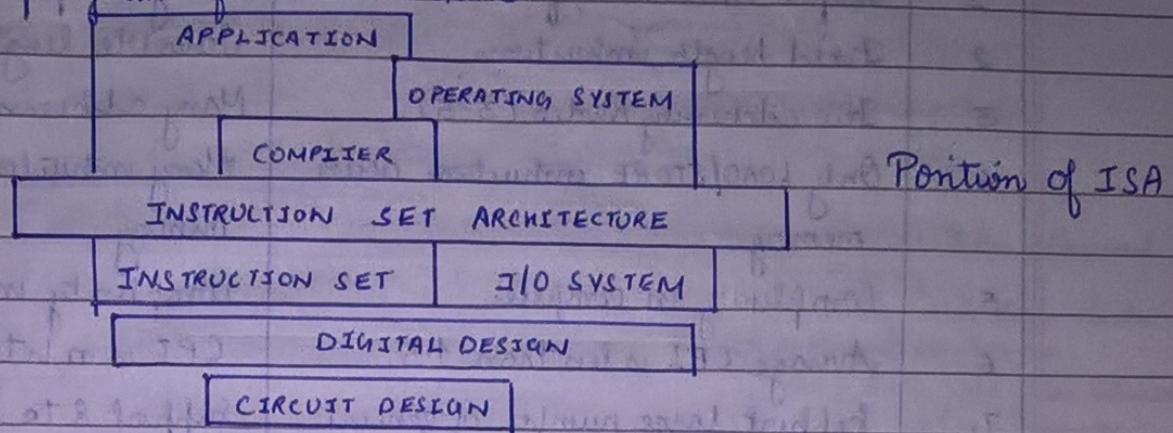


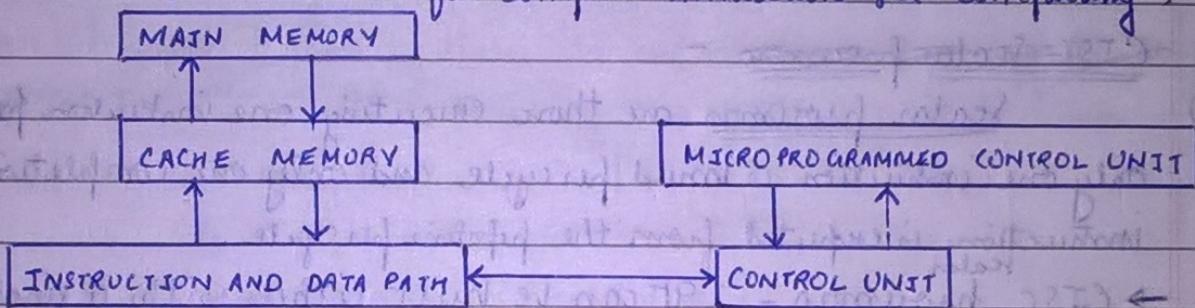
INSTRUCTION SET ARCHITECTURE

① (ISA) Instruction set architecture is the structure of a computer that a machine language programme (or a compiler) must understand to write a correct program for that machine.



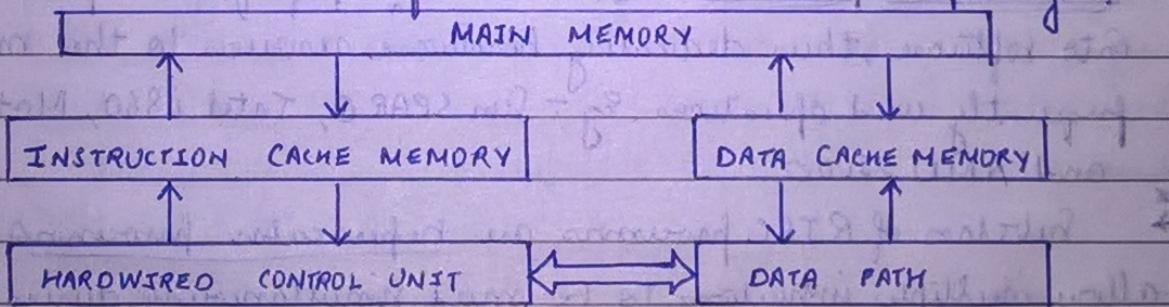
→ CISC architecture with microprogrammed control unit -

CISC stands for complex instruction set computing.

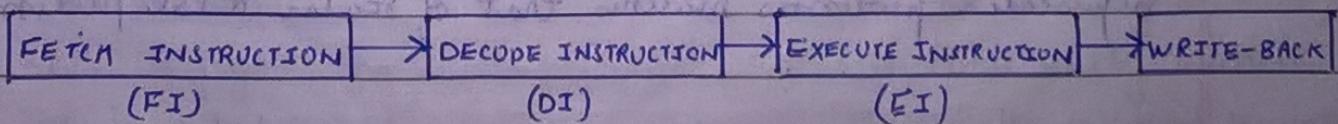


→ RISC architecture with hardware control unit -

RISC stands for reduced instruction set computing.



→ Stages in instruction pipeline



② Comparison between RISC and CISC -

S.No.	RISC	CISC
1.	Simple instructions, few in numbers	Many complex instructions
2.	Fixed length instructions	Variable length instructions
3.	Few addressing modes (3-5).	Many addressing modes (12-24)
4.	Only LOAD/STORE instructions access memory	Many instructions can access memory
5.	Complexity in compiler	Complexity in microcode
6.	Average CPI is less than 1.5	CPI is in between 2 and 15
7.	Support large number of general purpose registers	Support 8 to 24 number of general purpose registers
8.	Highly pipelined	Highly pipelined

③ CISC scalar processors -

Scalar processors are those executing one instruction per cycle, only one instruction is issued per cycle and only one completion of instruction is expected from the pipeline per cycle.

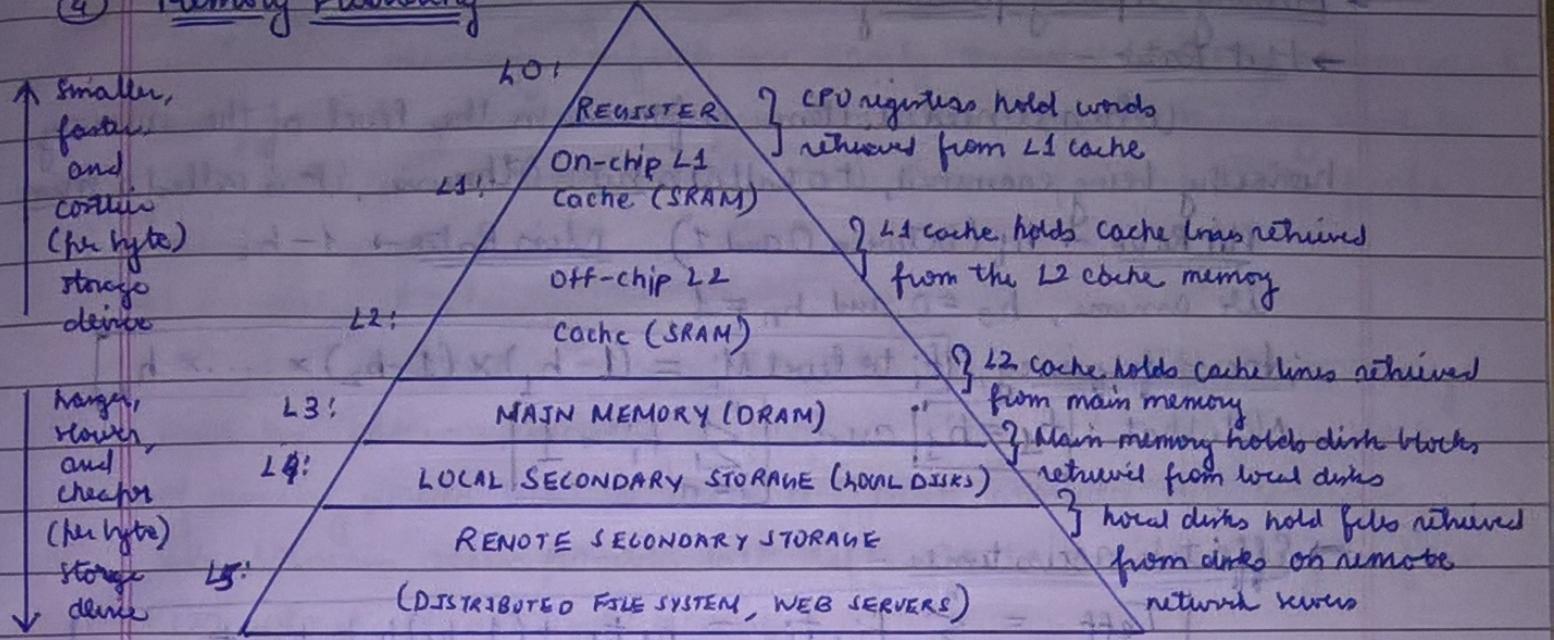
→ CISC processors - It can be built either with a single chip or with multiple chips mounted on a processor board.

Eg - Intel 386 and 486, M68040, VAX 8600 etc.

→ RISC scalar processors - It moves ^{least} frequent operations into software, thus dedicating hardware resources to the most frequently used operations. Eg - Sun SPARC, Intel i860, Motorola M8800 and AMD 29000.

→ Subclass of RISC processors are superscalar processors, which allow multiple instructions to be issued simultaneously during each cycle. Thus, the effective CPI of a superscalar processor should be lower than that of a scalar RISC processor. The clock rate of superscalar processors matches that of scalar RISC processors.

④ Memory Hierarchy



Organizing memory and storage systems is known as memory hierarchy.

→ Memory hierarchy stores the information it satisfies three properties -

(1) Inclusion - It is stated by the following set inclusion relations among n memory levels.

$$M_1 \subset M_2 \subset M_3 \subset \dots \subset M_n$$

$M_0 \rightarrow$ Cache memory, (lowest level)

$M_n \rightarrow$ contains all of the information words stored (highest level).

(2) Coherence - Copies of the same information item at higher levels of the memory hierarchy be consistent. Two strategies to maintain -

- Write-Through (WF) - Update data immediately to all higher levels
- Write-Back (WB) - Delay the update of copies at higher levels until the data being modified in the lower level is replaced or removed from that level.

(3) Locality - Three dimensions of the locality principle are -

→ Temporal Locality - If a program accesses one memory address, there is a good chance that it will access the same address again.

→ Spatial Locality - If a program accesses one memory address, there is a good chance that it will also access other nearby addresses.

→ Sequential Locality - Execution of program that follows certain sequential order.

→ Memory capacity planning -

→ Hit Ratio -

When a needed item is found in the level of the memory hierarchy being examined, it is called a hit. Otherwise, it is called a miss.

Hit Ratio $\rightarrow h_i$ (between 0 and 1) Miss Ratio $\rightarrow 1 - h_i$

We assume, $h_0 = 0$ and $h_n = 1$.

Access frequency, f_i to level $M_i = (1-h_1) \times (1-h_2) \times \dots \times h_i$

Note that $f_1 = h_1$ and $\sum_{i=1}^n f_i = 1$.

→ Effective access times -

$$T_{eff} = \sum_{i=1}^n f_i \cdot t_i = h_1 t_1 + (1-h_1) h_2 t_2 + \dots + (1-h_1)(1-h_2) \dots (1-h_{n-1}) h_n t_n$$

→ Hierarchy optimization

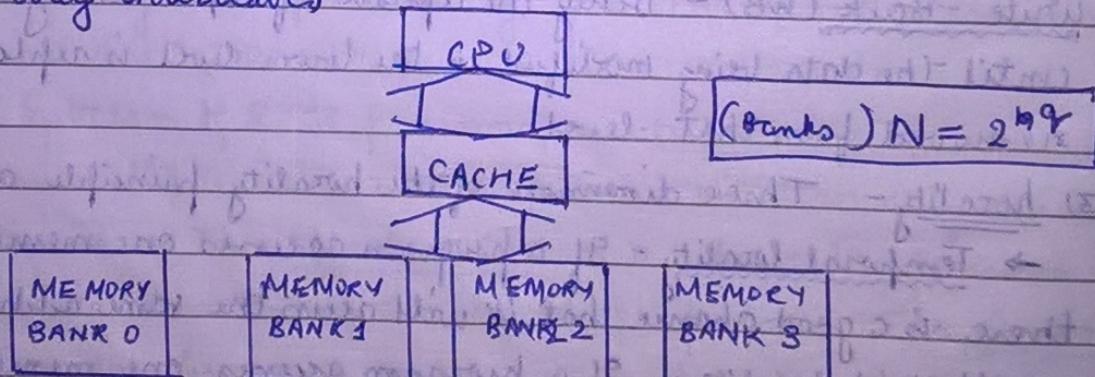
Total cost of a memory hierarchy is estimated as follows -

$$C_{total} = \sum_{i=1}^n C_i S_i \quad C \rightarrow \text{cost}, S \rightarrow \text{capacity}$$

Since $C_1 > C_2 > C_3 > \dots > C_n$, we have to choose $S_1 < S_2 < S_3 < \dots < S_n$

⑤ Interleaved memory organization -

In an interleaved memory, the memory is divided into a set of banks. An interleaved memory with n banks is said to be n -way interleaved.



There are two types of interleaving -

- (1) High-order (2) Low-order. (q -bits are lower)
- (q-bits are higher)

→ High-order interleaving -

Memory conflicts are easily avoided. Each processor executes a different program and programs stored in separate memory modules. Interconnection network is set to connect each processor to its proper memory module.

→ Low-order interleaving -

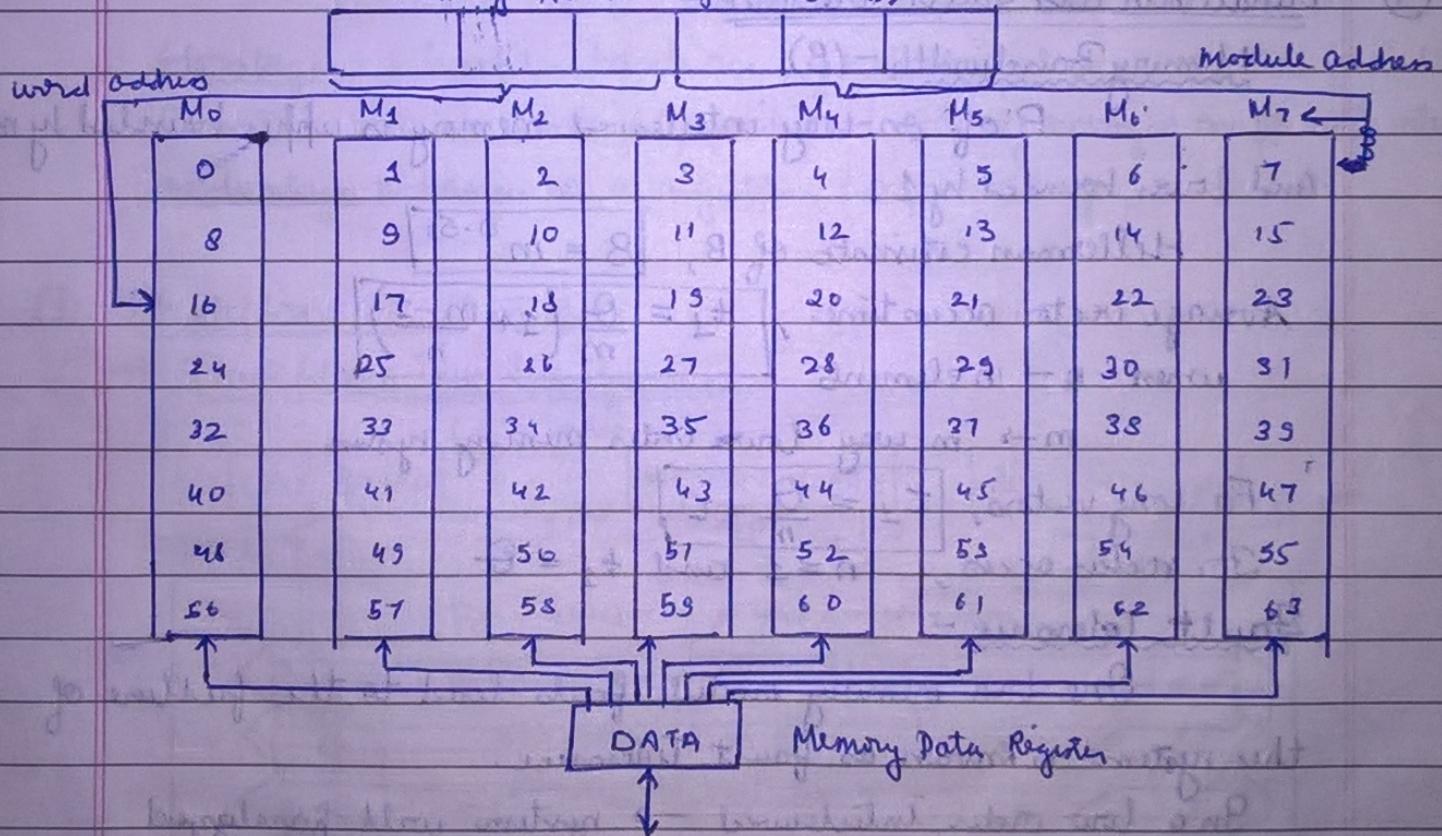
Consecutive memory locations reside in different memory modules. Processors executing a program stored in a contiguous block of memory would need to access different modules simultaneously. Simultaneous access possible but difficult to avoid memory conflicts.

→ Drawbacks of interleaved memory -

- (1) Involves complex design
- (2) Reduced fault-tolerance
- (3) Cannot be expanded conveniently.

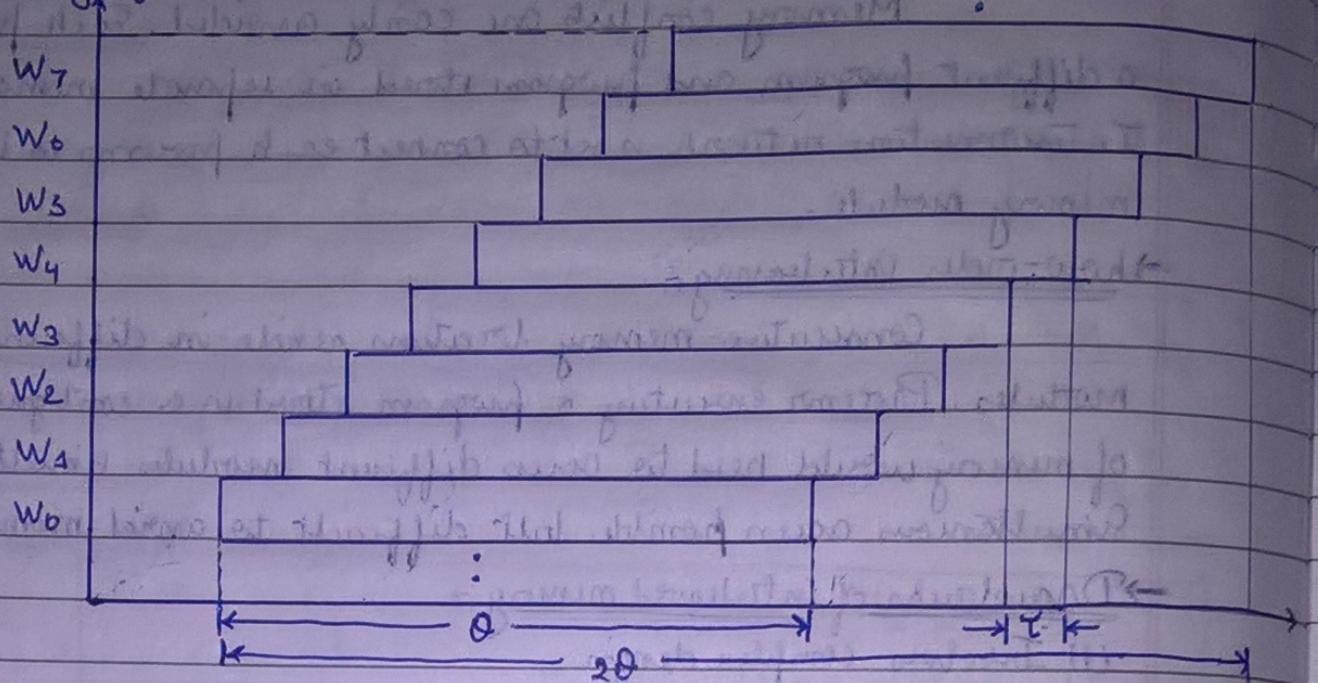
⑥ Pipelined Memory Access -

Memory Address Register (6 bits)



8-Way interleaved memory

Timing of the pipelined access of the eight contiguous memory words:



θ = Major Cycle

T = Minor Cycle

m = Degree of interleaving

$$T = \frac{\theta}{m}$$

⑦ Bandwidth and Fault Tolerance

Memory Bandwidth - (B)

B of m-way interleaved memory is upper bounded by m and lower bounded by 1

Hellerman estimate of B, $B = m^{0.56}$

Average vector access time, $t_1 = \frac{\theta}{m} \left(1 + \frac{m-1}{n} \right)$
where $n \rightarrow n$ elements

$m \rightarrow m$ -way lower order memory system

For long vectors, $t_1 = \frac{\theta}{m} = T$

For scalar access, $n=1$ and $t_1 = \theta$

Fault Tolerance -

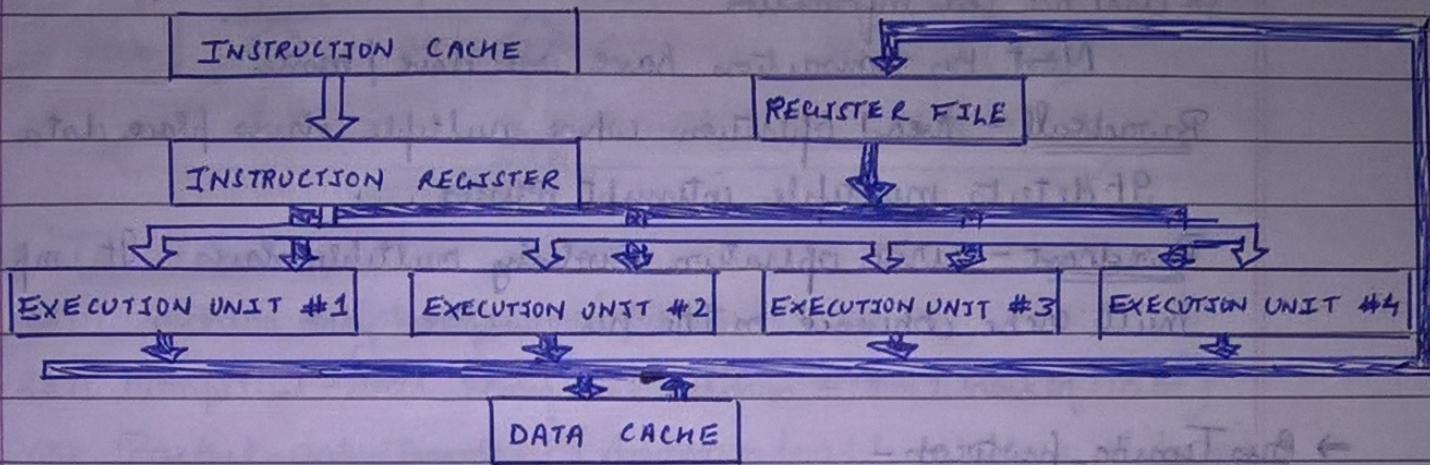
One memory module fails lead to the failure of the system is known as fault Tolerance.

In a low order interleaved \rightarrow system will fail

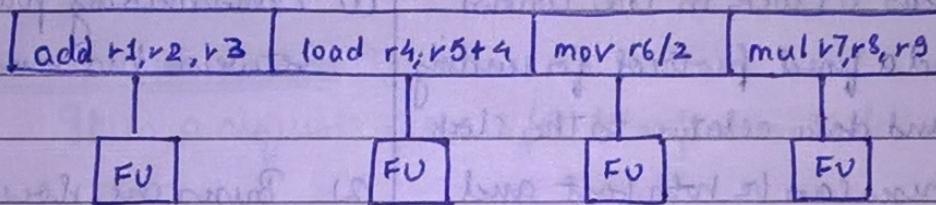
In a high order interleaved \rightarrow system can still be used.

⑧ VLIW architecture - (Very long Instruction word)

In a VLIW machine, many operations (instructions in a normal machine) are encoded in a single instruction. The word of the instruction is very long and can contain multiple operations independent



VLIW instruction execution -

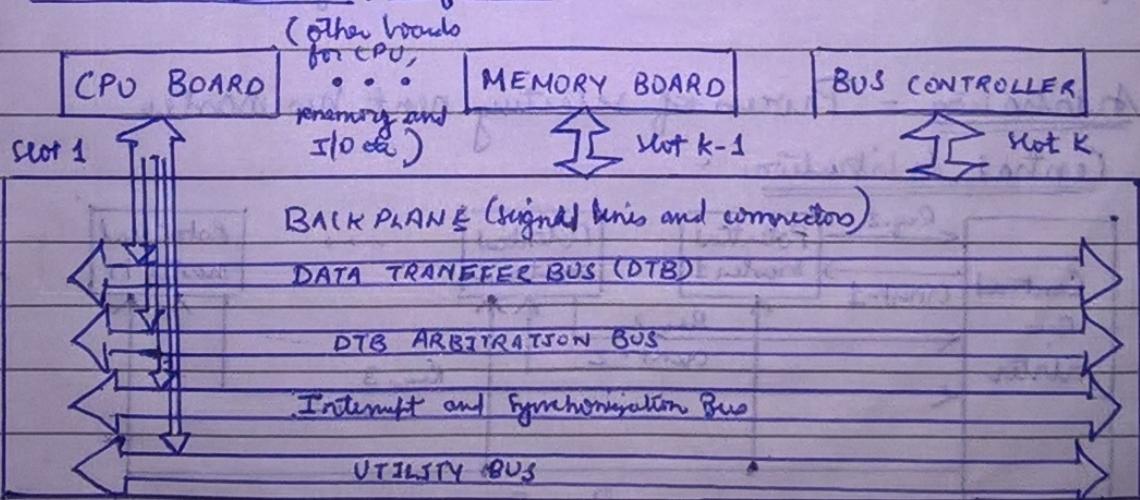


Advantages - simpler hardware, good compilers, compiler friendly, run-time behaviour is highly predictable, easily managed exceptions & interrupts

Disadvantages - large no. of registers, large code size

⑨ Backplane Bus System -

→ Backplane Bus Specification -



→ Addressing and Timing Protocols -

Bus Addressing -

The backplane bus is driven by a digital clock with a fixed cycle time called bus cycle. Backplane has limited physical size, so will not skew information.

Most bus transactions have one slave / master.

Roundcall - Read operation where multiple slaves place data on bus. It detects multiple interrupt sources.

Broadcast - Write operation involving multiple slaves. It implements multi-cache coherence on the bus.

→ Bus Timing protocols -

Synchronous Bus

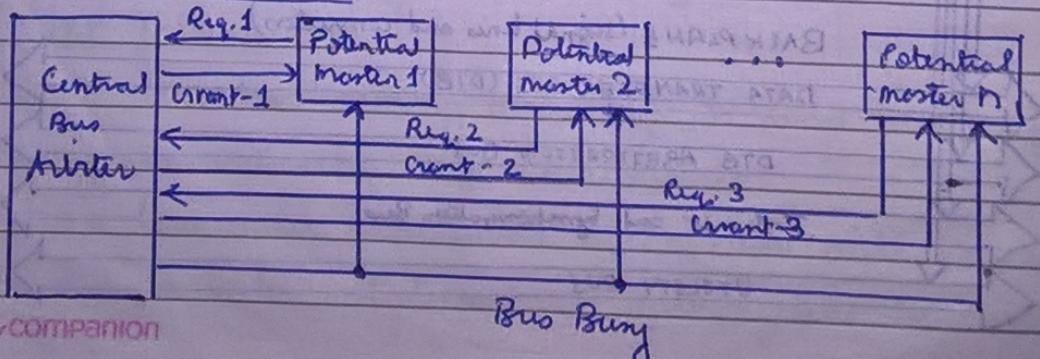
- (1) It includes a clock in the control lines and a fixed protocol for sending address and data relative to the clock.
- (2) These buses can be both fast and slow.
- (3) Buses cannot be long.
- (4) Everything on a bus must run at the same clock rate.
- (5) CPU-memory buses are typically synchronous.

Asynchronous Bus

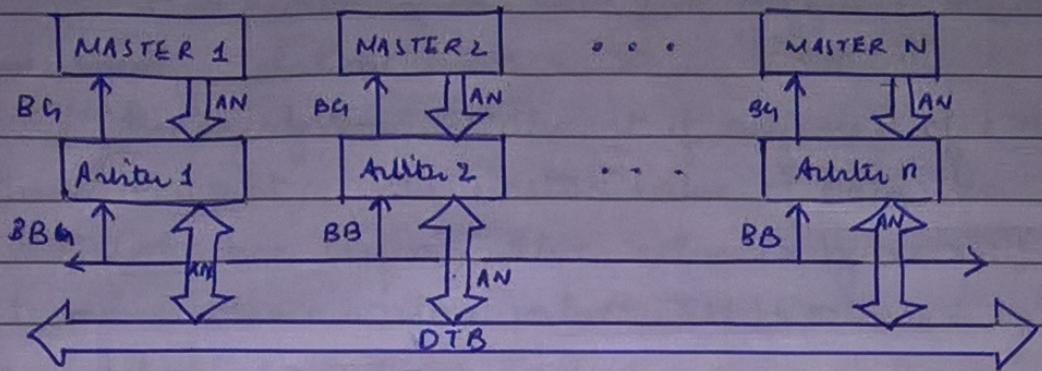
- (1) Handshaking protocols are used between bus sender and receiver.
- (2) Buses are slow.
- (3) Buses can be long enough.
- (4) Same clock is not used.
- (5) I/O buses are more likely to be asynchronous.

→ Arbitration - Process of selecting next bus master.

Central Arbitration -



Distributed (Decentralized) arbitration - priority based arbitration is used



→ Transaction modes -

- (1) Address only transfer (Address)
 - (2) Compelled-data transfer (Address + 1 or 2 block of Data)
 - (3) Packet-data transfer (Address + fixed length block of data)
- A bus transaction consists of a request followed by a response

→ Interrupt Mechanism -

It is a request from I/O to a processor for service or attention.