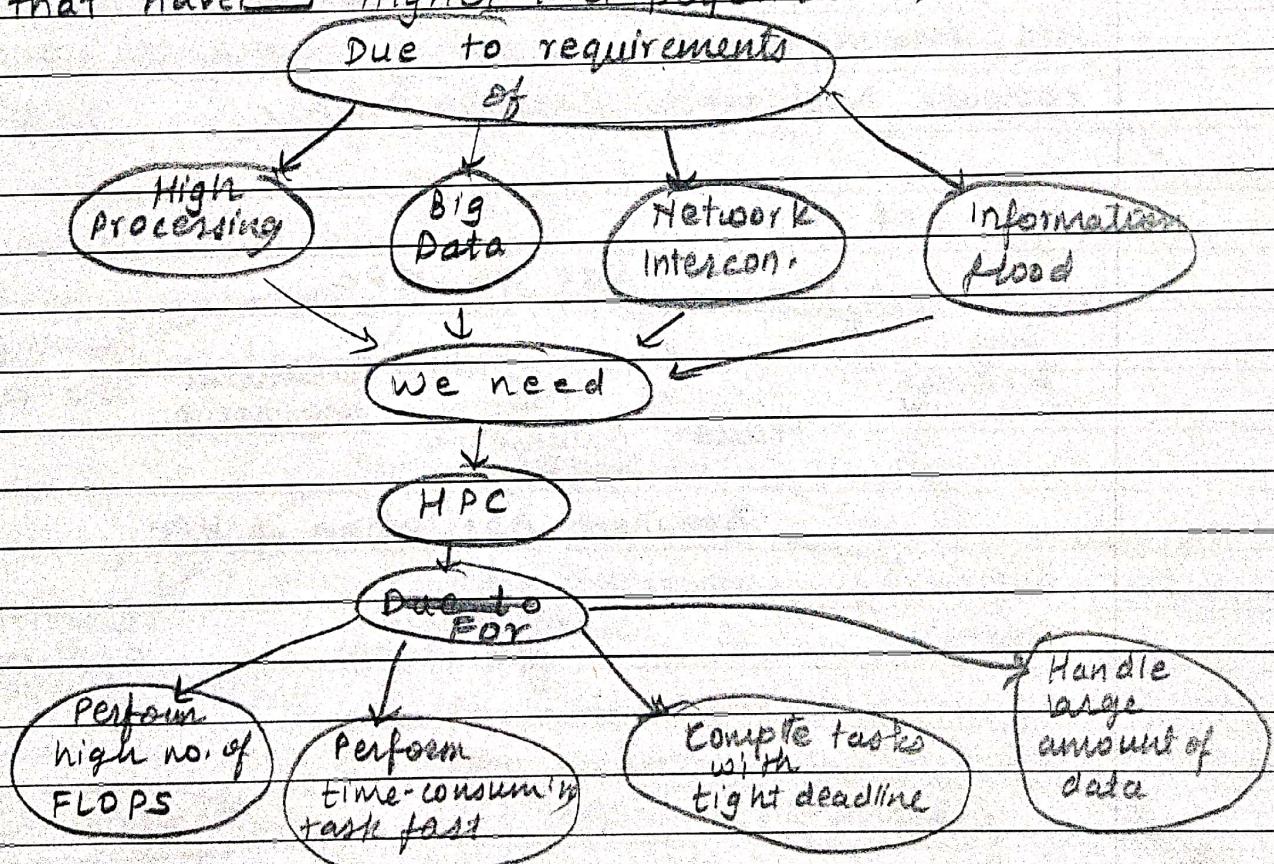


HIGH PERFORMANCE COMPUTING

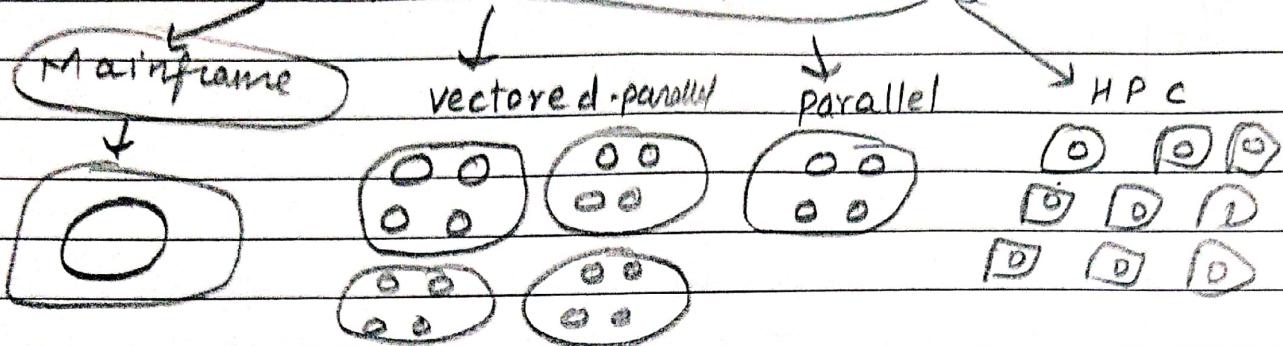
UNIT 1: INTRODUCTION TO PARALLEL COMPUTING

Q) What is HPC:-

- Refers to performing computation operations collaboratively on multiple computers that have higher level performance in terms of throughput.
- Keywords:- Computation collaboratively on multiple comps. that have [] higher level performance.

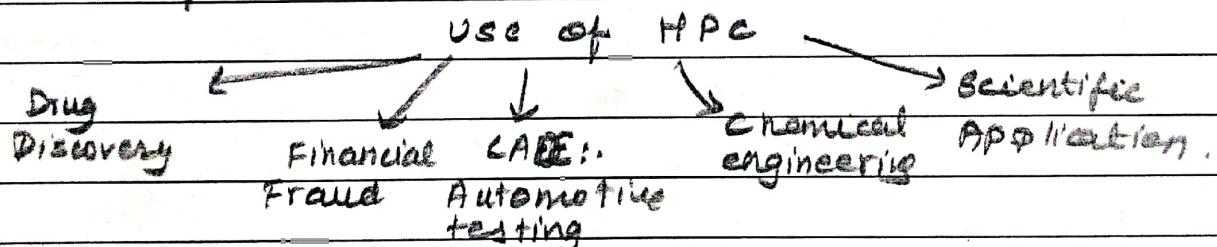
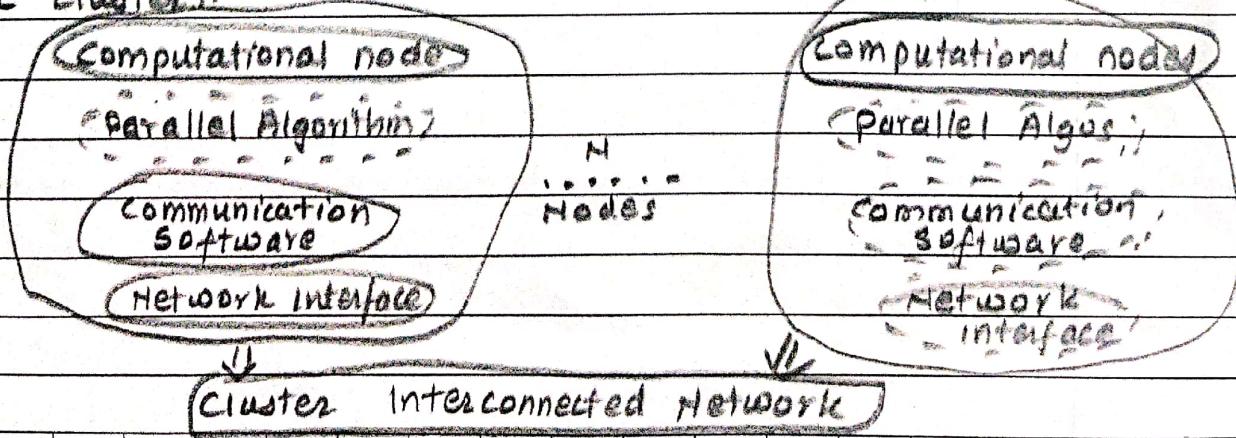


- HPC is mechanism of fast computation in parallel on computing nodes on a very fast computing network.

Different types of System

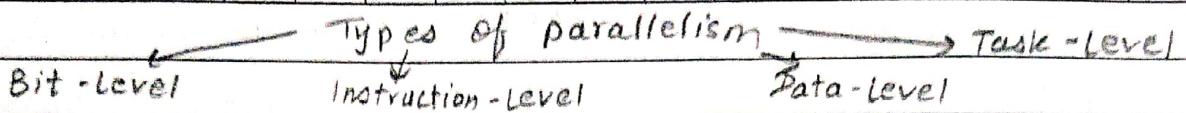
O : - is a computation Node.

- HPC uses large no. of high capacity efficiency computational node, whereas rest use less nodes but high end.
- HPC uses aggregate computing power for handling compute and data intensive tasks.
-
-
-

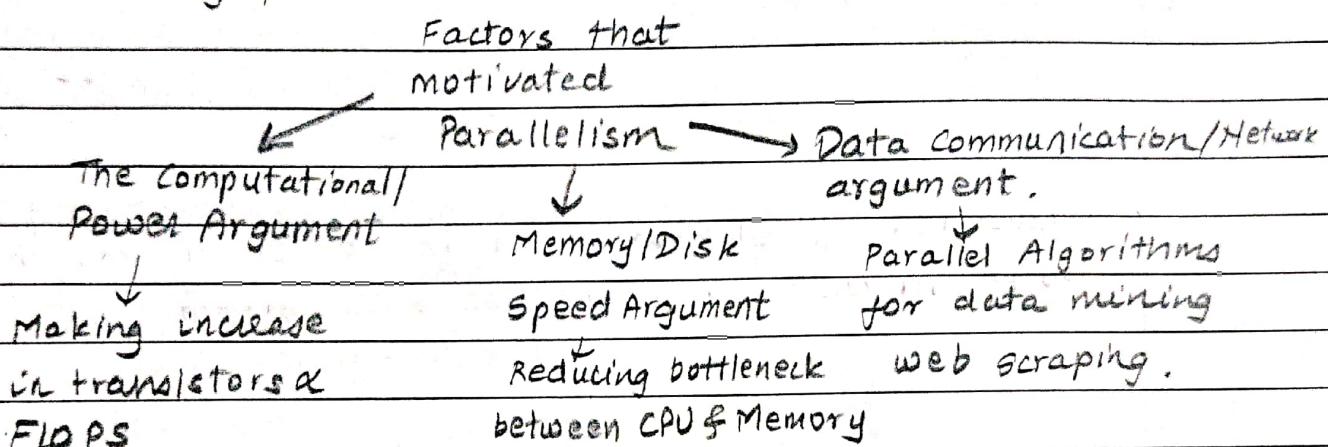
Uses of HPC:Prominent Application of HPCEngineering DesignCommercial ApplicationsScientific applicationHPC cluster:

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

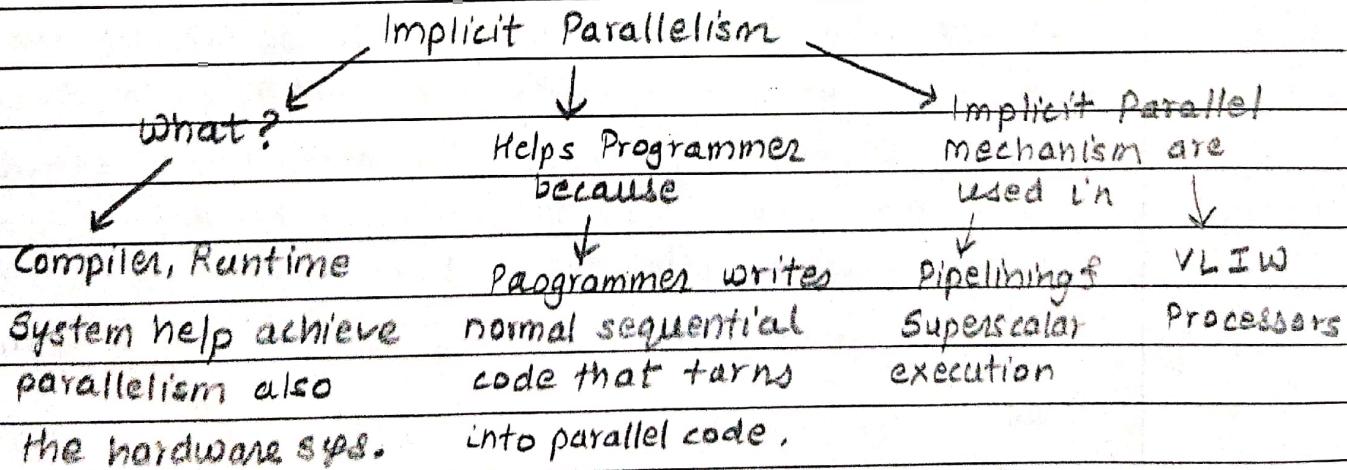


Q Motivating parallelism:-



Q Implicit Parallelism :-

- Helps programmer write programs without any concern for exploitation of parallelism.



Q Pipelining & Superscalar execution:-

Enhance Processor Performance by

Implementing
instruction-level parallelism

Subdivide pipeline
into multiple stages.

- In superscalar processor multiple instructions are issued per cycle and multiple results are generated by multiple pipelines per cycle.

- Superscalars designed to exploit instruction-level parallelism.

- Superscaling execution of code:

1) load R1, @1000

2) Load R2, @1008

3) add R1, @1004

4) add R2, @100C

5) add R1, R2

6) store R1, @2000

	IF	ID	OF	load R1, @1000
2)				load R2, @1008
3)				add R1, @1004
4)				add R2, @100C
5)				add R1, R2
6)				store R1, @2000

(i) Code

(ii) Execution schedule

IF: Instruction Fetch. ID: Instruction Decode. OF: Operand Fetch

E: Execution. WB: Write Back. NA: No Action

- [Describe Execution schedule :- (i) Instructions issued at same clock cycle ($t = L$) are mutually exclusive. (ii) Instructions issued after a cycle are not mutually exclusive.]

Superscalar execution depends on 3 things

True data dependency

Execution of one instr.
depends on result of prev.

Resource Dependency

Two instr. compete
for single resource

Branch /
Procedural
dependency

In exec. of
conditional statement we
don't know the flow of execn.

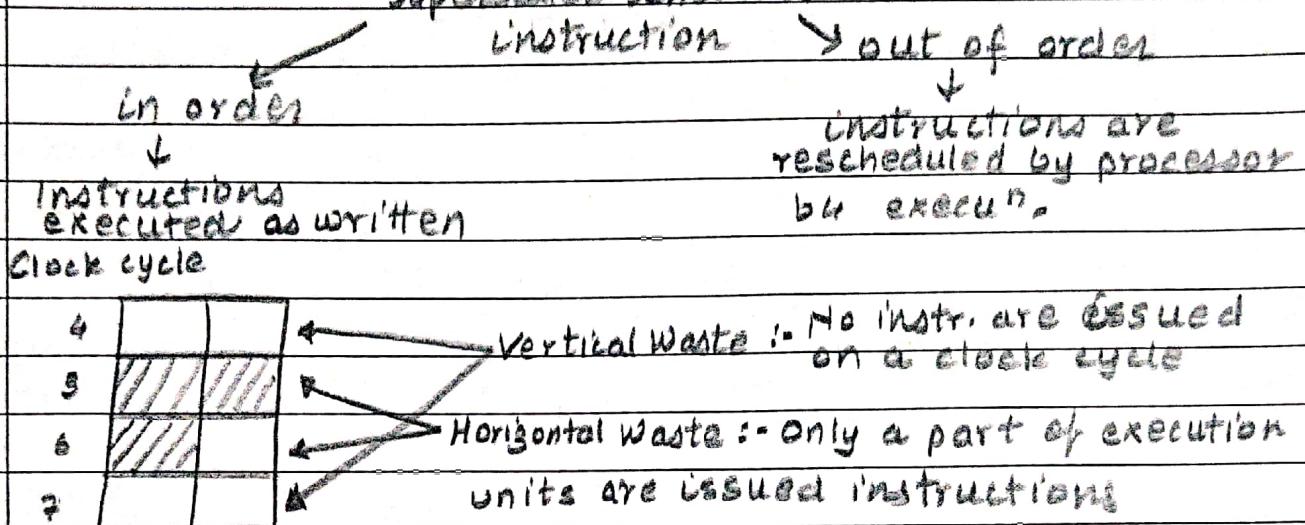
As a result scheduling instr. a priori
may lead to errors.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

- To tackle branch dependency speculative scheduling is done.
- Control independent instructions are moved up.
- Main superscaling job is to detect & schedule concurrent instructions

Superscalar schedules



Very Long Instruction Word Processors

- Helps exploit parallelism on instruction-level like superscalers.
- Compiler is the heart of ~~VLIW~~ VLWI processors.
- Resource availability and data dependency is checked during compile time.
- Techniques used :- BP SD LV [Acronym: ~~BP+LUSD~~]
 - Branch Predication
 - Speculative decomposition
 - Loop Unrolling
- Instructions that can run concurrently are combined into group or bundles and parceled off to processor as a single long instruction word.

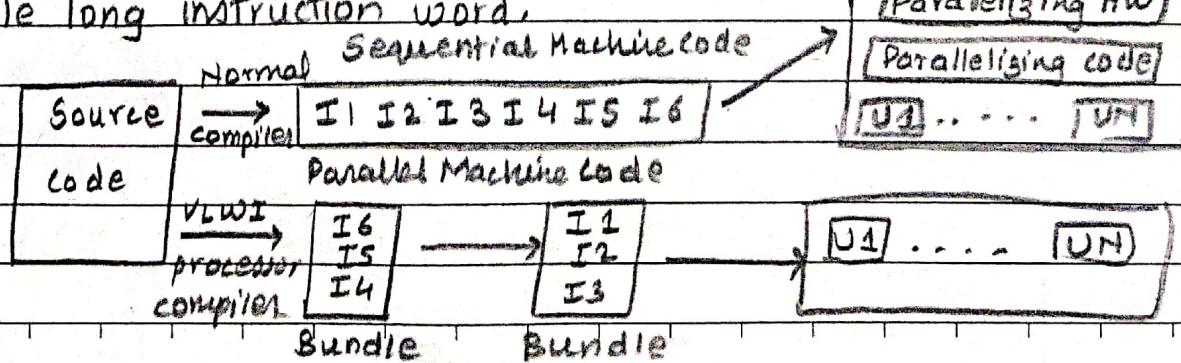
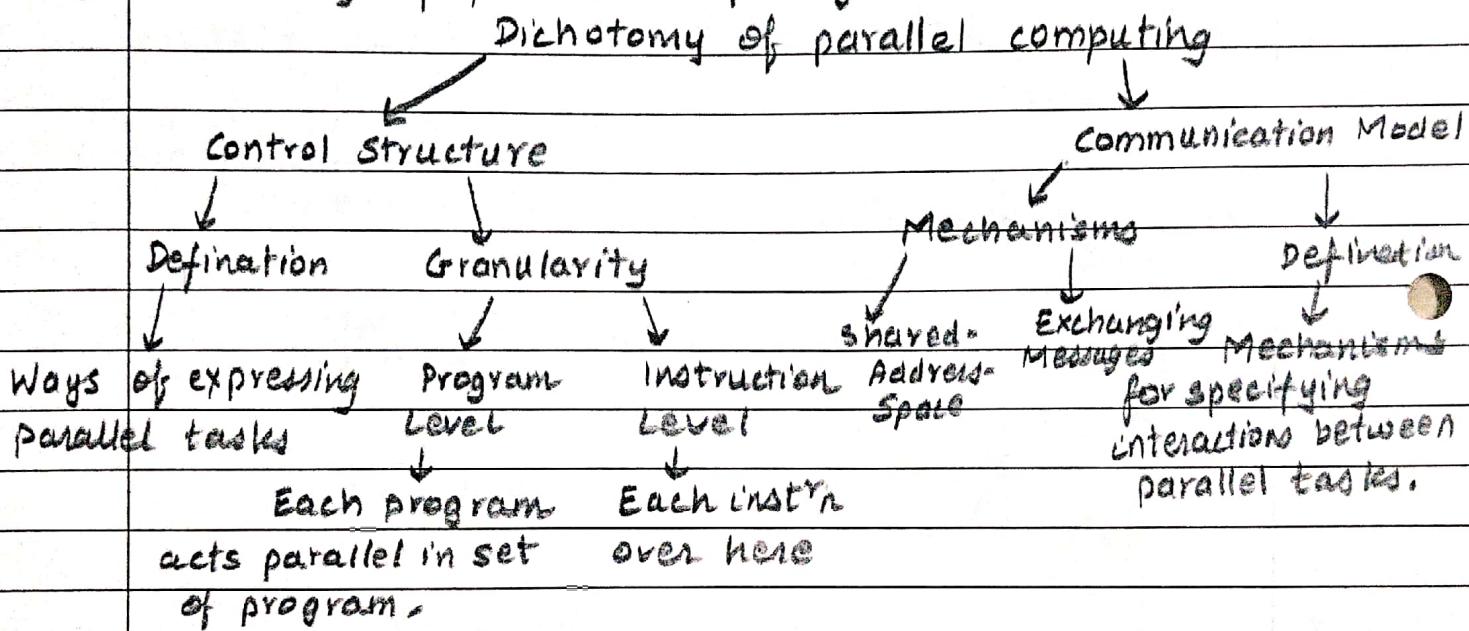


Fig: Working of VLWI processors

- Drawbacks :- (1) Dynamic's program unavailable to make schedules.
(2) Stalls on fetching data due to cache miss is unpredictable.

Q Dichotomy of parallel computing:-

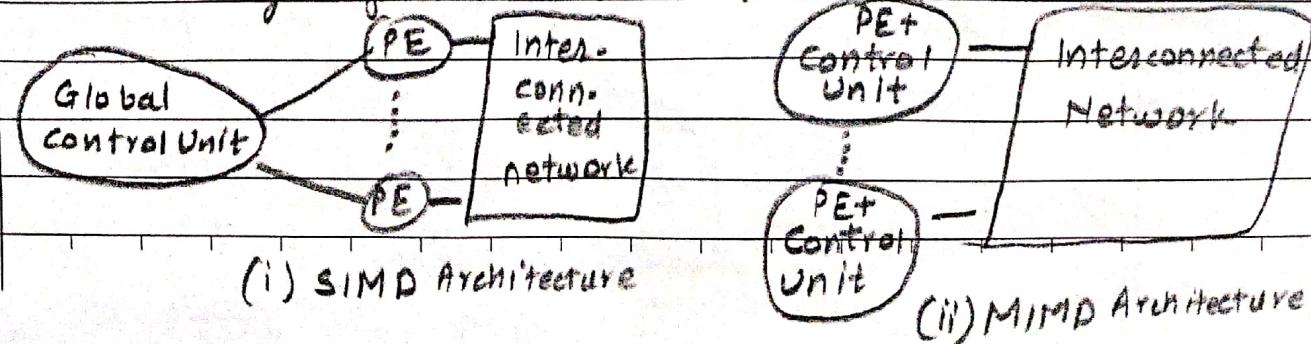


1 Control Structure:-

- Based on granularity various architectures are proposed
- | Sequential | Parallel |
|---|--|
| SISD | SIMD |
| Single Instruction,
Single Data | Single Instruction,
Multiple Data |
| MISD | MIMD |
| Multiple Instru ⁿ ,
Single Data | Multiple Instruction,
Multiple Data |

This is called Flynn's taxonomy.

Parallel sys. follow SIMD & MIMD.



- Communication Models:-

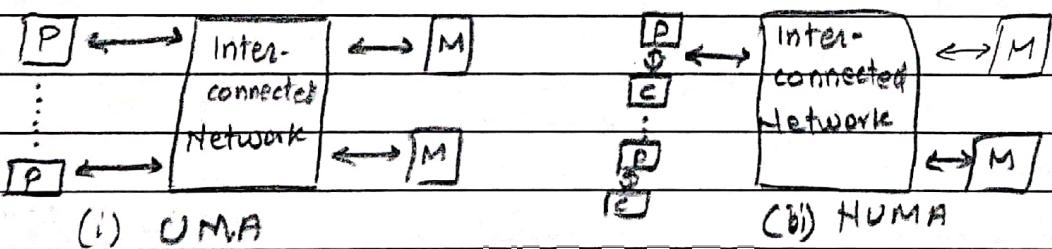
- Shared-address-space:-

Classified as

Uniform
Memory Access [UMA]
↓
Time taken to access
word is identical

Non-Uniform Memory
Access [NUMA]
↓
Some words require
different time to access.

- Set of all possible physical address is called address space of that processor.
- Shared-address-Space is a platform where all processors access common data space.



- 2) Messaging Platforms: Passing Platforms:-

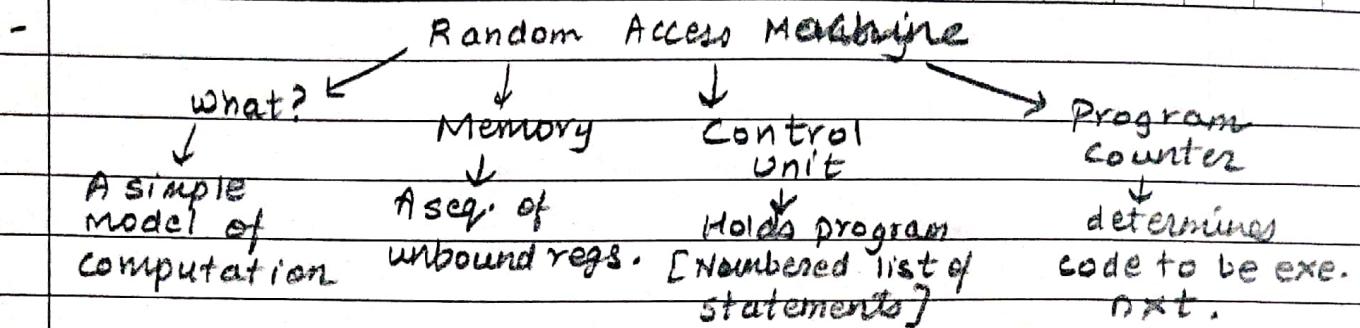
- P processors with separate address space communicate.
- Processors interact with each other through messages.
- Single node is complete in a sense. It is a single processor or shared-address-space multiprocessor.
- Operations consists of
 - (i) send
 - (ii) Receive
 - (iii) ID
 - (iv) Number of processes

- Physical organization of parallel platforms:-

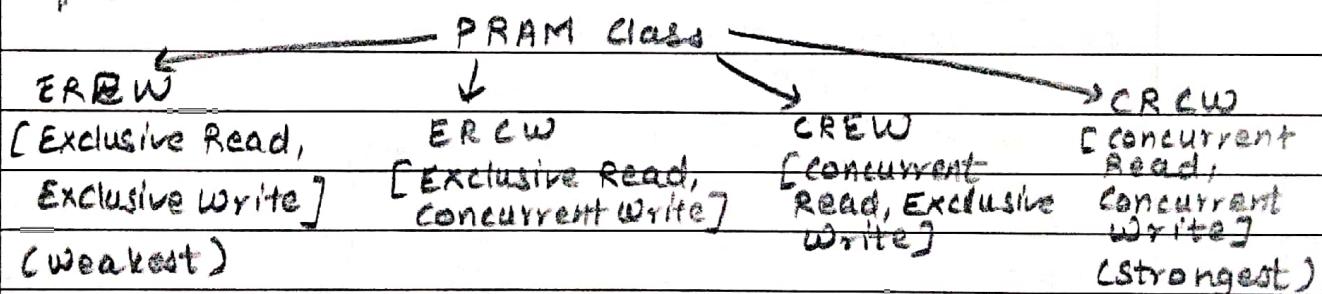
- Architecture of an ideal parallel computer:-

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI.

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.



- RAM can be made parallel by adding p processors and global memory of unbound size which is PRAM type.
- This is Parallel Random Access Machine [PRAM]
- Processors :- (i) work in same address space, (ii) Have different cycle, (iii) same clock.

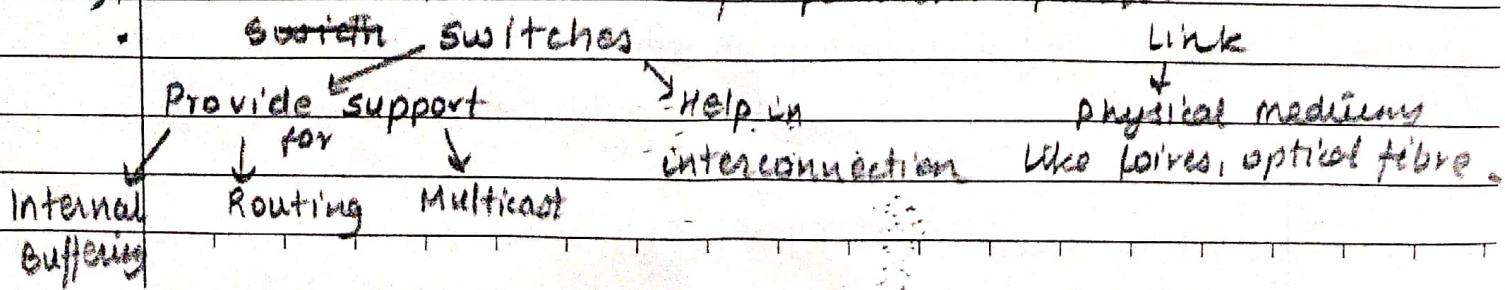


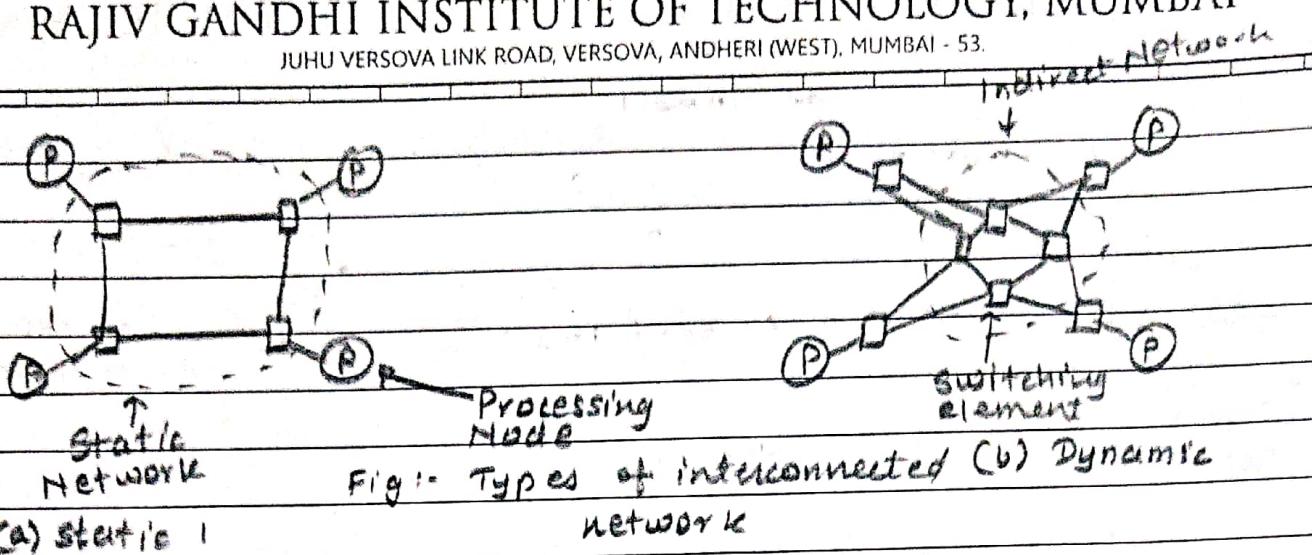
Exclusive :- Access is exclusive only to one resource

Concurrent :- Multiple Access.

- Consider a PRAM with p processors & m global memory.
- Set of switches connect processors to memory.
- A word cannot be accessed by more than 1 processor hence the total no. of switches required = $\Theta(mp)$
- This is not possible, hence PRAM cannot be implemented practically.

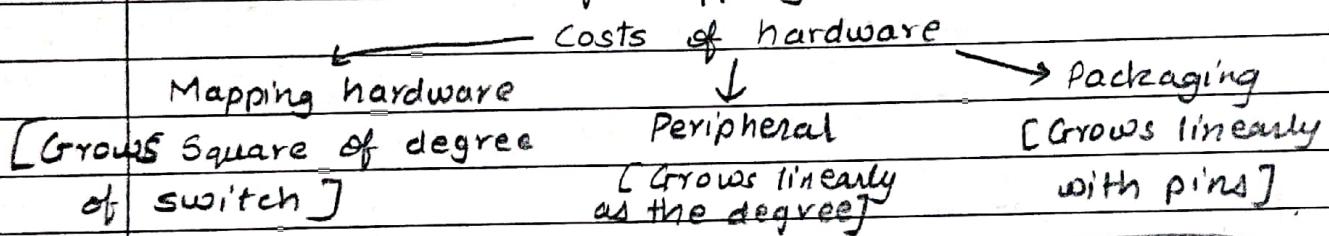
g) Interconnection network for parallel computers:-



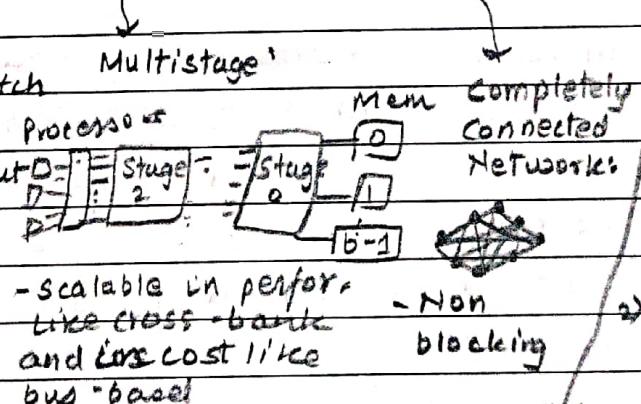
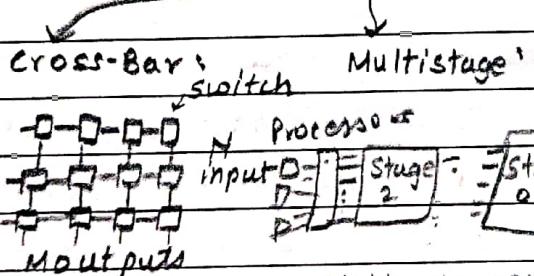
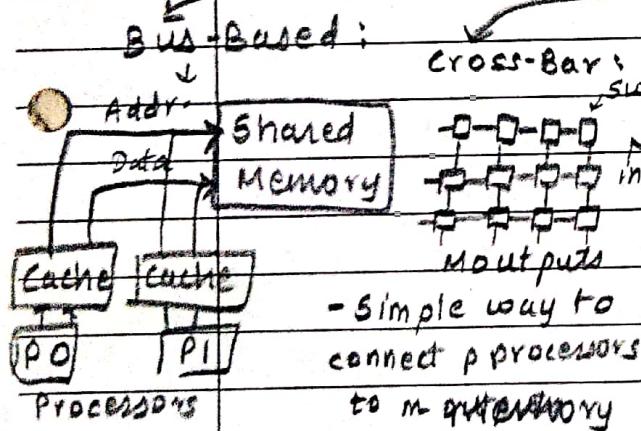


- Mapping of input to output is the basic functionality provided by switch.

Cost of mapping hardware



Network Topologies



Linear Array
Meshes,
k-d meshes

Linear array
with wraparound
Line,

2-D Meshes

Star Network 3) In
k-d
meshes
d-dimension



- Cost of Nodes, wires.

- Instance is J2 network

Bus interface - Non-blocking Let $p = \text{no. of I/Os}$
network then

[No one block
each other]

- Scalable in performance

but not in cost

$$j = f(2i) \quad 0 \leq i \leq p_1 - 1$$

$$2i + 1 = p_1 \quad p_1 \leq i \leq p - 1$$

If true then link betw

i & j exists,

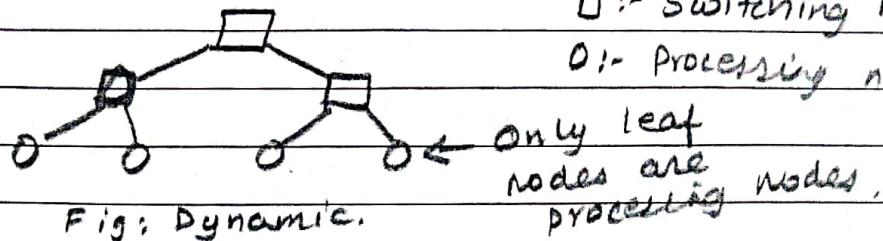
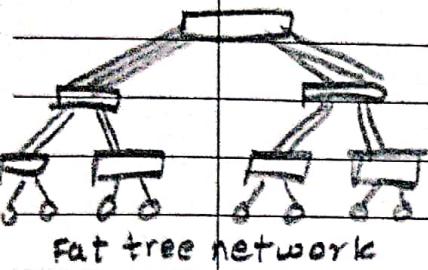
- central
x - no. of
processor nodes also
may
bottleneck
of connection

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI.

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Another one tried hard to fit in [last I tried so hard....]

7) Tree-based Networks:-



□ :- Switching nodes

○ :- Processing nodes,

Only leaf
nodes are
processing nodes.

- Static tree based shall have no switches.

② Static interconnected Network costs :-

Costs Associated with			Costs :-
Diameter			No. of wires required
- Max dist bet'n two processors.	Connectivity		Bisection - width and Bisection - bandwidth
	- Measure of		- Min. no. of links to remove
- Different Networks have different diameter.	multiplicity of paths between any two processors.	to divide network in equal parts is called bisection width of a network.	
- Some network diameter:-	- Reliability & connectivity	- Vol. of communication allowed between any 2 halves of network is bisection bandwidth.	
1) Fully-connected :- 1	- If network breaks in 2 disconnected network by removing min. no. of arcs		
2) Star-connected :- 2			
3) Ring-network:- $\lfloor P/2 \rfloor$	It is known as arc connectivity.		

③ Dynamic Interconnected networks:-

- Cost = Link cost + switch cost.
- Bisection determined by drawing various equal-partition and selecting min. no. of edges crossing the partition.

④ Cache coherence:-

- In case of shared address space processors additional hardware is required to keep multiple copies of data consistent with each other.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

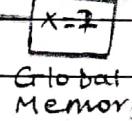
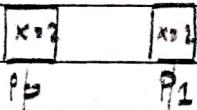
- This problem is known as cache coherence.
- To solve this additional protocols & hardware required.
- In Multiprocessor sys. multiple processors modify values making cache coherence more complex.

PROTOCOLS USED FOR COHERENCE

Invalidate [Invalidates all other cache]

loadX loadX copies

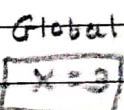
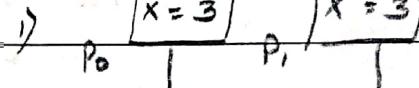
1>



[Updates all other cache copies]

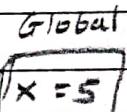
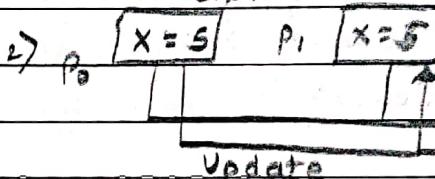
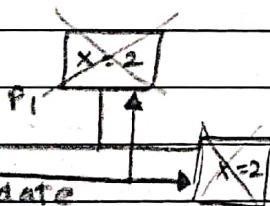
Update protocol with the same address line

loadX loadX



2>

write5,X



Update

- Whenever processors tries to modify var either invalidate or update proto. shall be fired to maintain cache coherence.

Q. Snoopy Cache System :-

- Coherence protocols are implemented using snoopy cache system.
- Coherency controller monitor/snoops bus transactions to maintain cache coherency protocols.
- Bus designed to constantly monitor caching events bet'n processor and memory m, they are called as snoopy coherence protocols.

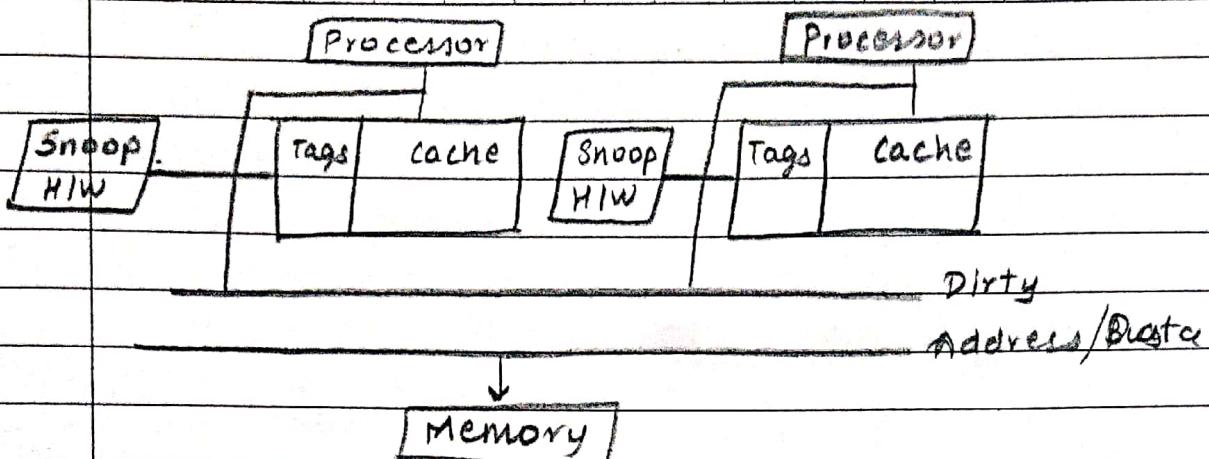
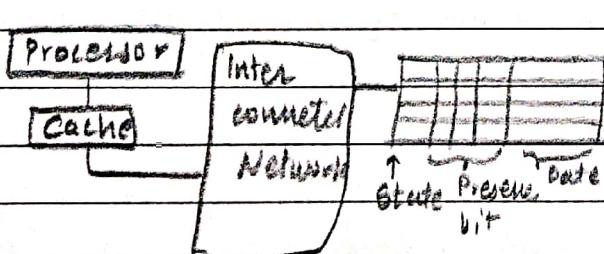


Fig: Snoopy Bus based coherence system.

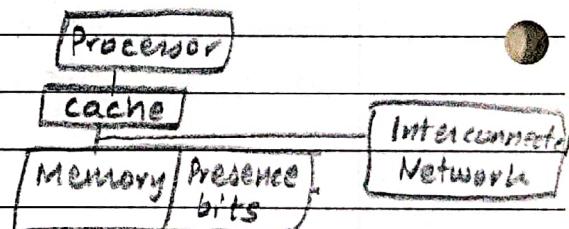
- Tags :- Determine the state of cache as per coherence protocol state diagram.
- Two cases :- (i) Snoop HW detects read ~~req~~^{req} to dirty data it then puts data outta bus.
(ii) If write req. on cache block then it invalidates it.

Q. Directory Based System :-

- Global memory connected to a directory that has bit-map representation of cache blocks and corresponding processor.

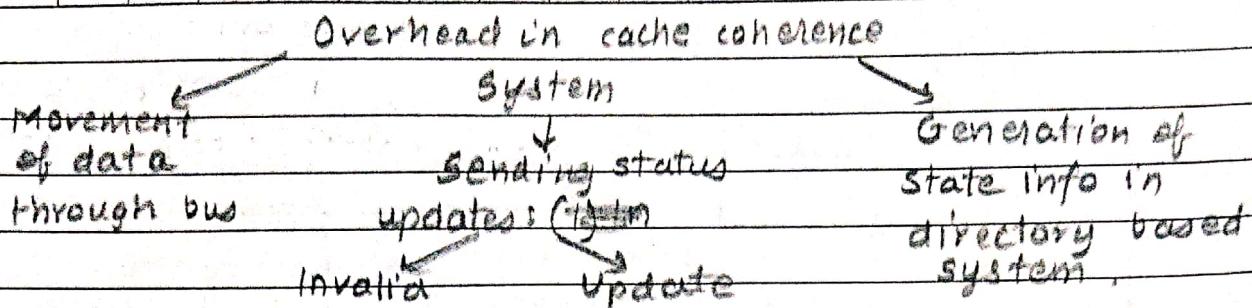


(i) centralized directory

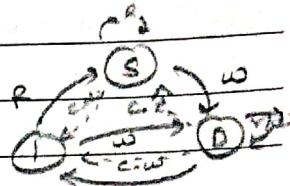
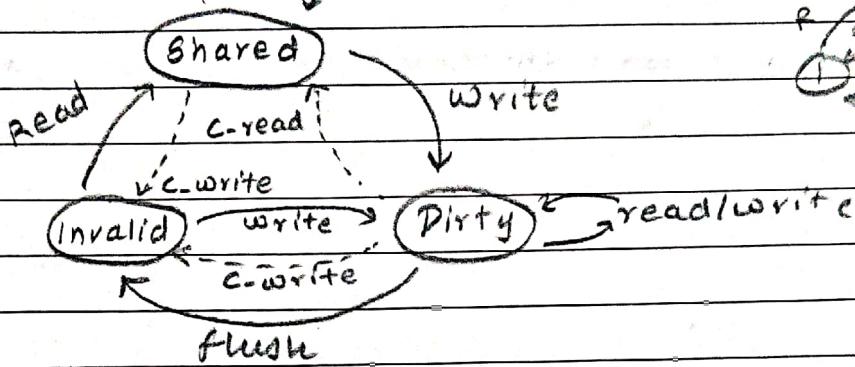


(ii) Distributed directory

- If any processor attempts to read a dirty value directory notices the dirty tag.
- Using the presence bit the request for the value is sent to the processor with the dirty value.



Q State-diagram of 3 state coherence protocol.



Q Communication cost in parallel machine:-

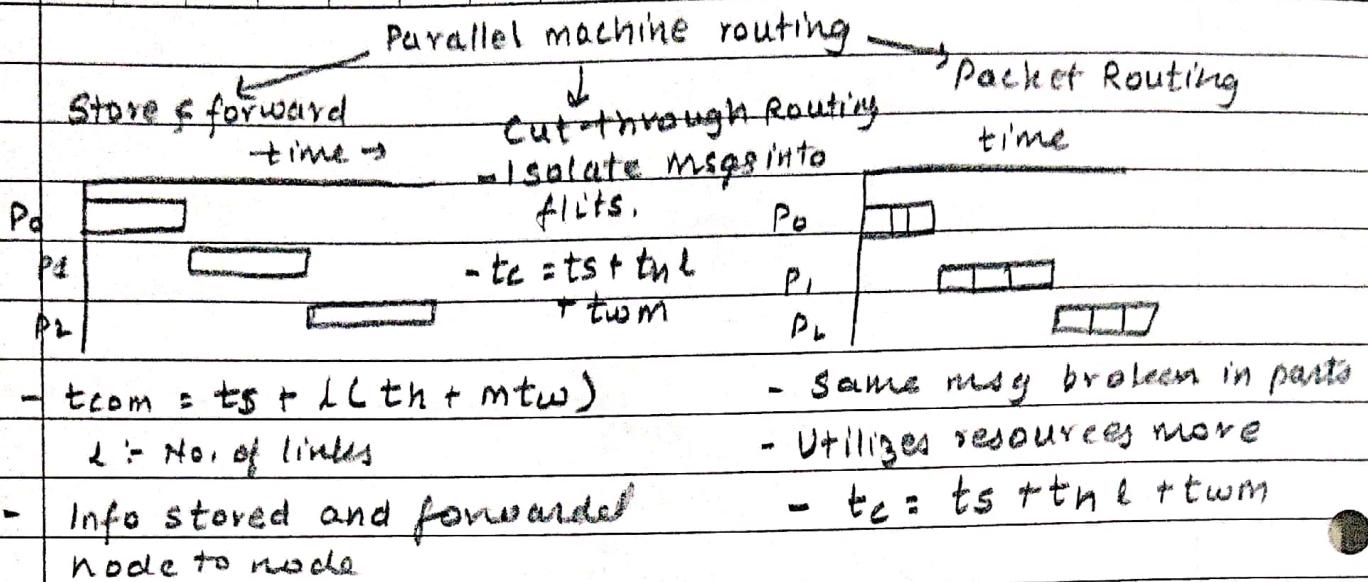
- cost of communication depends on :- (i) Programming modes semantic, (ii) Network topology (iii) Data handling and routing, (iv) Software protocols associated with a program.

Q Message passing costs:-

- Following params decide delay/latency in communication -
- i) Startup time(t_s) :- Time taken to handle message.
- ii) Per-hop time (t_h) :- Time taken by header to travel betⁿ nodes.
- iii) Per-word transfer time (t_w) :- Time taken by each word to traverse the link.

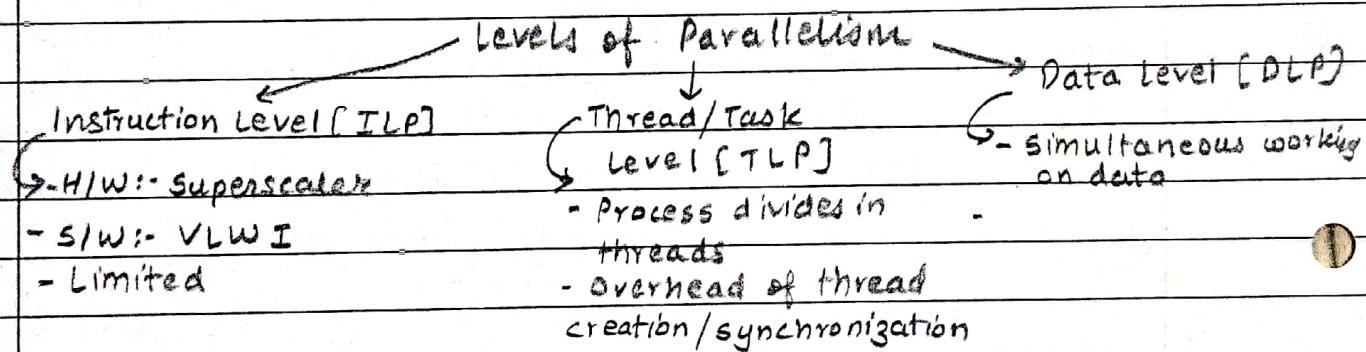
RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI.

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

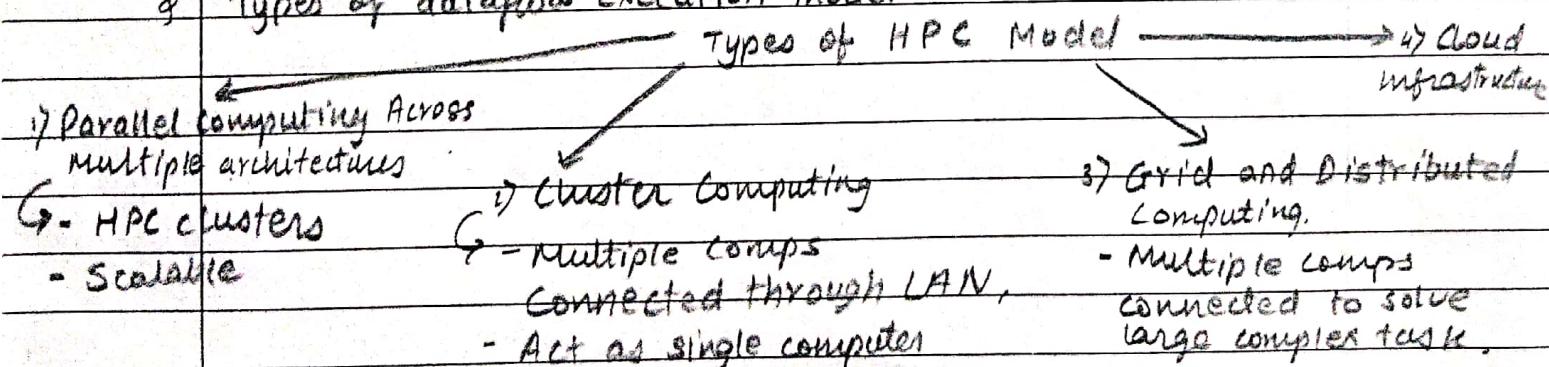


g) Communication costs in shared - Address space machine

g) Levels of Parallelism:-



g) Types of dataflow execution model:-



RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Point	SIMD	MIMD
Full Form	Single Instr. Multiple data	Multiple Instr. Multiple Data
Memory Req.	Small / Less	Large / More
Cost	Less	More
No. of Decoders	Single	Multiple
Synchronization	Facilit	Accurate
Type of programming	Synchronous	Asynchr.
Complexity	Low	High
Efficiency	Less	More

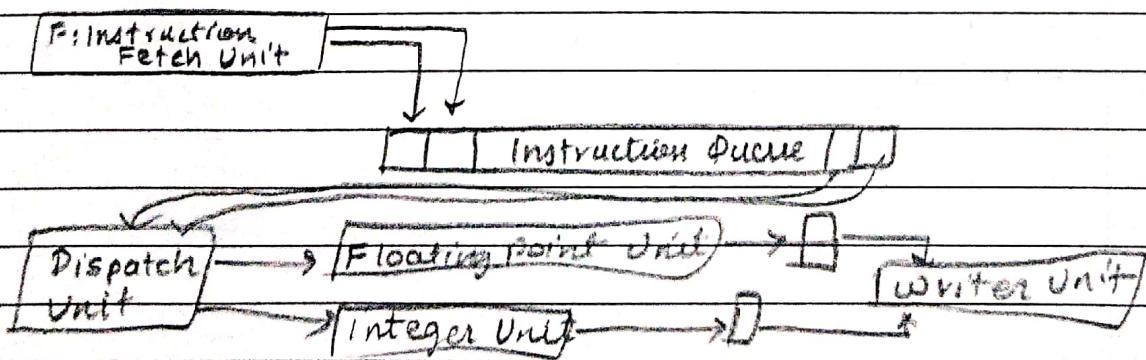
Q SIMT : [Single Instruction, Multiple Threads]

- SIMT = SIMD + Multi-threading
- Mostly implemented on GPUs.
- All threads are executed in lock-step.
- . SIMT is like a multi-core system where instructions are synchronized among the cores and not ran independently.

Q SPMD :: [Single Program, Multiple data]

- Single Technique to achieve parallelism.
- Sub-category of MIMD.
- Tasks split and run on multiple processors.

Q N-Wide Superscalar Architecture :-

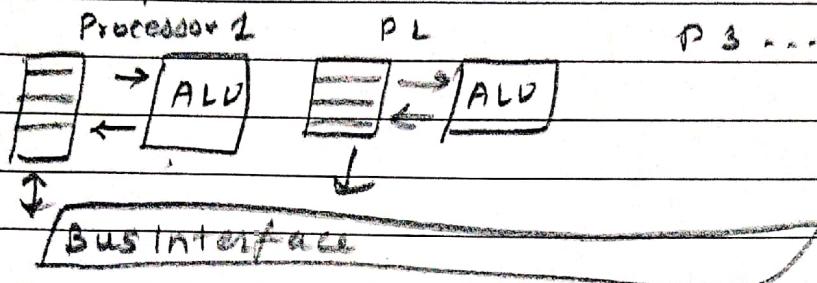


RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

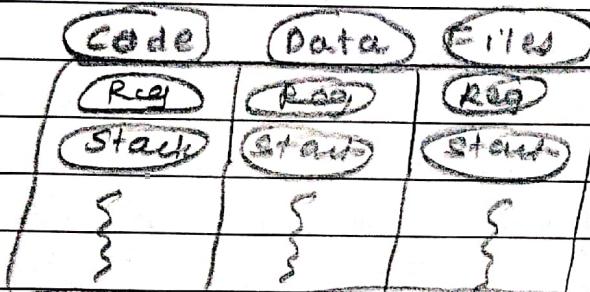
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

- Superscalar Archi. is N -wide if it supports fetch and dispatch of N instructions in every cycle.
- [Describe the diagram.]
- [Diagram for superscalar execution can be drawn]

Q Multi-core processors :-



Q Multi-thread :-



code, data & files are shared among threads.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

HPC UNIT 2 :: Parallel Algorithm Design

q. Decomposition, Tasks and Dependency Graph:-

- **Decomposition** :- Divide computations into sub-computation to execute them parallelly.
- **Tasks** :- Task is a programmer defined unit of computation.
- Tasks are generated by subdividing computation by decomposition.
- Tasks are indivisible parts of computation.
- Tasks don't need to be of same size they can be arbitrary size.

• Understanding tasks :-

$$\begin{array}{ccccccc}
 & 0 & 1 & \dots & A & \dots & n \\
 \text{Task } & \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{} \\
 : & \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{} \\
 \text{Task } & \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{} & \boxed{}
 \end{array} \times \begin{array}{c} b \\ \boxed{} \end{array} = \begin{array}{c} y \\ \boxed{} \end{array}$$

A: $n \times n$ matrixy: resultant vector $n \times 1$ b: $n \times 1$ vector

- Each $y(i)$ of the resultant vector is the dot product of entire i^{th} row of A and the entire vector b.
- We can consider computation of $y(i)$ as a task.
- Following points can be made about the tasks:-
- 1) All n tasks are mutually independent.
- 2) No task has to wait for another task to finish.
- 3) No data dependency exists betⁿ tasks so they can be executed in any order.

q. Task dependency graphs:-

- Acyclic graph where nodes are tasks and directed edges are dependencies amongst them.
- Possible dependencies among task and their order of execution

- can be represented pictorially by task dependency graphs.
- consider a query to fetch the following data with the following conditions :-

CIVIC AND 2001 AND (White OR Green)

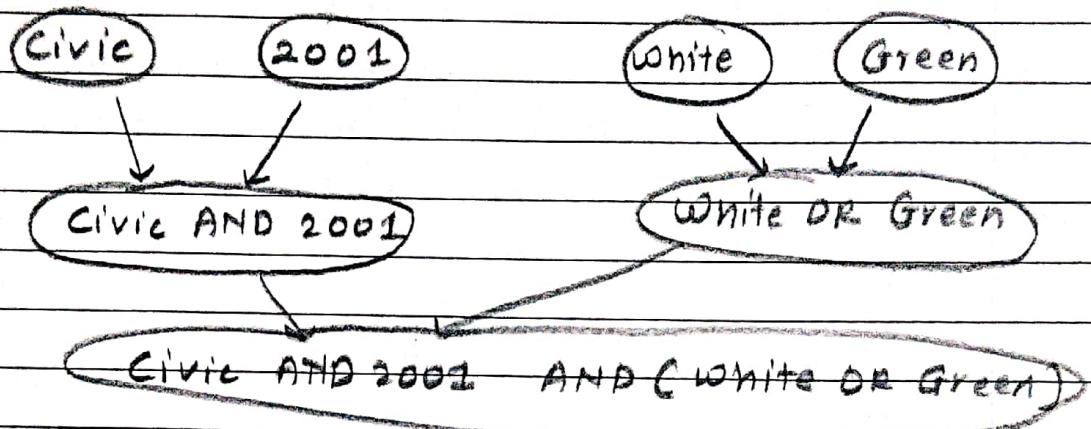


Fig : Task dependency graph.

- Granularity, concurrency and Task interaction :-
- Size of tasks = Granularity.
- Granularity in a parallel algo can be:- i) Fine or ii) Coarse
- i) Decomposing a computation into large number of small tasks is called fine-grained granularity.
- ii) If large subroutines of an algo are independent of each other, then they can be executed in parallel. This is coarse-grained parallelism. ∵ Coarse-grained decomposition of a computation into a ^{small} large no. of ~~one~~ large task.
- Degree of concurrency :-
- Number of tasks that can be executed in parallel.
- Max tasks that can be done is maximum degree of concurrency.

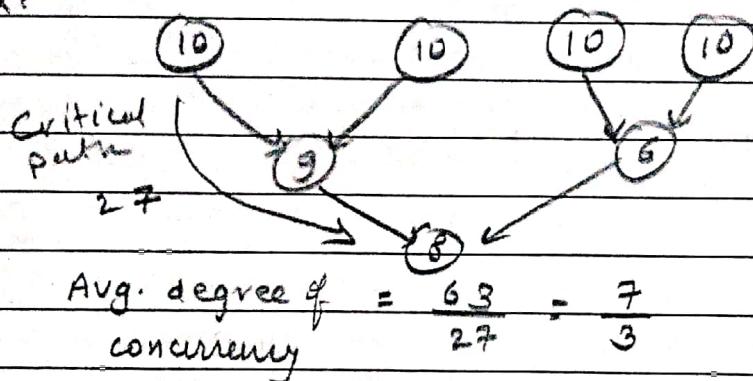
$$\text{Degree of concurrency} \propto 1$$

$$\text{Granularity}$$

- If dependency graph is a tree, maximum degree of concurrency = no. of leaves in tree.
- longest-directed path b/w start & finish node = critical path.
- The sum of weights of nodes along critical path is called critical path length.
- Ratio of total concurrency :-

$$\text{Avg. degree of concurrency} = \frac{\text{Total amount of work}}{\text{critical path length}}$$

Ex:-



- Dense matrix can't do more than n^2 concurrent tasks hence there's always a limit to granularity.

Q Task Interaction graphs:-

- Graphs of tasks and their data exchange is known as task interaction graph.
- Task interaction = Data Dependency whereas task dependency = control graph.
- This is required in case of sparse matrix computation.
- Processing & Mapping:-
- Process is the computation agent that performs tasks.
- Mapping of tasks to processes is called mapping.

Assignment

- The mapping scheme should exploit maximum concurrency.

Q Decomposition Techniques :-

- Job of dividing problems into sub-problems is called decomposition. Techniques to do this:-

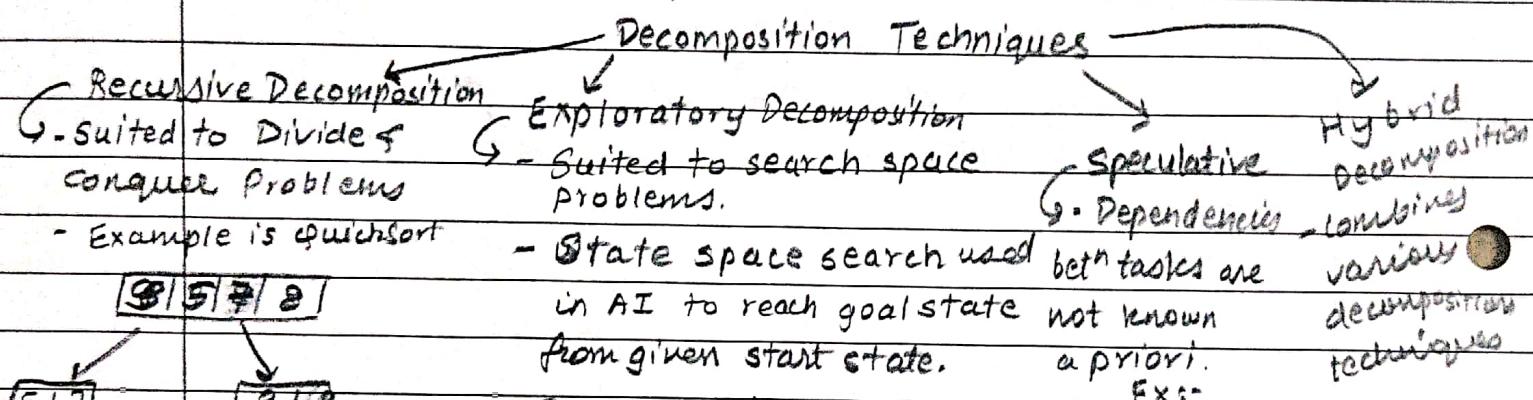
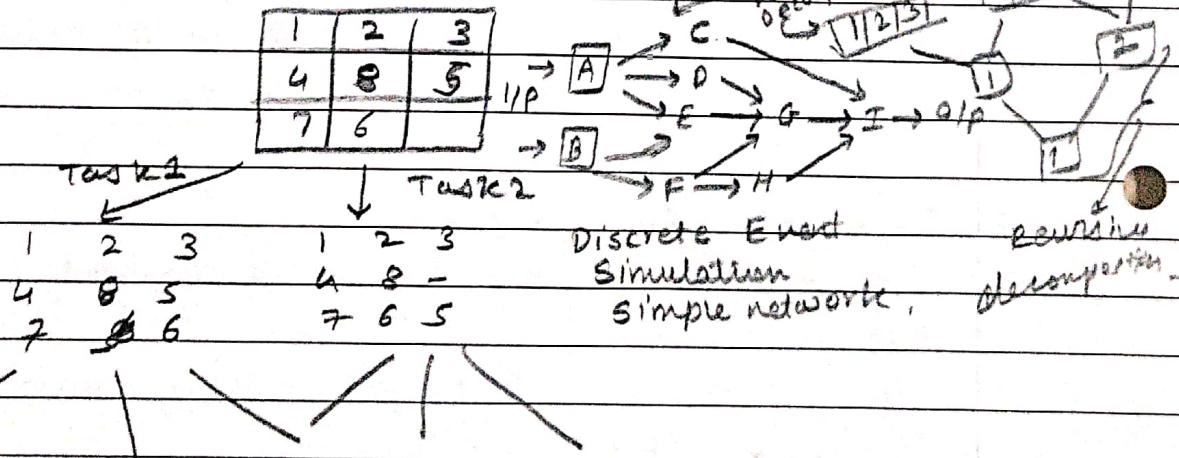


Fig Task dependency graph

- As we move down concurrency increases.



- Our problem is decomposed into 2 tasks

- In exploratory only one task can find a soln and as it finds it all tasks are terminated.

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Q Characteristics of Tasks & Interactions :-

- Good mapping depends on tasks & interactions.
- The effective parameters are :-

— Characteristics of tasks

Task Generation

Task Size

G-Time needed

Static

Dynamic

In decomposing
data centric apps
decomposition of size
of data and operation
is known apriori.

Tasks can be
generated by execu.
by algo.

This is static
task gen..

- Tasks not
known
apriori

- exploratory
decomposition

Examples.

- This is
dynamic

Task Gen.

to complete task
is the size of task.

- Matrix multiplication

partitioned is done

for each row.

for each task

as size is uniform.

- In case of quicksort

partitioned is done by

pivot which is

random. This is

non-uniform tasks.

size of data
associated with
a location
if size is known
or overhead of
data movement
is avoided.
knowledge of - size of
input.

Task sizes
- If knowledge of

size of data is
known or not.

- Its known

for matrix X

but not for

8-puzzle

problem.

* Characteristics of Inter-Task Interactions :-

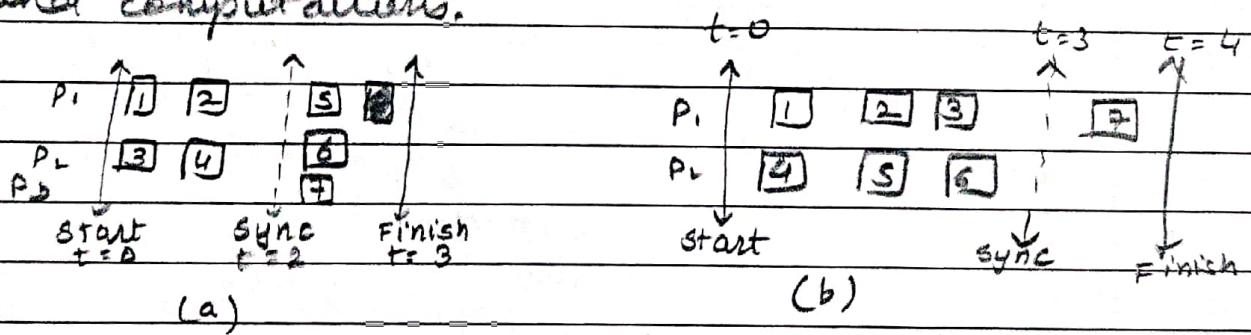
- i) Static Interaction :- Interaction happen at predetermined times. EX:- MX
- ii) Dynamic :- Interaction timings not known apriori. EX:- 8 puzzle.
- iii) Regular :- Interaction has a structure that can be exploited for efficient implementation. EX:- Image compression/ Dithering
- iv) Irregular :- No regular pattern exists. EX:- Sparse matrix X,
- v) Read-only interaction . vi) Read write Interaction ,
- vi) One-way interaction :- One pair of communicating tasks initiates and completes the task. vii) Two way interaction :- output of one is input to other .

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Q) Mapping technique for load balancing :-

- Mapping of tasks to process should be done in less time.
- Overheads that occur in decompositions are :-
 - (i) Inter-process interaction time. (ii) Time for which processes are idle.
- Good mapping shall aim to reduce this overheads.
- Due to load imbalancing some work is finished early.
- To balance load one may put tasks having interactions in the same process, however as some tasks may have different execution time pre optimum mapping is not obtainable.
- Good mapping schemes achieve balance bet'n interaction and computations.



Two configs for mapping (b's bad)

Types of Mapping

Static

- Tasks distributed by program exec.
- Task size known apriori
- Mapping schemes

Dynamic

- Task distribution during exec.
- Task gen. & map. done dynamically.
- Gets Ls in shared-addr space

Data partitioning

Task partitioning

Array distribution schemes :-

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

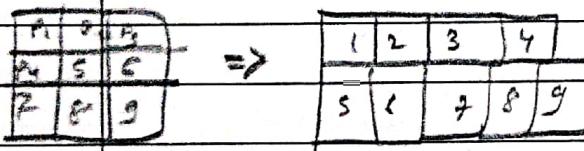
Schemes for static mapping

Data Partitioning

Array distribution schemes

- Tasks are responsible for exec. of data held by item.
- This is owner computes rule
- Techniques used

Block distributions:-
 - Uniform contiguous array portions distributed to different processes.



Task Partitioning

Randomized block distribution

- similar to cyclic but blocks randomly distributed.

Graph partitioning

- sparse data with highly irregular interactions

Cyclic & Block

- load balancing

mesh is divided in p parts, each part is assigned to a process p_i .



Assign blocks with repeating cycles as $(i \% p)$

- Such distributions are called cyclic distribution.

* Minimizing Interaction Overhead:-

- 1) Maximize data locality. 2) Minimize data Volume Exchange.
- 3) Minimize frequency of interactions. 4) Overlapping computation with interactions.
- 5) Replicating data communication computations.
- 6) Using group comm instead of point-to-point primitives.
- 7) Overlap interactions with other interactions.

MANJARA CHARITABLE TRUST
RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
 JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

Schemes of Dynamic Mapping

Centralized

Distributed

- Master / Slave model for processes.
- When slave ps run outta jobs they request master.
- With more ps master may bottleneck.
- To solve this processors pick chunk of tasks. This is chunk scheduling.
- Any process can send/receive work from other processes.

Parallel Algorithm Models :-

- Different ways to structure parallel algorithms to run on parallel system

Models

Data- Parallel

- Identical ops on different data concurrently.

- supports shared address-space

- message passing paradigm

- uses locality preserving

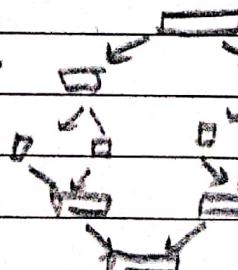
- decomposition and

- overlapping computation.

Task-Graph

- Task dependency graph is used to promote locality / reduce interaction costs.

- DFC algos.



Work-Pool

- Dynamic mapping
- Use message passing paradigm.

- Process generates work and add it to global pool.

Work Pool

Master-Slave

- 1 or > generate work and assign it to worker processes.

- Allocation can be static / dynamic.

Master

Slave

Slave

g) Sequential Computational Complexity :-

- $T(n)$:- maximum no. of ops performed by algo for a given input n .

g) Parallel Computational Complexity :-

- If p is the no. of processors then,
 - if $p \geq 1$ it is bounded parallelism.
 - if $p = \infty$ or $p \rightarrow \infty$ it is unbounded parallelism.
- Time complexity of a para algo to solve a problem of size n is given by a funcⁿ $T(p, n)$ which is the max time tht that elapses betⁿ start of algos execution by one processor and is terminated by 1 or ≥ 1 processor.
- Elementary ops are arithmetic or logical ops performed locally by a processor. They take $O(1)$ time.
- A problem is said to belong in NC (Nick's class) if it can be solved in polylogarithmic time using at most polynomial no. of processors.
- Problems in NC are thought to be parallel.

g) Anomalies :-

- If we write an unbounded parallel algorithm, the performance of the algo cannot be made better beyond a certain limit.
- This is so as some information needs to be known before further computation proceeds.
- The unbounded parallelism reflects this characteristic of a problem and is known as anomalies.
- In a balanced distribution each processors has $\lceil n/p \rceil$ or $\lceil P/n \rceil$ items

MANJARA CHARITABLE TRUST

RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI
JUHU VERSOVA LINK ROAD, VERSOVA, ANDHERI (WEST), MUMBAI - 53.

- + Types of anomalies:-
- If increase in P increases performance it is acceleration anomaly if it performs worse it is deceleration anomaly.