**Department of Computer Science and Engineering**

COURSE CODE – TITLE: 1151CS110 – COMPUTER
ORGANIZATION AND ARCHITECTURE

Course Instructor
Dr. M. Rajeev Kumar
Associate Professor (CSE)

# UNIT V- Memory and I/O Systems

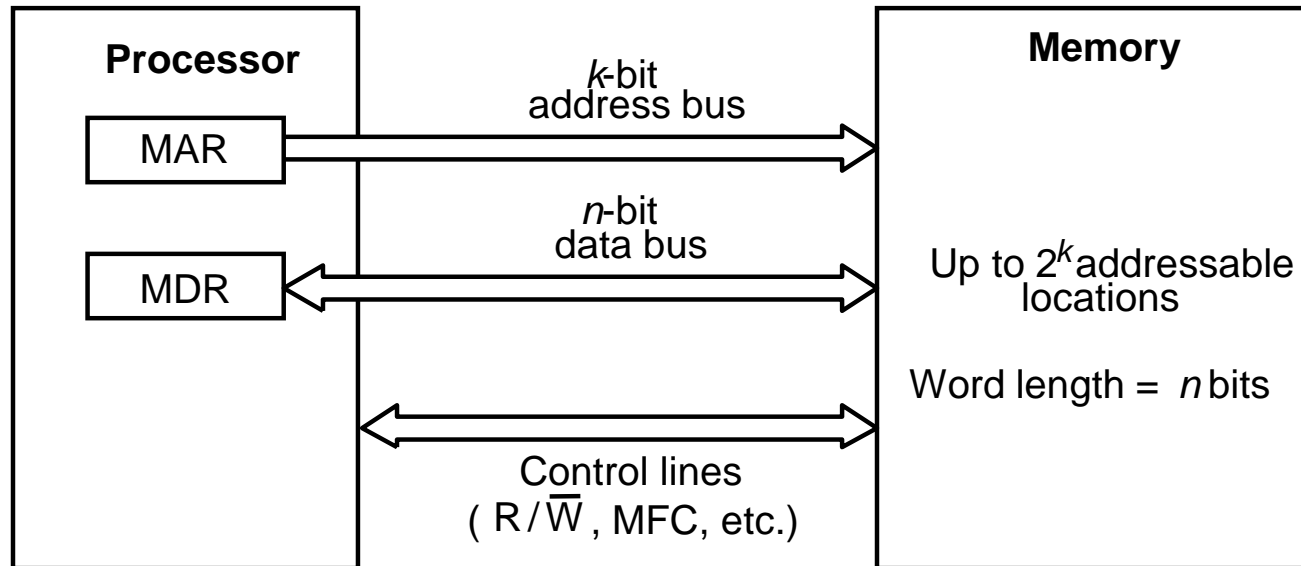| CO Nos. | Course Outcome(s) | Level of learning domain (Based on revised Bloom's) |
|---|---|---|
| CO5 | Identify performance issues in processor and memory design of a digital computer. | K3 |

# The Memory System

# Some basic concepts

- Maximum size of the Main Memory

- byte-addressable (little-endian & big-endian)

- CPU-Main Memory Connection



**Processor**

MAR

MDR

$k$-bit
address bus

$n$-bit
data bus

Control lines
( $R / \overline{W}$ , MFC, etc.)

**Memory**

Up to $2^k$ addressable
locations

Word length = $n$ bits

# Some basic concepts *(Contd.)*

- Measures for the speed of a memory:

  - memory access time.

  - memory cycle time.

- RAM

- 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).
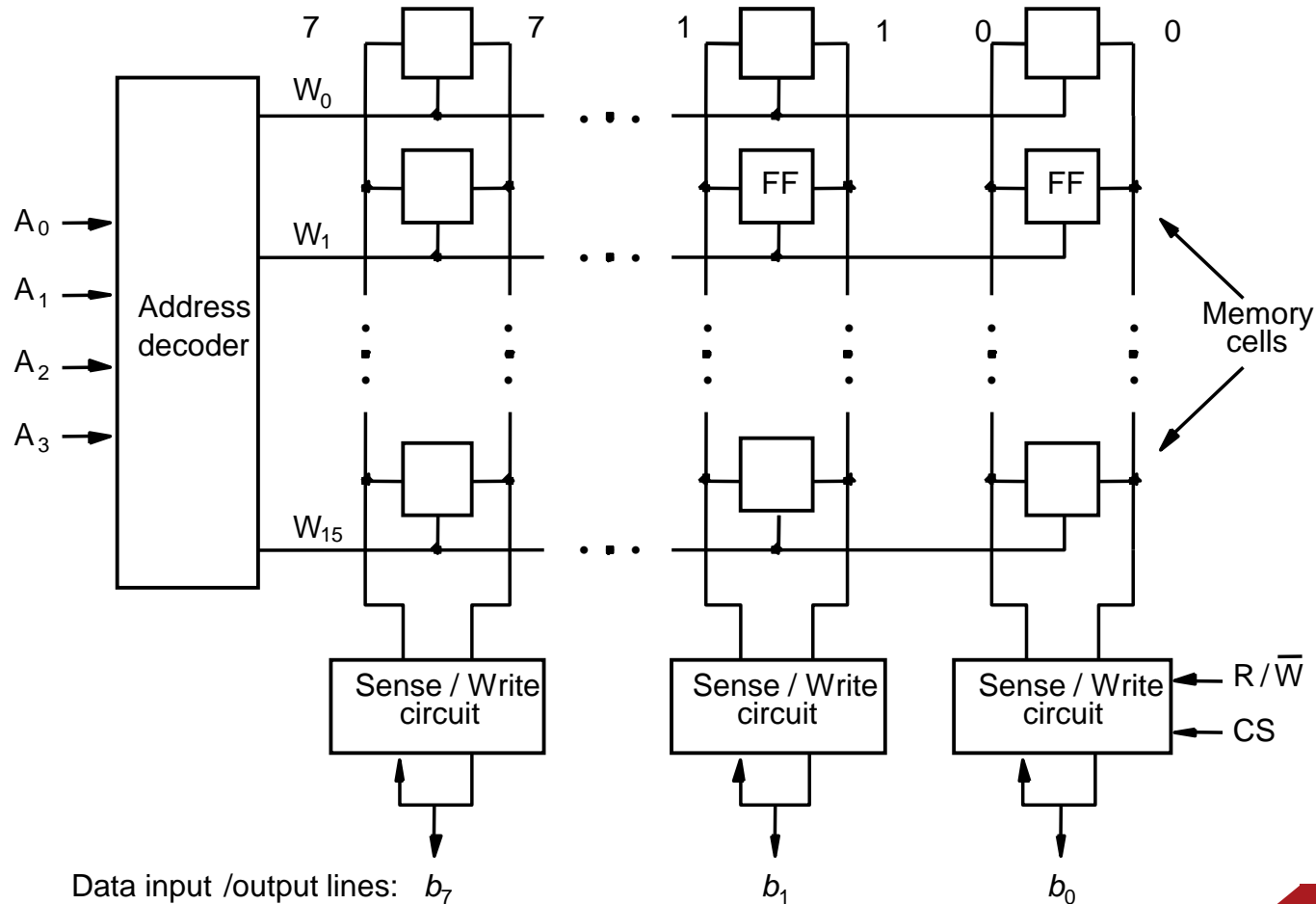
# Semiconductor RAM memories
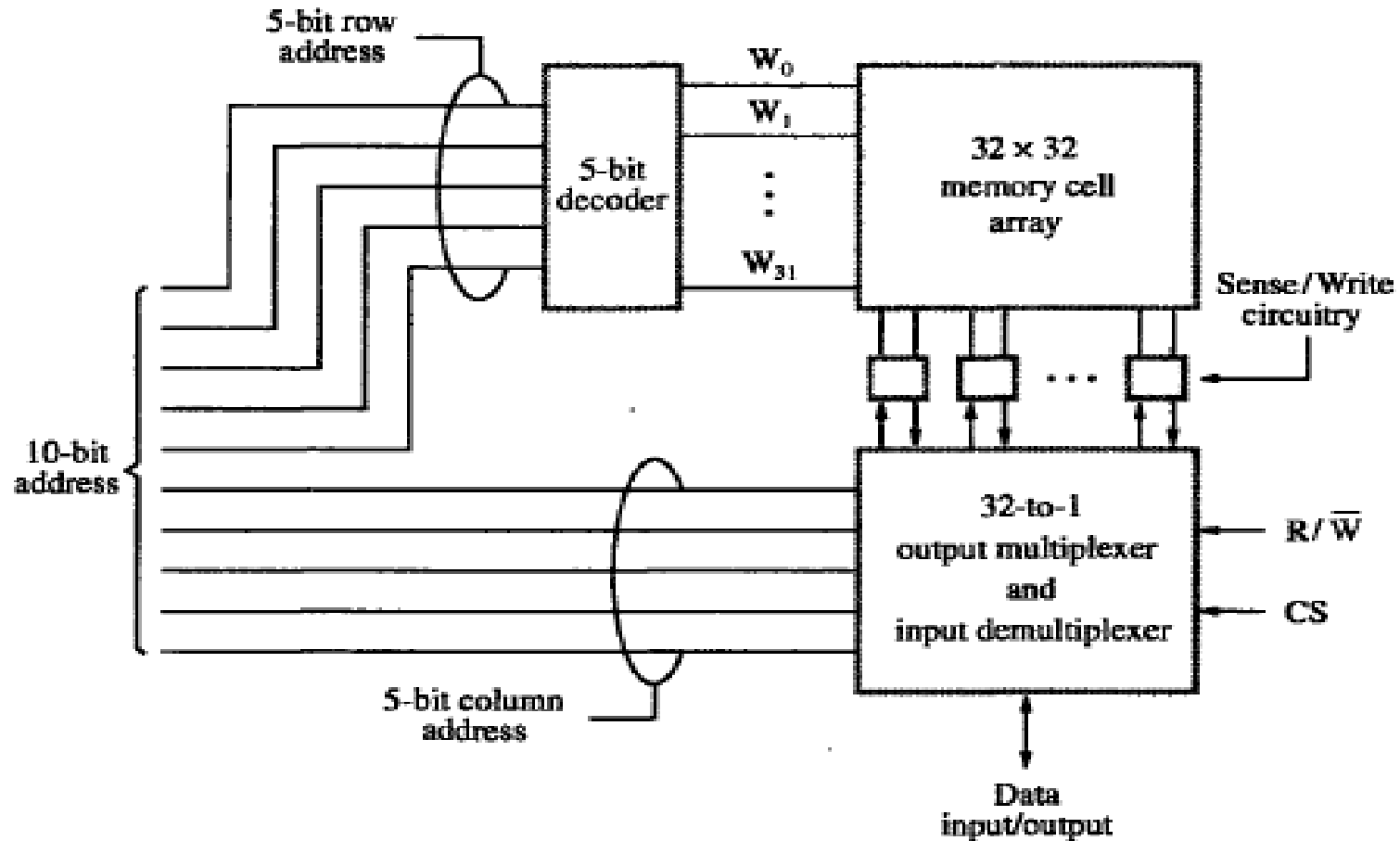
# Internal organization of memory chips

- 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.

# SRAM Cell

- Two transistor inverters are cross connected to implement a basic flip-flop.

- The cell is connected to one word line and two bits lines by transistors T1 and T2

- When word line is at ground level, the transistors are turned off and the latch retains its state

- Read operation: In order to read state of SRAM cell, the word line is activated to close switches T1 and T2. Sense/Write circuits at the bottom monitor the state of b and b'

# An example of CMOS memory cell

# 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.

# A single-transistor dynamic memory cell

# Asynchronous DRAMs



- *Each row can store 512 bytes. 12 bits to select a row, and 9 bits to select a group in a row. Total of 21 bits.*
- *First apply the row address, RAS signal latches the row address. Then apply the column address, CAS signal latches the address.*
- *Timing of the memory unit is controlled by a specialized unit which generates RAS and CAS.*
- *This is asynchronous DRAM*

# Fast Page Mode

- Suppose if we want to access the consecutive bytes in the selected row.
- This can be done without having to reselect the row.
  - Add a latch at the output of the sense circuits in each row.
  - All the latches are loaded when the row is selected.
  - Different column addresses can be applied to select and place different bytes on the data lines.
- Consecutive sequence of column addresses can be applied under the control signal CAS, without reselecting the row.
  - Allows a block of data to be transferred at a much faster rate than random accesses.
  - A small collection/group of bytes is usually referred to as a block.
- This transfer capability is referred to as the fast page mode feature.

# Synchronous DRAMs



- *Operation is directly synchronized with processor clock signal.*
- *The outputs of the sense circuits are connected to a latch.*
- *During a Read operation, the contents of the cells in a row are loaded onto the latches.*
- *During a refresh operation, the contents of the cells are refreshed without changing the contents of the latches.*
- *Data held in the latches correspond to the selected columns are transferred to the output.*
- *For a burst mode of operation, successive columns are selected using column address counter and clock. CAS signal need not be generated externally. A new data is placed during raising edge of the clock*
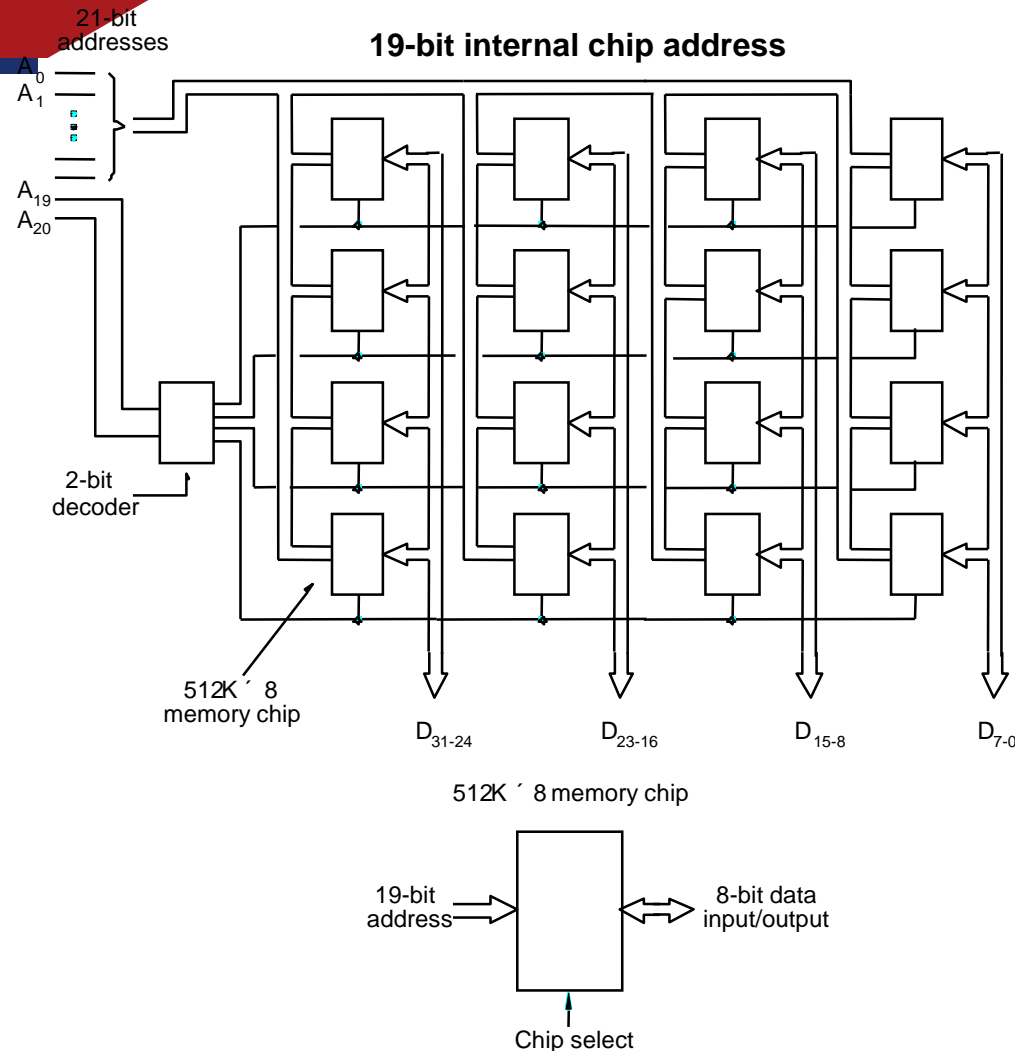
# Latency, Bandwidth, and DDRSDRAMs

- 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.

- DDRSDRAMs

  - Cell array is organized in two banks

# Static memories



A$_0$
A$_1$

A$_{19}$
A$_{20}$

2-bit decoder

512K ´ 8 memory chip

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

512K ´ 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 512x8 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.

    - Also, this arrangement is inflexible since the memory system cannot be expanded easily.

- 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.

    - Occupy less space on the motherboard.

    - Allows for easy expansion by replacement.

# 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.)*

# Read Only Memories (ROMs)

# Read Only Memories (ROMs)

- 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 Memories *(Contd.)*

- **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.

# Read Only Memories *(Contd.)*

- 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.
  - Flash devices have greater density.
    - Higher capacity and low storage cost per bit.
  - Power consumption of flash memory is very low, making it attractive for use in equipment that is battery-driven.
  - Single flash chips are not sufficiently large, so larger memory modules are implemented using flash cards and flash drives.

# 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



- *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. Much larger, but much slower than cache memory.*

- *Next level is magnetic disks. Huge amount of inexepensive 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.*

# Cache Memories

# Cache Memories

- Processor is much faster than the main memory.

  - As a result, the processor has to spend much of its time waiting while instructions and data are being fetched from the main memory.

  - Major obstacle towards achieving good performance.

- 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 "locality of reference".

# 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.

  - These instructions may be the ones in a loop, nested loop or few procedures calling each other repeatedly.

  - This is called "locality of reference".

- 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

```
┌─────────────┐        ┌─────────┐        ┌─────────────┐
│             │        │         │        │             │
│             │        │         │        │    Main     │
│  Processor  │ ◄────► │  Cache  │ ◄────► │   memory    │
│             │        │         │        │             │
│             │        │         │        │             │
└─────────────┘        └─────────┘        └─────────────┘
```

- 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 <u>Read or Write hit</u>.

- 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 <u>write-through</u> protocol.

    - Update the contents of the cache, and mark it as updated by setting a bit known as the <u>dirty bit or modified</u> bit. The contents of the main memory are updated when this block is replaced. This is <u>write-back or copy-back</u> 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 <u>cache coherence problem</u>.
- One option is to force a write-back before the main memory is updated from the disk.
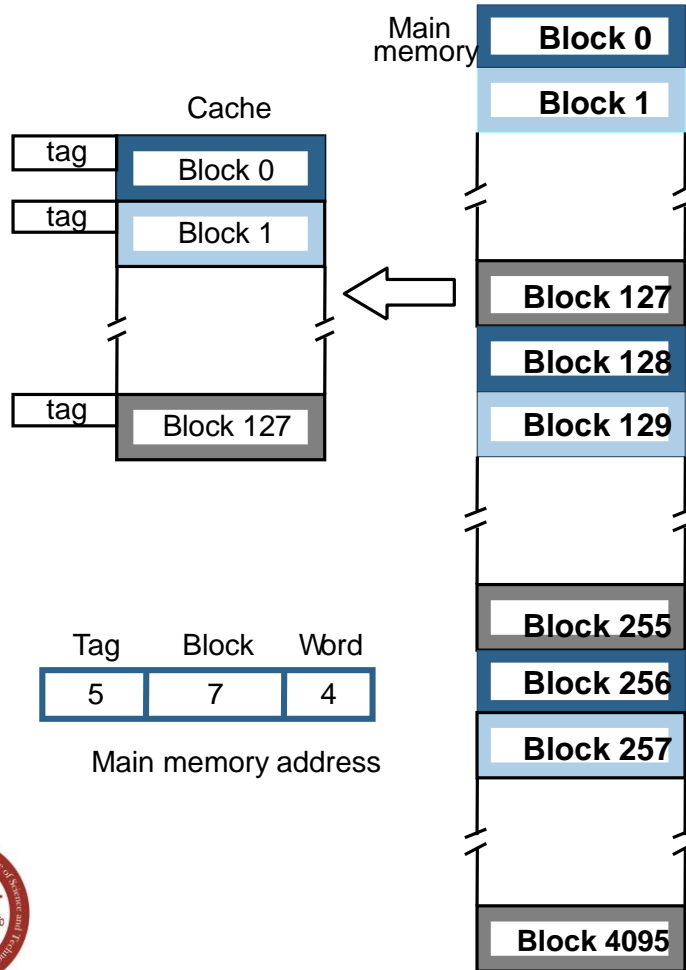
# Mapping functions

- 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
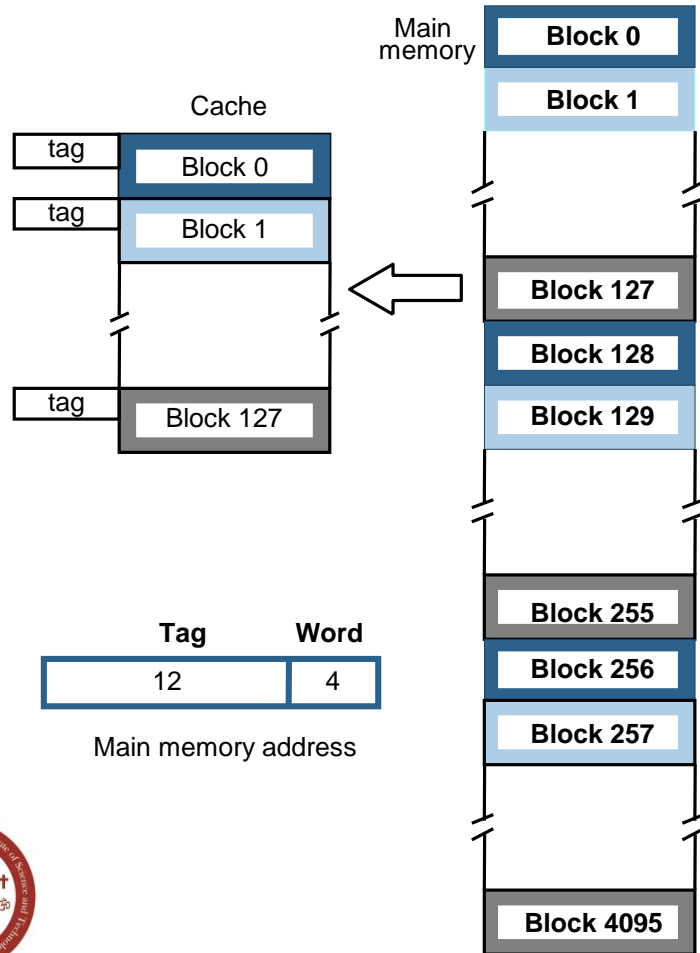
  - Set-associative mapping.

# Direct Mapping

**Cache**

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

**Main memory**

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.
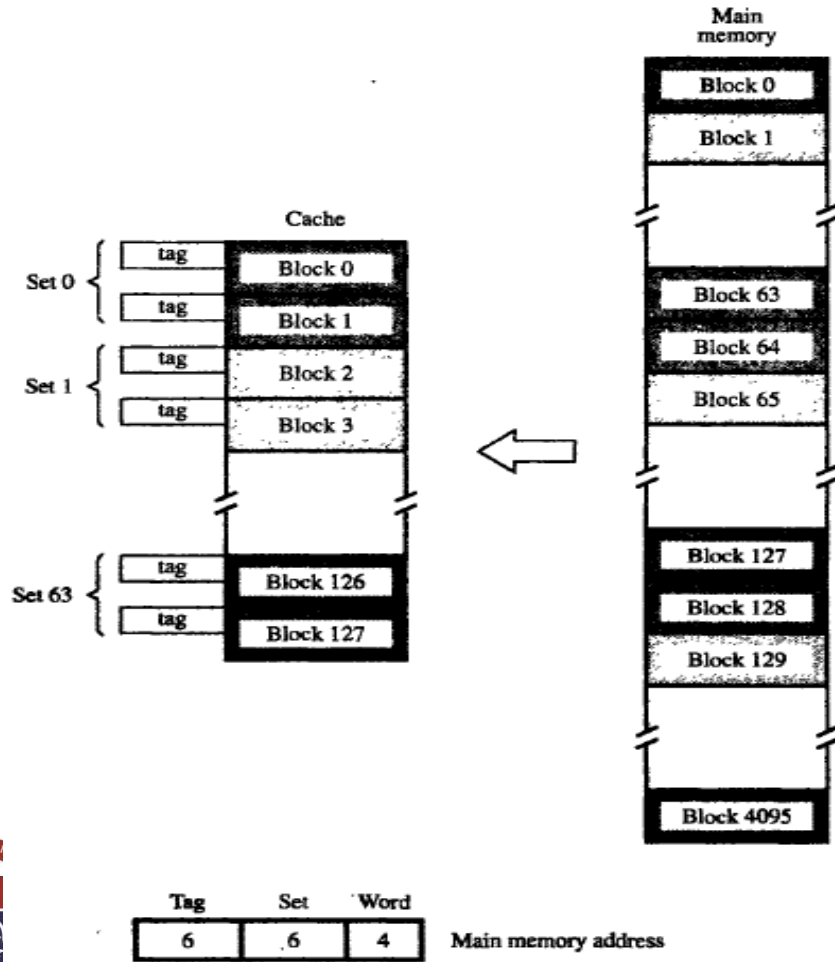- Simple to implement but not very flexible.

# Associative Mapping



Cache

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

Main memory

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

| 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

Main memory

Block 0
Block 1

Block 63
Block 64
Block 65

Block 127
Block 128
Block 129

Block 4095

Cache

| | | |
|---|---|---|
| tag | Block 0 | Set 0 |
| tag | Block 1 | |
| tag | Block 2 | Set 1 |
| tag | Block 3 | |
| tag | Block 126 | Set 63 |
| tag | Block 127 | |

- 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.

| Tag | Set | Word |
|---|---|---|
| 6 | 6 | 4 |

Main memory address

# Performance considerations

# Performance considerations

- A key design objective of a computer system is to achieve the best possible performance at the lowest possible cost.

  - Price/performance ratio is a common measure of success.

- Performance of a processor depends on:

  - How fast machine instructions can be brought into the processor for execution.

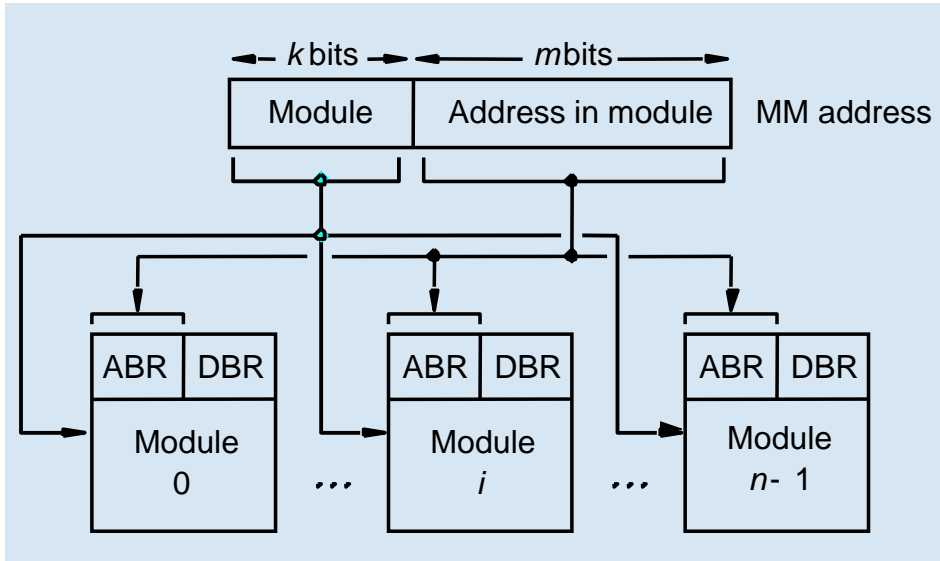  - How fast the instructions can be executed.

# Interleaving

- Divides the memory system into a number of memory modules. Each module has its own address buffer register (ABR) and data buffer register (DBR).

- Arranges addressing so that successive words in the address space are placed in different modules.

- When requests for memory access involve consecutive addresses, the access will be to different modules.

- Since parallel access to these modules is possible, the average rate of fetching words from the Main Memory can be increased.
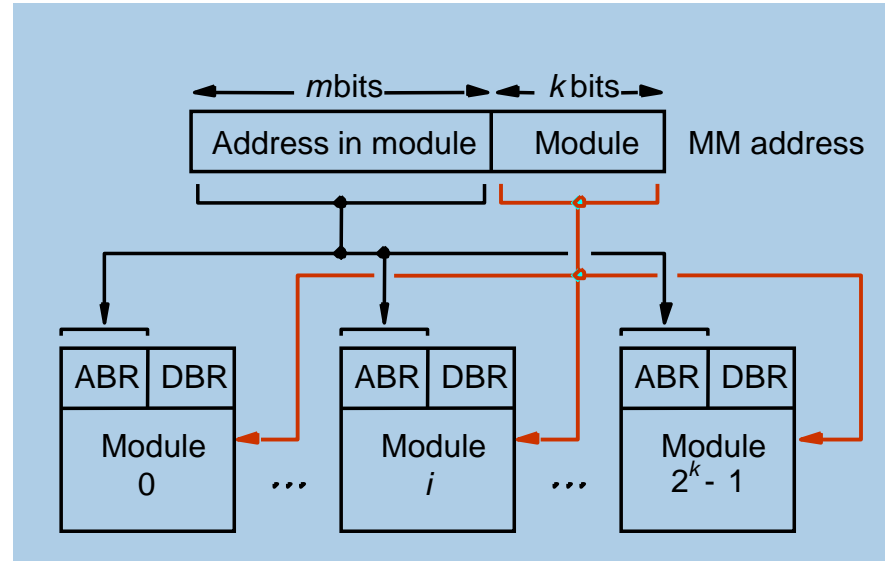
# Methods of address layouts



**Consecutive words in a module**

- Consecutive words are placed in a module.
- High-order k bits of a memory address determine the module.
- Low-order m bits of a memory address determine the word within a module.
- When a block of words is transferred from main memory to cache, only one module is busy at a time.

**Consecutive words in a consecutive module**

- Consecutive words are located in consecutive modules.
- Consecutive addresses can be located in consecutive modules.
- While transferring a block of data, several memory modules can be kept busy at the same time.

# Hit Rate and Miss Penalty

- Hit rate

- Miss penalty

- Hit rate can be improved by increasing block size, while keeping cache size constant

- Block sizes that are neither very small nor very large give best results.

- Miss penalty can be reduced if load-through approach is used when loading new blocks into cache.

# Caches on the processor chip

- In high performance processors 2 levels of caches are normally used.

- Avg access time in a system with 2 levels of caches is

    $T_{ave}$ = h1C1+(1-h1)h2C2+(1-h1)(1-h2)M

    Where
    > h1 is the hit rate in the L1 cache
    > h2 is the hit rate in the L2 cache
    > C1 is the time to access information in the L1 cache
    > C2 is the time to access information in the L2 cache
    > M is the time to access information in the main memory

# Other Performance Enhancements

Write buffer

- *Write-through:*
- *Each write operation involves writing to the main memory.*
- *If the processor has to wait for the write operation to be complete, it slows down the processor.*
- *Processor does not depend on the results of the write operation.*
- *Write buffer can be included for temporary storage of write requests.*
- *Processor places each write request into the buffer and continues execution.*
- *If a subsequent Read request references data which is still in the write buffer, then this data is referenced in the write buffer.*

- *Write-back:*
- *Block is written back to the main memory when it is replaced.*
- *If the processor waits for this write to complete, before reading the new block, it is slowed down.*
- *Fast write buffer can hold the block to be written, and the new*
  *block can be read first.*

## Prefetching

- New data are brought into the processor when they are first needed.

- Processor has to wait before the data transfer is complete.

- Prefetch the data into the cache before they are actually needed, or a before a Read miss occurs.

- Prefetching can be accomplished through software by including a special instruction in the machine language of the processor.

  - **Inclusion of prefetch instructions increases the length of the programs**.

- Prefetching can also be accomplished using hardware:

  - **Circuitry that attempts to discover patterns in memory references and then prefetches according to this pattern.**

## **Lockup-Free Cache**

- Prefetching scheme does not work if it stops other accesses to the cache until the prefetch is completed.

- A cache of this type is said to be "locked" while it services a miss.

- Cache structure which supports multiple outstanding misses is called a lockup free cache.

- Since only one miss can be serviced at a time, a lockup free cache must include  circuits that keep track of all the outstanding misses.

- Special registers may hold the necessary

  information about these misses.

# Virtual Memory

# Virtual Memories

- Recall that an important challenge in the design of a computer system is to provide a large, fast memory system at an affordable cost.

- Architectural solutions to increase the effective speed and size of the memory system.

- Cache memories were developed to increase the effective speed of the memory system.

- Virtual memory is an architectural solution to increase the effective size of the memory system.

# Virtual Memories *(Contd.)*

- Recall that the addressable memory space depends on the number of address bits in a computer.

  - For example, if a computer issues 32-bit addresses, the addressable memory space is 4G bytes.

- Physical main memory in a computer is generally not as large as the entire possible addressable space.

  - Physical memory typically ranges from a few hundred megabytes to 1G bytes.

- Large programs that cannot fit completely into the main memory have their parts stored on secondary storage devices such as magnetic disks.

  - Pieces of programs must be transferred to the main memory from secondary storage before they can be executed.

# Virtual Memories *(Contd.)*

- When a new piece of a program is to be transferred to the main memory, and the main memory is full, then some other piece in the main memory must be replaced.

  - Recall this is very similar to what we studied in case of cache memories.

- Operating system automatically transfers data between the main memory and secondary storage.

  - Application programmer need not be concerned with this transfer.

  - Also, application programmer does not need to be aware of the limitations imposed by the available physical memory.

# Virtual Memories *(Contd.)*

- Techniques that automatically move program and data between main memory and secondary storage when they are required for execution are called <u>virtual-memory</u> techniques.

- Programs and processors reference an instruction or data independent of the size of the main memory.

- Processor issues binary addresses for instructions and data.

  - These binary addresses are called logical or virtual addresses.

- Virtual addresses are translated into physical addresses by a combination of hardware and software subsystems.

  - If virtual address refers to a part of the program that is currently in the main memory, it is accessed immediately.

  - If the address refers to a part of the program that is not currently in the main memory, it is first transferred to the main memory before it can be used.

# Virtual Memory Organization

- Memory management unit (MMU) translates virtual addresses into physical addresses.

- If the desired data or instructions are in the main memory they are fetched as described previously.

- If the desired data or instructions are not in the main memory, they must be transferred from secondary storage to the main memory.

- MMU causes the operating system to bring the data from the secondary storage into the main memory.

| Processor |
|---|

Virtual address

Data

| MMU |
|---|

Physical address

| Cache |
|---|

Data    Physical address

| Main memory |
|---|

DMA transfer

| Disk storage |
|---|

# Address Translation

- Assume that program and data are composed of fixed-length units called pages.

- A page consists of a block of words that occupy contiguous locations in the main memory.

- Page is a basic unit of information that is transferred between secondary storage and main memory.

- Size of a page commonly ranges from 2K to 16K bytes.

  - Pages should not be too small, because the access time of a secondary storage device is much larger than the main memory.

  - Pages should not be too large, else a large portion of the page may not be used, and it will occupy valuable space in the main memory.

# Address Translation *(Contd.)*

- Concepts of virtual memory are similar to the concepts of cache memory.

- Cache memory:

  - Introduced to bridge the speed gap between the processor and the main memory.

  - Implemented in hardware.

- Virtual memory:

  - Introduced to bridge the speed gap between the main memory and secondary storage.

  - Implemented in part by software.

# Address Translation *(Contd.)*

- Each virtual or logical address generated by a processor is interpreted as a virtual page number (high-order bits) plus an offset (low-order bits) that specifies the location of a particular byte within that page.

- Information about the main memory location of each page is kept in the page table.

  - Main memory address where the page is stored.

  - Current status of the page.

- Area of the main memory that can hold a page is called as page frame.

- Starting address of the page table is kept in a page table base register.

# Address Translation *(Contd.)*

- Virtual page number generated by the processor is added to the contents of the page table base register.

  - This provides the address of the corresponding entry in the page table.

- The contents of this location in the page table give the starting address of the page if the page is currently in the main memory.

# Paging Hardware

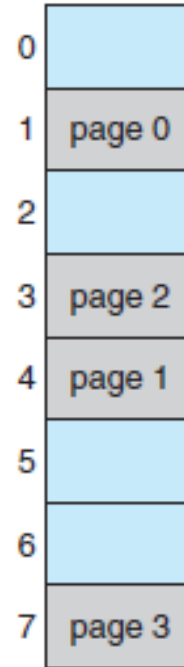# Paging model for logical and physical memory

# Paging example for a 32-bit memory with 4-Byte pages

# Address Translation *(Contd.)*

# Address Translation *(Contd.)*

- Page table entry for a page also includes some control bits which describe the status of the page while it is in the main memory.

- One bit indicates the validity of the page.

  - Indicates whether the page is actually loaded into the main memory.

  - Allows the operating system to invalidate the page without actually removing it.

- One bit indicates whether the page has been modified during its residency in the main memory.

  - This bit determines whether the page should be written back to the disk when it is removed from the main memory.

  - Similar to the dirty or modified bit in case of cache memory.

# Address Translation *(Contd.)*

- Other control bits for various other types of restrictions that may be imposed.

  - For example, a program may only have read permission for a page, but not write or modify permissions.

- Where should the page table be located?
- Recall that the page table is used by the MMU for every read and write access to the memory.
  - Ideal location for the page table is within the MMU.
- Page table is quite large.
- MMU is implemented as part of the processor chip.
- Impossible to include a complete page table on the chip.
- Page table is kept in the main memory.
- A copy of a small portion of the page table can be accommodated within the MMU.
  - Portion consists of page table entries that correspond to the most recently accessed pages.
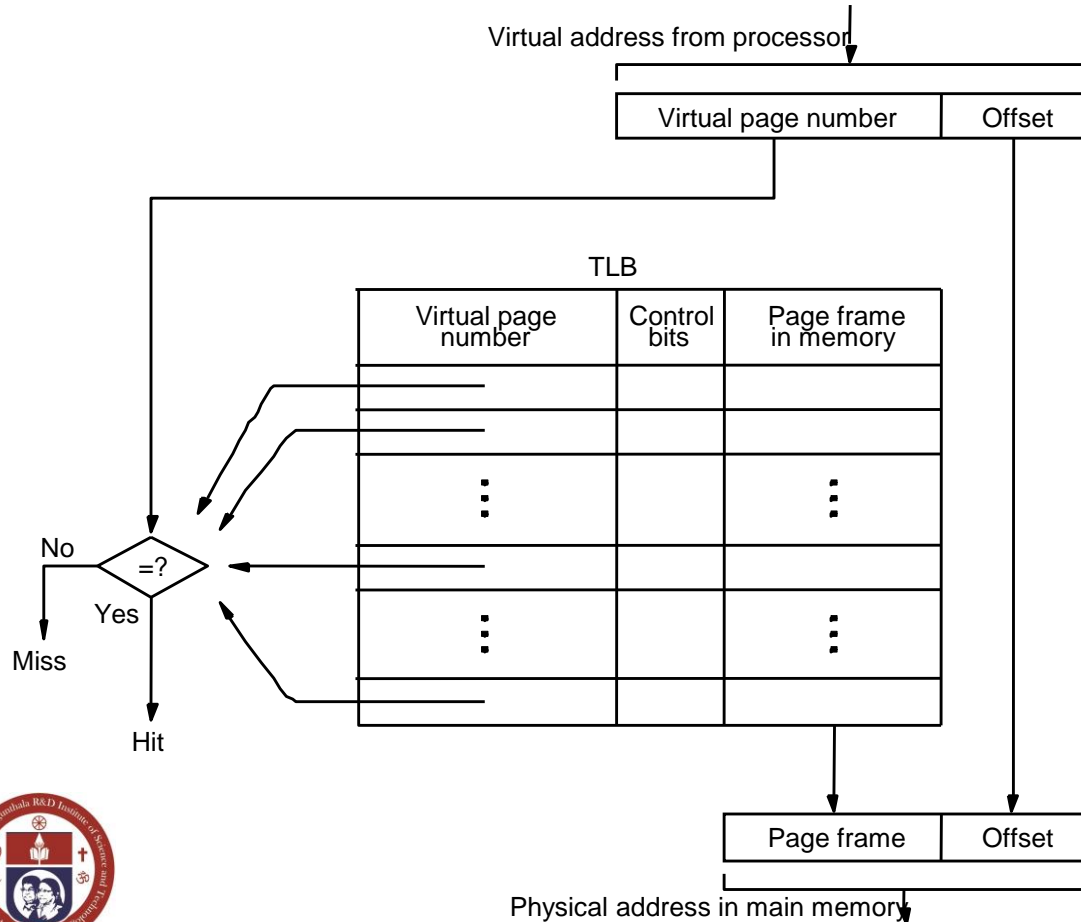
# Address Translation *(Contd.)*

- A small cache called as Translation Lookaside Buffer (TLB) is included in the MMU.

  - TLB holds page table entries of the most recently accessed pages.

- Recall that cache memory holds most recently accessed blocks from the main memory.

  - Operation of the TLB and page table in the main memory is similar to the operation of the cache and main memory.

- Page table entry for a page includes:

  - Address of the page frame where the page resides in the main memory.

  - Some control bits.

- In addition to the above for each page, TLB must hold the virtual page number for each page.

# Address Translation *(Contd.)*



## Associative-mapped TLB

- High-order bits of the virtual address generated by the processor select the virtual page.
- These bits are compared to the virtual page numbers in the TLB.
- If there is a match, a hit occurs and the corresponding address of the page frame is read.
- If there is no match, a miss occurs and the page table within the main memory must be consulted.
- Set-associative mapped TLBs are found in commercial processors.

# Address Translation *(Contd.)*

- How to keep the entries of the TLB coherent with the contents of the page table in the main memory?

- Operating system may change the contents of the page table in the main memory.

  - Simultaneously it must also invalidate the corresponding entries in the TLB.

- A control bit is provided in the TLB to invalidate an entry.

- If an entry is invalidated, then the TLB gets the information for that entry from the page table.

  - Follows the same process that it would follow if the entry is not found in the TLB or if a "miss" occurs.

# Address Translation *(Contd.)*

- What happens if a program generates an access to a page that is not in the main memory?

- In this case, a page fault is said to occur.

  - Whole page must be brought into the main memory from the disk, before the execution can proceed.

- Upon detecting a page fault by the MMU, following actions occur:

  - MMU asks the operating system to intervene by raising an exception.

  - Processing of the active task which caused the page fault is interrupted.

  - Control is transferred to the operating system.

  - Operating system copies the requested page from secondary storage to the main memory.

  - Once the page is copied, control is returned to the task which was interrupted.

# Address Translation *(Contd.)*

- Servicing of a page fault requires transferring the requested page from secondary storage to the main memory.

- This transfer may incur a long delay.

- While the page is being transferred the operating system may:

  - Suspend the execution of the task that caused the page fault.

  - Begin execution of another task whose pages are in the main memory.

- Enables efficient use of the processor.

# Address Translation *(Contd.)*

- How to ensure that the interrupted task can continue correctly when it resumes execution?

- There are two possibilities:

  - Execution of the interrupted task must continue from the point where it was interrupted.

  - The instruction must be restarted.

- Which specific option is followed depends on the design of the processor.

# Address Translation *(Contd.)*

- When a new page is to be brought into the main memory from secondary storage, the main memory may be full.

  - Some page from the main memory must be replaced with this new page.

- How to choose which page to replace?

  - This is similar to the replacement that occurs when the cache is full.

  - The principle of locality of reference (?) can also be applied here.

  - A replacement strategy similar to LRU can be applied.

- Since the size of the main memory is relatively larger compared to cache, a relatively large amount of programs and data can be held in the main memory.

  - Minimizes the frequency of transfers between secondary storage and main memory.

# Address Translation *(Contd.)*

- A page may be modified during its residency in the main memory.

- When should the page be written back to the secondary storage?

- Recall that we encountered a similar problem in the context of cache and main memory:

    - Write-through protocol(?)

    - Write-back protocol(?)

- Write-through protocol cannot be used, since it will incur a long delay each time a small amount of data is written to the disk.

# Input/ Output System - Accessing I/O Devices

# Accessing I/O devices



- Multiple I/O devices may be connected to the processor and the memory via a bus.
- Bus consists of three sets of lines to carry address, data and control signals.
- Each I/O device is assigned an unique address.
- To access an I/O device, the processor places the address on the address lines.
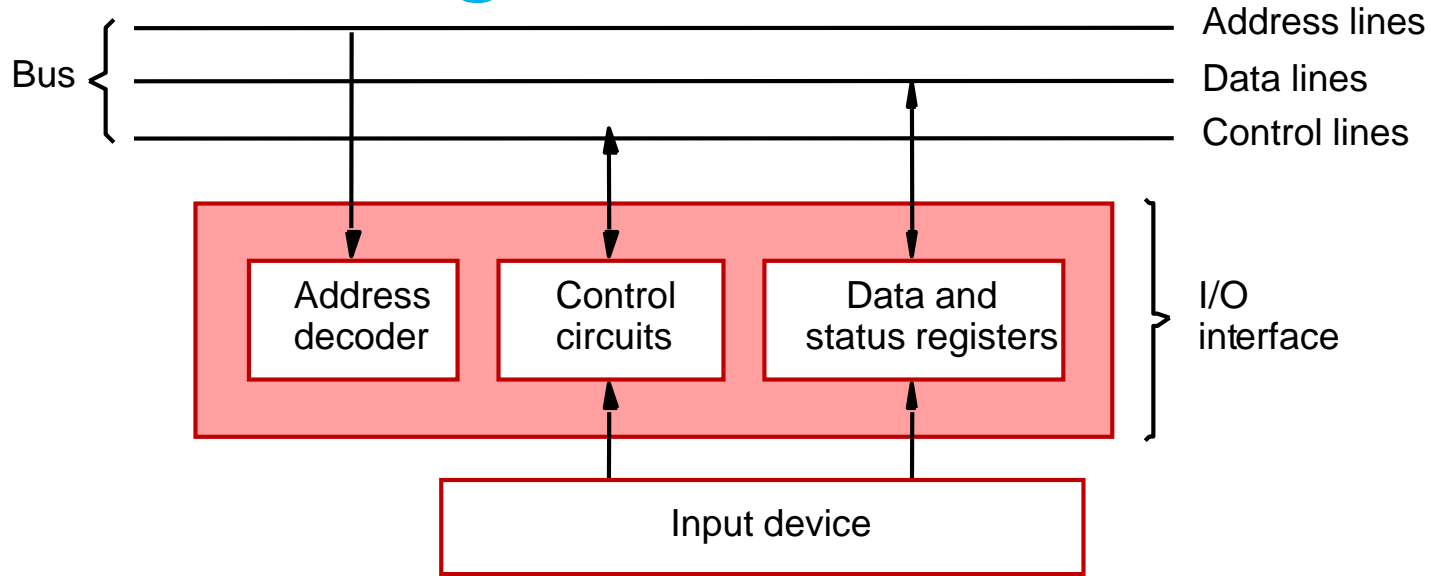- The device recognizes the address, and responds to the control signals.

# Accessing I/O devices *(Contd.)*

- I/O devices and the memory may share the same address space:
  - Memory-mapped I/O.
  - Any machine instruction that can access memory can be used to transfer data to or from an I/O device.
  - Simpler software.
- I/O devices and the memory may have different address spaces:
  - Special instructions to transfer data to and from I/O devices.
  - I/O devices may have to deal with fewer address lines.
  - I/O address lines need not be physically separate from memory address lines.
  - In fact, address lines may be shared between I/O devices and memory, with a control signal to indicate whether it is a memory address or an I/O address.

# Accessing I/O devices *(Contd.)*



- I/O device is connected to the bus using an I/O interface circuit which has:
  - Address decoder, control circuit, and data and status registers.
- Address decoder decodes the address placed on the address lines thus enabling the device to recognize its address.
- Data register holds the data being transferred to or from the processor.
- Status register holds information necessary for the operation of the I/O device.
- Data and status registers are connected to the data lines, and have unique addresses.
- I/O interface circuit coordinates I/O transfers.

# Accessing I/O devices *(Contd.)*

- Recall that the rate of transfer to and from I/O devices is slower than the speed of the processor. This creates the need for mechanisms to synchronize data transfers between them.

- Program-controlled I/O:

  - Processor repeatedly monitors a status flag to achieve the necessary synchronization.

  - Processor polls the I/O device.

- Two other mechanisms used for synchronizing data transfers between the processor and memory:

  - Interrupts.

  - Direct Memory Access.
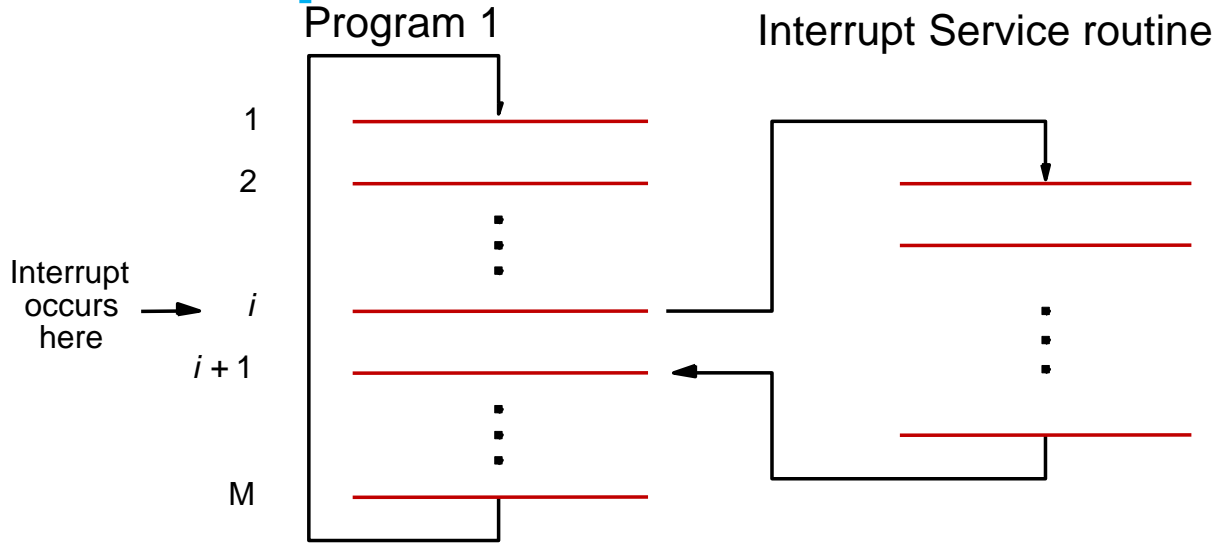
# Interrupts

# Interrupts

- In program-controlled I/O, when the processor continuously monitors the status of the device, it does not perform any useful tasks.

- An alternate approach would be for the I/O device to alert the processor when it becomes ready.

  - Do so by sending a hardware signal called an interrupt to the processor.

  - At least one of the bus control lines, called an interrupt-request line is dedicated for this purpose.

- Processor can perform other useful tasks while it is waiting for the device to be ready.

# Interrupts *(Contd.)*



Program 1

Interrupt Service routine

- Processor is executing the instruction located at address i when an interrupt occurs.
- Routine executed in response to an interrupt request is called the interrupt-service routine.
- When an interrupt occurs, control must be transferred to the interrupt service routine.
- But before transferring control, the current contents of the PC (i+1), must be saved in a known location.
- This will enable the return-from-interrupt instruction to resume execution at i+1.
- Return address, or the contents of the PC are usually stored on the processor stack.

# Interrupts *(Contd.)*

- Treatment of an interrupt-service routine is very similar to that of a subroutine.
- However there are significant differences:
  - A subroutine performs a task that is required by the calling program.
  - Interrupt-service routine may not have anything in common with the program it interrupts.
  - Interrupt-service routine and the program that it interrupts may belong to different users.
  - As a result, before branching to the interrupt-service routine, not only the PC, but other information such as condition code flags, and processor registers used by both the interrupted program and the interrupt service routine must be stored.
  - This will enable the interrupted program to resume execution upon return from interrupt service routine.

# Interrupts *(Contd.)*

- When a processor receives an interrupt-request, it must branch to the interrupt service routine.

- It must also inform the device that it has recognized the interrupt request.

- This can be accomplished in two ways:

  - Some processors have an explicit interrupt-acknowledge control signal for this purpose.

  - In other cases, the data transfer that takes place between the device and the processor can be used to inform the device.

# Interrupts *(Contd.)*

- Interrupt-requests interrupt the execution of a program, and may alter the intended sequence of events:
  - Sometimes such alterations may be undesirable, and must not be allowed.
  - For example, the processor may not want to be interrupted by the same device while executing its interrupt-service routine.
- Processors generally provide the ability to enable and disable such interruptions as desired.
- One simple way is to provide machine instructions such as *Interrupt-enable* and *Interrupt-disable* for this purpose.
- To avoid interruption by the same device during the execution of an interrupt service routine:
  - First instruction of an interrupt service routine can be Interrupt-disable.
  - Last instruction of an interrupt service routine can be Interrupt-enable.

# Interrupts *(Contd.)*

- Multiple I/O devices may be connected to the processor and the memory via a bus. Some or all of these devices may be capable of generating interrupt requests.
  - Each device operates independently, and hence no definite order can be imposed on how the devices generate interrupt requests?
- How does the processor know which device has generated an interrupt?
- How does the processor know which interrupt service routine needs to be executed?
- When the processor is executing an interrupt service routine for one device, can other device interrupt the processor?
- If two interrupt-requests are received simultaneously, then how to break the tie?

# Interrupts *(Contd.)*

- Consider a simple arrangement where all devices send their interrupt-requests over a single control line in the bus.

- When the processor receives an interrupt request over this control line, how does it know which device is requesting an interrupt?

- This information is available in the status register of the device requesting an interrupt:

  - The status register of each device has an *IRQ* bit which it sets to 1 when it requests an interrupt.

- Interrupt service routine can poll the I/O devices connected to the bus. The first device with *IRQ* equal to 1 is the one that is serviced.

- Polling mechanism is easy, but time consuming to query the status bits of all the I/O devices connected to the bus.

- The device requesting an interrupt may identify itself directly to the processor.

  - Device can do so by sending a special code (4 to 8 bits) to the processor over the bus.

  - Code supplied by the device may represent a part of the starting address of the interrupt-service routine.

  - The remainder of the starting address is obtained by the processor based on other information such as the range of memory addresses where interrupt service routines are located.

- Usually the location pointed to by the interrupting device is used to store the starting address of the interrupt-service routine.

# Interrupts *(Contd.)*

- Previously, before the processor started executing the interrupt service routine for a device, it disabled the interrupts from the device.

- In general, same arrangement is used when multiple devices can send interrupt requests to the processor.

  - During the execution of an interrupt service routine of device, the processor does not accept interrupt requests from any other device.

  - Since the interrupt service routines are usually short, the delay that this causes is generally acceptable.

- However, for certain devices this delay may not be acceptable.

  - Which devices can be allowed to interrupt a processor when it is executing an interrupt service routine of another device?

# Interrupts *(Contd.)*

- I/O devices are organized in a priority structure:

  - An interrupt request from a high-priority device is accepted while the processor is executing the interrupt service routine of a low priority device.

- A priority level is assigned to a processor that can be changed under program control.

  - Priority level of a processor is the priority of the program that is currently being executed.

  - When the processor starts executing the interrupt service routine of a device, its priority is raised to that of the device.

  - If the device sending an interrupt request has a higher priority than the processor, the processor accepts the interrupt request.

# Interrupts *(Contd.)*

- Processor's priority is encoded in a few bits of the processor status register.

    - Priority can be changed by instructions that write into the processor status register.

    - Usually, these are privileged instructions, or instructions that can be executed only in the supervisor mode.

    - Privileged instructions cannot be executed in the user mode.

    - Prevents a user program from accidentally or intentionally changing the priority of the processor.

- If there is an attempt to execute a privileged instruction in the user mode, it causes a special type of interrupt called as privilege exception.

- Each device has a separate interrupt-request and interrupt-acknowledge line.

- Each interrupt-request line is assigned a different priority level.

- Interrupt requests received over these lines are sent to a priority arbitration circuit in the processor.

- If the interrupt request has a higher priority level than the priority of the processor, then the request is accepted.

# Interrupts *(Contd.)*

- Which interrupt request does the processor accept if it receives interrupt requests from two or more devices simultaneously?.

- If the I/O devices are organized in a priority structure, the processor accepts the interrupt request from a device with higher priority.

  - Each device has its own interrupt request and interrupt acknowledge line.

  - A different priority level is assigned to the interrupt request line of each device.

- However, if the devices share an interrupt request line, then how does the processor decide which interrupt request to accept?

Polling scheme:
- If the processor uses a polling mechanism to poll the status registers of I/O devices to determine which device is requesting an interrupt.
- In this case the priority is determined by the order in which the devices are polled.
- The first device with status bit set to 1 is the device whose interrupt request is accepted.
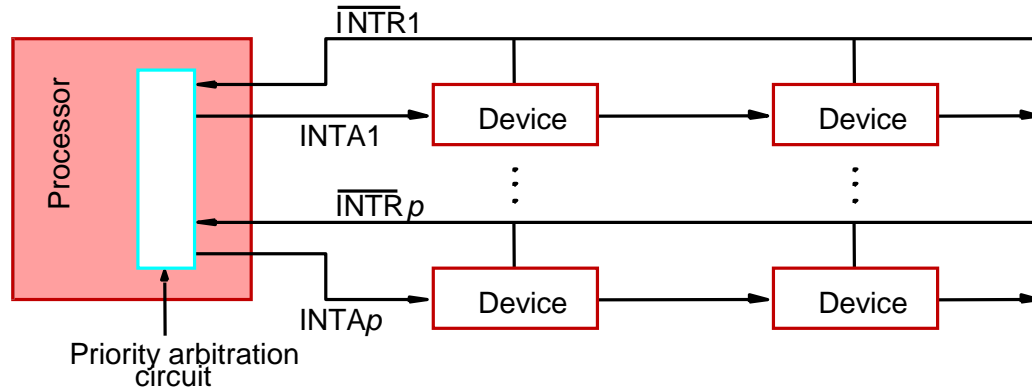
Daisy chain scheme:



- Devices are connected to form a daisy chain.
- Devices share the interrupt-request line, and interrupt-acknowledge line is connected to form a daisy chain.
- When devices raise an interrupt request, the interrupt-request line is activated.
- The processor in response activates interrupt-acknowledge.
- Received by device 1, if device 1 does not need service, it passes the signal to device 2.
- Device that is electrically closest to the processor has the highest priority.

# Interrupts *(Contd.)*

- When I/O devices were organized into a priority structure, each device had its own interrupt-request and interrupt-acknowledge line.
- When I/O devices were organized in a daisy chain fashion, the devices shared an interrupt-request line, and the interrupt-acknowledge propagated through the devices.
- A combination of priority structure and daisy chain scheme can also used.



- Devices are organized into groups.
- Each group is assigned a different priority level.
- All the devices within a single group share an interrupt-request line, and are connected to form a daisy chain.

# Direct Memory Access

# Direct Memory Access

- Direct Memory Access (DMA):

  - A special control unit may be provided to transfer a block of data directly between an I/O device and the main memory, without continuous intervention by the processor.

- Control unit which performs these transfers is a part of the I/O device's interface circuit. This control unit is called as a DMA controller.

- DMA controller performs functions that would be normally carried out by the processor:

  - For each word, it provides the memory address and all the control signals.

  - To transfer a block of data, it increments the memory addresses and keeps track of the number of transfers.
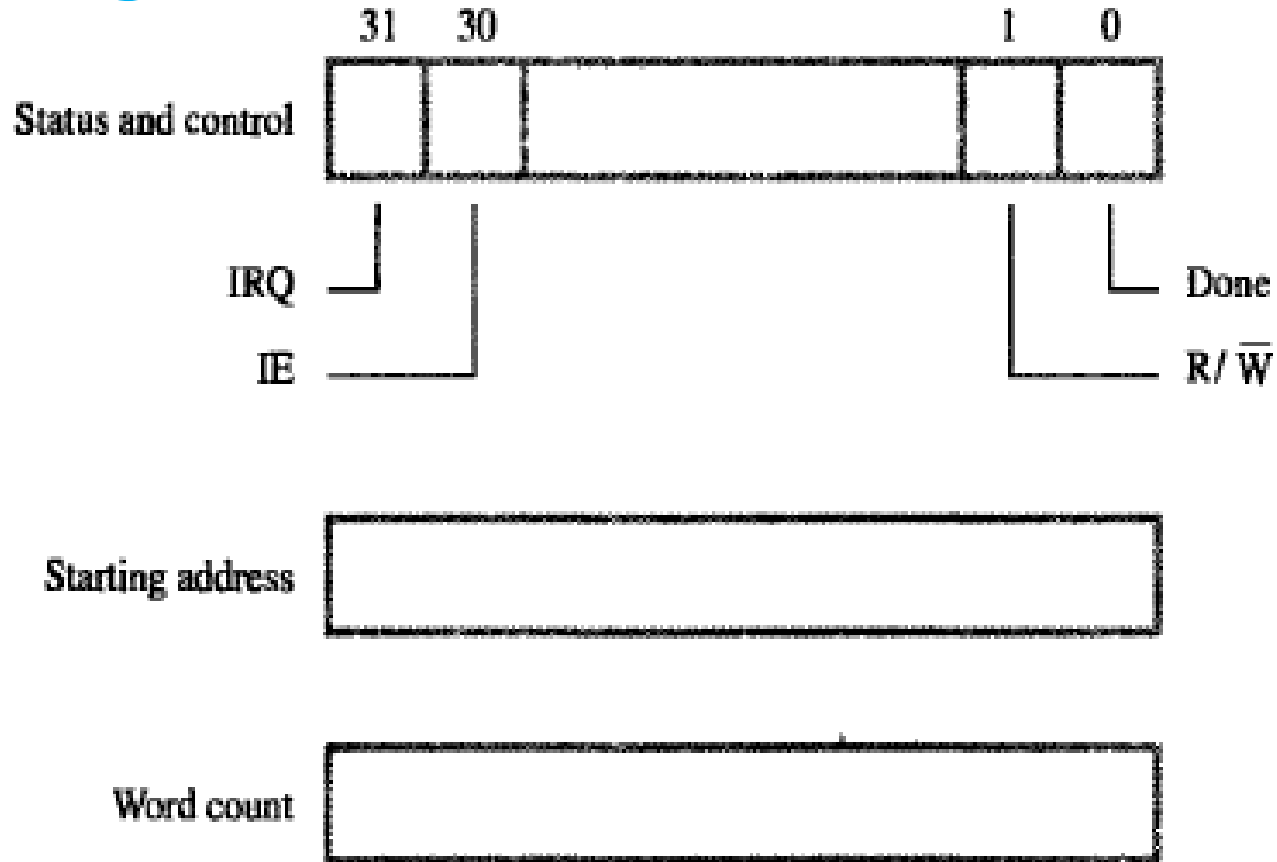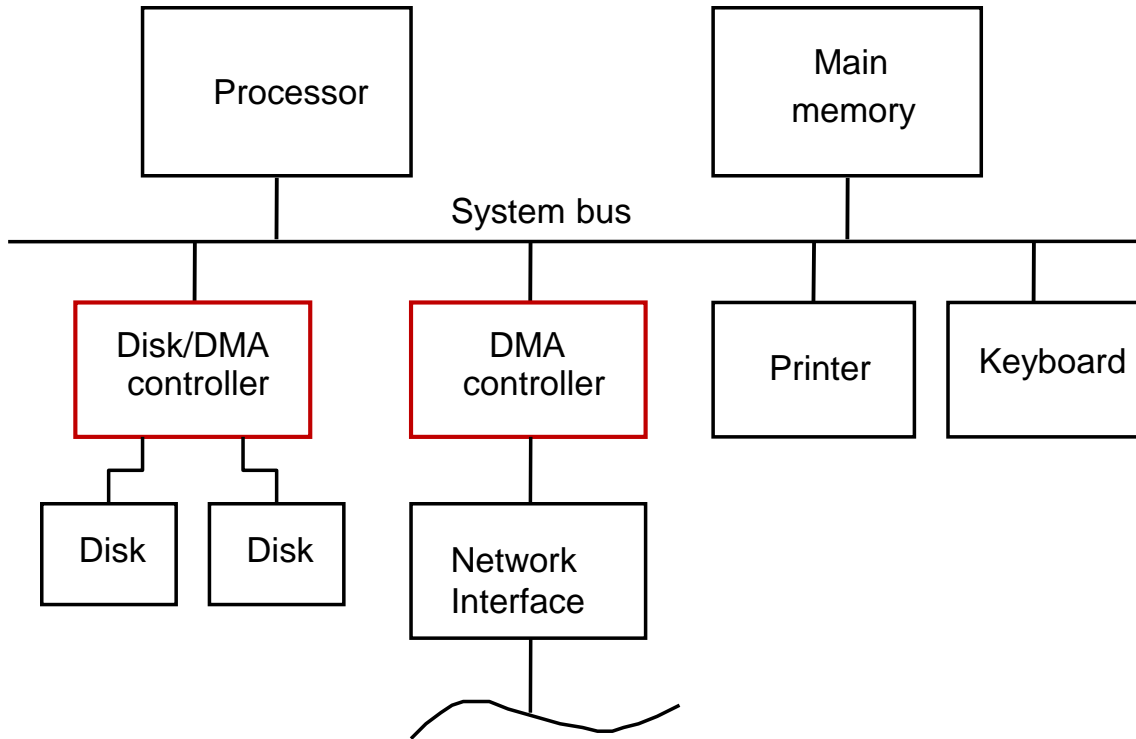
# Direct Memory Access *(Contd.)*

- DMA controller can transfer a block of data from an external device to the processor, without any intervention from the processor.

  - However, the operation of the DMA controller must be under the control of a program executed by the processor. That is, the processor must initiate the DMA transfer.

- To initiate the DMA transfer, the processor informs the DMA controller of:

  - Starting address,

  - Number of words in the block.

  - Direction of transfer (I/O device to the memory, or memory to the I/O device).

- Once the DMA controller completes the DMA transfer, it informs the processor by raising an interrupt signal.

# Registers in a DMA transfer

# Direct Memory Access *(Contd.)*



- DMA controller connects a high-speed network to the computer bus.

- Disk controller, which controls two disks also has DMA capability. It provides two DMA channels.

- It can perform two independent DMA operations, as if each disk has its own DMA controller. The registers to store the memory address, word count and status and control information are duplicated.
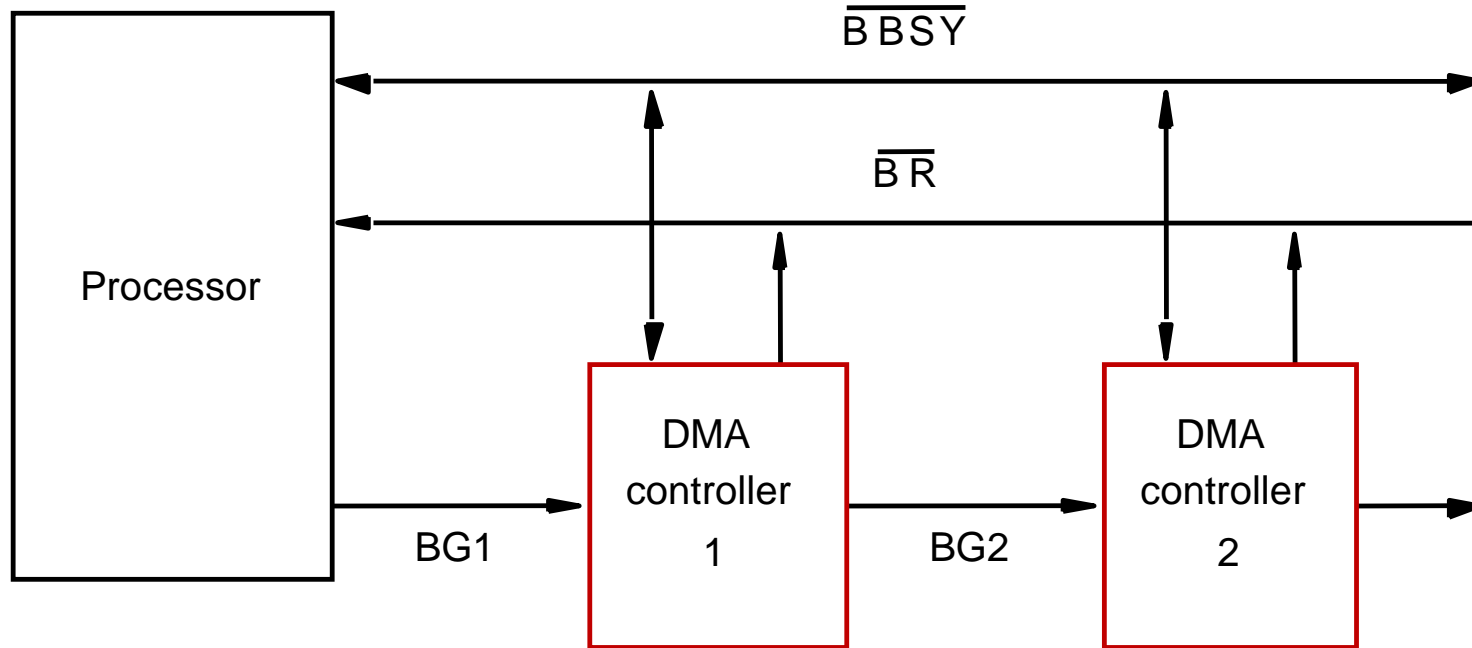
# Bus Arbitration

- Processor and DMA controllers both need to initiate data transfers on the bus and access main memory.

- The device that is allowed to initiate transfers on the bus at any given time is called the bus master.

- When the current bus master relinquishes its status as the bus master, another device can acquire this status.

  - The process by which the next device to become the bus master is selected and bus mastership is transferred to it is called bus arbitration.

- Centralized arbitration:

  - A single bus arbiter performs the arbitration.

- Distributed arbitration:

  - All devices participate in the selection of the next bus master.
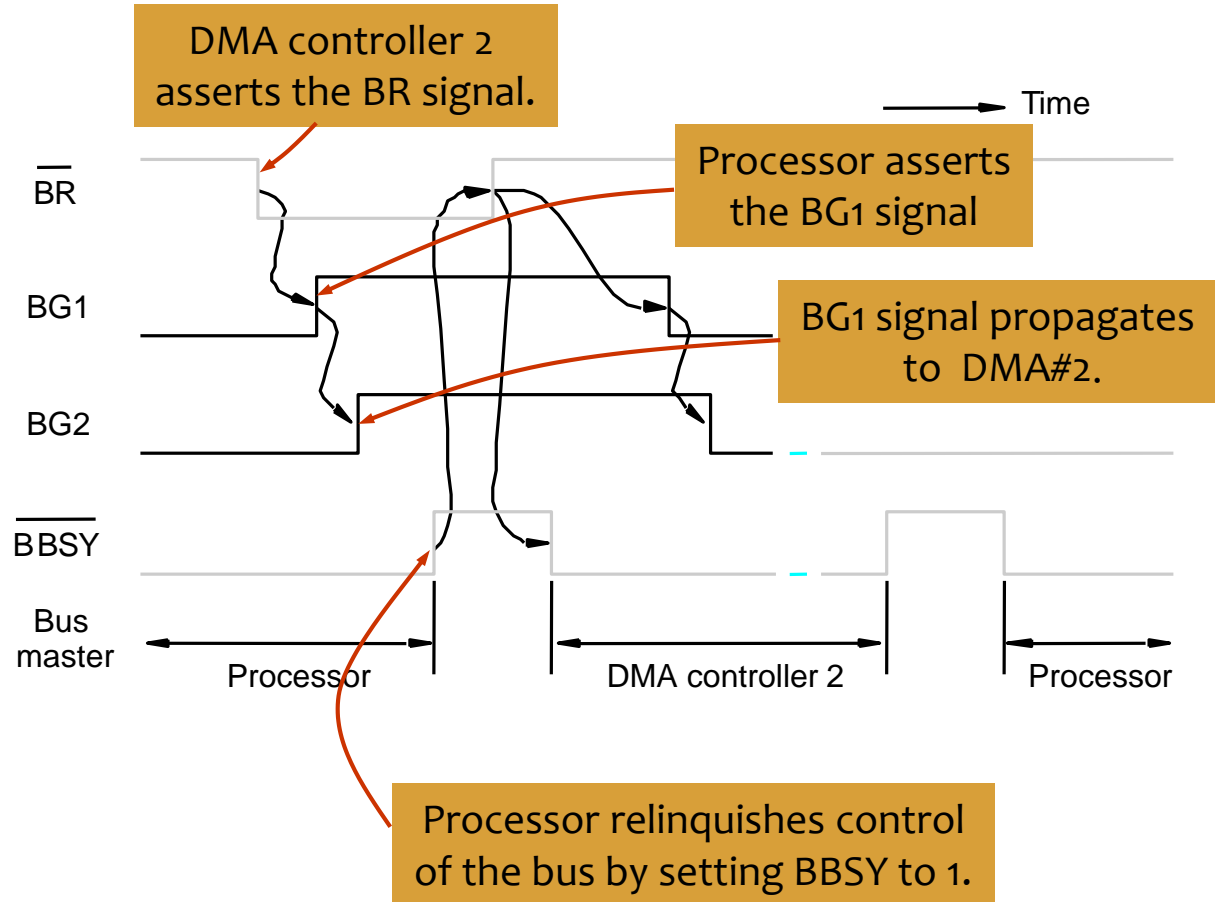
# Centralized Bus Arbitration

# Centralized Bus Arbitration *(Contd.)*

- Bus arbiter may be the processor or a separate unit connected to the bus.

- Normally, the processor is the bus master, unless it grants bus membership to one  of the DMA controllers.

- DMA controller requests the control of the bus by asserting the Bus Request (BR) line.

- In response, the processor activates the Bus-Grant1 (BG1) line, indicating that the  controller may use the bus when it is free.

- BG1 signal is connected to all DMA controllers in a daisy chain fashion.

- BBSY signal is 0, it indicates that the bus is busy. When BBSY becomes 1, the DMA  controller which asserted BR can acquire control of the bus.

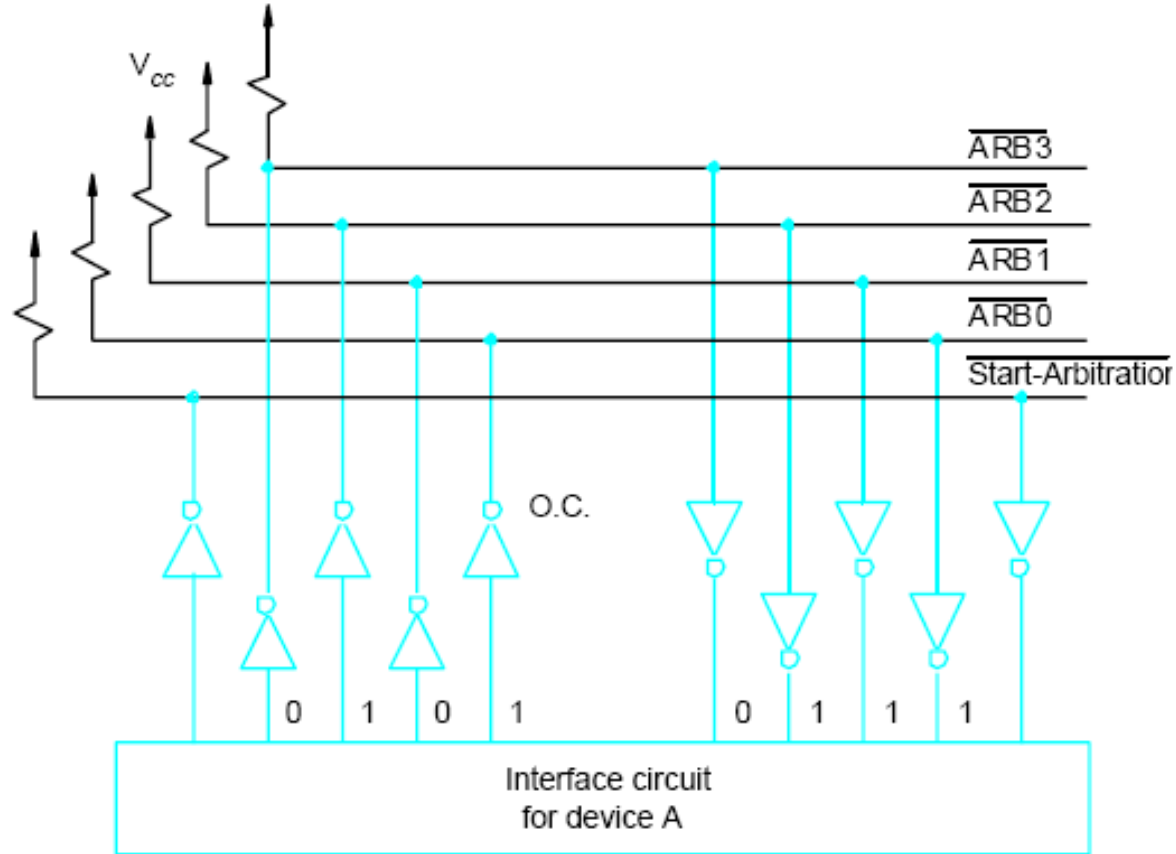# Centralized Bus Arbitration *(Contd.)*

# Distributed Arbitration *(Contd.)*

- All devices waiting to use the bus share the responsibility of carrying out the arbitration process.

  - Arbitration process does not depend on a central arbiter and hence distributed arbitration has higher reliability.

- Each device is assigned a 4-bit ID number.

- All the devices are connected using 5 lines, 4 arbitration lines to transmit the ID, and one line for the Start-Arbitration signal.

- To request the bus a device:

  - Asserts the Start-Arbitration signal.

  - Places its 4-bit ID number on the arbitration lines.

- The pattern that appears on the arbitration lines is the logical-OR of all the 4-bit device IDs placed on the arbitration lines.

Arbitration process:

- Each device compares the pattern that appears on the arbitration lines to its own ID, starting with MSB.

- If it detects a difference, it transmits 0s on the arbitration lines for that and all lower bit positions.

- The pattern that appears on the arbitration lines is the logical-OR of all the 4-bit device IDs placed on the arbitration lines.

# Distributed Arbitration *(Contd.)*

- Device A has the ID 5 and wants to request the bus:

  - Transmits the pattern 0101 on the arbitration lines.

- Device B has the ID 6 and wants to request the bus:

  - Transmits the pattern 0110 on the arbitration lines.

- Pattern that appears on the arbitration lines is the logical OR of the patterns:

  - Pattern 0111 appears on the arbitration lines.

Arbitration process:

- Each device compares the pattern that appears on the arbitration lines to its own ID, starting with MSB.

- If it detects a difference, it transmits 0s on the arbitration lines for that and all lower bit positions.

- Device A compares its ID 5 with a pattern 0101 to pattern 0111.

- It detects a difference at bit position 0, as a result, it transmits a pattern 0100 on the arbitration lines.

- The pattern that appears on the arbitration lines is the logical-OR of 0100 and 0110, which is 0110.

- This pattern is the same as the device ID of B, and hence B has won the arbitration.