

## Annexure No.:

↳ Communication Cost in Parallel M/c's.

Cost of communication depends upon.

↳ Model semantics.

↳ Network topology

↳ Data handling & routing

↳ Total time to transfer a message over n/w consists of following parameters.

- Start up time ( $t_s$ )

: time to spent for sending & receiving nodes. Executing routing protocols.

- Per hop time ( $t_h$ )

: time is a function of hop & factors  $\rightarrow$  no. of total routers.  
n/w delay, switch latency.

- Per word transfer time:

Overheads determined by the length of message.

## Unit 3

### Principles of Message Passing Paradigm

#### To Introduction

↳ The Basic Building blocks : Send / Rec. operation

#### To Send And Receive Protocols

#### To Message Passing Interface.

↳ min. set of MPI routines.  
Terminating

↳ starting the MPI library

↳ Communicators — MPI\_Comm, MPI\_com-  
world

#### To Every Informati

↳ MPI\_Comm\_Rank

MPI\_Comm\_size

#### To First MPI program.

↳ Sending & Receiving Messages.

↳ MPI datatypes -

→ Communication Model of Parallel platforms ←

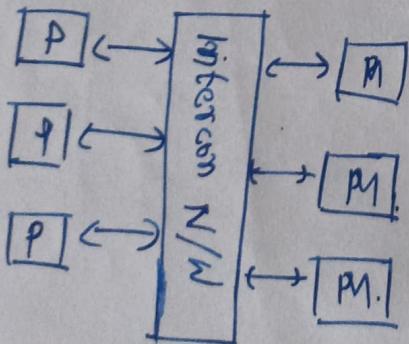


Shared data space

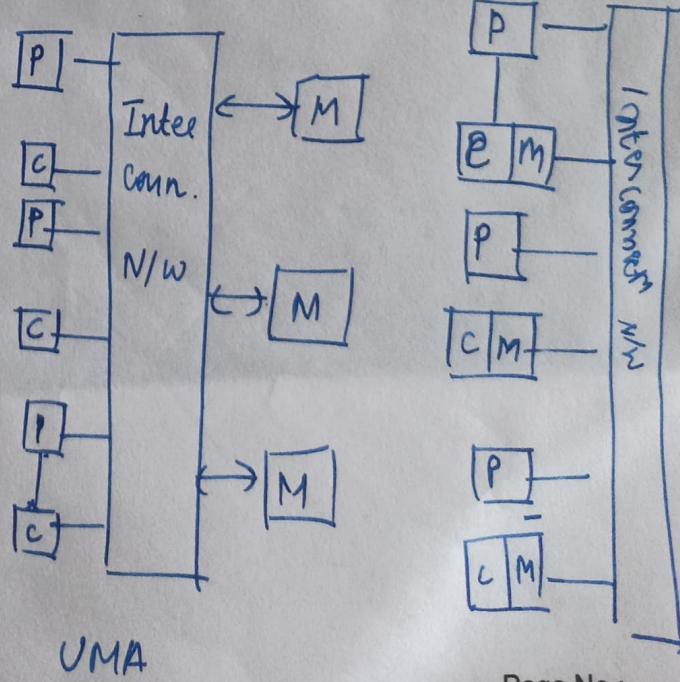


Shared Messages

- ↳ Part of all memory is accessible to all processor
- ↳ Time taken by a processor to access any memory word in the system global or local is identical the platform is the Uniform memory Access (UMA)
- Else is Non uniform memory Access (NUMA)
- ↳ All the processor share a common memory space and access the same data.
- ↳ Can be implemented using hardware & software.
- ↳ Adv: Exploit locality Reducing overhead latency.



Shared data space



## Annexure No.:

## → Message Passing

↳ communication method where processor has its own memory to send & receive message to and from processor.

↳ libraries: MPI (message passing interface)  
PVM (Parallel Virtual Machine)

Ad: Scalability, to distribute across different distributed machines.

### Libraries & function:

MPI send, MPI receive

Annexure No.:

→ Unit: 2 ←

## Decomposition Techniques :

### ① Task & Dependency Graphs

Step 1 to implement parallel computing is to break complex problem into tasks.

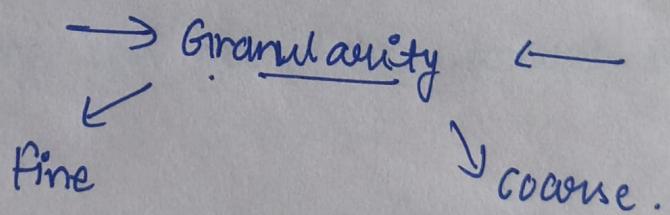
Task : Indivisible unit of computation.

Representation : Directed graph.

↳ It represents the control dependency.

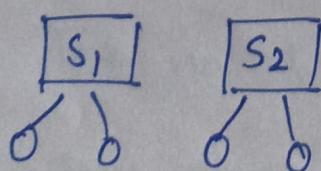
Task Dependency Graph

↳ Example: Data Query Processing.



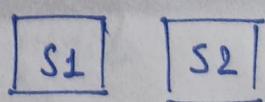
① smaller comp.

Ex: Op of comp. services.



① Large comp.

Ex: composite services.



② More objects holds less data.

Enrollment No. :

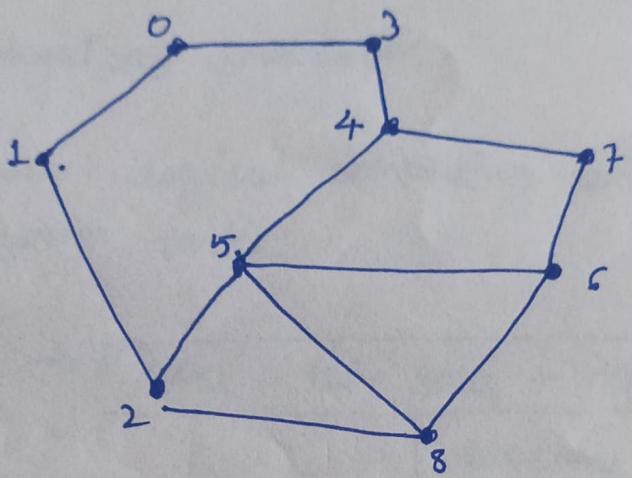
② Few object holds more data.

Page No. :

Annexure No.:

## → Task Interaction Graph ←

- ↳ Graph that captures the pattern of interactions among tasks.
- ↳ Represents data dependencies.



## → Process Limitations of Parallel Com ↵

- ① fine granularity can not hold parallel time arbitrariness small.
- ② Task Interaction graphs also a limiting factor.

## → Process Mapping ↵

- The no. of task > NO. of resources.
- Thus parallel algo. must provide the mapping of task into processes.

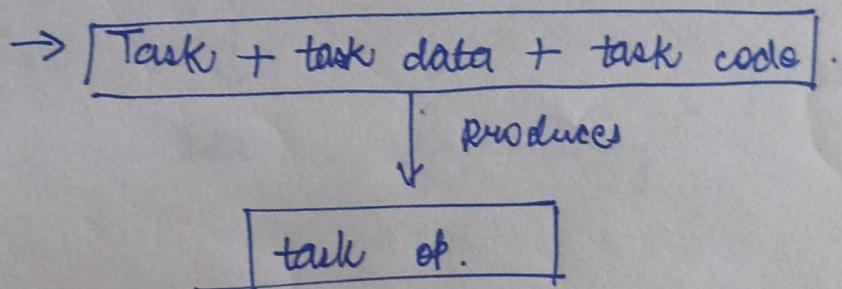
Mapping is done from task to process.

bcz: the typical prog. APIs do not allow the easy binding of task to physical processor.

↳ Task → Processes to map these process.

to physical processor.

↳ Process : logical computing sense / agents that performs tasks.



↳ Processors : Physical hardware unit that perform task

→ Mapping is done by observing the task dependency graphs.

# Mapping must minimize parallel time by :

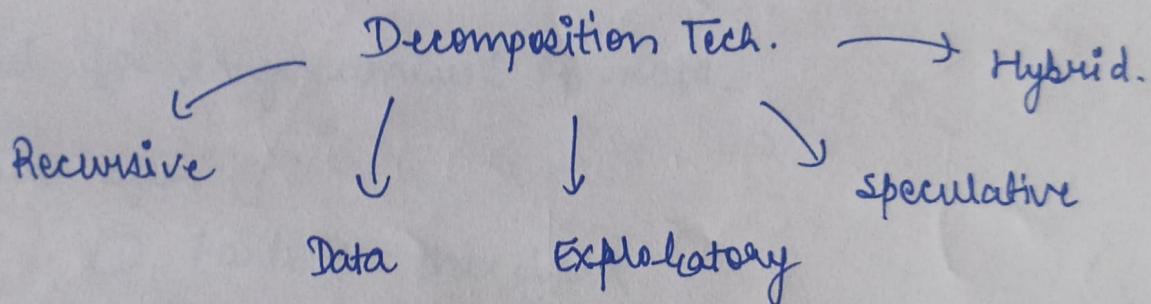
→ Properly mapping independent task to dif process.

→ Assigning task on critical path to processes accp. they available .

## Decomposition Techniques :

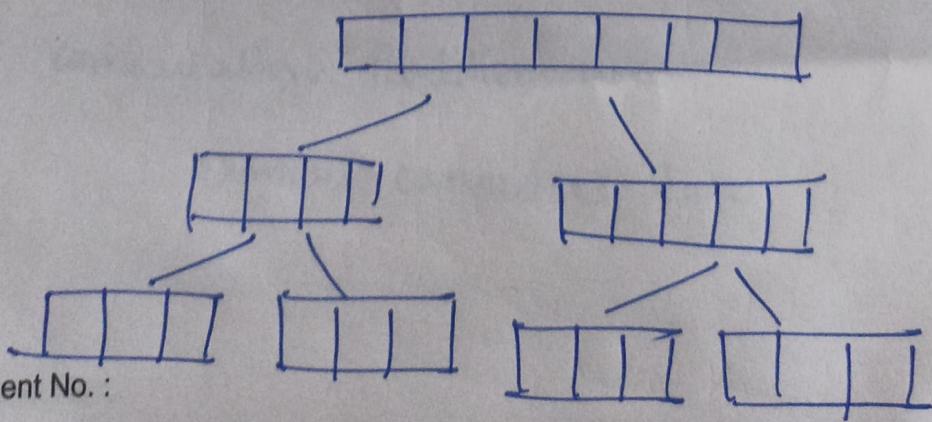
The procedure to break down the complex problem into simpler sub problems & task. which are indivisible.

↳ There is no concrete path for solving all subproblem



### ① Recursive Decomposition

- ↳ Use to solve problems like divide & conquer.
- ↳ Complex problems are broken into sub problems until a desired granularity reach.



Annexure No.:

$\min[A] = \infty$  min = A[0];

for i:1 to (n-1) do

If (A[i] < min) min := A[i];

endfor;

return min;

## ② Data Decomposition

- ↳ Used for large amount of data.
- ↳ Two steps
  - : ① Partition the data.
  - ② Induce comp. partition from data part.
- ↳ Type of Data to be partitioned.
  - ↳ Input
  - ↳ Output
  - ↳ Intermediate . } Uses different data decom. techniques.
- ↳ Computation Partition uses
  - : Owners computers rule

Example: Matrix Multiplication

## ① Partitioning O/P data.

Applied : O/p data can be comp. independently  
as a function of w/ input.

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \cdot \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix} \Rightarrow \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix}$$

$$C_{11} \rightarrow A_{11} B_{11} + A_{12} B_{12}$$

$$C_{12} \rightarrow A_{11} B_{12} + A_{12} B_{22}$$

$$C_{21} \rightarrow A_{21} B_{11} + A_{22} B_{12}$$

$$C_{22} \rightarrow A_{21} B_{12} + A_{22} B_{22}$$

Conclusion

- ↳ Database of transact<sup>n</sup> is replicated across process  
each task can be independently comp. w/o comm.
- ↳ Database of trans. is partitioned across process  
each task first computes partial counts -
- ↳ These task are then aggregated at app. task.

## Input Data Partitioning

- ↳ This induces a task decom in which each task generates partial count for all itemset.
- ↳ Combined subsequently for aggregate count.

## Intermediate

Sequence of transaction from i/p to o/p.

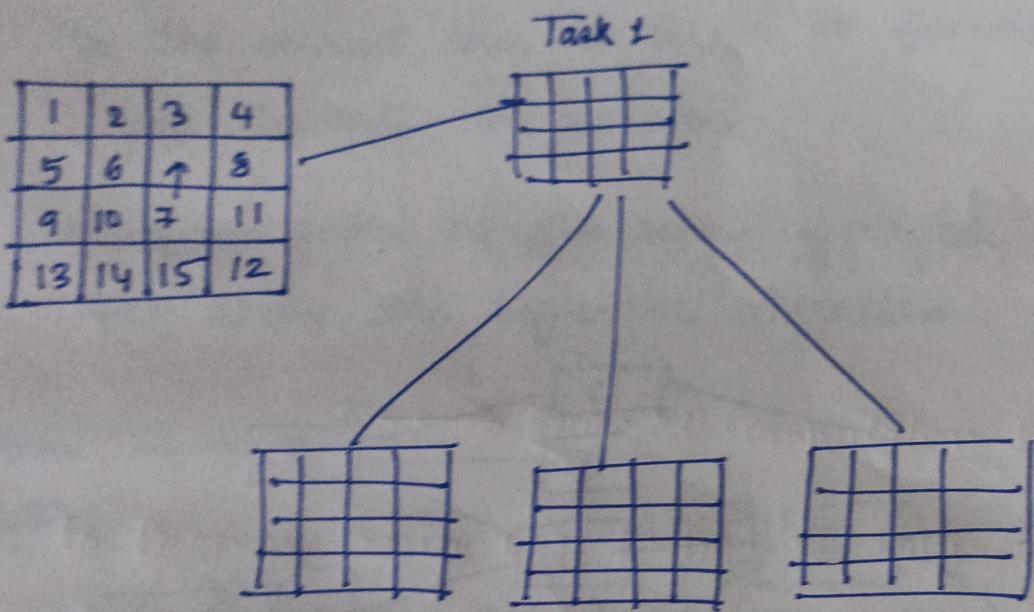
→ Owner's Computation Rule ←

- ↳ Process assigned a particular data item is responsible for all its computation.
- ↳ Input Partitioning : Computations that uses i/p data are performed by process.
- ↳ O/p partitioning: computation that uses o/p are computed by the process

## Exploratory Decomposition

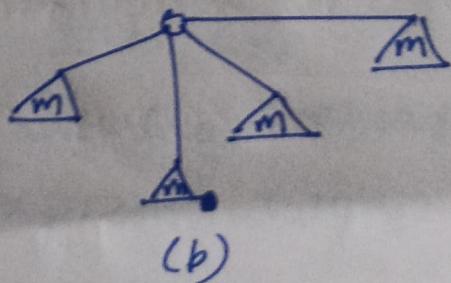
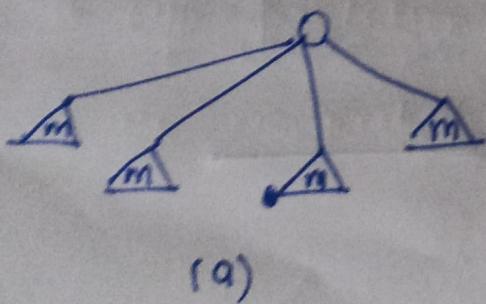
- ↳ used when computations that corresponds to search of space of solutions.

Ex: 15 - puzzle problem.



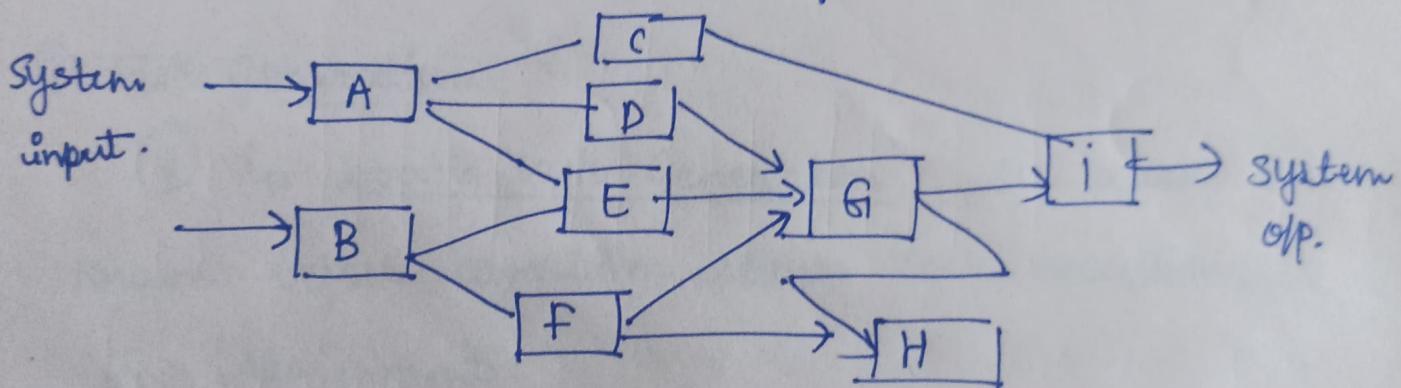
↳ Can cause change in amount of work done by the parallel formulation

↳ Change results in super or sub linear speedups.



## → speculative Decomposition

- ↳ This decomposition is based on such problems in which solutions are based on the current task finishes.
- ↳ The current task is based to predict the next task towards the solution.
- ↳ Perform more or the same aggregate work (but not less) than the sequential algorithm



## → If predictions are wrong

- ↳ Work is wasted.
- work is needed to be done
  - ↳ state - rastering overhead  
(memory / computations)

→ it may be the only way to extract concurrency.

## # The affect of good mapping on characteristics of Task

↳ The following are the four characteristics that can be have huge influence on the suitability of Task.

- ① Task Generation      ② Task Sizes
- ③ Knowledge of Task size      ④ Size of Data associated with task.

① Task Generation

① ↳ static task Generation: All tasks are known before algorithm starts task execution.

Ex: Assignments:

↳ Data Decomposition uses static task generat<sup>n</sup>.

Ex: Matrix Multiplication, LU factorizat<sup>n</sup>.

↳ Recursive Decomposition also uses static task generation (find min of a list of numbers)

② ↳ Dynamic Task Generation: All task are not available prior to before starting task execution algo.

↳ Recursive Algo. leads to dynamic task generat<sup>n</sup>.

Enrollment No.: (quick sort : the task are generated dynamically) Page No.

② Task size : Size of task is can be used to predict the the amount of time taken by a process

↳ Complexity often depends upon task complexity uniformity.

③ Knowledge of Task size: If the size of task is known, this information can be used in mapping of task to processes.

④ Size of Data associated with a task: data associated with the task must be available to the process performing that task.

→ Characteristics of Inter Task Interactions ←

① Static V/s Dynamic

↳ The task generation would be static if : for each task the interactions happened at predefined time.

↳ and, <sup>before</sup> algorithm execut<sup>n</sup> the tasks are known priori

↳ Static T.P. can be programmed easily in Message Passi the reason is interaction in message passing requir. active involvement of sender & receiver.

② Dynamic :

↳ The task generation would be dynamic if : each task the interact<sup>n</sup>s happened at non-uniform time.

↳ Tasks are not explicitly nor known priori before the execution of algorithm.

↳ Dynamic T.P. can not be easily programmed in Message Passing. because of inactive involvement of sender & receiver

## ② Regular v/s Irregular

→ Way of classifying interactions based on their spatial structure

### Regular

Patterns are regular

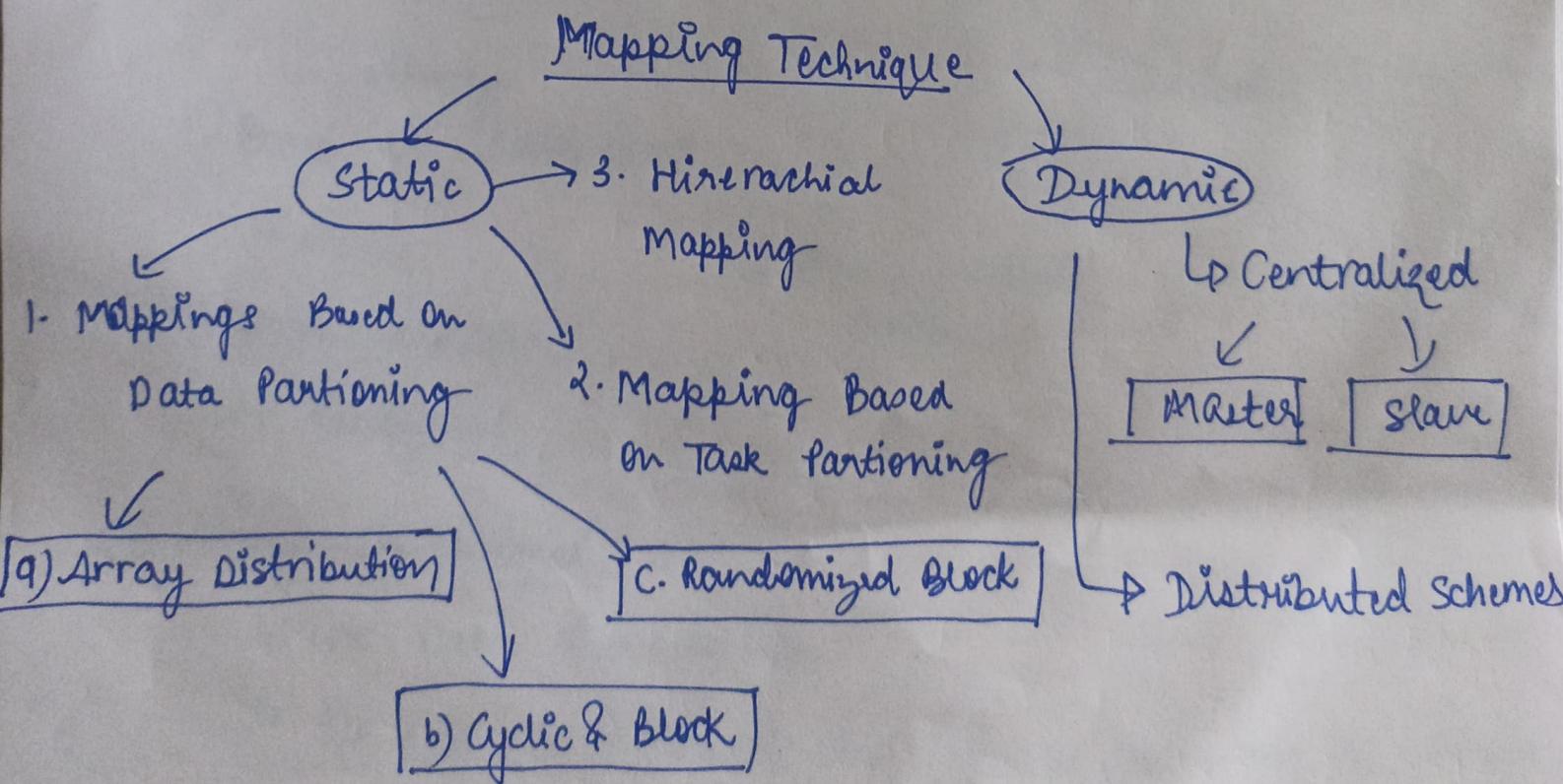
Ex: Image Dithering

### Irregular

Patterns are irregular

→ Mapping Techniques for Load Balancing ↗

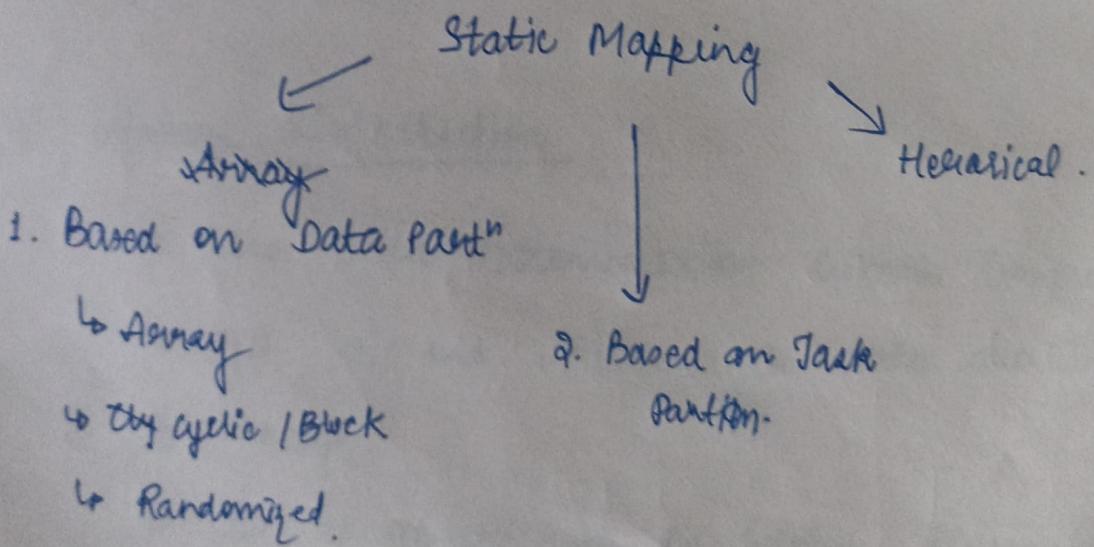
① Mapping techniques is used to map from task to process so that the execution task is done within small elapsed time



## Objective of Mapping

- ① Task should complete within shortest amount of time.
- ② Mapping must address following sources-
  - ↳ Load Imbalance.
  - ↳ Inter process communication.

Static Mapping: Task are mapped prior to process.  
 Task generated statically.  
 Have uniform computational requirement.



### ① Array Dis. Schemes

↳ Use Data Decomposition

↳ If, Of, If intermediate data are in the form of arrays

→ Based on Data Partitioning

- ↳ Combine partitioning with Owner's Computation Rule.
- ↳ Simplest Data Decompositn: 1-D block distribut.

Row wise distribut.

P <sub>0</sub>
P <sub>1</sub>
P <sub>2</sub>
P <sub>3</sub>
P <sub>4</sub>
P <sub>5</sub>
P <sub>6</sub>
P <sub>7</sub> .

Col. wise distribut.

P <sub>0</sub>	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub> -P <sub>8</sub> .

### (a) Array Distribution

- ↳ Use Data Decomposition: Owner's Computation Rule.
- ↳ S1P ; O1P and intermediate data are in form of arrays.
- ↳ Can be generalized to higher dimensions.

### [Array Block Distribution]

- ↳ Block is the simplest way to distribute an array.
- ↳ This leads to have uniform contiguous portions.

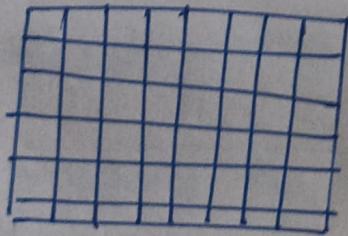
(b) Cyclic & Block. Ex: LU factorizat<sup>n</sup>

↳ Block-cyclic is variation of Block distribution, and can be used to alleviate the load balancing & idling problems.

↳ How to solve idling problems. :

→ All process have sampling of tasks from all parts of the matrix.

↳ As a result, even if different parts of the matrix require different amount of work, the overall work on each process balances out.



$t_{P01}$	$P_1$	$P_2$	$P_3$
$P_1$	$P_5$	$P_6$	$P_7$
$P_3$	$t_9$	$P_{10}$	$t_{11}$

## (c) Randomized Block.

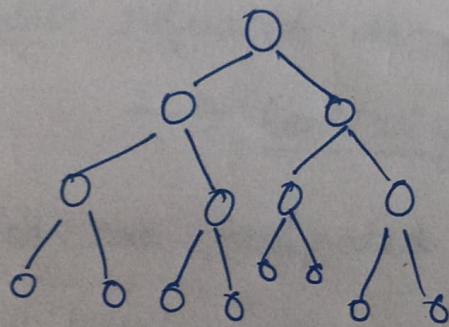
↳ Used where work has special patterns.

↳ Block could contained non-zero elements due to which the other processes may not get work / or remained idle.

Annexure No.:

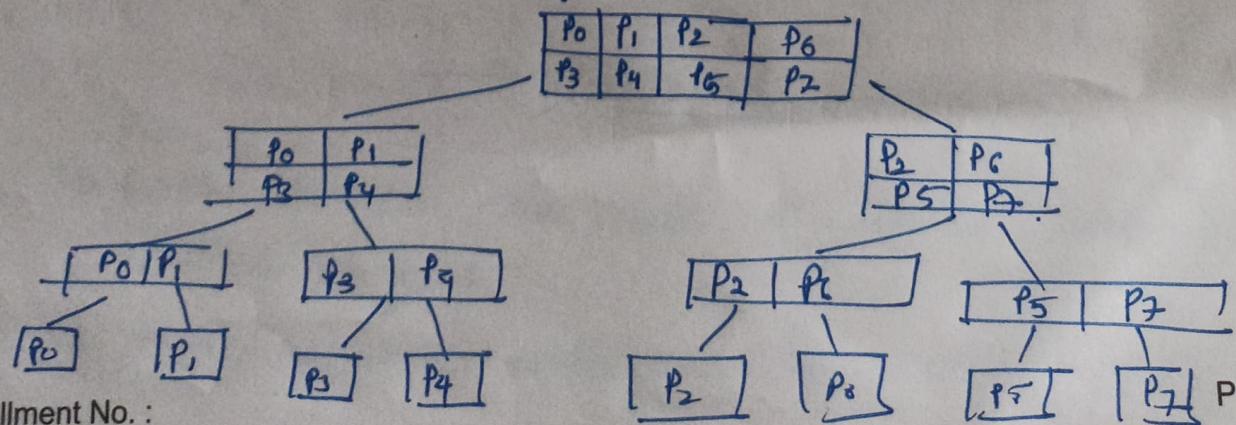
→ Here Based on Task Part " ←

- ↳ Partitioning a given task dependency across processes.
- ↳ Determining optimal mapping for task dependency graph is NP-complete problem.



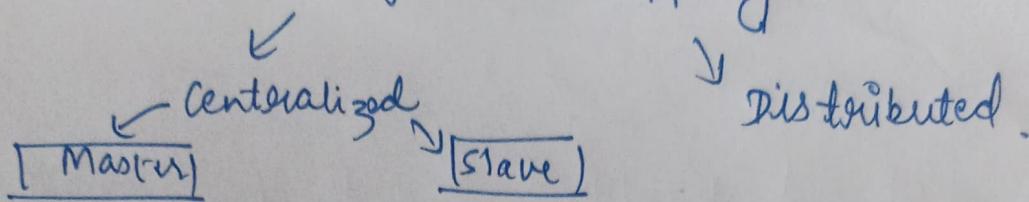
→ Hierarchical Mappings ←

- ↳ Single mapping have some limitations.
- ↳ task mapping of binary tree can not use a large numbers of process.
- ↳ For this reason task mapping is done at the top level and data partitioning within each level.



Annexure No.: \_\_\_\_\_

## Dynamic Mapping



- ↳ Dynamic Mapping do not know the task priori.
- ↳ It is also referred as dynamic load balancing.
- Centralized ←
- ↳ Processes are designated as 1) Master 2) Slaves.
- ↳ When a process runs out of work, it requests the master for more work.
- ↳ If the no. of process is increased.  
then, master becomes bottleneck.
- ↳ To solve this process may pick up no. of tasks called chunks at a time.

- ↳ This called chunk scheduling.
- ↳ Selecting large chunks = load imbalancing.

→ Distributed ←

- ↳ Each process can send or receive work from another process.

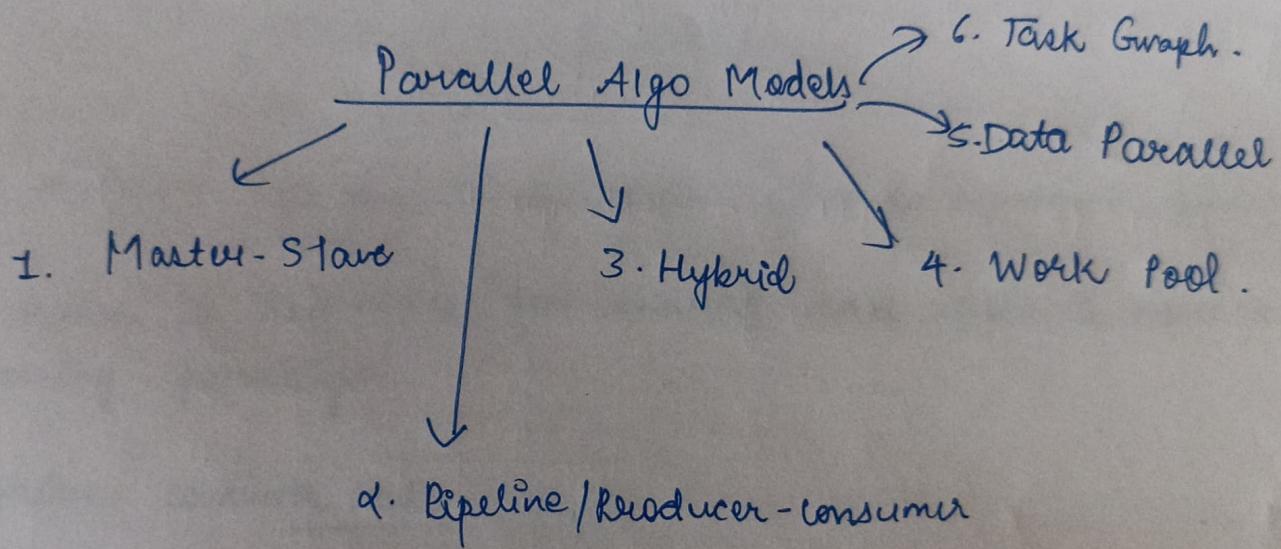
- ↳ This alleviates bottleneck in centralized system.



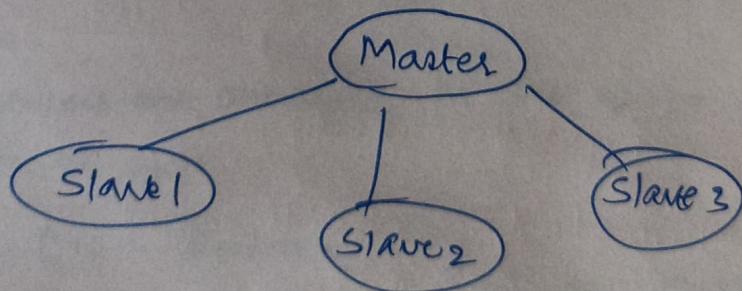
Annexure No.:

## → Parallel Algorithm Models ←

- ↳ It is a way of structuring a parallel alg. by
- Selecting: Decomposition & Mapping Techniques.
- Apply: Appropriate Strategy to minimize interact's.



### ① Master-slave Model



Annexure No.: \_\_\_\_\_

## Master-Slave Model

↳ One processor is Master; other units acts as slave.

### Responsibilities of Master

- ① Coordinate: activities among slave.
- ② Generates: works to be performed.
- ③ Assign: work to slaves.
- ④ Synchronize: Activities of slaves after each start. phase.

↳ Model is suitable for shared add space & message passing paradigm.

## # Pipeline: Consumer Model

↳ also called producer/consumer model.

↳ Set of data passed through series of process, which perform some task.

↳ The process are arranged in the queue in the form Multi-D array tree, graphs or cycle.

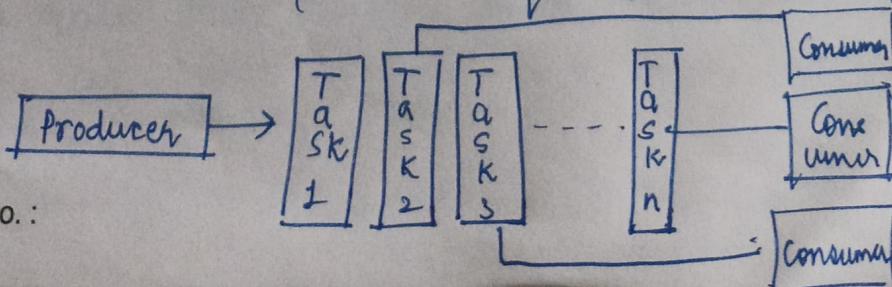
O/P = Consumer.

I/P = Producer

↳ Same as concept of pipeline

I/P of 1st phase = I/P for next phase

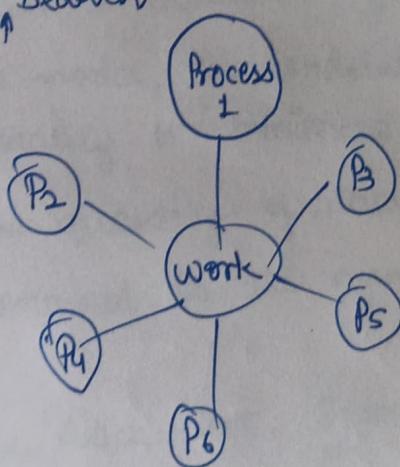
↳ Interact Minimization Technique used is Overlapping Interaction



# Work Pool :

- ↳ Task are dynamically assigned to processes for balancing load.
- ↳ Therefore any process may potentially execute any task.
- ↳ This model is used when the quantity of data associated with the task is comparatively smaller than the computation associated with task.
- ↳ If the task is generated dynamically & decentralized assigning of task is done, the termination detection algorithm is required.

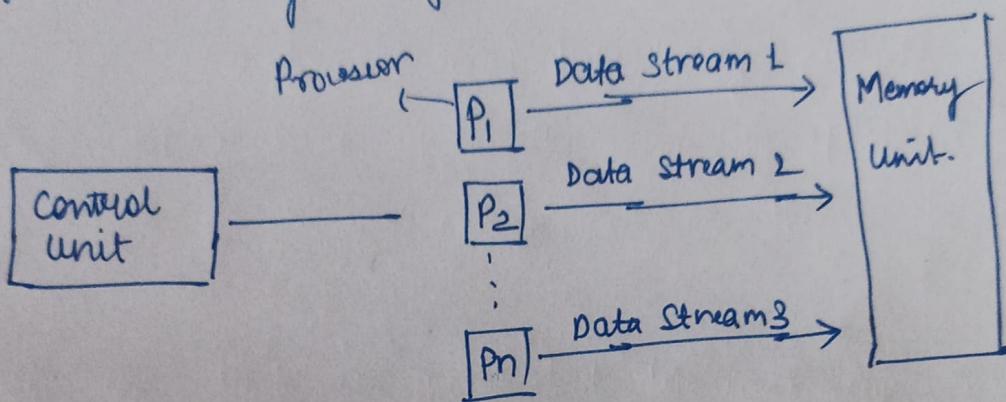
Ex: Parallel <sup>Tree</sup> Search

# Data Parallel :

- ↳ Task are assigned to process.
- ↳ Each task performs similar type of operations on diff data i.e. SIMD (single Inst. multi Data)

Characteristic: Intensity of Data parallelism increase with size of data.

use: for handling large data problems.



### # Task graph :

- ↳ Parallelism is expressed by task graph.
- ↳ In this model, the correlation of task are utilized to promote locality or minimize interact<sup>n</sup> cost.
- ↳ Use: the quantity of data associated with the task is huge compared to the number of computation associated with them.

Ex: Parallel Quick Sort, Sparse Matrix Multiplicat<sup>n</sup>

