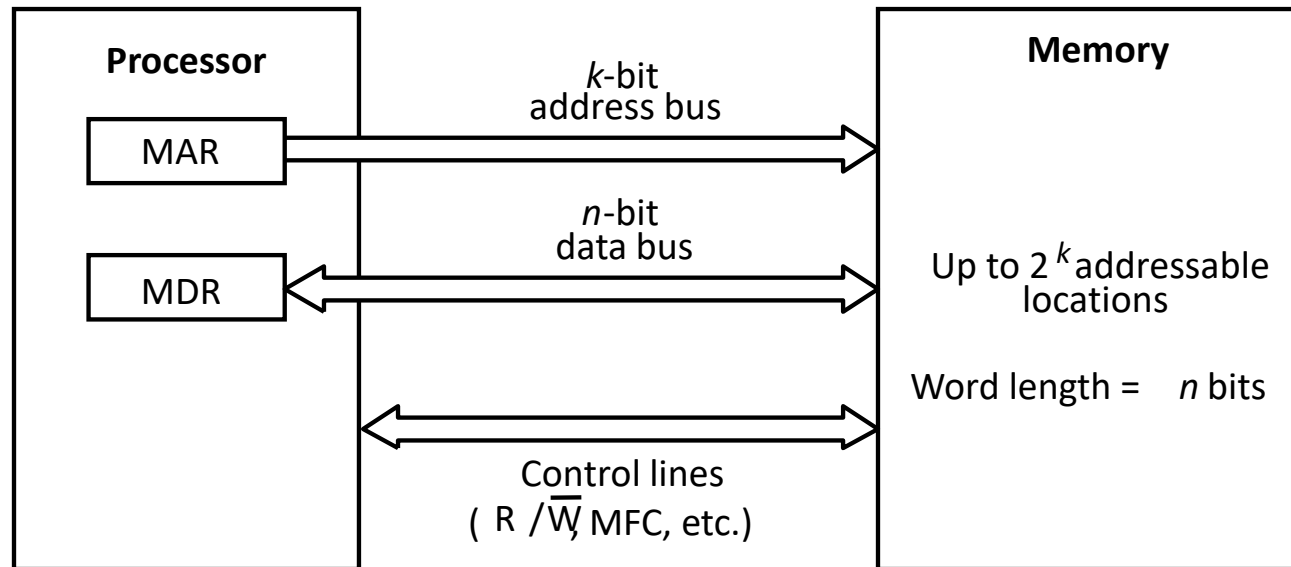# UNIT-4:
# The Memory System

# Memory Basic Concepts

- Maximum size of the Main Memory

- byte-addressable

- CPU-Main Memory Connection

# Memory Basic Concepts(Contd.,)

- Measures for the speed of a memory:

    - Memory access time.

    - Memory cycle time.

- An important design issue is to provide a computer system with as large and fast a memory as possible, within a given cost target.

- Several techniques to increase the effective size and speed of the memory:

    - Cache memory (to increase the effective speed).

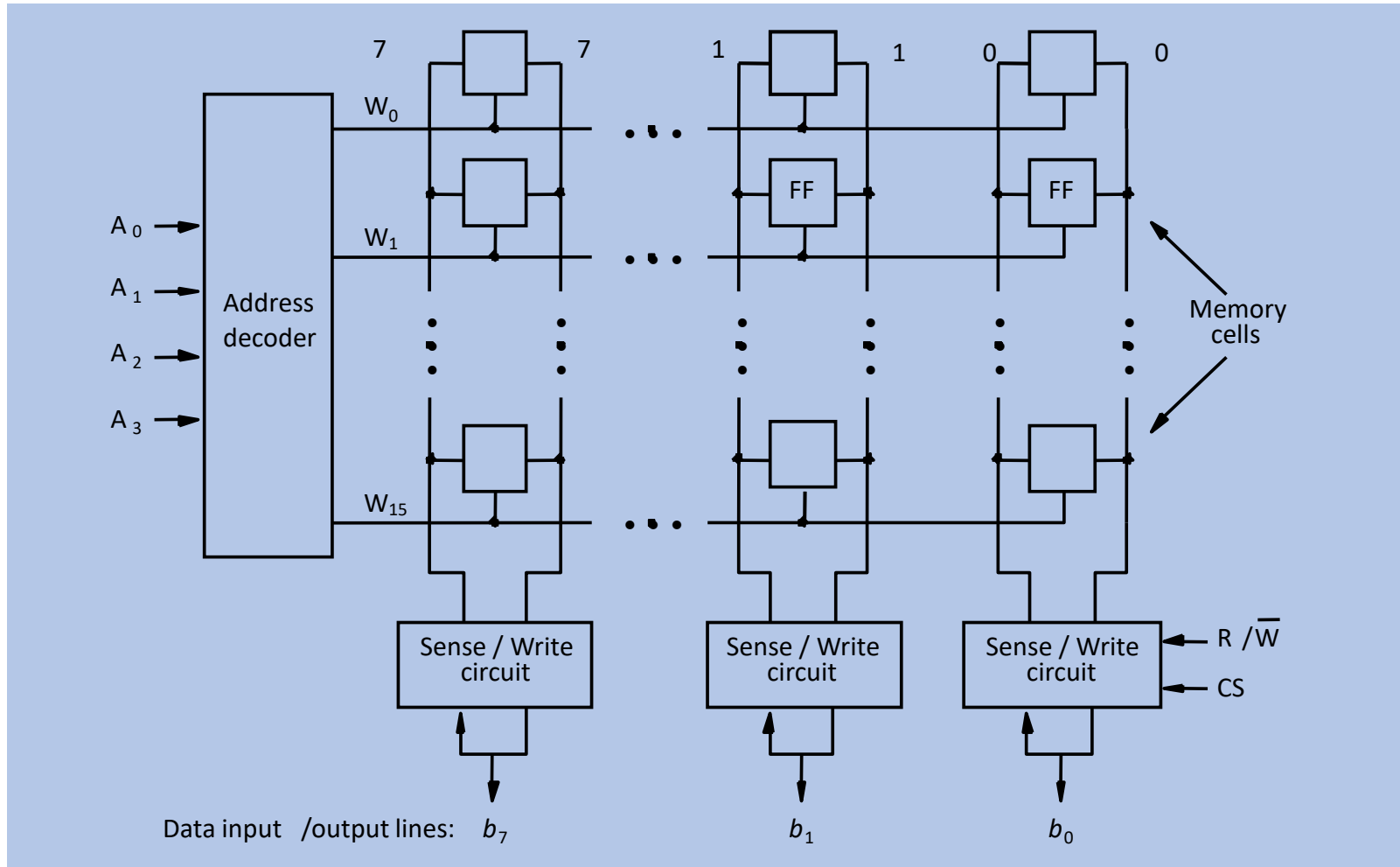    - Virtual memory (to increase the effective size).

# The Memory System

Semiconductor RAM memories

# Internal Organization of Memory Chips

- Each memory cell can hold one bit of information.

- Memory cells are organized in the form of an array.

- One row is one memory word.

- All cells of a row are connected to a common line, known as the "word line".

- Word line is connected to the address decoder.

- Sense/write circuits are connected to the data input/output lines of the memory chip.

# Internal Organization of Memory Chips (Contd.,)

# Asynchronous DRAMs

- ## Static RAMs (SRAMs):
  - Consist of circuits that are capable of retaining their state as long as the power is applied.
  - Volatile memories, because their contents are lost when power is interrupted.
  - Access times of static RAMs are in the range of few nanoseconds.
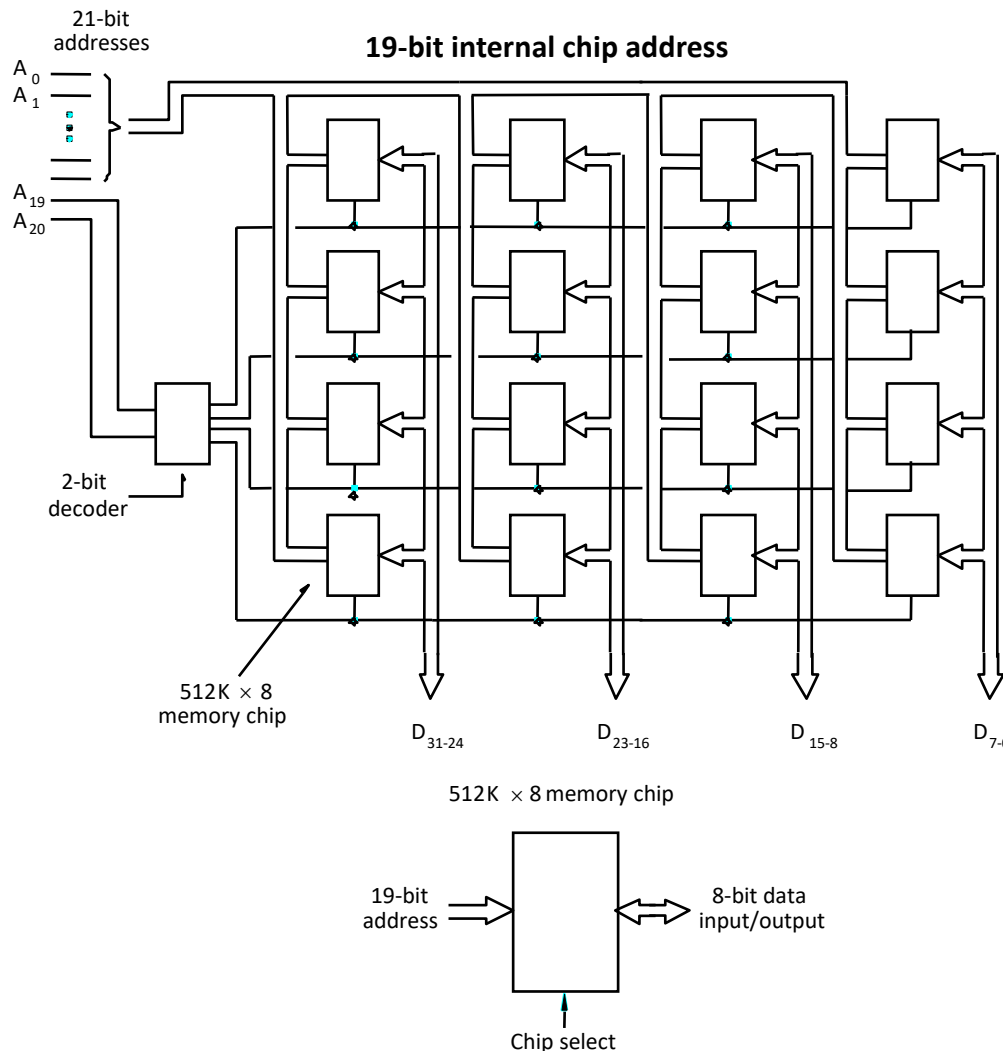  - However, the cost is usually high.

- ## Dynamic RAMs (DRAMs):
  - Do not retain their state indefinitely.
  - Contents must be periodically refreshed.
  - Contents may be refreshed while accessing them for reading.

# Latency, Bandwidth

- Memory latency is the time it takes to transfer a word of data to or from memory

- Memory bandwidth is the number of bits or bytes that can be transferred in one second.

# Static Memories



21-bit addresses

19-bit internal chip address

$A_0$
$A_1$
$A_{19}$
$A_{20}$

2-bit decoder

512K $\times$ 8 memory chip

$D_{31-24}$       $D_{23-16}$       $D_{15-8}$       $D_{7-0}$

512K $\times$ 8 memory chip

19-bit address

8-bit data input/output

Chip select

Implement a memory unit of 2M words of 32 bits each.
Use 512Kx8 static memory chips.
Each column consists of 4 chips.
Each chip implements one byte position.
A chip is selected by setting its chip select control line to 1.
Selected chip places its data on the data output line, outputs of other chips are in high impedance state.
21 bits to address a 32-bit word.
High order 2 bits are needed to select the row, by activating the four Chip Select signals.
19 bits are used to access specific byte locations inside the selected chip.
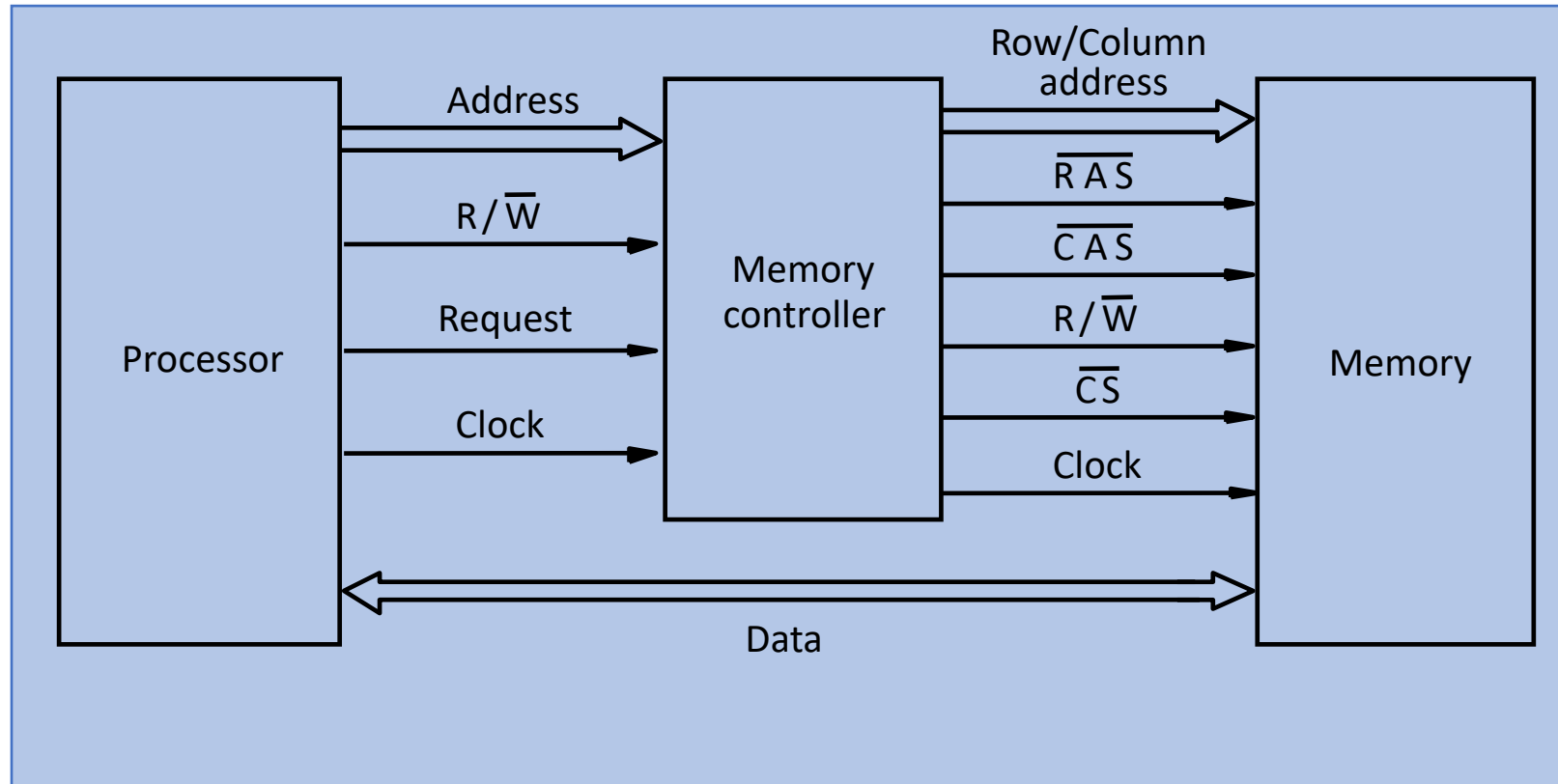
# Dynamic Memories

- Large dynamic memory systems can be implemented using DRAM chips in a similar way to static memory systems.

- Placing large memory systems directly on the motherboard will occupy a large amount of space.

- Packaging considerations have led to the development of larger memory units known as SIMMs (Single In-line Memory Modules) and DIMMs (Dual In-line Memory Modules).

- Memory modules are an assembly of memory chips on a small board that plugs vertically onto a single socket on the motherboard.

# Memory Controller

- Recall that in a dynamic memory chip, to reduce the number of pins, multiplexed addresses are used.
- Address is divided into two parts:
  - High-order address bits select a row in the array.
    - They are provided first, and latched using RAS signal.
  - Low-order address bits select a column in the row.
    - They are provided later, and latched using CAS signal.
- However, a processor issues all address bits at the same time.
- In order to achieve the multiplexing, memory controller circuit is inserted between the processor and memory.

# Memory Controller (contd..)

# The Memory System

- **SRAM and SDRAM chips are volatile:**
  - Lose the contents when the power is turned off.
- **Many applications need memory devices to retain contents after the power is turned off.**
  - For example, computer is turned on, the operating system must be loaded from the disk into the memory.
  - Store instructions which would load the OS from the disk.
  - Need to store these instructions so that they will not be lost after the power is turned off.
  - We need to store the instructions into a non-volatile memory.
- **Non-volatile memory is read in the same manner as volatile memory.**
  - Separate writing process is needed to place information in this memory.
  - Normal operation involves only reading of data, this type of memory is called Read-Only memory (ROM).

## Read-Only Memory:

- Data are written into a ROM when it is manufactured.

## Programmable Read-Only Memory (PROM):

- Allow the data to be loaded by a user.

- Process of inserting the data is irreversible.

- Storing information specific to a user in a ROM is expensive.

- Providing programming capability to a user may be better.

## Erasable Programmable Read-Only Memory (EPROM):

- Stored data to be erased and new data to be loaded.

- Flexibility, useful during the development phase of digital systems.

- Erasable, reprogrammable ROM.

- Erasure requires exposing the ROM to UV light.

- **Electrically Erasable Programmable Read-Only Memory (EEPROM):**

  - To erase the contents of EPROMs, they have to be exposed to ultraviolet light.

  - Physically removed from the circuit.

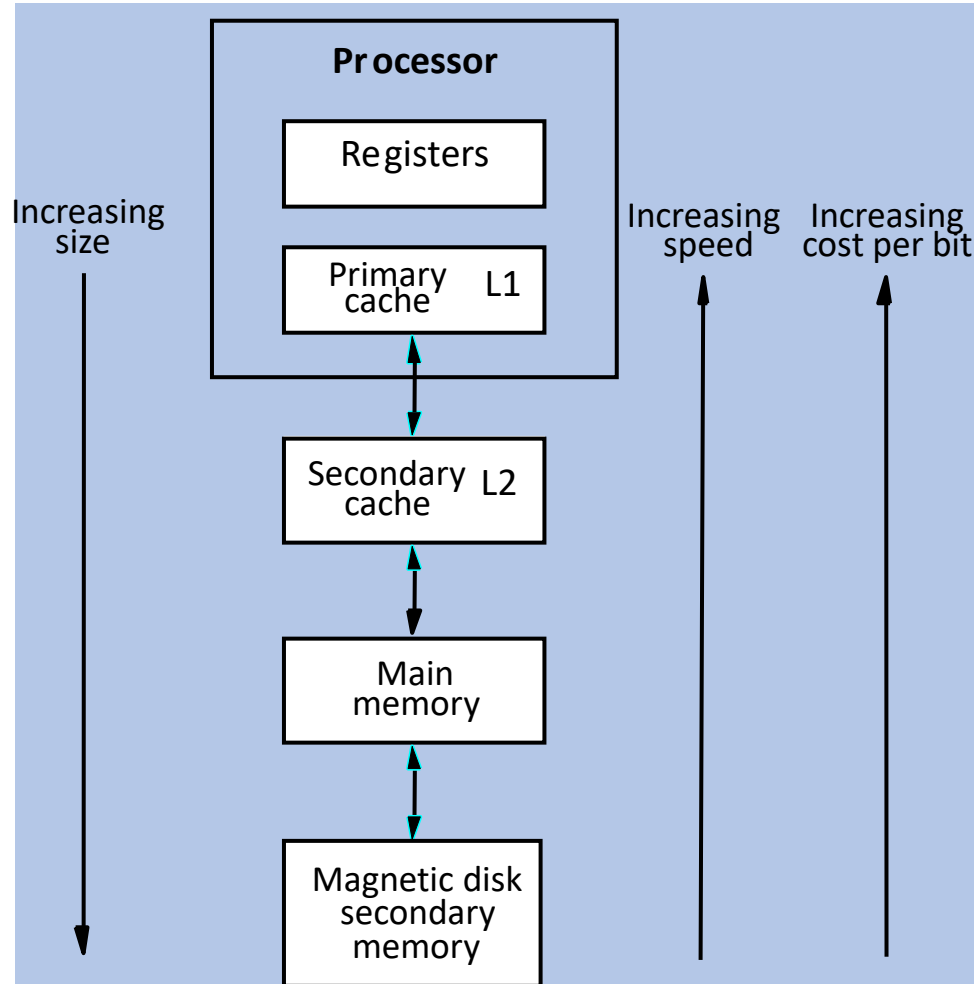  - EEPROMs the contents can be stored and erased electrically.

- **Flash memory:**

  - Has similar approach to EEPROM.

  - Read the contents of a single cell, but write the contents of an entire block of cells.

# Speed, Size, and Cost

- A big challenge in the design of a computer system is to provide a sufficiently large memory, with a reasonable speed at an affordable cost.

- Static RAM:
  - Very fast, but expensive, because a basic SRAM cell has a complex circuit making it impossible to pack a large number of cells onto a single chip.

- Dynamic RAM:
  - Simpler basic cell circuit, hence are much less expensive, but significantly slower than SRAMs.

- Magnetic disks:
  - Storage provided by DRAMs is higher than SRAMs, but is still less than what is necessary.
  - Secondary storage such as magnetic disks provide a large amount of storage, but is much slower than DRAMs.

# Memory Hierarchy



**Processor**

Registers

Primary cache   L1

Secondary cache   L2

Main memory

Magnetic disk secondary memory

Increasing size

Increasing speed

Increasing cost per bit

- *Fastest access is to the data held in processor registers. Registers are at the top of the memory hierarchy.*
- *Relatively small amount of memory that can be implemented on the processor chip. This is processor cache.*
- *Two levels of cache. Level 1 (L1) cache is on the processor chip. Level 2 (L2) cache is in between main memory and processor.*
- *Next level is main memory, implemented as SIMMs(Single in line memory module). Much larger, but much slower than cache memory.*
- *Next level is magnetic disks. Huge amount of inexpensive storage.*
- *Speed of memory access is critical, the idea is to bring instructions and data that will be used in the near future as close to the processor as possible.*
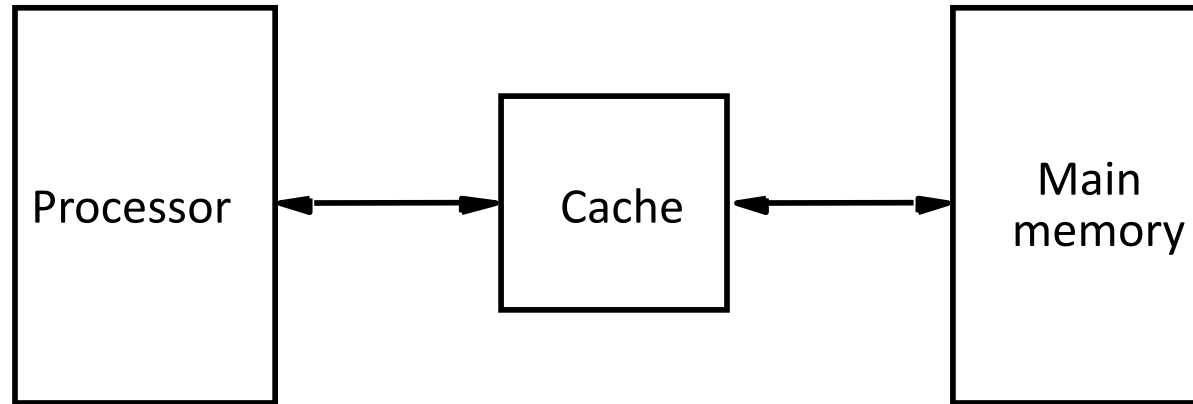
# The Memory System

Cache Memories

# Cache Memory

- Processor is much faster than the main memory.
- Speed of the main memory cannot be increased beyond a certain point.
- Cache memory is an architectural arrangement which makes the main memory appear faster to the processor than it really is.
- Cache memory is based on the property of computer programs known as <u>"locality of reference"</u>.

# Locality of Reference

■ Analysis of programs indicates that many instructions in localized areas of a program are executed repeatedly during some period of time, while the others are accessed relatively less frequently.

■ Temporal locality of reference:

- Recently executed instruction is likely to be executed again very soon.

■ Spatial locality of reference:

- Instructions with addresses close to a recently instruction are likely to be executed soon.

# Cache memories



- *Processor issues a Read request, a block of words is transferred from the main memory to the cache, one word at a time.*
- *Subsequent references to the data in this block of words are found in the cache.*
- *At any given time, only some blocks in the main memory are held in the cache. Which blocks in the main memory are in the cache is determined by a "mapping function".*
- *When the cache is full, and a block of words needs to be transferred from the main memory, some block of words in the cache must be replaced. This is determined by a "replacement algorithm".*

# Cache Hit

- *Existence of a cache is transparent to the processor. The processor issues Read and Write requests in the same manner.*

- *If the data is in the cache it is called a Read or Write hit.*

- *Read hit:*
    - *The data is obtained from the cache.*

- *Write hit:*
    - *Cache has a replica of the contents of the main memory.*
    - *Contents of the cache and the main memory may be updated simultaneously.    This is the write-through protocol.*
    - *Update the contents of the cache, and mark it as updated by setting a bit known as the dirty bit or modified bit. The contents of the main memory are updated when this block is replaced. This is write-back or copy-back protocol.*

# Cache Miss

- *If the data is not present in the cache, then a <u>Read miss or Write miss</u> occurs.*

- *Read miss:*
  - *Block of words containing this requested word is transferred from the memory.*
  - *After the block is transferred, the desired word is forwarded to the processor.*
  - *The desired word may also be forwarded to the processor as soon as it is transferred without waiting for the entire block to be transferred. This is called <u>load-through or early-restart.</u>*

- *Write-miss:*
  - *Write-through protocol is used, then the contents of the main memory are updated directly.*
  - *If write-back protocol is used, the block containing the addressed word is first brought into the cache. The desired word is overwritten with new information.*
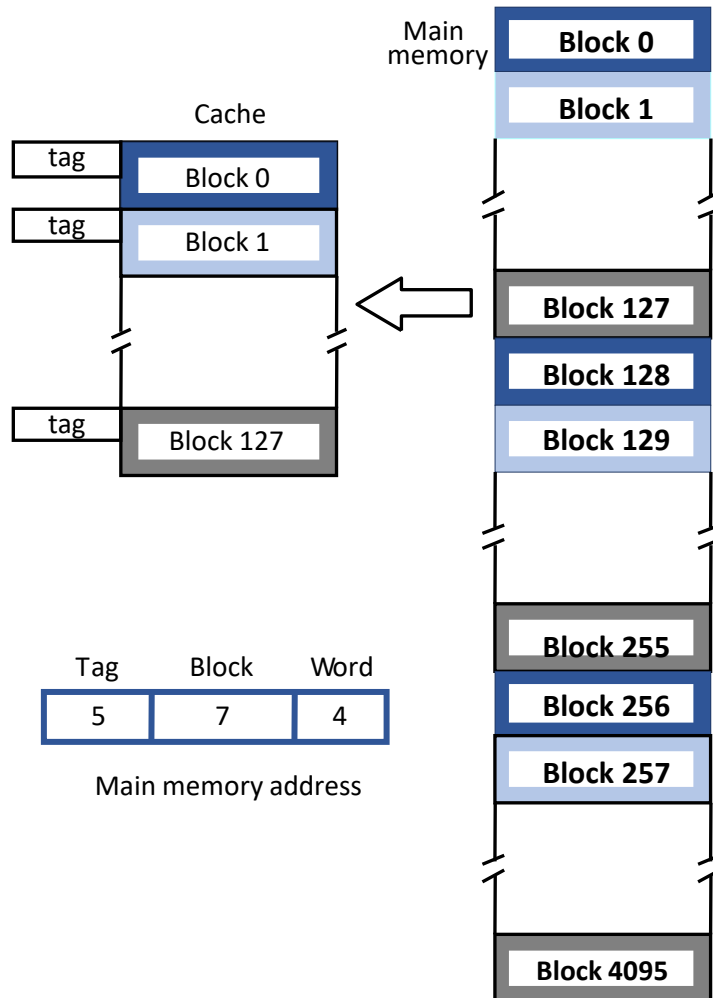
# Cache Coherence Problem

- *A bit called as "valid bit" is provided for each block.*
- *If the block contains valid data, then the bit is set to 1, else it is 0.*
- *Valid bits are set to 0, when the power is just turned on.*
- *When a block is loaded into the cache for the first time, the valid bit is set to 1.*

- *Data transfers between main memory and disk occur directly bypassing the cache.*
- *When the data on a disk changes, the main memory block is also updated.*
- *However, if the data is also resident in the cache, then the valid bit is set to 0.*

- *What happens if the data in the disk and main memory changes and the write-back protocol is being used?*
- *In this case, the data in the cache may also have changed and is indicated by the dirty bit.*
- *The copies of the data in the cache, and the main memory are different. This is called the cache coherence problem.*
- *One option is to force a write-back before the main memory is updated from the disk.*
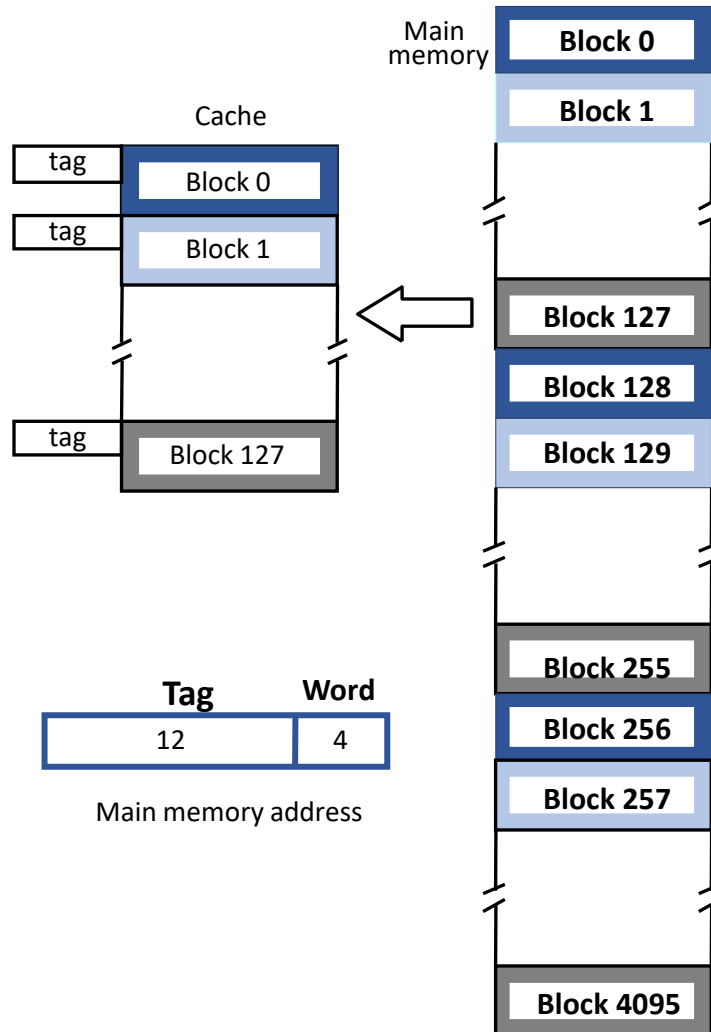
# Mapping Functions

- **Mapping functions** determine how memory blocks are placed in the cache.

- A simple processor example:

  - Cache consisting of 128 blocks of 16 words each.

  - Total size of cache is 2048 (2K) words.

  - Main memory is addressable by a 16-bit address.

  - Main memory has 64K words.

  - Main memory has 4K blocks of 16 words each.

- Three mapping functions:

  - Direct mapping

  - Associative mapping

  - Set-associative mapping.

# Direct Mapping



Main memory

Cache

| tag | Block 0 |
| tag | Block 1 |
| tag | Block 127 |

Block 0
Block 1
Block 127
Block 128
Block 129
Block 255
Block 256
Block 257
Block 4095

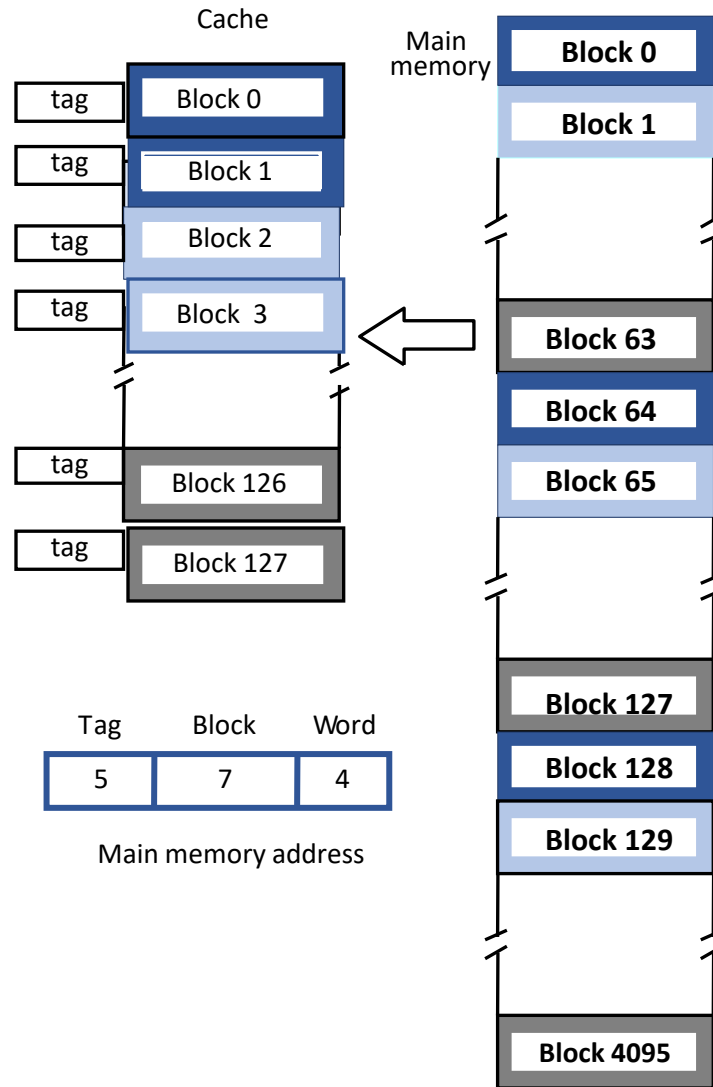| Tag | Block | Word |
|-----|-------|------|
| 5 | 7 | 4 |

Main memory address

- *Block j of the main memory maps to j modulo 128 of the cache. 0 maps to 0, 129 maps to 1.*
- *More than one memory block is mapped onto the same position in the cache.*
- *May lead to contention for cache blocks even if the cache is not full.*
- *Resolve the contention by allowing new block to replace the old block, leading to a trivial replacement algorithm.*
- *Memory address is divided into three fields:*
  - *Low order 4 bits determine one of the 16 words in a block.*
  - *When a new block is brought into the cache, the the next 7 bits determine which cache block this new block is placed in.*
  - *High order 5 bits determine which of the possible 32 blocks is currently present in the cache. These are tag bits.*
- *Simple to implement but not very flexible.*

# Associative Mapping



Main memory

| Block 0 |
|---|
| Block 1 |
| |
| Block 127 |
| Block 128 |
| Block 129 |
| |
| Block 255 |
| Block 256 |
| Block 257 |
| |
| Block 4095 |

Cache

| tag | | Block 0 |
|---|---|---|
| tag | | Block 1 |
| | | |
| tag | | Block 127 |

| Tag | Word |
|---|---|
| 12 | 4 |

Main memory address

- *Main memory block can be placed into any cache position.*
- *Memory address is divided into two fields:*
  - *Low order 4 bits identify the word within a block.*
  - *High order 12 bits or tag bits identify a memory block when it is resident in the cache.*
- *Flexible, and uses cache space efficiently.*
- *Replacement algorithms can be used to replace an existing block in the cache when the cache is full.*
- *Cost is higher than direct-mapped cache because of the need to search all 128 patterns to determine whether a given block is in the cache.*

# Set-Associative mapping



Cache

Main memory

| tag | Block 0 |
| tag | Block 1 |
| tag | Block 2 |
| tag | Block 3 |
| tag | Block 126 |
| tag | Block 127 |

| Tag | Block | Word |
|-----|-------|------|
| 5 | 7 | 4 |

Main memory address

Block 0
Block 1
Block 63
Block 64
Block 65
Block 127
Block 128
Block 129
Block 4095

*Blocks of cache are grouped into sets.*
*Mapping function allows a block of the main memory to reside in any block of a specific set.*
*Divide the cache into 64 sets, with two blocks per set.*
*Memory block 0, 64, 128 etc. map to block 0, and they can occupy either of the two positions.*
*Memory address is divided into three fields:*
*    - 6 bit field determines the set number.*
*    - High order 6 bit fields are compared to the tag fields of the two blocks in a set.*
*Set-associative mapping combination of direct and associative mapping.*
*Number of blocks per set is a design parameter.*
*    - One extreme is to have all the blocks in one set, requiring no set bits (fully associative mapping).*
*    - Other extreme is to have one block per set, is the same as* direct mapping.

# Replacement Algorithms

• For direct mapping where there is only one possible line for a block of memory, no replacement algorithm is needed.

• For associative and set associative mapping, however, an algorithm is needed.

•For maximum speed, this algorithm is implemented in the hardware. Four of the most common algorithms are:

1. Least Recently Used:- This replaces the candidate line in cache memory that has been there the longest with no reference to it.
2. First In First Out:- This replaces the candidate line in the cache that has been there the longest.
3. Least Frequently Used:- This replaces the candidate line in the cache that has had the fewest references.
4. Random Replacement:- This algorithm randomly chooses a line to be replaced from among the candidate lines. This yields only slightly inferior performance than other algorithms.

# FIFO (First In First Out)

- Pages in main memory are kept in a list

- First in first out is very easy to implement

- The FIFO algorithm select the page for replacement that has been in memory the longest time
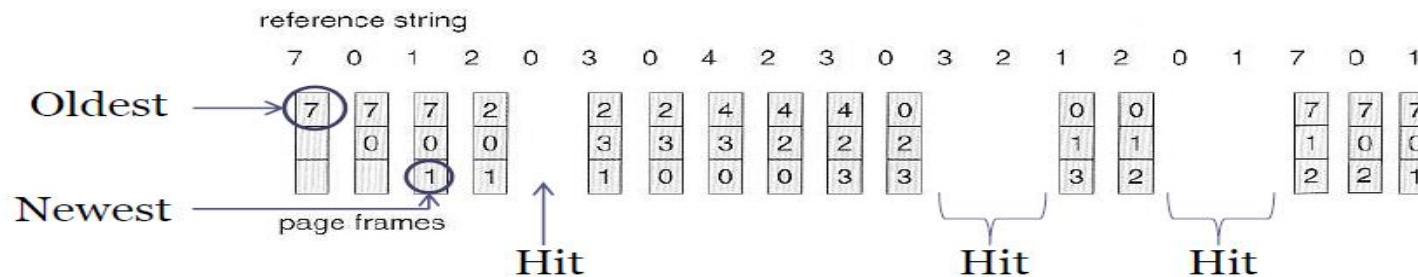


Fig: FIFO example

# FIFO (First In First Out)

- Advantages:
    - FIFO is easy to understand.
    - It is very easy to implement.

- Disadvantages:
    - The oldest block in memory may be often used.

# LRU (Least Recently Used)

- The least recently used page replacement algorithm keeps track page uses over a short period of time.
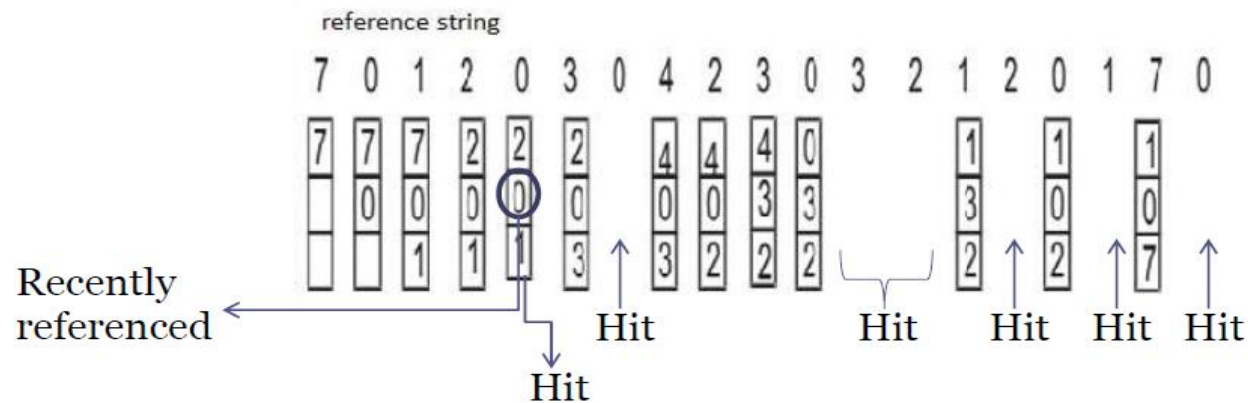


Fig: LRU example

# LRU (Least Recently Used)

- Advantages:
  - LRU page replacement algorithm is quiet efficient.
- Disadvantages:
  - Implementation difficult. This algorithm requires keeping track of what was used when, which is expensive if one wants to make sure
  - the algorithm always discards the least recently used item.

# Comparison of Clock with FIFO and LRU

# LFU (Least Frequently Used)

- The Least-Frequently-Used (LFU) Replacement technique replaces the least-frequently block in use when an eviction must take place.

- Software counter associated with each block, initially zero is required in this algorithm.

- The operating system checks all the blocks in the cache at each clock interrupt.

- The R bit, which is '0' or '1', is added to the counter for each block. Consequently, the counters are an effort to keep track of the frequency of referencing each block.

- When a block must be replaced, the block that has the lowest counter is selected for the replacement.

# LFU (Least Frequently Used)

Reference string:

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

| 7 | 7 | 7 | 2 | 2 | 2 | 2 | 4 | 4 | 4 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 3 | 3 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 |   |
|   |   | 1 | 1 | 1 | 3 | 3 | 3 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 7 | 7 | 7 |
| F | F | F | F |   | F |   | F | F | F | F |   |   | F |   | F |   | F |   |   |

Number of page faults = 12.
Number of page hits= 8.

# LFU (Least Frequently Used)

- Advantages:
  - Frequently used block will stay longer than (fifo)
- Disadvantages:
  - Older blocks are less likely to be removed , even if they are on longer frequently used because this algorithm never forgets anything.
  - Newer blocks are more likely to be replaced even if they are frequently used.
  - Captures only frequency factor.

# Random Replacement

- When we need to evict a page, choose one randomly

- Advantage:

  - Extremely simple

- Disadvantages:

  - Can easily make "bad" choices by swapping out pages right before

    they are needed.