# Distributed System.
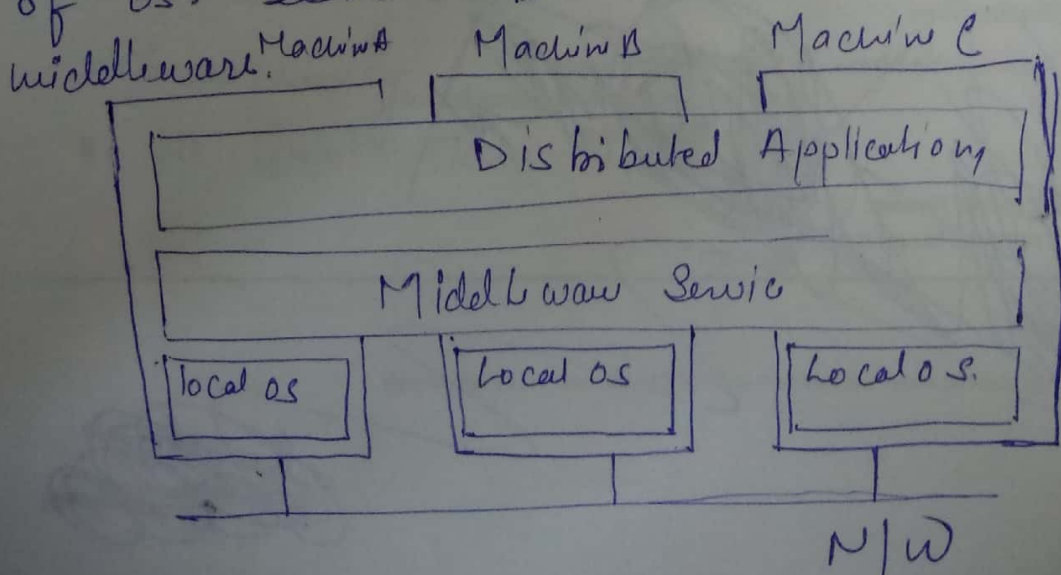
* A DS is a collection of independent computer that appears to the users of the system as a single computer

DS has Two Aspects → The 1st one deal with h/w
→ the machines are autonomous

* The 2nd one deal with s/w: the users think they are dealing with a single system.
  [Both are essential]

* Users are applications can interact with a DS in a consistent and uniform way,

* It is easy to expand & scale.

* To support heterogeneous. Computers & n/w while offering system view, DS are often organized by means of a layer of s/w that is logically placed b/w a higher level layer consisting of user & Application. and a layer underneath consisting of OS. whereas DS is sometimes called middleware.

Machine A    Machine B    Machine C

Distributed Application

Middleware Service

| local OS | Local OS | Local O.S. |

N/W

Examples of DS

* The Internet
* An Intranet
* Mobile Computing.

The Sharing of resources is a main motivation for Constructing Distributed Systems. Resources may be managed by servers and accessed by Clients or they may be encapsulated as object and accessed by other client objects.

## Challenges:

* Heterogeneity of its Components
* Openness. ( Which allow Components to be added or replaced.
* Scalability: the ability to work well when the no of Users increases.
* Failure Handling
* Concurrency of Components
* Transparency.
* No Global Clock

# Resource Sharing and Web Challenges.

1) **Heterogenity:** The internet enables users to access services and run applications over a heterogeneous collection of computer & n/ws
It applies to all of the following:

1) N/w
2) Computer H/w
3) OS
4) Programming language
5) Implementations by different Developer

* To solving the problems of heteroginity, middleware provide a uniform computational model foruse by the programmers of servers and distributed application. Ex: Remote Object Invocation, CORBA -) (Common Object Request Broker) provide remotes Remote object invocation. which allow an object in a program running on 1 computer to invoke a method of an object in a program running on another computer

# Openness:

1) Openness of a Computer system is the characteristic that determines whether the system can be extended and reimplemented in various ways.

2) The openness of D S is determined by the degree to which new resource sharing services can be added and be made available for use by a variety of client program.

3) Openness can not be achieved unless the specification and documentation of the key s/w interface of the components of a system are made available to s/w developer. The key interfaces r published

4) Systems that are designed to support resource sharing in this way are termed open D.S. to emphasize the fact that they are extensible. They may be extended at the h/w Level by the addition of Computers to the n/w. and at the s/w Level. by the introduction of new services and the re-implementation of old ones.

$\boxed{\text{RFC - Request for Comment}}$

## Security:

1) Many of info. resources that are made available and maintained in DS have a high intrinsic value to their users. Their security is importance.

2) Security for information resources has 3 components

   •) Confidentiality

   •) Integrity

   •) Availability

3) In DS client send request to access data Managed by servers which involves sending info. in message over a n/w

※ following two security challenges have not yet been fully met-

   ※ Denial of service Attack

   ※ Security of Mobile Mode

## Scalability:
DS operate effectively and efficiently at many different scales, ranging from small intranet to Internet

A system is scalable if it remain effective when there is a significant increase in the intranet no of resources and the no of user

   ※ The design of scalable DS presents the following challenges.

   ※ Controlling the cost of physical resources

   ※ Controlling the performance loss

# Failure Handling:

Computer System sometimes fail. When faults occur in h/w or s/w. programs may produce incorrect result or they may stop before they have computed the intended Computation

Handling of failures is difficult.

The techniques for dealing with failures are

1) **Detecting failures:** Some failures can be detected. for Example. CheckSum can be used to detect Corrupted data. in a msg or file.

The challenge is to manage in the presence of failure that can not be detected but may be suspected

2) **Masking Failure:** Some failure that have been detected can be hidden or made less severe.

Two examples of hiding failures:

1) Messages can be retransmitted when they fail to arrive

2) File data can be written to a pair of disk So that if one is corrupted the other may still be correct

Hiding failure are not guaranteed to work in the worst cases

## Tolerating failures:

**Recovery from failure:** Recovery involves the design of s/w so that the state of permanent data can be covered or rolled back after a server has crashed

**Concurrency :** The presence of Multiple users in a DS is a source of concurrent request to its resources. Each resource must be designed to be safe in a concurrent environment

**Transparency :** The aim is to make certain aspects of Distribution invisible to the application programmer so that they need only be concerned with the design of their particular application

for Ex : they need not be concerned with its location or the details of how its operations are accessed by other components.
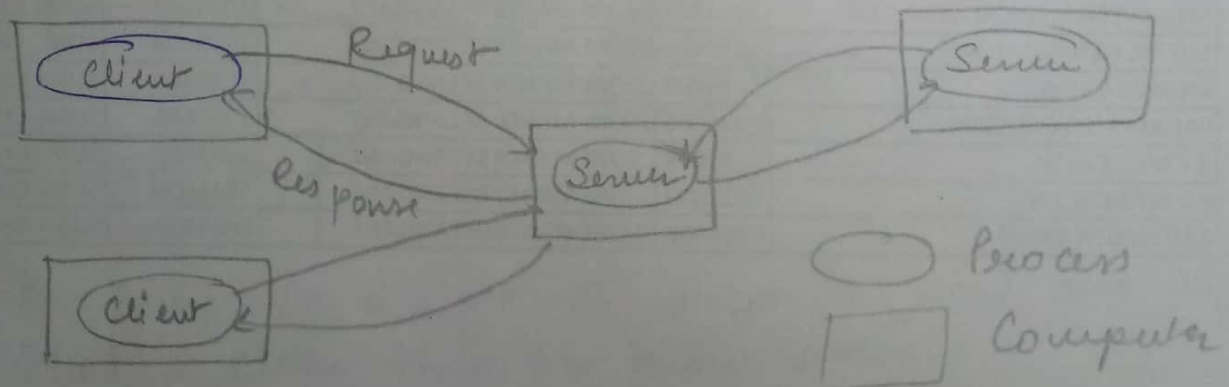
# Architectural Model

- Architectural Model is an abstract view of a D.S.
- Models are constructed to simplify reasoning about the System
- The Goal of A.M is to ensure that the structure will meet present & future demand on it
- Major Concern are to make the system reliable, manageable, adaptable and cost effective
- A Model of DS is expressed in term of

  Components

  Placement of Components

  Interaction among Components

- Component of a DS is a process (running program)

## System Architecture Models

1) Client - Server

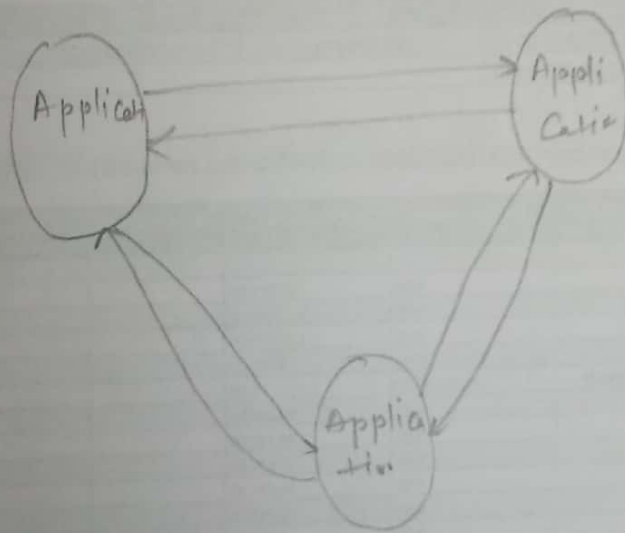2) Peer to Peer

3) Variations

### Client - Server Architecture

) Client Send Request to server
) Server Send Responses to client

Server may be client of other servers.

→ A web server is often a client of a file Server

→ An Internet Service is a client of a DNS server

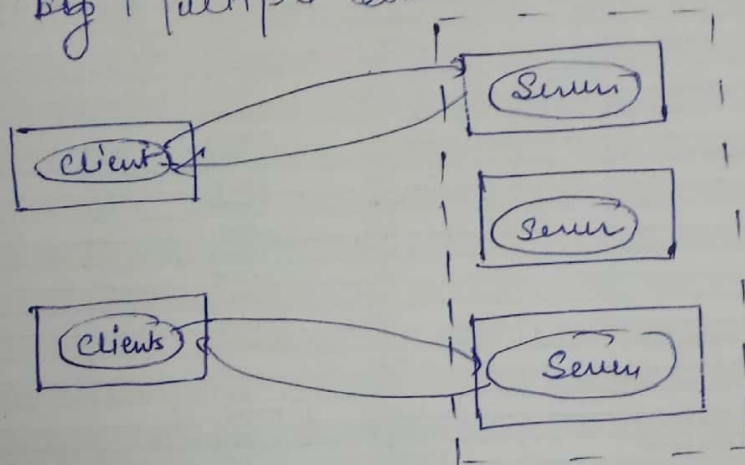- a server that translates DNS Name to IP address

Peer to Peer Architectures



•) All processes play similar role. i.e. they interact as peer

• No Central Component : Potentially better Scalability and resiliency to failure

•) Use the power of Modern desktop to implement a large Scale D.S

Ex. Skype. Bittorrent. Napster

Architectural Variations:

→ Services Provided by Multiple Server

⇒ Proxy Servers and Caches.

⇒ Mobile Agent
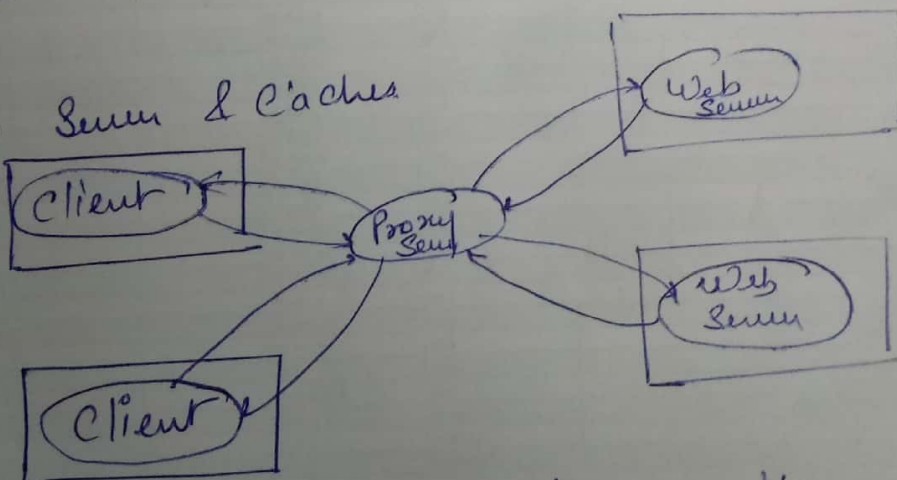
⇒ N/w Computer

→ Mobile devices

1) Services By Multiple Servers.



1) Multiple Servers provides services to client

2) Server May Partition the services object or replicate them

2) Proxy Server & Caches



1) A cache is a store of recently used data objects that is closer than the main store

A newly accessed object is added to the Cache. when that object is accessed again it is fetched from cache, if there is an up to date copy in the Cache.

1) Proxy Servers intercept communication with the real Server to provide faster Service. better Security.

## Mobile Agent

•) A running program (both Code and data) that transfers from 1 Computer to another

Ex. A worm
Used to attack Computer System

Mobile Device

→ Cellular Phone

⇒ Laptop's

→ Mobile Sensors

# Fundamental Model:

) understanding the characteristics that impact distributed System performance and operation.

Purpose of Model: To make generalization concerning what is possible or impossible

## Interaction: Model
Distributed Systems are comprised of entities that communicate and coordinate by passing messages. The following characteristics of communication channels impact the performance of the System.

- Latency - the time b/w the sending of a message at the Source and the receipt of the message at the destination.

- Bandwidth - the total amount of information that can be transmitted over a given time period.

- Jitter - the variation in the time taken to deliver a series of messages.

⇒ Coordination of actions of entities in a DS. is impacted by the fact that each entity will have a different clock drift rate.

⇒ Synchronous DS. that relay on certain action happening at the same time can only be built if you can guarantee bounds on System resources. and clock drift Rate

# Failure Model

It is important to understand the kind of failures that may occur in a system.

•) **Fail stop:** A process halts and remains halted. Other processes may not be able to detect this state

•) **Crash:** A process halts and remains halted. Other Processes ~~can~~ may not be able to detect this state.

•) **Omission:** A message inserted in an outgoing message buffer never arrives at the other ends incoming message buffer.

• **Send-Omission:** A process completes a send but the message is not put in its outgoing message buffer

• **Reciver-Omission:** A message is put in a process's incoming message buffer but that process does not receive it

➡ **Timing failure:** Clock drift exceeds allowable bounds.

⇒ **Byzantine failures**
   • Bug
   • Message carruption
   • Hardest to deal with

## Security Model:

There are several potential threat a system designer need be aware of.

* **Threat of Processes:** An attacker may send a request or response using a false identity (spoofing)

* **Threat to Communication channel:** An attacker may eavesdrop (listen to message) or inject new message into a Communication channel. An attacker can also save messages and reply them later

* **Denial of Service:** An attacker may over load a server by making excessive requests.

⇒ Cryptography and authentication are often used to provide adequacy security.

* Communication entities can use a shared Secret to ensure that they are communicate with one another and to end encrypt their message. So that they cannot be read by attackers.

# Limitation of Distributed System

1) A DS is a set of computers that communicate over a network, and do not share common memory and common clock.

2) Absence of a common clock (Global)

    i) No concept of Global time

    ii) Its difficult to reason about the temporal ordering of events.

       1) Cooperation b/w processes - (eg. Producer | Consumer, client | Server)

       2) Arrival of request to the OS (eg for resources)

       3) Collecting up-to-date global state.

3) It is difficult to design and debug algorithm in a D.S.

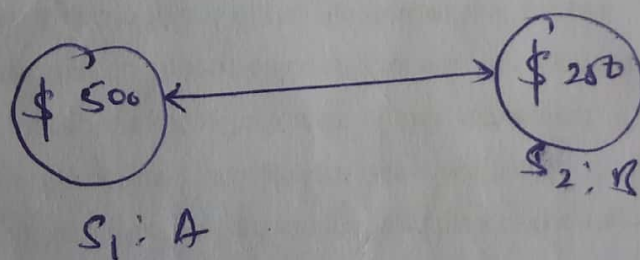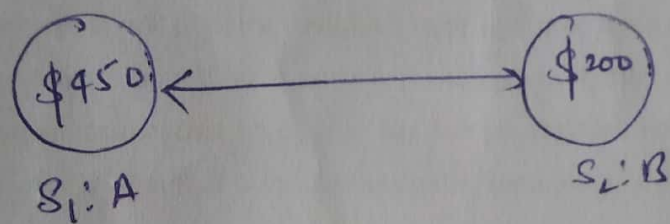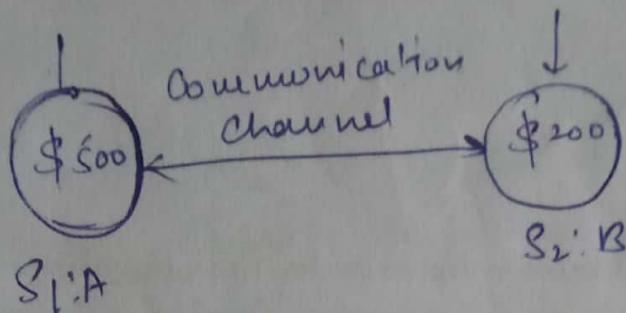    •) Mutual Exclusion

    •) Synchronization

    •) Deadlock

## Absence of Shared Memory

1) No up to date state of the entire system to any individual process as there's no shared Memory.

2) Coherent view: All observations of different processes (computers) are made at the same physical time.

    •) We can obtain a coherent but partial view of the system or incoherent view of the system.

    •) Complete view (Global state) — local view (local states) + Message in transit.

# Absence of shared Memory

local State A             local state B      Processes_100



$S_1 : A$             Communication channel      $S_2 : B$

$450 \leftarrow \longrightarrow $200

$S_1 : A$            $S_2 : B$

$500 \leftarrow \longrightarrow $250

$S_1 : A$            $S_2 : B$

**Logical Clock:** Lamport invented a simple mechanism by which the happened before ordering can be captured numerically called a logical clock

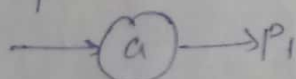*) A lamport logical clock is a monotonically increasing s/w counter.

## Drawback

1) Clocks are not strongly consistent : clock lose track of the timestamp of the event on which they are dependent on.

This is because we are using a single integer to store the local and logical time

2) if $a \to b$ then $c(a) < c(b)$. the reverse is not necessarily true if the events has occured in different processes.

That is. if a & b are event in different processes and $c(a) < c(b)$ then $a \to b$ is not necessarily true.

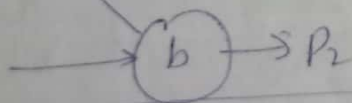events a & b may be causally related or may not be

# Logical clock:

- To order event across process, trying to sync clocks in one approach

- Another Approach = Assign timestamps to events

- As long as these time stamps obey causality, this would work

## Happen Before Relationship

(1)

$\boxed{a} \longrightarrow \boxed{b} \longrightarrow P$      $Ts(a) < Ts(b)$

if a and b are within the same process. and A occurred before B.

(2)

$\longrightarrow \boxed{a} \longrightarrow P_1$

$\longrightarrow \boxed{b} \longrightarrow P_2$   if a is the event of sending a msg M in own Process & b is the event of receiving M by another process.

$Ts(a) < Ts(b)$

## Properties Derived from (HBR)

### Transitive Relation

if
- $Ts(a) < Ts(b)$
- $Ts(b) < Ts(c)$

Then

$Ts(a) < Ts(c)$

### Causally Ordered Relation

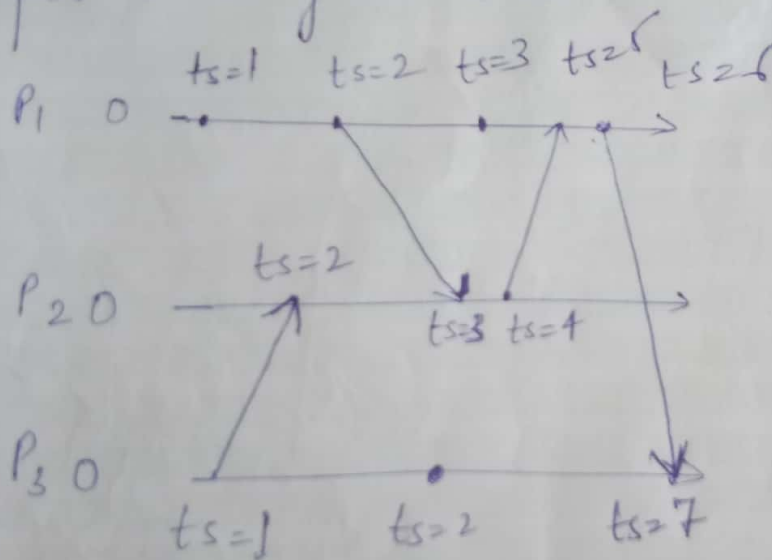$a \rightarrow b$

changes of a, affect b.

Concurrent Event

A → B    X

B → A    X

A || B

Lamport's @logical Clock

$P_1$  0  — • ts=1    ts=2  ts=3  ts=2$^5$ ts=6

ts Receiv = Max(own ts, Sender TS)

+1

$P_2$ 0 — ts=2 ↗    ts=3  ts=4

$P_3$ 0    ts=1    ts=2    ts=7

Event occur

ts = ts+1

Max (0,1) +1
1 +1 = 2

Max (2, 2) +1
2 +1 = 3

ts( 3, 4) +1
4+1 = 5

ts (2,6)+1
6+1 = 7

## Algo:

1) A process increment its ~~counter~~ Counter before each event in that process.

2) When a process sends a message it includes its Counter value with the msg

3) On receiving a msg. the counter of recipient is updated, if necessary. to the greater of its Current Counter & the timestamp. in the received msg. The counter is then incremented by 1 before the msg is considered received

algo for sending is

$$time = time + 1$$
$$time\ stamp = time$$
$$Send\ (Msg,\ time\_stamp);$$

Algo for receiving a msg

$$(msg,\ timestamp) = receive();$$
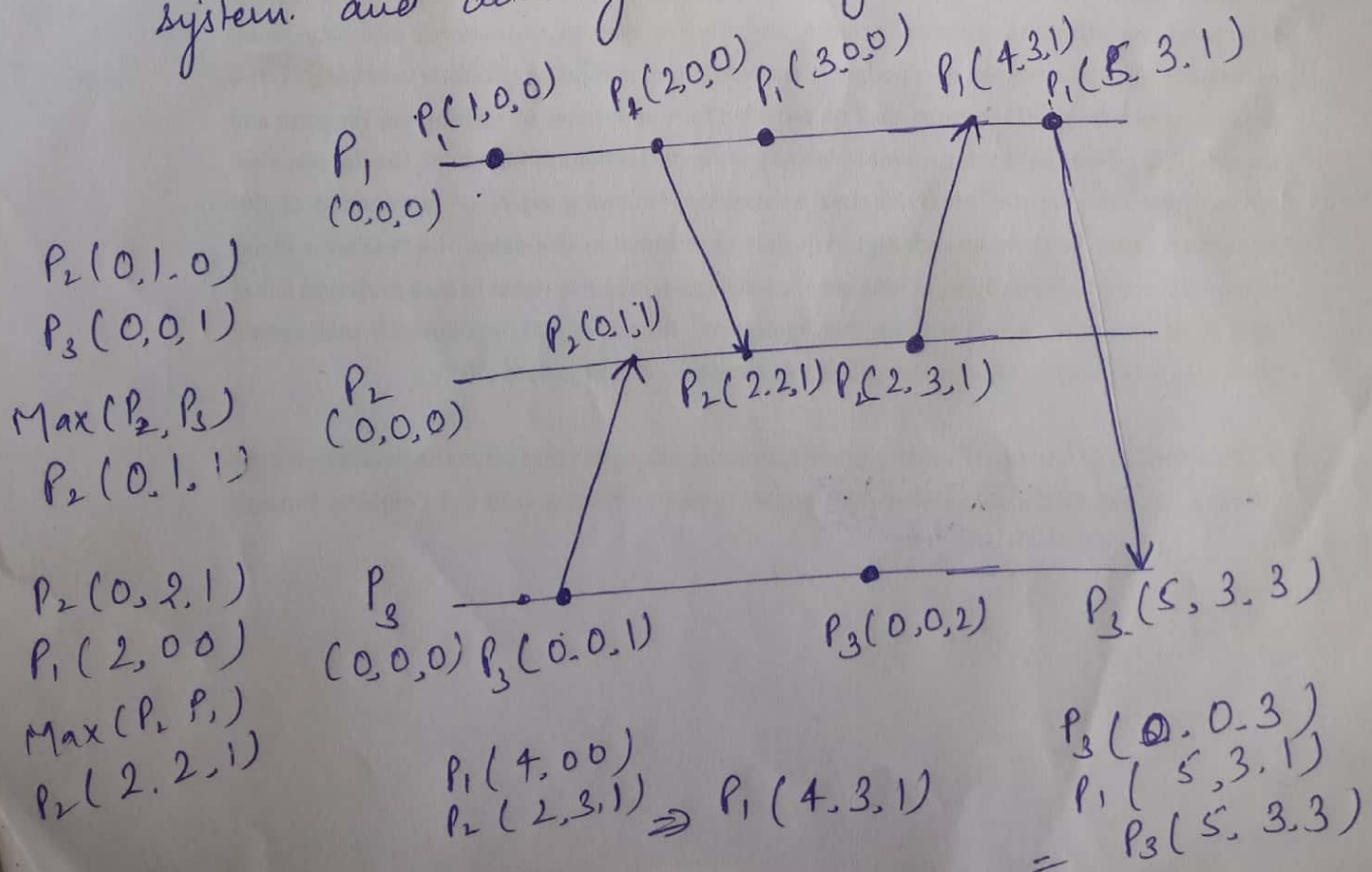$$time = max\ (time\_stamp,\ time\ ) + 1;$$

---

* A logical Clock is a mechanism for Capturing chronological and Causal relationships in a D.S.

* D.S. may have no physically synchronous global clock so a logical clock allow Global ordering on events from different processes in such system.

# Vector logical clocks:

1) Vector clocks overcomes the shortcoming of Lamports clock – The fact that $L(e) < L(e')$ we can not conclude the $e \rightarrow e'$

2) A Vector clock for a system of N processes is an array of N integers.

3) Each process keeps its own vector clock $V_i$, which it uses to timestamps on the message local event.

4) like Lamport timestamps, processes piggyback vector timestamps on the message they send to one another

5) A Vector clock is an algorithm for generating a partial ordering of events in a distributed system and detecting causality violations.

$P_1$ $\quad P_1(1,0,0)\quad P_4(2,0,0)\quad P_1(3,0,0)\quad P_1(4,3,1)\quad P_1(5,3,1)$
$(0,0,0)$

$P_2(0,1,0)$
$P_3(0,0,1)$

$Max(P_2, P_3)$ $\quad P_2$ $\quad P_2(0,1,1)$
$P_2(0,1,1)$ $\quad (0,0,0)$ $\qquad P_2(2,2,1)\ P_2(2,3,1)$

$P_2(0,2,1)$ $\quad P_3$
$P_1(2,0,0)$ $\quad (0,0,0)\ P_3(0,0,1)$ $\qquad P_3(0,0,2)$ $\quad P_3(5,3,3)$
$Max(P_2\ P_1)$
$P_2(2,2,1)$

$\qquad\qquad P_1(4,0,0)$
$\qquad\qquad P_2(2,3,1) \Rightarrow P_1(4,3,1)$ $\qquad P_3(0,0,3)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad P_1(5,3,1)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad = P_3(5,3,3)$

# Algorithm :

Step1 = Vector initialize to 0 for each process.

$$Vi[j] = 0 \quad \text{for } i, j = 1, 2, \cdots M.$$

Step 2 = Increment Vector before time stemp event

$$Vi[i] = Vi[i] + 1$$

Step 3 = Message is sent from Pi with Vi attached to it

Step 4 : when Pj receive message

$$Vj[i] = max(Vi[i], Vj[i]) \quad \text{for } i = 1 \cdots M.$$

Vector Timestamps have the disadvantage, compared with Lamport timestamps. of taking up an amount of storage and message payload that is proportional to M, the no of processes.

# Concept of Message Passing System
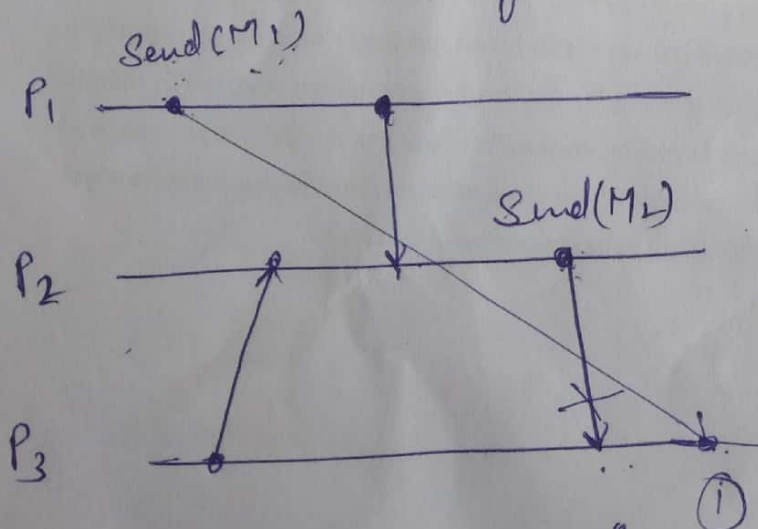
## Causal Ordering of Message

The purpose of Causal ordering of message is to insure, that the same Causal relationship for the "message send" event correspond with message receive event. i.e.

All the messages are processed in order that they were created.

In other words: if $Send(M_1) \rightarrow Send(M_2)$ (when Send(M)) is event sending message M). then every recipient of both messages $M_1$ & $M_2$ must receive $M_1$ before $M_2$.

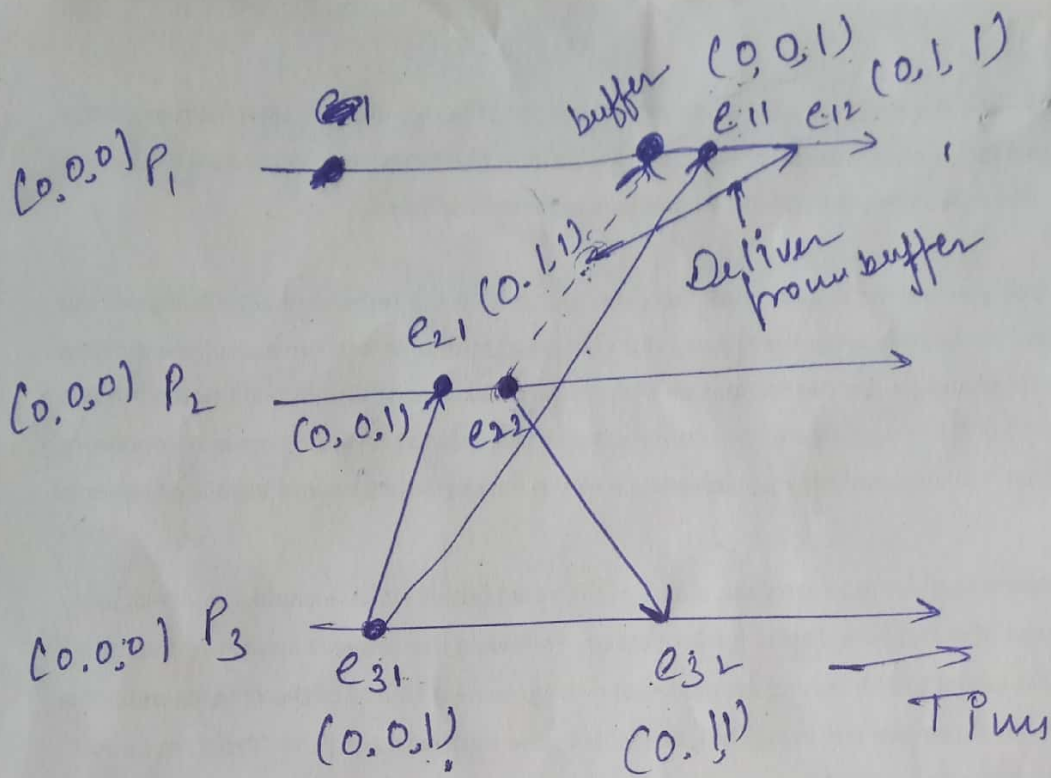There are Two protocols that make use of Vector clocks for the Causal ordering of messages in D.S.



An example of violation of Causal ordering of message

## BSS Algorithm: Birman-Schiper-Stephenson Protocol

(1) Broadcast based: A message sent is received by all other processes.

2) Deliver a message to a process only if the message preceding it immediately has been delivered to the process.

➔ otherwise, buffer the message

3) Accomplished by using a vector accompanying the message

## Algorithm:

1) Process $P_i$ increments the vector time $VT_{P_i}[i]$, time stamps, and broadcasts the message $m$. $VT_{P_i}[i]-1$ denotes the no of messages preceding $m$.

2) $P_j \neq P_i$ receives $m$, $m$ is delivered when:
   a) $VT_{P_j}[i] = VT_m[i] - i$
   b) $VT_{P_j}[k] >= VT_m[k]$ for all $k$ in $\{1,2,--n\}$
   - [i], $n$ is the total no of processes.
   Delayed message are queued in a sorted manner.
   c) Concurrent message are ordered by time of receipt.

3) When $m$ is delivered at $P_j$, $VT_{P_j}$ updated according Rule 2 of vector clocks.

$(0,0,0)\ P_1$     $e_{21}$     buffer $(0,0,1)$    $(0,1,1)$
                                $e_{11}$   $e_{12}$

$e_{21}\ (0,1,1)$     Deliver from buffer

$(0,0,0)\ P_2$     $(0,0,1)$    $e_{22}$

$(0,0,0)\ P_3$     $e_{31}$       $e_{32}$     Time
             $(0,0,1)$      $(0,1,1)$

There is no clock increment upon receiving.

Send $(M_1) \rightarrow$ Send $(M_2) \rightarrow$ Receive $(M_2) \rightarrow$ Send $(M_3)$

# Schiper-Eggli-Sandoz Protocol

* No need for broad cast Message
* Each process maintain a vector V_P of size M-1 (M the no of processes in the system)

* V_P is a vector of tuple ($P'$, t):
   $$P' = \text{Destination Process ID}$$
   $$t = \text{Vector timestamp}$$

* Tm: logical time of sending message m

* Tpi: Present logical time at Pi

* Initially, V_P. is empty.

## Algorithm

1) Sending a Message
   ⇒ Send Message M, timestamp tm, along with V_P₁ to P₂
   ⇒ Insert ($P_2$, tm) into V_P₁. Overwrite the previous value of ($P_2$, t) if any.
   ⇒ ($P_2$, tm) is not sent. Any further message carrying ($P_2$, tm) in V_P₁ can not be delivered to P₂ until tm < t P₂

# Delivering a Message:

1) if V-M (in the message) does not contain any pair $(P_2, t)$ it can be delivered.

2) /* $(P_2, t)$ exist */ if ~~t~~ $t \not> tp_2$ buffer the message. (Don't deliver)

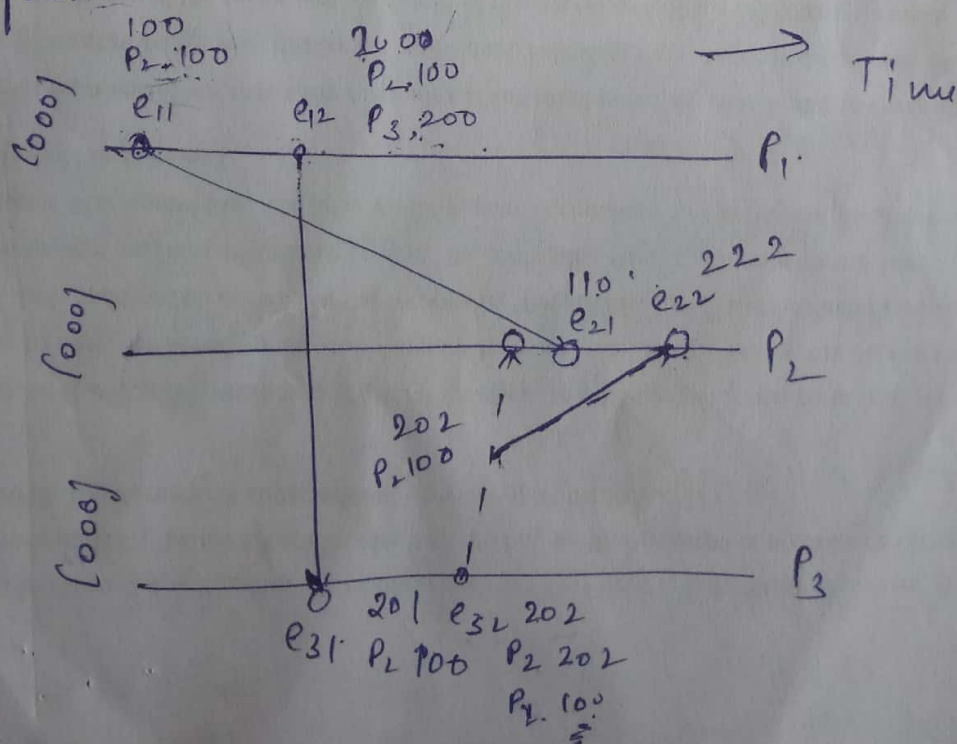3) else $(t < tp_2)$ deliver it

What does the condition $t \geq TP_2$ imply?

   t is message vector time stamp.

   $t > TP_2 \Rightarrow$ for all $j$, $t[j] > TP_2[j]$

This implies some events occurred without $P_2$'s knowledge in other processes. So $P_2$ decides to buffer the message

$\Rightarrow$ when $t < TP_2$ message is delivered & $TP_2$ is updated with help of V-$P_2$

# Global State of Distributed System.

⇒ Global state of DS is a collection of the local states of all the processes and the states of all the communication channels.

⇒ The state of process Pi is a collection of all events that happend at the process.
   — Internal message send and message receive

⇒ The state of channel is a collection of ~~the local~~ all the messages that have been sent and not yet received
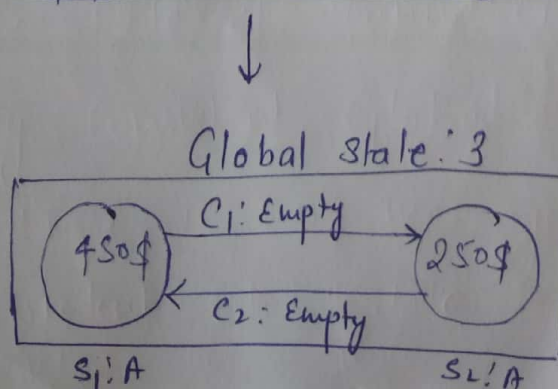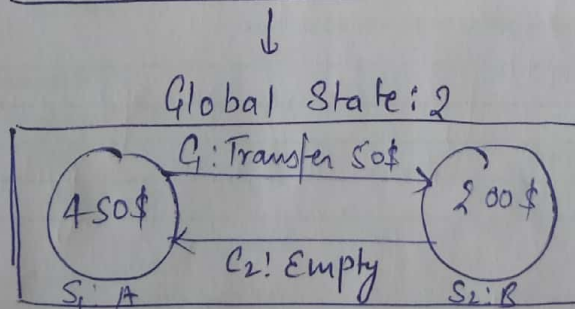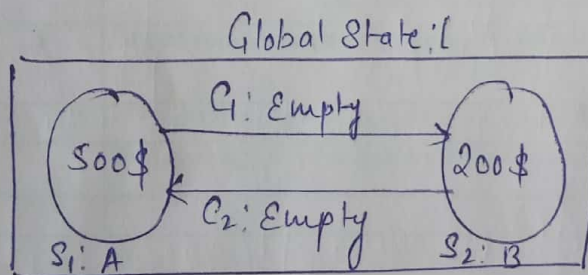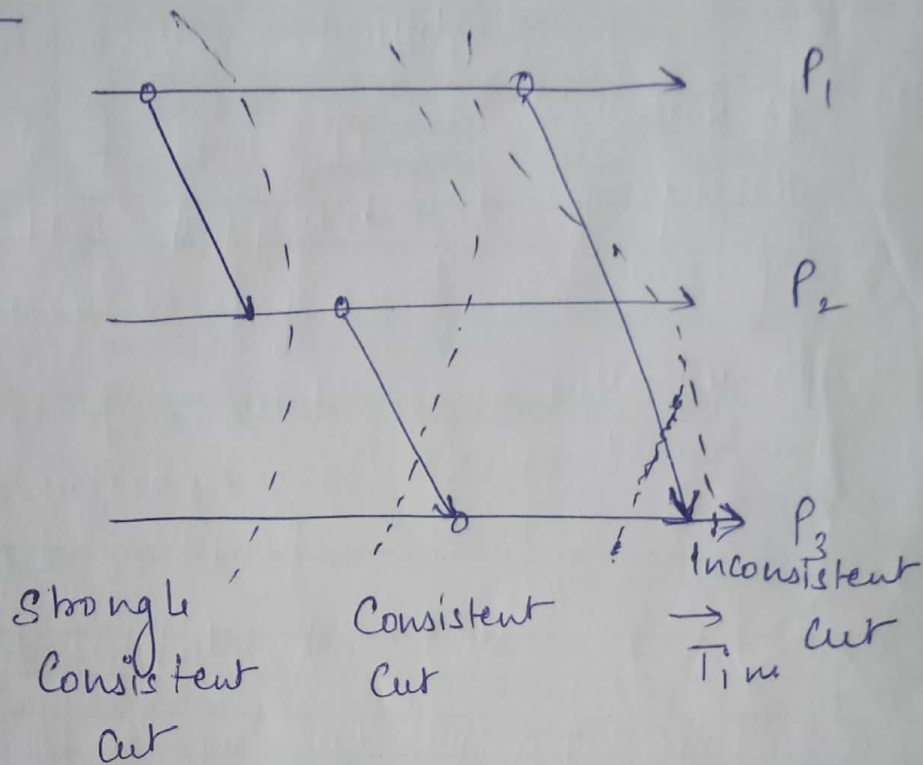
## Consistent Global State:

⇒ A message m sent by a process Pi to a process Pj and recorded in the local state of Pi should be either recorded in the local state of receiver Pj ~~and~~ or the ~~state~~ channel $C_{ij}$.

→ A message m sent by ~~the~~ a process Pi after recording its local state should not be recorded in the local state of receiver Pj and also not in the state of channel $C_{ij}$.

Strongly Consistent Global State: A consistent Global state in which there is no message in transit.

An Inconsistent Global state is the one that records the receiving of a message ~~and the sending of mess~~

in the receiving process. or as a transit message
in the channel. but the sending of message is
not recorded

## Example



Strongly
Consistent
Cut

Consistent
Cut

Inconsistent
→ Cut
Time

### Global State : 1

$C_1$ : Empty
$500\$$ → $200\$$
$C_2$ : Empty
$S_1$ : A        $S_2$ : B

↓

### Global State : 2

$C_1$ : Transfer $50\$$
$450\$$ → $200\$$
$C_2$ ! Empty
$S_1$ A        $S_2$ : B

↓

### Global state : 3

$C_1$ : Empty
$450\$$ → $250\$$
$C_2$ : Empty
$S_1$ : A        $S_2$ ! A

# Chandy Lamport's Algorithm

1) A distributed Algorithm to collect Global State

2) Assumption: All communication channel are FIFO (no overtaking of messages)

3) An initiator processor records its local state and send a marker to all its communication channels (before sending any message on the channels)

4) when a process recieves the marker
   → if it has already recorded its local state
      → Records the messages recieved after recording the local state (but before the marker was recieved) as the status of the channel.

   else:
      → Records the local state, record the state of channel as empty and send out the marker on all of its communication channels

5) All processes report their local state and the status of their associated incoming channels to the interior initiator process

# Termination Detection

→ A Distributed System generally consist of a set of cooperating processes which Communicate with each other by exchanging message. In this case, it is important to know when the Computation has terminated.

⇒ A process may either be in an active state or idle state.

→ Only active processes can send messages.

⇒ An active process may become idle at any time.

⇒ An idle process can become active on receiving a Computation message

→ Computation messages are those that are related to the underlying Computation being performed by the cooperating process

→ A Computation is said to have terminated if and only if all the processes are idle and there are no message in transit.

⇒ The message sent by the termination detection algo. are referred to as Control message

# Notations

- $B(Dw)$ = Computation message sent as a part of the computation and $Dw$ is the weight assigned to it

- $C(Dw)$ = Control message sent from the processes to the controlling agent and $Dw$ is the weight assigned to it.

# Huang's Termination Detection Algorithm

Rule 1: The controlling agent or an active process having weight $w$ may send a computation message to a process $P$ by doing:

Derive $w_1$ and $w_2$ such that

$$w_1 + w_2 = w, \qquad w_1 > 0, \quad w_2 > 0$$

$$w = w_1$$

Send $B(w_2)$ to $P$

Rule 2: On receipt of $B(Dw)$, a process $P$ having weight $w$ does.

$$w = w + \Delta w$$

if $P$ is idle, $P$ become active

Rule 3: An active process having weight $w$ may become idle at any time by doing

Send $C(w)$ to the controlling agent

$$w = 0$$

(The process becomes idle)

Rul 4. On succeiving e(Dw), the controlling agent
having weight w takes the following
actions:

$$w = w + Dw$$

if $w = 1$ Conclude that the computation
has terminated