

DISTRIBUTED SYSTEMS

① CLOCK SYNCHRONISATION

1. CLOCK synchronisation is the process of aligning the clocks of multiple devices or systems to ensure that they maintain a consistent notion of time.
2. In distributed systems, various computers and nodes work together to achieve a common goal. This is where clock synchronisation is crucial to ensure that operations are correct and coordinated.
3. In DS, certain processes may run concurrently on multiple computers located geographically apart from each other at larger distances.

In such situations, to run a process accurately across multiple nodes, the clocks of all nodes must be synchronised with each other.

4. To address these challenges, clock synchronisation protocols like NTP (Network time protocol) and PTP (Precision time protocol) are implemented employed in DS. These protocols ensure that different nodes have a consistent notion of time.

5. CLOCK synchronisation is a crucial aspect of DS. Proper synchronised clocks enable accurate event ordering, data consistency and coordinated operations, contributed to overall correctness and efficiency of distributed applications and services.

ISSUES WITH CLOCK SYNCHRONISATION

1. In real world, not two clocks are perfectly synchronised. But some specified time difference b/w clocks is tolerated.
↳ This specified time difference is a constant denoted by δ .

Difference b/w two clocks is called

CLOCK SKEW

and CLOCK SKEW should not be greater than δ .

TWO CLOCKS ARE SAID TO BE SYNCHRONIZED IF

$$\text{CLOCK SKEW} < \delta$$

2. In OS, errors occur due to unpredictable communication delay in message exchange to deliver a clock message and clock signals

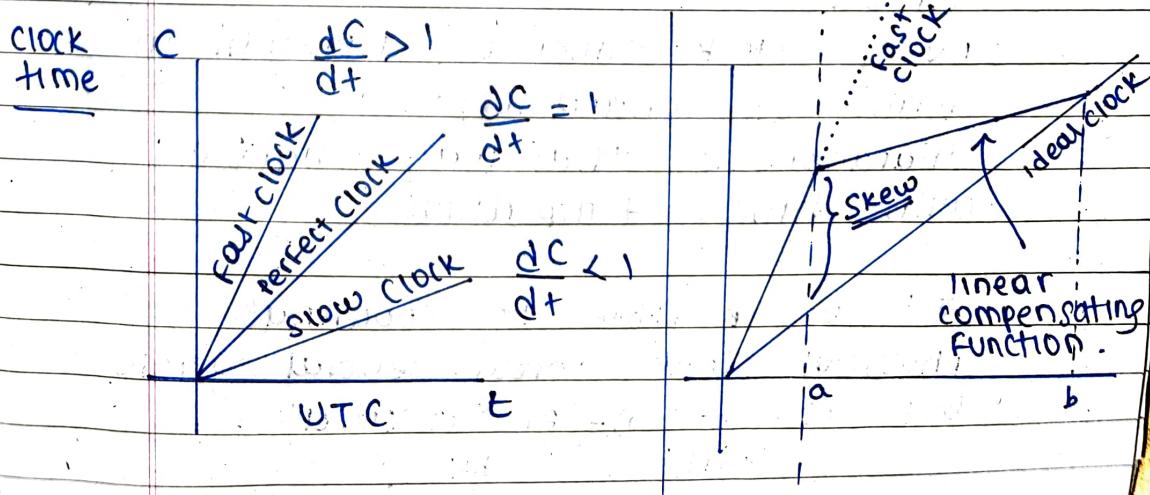
NOTE :

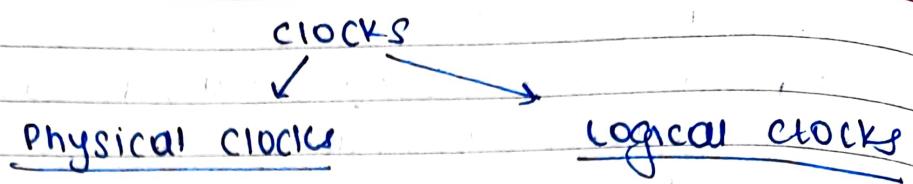
In clock synchronisation, clocks should never be adjusted backward as this can disrupt event ordering and system consistency.

Time should ALWAYS move forward

Instead, if a clock is ahead of the correct time, it should slow down slightly so that it gradually synchronises with correct time, without jumping backward.

Clock synchronisation algorithms are designed in such a way that the fast clocks are readjusted slowly instead of readjusting to actual time at once





In synchronisation, there are mainly two types of clocks.

① Physical clock :

- Time isn't a big issue in traditional centralised systems, where one or more CPU share the a common bus.
The entire system shares the same understanding of time, right or wrong, it is consistent.
- In distributed systems, that is not the case. Every system has its own timer that keeps the clock running.
These clocks are based on the oscillation of piezoelectric crystal or ~~sofa~~ similar integrated circuit.
Every clock timer is different in terms of characteristics - characteristics that may change with time, temperature
- It is possible to coordinate physical clock across several systems but it will never be accurate.

② logical clock

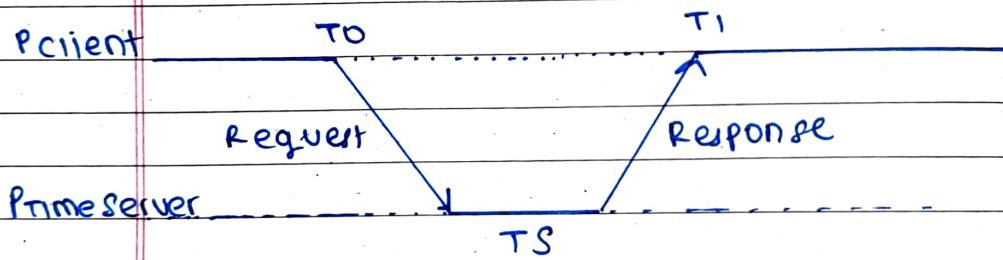
- logical clock means creating a protocol on computers in a distributed system so that computers can keep a uniform ordering of happenings inside some virtual time range.
- In a DS, a logical clock is a technique for recording temporal and causative links. Because distributed systems may lack a physically synchronized global clock, a logical clock provides for the global ordering of occurrences from various processes in certain systems.

Clock Synchronisation Algorithms

(I) Christian Algorithm

1. It is a synchronous centralized clock synchronisation algorithm used to synchronise time with time server by client processes.
2. This algorithm works well with a low latency network where the round trip time is short compared to accuracy.

Round trip time is the time start of request and end of corresponding response.



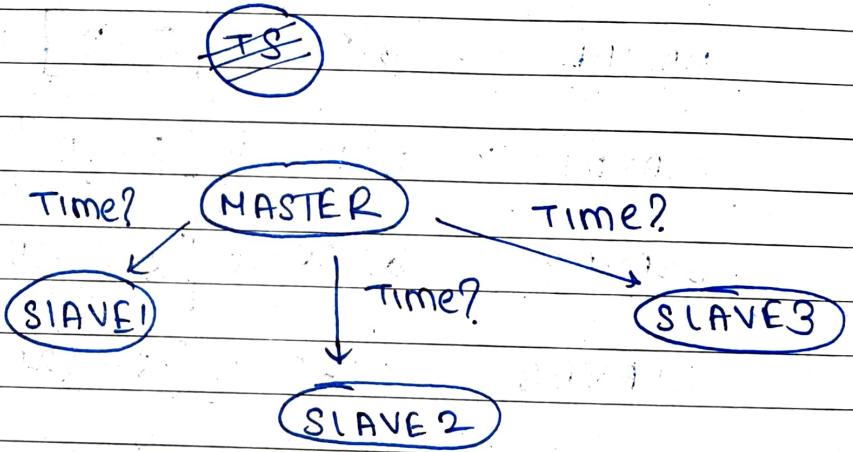
$$T_{\text{new}} = TS + \left[\frac{T_1 - T_0}{2} \right] \quad (1)$$

3. The client sends a request to a time server for its current value of UTC time (TS)
4. The client also records the time at which the request was submitted (T_0) and time it got the response (T_1)

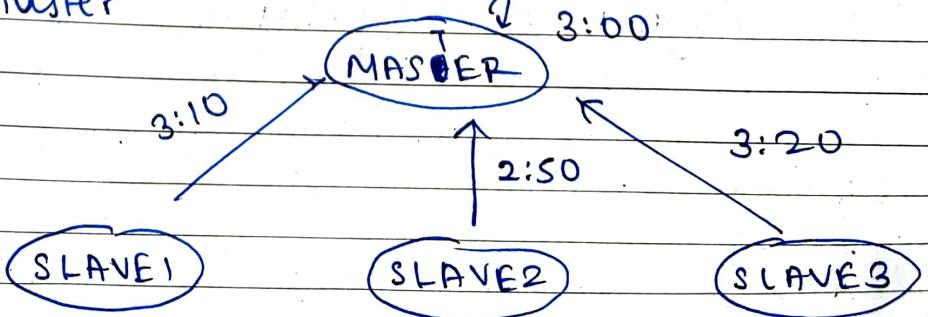
5. Then client changes the current time at T_1 with value received from T_S plus estimate of delay obtaining T_1 value.

(II) Berkeley Algorithm

1. It is a centralised clock synchronisation mechanism where in DS that implies that no computer has precise time timing source.
2. It is an example of active time server approach. The time server periodically sends a message to all computers in a group.



3. UPON receiving the request each slave node reads its local time and sends back current time to master



4. Master node calculates the time difference b/w all clock times received and clock time given by master system's clock ~~and~~

$$\text{Offset} = \frac{0 + (10) + (-10) + (20)}{4}$$

$$\text{Offset} = 5$$

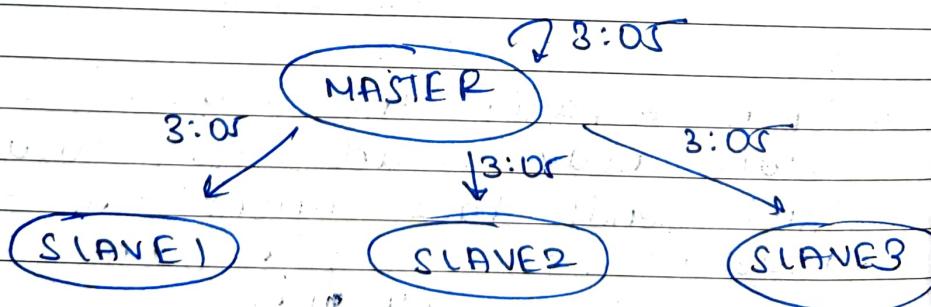
5. This average time difference is added to master system clock and broadcasted over the network.

$$\text{MASTER} : 3:00 + [5 - 0] = 3:05$$

$$\text{SLAVE1} : 3:10 + [5 - 10] = 3:05$$

$$\text{SLAVE2} : 2:50 + [5 - (-10)] = 3:05$$

$$\text{SLAVE3} : 3:20 + [5 - (20)] = 3:05$$

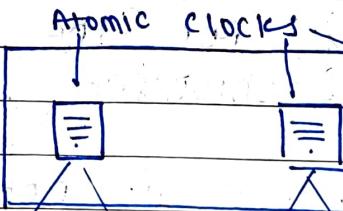


III) Network Time protocol (NTP)

- 1 - It is the standard followed by synchronisation clocks on internet.
It is a decentralised algorithm.
2. The main servers are directly linked to a precise and dependable UTC time source. They are the foundations of hierarchical time service with additional servers becoming operational as we go away from roots.

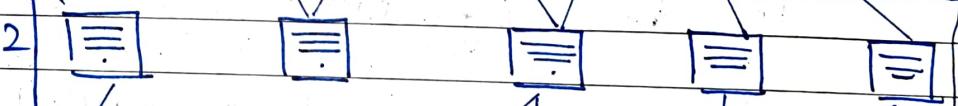
3.

Stratum 1

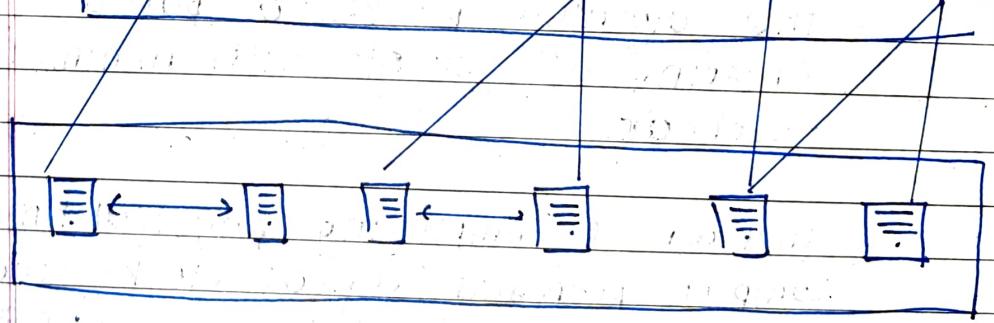


Stratum 1

Stratum 2



Stratum 3



4. The common configuration is :

Stratum 1: UTC servers at big governmental institutions

Stratum 2: Institutional time servers or ISP time server

Stratum 3: Most users linking to academic time servers

NTP may synchronise computers
in three modes

- { client - server mode / unicast
- multicast mode / multicast
- symmetrical (peer) mode / broadcast

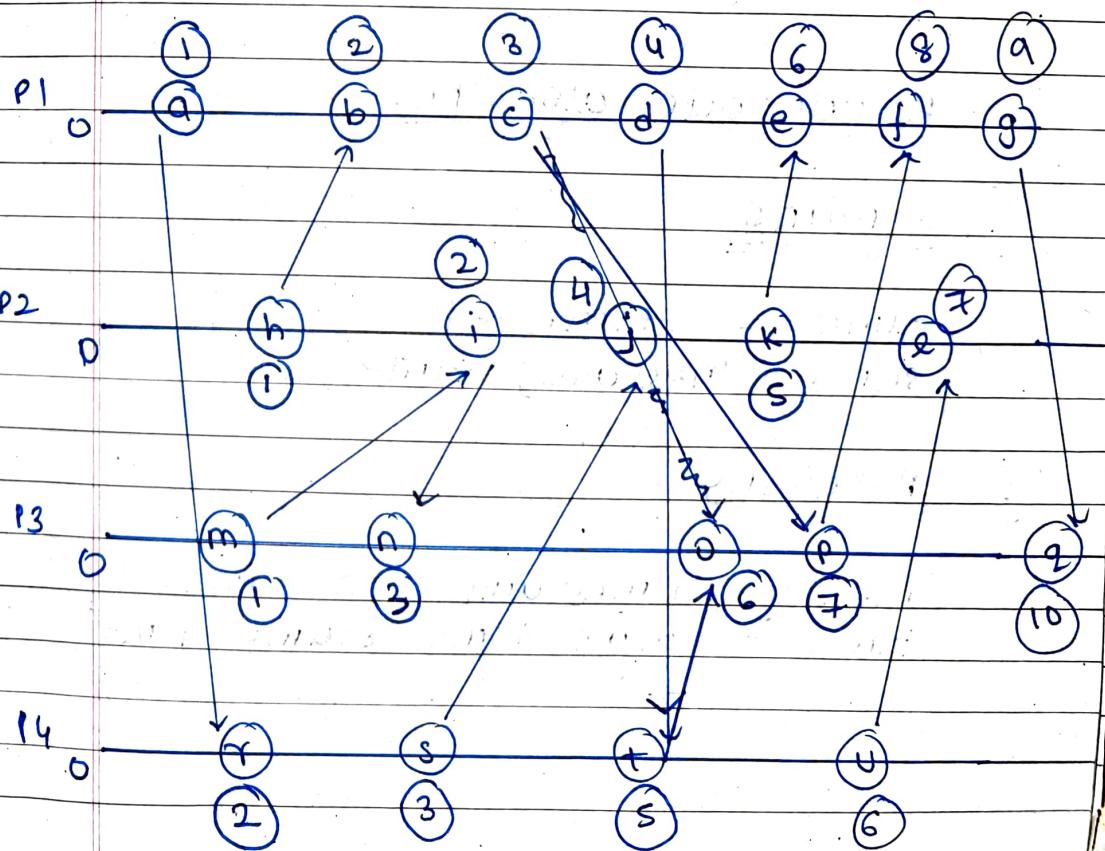


Lamport's logical clock

- In distributed systems, it is not necessary for the clocks to be synchronised.
- If two processes do not interact with each other, it is not necessary for their clocks to be synchronised.
- This algorithm provides a partial ordering of events with minimum overhead.
- Instead of employing physical time, Lamport proposed logical clocks that capture events' orderings through a 'happens-before' relationship.

or process

- Each machine has its own clock and the clock is just a 'counter' that increases every time an event happens.
- It is initially set to zero.
- The counter increments by 1 for each event it experiences.
- When a 'process' sends a message to another process, it includes its current clock value with message.
- When a process receives message, it updates its clock to be the maximum of its own clock and received clock value and increments by 1.



1	2	3	4	5	6	7	8	9	10
a	b	c	d		e	f	g		
h	i		j	k	l				
m	n			o	p				q
-	r	s		t	u				

conditions of Lamport logical clock

for any two events a and b

if $a \rightarrow b$

then $c(a) < c(b)$

But converse is not true. always

Pseudocode algorithm:

(1) sending

time = time + 1

send (message, time)

(2) Receiving

message, timestamp = receive()

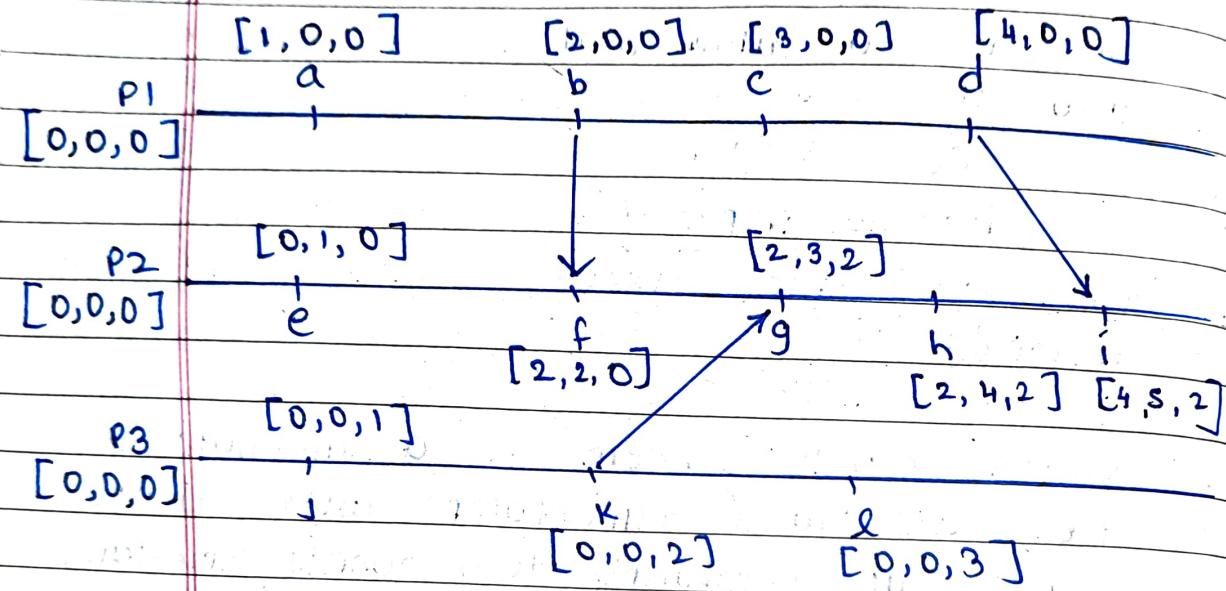
time = max(timestamp, time) + 1



vector clocks

- 1.- In Lamport's clocks, one cannot directly compare time stamps of two processes and determine their precedence relationship.
i.e. If a and b are two events and if $c(a) < c(b)$
 $\nrightarrow a \rightarrow b$
[does not imply]
- To solve this issue, vector clocks were designed
- 2- Vector clock system is a mechanism that associates timestamps with events such that comparing two events timestamps indicate whether those two events are causally related.
- 3- Vector clock is represented as a vector $[c_1, c_2, c_3, \dots, c_n]$ with an integer clock value of each process.
[c_i contains clock value of process i]
- 4- When a process generates an event or sends a message, it increments its own entry in the vector clock and includes its current vector clock with an event or message.
- 5- Upon receiving an event or message,
 - the recipient process updates its own vector clock by taking maximum value from each corresponding entry of its own vector clock and received vector clock.

Additionally the recipient increments its own entry in the vector clock to account for the received event.



MUTUAL EXCLUSION

- Mutual exclusion is an important concept in concurrent programming which ensures that simultaneous operations or node locations do not use common assets or crucial areas concurrently.

Thus it prevents conflicting and inconsistent results.

- ON A SINGLE COMPUTER, memory and other resources are shared by different processes. Thus status of shared resources and status of users is easily available. So with help of shared variable, mutual exclusion can be solved easily.

- IN DISTRIBUTED ENVIRONMENT

We have neither shared memory nor common physical clock and therefore mutual exclusion problem is difficult in DS.

- To eliminate the mutual exclusion problem in DS, new approach based on message passing is used.

Mutual exclusion 'mutex'

Requirements of Mutual exclusion Algorithm.

(1) No Deadlocks: Two or more site should not wait endlessly for any message that will never arrive.

(2) No Starvation:

Each site who wants to execute critical section should get an opportunity to execute it in finite time.

Any site should not execute critical section indefinitely.

Any site should not WAIT indefinitely to execute critical section while other sites are repeatedly executing critical sections.

(3) Fairness:

Each site should get a fair chance to execute critical section.

Any request made to critical section must be executed in the order they are made.

(4) Fault tolerant

In case of failure, it should be able to recognise it by itself in

order to function correctly without disruption

Performance metrics

- ① - synchronisation delay

previous site leaves
the critical section

following site enters
critical sections

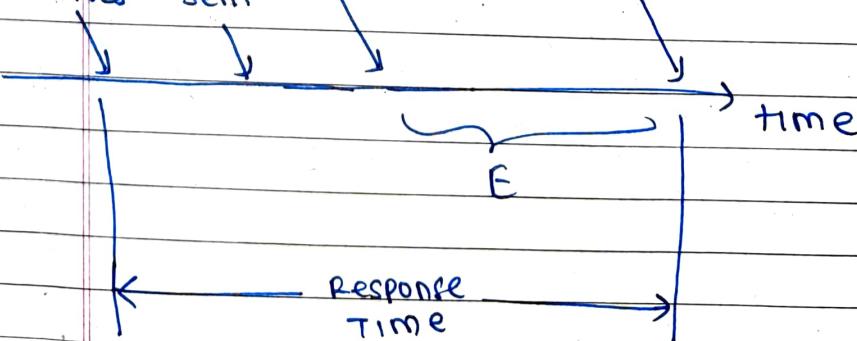


- ② - system throughput: rate at which request for critical section gets executed.

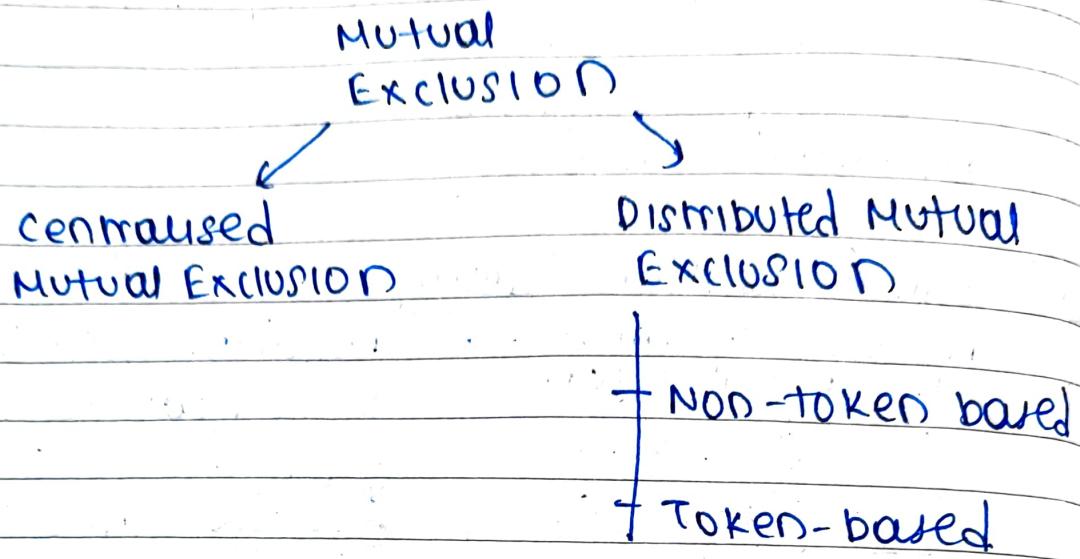
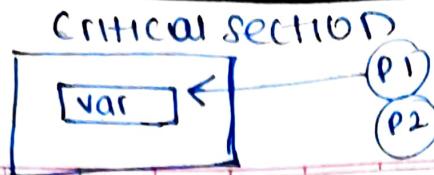
- ③ - Message complexity: Number of processes that where required by CS execution by a process.

- ④ - Response time: - Time interval from a request send to CS and execution completion.

CS req arrives Msg sent Enters CS Exits CS



one
process
at
a time



① Centralised Mutual Exclusion Algorithm.

1. In centralised algorithm, one process is elected as a coordinator which maybe the machine with highest network address.
2. The coordinator coordinates all processes and grant the permission to a process to enter in the critical section (CS) to use the shared resources.

3. When a process wants to enter the critical region, it sends a request message to the coordinator stating which critical region it wants to enter and asking for permission.
4. If no other process is currently in that critical region, the coordinator sends back a reply grant granting permission.

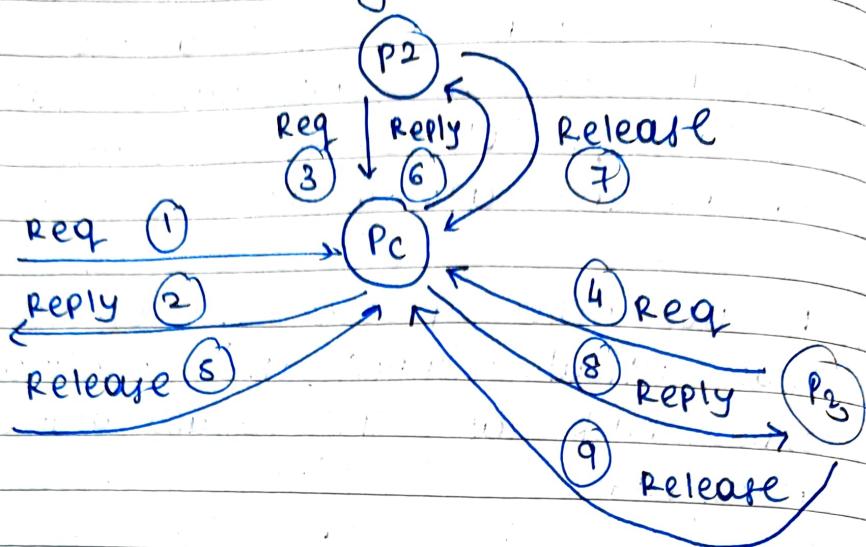
When the reply arrives, the requesting process enters critical region.

When the process exits the critical region, it sends a message to the coordinator to release its exclusive access.

If a process wants to access the same section of critical region that is currently occupied, then coordinator refrains from replying or replies with permission denied message.

If multiple processes concurrently seek permission to enter CS at some time, then coordinator grants permission to only one process to enter CS based on scheduling algorithms (FCFS).

Given following scenario



Queue

(1)

state



Initial status

(2)



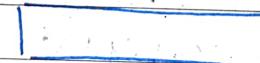
(3)



(4)



(5)



(6)



(7)



(8)



(9)



[Suzuki-Kasami]

Token based

Distributed Algorithm { - Quorum based [Nakayama]

Non token based

[Lamport, Ricart-Agarwala]
[Nakayama]

① Lamport's Algorithm for mutual exclusion in distributed systems.

- Lamport's Algorithm is a permission based algorithm proposed by Lamport as an illustration for the synchronisation scheme for distributed systems.
- In permission based, timestamp is used to order critical section requests and to resolve any conflict b/w requests.
- In this algorithm, critical section requests are executed in increasing order of timestamp. \Rightarrow FCFS.

① - When a process wants to enter critical section, it sends a time stamped request to all other processes in system.

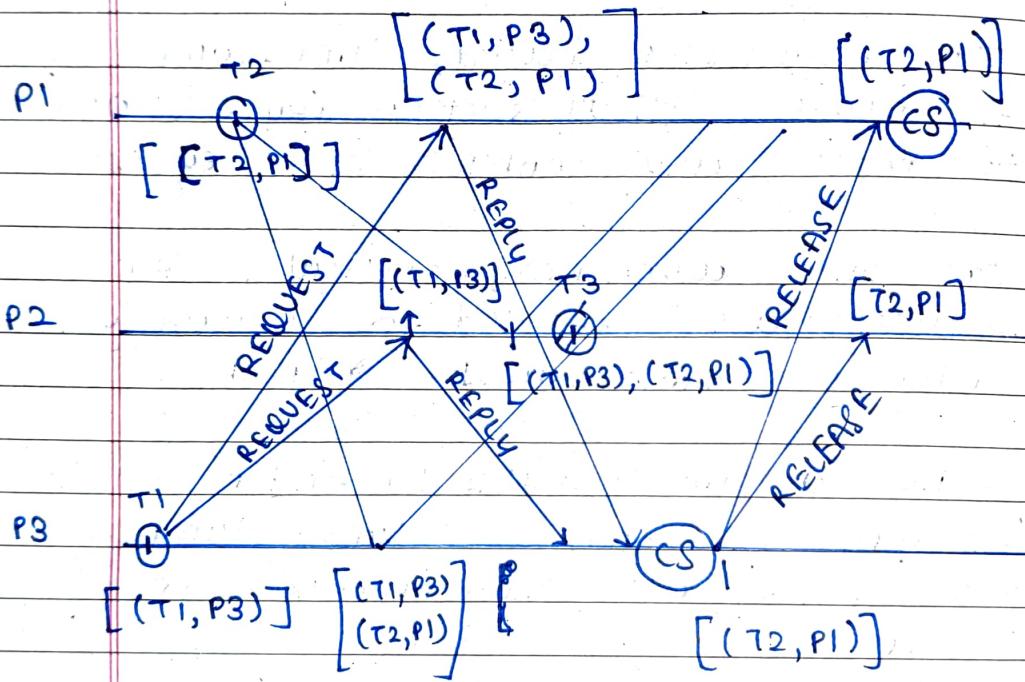
② - Upon receiving the request, each process sends a reply message with to requesting process.

③ - The requesting process can only enter critical section only after it has received messages from all processes and its own request has earliest time stamp among all outstanding request.

4) After executing critical section, the process sends a release message to all other processes informing them that the shared resource is now available.

— Types of messages in Lamport's Algo

{ REQUEST
REPLY
RELEASE }



(Timestamp, process-ID)

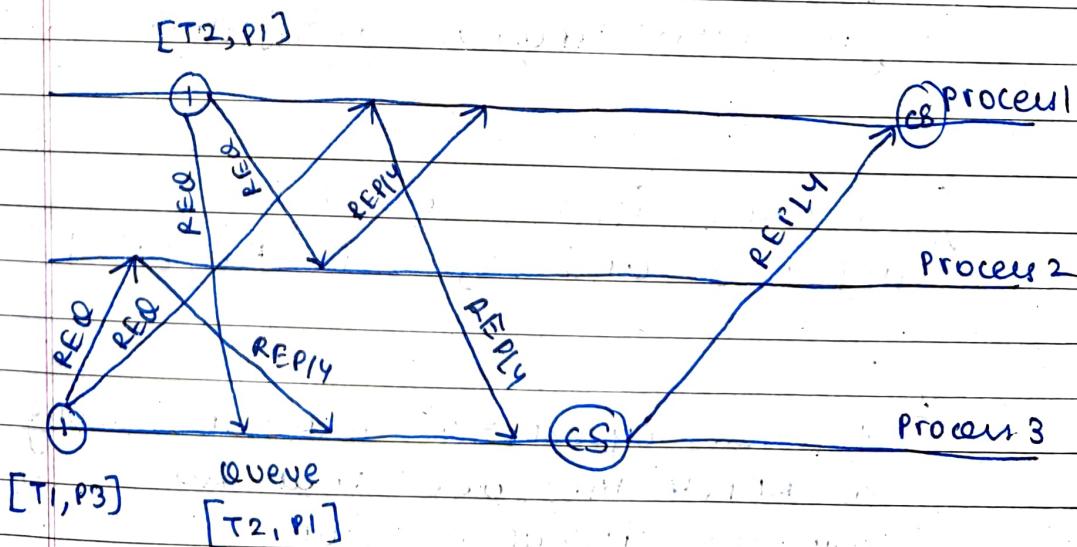
* — Performance

For each CS execution, Lamport's algo requires $(N-1)$ Req, $(N-1)$ reply, and $(N-1)$ release messages.

→ Lamport's algorithm requires $3(N-1)$ msg per CS invocation.

(2) Ricart - Agarwala Algorithm.

- Ricart - Agarwala algorithm is an optimization on Lamport algorithm.
- Ricart - Agarwala algorithm uses only two types of message: REQUEST and REPLY
- It is assumed that all processes keep a logical clock which is updated according to the clock rules
- The algorithm requires a total ordering of requests. Requests are ordered according to their global timestamps.
IF timestamps are equal, process identifiers are compared to their global timestamps
IF time stamps are equal, process identifiers are compared to other order them.



- A process sends a REQUEST message to all other processes to request their permission to enter critical section.
- A process sends a REPLY message to a process to give its permission to enter CS.

NOTE:

when process s_i requests process s_j

s_j sends a reply message
IF and only IF:

s_j is neither requesting nor
executing CS

OR

IF s_j is requesting and s_i 's
request timestamp is smaller than
 s_j 's own request timestamp

{ otherwise the request is
deferred and s_j sets $RP_j[i] = 1$

deferred request \rightarrow delay in sending
a REPLY message that has
requested critical section.

site S_i enters the CS after it has received a REPLY message from every mess. site it sent a REQUEST message to.

when site S_i exist exit the CS, it sends all the deferred REPLY messages

i.e $\forall j \text{ IF } RD_{ij} = 1 \text{ then }$

send Reply message to S_j and set $RD_{ij} = 0$

performance

for each CS execution, Ricart Agarwala algorithm requires $(N-1)$ request messages and $(N-1)$ reply messages.

thus, it requires $2(N-1)$ msgs per CS execution.

* QUORUM BASED MUTUAL EXCLUSION ALGORITHM

- Quorum based algorithm are different in two ways

- ① A site does not request permission from all sites, but only from a subset of the sites.

The request set of sites are chosen such that

$$\forall i \neq j : 1 \leq i, j \leq N$$

$$R_i \cap R_j \neq \emptyset$$

Consequently each pair of sites has a site that mediates conflicts b/w that pair.

- ② A site can send out only one REPLY message at any time. A site can only send a reply message only after it has received a RELEASE message for previous REPLY message

the

Naekawais Algorithm

- Naekawais Algorithm was the first quorum based mutual exclusion algorithm
- In this algorithm, a site does not request from every other site but from a subset of sites which are called quorum
- The request sets for sites i.e quorum in this algorithm are constructed as follows:

① $\forall i \forall j : 1 \leq i, j \leq N, i \neq j :: R_i \cap R_j \neq \emptyset$
i.e R_i and R_j must have overlap.
i.e there is atleast one common site
bw the request set of any two sites.

② $\forall i : 1 \leq i \leq N :: s_i \in R_i$

③ $\forall i : 1 \leq i \leq N :: |R_i| = k$

Every process s_i must be a member of its own quorum set R_i

This states that the size of each quorum R_i is exactly k where k is a constant.

④ $N = k(k-1) + 1$, and $|R_i| = \sqrt{N}$

where : $N \rightarrow$ total no. of processes in system

k is size of each quorum ($|R_i|$)

Tirthay Mahajan

$N = K(K-1) + 1$ ensures that each quorum R_i overlaps with all other quorums.

SUPPOSE : $K = 3$

$$N = 3(3-1) + 1 = 7$$

\therefore For 7 processes, quorum size = 3

$$|R_i| = \sqrt{N}$$

This ensures that size of each quorum is proportional to the total no. of processes.

If N is perfect square, then

$K = \sqrt{N}$ making quorum construction symmetric and balanced.

Algorithm

A site s_i executes the following steps to execute the CS.

①

To enter critical section

- when a site s_i wants to enter the critical section, it sends a REQUEST message to all other sites in request set R_i

- when a site s_j receives a REQUEST(i) message -

condition 1: If s_j has not already sent a reply message since the last RELEASE message, it immediately responds to s_i with a reply message!

Condition 2: If s_j has already already sent a reply message since the last release message

Condition 2: If s_j has already sent a reply message (e.g. another site s_k), s_j cannot send a reply msg immediately.

Instead it queues up the request from s_i .

(2) To execute the critical section:

A site s_i can enter the critical section IF it has received the REPLY message from all the sites in request set R_i

(3) To release the critical section:

- When a site s_i exits the CS, it sends a RELEASE(i) message to all the sites in request set R_i
- When a site s_j receives the RELEASE(i) message from site s_i , it sends a REPLY message to the next site waiting in queue and deletes the entry from queue

- in case queue is empty, site s_j update its status to show that it has not sent any REPLY message since the receipt of the last RELEASE message.

performance :

Maekawa's algorithm required invocation of $3IRI \approx 3\sqrt{N}$ messages per critical section execution as the size of a request set is \sqrt{N} .

$\sqrt{N} \rightarrow REQUEST$

$\sqrt{N} \rightarrow REPLY$

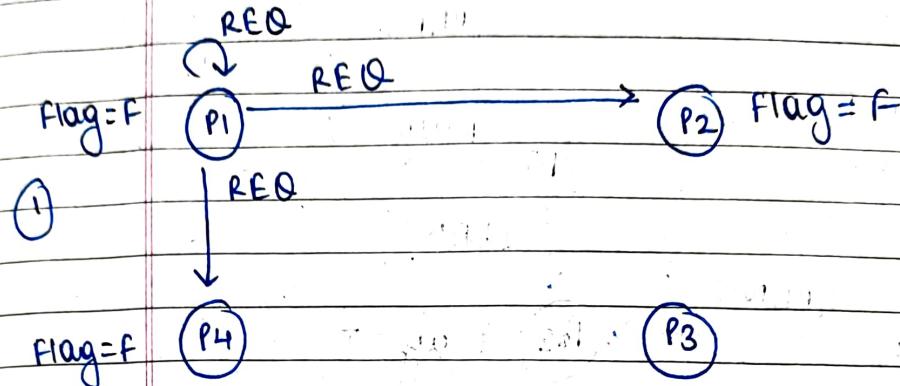
$\sqrt{N} \rightarrow RELEASE$.

NOTE

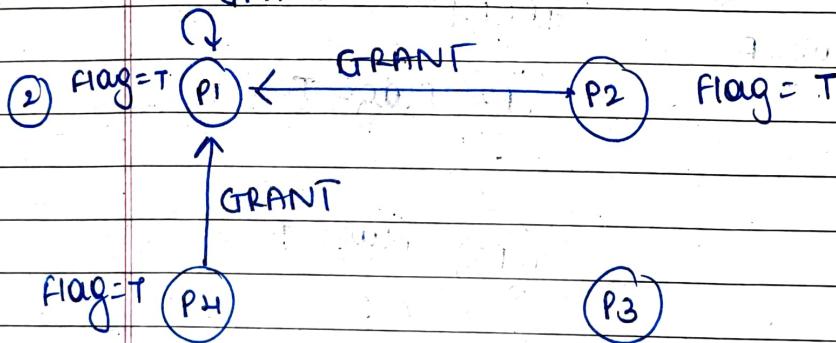
Maekawa's algorithm can deadlock because site is exclusively locked by other sites and requests are not prioritized by their timestamps.

(P1)

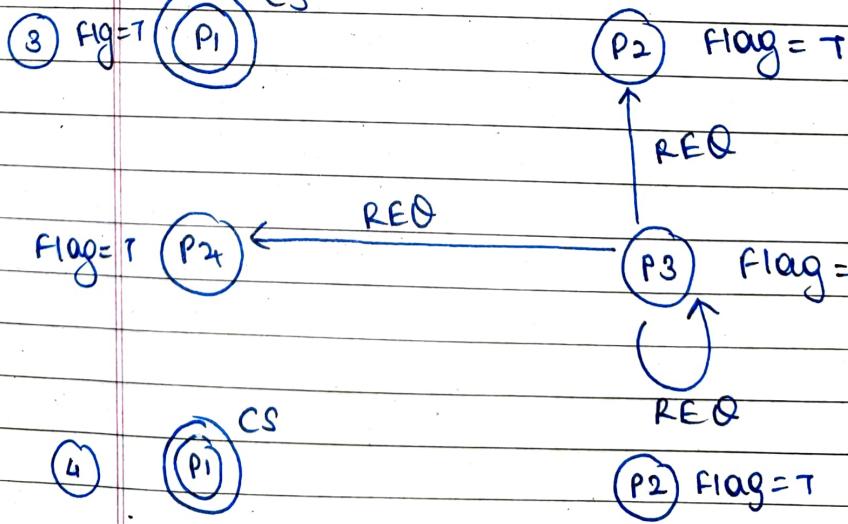
P1 → P1 P2 P4
P2 → P1 P2 P3
P3 → P2 P3 P4
P4 → P1 P3 P4



GRANT



CS



④

Flag=T

CS

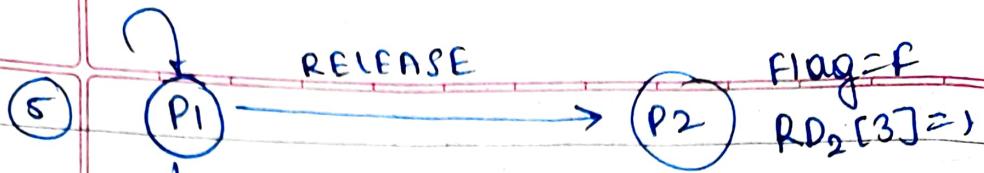
Flag=T

RD₄[3]=1

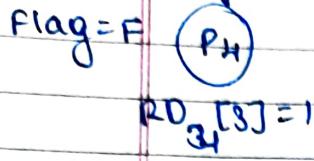
Flag=F

GRANT

RELEASE

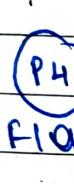


RELEASE



GRANT

Flag=F



Flag=T



Flag=T

GRANT



Flag=T

GRANT



Flag=F



Flag=T



Flag=T



Flag=T

C.S.

Token based algorithms:

- In token based algorithms, a token is shared among the sites.
- A site is allowed to enter its CS if it possesses the token.
- Token based algorithms use sequence nos instead of timestamps to distinguish b/w old and current req.
- Every site requesting to get a token contains a sequence no.
- On every token request, the sequence number is incremented by 1.

} TOKEN-RING ALGO
TOKEN BASED ALGORITHMS }
SUTUKI KARAMI
ALGORITHM.

① TOKEN RING algorithm

- TOKEN based
- TOKEN-RING algorithm is an example of token-based distributed mutual exclusion algorithm.
- It is based on concept of a circulating ring in a logical ring topology.

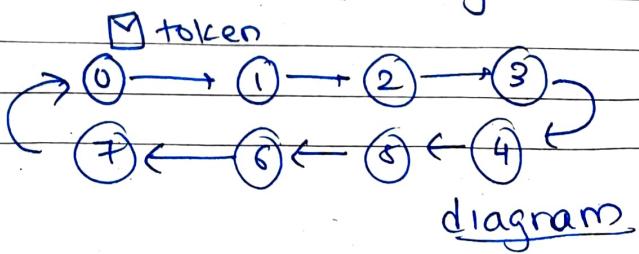
- * - Token creation and initial distribution.
- 1. - Initially there is only one token in the system. Each node is allowed to hold the token for a fixed amount of time. Once the node is done with the operation or does not need the token, it passes on to the neighbour ring.
- 2. - Suppose N processes, token is passed from P_k to $P_{(k+1) \bmod N}$ in point to point message.

* - critical section execution

- 1 - When a process acquires the token from its neighbour, it checks if it needs to access the shared resources.
- 2 - If so, the process goes ahead and enters CS.
- 3 - After it is finished, it exits the CS and passes the token along the ring.
- 4 - It is not permitted to immediately enter CS again for different event ^{using} FOR the same token.

* - TOKEN PASSING

1. - Depending on the order and distribution of token, certain processes might get more opportunity to access CS than others faster.
To ensure fairness, various strategies can be employed such as
2. - However, as the token circulates throughout the ring, it ensures each process gets a chance to enter CS.
3. - If the node holding the token crashes, we can detect the token disappearance by calculating the maximum time to complete one round.
Suppose there are N nodes, each node can hold token for maximum ' t ' time, total completion time = $NT_{(\max)}$
4. - Thus if the token does not return to a node by that time limit, we can consider token as lost.
5. - Recreation of token has to be done carefully so that only one machine ends up creating new token and we don't end up creating multiple tokens in ring.
6. - One strategy is to elect a leader who is responsible for detecting token loss and regenerating token.



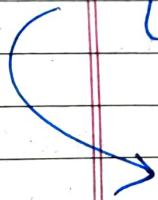
diagram

Election Algorithms

- 1- Election algorithms are used to elect a leader or coordinator among a group of distributed nodes.
- 2- The leader is responsible for creating, cooperating and managing the activities of the distributed system to ensure proper functioning and efficient communication.
- 3- Distributed election algorithms are essential for maintaining fault tolerance, load balancing and overall system reliability.
- 4- The choice of algorithm depends on the specific requirements of system and desired tradeoffs between simplicity, performance and fault tolerance.

common algos

{ RING
BULLY ..



learn from
GitHub

GOSSIP BASED COORDINATION

- 1 - Gossip based coordination is a decentralized approach used in DS for achieving coordination and information dissemination among nodes.
- 2 - Inspired from real world gossip spreads where information is passed from one person to another in random and probabilistic manner.
- 3 - In gossip based coordination system, each node communicates with few randomly selected neighbouring neighbour nodes at regular intervals of time.
- 4 - During these interactions, nodes exchange information they possess and this information spreads throughout nw gradually.
- 5 - The protocol continues till desired level of coordination or consensus is achieved among majority of nodes.
- 6 - It is important to note that while gossip protocols are scalable and fault tolerant, it may take time for message to propagate through the nw.
↳ eventual consistency rather than immediate consistency.

I Aggregation

- 1 - Aggregation is a gossip based coordination process in which we collect and combine information from multiple nodes to generate a summary or a consolidated view of data.
- 2 - Each node communicates only with small, randomly chosen subset of nodes in each round and they exchange local state information and converge them toward desired aggregate value.
- 3 - Each node maintains a partial ~~result~~ local state representing partial result.
- 4 - Each node maintains two different threads - active and passive thread.
- 5 - Active thread periodically sends the local state with random neighbour.
- 6 - Passive thread waits for messages by other nodes, and as soon as it receives a message, it will send back a reply containing its own local state.
- 7 - Both nodes update their local state on the aggregation function.
- 8 - Over successive rounds, the states across

nodes converge to a global aggregate

• Only need to do a single pass through the entire network

• Follows the same general steps as the local aggregate

• Local aggregate finds the best path to each destination

• Global aggregate finds the best path to each destination

• Local aggregate finds the best path to each destination

• Global aggregate finds the best path to each destination

• Local aggregate finds the best path to each destination

• Global aggregate finds the best path to each destination

• Local aggregate finds the best path to each destination

• Global aggregate finds the best path to each destination

• Local aggregate finds the best path to each destination

• Global aggregate finds the best path to each destination

• Local aggregate finds the best path to each destination

• Global aggregate finds the best path to each destination