

# UNIT - I

Computer Department  
Subject - High Performance Computing  
Sem - II  
gharun@gmail.com

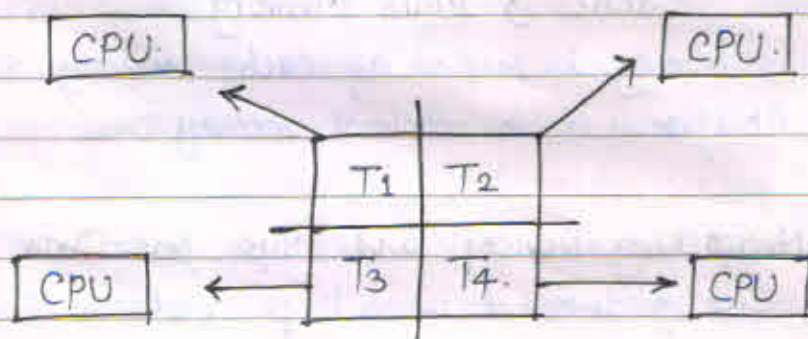
## 1. Parallel Processing Concepts.

### \* Introduction of Parallel Computing :

- parallel computing is basically used to describe about the solving of a single problem using two or more processor.  
e.g. Supercomputer

### - Definition :

"Parallel Computing is a type of computing/computation or the ex.e.g processes are carried out simultaneously. Large problems can often be divided into smaller ones, which can then be solved at the same time."



probl. to be solve.

1<sup>st</sup> Computation of tasks.

### \* Motivation for parallelism :

- 1] Increase nos of transistors in integrated chips enhances Computing power.
- 2] Improvements in storage technology (Memory/Disk.)
- 3] Improve n/wing devices and systems for Data communications.

- 1] Increase nos of transistors in integrated chips enhances Computing power.

→ - The power of CPU increases because of the additional transistors in the ICs.

- The modern processing elements are equipments with multiple units and these unit can be assigned to perform processing task individually.



### 2] Improvement in storage technology (Memory/Disk)

- As we may know, processor alone is not responsible for the enhancement in computation speed but primary memory as well as disk speed also play very significant role.
- As the no. of instr<sup>s</sup> executed in per clk cycle by the processor has increased but at the same time memory should be able to feed the required data for execution.
- The tech<sup>n</sup> called locality of references is used to manage the mismatch bet<sup>w</sup> the memory & the processor speed.



Cache & Main Memory

- The fastest memory is termed as cache memory is used to cope with situation. it is fastest memory.

### 3] Improved Networking Devices and Sys. for Data Comm<sup>n</sup>:

- Now a days, use of Internet as a large platform for IIC<sup>t</sup> and distributed Computation.
- Internet is used to provide the environment for performing large Computation.
- transfer information from one node to another node through use of Internet.

### \* Scope of Parallel Computing\*

- - The parallel computing is suitable for the problems require much more time for Computation completion.
- It is used based on high-end engg & Scientific problem.
- This includes Computer based simulations and Computational Fluid dynamics (CFD) as well as other computer based digital image processing and security algorithm.



## \* Application of H<sup>el</sup> Computing in Engg :

- - parallel computing is basically for speedup of computation.
- This phenomenon is adopted in several engg applications domain such as aerodynamics, optimization algo., like branch & bound and genetic programming, clustering etc.

## \* Application in Scientific area :

- To show the simulated behaviour of real world entities by using mathematics & mathematical formulas.
- The scientific appli<sup>s</sup> are major candidates for the H<sup>el</sup> computation e.g. weather forecasting and climate modeling, oil exploration and energy research, drug discovery & genomic research etc.

## \* Commercial Appli<sup>s</sup> for H<sup>el</sup> computing :

- It requires more processing power in the current market trends because performing many activities simultaneously.
- we can consider multimedia application as commercial application.

## \* H<sup>el</sup> computation application in Computer system :

- The computation perform in H<sup>el</sup>, by the use of collection of low power computing devices in the form of clusters.
- The group of computer connected together to solve complex problem is termed as cluster computing system.
- It is basically applied in the field of computer security.

## \* Parallel Programming Platforms :

- The basic components of the sequential computers are memory unit and the processor.
- In some system, memory unit is connected with the processor thro. the data path.
- The multiplicity is used for performance improvement.
- The multiplicity is achieved by introducing multiple elements of memory unit, processing element and data path.
- There are 2 variations in multiplicity.

One is in implicit H<sup>el</sup>ism where it is not seen by the programmer whereas it can be exposed to the programmer in diff. form.



## \* pipelining and Superscalar Execution

→ pipelining is the process ~~executing~~ <sup>fetching</sup> the ~~current~~ <sup>next</sup> instruction when current instruction is being executed.

- The pipeline processor is based on pipelining used in car assembly line. The whole task to be performed in an assembly line is divided into sub-tasks. It was invented by Henry Ford.
- Henry Ford got an idea that it took much more time to build a car physically but he could actually build a car in a minute.
- He did some step, car was divided into diff. stages on assembly line, one stage is assigned to put engine, another stage is assigned wheels, next stage is used to put cabinet & so on.

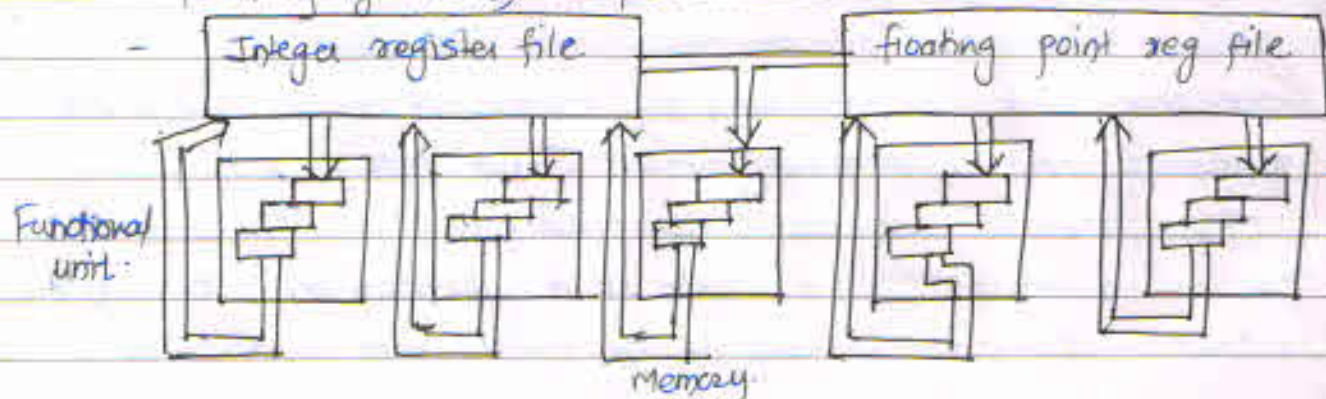
|   | 1        | Unit     | 3   | 4            |                                |
|---|----------|----------|-----|--------------|--------------------------------|
| 1 | Fetch N1 |          |     |              | instr <sup>n</sup> 1 completed |
| 2 |          | Fetch N2 |     |              | instr <sup>n</sup> 2 completed |
| 3 |          |          | ADD |              | instr <sup>n</sup> 3 completed |
| 4 |          |          |     | store result | instr <sup>n</sup> 4 completed |

Instruction Pipeline.

## \* Superscalar Execution

- A processor designed to execute more than one instr<sup>n</sup> at a time during single clock cycle is considered as superscalar execution.
- Superscalar processor fetches & decodes several instr<sup>n</sup> at a time.
- The superscalar architecture exploits the potential of ILP. (Instr<sup>n</sup> level parallelism).

- The typical behaviour of a superscalar archi. is described with the fetching of multiple instr<sup>n</sup> at a time and the attempt for fetching of nearby independent instr<sup>n</sup>.



Superscalar Organization.



- Superscalar m/c is based on von-neumann architecture But it can issue more than one instr<sup>y</sup> per clock cycle.
- A superscalar m/c uses multiple pipelines because with a single pipeline it is not possible to issue multiple instr<sup>y</sup>.
- The classifications of the superscalar processors are based upon the maximum number of instr<sup>y</sup> can be issued at the same time.
- This is depend on the pipeline structure used in the processor.

### \* Very Long Instr<sup>y</sup> Word processors (VLIW) §

- - VLIW (Very-Long Instr<sup>y</sup> word) Architectures are considered as one of the suitable alternatives to achieve Instr<sup>y</sup> level parallelism (ILP) in programs.
- VLIW architecture are used for exe. of more than one basic instr<sup>y</sup> at a time in prog.
- VLIW architecture can store multiple instr<sup>y</sup> in single word. In VLIW based system a H/L compiler is used to generate operation to be executed in H/L in the same word.
- The compiler is responsible for resolving dependancies among instr<sup>y</sup> at compile time.
- VLIW in this architecture indicates that the prog. to be executed in such processors is to be recompiled in a way that the instr<sup>y</sup> runs sequentially w/o existence of stall in the pipeline.
- In this architecture, it is not required for the h/w to examine the instr<sup>y</sup> stream about the instr<sup>y</sup> to be executed in H/L.
- The compiler determines which op<sup>r</sup> to be executed in H/L.
- The performance of VLIW architecture is very good when sequential prog. written in C or Fortran lang are executed after recompilation for such system.
- The VLIW compiler ensures that all op<sup>r</sup> in executing unit can perform simultaneously.
- The VLIW architecture (processors) are used in application areas where high performance is required with less cost.  
e.g. DSP (Digital signal processing)



## \* Basic Working of VLIW Processor : (Principle)

- VLIW processor's basic aim to speeding up computation by exploiting ILP.
- VLIW uses same h/w core as superscalar processors with multiple execution unit (EU) working in parallel.
- In VLIW processor an instr<sup>n</sup> consist of multiple operations in which a typical word length considered from 52 bits to 1Kbits.
- In this, all oper<sup>s</sup> in an instr<sup>n</sup> are executed in lock-step-mode.
- The VLIW processor relies on compiler to find parallelism and schedule dependency free prog. code.

## \* Advantages of VLIW :

- do not need complicated logic to check for dependencies.
- Eliminated Complicated Instr<sup>n</sup> scheduling.
- Compiler is critical to performance.

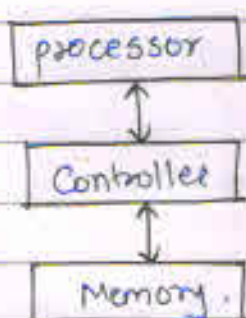
## \* Disadvantages of VLIW :

- Increased code size to empty "slots".
- Increased memory bandwidth.
- Compiler is critical to performance : must do all dependency resolution.
- cache misses in one pipeline will force all pipelines to stall in a "pure" VLIW machine.

## \* Limitation of Memory System Performance :

- As we know, performance of compute not only depend on processing element but memory also does matter.
- The use of cache memory shows a very useful impact on the overall performance of the system.
- A cache is a high-speed memory to be used as the buffer memory. It is logically placed bet<sup>n</sup> CPU & main memory.

- Controller resides bet<sup>n</sup> processor and memory
- The access of the memory is handled by the Controller as an intermediate ~~entry~~ entity.
- The processor send its read/write request to the ~~processor~~ <sup>Controller</sup>.
- The Controller perform the translation of memory addresses and request into appropriate signals for the underlying memory.
- Finally controller passes these signal to the memory chips.



### Controller bet<sup>n</sup> Memory and processor.

- The memory system performance can be assessed by measuring the speed at which a seq. of operation performed by the system.
- Latency is measure limitation of memory. because when we execute new instn<sup>s</sup> thro. memory, then processor should be able to provide speed as fast as needed. otherwise, there can be delay bet<sup>n</sup> processor & memory speed.
- In this, dependencies must be resolving while executing process.
- processor & memory access speed should be compatible.

### \* Use of Caches for Improvement of Latency :

\* Memory Bandwidth.

\* Memory Latency hiding tech<sup>n</sup>.

\* Effect of Multi-threading & Prefetching.



## \* Basic Working of VLIW Processor : (Principle)

- VLIW processor's basic aim to speeding up computation by exploiting ILP.
- VLIW uses same h/w core as superscalar processors with multiple execution unit (EU) working in parallel.
- In VLIW processor an instr<sup>n</sup> consist of multiple operations in which a typical word length considered from 52 bits to 1Kbits.
- In this, all oper<sup>n</sup> in an instr<sup>n</sup> are executed in lock-step-mode.
- The VLIW processor relies on compiler to find parallelism and schedule dependency free prog. code.

## \* Advantages of VLIW :

- do not need complicated logic to check for dependencies
- Eliminated Complicated Instr<sup>n</sup> scheduling
- Compiler is critical to performance.

## \* Disadvantages of VLIW :

- Increased code size to empty "slots"
- Increased memory bandwidth.
- Compiler is critical to performance : must do all dependency resolution.
- cache misses in one pipeline will force all pipelines to stall in a "pure" VLIW machine.

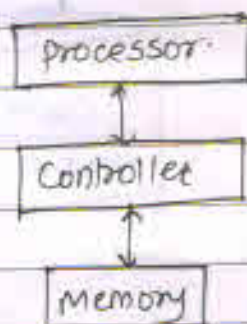
## \* Limitation of Memory System Performance :

- As we know, performance of compute not only depend on processing element but memory also does matter.
- The use of cache memory shows a very useful impact on the overall performance of the system.
- A cache is a high-speed memory to be used as the buffer memory. It is logically placed bet<sup>n</sup> CPU & main memory.



### \* Limitation of Memory System performance:

- The effective performance of a (computer) program on a computer relies not just on the speed of processor, but also on the ability of the memory system to feed data to the processor.
- Execution speed highly depend upon the speed at which instr<sup>s</sup> and data supplied to the processor by memory components.
- A cache is high speed memory to be used as buffer memory. It is logically placed bet<sup>n</sup> CPU & main memory.
- Latency:  
It is the time elapses bet<sup>n</sup> the start of op<sup>n</sup> and complet<sup>n</sup> of operation.  
Latency does not provide complete infor<sup>n</sup> about performance of the memory system.
- The controller resides bet<sup>n</sup> processor & physical memory. The access of memory by the processor is handled by the controller as an intermediate entity.



Controller bet<sup>n</sup> Processor & Memory.

### - Bandwidth:

The rate at which data can be pumped from the memory to the processor determines the bandwidth of memory system.

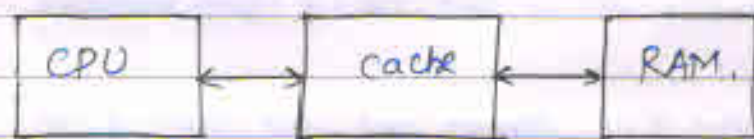
### \* Use of Cache for Improvement of latency:

- The cache memory is used to enhance the access speed of any storage devices  
e.g. disk drive, main memory, tape storage, web server, etc. for other cache also.
- The principle upon which cache works is called locality of references.

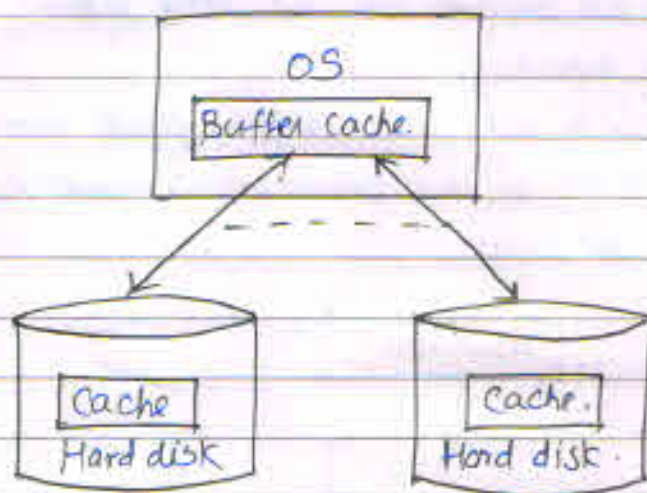


- The Cache are divided into 2 caches - hardware & Software.
- The Cache is memory with low latency & high bandwidth properties.

If one page/appl<sup>y</sup> is already fetched by processor then it is saved in cache memory so next time it can be fetched by cache mem so it improves latency of computer.



Hardware Cache.



Software Cache.

- There are 2 terms used here cache-hit and cache miss. When the requested data references is satisfied by the cache is called cache-hit otherwise cache miss.
- The cache memory comes with different System. Actually amount of data that can be stored in the cache is considered as the capacity of that cache, e.g. 32 KB Cache.

- 1) Cache block - contain multiple byte/words of data.
- 2) Cache set - row in the cache.
- 3) Tag - refer group of data uniquely.



## \* Memory Bandwidth $\frac{B}{s}$ (Double data Rate (DDR))

- Memory Bandwidth is the rate at which data can be read from or stored into a semiconductor memory by a processor.
- Memory Bandwidth usually expressed in units of bytes per second.
- The memory bus & memory unit are used to determine memory bandwidth.
- The memory Bandwidth can be improved by increasing size of memory block.

## \* Memory Latency hiding techniques:

- The increase in memory latency typically occurs when need arises to access remote memory.

e.g. Distributed Shared memory based system.

The important latency hiding techs are as follows:

### 1) Using Prefetching techniques:

- The prefetching is either software/hardware controlled.
- In sw controlled prefetching, explicit "prefetch" instructions are issued for data that is known to be remote.
- In hw controlled prefetching, it is done thro use of long cache line to capitalize on spatial locality or through the instr<sup>s</sup> lookahead.
- The prefetching tech is used for latency hiding because it brings instr<sup>s</sup> or data close to the processor before their actual requirement.
- The direct effect of this scheme is that the time duration bet<sup>h</sup> the issues of instr<sup>s</sup> and its actual references is increased. This is very significant impact when latencies are large.

### 2) Use of Coherent Caching tech

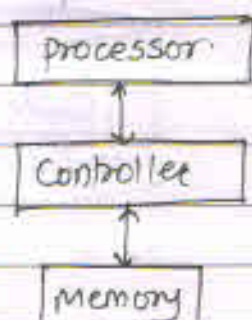
### 3) Relaxing the memory consistency Requirement.

### 4) Using Multiple Context to hide latency.



### \* Limitation of Memory System performance:

- The effective performance of a (computer) program on a computer relies not just on the speed of processor, but also on the ability of the memory system to feed data to the processor.
- Execution speed highly depend upon the speed at which instr<sup>s</sup> and data supplied to the processor by memory components.
- A cache is high speed memory to be used as buffer memory. It is logically placed bet<sup>n</sup> CPU & main memory.
- Latency:  
It is the time elapsed bet<sup>n</sup> the start of op<sup>s</sup> and complet<sup>n</sup> of operation.  
Latency does not provide complete info<sup>n</sup> about performance of the memory system.
- The controller resides bet<sup>n</sup> processor & physical memory.  
The access of memory by the processor is handled by the controller as an intermediate entity.



Controller bet<sup>n</sup> Processor & Memory.

### - Bandwidth:

The rate at which data can be pumped from the memory to the processor determines the bandwidth of memory system.

### \* Use of Cache for Improvement of latency:

- The Cache memory is used to enhance the access speed of any storage devices  
e.g. disk drive, main memory, tape storage, web server... for other cache also.
- The principle upon which cache works is called locality of references.



## \* Parallel Computing platforms :

- To facilitate HET platform, we need to study physical & logical org
- physical org means actual h/w org of the platform
- logical org means: programmer's view of the platform.

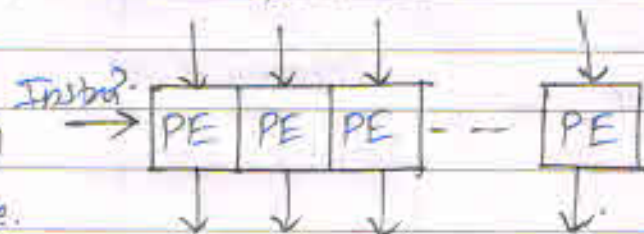
## \* Control structure of HET platform :

- In this, sys. is divided into multiple. Smaller part is granular in structure.
- There are 2 types 1) Coarse grained 2) Fine grained.
- When sys is divided into large nos of small part is called as fine grained.
- Sys. is divided into smaller nos of large parts is called as coarse grained.
- The term granularity is used to describe about the division of a task into nos of smaller subtasks.
- The granularity in HET prog. is considered as different level. e.g. prog. level and instr<sup>n</sup> level.
- There are 2 approaches for working of processing unit in a HET computer.
- First, Single Control unit is used to co-ordinate all processing unit centrally.
- Second, it is based on working of various processing units independently.

## \* SIMD (Single Instr<sup>n</sup> stream & Multiple Data stream) Architecture :

- In SIMD processor, one instruction works on several data item

Simultaneously by using several processing elements (PE's) all of which carry out the same operation as shown in fig



SIMD Architecture.

- The SIMD model sys. uses single control unit to dispatch multiple instr<sup>n</sup> to various processing element.



- In SIMD computing model, single control unit is used to read Instr<sup>s</sup> by single prog counter (PC), decode them and send control signal to the PEs.
- The nos of data paths are depend upon nos of PE.
- Data are supplied to & derived from PE by the memory.
- The diff. unit like PE and memory module are connected by interconnect<sup>n</sup> n/w.
- e.g. exe- g condal statements on SIMD architecture.

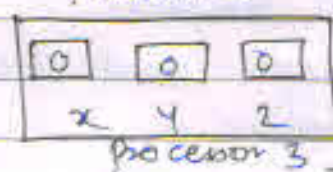
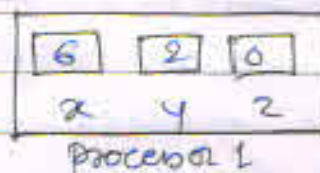
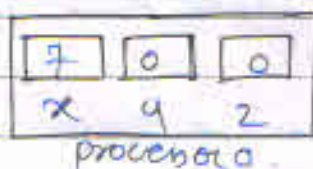
```

if (y == 0)
{
    z = x;
}
else
{
    z = x/y;
}

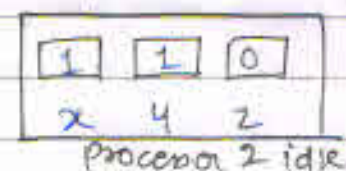
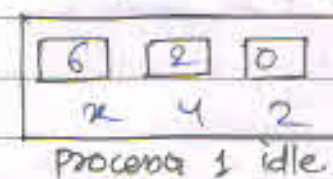
```

Conditional statement

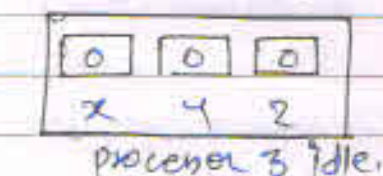
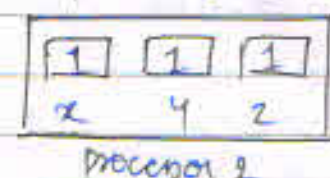
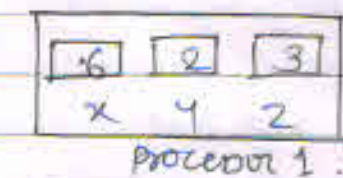
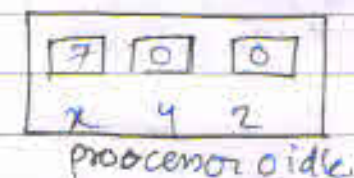
Initial value:



Step ÷ 1



Step ÷ 2

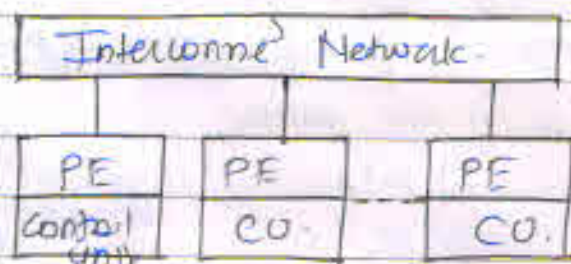


Execution of Conditional statement



### \* MIMD. (Multiple Instr<sup>?</sup> stream Multiple data stream) Architecture:-

- This model represent the sys. capable of executing multiple instr<sup>?</sup> on multiple data sets simultaneously.
- In fact, MIMD is used to describe a  $11^{\text{th}}$  mic able to perform independent computation at same time.
- The MIMD class of mic can execute independent program at the same time. The processing elements included in MIMD class of mic execute different progs at a time.
- This means in MIMD class of mic, each processor fetches its own instr<sup>?</sup> and operates on its own data.



MIMD Architecture.

### \* Comm<sup>?</sup> model for $11^{\text{th}}$ platform:-

- There are diff forms of transfer among nos of  $11^{\text{th}}$  task running concurrently.
- The basic 2 forms are - i) Shared data approach  
ii) message passing approach.

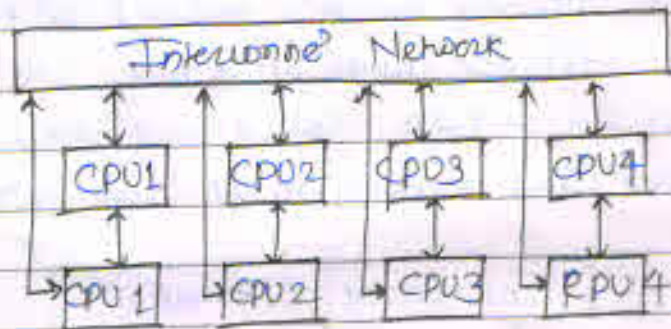
### \* Share-Address space:-

- In this, it provides common data space accessible by all processor included in the system.
- This shared space is used to provide interaction among nos of processor by modifying data objects.
- The sys. is also called as multiprocessor system, bcz it support the  $11^{\text{th}}$  prog<sup>?</sup> approach. term as single program, multiple data (SPMD) prog<sup>?</sup>.
- This type of platform uses separate mem. associated with the processor, or common memory unit globally available for all the processors included in multiprocessing environment.
- The share address space platform is classified as: NUMA & UMA.



## \* NUMA: Non Uniform Memory Access:

- It allows memory access to every processor w/o any restriction.
- A block of memory is attached to the processor and all blocks of memory can be accessed by the path provided by the system's interconnect network.



## NUMA:

- A processor has direct path to the block of memory attached to it.  
e.g. If accessing memory block MEM1 from CPU1 will be much faster than accessing block MEM2 from CPU1.
- It has significant implication:  
i.e. If we map addresses carefully it may be possible to keep most of the info required by a processor in the block attached to it.
- Therefore, the CPU can access that memory directly and reducing the contention for the common bus. Since the time to access a memory location depends on whether it is attached to invoking CPU or not.  
This model is called NUMA.



\* UMA : (Uniform Memory Access) :

- In this, each processor gets equal priority to access the main memory of the m/c is called as UMA.
- In fact, memory access by each processor is identical in UMA platform. as shown in fig.

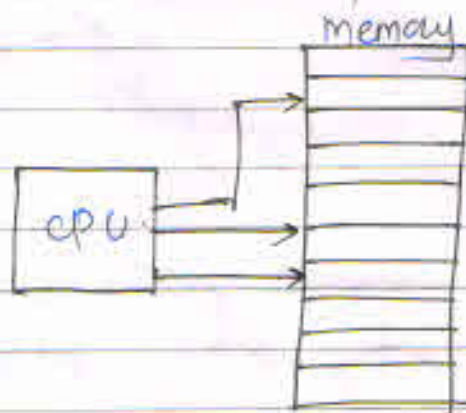


### UMA

- The programming is much easier in such platform because of availability of global memory space in the system.
- The coding for read only instr<sup>s</sup> among the nos of progs. running in different processors is not at all seen by the programmers.
- This happens bcoz coding for such platform is similar to coding usually done in serial prog. for single ~~pro~~ m/c.
- The interaction among the read/write operations are handled. requires the use of mutual exclusion or some other tools for synchronization.

\* physical Org<sup>n</sup> of I/I<sup>l</sup> platforms :

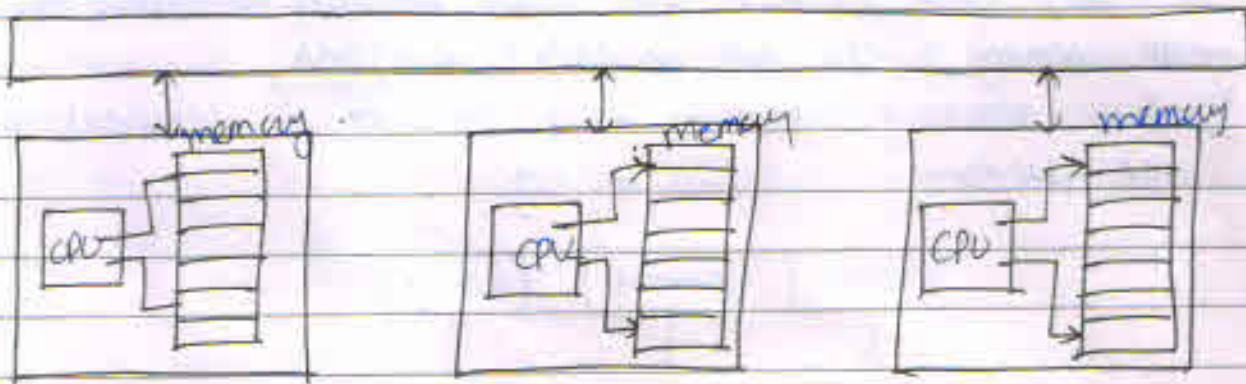
1) Architecture of Ideal I/I<sup>l</sup> computer :



Von Neumann Architecture Machine.

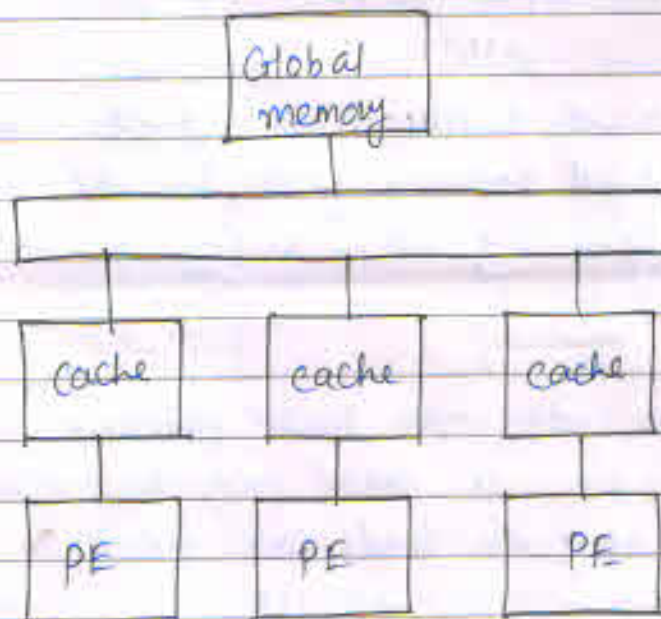


\* Multicomputer based IIC computing m/c model.



Multicomputer based IIC computing m/c model.

\* Ideal Model for IIC Computing:



Ideal Model for IIC Computing System.



## \* Interconnection n/w for IIC Computer :-

→ Def :-

" Parallel Computing system consist of more than one processor and these processing element is connected to memory unit directly or indirectly.

Interconn<sup>n</sup> n/w is needed to route data when processor needs to access memory structure.

- Classification of Interconnection Networks :-

- It is divided into static & dynamic classes.

- The conn<sup>n</sup> bet<sup>n</sup> i/p & o/p nodes are fixed for entire duration of comm<sup>n</sup> is called as static n/w.

e.g. linear array, ring, tree, star, mesh, hypercube etc.

- The conn<sup>n</sup> bet<sup>n</sup> the i/p and output is ~~static~~ not fixed (variable) is called as Dynamic n/w.

e.g. Buses, crossbar switches & multistage n/w, mesh n/w etc.

- In static interconnection n/w, nodes are connected using point to point comm<sup>n</sup> links.

Static n/w is also called as direct network.

- In the other hand, dynamic interconn<sup>n</sup> n/w, is built up using the switches and comm<sup>n</sup> link.

- Dynamic n/w are referred as Indirect n/w.

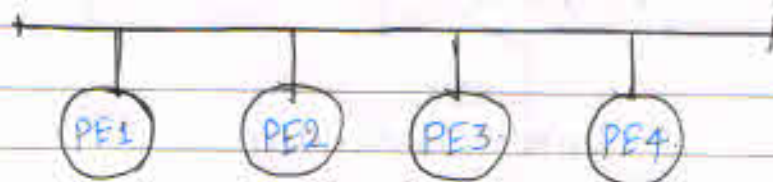
## \* Network topologies :-

- N/w topology are graphical representation of various arrangement of nodes and comm<sup>n</sup> links among nodes.

### \* Bus Based Network :-

- This is simplest type of n/w (interconn<sup>n</sup>) used and static by nature.

- In this n/w, all elements share common comm<sup>n</sup> link as shown in fig.





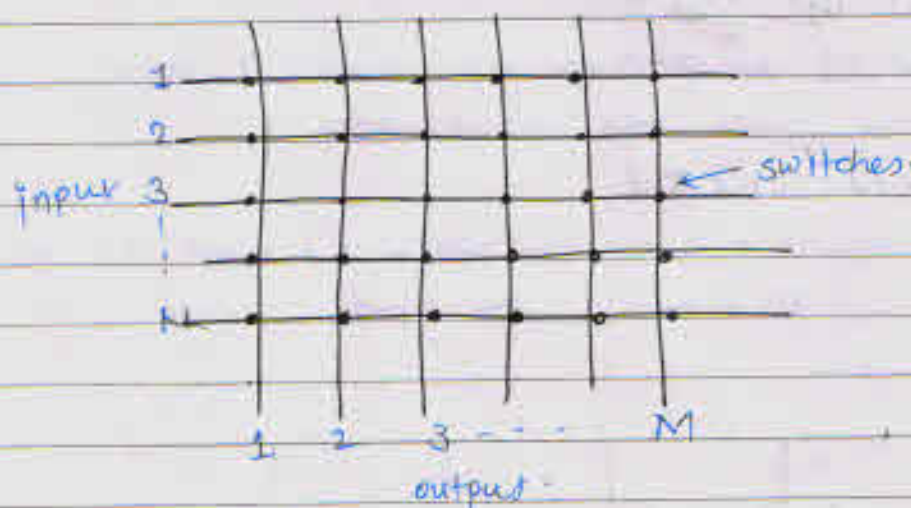
- this is most inexpensive n/w. the nodes can easily be added & deleted from this n/w.
- this n/w always has requirement for handling of conflict when several node elements request for the Bus at the same time.
- The Bus Controller mechanism is used to provide bus access on a FCFS basis.

#### \* Crossbar Network

- This is one of the simplest n/w topology used for interconnect among computing nodes.
- It consists of 2-dimensional grids of switches.
- this n/w provides non-blocking connectivity bet i/p and o/p.
- this type of interconnection n/w provides the characteristics in which any of the i/p can join to any of the o/p.



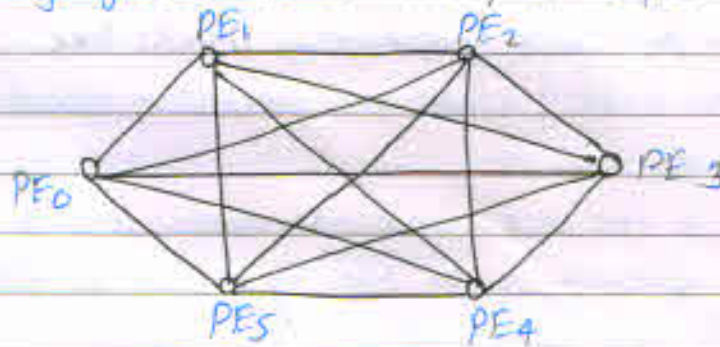
- The conn<sup>o</sup> establishment is done thro. a cross-point for particular row i/p and a particular column o/p.
- The no. of switch requirements for a n/w with N i/p & o/p is  $N^2$ .





### \* Fully-Connected Network†

- This is one of the most powerful interconnect topology used for providing connectivity among nodes.
- In this type of topology, each node is directly connected to all other nodes.
- The shortcoming of this n/w is that, it requires too many connections.



### - Linear Array†

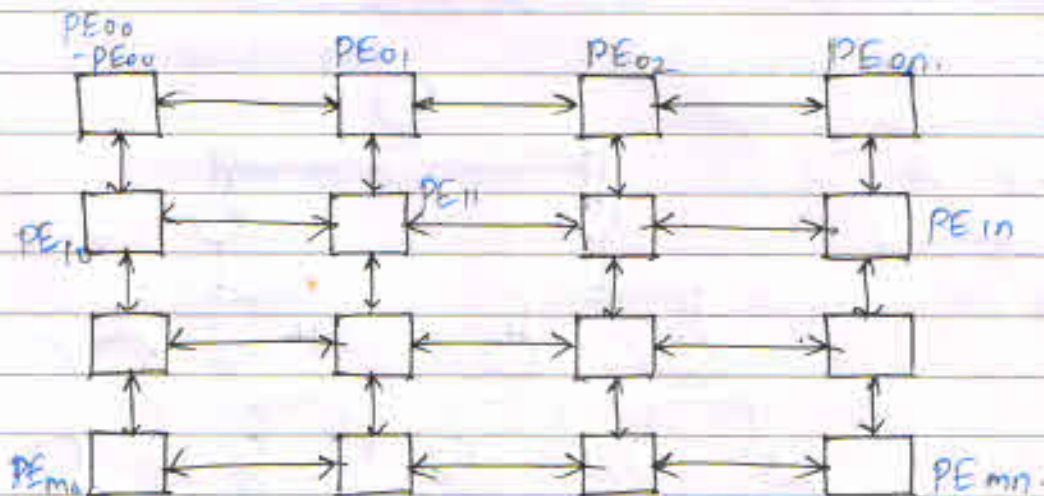
- It is one of the fundamental interconnect n/w pattern.
- In this type of n/w, processors are connected in 1-D linear array.



Linear Array

- This is 1-D n/w, where 1st & last node connected with only one adjacent processor and all the intermediate elements are connected with 2 adjacent processing element.

\* Meshes: A mesh is 2-D n/w. In which processing element are arranged in 2-D grid. The row & col. pos<sup>n</sup> are used to denote a particular processor in the mesh n/w.

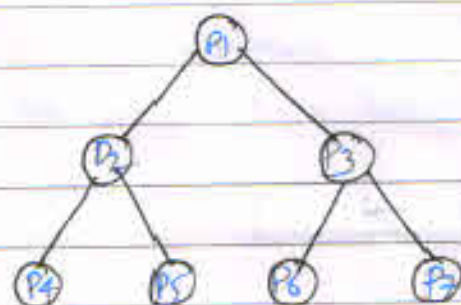


Mesh Network.

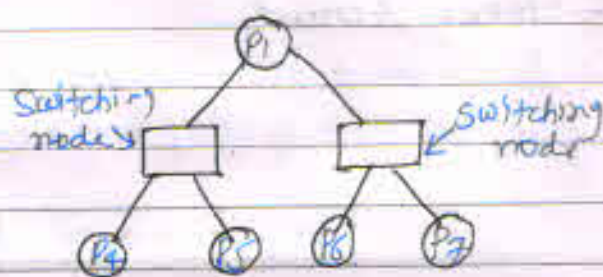


## \* Tree Based Network

- In this, only one path is used bet<sup>y</sup> any pair of nodes.
- In this, linear array of star connected n/w are considered as special case of tree n/w.
- In this n/w, processes are arranged in a completely like Binary-tree pattern.
- It is divided into 2 category. 1) Static tree 2) Dynamic tree.
- When each node of a tree are processing elements we call it as static tree n/w.
- Whereas, in the case, when the nodes included at the intermediate level are switching nodes the tree n/w is considered as a dynamic tree interconnect n/w.



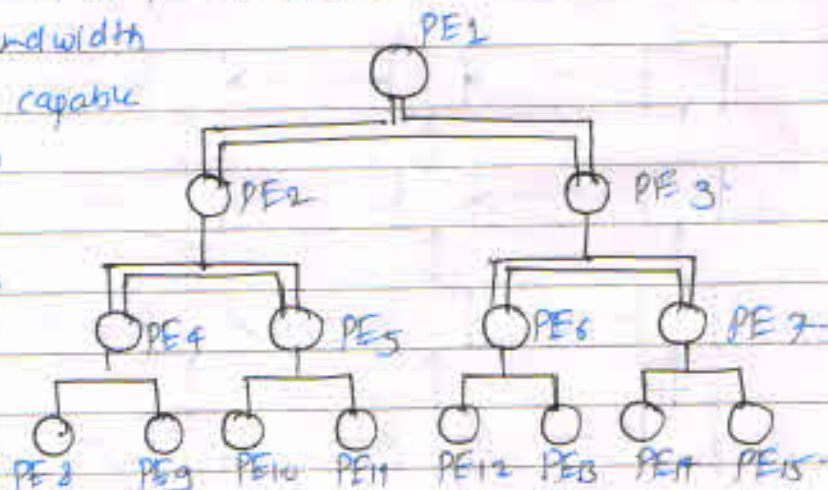
A tree n/w static in nature.



A tree n/w dynamic in nature.

## \* Fat Tree

- This type of n/w is modified form of the original tree n/w.
- This has increased bandwidth of edges into root direction.
- This simulated on the basis of actual trees where branches get thicker towards root. It provides more adaptability bcz practically more traffic occurs towards the root as compared to leaves.
- therefore, the larger bandwidth provided at root side is capable to handle the traffic in effective ways and no bottleneck can occur.

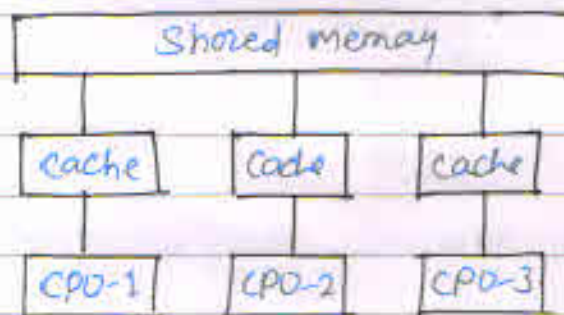


Fat Tree.



## \* Cache Coherence in Multi-processor Systems :

- The shared memory multiprocessor system equipped with a separate cache memory for each processor.
- Such sys keeps many copies of data and instr<sup>s</sup> in the fashion that one copy is kept in memory and one copy in each cache memory.
- when copy is changed by the processor other copies must also reflect with changes.
- Here cache coherency deals with the changes in the shared data should propagate throughout the sys in a timely fashion.



Multiprocessor system with shared memory.

## \* Cache Coherence Problem :

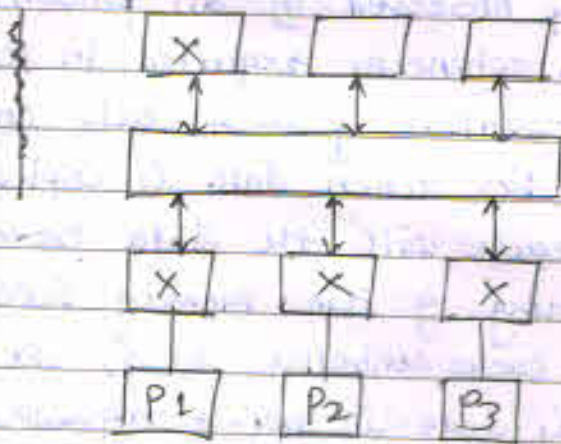
- In the shared memory, private cache is associated with each processor. In such situation each cache contains multiple copies of the same data.
- The problem occurs when all processors are allowed to update the data independently. This is called as cache coherence problem.
- The sys is called coherent only when every read op<sup>n</sup> results in the value which is updated by previous write operation, even by the process of any other processor of that system.
- The sys must be able to maintain a coherent view of memory. It is required to maintain consistency of data in shared resources of individual local cache memory in the system.
- The dealing with this consistency issues is called as cache coherence.
- The coherence problem may happen in multiprocessor system where each processor have local cache as well as can access global shared cache of the system.



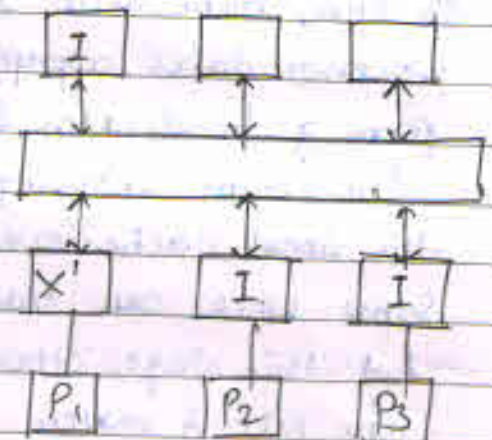
- The basic things included with coherency associated in multiprocessor system are migration and replication.
- The genuine drawback while accessing share data that is not available locally are latency & bandwidth.
- Cache coherence protocol are implemented at h/w level instead of the s/w bcz of the efficiency of the access.
- The cache coherence protocol is used to keep records about the shared data blocks and their associated states.
- Snooping and directory based are 2 types of coherence protocol.

\* Hardware-based protocol :- The basic policies used for maintaining Cache Consistency are Invalidate & Update.

\* Write - Invalidate.

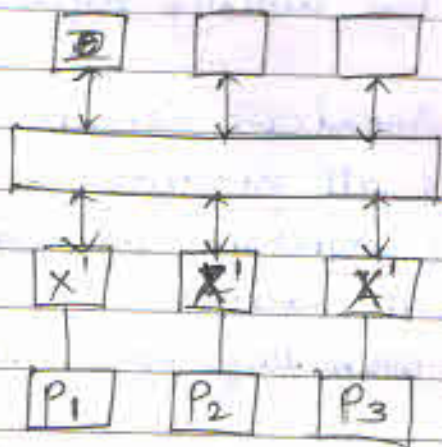


Four copies of block X



Copies of P2 & P3 are invalidated.

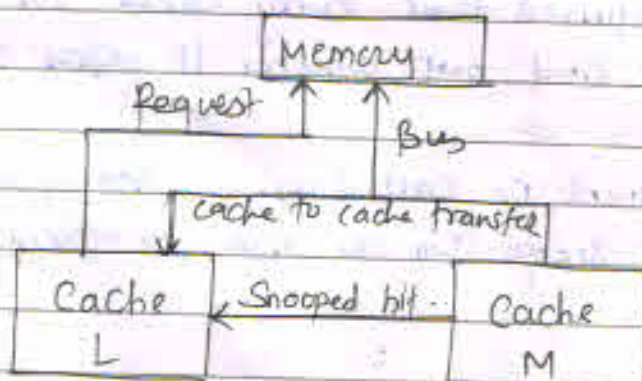
\* Write - Update :-



- 1) Memory copy is updated or not - it depends on imple<sup>2</sup> of protocol.
- 2) In this, if processor modifies data, it reflects to the all processor/system.
- 3) Actually processor cache has ability to broadcast changes to other.

Write update policy

\* Snoopy Cache protocol :-



Cache Coherence.

- It is suited for multiprocessor sys. with shared memory & interconn<sup>2</sup> among the nodes



- It is used in bus based system. It is suited for multiprocessor system.
- In this, mem. trans<sup>n</sup> is carefully observed. by all processor. processor takes appropriate action whenever required in the form of invalidation & update operation of local cache content.
- The term snooping is used bcz when data is copied to the local cache every other cache with the data from same block can track the sharing of the memory block.
- At this stage, other cache or cache controller snoops via the bus & watch about their data is being requested by other cache.
- In diag, cache L & M are connected via a bus.
- The data transfer occur b/w 2 cache. When cache M receive a Snoop hit<sup>n</sup> from cache L, which has actually requested it from memory.
- The shared bus has property to broadcast coherent info<sup>n</sup> in very fast manner among all processors.
- This multiprocessor system strictly maintain consistency of data by updating info<sup>n</sup> among all processors.
- The shared bus become bottleneck when large nos of processor is included.

This situation can be handled by increasing bandwidth.

- In some interconn<sup>n</sup> n/w, like Bus n/w it is feasible to broadcast consistency cmd to all cache.

In this it is required that every cache should execute consistency cmd to find out whether it refers to data in cache.

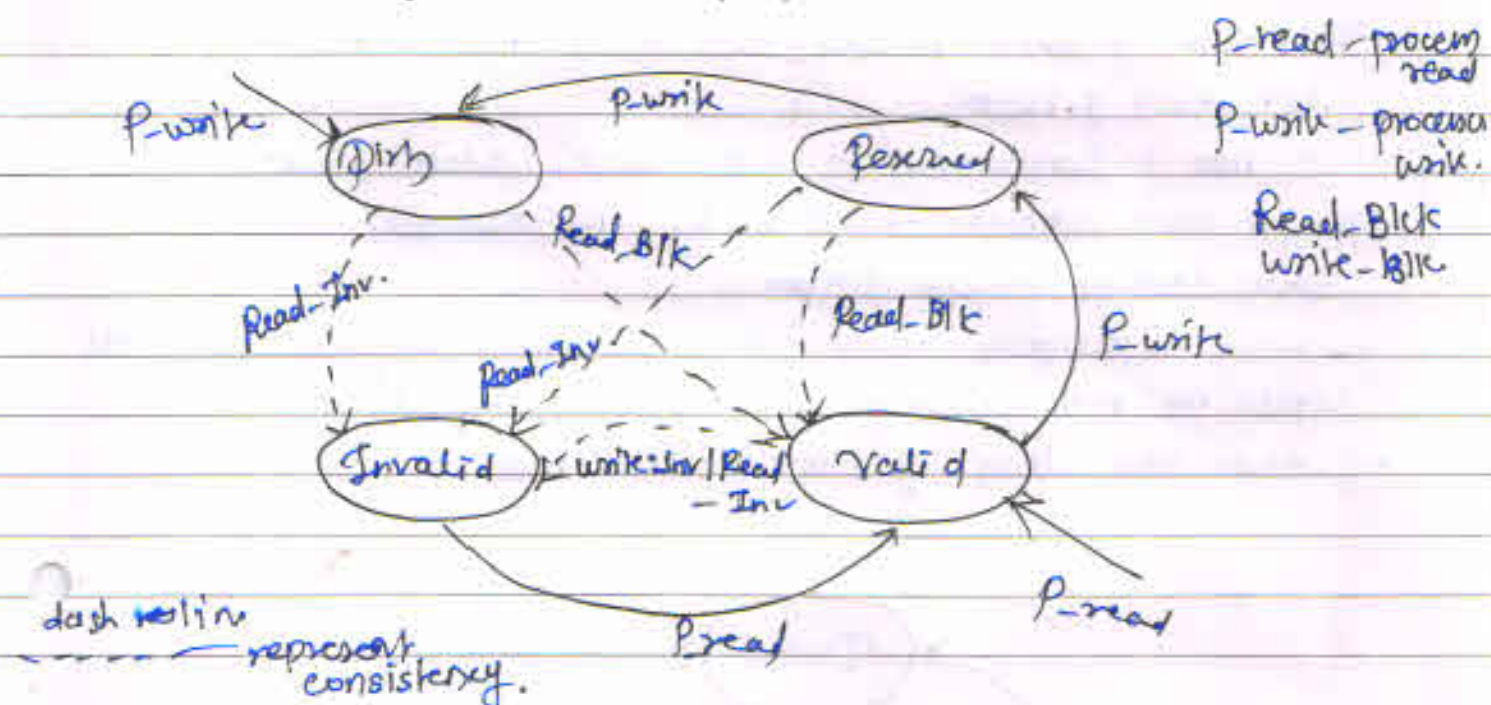
- This type of protocol is called as "Snoopy Cache Protocol" because each cache snoops on the n/w for every incoming consistency cmd.



\* Write-Invalidate Snoopy Cache protocol:

- 1<sup>st</sup> write-Invalidate Snoopy protocol is write-once protocol.
- It provides assumption of state with each copy of block.
  - These states are named as Invalid, Valid, Reserved & Dirty for copy of block.

- 1) Invalid - when copy is inconsistent
- 2) Valid - copy is consistent
- 3) Reserved - when data is written exactly once and copy is consistent with memory.
- 4) Dirty - Copy is only one in the system but data is modified more than once.



Write-Once Protocol state transition diagram of cache copies.

- The copy-back mem update policy is used in write-once.
  - It indicates that entire copy of data blocks must be written back to the memory when it is replaced.
  - This op<sup>n</sup> is performed only when the data has been modified during its cache residence time. (i.e. state is dirty)
  - The normal mem. cmd read block (Read-blk) and write block (write-blk) are applied & besides these cmds other consistency cmd required to maintain consistency are used.
- \* These cmds are:

- 1) Write-Inv - invalidates all other copies of block.
- 2) Read-Inv - used to read a block and all other copies are invalidated.



## \* Write-Update Snoopy Cache protocol

This protocol has 3 state.

### 1) valid-exclusive -

It indicates about only cache copy and this copy is consistent with the memory copy.

### 2) Shared -

The copy of data to be shared is consistent and there are existence of ~~three~~ other consistency copies also.

### 3) Dirty -

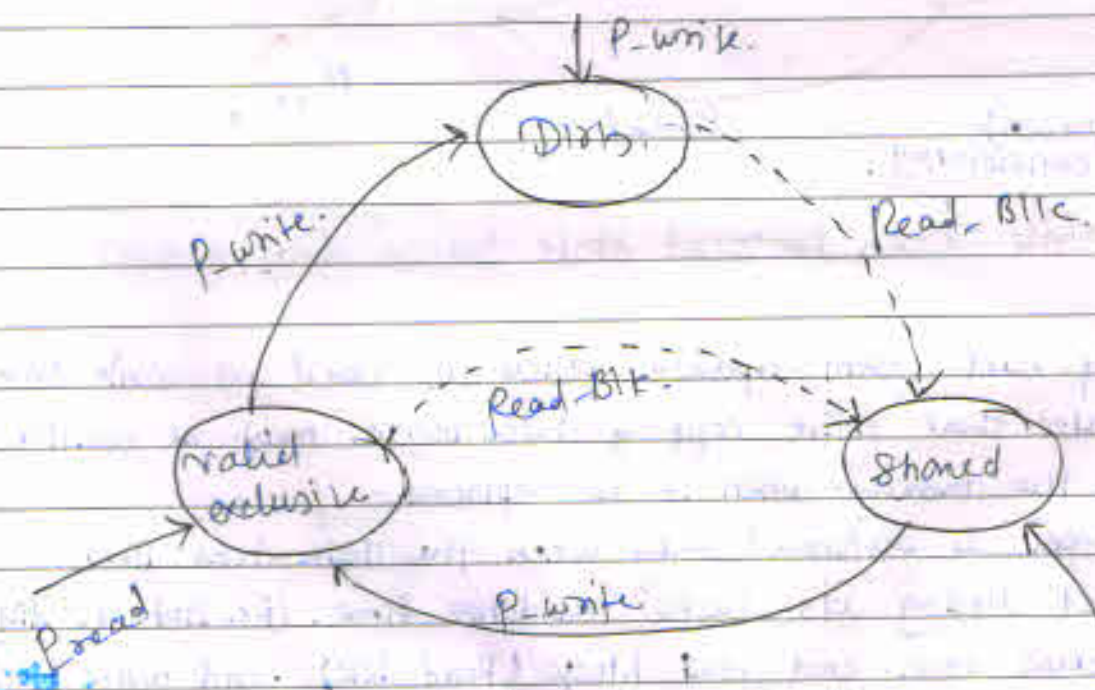
This is the only available copy of data & in this state the inconsistency of the memory copy of data is assumed.

There are 2 diff policies used by Firefly protocol for diff types of blocks.

- It has implementation of copy-back update policy for private block whereas shared blocks are supported during write-through policies option.

- write consistency update cmd is used to update all the copies for maintaining consistency in the system.

- state trans diag of this model is shown in fig



State trans diag. for Firefly protocol.



## 1) Instr<sup>n</sup> level parallelism (ILP) :

→ modern process has ability to execute several ops<sup>n</sup> of a prog simultaneously.

- ILP is the measures about how many ops<sup>n</sup> associated with a prog. can be performed simultaneously.
- This is handled by overlapping the instr<sup>n</sup> of prog during exe. and is actually termed as ILP.
- In this, multiple instr<sup>n</sup> from same instr<sup>n</sup> stream execute concurrently.
- Such instr<sup>n</sup> are generated & merged by h/w level e.g. Superscalar processor / compiler  
it generates instr<sup>n</sup> to be executed concurrently.  
as an example of VLIW.

## 2) Thread-level parallelism :

- Thread-level parallelism occurs when multiple thread or instr<sup>n</sup> sequences in a same application executed concurrently.
- In this, prog. is divided into independent parts in which each part is called thread. and these threads execute concurrently.
- The parallelism occurs here bcoz thread can have the effect of separate independent prog. running together and separate activities are being performed inside the same parallel prog.
- This type of parallelism is used in MIMD class of m/c.
- There are 2 types 1) fine grained & coarse grained multithreading.
- Thread level parallelism comes under task-level parallelism, bcoz it is based upon org<sup>n</sup> of prog. or computing sol<sup>n</sup> into set of thread for simultaneous exe. on different processor.

## \* Transaction level parallelism :

A trans<sup>n</sup> is seq. of info<sup>n</sup> exchange and related work considered as a unit in programming.

The term trans<sup>n</sup> level parallelism is used to describe concurrent exe. of multiple processors and threads from different transaction.



## \* Models :

1) SIMD 2) MIMD 3) NUMA 4) UMA.

We have already discussed this point

## \* Comparison of SIMD & MIMD.

### SIMD

- 1] SIMD archi. are simple.
- 2] low cost
- 3] Only one prog. copy is required to be stored
- 4] Control unit contains only one decoder.
- 5] provides scalability in terms of size & performance.
- 6] Sync<sup>n</sup> is required implicitly at prog. level.

### MIMD

- 1] MIMD archi. are complex.
- 2] medium cost
- 3] No. of prog. copies depends on the no. of PE. Each PE store its own copy of prog.
- 4] Control unit is available in every P.E.
- 5] It is complex in size but provides good performance.
- 6] Handling of sync<sup>n</sup> requires explicit data structure & operation.

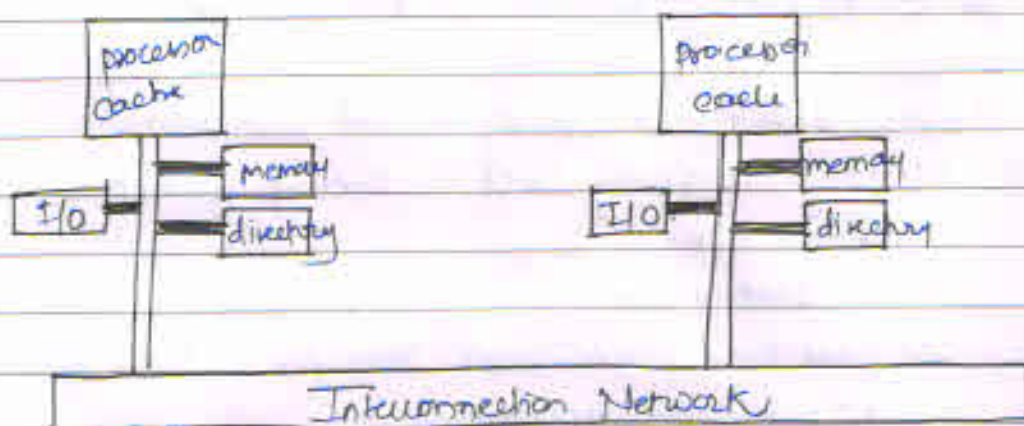
\* Dataflow model - Asyn<sup>d</sup> exe, no sequencing, Dataflow Repres<sup>n</sup>

\* Dataflow graph - operator, decoder, switch actor.

\* multithreaded archi - single, multi-queue.

### \* Directory Based Scheme:

- This is cache coherence protocol that is not based on broadcasting of data. This protocol stores the loc<sup>s</sup> of all cached copies of every block of shared data.
- The cache loc<sup>s</sup> can be kept at centralized place or at distributed manner and are called directories.
- used in scalable cache-coherent distributed memory multiprocessor system where cache directories are used to keep record on where copies of cache block reside.
- The directory based protocol is shown in fig.



### ① Directory Based Protocol Organization.

- In this two basic op<sup>s</sup> called handling read miss and handling a write to a shared cache block are implemented in this protocol.
  - The imple<sup>s</sup> of these op<sup>s</sup> are possible if the directory track the state of each cache block.
- The diff. state considered here are
- 1) Shared - The value in the mem: is up to date of the block cached by one/more processor.
  - 2) Uncached - processor do not have copy of cache block.
  - 3) Exclusive - Exactly one processor own copy of the cache block & it is designated as owner.
- In this state, the memory copy becomes out of date bcz the owner has updated the block.



## \* Comm<sup>n</sup> Cost in II<sup>e</sup>l machines:

- The II<sup>e</sup>l computing platform provides facility for execution of II<sup>e</sup>l progs.
- The various module of these II<sup>e</sup>l progs execute in diff. nodes of distributed computing based II<sup>e</sup>l system or in diff. processing elements of a shared memory II<sup>e</sup>l system.
- The module need to communicate with each other for completion of computing tasks.
- The comm<sup>n</sup> cost among these task depend on several factors like interconnect<sup>n</sup> n/w used for connecting nodes, modes used for II<sup>e</sup>l programming, protocol used in comm<sup>n</sup> etc.

## - Issues affect the Overall Comm<sup>n</sup>:

1) Message passing Cost: Theory you can write your own factors are:-

startup-time, sele<sup>n</sup> of routing algo., interface establishment per-hop time, per-word transfer, store & forward-routing and cut-through routing.

2) store & forward Routing

3) packet Routing

4) Cut-through routing.

## \* Level of II<sup>e</sup>lism:

- There are 2 types 1) h/w II<sup>e</sup>lism 2) s/w II<sup>e</sup>lism.

- In h/w, it includes activities provided at architecture level itself.

- In s/w, it is divided into II<sup>e</sup>lism based on task & data level.

- The types of II<sup>e</sup>lism in applications are -

1) ~~level~~ Instruction level

2) Thread level or Task level

3) Transaction level II<sup>e</sup>lism.

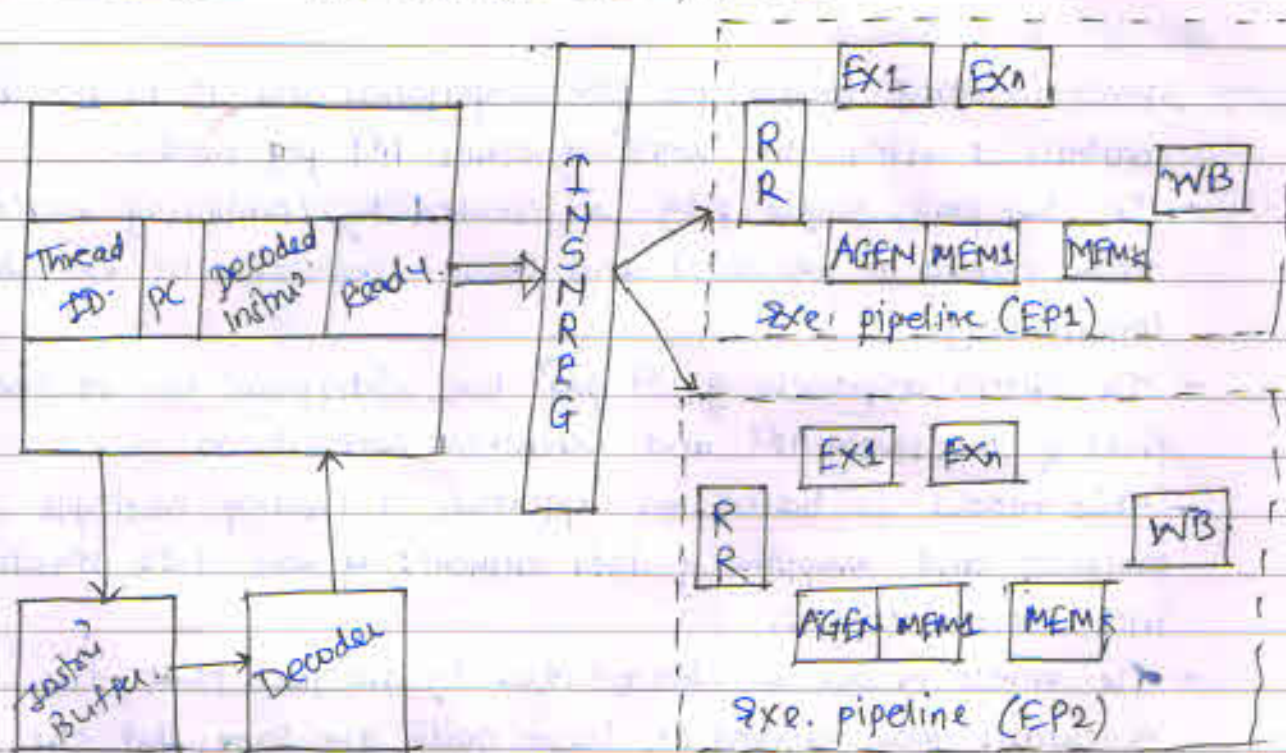


## \* SIMT :-

- The single instr<sup>2</sup> multiple thread archi. has been developed to achieve throughout Computing with high energy efficiency.
- In today's computing field, SIMT processors are used most of the times for Graphics Procence Units. (GPUs)
- The SIMT archi. is created with the combined support from h/w & s/w sides.
- The h/w archi. side in SIMT provides an approach for the conversion of scalar instr<sup>2</sup> into vector-style single-instr<sup>2</sup>-multiple data [SIMD] processing for the achievement of energy efficiency.
- On the other hand s/w side, the SIMT prog model for e.g. CUDA and OpenCL provides facility to achieve data lflism. to be implemented as task-level lflism.

## \* Micro-archi. of SIMT :-

- The micro archi. of an SIMT core is shown in fig. An instr<sup>2</sup> is selected & issues from ready warps by warp scheduler to the multiple exe. pipelines.



Micro-Archi. of an SIMT core.



- These multiple exe. pipelines are also called as streaming processors or processing elements (PEs)
- A warp is considered as a ready warp when all dependencies requires for the exe. g next instr<sup>s</sup> have been resolved.
- In the archi. each entry of warp scheduler contain the info<sup>s</sup> of warp. The instr<sup>s</sup> from instr<sup>s</sup> cache is read using Prog. Counter.

#### \* SIMT prog<sup>s</sup> model:

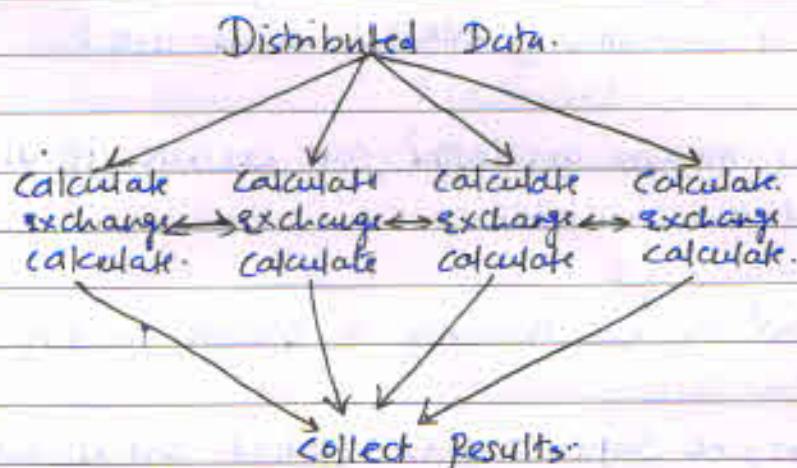
- The data parallelism is expressed as task level according to the SIMT prog<sup>s</sup> model.
- The prog<sup>s</sup> model here follows the - single prog. multiple thread paradigm this indicates that all threads shares the same prog.
- The scalar code commonly called as kernel code is written by an appli<sup>s</sup> developed to utilize the power of SIMT machines.
- All thread execute the same kernel fu<sup>s</sup> and are also referred as work items. Each thread identifies the data to be operated upon using the unique identifier.

#### \* SPMD (Single prog. Multiple Data):

- SPMD is prog<sup>s</sup> model for parallel computation and it is possible to construct it using the comb<sup>s</sup> of other parallel prog<sup>s</sup> model.
- In this prog, single prog. is created by combining multiple task (thread or processes) and these task can be executed concurrently.
- The SPMD approach of parallel prog<sup>s</sup> has widespread use in the field of massively parallel and scientific computation.
- This model is based on approach of having multiple processors and mapping of data elements of the data structure into these processors.
- The SPMD model is characterized by the fact that the executable prog running on these node are same. But the data upon which the step of these prog. operate are different.
- In this, it splits appli<sup>s</sup> data among various processors. This type of parallelism is referred as geometric parallelism, domain decomposition or data parallelism.

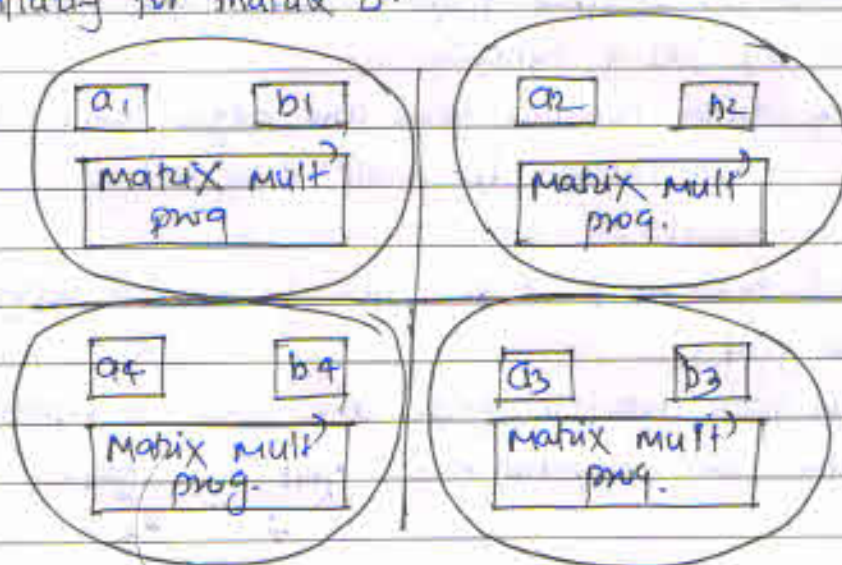


- The overall data used by appli<sup>n</sup> is divided and given to a set of processes included to solve the problem.
- One can easily understand SPMD by the use of an e.g. Consider matrix multi<sup>n</sup> prog on 4 nodes of 1K1 computing sys.



Base structure of SPMD model of computing.

- In this arrangement, 3 nodes are given exclusive responsibilities to act as slaves whereas one node act as master as well as involve in computations as a slave also.
- Also assume 2 matrixes A & B of size  $100 \times 100$  for multi<sup>n</sup>. The whole matrix can be divided into 4 chunks with 25 elements in each bcz overall 4 nodes are available for computation.
- matrix A is divided into 4 chunks as  $a_1[25][25]$ ,  $a_2[25][25]$ ,  $a_3[25][25]$ ,  $a_4[25][25]$  where  $a_1$  contains 1<sup>st</sup> 25 element,  $a_2$  next 25 & so on. Similarly for matrix B.



Matrix Multiplication as a SPMD example.



## \* Dataflow Model :-

- The dataflow model of comput<sup>n</sup> is based on graphical repres<sup>n</sup> of prog. in which oper<sup>s</sup> are represented by nodes and arcs are used to represent data dependencies.
- The ~~111~~ activities under dataflow model of computing provides various property to solve the computing prob. in efficient way.
- Some of properties of dataflow model are :-

### 1) Asynchronous Execution :-

- This means an instr<sup>n</sup> can execute if all its required operand are available.

### 2) No sequencing :-

- Instr<sup>s</sup> are not necessary to execute in any sequential fashion.

### 3) Dataflow Representation :-

- In this, no sequencing is required and it makes the possibility of dataflow repres<sup>n</sup> of a prog.
- The dataflow repres<sup>n</sup> of prog. provides the use of all forms of ~~111~~ prog. exe. w/o the assistance of any explicit tools of ~~111~~ exe.

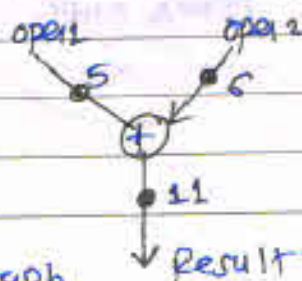
## - Dataflow Graph :-

- It is collec<sup>n</sup> of nodes & arcs. The dataflow computation works on dataflow graph. As we know, computer understand only m/c level lang. In similar, dataflow computer also works but m/c level lang repres<sup>n</sup> of dataflow computer in dataflow graph.
- Dataflow graph is directed graph which shows the data dependencies bet<sup>n</sup> nodes & functions.
- A dataflow graph contains node and edges where each node of dataflow graph contains i/p & o/p data ports.

$$\text{Result} = (\text{oper}_1) + (\text{oper}_2)$$

- The connec<sup>n</sup> bet<sup>n</sup> o/p & i/p ports are represented using the edges of dataflow graph.

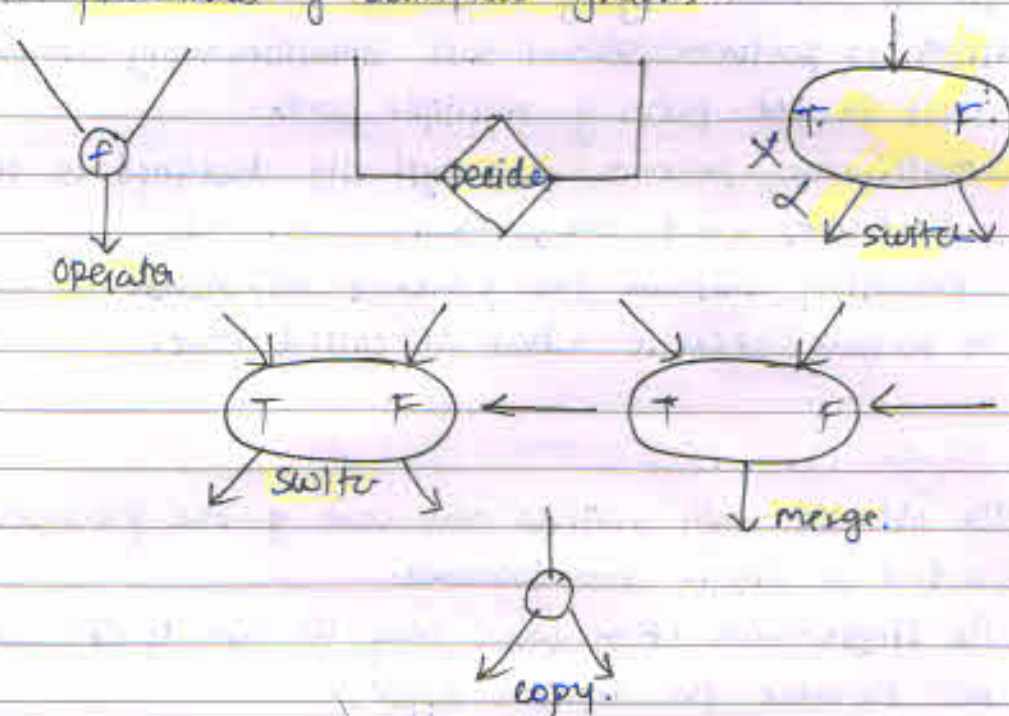
In this way dataflow graph are used to represent prog. for data flow computation.



Dataflow graph



## \* Basic primitives of dataflow graphs:



- operator: It perform some oper?
- Decider: It generates true/false value
- Switch (merge): used for directing data values. either by switch or merge
- copy: An identity operator called as copy is used to duplicate i/p token.

## \* Architectures - N-wide Superscalar architecture, multi-core and Multi-threaded:

### \* Multi-threaded Architecture:

- The thread in association with the h/w context is diff. than the j/w threads in multithreaded o/s. The h/w supported thread depends on specific form of multithreaded processor.
- The h/w based multi-threaded tech enables TLP at processor level. This is made up of duplication of architectural state on each processor while sharing one set of processor exe. states.
- The diff. archi. set are considered as separate logical processor by the OS during exe. of multiple threads.
- The OS schedules the diff. threads on unit of exe. on logical processor.
- The logical processor are created thro. h/w for sharing the processor's functional unit among diff. threads.
- The sharing of resources is crucial in multithreaded archi. bcoz max no. of sharing resources per logical processor provides more efficient oper. at multithreading level.



## \* Multi-core Architecture :-

- In the modern day of computing processor are constructed using to perform several task simultaneously at processor level itself in the form of multiple cores.

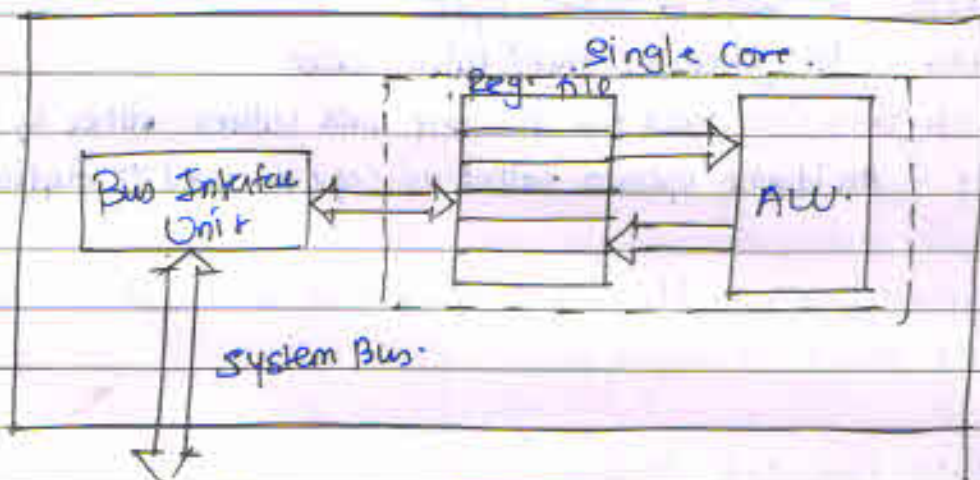
Multi-core processor are typically designed in the form of dual-core, quad-core, many-cores etc.

- Processing unit in the processor responsible to read instr<sup>n</sup> to perform specific action is called core.

### 1) Single-Core CPU :-

- The old CPU only utilizes one core of the processor is called as Single Core processor.

The single-core CPUs were very 1st type of CPU and not suitable for modern appl<sup>s</sup>.



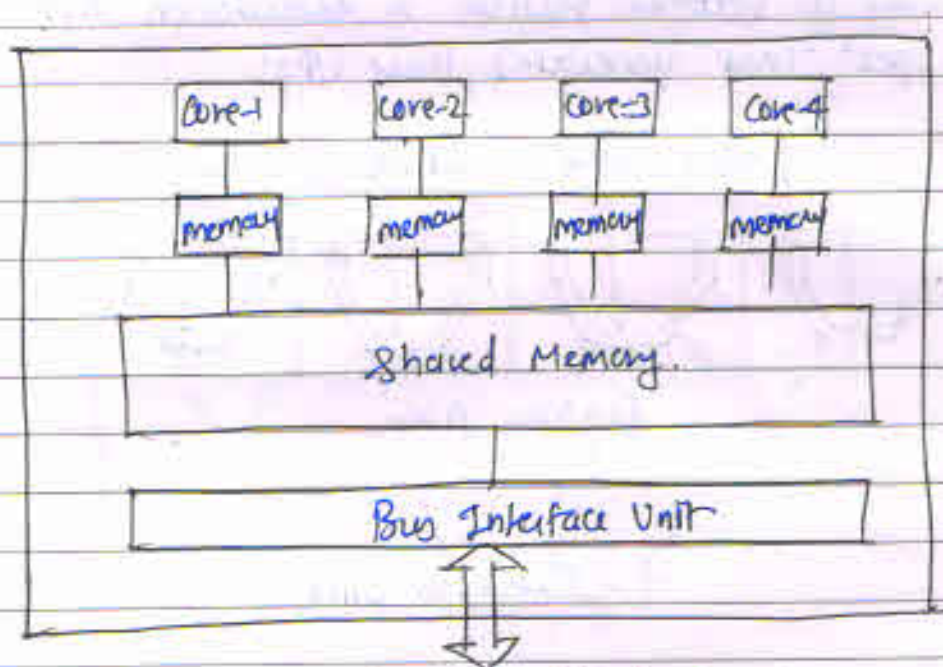
Single-Core chip.

### 2) Multi-core CPU :-

- The multi-core processor contain more than one distinct core in a chip. In this way, if a chip has 2-cores it is named as dual processor. Similarly quad core processor it contains 4 processors in it (core).
- The multi-core CPU design implements separate exe. pipelines for each of every cores.
- This means each core can run thread of exe. w/o requiring any resources from other cores.
- Multi-core processors are considered as MIMD category of p/c bcz multiple instr<sup>n</sup> are executed by diff cores.



- It also supports shared memory multiprocessor because all cores can share the same memory.



Other off-chip components.  
Multi-core processor.

- Intel xeon is multi-core processor packaged with the logic & circuitry for 2 or more Intel xeon processor.
- In fact, multi-core processor is based on packaging of more such processor in single physical processor.
- The advantage of multi-core processor is the improvement in performance by using multiple cores to run multiple tasks simultaneously.

#### \* N-wide Superscalar Architecture:

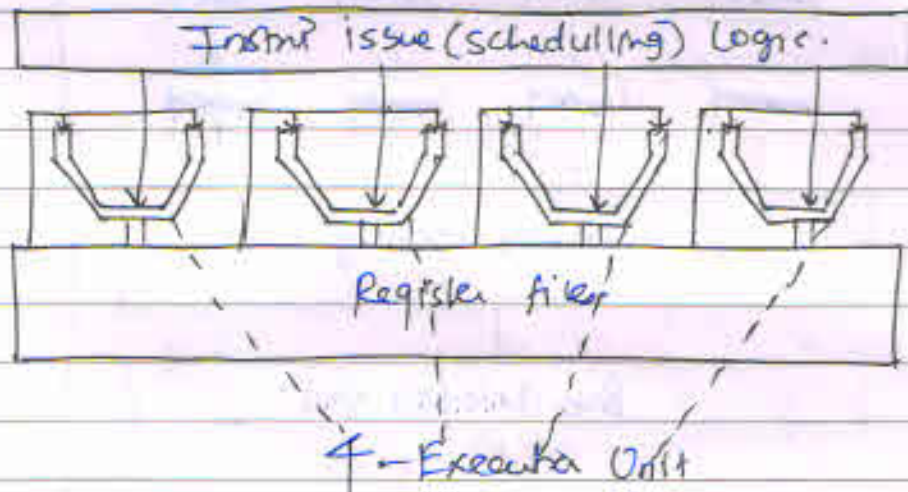
- The superscalar archit. implemented to support issue of one, two etc instr<sup>s</sup> and accordingly 4-wide, 6-wide. Superscalar architecture are designed by various vendors.

#### \* Superscalar archi. processor:

- The superscalar archi. contains multiple exe. unit for instr<sup>s</sup> exe.
- In this, single centralized reg. file is used to read operands from it and write results into it by each exe. unit.



- The result written back to the register file by the exe. unit becomes visible to all of the exe. unit on the next cycle. This way it becomes possible to execute on diff. units from opers that generates their i/p's.



### Instr. Issue logic & 4 Exe. Unit in a Superscalar.

#### \* ILP in superscalar processor:

- The superscalar arch. uses complex bypassing h/w to reduce delay bet dependent Instr.
- The delay is reduce bcz bypassing h/w forwards result in instr. exe. to all exe. units.
- The superscalar processor contains the Instr. Issue logic which is used to store Instr. of program.
- This instr. issue logic provides Instr. to the unit in parallel.

#### \* Instr. Issue logic in superscalar processor:

- The changes in control flow in prog. exe. due to branches occur simultaneously across all of the units.
- The prog. development under superscalar processor becomes much easier bcz the occurrences of simultaneous activities across all units.

• H/w used in superscalar processor

→ The ILP is extracted from a sequential prog from the h/w unit used in superscalar processor.

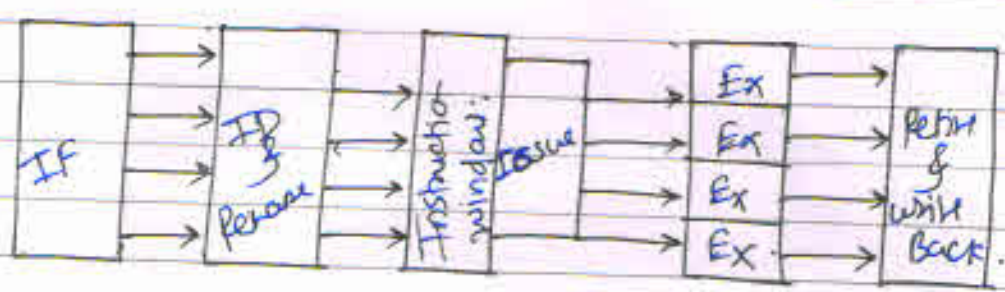
- The instr<sup>y</sup> issue logic of superscalar processor examines the instr<sup>y</sup> in the seq. prog. to determine which instr<sup>y</sup> may be issued on that cycle.

• ILP strength & weakness:-

- The superscalar processor executes instr<sup>y</sup> in parallel and in this way it is possible to achieve significant speedup while executing diff. types of progs.

- The max<sup>m</sup> improvement in the performance depend upon the limitation imposed by instr<sup>y</sup> dependancies.

• pipeline in superscalar:-



Ex - execute/addition calculation.

IF - Instr<sup>y</sup> fetch.

ID - Instr<sup>y</sup> decode/register fetch.

Superscalar Pipelines.

- The branch prediction logic used in superscalar pipelines provides control independancies for instr<sup>y</sup> in the instr<sup>y</sup> windows.

- The reg. renaming is used here to achieve name independancies in the pipeline. This way only true dependancies like data dependancies and structural conflict remain in the architecture to be handled.



\* Section of superscalar pipeline  
it is partitioned in distinct section based on the  
ability to issue & execute instr.

In-order section:

with the instr's fetch, decode, rename stages -  
the issue is also part of the in-order section.  
in case of in-order issue.

Out of order section:

Instr's can be executed in an order different from  
that specified in the prog.