

$$b \leq \theta \leq a.$$

$$b = \frac{a+b}{2} - \frac{a-b}{2} \leq \theta \leq \frac{a+b}{2} + \frac{a-b}{2} = a.$$

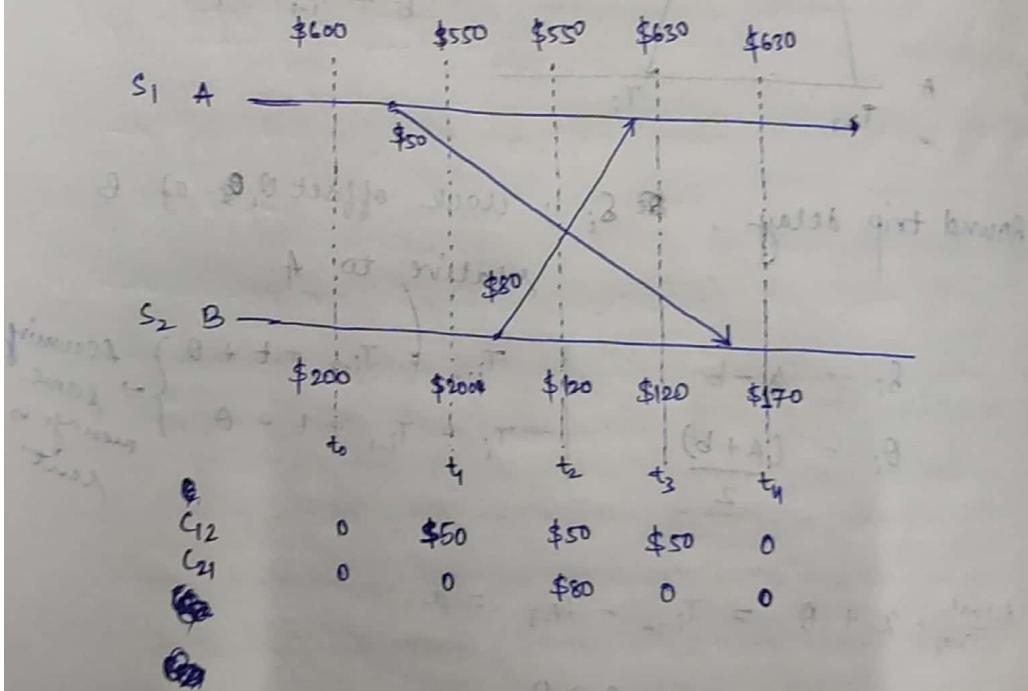
$$\theta_i - \frac{d_i}{2} \leq \theta \leq \theta_i + \frac{d_i}{2}$$

↳ Offset can be $-\frac{d_i}{2}$ or $\frac{d_i}{2}$.

Global State and Snapshot Recording

Algorithm

→ Recording global state is necessary for analyzing, testing or verifying properties associated with distributed execution.



‡ Determination of global states is non-trivial.

→ System Model

→ Transit

$$\hookrightarrow \text{Definition of } \xrightarrow{\text{transit}} (LS_i, LS_j) = \{m_{ij} \mid \begin{array}{l} \text{send}(m_{ij}) \in LS_i \\ \text{rec}(m_{ij}) \in LS_j \end{array}\}$$

→ A consistent global state,

$$GS = \{U, LS_i, V, SC_{ij}\}$$

if the following 2 conditions are satisfied:

i) $C_1: \text{send}(m_{ij}) \in LS_i \Rightarrow m_{ij} \in SC_{ij} \oplus \text{rec}(m_{ij}) \in LS_j$

(XOR)

ii) $C_2: \text{send}(m_{ij}) \notin LS_i \Rightarrow m_{ij} \notin SC_{ij} \wedge \text{rec}(m_{ij}) \in LS_j$

↳ C_1 : Indicates the law of conservation of messages.

↳ C_2 : States that, ~~unless there is a cause~~ for every effect, its cause must be present.

→ Issues in Recording a Global State

↳ I₁: How to distinguish between messages to be recorded in the snapshot (i.e., either in the channel state or a process state) from those not to be recorded.

Any message that is sent by a process before recording its snapshot must be recorded in the global state (from C₁).

Any message that is sent by a process after recording its snapshot must not be recorded in the global snapshot (from C₂).

→ I2: How to determine the instant when process takes its snapshot.

A process P_j must record its snapshot before processing a message that was sent by a process P_i after recording its snapshot.

→ 2 types of messages: computation messages and control messages.

→ Execution of snapshot algorithm is transparent to the underlying application, except for occasional delaying of some of the actions of the application.

Snapshot Algorithm for FIFO channels

1. Chandy - Lamport Algorithm

→ (control) message in this case is referred to as a marker.

→ After a site has recorded its state, it sends a marker along all of its outgoing channels, before sending out any ~~more~~ more messages.

Since the channels are FIFO, a marker separates the messages in the channel into those to be included in the snapshot from those not to be recorded in the snapshot.

→ Marker Sending Rule for process P_i :

i) Process P_i records its state, and

ii) for each outgoing channel c , on which a marker has not been sent P_i sends a marker

along c before p_i sends further message along c .

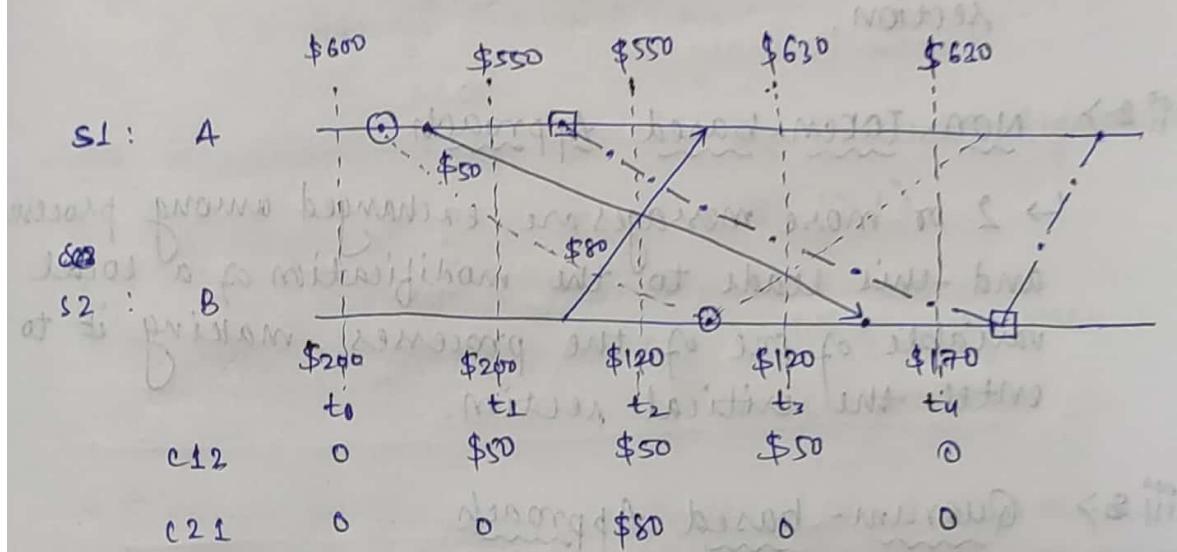
→ Marker Receiving Rule for process p_j

is on receiving marker along channel c :

↳ If p_j has not recorded its state, then record the state of c as the empty set.

• Execute the marker send up rule.

↳ Else, record the state of c as a set of messages received along c after p_j 's state was recorded and before p_j receive the marker along c .



□-- : $A = \$550 ; B = \$170 ; c_{12} = \$0 ; c_{21} = \80

○-- : $A = \$600 ; B = \$120 ; c_{12} = \$0 ; c_{21} = \80

These are consistent global states.

Distributed Mutual Exclusion Algorithm

↳ Sleeping Barber problem; Baker's algorithm; Dining Philosophers' problem; Dijkstra's algorithm; semaphores

→ Message passing is the sole means of implementing critical section system.

→ Three approaches:

i) Token-based Approach:

- ↳ Token is a privileged message
- ↳ The token is allowed to enter the critical section.

ii) Non Token-based Approach:

- ↳ 2 or more messages are exchanged among processes and this leads to the modification of a local variable of one of the processes, making it to enter the critical section.

iii) Quorum-based Approach

- ↳ A process will make a subset of processes (referred to as quorum), and not all the processes to enter the critical section.

- ↳ Effort required will be less than the above 2 approaches.

* In approaches (i) & (ii), the process wanting to enter the critical section makes request to all the processes.

→ System Model.

↳ A process can be in 3 states:

i) Requesting the critical section

ii) Executing the critical section

iii) Neither requesting nor executing the critical section (idle).

†. Time stamp is the measure for priority of requesting processes for critical section.
lower the timestamp, higher the priority.

Q: $N = \text{Total processes requesting for CS (critical section)}$
 $T = \text{Average Delay Time}$
 $E = \text{Average Execution Time}$

→ Requirements of Mutual Exclusion Algorithms (3 properties need to be satisfied)

↳ Safety Property: At any instant, only one process can execute the critical section.

↳ Liveness Property: Absence of deadlocks and starvation; there should be progress; No process should wait indefinitely to enter the critical section.

↳ Fairness Property: Each process gets a fair chance to execute the critical section.

→ Performance Metrics

↳ Message complexity

- No. of messages that are ~~needed~~ required per critical section execution by a site (assumption is that there is one process per site).

↳ synchronization Delay

↳ Response Time

- Time from request made by the process for the critical section till the execution is complete.

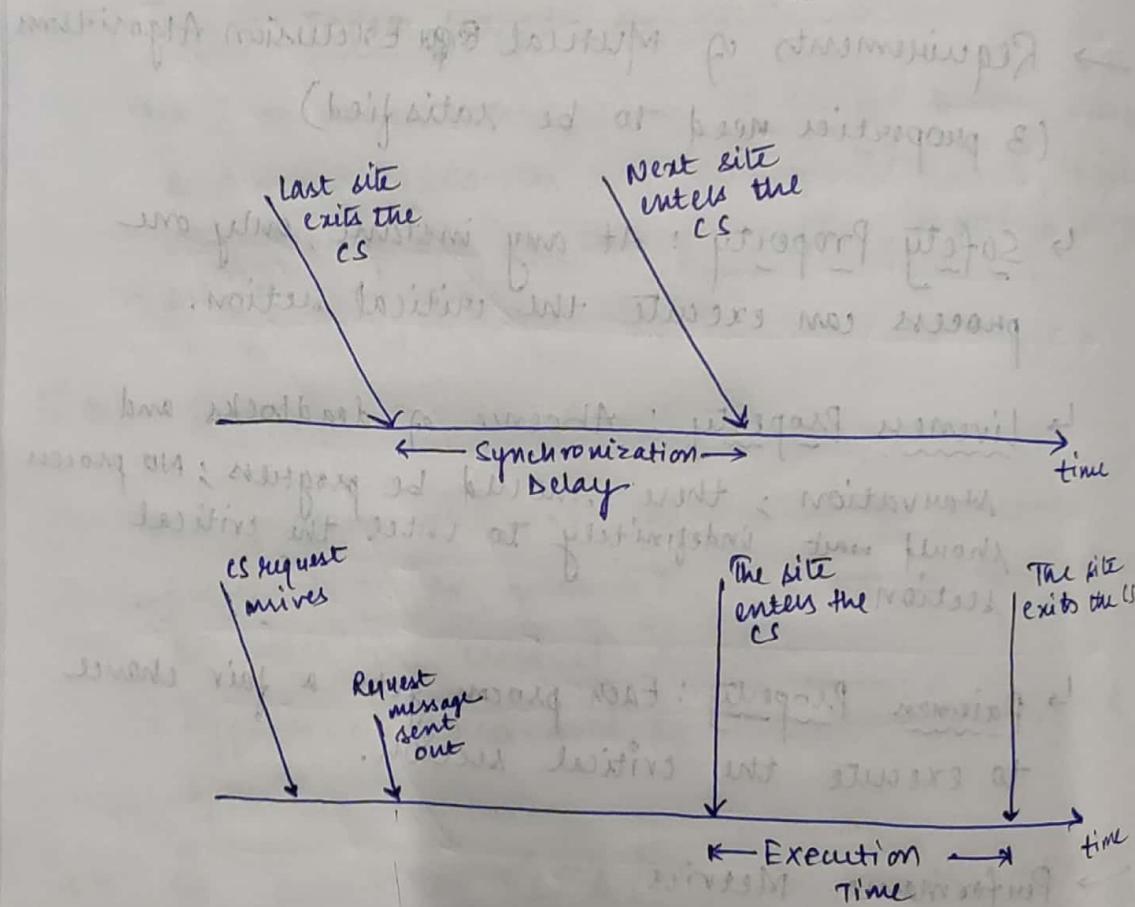
↳ System Throughput

- Rate at which system is able to execute the critical section.

SD : synchronization Delay

E : Average execution time

$$\text{System Throughput} = \frac{1}{SD+E}$$



- ↳ Low load and high load performance
 - Low load & high load mean how frequently the messages arrive.
- ↳ Best and worst case performance.
 - low load (less frequent messages), so best performance
 - In high load, processes may have to wait to get the message.

1. Lamport's Algorithm

→ Requests for CS are executed in order of their timestamp.

Time is defined by Lamport's logical clock.

i) → Requesting the critical section.

↳ When a site wants to enter the CS, it broadcasts a request message, i.e., (ts_i, i) message to all other sites and places the request on the request-queue_i. (Every process maintains a request queue).

↳ When a site s_j receives the REQUEST (ts_i, i) message from site s_i , it places site s_i 's request on its request-queue_j, and returns a timestamp REPLY message to s_i .

ii) → Executing the critical section.

↳ L1: s_i has received a message with timestamp larger than (ts_i, i) from all other sites.

↳ L2: s_i 's request is at the top of request-queue_i.

iii) → Releasing the Critical Section

↳ site s_i upon exiting the CS removes its request from the top of its queue and broadcasts a timestamp RELEASE message to all other sites.

↳ When a site s_j receives RELEASE message from site s_i , it removes s_i 's request from its queue.

→ $3(N-1)$ messages need to be sent for one CS exit of a process.

\Rightarrow

s_1

$(1,1)$

s_2

$(0,1,2)$

s_3

↳ site s_1 enters CS
↳ Execution complete, so site s_1 sends RELEASE message & removes $(1,1)$ from its queue

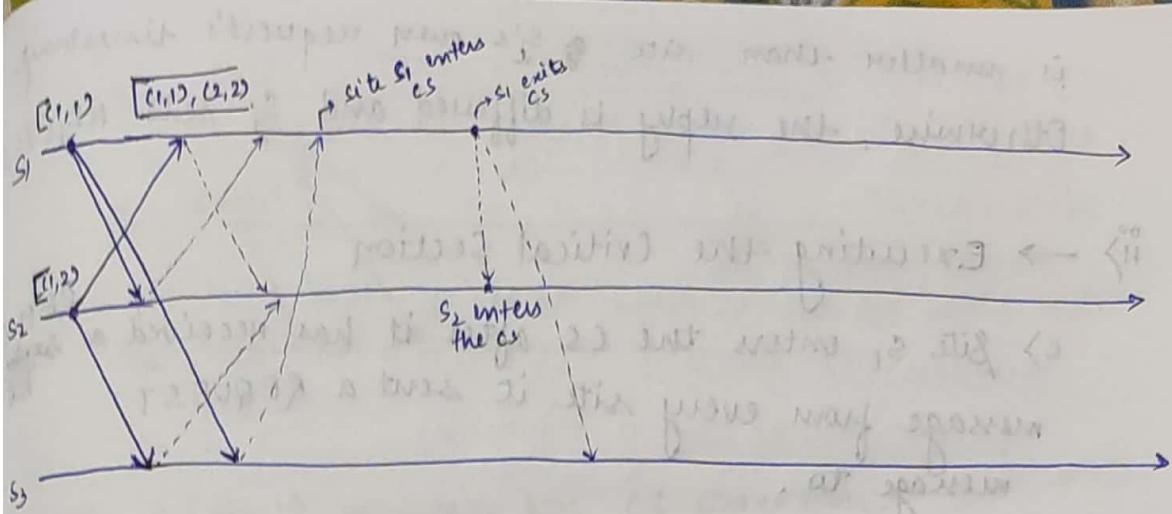
↳ site s_1 enters the CS

a: s_1 and s_2 are making requests for the CS.

b: site s_1 enters the CS

c: site s_1 exits the CS and sends RELEASE message.

(a, b, c are 3 separate diagrams, each containing the preceding one)



State S₂ enters the cs

→ Performance: For each cs execution, $3(N-1)$ messages for each execution.

→ Performance can be optimized from $3(N-1)$ to $2(N-1)$ by removing some of the reply messages.

→ Communication channel need to be FIFO.

2. Ricart - Agrawala Algorithm

→ REQUEST and REPLY messages.

→ Initially, $v_i, v_j : RD_i[j] = 0$.

RD_i : Request - differed array of site s_i .

→ Does not require communication channel to be FIFO.

i) → Requesting the critical section

a) When a site s_i wants to enter the CS, it broadcasts a timestamp ~~REQUEST~~ message to all other sites.

b) When a site s_j receives a REQUEST message from s_i , it sends a ~~REPLY~~ message to site s_i if site s_j is neither requesting nor executing the CS, or if the site s_j is requesting and s_i 's request timestamp

is smaller than site s_i 's own request's timestamp. Otherwise, the reply is differed and s_j sets $RD_{ij}[j] = 1$.

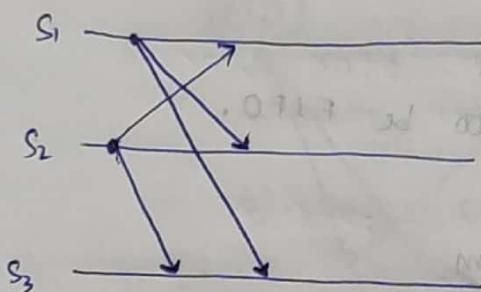
ii) → Executing the Critical Section

c) Site s_i enters the CS after it has received a REQUEST message from every site it sent a REQUEST message to.

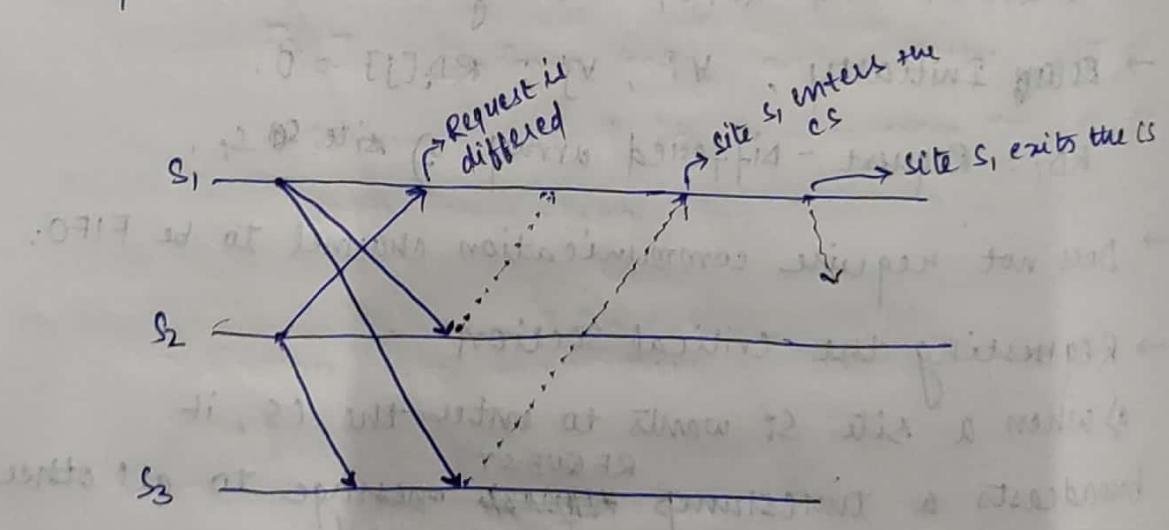
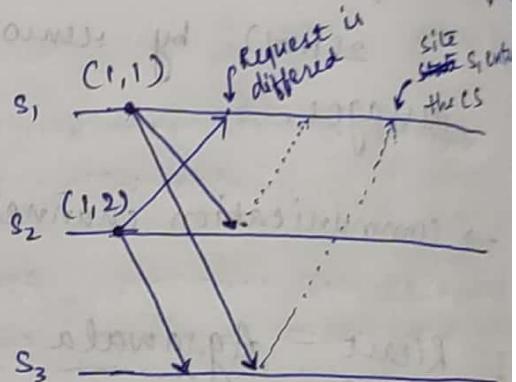
d) When site s_i exits the CS, it sends all the differed REPLY messages, i.e.

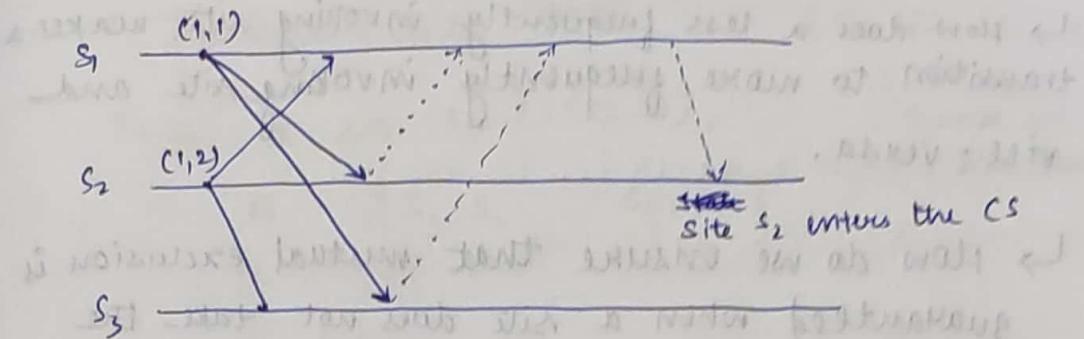
$$\forall j : RD_{ij}[j] = 1.$$

and, after sending the REPLY message, site s_i sets $RD_{ij}[j] = 0$.



Sites s_1 and s_2 make a request for the CS.





→ $2(N-1)$ messages per CS execution.

3. Singhal's Dynamic Information Structure Algorithm

→ Most mutual exclusion algorithms have a static approach.

→ Lack of efficiency: fail to exploit the changing condition in the system.

An algorithm can exploit changing condition to optimize the performance.

→ If a few sites are invoking mutual exclusion very frequently and other sites invoke mutual less frequently, then a frequently invoking site need not ask for the permission of less frequently invoking site everytime it requests an access for the CS.

→ Some of the design challenges:

↳ How does a site efficiently knows what sites are currently active, invoking mutual exclusion?

↳ When a less frequently invoking site needs to invoke mutual exclusion, how does it do it?

- ↳ How does a less frequently invoking site make a transition to more frequently invoking site and vice-versa.
- ↳ How do we ensure that mutual exclusion is guaranteed when a site does not take the permission of every other site?
- ↳ How do we ensure that a dynamic mutual exclusion algorithm does not waste resources and time in collecting the system state, offsetting any gain?

→ System Model

- ↳ Data structures are maintained by each process (they are basically sets):
 - i) R_i : Request set, contains the sites for which s_i must acquire permission before executing CS.
 - ii) I_i : Inform set, contains the sites to whom s_i must send its permission to execute after executing its CS.
- ↳ Every site maintains a logical clock c_i .
- ↳ Three boolean variables:
 - Requesting
 - Executing
 - My-Priority : will be true in case if requesting site s_i has a high priority

Initialization State

$\{s_n, s_{n-1}, \dots, s_i, \dots, s_2, s_1\}$

↳ For a site s_i ($i=1, \dots, N$),
Initial set is like this

$$R_i = \{s_1, s_2, \dots, s_{i-1}, s_i\}$$

$$I_i = \{s_i\}$$

$$C_i = 0$$

Requesting = Executing = false

The cardinality of R_i decreases in a step-wise manner from left to right.

This is referred to a staircase pattern based on the topological sense.

Step 1: (Request (cs))

↳ Requesting = true

$$\hookrightarrow C_i = C_i + 1$$

↳ send REQUEST (C_i, i) message to all sites in R_i

↳ Wait until ~~all sites in~~ $R_i = \emptyset$

(i.e., wait until all sites in R_i have sent a reply to s_i)

Step 2: (Executing (s))

↳ Executing = true

↳ Execute cs

↳ Executing = false

Step 3: (Release (s))

↳ for every site s_k in I_i (except s_i), do the following:
(on next page →)

begin

$$I_i = I_i - \{S_k\}$$

send a REPLY (C_i, i) message to S_k

$$R_i = R_i + \{S_k\}$$

end.

→ REQUEST Message Handler:

/* Site S_i is handling message REQUEST (C_j, j) */

$$C_i = \max(C_i, C_j)$$

↳ Case: Requesting = true.

begin

if My-Priority, then $I_i = I_i + \{S_j\}$

/* My-Priority is true if the pending request of S_i has priority over the incoming request

else

begin

send REPLY (C_i, i) message to S_j

if not S_j in R_i then,

begin

$$R_i = R_i + \{S_j\}$$

send REQUEST (C_i, i) message to S_j

end

end

end

↳ Case : Executing = true.

$$I_i^* = I_i^* + \{s_j\}$$

↳ Case : Executing = false \wedge Requesting = false.

begin

$$R_i := R_i + \{s_j\}$$

Send REPLY(c_i, i) message to s_j

end.

→ REPLY Message Handler:

* Site s_i is handling a message REPLY(c, j) */

begin

$$c_i = \max(c_i, c);$$

$$R_i := R_i - \{s_j\}$$

end.

e.g: 4 sites: $s_4 s_3 s_2 s_1$

$$\text{so, } R_1 = \{s_1\}; R_2 = \{s_2 s_1\}, R_3 = \{s_3 s_2 s_1\}$$

$$R_4 = \{s_4 s_3 s_2 s_1\}$$

NOW, s_3 requests.

$$\Rightarrow R_1 = \{s_3 s_1\}; R_2 = \{s_3 s_2 s_1\}$$

R_3 will be empty, and s_3 will become the rightmost thing, i.e., $s_3 s_4 s_2 s_1$.

Now, if s_4 requests after this, we will have the positions like: $s_2 s_1 s_3 s_4$.

* Stateless pattern is maintained.

4. Quorum-Based Mutual Exclusion Algorithm

$$\forall i, \forall j \quad 1 \leq i, j \leq N \quad R_i \cap R_j \neq \emptyset$$

→ A site can send only one REPLY message at a time.

→ A coterie C is defined as a set of sets where each set g is a member of C ($g \in C$), and g is referred to as the quorum.

→ Intersection Property.

$$\hookrightarrow \text{For every quorum, } g, h \in C, \quad g \cap h \neq \emptyset$$

→ Minimality Property.

\hookrightarrow There should be no quorums, $g, h \in C$, such that $g \supseteq h$ (g is superset of h).

5. Maekawa's Algorithm

→ 4 conditions.

$$\hookrightarrow M1: (\forall i, \forall j: i \neq j, 1 \leq i, j \leq N : R_i \cap R_j \neq \emptyset)$$

$$\hookrightarrow M2: (\forall i: 1 \leq i \leq N : s_i \in R_i)$$

$$\hookrightarrow M3: (\forall i: 1 \leq i \leq N : |R_i| = k)$$

Now, if there are N sites,

$$k = \sqrt{N}$$

$\hookrightarrow M4:$ Any site s_j is contained in k number of R_i ($1 \leq i, j \leq N$)

M1 & M2 : correctness

M3 : REQUESTS

M3 : Request sets of all sites must be equal.

M4 : Enforces that exactly the same number of sites should request permission from any site.

→ Requesting the critical section

(a) A site s_i requests access to the CS by sending REQUEST(i) message to all sites in its request set R_i .

(b) When a site s_j receives the REQUEST(i) message, it sends a REPLY(j) message to s_i , provided it has not sent the REPLY message to a site since its receipt of the last RELEASE message, otherwise, it queues up the REQUEST(i) message for later consideration.

→ Executing the critical section.

(c) Site s_i executes the CS only after it has received a REPLY message from every site in R_i .

→ Releasing the critical section

(d) After the execution of CS is over, site s_i sends a RELEASE(i) message to every site in R_i .

(e) When a site s_j receives a RELEASE(i) message from site s_i , it sends a REPLY message to the next site in the queue and deletes that entry corresponding to s_i from the queue. If the queue is

empty, then the state updates its state to reflect that it has not sent any REPLY message since the receipt of the last RELEASE message.

→ Performance

↳ \sqrt{N} REQUEST, \sqrt{N} REPLY and \sqrt{N} RELEASE messages.

↳ In total, $3\sqrt{N}$ messages.

* There is no notion of time, so requests can go in any order.

So, there is a possibility of deadlock.

(i) T2300239 INT division p 2 via a NMTU (d)

• → Problem of Deadlock:

$s_i s_j s_k$ enter the system simultaneously.

Suppose, $R_i \cap R_j = \{s_{ij}\}$

$R_j \cap R_k = \{s_{jk}\}$

$R_k \cap R_i = \{s_{ki}\}$

Suppose, s_{ij} has been locked by s_i , s_{jk} is locked by s_j , and s_{ki} has been locked by s_k .

Now, due to deadlock, we have 3 other type of commands/messages to avoid the deadlocks!

FAILED

INQUIRE

YIELD

i) Agarwal El Abbade Quorum-Based Algorithm

(Tree - quorum based approach)

→ Algorithm for constructing a tree-structured quorum:

↳ function GetQuorum (Tree : Network hierarchy):

Quorumset

↳ ~~variables~~ var : left, right, Quorumset

begin:

 if Empty(Tree) , then return ({}) ; end

 else if GrantPermission(Tree↑.Node) , then

 return (Tree↑.Node) ;

 return ((Tree↑.Node) ∪ (Tree↑.leftchild))

 return ((Tree↑.Node) ∪ (Tree↑.RightChild))

 else left ← GetQuorum (Tree↑.left);
 right ← GetQuorum (Tree↑.right);

 if (left = φ ∨ right = φ) then

 exit (-1)

 else return (left ∪ right)

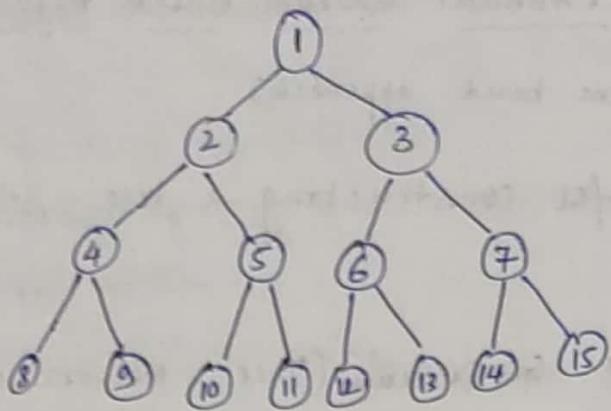
 end

end

end

end GetQuorum

-9:



$\{1, 2, 4, 8\}$ $\{1, 2, 4, 9\}$ $\{1, 2, 5, 10\}$ $\{1, 2, 5, 11\}$

$\{1, 3, 6, 12\}$ $\{1, 3, 6, 13\}$ $\{1, 3, 7, 14\}$ $\{1, 3, 7, 15\}$

Now, when node 3 fails, the possible paths are: $6-12$, $6-13$, $4-7-14$, $7-15$

$\{1, 2, 4, 8\}$ $\{1, 2, 4, 9\}$ $\{1, 2, 5, 10\}$ $\{1, 2, 5, 11\}$

$\{1, 6, 12, 7, 14\}$ $\{1, 6, 12, 7, 15\}$ $\{1, 6, 13, 7, 14\}$ $\{1, 6, 13, 7, 15\}$

If sites 1 and 2 are down, then the quorums must include either $\{4, 8\}$ or $\{4, 9\}$ and either $\{5, 10\}$ or $\{5, 11\}$.

$\{3, 6, 12\}$ $\{3, 6, 13\}$ $\{3, 7, 14\}$

$\{3, 7, 15\}$

$\{4, 5, 3, 6, 8, 10, 12\}$ ~~$\{4, 5, 3, 6, 8, 10, 12\}$~~

$\{4, 5, 3, 6, 8, 10, 13\}$

$\{4, 5, 3, 6, 8, 11, 12\}$

$\{4, 5, 3, 6, 8, 11, 13\}$

$\{4, 5, 3, 6, 9, 10, 12\}$

$\{4, 5, 3, 6, 9, 10, 13\}$

$\{4, 5, 3, 6, 9, 11, 12\}$

$\{4, 5, 3, 7, 8, 10, 14\}$

$\{4, 5, 3, 7, 8, 10, 15\}$

$\{4, 5, 3, 7, 8, 11, 14\}$

$\{4, 5, 3, 7, 8, 11, 15\}$

$\{4, 5, 3, 7, 9, 10, 14\}$

$\{4, 5, 3, 7, 9, 10, 15\}$

$\{4, 5, 3, 7, 9, 11, 14\}$

$\{4, 5, 3, 7, 9, 11, 15\}$

$\{4, 5, 3, 7, 9, 11, 16\}$

$\{4, 5, 3, 7, 9, 11, 17\}$

$\{4, 5, 3, 7, 9, 11, 18\}$

$\{4, 5, 3, 7, 9, 11, 19\}$

$\{4, 5, 3, 7, 9, 11, 20\}$

$\{4, 5, 3, 7, 9, 11, 21\}$

$\{4, 5, 3, 7, 9, 11, 22\}$

$\{4, 5, 3, 7, 9, 11, 23\}$

$\{4, 5, 3, 7, 9, 11, 24\}$

$\{4, 5, 3, 7, 9, 11, 25\}$

$\{4, 5, 3, 7, 9, 11, 26\}$

$\{4, 5, 3, 7, 9, 11, 27\}$

$\{4, 5, 3, 7, 9, 11, 28\}$

$\{4, 5, 3, 7, 9, 11, 29\}$

$\{4, 5, 3, 7, 9, 11, 30\}$

$\{4, 5, 3, 7, 9, 11, 31\}$

$\{4, 5, 3, 7, 9, 11, 32\}$

$\{4, 5, 3, 7, 9, 11, 33\}$

$\{4, 5, 3, 7, 9, 11, 34\}$

$\{4, 5, 3, 7, 9, 11, 35\}$

$\{4, 5, 3, 7, 9, 11, 36\}$

$\{4, 5, 3, 7, 9, 11, 37\}$

$\{4, 5, 3, 7, 9, 11, 38\}$

$\{4, 5, 3, 7, 9, 11, 39\}$

$\{4, 5, 3, 7, 9, 11, 40\}$

$\{4, 5, 3, 7, 9, 11, 41\}$

$\{4, 5, 3, 7, 9, 11, 42\}$

$\{4, 5, 3, 7, 9, 11, 43\}$

$\{4, 5, 3, 7, 9, 11, 44\}$

$\{4, 5, 3, 7, 9, 11, 45\}$

$\{4, 5, 3, 7, 9, 11, 46\}$

$\{4, 5, 3, 7, 9, 11, 47\}$

$\{4, 5, 3, 7, 9, 11, 48\}$

$\{4, 5, 3, 7, 9, 11, 49\}$

$\{4, 5, 3, 7, 9, 11, 50\}$

$\{4, 5, 3, 7, 9, 11, 51\}$

$\{4, 5, 3, 7, 9, 11, 52\}$

$\{4, 5, 3, 7, 9, 11, 53\}$

$\{4, 5, 3, 7, 9, 11, 54\}$

$\{4, 5, 3, 7, 9, 11, 55\}$

$\{4, 5, 3, 7, 9, 11, 56\}$

$\{4, 5, 3, 7, 9, 11, 57\}$

$\{4, 5, 3, 7, 9, 11, 58\}$

$\{4, 5, 3, 7, 9, 11, 59\}$

$\{4, 5, 3, 7, 9, 11, 60\}$

$\{4, 5, 3, 7, 9, 11, 61\}$

$\{4, 5, 3, 7, 9, 11, 62\}$

$\{4, 5, 3, 7, 9, 11, 63\}$

$\{4, 5, 3, 7, 9, 11, 64\}$

$\{4, 5, 3, 7, 9, 11, 65\}$

$\{4, 5, 3, 7, 9, 11, 66\}$

$\{4, 5, 3, 7, 9, 11, 67\}$

$\{4, 5, 3, 7, 9, 11, 68\}$

$\{4, 5, 3, 7, 9, 11, 69\}$

$\{4, 5, 3, 7, 9, 11, 70\}$

$\{4, 5, 3, 7, 9, 11, 71\}$

$\{4, 5, 3, 7, 9, 11, 72\}$

$\{4, 5, 3, 7, 9, 11, 73\}$

$\{4, 5, 3, 7, 9, 11, 74\}$

$\{4, 5, 3, 7, 9, 11, 75\}$

$\{4, 5, 3, 7, 9, 11, 76\}$

$\{4, 5, 3, 7, 9, 11, 77\}$

$\{4, 5, 3, 7, 9, 11, 78\}$

$\{4, 5, 3, 7, 9, 11, 79\}$

$\{4, 5, 3, 7, 9, 11, 80\}$

$\{4, 5, 3, 7, 9, 11, 81\}$

$\{4, 5, 3, 7, 9, 11, 82\}$

$\{4, 5, 3, 7, 9, 11, 83\}$

$\{4, 5, 3, 7, 9, 11, 84\}$

$\{4, 5, 3, 7, 9, 11, 85\}$

$\{4, 5, 3, 7, 9, 11, 86\}$

$\{4, 5, 3, 7, 9, 11, 87\}$

$\{4, 5, 3, 7, 9, 11, 88\}$

$\{4, 5, 3, 7, 9, 11, 89\}$

$\{4, 5, 3, 7, 9, 11, 90\}$

$\{4, 5, 3, 7, 9, 11, 91\}$

$\{4, 5, 3, 7, 9, 11, 92\}$

$\{4, 5, 3, 7, 9, 11, 93\}$

$\{4, 5, 3, 7, 9, 11, 94\}$

$\{4, 5, 3, 7, 9, 11, 95\}$

$\{4, 5, 3, 7, 9, 11, 96\}$

$\{4, 5, 3, 7, 9, 11, 97\}$

$\{4, 5, 3, 7, 9, 11, 98\}$

$\{4, 5, 3, 7, 9, 11, 99\}$

$\{4, 5, 3, 7, 9, 11, 100\}$

$\{4, 5, 3, 7, 9, 11, 101\}$

$\{4, 5, 3, 7, 9, 11, 102\}$

$\{4, 5, 3, 7, 9, 11, 103\}$

$\{4, 5, 3, 7, 9, 11, 104\}$

$\{4, 5, 3, 7, 9, 11, 105\}$

$\{4, 5, 3, 7, 9, 11, 106\}$

$\{4, 5, 3, 7, 9, 11, 107\}$

$\{4, 5, 3, 7, 9, 11, 108\}$

$\{4, 5, 3, 7, 9, 11, 109\}$

$\{4, 5, 3, 7, 9, 11, 110\}$

$\{4, 5, 3, 7, 9, 11, 111\}$

$\{4, 5, 3, 7, 9, 11, 112\}$

$\{4, 5, 3, 7, 9, 11, 113\}$

$\{4, 5, 3, 7, 9, 11, 114\}$

$\{4, 5, 3, 7, 9, 11, 115\}$

$\{4, 5, 3, 7, 9, 11, 116\}$

$\{4, 5, 3, 7, 9, 11, 117\}$

$\{4, 5, 3, 7, 9, 11, 118\}$

$\{4, 5, 3, 7, 9, 11, 119\}$

$\{4, 5, 3, 7, 9, 11, 120\}$

$\{4, 5, 3, 7, 9, 11, 121\}$

$\{4, 5, 3, 7, 9, 11, 122\}$

$\{4, 5, 3, 7, 9, 11, 123\}$

$\{4, 5, 3, 7, 9, 11, 124\}$

$\{4, 5, 3, 7, 9, 11, 125\}$

$\{4, 5, 3, 7, 9, 11, 126\}$

$\{4, 5, 3, 7, 9, 11, 127\}$

$\{4, 5, 3, 7, 9, 11, 128\}$

$\{4, 5, 3, 7, 9, 11, 129\}$

$\{4, 5, 3, 7, 9, 11, 130\}$

$\{4, 5, 3, 7, 9, 11, 131\}$

$\{4, 5, 3, 7, 9, 11, 132\}$

$\{4, 5, 3, 7, 9, 11, 133\}$

$\{4, 5, 3, 7, 9, 11, 134\}$

$\{4, 5, 3, 7, 9, 11, 135\}$

$\{4, 5, 3, 7, 9, 11, 136\}$

$\{4, 5, 3, 7, 9, 11, 137\}$

$\{4, 5, 3, 7, 9, 11, 138\}$

$\{4, 5, 3, 7, 9, 11, 139\}$

$\{4, 5, 3, 7, 9, 11, 140\}$

$\{4, 5, 3, 7, 9, 11, 141\}$

$\{4, 5, 3, 7, 9, 11, 142\}$

$\{4, 5, 3, 7, 9, 11, 143\}$

$\{4, 5, 3, 7, 9, 11, 144\}$

$\{4, 5, 3, 7, 9, 11, 145\}$

$\{4, 5, 3, 7, 9, 11, 146\}$

$\{4, 5, 3, 7, 9, 11, 147\}$

$\{4, 5, 3, 7, 9, 11, 148\}$

$\{4, 5, 3, 7, 9, 11, 149\}$

$\{4, 5, 3, 7, 9, 11, 150\}$

$\{4, 5, 3, 7, 9, 11, 151\}$

$\{4, 5, 3, 7, 9, 11, 152\}$

$\{4, 5, 3, 7, 9, 11, 153\}$

$\{4, 5, 3, 7, 9, 11, 154\}$

$\{4, 5, 3, 7, 9, 11, 155\}$

$\{4, 5, 3, 7, 9, 11, 156\}$

$\{4, 5, 3, 7, 9, 11, 157\}$

$\{4, 5, 3, 7, 9, 11, 158\}$

$\{4, 5, 3, 7, 9, 11, 159\}$

$\{4, 5, 3, 7, 9, 11, 160\}$

$\{4, 5, 3, 7, 9, 11, 161\}$

$\{4, 5, 3, 7, 9, 11, 162\}$

$\{4, 5, 3, 7, 9, 11, 163\}$

$\{4, 5, 3, 7, 9, 11, 164\}$

With the number of node failures greater than or equal to $\lceil \log n \rceil$, the algorithm may not be able to form tree-structured quorums.

If a failed site is a leaf node, the operation has to be aborted and a tree-structured quorum cannot be formed.

→ Suppose a site s wants to enter the critical section, the following events should occur:

1. site s sends a "Request" message to all other sites in the structured quorum it belongs.
2. Each site in the quorum stores incoming requests in a request queue ordered by their timestamp.
3. A site sends a "Reply" message indicating its consent to enter CS only to the request at the head of ~~the~~ its request queue, having the lowest timestamp.
4. If the site s gets a "Reply" message from all sites in the structured quorum it belongs to, it enters the CS.
5. After exiting the CS, s sends a "Relinquish" message to all sites in the structured quorum. On the receipt of a "Relinquish" message, each site removes s 's request from the head of its request queue.
6. If a new request arrives with a timestamp smaller than the request at the head of the queue, an "Inquire" message will be sent to the process whose request is at the head of the queue and the site waits for "Yield" or "Relinquish" message.

7. When a site receives an "Inquire" message, it does as follows:

- If s has acquired all of its ~~more~~ necessary to access the CS, then it simply ignores the message and proceeds normally and sends a "Relinquish" message after it exits the CS.

- If s has not yet collected enough replies from its quorum, then it sends a "Yield" message to the inquiring site.

8. When a site gets a "Yield" message, it puts the pending request (on behalf of which the "Inquire" message was sent) at the head of the queue and sends a "Reply" message to the requester.

Deadlock Detection in Distributed Systems

→ If the allocation sequence of process resources is not controlled, deadlocks can occur.

→ A deadlock can be defined as a condition where a set of processes request resources that are held by other processes in the set.

→ 3 Approaches to handle deadlocks:

↳ Deadlock Prevention

↳ Deadlock Avoidance

↳ Deadlock Detection

→ Deadlock detection requires the examination of the status of the process-resource receiver interaction.

→ To resolve a deadlock, we have to abort a deadlocked process.

(Generally the lowest priority process is aborted)

→ Assumption: The systems have only reusable resources.

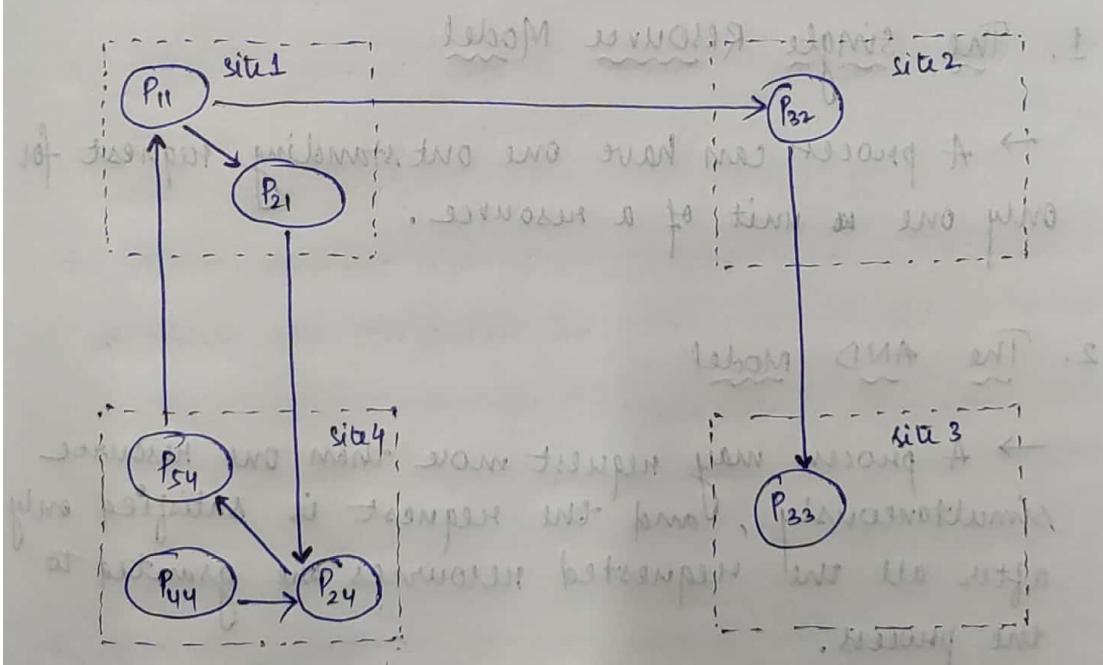
: Processes are allowed to make only exclusive access (not shared) to resources.

: There is only one copy of each resource.

→ Wait-for-Graph (WFG)

↳ Reflects the state of a system.

↳ Shows which process is waiting for which resource.



A WFG for 4 processes.

There is a cycle between sites 1 and 4, so it is a deadlock condition.

→ A system is deadlocked iff there exists a directed cycle or knot in the WFG.

Deadlock Handling Strategies

→ Prevention : inefficient

Avoidance : impractical

Detection : only thing that is practically efficient

→ Detection: two types of detection:

↳ Progress (no undetected deadlocks)

↳ Safety (no false deadlocks)

Models of Deadlock

1. The single Resource Model

→ A process can have one outstanding request for only one unit of a resource.

2. The AND Model

→ A process may request more than one resource simultaneously, and the request is satisfied only after all the requested resources are granted to the process.

3. The OR Model