

## Distributed Systems

↳ A system in which components are connected by a communication network, communicate and coordinates 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.

## HOSC!

### Challenges:

- 1) Heterogeneity: (DS Must be constructed from variety of diff. OS, Comp., HW, lang.)
- 2) Openness: (DS must be Extensible, implemented 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 Machines 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 same process , a sending msg , b recipient

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.

(Q1) for any 2 event , a & b in Process  $P_i$  if a happened before b ( $a \rightarrow b$ ) then.

$$[c_i(a) < c_i(b)]$$

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

(Q2) If a is event of sending msg in  $P_i$  and b is event of receiving msg in at process  $P_j$

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

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

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

Date \_\_\_\_\_

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

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

[IR1]

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

then

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

$d$  = drift time (generally 1).

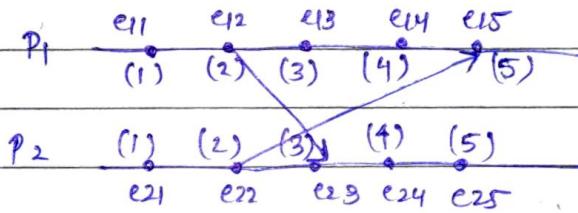
[IR2]

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



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

$$C_{P_2} = 0$$

casual ordering

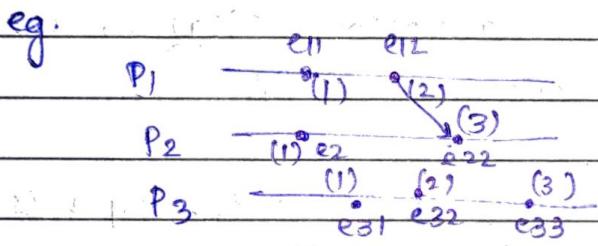
## Limitations of Lamport's clock

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

e.g.



$$c(e11) < c(e12)$$

$$c(e11) < c(e23)$$

$e11$  is causally related to  $e22$  but not  $e32$

- 3) since each clock independently advance due to  
occurrence of local events in process

Lamport  
clock  
system

LCS 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.

generates  
detects

partial ordering  
causality violation

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

### Vector Clock

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

Let n = 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<sup>th</sup> entry of  $c_i$ , corresponds to  
 $P_i$ '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  $P_i$

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

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

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

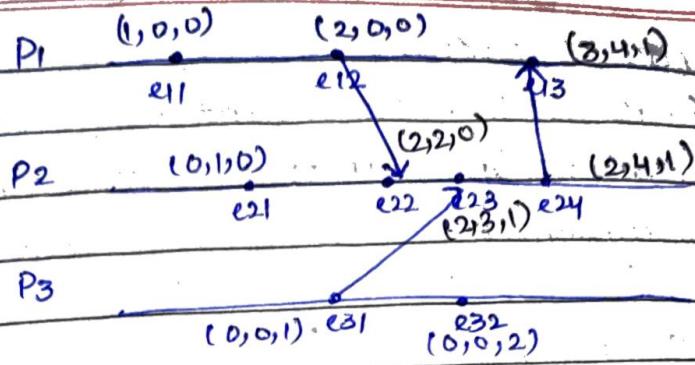
[IR2] If a is sending msg m by process  $P_i$   
then msg m is assigned a vector ts  
 $t_m = c_i(a)$

on receiving msg m by process  $P_j$

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

$$t_K \quad c_j[K] = \max(c_j[K], t_m[K])$$

eg.



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

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

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

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

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

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

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

$$e_{24} = (c_i^0 + d) = \underline{(2,4,1)} \checkmark$$

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

Advantages of vector clock over Lamport's

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

2) An system of vector clocks

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

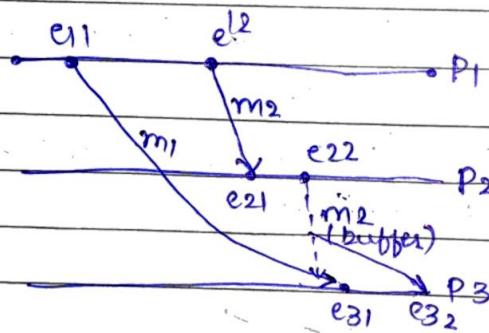
Thus, causal ordering can be determined.

Causal Ordering of Msg maintaining  
 Deals with the concept of "causal  
 relationship among 'msg sent' event  
 corresponding to 'msg received' event

If 1) send( $m_1$ )

2) Send( $m_2$ )

then every recipient of  $m_1$  &  $m_2$  msg.  
 must receive  $m_1$  before  $m_2$ .



Broadcast

Date \_\_\_\_\_

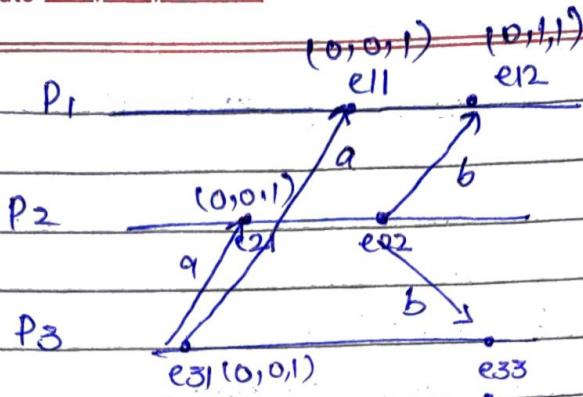
Bizman Sniper Stevenson Protocol:

Algo

- 1) Each process increases 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 clock for process.

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

- 1)  $P^i$  sending msg. to  $P^j$   
 $P^i$  increments  $c^i(i)$  and sets timestamp  
 $tm = c^i(i)$  for m.
- 2)  $P^j$  receiving msg. from  $P^i$
- 3) When  $P^j$ ,  $j \neq i$  receives m with ts  $tm$ , it delays msg. delivery until both are follows.
  - (1)  $c^j(k) = tm[i] - 1$  &
  - (2)  $\forall k < n \text{ s.t. } k \neq i \quad c^j(k) \leq tm[k]$
- 4) When msg. is delivered to  $P^j$ , update  $P^j$ 's v.c.
- 5) 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

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

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

 $e_{12}$  $P_1$  receives msg b as  $c_1 = (0, 0, 1)$ 

$$c_1[2] = t^b[2] - 1 = 1 - 1 = 0$$

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

$$c_1[3] = t^b[3] = 0$$

msg is accepted  $c_1 = (0, 1, 1)$

E32 P<sub>3</sub> receives msg b as, C<sub>3</sub> = (0, 0, 1)

$$c_3[2] = t^b[2] - 1 = 1 - 1 = 0$$

$$c_3[1] = t^b[1] = 0$$

$$c_3[2] = t^b[2] = 0$$

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

A(i)

### Transitless Global state

A Global state is transitless iff

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

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 computation

Need

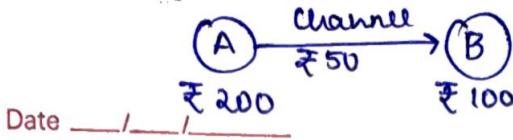
(1) Distributed Deadlock Detection

(2) Termination Detection

(3) Check point

& the state of comm. channels.

consistent  
Global  
state



Date \_\_\_\_\_

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

## Chandy Lampart'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 msgs.
- 4) Marker separates msg in the channel.

### Algorithm:

- 1) Marker sending rule process P :
  - a) P processes 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 msgs 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 marker sending rule.

else

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

## Termination Detection

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

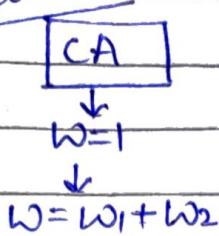
# example of use of Coherent view of DS.  
Notations:

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

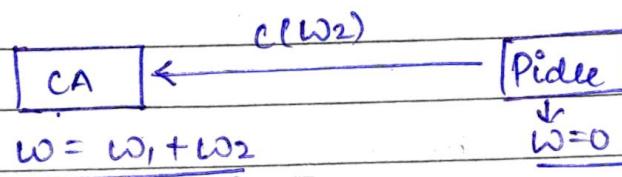
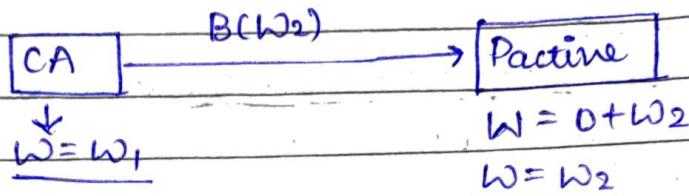
### Huang's Terminations Algo:

- 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.  $Dw = w_1$
  - ↳ send  $B(Dw)$  to process P.  $Dw = w_2$
- 2) On receiving B(DW) by process P:
  - ↳ Add Dw to weight of process P,  $w = w + Dw$
  - ↳ If P 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  $C(DW)$  to CA.
  - ↳ send control msg  $C(DW)$ ,  $Dw = w$ .
  - ↳ set weight of process as 0.  $[w=0]$

## Basic Implementation



$$\boxed{\text{Piddle}} \\ w=0$$



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

- ↳ Add weight through control msg to
- ↳ net of CA, i.e.  $w = w + DW$
- If After adding, CA weight  
 $w = 10 \Rightarrow$  process computation terminated

## 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.

## 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
- ↳ Rikant Agrawala
- ↳ Markowitz

Non token

## Non-Token Based Algo

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

### Example:

- ↳ Raymond's
- ↳ Suzuki-Kasami
- ↳ Singhal's Heuristic

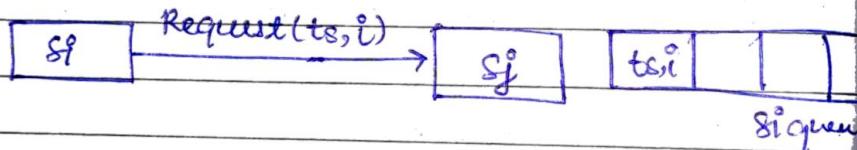
## Lamport's Non-Token Algorithm

### Requesting Set

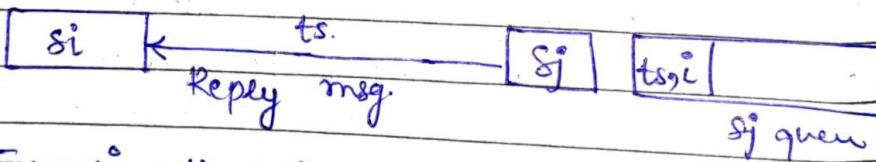
$R_i = \{S_1, S_2, \dots, S_n\}$ , every site has queue.

#### 1) Requesting the CS:

- a) Site  $S_i$  wants to enter the CS, it sends Request( $t_s, i$ ) msg to all the sites in  $R_i$  & places request in  $S_i$ 's req-queue.



- 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  $S_j$ 's queue.



#### 2) Executing the CS:

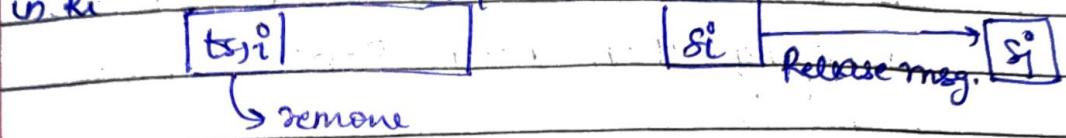
$S_i$  enters in CS if

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



3) Releasing the CS:

- a) On exiting SJ removes its request from top of RO & send the TS Release msg. to all sites in its queue.



- b) On receiving the release msg. SJ removes its request msg. from its RO.

Performance: Require  $3(N-1)$  msg / CS. involves

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.

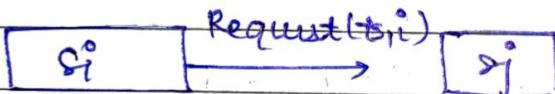
## Rickett-Agarwala Non-Token Algorithm

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

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  $t_s$ .  
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
- 2)  $(N-1)$  Request msgs
- 3)  $(N-1)$  Repay msgs.
- 4) synchronization Delay T.

Drawback :Unreliable approach

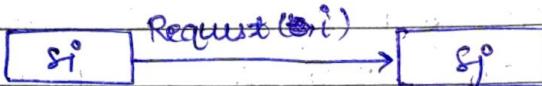
failure of any one node in system can halt the progress of system. In this situation process will stall.

## Markkawai Algorithms

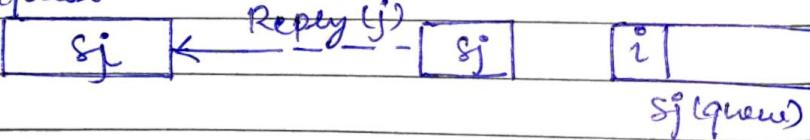
Requesting set,  $R_i = \{s_1, s_2, \dots, s_n\}$

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 & delete 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:

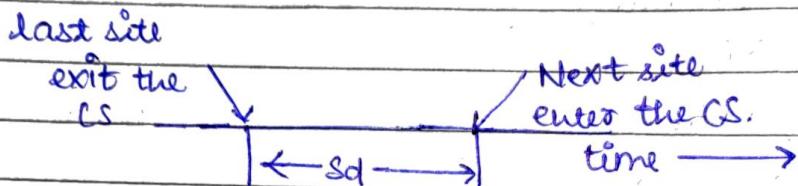
- 1) JN msgs per CS  
JN request msgs  
JN reply msgs  
JN release msgs.
- 2) synchronization Delay =  $\alpha T$  (msg. propagation Delay).

### Drawbacks:

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

## Performance Metric

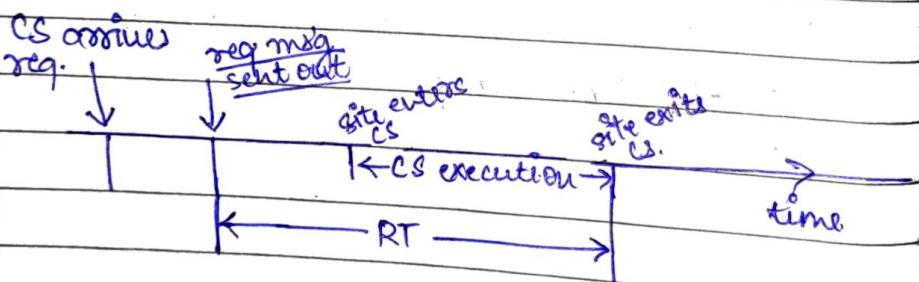
- 1) No. of msgs sent per CS:
- 2) Synchronization Delay



### 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.

$Sd$  = Synchronization delay.

$F$  = Avg. CS execution time

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

	Algorithm	RT(ll)	sd	Msg.(ll)	Msg(hl)
1)	Lamport	$2T+E$	T	$3(N-1)$	$3(N-1)$
2)	Ricart Agarwala	$2T+E$	T	$2(N-1)$	$2(N-1)$
3)	Makawa	$2T+E$	$2T$	$3\sqrt{N}$	$5\sqrt{N}$
4)	Suzuki Kasami	$2T+E$	T	N	N
5)	Singhal Heuristic	$2T+E$	T	$N/2$	N
6)	Raymond Tree	$T \log N + E$	$T \log N/2$	$\log N$	4

### Remember

- # Chandy-Lamport's Global state (Marker)
- # Casual Ordering Bismarck Schiper Stephenson
- # Termination Detection - Haeng's Termination CA
- # Lamport's logical clock  $c_j = \max(c_j, t_m + d_j)$  Rule
- # Vector clock  $c_j[k] = \max(c_j[k], t_m[k]) + 1$

## 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.
- 3) Security
- 4) 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 - Egeli - Sandor

#### Threads over processes.

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

### 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. Reg. made from Proxy server may help to hide client's IP from web server.

### Reason Middleware moved from distribution 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  $a, b, c$  be 3 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 & types of architecture exists that are usually used for distributed computing.

↓  
low level

Interconnection

of Multiple CPUs

↓  
high level

Interconnection of

processes running on the  
CPUs.

- a) Client Server Model. → Dependency
- b) 3-Tier Architecture. → No dependency
- c) Peer to Peer Model.
- d) Proxy server / cache
- e) 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  $n$  processors in DS, then only  $m$  processors out of them may be found as faulty processors.
- (ii) every processor in the system is free to communicate with each other in the system due to their logical connections with each other.
- (iii) A receiver processor always knows the identity of sender processor of msg.
- (iv) The communication medium is reliable and only processors 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) Intertemporal 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 arbitrary 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.

Protocols must meet following objectives:

(i) Agreement: All non-faulty processors agree on the same single value.

(ii) Validity: All 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 processors.

Protocol 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:

One Message Algo  $OM(m)$   $m \geq 0$  solves BAP  
 for  $3m+1$  processors in the presence of at most  
 $m$  faulty processors.

Let  $n =$  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 =

Algorithm  $OM(m)$ :

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

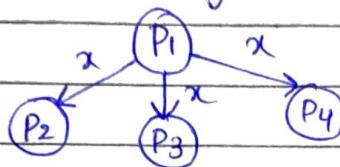
The majority function ( $v_1, v_2, \dots, v_{n-1}$ ) computes the majority value if it exists else it uses default value '0'.

BA cannot always be reached among 4  
Proof: processor if 2 processors are faulty.

4 processes: P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>, P<sub>4</sub>.

P<sub>1</sub> initiates the initial value.

P<sub>2</sub>, P<sub>4</sub> are faulty.



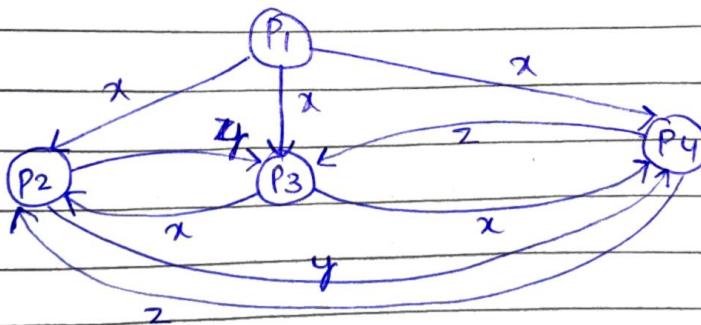
P<sub>1</sub> executes OM(1) and sends its value x to other processors.

After receiving the value x from P<sub>1</sub>,

P<sub>2</sub>, P<sub>3</sub>, P<sub>4</sub> execute OM(0).

• P<sub>3</sub> is non-faulty & sends x to P<sub>2</sub> & P<sub>4</sub>

P<sub>2</sub> & P<sub>4</sub> are faulty send y to P<sub>3</sub>, P<sub>4</sub>  
 and z to P<sub>2</sub>, P<sub>3</sub> resp.



Majority values for Byzantine solution:

processors	Received Majority	common
------------	-------------------	--------

P <sub>2</sub>	(x, x, z)	x
----------------	-----------	---

P <sub>3</sub>	(x, y, z)	0
----------------	-----------	---

P <sub>4</sub>	(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 Protocol:

- 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 \_\_\_\_\_ / \_\_\_\_\_ / \_\_\_\_\_

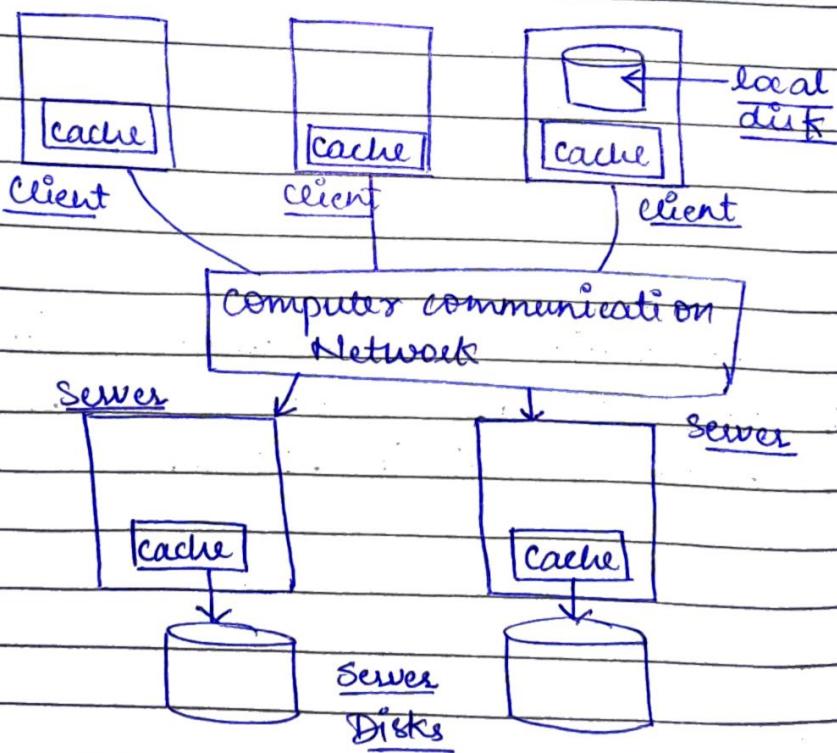
Page No.: \_\_\_\_\_

## 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 Service:

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

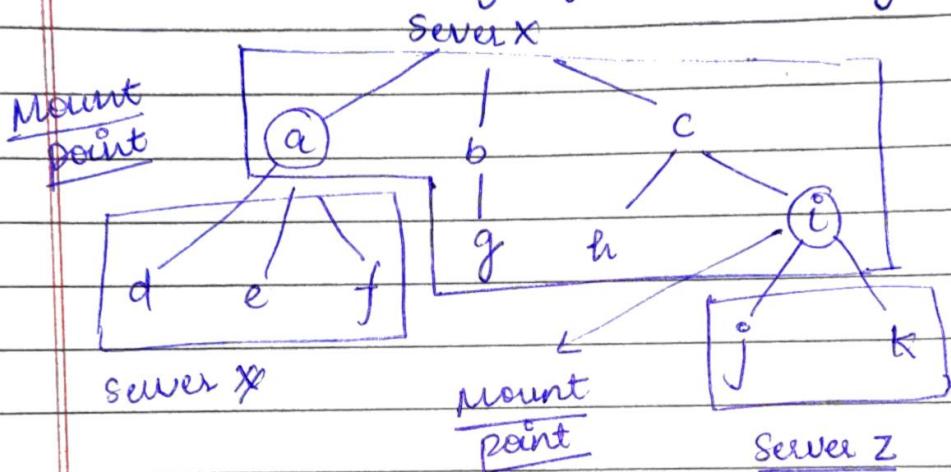
② Cache Manager:

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~~

- ↳ 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)

### Cache 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.

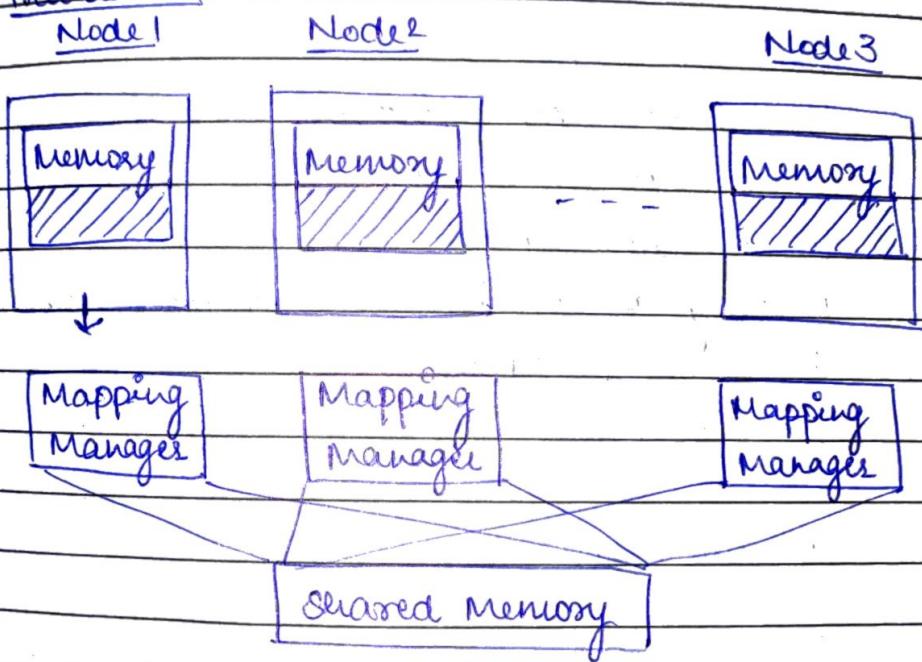
7)

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.
- Page Replacement:  
Replacement algo must take into acc. pg. access mode.  
↳ shared, private, readonly, writable.  
LRU: private pg. replaced before shared ones, private pg. map to dev. Shared pg sent over network to owner, 'Read only' maybe discarded.

## Memory Coherence

- ① Replicated shared data Obj.
- ② concurrent access to data obj at many nodes

① Sequential consistency:

② General consistency:

③ processes consistency:

opr. issued by processes 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.

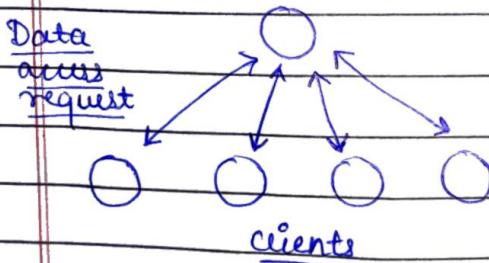
② Write update protocol.

## 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 associating seq. no. with write requests.
- ↳ If an application's req. to access shared data fails repeatedly, a failure condition is sent to the application.

Central  
Server



Clients

Send data request

Central Server

Receive Request  
perform data access  
send response

Receive response

## Migration Algorithm.

- ↳ In migration algo, data is migrated to the loc. of data access req, allowing subsequent accesses to data to be performed locally.
- ↳ Allow only one node to access shared data at a time.

# Read Replication Algorithm

## full Replication Algorithm