

UNIT III - DISTRIBUTED MUTEX & DEADLOCK

Distributed mutual exclusion algorithms: Introduction – Preliminaries – Lamport's algorithm – Ricart-Agrawala algorithm – Maekawa's algorithm – Suzuki-Kasami's broadcast algorithm. Deadlock detection in distributed systems: Introduction – System model – Preliminaries – Models of deadlocks – Knapp's classification – Algorithms for the single resource model, the AND model and the OR model.

DISTRIBUTED MUTUAL EXCLUSION ALGORITHMS

Mutual exclusion: Concurrent access of processes to a shared resource or data is executed in mutually exclusive manner.

- Only one process is allowed to execute the critical section (CS) at any given time.
- In a distributed system, shared variables (semaphores) or a local kernel cannot be used to implement mutual exclusion.
- Message passing is the sole means for implementing distributed mutual exclusion.
- Distributed mutual exclusion algorithms must deal with unpredictable message delays and incomplete knowledge of the system state.

Three basic approaches for distributed mutual exclusion:

- Token based approach
- Non-token based approach
- Quorum based approach

3.1 INTRODUCTION

Token-based approach:

- A unique token is shared among the sites.
- A site is allowed to enter its CS if it possesses the token.
- Mutual exclusion is ensured because the token is unique.

Non-token based approach:

- △ Two or more successive rounds of messages are exchanged among the sites to determine which site will enter the CS next.

Quorum based approach:

- △ Each site requests permission to execute the CS from a subset of sites (called a quorum).
- △ Any two quorums contain a common site.
- △ This common site is responsible to make sure that only one request executes the CS at any time.

3.2 PRELIMINARIES

System Model

- The system consists of N sites, S_1, S_2, \dots, S_N .
- We assume that a single process is running on each site. The process at site S_i is denoted by p_i .
- A site can be in one of the following three states: requesting the CS, executing the CS, or neither requesting nor executing the CS (i.e., idle).
- In the 'requesting the CS' state, the site is blocked and can not make further requests for the CS. In the 'idle' state, the site is executing outside the CS. In token-based algorithms, a site can also be in a state where a site holding the token is executing outside the CS (called the *idle token* state).
- At any instant, a site may have several pending requests for CS. A site queues up these requests and serves them one at a time.

Requirements of Mutual Exclusion Algorithms

- **Safety Property**: At any instant, only one process can execute the critical section.
- **Liveness Property**: This property states the absence of deadlock and starvation. Two or more sites should not endlessly wait for messages which will never arrive.
- **Fairness**: Each process gets a fair chance to execute the CS. Fairness property generally means the CS execution requests are executed in the order of their arrival (time is determined by a logical clock) in the system.

The first property is absolutely necessary and the other two properties are considered important in mutual exclusion algorithms.

Performance Metrics

The performance is generally measured by the following **four metrics**:

- **Message complexity**: The number of messages required per CS execution by a site.
- **Synchronization delay**: After a site leaves the CS, it is the time required and before the next site enters the CS (see Figure 1).

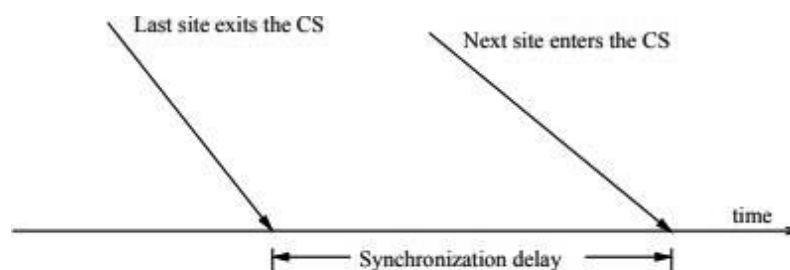


Figure 1: Synchronization Delay.

- **Response time:** The time interval a request waits for its CS execution to be over after its request messages have been sent out (see Figure 2).
- **System throughput:** The rate at which the system executes requests for the CS.

$$\text{system throughput} = 1 / (SD + E)$$

where SD is the synchronization delay and E is the average critical section execution time.

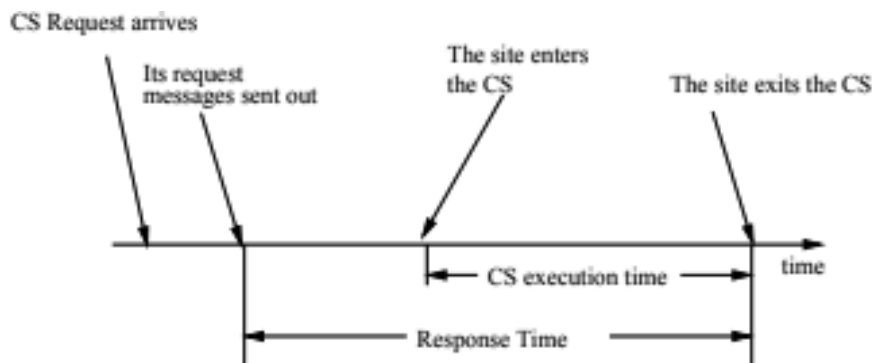


Figure 2: Response Time.

Low and High Load Performance:

- We often study the performance of mutual exclusion algorithms under two special loading conditions, viz., “low load” and “high load”.
- The load is determined by the arrival rate of CS execution requests.
- Under *low load* conditions, there is seldom more than one request for the critical section present in the system simultaneously.
- Under *heavy load* conditions, there is always a pending request for critical section at a site.

Best and worst case performance

- Generally, mutual exclusion algorithms have best and worst cases for the performance metrics. In the best case, prevailing conditions are such that a performance metric attains the best possible value. For example, in most mutual exclusion algorithms the best value of the response time is a round-trip message delay plus the CS execution time, $2T + E$.
- Often for mutual exclusion algorithms, the best and worst cases coincide with low and high loads, respectively. For examples, the best and worst values of the response time are achieved when load is, respectively, low and high; in some mutual exclusion algorithms the best and the worse message traffic is generated at low and heavy load conditions, respectively.

3.3 LAMPORT'S ALGORITHM

- Requests for CS are executed in the increasing order of timestamps and time is determined by logical clocks.
- Every site S_i keeps a queue, *request queue_i*, which contains mutual exclusion requests ordered by their timestamps.
- This algorithm requires communication channels to deliver messages the FIFO order.

The Algorithm

Requesting the critical section:

- When a site S_i wants to enter the CS, it broadcasts a REQUEST(ts_i, i) message to all other sites and places the request on *request queue_i*. ((ts_i, i) denotes the timestamp of the request.)
- When a site S_j receives the REQUEST(ts_i, i) message from site S_i , places site S_i 's request on *request queue_j* and it returns a timestamped REPLY message to S_i .

Executing the critical section: Site S_i enters the CS when the following two 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_i*.

Releasing the critical section:

- Site S_i , upon exiting the CS, removes its request from the top of its request queue and broadcasts a timestamped RELEASE message to all other sites.
- When a site S_j receives a RELEASE message from site S_i , it removes S_i 's request from its request queue.

When a site removes a request from its request queue, its own request may come at the top of the queue, enabling it to enter the CS.

Correctness

Theorem: Lamport's algorithm achieves mutual exclusion.

Proof:

- Proof is by contradiction. Suppose two sites S_i and S_j are executing the CS concurrently. For this to happen conditions L1 and L2 must hold at both the sites *concurrently*.
- This implies that at some instant in time, say t , both S_i and S_j have their own requests at the top of their *request queues* and condition L1 holds at them. Without loss of generality, assume that S_i 's request has smaller timestamp than the request of S_j .
- In condition L1 and FIFO property of the communication channels, it is clear that at instant t the request of S_i must be present in *request queue_j* when S_j was executing its CS. This implies that S_j 's own request is at the top of its own *request queue* when a smaller timestamp request, S_i 's request, is present in the *request queue_j* – a contradiction!

Theorem: Lamport's algorithm is fair.**Proof:**

- The proof is by contradiction. Suppose a site S_i 's request has a smaller timestamp than the request of another site S_j and S_j is able to execute the CS before S_i .
- For S_j to execute the CS, it has to satisfy the conditions L1 and L2. This implies that at some instant in time say t , S_j has its own request at the top of its queue and it has also received a message with timestamp larger than the timestamp of its request from all other sites.
- But *request queue* at a site is ordered by timestamp, and according to our assumption S_i has lower timestamp. So S_i 's request must be placed ahead of the S_j 's request in the *request queue*. This is a contradiction!

Performance

- For each CS execution, Lamport's algorithm requires $(N - 1)$ REQUEST messages, $(N - 1)$ REPLY messages, and $(N - 1)$ RELEASE messages. Thus, Lamport's algorithm requires $3(N - 1)$ messages per CS invocation.
- Synchronization delay in the algorithm is T .

An optimization

- In Lamport's algorithm, REPLY messages can be omitted in certain situations. For example, if site S_j receives a REQUEST message from site S_i after it has sent its own REQUEST message with timestamp higher than the timestamp of site S_i 's request, then site S_j need not send a REPLY message to site S_i .
- This is because when site S_i receives site S_j 's request with timestamp higher than its own, it can conclude that site S_j does not have any smaller timestamp request which is still pending.
- With this optimization, Lamport's algorithm requires between $3(N - 1)$ and $2(N - 1)$ messages per CS execution.

3.4 RICART-AGRAWALA ALGORITHM

- The Ricart-Agrawala algorithm assumes the communication channels are FIFO. The algorithm uses two types of messages: REQUEST and REPLY.
- A process sends a REQUEST message to all other processes to request their permission to enter the critical section. A process sends a REPLY message to a process to give its permission to that process.
- Processes use Lamport-style logical clocks to assign a timestamp to critical section requests and timestamps are used to decide the priority of requests.
- Each process p_i maintains the Request-Deferred array, RD_i , the size of which is the same as the number of processes in the system.
- Initially, $\forall i \forall j : RD_i[j]=0$. Whenever p_i defer the request sent by p_j , it sets
- $RD_i[j]=1$ and after it has sent a REPLY message to p_j , it sets $RD_i[j]=0$.

Description of the Algorithm

Requesting the critical section:

- (a) When a site S_i wants to enter the CS, it broadcasts a timestamped REQUEST message to all other sites.
- (b) When 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 CS, 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. Otherwise, the reply is deferred and S_j sets $RD_j[i]=1$

Executing the critical section:

- (c) Site S_i enters the CS after it has received a REPLY message from every site it sent a REQUEST message to.

Releasing the critical section:

- (a) When site S_i exits the CS, it sends all the deferred REPLY messages: $\forall j$ if $RD_i[j]=1$, then send a REPLY message to S_j and set $RD_i[j]=0$.

Notes:

- When a site receives a message, it updates its clock using the timestamp in the message.
- When a site takes up a request for the CS for processing, it updates its local clock and assigns a timestamp to the request.

Theorem: Ricart-Agrawala algorithm achieves mutual exclusion.

Proof:

- Proof is by contradiction. Suppose two sites S_i and S_j are executing the CS concurrently and S_i 's request has higher priority than the request of S_j .
- Clearly, S_i received S_j 's request after it has made its own request.
- Thus, S_j can concurrently execute the CS with S_i only if S_i returns a REPLY to S_j

(in response to S_j 's request) before S_i exits the CS.

- However, this is impossible because S_j 's request has lower priority. Therefore, Ricart-Agrawala algorithm achieves mutual exclusion.

Performance

- For each CS execution, Ricart-Agrawala algorithm requires $(N - 1)$ REQUEST messages and $(N - 1)$ REPLY messages.
- Thus, it requires $2(N - 1)$ messages per CS execution. Synchronization delay in the algorithm is T .

3.5 MAEKAWA'S ALGORITHM

Maekawa's algorithm was the first quorum-based mutual exclusion algorithm. The request sets for sites (i.e., quorums) in Maekawa's algorithm are constructed to satisfy the following conditions:

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

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

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

$$M4: \text{Any site } S_j \text{ is contained in } K \text{ number of } R_i\text{'s, } 1 \leq i, j \leq N.$$

Maekawa used the theory of projective planes and showed that $N = K(K - 1) + 1$. This relation gives $|R_i| = \sqrt{N}$.

- Conditions M1 and M2 are necessary for correctness; whereas conditions M3 and M4 provide other desirable features to the algorithm.
- Condition M3 states that the size of the requests sets of all sites must be equal implying that all sites should have to do equal amount of work to invoke mutual exclusion.
- Condition M4 enforces that exactly the same number of sites should request permission from any site implying that all sites have "equal responsibility" in granting permission to other sites.

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

Requesting the critical section

- (a) A site S_i requests access to the CS by sending REQUEST(i) messages 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 hasn't sent a REPLY message to a site since its receipt of the last RELEASE message. Otherwise, it queues up the REQUEST(i) 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

- (a) After the execution of the CS is over, site S_i sends a RELEASE(i) message to every site in R_i .
- (b) When a site S_j receives a RELEASE(i) message from site S_i , it sends a REPLY message to the next site waiting in the queue and deletes that entry from the queue. If the queue is empty, then the site updates its state to reflect that it has not sent out any REPLY message since the receipt of the last RELEASE message.

Correctness

Theorem: *Maekawa's algorithm achieves mutual exclusion.*

Proof:

- Proof is by contradiction. Suppose two sites S_i and S_j are concurrently executing the CS.
- This means site S_i received a REPLY message from all sites in R_i and concurrently site S_j was able to receive a REPLY message from all sites in R_j .
- If $R_i \cap R_j = \{S_k\}$, then site S_k must have sent REPLY messages to both S_i and S_j concurrently, which is a contradiction. \square

Performance

- Since the size of a request set is \sqrt{N} , an execution of the CS requires \sqrt{N} REQUEST, \sqrt{N} REPLY, and \sqrt{N} RELEASE messages, resulting in $3\sqrt{N}$ messages per CS execution.
- Synchronization delay in this algorithm is $2T$. This is because after a site S_i exits the CS, it first releases all the sites in R_i and then one of those sites sends a REPLY message to the next site that executes the CS.

Problem of Deadlocks

- Maekawa's algorithm can deadlock because a site is exclusively locked by other sites and requests are not prioritized by their timestamps.
- Assume three sites S_i , S_j , and S_k simultaneously invoke mutual exclusion. Suppose $R_i \cap R_j = \{S_{ij}\}$, $R_j \cap R_k = \{S_{jk}\}$, and $R_k \cap R_i = \{S_{ki}\}$.
- Consider the following scenario:
 - S_{ij} has been locked by S_i (forcing S_j to wait at S_{ij}).
 - S_{jk} has been locked by S_j (forcing S_k to wait at S_{jk}).
 - S_{ki} has been locked by S_k (forcing S_i to wait at S_{ki}).
- This state represents a deadlock involving sites S_i , S_j , and S_k .

Handling Deadlocks

Maekawa's algorithm handles deadlocks by requiring a site to yield a lock if the timestamp of its request is larger than the timestamp of some other request waiting for the same lock.

- A site suspects a deadlock (and initiates message exchanges to resolve it) whenever a higher priority request arrives and waits at a site because the site has sent a REPLY message to a lower priority request.

Deadlock handling requires three types of messages:

FAILED: A FAILED message from site S_i to site S_j indicates that S_i can not grant S_j 's request because it has currently granted permission to a site with a higher priority request.

INQUIRE: An INQUIRE message from S_i to S_j indicates that S_i would like to find out from S_j if it has succeeded in locking all the sites in its request set.

YIELD: A YIELD message from site S_i to S_j indicates that S_i is returning the permission to S_j (to yield to a higher priority request at S_j).

Maekawa's algorithm handles deadlocks as follows:

- When a REQUEST(ts, i) from site S_i blocks at site S_j because S_j has currently granted permission to site S_k , then S_j sends a FAILED(j) message to S_i if S_i 's request has lower priority. Otherwise, S_j sends an INQUIRE(j) message to site S_k .
- In response to an INQUIRE(j) message from site S_j , site S_k sends a YIELD(k) message to S_j provided S_k has received a FAILED message from a site in its request set or if it sent a YIELD to any of these sites, but has not received a new GRANT from it.
- In response to a YIELD(k) message from site S_k , site S_j assumes as if it has been released by S_k , places the request of S_k at appropriate location in the request queue, and sends a GRANT(j) to the top request's site in the queue. Maekawa's algorithm requires extra messages to handle deadlocks.

Maximum number of messages required per CS execution in this case is $5\sqrt{N}$.

Token-based algorithms

In token-based algorithms, a unique token is shared among the sites. A site is allowed to enter its CS if it possesses the token. A site holding the token can enter its CS repeatedly until it sends the token to some other site. Depending upon the way a site carries out the search for the token, there are numerous token-based algorithms. Next, we discuss two token-based mutual exclusion algorithms.

First, token-based algorithms use sequence numbers instead of timestamps. Every request for the token contains a sequence number and the sequence numbers of sites advance independently. A site increments its sequence number counter every time it makes a request for the token. (A primary function of the sequence numbers is to distinguish between old and current requests.) Second, the correctness proof of token-based algorithms, that they enforce mutual exclusion, is trivial because

an algorithm guarantees mutual exclusion so long as a site holds the token during the execution of the CS.

3.6 SUZUKI-KASAMI'S BROADCAST ALGORITHM

- If a site wants to enter the CS and it does not have the token, it broadcasts a REQUEST message for the token to all other sites.
- A site which possesses the token sends it to the requesting site upon the receipt of its REQUEST message.
- If a site receives a REQUEST message when it is executing the CS, it sends the token only after it has completed the execution of the CS.

This algorithm must efficiently address the following two design issues:

(1) How to distinguish an outdated REQUEST message from a current REQUEST message:

- Due to variable message delays, a site may receive a token request message after the corresponding request has been satisfied.
- If a site can not determine if the request corresponding to a token request has been satisfied, it may dispatch the token to a site that does not need it.
- This will not violate the correctness, however, this may seriously degrade the performance.

(2) How to determine which site has an outstanding request for the CS:

- After a site has finished the execution of the CS, it must determine what sites have an outstanding request for the CS so that the token can be dispatched to one of them.

The first issue is addressed in the following manner:

- A REQUEST message of site S_j has the form REQUEST(j, n) where n ($n=1, 2, \dots$) is a sequence number which indicates that site S_j is requesting its n^{th} CS execution.
- A site S_i keeps an array of integers $RN_i[1..N]$ where $RN_i[j]$ denotes the largest sequence number received in a REQUEST message so far from site S_j .
- When site S_i receives a REQUEST(j, n) message, it sets $RN_i[j] := \max(RN_i[j], n)$.
- When a site S_i receives a REQUEST(j, n) message, the request is outdated if $RN_i[j] > n$.

The second issue is addressed in the following manner:

- The token consists of a queue of requesting sites, Q , and an array of integers $LN[1..N]$, where $LN[j]$ is the sequence number of the request which site S_j executed most recently.
- After executing its CS, a site S_i updates $LN[i] := RN_i[i]$ to indicate that its request corresponding to sequence number $RN_i[i]$ has been executed.
- At site S_i if $RN_i[j] = LN[j] + 1$, then site S_j is currently requesting token.

The Algorithm

Requesting the critical section

- (a) If requesting site S_i does not have the token, then it increments its sequence number, $RN_i[i]$, and sends a REQUEST(i, sn) message to all other sites. (' sn ' is the updated value of $RN_i[i]$.)
- (b) When a site S_j receives this message, it sets $RN_j[i]$ to $\max(RN_j[i], sn)$. If S_j has the idle token, then it sends the token to S_i if $RN_j[i] = LN[i] + 1$.

Executing the critical section

- (c) Site S_i executes the CS after it has received the token.

Releasing the critical section

Having finished the execution of the CS, site S_i takes the following actions:

- (a) It sets $LN[i]$ element of the token array equal to $RN_i[i]$.
- (b) For every site S_j whose id is not in the token queue, it appends its id to the token queue if $RN_i[j] = LN[j] + 1$.
- (c) If the token queue is nonempty after the above update, S_i deletes the top site id from the token queue and sends the token to the site indicated by the id.

Correctness

Mutual exclusion is guaranteed because there is only one token in the system and a site holds the token during the CS execution.

Theorem: A requesting site enters the CS in finite time. Proof:

- Token request messages of a site S_i reach other sites in finite time.
- Since one of these sites will have token in finite time, site S_i 's request will be placed in the token queue in finite time.
- Since there can be at most $N - 1$ requests in front of this request in the token queue, site S_i will get the token and execute the CS in finite time.

Performance

- No message is needed and the synchronization delay is zero if a site holds the idle token at the time of its request.
- If a site does not hold the token when it makes a request, the algorithm requires N messages to obtain the token. Synchronization delay in this algorithm is 0 or T .

3.7 DEADLOCK DETECTION IN DISTRIBUTED SYSTEMS

- Deadlocks is a fundamental problem in distributed systems.
- A process may request resources in any order, which may not be known a priori and a process can request resource while holding others.
- If the sequence of the allocations of resources to the processes is not controlled, deadlocks can occur.
- A deadlock is a state where a set of processes request resources that are held by other processes in the set.

3.8 SYSTEM MODEL

- A distributed program is composed of a set of n asynchronous processes $p_1, p_2, \dots, p_i, \dots, p_n$ that communicates by message passing over the communication network.
- Without loss of generality we assume that each process is running on a different processor.
- The processors do not share a common global memory and communicate solely by passing messages over the communication network.
- There is no physical global clock in the system to which processes have instantaneous access.
- The communication medium may deliver messages out of order, messages may be lost garbled or duplicated due to timeout and retransmission, processors may fail and communication links may go down.

We make the following assumptions:

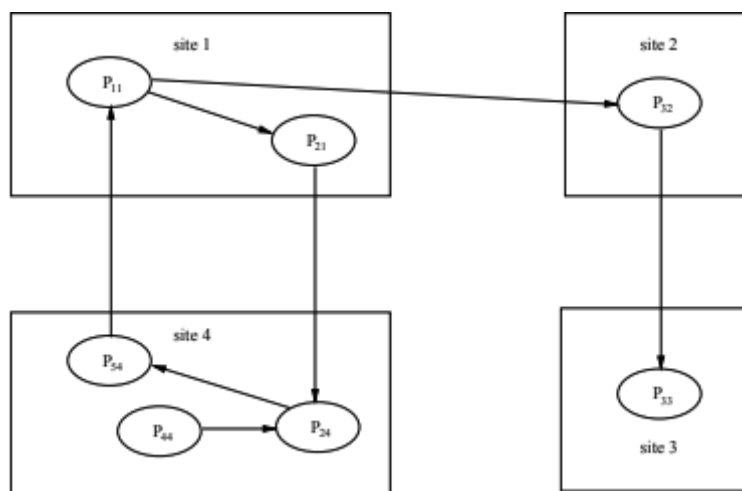
- The systems have only reusable resources.
- Processes are allowed to make only exclusive access to resources.
- There is only one copy of each resource.

Wait for Graph (WFG)

- A process can be in two states: *running* or *blocked*.
- In the running state (also called *active* 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 resource.

- The state of the system can be modeled by directed graph, called a *wait for graph* (WFG).
- In a WFG, nodes are processes and there is a directed edge from node P_1 to node P_2 if P_1 is blocked and is waiting for P_2 to release some resource.
- A system is deadlocked if and only if there exists a directed cycle or knot in the WFG.
- Figure 1 shows a WFG, where process P_{11} of site 1 has an edge to process P_{21} of site 1 and P_{32} of site 2 is waiting for a resource which is currently held by process P_{21} .
- At the same time process P_{32} is waiting on process P_{33} to release a resource.
- If P_{21} is waiting on process P_{11} , then processes P_{11} , P_{32} and P_{21} form a cycle and all the four processes are involved in a deadlock depending upon the request model.

An Example of WFG



3.9 PRILIMINARIES

Deadlock Handling Strategies

- There are three strategies for handling deadlocks, viz., deadlock prevention, deadlock avoidance, and deadlock detection.
- Handling of deadlock becomes highly complicated in distributed systems because no site has accurate knowledge of the current state of the system and because every inter-site communication involves a finite and unpredictable delay.
- Deadlock prevention is commonly achieved either by having a process acquire all the needed resources simultaneously before it begins executing or by preempting a process which holds the needed resource.

- This approach is highly inefficient and impractical in distributed systems.
- In deadlock avoidance approach to distributed systems, a resource is granted to a process if the resulting global system state is safe (note that a global state includes all the processes and resources of the distributed system).
- However, due to several problems, deadlock avoidance is impractical in distributed systems.
- Deadlock detection requires examination of the status of process-resource interactions for presence of cyclic wait.
- Deadlock detection in distributed systems seems to be the best approach to handle deadlocks in distributed systems.

Issues in Deadlock Detection

- Deadlock handling using the approach of deadlock detection entails addressing two basic issues: First, detection of existing deadlocks and second resolution of detected deadlocks.
- Detection of deadlocks involves addressing two issues: Maintenance of the WFG and searching of the WFG for the presence of cycles (or knots).

Correctness Criteria: A deadlock detection algorithm must satisfy the following two conditions:

(i) Progress (No undetected deadlocks):

- The algorithm must detect all existing deadlocks in finite time.
- In other words, after all wait-for dependencies for a deadlock have formed, the algorithm should not wait for any more events to occur to detect the deadlock.

(ii) Safety (No false deadlocks):

- The algorithm should not report deadlocks which do not exist (called *phantom or false* deadlocks).

Resolution of a Detected Deadlock

- Deadlock resolution involves breaking existing wait-for dependencies between the processes to resolve the deadlock.
- It involves rolling back one or more deadlocked processes and assigning their resources to blocked processes so that they can resume execution.

3.10 MODELS OF DEADLOCKS

Distributed systems allow several kinds of resource requests.

The Single Resource Model

- In the single resource model, a process can have at most one outstanding request for only one unit of a resource.
- Since the maximum out-degree of a node in a WFG for the single resource model can be 1, the presence of a cycle in the WFG shall indicate that there is a deadlock.

AND Model

- In the AND model, a process can request for more than one resource simultaneously and the request is satisfied only after all the requested resources are granted to the process.
- The out degree of a node in the WFG for AND model can be more than 1.
- The presence of a cycle in the WFG indicates a deadlock in the AND model.
- Since in the single-resource model, a process can have at most one outstanding request, the AND model is more general than the single-resource model.

Consider the example WFG described in the Figure 1.

- P_{11} has two outstanding resource requests. In case of the AND model, P_{11} shall become active from idle state only after both the resources are granted.
- There is a cycle $P_{11} \rightarrow P_{21} \rightarrow P_{24} \rightarrow P_{54} \rightarrow P_{11}$ which corresponds to a deadlock situation.
- That is, a process may not be a part of a cycle, it can still be deadlocked. Consider process P_{44} in Figure 1.
- It is not a part of any cycle but is still deadlocked as it is dependent on P_{24} which is deadlocked.

OR Model

- In the OR model, a process can make a request for numerous resources simultaneously and the request is satisfied if any one of the requested resources is granted.
- Presence of a cycle in the WFG of an OR model does not imply a deadlock in the OR model.
- Consider example in Figure 1: If all nodes are OR nodes, then process P_{11} is not deadlocked because once process P_{33} releases its resources, P_{32} shall become active as one of its requests is satisfied.
- After P_{32} finishes execution and releases its resources, process P_{11} can continue with its processing.
- In the OR model, the presence of a knot indicates a deadlock.

AND – OR Model

- A generalization of the previous two models (OR model and AND model) is the AND-OR model.
- In the AND-OR model, a request may specify any combination of *and* and *or* in the resource request.
- For example, in the AND-OR model, a request for multiple resources can be of the form x *and* (y *or* z).
- To detect the presence of deadlocks in such a model, there is no familiar construct of graph theory using WFG.
- Since a deadlock is a stable property, a deadlock in the AND-OR model can be detected by repeated application of the test for OR-model deadlock.

3.10.5 The $\binom{p}{q}$ model (called the P-out-of-Q model)

- The $\binom{p}{q}$ model (called the P-out-of-Q model) allows a request to obtain any k available resources from a pool of n resources.
- It has the same expressive power as the AND-OR model.
- However, $\binom{p}{q}$ model lends itself to a much more compact formation of a request.
- Every request in the $\binom{p}{q}$ model can be expressed in the AND-OR model and vice-versa.
- Note that AND requests for p resources can be stated as $\binom{p}{p}$ and OR requests for p resources can be stated as $\binom{p}{1}$.

3.10.6 Unrestricted model

- In the unrestricted model, no assumptions are made regarding the underlying structure of resource requests.
- Only one assumption that the deadlock is stable is made and hence it is the most general model.
- This model helps separate concerns: Concerns about properties of the problem (stability and deadlock) are separated from underlying distributed systems computations (e.g., message passing versus synchronous communication).

3.11 KNAPP'S CLASSIFICATION

Distributed deadlock detection algorithms can be divided into four classes:

- Path-Pushing
- Edge-Chasing
- Diffusion Computation
- Global State Detection.

Path-Pushing Algorithms

- In path-pushing algorithms, distributed deadlocks are detected by maintaining an explicit global WFG.
- The basic idea is to build a global WFG for each site of the distributed system.
- In this class of algorithms, at each site whenever deadlock computation is performed, it sends its local WFG to all the neighboring sites.
- After the local data structure of each site is updated, this updated WFG is then passed along to other sites, and the procedure is repeated until some site has a sufficiently complete picture of the global state to announce deadlock or to establish that no deadlocks are present.
- This feature of sending around the paths of global WFG has led to the term path-pushing algorithms.

Edge-Chasing Algorithms

- In an edge-chasing algorithm, the presence of a cycle in a distributed graph structure is verified by propagating special messages called probes, along the edges of the graph.
- These probe messages are different than the request and reply messages.
- The formation of cycle can be deleted by a site if it receives the matching probe sent by it previously.
- Whenever a process that is executing receives a probe message, it discards this message and continues.
- Only blocked processes propagate probe messages along their outgoing edges.
- Main advantage of edge-chasing algorithms is that probes are fixed size messages which is normally very short.

Diffusing Computations Based Algorithms

- In *diffusion computation* based distributed deadlock detection algorithms, deadlock detection computation is diffused through the WFG of the system.
- These algorithms make use of echo algorithms to detect deadlocks.
- This computation is superimposed on the underlying distributed computation. If this computation terminates, the initiator declares a deadlock.
- To detect a deadlock, a process sends out query messages along all the outgoing edges in the WFG.
- These queries are successively propagated (i.e., diffused) through the edges of the WFG.

- When a blocked process receives first query message for a particular deadlock detection initiation, it does not send a reply message until it has received a reply message for every query it sent.
- For all subsequent queries for this deadlock detection initiation, it immediately sends back a reply message.
- The initiator of a deadlock detection detects a deadlock when it receives reply for every query it had sent out.

Global state detection based deadlock detection algorithms

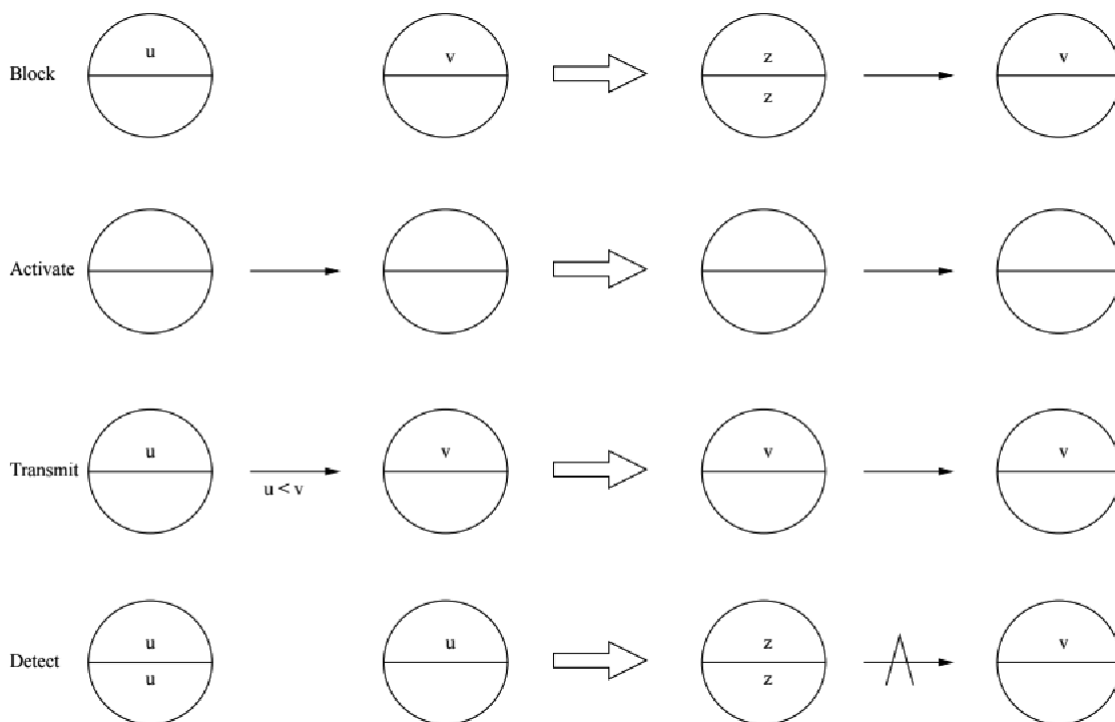
- Global state detection based deadlock detection algorithms exploit the following facts:
- A consistent snapshot of a distributed system can be obtained without freezing the underlying computation and
- If a stable property holds in the system before the snapshot collection is initiated, this property will still hold in the snapshot.
- Therefore, distributed deadlocks can be detected by taking a snapshot of the system and examining it for the condition of a deadlock.

3.12 MITCHELL AND MERRITT'S ALGORITHM FOR THE SINGLE-RESOURCE MODEL

- Belongs to the class of edge-chasing algorithms where probes are sent in opposite direction of the edges of WFG.
- When a probe initiated by a process comes back to it, the process declares deadlock.
- Only one process in a cycle detects the deadlock. This simplifies the deadlock resolution – this process can abort itself to resolve the deadlock.
- Each node of the WFG has two local variables, called labels:
 - **a private label, which is unique to the node at all times, though it is not constant, and**
 - **a public label, which can be read by other processes and which may not be unique.**
- Each process is represented as u/v where u and v are the public and private labels, respectively.
- Initially, private and public labels are equal for each process.
- A global WFG is maintained and it defines the entire state of the system.
- The algorithm is defined by the four state transitions shown in Figure 2, where $z = \text{inc}(u, v)$, and $\text{inc}(u, v)$ yields a unique label greater than both u and v labels that are not shown do not change.

- Block creates an edge in the WFG.
- Two messages are needed, one resource request and one message back to the blocked process to inform it of the public label of the process it is waiting for.
- Activate denotes that a process has acquired the resource from the process it was waiting for.
- Transmit propagates larger labels in the opposite direction of the edges by sending a probe message.

The four possible state transitions



- Whenever a process receives a probe which is less than its public label, then it simply ignores that probe.
- Detect means that the probe with the private label of some process has returned to it, indicating a deadlock.
- The above algorithm can be easily extended to include priorities where whenever a deadlock occurs, the lowest priority process gets aborted.

Message Complexity:

If we assume that a deadlock persists long enough to be detected, the worst-case complexity of the algorithm is $s(s - 1)/2$ Transmit steps, where s is the number of processes in the cycle.

3.13 CHANDY-MISRA-HAAS'S FOR AND MODEL

- Chandy-Misra-Haas's distributed deadlock detection algorithm for AND model is based on edge-chasing.
- The algorithm uses a special message called *probe*, which is a triplet (i, j, k) , denoting that it belongs to a deadlock detection initiated for process P_i and it is being sent by the home site of process P_j to the home site of process P_k .
- A probe message travels along the edges of the global WFG graph, and a deadlock is detected when a probe message returns to the process that initiated it.
 - A process P_j is said to be *dependent* on another process P_k if there exists a sequence of processes $P_j, P_{i1}, P_{i2}, \dots, P_{im}, P_k$ such that each process except P_k in the sequence is blocked and each process, except the P_j , holds a resource for which the previous process in the sequence is waiting.
 - Process P_j is said to be *locally dependent* upon process P_k if P_j is dependent upon P_k and both the processes are on the same site.
 - Data Structures
 - Each process P_i maintains a boolean array, *dependent_i*, where *dependent_i(j)* is true only if P_i knows that P_j is dependent on it.
 - Initially, *dependent_i(j)* is false for all i and j .

Algorithm

```

if  $P_i$  is locally dependent on itself
  then declare a deadlock
  else for all  $P_j$  and  $P_k$  such that
    (a)  $P_i$  is locally dependent upon  $P_j$ , and
    (b)  $P_j$  is waiting on  $P_k$ , and
    (c)  $P_j$  and  $P_k$  are on different sites,
      send a probe  $(i, j, k)$  to the home site of  $P_k$ 

```

On the receipt of a probe (i, j, k) , the site takes the following actions:

```

if
  (d)  $P_k$  is blocked, and
  (e) dependentk i is false, and
  (f)  $P_k$  has not replied to all requests  $P_j$ , then
  begin
    dependentk i = true; if  $k = i$ 

```

then declare that P_i is deadlocked else
for all P_m and P_n such that

- (a) P_k is locally dependent upon P_m , and
- (b) P_m is waiting on P_n , and
- (c) P_m and P_n are on different sites, send a probe (i, m, n) to the home site of P_n

end.

- A probe message is continuously circulated along the edges of the global WFG graph and a deadlock is detected when a probe message returns to its initiating process.

Performance Analysis

- One probe message (per deadlock detection initiation) is sent on every edge of the WFG which that two sites.
- Thus, the algorithm exchanges at most $m(n - 1)/2$ messages to detect a deadlock that involves m processes
- and that spans over n sites.
- The size of messages is fixed and is very small (only 3 integer words).
- Delay in detecting a deadlock is $O(n)$.

3.14 CHANDY-MISRA-HAAS DISTRIBUTED DEADLOCK DETECTION ALGORITHM FOR OR MODEL

Chandy-Misra-Haas distributed deadlock detection algorithm for OR model is based on the approach of diffusion-computation.

- A blocked process determines if it is deadlocked by initiating a diffusion computation.
- **Two types of messages are used in a diffusion computation:**
- **query(i, j, k) and reply(i, j, k)**, denoting that they belong to a diffusion computation initiated by a process P_i and are being sent from process P_j to process P_k .
- A blocked process initiates deadlock detection by sending query messages to all processes in its dependent set.
- If an active process receives a query or reply message, it discards it.
- When a blocked process P_k receives a query(i, j, k) message, it takes the following actions:
 - If this is the first query message received by P_k for the deadlock detection initiated by P_i (called the *engaging query*), then it propagates the query to all the processes in its dependent set and sets a local variable $num_k(i)$ to the number of query messages sent.
 - If this is not the engaging query, then P_k returns a reply message to it immediately provided P_k has been continuously blocked since it received the corresponding engaging query. Otherwise, it discards the query.

- Process P_k maintains a boolean variable $wait_k(i)$ that denotes the fact that it has been continuously blocked since it received the last engaging query from process P_i .
- When a blocked process P_k receives a $reply(i, j, k)$ message, it decrements $num_k(i)$ only if $wait_k(i)$ holds.
- A process sends a reply message in response to an engaging query only after it has received a reply to every query message it had sent out for this engaging query.
- The initiator process detects a deadlock when it receives reply messages to all the query messages it had sent out.

The algorithm works as follows:

Initiate a diffusion computation for a blocked process P_i :
 send query(i, i, j) to all processes P_j in the dependent set
 DS_i of P_i ;
 $num_i(i) := |DS_i|$; $wait_i(i) := true$;

When a blocked process P_k receives a query(i, j, k):
 if this is the engaging query for process P_i
 then send query(i, k, m) to all P_m in its dependent set DS_k ;
 $num_k(i) := |DS_k|$; $wait_k(i) := true$
 else if $wait_k(i)$ then send a $reply(i, k, j)$ to P_j .

When a process P_k receives a reply(i, j, k):
 if $wait_k(i)$
 then begin
 $num_k(i) := num_k(i) - 1$;

if $num_k(i) = 0$
 then if $i = k$ then **declare a deadlock**
 else send $reply(i, k, m)$ to the process P_m
 which sent the engaging query.

- In practice, several diffusion computations may be initiated for a process (A diffusion computation is initiated every time the process gets blocked), but, at any time only one diffusion computation is current for any process.
- However, messages for outdated diffusion computations may still be in transit.
- The current diffusion computation can be distinguished from outdated ones by using sequence numbers.

Performance Analysis

For every deadlock detection, the algorithm exchanges e query messages and e reply messages, where $e = n(n-1)$ is the number of edges.

1. Explain the Lamport's algorithm.
2. Discuss Ricart-Agrawala algorithm.
3. Explain Maekawa's algorithm.
4. Explain Suzuki-Kasami's broadcast algorithm.
5. How to detect Deadlock in distributed systems. Explain the system model
6. Discuss the Models of deadlocks

7. Explain the Knapp's classification.
8. Discuss the Algorithm for the single resource model. (**MITCHELL AND MERRITT'S ALGORITHM**)
9. Discuss the Algorithm the AND model . (**CHANDY-MISRA-HAAS'S**)
10. Discuss the Algorithm OR model. (CHANDY-MISRA-HAAS)