

UNIT - I

Distributed System

- A Distributed System is one in which components located at networked computers communicate & coordinate their actions only by passing messages.
- It is a no of autonomous processing element (not necessarily homogenous) that are interconnected by a computer Network & they cooperate in processing their assigned task transparently.
 - processing element → Computer / process / processor
 - Autonomous → Node / element must atleast be equipped with their private control.
 - Interconnected → Node must be able to exchange information.
 - Transparent → If implies existence of autonomous nodes is transparent to the other users.
- What actually is distributed in distributed system-
 - ① Processing logic - processing element with different logic are distributed
 - ② function → No single computer have designed with considerably lot of functionality. (Due to cost factor, Maintenance factor).
 - ③ Data → Data can be replicated & fragmented as it better to distribute Data for application such that it is nearby it is redundant.
 - ④ Control Control is not with the single node, but Control is distributed so that autonomy of control is not there.

Distributed System has the following significant consequences →

- ① concurrency.
- ② No global clock.
- ③ Independent failure.

http://www.abc.com/distrib/index.htm
→ protocol (HTTP)
→ Domain name of the web server Hosted at
www.abc.com.

Example of Distributed System :-

- ① web search-
- ② online games-
- ③ finance Trending
(market)
- ④ internet
- ⑤ Intranet
- ⑥ Mobile & ubiquitous computing.

(use of Network in selected domain)

- { - Finance & Commerce
- The Information Society
- Creative Industry & Entertainment
- Healthcare
- Education
- Transport & Logistics
- Science
- Environment Mgt

Challenges →

- ① Heterogeneity. (DS must be constructed from the variety of different Net^W, OS, Computer H/W, & programming languages.)
- ② openness. (DS must be extensible. (Develop interfaces of DS))
- ③ Security. (DDoS attack is a big problem) (using encryption)
- ④ Scalability. (DS need to be scalable & cost of adding a user in a constant amount of in terms the resources must be added.)
- ⑤ Failure Handling. (Any process, computer or net^W may fail independently)
- ⑥ Concurrency. (of the others).
- ⑦ Transparency. (Application programmer →).

Motivations

Resources sharing is the main Motivation for construction DS
Resources - Printers, files, web pages or database records are managed by some

System Models

②

- Systems that are intended for use in real-world environments should be designed to function correctly in the widest possible range of circumstances and in the face of many possible difficulties & threats.
- Each model type is intended to provide an abstract, simplified but consistent description of a relevant aspect of distributed system design.

① Physical Model → The most explicit way in which to describe a system, they capture the H/w composition of a system in terms of the computers, mobile phones, & their interconnecting networks.

② Architectural Model → Describes a system in terms of computational & communication tasks performed by its computational elements, the computational elements being individual computers or aggregates of them supported by appropriate Network interconnections.

③ Fundamental Model → Describes a system in terms of computational & communication tasks performed by its computational elements, the computational elements being individual computers or aggregates of them supported by appropriate Network interconnections.

Takes an abstract perspective in order to examine individual aspects of a distributed system.

Difficulties & Threats for distributed Models

- widely varying modes of use.
- wide range of system environment. (It accommodates heterogeneous H/w, OS & Netw.)
- Internal Problem. (No synchronized clocks, conflict data update many nodes of H/w & SW failure involving the individual system components.)
- External Problem. (Attack on data integrity & secrecy, DDoS attacks)

① Physical Model →

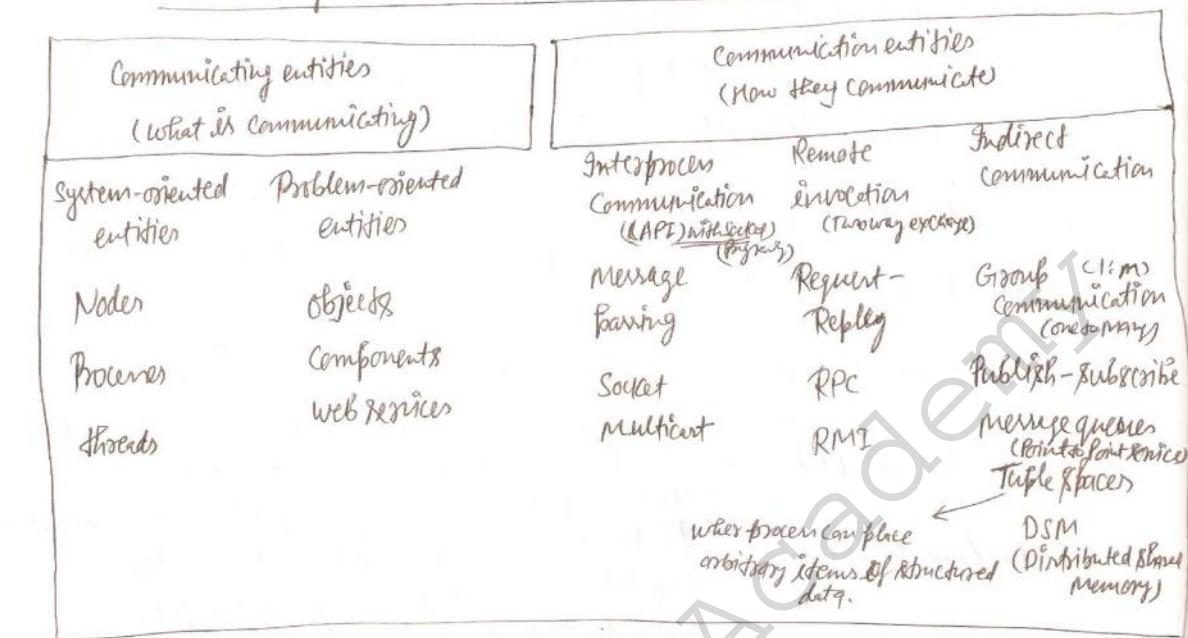
- It is a representation of the underlying H/w elements of a distributed system that abstracts away from specific details of the computer & networking technologies employed.

ex Barebone Physical Model -

where H/w or S/w component located at Networked Computers communicate & coordinate their actions only by passing messages.

#

Communicating entities & communication Paradigms



② Architectural Models

(a) Client Server Model (Search engine) -

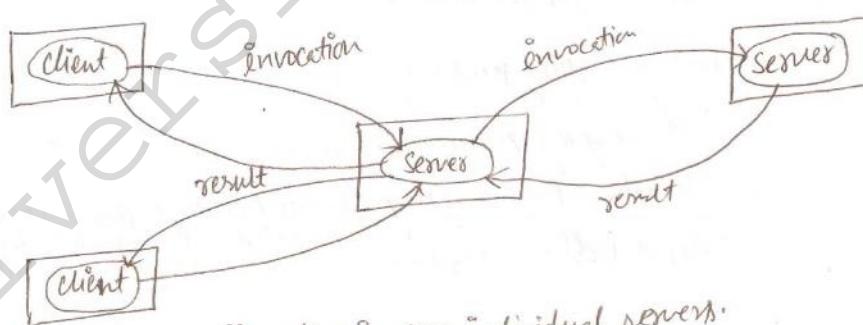
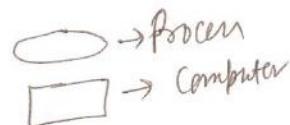
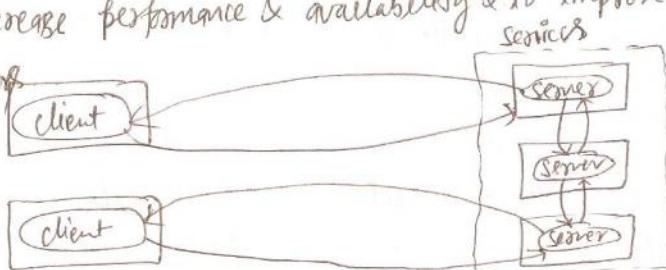


fig- client invokes individual servers.

- ③ Service provided by Multiple servers → Here data is partitioned & replicated to increase performance & availability & to improve fault tolerance.

fig- A service provided by multiple servers



③ Proxy Server & Caches

③

- A proxy server provides a shared cache of web resources for client machine at site or across several sites.
- A cache is a store of recently used data objects that is closer than the objects themselves.

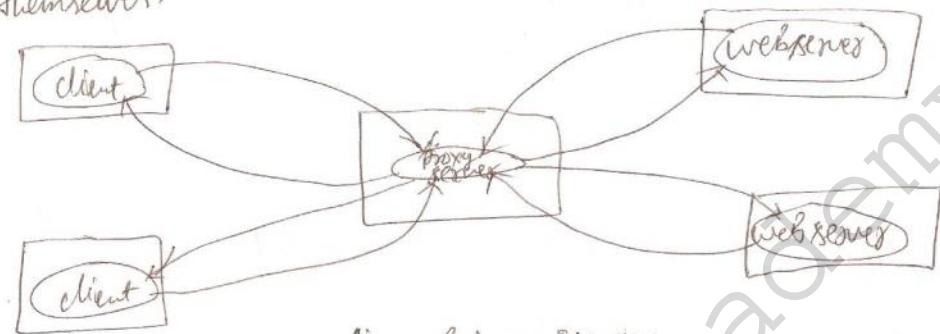


fig - web proxy server

④ Peer Processors

- All process play similar role, interacting cooperatively as peers to performs a distributed activity or computation without any distinction b/w client & server.

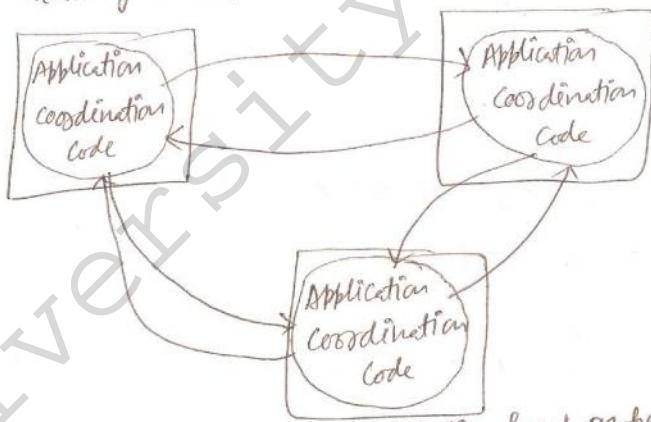
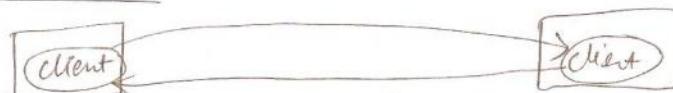


fig - Distributed application based on peer process

⑤ Variation of Client-Server Models

fig - webapplets



(a) client request results in the downloading of Applet Code



(b) client interact with the applet

- ① Mobile Code (~~Applet~~ Applet)
- ② Mobile agents: A mobile agent is running program that travels from one computer to another in network carrying out a task on someone's behalf such as collecting information eventually returning with results.
- ③ Network Computer → It downloads the operating system & any application sw needed by the user from remote file server. Application are run locally but the files are managed by remote file servers, thus user can migrate to another m/c easily.
- ④ Thin clients → The term thin client refers to sw ~~background~~ layer that support a window based user interface on a computer that is local to the user while executing applications programs on a remote computer.

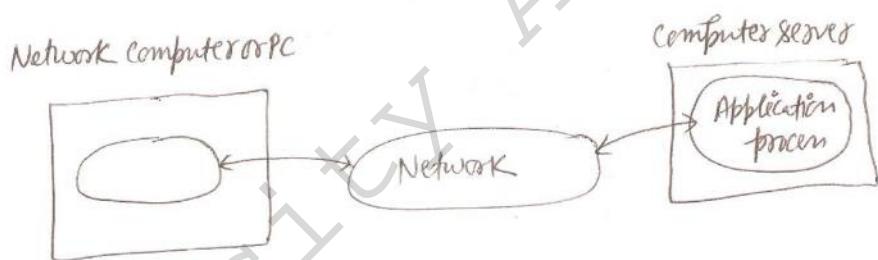


fig - Thin clients & Computer Servers

- ⑤ Mobile devices & Spontaneous networking →
 - Mobile devices - GSM, Bluetooth, IR etc.
 - Spontaneous net → Mobile devices + other device e.g. PDA
- ```

graph TD
 HWN([Hotel wireless Net]) <--> Internet
 HWN <--> DS[Discovery service]
 HWN <--> TV[TV/pc]
 HWN <--> Laptop[LAPTOP]
 HWN <--> PDA[PDA]
 HWN <--> Camera[Camera]
 HWN <--> GD[Guest device]
 HWN <--> MS((Music service))
 HWN <--> AS((Alarm service))

```
- fig - Spontaneous net work in Hotel.

### ③ Fundamental Model includes -

④

- ⑤ Interaction Model → It is concerned with performance of process communication channels and absence of global clock. Interaction Model is further classified as Synchronous systems.  
(ie based on process execution time, message delivery time & clock drift)
- ⑥ Asynchronous systems There is no bound on execution time, message delivery time & clock drift as Internet.
- ⑦ Failure Model → It classifies failure of processes & basic communication channels in a distributed system. It is based on marking, integrity & validity approach.
- ⑧ Security Models It identifies the possible threats to processes & communication channels in an open distributed system such as, integrity, auth, privacy, etc.

### Limitation Of Distributed System →

- Global Knowledge is not readily available, since it is impossible to collect up-to-date information of global state of distributed system.

Because of (limitation) :-

① Lack of Global Clock (common clock)

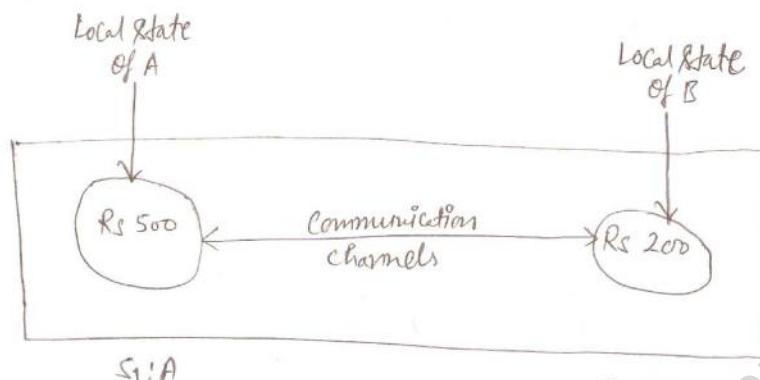
② unprecedented Message delay

③ Non-existence of Physically Shared Memory.

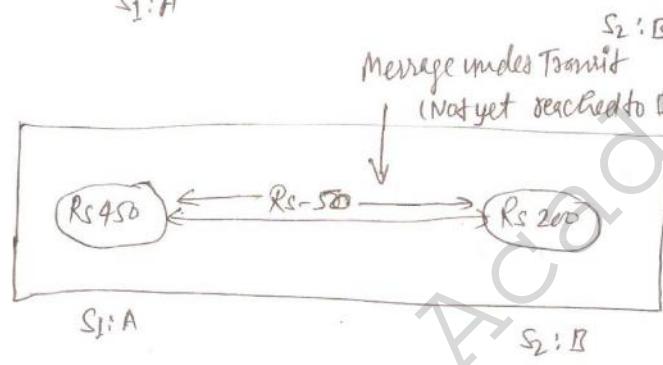
These limitation create difficulty in obtaining coherent global state.

Ex-

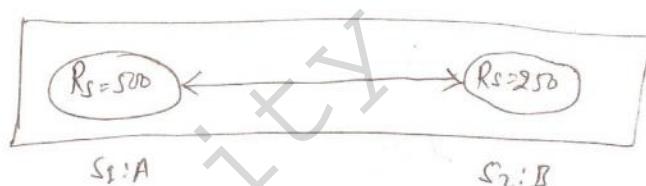
Initial  
State



Case A -



Case B



Case 1 - S<sub>1</sub> records its ~~global~~ local state (Rs 450) just after debit (-50) & S<sub>2</sub> records its location (200) before receiving. If commit merge is not taken care off -

Global State = local states<sub>1</sub> + local states<sub>2</sub>

$$= 450 + 200$$

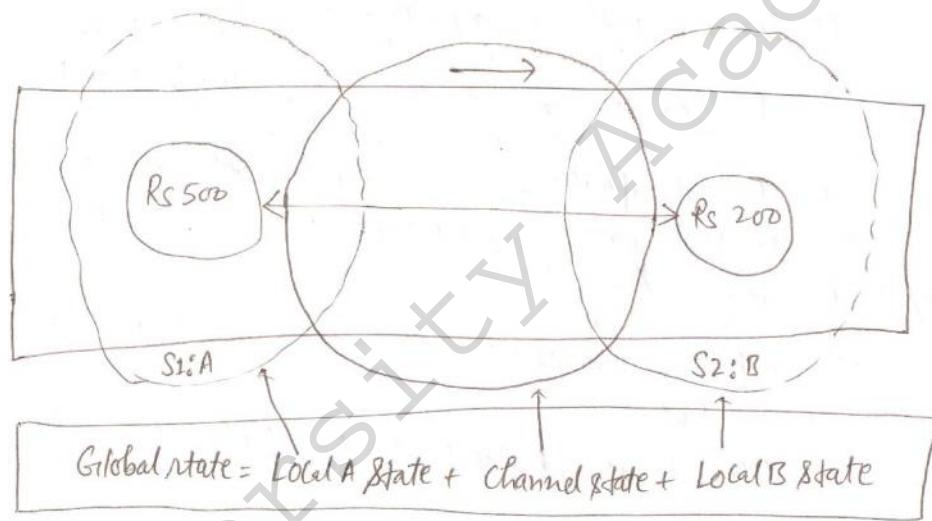
= 650 > 50 Rs missing i.e. Inconsistent system

Case 2 :  $s_1$  records its state (Rs-500) before debit &  $s_2$  records its state after credit hence there is surplus (Rs.50) which is also in-coherent.

From above two example, it is clear that communication channel cannot record its state by itself.

- Hence, nodes must have to coordinate their state recording activities in order to record the channel state or coherent state of system.

$$\text{Global State (GS)} = \text{Local State } s_1 + \text{Local State } s_2 + \text{Channel State between } s_1 \text{ & } s_2$$



\* Some types of H/w & S/w Resources that are distributed ~~are~~ are

Given as →

① H/w Resources -

- ① CPU - Computer Server, remote object servers, work programs, Cache servers use file servers, Remote servers their CPU is shared resources.
- ② Memory - Cache servers in shared holds recently access web pages.
- ③ Disk - file servers, virtual disk servers, video on demand servers.
- ④ Screen - Net & windows systems (x-11), allow processes in remote computers to update the content of windows.
- ⑤ Printer - Networked printers accept print jobs from many computers, managing them via a queue system.
- ⑥ Network Capacity - packet transmission enables many simultaneous communication channels (streams of Data) to be transmitted in the same circuit.

② Data/S/w Resources

- ① web pages - web servers enable multiple clients to share read-only page content.
- ② file - file servers enable multiple clients to share read-write files.
- ③ object - possibility for S/w objects are limitless.
- ④ database - Database are intended to record the definitive state of related sets of data.
- ⑤ Newsgroup Content - (News on internet)
- ⑥ Video/Audio Stream - server can store entire video on disk & deliver them at playback speed to multiple clients.
- ⑦ exclusive lock - a system level object provided by lock servers, enabling several clients to coordinate their use of a resource. (such as printer that does not include a queuing scheme)

⑥

### The impact of absence of global time or global clock $\rightarrow$

- It will be difficult to order event, that is temporal ordering of events. This will result in difficulty in design & development of distributed system.

$\rightarrow$  An OS is responsible for scheduling process. A basic criterion used in scheduling is the temporal order in which requests to execute process arrives. Due to the absence of global clock it is difficult to reason about the temporal order of events in a distributed system.

- Hence algorithms for DS is more difficult to design & debug as compared to algorithms for centralized system. In the absence of Global Clock makes it harder to collect up-to-date information on state of entire system.

### Lamport's Logical Clock $\rightarrow$

- An event changes the system state, which in turn influences the occurrence & outcome of future events. That is past events influence future events & this influences among causally related events referred as "causal effects".

So logic we define

#### ① Happened before relation ( $\rightarrow$ ) (causally related events of $a \rightarrow b$ )

$a \rightarrow b$ , means event "a happened before b".

- $a \& b$  are events in the same process & a occurred before b".
- a is sending message & b is receipt of same message.
- If  $a \rightarrow b$  &  $b \rightarrow c \Rightarrow a \rightarrow c$  (transitive)

#### ② Concurrent events ( $a || b$ ) Two distinct event $a \& b$ are said to be concurrent if $a \nrightarrow b$ & $b \nrightarrow a$ .

### Logical clocks

- Procen  $P_i$  has clock  $C_i$ .
- For an event  $a$ ,  $C_i(a)$  returns timestamp of event  $a$ .
- Timestamp monotonically increase.

clock works with following conditions  $\Rightarrow$

(C<sub>1</sub>) For any two events  $a & b$  in a Procen  $P_i$ , if  $a$  occurs before  $b$ , then

$$C_i(a) < C_i(b)$$

(C<sub>2</sub>) If  $a$  is the event of sending a message  $m$  in Procen  $P_i$  &  $b$  is the event of receiving the same message  $m$  at Procen  $P_j$ , then

$$C_j(a) < C_j(b)$$

The following implementation rules (IR) for the ~~clock~~ clock guarantee that the clocks satisfy the correctness condition C<sub>1</sub> & C<sub>2</sub>.

[IR<sub>1</sub>] clock  $C_i$  is incremented between any two successive events in Procen  $P_i$

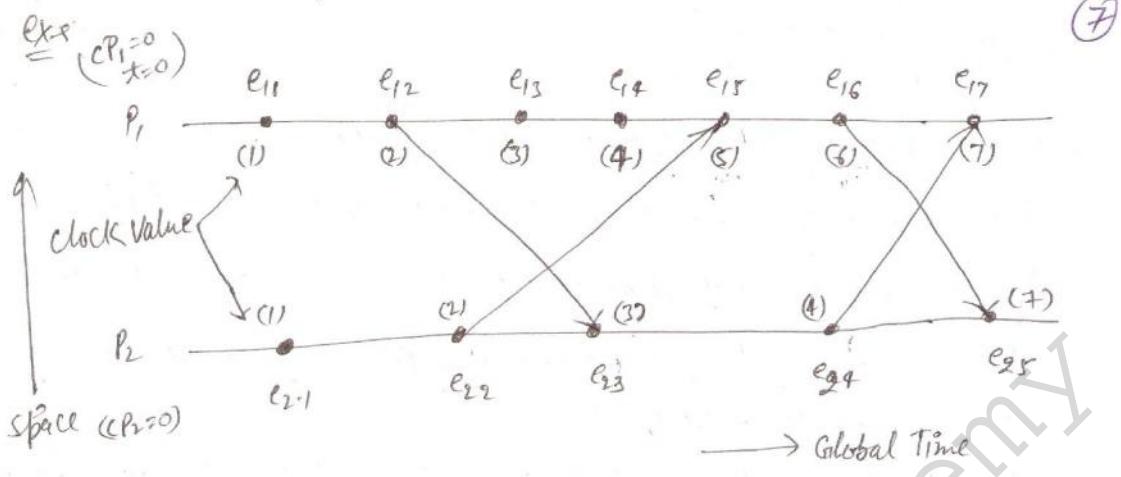
$$C_i^o = C_i + d \quad (d > 0)$$

i.e. If  $a & b$  are two successive events in  $P_i$  &  $a \rightarrow b$  then

$$C_i^o(b) = C_i^o(a) + d$$

[IR<sub>2</sub>] If event  $a$  is the sending of message  $m$  by Procen  $P_i$ , then message  $m$  is assigned timestamp  $t_m = C_i^o(a)$  (by IR<sub>1</sub>). On receiving the same message  $m$  by Procen  $P_j$ ,  $C_j^o$  is set to a value greater than or equal to its present value & greater than  $t_m$ .

$$C_j^o = \max(C_j^o, t_m + d) \quad (d > 0)$$



Let

$$CP_1 = 0 \quad \& \quad d = 1$$

$$\therefore CP_2 = 0$$

Q14  $\Rightarrow \text{Max}[C_{P_1}, t_m + d] \quad CP_2 = 0$   
 $\Rightarrow \text{Max}[1, 1+1]$

Q15  $e_{1.1} \Rightarrow$  is the internal event in Process 1 which causes  $CP_1$  to be incremented to 1 due to  $IR_1$ .

Similarly,  $e_{2.1}$  &  $e_{2.2}$  are two events in  $P_2$  resulting in

$$CP_2 = 2 \text{ due to } IR_1 \rightarrow C_2 = C_2 + d \quad (d > 0)$$

$d = 1$

$$e_{1.2} \Rightarrow \text{Max}(1, 1+1) = 2$$

$$e_{1.3} \Rightarrow \text{Max}(2, 2+1) = 3$$

$$e_{1.4} \Rightarrow \text{Max}(3, 3+1) = 4$$

$e_{1.6} \Rightarrow P_1 \rightarrow$  increment  $CP_1$  to 6 due to  $IR_1$ .

$$\text{Max}((4+1, 6+1))$$

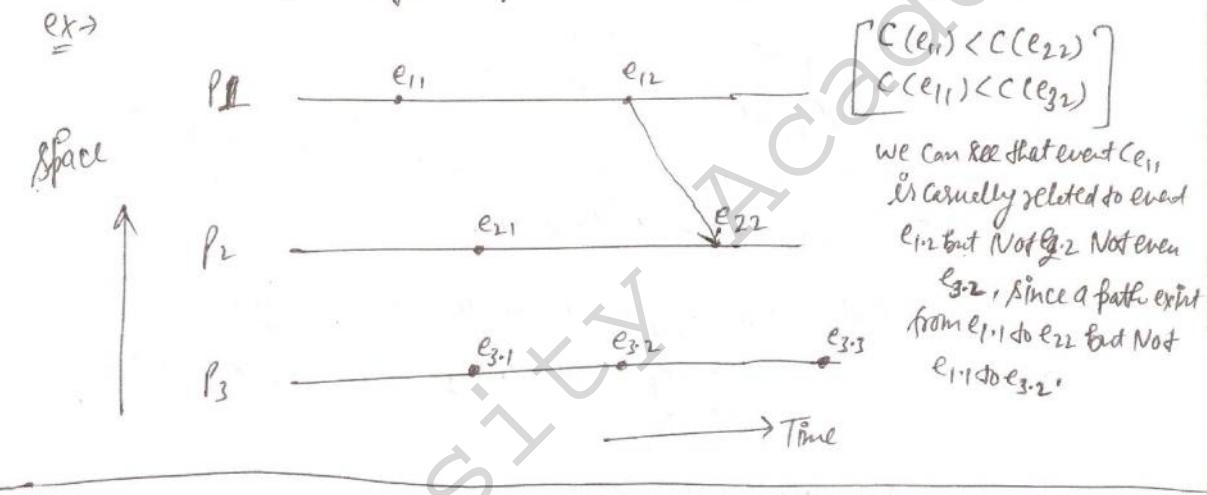
$$C_i = \text{Max}[C_i, t_m + d] \quad (d > 0)$$

$$\left\{ \begin{array}{l} e_{2.2} \Rightarrow \text{Max}(2+1, 2+1) = 3 \\ e_{2.4} \Rightarrow (2+1, 4+1) = 5 \\ e_{2.5} \Rightarrow \text{Max}(4+1, 6+1) = 7 \\ e_{1.7} \Rightarrow \text{Max}(4+1, 6+1) = 7 \end{array} \right.$$

due to  $IR_1$   
 $\underline{\underline{IR_2}}$

## Limitation of Lamport's Clocks

- Lamport system of logical clocks, if  $a \rightarrow b$  then  $C(a) < C(b)$  - However the reverse is not necessarily true if the events have occurred in different processes.
- Thus if  $a \& b$  are event at different processes &  $C(a) < C(b)$ , the  $a \rightarrow b$  is not necessarily true; event  $a \& b$  may be causally related or may not be causally related.
- Thus Lamport's system of clocks is not powerful to capture such situations.

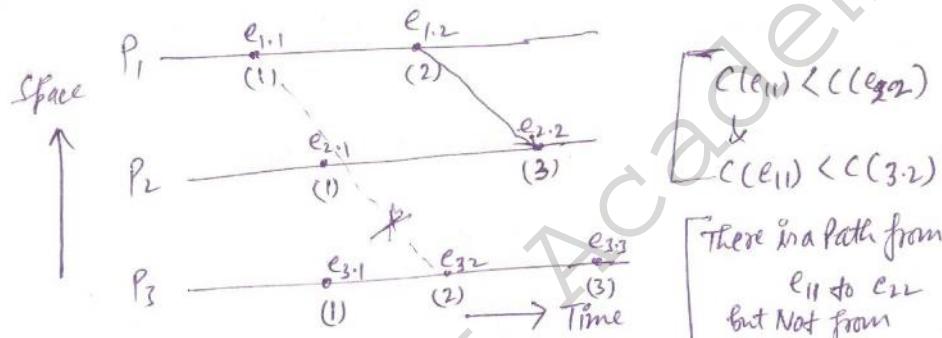


## Limitation of Lamport's Logical clock :-

(8)

- Lamport logical clock, if  $a \rightarrow b$  then  $c(a) < c(b)$ . However the reverse of this is not necessarily true if the events have occurred in different processes.
- That means, if  $a \& b$  are events in different processes &  $c(a) < c(b)$ , then  $a \rightarrow b$  is not necessarily true; events  $a \& b$  may be causally related or may not be causally related.

Ex -



We can't order events in two different process ( $P_1 \& P_2$  &  $P_1 \& P_3$ ) like  $\rightarrow_B$

By just comparing  $c(a) \& c(b)$  ( $c(a) < c(b)$ )

$$c(e_{1.1}) < c(e_{2.2})$$

$$c(e_{1.1}) < c(e_{3.2})$$

There is a path from  
 $e_{1.1}$  to  $e_{2.2}$   
but Not from  
 $e_{1.1}$  to  $e_{3.2}$

- Because each clock ~~can~~ can independently advance due to the occurrence of local events in a process & Lamport's logical clock system cannot distinguish between the advancement of clock due to local events, from those, due to exchange of message between processes.
- using the timestamps assigned by Lamport's clock, we cannot reason about the causal relationships of the events occurring in different processes by just looking the timestamp of the events.

## Vector clock $\rightarrow$

- A ~~concept~~ Let  $n$  be the no. of processes in a distributed system. each process  $P_i$  is equipped with a clock  $C_i$ , which is an integer vector of length  $n$ .  
The clock  $C_i$  can be thought of a function that assigns a vector  $C_i(a)$  to any event  $a$ .  $C_i(a)$  is referred to as the timestamp of event  $a$  at  $P_i$ .  
 $C_i[i]$ , in the  $i$ th entry of  $C_i$ , corresponds to  $P_i$ 's own logical time.  
 $C_i(j)$ ,  $j \neq i$  in  $P_i$ 's best guess of the logical time at  $P_j$ .

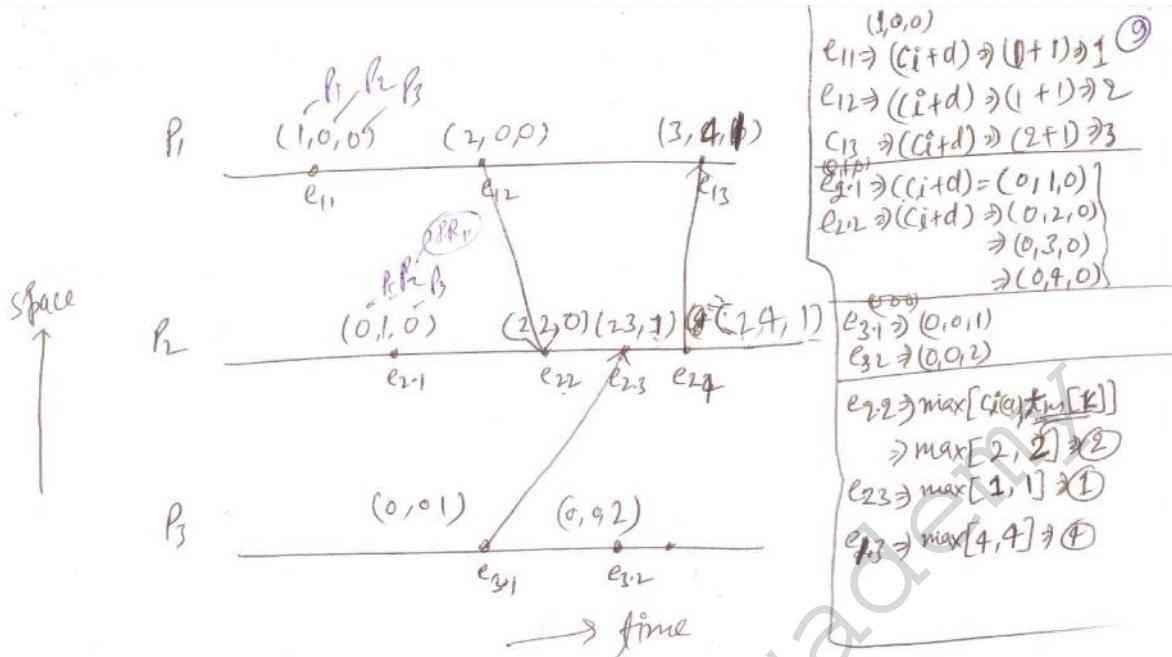
[IR<sub>1</sub>] Clock  $C^i$  is incremented b/w any two successive events in process  $P_i$ -

$$C_i[i] = C_i[i] + d \quad (d > 0) \quad - (I)$$

[IR<sub>2</sub>] If event  $a$  is sending of the message  $m$  by process  $P_i$ , then message  $m$  is assigned a vector timestamp  $t_m = C_i(a)$ ; on receiving the same message  $m$  by process  $P_j$ ,  $C_j$  is updated as follows-

$$\forall k, \quad C_j[k] = \max(C_j[k], t_m[k]) \quad - (II)$$

Note- On the receipt of messages, a process learns about the more recent clock values of the rest of the processes in the system.



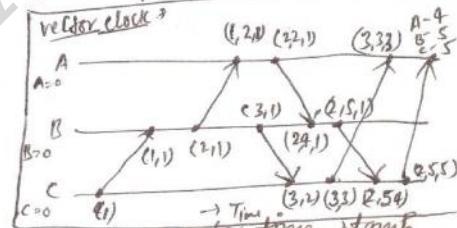
- In Rule IR1, we treat message send & receive by a process as events. In IR2 a message is assigned a timestamp after the sender process has incremented its clock due to IR1.

$$\boxed{A_i, A_j : c_i[i] \geq c_j[i]}$$

- Vector time stamps can be compared as follows. For any two vector time-stamps  $t^a$  &  $t^b$  of events  $a$  &  $b$ , respectively.

- Equal -  $t^a = t^b$  iff  $\forall i; t^a[i] = t^b[i]$ ;
- Not Equal -  $t^a \neq t^b$  iff  $\exists i; t^a[i] \neq t^b[i]$ ;
- Less Than or Equal -  $t^a \leq t^b$  iff  $\forall i; t^a[i] \leq t^b[i]$ ;
- Not Less Than or Equal -  $t^a \not\leq t^b$  iff  $\exists i; t^a[i] \not\leq t^b[i]$ ;
- Less Than -  $t^a < t^b$  iff  $(t^a \leq t^b \wedge t^a \neq t^b)$ ;
- Not less Than -  $t^a \not< t^b$  iff  $\neg(t^a \leq t^b \wedge t^a \neq t^b)$ ;
- Concurrent -  $t^a \parallel t^b$  iff  $(t^a \neq t^b \wedge t^b \neq t^a)$ ;

$\leftarrow$  Partial order  
 $\uparrow$  Not " "  
because  $\neq$  is not transitive.



- $a$  &  $b$  is causally related if  $t^a < t^b$ , even if  $a$  &  $b$  occurred at different processes.

- Then vector clocks allows to order events & decides whether two events are causally related or not by simply looking at the timestamp of the events.

### Causal ordering of messages →

- The causal ordering of message deals with the notation of maintaining the same causal relationship that holds among "message send" events with the corresponding "message receive" events.
- if  $\text{Send}(M_1) \rightarrow \text{Send}(M_2)$  (where  $\text{Send}(m)$  is the event sending message)
   
then every recipient of both message  $M_1$  &  $M_2$  must receive  $M_1$  before  $M_2$ .
- The causal ordering of message should not be confused with the causal ordering of events, which deals with the notation of causal relationships among the events.
- Causal ordering of message is not automatically guaranteed.

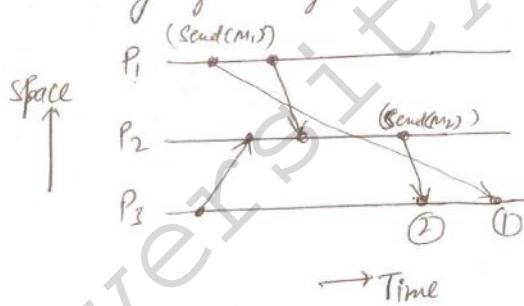
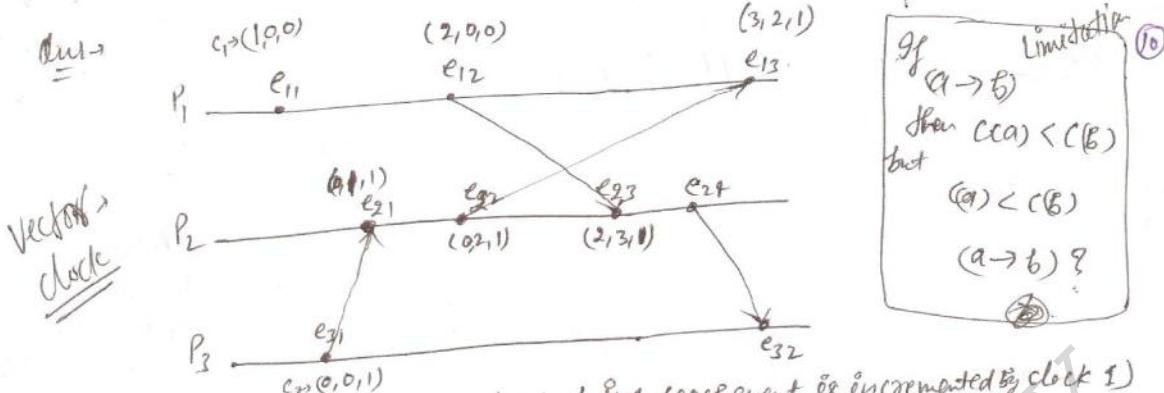


fig- example of violation of causal ordering of message.

Importance of Algorithm → Algorithm for causal ordering of message is important in developing (1) distributed algorithm (2) Implementation of above distributed algorithm.

Eg: for application such as (replicated database), it is important that every process in charge of updating a replicated receives the updates in the same order. Basic idea to maintain the consistency of database.

Alg is designed to deliver a message to a process only if the message immediately preceding it has been delivered to the process. otherwise the message is not delivered immediately but is buffered until the message immediately preceding it is delivered.



10

If  $(a \rightarrow b)$   
then  $C(a) < C(b)$   
but  
 $C(a) < C(b)$   
 $(a \rightarrow b) ?$

- Assume that all clock March at 0 & d by 1. (each event is incremented by clock 1)

At event  $e_{12}$ ,  $C_1(e_{12}) = 2$ , event  $e_{12}$  is sending of message to  $P_2$ .

when  $P_2$  receives the message ( $e_{23}$ ), if clock  $C_2 = 2$ , the clock is reset to 3.

- Event  $e_{24}$  is  $P_2$ 's sending a message to  $P_3$ , that message is received at  $P_3$  at  $e_{32}$ .
- Event  $e_{24}$  is  $P_2$ 's sending a message to  $P_3$ , that message is received at  $P_3$  at  $e_{32}$ .  
 $C_3$  is 1 (as one event happened). By rule TR2,  $C_3$  is reset to the max of  $C_2(e_{24}+1)$  & the current value of  $C_3$ , so  $C_3$  becomes 5.

$$C_3(e_{31}) = 1 < 2 = C_1(e_{12})$$

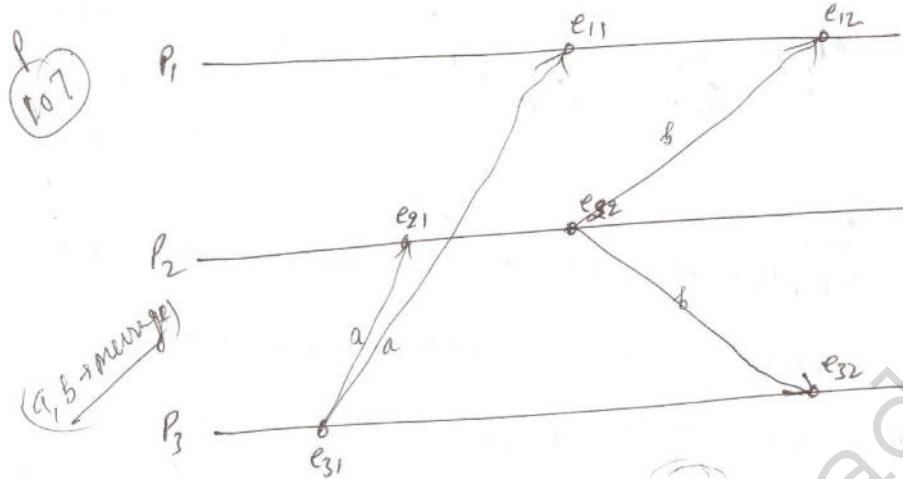
$e_{31}$  &  $e_{21}$  is causally unrelated,  $C_1(e_{21}) < C_3(e_{32})$

$e_{11} \rightarrow e_{32}$  Hence can not say one way to other.

So, The Vector clock, each process keeps track of what it believes the other processes' internal clocks are. The goal is to provide the ordering upon events within the system.

~~BSS~~ BSS →

Birman Schipper-Stephenson Protocols



- $n$  processes
  - $P_i$  is process
  - $c_i$  vector clock  
associated with  $P_i$
  - Jth element is  
 $[c_{ij}^1 c_{ij}^2 \dots c_{ij}^n]$  & contains  
 $P_j$ 's latest value  
for the current time  
in process  $j$
  - $\text{vectorTimeStamp}_m$  for  
message  $m$ .

$\ell_{3,1} \Rightarrow P_3$  send message  $a$ ,  $C_3 = \{0, 0, 1\}$ ;  $D^Q = \{0, 0, 1\}$ .

$e_{2,1} \Rightarrow p_2$  receive Murraya,  $A_8, C_2 = \{0, 0, 0\}, C_2[3] = t^9[3] - 1 = 1 - 1 = 0$  &

$$C_2[3] = \begin{bmatrix} C_2(3) & 0 \\ C_2(1) & 0 \\ C_2(1) & 0 \end{bmatrix} \quad C_2[1] = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad C_2[2] = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

So the message is accepted &  $C_2$  is sent to  $(0, 0, 1)$

$$C_{11} \rightarrow P_1 \text{ receives message } a, A \otimes C_1 = \{0, 0, 0\}, C_1[3] = t^a[3] - 1 \Rightarrow 1 - 1 = 0 \text{ by } t^a[3] = \{0, 0, 1\}$$

$C_1[1] \rightarrow t^a[1]$        $C_2[2] \rightarrow t^a[2]$       So, the Merger is accepted  
 $\& C_1$  is set to  $(0, 0, 1)$

$e_{22} \rightarrow P_2$  sends message  $b$ ,  $C_2 = \{0, 1, 1\}$ ,  $\neq^b = \{0, 1, 1\}$

$e_{12} \rightarrow p_1$  receives message  $t_1$ . As  $C_1 = \{0, 0, 1\}$ ,  $C_1[2] = t_1^6[2] - 1 \Rightarrow 1 - 1 \Rightarrow 0$

$c_1[1] = f^6[1]$  > 0, & the message is accepted  
 $c_1[3] = f^6[2]$  &  $c_1$  is set to (0, 1, 1)

$e_{32} \rightarrow p_3$  receives message 6,  $A \otimes C_3 = \{0, 0, 1\}$ ,  $C_3[2] = t^6[2]-1 = 1-1 = 0$

$c_1[1] = t^k[1] > 0$ , so the Merge is accepted  
 $c_3[2] = t^k[2] > 0$ , & G is set to  $(0, \underline{t}, 1)$

## Chandy-Lamport's Global State Recording Algorithm -

(1)

- It is a proposed model to capture a consistent global state. The algorithm uses a marker (a special message) to initiate the algorithm & the marker has no effect on the underlying computation. The communication channels are assumed to be FIFO.

### Marker sending Rule for a process P →

- Records its state
- For each outgoing channel C from P on which a marker has not been already sent, P sends a marker along C before P sends further message along C.

### Marker Receiving Rule for a process Q →

- On receipt of a marker along a channel C -

If Q has not recorded its state

then

begin

Record the state of C as an empty sequence.

Follow the "Marker sending Rule".

end.

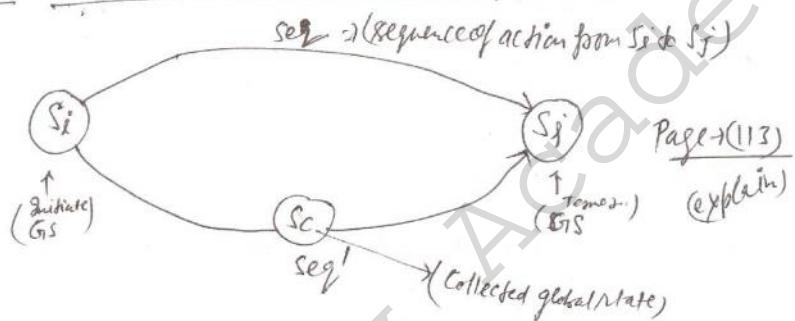
else

Record the state of C as the sequence of message received

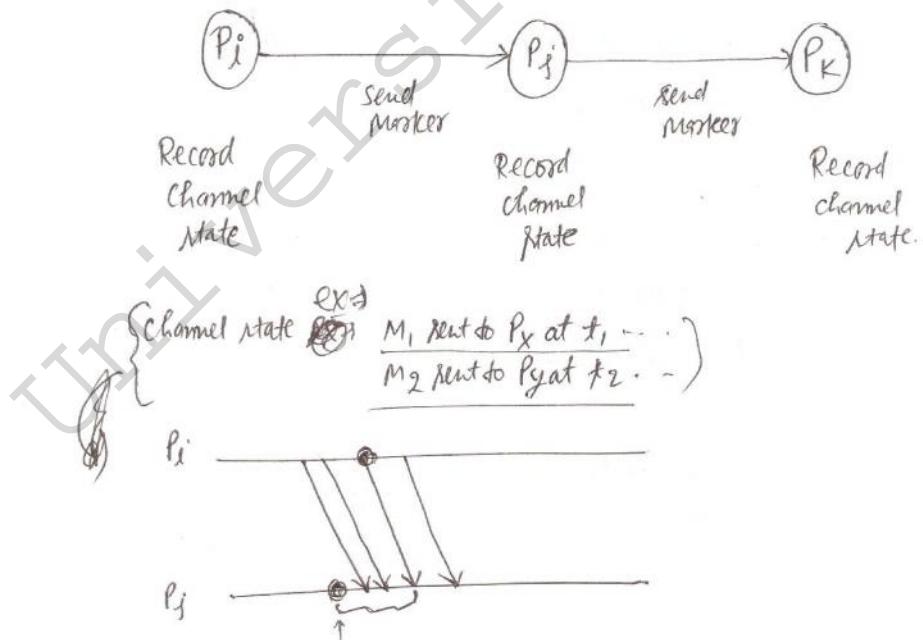
along C after Q's state was recorded & before Q received  
the marker along C.

- FIFO channel condition + markers help in satisfying consistency condition.

- Initiation of marker can be done by any process, with its own unique marker  $\Rightarrow$   $\langle \text{process id}, \text{sequence no} \rangle$ .
- Several processes can initiate state recording by sending markers. Concurrent sending of markers allowed.
- One possible way to collect global state - all initiator of markers
  - all process send the recorded state information to the initiator of marker. Initiator Process can sum up the global state.



ex



## \* Termination Detection :-

(12)

- A Distributed Computation generally consists of a set of cooperating processes which communicate with each other by exchanging messages. In this case of Distributed Computation, it is important to know when the computation has terminated.
- Termination detection, in fact, is an example of usage of the Coherent View (Consistent global state) of distributed system.

System Model → 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 message 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 & there are no messages in ~~transit~~ transit. The message sent by the termination detection algorithm ~~is~~ are known as the control messages.

Basic Idea →

- One of the cooperating processes monitoring the computation is called the controlling agent.
- Initially all processes are idle, the controlling agent's weight equals 1, & the weight of the rest of the process is zero. The computation starts when the controlling agent sends a computation message to one of the process.
- Any time a process sends a message, the process's weight is split between itself & the process receiving the message (the message carries the weight of the receiving process).
- The weight received along with a message is added to the weight of the process.
- Then also assigns a weight  $w$  ( $0 < w \leq 1$ ) to each active process (including agent) & to each message in transit.

- The weight assigned are such that at any time, they satisfy an invariant

$$\sum w = 1$$

- On receiving the computation, a process sends its weight to the controlling agent, which adds the received weight to its own weight.

- When the weight of the controlling agent is once again equal to 1, it concludes that the ~~computation~~ termination has terminated.

Notations -

✓  $B(Dw)$  → Computation message sent as a part of the computation &  $Dw$  is the weight assigned to it.

✓  $C(Dw)$  → Control message sent from the process to the controlling agent &  $Dw$  is the weight assigned to it.

Huang's Termination Detection algo.

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

Derive  $w_1$  &  $w_2$  such that

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

$$w = w_1$$

Send  $B(w_2)$  to  $P$ ;

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

$$w = w + Dw;$$

If  $P$  is idle,  $P$  becomes active.

Rule(3) → An Active process having weight  $w$  may become idle at any time by doing -

Send  $C(w)$  to the controlling agent;

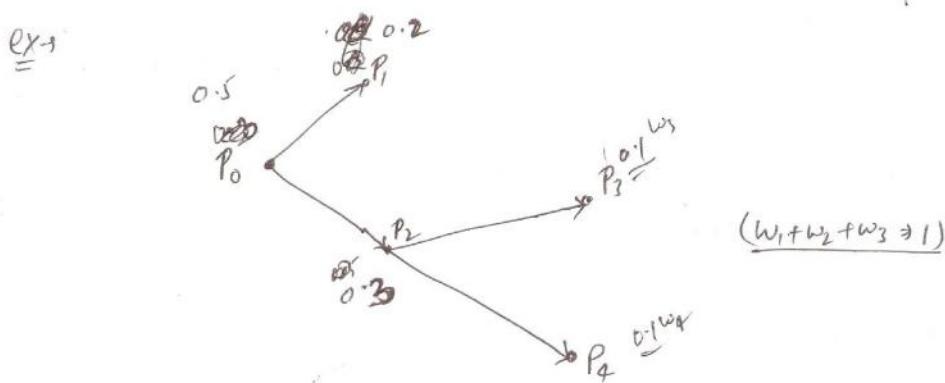
$$[w=0]$$

(The process becomes idle)

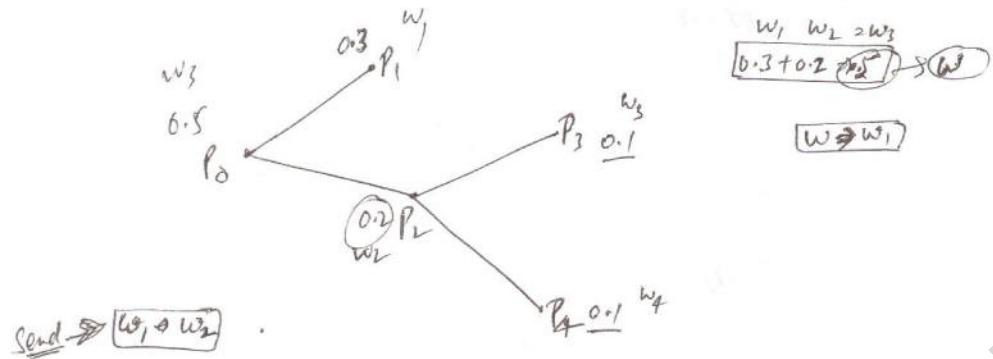
Rule(4) → On receiving  $C(Dw)$ , the controlling agent having weight  $w$  takes the following actions -

$$[w = w + Dw]$$

If  $w = 1$ , conclude that the termination has terminated.



- Process P<sub>0</sub>, designated the controlling agent, with w<sub>0</sub> = 1. It asks P<sub>1</sub> & P<sub>2</sub> to do some computation. It sets w<sub>1</sub> to 0.2 & w<sub>2</sub> to 0.3 & w<sub>3</sub> = 0.5.
  - P<sub>2</sub> in turn asks P<sub>3</sub> & P<sub>4</sub> to do some computations. It sets w<sub>3</sub> to 0.1 & w<sub>4</sub> to 0.1
  - When P<sub>3</sub> terminates, it sends c(w<sub>3</sub>) = c(0.1) to P<sub>2</sub>, which changes w<sub>2</sub> to (0.1 + 0.1) = 0.2.
  - When P<sub>2</sub> terminates, it sends c(w<sub>2</sub>) = c(0.2) to P<sub>0</sub>, which changes w<sub>0</sub> to 0.5 + 0.2 = 0.7
  - When P<sub>4</sub> terminates, it sends c(w<sub>4</sub>) = c(0.1) to P<sub>0</sub>, which changes w<sub>0</sub> to 0.7 + 0.1 = 0.8
  - When P<sub>1</sub> terminates, it sends c(w<sub>1</sub>) = c(0.2) to P<sub>0</sub>, which changes w<sub>0</sub> to 0.8 + 0.2 = 1  
P<sub>0</sub> thereupon concludes the computation is finished.
- Total No of Message Passed  $\rightarrow 8$  (One to start each computation & one to return the weight)



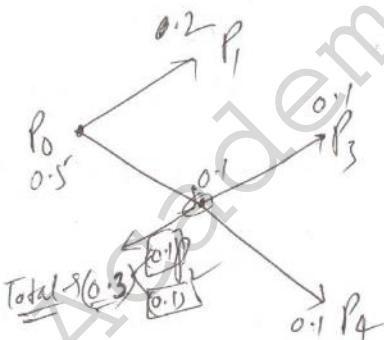
①

$$P_3 \rightarrow (0.4 + 0.1) \Rightarrow 0.2$$

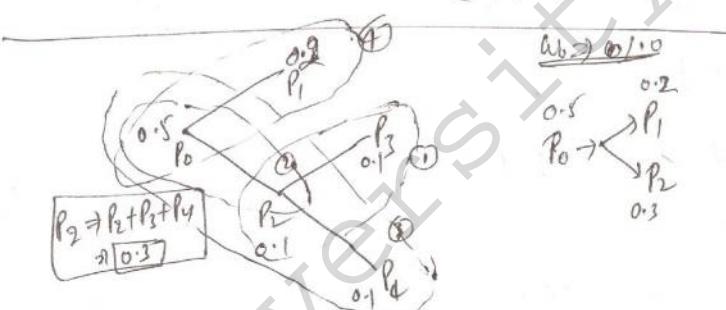
$$P_4 \rightarrow (0.2 + 0.1) \Rightarrow 0.3$$

$$P_0 \rightarrow (0.3 + 0.5) \Rightarrow 0.8$$

$$P_1 \rightarrow 0.8 + 0.3 \Rightarrow 1.1 \quad \times$$

②  $w_1, w_2$ 

$$(w_1 + w_2 = w_3)$$



→ When  $P_3$  terminates it sends -

$C(w) \rightarrow C(0.1)$  to controlling agent  $P_2$  which changes  $w_2$  to  $(0.1) + (0.1) = 0.2$

→ When  $P_2$  terminates it sends -

$$C(w) = C(0.2) \text{ to } " P_0 "$$

$$\text{ " } w_0 \text{ to } (0.5 + 0.2) \Rightarrow 0.7$$

→ When  $P_4$  "

$$C(w) = C(0.1) \text{ to } " P_0 \text{ which}$$

$$\text{ " } w_0 \text{ to } (0.7) + (0.1) \Rightarrow 0.8$$

→ When  $P_1$  "

$$C(w) = C(0.2) \text{ to } " P_0 "$$

$$\text{ " } w_0 \text{ to } (0.8) + (0.2) \Rightarrow 1.0 \text{ Ans}$$

Send B( $w_3$ ) to  $P_1$ 

$$(w \geq w_1) \iff$$

$$P_1 \rightarrow w_0 \Rightarrow 1$$

$$\begin{cases} w_1 \rightarrow 0.2 \\ w_2 \rightarrow 0.3 \\ w_3 \rightarrow 0.5 \end{cases}$$

$$\begin{matrix} w_1 + w_2 + w_3 = & 1.0 \\ 0.2 & 0.3 & 0.5 \end{matrix} = w_0$$

Global State  $\Rightarrow$

$$GS = \{LS_1, LS_2, \dots, LS_n\}$$

$rec(m_{ij}) \Rightarrow LS_i \text{ if } time(rec(m_{ij})) < time(LS_j)$

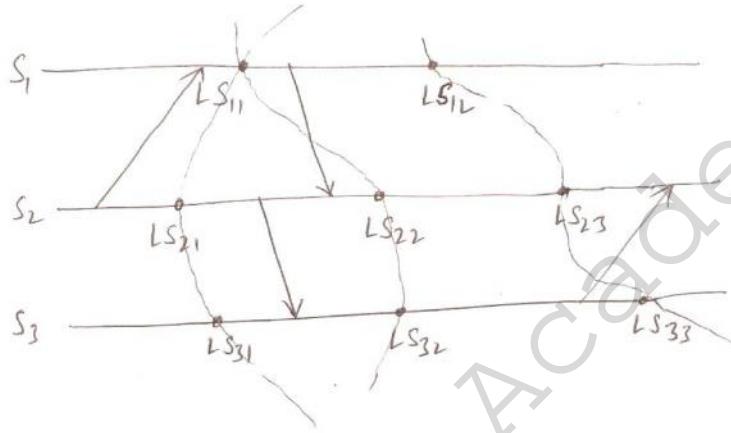
$send(m_{ij}) \Rightarrow LS_i \text{ if } time(send(m_{ij})) < time(LS_i)$

$$\text{transit}(LS_i, LS_j) = \{m_{ij} | send(m_{ij}) \in LS_i \text{ & } rec(m_{ij}) \notin LS_i\}$$

$$\text{inconsistent}(LS_i, LS_j) = \{m_{ij} | send(m_{ij}) \notin LS_i \text{ & } rec(m_{ij}) \in LS_j\}$$

$$\text{consistent}(LS_i, LS_j) = \forall i, \forall j, : 1 \leq i, j \leq n :: \text{inconsistent}(LS_i, LS_j) = \emptyset$$

Ex



①  $\{LS_{12}, LS_{23}, LS_{33}\}$  is consistent. ✓

②  $\{LS_{11}, LS_{22}, LS_{32}\}$  is inconsistent. ✓

③  $\{LS_{11}, LS_{21}, LS_{31}\}$  is strongly consistent. ✓

→ In a consistent GS, for every received message a corresponding send event is recorded in the Global state.

→ In an inconsistent global state, there is at least one message whose receive event is recorded but its send event is not recorded in the Global state.

→ In strongly consistent if it is consistent & transitive.

(Not only the send events of all the recorded received events are recorded, but the received events of all the recorded send events are also recorded.)

SEA Protocol  
Protocol

④ 15

P<sub>i</sub> sends a message to P<sub>j</sub>

① P<sub>i</sub> sends message m, timestamped t<sup>m</sup>, & V<sub>i</sub>, to process P<sub>j</sub>.

② P<sub>i</sub> sets V<sub>i</sub>[j] = t<sup>m</sup>.

P<sub>j</sub> receives a message from P<sub>i</sub>

① When P<sub>j</sub>, i ≠ j, receives m, it delays the message's delivery if both,

(a) V<sup>m</sup>[j] is set; &

(b) V<sup>m</sup>[j] < t<sup>j</sup>

② When the message is delivered to P<sub>j</sub>, update all set elements of V<sub>j</sub> with the corresponding elements of V<sup>m</sup>, except for V<sub>j</sub>[j], as follows.

(a) If V<sub>j</sub>[k] & V<sup>m</sup>[k] are uninitialized, do nothing.

(b) If V<sub>j</sub>[k] is uninitialized & V<sup>m</sup>[k] is initialized, set

$$V_j[k] = V^m[k]$$

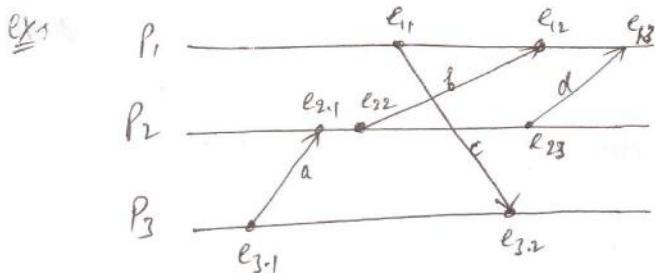
③ If both V<sub>j</sub>[k] & V<sup>m</sup>[k] are initialized, net

$$V_j[EK][EK'] = \max[V_j[EK][EK'], V^m[EK][EK']]$$

for all K' = 1, ..., n.

④ Update P<sub>j</sub>'s vector clock.

⑤ Check buffered messages to see if any one can be delivered.



$e_{31} \rightarrow P_3$  sends message  $a$  to  $P_2$ .

$$\begin{aligned} C_3 &= \{0, 0, 1\} \\ t^a &= \{0, 0, 1\} \\ V^a &= \{\_, \_, \_\} \\ V_3 &= \{\_, (0, 0, 1), \_\} \end{aligned}$$

$e_{21} \rightarrow P_2$  receives message  $a$  from  $P_3$ . As  $V^a[2]$  is uninitialized, the message is accepted.  
 $V_2$  is set to  $\{\_, \_, \_\}$  &  $C_2$  is set to  $\{0, 0, 1\}$ .

$e_{22} \rightarrow P_2$  sends message  $b$  to  $P_1$

$$\begin{aligned} C_2 &= \{0, 1, 1\} \\ t^b &= \{0, 1, 1\} \\ V^b &= \{\_, \_, \_\} \\ V_2 &= \{(0, 1, 1), \_, \_\} \end{aligned}$$

$e_{11} \rightarrow P_1$  sends message  $c$  to  $P_3$ .

$$\begin{aligned} C_1 &= \{1, 0, 0\} \\ t^c &= \{1, 0, 0\} \\ V^c &= \{\_, \_, \_\} \\ V_1 &= \{\text{scratched out}, \_, \_, \{1, 0, 0\}\} \end{aligned}$$

$e_{1.2} \rightarrow P_1$  receives message  $b$  from  $P_2$  -

As,  $V^b[1]$  is uninitialized the message is accepted.  
 $V_1$  is set to  $\{\_, \_, \_\}$   
 $C_1$  is set to  $\{1, 1, 1\}$

E<sub>3,2</sub> → P<sub>3</sub> receives message c from P<sub>1</sub>.

(2)  
16

As,  $V_{[3]}$  is uninitialized the message is accepted.

$V_3$  is set to  $\{?, ?, ?\}$

$V_C = \{1, 0, 1\}$ .

E<sub>2,3</sub> → P<sub>2</sub> sends merged to P<sub>1</sub>

$$\begin{cases} C_2 = \{0, 2, 1\} \\ f^d = \{0, 2, 1\} \\ Vd = \{\{0, 1, 1\}, ?, ?\} \\ V_2 = \{(0, 2, 1), ?, (0, 0, 1)\} \end{cases}$$

E<sub>1,3</sub> → P<sub>1</sub> receives message d from P<sub>2</sub>.

As,  $Vd_{[1]} < C_1_{[1]}$  no message is accepted,

$V_1$  is set to  $\{(0, 1, 1), ?, ?\}$  &  $C_1$  is set to  $\{1, 2, 1\}$

## Schiper - Eggli - Smidz Protocol

(17)

### Datastructure & Notation

- Each process  $P$  maintains a vector denoted by V-P of size  $(N-1)$ , where  $N$  is the No of processes in the System.
- An element of V-P is an ordered pair  $(P', t)$  where  $P'$  is the ID of the destination process of a message &  $t$  is ~~a~~ a vector timestamp.

$t_M$  = logical time at the sending of message  $M$ .

$t_{P_i}$  = Present/current logical time at process  $P_i$ .

### Protocols

#### Sending of Message $M$ from Process $P_1$ to Process $P_2$

- Send message  $M$  (timestamp  $t_M$ ) along with V-P<sub>1</sub> to process  $P_2$ .
- Insert Pair  $(P_2, t_M)$  into V-P<sub>1</sub>. If V-P<sub>1</sub> contains a pair  $(P_2, t)$ , it simply gets over written by new pair  $(P_2, t_M)$ . Note that the pair  $(P_2, t_M)$  was not sent to  $P_2$ . Any future message carrying  $(P_2, t_M)$  pair, can't be delivered to  $P_2$  until  $t_M \leq t_{P_2}$ .

#### Arrival of Message $M$ at Process $P_2$

- If V-M (the vector accompanying message  $M$ ) does not contain any pair  $(P_2, t)$  then the message can be delivered.
- else (\* A pair  $(P_2, t)$  exists in V-M\*)
  - If ( $t \leq t_{P_2}$  then)  
the message can't be delivered (\* it is buffered for later delivery)
  - else  
the message can be delivered.

If message M can be delivered at process  $P_2$ , then the following three actions are taken-

① Merge  $V_M$  accompanying M with  $V_{P_2}$  in the following manner.

\* If  $(\exists (p,t) \in V_M, \text{ such that } p \neq P_2) \& (\forall (p',t) \in V_{P_2}, p' \neq p)$ .  
then insert  $(p,t)$  at  $V_{P_2}$ .

**UNIT-2**Mutual Exclusion

(1)

- In the problem of mutual exclusion, concurrent access to a shared resource by several uncoordinated user requests is serialized to ensure the integrity of the shared resource.
- It requires that the actions performed by a user on a shared memory must be atomic. That is if several user concurrently access a shared resource then action performed by a user, as far as the other users are concerned, must be instantaneous & ~~indivisible~~ indivisible such that the net effect on the shared resources is the same as if user actions were executed serially, as opposed to in an interleaved manner.
- The Problem of ME frequently arises in distributed systems whenever concurrent access to shared resources by several sites is involved. For correctness ~~access to shared resources by several~~ it is necessary that the shared resources be accessed by a single site (or process) at a time.  
 Ex → Directory Mgt, where an update to a directory must be done automatically because if update & reads to a directory proceed to a directory must be done automatically because concurrently, reads may obtain inconsistent information. If an entry contains several fields, a read opn<sup>n</sup> may read some fields before the update & some after the update.

Mutual exclusion in Single Computer System Vs Distributed Systems

- The Problem of ME in Single-Computer System, where Shared memory exists.
- In Single-computer systems, the status of a shared resource & the status of work is readily available in the shared memory, & so the mutual exclusion problem can be easily implemented using shared variable (e.g. semaphores).

- However in DS, Both the shared resources & the writes may be distributed & Shared Memory does not exist. (Approach based on Shared variable are no applicable to DSs & approach based on Message Passing Must be used).
- The Problem of ME ~~becomes~~ becomes ~~more~~ much more complex in DS (as compared to single-computer systems) because of the lack of Both shared memory & a common physical clock & because of unpredictable Message delay.

### The classification of ME Algorithms

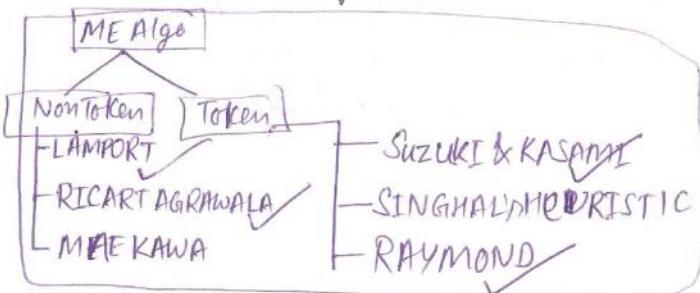
- They tend to differ in the communication topology (e.g tree, ring, etc.) & in the amount of info maintained by each site about other sites.
- These algorithms can be grouped into 2 classes.

~~(uses Token~~ *(uses message to arrive at ME)*

① Non Token Based - These algo require two or more successive rounds of message exchange among the sites. These algorithms are assertion Based because a site can enter its ~~critical~~ critical section(s) when an assertion defined on its local variables becomes true. ME enforced because the assertion becomes true only at one site at any given time.

*(uses Token*   
 *privileged*   
 *message to*   
 *arrive at*   
 *ME)*

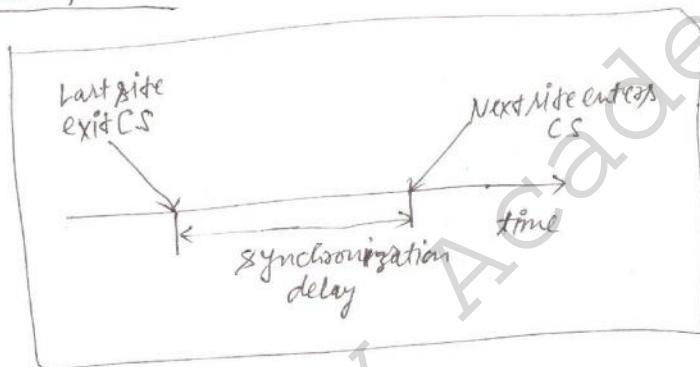
② Token based - In these algorithms, a unique token (also known as the Privilege message) is shared among sites. A site is allowed to enter in CS if it possesses the token & it continues to hold the token until the execution of the CS is over.



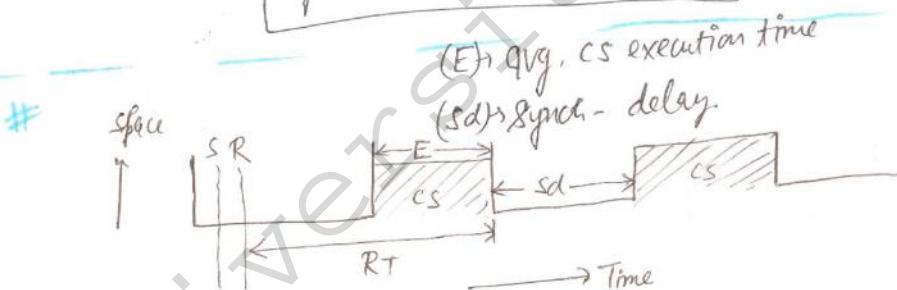
## Requirements of ME Algo's

- At most one process in critical section (Safety) (fault tolerance)
- if more than one requesting process, someone enters (liveness) (freedom from deadlock)
- a requesting process enters within a finite time (No starvation)
- requests are granted in order (fairness)

## How to measure the performance



$$\text{System throughput} = 1/(sd + E)$$



S → CS request arrives.

R → CS request message is sent out.

sd → Synchronous delay (time b/w two cs)

Performance matrix for ME algo can be easily depicted by time diagrams -

- ① Response time (RT) → It is the time b/w R & end cs. (time b/w reqt message is sent out & completion of critical section. For small RT the performance of ME algo is high.)
- ② Synchronous delay (sd) → It is time b/w two consecutive CS. That is time b/w END of CS & BEGIN of CS. In this period message are exchanged to arrive at ME decision.

Note - for small 'sd' performance of ME algo will be high.

### ③ No of Message Per CS →

- B/w two cs there are message exchange as discussed above. As no of message exchange reduces, the performance will improve.

(Note No of Message per cs must be less to improve the performance)

### ④ System throughput → It is derived from 'Sd' & 'E' as

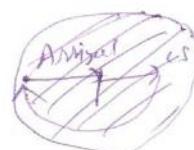
$$\text{System throughput} = \frac{1}{Sd+E}$$

- (Sd+E) time unit only one cs request is executing, so throughput is  $\frac{1}{(Sd+E)}$ . This can also be called as average time for each cs invocation or time at which cs request complete.

## # SIMPLE Sol<sup>n</sup> To Distributed ME →

- A site, called the control site, is assigned the task of granting permission for the cs execution. To request the cs, a site sends a Request message to the control site.
- ~~To request~~ The control site queues up the requests for the cs & grants them permission, one by one.
- This method to achieve mutual exclusion in distributed systems requires only three message per cs execution.

Central site has several drawbacks -

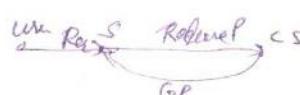


① Single Point of failure, the control site.

② The control site is likely to be swamped (flooded) with extra work.

③ Synchronous delay of algorithm  $\overset{\text{is}}{2T}$ , because a site should first release permission to the control site & then control site should grant permission to the next site to execute cs.

Thought →  $\frac{1}{(2T+E)}$  (If we reduce the synchronous delay the throughput will be max.)



(3)

## NON TOKEN BASED ALGORITHMS →

- In non-token based ME algo, a site communicates with a set of other sites to arbitrate who should execute the CS next.
- For a site  $s_i$ , request set  $R_i$  contains ids of all those sites from which the site  $s_i$  must acquire permission before entering the CS.

## LAMPORT'S ALGORITHM →

- LAMPORT was first to give a distributed mutual exclusion algorithm as an illustration of his clock synchronization scheme.
- In LAMPORT's algorithm,  $\forall i : 1 \leq i \leq N : R_i = \{s_1, s_2, s_3, \dots, s_N\}$ .
- Every site  $s_i$  keeps a queue, request\\_queue<sub>i</sub>, which contains mutual exclusion requests ordered by their timestamps.
- This Algorithm requires messages to be delivered in the FIFO order between every pair of sites.

Algo-

### Requesting the CS

- ① When a site  $s_i$  wants to enter the CS, it sends a Request( $t_{si}, i$ ) message to all the sites in its request set  $R_i$  & places the request on request\\_queue<sub>i</sub>:  
 $(t_{si}, i) \rightarrow$  timestamp of the request.
- ② When a site  $s_j$  receives the Request( $t_{si}, i$ ) message from site  $s_i$ , it returns a timed REPLY message to  $s_i$ , & places site  $s_i$ 's request on request\\_queue<sub>j</sub>.

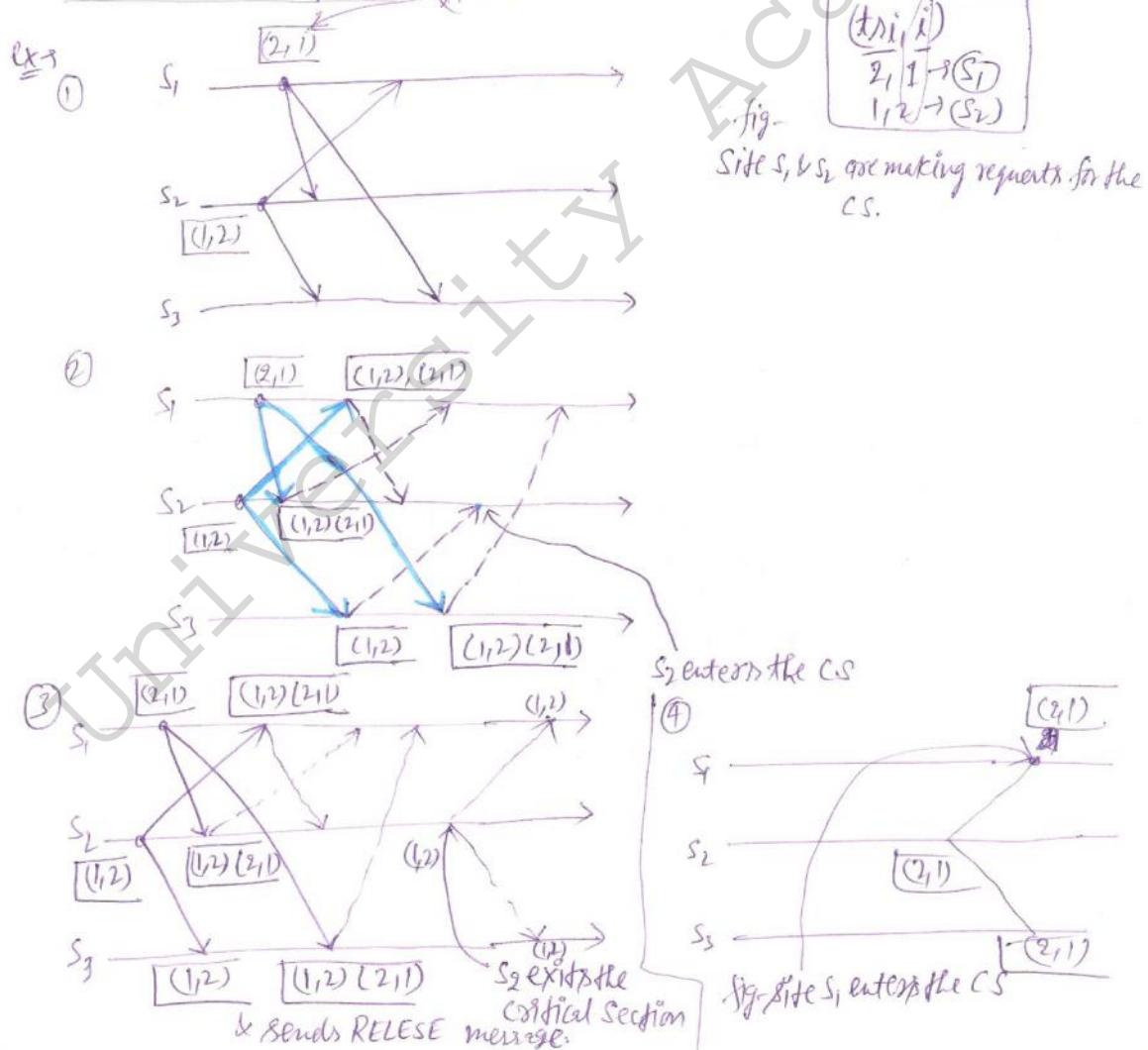
## Executing the CS

- Site  $S_i$  enters the CS when the two following conditions hold →

- [L<sub>i</sub>] -  $S_i$  has received a message with timestamp larger than  $(t_{SE}, i)$  from all other sites.
- [R<sub>i</sub>] -  $S_i$ 's request is at the top of request queue.

## Releasing the CS →

- ③ Site  $S_i$ , upon exiting the CS, removes its request from the top of its request queue & sends a timestamped RELEASE message to all the sites in its request set.
- ④ When a site  $S_j$  receives a Release message from site  $S_i$ , it removes  $S_i$ 's request from its request queue. (Timestamp of  $S_i$ )



## The Ricart-Agrawal Algorithm

(4)

- Ricart Agrawal is an optimization of Lamport's algorithm that dispenses with RELEASE messages by cleverly merging them with REPLY messages.
- Here,  $\forall i : 1 \leq i \leq N :: R_i = \{s_1, s_2, \dots, s_N\}$

### Algo

#### Requesting the CS

- ① When a site  $s_i$  wants to enter the CS, it sends a timestamped REQUEST message to all the sites in its request set.
- ② When site  $s_j$  receives a REQUEST message from site  $s_i$ , it sends a REPLY message to site  $s_i$ . If site  $s_j$  is neither requesting nor executing the CS or if site  $s_j$  is requesting &  $s_i$ 's request's timestamp is smaller than site  $s_j$ 's own request's timestamp, the request is deferred otherwise →  
(put off to a later time)

#### Executing the Critical Section

- ③ Site  $s_i$  enters the CS after it has received REPLY message from all the sites in its request set.

#### Releasing the critical section

- ④ When site  $s_i$  exits the CS, it sends REPLY message to all the deferred requests.

Ex:

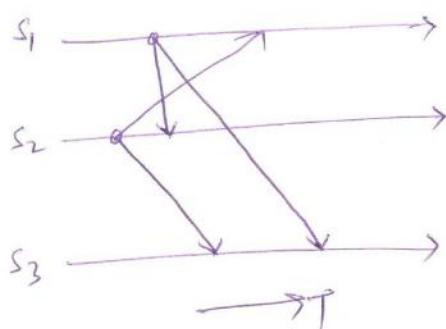
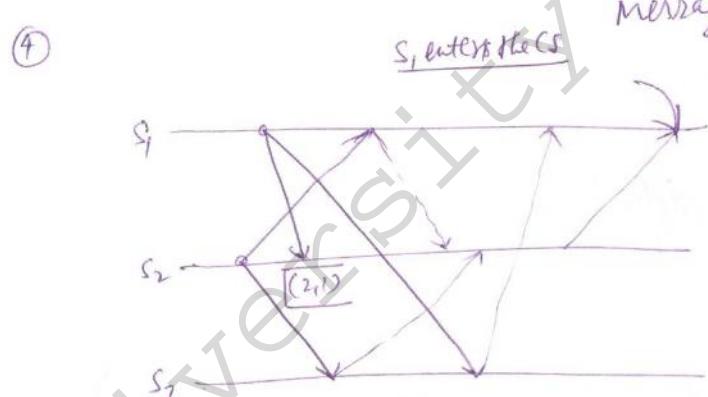
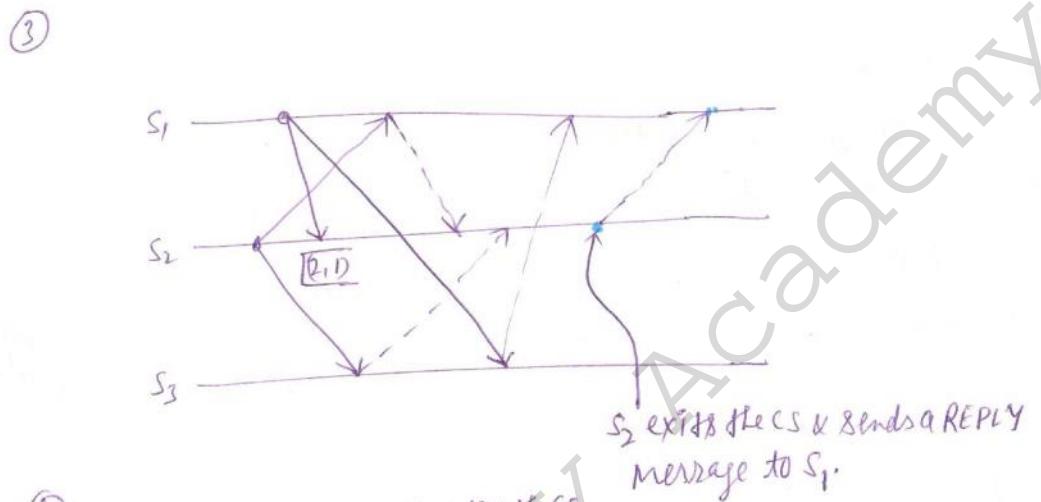
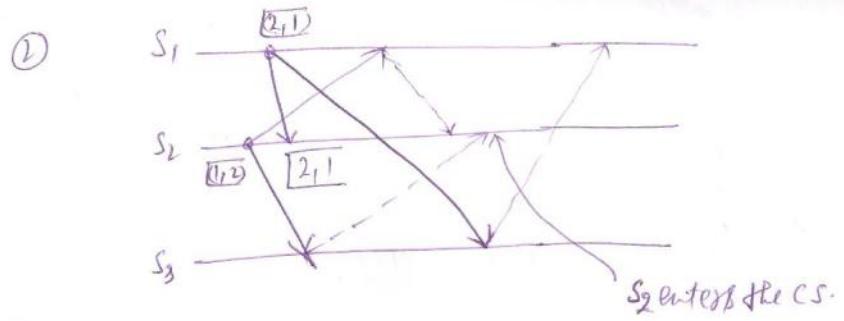


fig -  $s_1$  &  $s_2$  are making reg. for the CS.



## Token Based Algo

(5)

- A unique token is shared among all the sites. A site is allowed to enter its CS if it possesses the token. Depending upon the way a site carries out its search for the token, there are numerous token based algo.

① Token-Based algo use sequence no instead of timestamps.

(a site increments its seq no every time it makes a reqt for the token)

② (A Primary form of the sequence no is to distinguish bet<sup>n</sup> old & current request.)

③ → (37) correctness Proof of token-based algo.

## SUZUKI-KASAMI'S BROADCAST ALGORITHM

① <sup>Reqd CS</sup> If the requesting site  $s_i$  does not have the token, then it increments its sequence no,  $RN_i[i]$ , & sends a REQUEST( $i, s_m$ ) message to all other sites.  
 (  $s_n$  is the updated value of  $RN_i[i]$  )

② When a site  $s_j$  receives this message, it sets  $RN_j[i]$  to  $\max(RN_j[i], s_n)$   
 If  $s_j$  has the idle token, then it sends the token to  $s_i$ , if

$$RN_j[i] = LN[i] + 1$$

$LN[i] \rightarrow$  The sequence No of the request that sites  $\Rightarrow$  executed most recently.

Executing the CS:

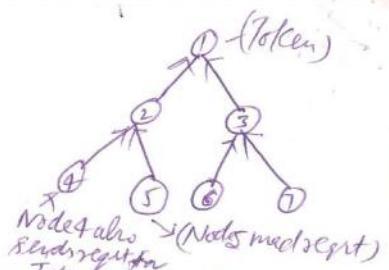
③ Site  $s_i$  executes the CS when it received the token.

Releasing the CS: After finishing the CS, site  $s_i$  makes the following actions-

④ It sets  $LN[i]$  element of the token array equals to  $RN_i[i]$ .

⑤ For every site  $s_j$  whose ID is not in the token queue it appends its ID to the token Queue if  $[RN_i[j] = LN[j] + 1]$

⑥ If token queue is nonempty after the above update, then it deletes the top site ID from the queue & sends the token to the site indicated by the ID.



Suzuki-Karumi exs

Step 1 :  $S_1$  has token,  $S_3$  is in queue  $\rightarrow$

| <u>Side</u> | <u>Seg Vector RN</u> | <u>Token Vector</u> | <u>Token Queue</u> |
|-------------|----------------------|---------------------|--------------------|
| $S_1$       | 10, 15, 9            | 10, 15, 8           | 3                  |
| $S_2$       | 10, 16, 9            | 10, 16, 8           |                    |
| $S_3$       | 10, 15, 9            |                     |                    |

Step 2 :  $S_3$  gets token,  $S_2$  in Queue  $\rightarrow$

|       |           |           |   |
|-------|-----------|-----------|---|
| $S_1$ | 10, 16, 9 |           |   |
| $S_2$ | 10, 16, 9 |           |   |
| $S_3$ | 10, 16, 9 | 10, 15, 9 | 2 |

Step 3 :  $S_3$  gets token (queue empty)

|       |           |           |         |
|-------|-----------|-----------|---------|
| $S_1$ | 10, 16, 9 |           |         |
| $S_2$ | 10, 16, 9 | 10, 16, 9 | (empty) |
| $S_3$ | 10, 16, 9 |           |         |

Tokens

- Queue (FIFO)  $\alpha$  of requesting Process.
- $LN[1 \dots n]$  :  $\rightarrow$  Sequence No of request that f executed most recently

The requested Message

- REQUEST ( $i, k$ ) - request message from node  $i$  for its  $K^{th}$  critical section.

Other Datastructures

- $RN_i[1 \dots n]$  for each node  $i$ , where  $RN_i[s]$  is the largest sequence no received so far by  $i$  in a Req<sup>nt</sup> message from  $f$ .

To Request CS  $\rightarrow$ 

$\rightarrow$  If  $i$  does not have token, increment  $RN_i[s]$  & send REQUEST

REQUEST ( $i, RN_i[s]$ ) to all nodes.

$\rightarrow$  If  $\boxed{i}$  has token  $i$  already, enter CS if the token is idle  
(No pending reqst), else follow rule to release CS.

~~Ques~~

on Receiving REQUEST ( $i, sn$ ) at  $f \Rightarrow$ 

- Set  $RN_f[i] = \max[RN_f[i], sn]$

$\rightarrow$  If  $f$  has the token & the token is idle, then send it to  $i$  ~~& f~~ if  
 $RN_f[i] = LN[i] + 1$ . If token is not idle, follow rule to

Release CS.

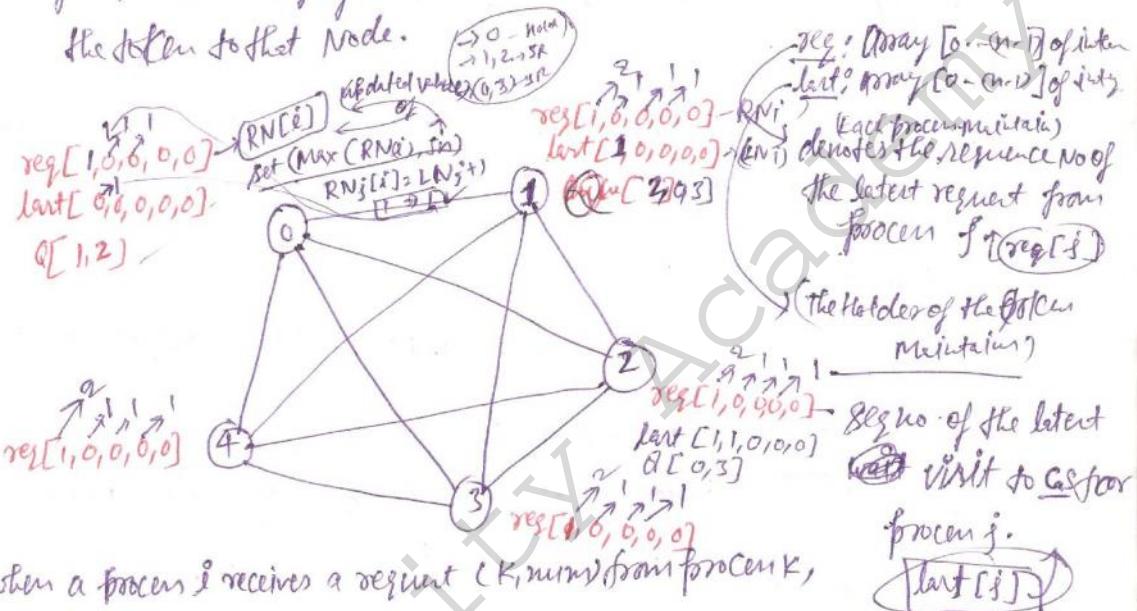
To enter the CS

- Enter the CS if token is present.

### To Relace the CS

- Set  $LN[i] = RN[i]$
- For every node  $j$  which is not in  $Q$  (in token), add node  $j$  to  $Q$  if ~~RN~~

$$[RN[j] = LN[j]+1]$$
- If  $Q$  is non empty after the above, delete first node from  $Q$  & send the token to that Node.



- When a process  $i$  receives a request ( $RN[i]$ ) from process  $k$ , it sets  $req[k]$  to  $\max(req[k], num)$ .
- $\rightarrow$  Queue of waiting process,  $token(Q, lant)$

### The Holder of the Token

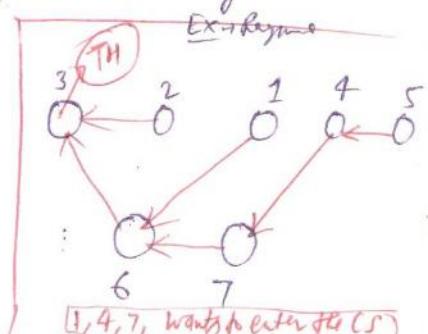
- complete CS
- sets  $lant[i] :=$  its own num
- update  $Q$  by retaining each process  $k$  only if -

$$1 + lant[k] \geq req[k] \rightarrow (\text{guarantees the fairness of requests})$$

- Sends the token to the head of  $Q$ , along with the array  $lant$  & the tail  $Q$ .

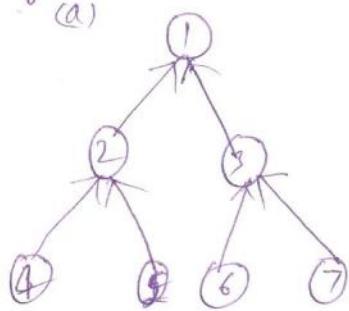
$$token = (Q, lant)$$

Initially -  $req[i] = \infty$ ,  $lant[i] = 0$

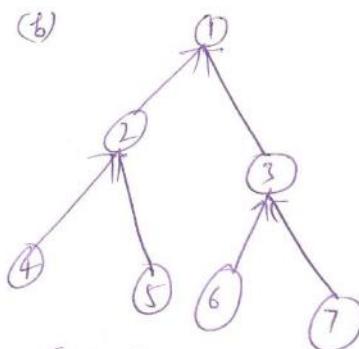


### Raymond's Tree Based Algo →

(7) (8)

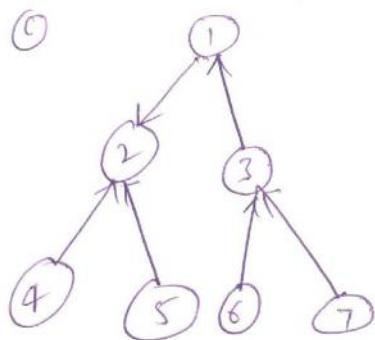


- ① Token is at Node 1  
node 5 made a request

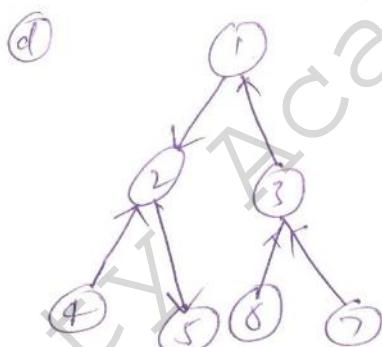


- ② Node 2 receives the Reqst,  
it sends the request to Node 1

- ③ Node 4 also sends a request, node 2 receives it.

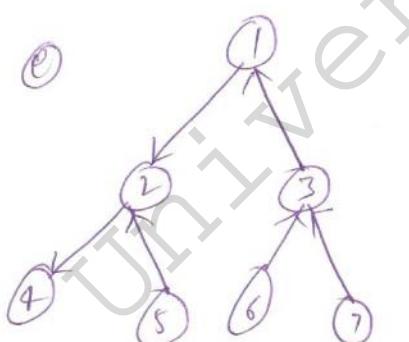


- ④ Token is at node 2 now  
node 2 becomes the root

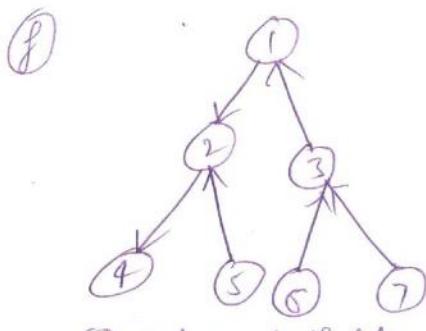


- ⑤ node 5 gets the token, it enters CS

- ⑥ node 2 sends a request to node 5

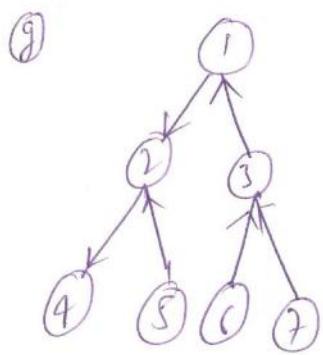


- ⑦ Node 5 sends the token to node 2

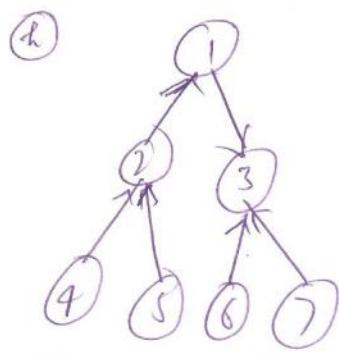


- ⑧ Node 4 gets the token, it enters CS

- ⑨ Node 3 sends a request



⑩ the resp from node 2  
comes to node 4



⑪ Node 3 gets the token, & becomes  
the root.

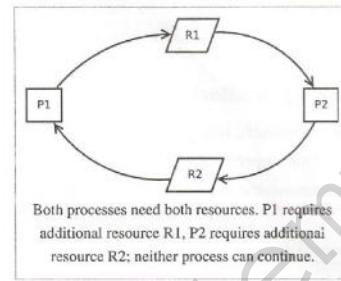
## Deadlock

A **deadlock** is a situation in which two or more competing actions are each waiting for the other to finish, and thus neither ever does.

In an operating system, a deadlock is a situation which occurs when a process enters a waiting state because a resource requested by it is being held by another waiting process, which in turn is waiting for another resource. If a process is unable to change its state indefinitely because the resources requested by it are being used by another waiting process, then the system is said to be in a deadlock.<sup>[1]</sup>

Deadlock is a common problem in multiprocessing systems, parallel computing and distributed systems, where software and hardware locks are used to handle shared resources and implement process synchronization.<sup>[2]</sup>

In telecommunication systems, deadlocks occur mainly due to lost or corrupt signals instead of resource contention.<sup>[3]</sup>



### Examples

Deadlock situation can be compared to the classic "chicken or egg" problem.<sup>[4]</sup> It can be also considered a paradoxical "Catch-22" situation.<sup>[5]</sup> A real world analogical example would be an illogical statute passed by the Kansas legislature in the early 20th century, which stated:<sup>[1][6]</sup>

"When two trains approach each other at a crossing, both shall come to a full stop and neither shall start up again until the other has gone."

A simple computer-based example is as follows. Suppose a computer has three CD drives and three processes. Each of the three processes holds one of the drives. If each process now requests another drive, the three processes will be in a deadlock. Each process will be waiting for the "CD drive released" event, which can be only caused by one of the other waiting processes. Thus, it results in a circular chain.

### Necessary conditions

A deadlock situation can arise if and only if all of the following conditions hold simultaneously in a system:<sup>[1]</sup>

1. **Mutual Exclusion:** At least one resource must be non-shareable.<sup>[1]</sup> Only one process can use the resource at any given instant of time.
2. **Hold and Wait or Resource Holding:** A process is currently holding at least one resource and requesting additional resources which are being held by other processes.
3. **No Preemption:** The operating system must not de-allocate resources once they have been allocated; they must be released by the holding process voluntarily.
4. **Circular Wait:** A process must be waiting for a resource which is being held by another process, which in turn is waiting for the first process to release the resource. In general, there is a set of waiting processes,  $P = \{P_1, P_2, \dots, P_N\}$ , such that  $P_1$  is waiting for a resource held by  $P_2$ ,  $P_2$  is waiting for a resource held by  $P_3$  and so on till  $P_N$  is waiting for a resource held by  $P_1$ .<sup>[1][7]</sup>

These four conditions are known as the **Coffman conditions** from their first description in a 1971 article by Edward G. Coffman, Jr.<sup>[7]</sup> Unfulfillment of any of these conditions is enough to preclude a deadlock from occurring.

|  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|
|  |  |  |  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|--|--|--|

## Deadlock handling

Most current operating systems cannot prevent a deadlock from occurring.<sup>[1]</sup> When a deadlock occurs, different operating systems respond to them in different non-standard manners. Most approaches work by preventing one of the four Coffman conditions from occurring, especially the fourth one.<sup>[8]</sup> Major approaches are as follows.

### Ignoring deadlock

In this approach, it is assumed that a deadlock will never occur. This is also an application of the Ostrich algorithm.<sup>[8][9]</sup> This approach was initially used by MINIX and UNIX.<sup>[7]</sup> This is used when the time intervals between occurrences of deadlocks are large and the data loss incurred each time is tolerable.

### ① Detection

Under deadlock detection, deadlocks are allowed to occur. Then the state of the system is examined to detect that a deadlock has occurred and subsequently it is corrected. An algorithm is employed that tracks resource allocation and process states, it rolls back and restarts one or more of the processes in order to remove the detected deadlock. Detecting a deadlock that has already occurred is easily possible since the resources that each process has locked and/or currently requested are known to the resource scheduler of the operating system.<sup>[9]</sup>

Deadlock detection techniques include, but are not limited to, *model checking*. This approach constructs a finite state-model on which it performs a progress analysis and finds all possible terminal sets in the model. These then each represent a deadlock.

After a deadlock is detected, it can be corrected by using one of the following methods:

1. **Process Termination:** One or more process involved in the deadlock may be aborted. We can choose to abort all processes involved in the deadlock. This ensures that deadlock is resolved with certainty and speed. But the expense is high as partial computations will be lost. Or, we can choose to abort one process at a time until the deadlock is resolved. This approach has high overheads because after each abortion an algorithm must determine whether the system is still in deadlock. Several factors must be considered while choosing a candidate for termination, such as priority and age of the process.
2. **Resource Preemption:** Resources allocated to various processes may be successively preempted and allocated to other processes until the deadlock is broken.

### ② Prevention

Deadlock prevention works by preventing one of the four Coffman conditions from occurring.

- Removing the **mutual exclusion** condition means that no process will have exclusive access to a resource. This proves impossible for resources that cannot be spooled. But even with spooled resources, deadlock could still occur. Algorithms that avoid mutual exclusion are called *non-blocking synchronization algorithms*.
- The **hold and wait** or **resource holding** conditions may be removed by requiring processes to request all the resources they will need before starting up (or before embarking upon a particular set of operations). This advance knowledge is frequently difficult to satisfy and, in any case, is an inefficient use of resources. Another way is to require processes to request resources only when it has none. Thus, first they must release all their currently held resources before requesting all the resources they will need from scratch. This too is often impractical. It is so because resources may be allocated and remain unused for long periods. Also, a process requiring a popular resource may have to wait indefinitely, as such a resource may always be allocated to some process, resulting in resource starvation.<sup>[1]</sup> (These algorithms, such as serializing tokens, are known as the *all-or-none algorithms*.)
- The **no preemption** condition may also be difficult or impossible to avoid as a process has to be able to have a resource for a certain amount of time, or the processing outcome may be inconsistent or thrashing may occur. However, inability to enforce preemption may interfere with a *priority* algorithm. Preemption of a "locked out"

(2)

resource generally implies a rollback, and is to be avoided, since it is very costly in overhead. Algorithms that allow preemption include lock-free and wait-free algorithms and optimistic concurrency control.

The final condition is the circular wait condition. Approaches that avoid circular waits include disabling interrupts during critical sections and using a hierarchy to determine a partial ordering of resources. If no obvious hierarchy exists, even the memory address of resources has been used to determine ordering and resources are requested in the increasing order of the enumeration.<sup>[1]</sup> Dijkstra's solution can also be used.

### Avoidance *(Safe/unsafe)*

Deadlock can be avoided if certain information about processes are available to the operating system before allocation of resources, such as which resources a process will consume in its lifetime. For every resource request, the system sees whether granting the request will mean that the system will enter an unsafe state, meaning a state that could result in deadlock. The system then only grants requests that will lead to safe states.<sup>[1]</sup> In order for the system to be able to determine whether the next state will be safe or unsafe, it must know in advance at any time:

- ✓ resources currently available
- ✓ resources currently allocated to each process
- ✓ resources that will be required and released by these processes in the future

It is possible for a process to be in an unsafe state but for this not to result in a deadlock. The notion of safe/unsafe states only refers to the ability of the system to enter a deadlock state or not. For example, if a process requests A which would result in an unsafe state, but releases B which would prevent circular wait, then the state is unsafe but the system is not in deadlock.

\* (1) One known algorithm that is used for deadlock avoidance is the Banker's algorithm, which requires resource usage limit to be known in advance.<sup>[1]</sup> However, for many systems it is impossible to know in advance what every process will request. This means that deadlock avoidance is often impossible.

\* (2) Two other algorithms are Wait/Die and Wound/Wait, each of which uses a symmetry-breaking technique. In both these algorithms there exists an older process (O) and a younger process (Y). Process age can be determined by a timestamp at process creation time. Smaller timestamps are older processes, while larger timestamps represent younger processes.

|                              | Wait/Die             | Wound/Wait           |
|------------------------------|----------------------|----------------------|
| O needs a resource held by Y | O waits              | Y dies $\frac{+}{-}$ |
| Y needs a resource held by O | Y dies $\frac{-}{+}$ | Y waits              |

### Livelock

A livelock is similar to a deadlock, except that the states of the processes involved in the livelock constantly change with regard to one another, none progressing.<sup>[10]</sup> Livelock is a special case of resource starvation; the general definition only states that a specific process is not progressing.<sup>[11]</sup>

A real-world example of livelock occurs when two people meet in a narrow corridor, and each tries to be polite by moving aside to let the other pass, but they end up swaying from side to side without making any progress because they both repeatedly move the same way at the same time.

Livelock is a risk with some algorithms that detect and recover from deadlock. If more than one process takes action, the deadlock detection algorithm can be repeatedly triggered. This can be avoided by ensuring that only one process (chosen randomly or by priority) takes action.<sup>[12]</sup>

|  |  |  |  |  |  |  |  |  |                     |  |  |
|--|--|--|--|--|--|--|--|--|---------------------|--|--|
|  |  |  |  |  |  |  |  |  | REMARKS /<br>OTHERS |  |  |
|--|--|--|--|--|--|--|--|--|---------------------|--|--|

### Distributed deadlock

Distributed deadlocks can occur in distributed systems when distributed transactions or concurrency control is being used. Distributed deadlocks can be detected either by constructing a global wait-for graph, from local wait-for graphs at a deadlock detector or by a distributed algorithm like edge chasing.

In a commitment ordering-based distributed environment (including the strong strict two-phase locking (SS2PL, or rigorous) special case) distributed deadlocks are resolved automatically by the atomic commitment protocol (like a two-phase commit (2PC)), and no global wait-for graph or other resolution mechanism is needed. Similar automatic global deadlock resolution occurs also in environments that employ 2PL that is not SS2PL (and typically not CO; see *Deadlocks in 2PL*). However, 2PL that is not SS2PL is rarely utilized in practice.

**Phantom deadlocks** are deadlocks that are detected in a distributed system due to system internal delays but no longer actually exist at the time of detection.

### References

- [1] Silberschatz, Abraham (2006). *Operating System Principles* ([http://books.google.co.in/books?id=WjvX0HmVTIMC&dq=deadlock+operating+systems&source=gbs\\_navlinks\\_s](http://books.google.co.in/books?id=WjvX0HmVTIMC&dq=deadlock+operating+systems&source=gbs_navlinks_s)) (7 ed.). Wiley-India. p. 237. . Retrieved 29 January 2012.
- [2] Padua, David (2011). *Encyclopedia of Parallel Computing* (<http://books.google.co.in/books?id=Ihm6LauVKFEC&lpg=PA1749&dq=computer%20networks%20deadlock%20definition&pg=PA524#v=onepage&q=deadlock&f=false>). Springer. p. 524. . Retrieved 28 January 2012.
- [3] Schneider, G. Michael (2009). *Invitation to Computer Science* (<http://books.google.co.in/books?id=gQK0pJONyhgC&lpg=PA271&dq=deadlock%20signal%20lost&pg=PA271#v=onepage&q=deadlock%20signal%20lost&f=false>). Cengage Learning. p. 271. . Retrieved 28 January 2012.
- [4] Rolling, Andrew (2009). *Andrew Rolling and Ernest Adams on game design* ([http://books.google.co.in/books?id=5SJHsm\\_PnUC&lpg=PA421&dq=deadlock%20chicken%20egg&pg=PA421#v=onepage&q=deadlock%20chicken%20egg&f=false](http://books.google.co.in/books?id=5SJHsm_PnUC&lpg=PA421&dq=deadlock%20chicken%20egg&pg=PA421#v=onepage&q=deadlock%20chicken%20egg&f=false)). New Riders. p. 421. . Retrieved 28 January 2012.
- [5] Oaks, Scott (2004). *Java Threads* ([http://books.google.co.in/books?id=mB\\_92VqJbsMC&lpg=PT82&dq=deadlock%20catch%22&pg=PT82#v=onepage&q=deadlock%20catch%22&f=false](http://books.google.co.in/books?id=mB_92VqJbsMC&lpg=PT82&dq=deadlock%20catch%22&pg=PT82#v=onepage&q=deadlock%20catch%22&f=false)). O'Reilly. p. 64. . Retrieved 28 January 2012.
- [6] A Treasury of Railroad Folklore, B.A. Botkin & A.F. Harlow, p. 381
- [7] Shibu (2009). *Intro To Embedded Systems* (<http://books.google.co.in/books?id=8hfntgwrR90MC&lpg=PA446&dq=coffman%20deadlock&pg=PA446#v=onepage&q=coffman%20deadlock&f=false>) (1st ed.). McGraw Hill Education. p. 446. . Retrieved 28 January 2012.
- [8] Stuart, Brian L. (2008). *Principles of operating systems* (<http://books.google.co.in/books?id=B5NC5-UfMMwC&lpg=PA112&dq=coffman%20conditions&pg=PA112#v=onepage&q=coffman%20conditions&f=false>) (1st ed.). Cengage Learning. p. 446. . Retrieved 28 January 2012.
- [9] *Distributed Operating Systems* (<http://books.google.co.in/books?id=l6sDRvKvCQ0C&lpg=PA177&dq=Tanenbaum%20ostrich&pg=PA177#v=onepage&q&f=false>) (1st ed.). Pearson Education. 1995. p. 117. . Retrieved 28 January 2012.
- [10] Mogul, Jeffrey C.; K. K. Ramakrishnan (1996). "Eliminating receive livelock in an interrupt-driven kernel" (<http://citeseer.ist.psu.edu/326777.html>). .
- [11] Anderson, James H.; Yong-jik Kim (2001). "Shared-memory mutual exclusion: Major research trends since 1986" (<http://citeseer.ist.psu.edu/anderson01sharedmemory.html>). .
- [12] Zobel, Dieter (October 1983). "The Deadlock problem: a classifying bibliography". *ACM SIGOPS Operating Systems Review* 17 (4): 6–15. doi:10.1145/850752.850753. ISSN 0163-5980.

### Further reading

- Kaveh, Nima; Emmerich, Wolfgang. *Deadlock Detection in Distributed Object Systems* (<http://www.cs.ucl.ac.uk/staff/w.emmerich/publications/ESEC01/ModelChecking/esec.pdf>). London: University College London.
- Bensalem, Saddek; Fernandez, Jean-Claude; Havelund, Klaus; Mounier, Laurent (2006). "Confirmation of deadlock potentials detected by runtime analysis". *Proceedings of the 2006 workshop on Parallel and distributed systems: Testing and debugging* (ACM): 41–50. doi:10.1145/1147403.1147412.
- Coffman, Edward G., Jr.; Elphick, Michael J.; Shoshani, Arie (1971). "System Deadlocks" ([http://www.cs.umass.edu/~mccorner/courses/691J/papers/TS/coffman\\_deadlocks/coffman\\_deadlocks.pdf](http://www.cs.umass.edu/~mccorner/courses/691J/papers/TS/coffman_deadlocks/coffman_deadlocks.pdf)). *ACM Computing Surveys* 3 (2): 67–78. doi:10.1145/356586.356588.
- Mogul, Jeffrey C.; Ramakrishnan, K. K. (1997). "Eliminating receive livelock in an interrupt-driven kernel". *ACM Transactions on Computer Systems* 15 (3): 217–252. doi:10.1145/263326.263335. ISSN 07342071.

- Havender, James W. (1968). "Avoiding deadlock in multitasking systems" (<http://domino.research.ibm.com/tchjr/journalindex.nsf/a3807c5b4823c53f85256561006324be/c014b699abf7b9ea85256bfa00685a38?OpenDocument>). *IBM Systems Journal* 7 (2): 74.
- Holliday, JoAnne L.; El Abbadi, Amr. "Distributed Deadlock Detection" ([http://www.cse.scu.edu/~jholliday/dd\\_9\\_16.htm](http://www.cse.scu.edu/~jholliday/dd_9_16.htm)). *Encyclopedia of Distributed Computing* (Kluwer Academic Publishers).
- Knapp, Edgar (1987). "Deadlock detection in distributed databases". *ACM Computing Surveys* 19 (4): 303–328. doi:10.1145/45075.46163. ISSN 03600300.
- Ling, Yibei; Chen, Shigang; Chiang, Jason (2006). "On Optimal Deadlock Detection Scheduling". *IEEE Transactions on Computers* 55 (9): 1178–1187.

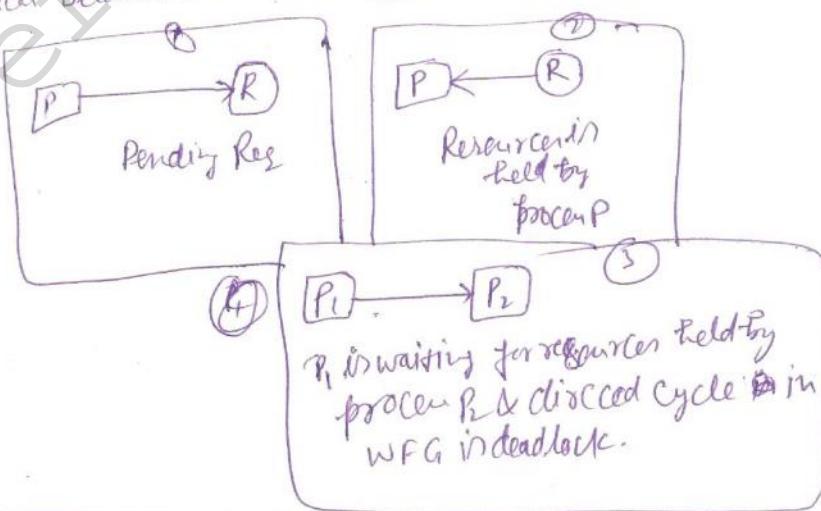
### External links

- "Advanced Synchronization in Java Threads (<http://www.onjava.com/pub/a/onjava/2004/10/20/threads2.html>)" by Scott Oaks and Henry Wong
- Deadlock Detection Agents (<http://www-db.in.tum.de/research/projects/dda.html>)
- DeadLock at the Portland Pattern Repository
- Etymology of "Deadlock" (<http://www.etymonline.com/index.php?term=deadlock>)
- ARCS - A Web Service approach to alleviating deadlock (<http://www.arcus.us>)

Deadlock Detection →

- ① Centralized deadlock detection - (use coordinator) → ① → ②  
- WFG
- ② Distributed deadlock detection - ②

- ③ Hierarchical Deadlock detection - ③



## DISTRIBUTED DEADLOCK DETECTION

(4)

Introduction →

- A distributed system is a network of sites that exchange information with each other by message passing.
- A site consists of computing & storage facilities & interface to local user & a communication network.
- In DS, a process can request & release resources (local or remote) in any order, which may not be known a priori & a process can request some resources while holding others.
- If the sequence of the allocation of resources to process is not controlled in such environments, deadlock can occur.

System Model →

- The problem of deadlock has been generally studied in DS under the following model -

- ① The systems have only reusable resources.
- ② Processes are allowed only exclusive access to resources.
- ③ There is only one copy of each resource.

- A process can be in two states - running or blocked.

Running state (active state) - A process has all the resources & is either executing or ~~not~~ ready to for execution.

Blocked state - Process is waiting to acquire some resources.

## Resource Vs Communication Deadlock

Two types of deadlock have been discussed -

① Resource deadlock - Process can simultaneously wait for several resources & can't proceed until they have acquired all those resources.

- A set of processes is resource-deadlock if each process in the set requests resources held by another process in the set & it must receive all of the requested resources before it can become unblocked.

② Communication Deadlocks Processes wait to communicate with other processes among a set of processes.

- A waiting process can be unblocked on receiving a communication from any one of these processes.

- A set of processes is communication-deadlock if each process in the set is waiting to communicate with another process in the set & no process in the set ever initiates any further communication until it receives the communication for which it is waiting.

"wait to communicate" → "wait to acquire a resource"

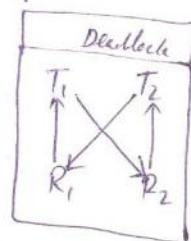
## A Graph-Theoretic Model

(5) (2)

- The state of Process-resource interaction in DS can be modeled by a bi-partite directed graph called a resource allocation graph.
- Nodes of this graph are processes & resources of a system, & the edges of the graph depict assignments or pending requests.
- A pending request is represented by a request edge directed from the node of a requesting process to the node of the requested resource.
- A resource assignment is represented by an assignment edge directed from the node of an assigned resource to the node of the assigned process.
- A system is deadlocked if its resource allocation graph contains a directed cycle or a knot.

## Wait-for-Graphs (WFG)

(See - Page No - (52) of Singhal & Shrivastri )



## Deadlock handling strategies in Distributed System →

- ① Prevention. (It works by preventing one of the 4 conditions from occurring)  
 $E_{NP} \rightarrow E_{HW}$
- ② Avoidance. ( Safe/unsafe state )
- ③ Detection. ( Here Deadlock are allowed to occur)  
If deadlock is detected, it can be corrected by using one of the following methods -
  - ① Process Termination.
  - ② Resource Preemption.

## Issues in Deadlock Detection & Resolution

- Deadlock detection & resolution entails addressing 2 basic issues -

① Detection of existing Deadlock.

② Resolution of Distributed deadlock.

### ① Detection

- The detection of Deadlock involves two issues -

\* Maintenance of the WFG & Search ~~Search~~ of the WFG for the presence of cycles (or knots).

- In Distributed system, a cycle may involve several nodes, so the search for cycles greatly depends upon how the WFG of the system is represented across the system.

- A correct deadlock detection algorithm must satisfy the following 2 conditions -

① Progress - No undetected deadlock -

- The algorithm must detect all existing deadlocks in finite time.

- Once a deadlock occurs, the deadlock detection activity should continuously progress until the deadlock ~~is~~ is detected.

② Safety - No false Deadlock -

- The algorithm should not report deadlocks which are non-existent (called phantom deadlock).

## Resolution

⑥

→ Deadlock Resolution involves breaking existing wait for dependencies in the System WFG to resolve the Deadlock.

- It involves rolling back one or more processes that are deadlocked & assigning their resources to blocked processes in the deadlock so that they can resume execution.
- When WFG dependency broken, the corresponding info should be immediately cleaned from the system. If the info is not cleaned appropriately in a timely manner, it may result in detection of Phantom Deadlock.

Control organization for Distributed Deadlock Detection → (for detail Read Page No. 155 - SKS)

① Centralized Control - (Central site has the responsibility of maintaining the Global WFG & searching it for cycles. (It has single point of failure))

② Distributed Control - Responsibility is shared equally. WFG is spread over many sites, several sites participate in deadlock detection also.

③ Hierarchical Control -

- This approach combines benefits of centralized & distributed approach in one
- Sites detect deadlock in only dependent sites (child) & resultant WFG or RAG is forwarded to site to calculate result of WFG or RAG.

Advantage - RAG & WFG is easy to be constructed.  
Disadvantage - Congestion (Bottleneck Problem) at single point failure in these.

Advantages Not susceptible to single point failure.

Disadvantages Difficult to design algorithms & sort of NP-Complete Problem, also proof correctness of algo is difficult.

## ① Centralized Deadlock Detection Algorithms

(a) Ho-Ramamoorthy's one & two phase algorithms.

## ② Distributed Deadlock Detection Algorithms

(a) Obermairk's Path Pursuing Algorithm.

(b) - Chandy - Misra - Haas Edge chasing algorithm  
(c) Diffusion Computation (d) Global State detection.

## ③ Hierarchical Deadlock detection -

(a) - Menasce - Mintz Algorithm.

(b) - Ho - Ramamoorthy's Algorithm.

### ① (a) HO-RAMAMOORTHY ALGO -

#### ① One-Phase algo -

- Each site maintains 2 status table -

- Process Table

- Resource Table

- one of the site becomes the central control site.

- the central control site periodically asks for the status tables.

- control site builds WFG using the status tables.

- control site analyzes WFG & resolves any present cycles.

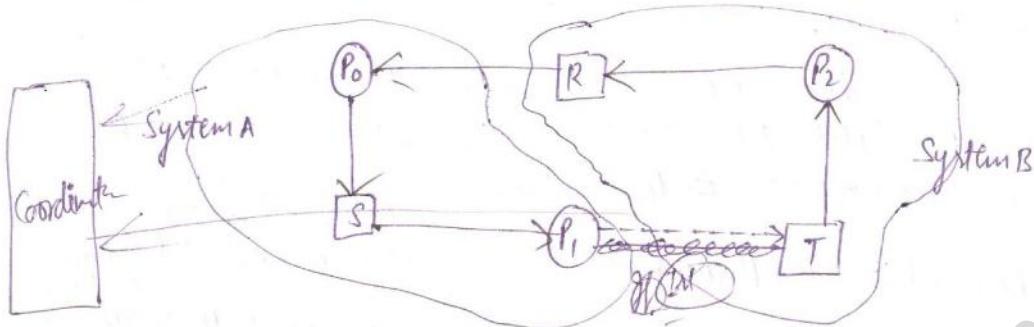
#### Shortcomings -

- Phantom Deadlocks.

- High Storage & Communication Costs.

## Phantom Deadlocks

(7) (4)



- System A has Process  $P_1 \rightarrow$  holding resource  $S$  & waiting for  $T$  which is in System B. Resource  $T$  is presently held by Process  $P_2$ .

- Here Process  $P_1$  sends a message to the Control Site stating the release of resources  $S$  & need of resources  $T$ .

- So in case the Message 2, that is waiting for  $T$  reaches the Control Site first then a cycle would be detected & it would be indicated that the system has a deadlock. (The false deadlock is called a phantom deadlock).



## 2-Phase Algorithm

- Each site maintains a status table for processes that initiated at that site.

- Resources Locked & Resources waited upon.

### ① Phase-1

- Control Site periodically asks for the locked & waited tables.  
- It then searches for presence of cycle in these tables. (WFG)

### ② Phase-2

- If cycle are found in Phase-1 search, Control Site makes 2nd request for the tables.  
- The details found common in both table requests will be analyzed for cycle confirmation.

Shortcoming - (phantom deadlock)

- The one phase algo is faster & requires fewer messages as compared to the two-phase algorithm.

- However, It requires (1-phase) more storage because every sites maintains two status tables & exchanges bigger messages because message contains 2 tables instead of one. so Process can define as TWF

② ① Obermark's Path-Pushing algo → (Designed for distributed database system)

- In Path Pushing deadlock detection algo, information about the wait for dependencies is propagated in the form of paths.

Algo

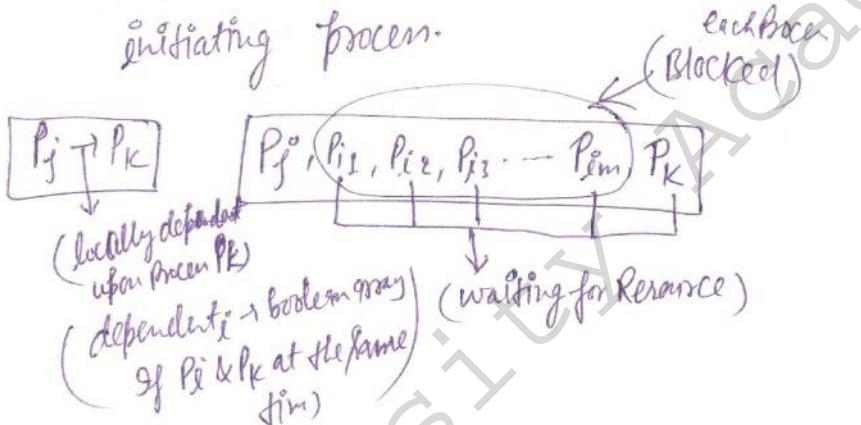
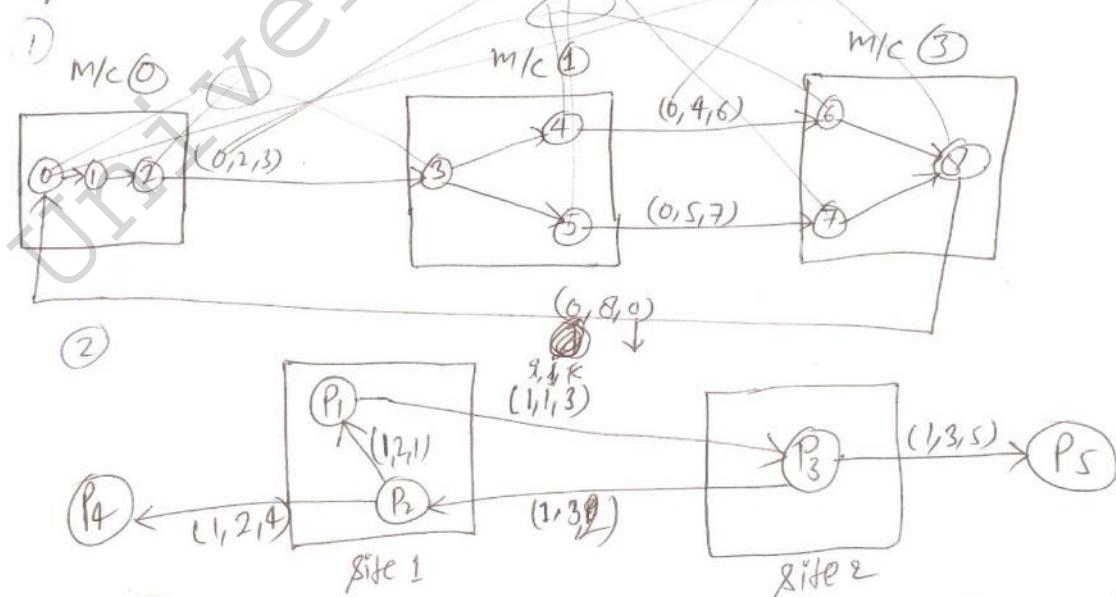
- The site waits for deadlock-related information from other sites.

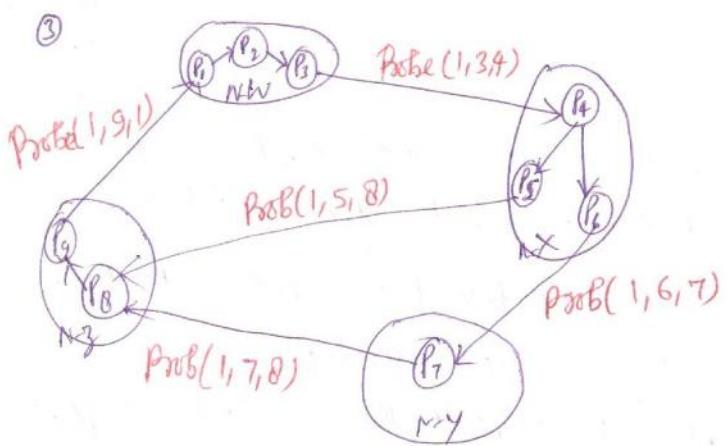
- The site combines the received information with its local TWF graph to build an updated TWF graph. (Transaction wait for G)

- For all cycles ' $Ex \rightarrow T_1 \rightarrow T_2 \rightarrow Ex$ ' which contain the node 'Ex', the site transmits them in string from 'Ex,  $T_1, T_2, Ex$ ' to all other sites where a sub-transaction of  $T_2$  is waiting to receive a message from the sub transaction of  $T_2$  at that site.

An Edge-Chasing Algorithm

- Chandy algo uses a special message called a probe. A probe is a triplet  $(i, j, k)$  denoting that it belongs to a deadlock detection initiated for process  $P_i$  & it is being sent by the home site of process  $P_j$  to the home site of process  $P_k$ .
- A Probe message travels along the edges of the global TWK graph, & a deadlock is detected when a probe message return to its initiating process.

Algos



## DISTRIBUTED DEADLOCK DETECTION

(9)

- In DDD algorithm, all sites collectively cooperate to detect a cycle in the state graph that is likely to be distributed over several sites of the system.
- A DDDAlg can be initiated whenever a process is forced to wait, & it can be initiated either by the local site of the process or by the site where the process waits.
- It is divided into 4 class -

① Path Pushing → The wait for dependencies information of the Global WFG is disseminated in the form of path (i.e. a sequence of wait-for dependency edges).

② edge-chaining → Special message called probes are circulated along the edges of the WFG to detect a cycle. When a blocked process receives a probe, it propagates the probe along its outgoing edges in the WFG.

- A process declares a deadlock when it receives a probe initiated by it. (probes are of a fixed size Normally very short)

③ Diffusion Computation types

- Make use of echo algorithm to detect deadlock. To detect a deadlock, a process sends out query messages along all the outgoing edges in the WFG.

- These queries are successively propagated through the edges of the WFG.

- Queries are discarded by a running process & are echoed back by block processes as following way - when a blocked process receives first query message for a particular deadlock detection initiation, it does not send a reply message until it has received a reply

Merge for every query it sent (so its successor in WFG), for all subsequent queries for this deadlock detection ~~algorithm~~ initiation, it ~~will~~ immediately sends back a reply merge. The initiator of a deadlock detection detects a deadlock when it receives a reply for every query it sent.

④ Global detection based deadlock detection,

- A consistent snapshot of a DS can be obtain without freezing the underlying computation.

### ③ Hierarchical Deadlock Detection

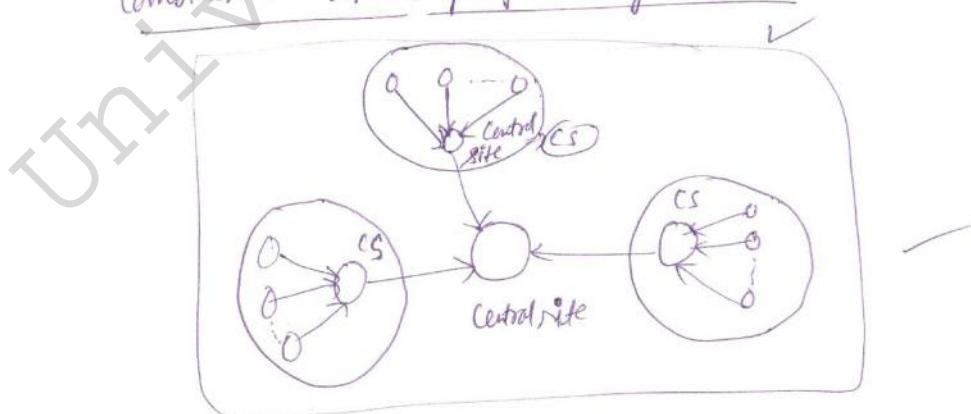
(10)

#### ① Menage-Mintz Algo

- Sites (controllers) organized in a tree structure.
- leaf controllers manage local WFGs.
- upper controllers handles Deadlock Detection.
- each parent node maintains a Global WFG, union of WFGs of its children. Deadlock detected for its children.
- Changes propagated Upwards in the tree.

#### ② Ho-Ramamoorthy Algo

- Site are Grouped into clusters.
- Periodically 1 site chosen as central control site -
  - Central Control Site chooses control site for other clusters.
- Control site for each cluster collects the status graph there -
  - Ho-Rammoorthy ~~algo~~ 1-Phase algo centralized DP used.
- All control sites forward status report to central control site which combines the WFGs & performs cycle search.



## UNIT-III

### UNIT-3

#### AGREEMENT PROTOCOL

(1)

##### Introduction :-

- \* In Distributed Systems, where sites (or processors) often compete as well as Cooperate to achieve a common goal, it is often required that sites reach Mutual agreement.
- \* For ex, in ~~Distributed~~ Data managers at sites must agree on whether to Commit or to abort a transaction.
- \* Reaching an agreement typically requires that sites have knowledge about the values of other sites.
  - ex- In ~~Do~~ Commit, a site should know the outcome of local Commit at each site.
- \* When the system is free from failures, an agreement can easily be reached among the processors (or sites).
  - ex- Processors can reach an agreement by communicating their values to each other and then by taking a majority vote ~~for~~ a minimum, Max<sup>m</sup>, mean etc, of those values.
- \* When a system <sup>tendency</sup> prone to failure this method does not work. This is because faulty processors can send conflicting values to other processors preventing them from reaching an agreement.
- \* In the presence of faults, processors must exchange their values with other processors & relay the values received from other processors several times to isolate the effects of faulty processors.
- \* A processor refines its value as it learns of the values of other processors.
  - this entire process of reaching an agreement is called an agreement Protocol.

- A very general model of faults is assumed. ex- a faulty processor may send spurious (ok but false) messages to other processors, may lie, may respond to received message correctly etc.
- Also non-faulty processors do not know which processors are faulty.
- In Agreement problems, nonfaulty processors in Distributed system should be able to reach a common agreement, even if certain components in the system are faulty.
- The agreement is achieved through an agreement protocol that involves several rounds of message exchange among the processors.

### The System Model :-

Agreement problems have been studied under the following system model.

- ① There are  $n$  processors in the system & at most  $m$  of the processors can be faulty.
- ② The processors can directly communicate with other processors by message passing. thus, the system is logically fully connected.
- ③ A receiver processor always knows the identity of the sender processor of the message.
- ④ The communication medium is reliable (i.e. it delivers all messages without introducing any errors) & only processors are prone to failures.

- For simplicity, we assume that agreement is to be reached between only two values, 0 & 1. Results can easily be extended to multivalue agreement. (2)

### Synchronous Vs Asynchronous Computations -

- In a Synchronous Computation, processes in the system run in lock step manner, where in each step, a process receives messages (sent to it in the previous step), performs a computation, & sends messages to other processes (received in the next step). A step of Synchronous Computation is also referred as a round.
- In Synchronous Computation, a process knows all the messages it expects to receive in a round. A message delay or a slow process can slow down the entire system or computation.
- In an asynchronous computation, on the other hand, the computation at processes does not proceed in lock steps. A process can read & received messages & perform computation at any time.  
"Here in this chapter, synchronous models of computation are assumed." the assumption of synchronous computation is critical to agreement protocols.

### Model of Processor Failures :-

- In Agreement Problem we consider a very general model of processor failures.
- \* A processor can fail in 3 modes: crash fault, omission fault, Malicious fault.
- \* In Crash fault  $\rightarrow$  A processor stops functioning & never resumes operation.
- \* In omission fault  $\rightarrow$  A processor "omits" to send messages to some processes. (there are the messages that the processor should have sent according to protocol or Algo it is executing.) (fail to perform some steps)
- \* In Malicious fault - a processor behaves randomly & arbitrarily. (Malicious faults are also known as Byzantine faults)

### Classification of Faults-

Based on components that faults

- program/ process
- Processor / Machine
- Link
- Storage
- Clock

Based on behaviors of faulty component

- Crash - Just Crash
- Faultstop - Crash with additional conditions
- omission - fails to perform some steps
- Byzantine - behaves arbitrarily
- Timing - violates timing constraints

### Authentication Vs Non-Authentication Merges -

- To reach an agreement, processors have to exchange their values & relay the received values to other processors several times. The capability of faulty processors to distort (make false by addition) what they receive from other processors greatly depends upon the type of underlying merges.
- There are two types of merge - authenticated & non-authenticated.
  - \* In an authenticated merge system, a (faulty) processor cannot forge a merge or change the contents of a received message.
  - \* In non-authenticated merge system, a (faulty) processor can forge a merge & claim to have received it from another processor or change the content of a received message before it relays the message to the other processors. A processor has no way to verify the authenticity of a received message.

Performance aspects → the performance or computational complexity of agreement protocol is generally determined by the following three metrics -

- ① time - (time taken to reach an agreement under a protocol)
- ② merge traffic - (measured by the no of messages exchanged to reach an agreement.)
- ③ storage overhead -

measures the amount of information that needs to be stored at processors during the execution of a protocol.

### Classification of Agreement Problems

(3)

- There are three well known Agreement Problem in distributed system-

- ① The Byzantine Agreement Problem.
- ② The Consensus Problem.
- ③ the interactive consistency Problem.

- In Byzantine agreement problem, a single value, which is to be agreed, is initialized by an arbitrary processor & all nonfaulty processors have to agree on that value.

- In Consensus Problem, every processor has its own initial value & all nonfaulty processors must agree on a single common value.

- In the interactive consistency problem, every processor has its own initial value & all nonfaulty processors must agree on a set of common values.

\* In all three problems, all nonfaulty processors must reach a common agreement.

- In Byzantine agreement & Consensus problems, the agreement is about a single value.

- whereas in interactive consistency problem, the agreement is about a set of common values.

Table 1 - Three Agreement Problem.

| Problem                | Byzantine Agreement | Consensus      | Interactive consistency |
|------------------------|---------------------|----------------|-------------------------|
| who initiate the Value | one processor       | All processors | All processors          |
| Final agreement        | Single Value        | Single Value   | A vector of values      |

## ① The Byzantine Agreement Problem :-

- In the Byzantine agreement problem, an arbitrarily chosen processor, called the source processor, broadcasts its initial value to all other processors.
- A sol<sup>n</sup> to the Byzantine agreement problem should meet the following two objectives -

Agreement → All nonfaulty processors agree on the same value.

Validity → If the source processor is nonfaulty, then the common agreed upon value by all nonfaulty processors should be the initial value of the source.

- Two points should be noted if ~~all~~ nonfaulty
  - ① If the source processor is faulty, then all nonfaulty processors can agree on any common value.
  - ② It is irrelevant what value faulty processors agree on or whether they agree on a value at all. (Termination) each nonfaulty processor must eventually decide on a value.

## ② The Consensus Problem :-

- Every processor broadcasts the initial value to all the other processors. Initial value of the processors may be different. A protocol for reaching consensus should meet the following conditions -
  - ① Agreement → All nonfaulty processors agree on the same single value.
  - ② Validity → If the initial value of every nonfaulty processor is  $v$ , then the agreed upon common value by all nonfaulty processors must be  $v$ .
- Note that if the initial values of nonfaulty processors are different, then all non-faulty processors can agree on any common value. again we don't care what value faulty processors agree on.
- ③ Termination → Each nonfaulty processor must eventually decide on ~~one~~ ~~any~~ a value.

### ③ The Interactive Consistency Problem →

(4)

- In the interactive consistency problem, every processor broadcasts its initial value to all other processors. The initial values of the processors may be different.
- A protocol for the interactive consistency problem should meet the following conditions:  
Agreement → All nonfaulty processors agree on the same vector  $(v_1, v_2, \dots, v_n)$ .

Validity If the  $i$ th processor is nonfaulty & its initial value is  $v_i$ ,  
then the  $i$ th value to be agreed on by all nonfaulty processors must be  $v_i$ .

~~Termination~~ Each nonfaulty processor must eventually decide on the array  $A$ .

### Relations Among the Agreement Problems →

- All three agreement problems are closely related.
- The Byzantine agreement problem is a special case of the interactive consistency problem, in which the initial value of only one processor is of interest.
- Conversely, if each of the  $n$  processors runs a copy of the Byzantine agreement protocol, the interactive consistency problem is solved.
- The consensus problem can be solved using the sol<sup>u</sup>t<sup>n</sup> of the interactive consistency problem, this is because all nonfaulty processors can compute the value that is to be agreed upon by taking the majority value of the common vector that is computed by an interactive consistency protocol, or by choosing a default value if a majority does not exist.
- Thus sol<sup>u</sup>t<sup>n</sup>s to the interactive consistency & consensus problems can be derived from the sol<sup>u</sup>t<sup>n</sup> to the Byzantine agreement problem.

## Ans<sup>n</sup> to the Agree Byzantine Agreement Problem

- The Byzantine agreement Problem was first defined & solved.
- (ex - Byzantine Agreement Problem in def<sup>n</sup>.)
- It is ~~also~~ obvious that all the processes must exchange the values through messengers to reach a consensus. Processes send their values to other processes & relay received value to other processes. During the execution of the protocol, faulty processes may confuse other processes by sending them conflicting values or by relaying to them fictitious values.  $\rightarrow$  (the)
- Byzantine Agreement Problem is referred as the Byzantine generals problem, because the Problem resembles a situation where a team of generals in an army is trying to reach agreement on an attack plan.
- The general are located at geographically distant & ~~bad~~ ~~change~~ ~~communication~~ communicate through messengers. Some of the generals are traitors (equivalently, faulty processes) & try to prevent loyal generals from reaching an agreement by deliberately transmitting erroneous information.

### (1) The upper bound on the No of Faulty Processors

- In order to reach an agreement on a common value, nonfaulty processes need to be free from the influence of faulty processes.
- If the faulty processes dominate in number, they can prevent non-faulty processes from reaching a consensus. Thus the no of faulty processes should not ~~not~~ exceed a certain limit if a consensus is to be reached.
- If it is impossible to reach a consensus (Agreement) if the no of faulty processes,  $m$ , exceeds  $\lceil \frac{1}{3}(n-1) \rceil$ . (Lamport)
- $\lceil \frac{1}{3}(n-1) \rceil$  rounds off  $\lceil \frac{1}{3}(n-1) \rceil$  messages exchanged ( $m$  = faulty process)

$\lceil \frac{1}{3}(n-1) \rceil$  is the lower bound on the No of messages exchanged to reach a Byzantine agreement in fully connected network.

(5)

## ② An Impossibility Result →

- Byzantine Agreement is impossible if  
 $m > \lfloor (n-1)/3 \rfloor$   
 $m \rightarrow \text{faulty Processors}$

$$\frac{n=3 \text{ process}}{\lfloor (3-1)/3 \rfloor = \lfloor 2/3 \rfloor = 0}$$

- Byzantine Agreement is impossible with  $< (m-1)$  message exchanged.

- Byzantine Agreement cannot be reached among three processors, where one processor is faulty.

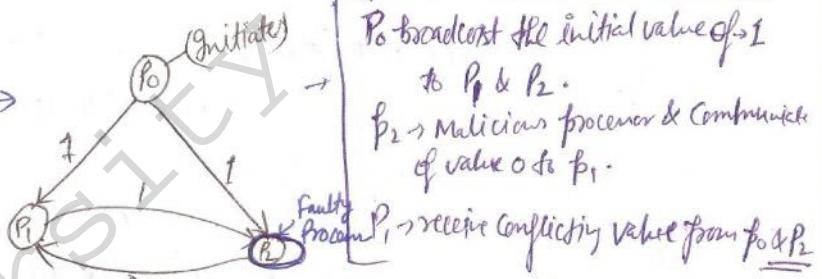
- Consider a system with three processors,  $p_0, p_1$  &  $p_2$ . we assumed that there are only 2 values 0, 1, on which processors agree & processor  $p_0$  initiates the initial value. There are two possibility -

- ①  $p_0$  is faulty or ②  $p_0$  is not faulty.

faulty processor oval (○)  
Nonfaulty → O

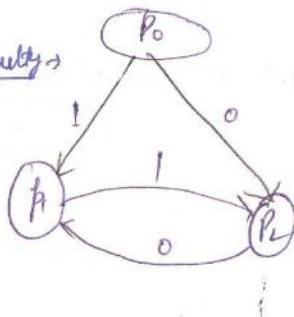
✓ ①  $p_0$  is not faulty →

Assume  $p_0$  is faulty →  
processor  $p_1$  must accept 1 as  
the agreed upon value if condition  
2 in be satisfied.



✓ ②  $p_0$  is faulty

Assume Initiator  $P_0$  is faulty →



$P_0$  send 1 to  $P_1$  & 0 to  $P_2$

$P_2$  communicate the value 0 to  $P_1$

Re-read

✓  $P_2$  → It must be 0 agreed upon value

& at  $P_1$  → It must be 1 agreed upon value

→ Therefore No sol<sup>n</sup> exists for Byzantine Agreement Problem for three processors, which can work under single processor failure.

### ③ Lamport-Shostak-Pease Algo (or message algo OM(m), m>0)

Solves Byzantine agreement Problem for  $3m+1$  or more processors in the presence of at most  $m$  faulty processors.

✓ If  $n$  be the total No of processors then  $n \geq 3m+1$ .

✓ The Algorithm is recursively defined as -

#### ① Algorithm OM(0) :-

- ① The source processor sends its value to every processor.
- ② each processor uses the value it receives from the source.  
(If receives no value then it uses a default value of 0).

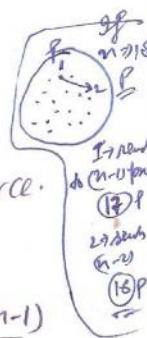
OM(0) (if there is 0 faulty processor in the system)

$$\text{if } m=1 \\ n \geq 3 \times 1 + 1 = 4 \text{ faulty}$$

3+1

#### ② Algorithm OM(m), m > 0 :-

- ① the source processor sends its initial value to every processor.
- ② for each  $i$ , let  $v_i$  be the value, processor  $i$  receives from the source.  
(If receives no value, then it uses a default value of 0).  
processor  $i$  acts as the new source & initiates Algorithm OM( $m-1$ )  
where it sends the value  $v_i$  to each of the  $(n-2)$  other processors.
- ③ for each  $i$  and each  $j \neq i$ , let  $v_{ij}$  be the value, processor  $i$  received from processor  $j$  in Step 2 using Algorithm OM( $m-1$ ).  
(If it receives no value, then it uses a default value of 0).  
Processor  $i$  uses the majority  $(v_1, v_2, \dots, v_{n-1})$ .



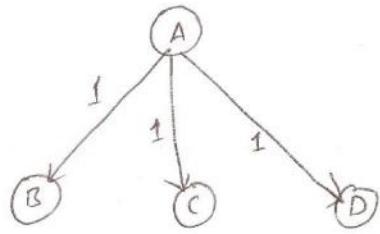
\* The execution of the Algorithm OM( $m$ ) invokes  $(n-1)$  separate executions of the Algorithm OM( $m-1$ ), each of which invokes  $(n-2)$  executions of the algorithm OM( $m-2$ ), & so on.

\* Therefore, there are  $(n-1)(n-2)(n-3) \dots (n-k)$  separate executions of algo OM( $m-K$ ),  $K=1, 2, 3, \dots, m+1$ .

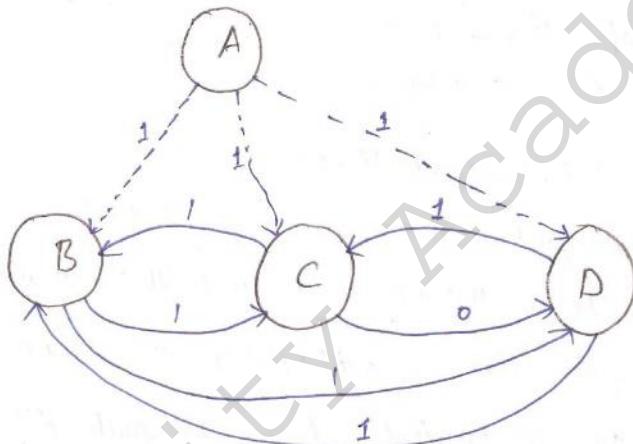
\* The Message complexity of the algo is  $O(n^m)$ .

Ex- 4 Processor ex- Nonfaulty Commander.

(6)



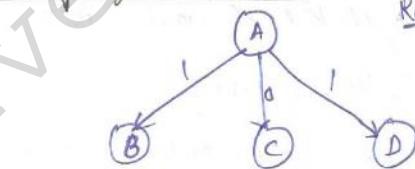
Round 1 - Processor A executes OM(1), where processor C is faulty



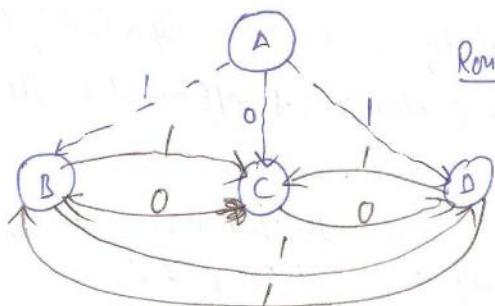
Round 2 : Processors B, C, & D execute OM(0). Dashed line shows that Memory Present during first phase.

Ex- 4 Processor - ex faulty commander

Round ① - Processor A executes OM(1), where processor A is faulty.



Round ② - Processor B, C, & D execute OM(0).



Merge Complexity  $\rightarrow$

(LIEUTENANT  
COMMANDER)  $\rightarrow$  (Problem)  
(Agreement)

$$* T(0, n) = (n-1)$$

$$* T(m, n) = (n-1) T(m-1, n-1), \text{ for } m > 0$$

$$* T(m, n) = (n-1)(n-2)(n-3) \cdots (n-(m+1)) \in O(n^m)$$

APPLICATIONS OF AGREEMENT ALGORITHMS  $\rightarrow$

- Algorithms for agreement Problems find applications in problems where processes should reach an agreement regarding their values in the presence of malicious failures.
- Two such applications - clock synchronization in distributed systems, atomic commit in distributed database.

(A) Fault Tolerance clock Synchronization  $\rightarrow$

- In Distributed Systems, it is often necessary that sites (or processes) maintain physical clocks that are synchronized with one another.
  - Since physical clocks have a drift problem, they must be periodically resynchronized.
  - Such periodic synchronization becomes extremely difficult if the Byzantine failures are allowed.
  - This is because a faulty process can report different clock values to different processes. (The description of clock synchronization is described in the paper of Lamport & Milligan-Smith. [H.W])
- The following assumptions regarding the system made-
- A<sub>1</sub>  $\rightarrow$  All clocks are initially synchronized to approximately the same value.
  - A<sub>2</sub>  $\rightarrow$  A Nonfaulty process's clock runs the approximately the correct rate.

A<sub>3</sub>  $\rightarrow$  A Nonfaulty process can read the clock value of another nonfaulty process with at most a small error  $\epsilon$ .

- A clock synchronization algorithm should satisfy the following two conditions. ⑦

- ① At any time, the values of the clocks of all nonfaulty processes must be approximately equal.
- ② There is a small bound on the amount by which the clock of a nonfaulty process is changed during each desynchronization.

\* There are 2 clock synchronization algs, namely ① Interactive Convergence algo.

② Interactive Consistency algo.

Pg No (190) - which are fault-tolerant to the Byzantine failures.

① Interactive Convergence algorithm  $\Rightarrow$  This alg. assumes that the clocks are initially synchronized & that they are resynchronized often enough so that two nonfaulty clocks never differ by more than  $\delta$ . The alg. works as follows:

Alg. Each process reads the value of all other processes' clocks & set its own clock to the average of these values - However, if a clock value differs from its own clock value by more than  $\delta$ , it replaces that value by its own clock's value when forming the average.

Let

$P, Q \rightarrow$  Non faulty process,

$r \rightarrow$  process

$C_{Pr}, C_{Qr} \rightarrow$  value used by  $P, Q$ .

When forming the average, respectively  $r$ 's values.

If  $r$  is nonfaulty - then  $C_{Pr} = C_{Qr}$

If  $r$  is faulty then  $C_{Pr} \& C_{Qr}$  will differ by  $\geq \delta$

$$\begin{cases} C_{Pr} \geq \delta \\ C_{Qr} \geq \delta \\ P \& Q \geq \delta \end{cases}$$

Let

$n \rightarrow$  total no of processes,  $m \rightarrow$  no of faulty processes,  $n > 3m$

processes  $P$  &  $Q$  set their clocks to the average of  $n$  values ( $C_{Pr}$ ) & ( $C_{Qr}$ ) respectively, for  $i=1, \dots, n$ .

We have,  $\rightarrow C_{Pr} = C_{Qr}$  for the  $[n-m]$  nonfaulty processes &  $|C_{Pr} - C_{Qr}| \leq 3\delta$

for the  $m$  faulty processes  $r$ .

- It follows from this that the average computed by P & Q will differ by at most  $(\frac{3m}{n})\delta$ .
- The assumption  $n > 3m$  implies  $(3m/n)\delta \ll \delta$ , so the algo succeeds in bringing the clocks closer together.
- Therefore we can keep the Nonfaulty processor's clocks synchronized to within  $\delta$  of one other by resynchronizing often enough so that clocks which are initially within  $(3m/n)\delta$  seconds of each other never drift further than  $\delta$  seconds apart.

- If the clock reading error is  $\epsilon$  then the difference in the clock values read by a process can be as long as  $\delta + \epsilon$ .  
 Therefore, only clock differences larger than  $\delta + \epsilon$  are replaced by 0 while computing the average increment.

$$\Delta_{QP} = \begin{cases} \Delta_{QP} & \text{if } \Delta_{QP} \neq \delta \\ 0 & \text{otherwise.} \end{cases}$$

(2) the Interactive Consistency Algorithm  
\* This algo add two improvements -

⑧ ~~Median~~ Median  $\frac{1}{2}(n+1)$ th value -  
Even M =  $\frac{(4th + 5th)}{2}$

- ① It takes the median of the clock values rather than the mean. The median provides a good estimate of the clock value, as the No of bad clocks will be low.
- ② It avoid the Problem of 2 faced clock (which reports differ values to differ processes) by using a more sophisticated technique to obtain clock values of the processes.

### (B) Atomic Commit in DDBS ->

- In the problem of atomic commit, sites of DBMS must agree whether to commit or abort the transaction.
- In the first Phase of the atomic commit, sites execute their part of a distributed transaction & broadcast their decisions (commit or abort) to all other sites.
- In the Second Phase, each site, based on what it received from other sites in the first phase, decides whether to commit or abort its part of distributed transaction.

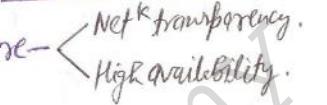
Since all sites receive an identical response from all other sites, they will reach the same decision. However, if some sites behaves maliciously, they can send a conflicting response to other sites, causing them to make conflicting decisions.

- In this situation, we can use algorithm for the Byzantine agreement to prove that all nonfaulty processes reach a common decision about a distributed transaction.
- It works as follow - In the first phase after a site has made a decision, it starts the Byzantine agreement, in the second phase processes determine a common decision based on the agreed vector of values.

University Academy

## Distributed Resource Management

①

- A distributed file system is a resource management component of a distributed operating system. It implements a common file system that can be shared by all the autonomous computers in the system.
  - NFS (Sun's Network File System)
    - Windows NT, 2000, XP
    - AFS (Andrew File System)
- Two important goals of distributed file systems are - 
  - ① Network transparency - The primary goal of a DFS is to provide the same functional capabilities to access files distributed over a network as the file system of a time-sharing mainframe system does to access files residing at one location.
    - Access
    - Location
    - Mobility
    - Performance
    - Scaling
    - Transparency
  - ② High availability → users should have the same easy access to files, irrespective of their physical location. System failures or regularly scheduled activities such as backups or maintenance should not result in the unavailability of files.
- In the recent years, several distributed file systems have been developed. Here we will understand that how common mechanisms & design aspects shared by today's distributed file systems.

## Architecture of Distributed file System

- Ideally, in a Distributed file system, files can be stored at any machine (or computer) and the computation can be performed at any machine (i.e. the machine or peers).
- When a machine needs to access a file stored on a Remote machine, the remote machine performs the necessary file access operations & return data if a read operation is performed.
- However, for higher performance, several machines, referred to as file servers, are dedicated to storing files and performing storage & retrieval operations. The rest of the machines in the system can be used solely for Computational Purpose. These machines are referred as clients & they access the files stored on the servers. (Client can use cache concept for fast access & serve)

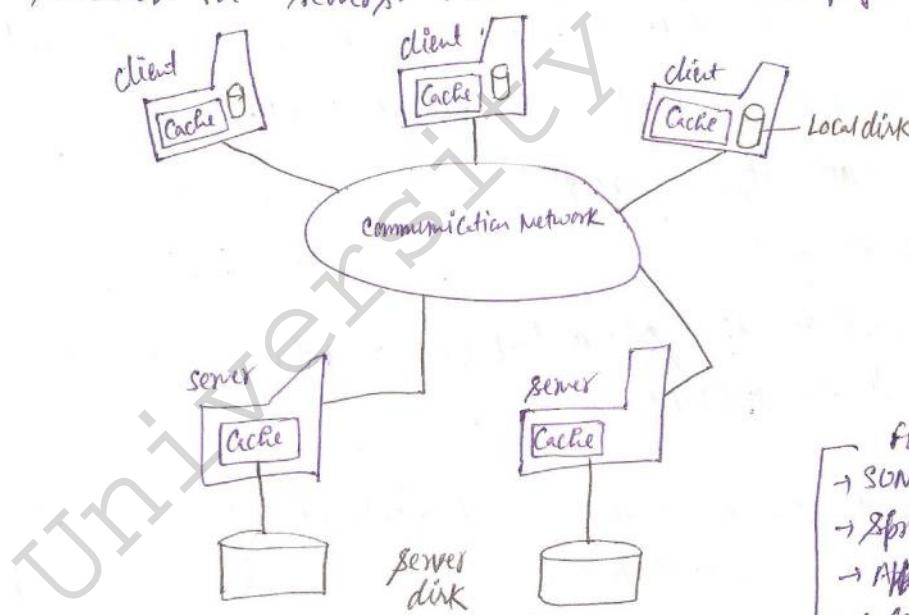
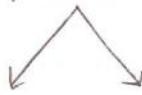


fig- Distributed file system Architecture

file system -  
→ SFS File System  
→ GFS File System  
→ Apollo Domain DFS  
→ Coda FS  
→ X-Kernel Logical FS  
→ Log Structured FS

(2)

Two most important services present in a distributed file system are the -



① Name Server

(Map Name to stored objects)

② Cache Manager (Perform file caching) & Can be present at both side Client & Server

① Name Server  $\Rightarrow$  It is a process that maps <sup>names</sup> specified by client to stored ~~object~~ object files & directories. The mapping (also referred as a name resolution) occurs when a process references a file or directory for the first time.

② Cache Manager  $\Rightarrow$  It is a process that implements file Caching.

- In file Caching, a copy of data stored at a remot file server, is brought to the client's m/c when ~~referenced~~ referenced by the client. Reducing the access delays due to Network latency , Cache managers can be present at both side of client's & servers
- Cache manager at the servers cache files in the main memory to reduce delays due to disk latency.

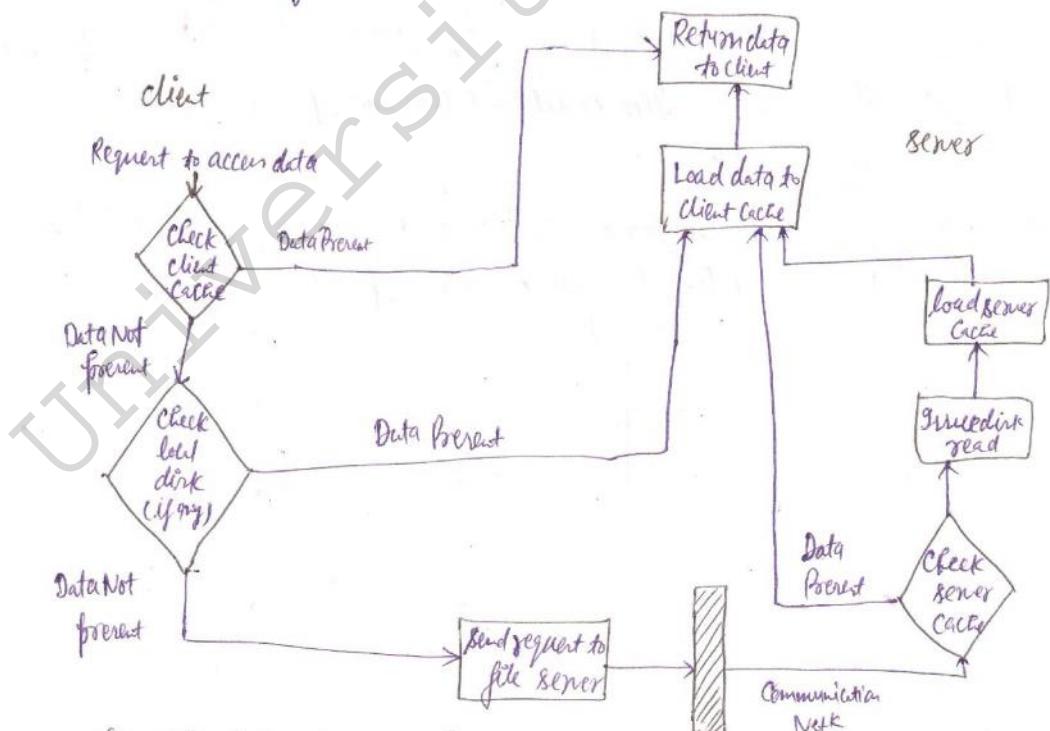
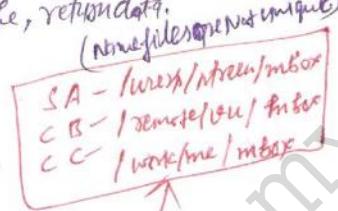


Fig- Typical Data access actions in DFS

Typical steps :-

- Check client cache, if present, return data.
- Check local disk, if present, load into local cache, return data.
- Send request to file server.

- Server checks cache, if present, load into client cache, return data.  
(name files are not unique)
- disk read
- load into server cache
- return data.

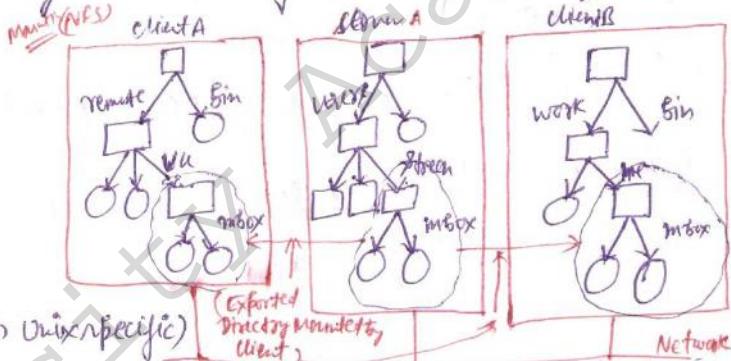


## Mechanisms for Building Distributed File Systems

- The Basic Mechanisms for building Distributed File Systems involves following steps:-

- ① Mounting.
- ② Caching.
- ③ Hints.
- ④ Bulk Data transfer.
- ⑤ Encryption.

### ① Mounting :- (Mounting is Unix specific)



- It allows the binding together of different filename space to form a single hierarchical structured name space.
- A name space (or a collection of files) can be bounded to or mounted at an internal node or a leaf node of a name space tree.

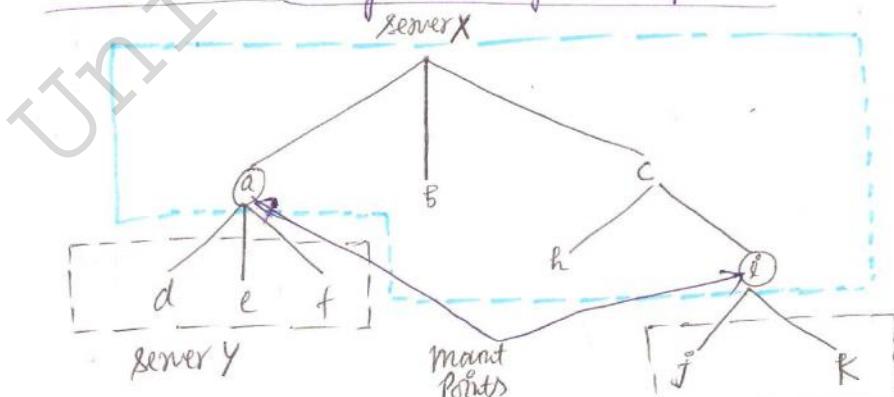


Fig - Namespace Hierarchy.

- A Node onto which a namespace is mounted is known as mount point. (3)
  - In this fig - Node a & Node i are mount points at which directories stored at Server Y & Server Z are mounted, respectively. Namespace (directory)
  - Node a & i are internal nodes in the namespace tree. The Kernel maintains a structure called the mount table, which maps mount points to appropriate storage devices.
  - In the DFS, file systems maintained by remote servers are mounted at the client.
  - There are two approaches to maintain the mount information-
    - ① At the client → (NFS) (Sun Net File system)
      - Different clients may see different files. (Every client needs to update its mount table)
    - ② At the servers (Sprite file system)
      - All clients may see the same filename space. good when file moves around the servers.
- ② Caching** → It implements in DFS to reduce delays in accessing data.
- In file caching, a copy of data stored at a remote file server is brought to the client when referenced by the client.
  - Subsequent access to the data is performed locally at the client, thereby reducing access delays due to network latency.
  - Caching exploits the (temporal locality of reference) exhibited by programs.
    - (Refers to the fact that a file recently accessed is likely to be accessed again in the near future.)
  - The data is cached ~~in the server~~ in the main memory (server cache) at the servers to reduce disk access latency.

### ③ Hints

- Caching improves file system performance by reducing the delay in accessing data.
- When multiple client cache & modify shared data, the problem of cache consistency arises.
- Specifically, it must be guaranteed that the cached data is valid (up-to-date) & that a copy of data - recently updated in some other client cache or in the file server - does not exist.
- Guaranteeing consistency is expensive in distributed file system as it requires elaborate cooperation between file servers & clients.
- An alternative approach is to treat the cached data as hints. In this case, cached data are not expected to be completely accurate.

Ex After the name of a file or directory is mapped to the physical object, the address of the object can be stored as a hint in the cache. If the address fails to map to the object in the following attempt, the cached address is purged (removed) from the cache. The file server consults filename server to determine the actual location of the file or directory & updates the cache.

### ④ Bulk Data Transfer

Transferring data in bulk reduces the protocol processing overhead at both servers & clients.

$$\text{Communication cost} = S + C \times B$$

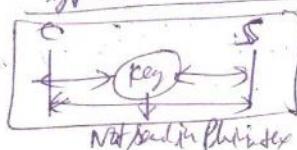
S → Startup cost, mostly fixed

C → Per Byte Cost

B → No of bytes to send.

### ⑤ Encryption

typical method to enforce security in distributed system.



## Design Issues

- (1) Naming & Name Resolution
- (2) Name Server
- (3) Cache Consistency
- (4) Writing Policy
- (5) Availability
- (6) Scalability
- (7) Semantics

(4)

### (1) Naming & name resolution

- \* A Name is associated with an object (such as file or directory).
- \* Name Resolution - Mapping a Name to an object.
- \* Name Space - (a collection of Names). three different methods to name files.

(1) Concentrate the Host name to the name of the file on that host.

- (A) Simple unique file Name.
- (B) No Network Transparency.
- (C) Moving a file from one machine to another may require modification to the application.

(D) Name Resolution is simple. Don't Need a Name Server.

(2) Mount remote directories onto local directory,

- (A) once Mounted, accessing the file becomes location transparent.
- (B) Name Resolution is simple.

(3) Have a single Global Directory (Sprite & Apollo FS)

- where all the files in the system belong to a single namespace.
- The main disadvantage of this scheme, however, is that it is mostly limited to one computing facility or to a few cooperating computing facilities.
- Not Scalable how to maintain the unique filenames.



## Name Servers

- In a centralized system, name resolution can be accomplished by maintaining a table that maps names to objects.
- In distributed system, name servers are responsible for name resolution.
- A name server is a process that maps names specified by clients to stored objects such as files & directories.
- The easiest approach to name resolution in DS is for all clients to send their queries to a single name server which maps names to objects.

- This Approach has the following serious drawback -

- ① If the Name Server crashes - the entire system is drastically affected.
- ② Name Server may become a bottleneck & seriously degrade the performance of the system. [a/c] → mapped objects may require Remote mapping for I/O
- ③ Cacher on Disk or Main Memory →
  - ✓ VirtualMemory (taken from Hard disk) (Real memory)
  - ✓ CacheMemory (External memory of CPU)

The advantages of having the cache in the Main memory are as following -

- ✓ Combined in Diskless workstation
- ✓ accessing a cache in main memory is much faster than accessing a cache on local disk.

- Same techniques for Server cache & Client cache.

The main disadvantage of having client-cache in main memory is that it competes with the Virtual memory system for Physical Memory space.

- Virtual Memory System more complex - data cannot be in both VM & Cache.
- Large files can't be cached completely - need to able to chop (cut in pieces) the file, server is more complex.

The Main Advantage of Caching on Disk -

- Large files can be cached.
- Virtual Memory <sup>not</sup> is simple.

### ③ ~~Writing Policy~~ <sup>→</sup>

⑤

- ✓ The writing Policy decides when a modified Cache at a client should be transferred to the server.
- ✓ The Simplest Policy is write-through, in this write-through, all writes requested by the applications at clients are also carried out at the servers immediately. (Synchronized way)
- ✓ The Main advantage of write-through is reliability. In the event of client crash little information is lost.
- ✓ A write-through policy does not take advantage of the Cache.
- ✓ An Alternative Writing Policy, delayed writing Policy, delays the writing at the server. In this case modifications due to a write are reflected at the server after some delay.
  - ✓ This approach can potentially take advantage of the Cache by performing many writes on a block present locally in the Cache.
- Another motivation for delaying the writes is that some of the data (ex- intermediate results) could be deleted in a short time, in which case data need not be written at the server at all.
- While the delayed writing Policy takes advantage of the Cache at a Client, it introduces the Reliability Problem.
- In the event of client crash a significant amount of Data can be lost.

#### ④ Cache Consistency :-

- The schemes that can guarantee consistency of the data cached at clients.
- There are two approaches to guarantee that the data returned to the client is valid.
  - \* In the server-initiated approach, servers inform cache managers whenever the data in the client cache becomes stale (old). Cache managers at clients can then retrieve the new data or invalidate the blocks containing the old data in their cache.
  - \* In the client-initiated approach, it is the responsibility of the cache managers at the clients to validate data with the server before returning it to the clients.
- Both of these ~~approaches~~ approaches are expensive & unattractive as they require elaborate cooperation bet<sup>n</sup> servers & cache managers. In both approaches communication cost is high.
  - The server-initiated approach requires the server to maintain reliable records on what data blocks are cached by which cache managers.
  - The client-initiated approach simply negates the benefits of having a cache by checking the server to validate data on every access.
- The third approach for cache consistency is simply not to allow file caching when concurrent write sharing occurs. In concurrent write sharing, a file is open at multiple clients & at least one client has it open for writing.
  - In this approach, the file server has to keep track of the clients sharing a file. When concurrent file sharing occurs for a file, the file server informs all the clients to purge their cached data items belonging to that file. (Alternatively, concurrent file sharing can be avoided by locking files. (Apollo FS))

- Another issue <sup>(4)</sup> that a cache consistency scheme needs to address is sequential-write sharing, which occurs when a client opens a file that has recently been modified & closed by another client. <sup>(5)</sup> <sup>(6)</sup>

- Two potential problems with sequential-write sharing are-

- ① When a client opens a file, it may have outdated blocks of the file in its cache, & (use timestamp associated with them).
- ② When a client opens a file, the current data blocks may still be in another client's cache waiting to be flushed (clean or remove).  
Server rewrites → Client flushed the modified blocks of a client's cache whenever a new client opens the file for writing.

## ⑤ Availability

- Availability is one of the important design issues of DFS. The failure of servers or the communication network can severely affect the availability of files.
- Replication is the primary mechanism used for enhancing the availability of files in DFS.

Replication Under replication, many copies or replicas of files are maintained at different servers. Replication is inherently expensive because extra storage space is required to store the replica & the overhead incurred in maintaining all the replicas up-to-date.

- The most serious problem of replication are-

- ① How to keep the replicas of a file consistent &
- ② How to detect inconsistency among replicas of a file & subsequent to recover these inconsistencies.

## Unit of Replication

- A fundamental design issue in replication is the unit of Replication.
- The most basic unit is a file. This unit allows the replication of only those files that need to have higher availability. It makes overall Replica Mgt easier. For ex - the protection rights associated with a directory have to be individually stored with each replica.
- Replication unit can be → a group of all the files of a single user or a files that in a server.
- The Group of files is referred as a Volume. (Hence of Volume Replication is that the Replica Mgt becomes easier)

## Replica Mgt

- Replica Mgt is concerned with the maintenance of replicas & in making use of them to provide increased availability.
- Replica Mgt depends on whether Consistency is guaranteed by the Distributed System.
- Here we are concerned with Consistency with Replica only. That is also known as Mutual Consistency, & not with the consistency within a file.
- To ensure mutual Consistency among replicas, a special weighted Voting Scheme can be used.

✓ Some No of (Voter + timestamp) are associated with → each replica.

Vote =  $r$  or  $w$  (Must be obtain through read, write respectively).

Condition →  $w > r$  &  $r+w > \text{total No of voters}$  of all the Replica  
(only Current update copies)  
( $r \rightarrow \text{read}$ )  
( $w \rightarrow \text{write}$ )

### (6) Availability

- Replication is the main mechanism.
- How to keep replicas consistent.
- How to detect inconsistencies among replicas.
- Consistency problem may decrease the availability.
- Replica mgt - Voting mechanism to read & write to replica.

### (7) Scalability

- \* How to meet the demand of growing system?
- \* the biggest problem - consistency issue
  - \* Caching - Reduce network load & server load.
  - \* Server-initiated Cache invalidation to maintain cache consistency: Complexity increases as the system size increases.
- \* Exploiting the mode for file use - (most shared files are readonly)

### (8) Semantics (Behaviour)

- what a user wants? strict consistency.
- users can usually tolerate a certain degree of errors in file handling - No need to enforce strict consistency.

#### ④ Writing Policy →

⑤

- The writing Policy defined when a modified Cache block at a client should be transferred to the Server.
- The Simpler Policy is write-through - In write through, all writes requested by the applications at clients are also carried out at the server immediately. (write is done synchronously both the cache & to the backing store)
- The main advantage of write-through is reliability. In the event of the client crash, little info is lost. A write-through does not take advantage of the cache.
- An alternating Policy, delayed writing Policy, delays the writing at the server. In this case, modifications due to a write are reflected at the server after some delay. This approach can potentially take advantage of the cache by performing many writes on a block present locally in the cache.

#### ⑤ Cache Consistency →

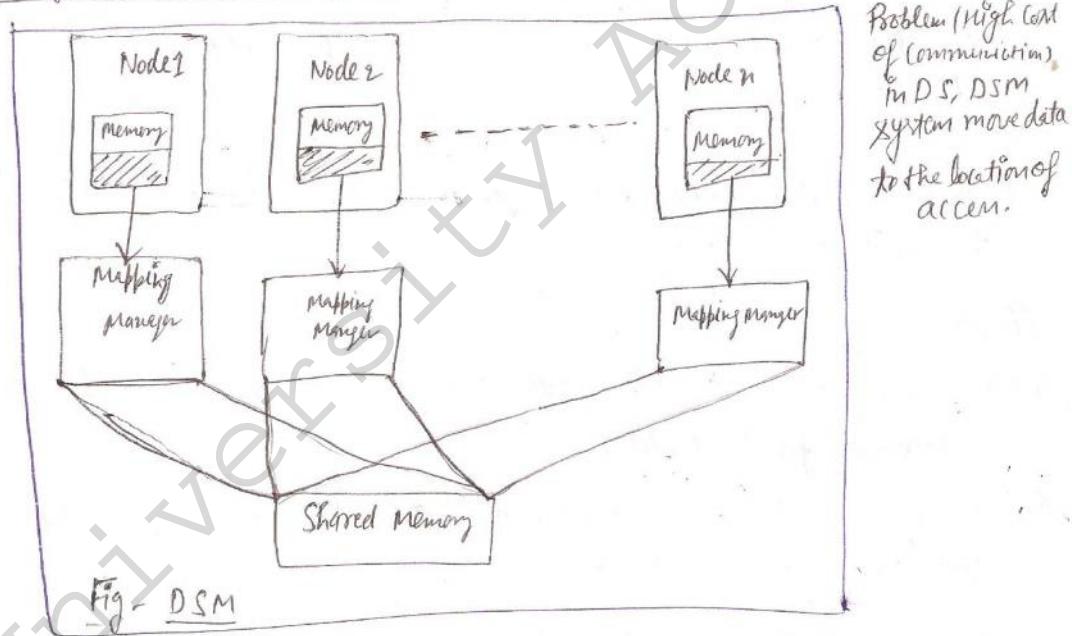
- Server-initiated - the server inform cache managers whenever the data in the client cache become stale (old).
- Client-initiated - the client makes sure the cache is clean before returning the data to the application.
  - ↓  
(validator)
- Both server-initiated & client-initiated ~~are~~ Cache Consistency Schemes are expensive.
- Not allow caching when concurrent-write sharing occurs.
- Sequential-write sharing may also cause problems - a client opens a file that has recently modified & closed by another client.
  - \* outdated blocks in the cache: timestamp of the cache copy. Contacting the server can find inconsistency.
  - \* update data in another client's cache not in the server (delay write) : flushing writes when other clients open.

## Distributed Shared Memory

①

### Introduction →

- Distributed Computing has been based on the Message Passing model in which processes interact & share data with each other by exchanging data in the form of messages. (Explicit Message Passing) (No Physical Shared Memory)
- DSM system is a resource mgmt component of Distributed Operating System that implements the shared memory Model in distributed systems, which have no physically shared memory.
- The shared memory model provides a virtual address space that is shared among all nodes (computers) in a distributed system.



- To overcome the problem (High cost of communication) in DS, DSM system moves data to the location of access.

### Architecture & Motivation →

- With DSM, programs access data in the shared address space just as they access data in traditional virtual memory.
- In systems that support DSM, data moves between Secondary Memory (HDD) & Main Memory (RAM) as well as between main memories of different nodes.

- Each node can own (A Node which creates a data object owns the data object initially)
  - ✓ data stored in Shared address space, and the ownership can change when data moves from one node to another.
  - ✓ When a process access data in the Shared address space, a mapping manager maps the Shared memory address to the Physical memory (which can be local or remote).
 

(RAM) Physical Blocks → PM → Logical Memory Unit  
Logical Blocks → LMSD → Page → CPU register
  - ✓ The Mapping Manager is a layer of SW implemented either in the OS Kernel or as a routine library routine.
  - ✓ To reduce delays due to communication latency, DSM may move data at the Shared memory address from a remote node to the node that is accessing data (when the Shared memory address maps to a Physical Memory location on a Remote Node).
  - ✓ In such case DSM make use of the communication services of the underlying communication system.
- Advantages of DSM are -
- ① Data sharing is implicit, hiding data movement (as opposed to "send/receive" in message passing Model).
  - ② Passing data structures containing pointers is easier (In message passing Model Data moves between different address spaces). (use Pass By reference).
  - ③ moving entire object to user <sup>DSM</sup> takes advantage of locality of reference.
  - ④ less expensive to build than tightly coupled multiprocessor systems off-the-shelf HW, no expensive interface to shared physical memory.
  - ⑤ very large total physical memory for all nodes - ~~so~~ Large programs can run more efficiently.

⑥ No serial access to common bus for shared Physical Memory like in multi-processor system. ②

⑦ Programs written for Shared Memory multiprocessors can be run on DSM system with minimum changes.

### ALGORITHMS FOR IMPLEMENTING DSM

- ① The central-server Algorithm.
- ② The Migration Algorithm.
- ③ The Read Replication Algorithm.
- ④ The Full Replication Algorithm.

- The Central Issue in the implementation of DSM are -

✓(A) How to keep track of the location of remote Data.

✓(B) How to Minimize Communication overhead when accessing Remote data.

✓(C) How to access concurrently remote data at several nodes (Computers).  
(Improve system performance)

① The central-server algorithm →

✓ Central server maintains all the shared data -

\* Read Request : returns data item.

\* Write Request : update data & returns acknowledgement message.

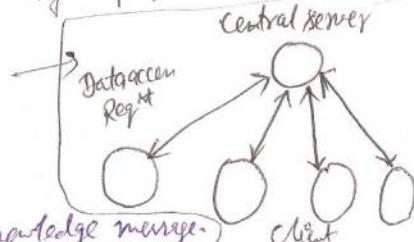


fig - The Central-server algo

\* A timeout is used to resend a request if acknowledgement fails.

\* Associated Sequence Numbers can be used to detect duplicate write requests.

\* If an Application's request to access shared data fails repeatedly, a failure condition is sent to the application.

✓ While the central server algo is easy to implement, the client server become a bottleneck.

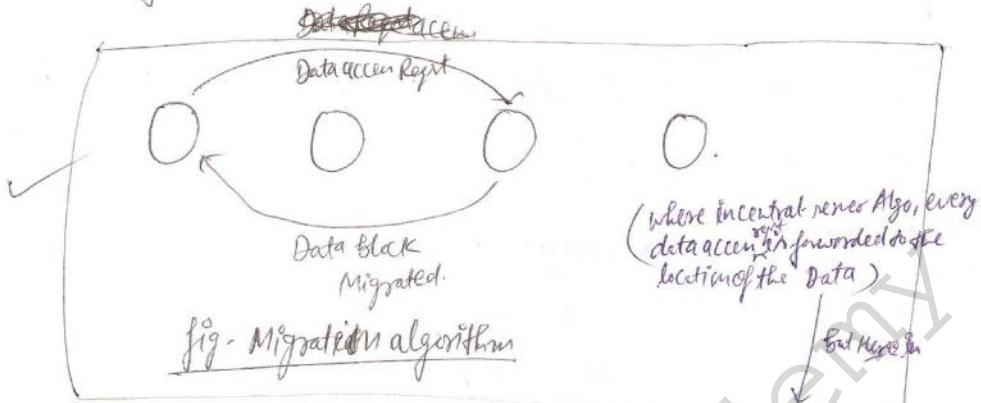
✓ To overcome this problem the shared data can be distributed over several servers.

✓ Issues → Performance & Reliability.

✓ Possible Solutions - Partition shared data over several servers.

- use a mapping function to distribute located data.

## ② The Migration Algorithms



### Operation -

- Ship (migrate) entire data object (page, block) containing data item to requesting location.
- Allow only one node to access a shared data at a time.

### Advantages

- Takes the advantage of locality of reference.
- DSM can be integrated with Virtual Memory provided by the operating system at individual nodes.
  - \* Make DSM page Multiple of VM page size.
  - \* A locally held shared memory page can be mapped into the VM page address space.
  - \* If page not local, fault handler migrates page & removes it from the address space at remote node.

### To locate a Remote Data object -

- \* use a location server.
- \* Maintain Rents at each node.
- \* Broadcast query.

### Disadvantages

- \* only one node can access a data object at a time.
- \* Thrashing can occur - to minimize it, set minimum time data object resides at a node.
- \* (causes system require more memory than it has)  
(Thrashing due to high page activity)

③ The Read - Replication Algorithm →

③

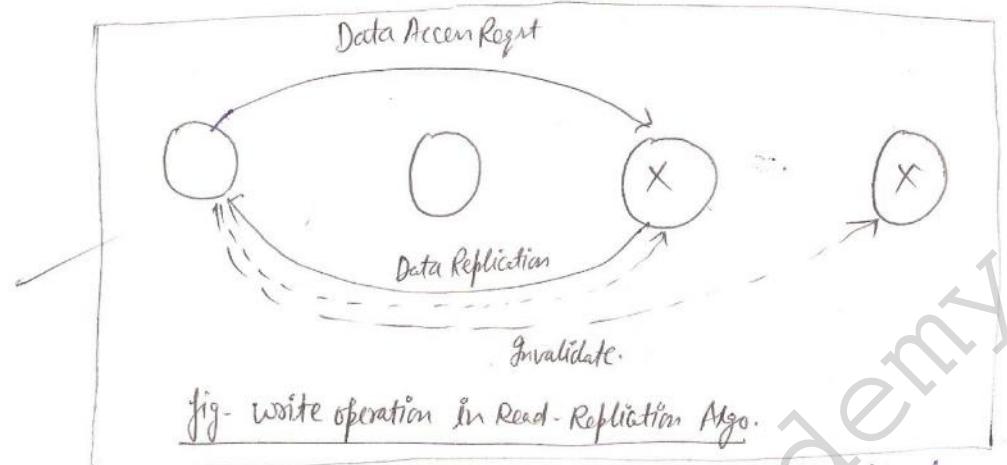


fig- write operation in Read-Replication Algo.

- In Previous approach, only processes on one node could access a shared data at any one moment. The read-replication algorithm extends the Migration algo by replicating data blocks & allowing multiple nodes to have read access or one node have read-write access (the Multiple readers-one writer protocol).

\*~~Readers~~

✓ Replicates Data objects to Multiple Nodes.

✓ DSM keeps track of location of Data objects.

✓ Multiple nodes can have read access or one node write access (Multiple Readers-one writer protocol)

✓ After a write, all copies are invalidated or updated the current value to maintain the consistency of the Shared Data Block.

✓ DSM has to keep track of locations of all Copies of data objects.

Example of implementation

(Integrate Shared Virtual Memory at Yale) ✓ \* IVY :> owner node of data object knows all nodes that have copies.

✓ \* PLUS :> Distributed link list tracks all nodes that have copies.

✓ Advantages

- The Read-Replication can lead to substantial performance improvements if the ratio of reads to write is large.

#### ④ The Full-Replication Algorithm ↗

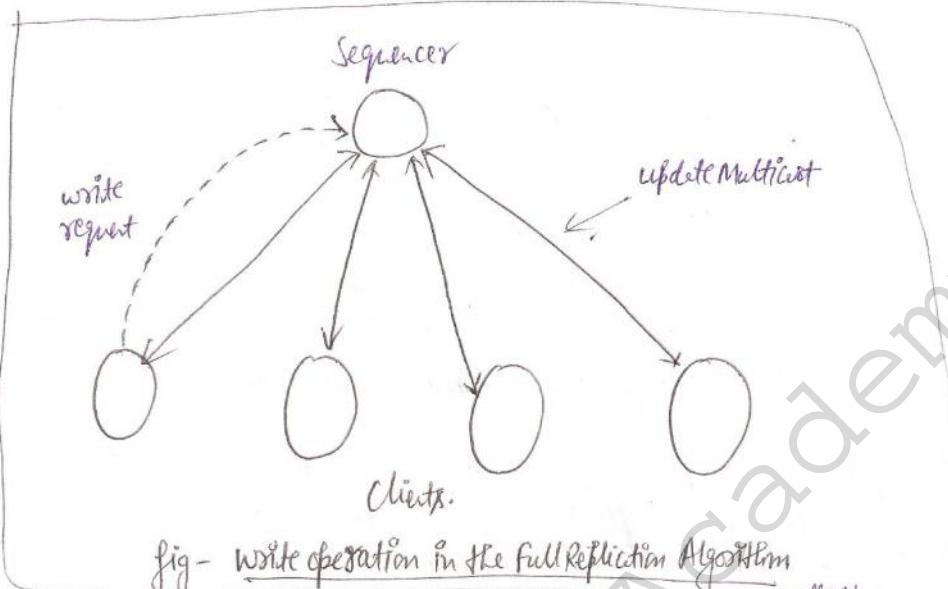


fig - Write operation in the Full Replication Algorithm

- ✓ This Algo is the Extension of the read-replication Algorithm. → Multiple Nodes can have both read & write access to Shared data blocks (Multiple readers, Multiple writers protocol). follow
- ✓ Because Many Node can write Shared data Concurrently, the access to the Shared data Must Be Controlled to Maintain its Consistency.
- ✓ Gives ~~multiple~~ Consistency of data for Multiple writers.
- ✓ All writers send to Sequencer.
  - ✓ The Sequencer will assign sequence No & sends write request to all the sites that have copies. (It multicasts the modification with the sequence No to all the nodes that have a copy of shared data item)
  - ✓ Each Node performs writer according to sequence Number.
    - A gap between the sequence No indicates that one or more modifications (written) missed.
    - Under such circumstances, Nodes will ask for retransmission of the modification it has missed.

-Two important issues to be considered in the Design of DSM system. They are important because the efficiency of DSM depends on the effectiveness of the size chosen for Granularity & the protocol used for Page replacement.

### Design Issues of Distributed Shared Memory-

1. Granularity- Refers to size of Shared Memory unit.
2. Page Replacement

#### 1. Granularity-

- A large page size for the shared memory unit will take advantage of the locality of reference.
  - By transferring large pages, less overhead is incurred due to paging activity and processing communication protocols.
- False sharing: occurs when two different data items, not shared but accessed by two different processes, are allocated to a single page.
  - More false sharing when the page size is large.
  - Smart compilers may partially solve the problem. However, if two processes share the same array, nothing can do about it.
- Another solution is to pre-fetch small pages.

#### 2. Page Replacement-

- When there is no free space in the memory, a page may need to be replaced
  - Traditionally, we use Least recently used (LRU).
- In DSM, LRU may need to be modified, since data may be accessed in different modes such as shared, private, read-only, writable, etc.
  - Private pages may be replaced before shared pages, as shared pages would have to be moved over the network, from their owner.
  - Read-only pages can simply be deleted as their owners have a copy.
- Once a page is selected for replacement, the DSM must ensure that the page is not lost forever.
  - One option is to swap the page onto disk.
  - Another option is to use reserved memory, wherein each node is responsible for certain portions of the global virtual space and reserves memory space for those portions.

#### IVY-

- IVY (Integrated Shared Virtual Memory at Yale) was implemented in the Apollo environment.
- The address space is divided into pages, with pages being spread over all the processors in the system
  - When a processor references an address that is not local, a trap occurs, and the DSM software fetches the page containing the address and restarts the faulting instruction, which now competes successfully.
  - Use replication to improve performance, especially read.
- Achieving sequential consistency by write-invalidation.

18

#### Sun Network File System-

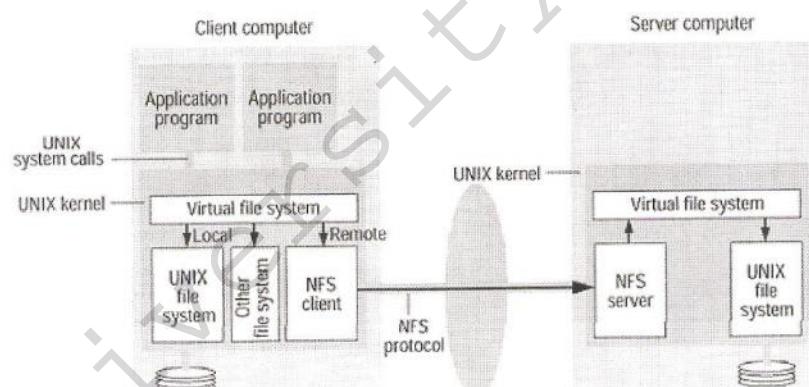
- Network file system (NFS) was developed by SUN Microsystem.
- It provides transparent access to remote files for client programs running on UNIX and other systems.
- Each computer has a client and server modules installed in its system kernel.
- Achieve a high level of support for hardware and OS heterogeneity.

### **Virtual File System-**

- The access transparency is achieved by a virtual file system (VFS) module.
- VFS keeps track of the file systems that are currently available both locally and remotely, and it passes each request to the appropriate system.
- V-node contains a reference to show whether a file is local or remote
  - If the file is local, the v-node contains a reference to the index of the local file (i-node).
  - If the file is remote, it contains the file handle of the remote file. The file handle contains fields uniquely identifying the file system type, the disk, the i-node number of the directory, and security information.

57

### **Network File System Architecture-**



**UNIT-4**

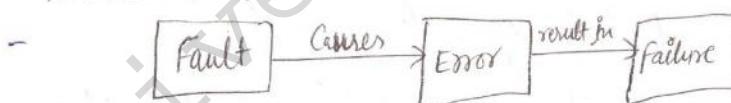
Introduction → Recovery in Computer Systems refers to restoring a system to its normal operational state. Recovery may be as simple as restarting a failed computer or restarting failed processes. Recovery is a very complicated process.

- In general, resources are allocated to executing processes in a computer.  
Ex- a process has memory allocated to it and a process may have locked shared resources, such as files & memory. Under such circumstances, if a process fails, it is imperative (requiring attention or action) that the resources allocated to the failed process be reclaimed (delivered from danger) so that they can be allocated to other processes.
- If a failed process has modified the database, then it is important that all the modifications made to the database by the failed process are undone.
- If a process has executed for some time before failure, it would be preferable to restart the process from the point of its failure & resume its execution.
- By restarting from the point of failure, the situation of having to reexecute the process from the beginning is avoided, which may be a time consuming & expensive operation.

- Distributed systems provide enhanced performance & increased availability.
- One way to realizing enhance performance is through the concurrent execution of many processes, which cooperate in performing a task.
- If one or more cooperating boxes fails, then the effect due to the interactions of the failed processes with the other processes must be evident, or every failed process would have to restart from an appropriate state.
- Increased the availability in DS is realized mainly through replication (e.g. data, processes, & HW, computers can be replicated).
- If site fails, copies of data stored at that site may miss updates, thus becoming inconsistent with the rest of the system when it becomes operational.
- Recovery in such cases involves the question of how not to expose the system to data inconsistencies & bring back the failed site to an up-to-date state lead to failure & the types of failure.

### Basic Concepts :-

- A system consists of a set of HW & SW components & is designed to provide a specific service.
- The components of a system may themselves be systems together with interrelationships.



Fault → is a defect within a system  
Error → is observed by a deviation from the expected behaviour of the system.

Failure → occurs when the system can no longer perform as required  
 (does not meet specification)

Fault Tolerance → is ability of system to provide a service, even in the presence of errors.

## QUESTION

- From the above definitions, it can be seen that an error is a manifestation of a fault in a system, which could lead to system failure. Failure recovery is a process that involves restoring an erroneous state (state which could lead to a system failure by sequence of valid state transition) to an error-free state.

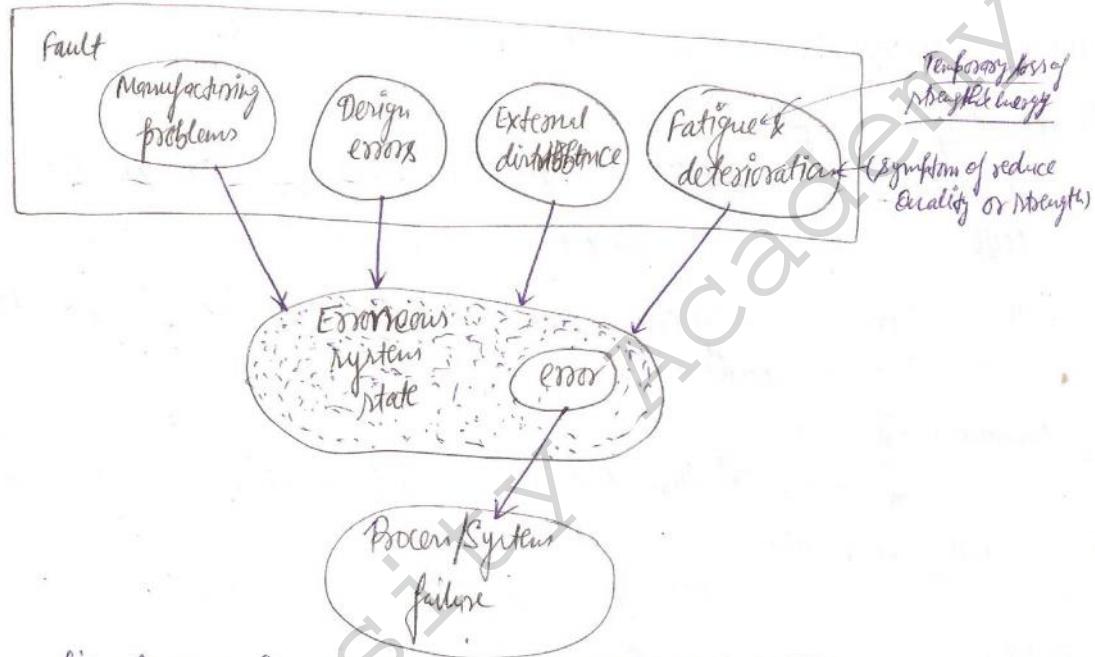


fig - An error is a Manifestation of fault & can lead to failure.

Classification of failure & failure in a computer system can be classified as follows-

- (1) Process failure → "The computation results in an incorrect outcome, the process causes the system state to deviate from specifications, the process may fail to progress. ex - (Deadlock)
  - (2) System failure → occurs when the processes fail to execute.
  - (3) Secondary Storage failure (SSF) → (ex - S/w error, H/w Problem)
  - (4) Communication Medium failures →
    - Tolerating Retries from SSF → (Archiving (Periodic Backup), Mirroring (continuous), Tolerating Comm. Medium failures → (Activity loggers))
    - ACK & retransmit
    - More complex fault tolerant Algs.
- "Recovering from system failure"

"Restart in predefined state"

"Reset the part of the state & dependent"

"Role back to before failure", "give up."
- ① Amnesia failure

② Partial-amnesia failure

③ Failure failure

④ Halting failure.

<sup>→ Partial or total loss of memory</sup>  
Ammena failure → occurs when a system restarts in a predefined state that does not depend upon the state of the system before its failure.

Partial - ammenia → occurs when a system restarts in a state wherein a part of state is the same as the state before the failure & the rest of the state is changed.  
(ex- file server has this type of failure)

Pause F → occurs when a system restarts in the same state it was before the failure.

Halting F → occurs when a crashed system never restarts.

### ③ Secondary storage failure

- A secondary storage failure is said to have occurred when the stored data (either some parts of it or in its entirety) cannot be accessed.
- This failure is usually caused by either parity error, head crash, or dust particles settled on the medium.

### ④ Communication Medium failures

- A communication medium failure occurs when a node cannot communicate with another operational node in the network.  
(ex- failure of switching nodes)

## BACKWARD AND FORWARD ERROR RECOVERY :-

- Error is that part of the state that differs from its intended value & can lead to a system failure, & failure recovery is a process that involves restoring an erroneous state to an error-free state.
- There are two approaches for restoring an erroneous state to an error-free state.  
✓ If the nature of errors & damages caused by faults can be completely & accurately accessed, then it is possible to remove those errors in the process's (system's) state & enable the process (system) to move forward. This technique is known as forward-error recovery. (using check point)  
(correct the error)

- It is not possible to foresee the nature of faults & to remove all the errors in the process (system's) state, then the process (system's) state can be restored to a previous error-free state of the process (system). This technique is known as Backward-error recovery. (Correct the error)  
 (use the check point) BR > FR

The major problem associated with the backward-error recovery approach are →

✓ Performance Penalty → The overhead to restore a process (system) state to a prior state can be quite high.

✓ There is no guarantee that faults will not occur again when processing begins from prior state.

✓ Some component of the system state may be unrecoverable.  
 ex - (Cash dispensed at an ATM can't be recovered)

### \* BACKWARD ERROR RECOVERY →

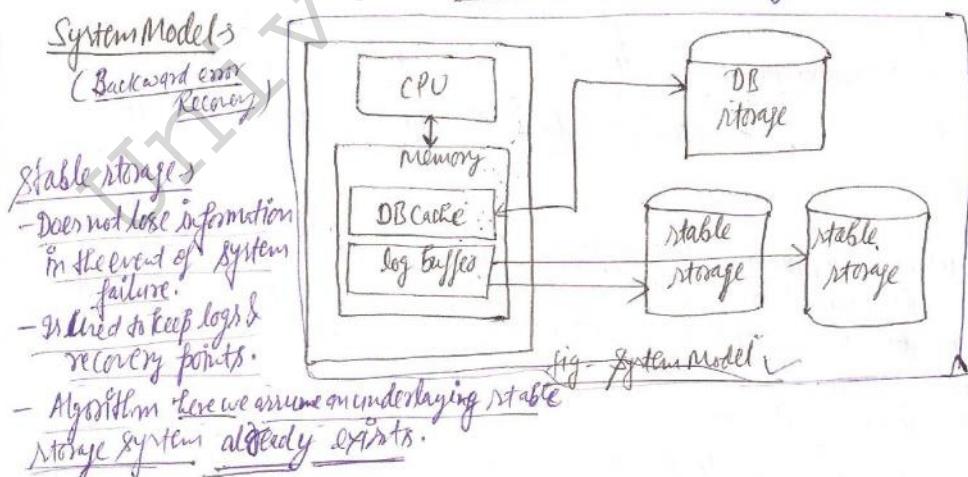


- It is based on Recovery Points.

- Two approaches are - (for implementing the backward Recovery)

✓ ① Operation-based recovery. ✓

✓ ② State-based Recovery. ✓



## Two Approaches to fault Tolerance :-

### ① Operation based Approach :-

- Record a log (audit trail) of the operations performed.

- Restore previous state by reversing steps.

### ② State based :-

- Record a snapshot of the state (checkpoint).

- Restore state by reloading snapshot (rollback).

- use of shadow pages can reduce size of checkpoints.

Updating in Place :- under this scheme every update (write)

operation to an object updates the object & results in a log to be recorded in a stable storage which has enough information to completely undo & redo the opn.  
(Name, location of object)

- The info recorded includes -

✓ the name of the object. (Do)

✓ the old state of the object (used for UNDO).

✓ the new state of the object (used for REDO).

(old state/value of object)

(New state/value of object)

operation which

DO - Does the action write a log

UNDO - gives log record written by

Do opn. Performed

By Do opn.

REDO -

DISPLAY -

Specified by the

Do opn.

Display the log

record.

② Write ahead Log Protocol :- In this a recoverable update opn is

implemented by the following opn -

✓ Update an object only after the undo log is recorded.

✓ before committing the update, redo & undo log are recorded.

- Both undo & redo opn performed properly in both updating in place & write ahead log protocol.

③ Shadow Pages :- whenever a process wants to modify an object the page containing the object is

Duplicated & is maintained on the stable storage.

✓ The ~~process~~ process only updates one copy. the unmodified copy is called the shadow page.

✓ If the process fails, the modified copy is discarded & the shadow page is used.

✓ If the process commits the shadow page is discarded.

## Recovery in Concurrent Systems $\rightarrow$ [IMP]

(4)

- In concurrent systems, several processes cooperate by exchanging information to accomplish a task. The information exchange can be through a shared memory in the case of shared memory machines (e.g. multiprocessor systems) or through messages in the case of a distributed system.
- In such systems, if one of the cooperating processes fails & resumes execution from a recovery point, then the effects it has caused at other processes due to the information it has exchanged with them after establishing the recovery point will have to be undone.
- To undo the effects caused by a failed process at an active process, the active process must also rollback to an earlier state.
- Thus in concurrent systems, all cooperating processes need to establish recovery points.
- Rolling back processes in concurrent systems is more difficult than in the case of single process.

The following discussion ~~discusses~~ shows how the rolling back of processes can cause further problems.

- ① Orphan Messages & the Domino Effect.
- ② Lost Messages.
- ③ Problem of livelocks.

### ① Orphan Messages & the Domino effect $\rightarrow$

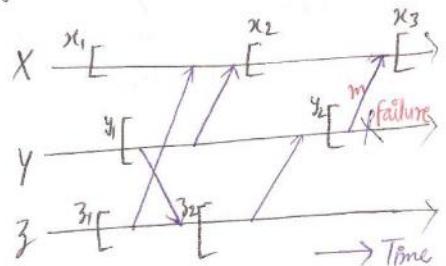
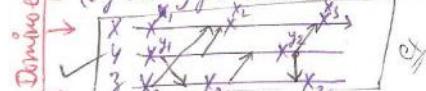


fig - Domino effect.

X, Y, Z  $\rightarrow$  Processes (Cooperate by exchanging messages)  
[  $\rightarrow$  Marks as a recovery point to which a process can be rolled back in the event of a failure.

Orphan message - whose receiving event is recorded in the checkpoint, but its sending event is not.

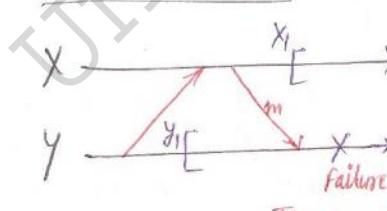
Domino effect  $\rightarrow$  one rollback causes others to rollback  
(e.g. when Y goes back to checkpoint  $y_2$ )



- If Process X is to be Rolled back, it can be rolled back to the recovery point  $x_3$  without affecting any other processes.
- Let, Y fails after sending message m & is rolled back to  $y_2$ . In this case, the receipt of m is recorded in  $x_3$ , but the sending of m is not recorded in  $y_2$ .
- Now we have a situation where X has received message m from Y, but Y has no record of sending it, which corresponds to an Inconsistent State.
- Under such circumstances, m is referred to as an orphan message & process X must also rollback.
- X must rollback because Y interacted with X after establishing its recovery point  $y_2$ . [Important]
- When Y is rolled back to  $y_2$ , the event that is responsible for the interaction is undone.
- Therefore, all the effects at X caused by the interaction must also be UNDONE.
- This can be achieved by rolling back X to recovery point  $x_2$ .
- Likewise, it can be seen that, if Z is rolled back, all three processes must rollback to their very first recovery points, namely  $x_1, y_1$  &  $z_1$ .
- This effect, where rolling back one process causes one or more other processes to roll back, is known as the domino effect, & the orphan messages are the cause.

## ② Lost Message $\rightarrow$

- A message whose sending event is recorded, but its receiving event is not recorded.



- Y fails & rolls back to  $y_1$ , & then X is in a state where it sent m but Y will never receive it.
- This situation can also happen if the communication channel is not reliable.

fig - message loss due to rollback recovery.

③ Livelock → Livelock is a situation in which a single failure can cause an infinite No. of rollbacks, preventing the system from making progress.

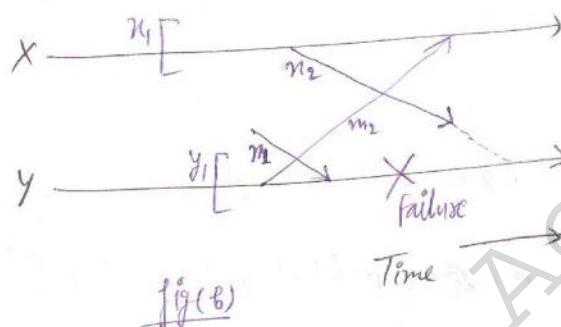
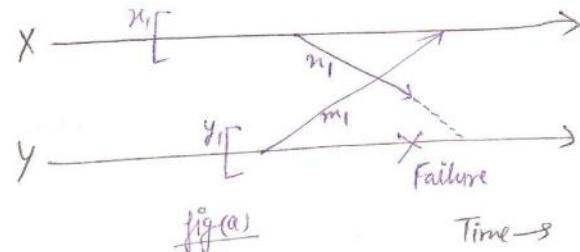


Fig - Livelock ✓

- It is a situation in the Distributed System.

Consistent Checkpoints → Strongly Consistent set of checkpoints (in distributed system)

- A process saves its local state on the stable storage, which is called a local checkpoint.
- The process of saving local states is called local checkpointing.
- All the local checkpoints, one from each node, collectively form a global checkpoint.

① Strongly Consistent set of checkpoints →

✓ The domino effect is caused by orphan messages, which themselves are due to rollbacks.

✓ To overcome the domino effect, a set of local checkpoints is needed (one for each process in the set) such that no info flow takes place (no orphan message) between any pair of processes in the set, as well as between any process in the set & any process outside the set during the interval formed by the checkpoints.

✓ A set of checkpoints is known as a recovery line or a strongly consistent set of checkpoints.

✓ A Global Check Point is a strongly Concurrent set of checkpoints if there is no orphan message & no lost message.

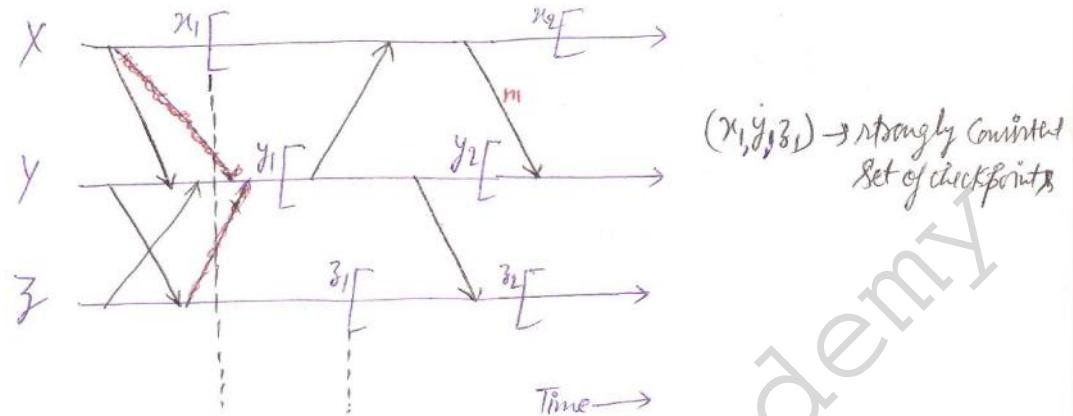


fig- Concurrent set of checkpoints

- A Global checkpoint is a strongly Concurrent set of checkpoints if there is no orphan message.  $\{n_2, y_1, z_1\}$
  - A concurrent set of checkpoints, set is similar to a concurrent global state in that it requires that each message recorded as received in a checkpoint (state) should also be recorded as sent in another checkpoint (state).
- \* A simple method for taking a concurrent set of checkpoints is
- ✓ Every process takes a checkpoints after sending every message, the set of the most recent checkpoints is always concurrent.
- ✓ However, it has a high overhead(expensive) that taking a checkpoint after every message has sent.
- ✓ So, one may attempt to reduce the overhead of in the above method by taking a ~~one~~ checkpoint after every  $K$  ( $K > 1$ ) message sent.
- ✓ This method suffers from domino effect.

- |                                                                                                                      |
|----------------------------------------------------------------------------------------------------------------------|
| ✓ Assumption - Checkpoint, send/receive are atomic.                                                                  |
| ✓ Take a checkpoint after every message.                                                                             |
| ✓ the set of the most recent checkpoint is always concurrent (why? Is it strongly concurrent?)                       |
| ✓ what is the main problem with this approach? If a checkpoint after every $K$ message sent? Is it still concurrent? |

Synchronous checkpoint & recovery  $\xrightarrow{\text{IMP}}$  (Koo & Toueg)  $\xleftarrow{① \text{The Checkpoint Algo.}} \xleftarrow{② \text{The Rollback Recovery Algo.}}$  ⑥

✓ A checkpointing & recovery technique proposed by Koo & Toueg that takes a consistent set of checkpoints & avoids livelock problems during recovery.

✓ The algorithm's approach is said to be synchronous, as the processes involved coordinate their local checkpoint actions such that the set of all recent checkpoints in the system is guaranteed to be consistent.

① The Checkpoint Algorithm  $\Rightarrow$  (Koo & Toueg)

The Checkpoint Algorithm assumes the following characteristics for the distributed system.

✓ Processors communicate by exchanging messages through channels.

✓ Channels are FIFO in nature. End-to-End protocols cope with message loss due to rollback recovery (SWP) + <sup>have</sup> <sub>(that is not yet sent)</sub> recording message in stable storage before sending them.

✓ Communication failures do not partition the network.

\* This Check Point Algorithm uses 2 kind of checkpoints on stable storage.

① Permanent  $\Rightarrow$  (A local checkpoint at a process & is a part of a consistent Global checkpoint)

② Tentative  $\Rightarrow$  (A temporary checkpoint set is made a permanent checkpoint on the successful termination of the checkpoint algorithm).

This Algorithm has 2 Phase  $\Rightarrow$

✓ ① Synchronous checkpoint Phase-1  $\Rightarrow$

(Initiator)  $\Rightarrow$  An initiating process  $P_i$  takes a tentative checkpoint & requests all the processes to take tentative checkpoints.

✓ Other Processes: can respond 'yes' or 'no'

Ans  $\Rightarrow$  Initiator  $\left\{ \begin{array}{l} \text{Decides to make checkpoints permanent if every one has} \\ \text{responded 'yes'} \end{array} \right\}$  otherwise  $P_i$  decides that all the tentative checkpoints should be discarded.

## ② Synchronous Checkpoint Phase 2 :-

- Initiator  $P_i$  informs all the processes of Phase-I (Decision at Phase-1) .  
(Commit or abort ~~according to~~ checkpoint)
- Others act accordingly.  
(Therefore either all or None of the processes takes permanent checkpoint.)
- Both tentative checkpoint & Commit/abort of checkpoint process must hold back messages.

### Synchronous checkpoint : Properties -

- ✓ All or none of the ~~processes~~ processes take permanent checkpoints.
- ✓ There is no record of a message being received but not sent.
- ✓ Checkpoint may be taken unnecessarily. e.g. ~~(ex)~~  $\rightarrow$  [Taking a checkpoint is an expensive operation]

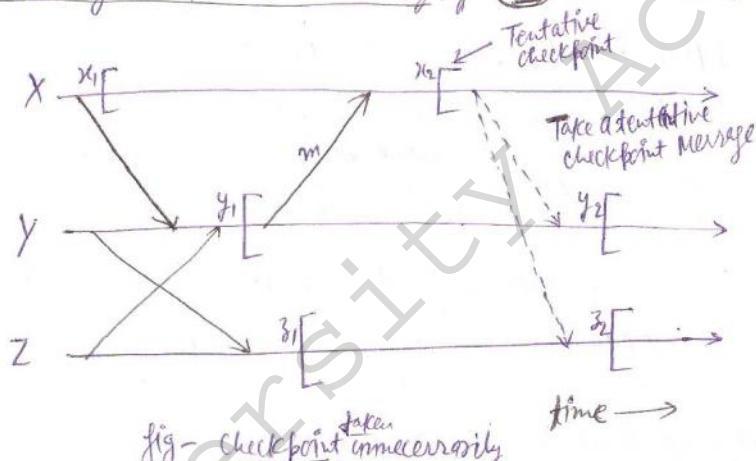


fig - Checkpoint <sup>taken</sup> unnecessarily

let,

- $\{x_1, y_1, z_1\} \leftarrow$  Committent set of checkpoint
- $X$  decides to initiate the checkpoint after receiving message  $M$ .
- It takes the tentative checkpoint  $x_2$  & sends "Take tentative checkpoint message" to process  $Y$  &  $Z$ , causing  $Y$  &  $Z$  to take checkpoint  $y_2$  &  $z_2$  respectively.
- Now, -  $\{x_2, y_2, z_2\}$  form a Committent set of checkpoint.
- $\{x_2, y_2, z_1\}$  can also form a Committent set of checkpoint.

Can this unnecessary checkpoints be avoided? A scheme is described as -

(7)

Main Ideas (Scheme)

\* Record all messages sent & received after the last checkpoint.

- last\_recv(x,y) -
- first\_send(x,y) -

\* When X request Y to take a tentative checkpoint, X send the last message received from Y with the request.

\* Y takes tentative checkpoints only if the last message received by X from Y was sent after Y sent the first message, after the last checkpoint  
 $(\text{last\_recv}(x,y) \geq \text{first\_send}(y,x))$

→ Page 310

\* When a process takes a checkpoints, it will ask all other processes that sent messages to the process to take checkpoint.

## ② The Rollback Recovery Algorithm :-

- It assumes that single process initiates the algorithm, as opposed to several processes concurrently invoking it to rollback & recover.
- It also assumes that checkpoint & rollback recovery algorithm are not currently involved.
- The rollback algo has two phases -  $\xrightarrow{\text{P-1}}$   $\xrightarrow{\text{P-2}}$

### ① Rollback Recovery : Phase(1) →

- (P-1)
- Initiator : check whether all processes are willing to restart from the checkpoints.  
→ If Yes → All processes restart otherwise P<sub>i</sub> will decide that the process should continue with their normal operations.
  - (A process may reply "no" to a restart reqd if it is already participating in a checkpointing or a recovering process initiated by some other process.)

### ② Rollback Recovery : Phase(2) →

- Initiator : Propagate go/no-go decision to all processes.
- others : carry out the decision of the initiator.

Between Request to Rollback & decision, no one sends other messages.

### Rollback Recovery Properties

- All or none the processes restart from checkpoints.
- After rollback, all process resume a consistent state.

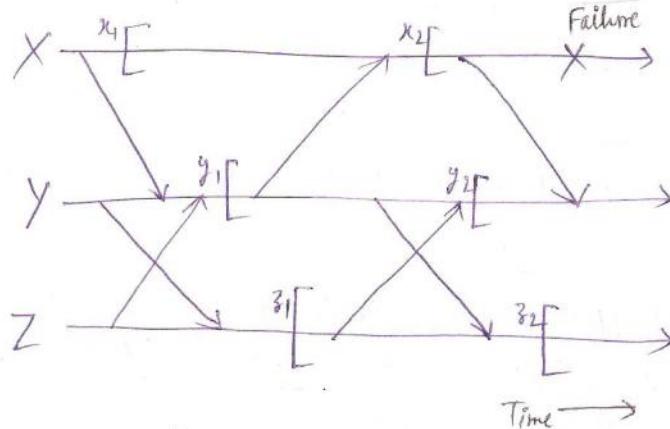


fig - Unnecessary rollback (problem)

- In the event of failure of process X, would require processes X, Y, Z to restart from checkpoints  $x_2, y_2, z_2$  respectively.
- However, that process Z need not have rolled back as there was no interaction betw Z & the other two processes.

To minimize the No of process (rollbacks), A scheme is defined as -

Main Idea - for Any two Processes X, Y (m, m<sub>1</sub>, m<sub>2</sub>)

- last\\_sent {x, y, z}  $\rightarrow$  (X restarts, Before X takes its latent permanent checkpoint) ✓
- last\\_recv {y, z, x}  $\rightarrow$  (Y will restart from its permanent checkpoint only if ~~the condition holds~~ the condition holds).
- last\\_recv {y, z}  $\rightarrow$  last\\_sent {x, y, z} when X restarts to checkpoint

Disadvantage of Synchronous Approach :-

- Checkpoint Algorithm generates merge traffic.

- Synchronization delays are introduced.

These cost may seems high if failed set of checkpoints are unlikely.

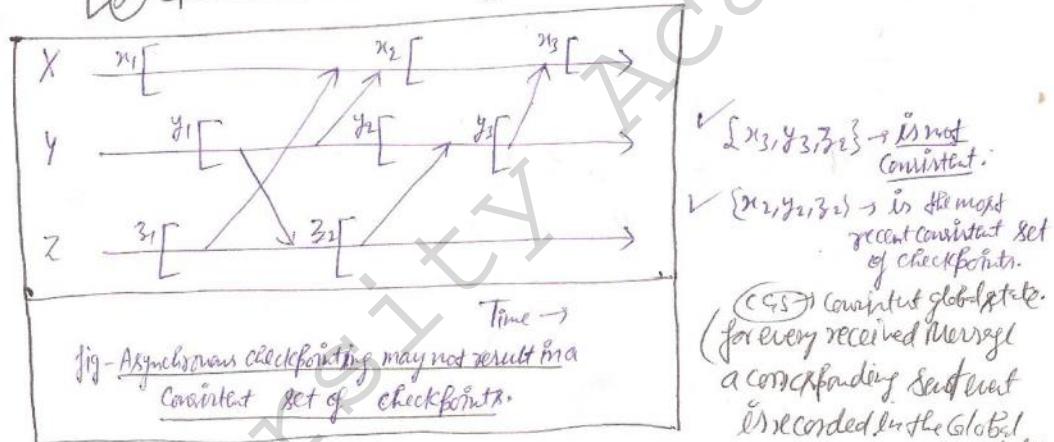
(X - Regt. Y to restart from the permanent checkpoint) It's not

(L  $\rightarrow$  Smallest level)  
(T  $\rightarrow$  largest level)  $\rightarrow$  Process's level

## Asynchronous Checkpointing And Recovery :-

(8)

- Takes multiple local checkpoints independently, (at each processor) without any synchronization among the processors. Because of the absence of synchronization there is no guarantee that a set of local checkpoints taken will be a consistent set of checkpoints.
- After a failure, try to find a consistent set of checkpoints among those that have been taken recently.
  - All incoming message sets' local checkpoints are logged. (Message received can be logged in 2 ways) →
    - ① optimistic approach: log each message before processing.
    - ② optimistic approach: Buffer messages & log in batches.



## A Scheme for Asynchronous Checkpointing & Recovery :- [Jiang & S. Venkatesan]

Algorithm makes the following assumptions →

- \* Communication channels are reliable.
- \* Communication channels are FIFO.
- \* Communication channels have no buffer limits (size).
- \* Message transmission delay is bounded. [finite].
- \* Underlying computation is event-driven, but locally timestamped events. Each event consists of the following - waiting for a message, process the message, change process state, & send a no of messages.

### Asynchronous Check Point :-

Two types of log storage are assumed to be available for logging in the system.

- ① Volatile log - (accessing volatile log takes less time than accessing stable log)
- ② Stable log - but contents of the volatile log are lost if corresponding processor fails  
(the contents of volatile log are periodically flushed to the stable storage & cleared.)

### Basic Idea :-

- \* At each event, a triplet  $\{s, m, \text{msgs\_sent}\}$  is put in the log.

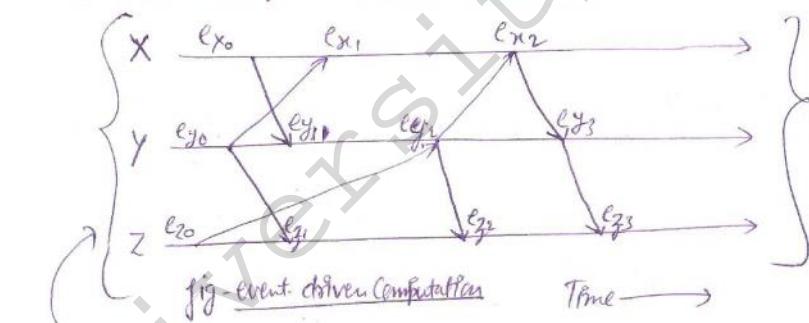
$\left\{ \begin{array}{l} s \rightarrow \text{state} \\ m \rightarrow \text{message causing the event} \\ \text{msgs\_sent} \rightarrow \text{in the set of message sent.} \end{array} \right\}$

- \* Two data structures used:

✓ RCVD(i, j, checkpoint) - the No of message received by processor  $i$  from processor  $j$  at checkpoint.

✓ SENT(i, j, checkpoint) - the No of message sent from  $i$  to  $j$  at checkpoint.

- \* Use the message send/receive counts to determine the point to rollback.



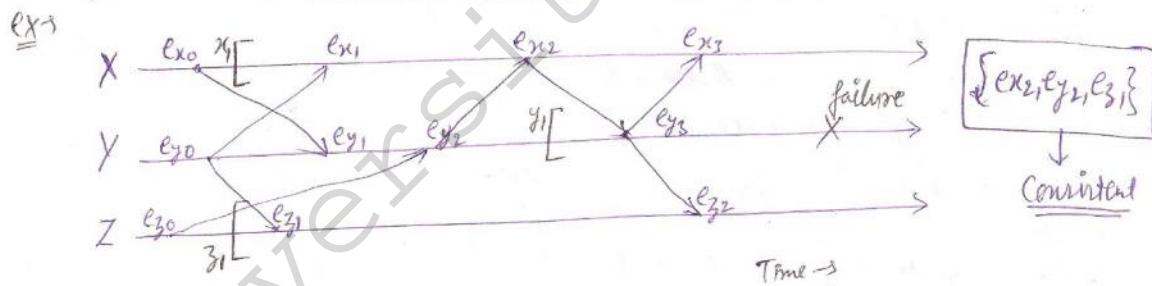
- If  $Y$  rolls back to a state corresponding to  $e_{y1}$ , then according to this state  $Y$  has sent only one message to  $X$ . According to  $X$ 's state, however, it has received two messages from  $Y$  thus far. Therefore,  $X$  has to roll back to a state preceding  $e_{x2}$  to be consistent with  $Y$ 's state. For the similar reason,  $Z$  will also have to roll back.

meanval(e<sub>Y</sub>)

Algorithm<sup>9</sup>

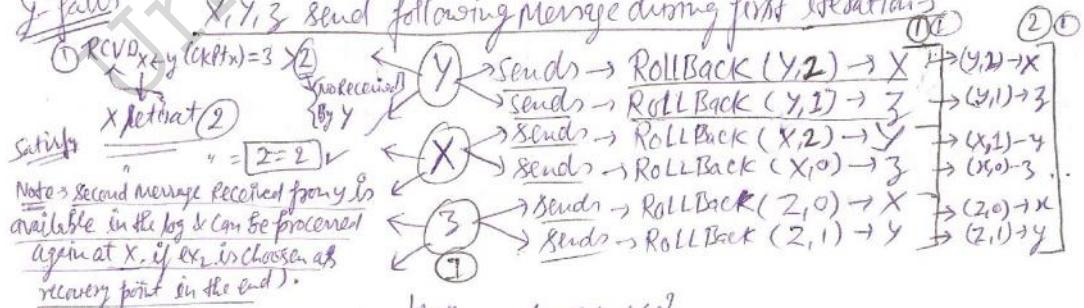
- ✓ All Node will be running the same recovering algorithm at each processor<sup>9</sup>
- ✓ If  $i$  is a processor that is recovering from failure, [checkpoint = the latest event logged in the stable storage.]
- \* else checkpoint = latest event that took place.
- ✓ for  $k=1$  to  $N$  do - (N → No of processors in the System)
- ✓ send [ROLLBACK( $i$ , SENT( $i, j$ , checkpoint))] to all neighbours  $j$ .
- ✓ wait for [ROLLBACK messages] from all neighbours
- ✓ for every ROLLBACK( $j, c$ ) received - (Preference of  
Log then Merge)  
- if  $RCVD(i, j, checkpoint) > c$  then  
\* find the latest event  $e$  such that  $RCVD(i, j, e) = c$   
\* checkpoint =  $e$

In each iteration, at least one processor will rollback to its final recovery point unless current recovery point is consistent.



If initially X, Y, Z set their Check Point at  $\rightarrow ex_3, ey_2, ez_2$

Y fails  
X, Y, Z send following message during first iteration



$$\begin{aligned} \text{① } RCVD_{Z \leftarrow Y}(\text{checkpoint}) &= 2 > 1, \\ \text{set } e_3 &= 1 \leq 1 \end{aligned}$$

## CheckPointing for Distributed Database Systems

(9)

- In this we discuss on a technique for taking checkpoints in a ~~sys~~ distributed Database System (DDBS) where a set of data objects is partitioned among several sites.

- A Checkpointing scheme for a DDBS should meet the following two basic objectives-

① Checkpoints should be highly desirable.

② All sites should take local checkpoints recording the state of the local database.

Issues →

① How sites decide or agree upon update of what transactions are to be included in their checkpoints.

② How each site can take a local checkpoint in a noninterfering way.

Son & Agrawala's Algorithm for checkpoint in a DDBS: → (Noninterfering & takes globally concurrent checkpoints).

Assumptions made by algo are→

- \* The basic unit of user activity is a transaction.
- \* The Transactions follow some Concurrency Control Protocol.
- \* Lamport logical clocks are used to associate each transaction with a timestamp.
- \* Site failure are detectable either by Net protocols or by time out mechanisms.
- \* Network partitioning never occurs.

Basic Idea →

- All participating sites agree upon a special timestamp known as the Global checkpoint Number (GCPN).
- The update of the transactions, which have timestamp  $\leq$  GCPN, are included in the checkpoint.
- These transactions are called before-checkpoint transaction (BCPTs).
- The update of the transaction which have timestamp  $>$  GCPN are not included in the checkpoint.
- These transactions are called after-checkpoint transaction (ACPTs).

→

- To avoid interfering with the normal operations while checkpointing is in progress, each site maintains multiple versions of data items in volatile storage that are being updated by ACPTs.
- Thus the state of the database is not disturbed once all the BCPTs terminate (at which time the database is consistent) until checkpointing completes.

#### DATA STRUCTURE:

This Algo requires that each site maintains following variables-

- \* LC: The local clock maintained as per Lamport's logical clock rules
- \* LCPN: A No Determined Locally for the current checkpoint.

#### Algorithm:

- The Algo is initiated by a special process known as the checkpoint coordinator (CC).
- It takes a consistent set of checkpoints with the help of processes known as checkpoint subcoordinators (CS), running at every participating site.

#### Phase-1 :-

- (1) The checkpoint coordinator broadcasts a checkpoint-Request message with local timestamp  $LC_{cc}$ .

$$(2) LCPN_{cc} = LC_{cc}$$

$$(3) CONVERT_{cc} = \text{false}$$
 (use of convert will become clear later)

- (4) The checkpoint coordinator waits for replies (obtaining LCPNs) from all the subordinate sites.

At all the checkpoint subcoordinators (CS) sites →

- (1) on receiving a checkpoint-Request message, a site m, updates its local clock as follows-

$$LC_m = \max\{LC_m, LC_{cc} + 1\}$$

$$(2) LCPN_m = LC_m$$

(3) site m informs LCPN<sub>m</sub> to the checkpoint coordinator.

$$(4) CONVERT_m = \text{false}$$

- (5) site m marks all the transactions with timestamp  $\neq LCPN_m$  as BCPT, & mark rest of the transactions as temporary - ACPT.

## Phase-2

At the checkpoint coordinator site :-

- Once all the replies for the checkpoint request messages have been received, the coordinator broadcast GCPN, which is -

$$GCPN = \boxed{\max \{ LCPN_1, LCPN_2, \dots, LCPN_n \}}$$

( $n \Rightarrow$  The No. of sites in the system)

At all sites -

- ① On receiving GCPN, a site m marks all temporary ACPTs which satisfy the following condition as BCPT -

$$\boxed{LCPN_m < transaction's timestamp \leq GCPN}$$

the update of those transactions, newly converted as BCPTs, are also included in the checkpoint.

- ② Convert<sub>m</sub> = true, when Convert<sub>m</sub> = true, ~~site m takes a local checkpoint~~. It indicates that GCPN is known and all BCPTs have been identified.
- ③ When all BCPTs terminates & Convert<sub>m</sub> = true, site m takes a local checkpoint by saving the state of data object.
- ④ When all the local checkpoint is taken, the database is updated with the committed temporary versions & then the committed temporary versions are deleted.

Note - if a site m receives a new "initiate transaction" message for a new transaction whose timestamp is  $\leq GCPN_m$  & the site m has already executed steps 1 & 2 of Phase 2, then site m reject the "initial transaction".

## Recovery in Replicated Distributed Database Systems (RDDBS)

(10)

- To enhance the performance & availability, a distributed database system is replicated where copies of data objects are stored at different sites.
- Such a system is known as RDDBS. In RDDBS, transactions are allowed to continue despite one or more site failure as long as one copy of the database is available.
- The copies of the database at the failed site may miss some updates while the sites are not operational. These copies will be inconsistent with the copies at the operating sites.
- The goal of recovery algorithm in RDDBS is to hide such inconsistencies from user transactions, bring the copies at recovering sites up-to-date with respect to the rest of the copies, & enable the recovering sites to start processing transactions as soon as possible.
- Two approaches have been proposed to recover failed sites-
  - ① Merge spoolers- Used to save all the updates ~~in log~~ directed towards failed sites.
  - ② Copies transactions- Read the up-to-date copies at the operational sites & update the copies at recovering sites. It runs concurrently with user transaction.
- The recovery algorithm should guarantee that-
  - ① The out-of-date replicas are not accessible to user transactions.
  - ② Once the out-of-date replicas are made up-to-date by copies transaction, they are also updated along the other copies by the user transactions.

## An Algorithm for Site Recovery $\Rightarrow$ [Bharagava & Ruan] $\Rightarrow$ Based on Copy transaction

### System Model $\rightarrow$

- The Database is assumed to be manipulated through transactions where access to the Database is controlled by Concurrency Control algorithm.
- Transactions either run to completion or have no effect on the database.
- The semantics of read & write operations on the Database are such that a read operation will read from any available copy & write operation updates all the available copies.
- All the out-of-date copies in the database are assumed to be marked as "unreadable".
- we also assume that database is fully replicated (every site has a copy of database).
- A site may be any one of the following states -

① operational/up  $\Rightarrow$  The site is operating normally & user transactions are accepted.

② Recovering  $\Rightarrow$  The Recovery is still in progress at the site & the site is not ready to accept user transaction.

③ Down  $\Rightarrow$  No RDBMS activity can be performed at the site.

④ Non operational  $\Rightarrow$  The site's state is either recovering or down.

- The operational session of a site is a time period in which the site is up. each operational session of a site is designated with a session no (an integer) which is unique in the site ~~history~~ history, but not necessarily unique systemwide.
- The session No are stored on nonvolatile storage so that a recovering site can use an appropriate new session No.

## DATA STRUCTURES

(12)

- Each site  $K$  maintains the following 2 data structures-

① The session No of site  $K$  is maintained in a variable  $AS_K$ .  $AS_K$  is set to zero when site  $K$  is nonoperational.

②  $PS_K$  is a vector of size  $n$ , where  $n$  is the No of sites in the system.

$PS_K[i]$  is the session no of site  $i$  as known to site  $K$ .  $PS_K[i]$  is set to zero whenever  $K$  learns that site  $i$  is down or some other site informs  $K$  that site  $i$  is down.

How the System functions under the Normal, failed & recovery Conditions

### User transaction

- Each request that originates at a site  $i$  for reading or writing a data item at site  $K$  carries  $PS_i[K]$ .
- If  $PS_i[K] \neq AS_K$  or  $AS_K = 0$ , then the request is rejected by site  $K$ . Otherwise there are 3 Possibilities -
  - ① the data item is readable - the request is processed at site  $K$ .
  - ② the data item is marked unreadable & the operation is a write operation - (the data item is modified & will be marked readable when the transaction commits.)
  - ③ the data item is marked unreadable & the operation is a read operation - (a copies transaction is initiated by site  $K$ ).

### Copies transaction

- Copies transaction <sup>may be</sup> initiated for all the data items marked as unreadable when a site starts recovering.
- On the other hand, a copies transaction may be initiated on a demand basis that is, whenever a read operation is received for individual data items marked unreadable.

### ③ Control transaction :-

- Control transactions are special transactions that update AS & PS at all sites (including any recovering sites).
- When a recovering site (say K) decides that it is ready to change its state from recovering to operational, it initiates a type-1 Control transaction.
- Type-1 transaction performs the following operations:-
  - It reads  $PS_i$  from some reachable site  $i$  & refreshes  $PS_K$ .
  - It chooses a new session No, sets  $PS_K[K]$  to this new session no, & writes  $PS_K[K]$  to all  $PS_i[K]$  where  $PS_i[i] \neq 0$ .  
[all the sites perceived up by site K]
- When a site discovers that one or more sites are down, it initiates a type-2 Control transaction. Ex- if site K learns that site m & n are down, then it initiates a type-2 Control transaction which performs following ops:-
  - It sets  $PS_K[m]$  &  $PS_K[n]$  to zero.
  - For all  $i$  such that  $PS_K[i] \neq 0$ , it sets  $PS_i[m]$  &  $PS_i[n]$  to zero.

### The SITE RECOVERY PROCEDURE :-

- When a site K restarts from failure, the recovery procedure at site K performs the following steps:-
  - It sets  $ASK$  to zero. That is site K is recovering & is not ready to accept user transactions.
  - It marks all the copies of data item unreadable.
  - It initiates a type-1 control transaction.
  - If the control transaction of step-3 successfully terminates, then the site copies the new session No from  $PS_K[K]$  to  $ASK$ . [one  $ASK$  is  $\neq 0$ , the site is ready to accept user transaction].
  - If step 3 fails because of discovering that another site has failed, site K initiates type-2 control transaction to exclude the newly failed site & then restarts from step-3.

System can be designed to be fault-tolerant in two ways →

- ① Mark failures, continue to perform, in degraded mode.

\* Basic mechanism is redundancy.

- ② or exhibit a well defined failure behaviour: may not perform the specified function, but facilitate actions suitable for recovery.

\* e.g. changes made to a database by a transaction are made visible to other transactions only if the transaction successfully commits.

Issues in Fault Tolerance → (we must study the implication of certain types of failures)

- ① - Process Death → { Server/client processes must inform each other before dying after sending a signal to other }.

{ Must recover resources. }

{ Must undo resultant ~~deadlock~~ blockages. } When a process dies it is important to free the resources allocated to that process so they may be recovered, otherwise they may be permanently lost.

- ② - Machine failure -

→ In the case of M/c failure all the processes running at the M/c will die.

→ The only difference between Process death & M/c failure, lies in how the failure is detected.

- ③ - Network failure -

{ A communication link failure can partition a network into subnets, making it impossible for a m/c to communicate with another M/c in a different subnet. }

\* usually a fault-tolerant algo assumes that M/c may be operating & processes on that M/c are active.

Atomic Actions →

- appear to other processes as if they were -

- indivisible.

- instantaneous.

- Serializability: the effect is the same as if atomic actions are interleaved, not concurrent.

Transaction →

- A set of sequence of actions that is implemented as if atomic.

- Commit → Guarantee transaction will be completed, so its effect is permanent.

\* Implemented by Commit protocols.

- Abort → guarantee to back out of the transaction, so none of its effect persist.

To back out  
Every write has  
Protocol  
Shadow copy

\* deadlock

\* timeout

\* protection violation

\* input data error

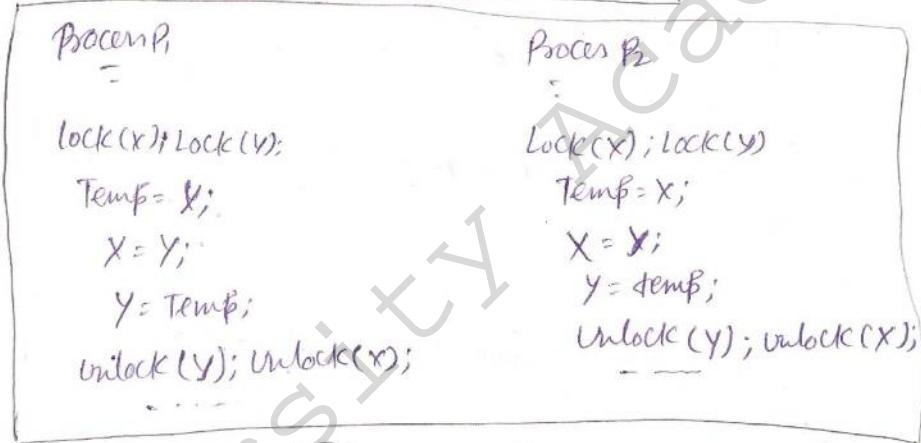
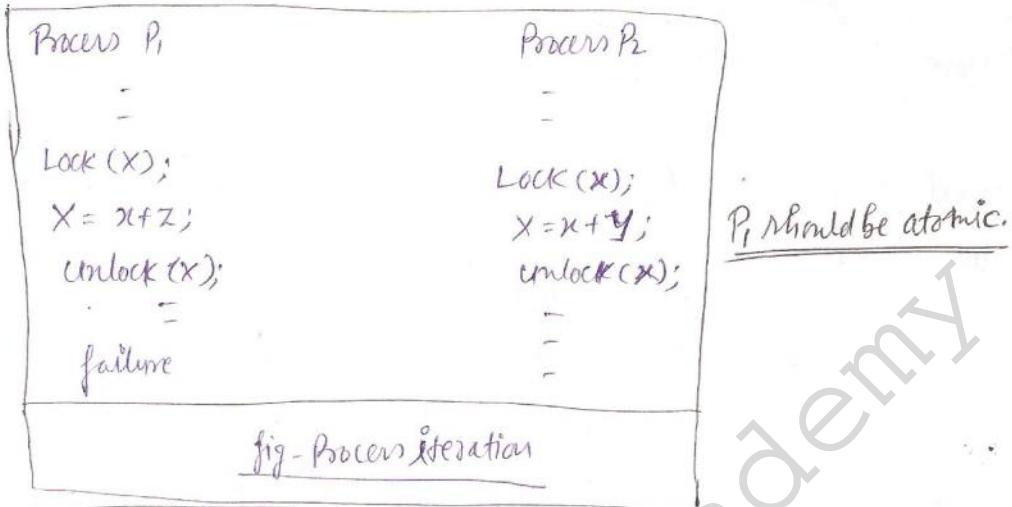
\* consistency violation.

A transaction may abort due

to following events.

Example :-

### X-Sharable Resources -



### Commit Protocols :- [Gray] (2 Phase Commit Protocol)

- Commit Protocols are used to ensure atomicity across sites.
  - A transaction which executes at multiple site must either be committed at all the sites, or aborted at all the sites.
  - Not acceptable to have a transaction committed at one site & aborted at another.
- The 2-Phase Commit (2PC) Protocol is widely used.
- The 3-Phase Commit (3PC) Protocol is more complicated & more expensive, but avoids some drawbacks of 2PC Protocol. (This Protocol is not used in Practice).

## 2-Phase Commit →

- A way to guaranteeing all-or-none agreement to commit a distributed transaction.

## General's Paradox →

This is an abstraction of a core problem in the design of commit protocols.

- Generals must attack simultaneously to capture hill.
- Some messengers may be killed or captured.
- The challenge is to use a protocol that allows the generals to agree on a time to attack, even though some messengers do not get through.

## 2-Phase Commit Protocol →

This protocol assumes as →

- \* One of the cooperating processes acts as a Coordinator.
- \* Other processes are referred to as Cohorts.
- \* Stable storage is available at each site.
- \* Each site uses write-ahead log protocol to achieve local fault recovery & rollback.
- \* At the beginning of the transaction, the Coordinator sends a start transaction message to every cohort.

### Phase-1 →

At the Coordinator →

- ① The Coordinator sends a COMMIT-REQUEST message to every cohort requesting the cohort to commit.
- ② The Coordinator waits for replies from all the cohorts.

At Cohorts →

- ① On receiving a COMMIT-REQUEST message, a cohort takes the following actions.
  - If the transaction executing at the cohort is successful, it writes UNDO & REDO log on the stable storage & sends an AGREED message to the Coordinator.
  - Otherwise, it sends an ABORT message to the Coordinator.

## Phase-2

### (a) At the Coordinator :-

- ④ If all the cohorts reply AGREED & the Coordinator also agrees, then the Coordinator writes a COMMIT record into the log. Then it sends a COMMIT message to all the cohorts. Otherwise the Coordinator sends an ABORT message to all the cohorts.
- ⑤ The Coordinator then waits for acknowledgements from each cohort. (Commit Message ACK)
- ⑥ If an acknowledgement is not received from any cohort within a timeout period, the Coordinator resends the Commit/Abort message to that cohort.
- ⑦ If all the acknowledgements are received, the Coordinator writes a COMPLETE record to the log (to indicate the completion of the transaction).

### (b) At the cohorts :-

- ① On receiving a COMMIT message, a cohort releases all the resources & locks held by it for executing the transaction, & sends an acknowledgement.
- ② On receiving a ABORT message, a cohort undoes the transaction using the UNDO log Record, releases all the resources & locks held by it for performing the transaction, & sends an acknowledgement.

- When there are no failure or message losses, it is easy to see that all nodes will commit only when all the participating (including the coordinator) agree to commit.
- In the case of lost messages (sent from either cohorts or the coordinator), the Coordinator simply resends messages after the timeout.

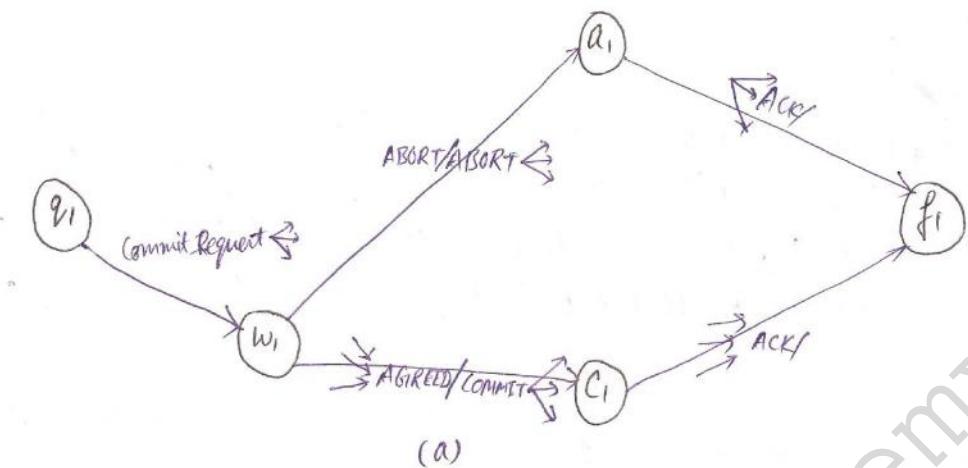
### Site failure :-

- ① The Coordinator crashes before the COMMIT record → on recovery, the Coordinator broadcast an ABORT message to all the Cohorts. All the Cohorts who had agreed to commit will simply undo the transaction using the UNDO log & abort, other cohorts will simply abort the transaction. All the cohorts are blocked until they receive an ABORT Message.
- ② - The Coordinator crashes after the COMMIT record but before COMPLETE → on recovery, the Coordinator broadcasts a COMMIT message to all the cohorts & waits for acknowledgements. In this case also the cohorts are blocked until they receive a COMMIT message.
- ③ - A Coordinator crashes after COMPLETE → on recovery, there is nothing to be done for the transaction.
- ④ - A cohort crash before reply (in Phase 1) → the Coordinator can abort the transaction because it did not receive reply from the crashed cohort.
- ⑤ - A cohort crash in Phase-2 (after UNDO/REDO) → on recovery, the cohort will check with coordinator whether to abort (ie perform an undo operation) or to commit the transaction. Note that committing may require a redo operation because the cohort may have failed before updating the database.

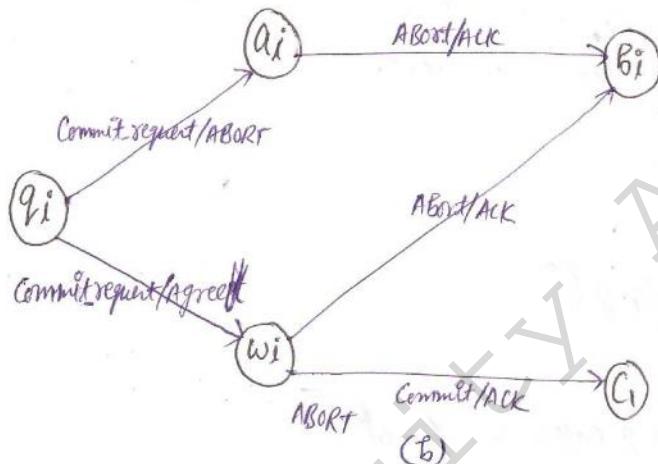
### 2-Phase Commit Protocol

- guarantees global atomicity (its biggest drawback is that it is a blocking protocol)

- if coordinator fails all sites block until coordinator recovers



(a)

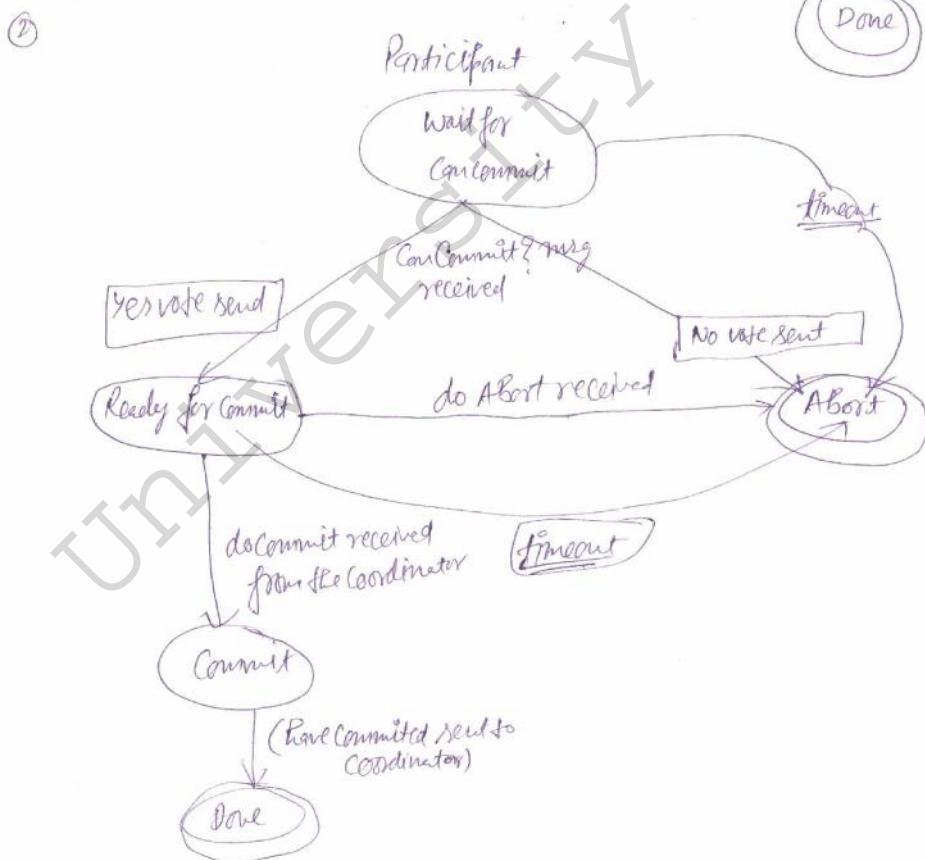
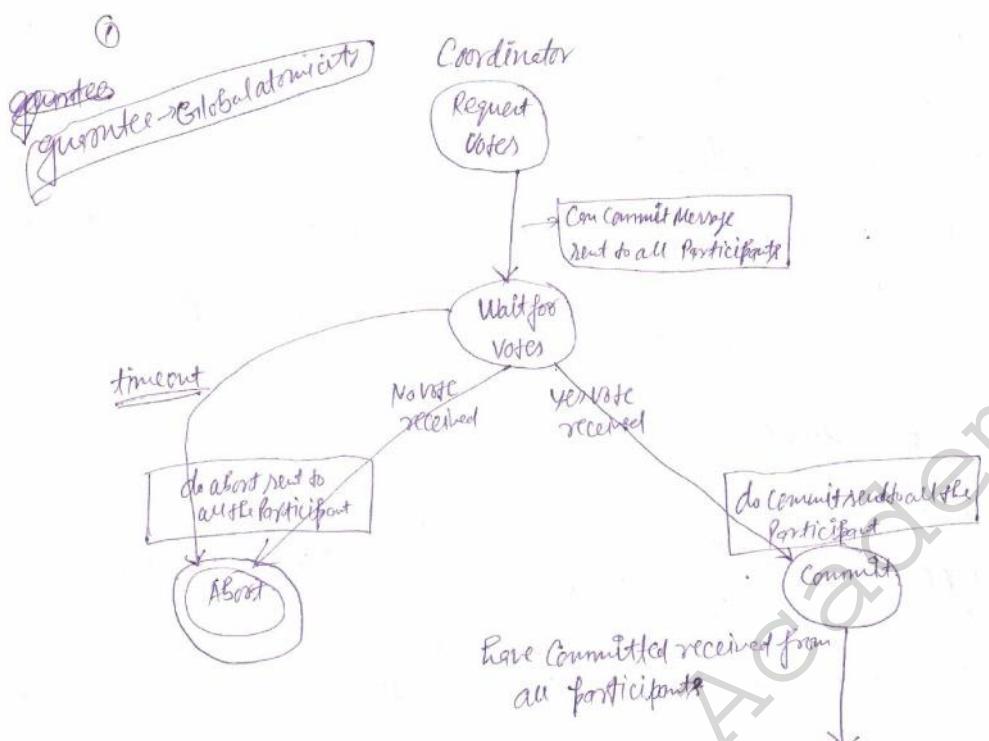


(b)

→ A label of the form "A/B" on a transition edge indicates that the transaction is taken when the set of merge received matches the specification A, & Causes the set of merge specified by B to be taken.

→ " $\xrightarrow{M}$ " represents the receipt of M from all other Rites.

→ " $M\xrightarrow{\cdot}$ " represents the sending of M to all other Rites.



## Voting Protocols

- Replicated data, at Multiple sites.
  - each site has some number of votes.
  - access to replicated data requires a majority of votes.
  - voters determine which version is the current one.
- ✓ Voting mechanism is more fault-tolerant than a Commit protocol in that it allows access to data under network partitions, site failures, & message losses without compromising the integrity of the data.

## Static Voting → (Gifford)

- replicas of files (data), are stored at Multiple sites.
- Each file access operation requires that an appropriate lock is obtained.
- reader-writer locks are supported.  
    { Some writers no reader } to access the file simultaneously.  
    { Multiple readers no writers }
- Every site has a Lock Manager, that performs the lock related operation.
- Every file has a Version No = No of time the file has been updated  
    (No of changes made)  
    Version No are stored in stable storage
- Every replica has some no of voters.
- Voter allocation is on stable storage.
- reader & writer require a quorum. (A minimum no of sites that must be present at a meeting to make it valid) <sup>Vote</sup>

## The Static Voting Algorithm

- When a process executing at site i issues a read or write request for a file, the following protocol is initiated.
- ① Give Lock Request to local lock manager.
  - ② Local lock manager eventually grants requests & then sends Vote Request to all other sites.

at site j

- ③ on receipt of Vote Request from i, issue Lock Request to local lock Manager.
- ④ If the lock request is granted by the local lock Manager, send the version No  $VN_j$  of its replica & the no of votes  $V_j$  of its replica to site i.

at site i

- ⑤ after votes are in, perform quorum test.
- \* Read quorum test →

$$V_r = \sum_{K \in P} V_k \geq r$$

Where  $P$  = the set of sites that Replicated.

If  $V_r \geq r$ , where  $r$  = read quorum, then site i has succeeded in obtaining the read quorum.

\* Write quorum test →

$$V_w = \sum_{K \in Q} V_k \geq w$$

If  $V_w \geq w$ , where  $w$  = write quorum, the site i has succeeded in obtaining the write quorum.

Set of site Q is determined as follows -

$M = \max \{ VN_j \mid j \in P \}$  is the largest version No reported in the Vote, &  
 $Q = \{ j \in P \mid VN_j = M \}$  includes only the votes that correspond to the version Number.

- ⑥ If quorum test fails, give Release\_lock to local manager & all sites in P that returned positive VOTE.
- ⑦ If quorum test succeeds, check whether local Copy is current, if not, obtain a fresh copy from another site.
- ⑧ for a read, just use the local copy.  
 for a write, update the local copy.  
 - then update  $VN_j$  & send the updates &  $VN_j$  to all the sites in Q.
- ⑨ Give a Release\_lock to local manager & all the sites in P.

at another sites  $\rightarrow$

- on receiving update messages, update own local copies.
- on Relock, release all local locks.

Vote assignment  $\rightarrow$

if  $v$  is the total No of Voter, we want to choose  $r$  &  $w$  such that

$$\begin{array}{l} r+w > v \\ w > \frac{v}{2} \end{array}$$

$\checkmark r \rightarrow$  Read Quorum.  
 $\checkmark w \rightarrow$  Write Quorum.

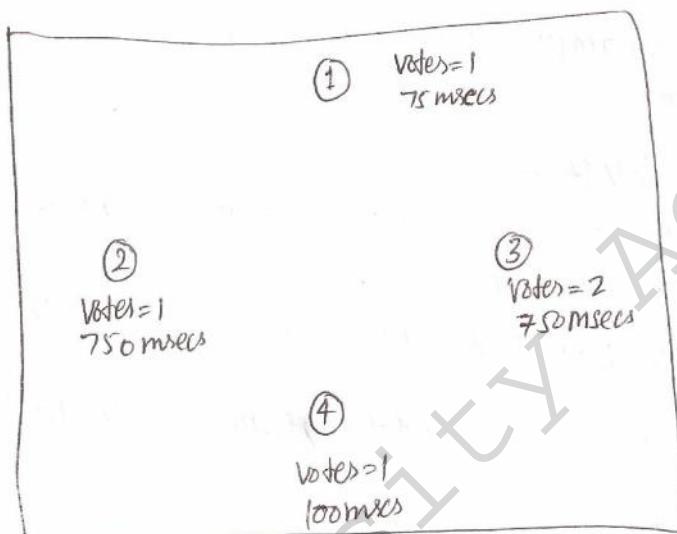


fig - vote assignment

$$\begin{cases} r=1, w=5 & \text{RAT} = 75 \text{ msec} \\ w=5 & \text{WAT} = 750 \text{ msec} \end{cases}$$

Read access time = 75 msecs (RAT)

$$\begin{cases} r=3 \\ w=3 \\ \text{RAT} = 75 \text{ msec} \end{cases}$$

| Site | Voter | Read Access Time |
|------|-------|------------------|
| 1    | 1     | 75 ms            |
| 2    | 1     | 750 ms           |
| 3    | 2     | 750 ms           |
| 4    | 1     | 100 ms           |

$$\begin{cases} r=1 \\ w=5 \\ \text{RAT} = 75 \text{ msec} \\ \text{WAT} = 750 \text{ msec} \end{cases}$$

If  $r=3, w=3$   
the time are unchanged,  
But writer are still parity.

If voter = 4, is more reliable we can further improve reliability by readjusting the voters following-

| S | V | RAT    |
|---|---|--------|
| 1 | 1 | 75 ms  |
| 2 | 1 | 750 ms |
| 3 | 1 | 750 ms |
| 4 | 2 | 100 ms |

Dynamic Voting Protocols → (The No of Vote rearrangements protocols, the No of votes assigned to a side changes dynamically).

- change the set of sides that can form a majority.
- change the distribution of votes.

Dynamic Vote Rearrangement →

Majority based approach.

- No of voters per side changes
- two kinds -

Dynamic Vote rearrangement:

① group consensus on new assignment

② autonomous rearrangement, ratified by majority of sides.

Autonomous Vote Rearrangement →

When side  $i$  wants to increase  $V_i[i]$  -

- send  $V_i$  &  $N_i$  along with new vote value  $x$  to all communicating sides.

- wait for a majority of sides to respond

- If a majority is collected, update  $V_i[i]$  to the new value of increment  $N_i[i]$ .

When side  $j$  receives a vote-increasing request from side  $i$  with  $V_i$ ,  $N_i$ , &  $x$  -

-  $V_j[i] = x$

-  $N_j[i] = N_i[i] + 1$

Vote Decreasing Protocols

When a side wants to decrease  $V_i[i]$

- Set  $V_i[i]$  to the New Value

- increment  $N_i[i]$

- Send  $V_i$  &  $N_i$  to other sides.

When side  $j$  receives a vote-decreasing report from side  $i$  with  $V_i$  &  $N_i$  -

-  $V_j[i] = V_i[i]$

-  $N_j[i] = N_i[i]$

## Vote Collecting Protocols

- for each reply  $V_j$  &  $N_j$  received by Site  $i$ ;
  - \*  $V_i[j] = V_j[j]$
  - \* if  $V_j[i] > V_i[i]$  or  $(V_j[i] < V_i[i] \& N_j[i] > N_i[i])$  then
    - $V_i[i] = V_j[i];$
    - $N_i[i] = N_j[i];$
  - end if;
- if Site  $j$  did not respond to Site  $i$  -
  - \* find  $K \in G$  &  $N_k[j] = \max \{N_p[j] : p \in G\}$ ,  
 $G$  = the set of all sites that replied to  $i$ , that is find the site that has the latest information on the votes assigned to Site  $j$ .
  - \*  $V_i[j] = V_k[j];$
  - \*  $N_i[j] = N_k[j];$
  - \*  $N_i[i] = N_k[i].$

## Deciding the outcome $\rightarrow$

- Let  $K$  be the set of all sites, &  $G$  be the set of sites that responded to the ballot.
- The total no of voters in the system is computed as -

$$TOT = \sum_{K \in K} V_i[K]$$

- The total no of votes received is computed as -

$$RCVD = \sum_{K \in G} V_i[K]$$

- Site  $i$  has a Majority iff  $RCVD > TOT/2$ .

# UNIT-5

(CS-7) 66 B

## Transactions and Concurrency Control

5

**Q1.** What do you mean by Transaction in context of Distributed Systems ? What are its major properties ?

**Ans.** In the database systems, multiple users concurrently access the data/records for various operations. This accessing is possible through the database transactions. The concurrency control is needed to maintain the correctness of database during concurrent access of data. In this way, the database system is assumed to be the set of data objects that are shared and accessed by the various users. The data objects include a data field value, a record, a file, etc. In database, every data value belongs to a specific domain set and comprises a data field. When the user access a database, it means, access to the database value for some operations. The operations and processing on the database are called the Database Actions or Transactions. Various processing required by a user are grouped together in the form of a program and this program works as a unit for program execution with database. This logical unit of the interaction between user and database is called as Transaction.

A transaction may be any read, write or compute operation on any data object of a database. The transaction process does not affect the consistent state of a database. If the database is in consistent state before the transaction, it will remain in the consistent state after the transaction also. It means, the transaction is atomic in nature. Either, the transaction unit executes successfully or it will not execute at all.

The result of a transaction are seen to the user only after the completion of the transaction, i.e., after transaction termination when a transaction is terminated, the user gets the notification of its success or failure through the resulting messages.

**Properties of Transactions :** Transactions have four essential properties, known as ACID properties :

- (1) **Atomicity:** Atomicity means, that the transaction completes its processing as an indivisible and atomic unit either successfully or unsuccessfully. After the transaction, the database changes consistency or not changed at all.
- (2) **Consistency:** Consistency means the database was in consistent state before transaction and after the termination of transaction, the database will also be in the consistent state.
- (3) **Isolation :** Isolation property indicates that every action processed by a transaction is kept isolated until the completion of transaction. During the transaction processing, the processing will be hidden from outside the transaction.
- (4) **Durability:** Durability means, any failure made after the commit operation will not affect the database. The commit action will reflect its results to the database after its termination.

**Q2.** Write down the States of Transaction.

**Ans.** The transaction itself is an independent unit which is processed by the user. But this atomic operation has several different life states during its processing. In a database system, a transaction may be found in several working states as modify, abort, commit, rollback, etc. When a transaction starts for a data object, it is called in the initiated state. This transaction moves forward in three different directions. The one way leads to the successful termination of the transaction where the commit operation is performed successfully. The second way forces a transaction to terminate. This is called as the Murderous Termination, which occurs on detecting an error during the transaction processing and the transaction is aborted at once. The another way for transaction termination is the Suicidal Termination where the transaction does not update the database but rollback itself and abort the operation. The following figure illustrates the various transaction states in a system.

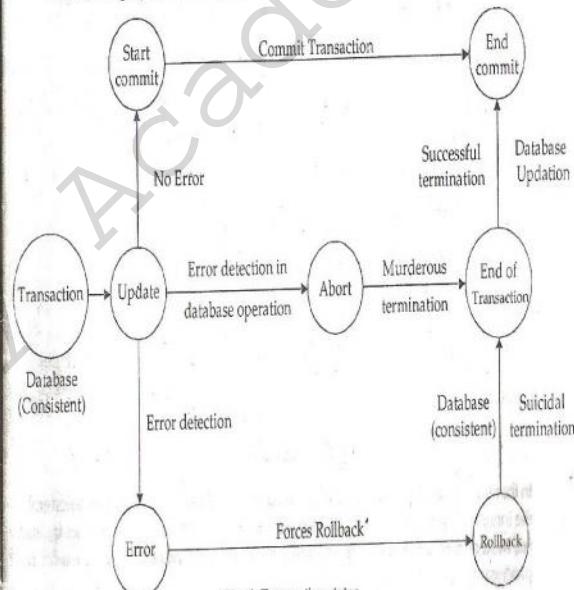


Fig. 1, Transaction states.

In the above figure, we can see that a transaction always reaches to its end either successfully or unsuccessfully. The database is assumed in consistent state before transaction state. When a transaction starts, it initiates the update operation, i.e., modify transaction. This state of the database is called as the modify state. Then the transaction is executed till its termination. When a transaction executes successfully, it completes the modify state of itself. The commit transaction is started if no error is detected out of transaction. When the transaction is terminated, it is in the consistent state. The change made by the transaction are then reflected to database. When all the changes, made by any transaction are reflected to the database, then the transaction is said to be in consistent state. The transaction may be successful or unsuccessful. Then transaction either be successful or murderous or suicidal.

**Q3.** What is Nested Transaction ? What are the advantages of nested transactions ?

**Ans.** In a database system, the transaction is said to be the query to the database management software program. Most often, the transaction is considered as the independent unit which

produces the result of one query at once. But sometimes, the result of a transaction may depend on the result of another transaction. In such case, the result is not independent. The database system uses a transaction to get a result database. Then this result value is used by any other transaction to get any other data value.

In this manner, the transaction (query) is based on other transaction (subquery). The subquery is the part of a transaction and so called nested transaction. These are called nested because, these become the integral part of a transaction process. The following figure illustrates the processing of the nested transaction:

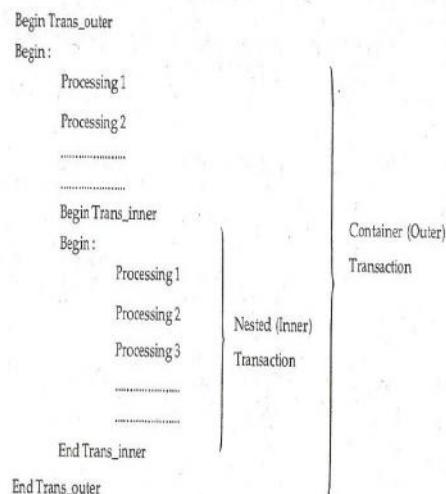


Fig. 2. Nested transaction.

In the above figure, two transactions are related with each other. The outer transaction embeds the inner transaction. When the transaction processing starts to execute, the inner transaction i.e., nested transaction executes first and pass its result to the outer transaction for further processing.

The overall transaction result is based on the accumulated result of the two transactions. Advantages: Nested transactions have the following main advantages :

- ✓ Subtransactions at one level may run concurrently with other subtransactions at the same level in the hierarchy. This can allow additional concurrency in a transaction.
- ✓ Subtransactions can commit or abort independently. In comparison with a single transaction, a set of nested sub transactions is potentially more robust.

#### Q. 4. What do you understand by concurrency control in distributed transactions ?

Ans. Each server manages a set of objects and is responsible for ensuring that they remain consistent when accessed by concurrent transactions. Therefore, each server is responsible for applying concurrency control to its own objects. The members of a collection of servers of distributed transactions are jointly responsible for ensuring that they are performed in a serially equivalent manner.

This implies that if transaction T is before transaction U in their conflicting access to objects at one of the servers then they must be in that order at all of the servers whose objects are accessed in a conflicting manner by both T and U.

#### Q. 5. What are Locks? Explain the locking algorithms for Distributed Database.

OR

What do you mean by two phase locking? How it is different from strict two phase locking ? Explain. [IUPTU 2010, Marks 10]

Ans. Every data object has a lock associated with the transaction with database can make the request for that lock as required. The transaction may hold or release the lock on the data object as required by the transaction. The transaction is said to have the locked data object, if it holds a lock.

In a transaction, the data object can be locked in two modes as:

- (i) Exclusive Lock : In case of Exclusive Lock, no other transaction can concurrently lock that data object at all.
- (ii) Shared Lock : If a data object is locked by a transaction in shared mode, the other transactions can also concurrently lock it (but only in shared mode).

Thus, we can say that locking the data object means restricting it for other transactions. The locked data object becomes inaccessible to other database transactions. Locking the data objects is helpful for concurrency control.

**Lock Based Locking Algorithms :** The lock based concurrency algorithms are used for concurrency control using lock based techniques. All these algorithms allow the locking of data object for transaction. A transaction is considered as the sequence of n actions. The lock based algorithms allow read, write, lock and unlock actions. A transaction can be locked for an action and this lock is removed by an unlock action. For a well-performed transaction, the following points are considered:

- The transaction locks a data object before accessing it.
- The transaction does not lock the data object more than once.
- The transaction unlocks all the locked data objects before the completion of the transaction. Various locking schemes are proposed for the distributed database. Some of them are :
- (i) **Static Locking Scheme :** According to static lock scheme, all the required data objects are locked before the execution of an action. All the data objects are predeclared. All these objects are unlocked by the transaction after the execution of all related actions. This type of locking is very simple, but useful for handling the concurrency in distributed database transaction. The drawback of this scheme is its requirement of previous knowledge about database.
- (ii) **Two-phase Locking Scheme :** This locking scheme is also called as '2PL'. This scheme is also called as Dynamic Locking Scheme, because the lock is applied to data object only when the transaction needs the data object. A specific constraint (Rule) is applied to the lock acquisition and lock release. A transaction cannot request a lock on the data object after its unlocking. So, all the required data objects must be locked before unlocking the data objects. The two-phase locking scheme has two phases of its working as :

Phase 1: Growing Phase : - *locks are acquired & no locks are released.*

- (i) Transaction may obtain locks.
- (ii) Transaction may not release locks.

Phase 2: Shrinking Phase : - *locks are released & no locks are acquired.*

- (i) Transaction may release locks.
- (ii) Transaction may not obtain locks.

| LOCK type  | Lock Compatibility |            |
|------------|--------------------|------------|
|            | read-lock          | write-lock |
| read-lock  | X                  |            |
| write-lock |                    | X          |

There is one lock point present in between these two phases. At the Lock Point Stage, the transaction holds locks on all the needed data objects. Figure below shows a two phase locking scheme.

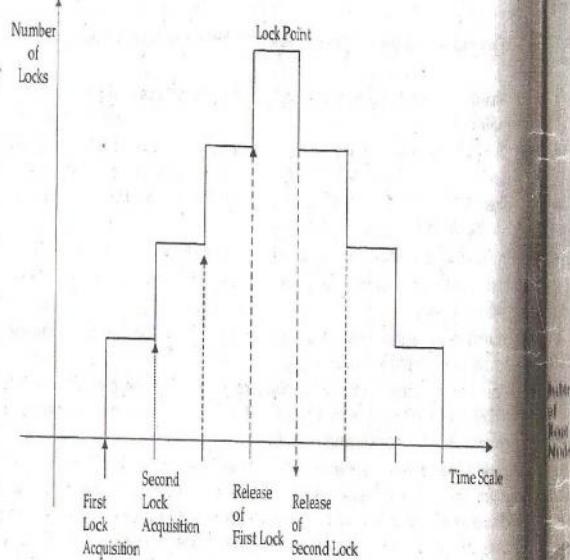


Fig. 3. Two Phase Locking Scheme.

Initially, a transaction is in the growing phase. The transaction acquires locks as needed. Once the transaction releases a lock, it enters the shrinking phase, and it can issue no new lock requests.

Two phase locking does not ensure freedom from deadlock. Cascading rollback may occur.

Cascading rollbacks can be avoided by a modification of two-phase locking called the Strict Two-Phase Locking Protocol. The strict two-phase locking protocol requires that in addition to locking being two-phase, all exclusive mode locks taken by a transaction must be held until that transaction commits.

(iii) **Timestamp Based Locking Scheme :** A transaction is assigned a unique timestamp when initiation. The timestamp is used to specify the order of the transactions and so helps in transaction ordering. The another use of the time stamp is to solve the conflict problem between the transactions and so the time stamps are helpful for deadlock prevention.

A conflict is resolved by Wait Action, Restart Action, Die Action, Wound Action.

- The die action and the wound action are associated with the Restart of transaction.
- The die action aborts the requesting transaction and starts it as a fresh transaction.
- The requesting transaction is preceded only when the wounded transaction aborts or completes.

(iv) **Non-Two Phase Locking Scheme :** This locking scheme is used for the database in which data objects are organized as an hierarchy, i.e., inverted tree like structure. The Non-Two Phase locking is helpful for deadlock free transaction in hierarchical database. In such schema,

a data object cannot be locked for more than one time. When a transaction needs some data object, it is provided by the system after locking it. Sometimes, a transaction may try to lock such an object which is already locked. In such situation, that transaction will not be able to lock that data object and the blocking of that transaction takes place. But any other transaction seeking for that data object will try to get it and lock the data object for the further processing. Look at the following hierarchical tree of the database system.

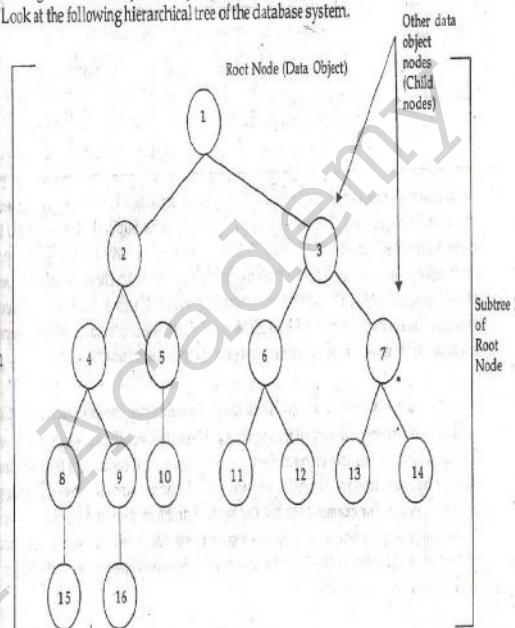


Fig. 4. Hierarchical Tree for Database System.

In such database system, the transaction initiates at a specific data object node and locks that data object. Furthermore, the transaction can only lock the data object in the subtree of that node. It means a data object node is locked first and then locks its child nodes only.

The Non-Two Phase locking scheme is more advantageous than the 2PL scheme, because it provides the freedom from deadlock in the system. One more advantage of this locking scheme is its dynamic nature, i.e., the lock is released after the use of data object by the transaction. By this concept, the data objects used by a transaction will be allowed to be used by other transaction very soon after the completion of first transaction.

#### Q. 6. Define Locking in Concurrency Control.

Ans. The oldest and most widely used concurrency control algorithm is locking. In a distributed transaction, the locks on an object are held locally (in the same server). The local lock manager can decide whether to grant a lock or make the requesting transaction wait. However, it cannot release any lock until it knows that the transaction has been committed or aborted at all the servers involved in the transaction. When locking is used for concurrency control, the objects remain locked and are unavailable for other transactions during the atomic commit protocol, although an aborted transaction releases its locks after phase 1 of the protocol.

(The transaction in conflict with the requesting transaction is tagged as) —————— wounded

As lock managers in different servers set their locks independently of one another, it is possible that different servers may impose different orderings on transactions. Consider the following interleaving of transactions T and U at servers X and Y:

|  | T                         | U                         |
|--|---------------------------|---------------------------|
|  | Write(A) at X locks A ①   |                           |
|  | Read(B) at Y wait for U ② | Write(B) at Y locks B ③   |
|  |                           | Read(A) at X wait for T ④ |

The transaction T locks object A at server X and then transaction U locks object B at server Y. After that, T tries to access B at server Y and waits for U's lock. Similarly, transaction U tries to access A at server X and has to wait for T's lock. Therefore, we have T before U in one server and U before T in the other. These different orderings can lead to cyclic dependencies between transactions and a distributed deadlock situation arises. When a deadlock is detected, a transaction is aborted to resolve the deadlock. In this case, the coordinator will be informed and will abort the transaction at the participants involved in the transaction.

Q.7. Why is Concurrency Control needed? Write down the main categories of Concurrency Control.

Ans. Need of Concurrency Control: If transactions are executed serially, i.e., sequentially with no overlap in time, no transaction concurrency exists. However, if concurrent transactions with interleaving operations are allowed in an uncontrolled manner, some unexpected, undesirable result may occur. For example: In the lost update problem, a second transaction writes a second value of a data-item (datum) on top of a first value written by a first concurrent transaction, and the first value is lost to other transactions running concurrently which need, by their precedence, to read the first value.

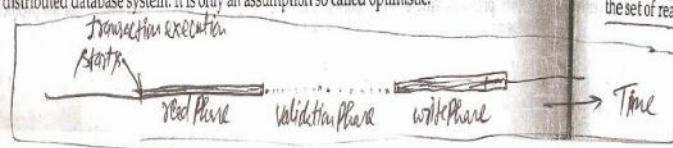
(Refer M&P Pg 1 Example)

Categories of Concurrency Control: The main categories of concurrency control mechanisms are:

- Optimistic: Delay the checking of whether a transaction meets the isolation and other integrity rules (e.g., serializability and recoverability) until its end, without blocking any of its (read, write) operations and then abort a transaction to prevent the violation, if the desired rules are to be violated upon its commit. An aborted transaction is immediately restarted and re-executed, which incurs an obvious overhead. If not too many transactions are aborted, then being optimistic is usually a good strategy.
- Pessimistic: Block an operation of a transaction, if it may cause violation of the rules, until the possibility of violation disappears. Blocking operations is typically involved with performance reduction.
- Semi-optimistic: Block operations in some situations, if they may cause violation of some rules, and do not block in other situations while delaying rules checking (if needed) to transaction's end, as done with optimistic.

Q.8. Explain Optimistic Concurrency Control.

Ans. Optimistic concurrency control states that the conflicts among the transactions are rare in distributed database system. It is only an assumption so called optimistic.



The optimistic concurrency algorithms are also based on this concept. These algorithms are developed by assuming that "the transaction conflicts are rare".

These algorithms apply a checking for the transaction. But this checking is not performed before transaction. When a transaction is executed, it is checked for its conflicts with other concurrent transactions. If any conflict is found, then the transaction will be aborted and will not forward the update operation. If no conflict is found, the transaction process to commit and do the updation in database as a result of transaction. The process of undone the transaction is also known as Rollback.

A transaction after starting leads to either commit or rollback. If the transaction is successfully executed, it leads to the transaction commit and causes the change in data of database as the result of transaction committed. On the other hand, the transaction leads to transaction rollback, if the transaction will not execute successfully due to some conflict with any other transaction.

The most common algorithm used for optimistic concurrency control was proposed by Kung and Robinson. They proposed that the transaction executes first with any synchronization and concurrency check. But it will not update the database finally. The transaction result is written to database only when it is not conflicted, otherwise not.

For example: Consider the following interleavings of transactions T and U, which access objects A and B at server X and Y, respectively.

| T            | U             |
|--------------|---------------|
| Read(A) at X | Read(B) at Y  |
| Write(A)     | Write(B) at Y |
| Read(B) at Y | Read(A) at X  |
| Write(B)     | Write(A)      |

The transactions access the objects in the order T before U at server X and in the order U before T at server Y. Now suppose that T and U start validation at about the same time, but server X validates T first and server Y validates U first.

Q.9. What do you mean by Timestamp Ordering?

Ans. Each transaction is assigned a unique timestamp value when it starts. The timestamp defines its position in the time sequence of the transaction. Requests from transaction can be totally ordered according to their timestamp. The basic timestamp ordering rule is based on operation conflicts and it is very simple. A transaction's request to write an object is valid only if that object was last read and written by earlier transaction. A transaction's request to read an object is valid only if that object was last written by an earlier transaction. The timestamp ordering rule is refined to ensure that each transaction accesses a consistent set of versions of the objects. The write operations may be performed after the close transaction operation has returned. Without making the client wait, but the client must wait when read operations need to wait for earlier transactions to finish. The time stamps are assigned from server's clocks.

The write operations are recorded in tentative versions of objects and are invisible to other transactions until a close transaction request is issued and the transaction is committed. The write timestamp of the committed object is earlier than that of any of its tentative versions and the set of read timestamps can be represented by its maximum members.

Whenever a transaction's write operation on an object is accepted, the server creates a new tentative version of the objects with write timestamp set to the transaction timestamp. Transactions read operation is directed to the version with the maximum write timestamp less than the transaction timestamp. Whenever a transaction's read operation on an object is accepted, the timestamp of the transaction is added to its set of read timestamp. When a transaction is committed, the value of the tentative versions becomes the value of the objects and the timestamp of the tentative versions becomes the timestamps of the corresponding objects.

**Q.10. Explain the methods for concurrency control in distributed transaction.** [UPTU 2010, Marks 5]

OR

What are the essential differences in the lock-based protocols and time stamp based protocols? [UPTU 2004, 2005; Marks 5]

**Ans.** The three types of concurrency control algorithms are: the locking scheme, timestamp method and the optimistic control. All of these are applied to distributed database system. Their comparison depends on the necessity of data objects, synchronization, data updation and data consistency.

- Timestamp method is similar to two-phase locking in that both use pessimistic approaches in which conflicts between transactions are detected as each object is accessed.
- Timestamp ordering decides serialization order statically when transaction start. On the other hand two-phase locking decides the serialization order dynamically.
- Timestamp ordering is better than strict two phase locking for read-only transaction, while two-phase locking is better when the operations in transactions are predominantly updated.
- Timestamp ordering aborts the transaction immediately, whereas locking makes the transaction wait but with a possible later penalty of aborting to avoid deadlock.
- When optimistic concurrency control is used, all transactions are allowed to proceed, but some are aborted when they attempt to commit or forward validation transaction are aborted earlier. This results in relatively efficient operation when there are few conflicts but a substantial amount of work may have to be repeated when a transaction is aborted.

**Q.11. Explain Multiversion Timestamp Ordering.**

OR

What are the advantages and drawbacks of Multiversion Timestamp Ordering in comparison with the Ordering Timestamp Ordering.

**Ans.** In multiversion timestamp ordering a list of old committed versions as well as tentative versions is kept for each object. This list represents the history of the values of the object. Each version has a read timestamp recording the largest timestamp of any transaction that has read from it in addition to a write timestamp. Whenever a write operation is accepted, it is directed to a tentative version with the write timestamp of the transaction. Whenever a read operation is carried out it is directed to the version with the largest write timestamp less than the transaction timestamp. If the transaction timestamp is larger than the read timestamp of the version being used, the read timestamp of the version is set to the transaction timestamp. When a read arrives late, it can be allowed to read from an old committed version, so there is no need to abort the read operations. In multiversion timestamp ordering, read operation is always permitted although they may have to wait for earlier transaction to complete (with commit or abort), which ensures that executions are recoverable.

**Advantages:** The main advantage of multiversion timestamp ordering is that a read request never fails and is never made to wait.

**Disadvantages:**

- (1) The reading of a data item also requires the updating of the R-timestamp field resulting in two potential disk accesses, rather than one.
- (2) The conflicts between transactions are resolved through rollbacks, rather than through waits. This alternative may be expensive.

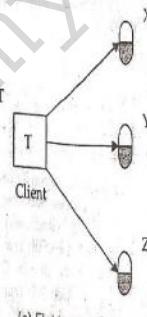
This multiversion timestamp ordering scheme does not ensure recoverability and cascadelessness. It can be extended in the same manner as the basic timestamp ordering scheme, to make it recoverable and cascadeless.

**Q.12. Discuss Flat and Nested Distributed Transactions.**

**Ans.** A client transaction becomes distributed if it invokes operations in several different servers. There are two different ways that distributed transactions can be structured as flat transactions and as nested transactions.

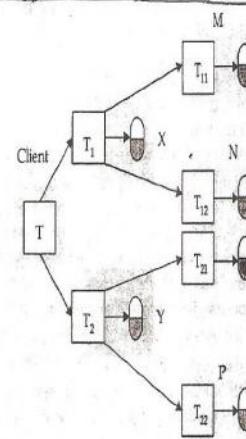
In a flat transaction, a client makes requests to more than one server. For example,

In the Fig. 5(a), transaction T is a flat transaction that invokes operations on objects in servers X, Y and Z. A flat client transaction completes each of its requests before going on to the next one. Therefore, each transaction can only be waiting for one object at a time.



(a) Flat transactions

In a nested transaction, the top-level transaction can open subtransactions, and each subtransaction can open further subtransactions down to any depth of nesting.



(b) Nested transactions

**Fig. 5.**  
The Fig. 5(b) shows a client's transaction T that opens two subtransactions T<sub>1</sub> and T<sub>2</sub>, which access objects at servers X and Y. The subtransactions T<sub>1</sub> and T<sub>2</sub> open further subtransactions

$T_{11}, T_{12}, T_{21}$  and  $T_{22}$ , which access objects at servers M, N and P. In the nested case, subtransactions at the same level can run concurrently, so  $T_1$  and  $T_2$  are concurrent, and as they invoke objects in different servers, they can run in parallel. The four subtransactions  $T_{11}, T_{12}, T_{21}$  and  $T_{22}$  also run concurrently.

Consider a distributed transaction in which a client transfers \$10 from account A to C and then transfers \$20 from B to D. Accounts A and B are at separate servers X and Y and accounts C and D are at server Z. If this transaction is structured as a set of four nested transaction, as shown in figure below, the four requests (two deposit and two withdraw) can run in parallel and the overall effect can be achieved with better performance than a simple transaction in which the four operations are invoked sequentially.

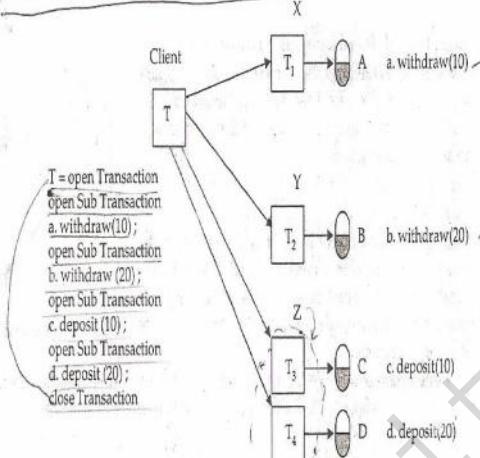


Fig. 6. Nested banking transaction.

### Q. 13. What are Atomic Commit Protocols?

Ans. Every transaction is atomic in nature. In the distributed system, the sequence of such atomic operations execute to perform a specific work. A transaction is said to be atomic, if it is indivisible, instantaneous and cannot be interrupted. Generally, the machine level instructions are assumed atomic because of having all such characteristics and continue process unless the system fails. Such instructions (transactions) are grouped together to perform a task. Atomic transactions are the building blocks for constructing the fault tolerant system. An atomic transaction must possess the following characteristics :

- A transaction is atomic if it does not affect by an existing process.
- A transaction is atomic if it does not communicate with other process.
- A transaction is atomic if its process can detect no state change.
- A transaction is atomic if it is considered as indivisible and instantaneous.

In the distributed system, a distributed transaction is a group of several transactions and these are processed in such a manner so that the database consistency is maintained. When a transaction starts, its execution leads to transaction commit or abort. There are various transaction states such as Abort, Commit, Rollback. The transaction may abort due to hardware failure, wrong input, data inconsistency, etc.

The distributed system have multiple sites together on which many processes execute to perform a task. The distributed system states that a transaction must be processed at every site or at no site to maintain the database integrity. In the same manner, all sites of the distributed system are responsible for the atomicity of the transaction. This type of transaction state is called as the Global Atomicity for that transaction.

Q. 14. What do you mean by two-phase commit protocol?

[IUPTU 2005, Marks 5]

OR

What is two-phase commit protocol? Explain its operation.

Ans. Two-phase commit protocol is designed to allow any participant to abort its part of transaction. Due to the requirement for atomicity, if one part of a transaction aborted then the whole transaction must also be aborted.

Phase 1 (voting phase):

1. The coordinator sends a canCommit? request to each of the participants in the transaction.
2. When a participant receives canCommit? request it replies with its vote (Yes or No) to the coordinator. Before voting Yes, it prepares to commit by saving objects in permanent storage. If the vote is No the participant aborts immediately.

Phase 2 (completion according to outcome of vote):

3. The coordinator collects the votes (including its own).
  - (i) If there are no failures and all the votes are Yes the coordinator decides to commit the transaction and sends a doCommit request to each of the participants.
  - (ii) Otherwise the coordinator decides to abort the transaction and sends doAbort requests to all participants that voted Yes.
4. Participants that voted Yes are waiting for a doCommit or doAbort request from the coordinator. When a participant receives one of these messages it acts accordingly and in the case of commit, makes a haveCommitted call as confirmation to the coordinator.

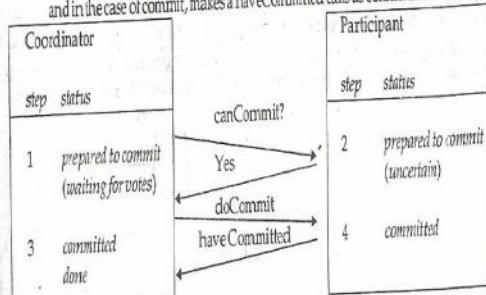


Fig. 7. Communication in two-phase commit protocol.

Operations for two-phase commit protocol:

canCommit ? (trans) → Yes/No

Call from coordinator to participant to ask whether it can commit a transaction.

Participant replies with its vote.

doCommit(trans)

Call from coordinator to participant to tell participant to commit its part of a transaction.  
doAbort(trans)

Call from coordinator to participant to tell participant to abort its part of a transaction.  
haveCommitted(trans, participant)

Call from participant to coordinator to confirm that it has committed the transaction.  
 $\text{getDecision}(\text{trans}) \rightarrow \text{Yes/No}$

Call from participant to coordinator to ask for the decision on a transaction after it has voted Yes but has still had no reply after some delay. Used to recover from server crash or delayed messages.

**Q. 15.** In the two-phase commit protocol, why can blocking be completely eliminated, even when the participants elect a new coordinator? [UPTU 2005, Marks 5]

Ans. The two-phase commit protocol is designed to allow any participant to abort its part of a transaction. Due to the requirement for atomicity, if one part of a transaction is aborted, then the whole transaction must also be aborted. In the first phase of the protocol, each participant votes for the transaction to be committed or aborted. Once a participant has voted to commit a transaction, it is not allowed to abort it. Therefore, before a participant votes to commit a transaction, it must ensure that it will eventually be able to carry out its part of the commit protocol, even if it fails and is replaced in the interim. A participant in a transaction is said to be in a prepared state for a transaction if it will eventually be able to commit it. To make sure of this, each participant saves in permanent storage of all the objects that it has altered in the transaction, together with its status-prepared.

In the second phase of the protocol, every participant in the transaction carries out the joint decision. If any one participant votes to abort, then the decision must be to abort the transaction. If all the participants vote to commit, then the decision is to commit the transaction.

**Q. 16.** Explain the Two-Phase Commit Protocol for Nested Transaction. Explain its operations.

Ans. The outermost transaction in a set of nested transactions is called the top-level transaction. Transactions other than the top-level transaction are called subtransactions. In Fig. 8,  $T$  is the top-level transaction,  $T_1$ ,  $T_2$ ,  $T_{11}$ ,  $T_{12}$ ,  $T_{21}$  and  $T_{22}$  are subtransactions.  $T_1$  and  $T_2$  are child transactions of  $T$ , which is referred to as their parent. Similarly,  $T_{11}$  and  $T_{12}$  are child transactions of  $T_1$ , and  $T_{21}$  and  $T_{22}$  are child transactions of  $T_2$ . Each subtransaction starts after its parent and finishes before it. Thus, for example,  $T_{11}$  and  $T_{12}$  start after  $T_1$  and finish before it.

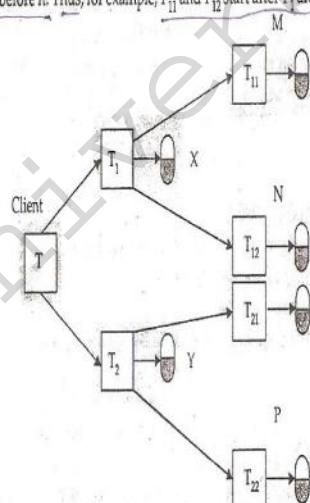


Fig. 8.

When a subtransaction completes, it makes an independent decision either to commit provisionally or to abort. A provisional commit is different from being prepared to commit: it is not backed upon permanent storage. If the server crashes subsequently, its replacement will not be able to commit. After all subtransactions have completed, the provisionally committed ones participate in a two-phase commit protocol, in which servers of provisionally committed subtransactions express their intention to commit and those with an aborted ancestor will abort. A prepared commit guarantees a subtransaction will be able to commit, whereas a provisional commit only means that it has finished correctly - and will probably agree to commit when it is subsequently asked to.

A coordinator for a subtransaction will provide an operation to open a subtransaction, together with an operation enabling the coordinator of a subtransaction to enquire whether its parent has yet committed or aborted.

A client starts a set of nested transactions by opening a top-level transaction with an  $\text{openTransaction}$  operation, which returns a transaction identifier for the top-level transaction. The client starts a subtransaction by invoking the  $\text{openSubTransaction}$  operation, whose argument specifies its parent transaction. The new subtransaction automatically joins the parent transaction, and a transaction identifier for a subtransaction is returned.

#### Operations in coordinator for nested transactions :

$\text{openSubTransaction}(\text{trans}) \rightarrow \text{subTrans}$

Open a new subtransaction whose parent is  $\text{trans}$  and returns a unique subtransaction identifier.

$\text{getStatus}(\text{trans}) \rightarrow \text{committed, aborted, provisional}$

Asks the coordinator to report on the status of the transaction  $\text{trans}$ . Returns values representing one of the following : committed, aborted, provisional.

**Q. 17.** Explain how the two-phase commit protocol for nested transaction ensures that if the top-level transaction commits all the right descendants are committed or aborted.

[UPTU 2008, Marks 10]

Ans. Consider the top-level transaction  $T$  and its subtransactions shown in Fig. 9. Each subtransaction has either provisionally committed or aborted. For example,  $T_{12}$  has provisionally committed and  $T_{11}$  has aborted, but the fate of  $T_{12}$  depends on its parent  $T_1$  and eventually on the top-level transaction,  $T$ . Although  $T_{21}$  and  $T_{22}$  have both provisionally committed,  $T_2$  has aborted and this means that  $T_{21}$  and  $T_{22}$  must also abort. Suppose that  $T$  decides to commit in spite of the fact that  $T_2$  has aborted, also that  $T_1$  decides to commit in spite of the fact that  $T_{11}$  has aborted.

When a top-level transaction completes, its coordinator carries out a two-phase commit protocol. The only reason for a participant subtransaction being unable to complete is if it has crashed since it completed its provisional commit. When each subtransaction was created, it joined its parent transaction. Therefore, the coordinator of each parent transaction has a list of its child subtransactions. When a nested transaction provisionally commits, it reports its status and the status of its descendants to its parent. When a nested transaction aborts, it just reports abort to its parent without giving any information about its descendants. Eventually, the top-level transaction receives a list of all the subtransactions in the tree, together with the status of each. Descendants of aborted subtransactions are actually omitted from this list.

Rules

- (1) A present can commit even if a subtransaction aborts.
- (2) If a present aborts, then its subtransaction must abort.

Transactions and Concurrency Control (CS-7) 80 B

Distributed Systems

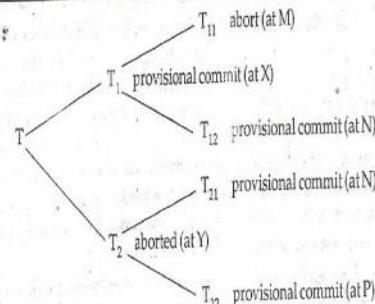


Fig. 9. Transaction T decides whether to commit.

#### Q. 18. Discuss Concurrency Control in Distributed Transactions.

Ans. The database system contains the three modules for concurrency control. These modules are the Transaction Manager, the Data Manager and the Scheduler. Any transaction from the user to database goes first of all to the transaction manager. Then the transaction manager interacts with data manager for further execution. The transaction manager assigns a timestamp to a transaction. It also assigns the request for lock and unlock of data objects as per the user's requirement.

The scheduler works between the transaction manager and data manager. It is used for concurrency control. The scheduler unlocks the data objects or release locks on them as used by the transaction. The data manager manages the database and read from it or write into it according to the user's requirement. The data manager is also responsible for failure recovery.

The transaction manager executes the transaction sequentially. It takes one request and sends it to data manager through scheduler and then data manager gives the appropriate result. After that the data manager executes the transaction actions in the form of execution.

For concurrency control the scheduler modifies the stream of actions towards the data manager.

#### Q. 19. What are the problems with Concurrency Control?

Ans. In the database system the probability of the consistency in database can be increased by executing transaction in the sequential manner. However, the serial execution of transaction also gives the poor response to the user requirements and also there is poor utilization of system resources. By executing transactions concurrently the efficiency can be increased. In concurrent execution, several accesses will be made to some data object and so a critical position can be made in the concurrent execution of transaction.

(i) **Inconsistent Retrieval of Data :** When many transaction read the data objects from the database through transaction manager, scheduler and data manager; but before any another transaction is completed, its modification to these data objects, then the former transaction undergoes in a problem of retrieving incorrect values of data objects. This is known as Inconsistent Retrieval of Data.

(ii) **Inconsistent Update of Data :** When many transaction work on the command set of data objects of any database, although leaving it in inconsistent state, the inconsistent update occurs. Since if the concurrency of the transactions is not controlled then because of many inconsistent states the database may be in inconsistent situation. To overcome this problem the concurrency control method is used. These mechanisms control the relative orders of the conflicting works, so that every transaction has consistent state and on the completion of all

#### Distributed Systems

#### (CS-7) 81 B Transactions and Concurrency Control

transactions the database is in consistent state. So the consistent state of the database system can be controlled by controlling the concurrency.

#### Q. 20. Explain Distributed Transaction.

OR

How does the Distributed Deadlock affect the transactions?

Ans. Deadlocks can arise within a single server when locking is used for concurrency control. Servers must either prevent or detect and resolve deadlocks. With deadlock detection schemes, a transaction is aborted only when it is involved in a deadlock.

Most deadlock detection schemes operate by finding cycles in the transaction wait-for-graph. In a distributed system involving multiple servers being accessed by multiple transactions, a global wait-for-graph can be constructed from the local ones. There can be a cycle in the global wait-for-graph that is not in any single local one; that is, there can be a distributed deadlock.

For example:

| U                                       | V                                      | W                                  |
|-----------------------------------------|----------------------------------------|------------------------------------|
| <i>a.deposit(10)</i><br>lock D          |                                        |                                    |
|                                         | <i>b.deposit(10)</i><br>lock B<br>at Y |                                    |
| <i>a.withdraw(20)</i><br>lock A<br>at X |                                        |                                    |
|                                         | <i>c.withdraw(20)</i><br>wait at Y     |                                    |
| <i>b.withdraw(30)</i><br>wait at X      |                                        |                                    |
|                                         | <i>c.deposit(30)</i><br>lock C<br>at Z |                                    |
|                                         |                                        | <i>a.withdraw(20)</i><br>wait at X |

Fig. 10. Interleavings of transactions U, V and W.

The above figure shows the interleavings of the transactions U, V and W involving the objects A and B managed by servers X and Y and objects C and D managed by server Z. The complete wait-for-graph, in Fig. 11(a), shows that a deadlock cycle consists of alternate edges, which represent a transaction waiting for an object and an object held by a transaction. As any transaction can only be waiting for one object at a time, objects can be left out of wait-for-graphs as shown in Fig. 11(b).

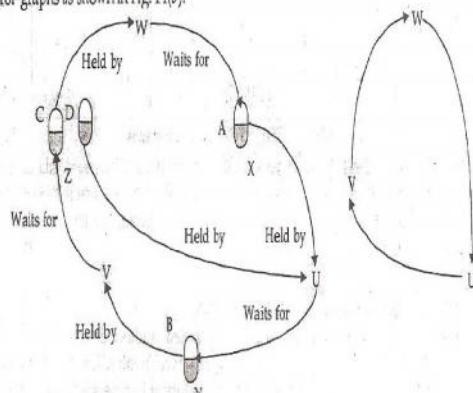


Fig. 11. Distributed deadlock

Detection of a distributed deadlock requires a cycle to be found in the global transaction wait-for-graph that is distributed among the servers that were involved in the transactions. Local wait-for-graphs can be built by the lock manager at each server. In this example, the local wait-for-graphs of the servers are:

✓ Server Y:  $U \rightarrow V$  (added when U requests b.withdraw(30))

✓ Server Z:  $V \rightarrow W$  (added when V requests c.withdraw(20))

✓ Server X:  $W \rightarrow U$  (added when W requests a.withdraw(20))

As the global wait-for-graph is held in part by each of the several servers involved, communication between these servers is required to find cycles in the graph.

A simple solution is to use centralized deadlock detection, in which one server takes on the role of global deadlock detector. From time to time, each server sends the latest copy of its local wait-for-graph to the global deadlock detector, which amalgamates the information in the local graphs in order to construct a global wait-for-graph. The global deadlock detector checks for cycles in the global wait-for-graph. When it finds a cycle, it makes a decision on how to resolve the deadlock and informs the servers as to the transaction to be aborted to resolve the deadlock.

#### Q. 21. What do you mean by Phantom Deadlock and how it occurs in Distributed System?

Ans. A deadlock that is 'detected' but is not really a deadlock is called a phantom deadlock. In distributed deadlock detection, information about wait-for relationships between transactions is transmitted from one server to another. If there is a deadlock, the necessary information will eventually be collected in one place and a cycle will be detected. As this procedure will take some time, there is a chance that one of the transactions that holds a lock will mean while have released it, in which case the deadlock will no longer exist.

For example : Consider the case of a global deadlock detector that receives local wait-for graphs from servers X and Y, as shown in Fig. 12. Suppose that transaction U then releases an object at server X and requests the one held by V at server Y. Suppose also that the global deadlock detector receives server Y's local graph before server X's. In this case, it would detect a cycle  $T \rightarrow U \rightarrow V \rightarrow T$ , although the edge  $T \rightarrow U$  no longer exists.

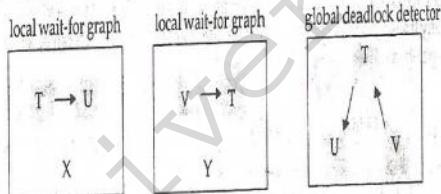


Fig. 12. Local and global wait-for graphs.

A phantom deadlock could be detected if a waiting transaction in a deadlock cycle aborts during the deadlock detection procedure. For example, if there is a cycle  $T \rightarrow U \rightarrow V \rightarrow T$  and U aborts after the information concerning U has been collected, then the cycle has been broken already and there is no deadlock.

#### Q. 22. Define Edge Chasing in Distributed Deadlock.

Ans. Edge chasing or path pushing is used in distributed deadlock detection. In this, the global wait-for graph is not constructed, but each of the servers involved has knowledge about some of its edges. The servers attempt to find cycles by forwarding messages called probes, which follows the edges of the graph throughout the distributed system. A probe message consisting of transactions wait-for relationship representing a path in the global wait-for graph.

Edge chasing has three steps:

- Initiation : The server initiates to detect deadlock.
- Detection : Detection consists of receiving probes and deciding whether deadlock has occurred and forward probes.
- Resolution : When a cycle detected. A transaction in the cycle is aborted to break deadlock.

#### Q. 23. Discuss Transaction Recovery and different methods of Transaction Recovery.

Ans. The atomic property of transactions requires that the effects of all committed transactions and none of the effects of incomplete or aborted transactions are reflected in the objects they accessed. This property can be described in terms of two aspects : durability and failure atomicity. Durability requires that objects are saved in permanent storage and will be available indefinitely thereafter. Therefore, an acknowledgement of client's commit request implies that all the effects of the transaction have been recorded in permanent storage as well as in the server's (volatile) objects. Failure atomicity requires that effects of transactions are atomic even when the server crashes. Recovery is concurrent with ensuring that a server's objects are durable and that the service provides failure atomicity.

Although file servers and database servers maintain data in permanent storage, other kinds of servers of recoverable objects need not do so except for recovery purposes. When a server is running it keeps all of its volatile memory and records its committed objects in a recovery file or files. Therefore, recovery consists of restoring the server with the latest committed versions of its objects from permanent storage. Databases need to deal with large volumes of data. They generally hold the objects in stable storage on disk with a cache in volatile memory. The two requirements for durability and for failure atomicity are not really independent of one another and can be dealt with by a single mechanism—the recovery manager. The task of a recovery manager is:

- ✓ To save objects in permanent storage (in a recovery file) for committed transactions;
- ✓ To restore the server's objects after a crash;
- ✓ To reorganize the recovery file to improve the performance of recovery;
- ✓ To reclaim storage space (in the recovery file).

#### Q. 24. Write short note on:

- Logging.
- Shadow versions.

Ans. (i) Logging : In the logging technique, the recovery file represents a log containing the history of all the transactions performed by a server. The history consists of values of objects, transaction status entries and intentions lists of transactions. The order of the entries in the log reflects the order in which transactions have prepared, committed and aborted at that server.

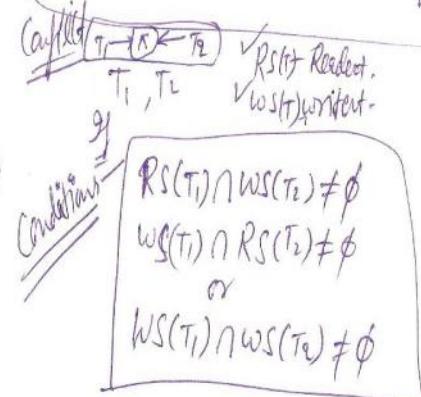
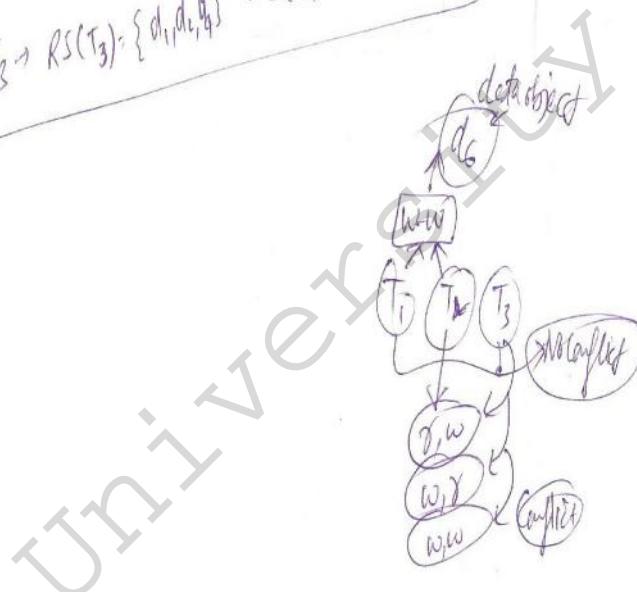
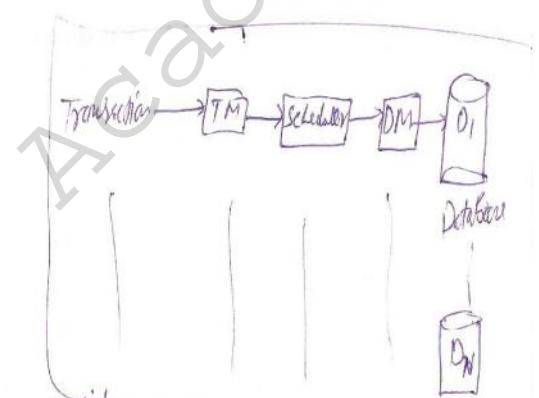
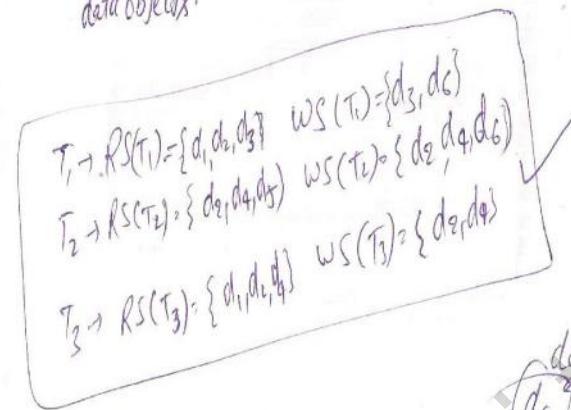
During the normal operation of a server, its recovery manager is called whenever a transaction prepares to commit, commits or aborts a transaction. When the server is prepared to commit a transaction, the recovery manager appends all the objects in its intentions list to the recovery file, followed by the current status of that transaction (prepared) together with its intentions list. When a transaction is eventually committed or aborted, the recovery manager appends the corresponding status of the transaction to its recovery file.

(ii) Shadow Versions : The shadow versions technique is an alternative way to organize a recovery file. It uses a map to locate versions of the server's object in a file called a version store. The map associates the identifiers of the server's versions in the version store. The versions written by each transaction are shadows of the previous committed versions.

## DATABASE SYSTEM

(19)

- In database system users concurrently access data object by executing transactions. Concurrency control is the process of controlling concurrent access to a database to ensure that the consistency of the database is maintained.
- A Database system consists of a set of shared data objects that can be accessed by users. A data object can be a page, a file, a segment or a record.
- For the purpose of concurrency control, we will view a database as a collection of data objects ( $d_1, d_2, d_3, \dots, d_n$ ).
- Each data object takes values from a specified domain. The state of database is given by the values of its data objects.
- In database contain semantic relationship, called consistency assertions or integrity constraints must hold among its data objects.



## S/wo Model - for the purpose of Concurrency Control

✓ Transaction Manager (TM)

✓ Data Manager (DM)

Transaction - TM - Scheduler - DM - Database

✓ Scheduler

$T_1$   
lock A  
 $\underline{A \neq 0 \rightarrow 1}$

$T_2$   
lock B  
 $\underline{B \neq 0 \rightarrow 1}$

Data object

A > 0

B > 0

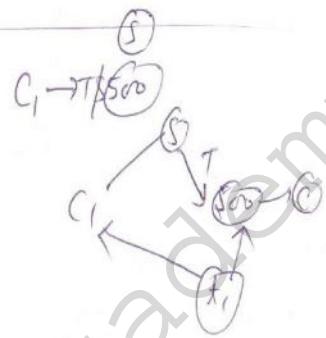
$T_1 = \text{if } A > 0 \text{ then } B := B + 1$   
 $T_2 = \text{if } B > 0 \text{ then } A := A + 1$

$S = U_{i=0}^n S_i$

$D_i = \{d_j | j \in [1..M]\}$

$D_i \cap D_j = \emptyset$

$D_1 \cup D_2 \cup D_3 \cup D_N = \{d_1, d_2, \dots, d_M\}$



$S = 1000$   
 $C = 500$

$C_1 \rightarrow 500$  (1000 - 500) at ①  
 $T_1 \rightarrow S (500) \& C (500)$  at ②  
 $\rightarrow 1000 - T_1$

When a transaction is prepared to commit, any of the objects changed by the transaction are appended to the version store, leaving the corresponding committed versions unchanged. These new as yet tentative versions are called shadow versions. When a transaction commits, a new map is made by copying the old map and entering the positions of the shadow versions. To complete the commit process, the new map replaces the old map.

To restore the objects when a server is replaced after a crash, its recovery manager reads the map and uses the information in the map to locate the objects in the version store. The shadow version method provides faster recovery than logging because the positions of the current committed objects are recorded in the map, whereas recovery from a log requires searching throughout the log for objects.

Shadow versions on their own are not sufficient for a server that handles distributed transactions. Transaction status entries and intentions lists are saved in a file called the transaction status file. Each intentions list represents the part of the map that will be altered by a transaction when it commits.

**Q. 25. Explain Recovery of the Two-Phase Commit Protocol.**

**Ans.** In phase 1, when coordinator is prepared to commit, its recovery manager adds a coordinator entry to its recovery file. Before participants votes Yes, it must have already prepared to commit. When it votes Yes, its recovery manager records a participant entry and add uncertain transaction status to its recovery file as a forced write. When votes No, it adds abort transaction status to its recovery file.

In phase 2, of the protocol, the recovery manager of the coordinator adds either a committed or aborted transaction status to its recovery file, according to the decision. This must be a forced write. Recovery manager of the participants add a commit or abort transaction status to their recovery file according to message received from coordinator. When a coordinator has received a confirmation from all of its participants, its recovery manager adds a done transaction status to its recovery file—this need not to be forced.

The done status entry is not part of the protocol but is used when the recovery file is recognized.

**Q. 26. What do you understand by Replication?**

**Ans.** Replication is a key to providing high availability and fault tolerance in distributed systems. High availability is of increasing interest with the tendency towards mobile computing and consequently disconnected operation. Fault tolerance is an abiding concern for services provided in safety-critical and other important systems.

The replication is "the maintenance of copies of data at multiple computers. Replication is a key to the effectiveness of distributed systems in that it provides enhanced performance, high availability and high fault tolerance. For example, the caching of resources from web server, from web servers in browser and web proxy servers is a form of replication, since the data held in cache and at servers are replicas of one another."

Replication is a technique for enhancing a service. The motivations for replication are to improve a service's performance, to increase its availability, or to make it fault tolerant.

When data are replicated the replication transparency is required, i.e., clients should not normally have to aware that multiple copies of data exist. The consistency of data is another requirement.

**Q. 27. Explain the concept of System Models and Group Communication.**

**Ans.** System model : The model involves replicas held by distinct replica managers, which are components that contain the replicas on a given computer and perform operations upon them.

directly. This general model may be applied in a client-server environment, in which case a replica manager is a server. We shall sometimes simply call them servers instead. Equally, it may be applied to an application and application processes can in that case act as both clients and replica managers. For example, the user's laptop on a train may contain an application that acts as a replica manager for their diary.

We shall always require that a replica manager applies operations to its replicas recoverably. This allows us to assume that an operation at a replica manager does not leave inconsistent results if it fails part-way through. The set of replica managers may be static or dynamic.

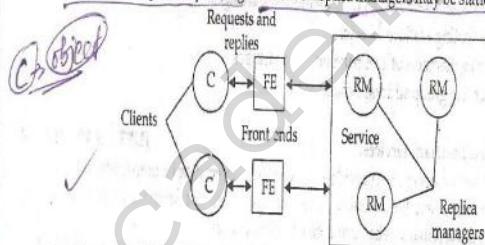


Fig. 13. Architectural model for the management of replicated data.

The general model of replica management is shown in above figure. A collection of replica managers provides a service to clients. The client see a service that gives them access to objects (for example, diaries or bank accounts), which in fact are replicated at the managers. Clients each request a series of operations—invocations upon one or more of the objects. An operation involves no updates are called read-only requests; requested operations that update an object are called update requests (these may also involve reads).

**Group Communication :** Multicast communication is also known as group communication because process groups are the destinations of multicast messages. Groups are useful for managing replicated data and in other systems where processes cooperate towards a common goal by receiving and processing the same set of multicast messages. They are also useful where the group members independently consume one or more common streams of messages, such as messages carrying events to which the processes react independently.

In a service that manages replicated data, for example, users may add or withdraw a replica manager, or a replica manager may crash and thus need to be withdrawn from the system's operation. A full implementation of group communication incorporates a group membership service to manage the dynamic membership of groups, in addition to multicast communication.

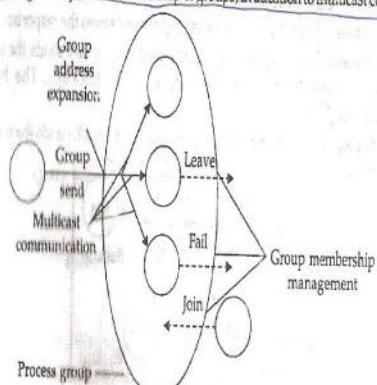


Fig. 14. Services provided for process groups.

Multicast and group membership management are strongly interrelated. Fig. 14 shows an open group, in which a process outside the group sends to the group without knowing the group's membership. The group communication services has to manage changes in the group's membership while multicasts take place concurrently.

**Q. 28. What are Roles of the Group Membership Service?**

Ans. A group membership service has four main tasks which are as follows:

- Providing an interface for group membership changes.
- Implementing a failure detector.
- Notifying members of group membership changes.
- Performing group address expansion.

**Q. 29. Explain fault-tolerant services.**

[UPTU 2010, Marks 5]

Ans. For fault tolerant services we provide a service that is correct despite upto/ process failures, by replicating data and functionality at replica managers. For simplicity, we assume that communication remains reliable and no partitions occurs.

Each replica manager is assumed to behave according to a specification of the semantics of the objects it manages, when they have not crashed. For example, a specification of bank accounts would include an assurance that funds transferred between bank accounts can never disappear, and that only deposits and withdrawals affect the balance on any particular account.

(i) **Passive (primary-backup) Replication:** In a passive model of replication for fault tolerance, there is at any time a single primary replica manager and one or more secondary replica managers—'backups' or 'slaves'. In the pure form of the model, front ends communicate only with the primary replica manager to obtain the service. The primary replica manager executes the operations and sends copies of the updated data to the backups. If the primary fails, one of the backups is promoted to act as the primary.

The sequence of events when a client requests an operation to be performed is as follows:

- a. **Request:** The front end issues the request, containing a unique identifier, to the primary replica manager.
- b. **Coordination:** The primary takes each request atomically, in the order in which it receives it. It checks the unique identifier, in case it has already executed the request and if so it simply re-sends the response.
- c. **Execution:** The primary executes the request and stores the response.
- d. **Agreement:** If the request is an update then the primary sends the updated state, the response and the unique identifier to all the backups. The backups send an acknowledgement.
- e. **Response:** The primary responds to the front end, which hands the response back to the client.

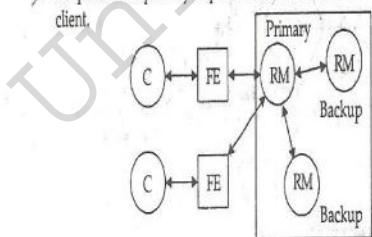


Fig. 15. The passive (primary-backup) model for fault tolerance.

(2) **Active Replication:** In the active model of replication for fault tolerance, the replica managers are state machines that play equivalent roles and are organized as a group. Front ends multicast their requests to the group of replica managers and all the replica managers process the request independently but identically and reply. If any replica manager crashes, then this need have no impact upon the performance of the service, since the remaining replica managers continue to respond in the normal way.

Under active replication, the sequence of events when a client requests an operation to be performed is as follows:

- a. **Request:** The front end attaches a unique identifier to the request and multicasts it to the group of replica managers, using a totally ordered, reliable multicast primitive.
- b. **Coordination:** The group communication system delivers the request to every correct replica manager in the same (total) order.
- c. **Execution:** Every replica manager executes the request.
- d. **Agreement:** No agreement phase is needed, because of the multicast delivery semantics.
- e. **Response:** Each replica manager sends its response to the front end.

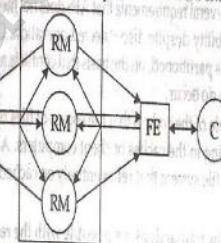


Fig. 16. Active replication.

**Q. 30. Discuss Highly Available Services.**

[UPTU 2010, Marks 5]

OR

Write a short note on:

- (i) Gossip architecture.
- (ii) Bayou and operational transformation approach.
- (iii) Coda file system.

Ans. The systems that provide highly available services are:

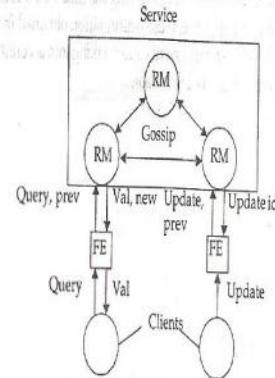


Fig. 17. Query and update operations in a gossip service.

(i) **Gossip architecture**: Gossip architecture is a framework for implementing highly available services by replicating data close to the points where groups of clients need it. The name reflects the fact that the replica managers exchange 'gossip' messages periodically in order to convey the updates they have each received from clients.

A gossip service provides two basic types of operation: queries are read-only operations and updates modify but do not read the state. A key feature is that front ends send queries and updates to any replica manager they choose—any that is available and can provide reasonable response times.

(ii) **Bayou and the operational transformation approach**: The Bayou system provides data replication for high availability with weaker guarantees than sequential consistency. Bayou replica managers cope with variable connectivity by exchanging updates in pairs, in what the designers also call an anti-entropy protocol. But Bayou adopts a markedly different approach in that it enables domain-specific conflict detection and conflict resolution to take place.

(iii) **Coda file system**: The coda file system is a descendent of Andrew File System (AFS) that aims to address several requirements that AFS does not meet, particularly the requirement to provide high availability despite disconnected operation. Coda allows clients to update data while the system is partitioned, on the basis that conflicts are relatively unlikely and that they can be fixed if they do occur.

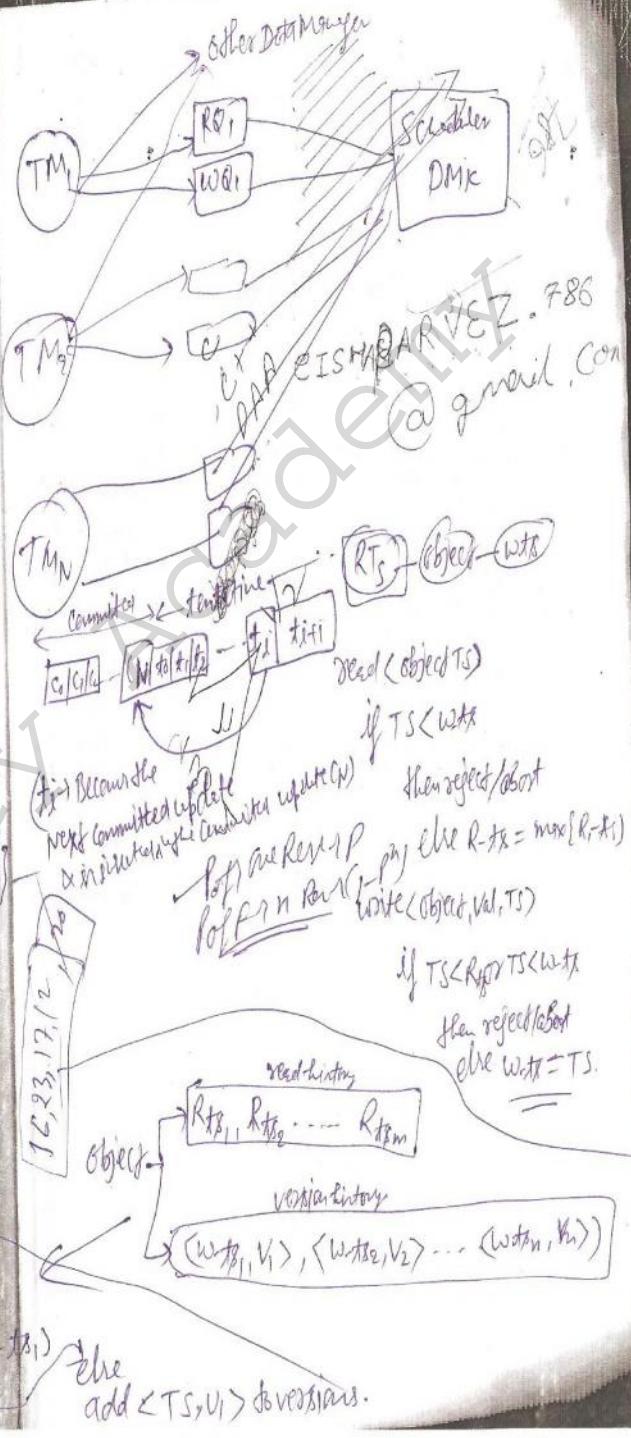
The principle of the coda is that the copies of files residing on servers are more reliable than those residing in the caches of client computers. Although it might be possible logically to construct a file system that relies entirely on cached copies of files in client computers.

Q.31. Explain, how transactions are possible with the replicated data?

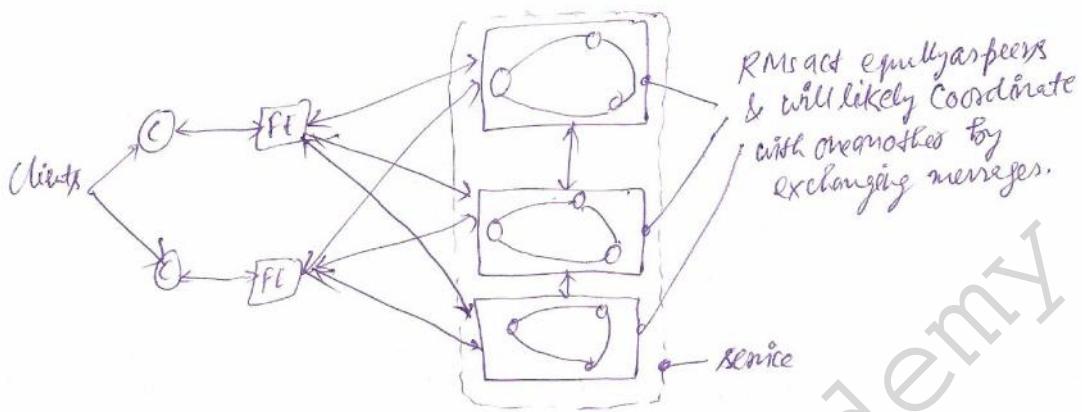
Ans. Objects in transactional systems may be replicated to increase both availability and performance. From a client's view point, a transaction on replicated objects should appear the same as one with non-replicated objects. In a non-replicated system, transactions appear to be performed one at a time in some order. This is achieved by ensuring a serially equivalent interleaving of client's transactions. The effect of transactions performed by clients on replicated objects should be the same as if they had been performed one at a time on a single set of objects. This property is called one-copy serializability.

Each replica manager provides concurrency control and recovery of its own objects. Recovery is complicated by the fact that a failed replica manager is a member of a collection and that the other members continue to provide a service during the time that it is unavailable. When a replica manager recovers from a failure, it uses information obtained from the other replica managers to restore its objects to their current values, taking into account all the changes that have occurred during the time it was unavailable.

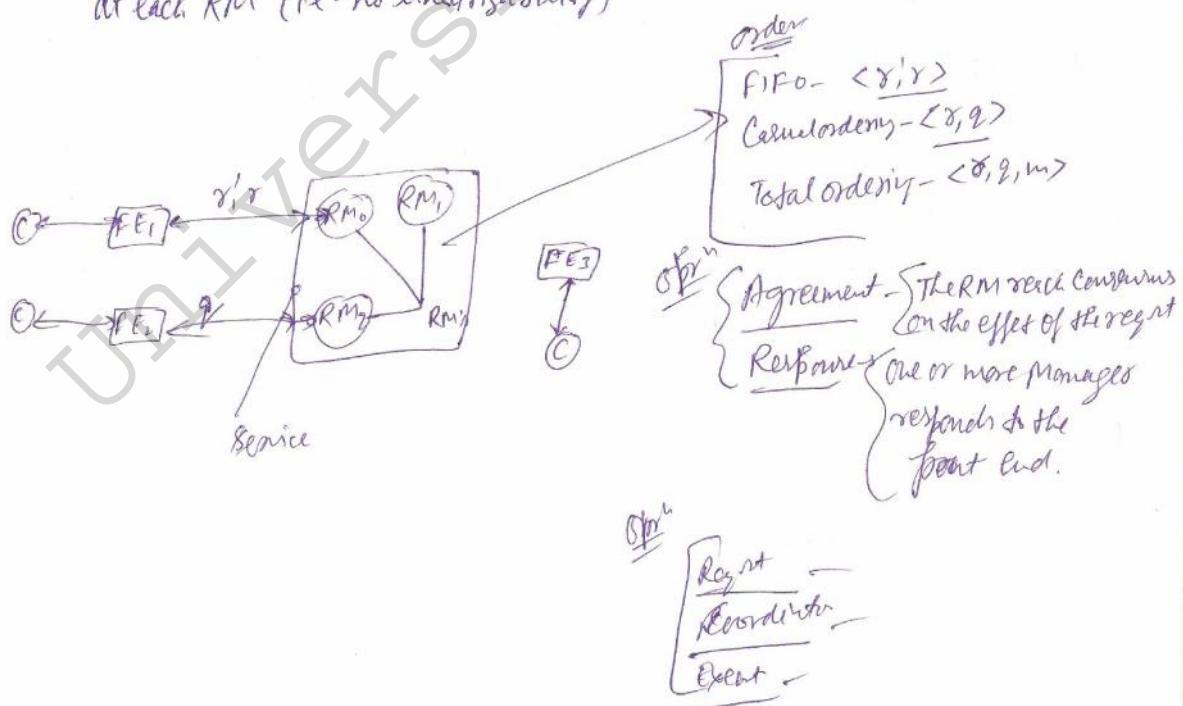
**read (Object, TS)**  
 read  $V_j$  where  $j = \max\{i | w_{ti} < TS\}$   
**add (TS) to read history**  
**written object, Val, TS**  
 If there is a K such that  
 $TS < R_{tk} < W_{tj}$  where  $j = \min\{i | w_{ti} < TS < R_{tk}\}$   
 Then reject op.



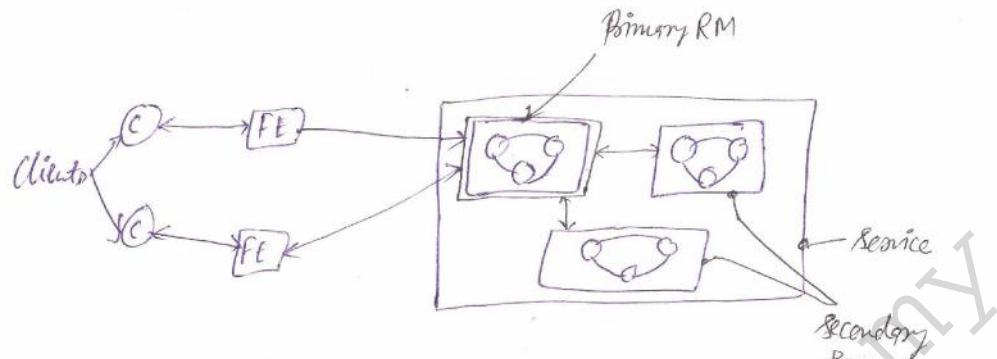
## Active Replication



- Client Requests are Multicast (unidirectionally ordered multicast) by the FE to All RMs in the Service Group.
- FEs collect responses from all non-failed RMs & form a single response to the client.
- Although Client requests are handled in order (i.e. with Sequential Consistency), due to arbitrary delays, they may interfere with other clients' requests differently at each RM (i.e. no linearizability)

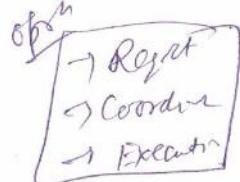


## Passive (Primary Backup) Replication



- Clients communicate, via their FE, with the Primary RM (e.g. read or update to objects)
- Primary sends updates to the Secondary RMs so that they can replace the Primary if it fails, perhaps through some pre-arrangement or an election.
- To maintain consistency during primary replacement, RMs can implement view synchronization, so they know which options have been processed up until the failure.
- Since all of the client operations are handled by the primary RM, if concurrent, that exactly the same sequencing is executed by its secondary RMs, so this technique can achieve linearizability.

### Replica Manager (RM)



- one replica manager per replica
- Receive FE's reqt apply op<sup>n</sup> to its replica atomically.

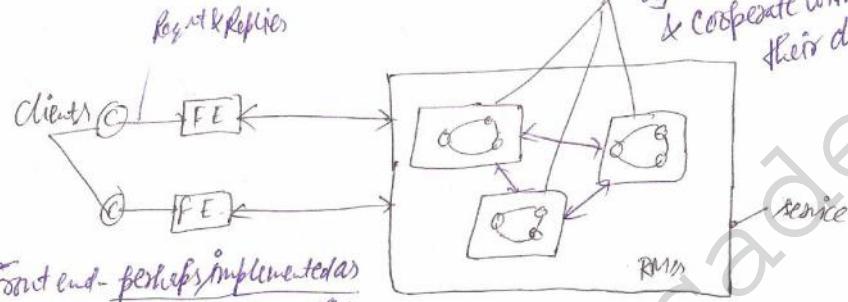
### Frontend (FE)

- one front end per client
- Receive client's reqt communicate with RM by message passing

## Replication

Replicated Service -

- Model of a Replicated Service.
- View Synchronization Communication.
- Linearizability & Sequential Consistency.



Front end - perhaps implemented as  
Middleware at the Client, mediates  
bet'w the client's opns & the replicated  
Services to provide transparency.

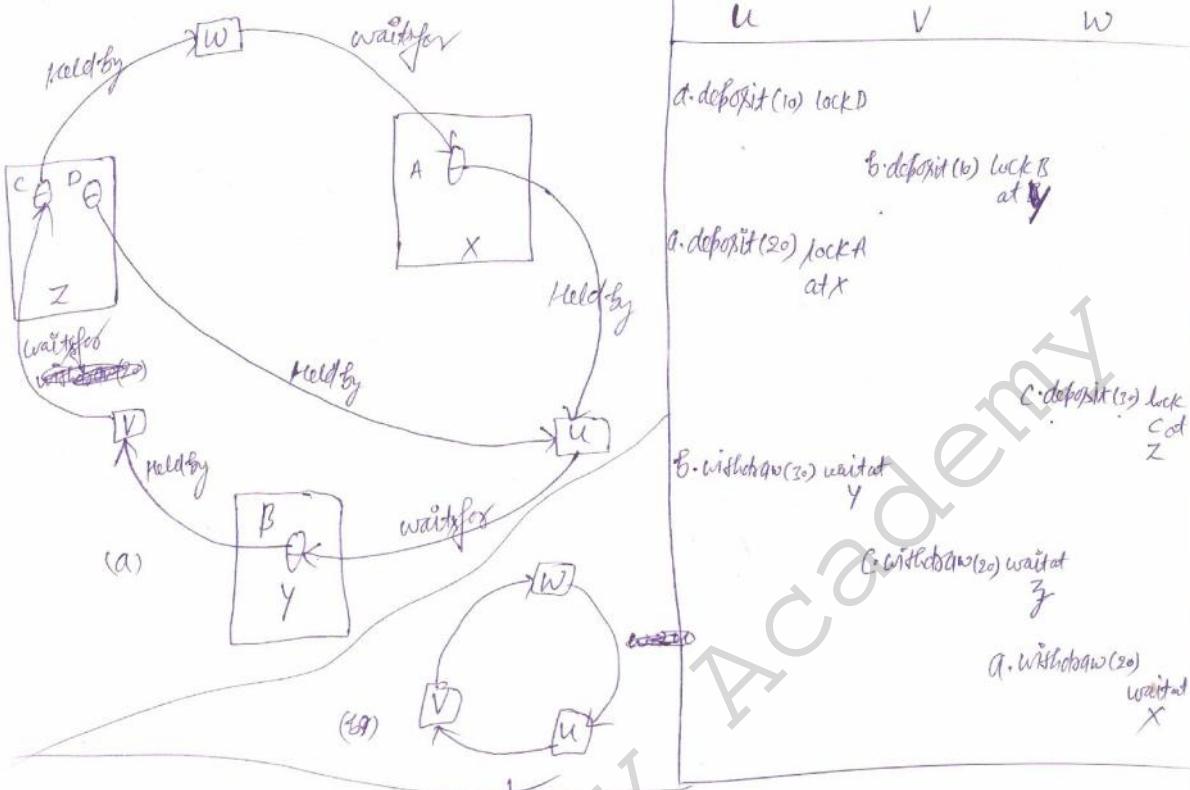
Replica Manager(s) — Maintains Recoverable  
object (i.e. that support tentative opns)  
& cooperate with one another to keep  
their data consistent.

In general, there are 5 phases involved in handling the client request -

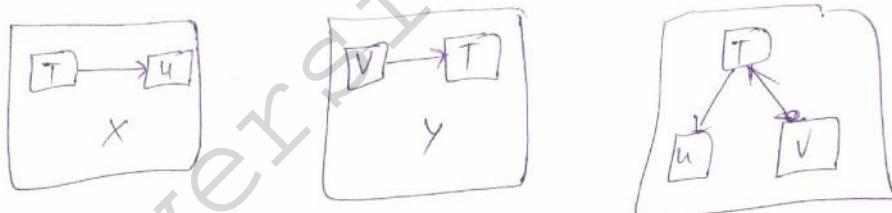
- ① Request - the FE communicates the reqt to one or more of the RM's (depending on the particular replication system) on behalf of the client.
- ② Coordination - the RM's coordinate to handle the request consistently (e.g. agreement on ordering of requests, whether or not the request is valid.)
- ③ Execution - the RM's execute the request, perhaps tentatively (i.e. to allow for recovery).
- ④ Agreement - if appropriate (e.g. for transactions), the RM's will reach a consensus on whether to commit or abort the update.
- ⑤ Response - one or more of the RM's return a response to the FE.  
The FE will then deliver the Response to Client.

### Group Membership View

- An RM may join or leave or fail.
- RMs need to agree on the current membership of the group so that they can be confident not only of who they should be communicating with but also who the other RMs are communicating with.
- This can be achieved by the notion of a (group membership) view, which is simply an agreed list of current RMs (e.g. IP addresses, process IDs, etc.).
- When a RM notices that another RM has failed (or receives a join request from a new RM), it multicasts a new view to the other RMs.
- Importantly, through, the RMs do not immediately treat the new view as the current view until all RMs in the new view have acknowledged to other RMs that they got it.



Deadlock wait for graphs



WFG of these trans.

At server Y:  $U \rightarrow V$  (added when U reqt B withdraw)

$Z: V \rightarrow W$  (" " " "  $V = c.\text{withdraw}(3)$ )

$X: W \rightarrow U$  (" " " "  $W = a.\text{withdraw}(2)$ )

$\rightarrow$  Transaction  $U, V, W$  involves objects  $A, B, C, D$ .

- ①  $W \rightarrow U \quad (\exists X \rightarrow Y)$
- ②  $W \rightarrow U \rightarrow V \quad (Y - 3)$
- ③  $W \rightarrow U \rightarrow V \rightarrow W$

~~① Server X initiate detection of deadly probe  $\langle W \rightarrow U \rangle$ , notes that B is held by V & appends V to the probe.~~

① Server X initiates detection by sending probe  $\langle W \rightarrow U \rangle$  to the server of B [Server Y]  
 ② Server Y receives probe  $\langle W \rightarrow U \rangle$ , notes that B is held by V & appends V to the probe to produce  $\langle W \rightarrow U \rightarrow V \rangle$ . ~~If Notes~~ that V is waiting for C at Server Z.  
 This probe is forwarded to Z.

③ Server Z receives probe  $\langle W \rightarrow U \rightarrow V \rangle$ , notes C is held by W & appends W to the probe to produce  $\langle W \rightarrow U \rightarrow V \rightarrow W \rangle$ .  
 → This path contains cycle. The server detects a deadlock. One of the transactions in the cycle must be aborted to break the deadlock.

| Coordinator of Transaction | Child transactions | Participant | Provisional Commit list           | Abort list                    |
|----------------------------|--------------------|-------------|-----------------------------------|-------------------------------|
| T                          | $T_1, T_2$         | yes         | $T_1, T_{12}$                     | <del>A</del> $T_{11}, T_2$    |
| $T_1$                      | —                  | —           | —                                 | —                             |
| $T_2$                      | —                  | —           | —                                 | —                             |
| $T_{11}$                   | —                  | no (aborts) | —                                 | —                             |
| $T_{12}, T_{21}$           | —                  | —           | $T_{12} \text{ by } Not T_{21}^*$ | <del>A</del> $T_{21}, T_{12}$ |
| $T_{22}$                   | —                  | —           | No (Presentation)                 | —                             |

$\times (T_2, \text{ Present but aborted})$

### Replica Managers

- one Replica Manager per replica
- Receives replica request, apply operations to its replicas automatically.

Front End's Request Communication (Request can be made to a single RM or to Multiple RMs).

- One FE per client
- Receives client's request, communicate with RM by message passing.

### Coordinator's RM's decide

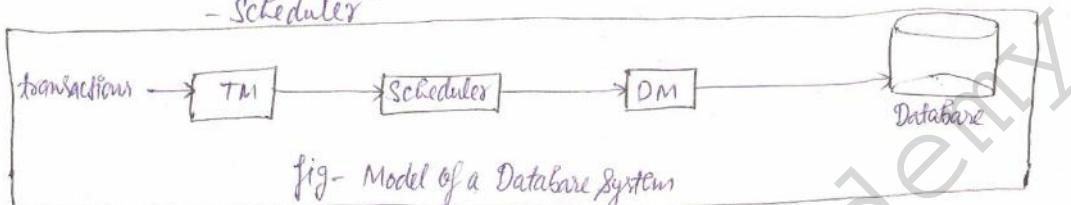
- The order of request
- FIFO -  $[r, r']$  first  $r$  then  $r'$  then good & timely
- Camel ordering -  $[r, r'] \rightarrow r$  happens before  $r'$
- Total ordering -  $[r, r'] \rightarrow$  correct  $r$  & then  $r'$

Unit-5

## A Concurrency Control Model of Database Systems

(1)

- For the purpose of concurrency control, we can view a database system as consisting of three S/w modules:-
  - Transaction Manager (TM)
  - Data Manager (DM)
  - Scheduler



- The transaction Manager supervises the execution of a transaction. It intercepts & executes all the submitted transactions.
- A TM interacts with the DM to carry out the execution of a transaction.
- It is the responsibility of the TM to assign a timestamp to a transaction or issue requests to lock & unlock data objects on behalf of a user.
- Thus we can say that TM is an interface bet<sup>1</sup> user & the database system.
- The Scheduler is responsible for Concurrency Control. It grants or releases locks on data objects as requested by a transaction.
- The Data Manager Manages the database. ~~be operating them~~ It carries out the read-write requests issued by the TM on behalf of a transaction by operating them on the database.
- DM is an interface bet<sup>2</sup> the Scheduler & the database. DM executes a stream of transaction actions, directed toward it by the TM.

## The Problem of Concurrency Control :-

- Typically, in a Database System, several transactions are under execution simultaneously.
- Since a transaction preserves database consistency, a database system can guarantee consistency by executing transactions serially, i.e. one at time.
- Such a serial execution of transactions is inefficient as it results in poor response to user requests & poor utilization of system resources.
- Efficiency can be improved by executing transactions concurrently, that is, by executing read & write actions from several transactions in an interleaved manner.
- Because the actions of concurrently running transactions may access the same data objects, several anomalous situations may arise if the interleaving of actions is not controlled in some orderly way.

Some Situations are -

### ① Inconsistent Retrieval -

- It occurs when a transaction reads some data objects of a database before another transaction has completed with modification of those data objects.
- In such transaction, the former transaction faces the risk of retrieving incorrect values of the data objects.

Ex -

$$S = 1000 \quad (-500) \quad C = 500$$

- customer  $t_1$  → reads  $S$  into its workspace & subtract 500 from it & writes  $S$ .
  - teller  $t_2$  → reads  $S$  ( $= 500$ ) &  $C$  ( $= 500$ ) into its workspace
  - $t_1$  → reads  $C$  ( $= 500$ ) into its workspace & add 500 to it, & write back  $C$  ( $= 1000$ ).
  - $t_1$  → Output 1000 as the balance
- Here,  $t_2$  reads  $S$  after  $t_1$ , has modified it & reads  $C$  before  $t_1$  has modified it resulting the incorrect retrieval of the total balance.

### ② Inconsistent Updates occurs when many transactions read & write onto a common set of data objects of a database, leaving the database in an inconsistent state.

ex (inconsistent state) →

Let, data object - A, B

Consistent condition → "(A=0) OR (B=0)"

Concurrent Modified transaction -

" $T_1$  : if  $A=0$  then  $B=B+1$ "

" $T_2$  : if  $B=0$  then  $A=A+1$ ".

Now,

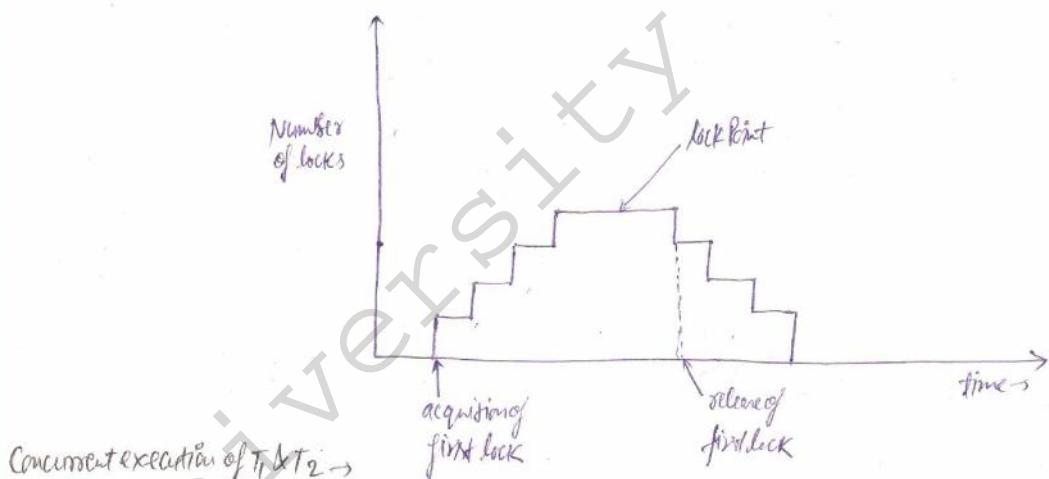
- A Possible execution trace is -

initially  $(A=0 \wedge B=0)$

$T_1$  reads  $A(=0) \wedge B(=0)$ , since  $A=0$  in  $T_1$  so it increments  $B$  by 1 & writes in the Database ( $B=1$ )

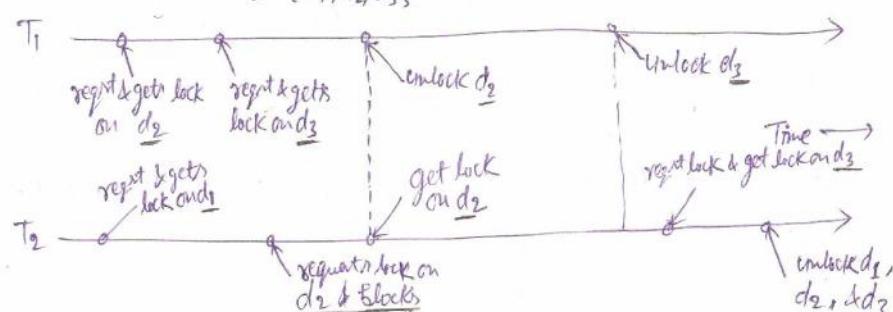
$T_2$  reads  $A(=0) \wedge B(=0)$ , since  $B=0$ , in  $T_2$  — A "1" — ( $A=1$ ).

The final Database state " $(A=1) \wedge (B=1)$ " is inconsistent.



$$RS(T_1) = \{d_2, d_3\}, WS(T_1) = \{d_3\}$$

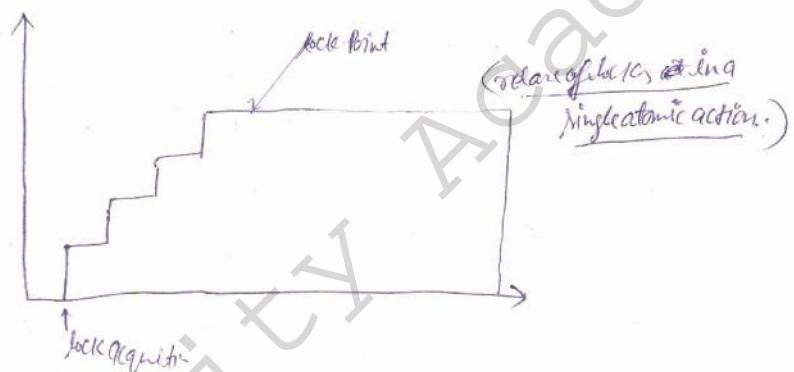
$$RS(T_2) = \{d_1, d_2, d_3\}, WS(T_2) = \{d_1, d_2, d_3\}$$



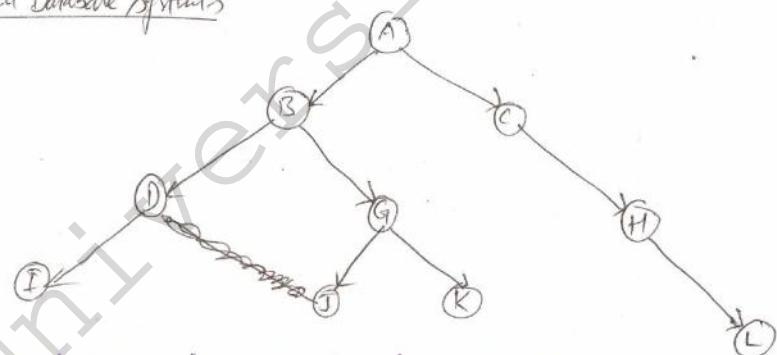
Cascaded Roll-backs → when a transaction Roll-back (for any reason - user kills it, the system crashes, or it becomes deadlocked), all the Data objects modified by it, ~~the system~~ are restored their original states.

- In this case, all transactions that have read the backed up data objects must also be rolled back & the data objects modified by them must also be restored & so on.
- This phenomenon is called the Cascaded Roll-back.
- 2PL suffers from the problem of cascaded roll-back.

Strict 2PL →



Hierarchical Database Systems



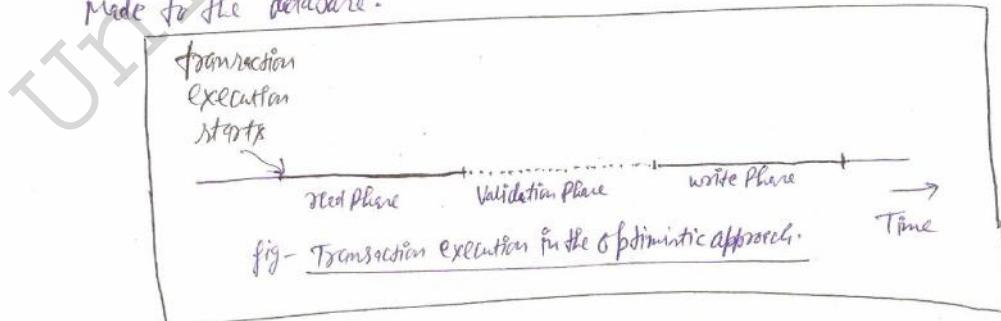
- T<sub>1</sub> must lock node B first so that it can lock any other data object (D, G, & I).
- T<sub>2</sub> must lock node A before accessing any other node. Before T<sub>2</sub> can lock H, it must first lock node C.

- |                                                                                                                                                                                                                                                                                                                                            |             |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------|
| <ul style="list-style-type: none"> <li>→ T<sub>i</sub> can lock data object R ≠ E(T<sub>i</sub>) iff T<sub>i</sub> is holding a lock on R's ancestor.</li> <li>→ After unlocking a data object, T<sub>i</sub> cannot lock it again.</li> <li>→ T<sub>i</sub> can only access those data objects for which it is holding a lock.</li> </ul> | (Root Node) |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------|

## # OPTIMISTIC CONCURRENCY CONTROL $\Rightarrow$ (KUNG-ROBINSON)

- In Optimistic algo, no synchronization is performed when a transaction is executed, but at the end of transaction's execution, a check is performed to determine if the transaction has conflicted with any other concurrently running transaction.
- In the case of conflict, the transaction is aborted, otherwise it is committed.  
(Comments)
- In this technique, a transaction always executes (tentatively) concurrently with other transactions without any synchronization check, but before its writes are written in the database (& become accessible to other transactions), it is validated.
- In Validation Phase, it is determined whether actions of the transaction have conflicted with those of any other transaction. If found in conflict, then the tentative writes of the transaction are discarded & the transaction is restarted.
- The execution of transaction is divided into 3 phases:  
① Read Phase ② Validation Phase ③ Write Phase

- In Read Phase, appropriate data objects are read, the intended computation of the transaction is done, & writes are made on a temporary storage.
- In Validation Phase, it is checked if the writes made by the transaction violate the consistency of the database.
- If check passes, then in the Write Phase, all the writes of the transaction are made to the database.



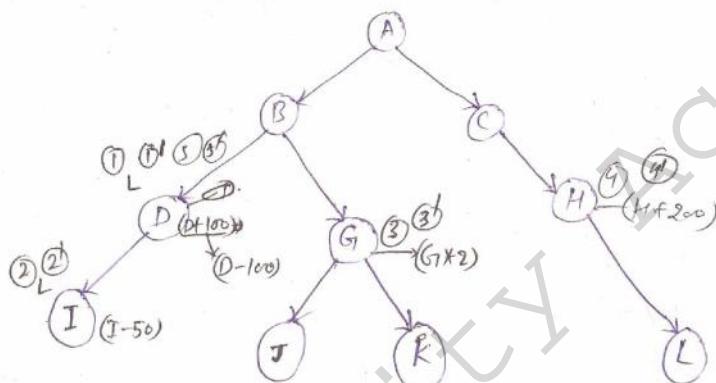
(3)

 $T_1$ 

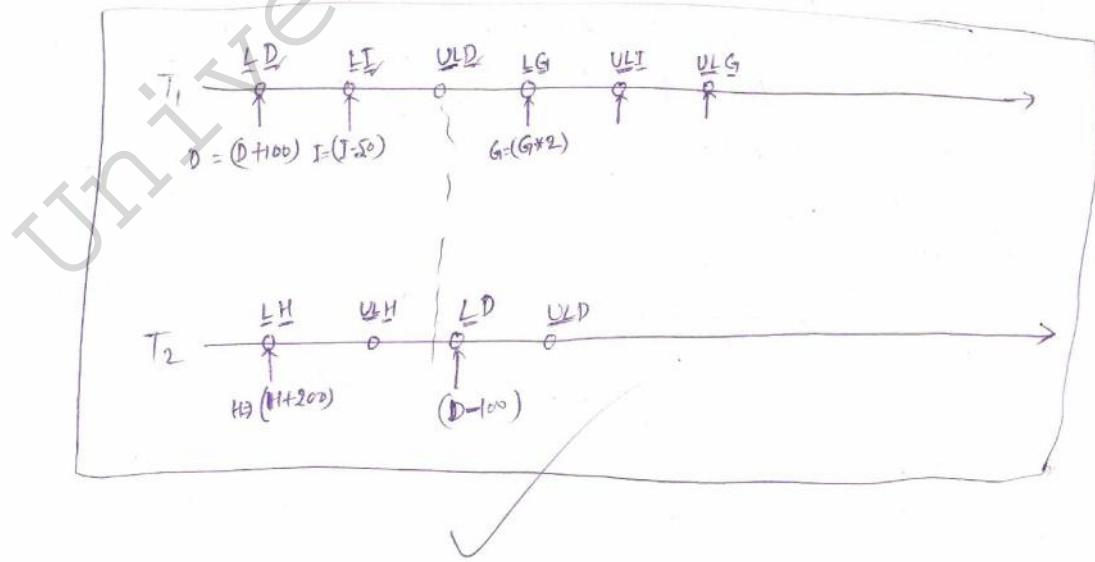
lock D  
 $D+100 \rightarrow D$   
 lock I  
 $I-50 \rightarrow I$   
 unlock D  
 lock G  
 $G_1 * 2 \rightarrow G_1$   
 unlock I  
 unlock  $G_1$

 $T_2$ 

LOCK H  
 $H+200 \rightarrow H$   
 unlock H  
 lock D  
 $D-100 \rightarrow D$   
 unlock D



~~①, ②, ③, ④, ⑤, ⑥, ⑦, ⑧, ⑨, ⑩, ⑪~~  
unlock UL UL



## Timestamp Based Algorithm

TM → attaching timestamp to each transaction  
DMs prevent conflicting opr in timestamp order

- timestamp based concurrency control algorithm, every site maintains a logical clock that is incremented when a transaction is submitted at that site & updated whenever the site receives a message with a higher clock value (every message contains the current clock value of its sender site).

- Timestamps can be used in two ways -

- determine the outdatedness of a request with respect to data objects.
- used to order events (read-write reqt) with respect to one another.

### ① Basic Timestamp ordering Algo's (BTO)

- The Scheduler at each DM keeps track of the largest timestamp of any read & write processed thus far for the each data object.

Let  
 $\textcircled{1} - \begin{cases} R-ts(\text{object}) - \text{read-timestamp.} \\ W-ts(\text{object}) - \text{write timestamp.} \end{cases}$

$\textcircled{2} - \begin{cases} \text{read}(x, TS) \\ \text{write}(x, v, TS) \end{cases}$  {read & write request with timestamp TS on a data object x. In write operation, v is the values to be assigned to x.}

- A read  $(x, TS)$  request is handled in the following manner -
  - If  $TS \leq W-ts(u)$ , then the read reqt is rejected & the corresponding transaction is aborted, otherwise it is executed &  $R-ts(u)$  is set to  $\max\{R-ts(u), TS\}$ .
- A write  $(x, v, TS)$  request is handled in the following manner -
  - If  $TS \leq R-ts(u)$  or  $TS < W-ts(u)$ , then write request is rejected, otherwise it is executed &  $W-ts(u)$  is set to TS.
  - If a transaction aborted it is started with new timestamp.

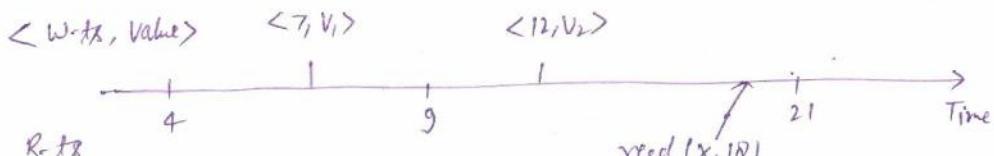
## ② Thomas Write Rule (TWR) $\rightarrow$ (Suitable for only write actions)

- for a write  $(X, V, TS)$ , if  $TS < W.ts(w)$ , the TWR says that instead of rejecting the write simply ignore it.
- this is sufficient to enforce synchronization among writers because the effect of ignoring an obsolete (No longer used) write request is the same as executing all writers in their timestamp ordering.
- TWR takes care of only write-write synchronization.
- TWR is an improvement over the BTO algorithm because it reduces the No of transaction aborts.

## ③ Multiversion Timestamp ordering Algo (MTO)

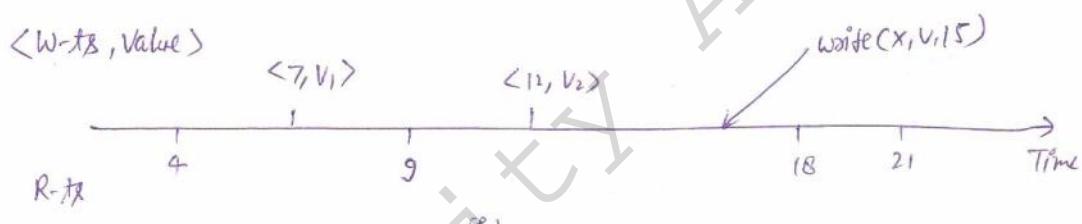
- A history of a set of R-txs & W-txs, value > pairs (called versions) is kept for each data object at the respective DM's.
- The R-txs of a data object keep track of the timestamp of all the executed read operations, & the versions keep track of the timestamp & the value of all the executed write ops.
- Read & write reqt are executed in following manner.
  - A read  $(X, TS)$  request is executed by reading the version of  $X$  with the largest timestamp less than  $TS$  & adding  $TS$  to the  $X$ 's set of R-txs. A Read request is never rejected.
  - A write  $(X, V, TS)$  request is executed in the following way.
    - If there exist a R-tx( $x$ ) in the interval from  $TS$  to the smallest ~~ts~~  $W.ts(w)$  that is larger than  $TS$ , then the write is rejected,
    - Otherwise it is accepted & a new version of  $X$  is created with timestamp  $TS$ .

Ex- (a) A read ( $x, 18$ ) is executed by reading the version  $\langle 12, v_2 \rangle$  & the resulting history is shown as -



(a)

Ex- (b) write ( $x, v, 15$ ) is rejected because a read with timestamp 18 has already been executed. However, a write ( $x, v, 22$ ) is accepted & is executed by creating a version  $\langle 22, v \rangle$  in the history.



(b)

#### ④ Commutative timestamp ordering Algorithms (CTO)

- It eliminates the abort & restarts of transactions by executing the requests in strict timestamp order at all DM's.
- A Scheduler processes a request when it is sure that there is no other request with smaller (older) timestamp in the system.
- Each Scheduler maintains 2 queues -
  - R-queue } per TM, these queue hold read & write requests.
  - W-queue }
- ATM sends requests to Schedulers in timestamp order & the communication medium is order preserving.
- A Scheduler puts a new read or write request in the corresponding queue in timestamp order.

This algo executes read & write actions in the following way -

① A  $\text{read}(x, TS)$  request is executed in following way -

\* If every w-queue is nonempty & the first write on each w-queue has a timestamp greater than TS, then the read is executed, otherwise the  $\text{read}(x, TS)$  request is buffered in the R-queue.

② A  $\text{write}(x, v, TS)$  request with timestamp TS is executed in following manner.

\* If all R-queues & all w-queues are nonempty & the first read on each R-queue has a timestamp greater than TS & the first write on each w-queue has a timestamp greater than TS, then the write is executed, otherwise the  $\text{write}(x, v, TS)$  request is buffered in the w-queue.

③ When any read or write reqt is buffered or executed, buffered reqt is checked to see if any of them can be executed.

Problem with CTOs

① Termination is not guaranteed.

② The algorithm is overly conservative.

