

Course Title : Parallel and Distributed Computing

Course No: C.Sc. 544

Nature of the Course : Theory + Lab

Full Marks : 45+30

Pass Marks : 22.5+15

Credit Hrs: 3

Course Description:

- This course covers a broad range of topics related to parallel and distributed computing, including foundations, models, and algorithms.

Course Contents:

Unit 1: Foundations

1.1 Parallel and Distributed Computing: the Scene, the Pros, the Players (2 hrs)

A Perspective, Parallel Processing Paradigms, Modeling and Characterizing Parallel Algorithms, Cost vs. Performance Evaluation, Software and General Purpose PDC

1.2 Semantics of Concurrent Programming (2 hrs)

Models of Concurrent Programming, Semantic Definitions, Axiomatic Semantic Definitions, Denotational semantics Definition Operational Semantic Definitions

1.3 Formal Methods: A Petri Nets Based Approach (2 hrs)
Process Algebras, PETRI Nets, High-Level Net Models

1.4 Complexity Issues in Parallel and Distributed Computing (3 hrs)
Introduction, Turing Machines as the Basis and Consequences,
Complexity Measures for Parallelism, Parallel Complexity Models
and Resulting Classes, VLSI computational Complexity,
Complexity Measures for Distributed Systems, Neural Networks
and complexity Issues.

1.5 Distributed Computing Theory (3 hrs)

The computational Model, A simple Example, Leader Election,
Sparse Network Covers and Their Applications, Ordering of
Events, Resource Allocation, Tolerating Processor Failures in
Synchronous systems, Tolerating Processor Failures in
Asynchronous Systems, Other Types of Failures, Wait-Free
Implementations of shared Objects.

Unit 2 : Models

2.1 PRAM Models (5 hrs)

Introduction; Techniques for the Design of Parallel Algorithms;
The PRAM Model; Optimality and Efficiency of Parallel
Algorithms; Basic PRAM Algorithms; The NC-class; P-completeness
Hardly Parallelizable Problems, Randomized Algorithms and
Parallelism; List Ranking Revisited; Optimal $O(\log n)$
Deterministic List Ranking; Taxonomy of Parallel Algorithms
Deficiencies of the PRAM Model.

2.2 Broadcasting with Selective Reduction : A Powerful Model of Parallel Computation (3 hrs)

Introduction; A Generalized BSR Model; One Criterion BSR Algorithms; Two Criteria BSR Algorithms; Three Criteria BSR Algorithms; Multiple Criteria BSR Algorithms

2.3 Dataflow Models (3 hrs)

Kinds of Dataflow; Data-Driven Dataflow Computing Models; Demand-driven Dataflow Computing Models; Unifying Data-Driven and Demand-Driven

2.4 Partitioning and Scheduling (4 hrs)

Program Partitioning; Task scheduling; scheduling system model; communication models; Optimal scheduling Algorithms; Scheduling Heuristic Algorithms; Scheduling Nondeterministic Task Graphs; Scheduling Tools; Task Allocation; Heterogeneous Environments

2.5 Checkpointing in Parallel and Distributed System (3 hrs)

Introduction; Checkpointing using Task Duplication; Techniques for consistent checkpointing

2.6 Architecture for open Distributed software System (2 hrs)

Introduction to Open Distributed System Architecture; Computational Model; Engineering Model; ODP Application

Unit 3 : Algorithms

3.1. Fundamentals of Parallel Algorithms (3 hrs)

Introduction; Models of Parallel Computation; Balanced Trees;
Divide and conquer, Partitioning; Combining

3.2. Parallel Graph Algorithms (5 hrs)

Graph-Theoretic Concepts and Notation; Tree Algorithms; Algorithms
for General Graphs; Algorithms for Particular Classes of Graphs

3.3. Data Parallel Algorithms (5 hrs)

Chapter Overview; Machine Model; Impact of Data Distribution;
CU/PE Overlap; Parallel Reduction Operations; Matrix and Vector
Operations; Mapping Algorithms onto Partitionable Machines;
Achieving Scalability Using set of Algorithms

Recommended Books

Unit-1

Date |

Page |

Unit 1.1 : Parallel and Distributed Computing

Introduction

- In parallel computing, all processors may have access to a shared memory to exchange information between processors.
- In distributed computing, each processor has its own private memory, where information is exchanged by passing message between processors (RPC, RMI)

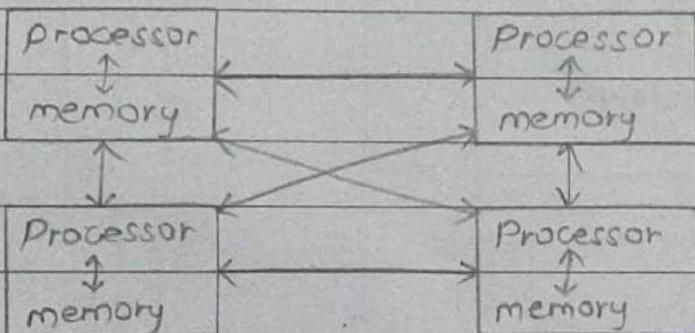


Fig: Distributed computing

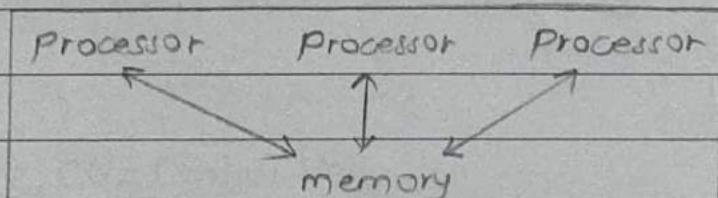


Fig: Parallel computing

§ Why do we need Parallel Computing?

- Serial computing is too slow
- need for more computing power (e.g. cryptography, data mining, search engine)

Problem

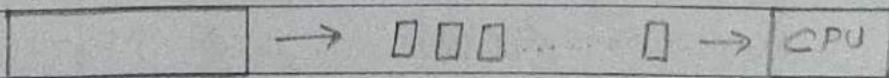


Fig: Traditional approach

Problem

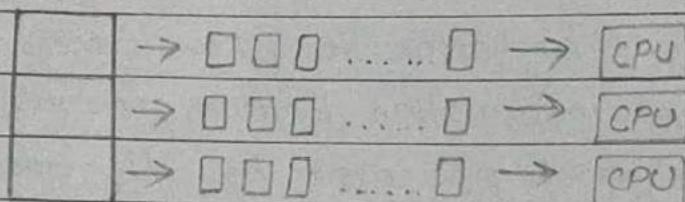


Fig: Parallel computing approach

Issues in sharing resources

Spooler Directory

Printer	:	
	:	
	:	
	:	
Process A	4 5 6 7	[out 4] → next to be printed
Process B		[in = 7] → next free slot to store the document
	:	
	:	
	:	
	:	

Parallel Processing Paradigms

- Flynn's Taxonomy

1. SISD (Single Instruction Single Data)
2. SIMD (Single Instruction Multiple Data)
3. MISD (Multiple Instruction Single Data)
4. MIMD (Multiple Instruction Multiple Data)

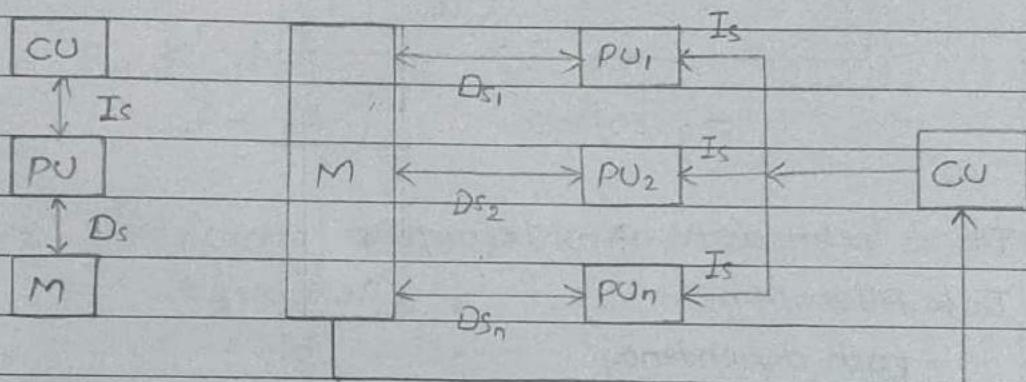


Fig: SISD

IS

Fig: SIMD

Where, CU = Control Unit

PU = Processing Unit

M = Memory

Is = Instruction Stream

Ds = Data Stream

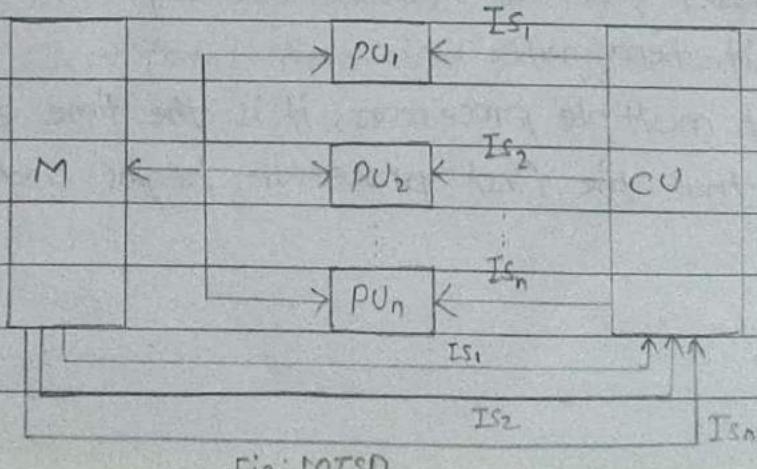


Fig: SIMD

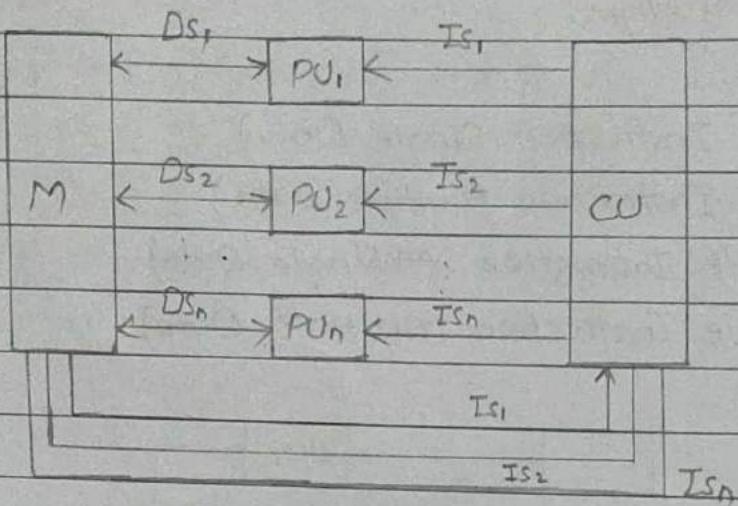


Fig: MIMD

Issues related to parallelization

1. Task Allocation

- task dependency
- Load balancing

2. Communication Time

- Communication time between the processors.
- Serious, when the number of processors increases to hundred or thousands.

3. Execution Time

- time elapsed from the moment the algorithm starts to the moment it terminates
- In case of multiple processors, it is the time elapse between the time that the first processor begins and the last one terminates.

4. Speedup

- Speedup (S) = $\frac{\text{Running time of the best sequential algorithm}}{\text{Running time of the parallel algorithm}}$

- it is not always possible to decompose the task
- so maximum speed of N processor system in executing algorithm is:

$$S_N \leq \frac{1}{f + \frac{(1-f)}{N}} \leq \frac{1}{f}$$

Where, $f \rightarrow$ fraction of computation task must be done sequentially.

^{no speedup}

$\hat{S}_{\max} = 1$, where $f = 1 \rightarrow$ every job sequentially

$\check{S}_{\max} = 0$, where $f = 0 \rightarrow$ every job parallelly

^{full speedup}

5. Communication Penalty

$$CP_i = \frac{E_i}{C_i}$$

Where, $E_i \rightarrow$ total execution time spent by P_i

$C_i \rightarrow$ total time used for communication by P_i .

Unit 1.2: Semantics of Concurrent Programming

- concurrent programming refers to activation of more than one instruction at a time
- semantics is the field concerned with the mathematical study of the meaning of the programming language.
- Syntactic vs semantic structure.

e.g.

English language

programming language

Categories

1. Axiomatic semantic Definition
2. Operation semantic Definition
3. Denotational Semantic Definition

Models of concurrent programming

1. Level of granularity
2. Sharing the clock
3. Sharing the memory (mutually exclusive)
4. Pattern of interaction
 - 4.1 Synchronization → mutual exclusion/admission
 - 4.2 Communication → synchronous / Asynchronous

1. Axiomatic Semantic Definition

- defines the meaning of language construct by making statement in the form of axioms or inference rules.
- Based on the Hoare's program (Hoare Triple), which is in

the form

$\{P\} \subseteq \{Q\}$ where, $P \rightarrow$ pre-condition
 $S \rightarrow$ statement
 $Q \rightarrow$ Post-condition

e.g. $\{x=5\} \quad x=x+1 \quad \{x=6\}$

- eg. 1. $\{j=3 \text{ AND } k=4\} \quad j=j+k \quad \{j=7 \text{ AND } k=4\}$
- 2. $\{x > 0\}$ while $(x \neq 0)$ do $x=x-1 \quad \{x=0\}$
- 3. $\{x \leq N-1\} \quad x=x+1 \quad \{x \leq N\}$

Partial Correctness

- if answer is returned, it will be correct

Total Correctness

- it requires that the algorithm terminates

Sequence Statement Rule

- It has the form

$$\frac{H_1, H_2, \dots, H_n}{H}$$

i.e if H_1, H_2, \dots, H_n have all been verified then H is also valid.

e.g. $x > 1 \quad \{x=n+1\} \quad x > 2$

$x > 2 \quad \{x=n+1\} \quad x > 3$

$x > 1 \quad \{x=n+1, x=n+1\} \quad x > 3$

i.e $\{P\} \rightarrow S_1 \{Q\}, \{Q\} \rightarrow S_2 \{R\}$
 $\{P\} \rightarrow S_1, S_2 \{R\}$

Alternation Statement Rule

IF - THEN

$\{P \text{ and } B\} \rightarrow \{Q\}$
 $\{P \text{ and NOT } B\} \rightarrow \{Q\}$
 $\{P\} \text{ IF } B \text{ THEN } S \text{ END IF } \{Q\}$

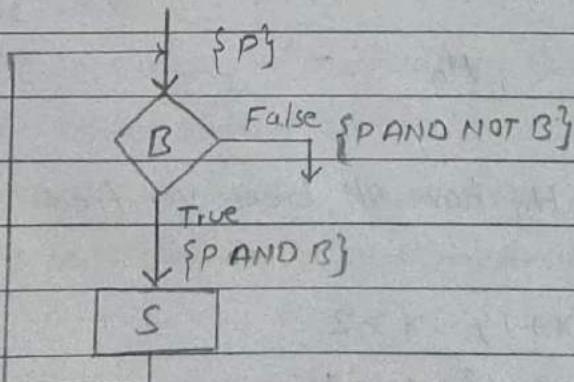
IF - ELSE

$\{P \text{ and } B\} \rightarrow S_1 \{Q\}$
 $\{P \text{ and NOT } B\} \rightarrow S_2 \{Q\}$
 $\{P\} \text{ IF } B \text{ THEN } S_1 \text{ ELSE } S_2 \text{ END IF } \{Q\}$

Loop

While

$\{P \text{ and } B\} \rightarrow \{P\}$
 $\{P\} \text{ while } B \text{ do } S \text{ END while } \{P \text{ and NOT } B\}$



e.g.

$x > 0 \wedge x < 10 \quad \{x = x + 1\} \quad x > 0$
 $x > 0 \{ \text{while } x < 10 \text{ do } x = x + 1 \} \quad x \geq 10 \wedge x > 0$

Disjoint Parallel Program

- Two programs s_1 and s_2 are said to be disjoint if none of them can change the variables accessed by the other.
- Denoted by $[s_1 \parallel s_2]$

Example:

$$\begin{array}{ll} s_1 : x = z & \xrightarrow{\text{but not}} s_1 : x = y + 1 \\ s_2 : y = z & s_2 : x = z \end{array}$$

- Semantics of this statement can be denoted as

$$\begin{array}{ll} \{P_1\} & s_1 \{Q_1\} \\ \{P_2\} & s_2 \{Q_2\} \\ \{P_1 \text{ AND } P_2\} & [s_1 \parallel s_2] \{Q_1 \text{ AND } Q_2\} \end{array}$$

Await then Rule

- Await T then B , where T is a condition and B is a block of code.
- If condition T holds, then block B is executed else the process is suspended until T becomes true.
- Semantically, it can be denoted as,

$$\begin{array}{l} \{P \text{ AND } T\} B \{Q\} \\ \therefore \{P\} \text{ Await } T \text{ then } B \{Q\} \end{array}$$

Operational Semantics

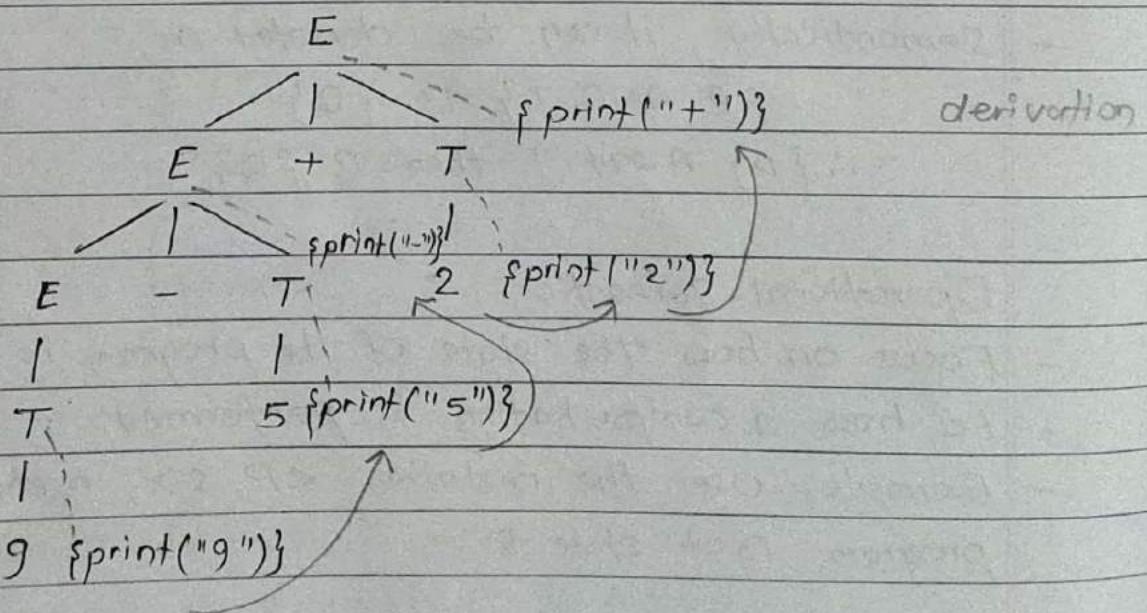
- Focus on how the state of the program is affected.
- i.e. how a computation is performed.
- Example, use the notation $\langle P, S \rangle$, means semantic of program P at state S

$\langle z = x, \quad x = y, \quad y = z, \quad [x \rightarrow 5, \quad y \rightarrow 7, \quad z \rightarrow 0] \rangle$
 $\langle x = y, \quad y = z, \quad [x \rightarrow 5, \quad y \rightarrow 7, \quad z \rightarrow 5] \rangle$
 $\langle y = z, \quad [x \rightarrow 7, \quad y \rightarrow 7, \quad z \rightarrow 5] \rangle$
 $\langle, \quad [x \rightarrow 7, \quad y \rightarrow 5, \quad z \rightarrow 5] \rangle$

- Example, Syntax directed translation can also be represented as operational semantics

Production	semantic Rule
$E \rightarrow E + T$	{printf("+)}
$E \rightarrow E - T$	{printf("-")}
$E \rightarrow T$	—
$T \rightarrow 0$	{printf("0")}
$T \rightarrow 1$	{printf("1")}
:	:
$T \rightarrow 9$	{printf("9")}

input string $\Rightarrow 9 - 5 + 2$



Denotational semantics

- Each construct of the language is mapped into a mathematical object that defines its meaning.

Example

$$3 * 5 + 4 \Rightarrow \begin{array}{c} + \\ / \quad \backslash \\ * \quad 4 \\ / \quad \backslash \\ 3 \quad 5 \end{array}$$

- Example; also the ordered pair can be used to represent the program
 Eg. $\text{fact}(n) = \text{if } (n == 0) \text{ then } 1$
 $\qquad\qquad\qquad \text{else } n * \text{fact}(n-1)$
- can be viewed as in ordered pair for denotational semantics
 as $\langle n, n! \rangle$
- $\langle 0, 1 \rangle, \langle 1, 1 \rangle, \langle 2, 2 \rangle, \langle 3, 6 \rangle, \langle 4, 24 \rangle$.

Example: Communicating Sequential Process

- A process P can be represented by the notation
 $n \rightarrow Q$, where n is an event and Q is a process
- i.e P as process which first engage in that event n and behaves exactly as described by Q
- Example: a vending machine
 1. $\text{VM}_0 \stackrel{\text{def}}{=} (\text{coin} \rightarrow (\text{candy} \rightarrow \text{stop}))$
 2. $\text{VM}_1 \stackrel{\text{def}}{=} (\text{coin} \rightarrow (\text{candy} \rightarrow \text{VM}_1))$
 3. $\text{VM}_2 \stackrel{\text{def}}{=} (\text{coin} \rightarrow (\text{candy} \rightarrow \text{VM}_2)) // (\text{note bill} \rightarrow (\text{toffee} \rightarrow \text{VM}_2))$

Unit 1.3 Process Algebra

Transition system

- Formally, it can be defined as

$A' = (A, S, \rightarrow, s_0)$, where $A \rightarrow$ set of alphabet

$S \rightarrow$ set of states

$\rightarrow \rightarrow$ transition $S \times A \rightarrow S$

$s_0 \rightarrow$ starting state

Conventions

- a, b, c, \dots are used to indicate elements of A
- u, v, \dots are used to indicate $A \cup \{\tau\}^*$
- w is used to indicate $(A \cup \{\tau\})^*$
- ow is used to indicate A^*

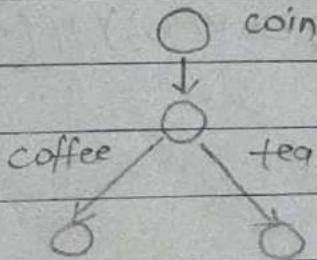
Example:

$s \xrightarrow{ow} s'$ means $s = s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} s_3 \xrightarrow{a_3} \dots \xrightarrow{a_n} s'$

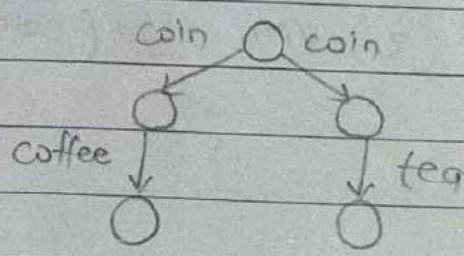
Reachable states

- Let $A' = \{A, S, \rightarrow, s_0\}$ be a transition system
then $s' \in S$ is a reachable state of A'
if there exists $w = u_1, u_2, \dots, u_n \in (A \cup \{\tau\})^*$
such that $s_0 \xrightarrow{w} s'$
- set of all reachable states of A' is denoted by $\text{reach}(A')$
- transition system can also be represented through graphs.

Example:



fig(1) Non-deterministic

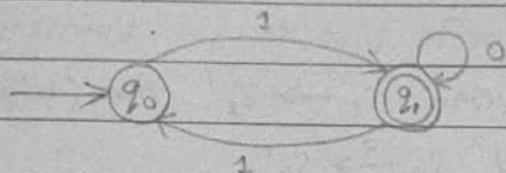


fig(2) Deterministic

Process Algebra

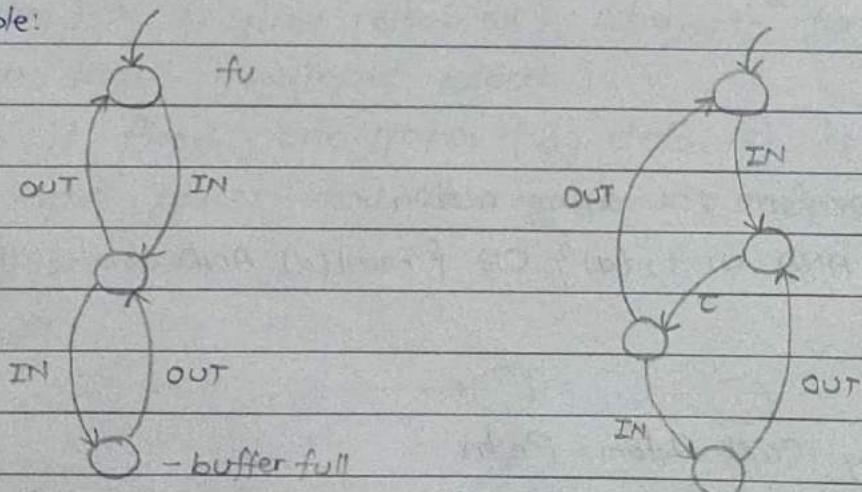
- System behavior generally consists of process and data
- Process are the control mechanism for the manipulation of data.

Eg. DFA



- is a labelled transition system in which one state is considered as root state
- for the mathematical reasoning, process graphs are expressed algebraically called Process Algebra.
- consists of two types of actions
 1. Observable action \rightarrow explicit
 2. Unobservable action \rightarrow implicit and denoted by τ .

Example:



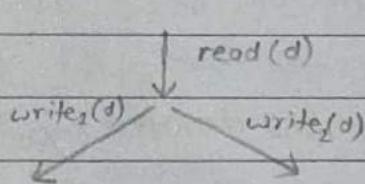
Fig(i) Transition system
representing two place buffer

Fig(ii) Transition system representing
two one place buffer

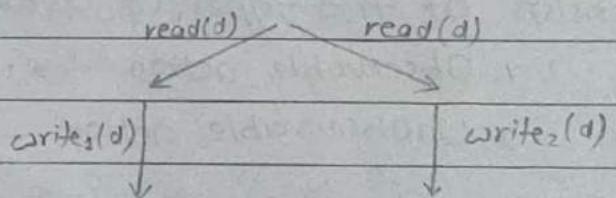
Observation bisimilarity

- Let $A'_i = (A_i, S_i, \rightarrow_i, S_{0i})$, $i=1, 2, \dots$ be two transition system
- A'_1 and A'_2 are observation bisimilar, denoted by $A'_1 \cong A'_2$, if there is a relation $B \subseteq S_1 \times S_2$ such that,
 1. $(S_{01}, S_{02}) \leftarrow B$
 2. If $(S_1, S_2) \leftarrow B$ and $S_1 \xrightarrow{a} S_1'$ then there is S_2' , such that $(S_1', S_2') \leftarrow B$ and $S_2 \xrightarrow{a} S_2'$
 3. If $(S_1, S_2) \leftarrow B$ and $S_2 \xrightarrow{a} S_2'$ then there is S_1' , such that $(S_1', S_2') \leftarrow B$ and $S_1 \xrightarrow{a} S_1'$

Example: Consider the following two processes



Fig(1)



Fig(II)

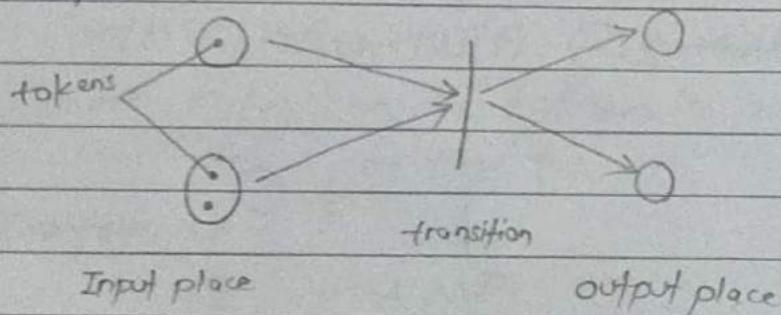
- Both system perform the same action
- i.e {read(d) AND write₁(d)} OR {read(d) AND write₂(d)}

Petri Nets

- Developed by Carl Adam Petri
- is the mathematical modelling language for the description of the system.
- Petri Nets (PN) are represented through directed graph (bi-partite graph)

- consists of two types of nodes
 1. Place (Buffer that holds something (sth))
 - represented by circle
 2. Transition (event)
 - represented by straight line (|) or rectangle (\square) or box (\blacksquare)

Example:



- a transition is fired (enabled), whenever there is at least one token in all its input places.
- When it fired, one token (by default) is removed from all the input places and one token is added to all its output places.

Example:

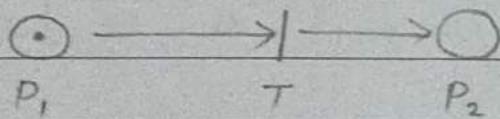


Fig: Before firing T

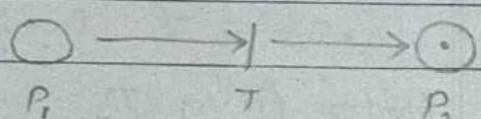


Fig: After firing T

Example:

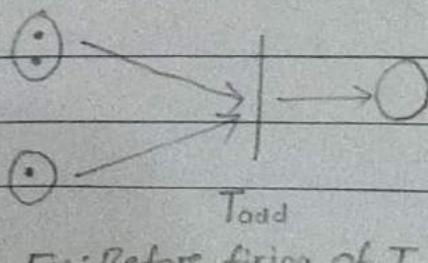


Fig: Before firing of T

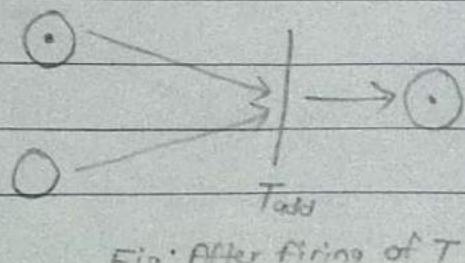
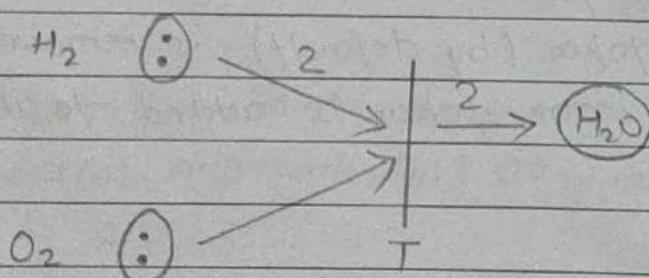
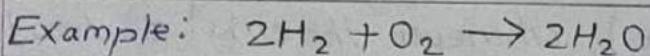
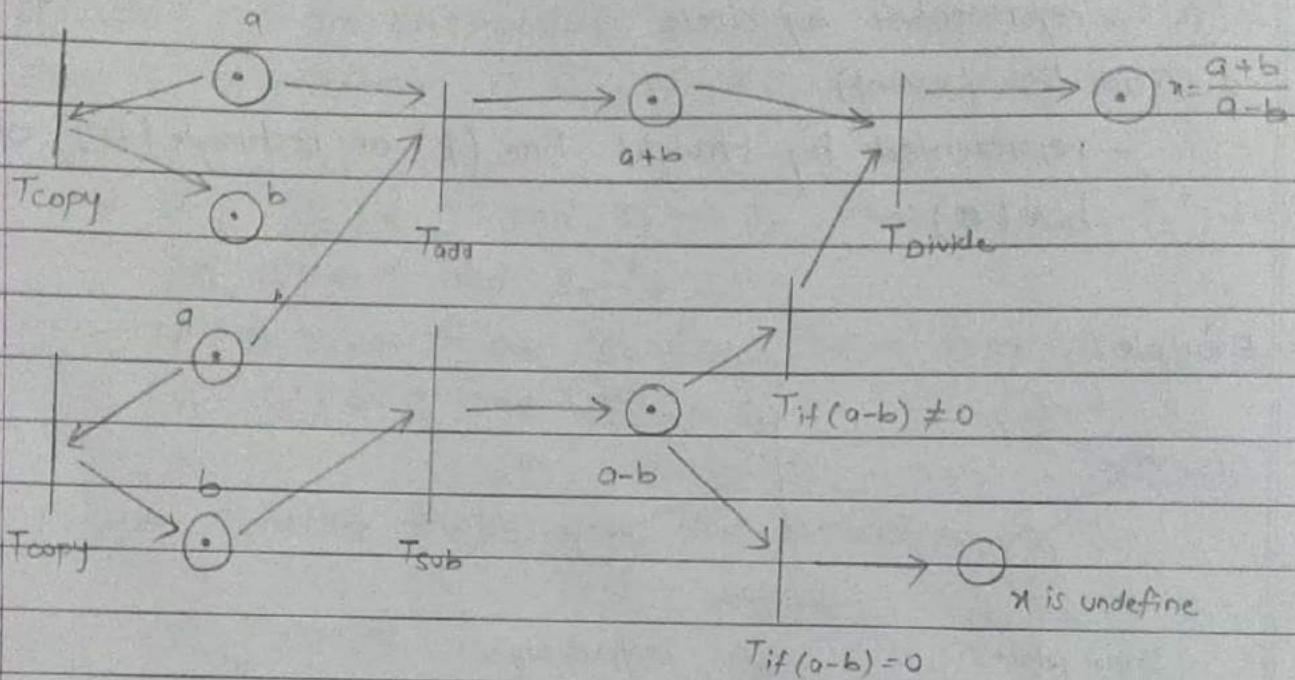
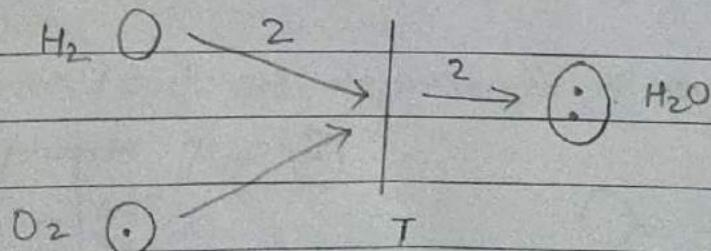


Fig: After firing of T

Example: $x = \frac{a+b}{a-b}$



After firing T ,



Formal Definition of Petri nets

- $PN = (P, T, I, O, M_0)$ where

$P \rightarrow$ set of places = $\{P_1, P_2, \dots, P_m\}$

$T \rightarrow$ set of transitions = $\{T_1, T_2, \dots, T_n\}$

$I \rightarrow$ is an input matrix of $n \times m$

$O \rightarrow$ is an output matrix of $n \times m$

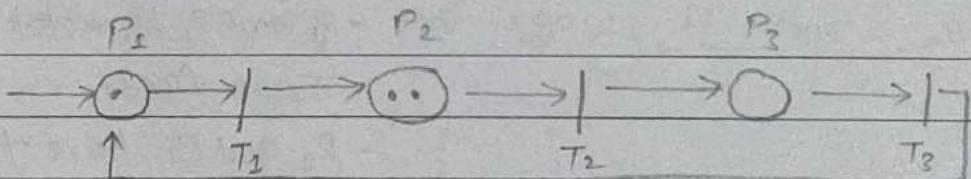
$M_0 \rightarrow$ represents initial marking of ~~matrix~~ machine.

- A marking of a PN at time t ,

$$M(t) = \{M_1(t), M_2(t), \dots, M_m(t)\}$$

where $M_i(t) =$ no. of tokens in place i ,
at time t .

Example:



$$P = \{P_1, P_2, P_3\}$$

$$T = \{T_1, T_2, T_3\}$$

$$I = T_1 \begin{bmatrix} P_1 & P_2 & P_3 \\ 1 & 0 & 0 \end{bmatrix} \\ T_2 \begin{bmatrix} P_1 & P_2 & P_3 \\ 0 & 1 & 0 \end{bmatrix} \\ T_3 \begin{bmatrix} P_1 & P_2 & P_3 \\ 0 & 0 & 1 \end{bmatrix}_{3 \times 3}$$

$$O = T_1 \begin{bmatrix} P_1 & P_2 & P_3 \\ 0 & 1 & 0 \end{bmatrix} \\ T_2 \begin{bmatrix} P_1 & P_2 & P_3 \\ 0 & 0 & 1 \end{bmatrix} \\ T_3 \begin{bmatrix} P_1 & P_2 & P_3 \\ 1 & 0 & 0 \end{bmatrix}_{3 \times 3}$$

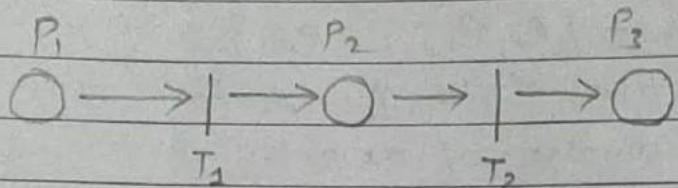
$$M_0 = \begin{bmatrix} P_1 & P_2 & P_3 \\ 1 & 1 & 0 \end{bmatrix}$$

Suppose T_2 is fired

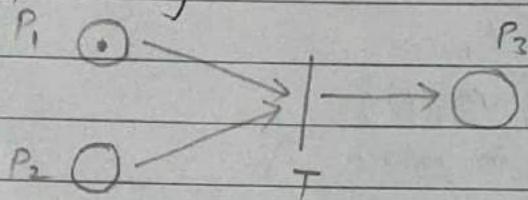
$$M_{T_2} = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$$

Variations in PN

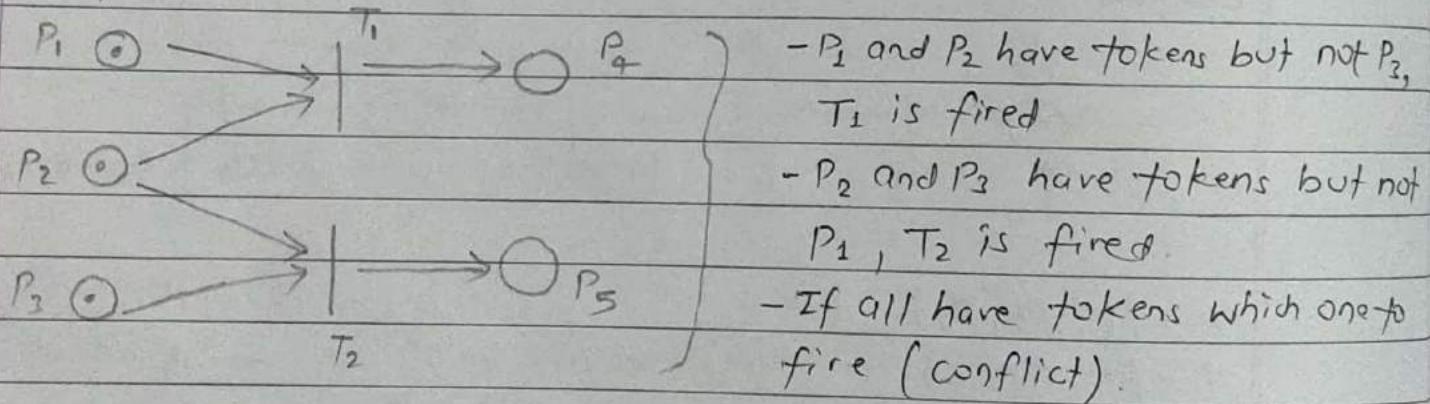
1. Sequential action



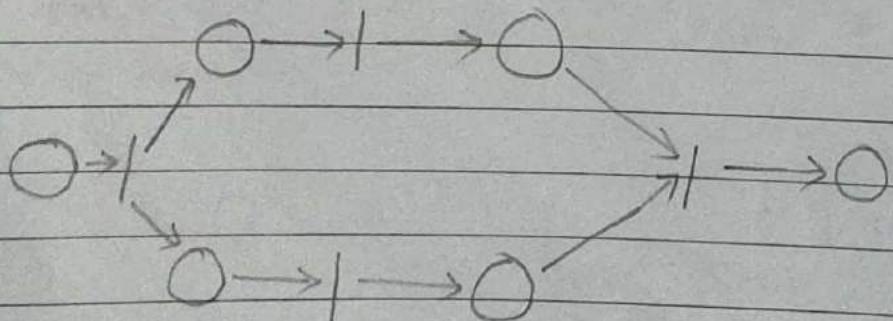
2. Dependency



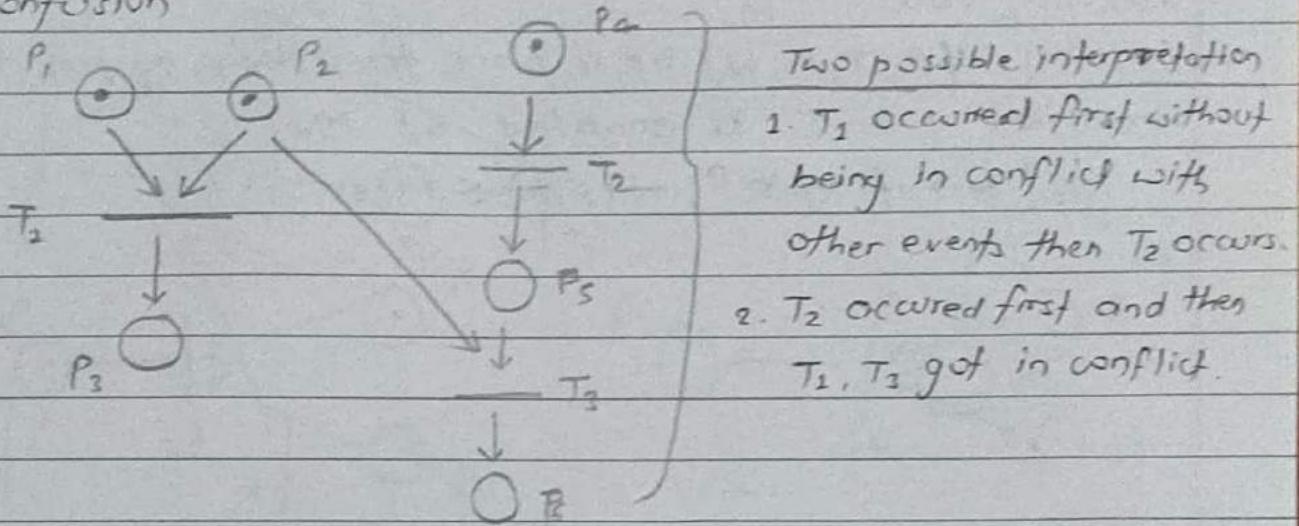
3. Conflict



4. Concurrency



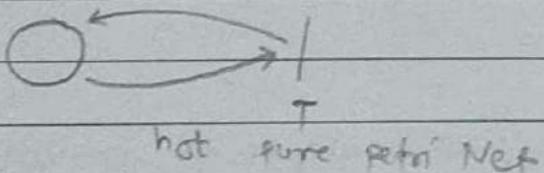
5. Confusion



Pure Petri Net

- Let $PN = \{P, T, F\}$, PN is pure if $\forall x \in X : n \cap n_0 = \emptyset$
- i.e no self loop

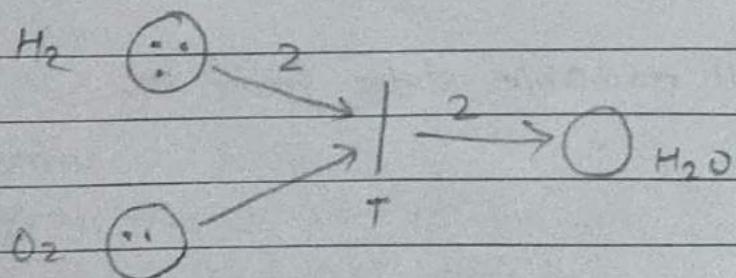
Eg.



Place/Transition Net (P/T net)

- defined by 4 tuples (P, T, F, W) where (P, T, F) is a Petri net and W is a weight function.

Eg.



Transition rule

- Let $M = (P, T, F, W)$ be a place transition system, then a transition $t \in T$ is enabled at M_0 ,

$$M_0 [t > \text{ if } \forall p \in P, W(p, t) \leq M_0(s)]$$

Eg.

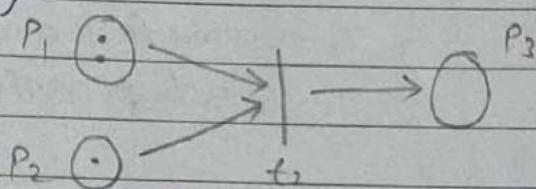


Fig: Before firing t_1

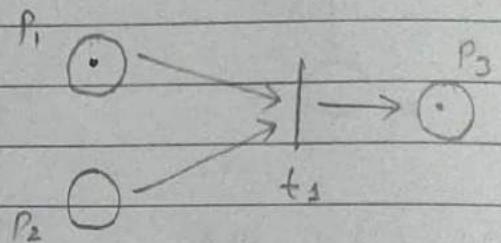


Fig: After firing t_1

- $M_0 [t_1 > M_1, M_0 = (2, 1, 0), M_1 = (1, 0, 1)]$

Properties of Petri Net

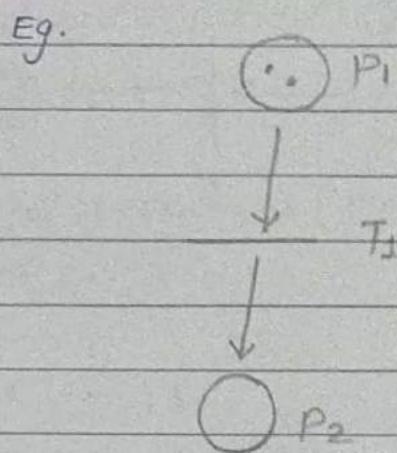
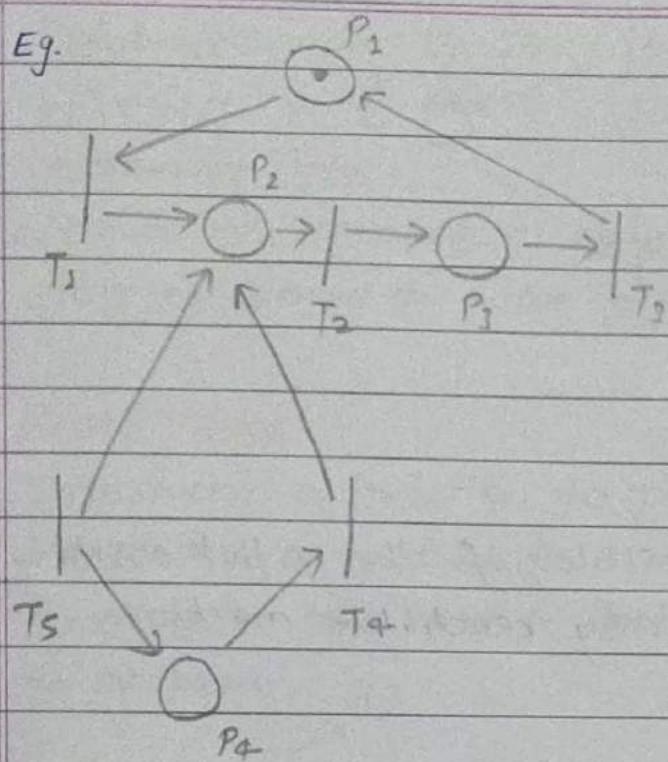
1. Reachability
2. Boundedness
3. Liveness
4. Reversibility

1. Reachability

- A marking M is said reachable from another marking M' , if there exist a sequence, σ of transitions such that

$$M' [\sigma > M]$$

- $R(M_0) \rightarrow$ set of all reachable states from M_0 .



$$M_0 = (2, 0)$$

$$R(M_0) = \{(0, 1, 0, 0), (0, 0, 1, 0)\}$$

$$M_0 = (2, 0)$$

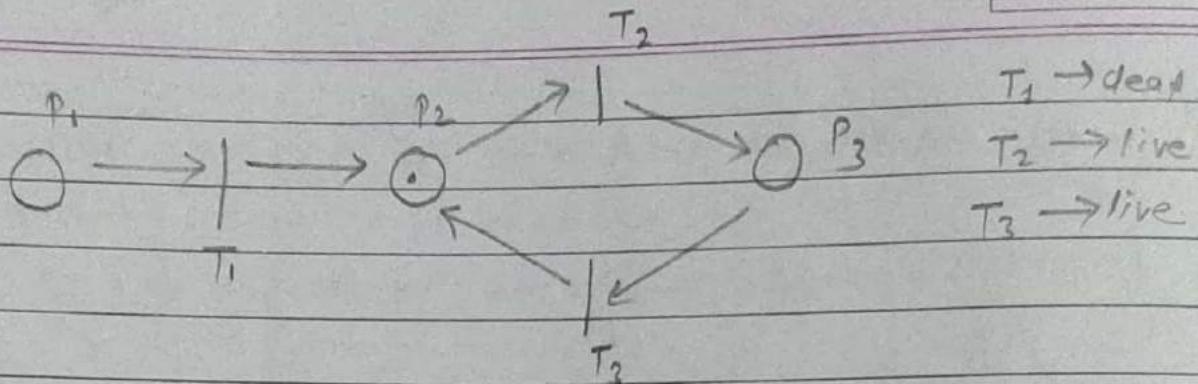
$$R(M_0) = \{(1, 1), (0, 2)\}$$

2. Boundness

- no. of tokens in a place is bounded
- a place is said to be k -bounded, if the ~~the~~ number of tokens in P never exceed k .
- a petri net is k -bounded if all of its places are k -bounded.

3. Liveness

- a transition is said to be live if it can always be made from any reachable marking
- a transition is deadlock if it can be never fire

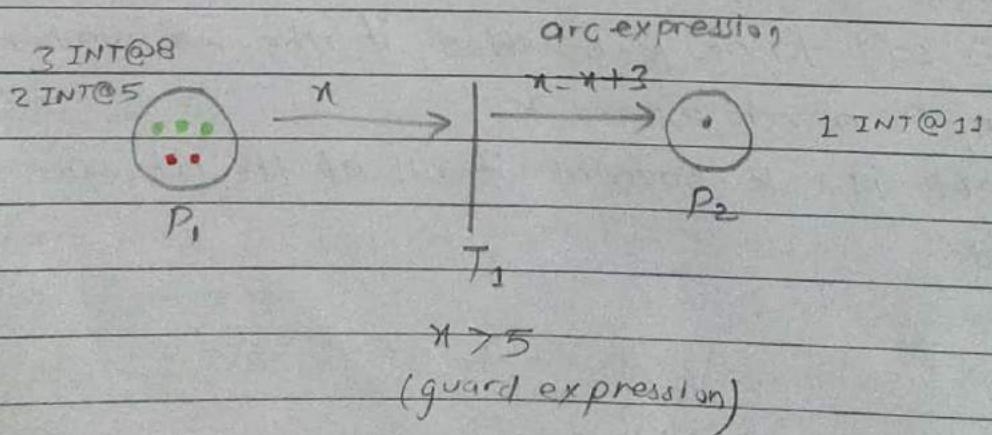


4. Reversibility

- a petri net is said reversible, if the initial marking (M_0) remains reachable from any reachable marking.

High Level Petri Net

- consists of following properties
 1. Petri Net
 2. Marking Scheme (color)
 3. Transition Guard
 4. Arc expression



Assignment-1

Unit 1.4: Complexity Issues in Parallel and Distributed Computing

Complexity Theory

- estimate the amount of computational resources (time and storage) needed to solve a problem.

Basic Terms

1. Introducing a notation to specify complexity
- introduce a mathematical notation to specify the functional relationship between the input and size of the problem.
2. Choice of machine model to standardize the measure
 - specifying the underlying machine
 - RAM, PRAM
3. Refinement of the measure of parallel computation
 - how fast can we solve the problem, when a number of processors are put together to work in parallel.

Turing machine as a basis on consequences

1. Simplicity of abstraction
2. Dual nature of TM
3. Inclusion of non-determinism.
4. Realization of unbounded parallelism
5. Provision of an easy abstract measure for resources

1. Simplicity of abstraction

- simple model of an abstract machine that consists of three components.

a. Finite state machine

$DFA \rightarrow RE$

b. A read/write head

$x \uparrow$
 $CFG \rightarrow PDA$

c. Infinite tape memory

- formally, it can be defined as

$$TM = (Q, \Sigma, \Gamma, \delta, q_0, B, F), \text{ where}$$

Q = set of states

Σ = set of input alphabets

Γ = set of tape symbols

δ = Transition function

q_0 = starting state

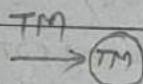
B = Blank symbol

$$S \times \Sigma = S \times \Gamma \times \{S, L, R\}$$

F = set of final states

2. Dual nature of TM

- Universal Turing Machine (UTM)



3. Inclusion of non-determinism

- NDTM \rightarrow alternative move,

- branching tree called OR tree

- a computation halts, if there is a sequence that leads to an accepting state

- rejects an input if and only if there is no possible sequences of moves.

Realization of unbounded parallelism

- NDTM can be simulated using DTM by allowing unbounded parallelism in computation
- DTM that reaches a solution early, halts all other DTM, if it fails, it only terminates.

Provision of easy abstract measure for resources

- DTM \Rightarrow By counting each move as one unit time and each cell as scanned as one space time.
- NDTM \Rightarrow Complexity is the number of moves required for an accepting sequence.

ATM (Alternating Turing Machine)

- States of ATM can be partitioned as
 1. Universal state \rightarrow AND tree
 2. Existential state \rightarrow OR tree
 3. Normal state \rightarrow computation lead to acceptance if and only if its unique descendent is accepting.

Parallel complexity Models

- Synchronous parallel computation can be used to solve problem like matrix multiplication
- It is because all identical processors work under the control of centralized clock.

VLSI computational complexity

- practical use of algorithm depends upon the design of VLSI.
- two fundamental parameters deciding the efficiency of VLSI is
 1. Area of circuit (space complexity)
 2. Time taken to provide an output for a given n input
(Time complexity)

Unit 1.5

Distributed Computing

- Major obstacles in distributed computing is uncertainty due to differing processor speed and occasional failure of components.
- A system consists of n processors P_1, P_2, \dots, P_n and the network net.
- Each processor P_i is modeled as state machine with state sets Q_i .
- The state set Q_i contains a distinguished initial state $q_{0,i}$.
- Each processor P_i contains $buff_i$.

Leader Election

- In distributed computing, leader election is the process of designating a single processor as the leader of some task distributed among several computers.
- The existence of leader can simplify co-ordination among the processors and is helpful in achieving fault tolerance and saving resources.
- Election Algorithm
 - 1. Bully Algorithm
 - 2. Ring Algorithm

1. Bully Algorithm

- The biggest guy in town always win.
- When a process notices that the co-ordinator is not

responding to request, if initiates an election.

- Procedure

1. Each process has a unique processor ID
2. Processes know the IDs and address of every other processes
3. Communication is assumed reliable
4. Process initiates election if the co-ordinator failed
5. Select the processor with highest ID.

Example:

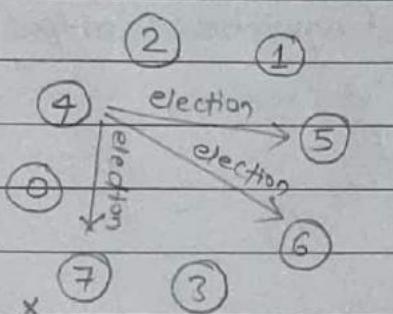


Fig (a)

- currently the leader is 7 and is not responding
- suppose 4 initiates the election

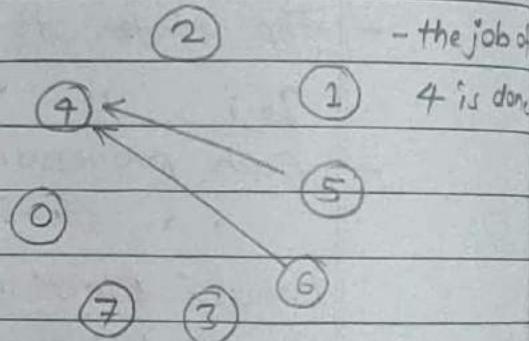


Fig (b)

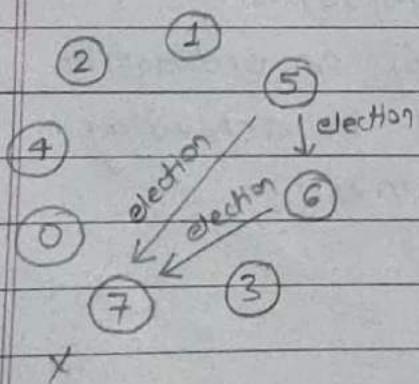


Fig (c)

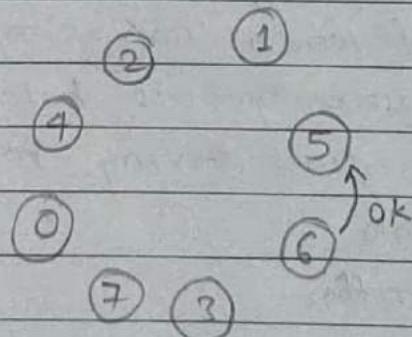


Fig (d)

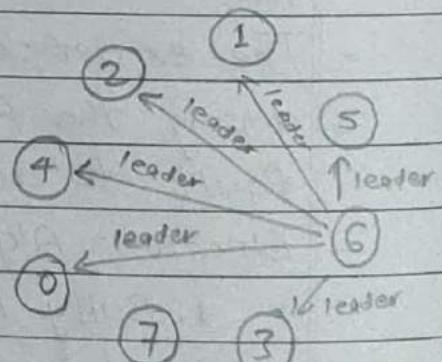


Fig (e)

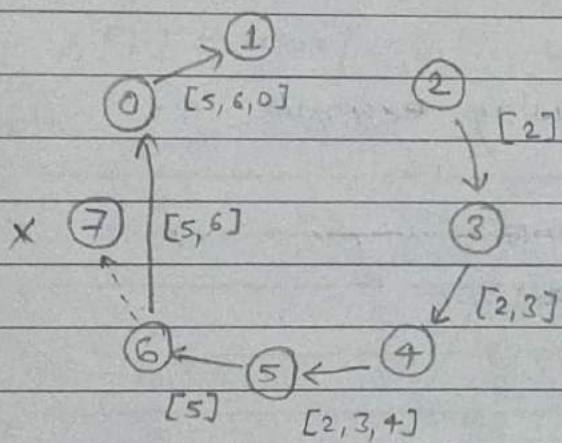
- both 5 & 6 Initiates election

- 6 is winner

- 6 is announcing

Ring Algorithm

- If any process notices that the current co-ordinator has failed, it starts an election by sending message to the first neighbor of the ring.
- The election message contains the node's identifier and is forwarded around the ring.
- When the election message reaches the originator, the election is complete
- Leader is chosen based on the highest numbered process.



Step 1: 2 and 5 are initiating the election

Step 2: they forward the election message to its neighbor until the circuit builds.

Step 3: $2 \rightarrow [2, 3, 4, 5, 6, 0, 1]$

$5 \rightarrow [5, 6, 0, 1, 2, 3, 4]$

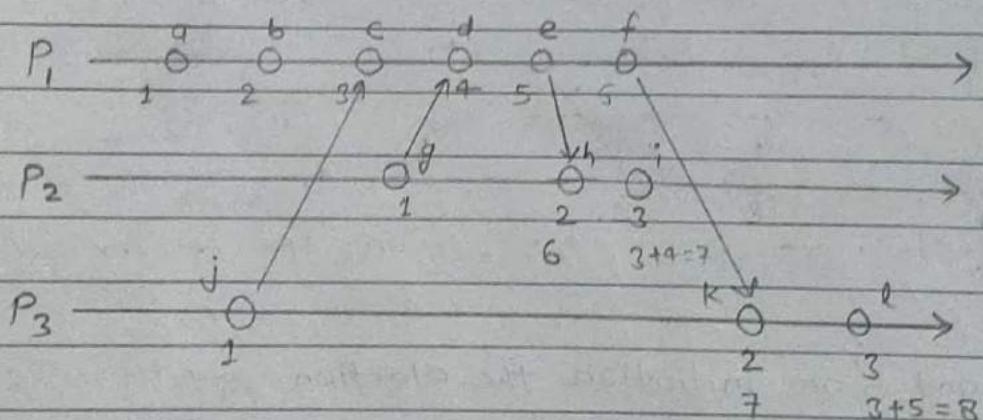
Both 2 and 5 elects 6 as new leader.

Ordering of Events

- Assign logical time stamps to events based on ordering among them
- if event a happened before event b , then it can be represented as
 $a \rightarrow b$ OR $a \leq_h b$
- if $a \rightarrow b$ and $b \rightarrow c$ then $a \rightarrow c$
- if $a \leq_h b$ then $T(a) < T(b)$ where T is time stamp
- but if $T(a) < T(b)$, can we guarantee $a \leq_h b$.

Lamport's Logical clock

- Let us assume the following example.



- i) $(j, c) \rightarrow (1, 3) \rightarrow \text{do nothing}$
 - ii) $(g, d) \rightarrow (1, 4) \rightarrow \text{do nothing}$
 - iii) $(e, h) \rightarrow (5, 2)$
 - iv) $(f, k) \rightarrow (6, 2)$
- } good ordering
} bad ordering

- Each message carries a time stamp of the sender's clock
- When a message arrives
 1. If receiver's clock $<$ message-time-stamp
set system clock to message-time-stamp + 1
 2. else do nothing

Vector Clock

Rules

1. Vector initialization to 0 at each process

$$v_i[j] = 0, \text{ for } i, j = 1, 2, \dots, N$$

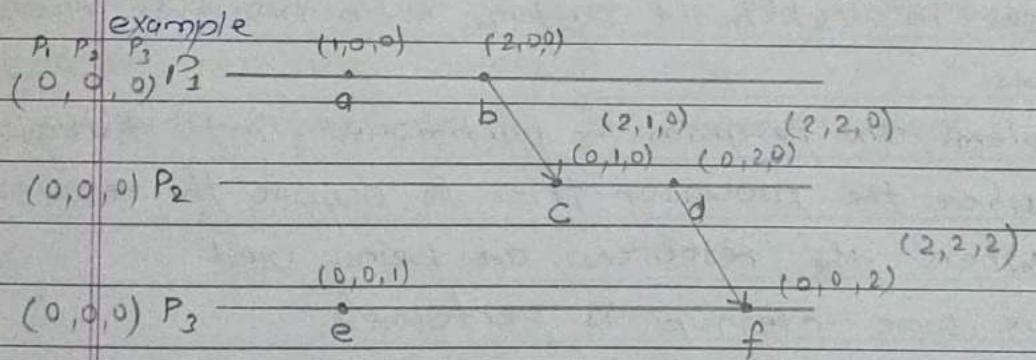
2. Process increments its elements of the local vector before time stamping event.

$$v_i[i] = v_i[i] + 1$$

3. Message is sent from process P_i with v_i attached to it.

4. When P_j receives message, compare vector elements and set local vector to higher of two values.

$$v_j[i] = \max(v_i[i], v_j[i])$$



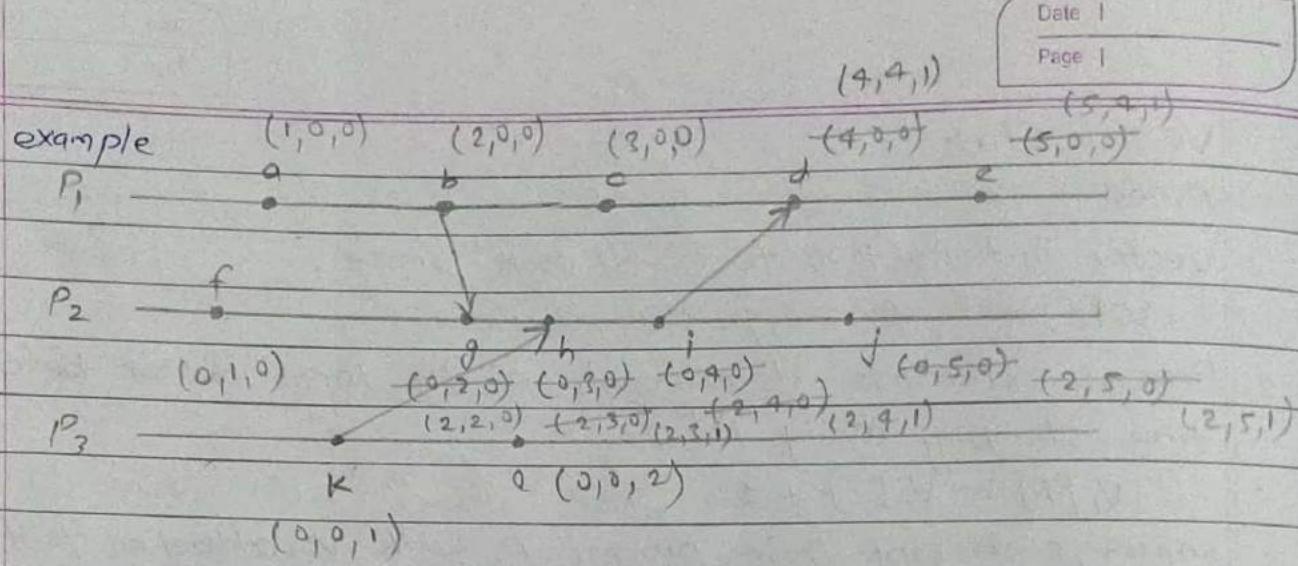
Concurrent events

- Two events e and e' are concurrent iff neither $v(e) < v(e')$ is true nor $v(e') < v(e)$ is true.

Example

Events	Timestamp
a	(1, 0, 0)
b	(2, 0, 0)
c	(2, 1, 0)
d	(2, 2, 0)
e	(0, 0, 1)
f	(2, 2, 2)

- which are concurrent events.
(a,e), (b,e)



Resource Allocation

- Problem arises in distributed systems when processor needs to share resources
- In such problem, the program is partitioned into 4 regions.
 1. Trying : where the processor tries to acquire the resources
 2. Critical : where the resources are being used.
 3. Exit : where some cleanup is performed.
 4. Remainder : the rest of the code.

Tolerating processor failure in synchronous system

- In real systems, the various components may not operate correctly all the time.
- It is assumed that communication links are reliable but processors are not reliable.
- Two types of failure
 1. Crash Failure
 2. Byzantine Failure

Crash Failure (Technical problem)

- A processor crashes in round r , it is correctly operates in round $1, 2, \dots, r-1$, and does not take a step in round $r+1$.

Byzantine Failure (logical problem)

- Faulty processor (traitor) is called byzantine.
- It can send different message to different processors, that is not supposed to send by it.

Assignment 2:

- Tolerating processor failure in Asynchronous system
- sparse network covers
- wait free implementation of shared object. - eq busy wait

Unit 2: Models

Unit 2.1. PRAM Models

PRAM Models

- Parallel Random Access Machine models
- Parallel can be defined as a technique for increasing the computation speed by allocating multiple processors to execute sub tasks simultaneously.

e.g.

Finding the smallest value in the set of N values.

Algorithm

- no. of processors needed (P) = $\frac{n(n-1)}{2}$

Do in Parallel

- P_i derives the pair (i_1, i_2) of indices that correspond to different i .
- P_i reads values $L(i_1) \neq L(i_2)$
- If $L(i_1) > L(i_2)$ then P_i sends negative outcome to P_{i_1} else to P_{i_2}
- Only one processor, P_j , $1 \leq j \leq n$, did not receive negative outcome
- P_j reads the value $L(j)$ and write it onto output cell

1 2 3 4 → index

eg. $L = \{2, 6, 4, 8\}$, $n = 4$

$$P = \frac{n(n-1)}{2} = \frac{4 \times 3}{2} = 6$$

- $\rightarrow P_1 \rightarrow (1, 2) \quad P_2(\text{-ve}) \quad (i_1, i_2)$
 $\rightarrow P_2 \rightarrow (1, 3) \quad P_3(\text{-ve}) \quad \downarrow P_1$
 $\rightarrow P_3 \rightarrow (2, 3) \quad P_2(\text{-ve}) \quad P_1 = (i_2 - 1)(i_2 - 2) + i_1$
 $\rightarrow P_4 \rightarrow (1, 4) \quad P_4(\text{-ve}) \quad 2$
 $\rightarrow P_5 \rightarrow (2, 4) \quad P_4(\text{-ve})$
 $\rightarrow P_6 \rightarrow (3, 4) \quad P_4(\text{-ve})$

$$P_j = P_1 = 2$$

Basic PRAM algorithms

1. Computing prefix sums
2. List Ranking
3. Parallel sorting algorithm

Computing Prefix sums

- given $A = \{a_0, a_1, a_2, \dots, a_{n-1}\}$
- we have to compute $\{a_0, (a_0 + a_1), (a_0 + a_1 + a_2), \dots, (a_0 + a_1 + \dots + a_{n-1})\}$

e.g.

$$A = \{5, 3, -6, 2, 7, 10, -2, 8\}$$

- the output is

$$\{5, 8, 2, 4, 11, 21, 19, 27\}$$

Algorithm

- input an array $\{a_1, a_2, \dots, a_n\}$
- if ($n=1$) then $s_1 \leftarrow a[1]$
- Else
 - For $j = 0$ to $\log n - 1$ do
 - For $i = 2^j + 1$ to n do in parallel
 - 1. obtains $a[i - 2^j]$ from P_{i-2^j} through shared memory
 - 2. $a[i] = a[i - 2^j] + a[i]$

Eg. A $\{5, 3, -6, 2, 7, 10, -2, 8\}$

Here,

$$n = 8, \log 8 = 3 \quad j = 0, 1, 2$$

For $j = 0$

$$i = 2^j + 1 \text{ to } n = 2 \text{ to } 8$$

$$i - 2^j = 2 - i = 1$$

$$P_1 \rightarrow a[1] = 5$$

$$P_2 \rightarrow a[2] = a[1] + a[2] = 5 + 3 = 8$$

$$P_3 \rightarrow a[3] = a[2] + a[3] = 3 + (-6) = -3$$

$$P_4 \rightarrow a[4] = a[3] + a[4] = (-6) + 2 = -4$$

$$P_5 \rightarrow a[5] = a[4] + a[5] = 2 + 7 = 9$$

$$P_6 \rightarrow a[6] = a[5] + a[6] = 7 + 10 = 17$$

$$P_7 \rightarrow a[7] = a[6] + a[7] = 10 + (-2) = 8$$

$$P_8 \rightarrow a[8] = a[7] + a[8] = (-2) + 8 = 6$$

$\{5, 8, -3, -4, 9, 17, 8, 6\}$

or, $j = 1$

$$i = 2^j + 1 \rightarrow n \\ = 3 \rightarrow 8$$

$$i - 2^j = 3 - 2^1 = 1$$

$$P_1 \rightarrow a[1]$$

$$P_2 =$$

$$P_3 \rightarrow a[3] = a[1] + a[3] = 5 + (-3) = 2$$

$$P_4 \rightarrow a[4] = a[2] + a[4] = 8 + (-4) = 4$$

$$P_5 \rightarrow a[5] = a[3] + a[5] = (-3) + 9 = 6$$

$$P_6 \rightarrow a[6] = a[4] + a[6] = (-4) + 17 = 13$$

$$P_7 \rightarrow a[7] = a[5] + a[7] = 9 + 8 = 17$$

$$P_8 \rightarrow a[8] = a[6] + a[8] = 17 + 6 = 23$$

$\{5, 8, 2, 4, 6, 13, 17, 23\}$

For $j = 2$

$$i = 2^j + 1 \rightarrow n \\ = 5 \rightarrow 8$$

$$i - 2^j = 5 - 2^2 = 1$$

$$P_5 \rightarrow a[5] = a[1] + a[5] = 5 + 6 = 11$$

$$P_6 \rightarrow a[6] = a[2] + a[6] = 8 + 13 = 21$$

$$P_7 \rightarrow a[7] = a[3] + a[7] = 2 + 17 = 19$$

$$P_8 \rightarrow a[8] = a[4] + a[8] = 4 + 23 = 27$$

$\{5, 8, 2, 4, 11, 21, 19, 27\}$ output

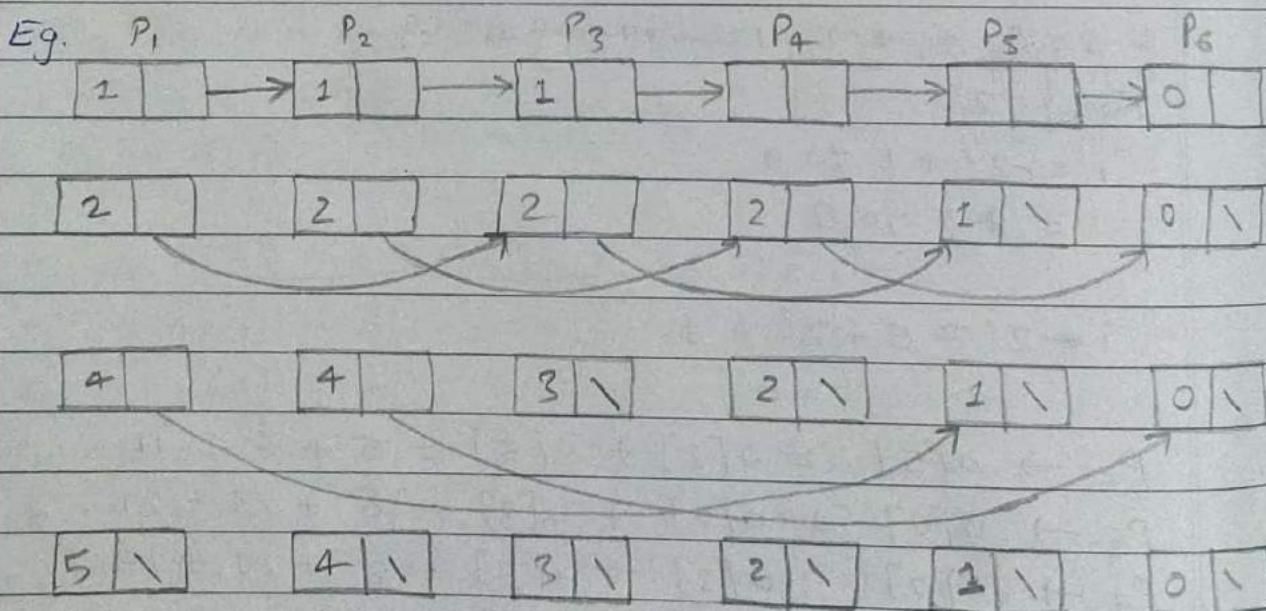
List Ranking

- The problem is that, given a singly linked list L , with N objects, for each node compute the distance to the end of the list.
- if d denotes the distance

$$\text{node}.d = \begin{cases} 0 & \text{if } \text{node}.next = \text{nil} \\ \text{node}.next.d + 1, & \text{otherwise} \end{cases}$$

Parallel Algorithm

- assign one processor for each node
- for each node i , do in parallel
 1. $i.d = i.d + i.next.d$
 2. $i.next = i.next.next$



Parallel Sorting Algorithm

- bitonic sort
- consists of two sequences, one increasing and one decreasing

5, 6, 7, 8, 4, 3, 2, 1

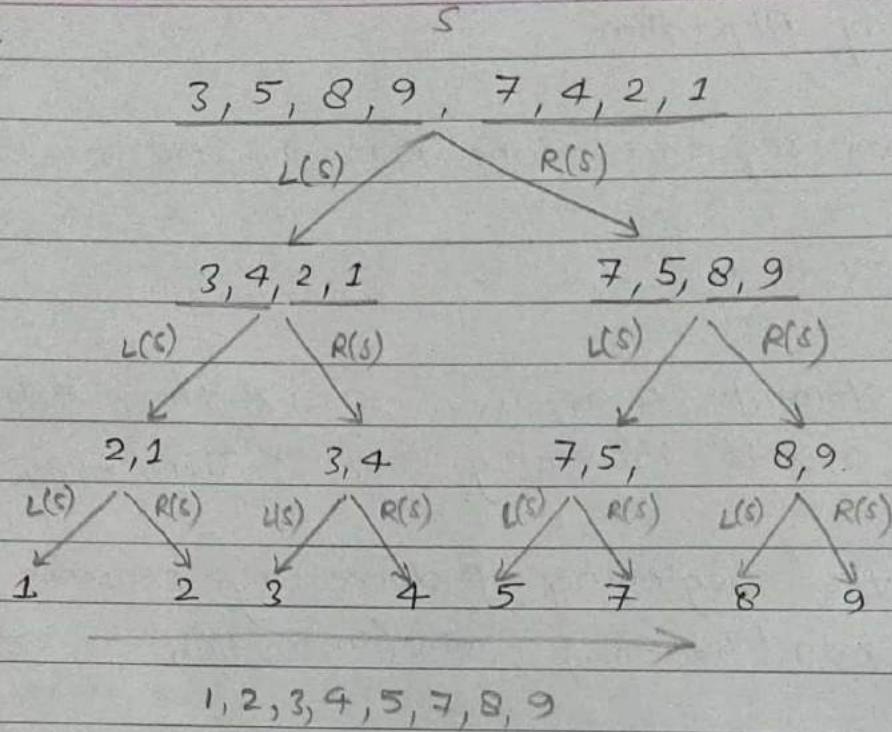
- sequence of elements $\langle a_0, a_1, \dots, a_{n-1} \rangle$ there exists $\langle a_0, a_1, \dots, a_i \rangle$ is increasing and $\langle a_{i+1}, a_{i+2}, \dots, a_{n-1} \rangle$ is decreasing.
- if $a_{n/2}$ is the beginning of decreasing sequence s, then
 $L(s) = \{ \min(a_0, a_{n/2}), \min(a_1, a_{n/2+1}), \dots, \min(a_{n/2-1}, a_{n-1}) \}$

$$R(s) = \{ \max(a_0, a_{n/2}), \max(a_1, a_{n/2+1}), \dots, \max(a_{n/2-1}, a_{n-1}) \}$$

Algorithm

- Input \rightarrow a bitonic sequences
- If s is the length 1 then STOP
- Else
 - \rightarrow Form $L(s)$ and $R(s)$
 - \rightarrow Do in parallel
 - $L(s) \leftarrow$ recursive bitonic merge $L(s)$
 - $R(s) \leftarrow$ recursive bitonic merge $R(s)$
 - concatenate $L(s)$ and $R(s)$

Eg.



Randomized Algorithm

- is the algorithm where execution is controlled at one or more points by randomly made choices.

e.g. Fermat's little theorem

- Primality testing $\rightarrow 2^7 \equiv 2 \pmod{7}$

$$b^p \equiv b \pmod{p} \quad p = \text{prime}$$

where $b \in \mathbb{Z}_p$ $\mathbb{Z}_n = \text{set of residue}$

Deficiencies of PRAM algorithm

- no mechanisms for representing communication between the processors.
- storage management and communication issues are hidden from the algorithm designer.

Unit : 2.2

Broadcasting with Selective Reduction (BSR)

- Broadcast instruction consists of three phases
 1. Broadcast Phase
 2. Selection Phase
 3. Reduction Phase

Broadcast Phase

- allows all of the N processors to write concurrently to all of the M memory locations.
- each processor P_i produced a record containing two fields, a tag t_i and datum d_i .

Selection Phase

- after the data are received and each memory locations U_j , $1 \leq j \leq m$, a switch S_j will select the receiving data d_i by comparing tag value t_i by using a selection rule.

$$\sigma \{ <, \leq, =, >, \geq, \neq \}$$

Reduction Phase

- selected data are reduces to a single value using reduction rule, $R \{ \sum, \pi, \wedge, \vee, \oplus, \cap, \cup \}$

Formal Definition of RSR.

- Mathematically it can be expressed as

$$v_j \leftarrow R(d_i) \mid \wedge t(i, h) \sigma_h L(j, h)$$

$1 \leq i \leq N \quad 1 \leq h \leq k \quad 1 \leq j \leq M$

where,

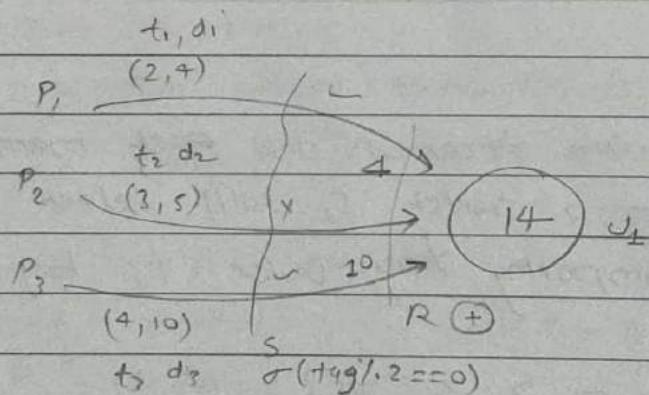
$N \Rightarrow$ no. of processor

$d_i \rightarrow$ datum broadcast by P_i

$\sigma_h \rightarrow$ selection operation

$t(i, h) \rightarrow$ tag broadcast by P_i for criteria σ_h

$L(j, h) \rightarrow$ limit value j for criteria σ_h .



Parenthesis Matching

- One criterion RSR algorithm.
- The problem is to find the pairs of matching parenthesis in a given legal sequences l_1, l_2, \dots, l_n of parentheses.
- By legal means that every parenthesis has its matching parenthesis in a sequence
- Example:

output for the input sequence,

$(((())))(())$ is

()	((())	(())	
1	2	3	4	5	6	7	8	9	10	11	12
2	1	8	7	6	5	4	3	12	11	10	9

- can be solved by N processors

Algorithm

- For each processor j do in parallel
- If $I_j == "("$ then $b_j = 1$ else $b_j = -1$
- $P_j = \sum b_j \mid i \leq j$
- If $b_j == -1$ then $P_j = P_j + 1$
- $P'_j = P_j - \frac{1}{j}$
- $q_j = -1, t_j = 0, r_j = 0$
- $q_j = \cap P'_i \mid P'_i < P'_j$
- $t_j = \cap i \mid P'_j = q_j$
- $r_j = \cup i \mid t_j = i$
- If $I_j == ")"$ then $m_j = t_j$ else $m_j = r_j$

	P_1	P_2	P_3	P_4	P_5	P_6	P_7	P_8	P_9	P_{10}	P_{11}	P_{12}	
Sequence	()	((()))	(())	
Index(j)	1	2	3	4	5	6	7	8	9	10	11	12	
b_j	1	-1	1	1	1	-1	-1	-1	1	1	-1	-1	
$\sum b_i \mid i \leq j$	1	0	1	2	3	2	1	0	1	2	1	0	
Updated P'_j	1	1	1	2	3	3	2	1	1	2	2	1	
$1 - \frac{1}{1}, 1 - \frac{1}{2}$	P'_j	0	0.5	0.66	1.75	2.8	2.83	1.86	0.87	0.89	1.9	1.91	0.92
default value	q_j	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	
	q'_j	-1	0	0.5	0.92	1.91	2.8	1.75	0.66	0.87	1.86	1.9	0.89
	t_j	0	1	2	12	11	5	4	3	8	7	10	9
	r_j	2	3	8	7	6	0	10	9	12	11	5	4
	m_j	2	1	8	7	6	5	4	3	12	11	10	9

Maximal Sub Sum Set

Maximal Sum Sub Segment

- Given an array of numbers d_1, d_2, \dots, d_m , it is required to find a contiguous sub array of maximal sum.

Example:

$$A = \{-2, -3, 4, -1, -2, 1, 5, -3\}$$

$$\text{Maximal Sum} = 4 + (-1) + (-2) + 1 + 5 = 7$$

Algorithm

$$s_j = \sum d_i \mid i \leq j$$

$$m_j = \max s_i \mid i \leq j$$

$$e_j = \min i \mid s_i = m_j$$

$$m_j = m_j - s_j + d_j$$

$$t = n m_i$$

$n = n i \mid m_i = t \rightarrow$ starting point

$y = e_x \rightarrow$ ending point

$$A = \{31, -41, 59, 26, -53, 58, 97, -93, -23, 84\}$$

s_j	31	-10	49	75	22	80	177	84	61	145
m_j	177	177	177	177	177	177	177	145	145	145
e_j	7	7	7	7	7	7	7	10	10	10
m_i	177	146	187	128	102	155	97	-32	61	84

$$t = n m_i = 187$$

$n = 3 \rightarrow$ starting point

$y = 7 \rightarrow$ ending point

$$\therefore 59 + 26 - 53 + 58 + 97 = 187$$

Two criterion RSR algorithm

- RSR algorithm where broadcast data are tested for their satisfaction of two criteria before being allowed to take part in reduction process.

Example:

- Counting inversions in a permutation
- Given a set S , a permutation π of S is a set S' , containing all elements of S , but in different order.
- Inversions in permutations are those numbers of pairs, which are in disorder.

Eg.

$\begin{matrix} 1 & 2 & 3 & 4 & 5 & -i \end{matrix}$
NO. of inversions in $\{1, 6, 2, 9, 5\} = 3$

They are $\{6, 2\}, \{6, 5\}, \{9, 5\}$

Rules

$$\bullet y_j = \sum_{i=1}^j \{ i < j \wedge \pi(i) > \pi(j) \}$$

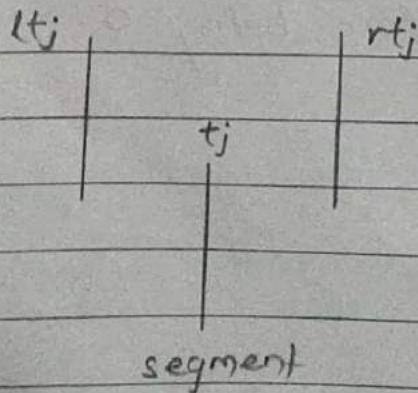
$$\begin{array}{cccc} (1, 6) & (6, 2) \cancel{\sim} & (2, 9) & (9, 5) \cancel{\sim} \\ (1, 2) & (6, 9) & (2, 5) & \\ (1, 9) & (6, 5) \cancel{\sim} & & \\ (1, 5) & 2 < 5 \wedge 6 > 5 \text{ true} & & \end{array}$$

Three criterion BSR algorithm

- Have to pass three criterion to go to reduction process.

Eg:

vertical segment visibility



$$\{t_j = \gamma_i \mid x_i < x_j \wedge t_i \geq t_j \wedge b_i \leq t_j\} \text{ (left top)}$$
$$\{t_j = \gamma_i \mid x_i > x_j \wedge t_i \geq t_j \wedge b_i \leq t_j\} \text{ (right top)}$$

Unit : 2.3

Data Flow Models

- Flow of information between data processing entities
- Control flow and data flow
- Control flow assumes that a program in a series of instructions, each of which specifies
 - Either an operation with memory location
 - OR specifies transfer of control to another instruction

Example:

$$C = \text{if } (n == 0) \text{ then } a + b \text{ else } a - b$$

- But if in data flow computing model, it executes any operation as soon as necessary operations are available

Basic primitives of data flow

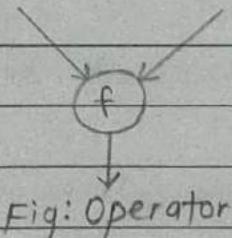


Fig: Operator

→ a data value is produced by an operator as a result of some operation f

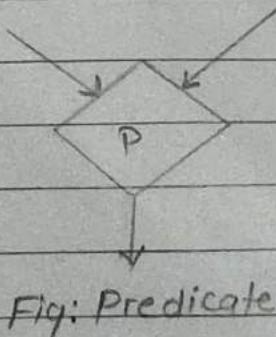
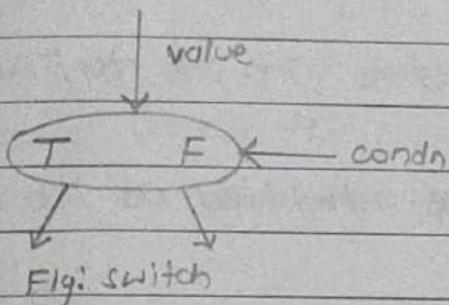
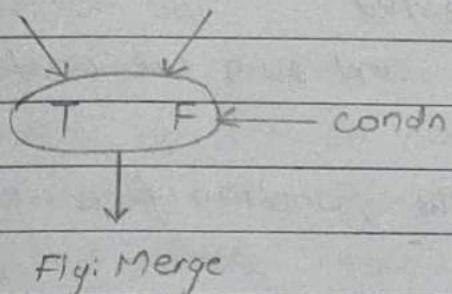


Fig: Predicate

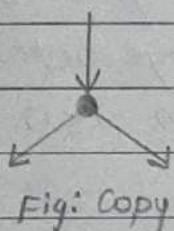
→ a TRUE or FALSE control value is generated by a decider (predicate) depending on its input tokens.



→ a switch actor directs an input data token to one of its output depending upon the control input

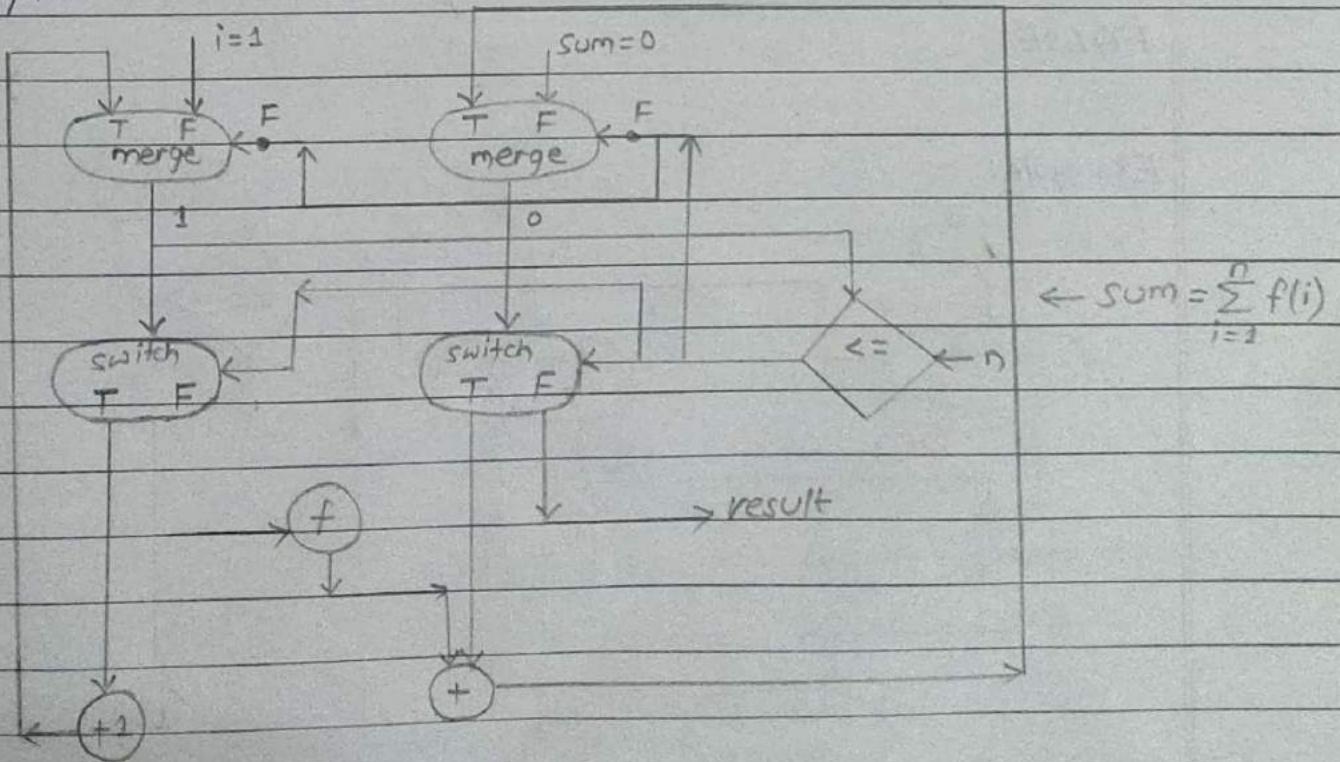


→ A merge actor passes one of its input tokens to the output depending on the input control token



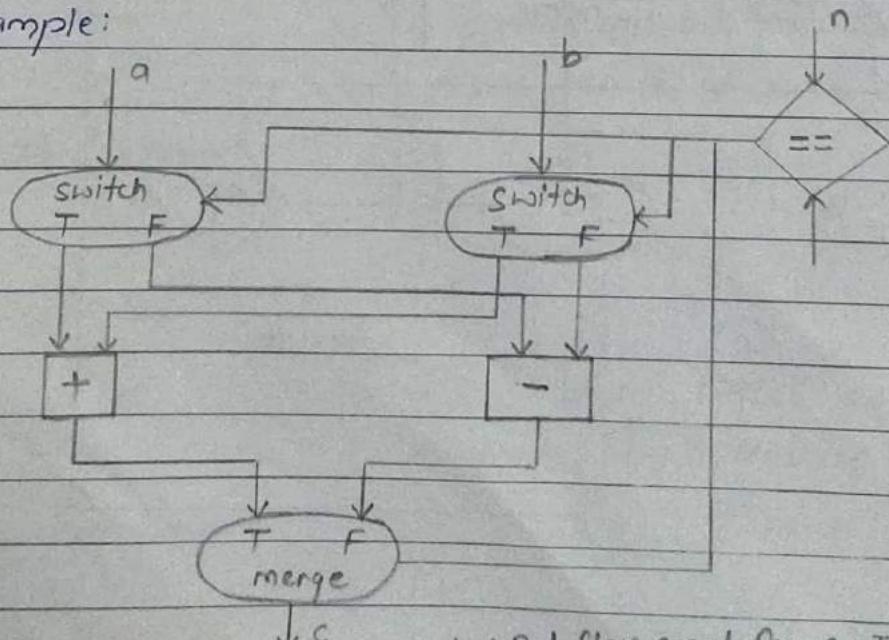
→ A copy is an identity operator which duplicates input tokens.

Example:



Description:

- Control units to the merge actors are set to FALSE for initialization.
- The input values i and sum are admitted as the initial values of the iteration.
- The predicate $i \leq n$ is then tested
- If it is true, the values of i and sum are routed to the TRUE sides of switch.
- This initiates the firing of the function f as well as increment of i .
- Once the execution of the body of the loop completes, the next iteration initiated at the merge actors.
- The initiation continues until the condition $i \leq n$ is no longer true.
- The final sum is then routed out of the loop and the initial Boolean values at merge actors are restored to FALSE.

Example:Fig: Dataflow graph for $C = \text{IF } n == 0 \text{ THEN } a+b \text{ ELSE } a-b$

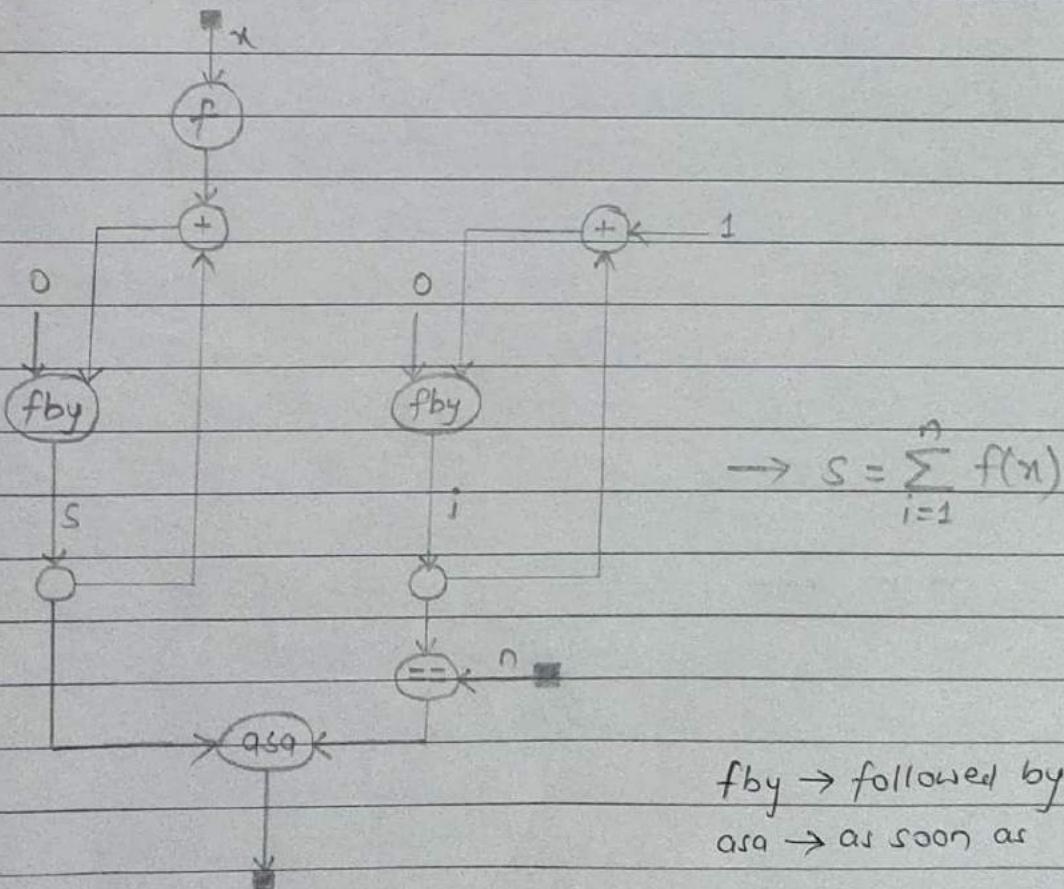
Demand Driven Data Flow Computing Model

- An operation of a node will be performed only if there are tokens on all the input edges and there is a demand for the result of applying the operation

Operator net

- Programming model for demand driven
- consists of set of nodes, a set of arcs and set of equations that relate the output arcs of nodes to functions/operators applied to the input arcs.

Example:

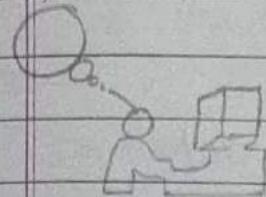


Unit 2.4: Partitioning and scheduling

Partitioning and scheduling

- Partitioning deals with parallelism

Explicit Approach



Implicit Approach

Sequential Program

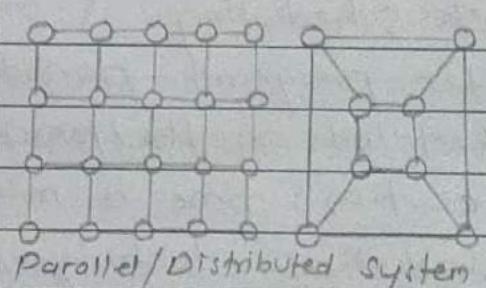
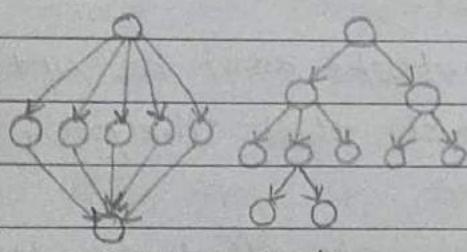
Dependence Analyzer

Parallel
Program
Tasks

Gains of
sequential code

Partitioner

Ideal
Parallelism



Scheduler

Processors

Time

schedule

Program Partitioning

1. Data Partitioning

- same computation on different set of data concurrently
- Eg. matrix multiplication

2. Function Partitioning

- Different computation on same set of data concurrently
- Eg. flight simulation

Task scheduling

- After program partitioning, tasks must be optimally scheduled on the processors such that the program execution time is minimized.
- Scheduling techniques can be classified as deterministic and non-deterministic
- In deterministic scheduling, all the information about tasks to be scheduled is entirely known prior to execution time.
- In non-deterministic, some information may not be known before program execution
(Eg. conditional branches)

Scheduling system Models

- consists of

1. Parallel program tasks
2. Target machine
3. Schedule
4. Performance criteria

1. Parallel program tasks

Parallel program tasks can be defined as the system (T, \leq, D_{ij}, A_i) , where

1. $T = \{t_1, t_2, \dots, t_n\}$ is a set of tasks to be executed
2. \leq is a partial order defined on T , i.e. $t_i \leq t_j$ means t_i must be completed before t_j can start execution,
3. D_{ij} is an $n \times n$ matrix of communication data, where $D_{ij} \geq 0$ is the amount of data required to be transmitted from task t_i to t_j .
4. A_i is a vector of the amount of computations i.e $A_i > 0$, is the number of instructions required to execute t_i .

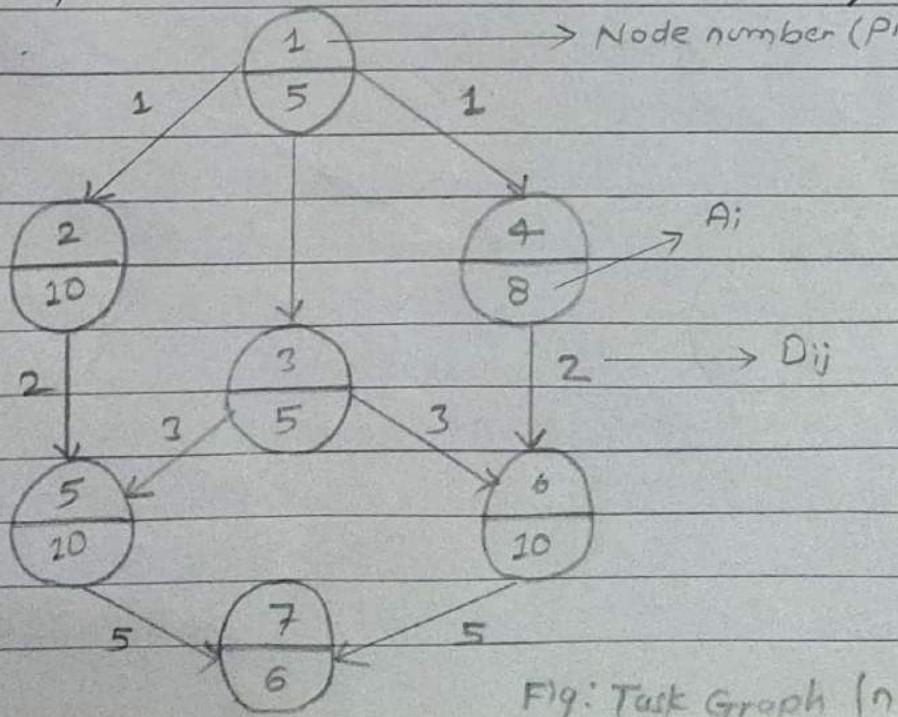


Fig: Task Graph ($n=7$)

1. $T = \{t_1, t_2, t_3, t_4, t_5, t_6, t_7\}$

2. $L = \{t_1 < t_2, t_1 < t_3, t_1 < t_4,$
 $t_2 < t_5, t_3 < t_5, t_6 < t_3,$
 $t_6 < t_4, t_5 < t_7, t_6 < t_7\}$

3.. D_{ij}

	1	2	3	4	5	6	7
1	0	1	2	1	0	0	0
2	0	0	0	0	2	0	0
3	0	0	0	0	3	3	0
4	0	0	0	0	0	2	0
5	0	0	0	0	0	0	5
6	0	0	0	0	0	0	5
7	0	0	0	0	0	0	0

4. $A_1 = 5$

$A_2 = 10$

$A_3 = 5$

$A_4 = 8$

$A_5 = 10$

$A_6 = 10$

$A_7 = 6$

2. Target Machine

- Assumed to be made up of "m" heterogeneous processing elements connected using an arbitrary interconnection network
- Formally target machine can be defined as $(P, P_{ij}, S_i, T_i, \beta_i, R_{ij})$, where
 1. $P = \{P_1, P_2, \dots, P_m\}$ is a set of m processors forming the parallel architecture
 2. P_{ij} is an $m \times m$ interconnection topology matrix
 3. S_i specifies the speed of processor P_i .
 4. T_i specifies the startup cost of initiating message m on processor P_i .
 5. β_i specifies the startup cost of initiating a process on processor P_i .
 6. R_{ij} is the transmission rate over link connecting P_i and P_j .

- Now execution time of task i on processor j (T_{ij}) can be computed as,

$$T_{ij} = \frac{A_i}{S_j} + \beta_j$$

↑ → start(process j)
 task processor speed

Schedule

Formally defined as

$$f: v \rightarrow \{1, 2, \dots, m\} \times \{0, 1, \dots, \infty\}$$

- i.e $f(v) = (i, t)$ for some v means task v is scheduled to be processed by processor P_i starting at time t
- Task 1 starts on processor P_2 at time 0 and finishes at time 10
- Task 5 starts on processor P_3 at time 0 and finishes at time 30
- Shaded area is communication delay.

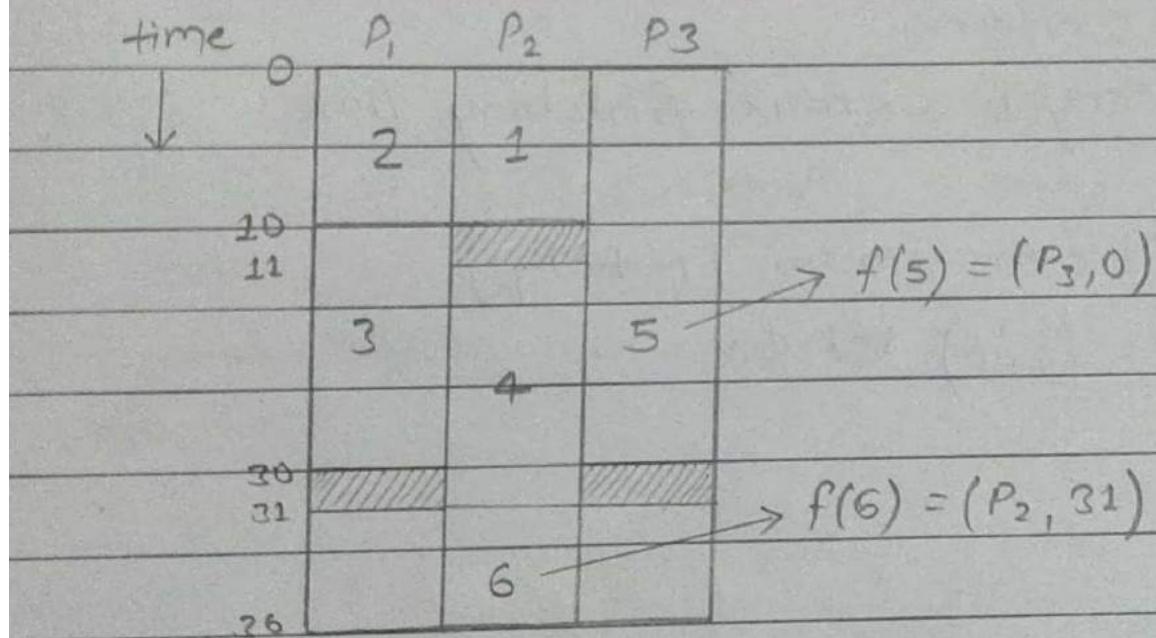


Fig: Gantt chart

Performance Criteria

- Scheduling length or max finishing time
- Formally
 - length(f) = maximum $\{f_i + T_{ij}\}$
 - Where $f(i) = (j, t) \forall i \in V$

Communication Model

- Two key component that contribu to the total completion cost is
 1. Execution Time
 2. Communication Delay

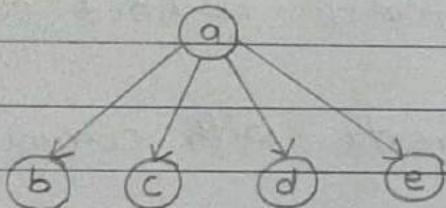
Model A

- Total cost = communication cost + execution cost
- Where
 - Execution cost = schedule length
 - Communication cost = no. of message \times cost per message
 - No. of message = the no. of node pairs (u, v) such that $(u, v) \in A$ and $\text{proc}(u) \neq \text{proc}(v)$

Model B

- If tasks u, v and $w \in V$ such that $(u, v), (u, w) \in A$, then the result of computing u must be communicated to the processor v and w
- This communication is counted twice in model A but one in model B

Example:



- Suppose that each of given task takes one unit of execution time in either P_1 or P_2
- Task a sends message to b, c, d, e
- Suppose $a, b, d \rightarrow$ allocated to P_1
 $c, e \rightarrow$ allocated to P_2

Model A

$$\text{No. of message} = |(a, c), (a, e)| = 2$$

Model B

$$\text{No. of message} = |(a, P_2)| = 1$$

Model C :

- This model assumes that the existence of an I/O processor that is associated with every processor in the system.
- communication during execution period happens.

Optimal Scheduling Algorithm

- Deal in 3 cases

1. When the task graph is a tree
2. When the task graph is an interval order
3. When there are only 2 processors available

Scheduling tree structure task graphs

- A. Scheduling in-forests / out-forests without communication
- B. Scheduling in-forests / out-forests with communication

Scheduling in-forests / out-forests without communication

- Task graph is either

- In-forest \rightarrow each node has at most one immediate successors OR
- Out-forest \rightarrow each node has at most one immediate predecessors

- General strategy \rightarrow used the highest level first

Algorithm

1. The level of each node in the task graph is calculated and used as each node priority
2. Whenever a processor becomes available, assign the unexecuted ready task with the highest priority

Example :

Level 5

Level 4

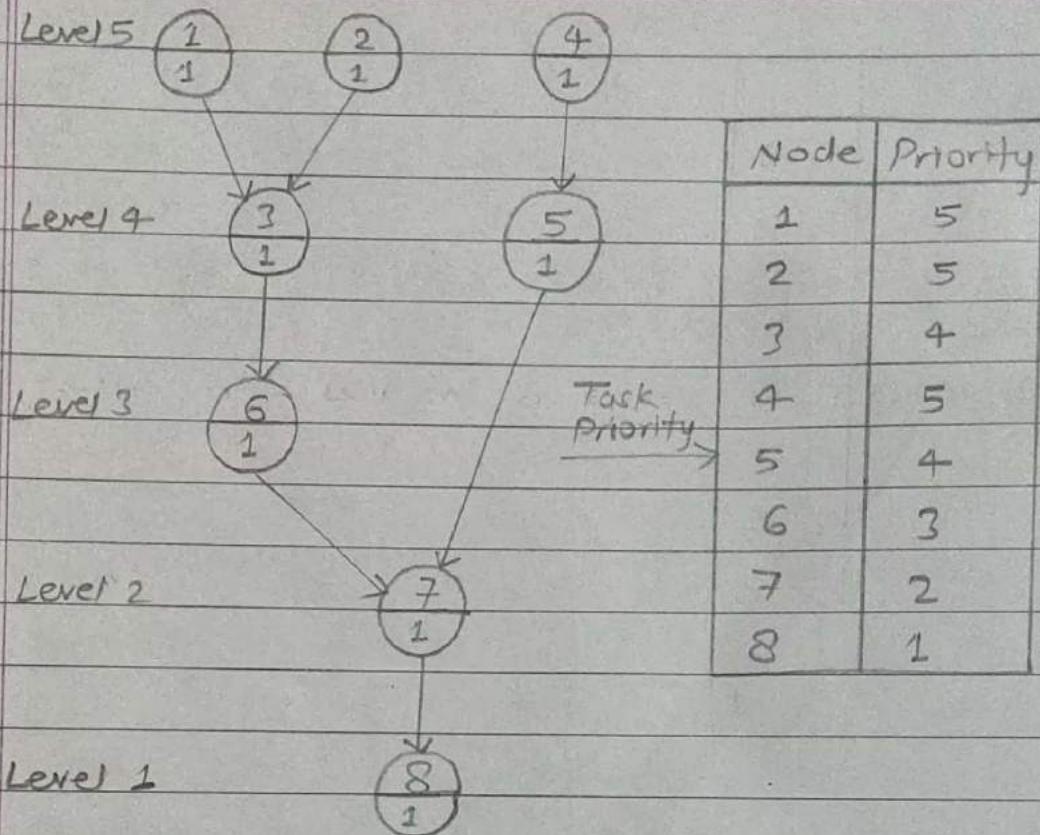
Level 3

Level 2

Level 1

Node Priority

	14	5		P ₁	P ₂	P ₃	
Task Priority →	13	5		0	14	13	11
	12	4		1	12	10	9
	11	4		2	8	6	7
	10	4	Task scheduling →	3	5	3	4
	9	4		4	2		
	8	4		5	1		
	7	3					
	6	3					
	5	3					
	4	2					
	3	2					
	2	2					
	1	1					



Task scheduling →

	P_1	P_2	P_3
0	1	2	4
1	3		5
2	6		
3			7
4			8
5			

B. Scheduling in-forests / out-forests with communication

- Based on the idea of adding new precedence relations to the task graph
- Scheduling the augmented task graph without communication is equivalent to scheduling the original task graph with communication.

Task graph (T) \rightarrow Aug. task graph (T') \rightarrow without C

Node-depth with C

- The depth of a node is defined as the length of the longest path from any node with depth zero to that node.

operation swap-all

- given a schedule f , the operation swap-all (x, y, t) where x and y are two tasks in f scheduled to start at time t , on processors i and j .
- The effect of this operation is to swap all the task pairs scheduled on processors i and j is the scheduled at time t_1 , $\forall t_2$, $t_1 \geq t$.

Algorithm

1. Given an in-forest $G = (V, E)$, identify the set of siblings s_1, s_2, \dots, s_k where s_i is the set of all nodes in V with a common child (s_i)
2. $A_1 \leftarrow A$
3. For every set s_i
 - 3.1 Pick node $u \in s_i$ with the maximum depth.

3.2 $A_2 \leftarrow A_2 - (v, \text{child}(s_i))$, $\forall v \in s_i$ and $v \neq u$

3.3 $A_2 \leftarrow A_2 \cup (v, u) \quad \forall v \in s_i$ and $v \neq u$.

4. Obtain the schedule f on A_2 without communication.
5. For every set s_i in the original task graph G , if node v is scheduled in the time slot immediately before $\text{child}(s_i)$ but in a different processor, apply swap-all.

Example:

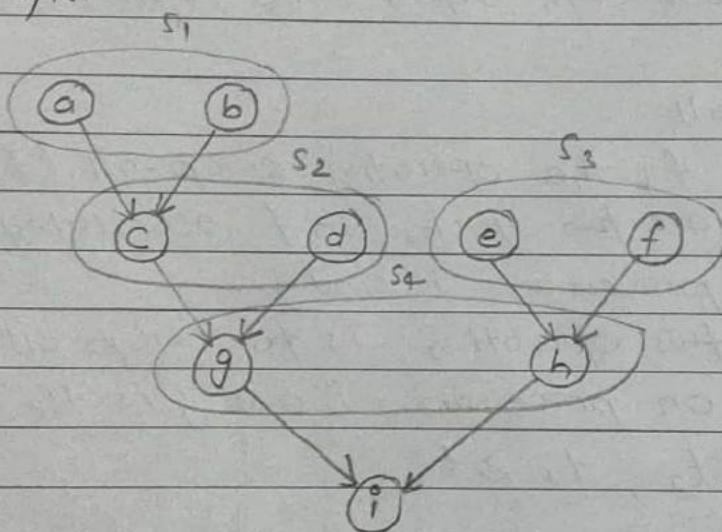
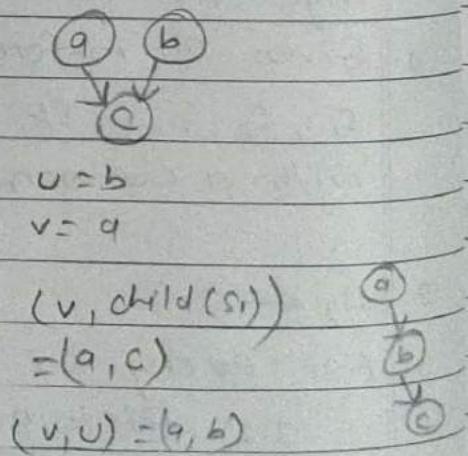


Fig 1: Original task graph

max depth (v)	
$s_1 = \{a, b\} \Rightarrow b$	$v = b$
$s_2 = \{c, d\} \Rightarrow c$	$v = d$
$s_3 = \{e, f\} \Rightarrow f$	$v = e$
$s_4 = \{g, h\} \Rightarrow g$	$v = h$



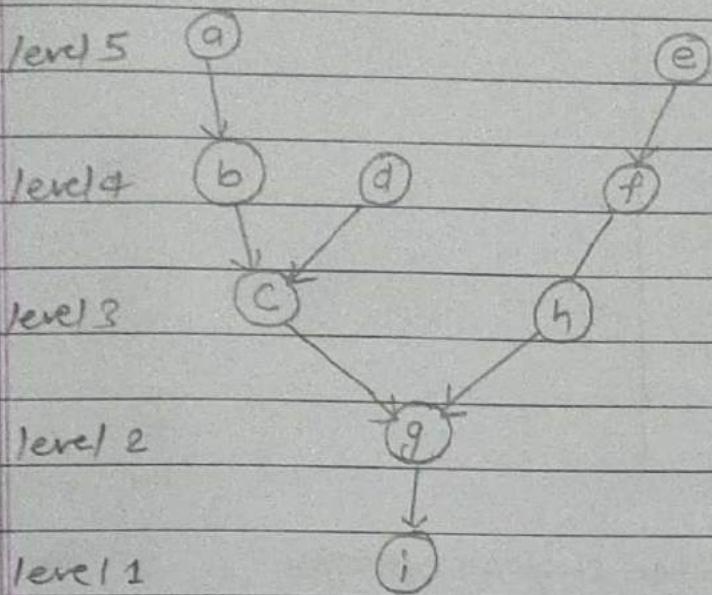


Fig2: Augmented task graph

	Node	Level				
	a	5				
	b	4		Time	P ₁	P ₂
	c	3		0	a	e
<u>Task priority</u> →	d	4	Task scheduling	1	b	d
	e	5	→	2	f	c
	f	4		3	h	.
	g	2		4	g	
	h	3		5	i	
	i	1				

Fly 4: Task scheduling

Fig3: Task priority

Time	P_1	P_2
0	a	e
1	b	d
2	c	f
3		h
4		g
5		i

Fig5: Task schedule after swap all operation