

06/02/2023

Module - 1

Syllabus

Distributed systems - definitions, relation to computer systems components, motivation, primitives for distributed communication, design issues, challenges and applications, a model of distributed computation, distributed programs, model of distributed execution, model of communication networks, global state of distributed systems, cuts of distributed computation, past and future comes to an end, models of process communication.

⇒ Distributed Systems.

A distributed system is a collection of independent entities that cooperate to solve a problem that cannot be individually solved.

A distributed system can be characterized as a collection of mostly autonomous processes communicating over a communication network and have following features:

- * No common physical clock
- * No shared memory

papergrid

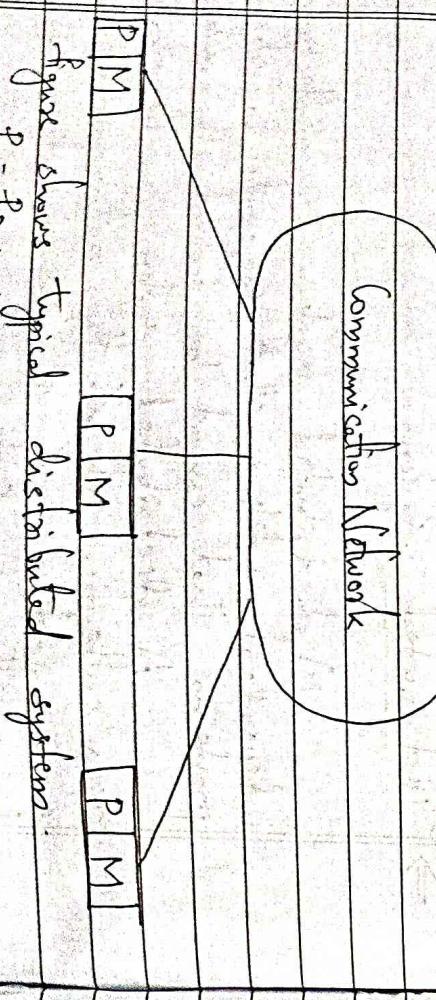
Date: / /

- * Geographical Separation
 - * Autonomy and heterogeneity

→ Distributed Computation

Field of computer science that deals with the study of distributed systems is called as distributed computation.

Relation to Computer Systems Responses -



Shows typical distributed system
P - Processor
M - Memory.

papergrid

Date: / /

- * Each computer has a memory processing unit and computers are connected by communication network
 - * A distributed software is also termed as

- A distributed execution is a execution of processes across the distributed systems to collaboratively achieve a common goal.

2. *Scutellaria* *lutea* L.

- * a two
Distributed systems uses layered architecture
to break down complexity of system designs.
The middleware is the distributed software that
use distributed systems.

Fig: Interaction of software components at each process

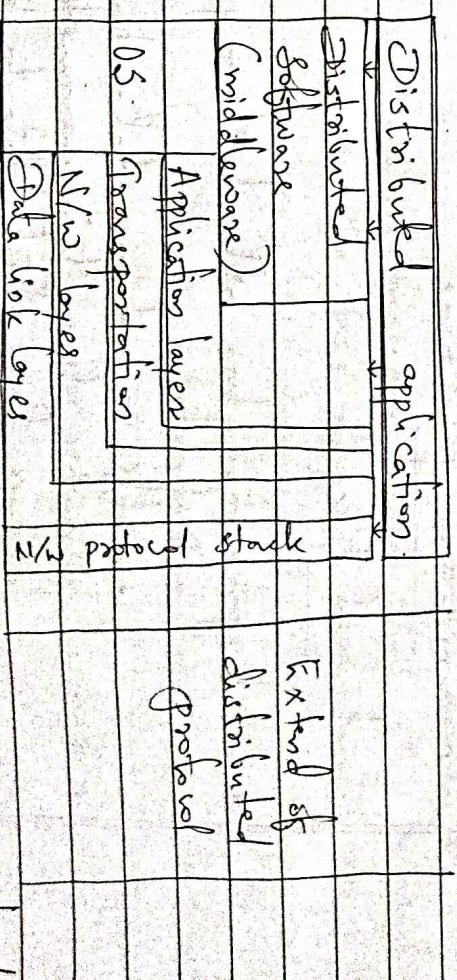


Fig. shows the interaction of this software base systems components with each other.

papergrid

Date: / /

to be received.

→ Motivation

- * Globally distributed computing, e.g.: money transfer
- * Resource sharing
- * Access to geographically remote data and resources
- * Enhance reliability
- * Increased performance / cost ratio
- * Scalability
- * Modularity and incremental expandability

⇒ Primitives for distributive computation

Blocking / Non-blocking / Synchronous / Asynchronous primitives.

→ Communication primitives.

Send and Receive

* Send primitive has at least two parameters:

- Destinates
- Buffer - containing data to be sent.

* Receive primitive has two parameters :

- Source from which data to be received
- We use a buffer into which data is

→ Synchronous Primitive

A send or receive primitive is synchronous if both send and receive handshake with each other. Processing for send primitive complete only after invoking processor. Same that other corresponding receive primitive has also been invoked and receive operation has been completed. The processing for receive primitive complete when data to be received is copied into receiver user buffer.

Asynchronous Primitive

A send primitive is said to be asynchronous

papergrid

Date: / /

papergrid

Date: / /

It waiting until the invoking process after data items to be sent have been copied out of user specified buffer.

→ Blocking primitive

A primitive is blocking if control returns to invoking process after processing for primitive complete.

→ Non-blocking primitive

A primitive is non-blocking if control returns back to invoking process immediately after invocation even though operations has not completed. For non-blocking send control released process over before data is copied out of user buffer. For non-blocking receive control is released from process before data may have arrived from buffer.

For non-blocking primitive a return parameter on the primitive call between a system generated handle which can be later used to check status of the completion of call. The process can block for the completion of call in two ways:

(i) It can keep checking if no handle has been first or packed.
(ii) It can issue a wait with a list of handles as parameter.

papergrid

Date: / /

Send | Receive

Synchronous	Asynchronous	Synchronous
Send	Send	Receive

Blocking	Non-Blocking	Blocking	Non-Blocking	Blocking	Non-Blocking
Synchronous	Synchronous	Asynchronous	Asynchronous	Synchronous	Synchronous
Send	Send	Receive	Receive	Send	Send

Duration to copy data from or to user buffer.

Duration in which process issuing send receive primitive is blocked.

S - send primitive issued

S-C - processing for send complete.

R - receive primitive issued

R-C - processing for receive completes.

P - The completion of previously initiated blocking



W - process may issue wait to check completion of non-blocking operation.

~~8/02/2023~~ Blocking synchronous send, Blocking synchronous receive.

Blocking synchronous send

The data is copied from blocking buffer to kernel buffer and then send over via After data is copied to receiver buffer system and receive call has been issued. An acknowledgement back to sender comes back to the process that issued the send operation and complete the send.

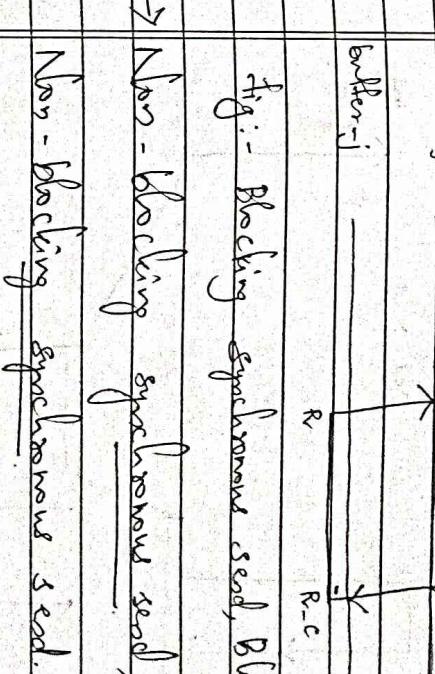


Fig:- Blocking synchronous send, Blocking synchronous receive

→ Non-blocking synchronous send, non-blocking receive

Non-blocking synchronous send.

Blocking receive

The receive call blocks until the data expected arrives and is returned to the specified user buffer. This control is returned to user buffer.

Non-blocking receive

The receive call will cause the kernel to

control return back to invoking process as soon as copy of data from user buffer to kernel buffer is initiated

register the call and return handle of the user process can later check location that the user process can locate check for the completion of non-blocking receive operation. This location gets posted by the kernel after the expected data arrives and is copied to user specified buffer. The user process can check for the completion of non-blocking receive by invoking wait operation on sevices handle.

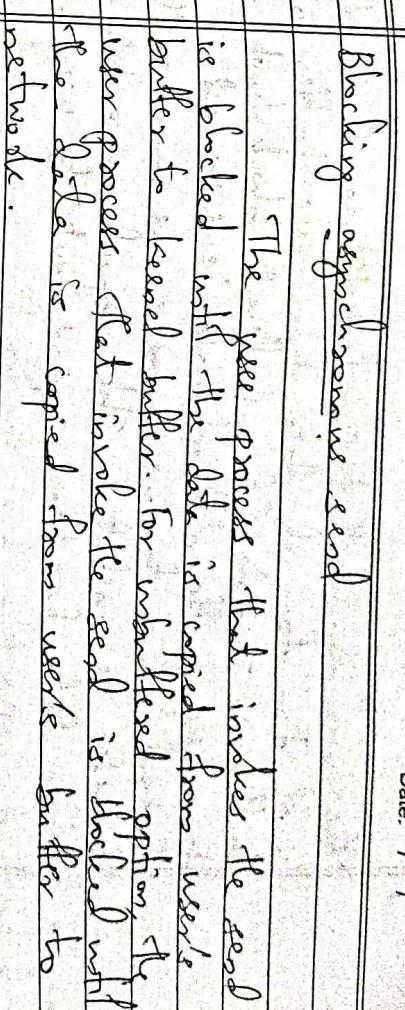


Fig:- Blocking asynchronous send.

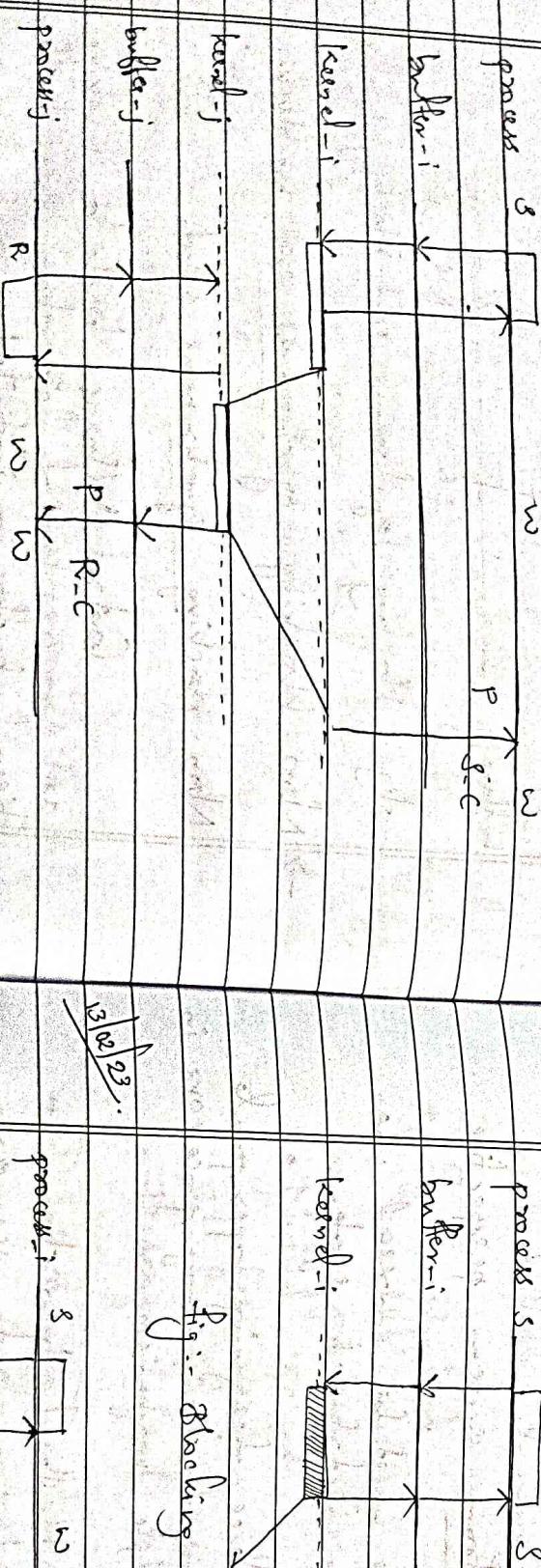


Fig :- Non-Blocking Asynchronous Send,
Non-Blocking Asynchronous Receive.

The user process that invokes the send is blocked until transfer of data from user buffer to local buffer is initiated (for transferred block with transfer buffer is initiated). User buffer is copied to network is initiated. Local server to user process as soon as this command has been copied out of user buffer.

\rightarrow Design issues and challenges:

Design challenges from systems perspective

The following questions must be addressed when designing and building distributed systems.

* Communications

This stage involves designing appropriate mechanisms for communication among processors in network. Some example mechanisms are :-

. RPC (Remote

. RMI (Remote Method Invocation)

- Message oriented communication vs stores oriented communication.

* Processes

Some of the issues involved are

- Management of processes and threads - at client/server
- Code migration
- Design of software and mobile engine

* Naming

Device easy to use and robust scheme for names, identities and addresses is essential for locating resources and processes in a transparent and scalable manner.

* Synchronization

Mechanisms for synchronization or coordination among processes are essential.

- Mutual exclusion
- Leader election
- Synchronizing physical clocks
- Designing logical clock est

- Data storage and access schemes for data storage and access
- Data is a fast and scalable manner across network are important for efficiency. File systems designs have to be considered in setting of distributed system.

papergrid

Date: / /

connected

Date: / /

papergrid

Date: / /

* Consistency and replication.

To provide fast access of data and to provide scalability replication of data is highly desirable. This leads to issues of managing replicas and dealing with consistency among replicas/cohorts in distributed setting.

* Fault Tolerance

Fault tolerance requires mainly to correct and efficient operation in case of failure in links, nodes and processes. Some mechanisms to provide fault tolerance like tolerance zone:

- Process resilience
- Reliable communication
- Distributed constraint
- Check pointing and recovery.

* Security

Distributed systems security involves various aspect of crypto graphy, secure channels, access control, key management etc.

* API and transparency

API for communication and other specialized services is important for ease of use and wider

* Adoption of distributed systems services by users.

Transparency hide implementation policies from user. Different types of transparencies:-

- Access transparency
- If hides difference in date representation of different system
- Location transparency
- Makes locations of resources transparent to user

- Migration transparency.
- Allows relocating resources without changing names.

- Relocation transparency Ability to relocate resources as they are being accessed.

- Replication transparency Does not let user become aware of a replication.

- Consistency transparency Deals with masking concurrent use of shared resources from the user.

- Failure transparency
- Allowing user and application programs to complete their task despite the failure of hardware or software component.

* Scalability and modularity
A system is described as scalable if it will remain effective when there is a significant increase in number of resources and number of users. Various techniques such as

- Replication
- Load balancing
- Decentralized processing
- Help to achieve scalability.

→ Model

A model of distributed computation

(i) Message send event : A send event changes the state of process that sends the message and

(ii) Message receive event : A receive event changes the state of channel on which message is

sent.

(iii) Message receive event : A receive event changes the state of process that receives the message and state of channel on which message is received.

Distributed program

A distributed program consist of set of asynchronous processes. P_1, P_2, \dots, P_n that communicate by message passing over communication links. We assume each process runs on different processor. Processor does not share global memory.

C_{ij} denotes channel from process P_i to process P_j , may denote message send by P_i to P_j . These process do not share global memory. Communication delay is finite. Process executions and message transfers are asynchronous. The global state of distributed computation consist of state

of channel or computation channel.

Let e_i^x denote n^{th} event at process P_i . Then send(m) and receive(m) denotes send and receive events. The

occurrence of events changes state of respective process and channels.

The events at a process are linearly ordered by their order of occurrence. The execution of process P_i produces a set of events $e_i^1, e_i^2, \dots, e_i^n$ and it is denoted

by H_i
 $H_i = (b_i \rightarrow i)$

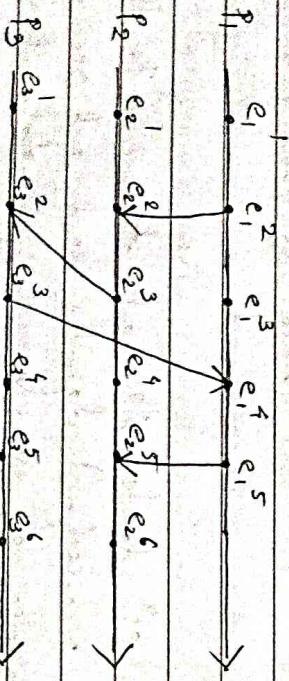
b_i - set of events produced by P_i
 \rightarrow i - binary relation defines linear order on these events, relation express causal dependencies

events of P_i .

For every message m that is exchanged between 2 process, we have

$send(m) \rightarrow msg\ rec(m)$

Relation \rightarrow msg defines causal dependencies between pairs of corresponding send and receive event. The figure shows space-time diagrams of distributed execution involving 3 process.



Time \rightarrow

For any 2 event $e_i \rightarrow e_j$ denotes the fact that event e_j is directly or transitively caused by e_i .

Models of communication network

These are several models of the service provided by communication network.

(i) FIFO Model

In FIFO model, each channel act as first in first out message queue and thus message ordering is preserved.

(ii) Non-FIFO Model

In Non-FIFO model, channel act like a set in which write process act messages and receive process and remove messages in random model.

(iii) Causal Ordering.

This model is based on Lamport's happens before relation. A system that supports causal ordering need satisfies the following properties.

For any two messages m_{ij} and m_{kl}

$$\text{send}(m_{ij}) \rightarrow \text{send}(m_{kl})$$

$$\text{rec}(m_{ij}) \rightarrow \text{rec}(m_{kl})$$

This property ensures that those causally related messages destined to same destination are connected.

CO C FIFO C No, FIFO

→ Models of process communication.

Two basic models of process communication.

(i) Synchronous communication model.

It is a blocking type where on message

sent, the sender process blocks until the message has been received by receiver process. The

sender process executes only after it receives that receiver process has accepted the

message. Thus sender and receiver must synchronize to exchange message. The synchronous is simpler to handle and implement. poor performance.

(ii) Asynchronous communications model.

* It is a non-blocking type.

* Sender and receiver don't synchronize to exchange message. After sending a message the sender process doesn't wait for message to be delivered to receive a process. The message is buffered by the system and is delivered to receiver process when it is ready to accept a message.

* A buffer overflow may occur. It provides light

parallelism.

* It requires more complex buffer management.

* Difficult to design, verify and implement distributed algorithms for asynchronous communication.

→ Cut off distributed computing

for the space-time diagram, of distributed

computation zigzag line joining 1 arbitrary point on each process line is termed as cut in computation. Such line slices space-time diagrams and thus set of events is distributed computation into past and future PAST and FUTURE.

The PAST contains all events to left of cut and FUTURE contains all events to right of cut.

Every cut corresponds to a global state and every global state can be graphically represented as cut in computation space-time diagram.

→ Two types of cut

(i) Consistent cuts

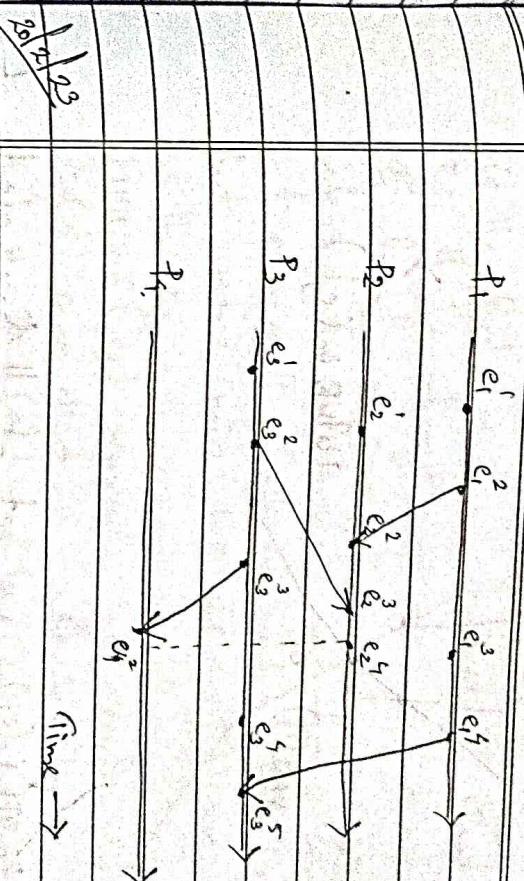
A consistent global state corresponds to a cut so which every message received in the past of cut was send in the PAST of cut. All messages that goes from PAST to FUTURE are in consistent global state cut.

$e_i \rightarrow e_j$
and all information available at e_i could be made accessible at e_j . All such events e_i belong to PAST of e_j .

e_j (ej) do not all events in the PAST of e_j is the computation ($H \rightarrow$) then

$$PAST(e_j) = \{e_i | e_i \in H, e_i \rightarrow e_j\}$$

If a message changes the consistency of cut from FUTURE to PAST



$\max(\text{PAST}(e_j))$
 $\min(\text{FUTURE}(e_j))$

is the first event in process P_i . This is affected by e_j .

P_i is always message receiving event.

If denotes a consistent cut in computation. It is referred to as interface of future view. All events at process P_i that occurred after $\max(\text{PAST}(e_j))$ but before $\min(\text{FUTURE}(e_j))$ are concerned with e_j .

$\text{FUTURE}(e_j)$ represents all events in the FUTURE light cone that is affected by e_j . Set $\text{H-PAST}(e_j) - \text{Future}(e_j)$ are concurrent with each other.

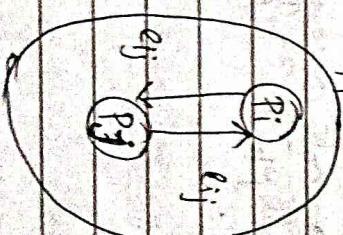
Global State of distributed systems is a collect of local state of its components namely processes and communication channel. The state of process at any time is defined by context of processor register, processor stack (local memory) increase in local context of distributed application.

FUTURE of an event e_j denoted by $\text{FUTURE}(e_j)$ contains all events e_i that are causally affected by e_j .

for a computation, ($H \rightarrow$). FUTURE of e_j is defined as

$$\text{FUTURE}(e_j) = \{ e_i | e_i \in H, e_j \rightarrow e_i \}$$

$\text{FUTURE}(e_j)$ is set of those events of $\text{FUTURE}(e_j)$ that are in process P_i and $\min(\text{FUTURE}(e_j))$



e_j is global state.

27/2/23
Page

MODULE - 2

Leader Election Algorithms

State of channel is given by set of message in transit to the channel. The sequence of event changes state of respective processes and channels thus causes transitions in global system state. Let L_i denote state of process P_i after the occurrence of event e_i^n and before event e_{i+1}^n .

L_i^n = initial state of process P_i

$\text{send}(m) \leq L_i^n$

$e_i^n \neq L_i^n$

$\exists j \quad 1 \leq j < n$

$e_j^n = \text{send}(m)$

$\forall j \quad 1 \leq j \leq n \quad e_j^n \neq \text{rec}(m)$

$S_{ij}^{n,y} \text{ denote state of channel } C_{ij} \text{ if }$

$S_{ij}^{n,y} = \{ w_{ij} \mid \text{send}(w_{ij}) \leq L_i^n \wedge \text{rec}(w_{ij}) \leq L_j^n \}$

State of channel depends upon state of processes involved. Thus channel state $S_{ij}^{n,y}$ denotes all messages that P_i send upto event e_i^n and which process P_j had not received until event e_j^n

$C_n = \{ U_i L_i \cup_j S_{ij}^{n,y} \}_{i=1}^n$

→ Issues in consistency of local states

- * If a global physical clock available the following simple procedure could be used to record a consistent global

Global state of distributed systems is a collection of local state and its consistency. Recording of global state of distributed system is an important paradigm and it find application in several aspect of distributed system design. For example in defining of stable properties such as deadlock and termination.

Global state of system is examined for certain properties:

- (i) for fairness property
- (ii) for the voting list and total software

Presence of global memory necessitates way of getting coherent and complete view of system based on local state of individual process.

snapshot of distributed system. So this algorithm snapshot collector decides a future time at which snapshot is to be taken and broadcast this time to every processes. All process take their local snapshots at that instance in the global time.

A global physical clock is not available in distributed system and the following two issues need to be addressed in reasoning of a consistent global snapshot of distributed system.

- How to distinguish the message to be recorded in the snapshot (in a channel state space state) from those node to be recorded.
- How to determine the instant when a process takes its snapshot.

SnapShot Algorithm for FIFO channels

Chandy-Lamport Algorithms

- * This algorithm uses a control message call a marker.
- After a right has recorded its snapshot, it sends a marker along all of its outgoing channels before sending out any more messages. Since

channels are FIFO, a marker separates the messages in the channel into those to be included in the snapshot (channel state/process state) from those not to be recorded in the snapshot.

Chandy-Lamport algorithms

- Marker sending rule for process p_i .
- 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 messages along C .

Marker receiving rule for process p_j .
On receiving a 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 sending rule".

else

Record the state of C as the set of messages received along C after p_j 's state was recorded and before p_j received the marker along C .

Election algorithms choose a process from group of processor to act as a coordinator. If the coordinator process crashes due to some reasons, then a new coordinator is elected as other processor. Election algorithm basically determines where a new copy of coordinator should be selected.

Election algorithm assumes that every active process in the system has a unique priority number. The process with highest priority will be chosen as a new coordinator.

Hence, when a coordinator fails, this algorithm checks what active process which has highest priority number. Then this number is send to every active process in the distributed system.

We have two election algorithms for two different configurations of distributed system.

1. The Bully Algorithm
2. The Ring Algorithm.

This algorithm was proposed by Chandra - Toueg. When the process notices that the coordinator is no longer responding to request, it initiates an election. A process, P, holds an election as follows:

- (i) P sends an ELECTION message to all processes with higher number.
- (ii) If no one responds, P wins the election and becomes the coordinator.
- (iii) If one of the higher-ups answers, it takes over P's job is done.

Distributed algorithms require one process to act as a coordinator or initiator. To decide which process becomes the coordinator different types of algorithms are used.

Election algorithms are meant for electing a coordinator process from among the currently running processes so such a manner that at any instance of

time there is a single coordinator for all processes in the systems.

The goal of an election algorithm is to ensure that when an election starts it concludes with all the processes agreeing on who the coordinator should be. Therefore, unless initiated, an election algorithm basically finds out which of the currently active processes has the highest priority number and then informs this to all other active processes.

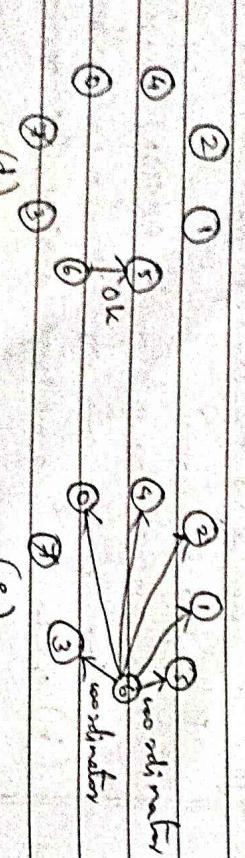
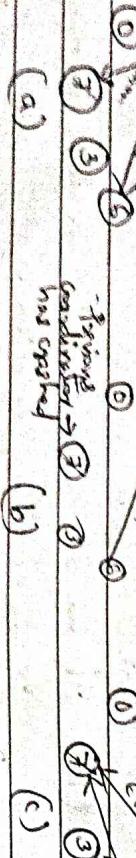
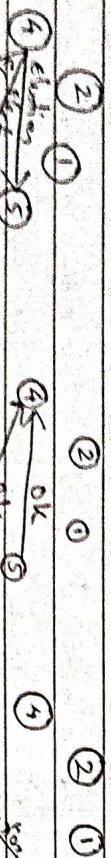
Bully Algorithm

Ring Algorithms

OK message back to the sender to indicate that he is elected and will take over. The receiver then becomes elected, unless it is already holding one.

(c) All process give up except one that is the new coordinator. It announces its victory by sending all process a message telling them that starting immediately, it is the new coordinator.

(d) A process that was previously down comes back up. It looks at election. If it happens to be highest numbered process currently running, it will send messages to all the other coordinators who are still up and tells them the coordinates of the job. That's the "biggest guy in town always wins" rule. He can name "Bully algorithm".



Ring Algorithms

This algorithm uses a ring for its election but does not use any tokens. In this algorithm, it is assumed that the processes are physically or logically ordered so that each processor knows its successor.

1) When any process notices that a coordinator is not functioning, it builds an ELECTION message containing its own process number and sends the message to its successor. If the successor is down, the sender skips over the successor and goes to the next process along the ring until a process is located.

A ring based election algorithm

Call process P_i has a communication channel to the next process in the ring $P(i+1)$ in. All messages are sent clockwise around the ring. The goal of this algorithm is to elect a single process called the coordinator. Initially, every process is marked as a non-participant in an election. Any process can begin an election, the

Termination Detection

- Starting process marks itself as a participant and place its identifier in a message to its neighbours.
2. A process receives a message and compare it with its own. If the arrived identifier is larger, it passes on the message.
 3. If arrived identifier is smaller and receiver is not a participant, substitute its own identifier in the message and forward it. It does not forward the message if it is already a participant.
 4. On forwarding of any case, the process marks itself as a participant.
 5. If the received identifier is that of receiver itself then this process identifier must be the greatest and it becomes the coordinator.

The coordinator marks itself as non-participant and it becomes the coordinator.

6. The coordinator marks itself as non-participant and delete; and seize the deleted message to its neighbours echoing its id.
7. When a process receives deleted message, it marks itself as non-participant and sets its variable deleted; and forwards the message.

A termination detection algorithm must ensure the following:

- (i) Execution of termination detection algorithms can't indefinitely delay underline computation i.e. execution of termination detection algorithms must not freeze underline algorithm.
- (ii) Termination detection algorithm must not require addition of new communication channels by processes.

A distributed problem in distributed system is to be determined if a distributed computers has terminated.

A distributed computation is considered to be globally terminated if every process is locally terminated and there is no message is transmit b/w any processes.

A locally terminated state is a state to which process has finished its computations and will not restart any action unless it receives a message.

Two distributed computations taking place in distributed system are

1. Undelived computation :- basic message
2. Termination detection algorithm - control message.

→ System model of distributed computing

* Communication is asynchronous.

* Message send over same communication channel

* May not obey FIFO ordering - following three

A distributed computation has following three

tenants

(a) Process can be in one or two states either

(i) Active (Busy) - unless it is doing

local computation.

(ii) Idle (Passive) - when process has tempor-

arily finish execution of local computation and

activated only on the receipt of message

from other process.

(b) An active process can become idle at any time (on

finishing local computation & process all received

messages).

(c) An idle process can become active only on the

receipt of message from another process.

(d) Only active processes can send messages.

(e) A message can be received by a process when

The process is in active or idle state.

(f) The sending of a message and receipt of a message occur as atomic actions.

Assignment topic

6 | 3 | 23

Version time

⇒ Logical time

→ A framework for a system of logical clocks.

* A system of logical clock consists of the time domain T and logical clock, C .

* Elements of T form a partially ordered set over a relation \leq . This relation is usually called happens before or causal ordering.

* Distributed systems may have no physically synchronous global clock. So a logical clock allows global ordering on events from different processes in such systems.

The logical clock C is a function that map an event in a distributed system to an element in the domain T , denoted as $C(e)$ and called timestamp

of e and is defined as
 $C:H \rightarrow T$

such that following property is satisfied.

such that following property is satisfied.
 for two events e_i and e_j , $c(e_i) < c(e_j)$, this is called
 $e_i \rightarrow e_j \Rightarrow c(e_i) < c(e_j)$, this is called
 local consistency condition.

when T and C satisfy following conditions for two
 events $e_i \neq e_j$
 $c(e_i) < c(e_j)$

The system of clock is said to be stronger
 consistent.

\rightarrow Implementing logical clock

It requires addressing two issues:
 (i) Data structure local to every process to represent
 logical time.

(ii) Protocol to update data structure to ensure
 consistency condition.

Each process P_i maintains data structure that
 allows following two capabilities.

(1) A local logical clock denoted by L_i , that
 helps process P_i measure its own progress.

(2) A logical global clock denoted by G_i , that
 is representation of process P_i 's local view of

logical global time. This allows the process to assign
 consistent time stamp to its local events.

The protocol consists of following rules:

R₁: The rule governs how process updates its
 updated by a process when it executes an
 event [read, receive or internal]
 R₂: This rule governs how process updates its
 global logical clock to update its view of
 global time and global progress.

Scales time

This representation was proposed by Lamport
 as an attempt to totally order events in distributed
 systems.

Time domain in this representation is the set of
 non-negative integers.

Rules R₁ & R₂ to update clock are as follows

$$C_i = C_i + d \quad (d > 0)$$

Every time R₁ is executed, it can have different
 values. This value may be application dependent.
 R₂: Each message piggybacked. The clock value of its
 receiver sender at sending time. When a

process P_i will exchange with the following message if executing the following actions.

$$C_{\text{message}} = \max(C_1, C_{\text{msg}})$$

2. execute R_1

3. deliver the message.



fig: Shows the evolution of idle time width $d=1$.

\rightarrow Termination detection by weight throwing

In termination detection by weight throwing, a process

called controlling agent monitors the computation. It communicates channel exists between each of

the processes and the controlling agent and also for every pair of processes.

Initially, all processes are in the idle state. The weight of each process is zero and the weight at the bootstrapping agent is 1.

The computation starts when the controlling agent

- * sends a basic message to one of the processes.
- * The process becomes active and the computation starts.
- * If the non-zero weight $w(\text{proc} \leq 1)$ is assigned, each process in the active state adds to each message in transit in the following manner.

- * If the process sends a message, it sends a part of its weight in the message.
- * When the sum of weight on all the processes and received in the message is transmitted, it is always 1.

- * When a process becomes passive, it sends its weight to the controlling agent in a control message, which the controlling agent adds to its weight.

- * The controlling agent concludes termination of its weight becomes 1.
- * Notation:

- * The weight on the controlling agent and a process is represented by w .

- * $B(DW)$: A basic message B is sent as a part of the computation, where DW is the weight assigned to it.
- * $C(DW)$: A control message C is sent from a process

to the controlling agent where DW is the weight assigned to it.

Formal description

The algorithm is defined by the following four rules:

Rule 1: The controlling agent or an active process may send a basic message to one of the processes, say P , by splitting its weight W into W_1 and W_2 such that $W_1 + W_2 = W$, $W_1 > 0$ and $W_2 > 0$. It then assigns its weight $w := W_1$ and sends a basic message $B(DW := w_2)$ to P .

Rule 2: On the receipt of the message $B(DW)$, process P adds DW to its weight $(w := w + DW)$. If the receiving process is in the idle state, it becomes active.

Rule 3: A process transitions from the active state to the idle state at any time by sending a control message $C(DW := w)$ to the controlling agent and making its weight $w := 0$.

→ 8/3/23

A spanning tree based termination detection algorithm

The algorithm description.

The algorithm works as follows:

Initially, each leaf process is provided with a token. The set S is used for book-keeping to know which processes have the token. There S will be the set of all leaves in the tree.

Initially, all processes and tokens are idle. As explained above, coloring helps the agent know if a message-passing process involved in one of the subtrees has already terminated, it sends the token to its parent process.

4. A parent process will collect the tokens sent by each of its children. After it has received a token from all of its children and after it has terminated the parent process sends a token to its parent.

5. A process turns back when it sends a message to some other process. This coloring scheme helps a process remember that it has sent a message

Rule 4: On the receipt of a message $C(DW)$, the controlling agent adds DW to its weight ($w := w + DW$). If $w = 1$, then it concludes that the computation has terminated.

When a process terminates, if it's black, it sends a black token to its parent.

6. A black process turns back to white after it has sent a black token to its parent.

7. A parent process holding a black token (from one of its children), sends only a black token to its parent, to indicate that a message-passing was correct, involved in the subtree.

8. Process are propagated to the root in this fashion.

The root, upon receiving a black token, will know that a process in the tree had sent a message to some other process. Hence, it restarts the algorithm by sending a Repeat signal to all its children.

9. Each child of the root propagates the Repeat signal to each of its children and so on, until the signal reaches the leaves.

10. The leaf nodes restart the algorithm on receiving the Repeat signal.

11. The root concludes that termination has occurred if:

- it is white;
- it is idle; and
- it has received a black token from each of its children.

All leaf nodes have tokens

$S = \{3, 4, 5, 6\}$

T_1, T_2, T_3

\Rightarrow

$0, 1, 2$

T_4, T_5, T_6

\Rightarrow

$3, 4, 5, 6$

T_7

\Rightarrow

0

T_8

\Rightarrow

1

T_9

\Rightarrow

2

T_{10}

\Rightarrow

3

T_{11}

\Rightarrow

4

T_{12}

\Rightarrow

5

T_{13}

\Rightarrow

6

T_{14}

\Rightarrow

7

T_{15}

\Rightarrow

8

T_{16}

\Rightarrow

9

T_{17}

\Rightarrow

10

T_{18}

\Rightarrow

11

T_{19}

\Rightarrow

12

T_{20}

\Rightarrow

13

T_{21}

\Rightarrow

14

T_{22}

\Rightarrow

15

T_{23}

\Rightarrow

16

T_{24}

\Rightarrow

17

T_{25}

\Rightarrow

18

T_{26}

\Rightarrow

19

T_{27}

\Rightarrow

20

T_{28}

\Rightarrow

21

T_{29}

\Rightarrow

22

T_{30}

\Rightarrow

23

T_{31}

\Rightarrow

24

T_{32}

\Rightarrow

25

T_{33}

\Rightarrow

26

T_{34}

\Rightarrow

27

T_{35}

\Rightarrow

28

T_{36}

\Rightarrow

29

T_{37}

\Rightarrow

30

T_{38}

\Rightarrow

31

T_{39}

\Rightarrow

32

T_{40}

\Rightarrow

33

T_{41}

\Rightarrow

34

T_{42}

\Rightarrow

35

T_{43}

\Rightarrow

36

T_{44}

\Rightarrow

37

T_{45}

\Rightarrow

38

T_{46}

\Rightarrow

39

T_{47}

\Rightarrow

40

T_{48}

\Rightarrow

41

T_{49}

\Rightarrow

42

T_{50}

\Rightarrow

43

T_{51}

\Rightarrow

44

T_{52}

\Rightarrow

45

T_{53}

\Rightarrow

46

T_{54}

\Rightarrow

47

T_{55}

\Rightarrow

48

T_{56}

\Rightarrow

49

T_{57}

\Rightarrow

50

T_{58}

\Rightarrow

51

T_{59}

\Rightarrow

52

T_{60}

\Rightarrow

53

T_{61}

\Rightarrow

54

T_{62}

\Rightarrow

55

T_{63}

\Rightarrow

56

T_{64}

\Rightarrow

57

T_{65}

\Rightarrow

58

T_{66}

\Rightarrow

59

T_{67}

\Rightarrow

60

T_{68}

\Rightarrow

61

T_{69}

\Rightarrow

62

T_{70}

\Rightarrow

63

T_{71}

\Rightarrow

64

T_{72}

\Rightarrow

65

T_{73}

\Rightarrow

66

T_{74}

\Rightarrow

67

T_{75}

\Rightarrow

68

T_{76}

\Rightarrow

69

T_{77}

\Rightarrow

70

T_{78}

\Rightarrow

71

T_{79}

\Rightarrow

72

T_{80}

\Rightarrow

73

T_{81}

\Rightarrow

74

T_{82}

\Rightarrow

75

T_{83}

\Rightarrow

76

T_{84}

\Rightarrow

77

T_{85}

\Rightarrow

78

T_{86}

\Rightarrow

79

T_{87}

\Rightarrow

80

T_{88}

→ Termination detection using distributed snapshot.

The algorithm assumes that there is a logical communication channel b/w every pair of processes. Communication channels are reliable but non-FIFO message delay is arbitrary but finite.

The main idea behind the algorithm is as follows:

- * Idea a computation terminates there must exist a unique process which became idle last.
- * When a process goes from active to idle it issues a request to all other processes to take a local snapshot and also request itself to take a local snapshot.
- * When a process receives the request if it agrees that requester becomes idle before itself it grants the request by taking a local snapshot for the request.
- * Request is said to be successful if all processes have taken a local snapshot for it.
- * The requester or any external agent may collect all the local snapshots for the request.
- * If the request is successful, a global snapshot can be obtained and the desired state will indicate termination of the

communication. In the recorded snapshot all the processes are idle and there is no message in transit to any of the message.

13/3/23

Module - 3Distributed Mutual Exclusion Algorithms.

→

Distributed Mutual ExclusionAlgorithms.

- * Mutual exclusion is a distributed system where at most only one process is allowed to execute critical section at any given time.

- * Message passing is the means for implementing distributed mutual exclusion.

- * Three basic approaches for implementing distributed mutual exclusion-
 - Token based approach
 - Non-token based approach
 - Anonymous based approach.

Asume channels reliably deliver messages do not crash and network doesn't get partitioned.

Requirements of mutual exclusion algorithm

→

System ModelSystem Model

- * N sites S_1, S_2, \dots, S_n
- * Single process running on each site, the processes at site i denoted by P_i . Processes communicate asynchronously.
- * Process send request message to other critical section and wait for reply before next entering critical section.
- * A site can be in any of the three states.

Performance MatrixForward

- (i) Requesting the critical section - site block and cannot make further request
- (ii) Executing critical section (CS)
- (iii) Neither requesting nor executing CS (idle)
 - Site is executing outside CS.

- (i) At any instant site may have several pending requests for CS, site queues up this request and serves them one at a time. No assumptions regarding communication channel if they use FIFO or not.
- (ii) Assume channels reliably deliver messages do not crash and network doesn't get partitioned.

- (i) It should satisfy the following properties
 - (a) Safety property - At any instant only one process can execute CS
 - (b) Liveness property - This property states absence of deadlock & starvation
 - (c) Fairness - each process must get a fair chance to execute CS.

(i) Message complexity - No. of messages that are required per critical section execution by a site.

(ii) Synchronization delay - After a site leaves the critical section, it is the time required before next site enters critical section.

(iii) Response time - This is the time interval a request wait for its critical section execution to be over after its request messages have been send out.

$\text{U request} \xrightarrow{\text{out}} \text{the site enters CS} \xrightarrow{\text{in}} \text{U exit CS}$

\Rightarrow Lamport's algorithm for mutual exclusion based algorithm

* Requesting the critical section

- When a site S_i wants to enter the CS, it broadcast a REQUEST(t_{si}, i) message to all other sites and places the request on request-queue: $((t_{si}, i))$ and denies the timestamp of the request.
- When a site S_j receives the REQUEST(t_{si}, i) message from site S_i , it places site S_i 's request on request-queue and removes a timestamped REPLY message to S_i .

* Executing the critical section

site S_i enters the CS unless the following two conditions hold:

- i : If has received a message with timestamp larger than (t_{si}, i) from all other sites.
- ii : Site's request is at the top of request-queue

$$(iv) \text{System throughput} = \frac{1}{SD + C}$$

This is the rate at which systems execute request for critical section.

SD - Synchronization Delay.

C - Average Critical Section executing time

- When a site S_j receives a RELEASE message

from site S_i , it removes S_i 's request from its request queue.

\rightarrow Ricart - Agrawala algorithm

If

- * Each process p_i maintains the request-deferred array RDi , the size of which is the same as the number of processes in the system.
- * Initially, $V_i V_j : RDi[j] = 0$.

* Requesting the critical section.

- (a) If 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 $RDi[i] := 1$.

- * If a process p_i that is waiting to execute the critical section receives a REQUEST message from process p_j ,

* Executing the critical section

- * Then if the priority of p_j 's request is lower, p_i defers the REPLY to p_j and sends a REPLY message to p_j only after executing the CS for its pending request.
- * Otherwise, p_i sends a REPLY message to p_j immediately provided it is currently not executing the CS.

- (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

(d) When site S_i exits the CS, it sends all the delivered REPLY messages: R_j if $RD_{[j]} = 1$, then sends a REPLY message to S_j and sets $RD_{[j]} := 0$.

~~14/3/23~~
Process Based Mutual Exclusion

* This algorithm reduces the complexity of handling mutual exclusion by having sites as permissions from only a subset of sites.

→ Mackawa's Algorithm

It was the first process based mutual exclusion algorithm. This algorithm organizes delivery of message to be handled in order they are sent along every pair of sites.

Requesting the critical section:

- A site S_i requests access to the CS by sending REQUEST (i) messages to all sites in its request set R_i .
- When a site S_j receives the REQUEST (i) message,

it sends a REPLY (i) 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.

Exiting the critical section:

Site S_i executes the CS only after it has received a REPLY message from every site in R_i .

Releasing the critical section:

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

~~15/3/23~~

Process on subset of each process:

- $P_1 = (3, 5)$ // process 3 and 5 are in queue
- $P_2 = (4, 5)$ // " 4 & 5 " of P_1
- $P_3 = (2, 3, 4)$ // " 2, 3 & 4 " of P_2

$$P_4 = (3, 4)$$

$$P_5 = (2, 3)$$

$$P_6 = (1, 3)$$

→ Rules

) The intersection of two group ≠ Null
ie; atleast one process should be present in both subsets.

2) A process can be part of any group

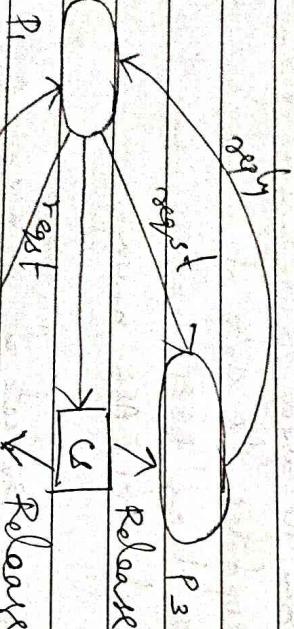
3) A single process cannot be considered as a group.
If a process want to access CS then it should send request msg to all members in its group if the it should receive reply from all members is quorum.

→ Token based algorithm

Suzuki - Kawano's broadcast algorithm

* Requesting the CS

- (a) If requesting the site i does not have token then it increases its sequence number $RN[i]$ and sends a REQUEST(i, s_n) msg to all other sites ("s_n" is updated value of $RN[i]$).
 When a site j receives this message it sets $RN[j][i]$ to $\max[RN[j][i], s_n]$. If j has the idle token, then it sends the token to site i if $RN[j][i] > LN[i] + 1$.



P₁ → req → C → Release or

P₂ → req → C → Release or

P₃ → Release or

(c)

Site s_i encounters the C after it has received the token.

→ Drawback

- * Chance of deadlock
- * Requests are not prioritized by timestamp.