memory block is assigned to a specific line in the cache. Mapping is expressed as

$$i = j \bmod C$$

Where i is the cache line number assigned to main memory block j .

If $M = 64$, $C = 4$

line 0 can hold blocks 0, 4, 8, 12,.....

line 1 can hold blocks 1, 5, 9, 13,.....

2⁰t

2¹

line 2 can hold blocks 2, 6, 10, 14,

line 3 can hold blocks 3, 7, 11, 15,.....

Direct Mapping cache treats a main memory address as 3 distinct fields -

Tag identifier- The tag is stored in the cache along with the data words of the line.

Line number identifier- Line identifier specifies the physical line in cache that will hold the referenced address.

Word identifier (offset)-Word identifier specifies the specific word (or addressable unit) in a cache line that is to be read.

For every memory reference that the CPU makes, the specific line that would hold the reference (if it has already been copied into the cache) is determined. The tag held in that line is checked to see if the correct block is in the cache.

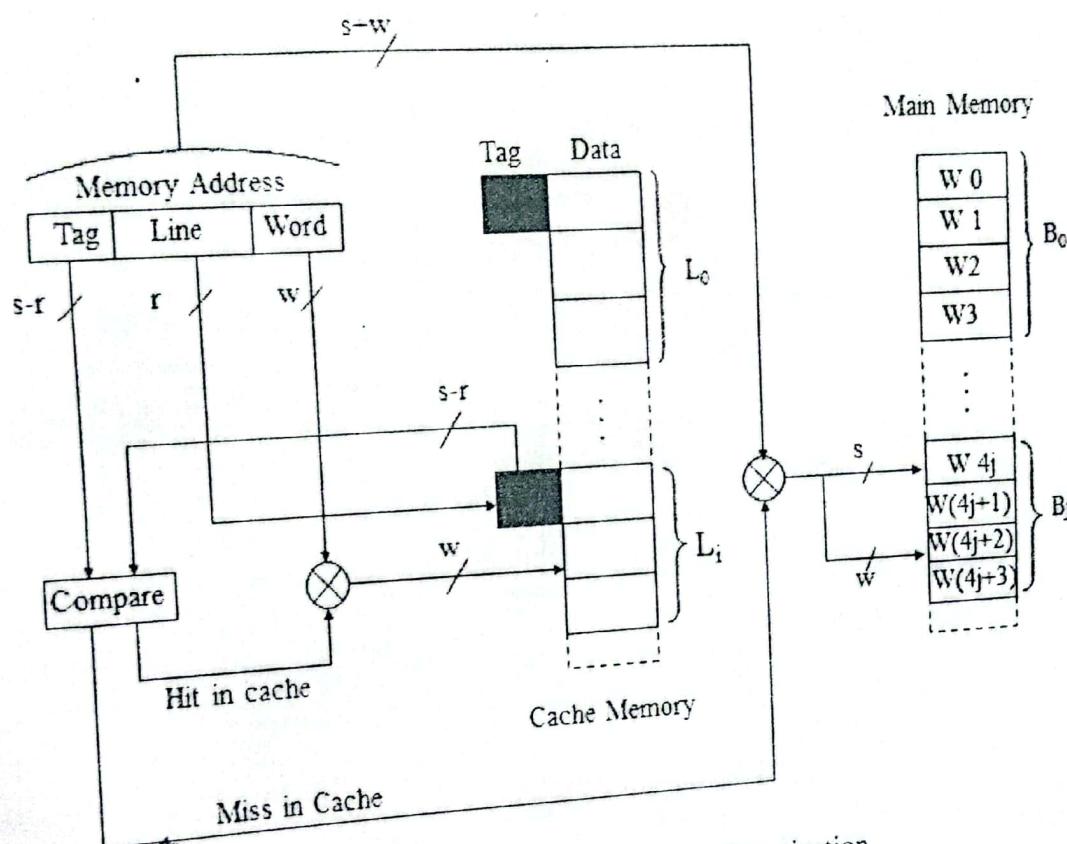


Fig: Direct Mapping Cache Organization

- **Advantages:**

- Easy to implement.
- Relatively inexpensive to implement.
- Easy to determine where a main memory reference can be found in cache.

- **Disadvantages:**

- Each main memory block is mapped to a specific cache line.

- Through locality of reference, it is possible to repeatedly reference to blocks that map to the same line number.
- These blocks will be constantly swapped in and out of cache, causing the hit ratio to be low.

Associative Mapping

Associative mapping removes the drawback of direct mapping by providing main memory block to be loaded into any line of cache. Memory address logic interprets as a tag and a word field. Tag field identifies the block of main memory. to find whether the block is in cache, the cache control logic examines every line's tag for a match.

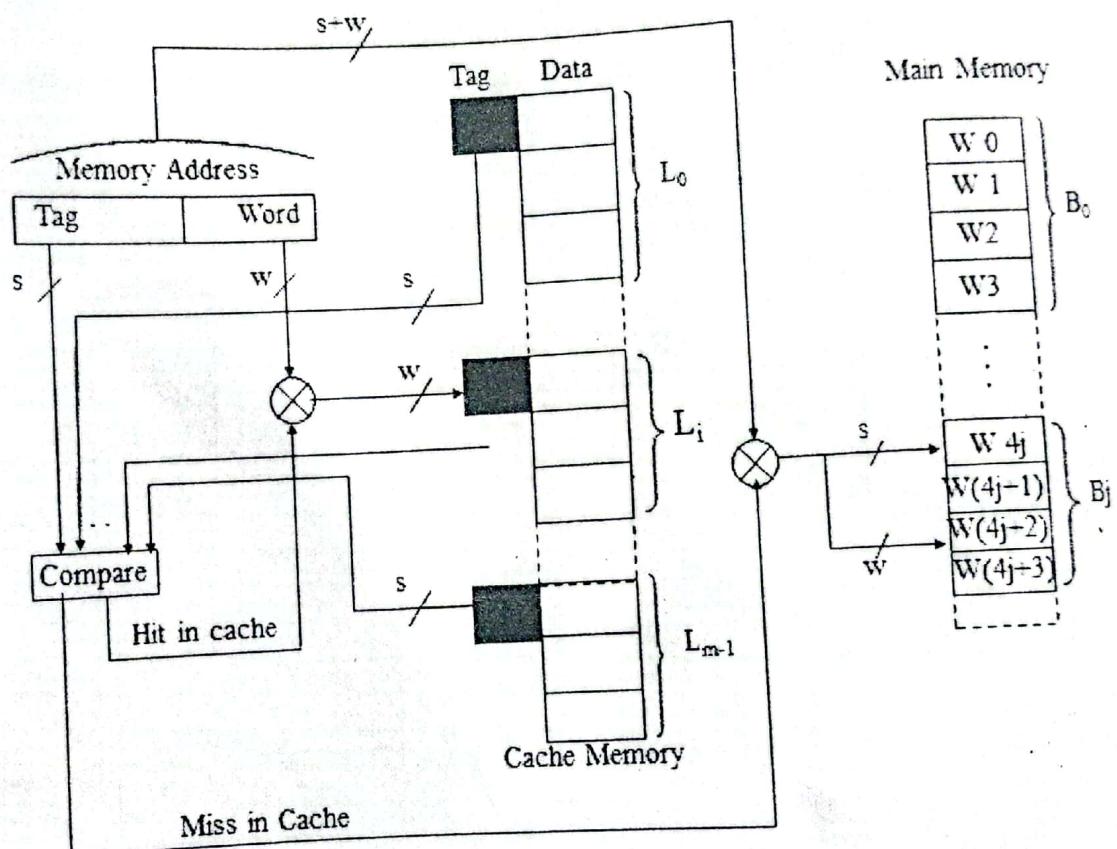


Fig: Associative Mapping

Advantage

Fast due to associative memory. Flexible to replace block when new block is read into cache.

Disadvantage

Need of complex circuit to examine the tags of all cache lines in parallel.

Set-Associative Mapping

Compromise between direct and associative mappings that builds on the strengths of both. Divide cache into a number of sets (v), each set holding a number of lines (k). A main memory block can be stored in any one of the k lines in a set such that

$$\text{set number} = j \bmod v$$

Where, j is the assigned main memory block.

If a set can hold X lines, the cache is referred to as an X -way set associative cache. Most cache system today that use set associative mapping are 2- or 4-way set associative.

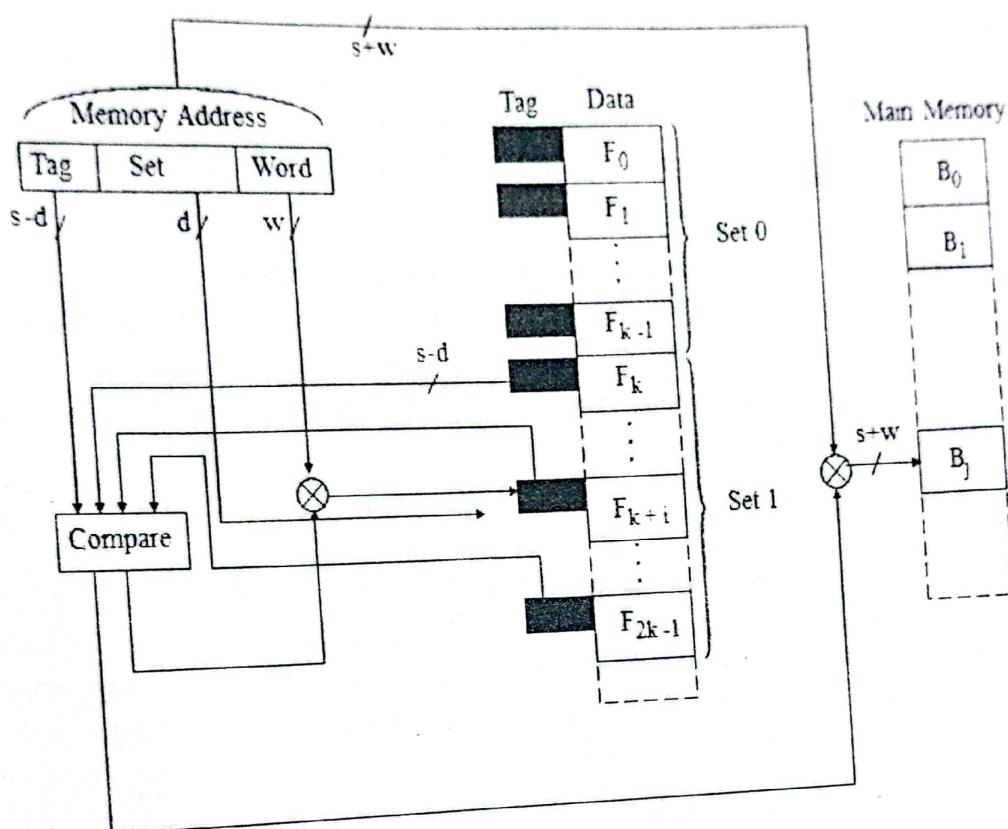


Fig: Set-Associative Cache organization

Replacement Algorithm

When a new block is brought into cache, one of the existing block must be replaced. There are four commonly used algorithms:

1. Least recently used (LRU)
2. Least frequently used (LFU)
3. First-in -First out (FIFO)
4. Random

With LRU, cache replaces the block in the set that has been in the cache for longest with no reference to it. With least frequently used that block is replaced in the set that has fewest references. It can be implemented by associating the counter with each line. FIFO replaces the

block that has been in the cache for longest. it can be implemented by round robin or circular buffer technique. Random picks a line randomly and replaces that block in the set.

Write Policies:

Before replacement of block from cache it is necessary to examine if it has been changed in cache. if no, the block can be overwritten. If yes, then there is write operation in cache and it should be updated in main memory. Write policies are:

1. Write through
2. Write back

With write through, write operations are made to main memory as well as to the cache insuring that main memory is valid. Cache memory module can monitor traffic to main memory to maintain its consistency. Main disadvantage is it creates memory traffic and bottleneck.

An alternative is write back policies. In this technique update is made only in the cache. When an update occurs UPDATE bit associated with the line is set. Then when a block is replaced it is written back to main memory if and only if the update bit is set. Disadvantage of write back policy is portion of main memory is invalid. Thus access from I/O module is allowed only through cache.

Associative memory page table:

In above mapping technique memory accommodation depends upon the blocks allocation in main memory. In this case some pages remain unused. Thus efficient way to organize page table is to construct it with a number of words equals to the number of blocks in main memory. In this way the size of memory is reduced and each location is fully utilized. This method can be implemented by means of an associative memory.

Each word in main memory contains page number with its corresponding block number. Page field in each word compared with page no. in virtual address. If a match occurs the word is read from memory and its corresponding block number is extracted.

101	Line number
-----	-------------

111	00
-----	----

001	11
010	00
101	01
110	10

Associative memory consists of two fields. The first three bits specify page number and last two bits specify block number. Virtual address is in argument register. Page no. bits are compared with all pages nos. in page field of associative memory. If page no. id found the 6-bit word is read out from memory. The corresponding block number is transferred to main memory address register. If no match occurs, a call is generated to bring required page from auxiliary memory.

The most common replacement and least recently used (LRU).

Write-through and Write-back cache write method.

Write Through

- The simplest and most commonly used procedure is to update main memory with every memory write operation.
- The cache memory being updated in parallel if it contains the word at the specified 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 updated data.

Write-Back (Copy-Back)

- 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 need be rewritten into main memory.

Virtual Memory

Virtual memory is used to give programmers the illusion that they have a very large memory at their disposal, even though the computer actually has a relatively small main memory.

A virtual memory system provides a mechanism for translating program-generated addresses into correct main memory locations.

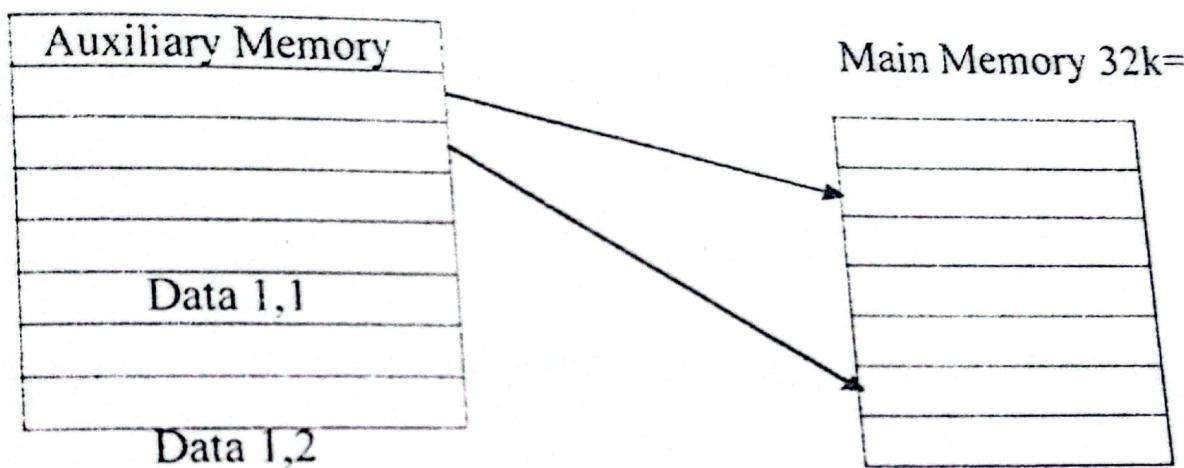
Address space

An address used by a programmer will be called a virtual address, and the set of such addresses is known as address space.

Memory space

An address in main memory is called a location or physical address. The set of such locations is called the memory space.

Program 1



Program 2

Data 2,1

Address space $1024k = 2^{10}$

As an illustration, consider a computer with a main-memory capacity of 32K 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 $2^{20} = 1024K$ 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 $N = 1024K$ and $M = 32K$.

- In a multiprogramming computer system, programs and data are transferred to and from auxiliary memory and main memory based on demands imposed by the CPU.
- Suppose that program 1 is currently being executed in the CPU. Program 1 and a portion of its associated data are moved from auxiliary memory into main memory as shown in figure 9.9.
- Portions of programs and data need not be in contiguous locations in memory since information is being moved in and out, and empty spaces may be available in scattered locations in memory.
- In our example, the address field of an instruction code will consist of 20 bits but physical memory addresses must be specified with only 15 bits.

Thus CPU will reference instructions and data with a 20-bit address, but the information at this address must be taken from physical memory because access to auxiliary storage for individual words will be too long.

Address mapping using pages.

- A The table implementation of the address mapping is simplified if the information in the address space and the memory space are each divided into groups of fixed size.
- The physical memory is broken down into groups of equal size called blocks, which may range from 64 to 4096 words each.
- The term page refers to groups of address space of the same size.
- 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 as shown in figure 9.9

At any given time, up to four pages of address space may reside in main memory in any one of the four blocks.

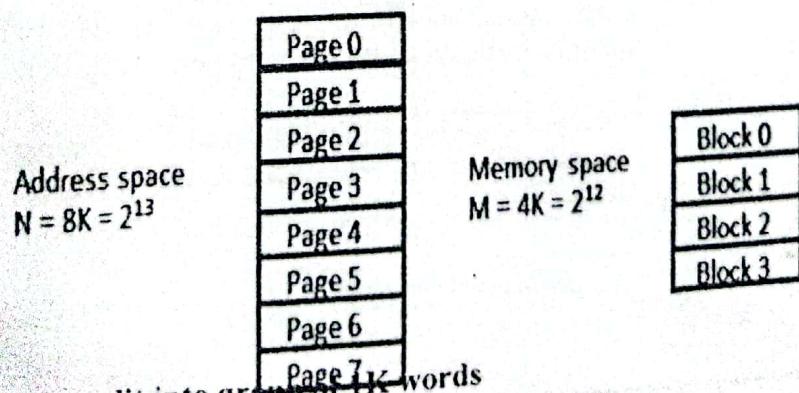


Figure 9.10 Address and Memory space split into group of 1K words

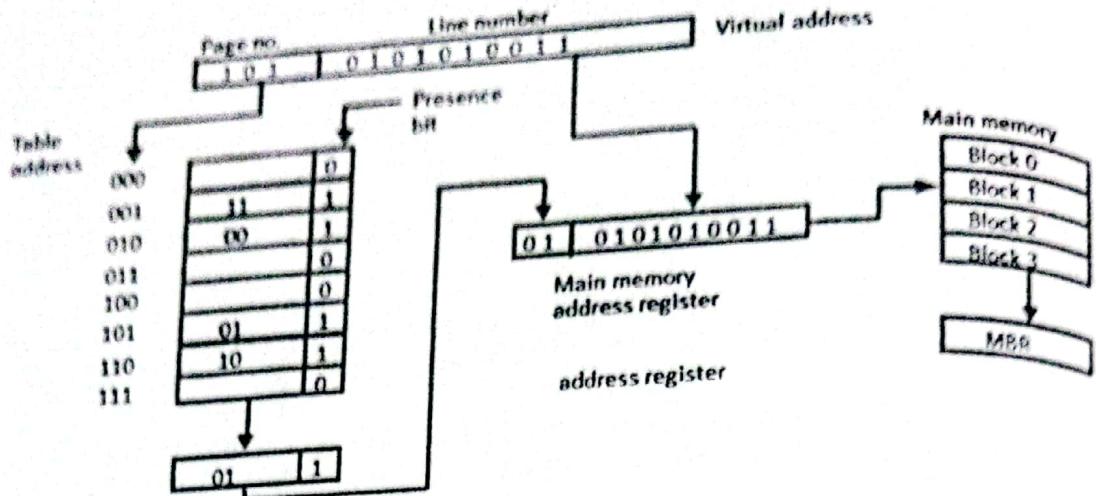


Figure 9.11: Memory table in paged system

The organization of the memory mapping table in a paged system is shown in figure 9.10.

- The memory-page table consists of eight words, one for each page.
- The address in the page table denotes the page number and the content of the word give the block number where that page is stored in main memory.
- The table shows that pages 1, 2, 5, and 6 are now available in main memory in blocks 3, 0, 1, and 2, respectively.
- A presence bit in each location indicates whether the page has been transferred from auxiliary memory into main memory.
- A 0 in the presence bit indicates that this page is not available in main memory.
- The CPU references a word in memory with a virtual address of 13 bits.

The three high-order bits of the virtual address specify a page number and also an address for the memory-page table.

The content of the word in the memory page table at the page number address is read out into the memory table buffer register.

If the presence bit is a 1, the block number thus read is transferred to the two high-order bits of the main memory address register.

The line number from the virtual address is transferred into the 10 low-order bits of the memory address register.

A read signal to main memory transfers the content of the word to the main memory buffer register ready to be used by the CPU.

If the presence bit in the word read from the page table is 0, it signifies that the content of the word referenced by the virtual address does not reside in main memory.

A segment is a set of logically related instructions or data elements associated with a given name.

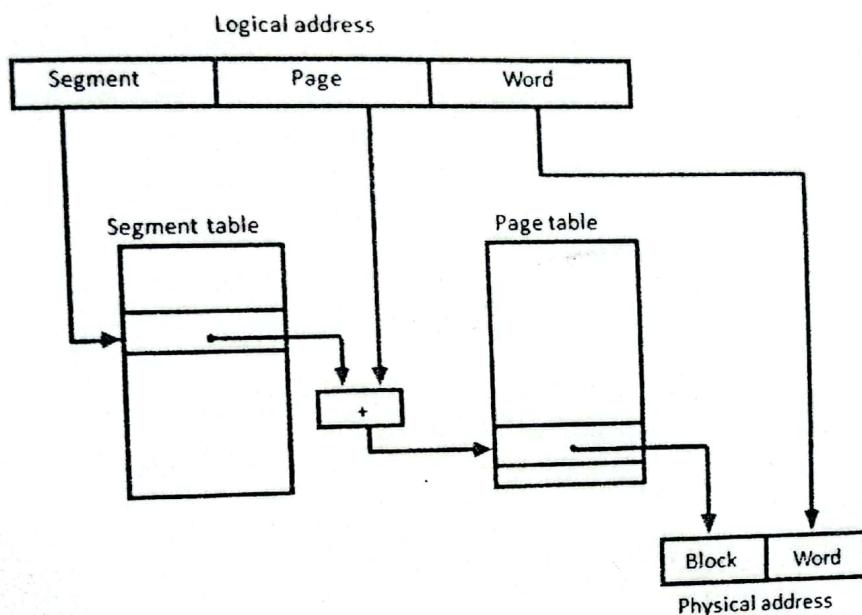
Logical address

The address generated by segmented program is called a logical address.

Segmented page mapping

- The length of each segment is allowed to grow and contract according to the needs of the program being executed. Consider logical address shown in figure 9.12.

Figure 9.12: Logical to physical address mapping



- The logical address is partitioned into three fields.
- The segment field specifies a segment number.
- The page field specifies the page within the segment and word field gives specific word within the page.
- A page field of k bits can specify up to 2^k pages.
- A segment number may be associated with just one page or with as many as 2^k pages.
- Thus the length of a segment would vary according to the number of pages that are assigned to it.

The mapping of the logical address into a physical address is done by means of two tables, as shown in figure 9.12.

The segment number of the logical address specifies the address for the segment table.

The entry in the segment table is a pointer address for a page table base.

The page table base is added to the page number given in the logical address.

The sum produces a pointer address to an entry in the page table.

- The concatenation of the block field with the word field produces the final physical mapped address.
- The two mapping tables may be stored in two separate small memories or in main memory.
- In either case, memory reference from the CPU will require three accesses to memory: one from the segment table, one from the page table and the third from main memory.
- This would slow the system significantly when compared to a conventional system that requires only one reference to memory.