

# Distributed Systems

Notes by Rohan Khurana

## Distributed Systems

- ↳ A system in which components are connected by a communication network, communicate and coordinate their actions by passing msg.

Definition: A Distributed System is defined as the no of autonomous processing elements (not necessarily homogeneous) that are interconnected by a computer Network, cooperate in processing their assigned tasks transparently.

### HSSC

#### Challenges:

- 1) Heterogeneity: (DS Must be constructed from variety of diff. OS, Comp., H/w, lang.)
- 2) Openness: (DS must be Extensible, reimplemented in diff. ways)
- 3) Security: (CIA of Info. Resources, using Encryption)
- 4) Scalability: (DS remain effective when increase in Resources and Users)
- 5) Concurrency: (Diff. users accessing shared Resources)
- 6) Transparency: (Existence is transparent to others)

# Motivation: Resource Sharing

Behind DS

# Examples: Internet, Intranet, Finance Trading, Online Games, Mobile & Ubiquitous Computing.

### Limitations of DS:

- 1) Absence of Global clock
- 2) Unprecedented message passing
- 3) Non-existence of physical shared memory.
- 4) Initial Deployment cost is high

Logical clock

- (1) Logical clock refer to implementing a protocol on all Machine within the DS, so that the machines are able to maintain consistent ordering of Event within some virtual ts.
- (2) A logical clock is a mechanism of capturing chronological & causal relationship in DS.

Happened Before Relation ( $\rightarrow$ ):

$a \rightarrow b$ , 'a' happened before b.

a, b e. same process, a sending msg, b receiving  
 If  $a \rightarrow b$ ,  $b \rightarrow c$ .  
 $\Rightarrow a \rightarrow c$  (transitive)

Logical clock

Let  $P_i$  is process,  $c_i$  is vector clock associated with it

for event a,  $c_i(a) = TS$  of event a.

[01] for any 2 event, a & b in Process  $P_i$  if a happened before b ( $a \rightarrow b$ ) then.

$$c_i(a) \prec c_i(b)$$

$$c_i(a) < c_i(b)$$

[02] if a is event of sending msg in  $P_i$  and b is event of receiving msg in at process  $P_j$ .

$$c_i(a) \prec c_j(b)$$

$$c_i(a) < c_j(b)$$

$$C_j = \max(c_j, t_m + d)$$

Date \_\_\_\_\_

Page No.: \_\_\_\_\_

$C_i$  is incremented b/w any 2 successive events

[IRD]

If  $a \rightarrow b$  in process  $P_i$

then

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

$d$  = drift time (generally 1).

[IR2]

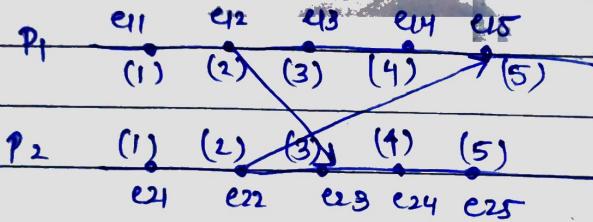
If there are more no. of processes & event  $a$  is sending of msg  $m$  in  $P_i$   
then  $T_s - t_m = C_i(a)$

On receiving same msg  $m$  in process  $P_j$   
then

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

Example

space ↑



$$C_{P_1} = 0 \quad d = 1$$

$$C_{P_2} = 0$$

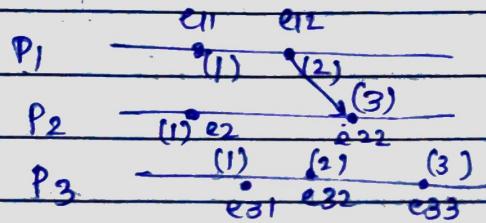
$$C_j = \max(c_j, t_m + d)$$

## Limitations of Lamport's clock

- 1) Lamport's system of logical clock states, if  $a \rightarrow b$  then  $c(a) < c(b)$

However the reverse is not necessarily true if events have occurred in diff processes

- 2) thus if  $a \leq b$  are events in different processes &  $c(a) < c(b)$  then  $a \rightarrow b$  is not always true.  
eg.



$$c(e_{11}) < c(e_{12})$$

$$c(e_{11}) < c(e_{21})$$

e<sub>11</sub> is causally related to

e<sub>22</sub> but not e<sub>32</sub>

- 3) since each clock independently advance due to occurrence of local events in process  
 Lamport clock system cannot distinguish b/w advancement of clocks due to local events from those due to exchange of msg. b/w processes.

∴ By using TS assign by LCS we cannot reason about the causal ordering of 2 events occurring in diff processes.

generate  
detect

partial ordering  
causality violation

Page No.: \_\_\_\_\_

### Vector Clock

An algorithm that generates partial ordering  
of events and detects causality violations.  
in DS. (causality)

Let  $n = \text{no. of processes in DS}$ .

each process  $P_i$  equipped with clock  $c_i$   
 $c_i$  can be assumed to be a vector function  
that assigns a vector  $c_i(a)$  to event  $a$ .  
 $c_i(a)$  TS of event  $a$  at  $P_i$ .

$c_i[i]$   $\rightarrow$   $i$ th entry of  $c_i$ , corresponds to  
Pi's own logical time.

$c_i[j]$   $\rightarrow$   $p_i$ 's best guess of logical time  
at  $P_j$

[IR1]

Clock  $c_i$  is incremental b/w 2 successive  
events in process  $P_i$

$$c_i[i] = c_i[i] + d$$

$$c_i[i] = c_i[i] + d$$

$d = \text{drift time} = 1$

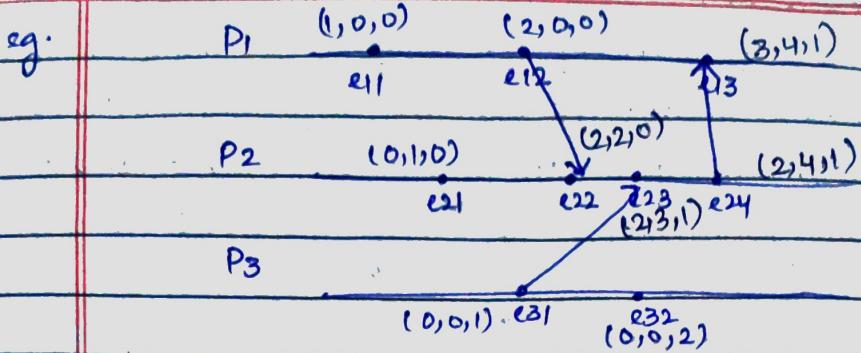
[IR2]

If  $a$  is sending msg  $m$  by process  $P_i$   
then msg.  $m$  is assigned a vector ts  
 $t_m = c_i(a)$

On receiving msg  $m$  by process  $P_j$

$$g_j[K] = \max(g_j[K], t_m[K])$$

$$\forall K \quad g_j[K] = \max(g_j[K], t_m[K])$$



$$e_{11} = (c_i + d) = (0+1) = 1 \Rightarrow (1, 0, 0) \checkmark$$

$$e_{12} = (c_i + d) = (1+1) = 2 \Rightarrow (2, 0, 0) \checkmark$$

$$e_{21} = (c_i + d) = (0+1) = 1 \Rightarrow (0, 1, 0) \checkmark$$

$$e_{31} = (c_i + d) = (0+1) = 1 \Rightarrow (0, 0, 1) \checkmark$$

$$e_{32} = c_i + d = 1+1=2 \Rightarrow (0, 0, 2) \checkmark$$

$$e_{22} = \max \begin{bmatrix} 0, 2, 0 \\ 2, 0, 0 \end{bmatrix} = (2, 2, 0) \checkmark$$

$$e_{23} = (c_i + d) = \max \begin{bmatrix} (2, 3, 0) \\ 0, 0, 1 \end{bmatrix} = \underline{(2, 3, 1)} \checkmark$$

$$e_{24} = (c_i + d) = \underline{(2, 4, 1)} \checkmark$$

$$e_{13} = \max \begin{bmatrix} 3, 0, 0 \\ 2, 4, 1 \end{bmatrix} = \underline{(3, 4, 1)} \checkmark$$

Advantages of vector clock over Lamport's

1) Events a & b are causally related if  $t^a < t^b$  or  $t^b < t^a$  else they are concurrent

2) One system of vector clocks

$$a \rightarrow b \text{ iff } t^a < t^b$$

Thus, causal ordering can be determined.  
(causal)

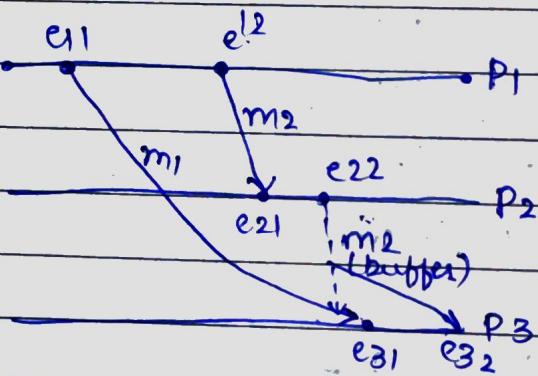
(causal)

Causal ordering of Msg maintaining  
 Deals with the concept of causal (ausal)  
 relationship among 'msg sent' event  
 corresponding to 'msg received' event.

If 1) send(m<sub>1</sub>)

2) Send(m<sub>2</sub>)

then every recipient of m<sub>1</sub> & m<sub>2</sub> msg.  
 must receive m<sub>1</sub> before m<sub>2</sub>.



Birman Sniper Stevenson Protocol:algo

- 1) Each process increase its vector clock by 1 upon sending the msg.
- 2) Msg is delivered to the process if the process has received all the prev. msg. use buffer the msg.
- 3) Update the vector clocks for process.

Reference: $P^o$  = process $e^{ij}$  =  $i \rightarrow$  process no,  $j \rightarrow$  event no. $c^i$  = v.c associated with  $P^o$  $c^i(j)$  =  $j$ th elementProtocol

- 1)  $P^o$  sending msg. to  $P^j$

$P^o$  increments  $c^o(i)$  and sets timestamp  
 $tm = c^o(i)$  for m.

- 2)  $P^j$  receiving msg. from  $P^o$

a) When  $P^j$ ,  $j^o = i$  receives m with ts  $tm$ , it delays msg. delivery until both are as follows.

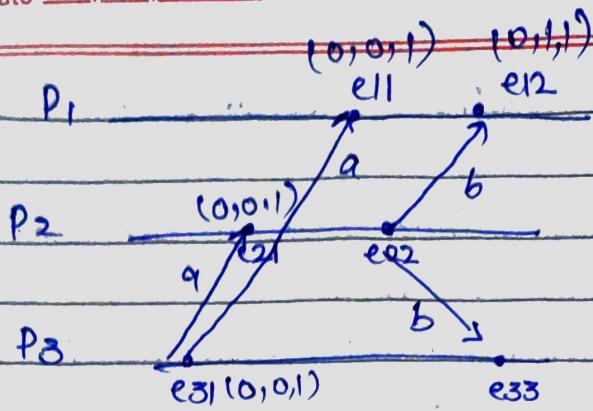
$$\textcircled{1} \quad c^j(k) = tm[i] - 1 \quad \forall k$$

$$\textcircled{2} \quad \forall k < n \quad \exists k = i \quad c^j(k) \leq tm[k]$$

- b) When msg. is delivered to  $P^j$ , update  $P^j$ 's v.c.

- c) Check buffer if anything can be delivered.

eg:

 $e_{31} \Rightarrow P_3$  sends msg a to

$$e_3 = (0, 0, 1)$$

$$t^a = (0, 0, 1)$$

 $e_{21} \Rightarrow P_2$  receives msg a

$$\text{as } c_2 = (0, 0, 0)$$

$$c_2[3] = t^a[3] - 1 = 1 - 1 = 0$$

$$c_2[1] = t^a[1] = 0$$

$$c_2[2] = t^a[2] = 0$$

msg is accepted &  $c_2 = (0, 0, 1)$  $e_{11}$  $P_1$  receives msg a as,

$$c_1 = (0, 0, 0)$$

$$c_1[3] = t^a[3] - 1 = 1 - 1 = 0$$

$$c_1[1] = t^a[1] = 0$$

$$c_1[2] = t^a[2] = 0$$

msg is accepted &  $c_1 = (0, 0, 1)$  $e_{12}$  $P_1$  sends msg b

$$c_2 = (0, 1, 1) \quad t^b = (0, 1, 1)$$

 $e_{22}$  $P_2$  receives msg b as  $c_2 = (0, 0, 1)$ 

$$c_2[2] = t^b[2] - 1 = 1 - 1 = 0$$

$$c_2[1] = t^b[1] = 0$$

$$c_2[3] = t^b[3] = 0$$

msg is accepted  $c_2 \neq (0, 1, 1)$

Q32 P<sub>3</sub> receives msg b as, C<sub>3</sub> = (0, 0, 1)

$$C_3[2] = t^b[2] - 1 = 1 - 1 = 0$$

$$C_3[1] = t^b[1] = 0$$

$$C_3[2] = b^b[2] = 0$$

So, msg is accepted  $\rightarrow$  C<sub>3</sub> is set to (0, 1, 1)

A(i)

### Transitless Global state

A Global state is transitless iff

$$\forall i \forall j : (1 \leq i, j \leq n) : \text{transit}(L_i, L_j) = \emptyset$$

Thus all communication channels are empty in transitless Global state.

A(ii)

### Strongly Consistent Global state:

(consistent + transitless)

(i) All channels are empty.

(ii) Send events of all recorded received events are recorded, receive events of all recorded send events are recorded.

(A)

Global state It is a set of local state of all local individual processes involved in

Need

i) Distributed Deadlock Detection.

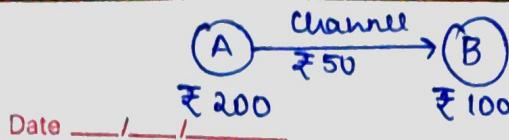
ii) Termination Detection

iii) Checkpoint

computation

& the state of comum. channels.

consistent  
Global  
state



Date \_\_\_\_\_

Page No.: \_\_\_\_\_

## Chandy Lampert's Global State Recording Algorithm

- 1) Model to capture a consistent Global state.
- 2) It uses a control msg called Marker whose role in a FIFO system is to separate msgs in channel.
- 3) After a site has recorded its local state it sends marker along all its outgoing channel before sending more msg.
- 4) Marker separates msg in the channel.

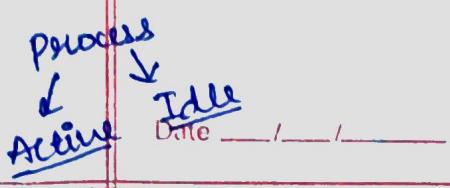
### Algorithm:

- 1) Marker sending rule process P :
  - a) P process records its state
  - b) for each outgoing channel C from P on which marker has not been already sent, P sends marker along C before sending more msg along C.
- 2) Marker receiving rule process Q :

On receiving marker along channel C if (Q hasn't recorded its state)  
↳ record its state  
↳ record the state of C as empty set  
↳ follow markers sending rule.

else

record the state of C as the set of msg received along C after Q's state was recorded & before Q received the marker along C.



problem comes  
when DC terminates  
implicitly Page No.: \_\_\_\_\_

## Termination Detection

Detecting whether an ongoing distributed computation has finished all its activities.

# example of use of Coherent view of DS.  
Notations:

- ① B(DW): computation msg sent as part of comp.  
DW is weight assigned to it.
- ② C(DW): control msg sent from process to CA  
DW is weight assigned to it.

### Huang's Termination Algo:

1) Rule to send B(DW):

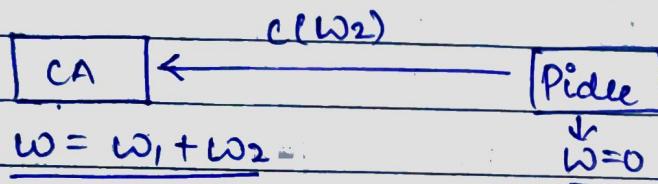
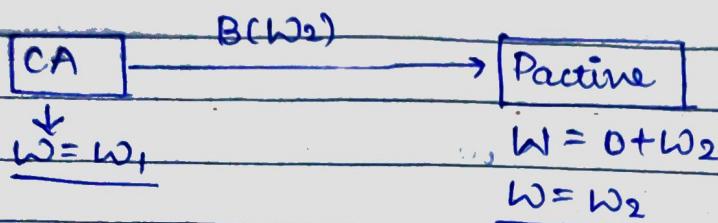
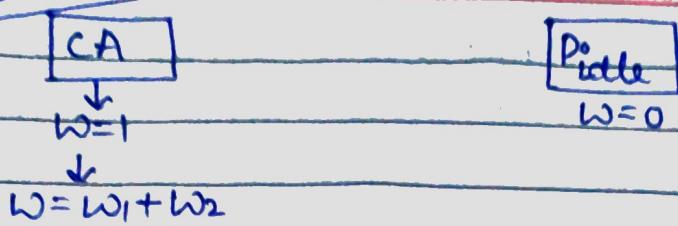
- ↳ CA having weight  $w$  may send a computation msg to process P.
- ↳ split the weight,  $w = w_1 + w_2$ ,  $w_1 > 0, w_2 > 0$
- ↳ set weight of P.  $w = w_1$
- ↳ send  $B(w_2)$  to process P.  $DW = w_2$

2) On receiving  $B(DW)$  by process Q:

- ↳ Add DW to weight of process Q,  $w = w + DW$
- ↳ If B was idle it becomes active on receiving  $B(DW)$ .

3) Rule to send C(DW):

- ↳ Any active process with weight  $w$  becomes idle by sending  $C(w)$  to CA.
- ↳ send control msg  $C(w)$ ,  $DW = w$ .
- ↳ set weight of process as 0.  $[w = 0]$

Basic Implementation

4) On Receiving  $c(DW)$  by CA:

↪ Add weights through control msg to

↪ Net of CA, i.e.  $w = w + DW$

If After adding, CA weight

$$w = 10 \Rightarrow$$

process  
computation  
terminated

*Johann Aurora*

## Mutual Exclusion

- ↳ Program obj that prevents simultaneous access to a shared resource.
- ↳ This concept is used in concurrent programming with a critical section. (A piece of code in which process access shared resource.)

## Requirement

- 1) freedom from Deadlock
  - 2) freedom from Starvation
  - 3) Fairness
  - 4) Fault Tolerance
- enters wait due to circular wait Relationship
- unbounded waiting due to order of service policy
- requests are not served in order they are made
- algo breaks if process die/mig are late/garbled

## Token based Algo

- 1) Token is shared among all the sites.
- 2) Token contains seq. of No. of sites in order to execute CS.
- 3) A site having token can only enter CS.

### ↳ Example:

- ↳ Lamport's
- ↳ Rickart Agarwala
- ↳ Maekawa's

Non-token

- 5) req. are executed exactly in order as they made it
- 6) scalable
- 7) Easy Authentication

## Non-Token Based Algo

- 1) Central site will communicate with all sites
- 2) Site uses TS value in order to request for CS.
- 3) A site with smaller TS will enter the CS.

### ↳ Example:

- ↳ Raymond Tree
- ↳ Suzuki-Tessman
- ↳ Singhal's Heuristic

- 5) No strict order of exec
- 6) less scalable

## Lamport's Non-Token Algorithm

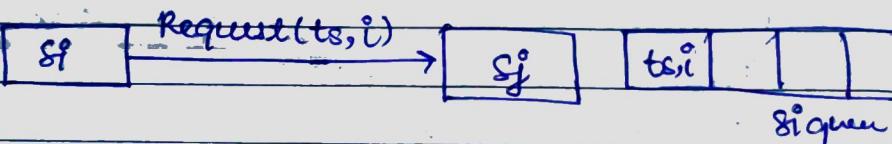
### Requesting set

$R^i = \{s_1, s_2, \dots, s_n\}$ , every site has queue.

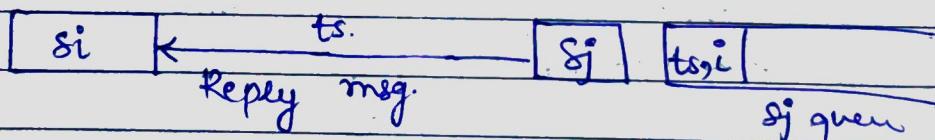
#### 1) Requesting the CS:

a) Si wants to enter the CS, it sends

Request( $t_s, i$ ) msg to all the sites in  $R^i$   
& places request in seq-queue<sup>i</sup>



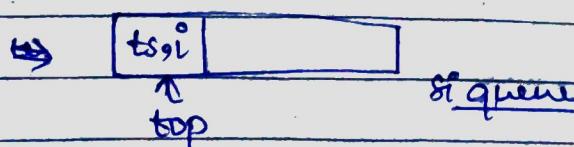
b) On receiving Request( $t_s, i$ ) msg by  $s_j$   
it returns a timestamp reply msg to  $s_i$   
& places  $s_i$ 's req. in seq-queue.



#### 2) Executing the CS:

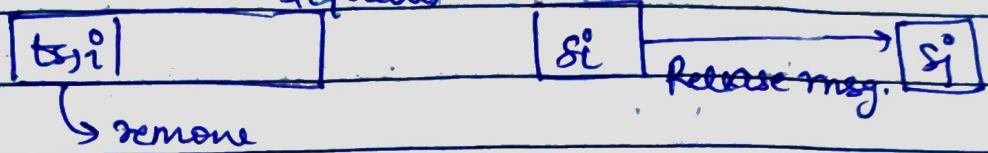
Si enters in CS if

- (i) Si has received msg with  $(t_{sj}) > (t_s, i)$  from all sites. AND
- (ii) Si's request is at top of seq-queue<sup>i</sup>.



### 3) Releasing the CS:

- a) On exiting si removes its request from top of RQ & send the ts Release msg. to all sites in fi



- b) On receiving the release msg. si removes its request msg. from its RQ.

Performance: Require  $3(N-1)$  msg / cs. insertion

1)  $(N-1)$  Request

$(N-1)$  Reply

$(N-1)$  Release

2) Synchronization Delay = T

3) can be optimized to  $2(N-1)$  msg / cs.

by suppressing Reply msg. in certain situation

Drawbacks:

1) Unreliable Approach:

failure of any process will halt the progress of entire system.

2) High Message Complexity:  $3(N-1)$  / cs.

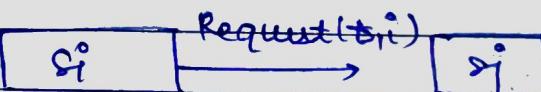
## Ricart-Agarwala Non-Token Algorithm

# Optimization of Lamport's Algo,  
it dispenses with release msg by  
merging them with Reply msg.

Request set:  $R_i = \{s_1, s_2, \dots, s_n\}$   
every site has a queue.

1) Requesting the CS:

a)  $s_i$  wants to enter the CS, it sends its request to all sites



b) When  $s_j$  receives request msg from  $s_i$ ,  
it sends a reply msg to  $s_i$ ,  
if  $s_j$  is neither requesting nor executing  
the CS

or if  $s_j$  is requesting but  $(ts, i) < (ts, j)$   
else request is deferred.

2) Executing the CS:

$s_i$  enters the CS after receiving reply msg  
from all sites in  $R_i$ .

3) Releasing the CS:

$s_i$  exits the CS, it sends reply msg to  
all the deferred requests.

### Performance :

- 1)  $2(N-1)$  msgs / CS execution  
 $(N-1)$  Request msgs.  
 $(N-1)$  Reply msgs.
- 2) synchronization Delay T.

### Drawbacks :

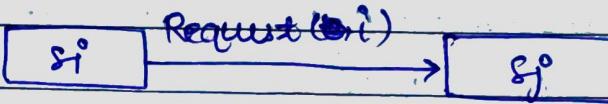
#### Unreliable approach

failure of any one node in system can  
halt the progress of system. In this  
situation process will stall.

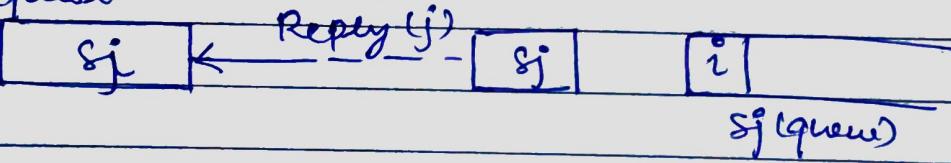
## Markovian Algorithms

Requesting set,  $R_i = \{s_1, s_2, \dots, s_n\}$

- 1) Requesting the CS:
- a)  $s_i$  wants to enter the CS, it sends request msg to all the sites in  $R_i$



- b) When  $s_j$  receives the request msg from  $s_i$  it will send a reply msg to  $s_i$  only if  $s_j$  hasn't sent a reply msg to a site from the time it received the last Release msg. else it queues up the request



- 2) Executing the CS:

$s_i$  enters the CS only after receiving reply msg from all the sites in  $R_i$ .

- 3) Releasing the CS:

- a)  $s_i$  exits the CS & sends Release(i) msg to all the sites in  $R_i$
- b) On receiving Release(i),  $s_j$  sends reply to next site waiting in queue & deletes that entry from queue.

Q) If queue is empty, sj update its status to show that it hasn't sent any reply msg since the receipt of last release msg.

### Performance:

- 1)  $3JN$  msgs per CS  
 $JN$  request msgs  
 $JN$  reply msgs  
 $JN$  release msgs.
- 2) Synchronization Delay =  $\alpha T$  (msg. propagation delay).

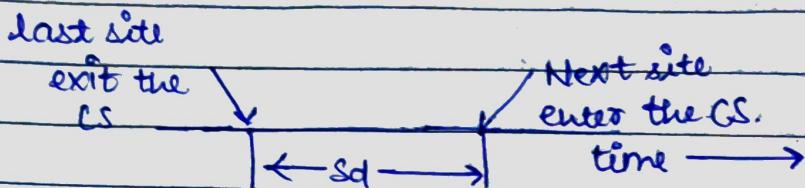
### Drawback:

Deadlock prone  $\rightarrow$  A site is exclusively locked by other sites & request are not prioritized by the timestamp.

## Performance Metric

1) No. of mega users per CS:

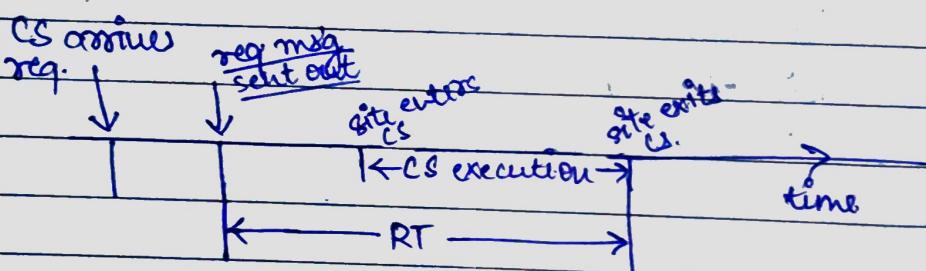
2) Synchronization Delay



Time between 2 consecutive CS:

3) Response Time (RT):

Time between request msg sent out and completion of Critical section.



4)

System throughput: Rate at which system executes for CS.

$S_d$  = synchronization delay

$F$  = Avg. CS execution time

$$\boxed{\text{System throughput} = \frac{1}{S_d + F}}$$

ll  $\rightarrow$  light load  
ll  $\rightarrow$  heavy load

Date \_\_\_\_\_

Page No.: \_\_\_\_\_

	<u>Algorithm</u>	RT(ll)	sd	Meg. (ll)	Msg (ll)
1)	Lamport	QT+E	T	$3(N-1)$	$3(N-1)$
2)	Ricart	2T+E	T	$2(N-1)$	$2(N-1)$
	Agarwala				
3)	Markkula	2T+E	2T	$3\sqrt{N}$	$5\sqrt{N}$
4)	SuzukiKasami	2T+E	T	N	N
5)	Singhal Hechtman	2T+E	T	$N/2$	N
6)	Raymond	$T \log N + E$ $T \log N/2$	$T \log N/2$	$\log N$	4
	$T_{ll}$				

### Remember

- # Chandy-Lamport's Global state (Marker)
- # Casual Ordering Birman Sevcik Stephenson
- # Termination Detection - Huang's Termination CA
- # Lamport's logical clock  $c_j = \max(c_j, t_m[j])$  Rule
- # Vector clock  $c_j[k] = \max(c_j[k], t_m[i[k]] + k)$

**Rohan**

## Goals of DS

- 1) Resource sharing
- 2) Improved system performance
- 3) Improve Reliability & Availability.
- 4) Modular ~~complexity~~ (expandability).

## Web challenges:

- 1) Naming
- 2) Access control.
- 2) Security
- 3) Availability
- 5) Performance management.

System Models: Common properties & design choices for DS in single discipline Model.

- (i) Architectural
- (ii) Fundamental.

## Causal Ordering Algo:

- 1) Birman - Schiper - Stephenson.
- 2) Schiper - Eggers - Sandor

## Threads over processes.

- ↳ Created & Destroyed much faster.
- ↳ faster to switch b/w threads.
- ↳ Threads share data easily.

QUESTION

### Proxy Server:

- ↳ A server that acts as a gateway between user's computer & Internet.
- ↳ It verifies and fwd incoming client req. to other servers for further communication.
- ↳ Enhance privacy.
- ↳ e.g. Req. made from Proxy server may help to hide client's IP from web server.

Reason Middleware moved from distributed objects to distributed components.

- ① Implicit Dependencies
- ② Interaction with Middleware.
- ③ Lack of separation of distributed concern.
- ④ No support for deployment.

### Consistent cut:

A cut is consistent if for each event that it contains, it also includes all events causally ordered before it.

Let a, b be 2 events in DS

$$\begin{array}{|c|} \hline (a \in \text{consistent cut } C) \wedge (B \leq a) \\ \hline \Rightarrow b \in C \\ \hline \end{array}$$

Inconsistent cut:

A cut  $c$  is inconsistent iff any of the following conditions are true.

1)  $e \rightarrow$  sending event of msg.

$e' \rightarrow$  receiving event of msg.

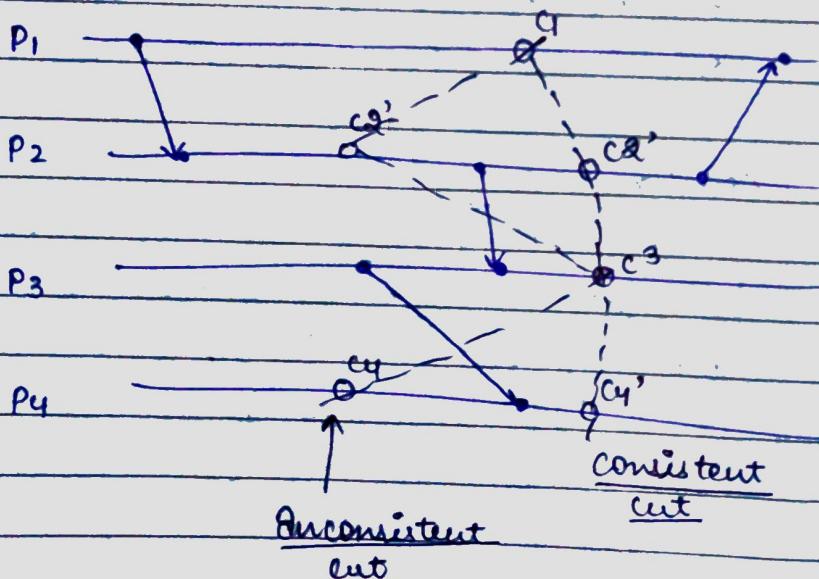
$$e' \in C \wedge e \notin C$$

2)  $e$  and  $e'$  are on events and

$$e \leftrightarrow e' \wedge e' \in C \wedge e \notin C$$

3)  $e$  and  $e'$  are off events and

$$e \leftrightarrow e' \wedge e' \in C \wedge e \notin C$$



## Architectural Model

various types of architectures exist that are usually used for distributed computing.

low level

Interconnection

of Multiple CPUs

high level

Interconnection of

processes running on those  
CPUs.

- a) Client Server Model. → Dependency
- b) 3-Tier Architecture. → No dependency
- c) Peer to Peer Model
- d) Proxy server / cache.      c) tightly coupled (clustered)

Fundamental Model: Deals with communication that can be affected by delays, failure & security attacks.

- a) Interaction Model

↳ performance of comm. channel,  
↳ sync. & async

- b) Failure Model:

↳ Masking, Integrity & Validity approach

↳ Specification of faults.

- c) Security Model:

↳ Identify possible threats to process domain

## Agreement Protocols

- ↳ process of sending and reaching the agreement to all sites
- ↳ AP are useful for error free communication among various sites.

System Model: (Where agreement protocols are used)

- (i) If there are n processes in DS. then Only m processes out of them maybe found as faulty processes.
- (ii) every process in the system is free to communicate with each other in the system due to their logical connections with each other.
- (iii) A receiver process always knows the identity of sender process of msg.
- (iv) The communication medium is reliable and only processes are prone to failures.

## Classification Of Agreement Problems

1)

Byzantine Agreement Problem:

↳

A single value is to be agreed upon.

↳

Agreed value is initialized by an arbitrary process and all non-faulty processes have to agree on that value.

2)

Consensus Problem:

↳

every process has its own initial value and all non-faulty processes must agree on a single common value.

3)

Intentional Consistency Problem:

↳

every process has its own initial value and all non-faulty processes must agree on a set of common values.

## Byzantine Agreement Problem

An arbitrarily chosen process (source process) broadcasts its value to others

A solution to BAP should meet the following objectives:

- (i) Agreement: All non-faulty processes agree on the same value.
- (ii) Validity: If source is nonfaulty, then the common agreed value must be the value supplied by the source process.  
If source is faulty, then all non-faulty processes can agree on any common value.
- (iii) Termination: each non-faulty process must eventually decide on a value.

## Consensus Problem

Every processor broadcasts its initial value to other processors  
Initial value may be different for different processors.

Protocol must meet following objectives:

- (i) Agreement: All non-faulty processors agree on the same single value
- (ii) Validity: At NFP  
If the initial value of every NFP is  $v$  then the agreed upon common value by all NFP must be  $v$ .
- (iii) Termination: Each NFP must eventually decide a value.

## Iterative Consistency Problem

Every processor broadcasts its initial value to all other processors.

Initial value must be different for different processor

Protocols must meet the following objectives:

- (i) Agreement: All NFP agree on the same vector ( $v_1, v_2, v_3, \dots, v_n$ )
- (ii) Validity: If  $i$ th processor is NF and its initial value is  $v_i$  then the  $i$ th value agreed by all NFP must be  $v_i$ .
- (iii) Termination: Each NFP must eventually decide on different value of vectors

Lamport's Shostak Pease Algorithm:Lamport's Shostak Pease Algorithm:

One Message Algo  $OM(m)$   $m > 0$  solves BAP  
 for  $3m+1$  processors in the presence of atmost  
 $m$  faulty processors

Let  $n = \text{total no. of processors}$  ( $n \geq 3m+1$ )

↳ 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 it receives no value, default value = 0.

Algorithm  $OM(m)$ :

- 1) The source processor sends its value to every processor.
- 2) For each  $i$ , let  $v_i^0$  = value processor  $i$  receives from the source.
- 3) Processor  $i$  acts as the new source and initializes algo  $OM(m-1)$  wherein, it sends the value  $v_i^0$  to each of other  $n-2$  processors.
- 4) For each  $i$  and  $j$  ( $\neq i$ ), let  $v_{ij}$  (value processor  $i$  received from processor  $j$ ) in step ③.  
 Processor  $i$  uses the value majority( $v_i^0, v_{i1}, v_{i2}, \dots, v_{i(n-1)}$ )

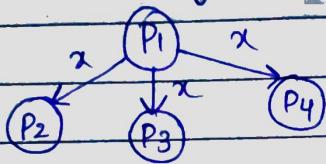
The majority function ( $v_1, v_2, \dots, v_{n-1}$ )  
 computes the majority value if it exists else  
 it uses default value '0'.

~~BA cannot always be reached among 4 processors if 2 processors are faulty.~~

4 processors :  $P_1, P_2, P_3, P_4$

$P_1$  initiates the initial value.

$P_2, P_4$  are faulty.



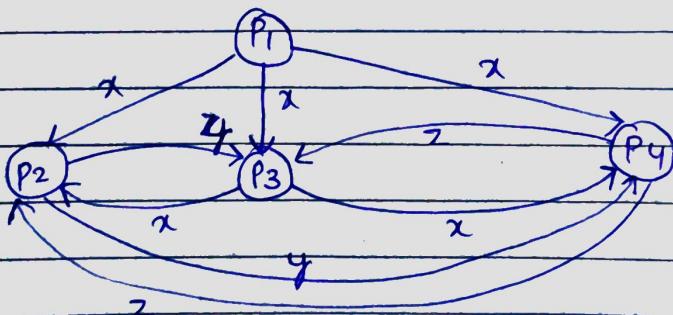
$P_1$  executes  $DM(1)$  and sends its value  $x$  to other processors.

After receiving the value  $x$  from  $P_1$ ,

$P_2, P_3, P_4$  execute  $DM(0)$ .

$\bullet$   $P_3$  is non-faulty & sends  $x$  to  $P_2$  &  $P_4$

$P_2$  &  $P_4$  are faulty sends  $y$  to  $P_3, P_4$  and  $z$  to  $P_2, P_3$  resp.



Majority values for Byzantine solution:  
processors      Required Majority      Common

Processors	Required Majority	Common
$P_2$	$(x, x, z)$	$x$
$P_3$	$(x, y, z)$	$0$
$P_4$	$(x, x, y)$	$x$

According to Majority value table, processor doesn't agree on single common majority value, which violates the condition of BAP.

### Applications of Agreement Protocols:

- 1) Fault-Tolerance Clock synchronization:
  - ↳ DS require physical clocks to synch.
  - ↳ Physical clocks have drift problems.
  - ↳ AP may help to reach a common clock value.
- 2) Atomic Commit in DDBS:
  - ↳ DDBS sites must agree whether to commit or abort the transactions.
  - ↳ AP may help to reach a consensus.

### Application of Agreement Protocol:

- 1) Fault-Tolerance clock synchronization
- 2) Atomic Commit in DDBMS.

# Resource Management Component

Date   /  /  

Page No.:   

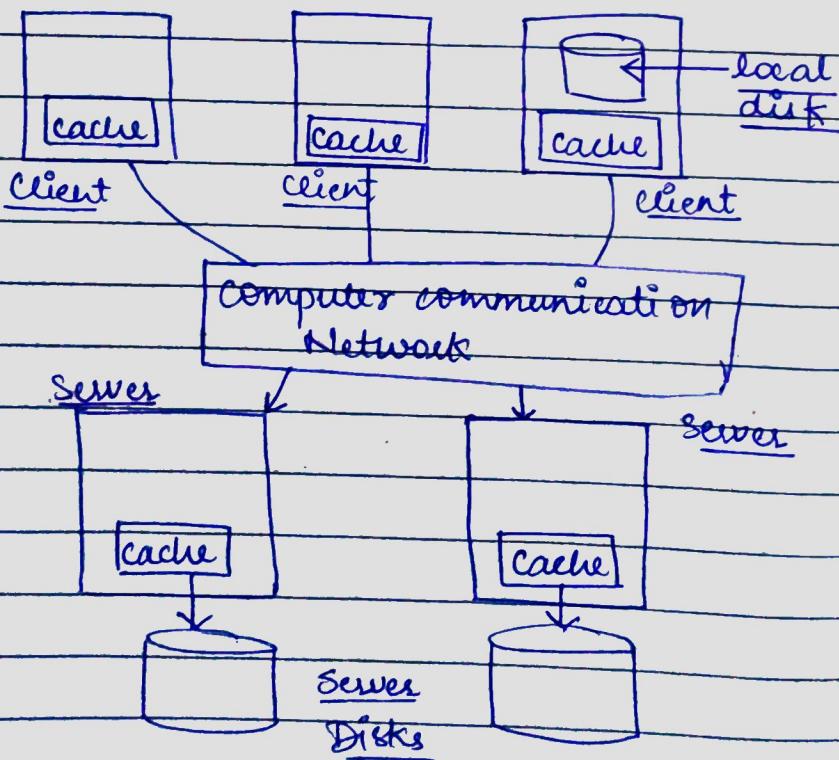
## Distributed file system (DFS)

- ↳ DFS is a resource management component of Distributed Operating System.
- ↳ It implements a common file system that can be shared by all autonomous computer in system.

Goals → ① Network Transparency

→ ② High Availability

### Architecture:



(1) Name Server:

A process that maps names specified by the clients to stored objects such as files & directories.

(2)

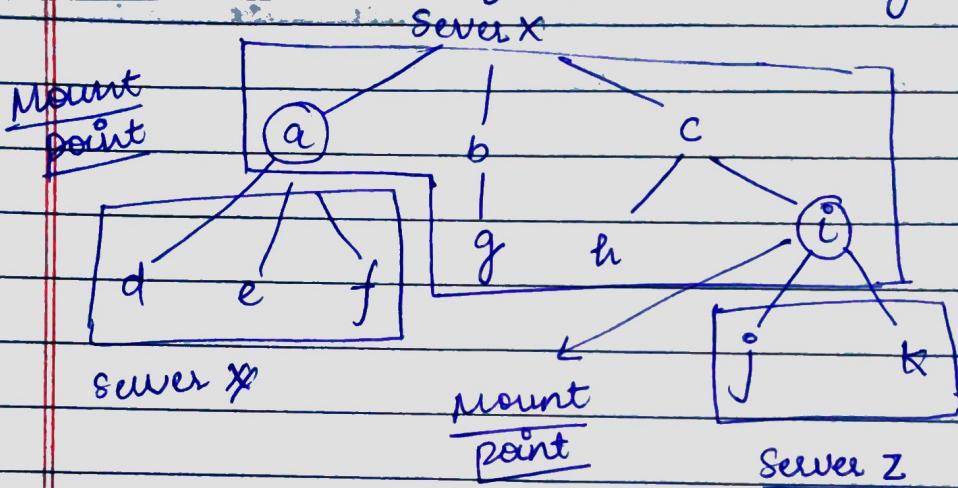
Cache Manager:

Implements file caching  
The mapping occurs when process references a file/directory for the first time.

## Mechanisms for Building DFS:

①

Mounting: It means assigning the root directory of the new file system to a subdirectory of root directory.



2)

Caching:

- ↳ Improves file system performance by exploiting the locality of reference.
- ↳ When client references a remote file, the file is cached in Main Memory of server and at client side.
- ↳ When multiple clients modify shared data, cache consistency becomes problem.
- ↳ It is difficult to implement a solution that guarantees consistency.

3) Hints:~~Timestamp~~  
~~+ cache~~~~Caching without cache~~  
~~coherence~~

- ↳ Treat cached data as hints:  
ie cached data may not be completely accurate.
- ↳ can be used by application that can discover that the cached data is invalid and can recover

4) Bulk Data Transfer:

- ↳ Overhead introduced by protocols doesn't depend on amt of data transferred in one transaction.
- ↳ Most files are accessed in their entirety
- ↳ When client req one block of data, multiple consecutive blocks are transferred

5) Encryption:

- ↳ Needed to provide security in DS.
- ↳ Entities that need to communicate send req to authentication server.
- ↳ Authentication server provide key for conversation.

## Design Issues

1)

### Naming and Name Resolution.

A name in file system is associated with an object. Name Resolution is the process of mapping a name to an object or in case of replication to multiple objects.

2)

### Caches on Disk or Main Memory

Whether the data required by a client should be in main memory at the client or on a local disk at the client.

3)

### Writing Policy:

It decides when a modified cache block at client side should be transferred to server.

4)

### Cache consistency:

5)

### Availability:

- a) Replication
- b) Unit of Replication
- c) Replica Management.

6)

Scalability: It deals with the suitability of design of system to cater the demands of growing system.

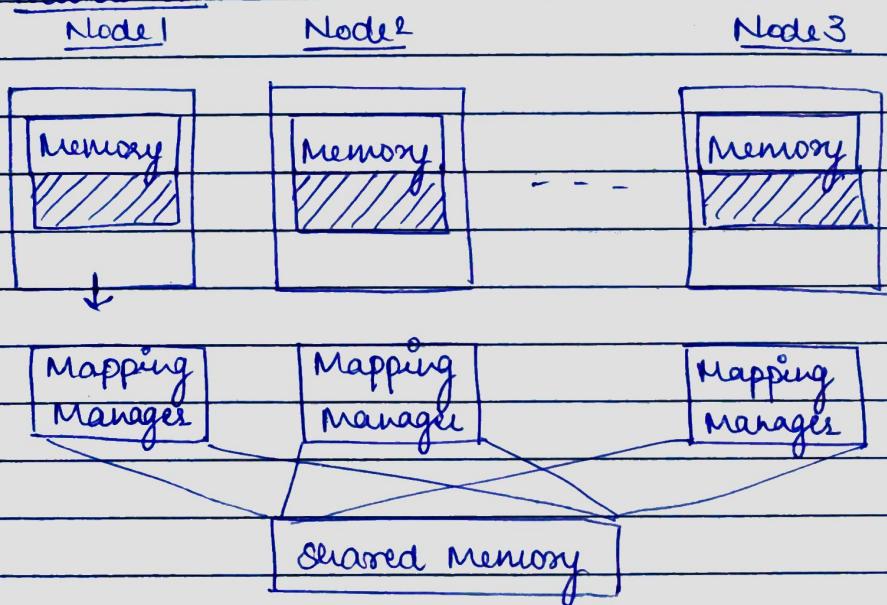
7)

Semantics: It characterizes effects of access on files.

## Distributed Shared Memory

- DSM implements shared memory model in DS, which has no physical shared memory
- Provides virtual address space shared b/w all nodes

Architecture:



Design Issues:

- Granularity:
  - size of shared memory unit.
  - large pg size, more locality, less overhead commun'
  - more contention
  - more false sharing
- Page Replacement:
  - Replacement algo must take into acc. pg. access modes
  - ↳ shared, private, readonly, writable
  - ↳ private pg. replaced before shared ones, private pg. swap to diff. Shared pg sent over network to owner, 'Read only may be discarded'

## Memory Coherence

① Replicated shared data Obj.

② concurrent access to data obj at many nodes

① Sequential consistency:

→ result of exec. of opr. of all prcs.  
is same as if they are  
executed seq.

② General consistency: → All copies of Mem loc. contain same data

③ Process consistency:

→ Opr. issued by processor are performed in order they are issued.

④ Weak consistency:

Memory is consistent only after sync opr.

⑤ Release consistency:

Sync. opr. must be consistent.

## Coherence Protocols:

① Write-invalidate

Protocol.

A write to shared data causes the invalidation of all copies except one before the write.

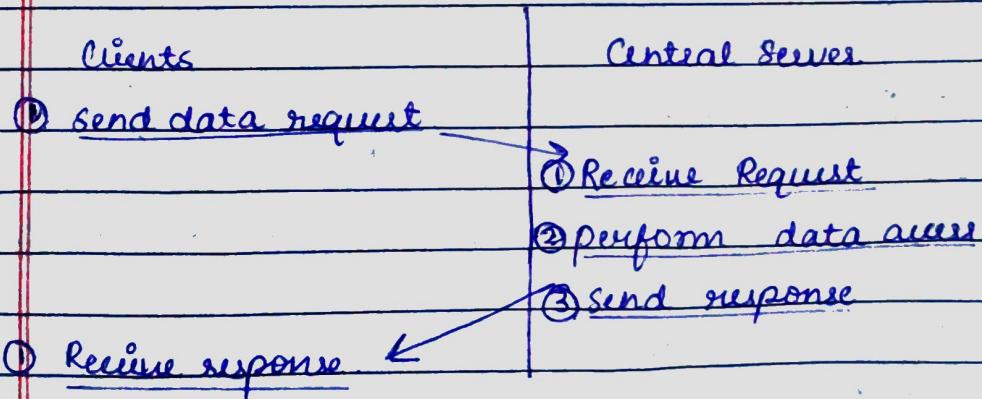
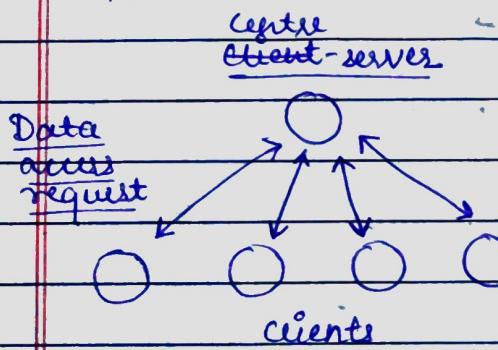
② Write update protocol.

A write to a shared data causes all the copies of that data to be updated.

## Algo to implement DSM

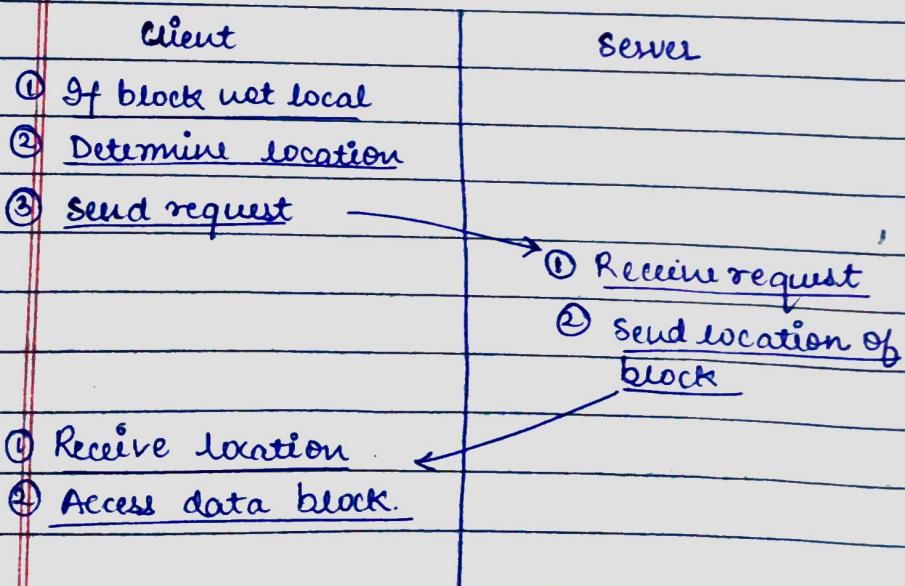
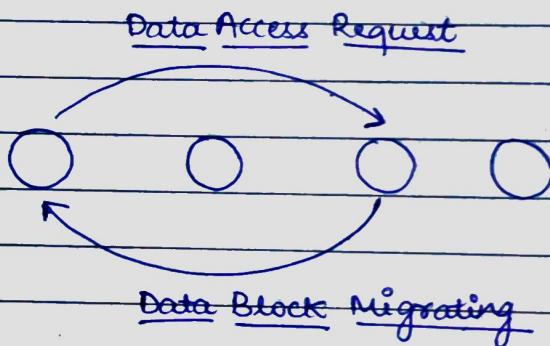
### Central Server Algorithms:

- ↳ A central server maintains all the shared data
- ↳ Read Request : return data item
- ↳ Write Request : update data & return ack. msg.
- ↳ A timeout is used to resend a request if ack fails
- ↳ Duplicate write req. can be detected by associated seq. no. with write requests.
- ↳ If an application's req. to access shared data fails repeatedly, a failure condition is sent to the application.



## Migration Algorithm

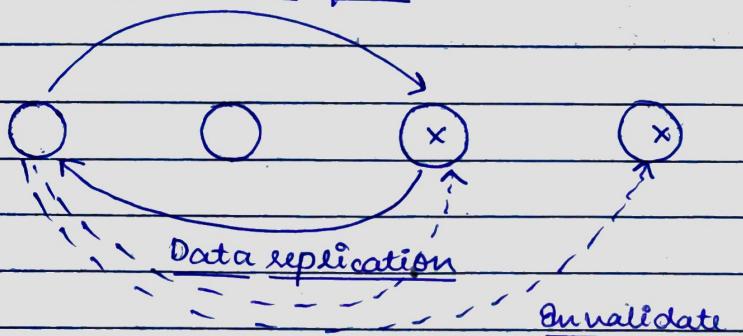
- ↳ In migration algos, data is migrated to the loc. of data access req., allowing subsequent access to data to be performed locally.
- ↳ Allows only one node to access shared data at a time.
- ↳ Keeping track of memory location (location service, home machine for each page, broadcast)



## Read Replication Algorithm

- 1) Allow Multiple nodes to have read access and one node to have write access
- 2) A write to a copy causes all copies of the data to be updated or invalidated
- 3) All locations must be kept track of: location service/home Machine

### Data access request

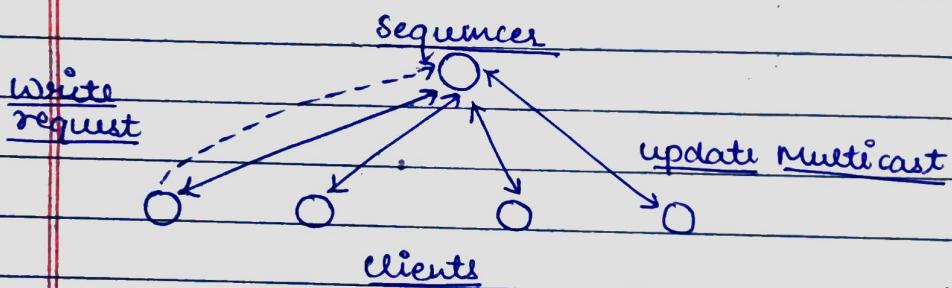


Clients	Remote Host
<ol style="list-style-type: none"> <li>① If block not local</li> <li>② determine location</li> <li>③ send request.</li> </ol>	<ol style="list-style-type: none"> <li>① Receive Request</li> <li>② send block</li> </ol>
<ol style="list-style-type: none"> <li>① Receive block</li> <li>② Multicast Invalidate</li> </ol>	<ol style="list-style-type: none"> <li>① Receive Invalidate</li> <li>② Invalidate Block</li> </ol>
① Access data	

## Full Replication Algorithm

- 1) Allow multiple read/write concurrently
- 2) Must control access to Shared Memory
- 3) Need to maintain consistency

use gap-free sequences, all nodes wishing to modify shared data will send the modification to sequence which Multicast the sequences.



Clients	Sequencer	Hosts
① <u>if write send data</u>	① <u>receive data add</u> ② <u>Sequence No.</u> ③ <u>multicast</u>	
① <u>receive ack update</u> <u>local memory</u>		① <u>receive data update</u> <u>local memory</u>

## Backward Error Recovery

- operation Based
- state-Based

Page No. \_\_\_\_\_

## Failure Recovery

Process that involves restoring an erroneous state to an error-free state

### Forward Error Recovery:

- ↳ Remove errors in process/system state if errors can be completely accessed
- ↳ Continue process/system forward execution

### Backward Recovery

#### Advantages:

- 1) simple to implement
  - 2) can be used as general fault does not occur recovery mechanism again
- Disadvantages:
- 1) No guarantee that
  - 2) Some components can't be recovered.
  - 3) perform penalty

## Backward Error Recovery

Restore process/system to previous error-free state and restart from there.

### Forward Recovery

#### Advantages:

- 1) less overhead

#### Disadvantages:

- 1) limited to use
- 2) cannot be used as general mechanism for error recovery

~~Data  
Application /  
of Agreement protocol~~

indivisible  
instantaneous  
uninterrupted  
No.

} Atomic  
Transactions

## Atomic Commit

Phase 1: Every site executes their part of distributed transactions & broadcasts their decisions (commit / Abort) to all other sites.

Phase 2: Each site decides whether to commit or abort its part of distributed transactions based on what it has received from other sites in first phase.

### Advantages of DSM:

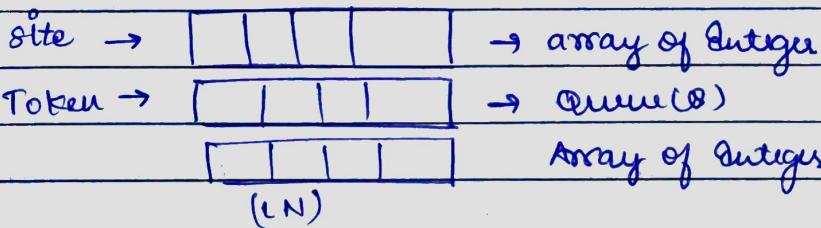
- 1) Cheaper than tightly coupled processor system.
- 2) Can form a large physical memory.
- 3) Easier programming.
- 4) DSM takes advantages of Memory reference locality.

Token Based Algo:

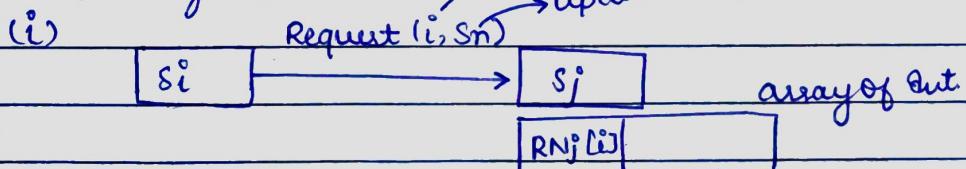
In token based algorithm, a unique token is shared among all sites. If the sites possess token then it is allowed to access the CS.

Suzuki-Kasami Algorithm: (Broadcast)Design Issues:

- Distinguish outdated Request Msg from Current Request Msg
- Determine which site has outstanding req. (RN)



- 1) Requesting the CS:  $\xrightarrow{\text{items}}$  updated seq.no.



- (ii)  $RN_j[i] \leftarrow \max(RN_j[i], s_n)$ .

$s_j$  send idle token if  $RN_j[i] = LN[i] + 1$

- 2) Executing the CS:

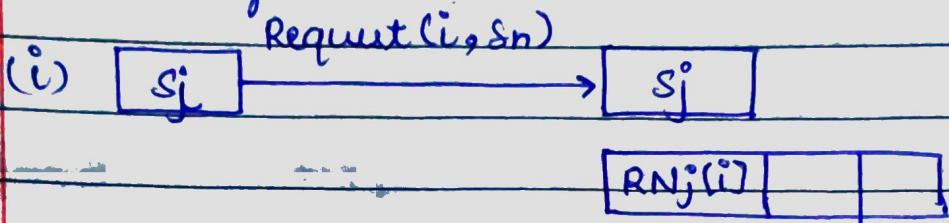
Site  $s_i$  execute the CS when it received token.

- 3) Releasing the CS: (i)  $LN[i] = RN_i[i]$

(ii) updated ID of site not in token queue if  $RN_j[i] = LN[j] + 1$

(iii) Delete top site ID from queue

## 1) Requesting the CS:



(ii)  $RN_j^i[i] \leftarrow \max(RN_j^i[i], s_n)$

$s_j$  send idle token if

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

## 2) Executing the CS:

Site  $s_i$  executes the CS when it received token

## 3) Releasing the CS:

(i) set  $LN[i] = RN_i^i[i]$

(ii) Updated ID of site not in token only  
if  $RN_j^i[j] = LN[j] + 1$

(iii) Pop queue.

## Singhal's Heuristic Algorithm

$$R = \{S_1, \dots, S_n\}$$

Data structures: state

site  $\rightarrow$   $SV[1-N]$ ,  $SN[1-N]$

Token  $\rightarrow$   $TSV[1-N]$ ,  $TSN[1-N]$

A site can be in following states:

R  $\rightarrow$  req. the cl.

E  $\rightarrow$  execute the cl.

H  $\rightarrow$  Hold the idle token.

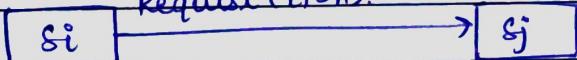
N  $\rightarrow$  None

### 1) Requesting the CS:

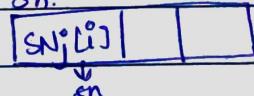
(i) set  $SV[i] = R$

(ii) increment  $SN[i] = SN[i] + 1$

Request( $i, SN[i]$ ).

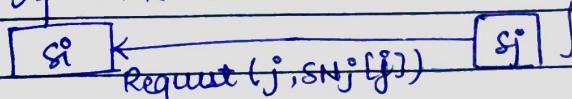


(iii) Discard if  $SN[j] \geq n$ .



state  
of site  $\{ j \rightarrow N, \text{ set } i \rightarrow R \}$

$\{ j \rightarrow R, \text{ if } i \neq R \text{ set } i = R \}$



$\{ j \rightarrow E, \text{ set } i \rightarrow R \}$

$\{ j \rightarrow H, \text{ set } i \rightarrow R \}$

Token  $\Rightarrow$  state  $\rightarrow R$

sequence NO  $\rightarrow SN$ .

2)

Executing the cl:Si execute cl and set  $SV_i[i] = E$ 

3)

Releasing the cl:Si  $\rightarrow SV_i[i] = N, TSV_i[i] = N$ for all  $s_j$ If  $SN_i[j] > TSN_i[j]$ 

then

$$TSV_i[j] = SV_i[j]$$

$$TSN_i[j] = SN_i[j]$$

else

$$SV_i[j] = TSV_i[j]$$

$$SN_i[j] = TSN_i[j]$$

If  $SV_i[j] = N$  set  $SV_i[i] = N$ 

Tokens	RT	Sd.	Message (u)	Message (u)
Suzuki-Karaii	$2T+E$	$\frac{T}{2}$	$N$	$N$
Singhal-Heuristic	$2T+E$	$\frac{T}{2}$	$N/2$	$N$
Raymond	$T(\log N)+E$	$T\log N/2$	$\log(N)$	$4$

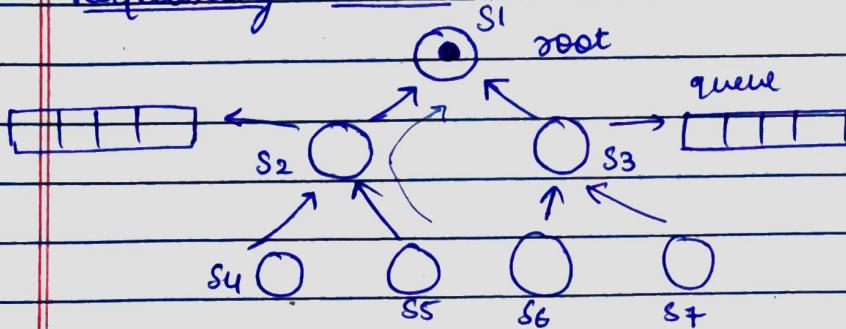
~~logical  
variable  
node~~

## Raymond Tree Algorithm

### Basic Idea:

- 1) The sites are logically arranged as directed tree.
- 2) Every site has local variable holder that points to an immediate neighbour node.

### Requesting the C:



### Executing the C:

A site enters C when it receives the token & deletes the top entry of its requesting queue.

### Releasing the C:

If request queue is nonempty it deletes top entry & send token to that site and set holder.

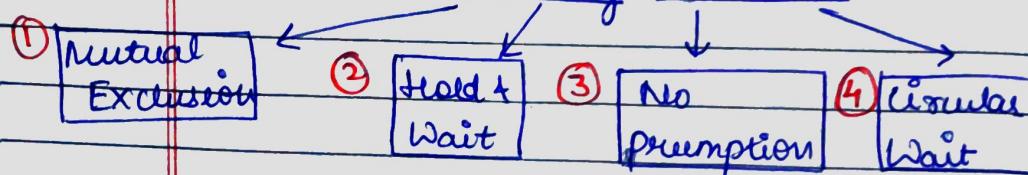
### Properties

- 1) free from deadlock, starvation
- 2) Average case  $O(\log N)$  msg complexity  
 $O(C \log N/2)$  sd.
- 3) Worst case  $\rightarrow$  Balanced Binary Tree.

## Deadlock

Defined as the permanent blocking of the process ie a set of processes is waiting for an event need by some other process.

### Necessary conditions



### Communication Deadlock

① Waiting is done for messages

② A process cannot proceed with the execution until it can communicate with one of the processes for which it was waiting.

③ A process can know the identity of those processes on the actions of which it depends.

④ Diagram can't be presented

### Resource Deadlock

① Waiting is done for resources

② A process cannot proceed with the execution until it receives all the resources for which it was waiting.

③ The dependence of one transaction on actions of other transaction is not directly known.

④ Can be presented by safestate by safestate

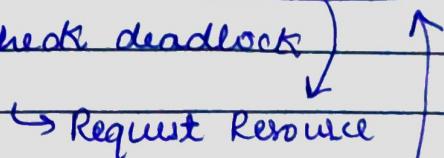
## Centralized Deadlock Detection

(i)

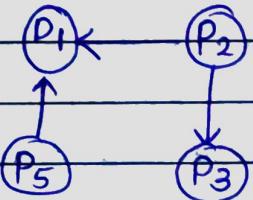
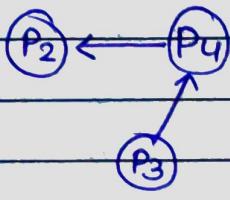
### Complete centralized Algo

$$R = \{S_1, S_2, \dots, S_n\}$$

In this a designated site called control site maintains a WFG to check deadlock.

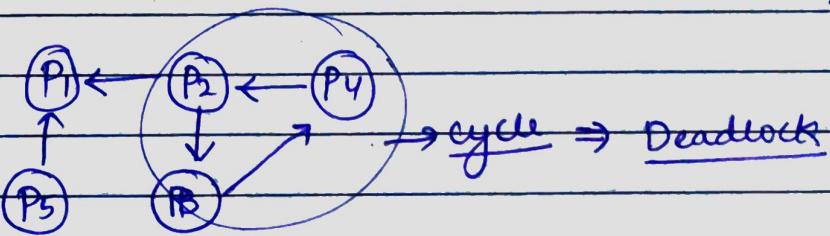


eg.

 $S_1$  $S_2$ 

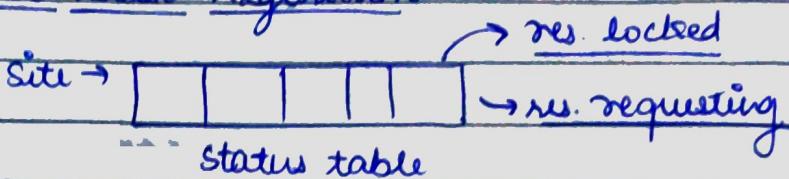
local WFG

Control  
site  
WFG

 $S_3$ 

### Disadvantages

↳ Dependent only on control site.

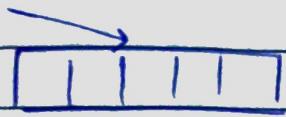
(ii) Ho-Ramamoorthy AlgorithmA) Two PhaseB) One PhaseA) Two Phase Algorithm

- 1) every site maintains a status table that contains the status of all processes initiated at that site
- 2) status includes → all resources locked → all resources waited upon.
- 3) A designated site req. the status table from all sites & constructs a WFG from the info. received & searches it for cycle.
- 4) If no cycle, system is free from deadlock else, designated site again requests status table from all the sites & again constructs a WFG using only those transactions which are common to both reports.
- 5) If same cycle is detected, it is deadlocked

Limitations It may indeed report false deadlock.

### B) One Please Algorithm:

- ↳ Only one status report from each site
- ↳ each site maintains 2 status tables.



Resource status Table

Keeps track of the transactions that have locked by or waited for by all the transactions at that site.

Process status Table

Keeps track of the resources locked by or waited for by all the transactions at that site.

- ↳ A designated site requests both the tables from every site, constructs LFG using only those transactions for which entry in the resource table matches with process table & searches for the cycle.
- ↳ If no cycle is detected so no deadlock.
- # It does not detect false deadlock b'coz it eliminates the inconsistency in state info by using only the info. That is common to both tables.
- # Faster & few nrgs, more storage  $\rightarrow$  2 tables.

## Distributed Deadlock Detection

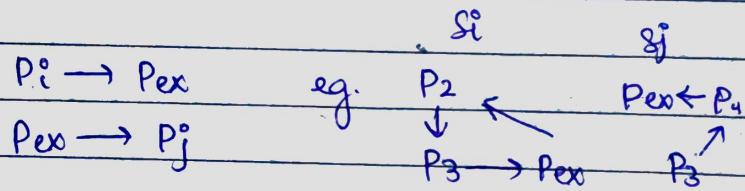
All the sites cooperate to detect a cycle in the state graph that likely to be distributed.

$$R = \{S_1, S_2, \dots, S_n\}$$

Obermark's

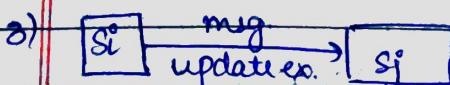
### Path-Pulling Algorithm

- 1) In this algo, info about the waitfor dependencies is propagated in the form of paths.
- 2) In this, each site maintains its local WFG that includes Pex.



### Algorithm:

- 1) If local WFG contains cycle that doesn't involve Pex then system is in deadlock.
- 2) If Pex exists suppose  $\text{Pex} \rightarrow P_1 \rightarrow P_2 \rightarrow P_3 \rightarrow \text{Pex}$



- 3)  $S_j$  forms WFG with new info to detect.

## Problem

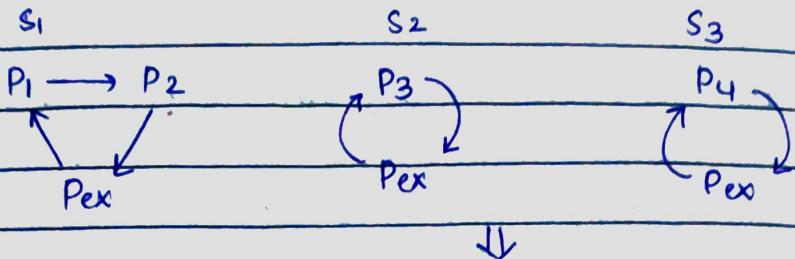
↳ false deadlock

Date \_\_\_\_\_

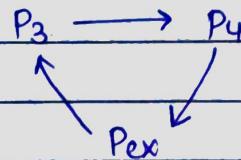
(due to anyone snapshots)

Page No.: \_\_\_\_\_

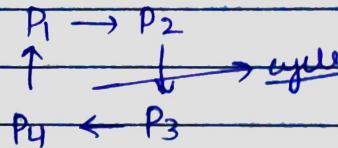
## example



$$S_4 = S_{2,3}$$



$$\Downarrow S_4, S_1$$



## Performance:

- 1)  $O(n(n+1)/2)$  messages
- 2)  $O(n)$  message size
- 3)  $O(n)$  delay to detect deadlock  
 → completely centralized  
centralized → Ho-Ramamoorthy

Distributed → Edge chasing

Path Pushing

Hierarchical → Manas - Munz

Ho-Ramamoorthy

## Edge chasing Algorithm

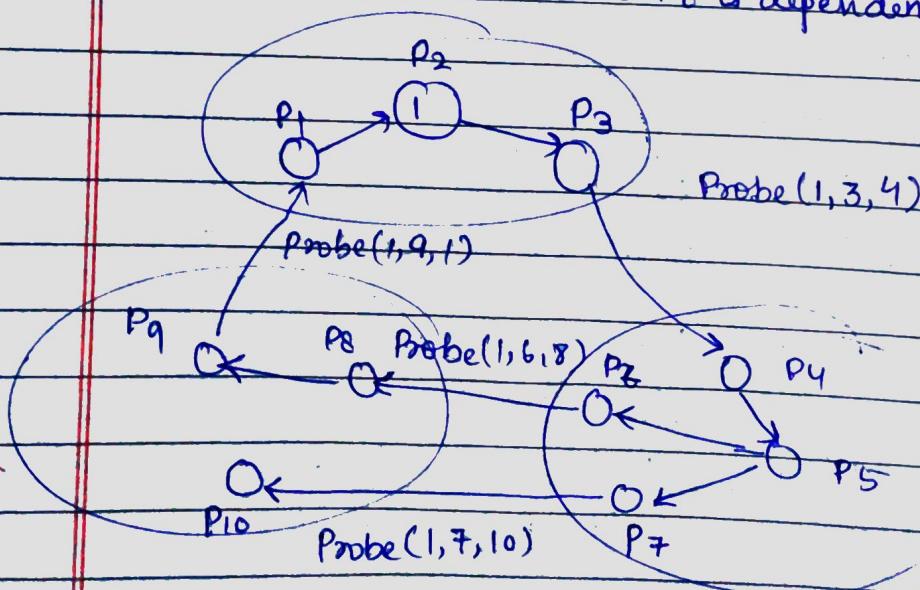
1) In this, we use a special msg. called probe.  
A probe is a triplet  $(i, j, k)$  denoting A' that it belongs to a deadlock detection initiated for process  $P_i, P_j, P_k$ .

a) Here we construct ~~GTEFG~~, GTWFG.  
(Global Transition waitfor graph)

Data structure: Process  $(P_i, P_j, P_k)$

$P_i$  is dependent on  $P_k$  if it exists a seq. of processes  $\{P_i, P_{j_1}, P_{j_2}, P_{j_3}, P_k\}$

the system maintains a boolean array like dependent j {if dependent;  $i\}$   
ie  $P_i$  is dependent on  $P_j$



Edge chasing Algorithm → phantom Deadlock Removal  
 3 steps

- 1) Initialization: Server initiates to detect deadlock
- 2) Detection: Detection consists of receiving probes and deciding whether deadlock has occurred & forward probes.
- 3) Resolution: When a cycle is detected, a transaction in cycle is aborted to break deadlock.

### Performance

- 1)  $(m-1)/2$  exchange of msg:  
 $m \rightarrow$  no. of processes  
 $n \rightarrow$  no. of sites
- 2) Delay  $\rightarrow$  O(n)  
to detect  
deadlock
- 3) 3-word message length

Centralized control

Distributed control

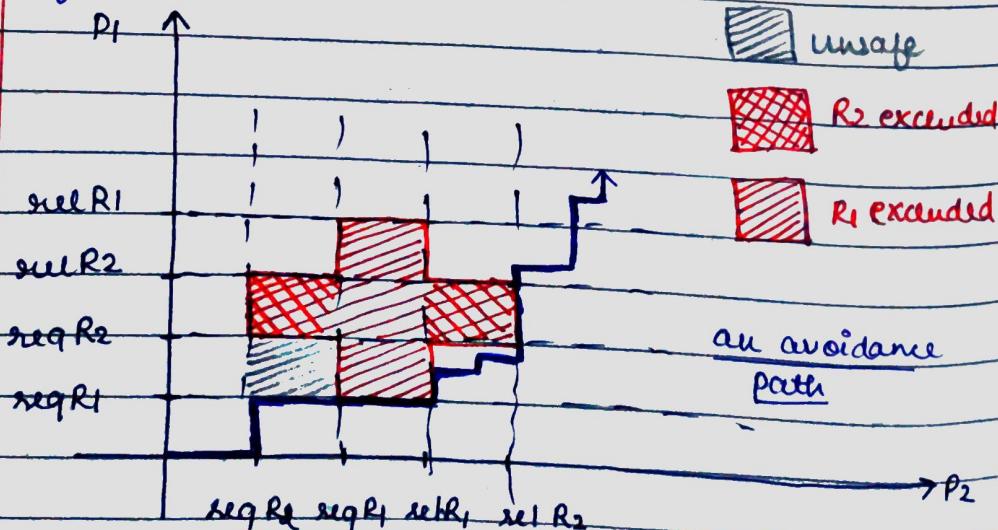
Hierarchical control

- 1) A control site has the resp. to detect Global WFG.
- All sites have the resp. to detect a GWFG.
- Descendent site can detect a GWFG.

- 2) Have single point of failure
- No single point of failure
- No single point of failure.

- 3) Easy to implement
- Difficult to implement
- Simple to implement

- 4) Completely Centralized Algo.  
Ho-Romanovsky Algo.
- Path flushing Edge-chasing Algo.
- Menasci Muntz Ho-Rama-moorthy Algo.



## Deadlock Handling strategies

A)

Prevention:1) Never satisfyMutual exclusionHold & waitNo preemptioncircular wait

- 2) Methods
- 1) Collective Requests
  - 2) Ordered Requests
  - 3) Preemption.

B)

Avoidance:Banks's Algorithm

A resource is assigned to a process if the state of global system is safe.

C)

Detection:

maintaining WFG and searching the cycles.

### Deadlock prevention

(i)

Ordered Allocation:

P<sub>1</sub> req(R1), req(R2), release(R2), release(R1)

P<sub>2</sub> req(R1), req(R2), release(R2), release(R1)

(ii)

Release before Request

P<sub>1</sub> req(R1), release(R1), req(R2), release(R2)

P<sub>2</sub> req(R2), release(R2), req(R1), release(R1)

(iii)

Request all at once:

P<sub>1</sub> request(R1, R2), release(R1, R2)

P<sub>2</sub> request(R1, R2), release(R1, R2)

## Edge chasing Algorithm

- 1) In this, we use a special msg called probe.
- 2) Probe is a triplet  $(i, j, k)$ , sent for detection initiated by  $P_i$  by site  $P_j$  to site  $P_k$ .
- 3) Deadlock is detected when probe returns to its initiator.
- 4) Here, we construct a global transition wait for graph. [GTWFG]

### Terminology & DS:

- (i)  $P_i$  is dependent on  $P_k$  if there exist a sequence of processes  $(P_i, P_{j_1}, P_{j_2}, P_{j_3}, P_k)$
- (ii) Boolean Array  
 $\text{dependent}_j(i) = \begin{cases} \text{True} & P_i \text{ is dependent on } P_j \\ \text{False} & \text{otherwise} \end{cases}$

- A) Algorithm Initiation by  $P_i$
- 1) If  $P_i$  is locally dependent on itself, then declare deadlock.
  - 2) else send probe  $(i, j, k)$  to homesite of  $P_k$  if
    - (i)  $P_i$  is locally dependent on  $P_j$
    - (ii)  $P_j$  is waiting for  $P_k$
    - (iii)  $P_j$  &  $P_k$  are on diff sites

### 8) Algorithm on receipt of probe (i,j,k)

1) Check the following

(i)  $P_k$  is blocked

(ii)  $\text{dependent}_K(i) = \text{false}$ .

(iii)  $P_K$  hasn't replied to all req. of  $P_j$

2) If all are true.

(i) set  $\text{dependent}_K(i) = \text{true}$

(ii) If  $K=i \Rightarrow P_i$  is deadlocked

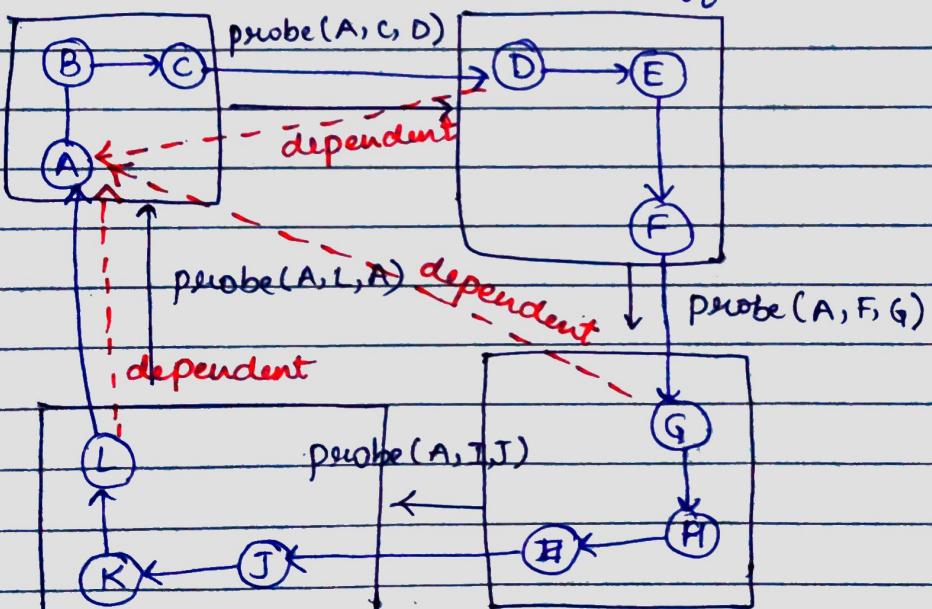
(iii) else send  $\text{probe}(i,m,n)$  to home site of

$P_n$  for every  $m \neq n$  following condition hold.

↳  $P_K$  is locally dependent on  $P_m$

↳  $P_m$  is waiting on  $P_n$

↳  $P_m$  &  $P_n$  are on diff. sites



Performance:

$n = \text{sites}$

1)  $m(n-1)/2$  messages  $\rightarrow m = \text{processors}$

2) 3-word msg length

3)  $O(n)$  delay

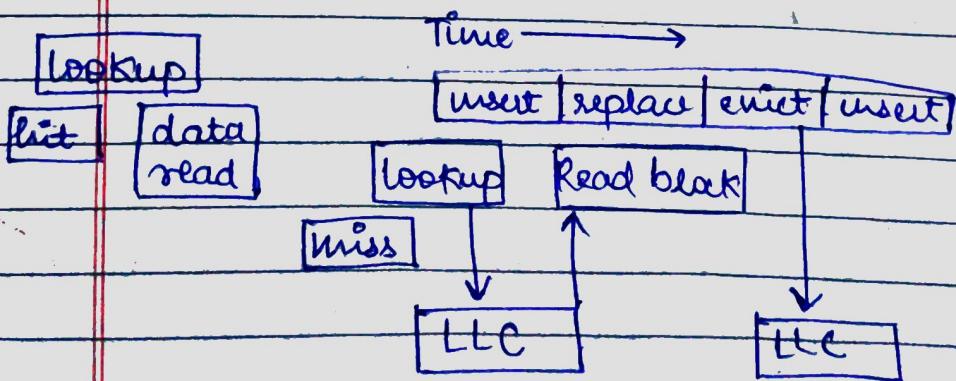
"In a gentle way, you can shake the world." -Mahatma Gandhi

## Cache

Data storage technique that provide the ability to access data or files at higher speed.

### Read Operation

- 1) Starts with lookup operation.
- 2) If cache hit, then cache return value to processor or higher level cache.
- 3) If cache miss, then cancel read op & send request to lower level cache.
- 4) Lower level cache will perform the same seq. of access & return entire cache block.
- 5) Cache controller then extract the requested data from block & send it to processor.
- 6) simultaneously, cache controller initiates the insert operation to insert the block into the cache.



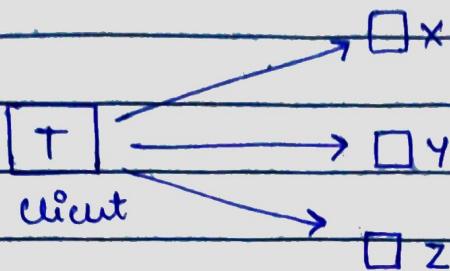
# Distributed Transactions

Date \_\_\_\_\_ / \_\_\_\_\_ / \_\_\_\_\_

Page No.: \_\_\_\_\_

## Flat Transactions

- ↳ A flat transaction has a single initiating point and a single end point.
- ↳ for short activities.

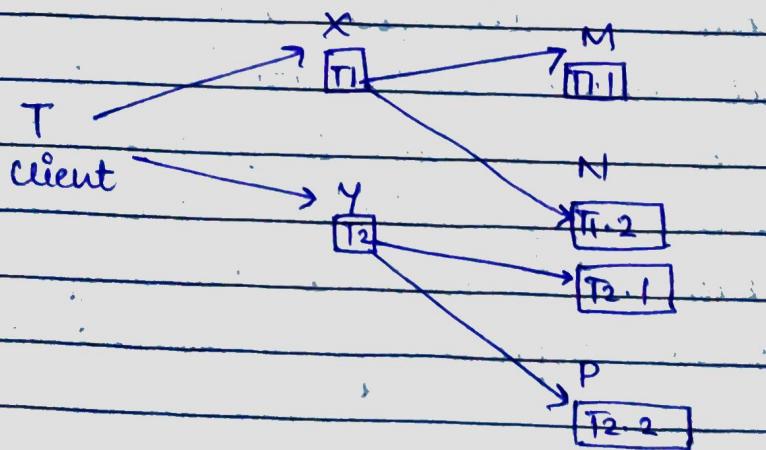


## Limitations:

- ↳ All work is lost in the event of crash
- ↳ only one DBMS is used at a time
- ↳ No partial rollback.

## Nested Transactions:

- ↳ A transaction that includes other transactions within its initiating point & end point

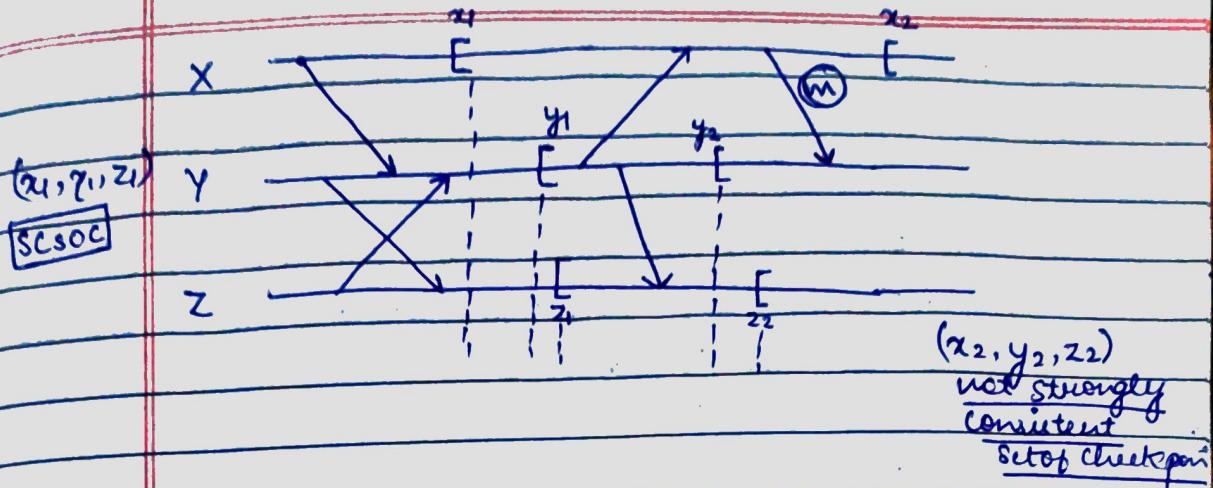


### Checkpointing:

- ↪ A mechanism that enables transactions to recover from inconsistent state using backward error recovery.
- ↪ In DS, all the sites save their local ~~inter~~ states which are known as local checkpoints, and the process of saving local state is called local checkpointing.
- ↪ All the local checkpoints, one from each site, collectively form a global checkpoint.

Strongly consistent set of checkpoints

- ① Domino effect is caused by orphan mgs which themselves are due to rollbacks.
- ② To overcome domino effect, a set of local checkpoints are needed such that no info. flow starts place b/w any set of processes without saving checkpoints.
- ③ A set of checkpoints is known as Recovery line or strongly consistent set of checkpoints.
- ④ A Global checkpoint is SCSOL if there is no orphan msg & no lost msg.



Method of taking consistent set of cp:

- ① If every process takes a checkpoint after sending every msg, the set of most recent cps are always consistent.
- ② However it has high overhead, so we can resolve this, by taking cp after every k, K71

### Recovery in DDBS:

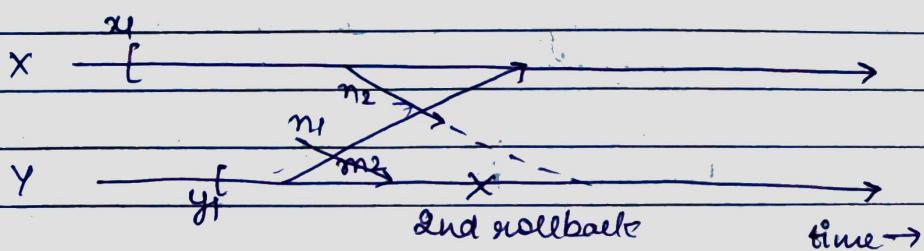
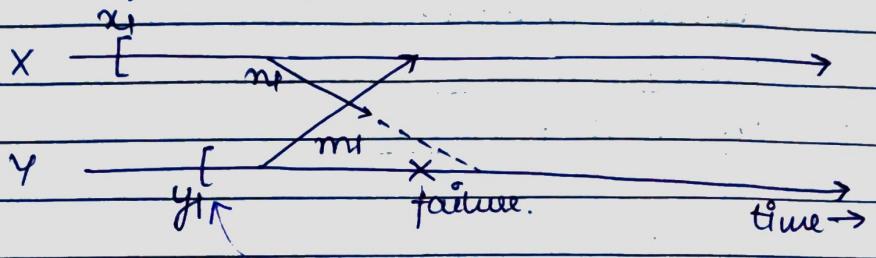
- 1) To enhance performance, availability, a DDBS is replicated whose copies are stored in diff sites.
- 2) The copies of DB at the failed site may miss some updates whereas sites which are not operational have inconsistent copies of data.
- 3) Proposed Approach
  - a) Message spooler
  - b) Copies transactions

Save all the update Read up-to-date copies at all directed towards operational sites & update failed sites the copies at recovering site

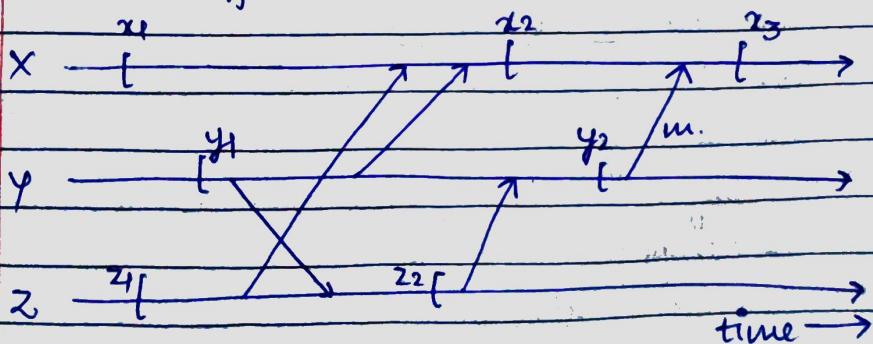
- 4) Recovery algorithms should guarantee that-
- ↳ The out of date replicas are not accessible to user transactions.
  - ↳ Once outofdate replicas are updated by copies transactions, they are also updated along the other copies of user transactions.

Livelocks:

In rollback recovery, livelock is a situation in which single failure can cause an infinite no. of rollbacks preventing the system from making progress.

Domino Effect:

When rolling back one process causes one or more other processes to rollback it is known as Domino effect.



### failure Resilient Systems:

- ↳ A process is said to be resilient if it makes failures and guarantees progress despite a certain no. of system failure.
- Appropriate proposed to implement Resilient process:

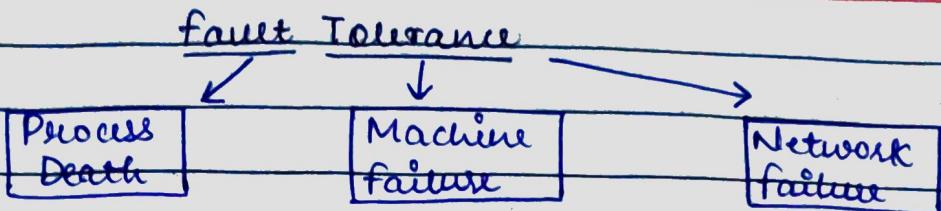
  - Backup processes
  - Replicated Execution

### Backup process:

- ↳ ① primary process & ② backup process.
- ↓  
terminates because of failure → backup process becomes active & takes over the func of primary process.
- ↳ To minimize computation that was to be redone by the backup process, the state of primary process is stored at appropriate intervals.

### Replicated Execution:

- ↳ several processes execute same program concurrently
- ↳ As long as one of the processes survive failures, the computation/ service continues.
- ↳ Increases reliability & availability.

Process Death:

- ↳ If process has reached a deadlocked state or had terminated abruptly, then all the resources allocated to the process must be recouped
- ↳ In client-server architecture, if client dies then server must be made aware about it so that all the resources can be released & vice-versa

Machine failure:

- ↳ In this case, all processes running in the machine will die.

Network failure:

- ↳ A communication link failure can partition a Network into subnet, making it impossible for a machine to communicate with another machine in a different subnet.

## Commit Protocols

- These protocols are used to ensure atomicity across sites.
- A transaction which executes at multiple sites must either be committed at all sites or aborted at all the sites.

### 2-Phase Commit Protocol

#### Assumptions:

- ↳ one of the cooperating process acts as a coordinator other is referred as cohorts.
- ↳ stable storage is available at each site.
- ↳ each site uses write-ahead protocol to achieve local fault recovery & rollback.
- ↳ At the beginning of transaction, the coordinator sends start transaction msg to every cohort.

#### A) PHASE 1:

- ↳ At coordinator.

- a) coordinator → COMMITREQUEST → cohorts
- b) wait for reply.

- ↳ At cohorts.

- a) If transaction executing at cohort is successful, it writes UNDO + REDO log on stable storage & sends AGREED msg to Coordinator else send ABORT msg to the coordinator.

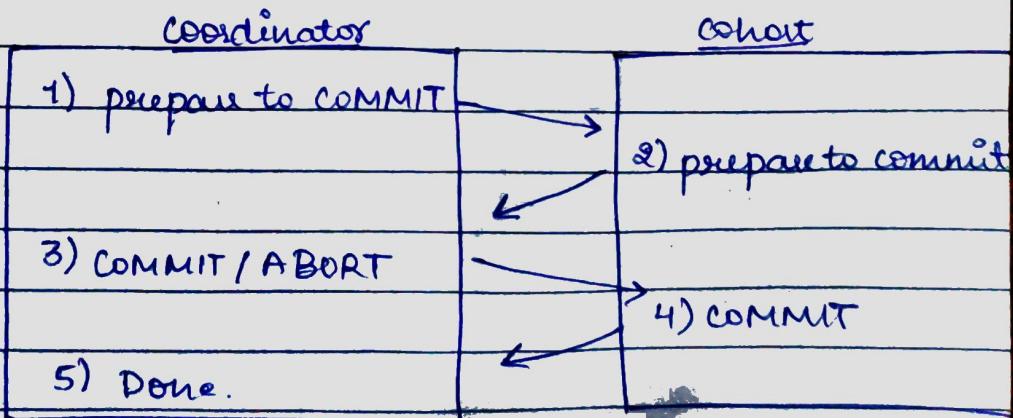
B) PHASE 2:

## ↳ at coordinator :

- If all cohorts agreed & coordinator agrees, then coordinator writes a COMMIT record into the log. & sends COMMIT msg to all the cohorts else ABORT msg.
- Wait for ACK from each cohort, if ACK is not received from any cohort within timeout, coordinator resends the COMMIT/ABORT msg to that cohort.
- If all ACK's are received then coordinator writes a COMPLETE record to the log.

## ↳ at cohort :

- On receiving COMMIT msg, a cohort releases all the resources & locks held by it for executing transactions & send an ACK.
- On receiving ABORT msg, cohort undoes the transaction using UNDO tag record, releases all the resources & locks held & send an ACK.



## Voting Algorithm

- ↳ More fault Tolerant
- ↳ Allow data access under
  - ① Network partition
  - ② Site failure
  - ③ Net losses
- ↳ Every replica is assigned certain no. of votes.  
This info is stored in stable storage.
- ↳ R/W opr is permission when a certain no. of votes are collected by the req. process.

### Assumptions / Basic Idea

- ① each file access requires that an appropriate lock access is obtained (R/W lock)
- ② each site has lock manager
- ③ each file has a version no. =  $\frac{\text{no. of time it has been updated}}{\text{}}$
- ④ each replica has certain no. of votes.
- ⑤ vote allocation  $\rightarrow$  stable storage.
- ⑥ R/W opr requires a quorum

stable storage

lock manager

Static Voting AlgoRequesting set,  $R = \{S_1, S_2, \dots, S_n\}$ 

- 1)  $S_i$  send LOCK REQ to its local lock manager
- 2) If request granted,  $S_i$  sends VOTE REQ to all sites in  $R$ .
- 3)  $S_j$  send LOCK REQ, if granted then it return Version No. & No. of votes assigned to replica
- 4)  $S_i$  decides whether it has quorum or not.
  - ↳ If  $S_i$  does not have quorum, issue RELEASE LOCK
  - ↳ If  $S_i$  is successful, checks if its copy of file is current
  - ↳ If request is read, read local copy
  - ↳ If request is write, update copy & Version No. and send RELEASE LOCK request.

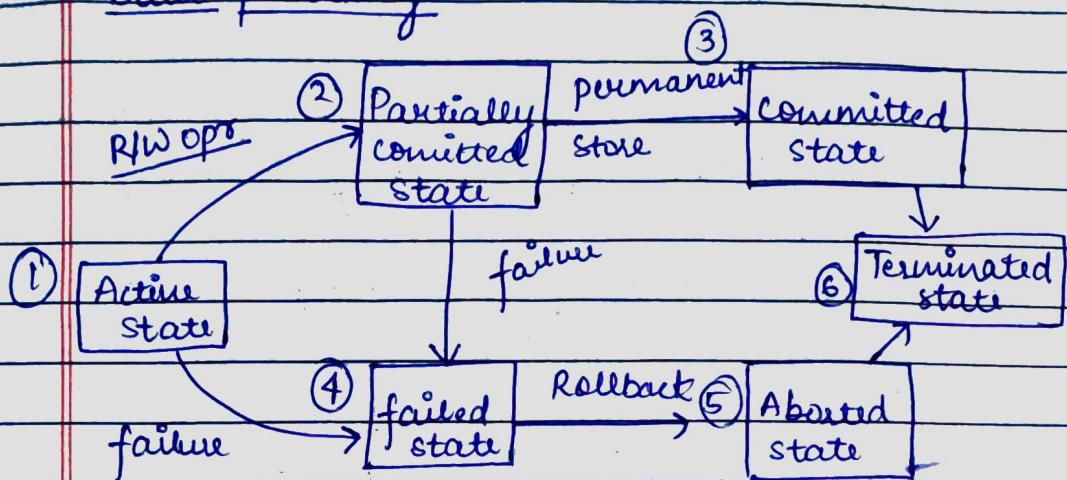
Dynamic Voting Protocols

- ↳ change set of sites that can form majority
- ↳ change distribution of votes

Static Node ProtocolDynamic Node protocol

- |  |   |
|--|---|
| 1) Non-adaptive in nature                                    | 1) Adaptive.  |
| 2) can cause comm. failure                                   | 2) prevent comm. failure                                      |
| 3) may / may not ensure availability                         | 3) ensure <u>availability</u> .                               |
| 4) Selection time of replica doesn't depend on system state. | 4) selection time of replica depends on <u>system state</u> . |

Transaction : A program including collection of DB oprs, executed as a logical unit of Data processing



### Properties of Transactions (ACID)

- 1) Atomicity:
  - ↳ either transaction executes completely or it does not occur at all.
  - ↳ Responsibility of Transaction Control Manager
- 2) Consistency:
  - ↳ ensures integrity constraints are maintained, system remains consistent before & after the transaction.
  - ↳ Responsibility of DBMS & Application programs
- 3) Isolation:
  - ↳ ensure multiple transactions can occur simultaneously without causing inconsistency, each transaction exec. is isolated.
  - ↳ responsibility of Concurrency Control Manager

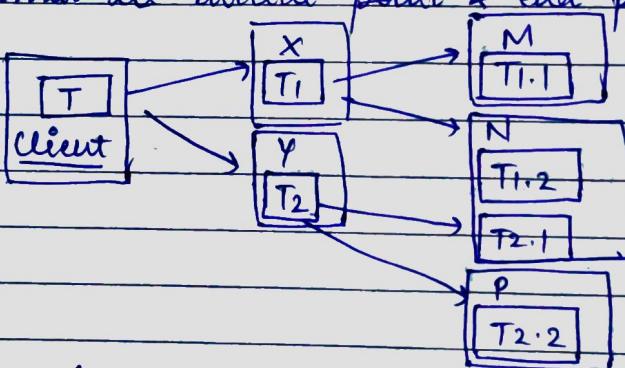
4)

Durability:

- ↳ ensures all the changes made by transaction after successful execution are written successfully in disk.
- ↳ Changes never lost in case of failure.
- ↳ Responsibility of Recovery Manager.

Nested Transactions

- ↳ A transaction that includes other transaction within its initial point & end point



4

Rules:

- ① When a parent transaction commits then all the sub transactions that have provisionally commit can commit.
- ② When a parent aborts, all of its subtransactions are aborted.
- ③ When a child completes, it makes an independent decision either to commit provisionally or abort.
- ④ When a child aborts, parent can decide whether to abort or not.

Locks : Variable associated with shared resource that ensures serializability, by providing access to a data item in natively exclusive manner

Types of locks:

1) Binary locks

↳ locked (1)

↳ unlocked (0)

2) Shared / Exclusive locks

(i) Shared: If  $T_i$  locks  $X$  in shared mode then before  $T_i$  unlock  $X$ , no other transaction  $T_j$  can write into  $X$ , only  $T_j$  can read.

(ii) Exclusive: If  $T_i$  lock  $X$  in exclusive mode then before  $T_i$  unlocks  $X$ , no other transaction  $T_j$  can read or write into  $X$ .

Strict  $\rightarrow$  2PL  $\rightarrow$  NID shrinking phase  
 $\rightarrow$  Release all the locks when whole transaction commit.

Date \_\_\_\_\_ / \_\_\_\_\_ / \_\_\_\_\_

Page No.: \_\_\_\_\_

### Two Phase locking Protocol

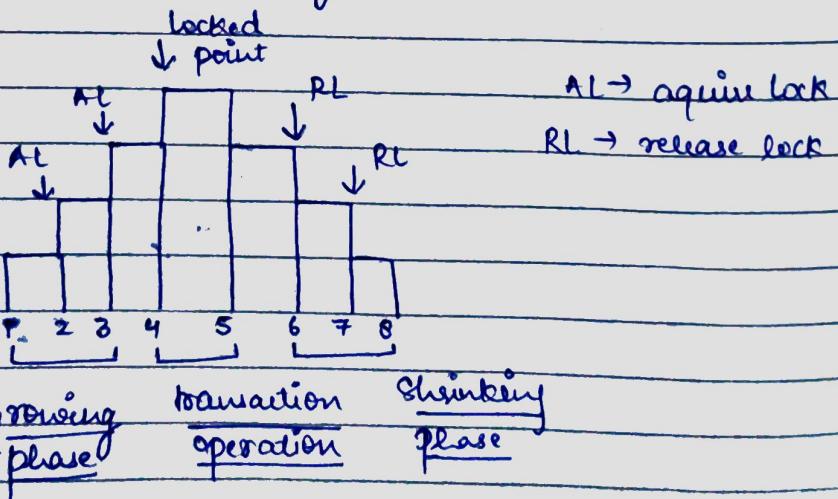
2PL is a method of concurrency control in DB that ensures serializability by applying lock to transaction Data which blocks other trans. to access same data simultaneously.

#### 1) Growing Phase:

- When transaction begins to execute, it requires permission of locks <sup>on DI</sup> it needs.
- In this phase, transaction may obtain any locks but may not release any lock.

#### 2) Shrinking phase:

- After executing transaction opr., it releases all the required locks on DI
- In this phase, transaction can only release locks but not obtain any locks



## Timestamp Ordering Protocol

- $TS(T_i)$  Timestamp : unique identifier created by DS to identify a transaction for two transaction,  $T_i$  &  $T_j$
- If  $TS(T_i) < TS(T_j)$   
 $\Rightarrow T_i$  started before  $T_j$

with every shared data item, 2 timestamps are associated.

- (i)  $R\_TS(X)$   $\rightarrow$  largest TS value that execute successful read operation  $Read(X)$
- (ii)  $W\_TS(X)$   $\rightarrow$  largest TS value that execute successful write operation  $Write(X)$

### Conditions :

- 1)  $T_i$  issues  $Read(X)$  opr
  - $\hookrightarrow$  If  $W\_TS(X) > TS(T_i)$  then opr is rejected
  - $\hookrightarrow$  If  $W\_TS(X) \leq TS(T_i)$  then opr is executed.
  - $\hookrightarrow$  timestamps of all data items updated.
- 2)  $T_i$  issues  $Write(X)$  opr
  - $\hookrightarrow$  If  $TS(T_i) < R\_TS(X)$  then opr is rejected
  - $\hookrightarrow$  If  $TS(T_i) < W\_TS(X)$  then opr is rejected &  $T_i$  is rolled back.
  - $\hookrightarrow$  else operation is executed.

## Validation-Based Protocol (Optimistic Concurrency Control Technique)

In VBP, transaction is executed in 3 phases

### 1) Read Phase:

↳ Transaction  $T_i$  is read & executed

↳ It is used to read the value of various data items and store them in temp variable. It can perform all write operations on temp variable without update of actual database.

### 2) Validation Phase:

In this phase, temp variable is validated against the actual data to see if it violates the serializability.

### 3) Write Phase:

If validation of transaction is validated, then the temporary results are written in db/system otherwise transaction is rolled back.

Each Phase has the following Timestamps:

(i) start( $T_i$ ): time when  $T_i$  starts its execution

(ii) validation( $T_i$ ): time when  $T_i$  finishes its read phase and starts validation phase.

(iii) finish( $T_i$ ): time when  $T_i$  finishes its write phase.

## Comparison of Methods of Concurrency Control

### Locking

The locking schemes are used to restrict the availability of DI for other transactions.

### Timestamp

It works on checking of TS. TS is the start time of transaction generated by logical clock.

### OCC

All transactions are allowed to proceed but some are aborted when they attempt to commit.

### Timestamp ordering

1) Decide serialization order statistically  
2) Better than locking for Read dominated trans.

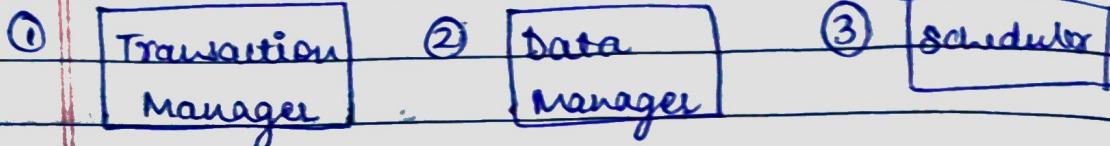
### Two phase lock Protocol

1) decide serialization order dynamically  
2) Better than TS ordering for update-dominated trans.

Both are pessimistic Methods

## Concurrency Control in Distributed Transaction

DBS contains 3 modules



### 1) Transaction Manager:

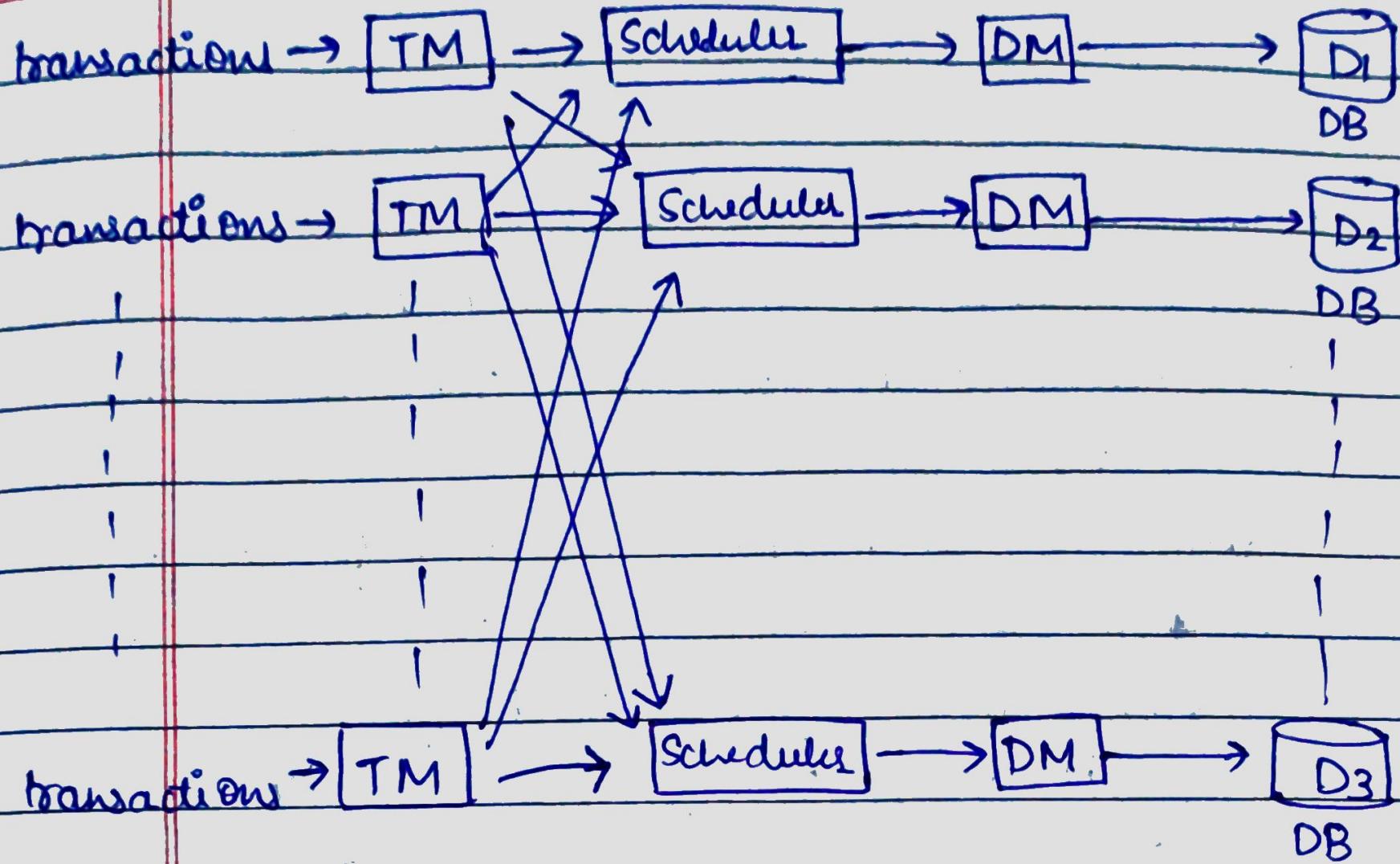
- ↳ supervise the execution of transaction
- ↳ interact with DM to carry out the est.
- ↳ It assign a TS to a transaction or issue request to lock/unlock data obj. on behalf of user.
- ↳ Acts as interface b/w user & DBS

### 2) Scheduler:

- ↳ used to enforce concurrency control.
- ↳ It grant/ release locks on data objects as requested by a transaction.

### 3) Data Manager:

- ↳ Manage the DB
- ↳ It carries out R/W request issued by TM on behalf of a transaction by operating them on DB
- ↳ DM is interface between scheduler & DB
- ↳ Execution of transaction at TM results in execution of its actions at the DM.

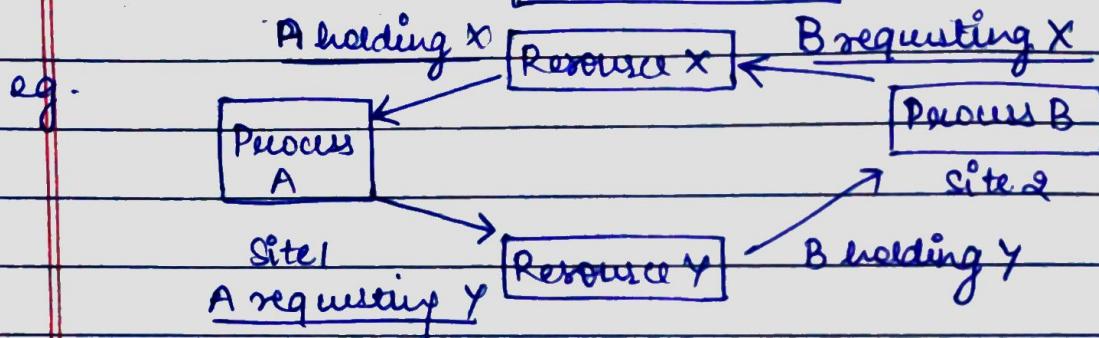


## Distributed Deadlock

DS consists of a no. of sites connected by a network  
 each site with globally unique identifier,  
 makes resource req. to a controller  
 The controller at each site could maintain  
 WFG on a process request.

Each site WFG could be cycle free & yet DS  
 could be in Deadlock.

This is called Global Deadlock.



## Transaction Recovery : **REPLICATION**

- Replication is the maintenance of copies of data multiple sites.
- Enhance performance, Higher Availability, Fault Tolerance, Increase Reliability, Less Access time, Easy load Balancing

Problem : Data Inconsistency

## Group communication

- ↳ It occurs when source process sends a msg to a group of processes (Multicast Commn)
- ↳ Groups are useful for managing Replicated data and in other system where processes cooperate towards a common goal by receiving & processing the same set of multicast msgs.

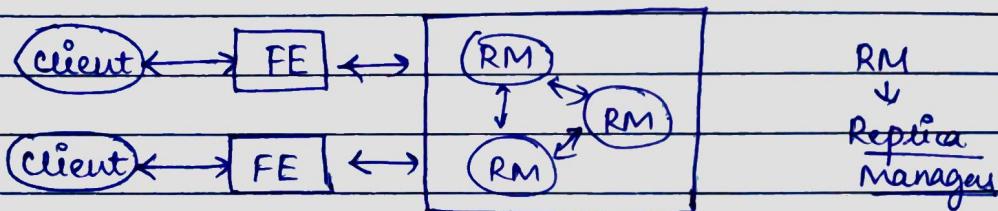
## Highly Available Services:

Replication technique is used to make service highly available.

i.e. client access the service with reasonable response time.

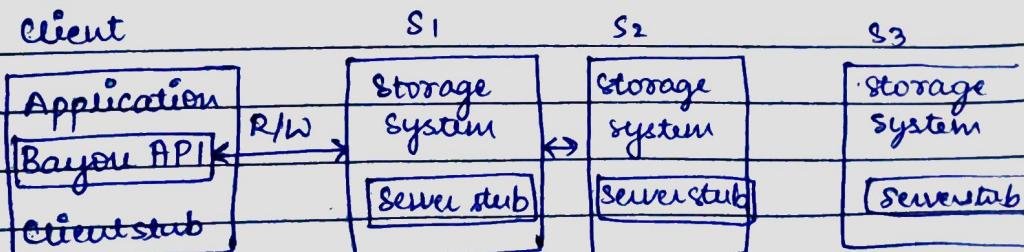
## (i) Group Architecture:

Client request service operation which are processed initially by front-end.



→ Replica manager update one another by exchanging gossip msg which contain the most recent update.

## (ii) Bayesian Basic System Model:



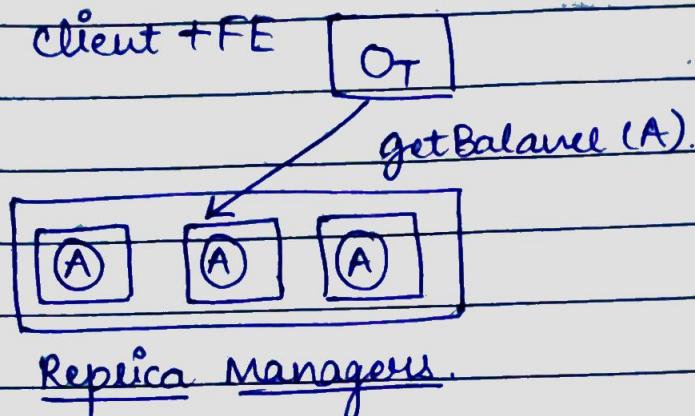
↳ It is a replicated, weakly consistent storage designed for mobile computing envt.

### (iii) Coda file system:

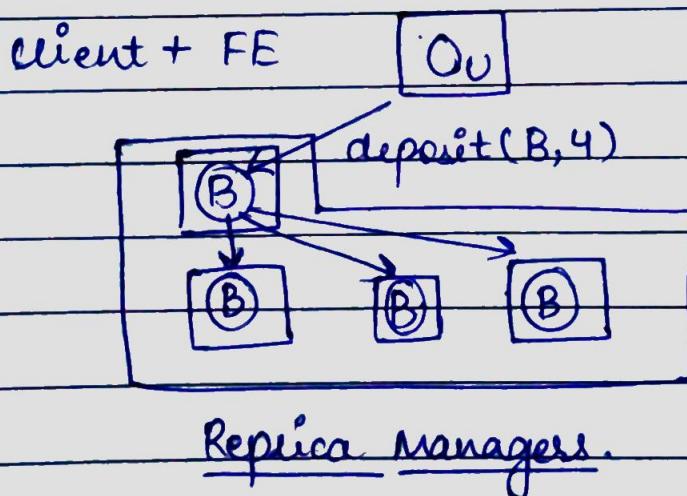
- ↳ Purpose of CFS is enabling disconnected operations.
- ↳ Client communicate over high bandwidth network with server.
- ↳ Coda use optimistic replica control for high availability.

## Transaction with Replicated Server

(i) Read operation:



(ii) Write Operation:

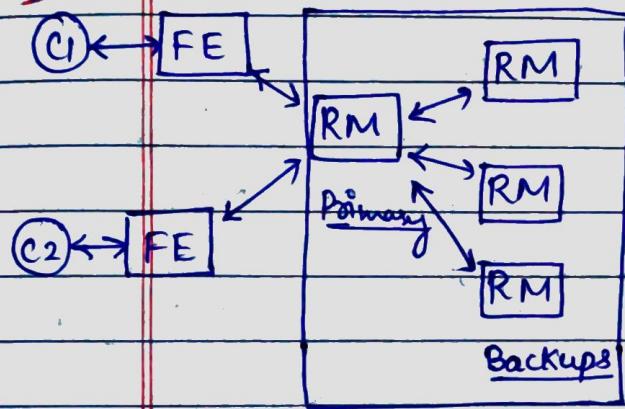


## Fault Tolerant services

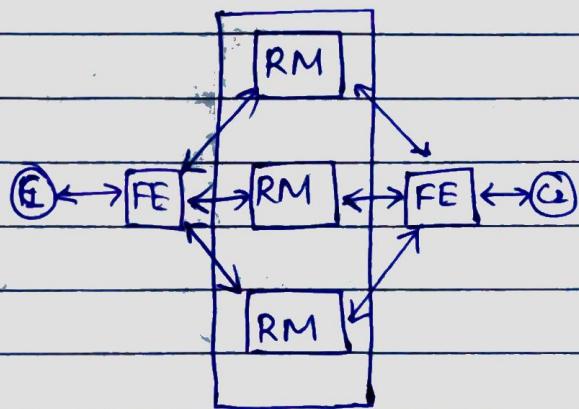
- ↳ FTS are obtained by using replication.
- ↳ By using Multiple Independent Server replica Managers each managing replicated data.

Primary  
Backup  
Replication

### Passive Replication



### Active Replication



RM acts as state machines

### Sequence of Events in Both Replications

- ① Request
  - ② Coordination
  - ③ Execution
  - ④ Agreement
  - ⑤ Response
- unique identifier

## Passive Replication :

- ↳ single primary replica Manager, one or more secondary RM ie slaves or Backups'
- ↳ Front end communicates with primary RM to obtain the service.
- ↳ Primary RM executes the operations & sends copies of updated data to backups
- ↳ If primary RM fails, one of the backup is promoted.

- 1) Request: FE issues request, containing Unique Identifier to primary RM.
- 2) Coordination: PRM takes each request atomicity. It checks the UI in case it has already executed the request.
- 3) Execution: PRM executes request & stores the response.
- 4) Agreement: If request is an update then primary send the updated state, response and UI to all the backups. The backup send an agreement acknowledgement.
- 5) Response: PRM respond to FE which sends response back to client.

## Active Replication:

- ↳ RM are state machines and play equivalent roles and are organized as group.
- ↳ FE multicast their requests to Group of RM and all RM process the req. independently but identically & reply.

- 1) Request: FE attach UI to request & multicast to Group of RM
- 2) Coordination: The Group Communication system delivers request to every RM in same order.
- 3) Execution: Every RM executes the request
- 4) Agreement: Not needed
- 5) Response: Each RM sends its response to FE

## Multiversion Timestamp Ordering Protocol

- Multiversion schemes keep old version of data item to increase concurrency.
- each successful write results in creation of new version of data item written.

Each data item  $Q$  has seq. of versions

$\langle Q_1, Q_2, \dots, Q_m \rangle$  Each version  $Q_k$  contains 3 data fields:

- ① [content]
- ② [W-TS( $Q_k$ )]
- ③ [R-TS( $Q_k$ )]

Value of version OK	TS of transaction that created OK	Largest TS of transaction that successfully read OK.
------------------------	--------------------------------------	--

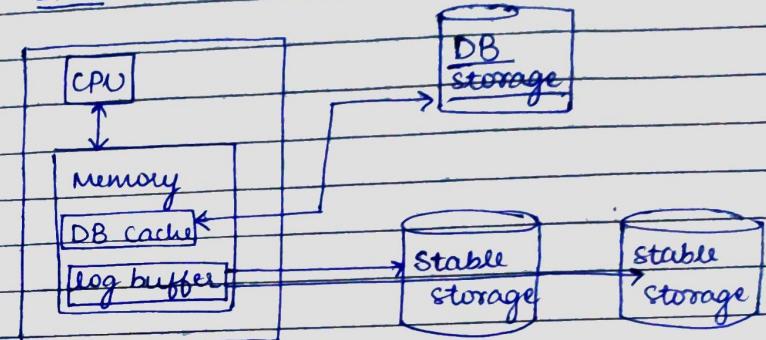
- ↳ When  $T_i$  creates new version  $Q_k$  of  $Q$   
 $Q_k$ 's W-TS and R-TS are initialized to TS( $T_i$ )
- ↳ Suppose  $T_i$  issues read( $Q$ ) or write( $Q$ )
  - 1) If  $T_i$  issues read( $Q$ ) then value returned is content of version  $Q_k$ .
  - 2) If  $T_i$  issues write( $Q$ )
    - (i) if  $TS(T_i) < R-TS(Q_k)$ ,  $T_i$  is rolled back.
    - (ii) if  $TS(T_i) = W-TS(Q_k)$ , content of  $Q_k$  overwritten.
    - (iii) else new version of  $Q$  created.

Advantages :

- 1) It allows more concurrency in D.S.
- 2) Read request never fail & never made to wait
- 3) Improved system responsiveness by providing multiple versions
- 4) Reduces the portability of conflict transactions.

Disadvantages :

- 1) Conflict byz transactions are resolved through rollbacks rather than waits
- 2) Requires Huge amount of storage for storing multiple versions of data objects.
- 3) Does not ensure recoverability & correctness.

Backward Recovery  
Error

- ↳ does not lose info in case of sys failure
- ↳ keep logs & recovery points

Operation Based

- ↳ Record a log (audit trail) of operations performed.
- ↳ Restore previous state by reversing steps

state based

- ↳ Record a snapshot of state (checkpoint)
- ↳ Restore state by reloading snapshot (rollback)

### forward Error Recovery:

- 1) Remove errors in the process or system state if errors can be easily and completely accessed
- 2) Continue process or system forward execution.

### Advantages :

- 1) less overhead

### Disadvantages :

- 1) limited to one
- 2) cannot be used as general mechanism for error recovery.

### Backward error Recovery:

- 1) Restore process or system to previous error-free state and restart from there
- 2) Two approaches to implement this
  - (i) operation Based
  - (ii) state Based

### Advantages:

- 1) Simple to implement
- 2) Can be used as a general mechanism for error recovery.

### Disadvantages:

- 1) No guarantee that fault doesn't occur again
- 2) Some components can't be recovered
- 3) performs penalty.

Fault is anything that deviates from what we expect when something unexpected happens in system.

Errors are manifestation of faults in system and when propagate through system can lead to failure.

### Types of faults

- |   |   |  |
|---|---|--|
| ① <u>Transient</u><br>occur once &<br>goes away | ② <u>Intermittent</u><br>occurs, disappear<br>again | ③ <u>Permanent</u><br>occurs once<br>continues<br>until fixed. |
|---|---|--|

### Types of failures

- |   |                                   |   |  |
|---|-----------------------------------|---|--|
| ① <u>Method</u><br><u>failure</u><br>( <u>process</u><br><u>failure</u> ) | ② <u>System</u><br><u>failure</u> | ③ <u>Disk</u><br><u>failure</u><br>( <u>secondary</u><br><u>storage</u><br><u>failure</u> ) | ④ <u>Media</u> /<br><u>Network</u><br><u>failure</u> |
|---|-----------------------------------|---|--|