

Distributed Mutual Exclusion

Mutual Exclusion in Single Computer System Vs Distributed System

= In single computer systems, the status of a shared resource and the status of user is available in the shared memory. and solutions to the Mutual Exclusion problem can be easily implemented using shared variables (eg. Semaphores)

⇒ In D.S. both the shared resources and the users may be distributed and shared memory does not exist. Approaches based on shared variables are not applicable to D.S. and approaches based on message passing must be used

⇒ The problem of M.E. becomes much more complex in D.S. because of the lack of both shared memory and a common physical clock and because of unpredictable message delay

Classification of Mutual Exclusion Algo.

(2)

Non Token
Based

⇒ These algo. requires 2 or more successive round of msg exchanges among the site

Token Based

Algorithm.

= In these algorithm, a unique token is shared among the site

A site allowed to enter the cs if it possesses the token. It continues to hold the token until the execution of cs is over.

Requirement of Mutual exclusion

The primary objective is to guarantee that only one request accesses the cs at a time.

The following characteristics are considered important in a Mutual exclusion Algo.

- (1) Freedom from Deadlocks: 2 or more sites should not endlessly wait for msg. that will never arrive
- (2) Freedom from Starvation: A site should not be forced to wait indefinitely to execute cs while other sites are repeatedly executing cs.

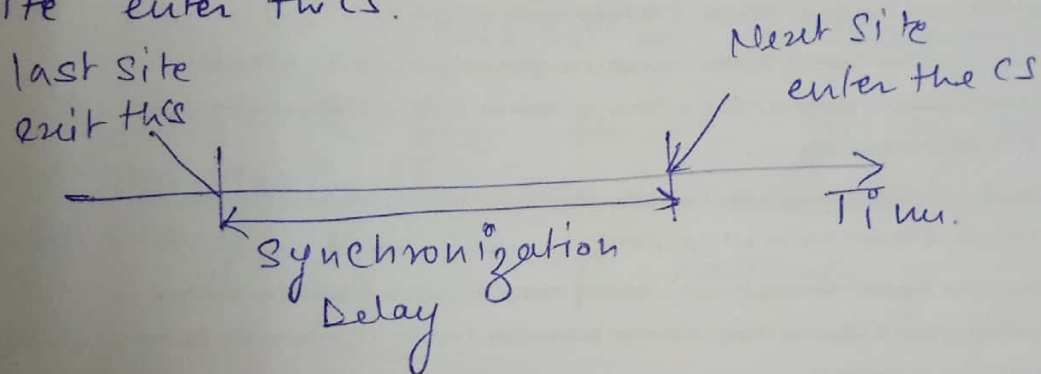
3) Fairness: Requests must be executed in the order they are made. (3)

4) Fault Tolerance: A mutual exclusion algo. is fault tolerant if in the wake of failure it can reorganize itself so that it continues to function without any disruption

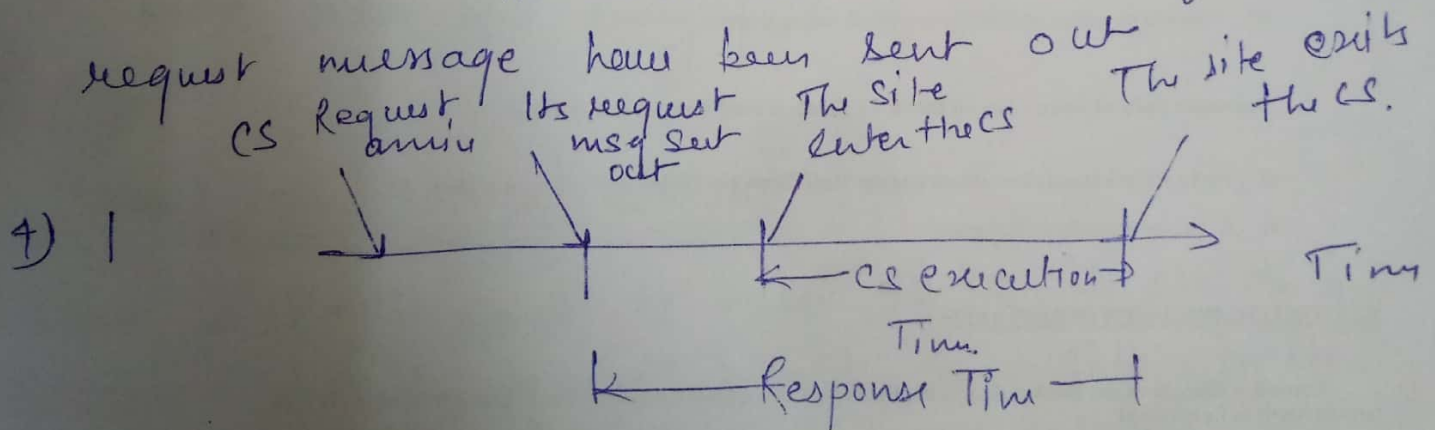
Performance Measures:

(1) Synchronization Delay: Time required after a site leaves the CS and before the next

site enters the CS.



(2) Response Time: Time interval a request waits for its CS execution to over after its request message has been sent out. The site enters the CS. The site exits the CS.



3) Throughput = the rate at which the system executes the requests for the CS.

if S_d = synchronization Delay
 E = avg CS execution time

$$\text{Throughput} = 1/(S_d + E)$$

17

Non Token Based Algorithm

(4)

In non token Based Algorithm, a site communicate with a set of other sites to arbitrate who should execute the CS next.

⇒ for a site s_i , request set R_i contains ids of all those sites from which site s_i must acquire permission before entering the CS.

⇒ Non Token Based mutual exclusion algo use timestamps to order requests for the CS. ~~and to~~.

In all these algo. logical clocks are maintained and updated according to Lamport's Scheme

Each request for the CS gets a timestamp and smaller timestamp requests have priority over larger timestamp requests

Lamport's Algorithm

This algorithm requires messages to be delivered in the FIFO order b/w every pair of sites.

Algo:

Request the critical Section:

1) When a site s_i wants to enter the CS, it sends a $REQUEST(t_{s_i}, i)$ message to all the sites in its request set R_i , and places the request on request-queue(t_{s_i}, i) is the timestamp of the request)

Executing out

(2) When a site S_j receives the REQUEST (ts_i, i) message from site S_i , it returns a timestamped REPLY message to S_i and places site S_i 's request on request-queue.

9

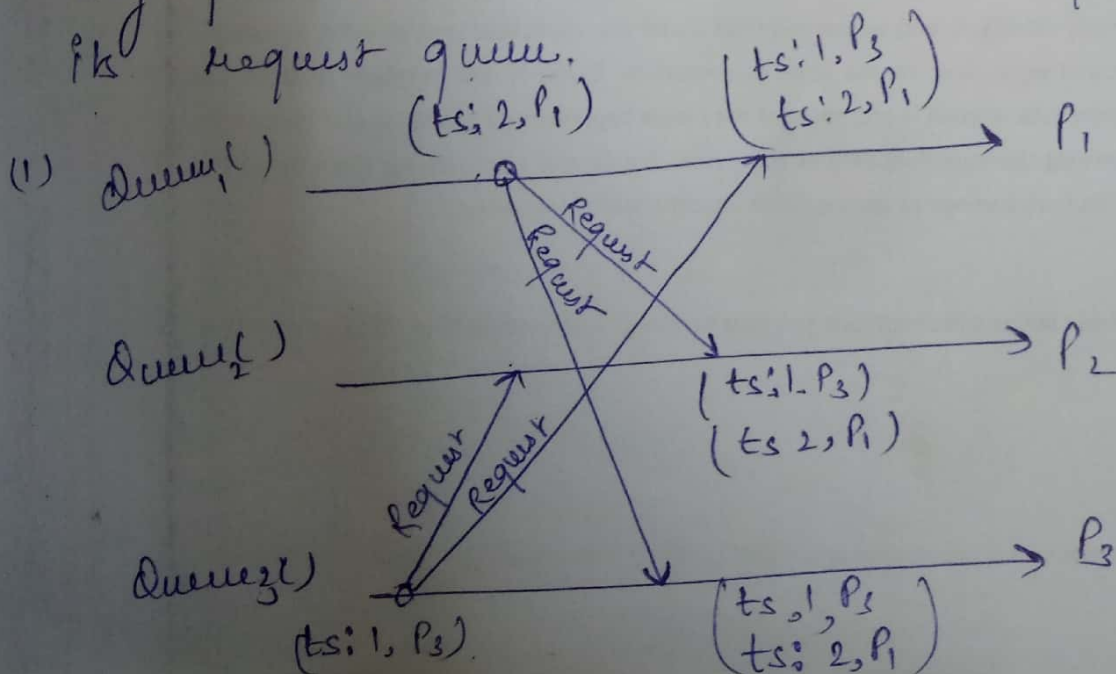
Executing the CS.

site S_i enters the CS when the 2 following conditions hold.

- (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.

Releasing the Critical Section

- (1) Site S_i upon exiting the CS, removes its request from top of its request queue and sends a timestamped RELEASE message to all the sites in its request set.
 (2) When a site S_j ~~release~~ receives a RELEASE msg from site S_i , it removes S_i 's request from its request queue.



The RICART AGRAWALA Algorithm

(7)

⇒ It is an optimisation of Lamport's algorithm

Algo

Requesting the c.s.

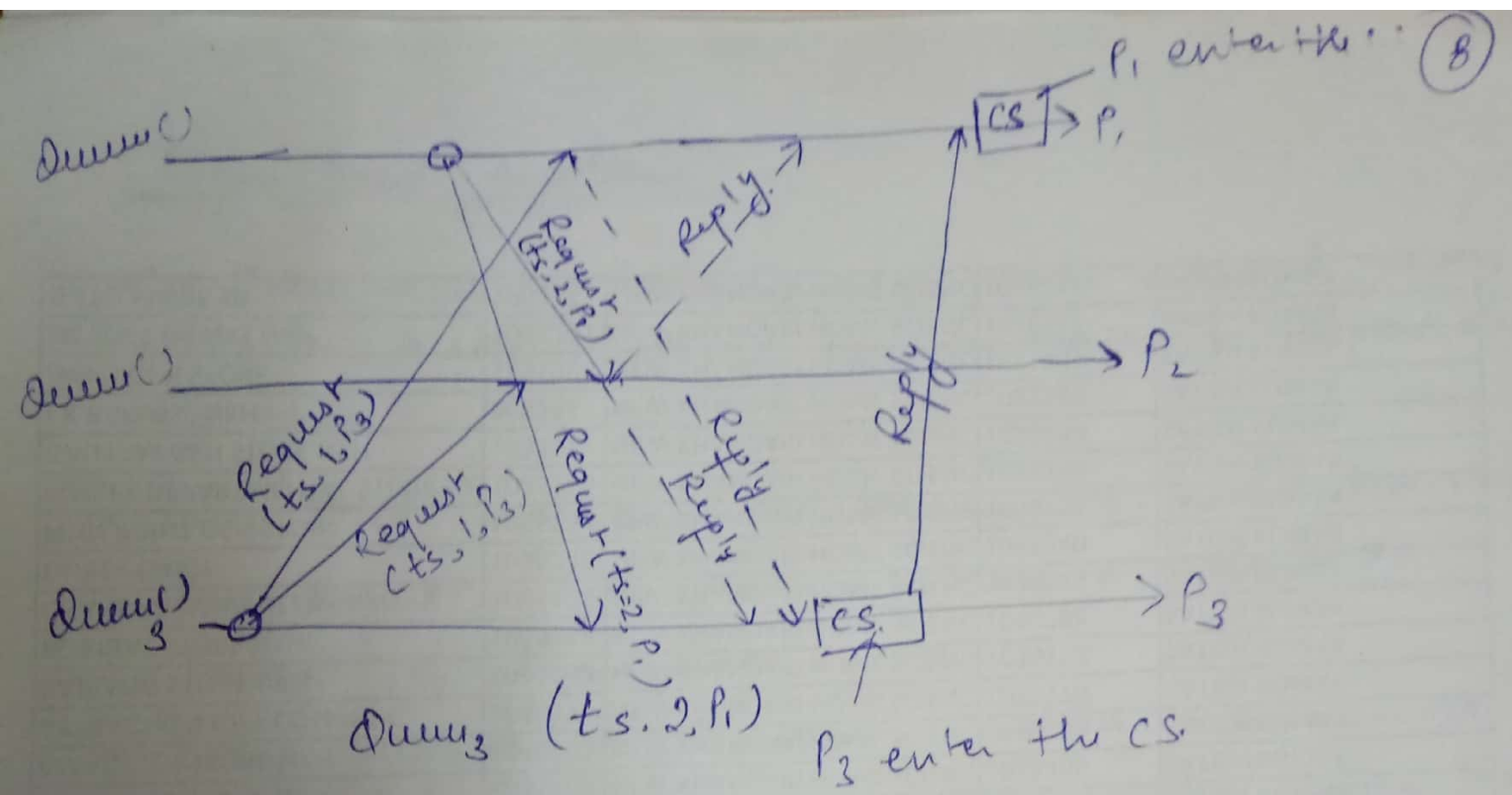
- 1) When a site S_i wants to enter the c.s., it sends a timestamped REQUEST message to all the sites in its request set.
- 2) When a site S_j receives a REQUEST message from site S_i , it sends a REPLY message to site S_i .
If site S_j is neither requesting nor executing the c.s. or if the site S_j is requesting and S_i 's request's timestamp is smaller than site S_j 's own request's timestamp. The request is deferred otherwise.

Executing the c.s.

- 3) site S_i enter the c.s. after it has received REQUEST message from all the sites in its request set.

Releasing the C.S.

- 4) When a site S_i exits the c.s. it sends REPLY message to all the ~~diff~~ deferred request



The Ricant Agrawala algo requires $2(N-1)$ messages
 per CS. $(N-1)$ Request, $(N-1)$ Reply msg.
 Synchronization delay = T .

Token Based Algorithm

- * In Token based Algorithm, a unique token is shared among all sites. A site is allowed to enter the cs if it possesses the token.
- * Token based algorithm use sequence no. instead of timestamps. Every request for the token contains a sequence no. instead of timestamps.
- * Every request for the token contain a sequence no.
- * A site increments its sequence no. counter every time it makes a request. for the token
A primary function of the sequence no is to distinguish b/w old and current requests.

Suzuki Kasami Broadcast Algorithm

- * If a site attempting to enter the cs. does not have the token. It broadcast a REQUEST message. for the token to all the other sites.
- * A site that possesses the token send it to the requesting site upon receiving its REQUEST. message.
- * If a site receives a REQUEST. message when it is executing the cs. it sends the token only after it has exited the cs.
- * A site holding the token can enter its cs. repeatedly until it sends the token to some other site.

Algorithm

The token: \rightarrow Queue (FIFO) of requesting processes
 $LN[1..n]$: Sequence no of request that j executed most recently.

Requesting the C.S.

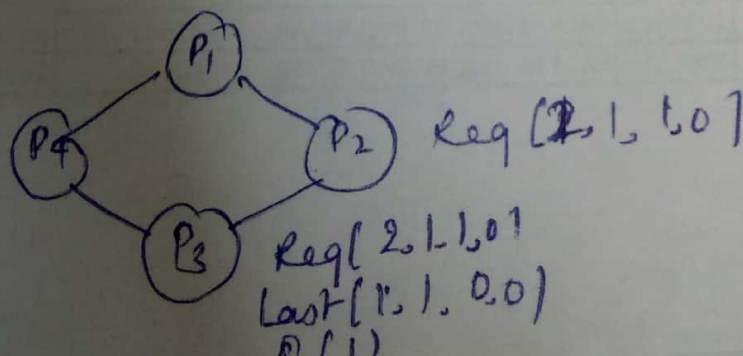
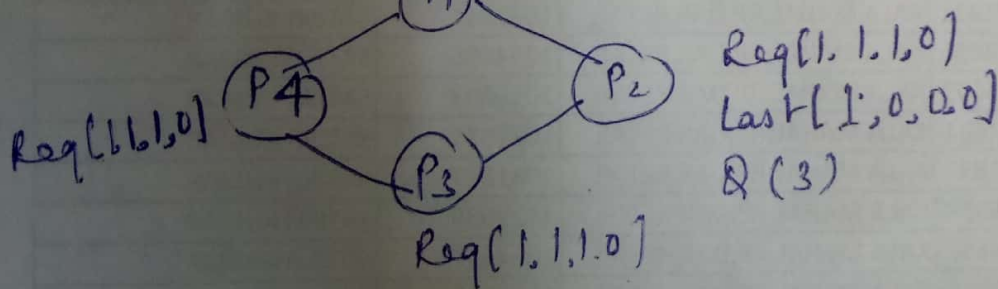
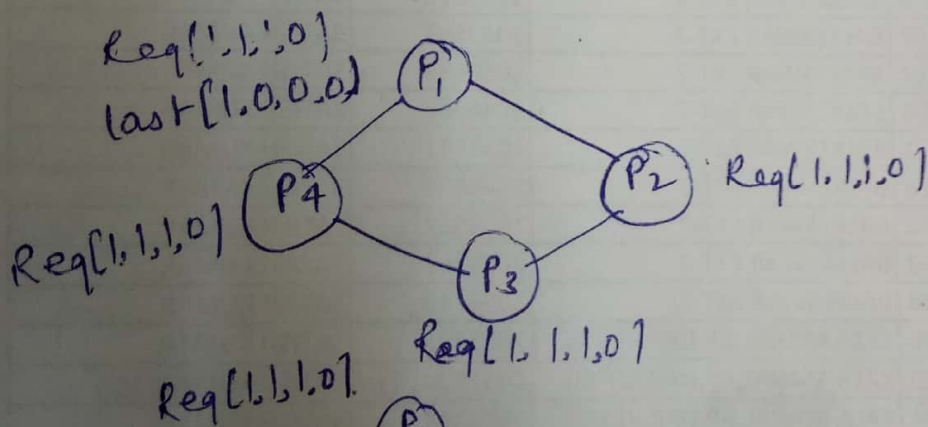
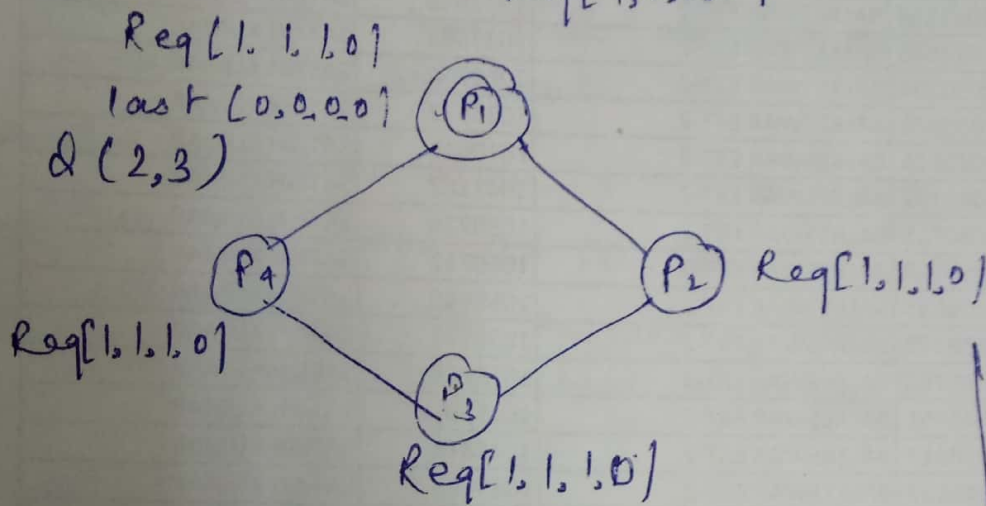
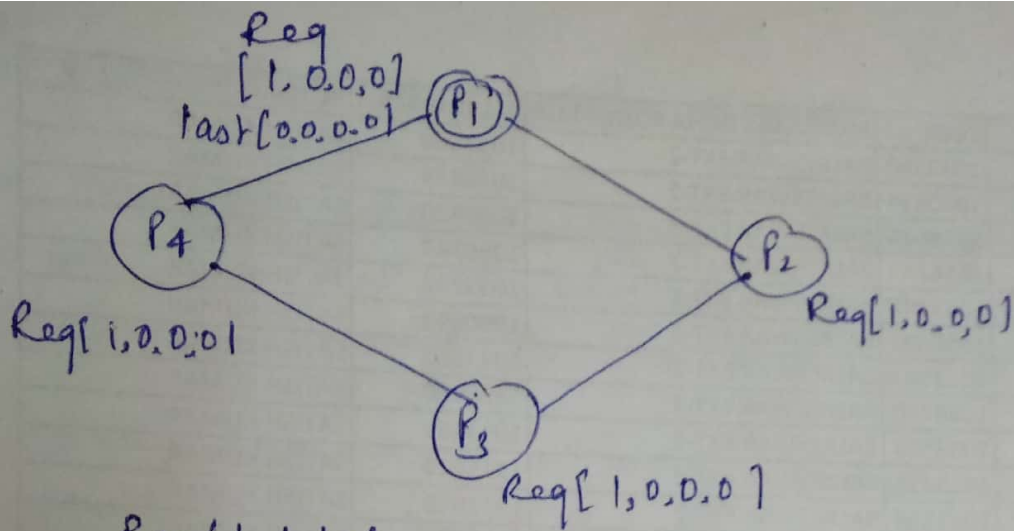
- (1) If the requesting site S_i does not have the token then it increments its sequence no. $RN[i]$, and send a ~~send~~ REQUEST(i, sn) message to all other sites.
- (2) When a site S_j receives this message it sets $RN[j]$ to $\max(RN[j], sn)$. If S_j has the idle token, then it send the token to S_i if $RN[j] = LN[j+1]$.

Executing the C.S.:

- 3) Site S_i executes the C.S. when it has received the token.

Releasing the critical section:

- 4) It sets $LN[i]$ element of the token array equal to $RN[i]$.
- 5) For every site S_j whose ID is not in the token queue it appends its ID to the token queue if $RN[i][j] = LN[j] + 1$.
- 6) If token queue is non empty after the above update then it delete the top site ID from the queue and send the token to the site indicated by the ID.



Synchronization
 Delay $\Rightarrow 0$
 The algo requires
 O or M message
 per CS invocation

Raymond's Tree Based Algorithm:

- ⇒ In this algo. sites are logically arranged as a directed tree. Such that the edge of the tree are assigned directions toward the site that has token.
- ⇒ Every site has a local variable holder that points to an immediate neighbor node on a directed path to the root node.
- ⇒ At root site, holder points to itself.
- ⇒ Every site keeps a FIFO Queue, called request-q, which stores the requests of those neighboring sites that have sent a request to this site.

Algo:

Requesting the C.S.

- (1) When a site wants to enter the C.S. It send a Request Message to the node along the directed path to the root, provided it does not hold the token and its request-q is empty. It then add its request to its request-q.
- (2) When a site on the path receives this message it places the REQUEST in request-q and send a REQUEST msg along the directed path to the root provided it has not sent out a REQUEST msg on its outgoing edge.

3) When a root site receives a REQUEST message. it sends the token to the site from which it received the request message ~~on its outgoing edge~~ and sets its holder variable to point at that site

4) When a site receives the token it deletes the top entry from its request-q, sends the token to the site indicated in this entry, and sets its holder variable to point at that ~~to~~ self

Executing the CS

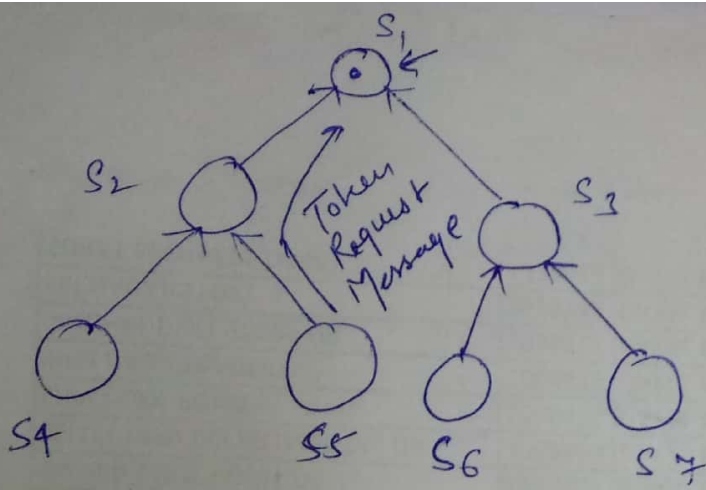
5) A site enters the CS when it receives the token and its own entry is at the top of its queue. then, the site deletes the top entry from its request-q and enters the CS.

6) Releasing the CS After a site has finished execution of the CS it takes the following actions:

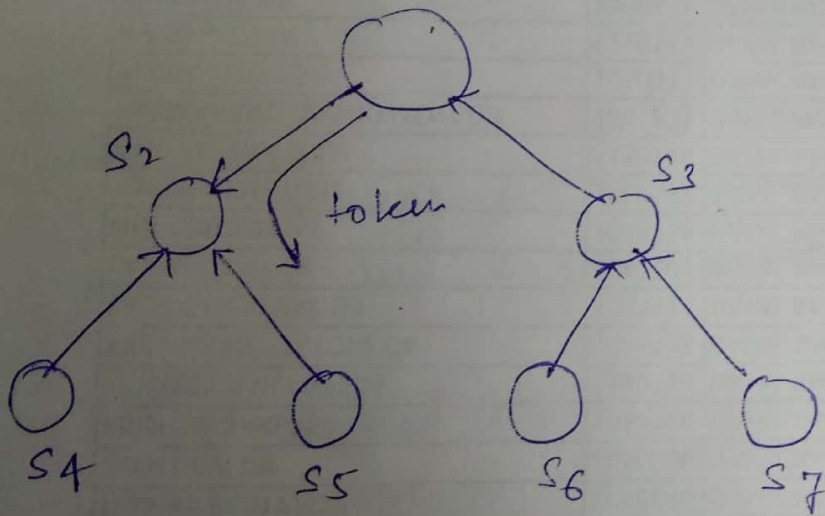
(b) If its request-q is nonempty then it deletes the top entry from its request-q, sends the token to that site and set its holder variable to point at that site

~~(7) If its request-q is empty~~

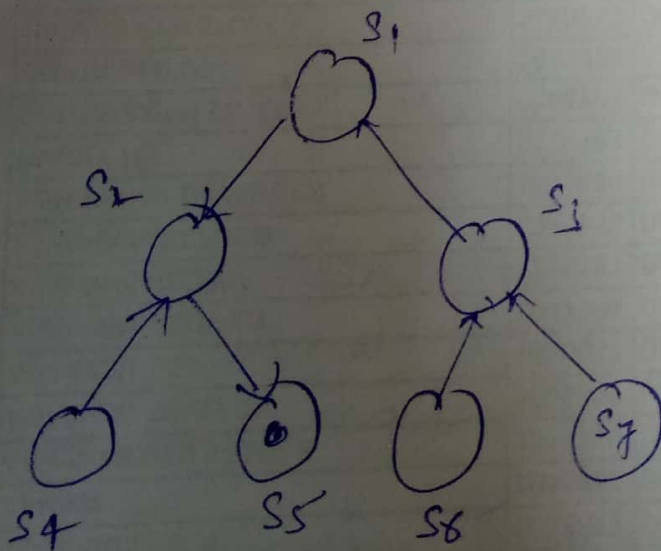
⇒ Algo is free from deadlocks



Site S_5 is requesting the token



The token is in transit to S_5



S_5 received the token

System Model: The problem of deadlocks has been generally studied in D.S. under the following Model

- The system has only reusable resources.
- Processes are allowed only exclusive access to resources
- There is only one copy of each resource

A process can be in 2 states - Running or Blocked. In the running state, a process has all the needed resources and is either executing or is ready for execution.

In the blocked state, a process is waiting to acquire some resources.

Resource Vs. Communication Deadlock:

Two types of deadlock $\begin{cases} \text{Resource Deadlock} \\ \text{Communication Deadlock} \end{cases}$

In Resource Deadlock, processes can simultaneously wait for several resources and can not proceed until they have acquired all those resources.

In Communication Deadlock; processes wait to communicate with other processes among a set of processes. A waiting process can be unblocked on receiving a communication from any other one of these processes.

Deadlock Handling Strategies in distributed System

Deadlock Prevention : Deadlock prevention is commonly achieved by either having a process acquire all the needed resource simultaneously before it begins execution, or by preempting a process that hold the needed resource.

Deadlock Avoidance In the deadlock avoidance approach to D.S., a resource is granted to a process if the resulting global ~~state~~ system state is safe.

Beoz of the following problems, deadlock avoidance can be seen ~~is~~ as impractical in D.S.

- (1) Every site has to maintain information on the global state of system, which translate into huge storage ~~information~~ ~~can be global~~ ~~data~~ ~~of the system~~ requirements and extensive communication cost.
- (2) The process of checking for a safe global state must be mutually exclusive.
- (3) Due to the large no of processes and resources it will computationally expensive to check for a safe state.

Deadlock Detection : Deadlock Detection in D.S has 2 favorable Condⁿ.

- (1) Once a cycle is formed in the wFG, it persist until it is detected and broken
- (2) Cycle detection can proceed, concurrently with the normal activity of a system.

Resolution: Deadlock resolution involves breaking existing ⁽²⁾ wait for dependencies in the system wfg to solve the deadlock.
It involves rolling back one or more processes that are deadlocked and assign their resources to blocked processes in ~~the~~ the deadlock so that they can resume execution.

Distributed Deadlock Detection:

⇒ Occurrence of a cycle b/w Resource and process in wfg is known as deadlock Detection
⇒ Deadlock Detection is the most popular strategy for handling deadlocks in D.S.

We can control Deadlock detection in 3 types

- (1) Centralized
- (2) Distributed
- (3) Hierarchical

(1) Centralized Deadlock Detection Algo.

- In this only one node is responsible for building and ~~any~~ analyzing a real wfg for a cycle
- Each site maintain its wfg and update control site periodically or on request
- A central control site constructs the global wfg and searches for cycle.
- Control site maintain wfg continuously
- All site request resources and release resources by sending corresponding message to control site

-) Control site update WFG for each request and release of resources
-) for every request edge added to WFG and control site checks WFG for deadlock

Disadvantages: single point of failure and congestion

Distributed Deadlock Detection Algorithm:

In this each node participates equally in detecting the deadlock --- WFG

Principle • All sites are responsible for detecting a global deadlock

- Global state graph distributed over many sites; several of them participate in detection
- Detection initiated when a process suspected to be deadlocked

Advantage: No single point of failure and no congestion

Disadvantage: Difficult to implement

Hierarchical Deadlock Detection Algorithm:

In this, nodes are organized or arranged in a tree pattern

Principle: In this the sites are organized in a tree structure with one site at root of the tree.

Each node has information about the dependent nodes.

Deadlock is detected by the node that is the common ancestor of all sites which have resource allocations

in conflict

- Deadlock is detected at the lowest level

In this sets of nodes are required to report periodically to a control site node but control sites are organized in a tree structure

- The Master control site forms the root of the tree with leaf nodes having no control responsibility and interior nodes serving as controllers for their branches

Centralized Deadlock Detection Algorithm

The Completely Centralized Algorithm:

- ⇒ The Completely Centralized algo. is the simplest centralized deadlock detection algorithm, wherein a designed site called the control site maintains the WFG of the entire system and checks it for the existence of deadlock cycles.
- ⇒ All sites request and release resources by sending request resource and release resource message to control site.
- ⇒ When a control site receives a request resource or a release resource message it updates its WFG. The control site checks the WFG for deadlocks whenever a request edge is added to the WFG.

- ⇒ This algo is simple and easy to implement.
- ⇒ but it is also highly inefficient bcoz all resources acquisition and release request must go through the control site. even when the resource is local
 - ⇒ large delays.
 - ⇒ large communication overhead
 - ⇒ Reliability is poor because of control site

The Ho Ramamoorthy Algorithm

~~Algorithm~~. Two phase Model (Can be for AND or OR Model)

- ⇒ In this each site has a status ~~table~~ Table of locked and waited resources
- ⇒ The Control site will periodically ask for this table from each node
- ⇒ The Control node will search for cycles and if found will request the table again from each node
- ⇒ Only common info. in both reports will be analyzed for confirmation of a cycle

One Phase (Can be for AND or OR Model)

- In this each site keeps 2 tables for all local process:
 - Process Status
 - Resource Status Table
- The Control site will periodically ask for these tables via msg from each node
- The Control site will build and analyze the WFG, looking for cycles and resolving them when found

In One phase algo. does not detect the ~~dead~~ false deadlock because it eliminates the inconsistency in state info.

⇒ The 1 phase algo is faster and requires fewer msg as compared to 2 phase algo

Distributed Deadlock Detection Algo.

In Distributed deadlock detection algo. all sites collectively cooperate to detect a cycle in the state graph

A Path Pushing Algorithm:

⇒ In this path information sent from waiting node to blocking node

⇒ This approach deal with Transactions
Each transaction may have sub transactions but they executes in sequentially manner

⇒ In this all the transactions are totally ordered or sequentially.

⇒ Path info sent from waiting node to blocking node

⇒ In this WFG. Constructed by disseminating dependency Sequences.

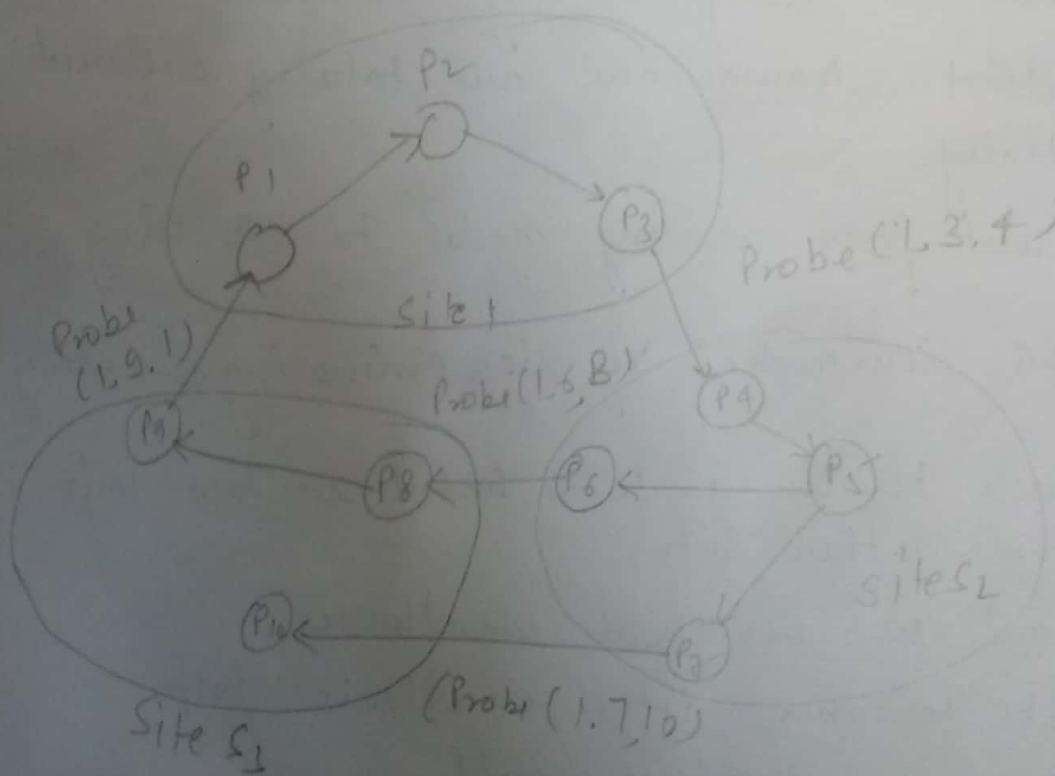
— Each node builds a WFG based on local info and info. from other site.

— Detect and resolves local deadlocks

— Transmit to other sites deadlock info in form of waiting path

Edge Chasing Algorithm:

- ⇒ In this message are sent along graph edges.
 - Special message (probe) sent along edges of WFG to detect a cycle
 - When Blocked process receive msg. resends it on its outgoing edge of WFG
 - When a process receive a msg it initiated declare deadlock
- ⇒ Process send special msg
- ⇒ A blocked process receive msg and circulates it along its outgoing dependency edge



Hierarchical Deadlock Detection Algo.

⇒ Sites are arranged in hierarchical fashion and a site is responsible for detecting deadlock involving only its children sites.

The Menasce Muntz Algorithm

⇒ In the hierarchical deadlock detection algorithm of Menasce Muntz, all the controllers are arranged in tree fashion. (Controller manages a resource or it is responsible for deadlock detection)

⇒ The Controller at the bottom most level (Leaf Controller) manages resource and others. are responsible for deadlock detection.

= A Leaf Controller maintains a part of the global TWF Graph. concerned with the allocation of resources at that Leaf Controller

⇒ Whenever a change occurs in a controller TWF Graph due to a resource allocation, wait or release, it is propagated to its parent controller

The parent Controller makes changes in its TWF Graph searches for cycle and propagate the changes upward if necessary.

A non leaf controller can receive up to date info. concerning the TWF Graph of its children continuously.

The Ho Ramamoorthy Algo.

- ⇒ In the hierarchical algo. of Ho and Ramamoorthy sites are grouped into several disjoint clusters. Periodically a site is chosen as a central control site, which dynamically chooses a control site for each cluster.
- ⇒ The Central Control site requests from every control site their intercluster transaction status info. and wait for replies.
- ⇒ Control site collects status table from all the sites in its cluster and applies the two phase deadlock detection algo. to detect all deadlock involving only intercluster transactions.
- ⇒ It sends intercluster transaction status info. and wait for replies, to the Central Control site.
- ⇒ The Central Control site ~~separately~~ detects all deadlocks located in its cluster and the Central Control site detects all intercluster deadlocks.

