# Large Computer System

# Large Computer System:-

Forms of parallel processing, Array Processor, The Structure of General-purpose, Inter Connection Networks.

## Forms of Parallel Processing:-

Many oppurtunities are available for parts of a given Computational task to be excecuted in parallel. For ex:- in handling I/O operations, most Computer Systems have hardware for that performs direct memory access (DMA) between an I/o device and main memory. The transfer of data in either direction between the main memory and a magnetic disk under the direction of DMA Controller that Operates in parallel with the processor.

When a block of data is to be transferred from disk to main memory, the processor indicates the transfer by sending instructions to the DMA Controller While the Controller transfers the required data using Cycle Stealing, the processor Continues to perform Some Computation that is unrelated to the data transfer. When the Controller Completes the transfer it sends an interrupt request to the processor to

Signal that the requested data are available in the S
main memory. In response, the processor switches to a
computation that uses the data.

## Classification of Parallel Structures :-

A general Classification of parallel processing has
been proposed by "Flynn". In this classification. a
Single -Processor Computer System is called 'Single
Instruction Stream', Single Data Stream [SISD] System.
A program excecuted by the processor Constitutes the.
Single instruction Stream. In the Second Scheme, a
Single Stream of instruction is broad cast to a no.of
Processors.

Each processor Operates on its Own data. This
Scheme in which all processors excecute the Same progra
but operates on different data is Called "Single
Instruction Stream Multiple Data Stream" [SIMD] System.

The multiple data streams are the Sequences of
data items accessed by the individual processor in
their Own memories. The third Scheme involves a no.of
independent processors, each excecuting a different
program and accessing in Own Sequence of data
items. Such machines are Called "Multiple Instruction

Stream, Single Data Stream"[MISD]. System. In such a system, a Common data structure is manipulated by Seperate processors, each excecuting a different program.

## Array Processors :-

The SIMD form of parallel processing also Called "array processing" was the first form of parallel processing to be studied and Implemented. In early 1970's a system named ILLIAC-IV was designed at the university of Illinois using the approach and was later built by Bro Burroughs Corporation. A two-dimensional grid of processing elements excecutes an instruction stream that is broad cast from Central Processor. Control processor. As each instruction is broad cast, all elements excecute it Simultaneously. Each processing element is Connected to its four nearest neighbours for purpose of exchanging data.
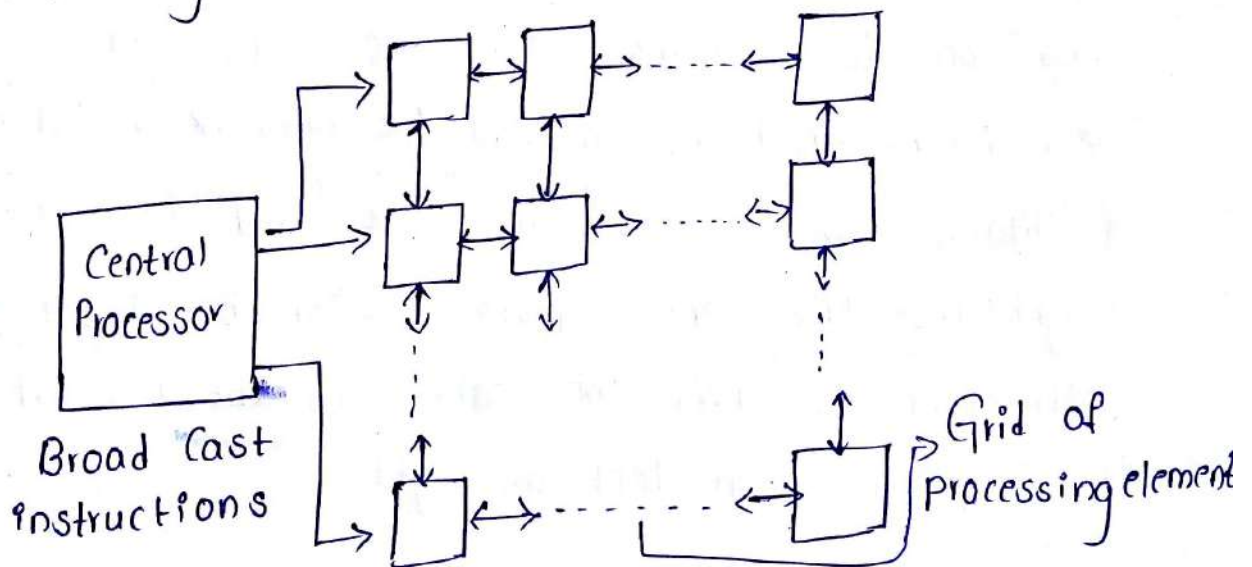


fig :- An array Processor

End-arround connections may be provided in both rows and columns. The grid of processing elements can be used to solve two-dimensional problems. Assume that the edges of the plane are held at some fixed temperatures. The outer edges are initialized to the specified temperatures. All interior points are initialized to arbitary values, not necessarily the same. Iterations are then excecuted in parallel at each element. Each iteration consists of calculating an improved estimate of the temperature at a point by averaging the current values of its four nearest neighbours. The process stops, when changes in the estimate during successive iterations are less than some predefined small quantity.

In array processor each element must be able to exchange values with each of its neighbours. Over the parts. Each processing element has a few registers and some local memory to store data. It also has a register, which we can call the network register, that facilitates movement of values to and from its neighbours. The Central processor can broadcast an instruction to shift the values in the network register one step up, down, left or right.

Each processing element also contains an ALU to excecute arithmetic instructions broad cast by the Control processor. The Control processor must be able to determine which when each of the processing elements has developed its component of the temperature to the required accuracy.
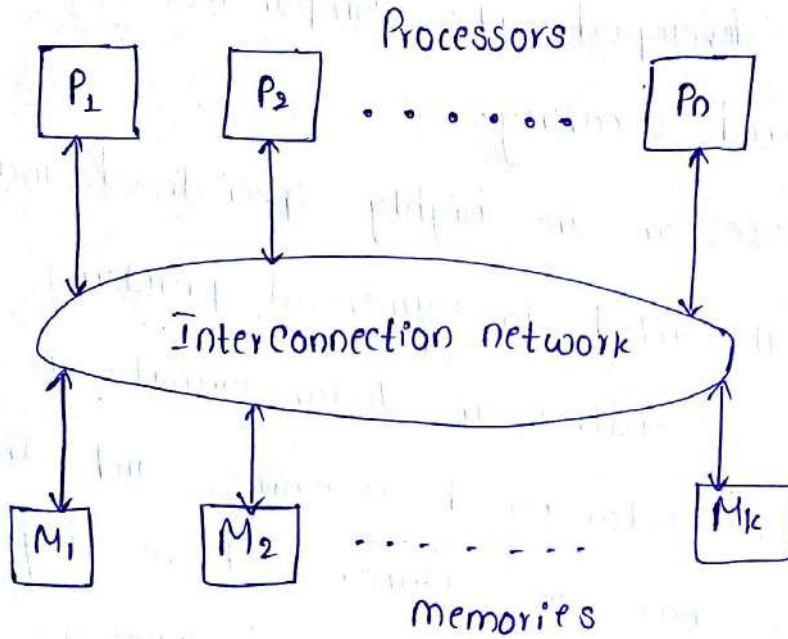
Array processor are highly specialized machines. They are well-suited to numerical problems that can be expressed in matrix or vector format. A key difference b/w vector-based machines and array processors is that the former achieve high performance through a heavy use of pipelining, whereas the latter provide extensive parallelism by replication of Computing modules. Neither array processors nor vector-based machines are particularly useful in speeding up general Computations.
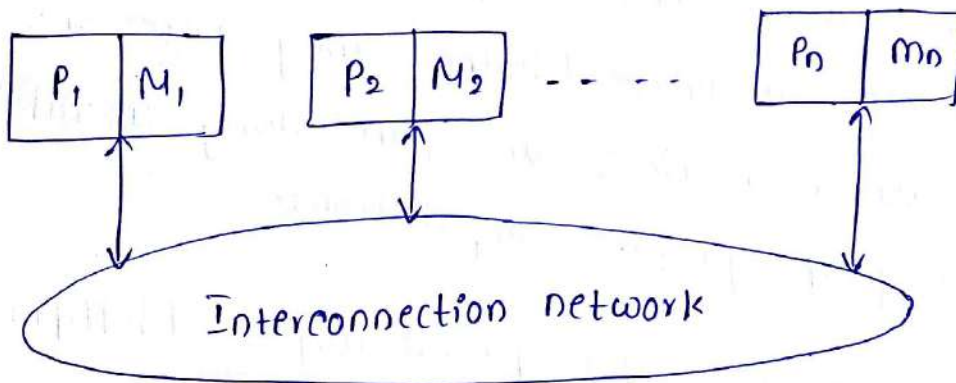
## The structure of General-Purpose Multiprocessors :-

The array processor architecture described in the peceding Section is a design for a Computer System that Corresponds directly to a class of Computational problems that exhibit an obvious form of data parallelism. In more general Cases in

which parallelism is not so obvious, it is useful
to have an MIMD architecture, which involves a no. of
processors capable of independently execcuting differen
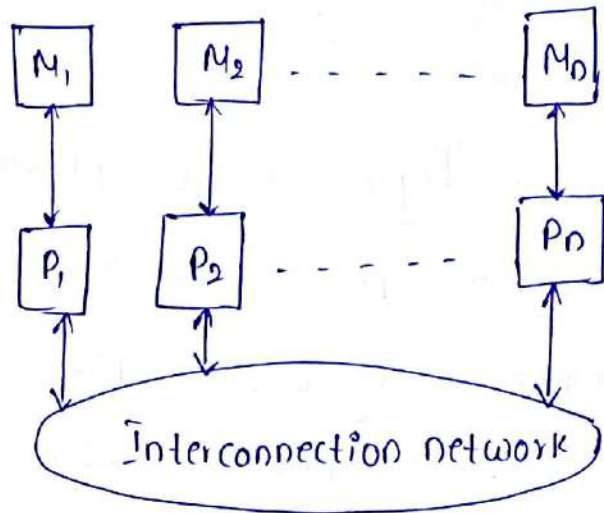routines in parallel.

Processors



Interconnection network

Memories

(a) A UMA multiprocessor



Interconnection network

(b) A NUMA multiprocessor

An interconnection network permits n processor to
access k memories so that any of the processors
can access any time of memories. The inter
connection network and may introduce considerable

delay b/w a processor and memory. If this delays is same for all access to memory which is common for this Organization, then Such machine is called a "Uniform Memory Access" (UMA) multiprocessor.



(c) Distributed memory System

Because of extremely short instruction excecution times achievable by processors, the network delay in fetching instructions and data from the memories in unacceptable if it is too long. Unfortunately interconnection network with very short delays are costly and Complex to implement.

An attractive alternative, which alows a high Computation in all processors, is to attach the memory modules directly to the processor. In addition to accessing its local memory, each processor Can also access, other

memories over network. Since, the remote access process through the network, the accesses take considerably longer than access to the local memory. Because of this difference in access times, such multiprocessors are called "Non-Uniform Memory Access" (NMUA) multiprocessors.

The Organization of figures a & b provides a "global memory" where any processor can access any memory module without intervention by another processor. Here all memory modules serve as private memories for the processors that are directly connected to them. A processor cannot access a remote memory without the cooperation of the remote processor. This cooperation takes place in the form of messages exchanged by the processor. Such systems are called "distributed-memory System" with a message-passing protocol.

Figures a,b,c depict a high-level view of possible multiprocessor Organizations. The performance and cost of these machines depend greatly on on implementation details.

# Interconnection Networks:-

In general, the network must allow information transfer between any pair of modules in the system. The network may also be used to broadcast information from one module to many other modules. The traffic in the network consist of requests (such as R/W) data transfers and various commands.

The suitability of a particular network is judged in terms of cost bandwidth, effective p. throughput and ease of implementation. The term "band width" refers to the capacity of a transmission link to transfer data and is expressed in bits or bytes per second. Then "effective throughput" is the actual rate of data transfer. This rate is less than the available bandwidth because a ginven link usually doesnot carry data all of the time.

Information transfer through the network usually takes place in the form of packet of fixed length and specified format. Longer messages may require many packets.

Ideally, a complete packet would be handled in parallel in one clock cycle at any node or switch

in the network. However, to reduce cost and complexity, the links are often considerably narrower. In such cases, a packet must be divided into smaller piece each of which can be transmitted in one clock cycle.

# Single Bus :-

The simplest and most economical means for interconnecting a no.of modules is to use a single bus. The Since, several modules are connected to the bus and any module can request a data transfer any time, it is essential to have an efficient bus arbitration scheme.

In a simple mode of operation, the bus is dedicated to a particular source-destination pair for the full duration of the requested transfer. The memory module needs a certain amount t of time to access the data, the bus will be idle until the memory is ready to respond with the data.

Suppose the bus tranner takes T times units, and the memory access time is 4T units. It then takes 6T units to complete a read request. Thus the bus is idle for two thirds of time. A scheme

known as "Split-transaction protocol" makes it possible to use the bus during the idle period to serve another request.

Consider the following method of handling a series of read requests. Possibly from different processors. After transferring the address involved in the first request, the bus may be reassigned to transfer the address for the second request. Assuming that this request is to a different memory module. We now have two modules proceeding with read access cycles in parallel. If neither module has finished with its access, the bus may reasigned to a third request and so. on. Address and data transfers for different requests represent independent use of the bus that can be interleaved in any order.

The split-transaction protocol allows the bus and the available bandwidth to be used more efficiently. The performance improvement achieved with this protocol depends on the relationship b/w the bus transfer time and the memory access time. Performance is improved at the cost of increased bus complexity. There are two reasons why complexity increases. Since a memory module needs to know which source

initiated a given read request, a source identification
tag must be attached to the request. Complexity also
increases because all modules, not just the processors,
must be able to act as bus masters.

Multiprocessor that use the split transaction bus
vary in size from 4-32 processors. In larger sizes, the
bandwidth of the bus can become a problem. The
bandwidth can be increased if a wider bus, i.e., a bus
that has more wires is used. Most of the data transferre
between the processor and memory modules consists of
a cache blocks. where a block consists of no. of words.
The Challenge multiprocessor from silicon Graphics
Corporation uses a bus that allows parallel transfer
of 256 bits of data.

The main limitation of single bus that is no. of
modules that can be connected to the bus is not large.
An ordinary bus functions well if no more than 10
to 15 modules are connected to it. Using a wider
bus to increase the bandwidth allows the no. of modules
to be doubled. Networks that allow multiple independent
transfer operations to proceed in parallel can provide
significantly increased data transfer rates.

# Crossbar Networks:-

A Versatile Switching arrangement (is shown in fig). It is known as "Crossbar Switch" which was Originally developed for use in telephone networks. Any module $Q_i$ can be Connected to any other module $Q_j$ by Closing an appropriate Switch. Such networks, where there is a direct link between all pairs of nodes. are Called "fully Connected networks". many Simultaneous transfers are possible. If n Sources need to send data to n distinct destinations, then all of these transfers Can take place Concurrently. Since no transfer is prevented by the lack of a Communication path, the Crossbar is called "non blocking Switch".
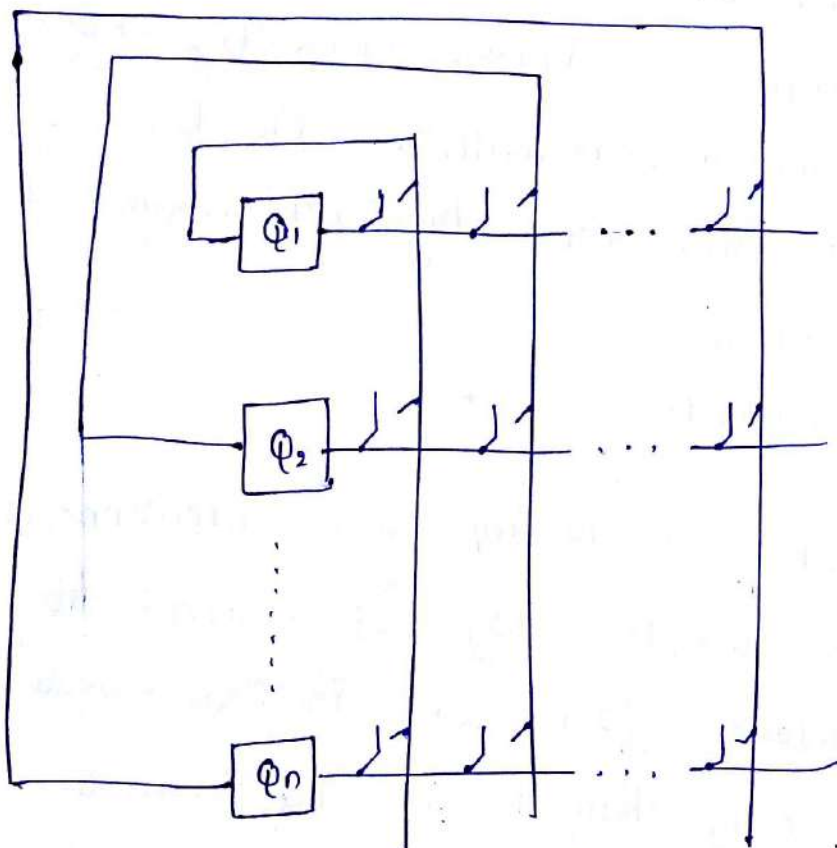


fig. Cross bar interconnection network.

If fig. we show just a single switch at each crosspoint. In actual multiprocessor. however, the paths through the Crossbar network are much wider. The no·o f crosspoints is $n^2$ in a network used to interconnect $n$ modules, the total no·of switches becomes large as $n$ increases. This results in high cost and Cumbersome implementation.

One of the large Crossbar switches is found in Sun's E10000 System in which 16 four-processor nodes are connected by a 16×16 crossbar switch. It is also possible to use a multilevel Crossbar switch, where a Crossbar switch at level 1 connects to a Crossbar swith at level 2 and so on. Such schemes are found in Fujits'u's VPP5000, Hitachi's SR8000, and NEC's Sx-5 machines. A multilevel Crossbar has become a popular choice for a high performance interconnection medium.

## Multistage Networks :-

It is also possible to implement interconnection network that use multiple stages of switches to set up paths between sources and destination. Such networs are less Costly than the Crossbar structure. Yet, they provide a reasonable large no·of parallel

paths b/w Sources and destinations.

The (fig) shows a three-stage network called a "Shuffle network that interconnects eight modules.
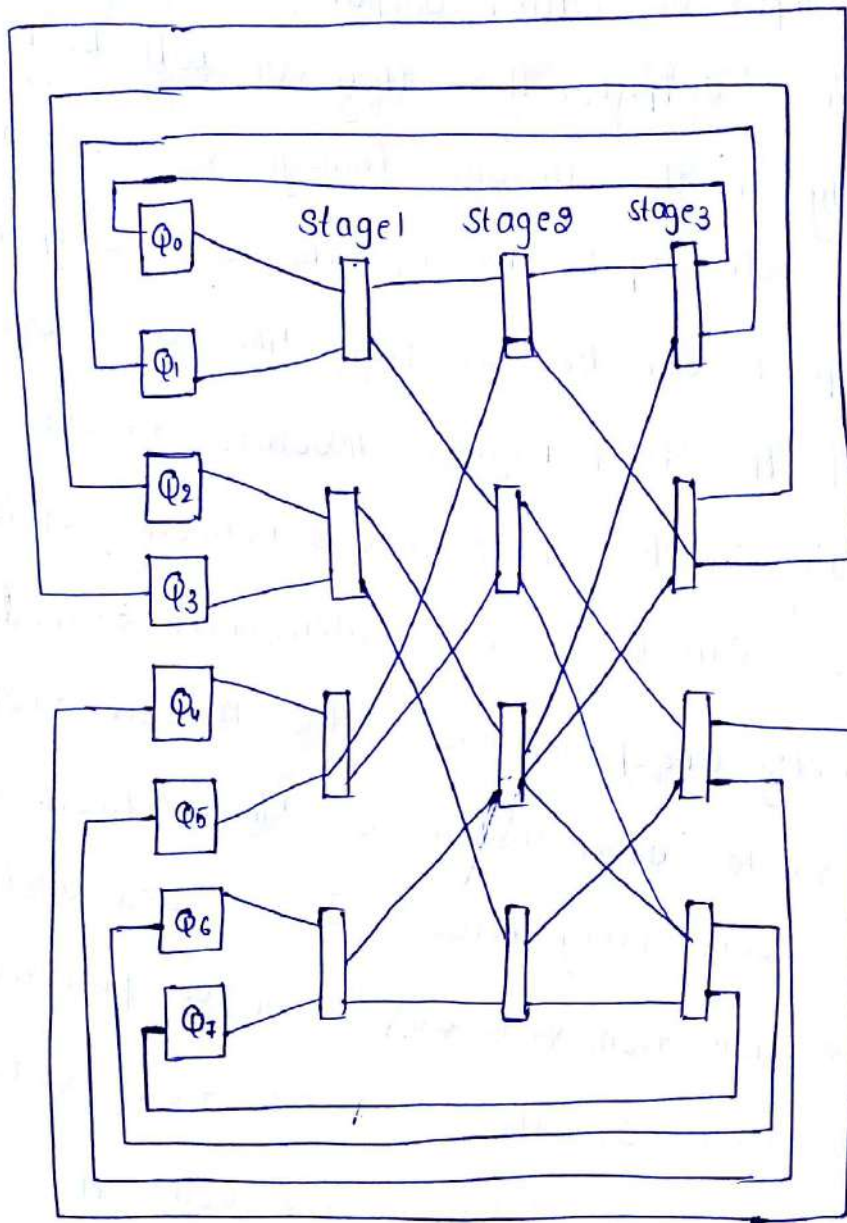


fig: Multistage Shuffle network

The term shuffle describes the pattern of Connections from the outputs of one stage to the inputs of next stage. This pattern is identical to the repositioning of playing Cards in a deck that is shuffled by

Splitting the deck into two halves and interleaving the cards in each half.

Each Switchbox in the fig is a $2 \times 2$ Switch that can route either input or either output. If the inputs request distinct outputs, then they can both be routed simultaneously in the straight through or crossed pattern. If both inputs request the same output. Only one request can be satisfied. The other one is blocked until the first request involves finishes using the switch. It can be shown that a network consisting of 's' stages can be used to interconnect $2^s$ modules.

There is exactly one path through the network from any module $Q_i$ to other module $Q_j$. This network provides full connectivity between sources and destinations

Eg: The connection from $Q_n$ to $Q_u$ cannot be provided at the same time as the connection from $Q_1$ to $Q_5$

If n nodes are to be interconnected using the scheme (in fig) then we must use $s = \log_2 n$ stages with n/2 Switch boxes per stage. Since each switchbox contains four switches, the total no. of switches is

$$4 \times \frac{n}{2} \times \log_2 n = 2n \times \log_2 n$$

which for large networks, is considerably less than $n^2$ switches needed in a Crossbar network.

Multistage networks are less capable of providing concurrent connections than Crossbar switches, but they are also less costly to implement.

## Hypercube Networks :-

The interconnection network imposes the same delay for paths connecting any two modules. Such schemes can be used to implement UMA multiprocessors. The network topologies that are suitable only for NUMA multiprocessor. The first such scheme that gained popular it uses the topology of n-dimensional cube, called a "hypercube". to implement a network that interconnects $2^n$ nodes.

$N_3(011)$   $N_7(111)$

$N_2(010)$

$N_6(110)$

$N_1(001)$   $N_5(101)$
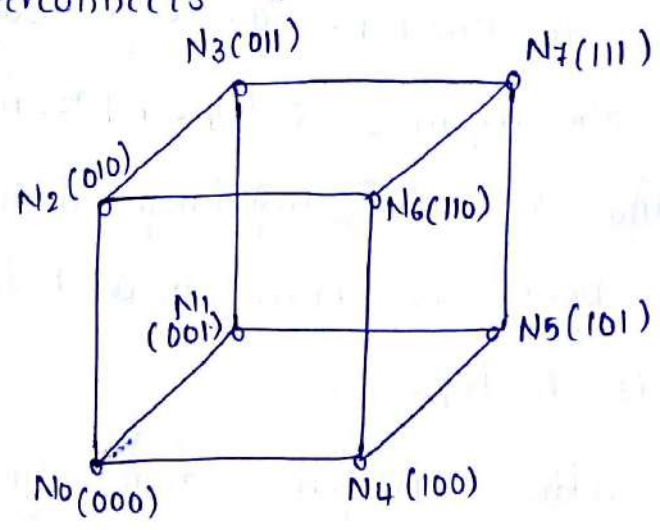
$N_0(000)$   $N_4(100)$

fig: A 3-dimensional hypercube network

The small circles represent the communication circuits in the nodes. The functional units attached to each node. The edges of the cube represent bidirection Communication links between neighbouring nodes. In an n-dimensional hypercube, each node is directly connected to n neighbours. A useful way to label the nodes is to assign binary addresses to them in such way that the address of any two neighbours differ in exactly one-bit position.

Routing messages through the hypercube is particularly easy. The processor at node Ni wishes to send a message to node Nj. The binary addresses of the source i and the destination j are compared from least to most significant bits. The message gets closure to destination node Nj with each of these hops from one node to another. Eg: a message from node N2 to node N5 requires 3 hops, passing through nodes N3 and N1. The maximum distance that any message needs to travel in a n-dimensional hypercubes is n hops.

Scanning address patterns from right to left is only one of the methods that can be used to

determine message routing. The existence of multiple paths b/w two nodes means that when faulty links are encountered, they can usually be avoided by simple, local routing decision. If one of the shortest routes is not available, a message may be sent over a long path.

Hypercube interconnection networks have been used in a no.of machines. Eg's include Intel iPSC, which used a 7-dimensional cube to connect upto 128 nodes and NCUBE's NCUBE/ten which had upto 1024 nodes in a 10-dimensional cube.
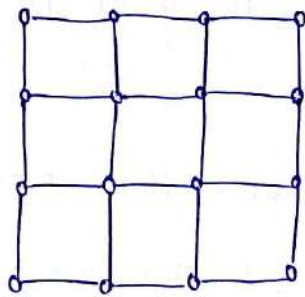
## Mesh Networks :-



Fig: 2-dimensional mesh network

Interconnecting a large no.of nodes is by means of a mesh. Again the links b/w the nodes are birectional. Mesh gained popularity in the early 1990's and essentially displaced hypercubes for interconnection

networks in large multiprocessors.

One of the simplest and most effective possibilities is to choose a path b/w a source node $N_i$ and destination node $N_j$ such that the transfer first takes place in the horizontal direction from $N_i$ toward $N_j$.

Eg's: Intel's paragon and the experimental machine Dash and Flash at stanford University and Alewife at MIT.

If a waparound connection is made b/w the nodes at the opposite edges in the result in a network that consists of a set of bidirectional rings in the x torus, the average latency of information is reduced, but the cost of greater complexity. Such an interconnection network is used in Fujitsu's AP 3000 machines.

Both the regular mesh and the torus schemes can also be implemented as 3-D networks, in which the links are b/w neighbors in the x, y & z directions. Eg's of 3-D torus is found in Cray's T3E multiprocessor.

# Tree Networks :-

A hierarchically structured network implemented in the form of tree is another interconnection topology. The (fig a) shows a 4-way tree that interconnects 16-modules. In this tree, each pattern node f allows communication b/w two of its children at a time. An intermediate-level node can provide a connection from one of its child nodes to its parent. This enables two 'leaf' nodes that are any distance apart to communicate.
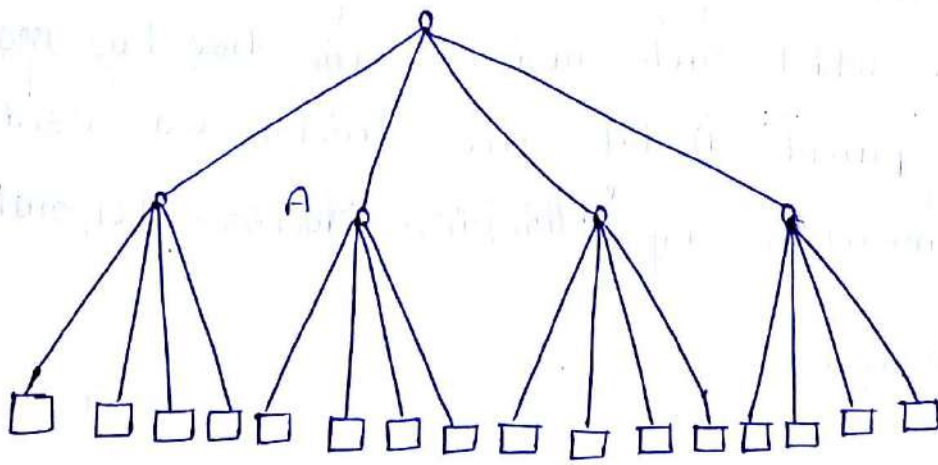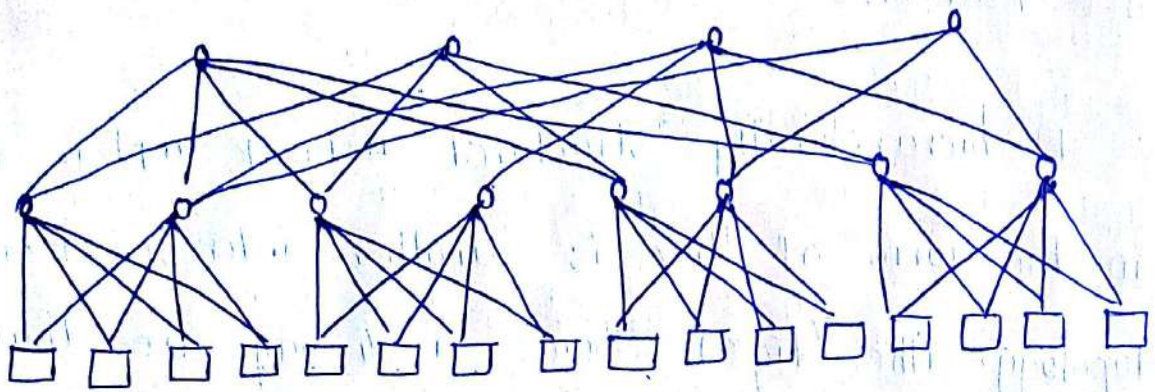


fig: A four-way tree.

A tree network performs well if there is a large amount of locality in communication that is only a small portion of network traffic goes through the single rote noode.
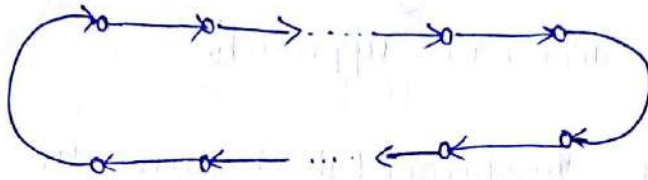
(b) fat Tree

If this is not these case, performance detoriorates rapidly because the root node becomes a bottleneck.

To reduce the possibility of a bottleneck the no.of links in the upper levels of a tree hierarchy can be increased. This is done in a "fat tree" network, in which each node in the tree has more than one parent. A fat tree structure was used in CM-5 meachine by ".Thinking Machines Corporation."
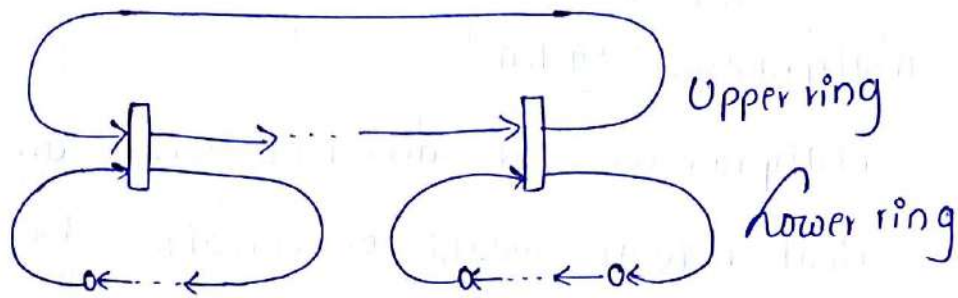
## Ring Networks:-

One of the simplest network topologies is uses a ring to interconnect the node in the system. The main advantage of this arrangement is that the ring is easy to implement. Links in the ring can be wide, usually accommodating a complete packet in parallel, because each node is connected to only

two neighbours. However, it is not useful to construct a very long ring to connect many nodes because the latency of information transfer would be unacceptably large.



(a) Single Ring



Upper ring

Lower ring

(b) Hierarchy of rings

Having short rings reduces substantially the latency of transfers that involves nodes on the same ring. The drawback of this scheme is that the highest-level ring may become a bottleneck for traffic.

Commercial machines that feature ring network include Exemplar v2600 by Hewlett-Packard and KSR-2 by Kendal Square Research.

# Practical Considerations:-

Several different topologies can be used to implement the interconnection network in a multiprocessor system. Each has certain advantages and disadvantages when comparing different approaches.

The most fundamental requirement is that to communication network be fast enough and have sufficient throughput to satisfy the traffic demand in multiprocessor system.

Multiprocessor of different sizes are needed. The ideal network would be suitable for all sizes ranging from just a few processor to possibly thousands of processors. The term "Scalability" is often used to describe the ability of multiprocessor architecture to provide increased performance as the size of system increases, while the increase in cost is proportional to the increase in size.

In addition to providing the basic communication b/w sources and destinations it is useful to have broadcasting capability where a message traverses the entire network and is received by all nodes.

The ability to send a message to only a subset of the network nodes is to also beneficial. Such transfers are called multicasting.

Ideally, the machine could continue to function even if some link in the network fails.

## Meshes and Rings:-

Both mesh and ring network are characterized by point-to-point links which can be driven at high clock rates. Both are viable in small configura- -tions and can be expanded without difficulty. Incremental expansion is smaller in a ring network than in mesh network.
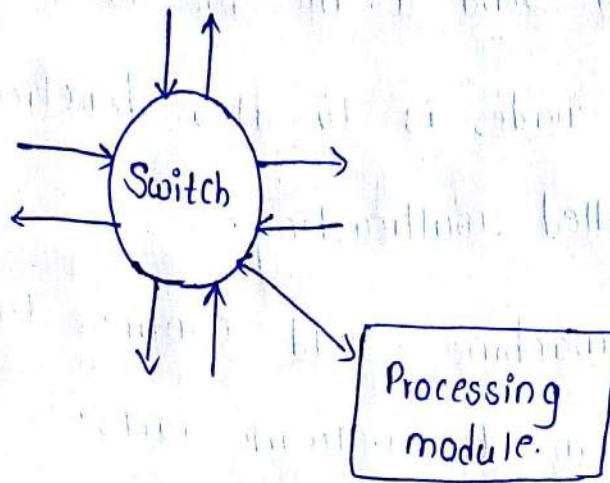
The switch block includes both the circuitry that selects the path for transfer and the buffers needed to hold the data being transferred. Data are transferred from the buffer in one node to the buffer in the next node in one clock cycle.
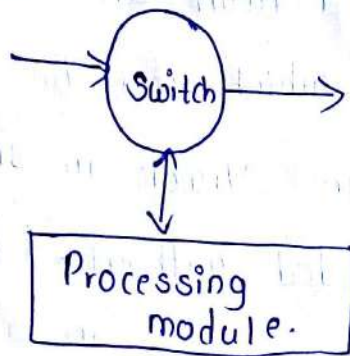
The bidirectional communication is needed in both x & y directions, eight distinct network link must be connected to the node.

(a) node in a mesh



(b) Node in a ring

The width of the links are limited by total no. of wires that can be used, taking into the account the cost and packaging. The term flit [Flow Control digIT] is often used to refer to a portion of the packet that can be accepted by Switching circuitry in the node for forwarding, or buffering in case the forward path is blocked by another transfer.

A straight forward scheme known as the

"store-and-forward" method is to provide a large enough buffer in each node to hold all flits of packets. Thus, an entire packet is transferred from one module node to another where it is stored until it can be forwarded to the next node.

An attractive alternative is the "wormhole" routing scheme in which sequence of flits that constitute a packet can be viewed as a worm that moves through the network. The first flit in a worm contains a header that includes the address of the destination node.

Wormhole routing is an application of a strategy known as Circuit Switching which is a familiar concept from telephone networks, where a path through the network is established when a number is dialed. The Conservation takes place along this path called Circuit. Once a Circuit is established, however, the remaining flits of the packet move toward the destination without experiencing any contention. The stratagies in which an entire packet is buffered at each node as in the store-and-forward method called "packet Switching".

The main disadvantage of a hierarchical ring network that the ring at the top of the hierarchy may become a bottleneck if too many packets need to be transfered over it. This will occur if the locality in comminication is low. The ring base system has hundreds of processors. In Contrast mesh-based system scale well to thousands of processors.

The Ring-based System are easier to implement but do not scale as well as mesh-based System. The mesh Systems are Suitable for use in both Small and very large Systems.

Mixed Topology Networks :-
= = = = = = = = = = = = = =

Designers of multiprocessors System Strive to achieve Superior performance at a reasonable cost. In an effort to exploit most advantageous Characteristics of different topologies, many Sucessful machines feature mixed topologies. Bus and Crossbar are excellent Choices for Connecting a few processors together.

Data General's AV25000 System uses nodes where processors are Connected by the bus.

These nodes are interconnected using a ring network. Hewlett-packard's Exemplar v2600 also uses a ring network to interconnect nodes. Where each nodes has a Crossbar Switching Connect the processors. Compaq's Alpha Server SC uses a fat tree to interconnect the node that Comprice processors Connected by Crossbar Switch.

## Symmetric Multiprocessors :-

A multiprocessor System in which all Processors have identical access to all memory modules and all I/O devices. If any processor can excecute either the Operating System kernel or use programs, the machine is Called a "Symmetric Multiprocessor "[SMP]". This also implies that any processor can initiate an I/O Operation to on any I/O device and it can handle any external interrupt.

SMP's are usually implemented using either a bus or Crossbar network. Often. an smp is used as an node in much large multiprocessor system.