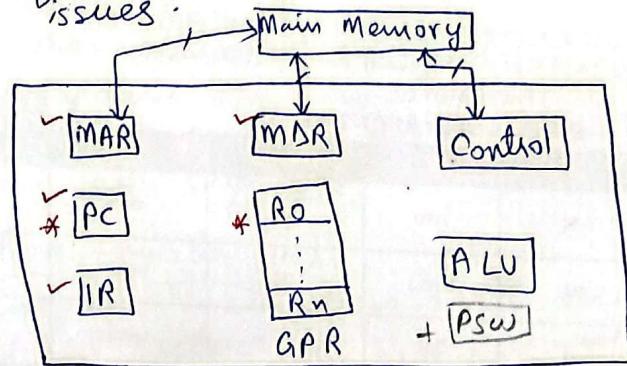


## REGISTERS

- length of various registers and which registers to include are design issues.



user visible \*

not user visible ✓

read only +

ALU contains storage location called registers.

- \* User Visible Registers: Enable the machine or assembly-language programmer to minimize main memory references by optimizing use of registers.
- They can be referenced by means of machine language programs.
- \* Control and Status Registers: Used by control unit to control the operation of the processor and by privileged, operating system programs to control the execution of programs.

There is not a clean separation of registers into these two categories. For example, on some machines program counter is user visible (e.g.: Pentium), but on many is not.

### USER VISIBLE REGISTER required by machine language

#### 1- General Purpose Registers:

- used to temporary store the operands.

- can contain address also in some cases (indirect addressing mode).

#### 2- Dedicated Data and Address Registers:

- Data registers hold data only and cannot be implied in the calculation of an operand address.

- Address registers may be devoted to a particular add. mode.

- Some address registers are:

a- Segment pointers: holds the address of the base of the segment.

b- Index registers: used for indexed addressing and may be auto-indexed.

c- Stack pointer: If there is user visible stack addressing, then this register points to the top of the stack.

In some machines subroutines call result in automatic saving of all user-visible registers, to be restored on return. On other machines, it is the responsibility of the programmer to save the relevant contents of regs.

## CONTROL & STATUS REGISTERS:

Most commonly found registers are:  
Program Counter (PC):

Contains the address of an instruction to be fetched.  
Processor updates the PC after each instruction fetch, so that the PC always points to the next instruction to be executed.  
A branch or skip instruction also modify the contents of PC.

### Instruction Register (IR)

- Contains the instruction most recently fetched.
- The fetched instruction is loaded into an IR, where the opcode and operand specifiers are analyzed.

### Memory Address Registers (MAR):

Contains the address of an instruction in memory.  
Connects directly to the address bus.

### Memory Buffer Register (MBR):

- Contains a word of data to be written to memory or word most recently read.
- Connects directly to the data bus.

MAR and MBR are not provided in a machine then some equivalent buffering mechanisms must be there whereby the bits to be transferred on the system bus are staged and bits to be read from the data bus are temporarily stored.

The four registers mentioned above are used for the movement of data between the processor and memory.

## ACCESSORS STATUS

### PROGRAM STATUS WORD (PSW):

It could be a register or a set of registers.

- It could be a register or a set of registers.
- Contains program status information.
- Typically contains condition codes (flag register) and other information (starting address of program in memory, etc).

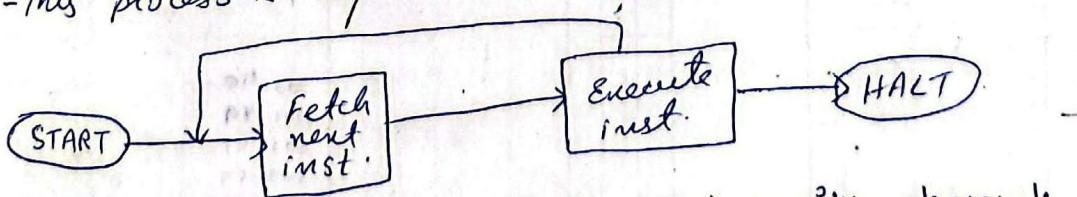
### Condition Codes (Flag Register):

- Partially visible to the user. These are read only.
- Condition codes are set by the processor hardware as the result of operations. For e.g. an arithmetic operation may produce a positive, negative or zero result.
- In addition to the result itself being stored in memory, a condition code is also set. Subsequently the code may be tested as a part of some conditional branch operation.

- they make the circuitry complex, so some machines don't include condition codes at all. (eg IA-64 architecture)
- common fields or flag include the following:
  1. sign : contains the sign bit of the result of the last arithmetic operation.
  2. zero: set when the result is 0.
  3. carry: set if an operation resulted in a carry into, or borrow out of a higher order bit. Used for multi-word arithmetic operations.
  4. equal: set if a logical compare result is equality.
  5. overflow: used to indicate arithmetic overflow.
  6. interrupt enable/disable: used to enable or disable interrupts.
  7. supervisor: indicates whether the processor is executing in supervisor or user mode. Certain privileged instructions can be executed only in supervisor mode, and certain areas of memory can be accessed only in supervisor mode.

### BASIC INSTRUCTION CYCLE:

- A program is a set of instructions stored in memory (mostly in consecutive locations).
- The processing required for a single instruction is called an instruction cycle.
- In simplest form, an instruction cycle consists of two steps:
  1. instruction fetch: fetches the instruction from memory (one at a time).
  2. instruction execute: executes the instruction.  
Instruction execution may involve several operations and depends on the nature of instruction.
- This process is repeated until all instructions are executed.

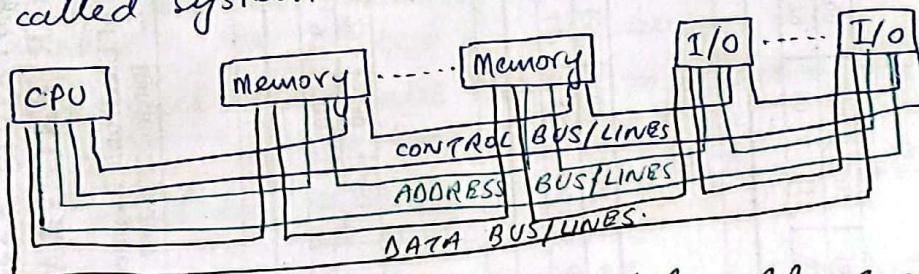


Accumulator: Employed to hold temp orally operands and results of ALU operations.

## BUS INTERCONNECTION: 34

(15)

- The most commonly used interconnection structure so far is bus and various multiple-bus structures.
- A bus is a communication pathway connecting two or more devices.
- A bus is a shared transmission medium.
- Only one device can successfully transmit the data at a time, otherwise the signals will be garbled.
- Signal transmitted by one device is available for reception by all other devices attached to the bus.
- A bus consists of multiple lines, each line is capable of transmitting signals representing binary 1 and binary 0.
- A group of bits can be sent simultaneously, thus a reasonable speed is achieved. 8 bits  $\rightarrow$  8 lines transmit them in parallel.
- Computer systems contains a number of different buses that provide pathways between components at various level of computer system hierarchy.
- A bus that connects major computer components (CPU, MM, I/O) is called system bus.



# of lines  
= width of bus.

- There are three functional groups: data, address, control.
- there may be power distribution lines that supply power to the attached modules.

### 1-DATA BUS:

- provides path for memory data between system modules.
- if width of data bus increases, performance increases.
- e.g: if data bus is 8-bits wide and instruction is 16-bits wide, then two cycles will be needed to fetch each instruction.
- width is usually in power of 2. (8/16/32/64...).

### 2-ADDRESS BUS:

- used to designate source or destination of data on data bus.
- also used to address I/O module and I/O ports.
- width of the address bus determines the maximum possible memory capacity of the system.

## INTERCONNECTION STRUCTURES

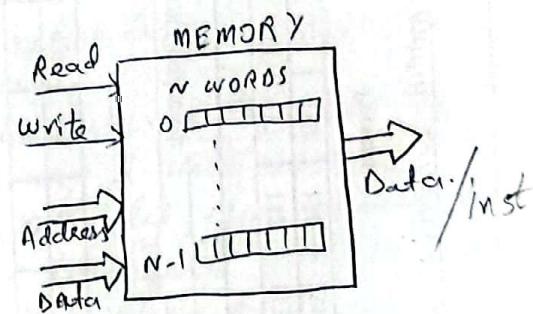
(14)

- The collection of paths connecting the various modules (processors, memory, I/O) is called interconnection structures.
- Major forms of input and output (exchanges) are as follows:

### MEMORY:

- two operations possible:
  - \* read
  - \* write

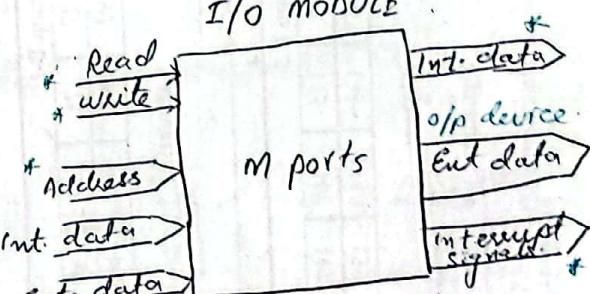
- sometimes the memory also sends acknowledgement signals to processor (asynchronous mode).



### I/O MODULE:

- from an internal point of view (from/to processor), I/O is functionally similar to memory.
- 2 I/O operations, read and write.
- may control more than one I/O device.
- interface to an external device is called port. Each port has a unique address.
- there are external data paths for the input and output of data with an external device.
- it may send INT signal to processor.

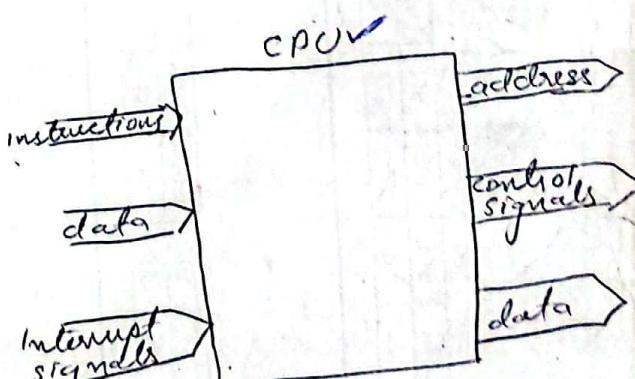
### I/O MODULE:



### PROCESSOR:

- reads instruction and data
- writes out data after processing.
- uses control signals to control the overall operation of system.
- receives interrupt signals.

Thus the interconnection structure must support the following types of transfers:



- 1- Memory to processor: Processor reads instruction/data from memory.
- 2- Processor to memory: Processor writes data to memory.
- 3- I/O to processor: Processor reads (via I/O module) from a device.
- 4- Processor to I/O: Processor sends data to I/O (via I/O module).
- 5- I/O to or from memory: an I/O module is allowed to exchange data directly with memory, without processor intervention — DMA (direct memory access)

## CONTROL BUS:

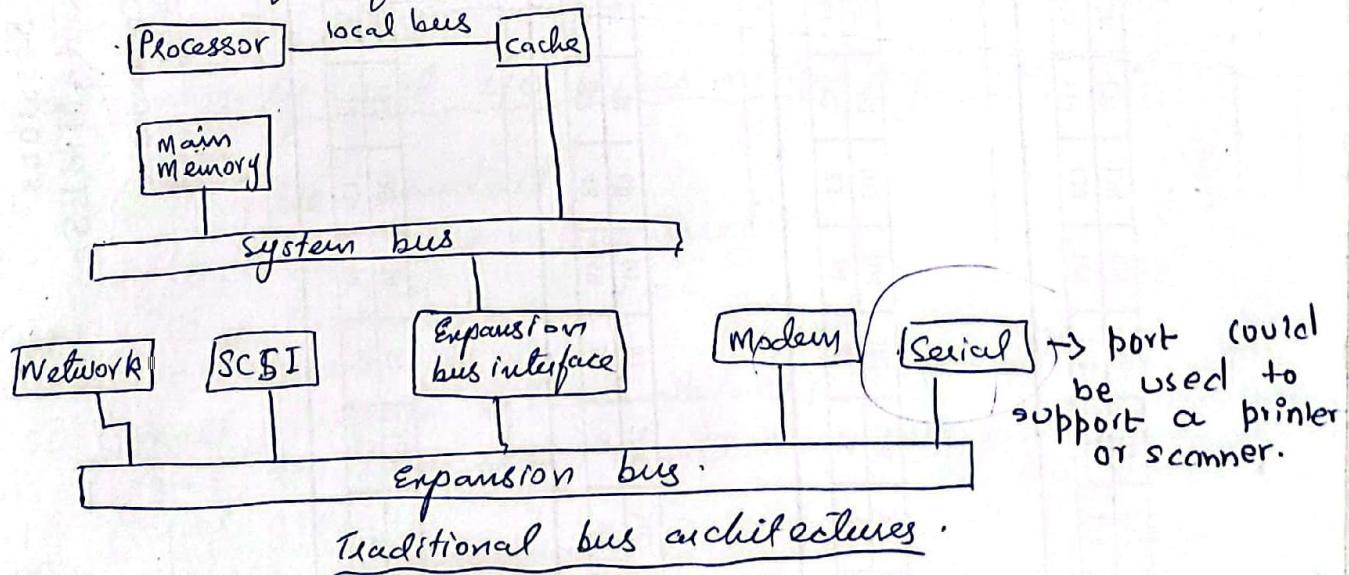
(16)

- used to control the access to and use of data and address buses.
- because the data and address lines are shared by all components, there must be a means of controlling their use.
- Control signals/lines transmit both:
  - a - timing signals: indicate the validity of data and address information.
  - b - command signals: specify operations to be performed.
- typical control lines include:
  - 1 - Memory write: causes data on the bus to be written into the addressed location.
  - 2 - Memory read: causes data from the addressed location to be placed on the bus.
  - 3 - I/O write: causes data on the bus to be output to the addressed output port.
  - 4 - I/O read: causes data from the addressed I/O port to be placed on the bus.
  - 5 - Transfer ~~over~~ ACK: Indicates that data have been accepted from or placed on the bus.
  - 6 - Bus request: Indicates that a module needs to gain control of the bus.
  - 7 - Bus grant: Indicates that the requesting module has been granted control of the bus.
  - 8 - Interrupt request: Indicates that an interrupt is pending.
  - 9 - Interrupt ACK: Acknowledges that the pending interrupt has been recognized.
  - 10 - Clock: Used to synchronize operations.
  - 11 - Reset: Initializes all modules.
- Before starting any bus operation, a module needs to get control of the bus.
- system bus is usually etched on the board or the components are on the same chip.
- Nowadays, for increasing number of modules ~~to share the~~ bus, bus architectures are replaced by Icn (Interconnection network).
- If a great number of devices are connected to the bus, performance will suffer because of the following two reasons:

The more devices attached to the bus, the greater the bus length and hence the greater the propagation delay, <sup>more contention.</sup> (17) The bus may become a bottleneck as the aggregate data transfer demand approaches the capacity of the bus. This problem can be countered to some extent by increasing the data rate that the bus carry and by using wider buses. But still not sufficient for rapidly growing devices (e.g: graphics and video controllers, network interfaces).

### BUS HIERARCHIES Page No: 88

- most computer systems use multiple buses, laid out in a hierarchy.
- It is possible to connect I/O controllers directly onto the system bus. A more efficient solution is to use expansion buses.
- An expansion bus interface buffers data transfers between the system bus and the I/O controllers on the expansion bus.
- This arrangement allows the system to support a wide variety of I/O devices and at the same time insulate memory-to-processor traffic from I/O traffic.



- \* shows typical examples of I/O devices that might be attached to the expansion bus.
- \* Network connections include (LAN) and (WAN)
- \* SCSI (Small Computer System Interface) is itself a type of bus used to support local disk drivers and other peripherals.

## Mezzanine Architecture

Some commonly used I/O buses nowadays are:

[See handout]

- PCI (Peripheral component interconnect) 32-bit / 64-bit
- PCI-X (PCI extended) 64-bits
- PCI-e (PCI express) serial
- SCSI (small computer system interface). 16-bit
- SAS (serial attached SCSI) [faster than SCSI & SATA].  
[SCSI - ethernet port]
- SATA (Serial ATA) AT attachment
- Gb PU (graphics processing unit)  
[SAN-connected to fibre cables]
- AGP (accelerated graphics port)
- USB (universal serial bus)
- Firewire.

## ELEMENTS OF BUS DESIGN

(19)

### Bus types

- 1) Method of arbitration (centralized or distributed)
- 2) Timing.
- 3) Bus width.
- 4) Data transfer type.
- 5) Addressing.

### ① BUS TYPES:

\* Dedicated bus lines: A dedicated bus line is permanently assigned either to one function or to a physical subset of computer components. e.g. separate buses for data & address.

#### Time multiplexed bus lines:

- address and data may be transmitted over the same set of lines using Address Valid control line.
- Address valid control line is activated when address is placed on the bus and vice versa when data is placed.
- adv: save space and cost.
- disadv: requires complex circuitry.
- + potential reduction in performance because certain events that share the same lines, cannot take place in parallel. In 8086/88: half of the 16-bit address lines are multiplexed with data and upper half is multiplexed with status line. In 8085 lower 8-bits are multiplexed with data.

### ② METHOD OF ARBITRATION:

Arbitration → problem solving.

- In all but simplest systems, more than one module may need control of the bus. Bus is shared among various module.
- e.g.: bus can be needed simultaneously by the PR and by I/O module (DMA operation).
- ✓ As only one module ~~can~~ at a time can successfully transmit data, so some method of arbitration is needed (to allocate bus to be initiator of the transfer).

#### Centralized scheme:

- \* Centralized scheme:
  - a single hardware device, called bus controller or arbiter is responsible for allocating time on the bus.
  - the device may be a separate module or part of the controller.

- \* Distributed scheme:
  - there is no central controller, rather each module contains access control logic and the modules act together to share the bus.
  - transfer initiator, gets control of the bus → master.
  - other device involved in the exchange → slave.

## BUS WIDTH

- \* wider the data bus, greater the number<sup>(20)</sup> of bits transferred at one time.
- \* wider the address bus, greater the range of location that can be directly referenced.

## DATA TRANSFER TYPES

- \* Read operation: For dedicated as well as time multiplexed lines:
  - delay occurs for bus arbitration.
  - some delay occurs for data onto the data bus as soon as it has loaded the data.

# SEMICONDUCTOR MEMORIES

45

RWM  
read/write memory

ROM  
read only memory

RMM  
read mostly memory.

## \* RWM (read/write memory):

- data can be easily and rapidly read from or written into RAMs.
- reading and writing is accomplished through electrical signals.
- volatile memory.
- could be DRAM ~~or SRAM~~ or SRAM.

## \* DRAM (Dynamic RAM):

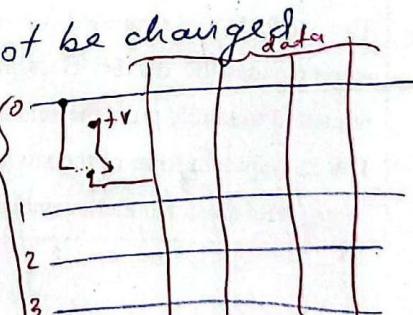
- made with cells that store data as charge on capacitor
- capacitors have natural tendency to discharge, they require periodic charge refreshing to maintain charge storage.
- they are DRAM memories
- used in RAMs.

## \* SRAM (Static RAM):

- data is stored on flip flops (a logical element)
- no refreshing is required, therefore fast
- expensive - complex circuit (relatively)
- used in most cache's.

## \* ROM (read only memory):

- contains permanent patterns of data that cannot be changed
- read operations are possible any time
- write operations is possible only once
- non-volatile and non erasable
- used in microwave oven, older BIOS, etc.
- data is actually wired into the chip
- data is actually wired into the chip
- MASK ROMs - factory programmed ROM.
- \* PROM (programmable ROM) - can be programmed by the user once.



## \* RMM (read mostly memory):

- read operations are more frequent than write operations.
- read operation is as simple as in RAM but write operation requires special hardware called ROM programmer.
- non-volatile and erasable.
- these ROMs should be erased before being rewritten.
- erasure also requires special hardware.
- types:

\* EPROM

\* FLASH

\* EEPROM



## EEPROM

## FLASH

## EEPROM

(146)

erase process	uses UV light	→	electrical voltages are used - in circuit erasure
size	smaller in size than EEPROM (more dense)	← more dense than EEPROM.	larger in size.
erase time	20 min or more	1 on a few seconds	~10 ms
endurance	100 programming cycles before chip is damaged.	100 to 10000 cycles	100000 cycles.
erase unit	whole contents are erased	section/block level erase uses: USB flash memories, cameras, video cards, etc.	byte level erase is possible.

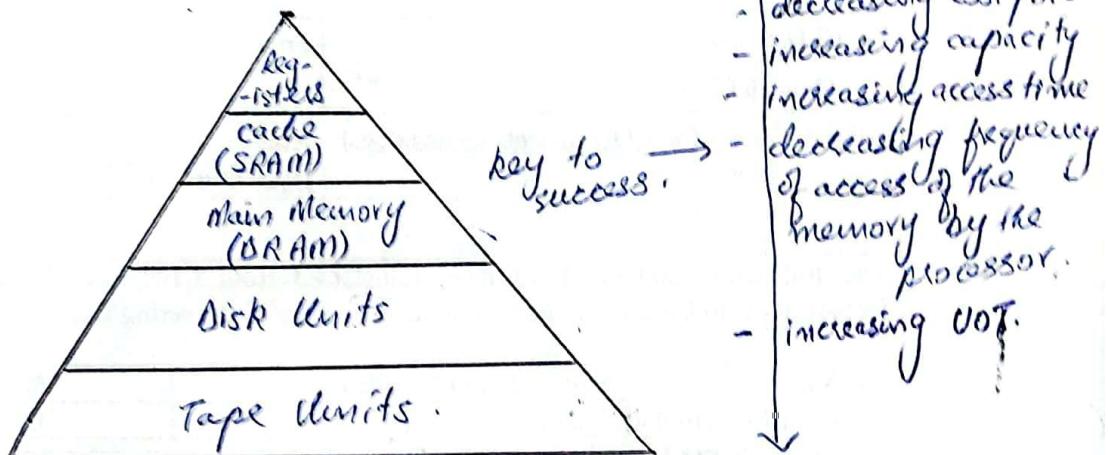
LOCALITY OF REFERENCE states that a program tends to access any small portion of logical address space at any point of time. Logical address space is the set of all logical addresses generated by the CPU.

Temporal locality: It is the locality over time. Recently referenced items are likely to be referenced again in the near future.  
e.g.: loops and subroutines.

Spatial locality: It is the locality over space. Tending to access items whose addresses are near to one another.  
e.g.: array traversal.

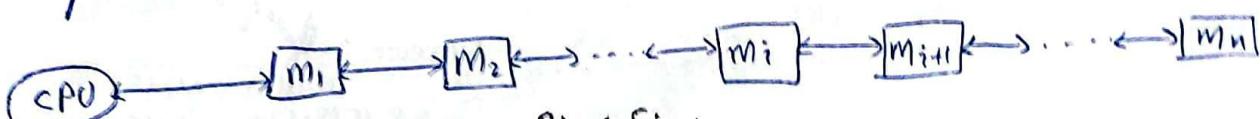
## MEMORY HIERARCHY

- There is a tradeoff between cost, capacity and speed of a memory.
- • Faster access time, greater cost per bit
- • Greater capacity, (if) smaller cost per bit
- • Greater capacity, slower access time.



Rationale behind the design of a memory hierarchy:

- cost of a memory varies inversely with the access time.  
 $t \propto \frac{1}{\sqrt{\text{cost}}}$
- It is therefore impossible to design an entire system using just one technology that offers both low cost and high operating speed (i.e. low access time). Hence a memory system is designed using different physical devices to achieve optimum performance.
- memory hierarchy reduces the speed gap between processor and memory system at reasonable cost.



size:  $S_i < S_{i+1}$   
access time:  $t_i < t_{i+1}$   
cost:  $C_i > C_{i+1}$   
unit of transfer:  $R_i < R_{i+1}$   
access frequency:  $f_i > f_{i+1}$

Hit = a reference generated by CPU is found in  $M_i$   
Miss = a reference generated by CPU is not found in  $M_i$