

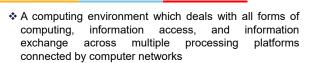


Modular Overview

- M1: Introduction to Distributed Computing
- M2: Logical Clocks & Vector Clocks
- M3: Global state and snapshot recording algorithms
- M4: Message ordering and Termination detection
- M5: Distributed Mutual Exclusion
- M6: Deadlock Detection
- M7: Consensus and Agreement Algorithm
- M8: Peer to Peer Computing and Overlay graphs
- M9: Cluster computing, Grid Computing
- M10: Internet of Things

BITS Pilani, Hyderabad Campus

What is a distributed system?



Ex: Banking systems; Communication systems, Information systems (WWW),Manufacturing and process control, Inventory Systems, General purpose automation systems etc,

- So, what is a good definition of a distributed system?
- A distributed system is a collection of independent entities that cooperate to solve a problem that cannot be individually solved.
- Distributed systems have been in existence since many years

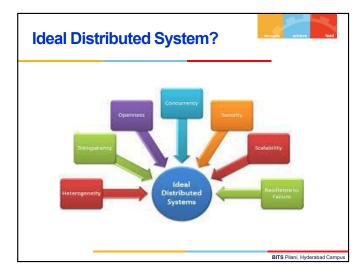
BITS Pilani, Hyderabad Campus

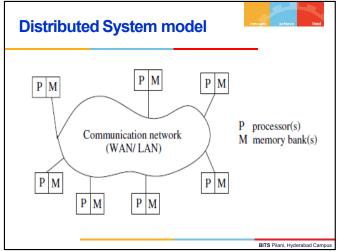
Characteristics of a distributed system

- Autonomous processors communicating over a communication network make a distributed system
 - √ No common physical clock
 - √ No shared memory
 - √ Geographical separation
 - ✓ Autonomy and heterogeneity

BITS Pilani, Hyderabad Campus

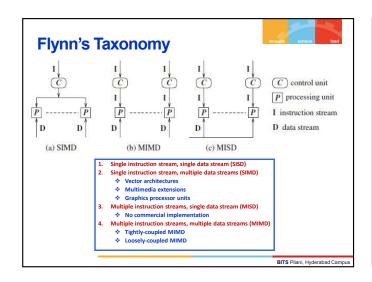
BITS Pilani, Hyderabad Campus

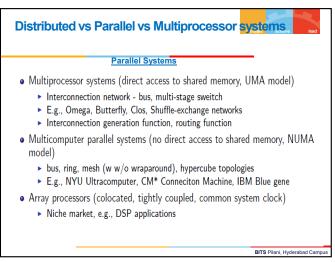


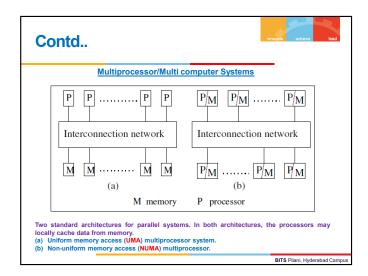


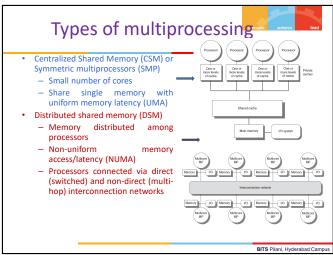
Software Components and their interaction in **Distributed environment** Distributed application Extent of distributed protocols Distributed software (middleware libraries) protocol Application layer Transport layer Operating Network system Network layer Data link layer BITS Pilani, Hyderabad Campus

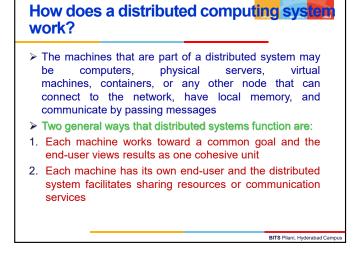
Why do we need distributed systems? Inherently distributed computation Resource sharing Access to remote resources Increased performance/cost ratio Reliability Availability, integrity, fault-tolerance Scalability, Modularity and incremental expandability

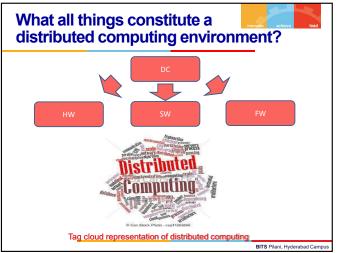












Interconnection networks in Distributed Computing



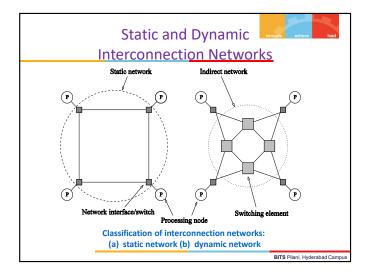
- Interconnection networks are composed of switching elements
- Topology is the pattern to connect the individual switches to other elements, like processors, memories and other switches
- A network allows exchange of data between processors in the distributed computing system

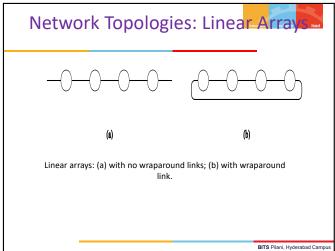
BITS Pilani, Hyderabad Campus

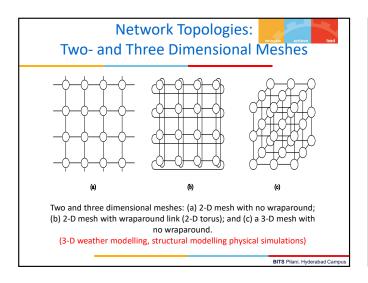
Types of Interconnection networks

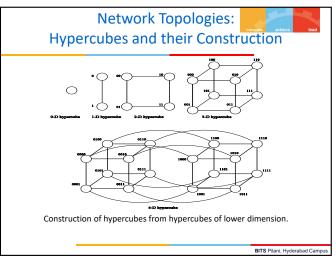
- Direct connection networks Direct (static) networks have point-to-point connections between neighboring nodes. Ex:
 - ≻Rings
 - **≻**Meshes
 - **≻**Cubes
- Indirect connection networks Indirect (dynamic) networks have no fixed neighbors. The communication topology can be changed dynamically based on the application demands. Ex:
 - **➢**Bus networks
 - **≻**Multistage networks
 - ➤ Cross bars

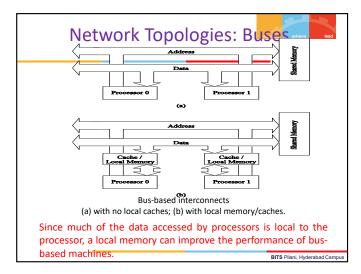
BITS Pilani, Hyderabad Campus

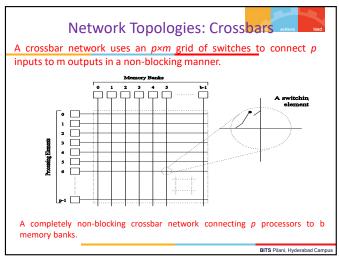


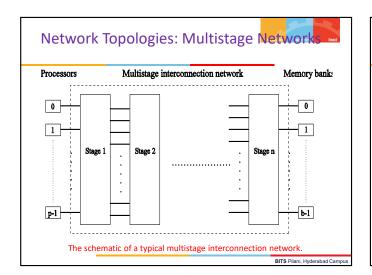




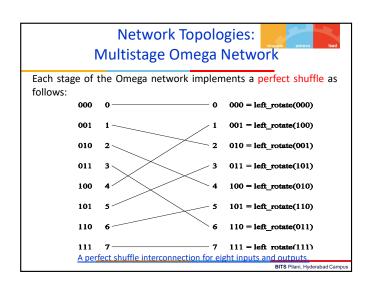


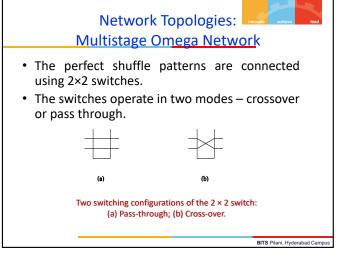


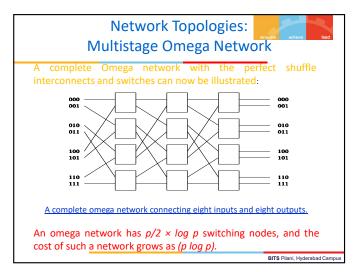


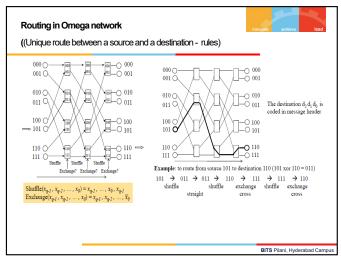


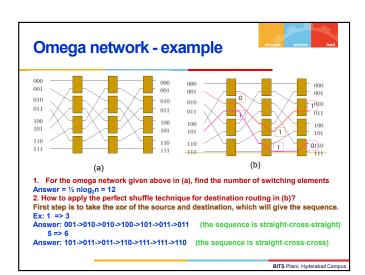
Network Topologies:

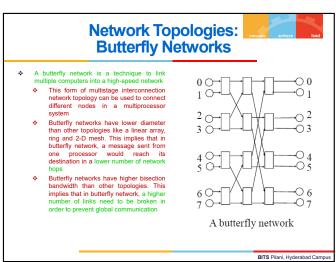












Coupling



- ✓ The degree of coupling among a set of modules, whether hardware or software, is measured in terms of the interdependency and binding and/or homogeneity among the modules
- ✓ When the degree of coupling is high (low), the modules are said to be tightly (loosely) coupled
- ✓ SIMD and MISD architectures generally tend to be tightly coupled because of the common clocking of the shared instruction stream or the shared data stream

Coupling in MIMD Architectures



- Tightly coupled multiprocessors (with UMA shared memory)
- 2. Tightly coupled multiprocessors (with NUMA shared memory or that communicate by message passing)
- 3. Loosely coupled multi computers (without shared memory) physically collocated
- Loosely coupled multi computers (without shared memory and without common clock) that are physically remote

BITS Pilani, Hyderabad Campus

BITS Pilani, Hyderabad Campus

Parallelism



- This is a measure of the relative speedup of a specific program, on a given machine
- The speedup depends on the number of processors and the mapping of the code to the processors
- Speedup (S) is the ratio of the time taken to solve a problem on a single processor to the time required to solve the same problem on a parallel computer with p identical processing elements
- It is expressed as the ratio of the time T(1) with a single processor, to the time T(n) with n processors.
- Parallelism within a parallel/distributed program is an aggregate measure of the percentage of time that all the processors are executing CPU instructions productively, as opposed to waiting for communication (either via shared memory or message-passing) operations to complete
- \triangleright S = $T_{\text{serial}}/T_{\text{parallel}}$

BITS Pilani, Hyderabad Campus

Concurrency: Definition



- ✓ A broader term that means roughly the same as parallelism of a program, but is used in the context of distributed programs
- √ The parallelism/concurrency in a parallel/distributed program can be measured by the ratio of the number of local (non-communication and non-shared memory access) operations to the total number of operations, including the communication or shared memory access operations

BITS Pilani, Hyderabad Campu

Concurrency



- The maximum number of tasks that can be executed simultaneously at any time in a parallel algorithm is called its degree of concurrency
- ❖If C(W) is the degree of concurrency of a parallel algorithm, then for a problem of size W, no more than C(W) processing elements can be employed effectively

BITS Pilani, Hyderabad Campu

Granularity in Distributed Computing



- The ratio of the amount of computation to the amount of communication within the parallel/distributed program is termed as granularity
- If the degree of parallelism is coarse-grained (fine-grained), there are relatively many more (fewer) productive CPU instruction executions, compared to the number of times the processors communicate either via shared memory or message passing and wait to get synchronized with the other processors
- Programs with fine-grained parallelism are best suited for tightly coupled systems
 - These typically include SIMD and MISD architectures, tightly coupled MIMD multiprocessors (that have shared memory), and loosely coupled multi computers (without shared memory) that are physically co located

BITS Pilani, Hyderabad Campu

Communication Paradigms for Distributed Computing



- Recap on interaction of software components in Distributed Computing (CORBA, RPC, DCOM,RMI etc.,)
- A distributed computation involves a number of processes communicating with one another and inter process communication mechanism is fairly complex
- many dimensions of variability in distributed systems like network topology, inter process communication mechanisms, failure classes, and security mechanisms
- * Two major paradigms are:
- 1. Shared Memory Process Communication Model
- 2. Message Passing Process Communication Model

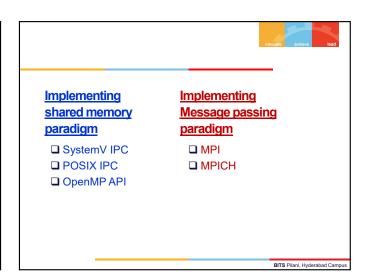
BITS Pilani, Hyderabad Campus

Shared Memory vs Message Passing Process P1 Process P2 Message Queue In In 2 Message Passing Model Shared Memory Process Communication Model BITS Plant, Hyderabad Campus

Message-passing and **Shared Memory - emulation**

- Emulating MP over SM:
 - ▶ Partition shared address space
 - ▶ Send/Receive emulated by writing/reading from special mailbox per pair of processes
- Emulating SM over MP:
 - ▶ Model each shared object as a process
 - ▶ Write to shared object emulated by sending message to owner process for the
 - ▶ Read from shared object emulated by sending query to owner of shared object

BITS Pilani, Hyderabad Campus



Classification of primitives



Synchronous (send/receive)

- Handshake between sender and receiver
- Send completes when Receive completes
- Receive completes when data copied into buffer

Asynchronous (send)

Control returns to process when data copied out of user-specified buffer

Blocking (send/receive)

Control returns to invoking process after processing of primitive (whether sync or async) completes

Nonblocking (send/receive)

- Control returns to process immediately after invocation
- Send: even before data copied out of user buffer
- > Receive: even before data may have arrived from sender

Non-blocking Primitive

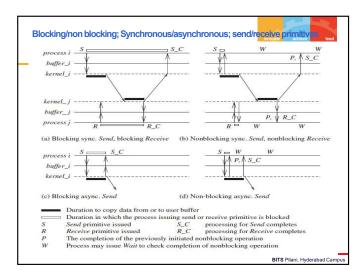


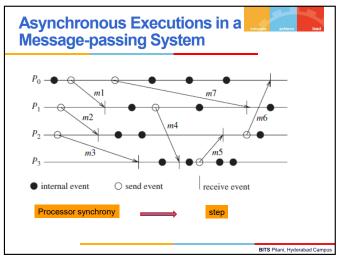
nonblocking send primitive

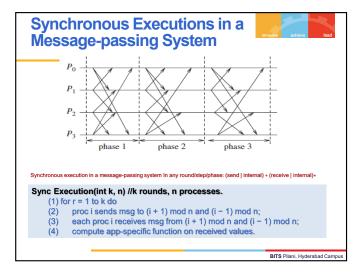
 $Send(X, destination, handle_k)$ //handle_k is a return parameter $Wait(handle_1, handle_2, ..., handle_k, ..., handle_m)$ //Wait always blocks

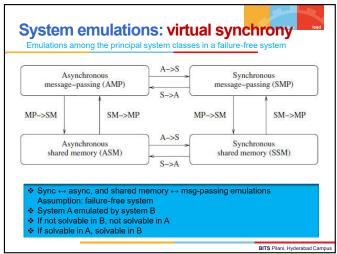
- When the Wait call returns, at least one of its parameters is posted
- Return parameter returns a system-generated handle
 Use later to check for status of completion of call
 - - Keep checking (loop or periodically) if handle has been posted

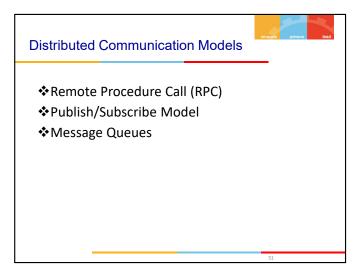
 - Issue Wait(handle1, handle2, . . .) call with list of handles
 Wait call blocks until one of the stipulated handles is posted

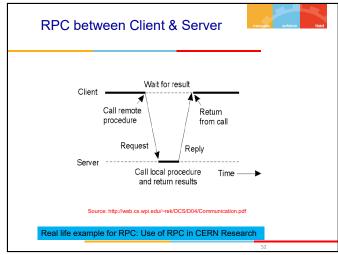


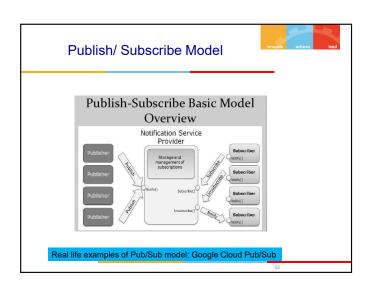


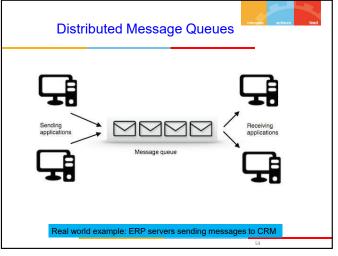












Design Issues & Challenges System Perspective

- Communication mechanisms
- Processes
- Synchronization
- ❖ Data storage and access
- Consistency and replication
- **❖** Fault-tolerance
- **❖** Security
- ❖ Scalability and modularity
- ❖ API and transparency

Algorithmic Challenges



- Designing useful execution models and frameworks: to reason with and design correct distributed programs
- Dynamic distributed graph algorithms and routing algorithms
- Load balancing: to reduce latency, increase throughput, done dynamically
- * Real-time scheduling: difficult without global view, network delays make task harder
- Performance modeling and analysis: Network latency to access resources must be reduced
- Synchronization/coordination mechanisms
- Group communication, multicast, and ordered message delivery
- Monitoring distributed events and predicates
- Distributed program design and verification tool
- ❖ Time and global state
- ❖ Debugging distributed programs
- Data replication, consistency models, and caching
- * World Wide Web design: caching, searching, scheduling
- Distributed shared memory abstraction
- * Reliable and fault-tolerant distributed systems

Recap Quiz



1. Which of the following is not one of the communication mechanisms in distributed computing?

(a) RPC

(c) ROI

Let Ts = 10 time units be the sequential time of execution in a distributed system and let
 Tp = 8 time units be the parallel time. What is the speedup? Assume other delays are
 negligible.

(a) 1.25

(b) 1.0

(c) 0.8

execution called (a) Task

(b) step

3. In distributed computing processor synchrony is achieved through units of (d) thread

(c) process

4. What is the granularity of a distributed computing environment if the amount of computation is 100 units and the amount of communication is 36 units?

(b) 2.78

(c) 1.36

5. The maximum number of tasks that can be executed simultaneously at any time in a distributed computing environment is called the degree of

(a) Efficiency

(b) granularity (c) consistency (d) concurrency

Recap Quiz key



Q1	Q2	Q3	Q4	Q5
d	а	b	b	d

BITS Pilani, Hydera

Major references



- Ajay D. Kshemkalyani, and Mukesh Singhal, Chapter 1, "Distributed Computing: Principles, Algorithms, and Systems", Cambridge University Press, 2008.
- http://www.ois.com/Products/what-is-corba.html
- http://wiki.c2.com/?PublishSubscribeModel
- https://www.slideshare.net/ishraqabd/publish-subscribe-modeloverview-13368808/5
- https://www.cloudamqp.com/blog/2014-12-03-what-is-messagequeuing.html
- https://en.wikipedia.org/wiki/Message_queue

SS ZG 526: Distributed Computing (CS2) Prof. Geetha ssor, Dept. of Computer Sc. & Information Systems BITS Pilani Hyderabad Campus **BITS Pilani** geetha@hvderabad.bits-pilani.ac.in

Introduction

- ❖ concept of causality between events → important for design and analysis of parallel and distributed computing and operating systems
- causality is tracked using physical time
- ❖ But in distributed systems
 - not possible to have global physical time
 - possible to realize only an approximation of physical time
- way of measuring time
 - asynchronous distributed computations progress in spurts
 - logical time also advances in jumps
 - sufficient to capture the fundamental monotonicity property associated with causality in distributed systems

Course ID: SS ZG526, Title: Distributed Computin

BITS Pilani, Hyderabad Camp

Introduction

- monotonicity property: if event a causally affects event b, then timestamp of a is smaller than timestamp of b
- Ways to represent logical time:
 - scalar time
 - vector time

Course ID: SS ZG526. Title: Distributed Computing

BITS Pilani, Hyderabad Ca

Introduction

- Causality / Causal Precedence Relation
- described among events in a distributed system
- useful in reasoning, analyzing, and drawing inferences about a computation
- helps to solve several problems in distributed systems
 - distributed algorithms design
 - tracking of dependent events
 - knowledge about progress
 - concurrency measure

SS 2G526, Title: Distributed Computing 4 BITS Pilani, Hy

Introduction

➤ Distributed algorithms design

knowledge of the causal precedence relation among events helps to:

- •ensure liveness and fairness in mutual exclusion algorithms •maintain consistency in replicated databases
- •design correct deadlock detection algorithms

Tracking of dependent events

- •helps construct a consistent state for resuming re-execution
- •in failure recovery, helps build a checkpoint
- •in replicated databases, aids in detection of file inconsistencies

Course ID: SS ZG526, Title: Distributed Computin

BITS Pilani, Hyderabad Ci

Introduction



- Knowledge about progress knowledge of causal dependency among events
 - •helps measure progress of processes in distributed computation
 - •helps to discard obsolete information, garbage collection, and termination detection
- Concurrency measure knowledge regarding how many events are causally dependent
 - $\hbox{-useful in measuring amount of concurrency} \\$
 - •events not causally related are concurrent
 - •analysis of causality tells about concurrency in program

Capturing Causality between Events



- •in distributed systems
 - •rate of event occurrence is of high magnitude
 - •event execution time is of low magnitude
- •physical clocks must be precisely synchronized
- •in distributed computation
 - •clocks accurate to a few tens of milliseconds are not sufficient
 - •progress occurs in spurts
 - •interaction between processes occurs in spurts
- •logical clocks can accurately capture
 - •causality relation between events produced by a program execution
 - •fundamental monotonicity property

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Ci

Logical Clocks

- •every process has a logical clock
- •logical clock is advanced using a set of rules
- •each event is assigned a timestamp
- •causality relation between events can be inferred from their timestamps
- •timestamps follow the monotonicity property

Representing Logical Clocks

- 1. Lamport's scalar clocks
 - time is represented by non-negative integers
- 2. Vector clocks (Fidge, Mattern and Schmuck)
 - time is represented using a vector of non-negative integers

Real world applications of Logical clocks



❖ Vector clocks => Voldemort (Linkedin) Amazon Dynamo, Version control systems etc.

Self assessment Question: • Why Lamport's clock cannot be used in blockchain technology?

Definition of Logical Clocks



- •system of logical clocks consists of
 - •time domain T
 - •logical clock C
- •elements of T form a partially ordered set over a relation <
- - •happened before or causal precedence relation
 - •analogous to 'earlier than' relation provided by physical time

Definition of Logical Clocks



•is a function that maps an event e in a distributed system to an element in the time domain T

- •denoted as C(e)
- •C(e): timestamp of e
- •C is defined as: $C: H \mapsto T$
- •satisfies the following property
 - ullet for any 2 events e_i and e_j , $e_i
 ightarrow e_j \Longrightarrow \mathsf{C}(e_i) < \mathsf{C}(e_j).$
 - monotonicity property
 - •known as clock consistency condition

Definition of Logical Clocks



•system of clocks is *strongly consistent* if T and C satisfy the following:

•for 2 events e_i and e_i ,

$$e_i \rightarrow e_j \Leftrightarrow \mathsf{C}(e_i) < \mathsf{C}(e_j)$$

Implementing Logical Clocks

- Two issues need to be addressed
 - data structures local to every process to represent logical time
 - a protocol to update the data structures to ensure the consistency condition

Course ID- 55 70576. Titles Distributed Committee

BITS Pilani, Hyderabad Campu

Implementing Logical Clocks

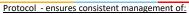
Data Structures

- •Each process p_i maintains two data structures
- •A local logical clock
 - •denoted by Ic;
 - •helps process p_i measure its own progress
- •A logical global clock
 - •denoted by gc_i
 - •represents p_i 's local view of the logical global time
 - •allows p_i to assign consistent timestamps to its local events
 - •lc; is a part of gc;

Course ID: SS ZG526. Title: Distributed Computing

BITS Pilani, Hyderabad Ca

Implementing Logical Clocks



•a process's logical clock

process's view of global time

- •consists of the following 2 rules:
- R1: governs how the local logical clock is updated by a process when it executes an event (send, receive, or internal)
- R2
 - •governs how a process updates its global logical clock to update its view of the global time and global progress
 - •dictates what information about logical time is piggybacked in a message
 - •how this information is used by the receiving process to update its view of the global time

Course ID, 66 30536. Tales Distributed Course the

BITS Pilani, Hyderabad Campu

Implementing Logical Clocks

- •systems of logical clocks differ in
 - •representation of logical time
 - •the protocol to update the logical clocks
- •all logical clock systems
 - •implement rules R1 and R2
 - •ensure the fundamental monotonicity property associated with causality
 - •provide users with some additional properties

Course ID: SS ZG526, Title: Distributed Computing

Scalar Time - Definition

•representation was proposed by Lamport in 1978 to totally order events in distributed system

•in this representation, time domain is the set of non-negative integers

•C;

- •integer variable
- •denotes logical local clock of p_i and its local view of global time

Scalar Time - Definition

Rules for updating clock:

*R1: Before executing an event (send, receive, or internal), process p_i executes $C_i := C_i + d \qquad (d > 0)$

For each execution of R1,

- •d can have different value
- •value of d may be application-dependent

Usually d = 1

- •helps to identify the time of each event uniquely at a process
- •keeps rate of increase of d to lowest level

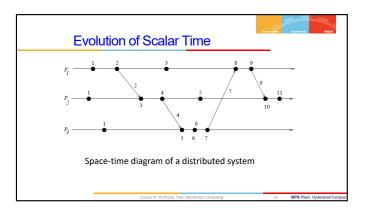
19 BITS Pilani, Hyderabad 0

Scalar Time - Definition

Rules for updating clock (contd. from previous slide):

R2: When process p_i receives a message with timestamp C_{msg} , it executes the following actions:

- 1. $C_i = max(C_i, C_{msg});$
- 2. execute R1;
- 3. deliver the message.



Scalar Time - Basic Properties

- Consistency
- •Total Ordering
- •Event Counting
- •No Strong Consistency

Scalar Time – Basic Properties

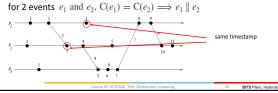
Consistency Property

- for two events e_i and e_j , $e_i \rightarrow e_j \Longrightarrow C(e_i) < C(e_j)$
- scalar clocks satisfy monotonicity property
- hence, they also satisfy consistency property

Scalar Time - Basic Properties

Total Ordering

- Scalar clocks are used to totally order events in a distributed system
- Problem in totally ordering events 2 or more events at different processes may have same timestamp



Scalar Time - Basic Properties

Total Ordering

- Tie breaking mechanism is needed
- Tie breaking procedure:
 - process identifiers are linearly ordered
 - break tie among events with identical scalar timestamps on the basis of their process identifiers
 - lower process identifier implies higher priority
 - timestamp of an event is denoted by a tuple (t, i)
 - $t \rightarrow time of occurrence$
 - $i \rightarrow identity of the process where it occurred$



Total Ordering

- Total order relation < on 2 events x and y with timestamps (h, i) and (k, j) respectively, is defined:
 - $x \prec y \Leftrightarrow h < k \text{ or } (h = k \text{ and } i < j)$
- events that occur at the same logical scalar time are independent (i.e., they are not causally related)
- such events can be ordered using any arbitrary criterion without violating the causality relation \Rightarrow
- so, total order is consistent with the causality relation "→"
- Note:- x < y ⇒ x → y ∨ x | | y

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campu

Scalar Time - Basic Properties

Total Ordering

- Utility
 - total order is generally used to ensure liveness properties in distributed algorithms
 - · Liveness property:
 - something good eventually happens
 - system makes progress, no starvation, programs terminate
 - requests are timestamped and served according to the total order based on these timestamps

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Ca

Scalar Time - Basic Properties

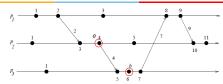
Event Counting

- If the increment value d=1, scalar time has the following property:
 - if event e has a timestamp h, then h 1
 - is the minimum logical duration required before producing event e
 - counted in units of events
 - called the *height* of event *e*
 - implies that h 1 events have been produced sequentially before event e, irrespective of the processes that produced these events

Course ID: SS ZG526, Title: Distributed Computing

BIT'S Pilani, Hyderabad Campu

Scalar Time - Basic Properties



- height of event a = 4 1 = 3, i.e., 3 events precede a on the longest causal path ending at a
- height of event b = 6 1 = 5, i.e., 5 events precede b on the longest causal path ending at b

Course ID: SS ZG526, Title: Distributed Computing

Scalar Time - Basic Properties

No strong consistency

- system of scalar clocks is not strongly consistent
- for 2 events e_i and e_j , $C(e_i) < C(e_j) \not\Rightarrow e_i \rightarrow e_j$

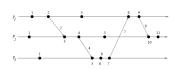


- scalar timestamp of 3rd event of p_1 < scalar timestamp of 3rd event of p_2
- but the former did not happen before the latter

Scalar Time - Basic Properties

No strong consistency

- scalar clocks are not strongly consistent because logical local clock and logical global clock of a process are combined into one
- results in the loss causal dependency information among events at different processes



- when p₂ receives the 1st message from p₁, p₂ updates its clock to 3
- forgets that the timestamp of the latest event at p₁ on which it depends is 2

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Cam



- time domain is represented by a set of n-dimensional nonnegative integer vectors
- each process p_i maintains a vector vt_i[1..n]
 - vt_i[i]: local logical clock of p_i
 - describes progress of logical time at p_i
 - vt_i[j]: represents process p_i's latest knowledge of process p_j's local time
 - if $vt_i[j] = x$
 - p_i knows that local time at p_j has progressed till x
- vector vt_i constitutes p_i 's view of the global logical time
- vt_i is used to timestamp events

urse ID: SS ZG526, Title: Distributed Computing 32

Vector Time - Definition

Rules used by p; for clock updating

• R1: Before executing an event, process p_i updates its local logical time as follows:

$$vt_i[i] = vt_i[i] + d, d > 0$$

urse ID: SS ZG526, Title: Distributed Computing 33 BITS Pilan

Vector Time - Definition

Rules used by p_i for clock updating

- R2:
 - each message *m* is piggybacked with the vector clock *vt* of the sender process at sending time
 - when p_i receives message (m, vt), it executes the following sequence of actions:
 - 1. update its global logical time as: $1 \le k \le n$: $vt_i[k] = max(vt_i[k], vt[k])$;
 - 2. execute R1;
 - 3. deliver the message *m*.

rse ID: SS ZG526, Title: Distributed Computing 34 BITS Pilar

Vector Time - Definition Rules used by p, for clock updating • R2: • timestamp associated with an event is the value of the vector clock of its process when the event is executed • initially, a vector clock is [0, 0, 0,, 0]

Vector Time - Definition

Relations for comparing 2 vector timestamps vh and vk

- $vh = vk \Leftrightarrow \forall x : vh[x] = vk[x]$
- $vh \le vk \iff \forall \ x : vh[x] \le vk[x]$
- $vh < vk \Leftrightarrow vh \leq vk \text{ and } \exists x : vh[x] < vk[x]$
- $vh \mid \mid vk \Leftrightarrow \neg(vh < vk) \land \neg(vk < vh)$

Vector Time - Basic Properties

- Isomorphism
- Strong Consistency
- Event Counting

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Car



Isomorphism

- if events in a distributed system are timestamped using a system of vector clocks, following property holds:
 - if 2 events x and y have timestamps vh and vk, respectively, then $x \to y \Leftrightarrow vh < vk$ $x \mid | y \Leftrightarrow vh | | vk$
- there is an isomorphism between the set of partially ordered events produced by a distributed computation and their vector timestamps
- very powerful, useful, and interesting property of vector clocks

ID: SS ZGS26. Title: Distributed Computing 38 BITS Pilani. Hyderabad

Vector Time - Basic Properties

Isomorphism

- if the process at which an event occurred is known, test for comparing 2 timestamps can be simplified as:
 - if events x and y respectively occurred at processes p_i and p_j and have timestamps vh and vk respectively, then

$$x \rightarrow y \Leftrightarrow vh[i] \le vk[i]$$

 $x \mid | y \Leftrightarrow vh[i] > vk[i] \land vh[j] < vk[j]$

Course ID: SS 7G526 Title: Distributed Computing

BITS Pilani, Hyderabad Campi

Vector Time – Basic Properties

Strong Consistency

- system of vector clocks is strongly consistent
- examination of vector timestamps of 2 events can determine if the events are causally related
- dimension of vector clocks can't be less than n for this property to hold

(n = total no. of processes in the distributed computation)

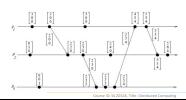
Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campu

Vector Time - Basic Properties

Event counting

- if d is always 1 in rule R1 and vt_i[i] is the ith component of vector clock at process p_i
 - then $vt_i[i]$ = no. of events that have occurred at p_i until that instant

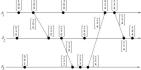


Vector Time – Basic Properties

Event counting

•if an event e has timestamp vh

- vh[j] = number of events executed by process p_j that causally precede e
- Σ vh[j] -1: total no. of events that causally precede e in the distributed computation



Efficient Implementation of Vector Clocks

Why is efficient implementation necessary?

- when no. of processes in a distributed computation is large
 - vector clocks will require piggybacking of huge amount of information in messages
 - messages are required for disseminating time progress and updating clocks
 - message overhead grows linearly with the no. of processes
 - message overhead is not affected by the no. of events occurring at the processors

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Cam

Singhal–Kshemkalyani's Differential Technique

- ☐ based on the observation between successive message sends to the same process, only a few entries of the vector clock at the sender process are likely to change
- $\hfill \square$ this is more likely when the number of processes is large lacktriangledown reason: only a few of them will interact frequently by passing messages
- ☐ fundamental idea behind the technique
 - \square when process p_i sends a message to process p_i , it piggybacks only those entries of its vector clock that differ since the last message sent to p;

Singhal-Kshemkalyani's Differential Technique

•if entries i_1 , i_2 ,, i_{n_1} of the vector clock at p_i have changed to v_1 , v_2 ,, $v_{n,i}$, respectively, since the last message sent to p_i , •then p_i piggybacks a compressed timestamp of the form $\{(i_1,\, v_1),\, (i_2,\, v_2),....,(i_{n_1},\, v_{n_1})\}$ to the next message to p_i

• on receiving this message, p_i updates its vector clock as follows: $vt_i[i_k] = max(vt_i[i_k], v_k)$ for $k = 1, 2,, n_1$

Singhal-Kshemkalyani's Differential Technique



Benefit:

> reduces message size, communication bandwidth and buffer (to store messages) requirements

Worst case:

- \triangleright every element of the vector clock has been updated at p_i since the last message sent to p_i
- \succ next message from p_i to p_j will need to carry the entire vector timestamp of size n
- > Average case: size of the timestamp on a message will be less than

Singhal-Kshemkalyani's Differential Technique



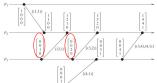
•Requirements for implementation:

- •each process must remember the vector timestamp in the message last sent to every other process
- •communication channels follow FIFO discipline for delivery of messages

Singhal–Kshemkalyani's Differential Technique



- in general, the timestamp is of the form:
 - $\{(p_{\scriptscriptstyle 1}, latest_value), (p_{\scriptscriptstyle 2}, latest_value), ...\}$
 - p_i indicates p_i^{th} component of the vector clock has changed



- 2nd message from ρ_3 to ρ_2 contains a timestamp $\{(3,2)\}$ informs ρ_2 that the 3rd component of the vector clock has been modified and the new value is 2

Singhal–Kshemkalyani's Differential Technique



- cost of maintaining vector clocks in large systems can be substantially reduced
- especially if the process interactions exhibit temporal or spatial localities
- useful in applications
 - · causal distributed shared memories
- · distributed deadlock detection
- enforcement of mutual exclusion and localized communications

Fowler–Zwaenepoel's Direct-Dependency Technique

- Innovate achieve lead
- · reduces the size of messages by transmitting only a scalar value
- · no vector clocks are maintained on-the-fly
- process only maintains information regarding direct dependencies on other processes
- vector time for an event represents transitive dependencies on other processes
- vector time for event is constructed off-line from a recursive search of the direct dependency information at processes

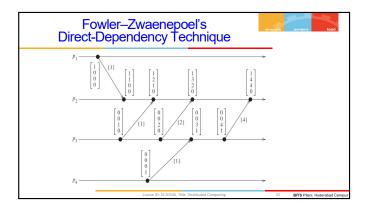
ID-SS-7GS26. Title: Distributed Computing 50 RITS Pilani Hyderabac

Fowler–Zwaenepoel's Direct-Dependency Technique

- p_i maintains a dependency vector D_i
- initially, $D_i[j] = 0$ for j = 1,....,n
- D, is updated using the following rules:
- 1. When an event occurs at p_p $D_i[i] = D_i[i] + 1$ (i.e., increment the vector component corresponding to its own local time by one)
- 2. When p_i sends a message to p_j , it piggybacks the updated value of $D_i[i]$ in the message
- 3. When p_i receives a message from p_j with piggybacked value d, p_i updates its dependency vector as follows: $D[j] = max\{D[j], d\}$

Course ID-00 700726 Title-Distributed Course time

BITS Pilani, Hyderabad Can



Fowler–Zwaenepoel's Direct-Dependency Technique

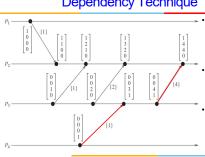


- D_i reflects only direct dependencies
- at any instant, D_i[j] denotes the sequence no. of the latest event on
 p_i that *directly* affects the current state
- this event may precede the latest event at ρ_j that $\it causally$ affects the current state

Course ID: SS ZG526, Title: Distributed Computin

BITS Pilani, Hyderabad Camp

Fowler–Zwaenepoel's Direct-Dependency Technique



- when p₄ sends a message to p₃, it piggybacks a scalar that indicates the direct dependency of p₃ on p₄ because of this message
- then p_3 sends a message to p_2 piggybacking a scalar to indicate the direct dependency of p_2 on p_3 because of this message
- p_2 is indirectly dependent on p_4 since p_3 is dependent on p_4

Fowler–Zwaenepoel's Direct-Dependency Technique



- transitive (indirect) dependencies are not maintained by this method
- transitive dependencies
 - can be obtained only by recursively tracing the direct dependency vectors of the events off-line
 - involves computational overhead and latencies
- this method is ideal only for those applications
 - that do not require computation of transitive dependencies on the fly
 - eg. applications: causal breakpoints, asynchronous checkpoint recovery

a was assessed with a later a later at

BITS Pilani, Hyderabad Car

Fowler-Zwaenepoel's Direct-Dependency Technique

- saves cost considerably
- not suitable for applications that require on-the-fly computation of vector timestamps

Physical Clock Synchronization

- no global clock or common memory
- each processor has its own internal clock and its own notion of time
- clocks can drift apart by several seconds per day, accumulating significant errors over time
- clock rates are different, may not remain always synchronized
- for most applications and algorithms that run in a distributed system, need to know time in one or more contexts:
 - time of the day at which an event happened on a specific machine in the network
 - time interval between two events that happened on different machines in
 - network

relative ordering of events that happened on different machines in network

Physical Clock Synchronization



- Clock synchronization -
- ensuring that physically distributed processors have a common notion of time
- · performed to correct clock skew in distributed systems
- synchronized to an accurate real-time standard like UTC (Universal Coordinated Time)
- Offset -
 - difference between the time reported by a clock and the real time
 - offset of the clock C_a is given by C_a(t) -t (C_a(t) time of clock)
 - offset of clock C_a relative to C_b at time $t \ge 0$ is given by $C_a(t) C_b(t)$

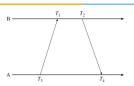
Network Time Protocol



- uses offset delay estimation method
- involves a hierarchical tree of time servers
- primary server at the root synchronizes with the UTC
- next level contains secondary servers, which act as a backup to the primary server
- at the lowest level is the synchronization subnet which has the clients

Network Time Protocol





- Each NTP message includes the latest three timestamps T_1 , T_2 , and T_3 T_4 is determined upon arrival Peers A and B can independently calculate delay and offset using a bidirectional message stream
- T₁, T₂, T₃, T₄ values of the 4 most recent timestamps as shown
- assume that clocks A and B are stable and running at the same speed
- $a = T_1 T_3$, $b = T_2 T_4$
- if the network delay difference from A to B and from B to A, called differential delay, is small, the clock offset $\boldsymbol{\theta}$ and roundtrip delay δ of B relative to A at time $\rm T_4$

 $\theta = (a+b)/2$ δ= a-b

Network Time Protocol

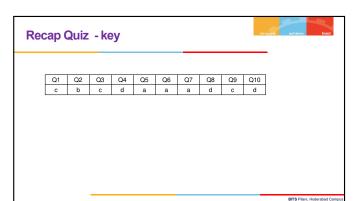


- A pair of servers exchange pairs of timing messages.
- A store of data is then built up about the relationship between the two
- Servers (pairs of offset and delay).
- Each peer maintains pairs (O_i, D_i) , where:
- O_i measure of offset (θ)
- D_i transmission delay of two messages (δ) Offset corresponding to the minimum delay is chosen
- Assume that message m takes time t to transfer and m' takes t' to transfer
- Offset between A's clock and B's clock is O
- If A's local clock time is A(t) and B's local clock time is B(t)
- A(t) = B(t)+O
- Then, $T_{i-2} = T_{i-3} + t + O$, $T_i = T_{i-1} O + t$

Network Time Protocol

- Assume t = t'
- Offset O_i can be estimated as O_i = $(T_{i-2} T_{i-3} + T_{i-1} T_i)/2$
- Round-trip delay is estimated as $D_i = (T_i T_{i-3}) (T_{i-1} T_{i-2})$
- 8 most recent pairs of (O_i, D_i) are retained Value of O_i that corresponds to minimum D_i is chosen to estimate O

Recap Quiz





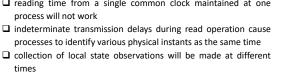




Global state and snapshot contd.. Every distributed system component has a local state state of a process is characterized by state of its local memory history of process activity channel state is characterized by the set of messages sent along the channel less the set of messages received along the channel Global state of a distributed system is a collection of the local states of its components Snapshot is the state of a system at a particular point in time

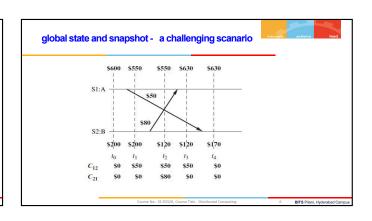
Global state and snapshot contd.. | benefit of shared memory - up-to-date state of the entire system is available to the processes sharing the memory | absence of shared memory requires ways of getting a coherent and complete view of the system based on the local states of individual processes | meaningful global snapshot can be obtained | if the components of the distributed system record their local states at the same time | requires local clocks at processes to be perfectly synchronized | existence of global system clock that could be instantaneously read by the processes

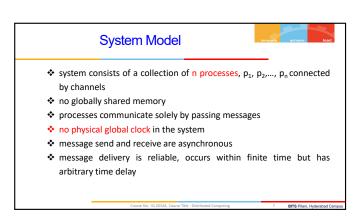
Global state and snapshot contd.. technologically infeasible to have perfectly synchronized clocks at various sites clocks are bound to drift reading time from a single common clock maintained at one process will not work

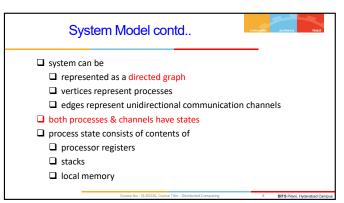


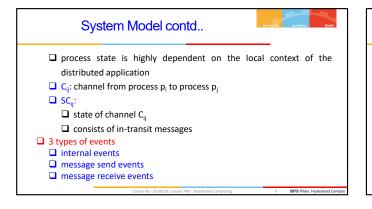
☐ may not be meaningful

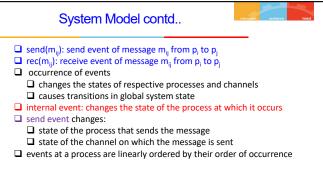
BITS Pilani, Hyderabad Camp











System Model contd..

Innovate achieve load

- receive event:
 - > changes state of the receiving process
 - > state of the channel on which the message is received
- LS_i: state of process p_i
- \succ at an instant t, LS $_{i}$ state of p_{i} as a result of the sequence of events executed by p_{i} up to t
- ightharpoonup an event $e \in LS_i$ iff e belongs to the sequence of events that have taken p_i to LS_i
- ightharpoonup e $otin LS_i$ iff e does not belong to the sequence of events that have taken p_i to LS_i

Course No.- SS ZG526, Course Title - Distributed Computing

BITS Pilani, Hyderabad Campu

System Model contd..



- $\ \, \mbox{\ \, } \mbox{\ \, for a channel } \mbox{\ \, C}_{ij}$, in-transit messages are:
 - **♦** transit(LS_i, LS_i) = { m_{ij} | send(m_{ij}) ∈ LS_i \land rec(m_{ij}) \notin LS_i}
- if a snapshot recording algorithm records the states of p_i and p_j as LS_i and LS_i, respectively,
 - ★ it must record the state of channel C_{ij} as transit(LS_i, LS_j)

Course No. - SS ZG526, Course Title - Distributed Computing

System Model contd..



- ☐ Several models of communication among processes exist
- ☐ FIFO model:
 - ☐ each channel acts as a first-in first-out message queue
- ☐ message ordering is preserved by a channel
- non-FIFO model:
 - ☐ channel acts like a set
 - sender process adds messages to the channel in a random order
 - ☐ receiver process removes messages from the channel in a random order

Course No. - SS 2G526, Course Title - Distributed Computing

BITS Pilani, Hyderabad Campus

A Consistent Global State



- ☐ global state GS is a *consistent global state* iff it satisfies the following two conditions:
 - \square C1: send(m_{ii}) \in LS_i \Rightarrow m_{ii} \in SC_{ii} \bigoplus rec(m_{ii}) \in LS_i (\bigoplus : XOR)
 - $\begin{tabular}{ll} \blacksquare $\textbf{C2}$: send(m_{ij}) \not\in LS_i $\Rightarrow m_{ij} \not\in SC_{ij} \land rec(m_{ij}) \not\in LS_j$ \end{tabular}$

Course No.- SS ZG526, Course Title - Distributed Computing

BITS Pilani, Hyderabad Camp

A Consistent Global State



☐ Condition C1 states the law of conservation of messages:

- \square every message m_{ij} that is recorded as sent in the local state of sender p_i must be captured
 - \Box in the state of the channel C_{ij}
- or in the collected local state of the receiver p_i

☐ Condition C2 states that:

- ☐ in the collected global state, for every effect, its cause must be present
- \square if a message m_{ij} is not recorded as sent in the local state of p_i , then it must neither be present in the state of the channel C_{ij} nor in the collected local state of the receiver p_i

Course No.- SS ZG526, Course Title - Distributed Computin

BITS Pilani, Hyderabad Campu

Chandy-Lamport Approach

- uses a control message called marker
- ☐ after a process has recorded its snapshot, it sends a marker along all of its outgoing channels before sending out any more messages
- ☐ all messages that follow a marker on a channel have been sent by the sender after it took its snapshot
- channels are FIFO
- marker separates the messages in the channel into those to be included in the snapshot from those not to be recorded
- ☐ a process must record its snapshot no later than when it receives a marker on any of its incoming channels

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Cam

Chandy-Lamport Algorithm

Marker sending rule for process p

(1) Process p_i records its state.

(2) 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_i

On receiving a marker along channel C:

if p_i 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 ρ_{js} state was recorded and before ρ_{l} received the marker along C

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Camp

Chandy-Lamport Algorithmic approach

Putting together recorded snapshots

- global snapshot creation at initiator:
 - \Box each process can send its local snapshot to the initiator of the algorithm
- ☐ global snapshot creation at all processes:
 - $\hfill \Box$ each process sends the information it records along all outgoing channels
 - each process receiving such information for the first time propagates it along its outgoing channels
 - ☐ all the processes can determine the global state

Course ID: SS ZG526, Title: Distributed Computin

BITS Pilani, Hyderabad Ca

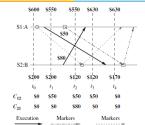
Chandy-Lamport Algorithm contd..

- ☐ algorithm can be initiated by any process by executing the marker sending rule
- ☐ **termination criterion** each process has received a marker on all of its incoming channels
- ☐ if multiple processes initiate the algorithm concurrently
 - each initiation needs to be distinguished by using unique markers

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campus

Chandy-Lamport Algorithm-scenario1



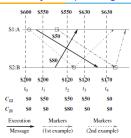
- ❖ Markers shown by dashed-and-dotted
- \$1 initiates the algorithm just after t.
- \$1 records its local state (account A=\$550) and sends a marker to \$2
- and sends a marker to 52
 marker is received by S2 after t₄
- when S2 receives the marker, it records its local state (account B=\$170), state C₁₂ as \$0, and sends a marker along C₂₂
- and sends a marker along C₂₁

 ❖ when S1 receives this marker, it records the state of C₂₁ as \$80
- ❖ \$800 amount in the system is conserved in the recorded global state
 ❖ A = \$550, B = \$170, C₁₂ = \$0, C₂₁ = \$80

se ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad (

Chandy-Lamport Algorithm-scenario2



- Markers shown using dotted arrows
- \$1 initiates the algorithm just after t₀ and before sending the \$50 for \$2
- ❖ S1 records its local state (account A = \$600) and sends a marker to S2
- * marker is received by S2 between t₂ and t₃
- when S2 receives the marker, it records its local state (account B = \$120), state of C_{12} as \$0, and sends a marker along C₂₁
 when S1 receives this marker, it records the
- state of C₂₁ as \$80
- \$\$800 amount in the system is conserved in the recorded global state

❖ A = \$600, B = \$120, C₁₂ = \$0, C₂₁ = \$80

Snapshot Algorithms for



- FIFO system
- ensures that all messages sent after a marker on a channel will be delivered after the marker
- in a non-FIFO system
- a marker cannot be used to differentiate messages into those to be recorded in the global state from those not to be recorded in the global state
- non-FIFO algorithm by Lai and Yang
 - use message piggybacking to distinguish computation messages sent after the marker from those sent before the marker



- fulfills the role of marker using a coloring scheme ☐ Coloring Scheme:
- every process is initially white
- process turns red while taking a snapshot
- ☐ equivalent of the "marker sending rule" is executed when a process turns red
- $\hfill \Box$ every message sent by a white process is colored white
 - $oldsymbol{\square}$ a white message is a message that was sent before the sender of that message recorded its local snapshot

Lai-Yang Algorithm



every message sent by a red process is colored red

- a red message is a message that was sent after the sender of that message recorded its local snapshot
- every white process takes its snapshot no later than the instant it receives a red message

Lai-Yang Algorithm contd..

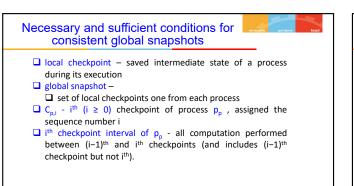
☐ Coloring Scheme:

- $\hfill \square$ when a white process receives a red message, it records its local snapshot before processing the message
- $\ \square$ ensures that
 - lacksquare no message sent by a process after recording its local snapshot is processed by the destination process before the destination records its local snapshot
- ☐ an explicit marker message is not required
- marker is piggybacked on computation messages using a coloring scheme

Lai-Yang Algorithm contd..

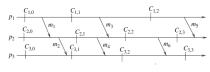


- ❖ Lai-Yang algorithm fulfills this role of the marker for channel state computation as follows:
 - every white process records a history of all white messages sent or received by it along each channel
 - ❖ when a process turns red, it sends these histories along with its snapshot to the initiator process that collects the global snapshot
 - initiator process evaluates transit(LS_i, LS_i) to compute the state of a channel C_{ij} as:
 - $SC_{ij} = \{\text{white messages sent by } p_i \text{ on } C_{ij}\} \{\text{white messages } \}$ received by p_j on C_{ij}}



Necessary and sufficient conditions for consistent global snapshots a causal path exists between checkpoints C_{i,x} and

- \square a causal path exists between checkpoints $C_{i,x}$ and $C_{j,y}$ if a sequence of messages exists from $C_{i,x}$ to $C_{j,y}$ such that each message is sent after the previous one in the sequence is received
- ☐ a zigzag path between two checkpoints is a causal path, however, allows a message to be sent before the previous one in the path is received



Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campus

Necessary and sufficient conditions for consistent global snapshots



- necessary condition for consistent snapshot absence of causal path between checkpoints in a snapshot
- necessary and sufficient conditions for consistent snapshot absence of zigzag path between checkpoints in a snapshot

Course ID: SS ZG526, Title: Distributed Computing 29 BITS Pilani, Hyde

Recap Quiz



- 1. Which of the following is not a type of event in distributed computing environments? (a) Internal (b) external (c) message send (d) message receive
- 2. If the message ordering is preserved in the distributed computing environment, then this system model is called (a)non-FIFO (b) LIFO (c) queue (d) FIFO
- 3. The control message used in Chandy-Lamport algorithm is called
 (a) Master (b) snapshot (c) marker (d) checkpoint
- 4. Message piggybacking is used in snapshot algorithms for ____ channels (a)non-FIFO (b) LIFO (c) queue (d) FIFO
- 5. The Lai-Yang algorithm for non-FIFO channels uses ______ instead of a marker (a) Checkpointing (b) colouring (c) creating (d) initiating

Course No. 15 TOTAL Course Title Distributed Courseling

BITS Pilani, Hyderabad Camp

Recap Quiz - key Q1 Q2 Q3 Q4 Q5 b d c a c Course No.-15/20256, Course Title - Distributed Computing 31 BTS Plant, Hydenslad Campu

Reference

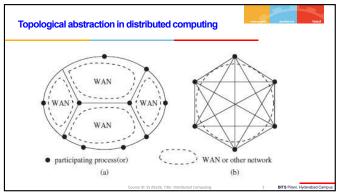


□Ajay D. Kshemkalyani, and Mukesh Singhal, Chapter 4, "Distributed Computing: Principles, Algorithms, and Systems", Cambridge University Press, 2008.

Course No. - SS ZG526, Course Title - Distributed Computing

BITS Pilani, Hyderabad Campu





Centralized and distributed algorithms



- centralized algorithm
 - predominant amount of work is performed by one (or possibly a few) processors
 - other processors play a relatively smaller role in accomplishing the joint task
 - other processors usually request or supply information
- distributed algorithm each processor plays an equal role in sharing the message overhead, time overhead, and space overhead

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campi

Symmetric and asymmetric algorithms



- **symmetric algorithm** algorithm in which all the processors execute the same logical functions
- asymmetric algorithm algorithm in which different processors execute logically different functions

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Camp

Synchronous and asynchronous systems



- synchronous system satisfies the following properties:
- a known upper bound on the message communication delay
- a known bounded drift rate for the local clock of each processor with respect to real-time
- a known upper bound on the time taken by a process to execute a logical step in the execution
- asynchronous system a system in which none of the above 3 properties of synchronous systems are satisfied

Failure models



- Failure model specifies the manner in which the component(s) of the system may fail
- Process failure models:
- > Fail-stop a properly functioning process may fail by stopping execution from some instant, other processes can learn that the process has failed
- Crash a properly functioning process may fail by stopping to function from any instance, other processes do not learn of this crash
 Byzantine or malicious failure, with authentication a process may exhibit any
- arbitrary behavior, authentication can be used to detect forgery
- Byzantine or malicious failure a process may exhibit any arbitrary behavior, no authentication techniques are applicable

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campus

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Camp



- ❖ undirected unweighted graph G = (N, L) represents topology
- ❖ n = |N|
- **❖** /= |L|
- ❖ diameter of a graph
 - minimum number of edges that need to be traversed to go from any node to any other node
 - ❖ diameter = $max_{i, j \in N}$ {length of the shortest path between i and j}
 - For a tree embedded in the graph, its depth is denoted as h

Course ID: SS ZG526, Title: Distributed Computi

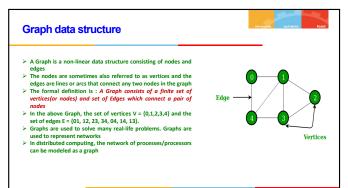
BITS Pilani, Hyderabad Campu

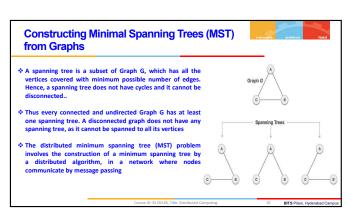
• Node Information: • Adjacent links • Neighboring nodes. • Unique identity. • Graph model: undirected graph.

Nodes V

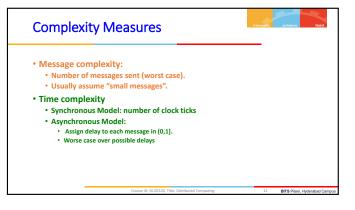
• Edges E
• Diameter D

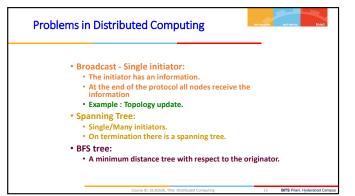
urse ID: SS ZG526, Title: Distributed Computing 8 BITS Pilani, Hy

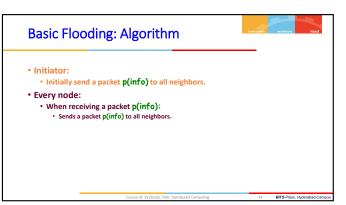


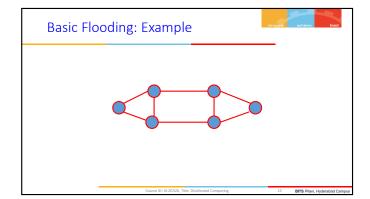


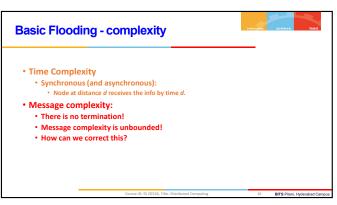
Synchronous Model: There is a global clock. Packet can be sent only at clock ticks. Packet sent at time t is received by t+1. Asynchronous model: Packet sent is eventually received.

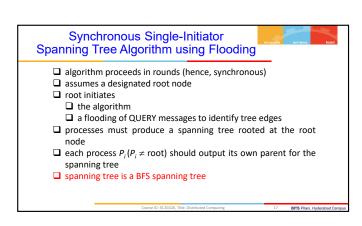


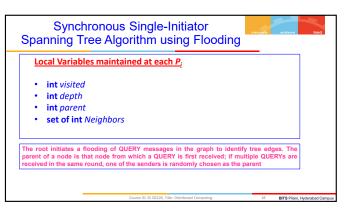


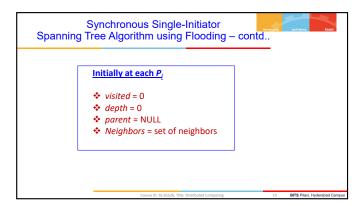


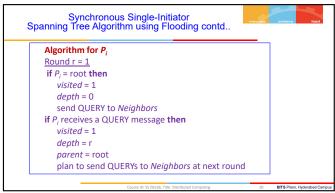








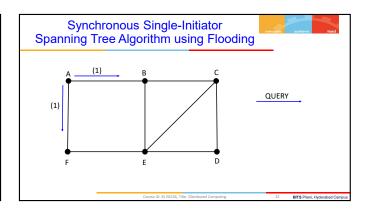


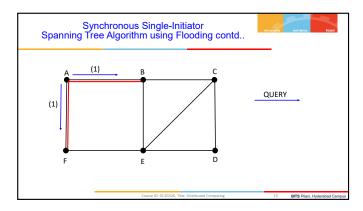


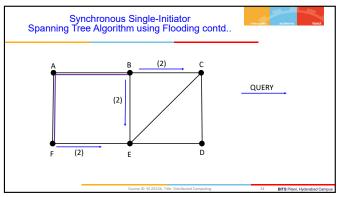
Synchronous Single-Initiator Spanning
Tree Algorithm using Flooding

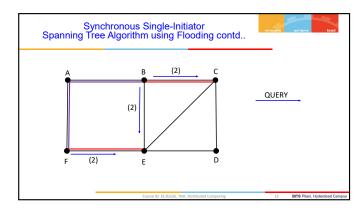
Algorithm for P_i
Round r > 1 and r <= diameter

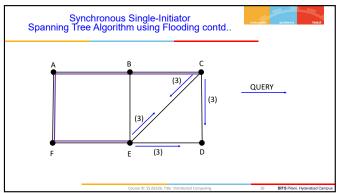
if P_i planned to send in previous round then
P_i sends QUERY to Neighbors
if visited = 0 then
if P_i receives QUERY messages then
visited = 1
depth = r
parent = any randomly selected node from which QUERY was received plan to send QUERY to Neighbors \ {senders of QUERYs received in r}

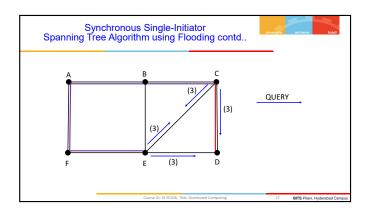


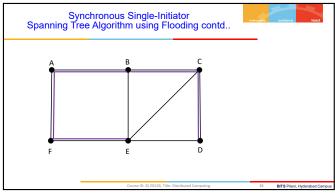


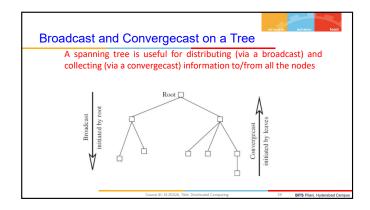


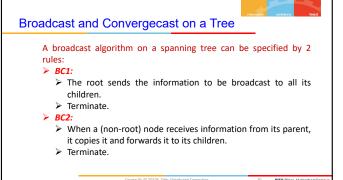












Broadcast and Convergecast on a Tree $\hfill \square$ A convergecast algorithm collects information from all the nodes at the root node in order to compute some global function ☐ It is initiated by the leaf nodes of the tree, usually in response to receiving a request sent by the root using a broadcast

Broadcast and Convergecast on a Tree

The convergecast algorithm is specified as follows:

- **❖** CVC1:
 - Leaf node sends its report to its parent.
 - Terminate.
- CVC2:
 - At a non-leaf node that is not the root: When a report is received from all the child nodes, the collective report is sent to the parent.
 - Terminate.
- CVC3:
 - At the root: When a report is received from all the child nodes, the global function is evaluated using the reports.
 - Terminate.

Broadcast and Convergecast on a Tree

Complexity

- \square each broadcast and each convergecast requires n-1 messages
- lacksquare each broadcast and each convergecast requires time equal to the maximum height h of the tree, which is O(n)

Applications:

- $\hfill \square$ computation of minimum of integer variables associated with the nodes of an application
- □ leader election

Synchronizers

- It is difficult to design algorithms for asynchronous systems
- ❖ solution simulate synchronous behavior on an asynchronous system
- synchronizers transformation algorithms to run synchronous algorithms on asynchronous systems
- synchronizer signals to each process when it is sure that all messages to be received in the current round have arrived and it is safe to proceed to the next round

Synchronizers

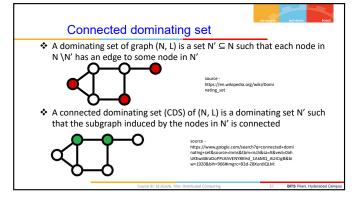
- requires each process to send every neighbor only one message in each round
- $\hfill \square$ if no message is to be sent in the synchronous algorithm, an empty dummy message is sent in the asynchronous algorithm
- $\hfill \square$ if more than one messages are sent in the synchronous algorithm, they are combined into one message in the asynchronous algorithm in any round, when a process receives a message from each
- neighbor, it moves to the next round α, β, and γ synchronizers: (rounds, rooted spanning tree; n/w of clusters)
 use the notion of process safety
 - a process i is safe in round r if all messages sent by i in round r have
 - been received

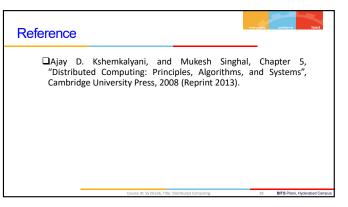
Maximal Independent Set

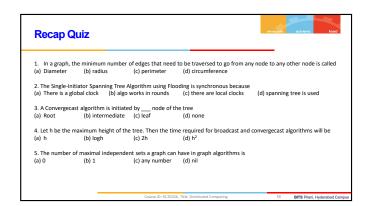
- \triangleright For a graph (N, L), an independent set of nodes N', where N' \subset N, is such that for each i and j in N', $(i, j) \notin L$
- An independent set N' is a maximal independent set if no strict superset of N' is an independent set
- A graph may have multiple maximal independent sets
- Adjacent nodes must not be chosen

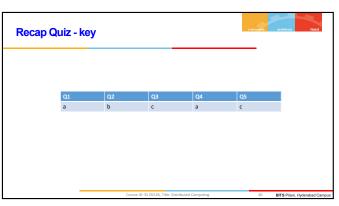


http://www.martinbroadhurst.com/ greedy-max-independent-set-in-c.html

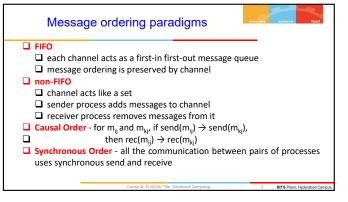












Asynchronous executions (A-executions)





(a) An A-execution that is not a FIFO execution (b) An A-execution that is also a FIFO

FIFO execution

for all (s, r) and $(s', r') \in \mathcal{T}$, $(s \sim s')$ and $(s', r') \implies r \prec r'$.

- ❖FIFO logical channel over a non-FIFO channel would use a separate numbering scheme to sequence the messages on each logical channel
- ❖The sender assigns and appends a <u>sequence_num</u>, <u>connection_id</u> tuple to each message
- The receiver uses a buffer to order the incoming messages as per the sender's sequence numbers, and accepts only the "next" message in sequence

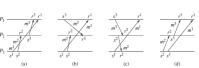
Causally Ordered (CO) executions



for all (s, r) and $(s', r') \in T$, $(r \sim r')$ and $s \prec s' \implies r \prec r'$.

- \blacktriangleright If two send events s and s' are related by causality ordering (not physical time ordering), then a causally ordered execution requires that their corresponding receive events \boldsymbol{r} and \boldsymbol{r}' occur in the same order at all common destinations.
- ▶If s and s' are not related by causality, then CO is vacuously satisfied because the antecedent of the implication is false

Illustration of causally ordered executions



(a) Not a CO execution. (b),(c), and (d) CO executions.

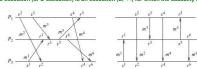
Group Communication > Group communication needs to be supported

group, of the processes in the system unicasting - point-to-point message communication

distributed system

Synchronous (SYNC) execution - S executiom



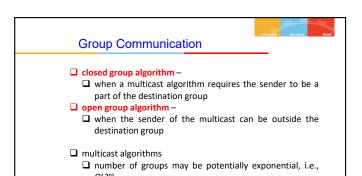


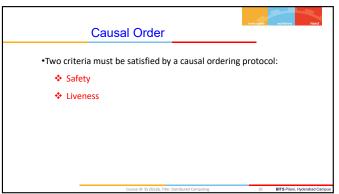
(a) Execution in an asynchronous system (b) Equivalent instantaneous communication

- ☐ When all the communication between pairs of processes uses synchronous send and receive primitives, the resulting order is
- the synchronous order
 □ As each synchronous communication involves a handshake between the receiver and the sender, the corresponding send and receive events can be viewed as occurring instantaneously and atomically
 □ In a timing diagram, the "instantaneous" message communication can be shown by bidirectional vertical message lines

> message broadcast - sending a message to all members in the

multicasting - a message is sent to a certain subset, identified as a







Safety

- > a message M arriving at a process may need to be buffered until all system-wide messages sent in the causal past of the send(M) event to the same destination have already arrived
- > distinction is made between
 - > arrival of a message at a process
 - > event at which the message is given to the application process

> A message that arrives at a process must eventually be delivered to the process

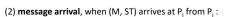
Raynal-Schiper-Toueg Algorithm seach message M should carry a log of all other messages or their identifiers, sent causally before M's send event, and sent to the same destination dest(M) $\ensuremath{\clubsuit}$ log can be examined to ensure whether it is safe to deliver a message channels are FIFO each send event to unicast, multicast, or

Raynal-Schiper-Toueg Algorithm contd.

local variables:

- array of int SENT[1 n, 1 n] (n x n array)
 - SENT_i[j, k] = no. of messages sent by P_i to P_k as known to P_i
- array of int DELIV [1 n]
 - DELIV_i[j] = no. of messages from P_j that have been delivered to P_i
 - DELIV[j] = no. of messages sent by P_j that are delivered locally
 - (1) send event, where P_i wants to send message M to P_i: (1a) send (M, SENT) to Pi (1b) SENT[i, j] = SENT[i, j] + 1

Raynal-Schiper-Toueg Algorithm



(2a) deliver M to P_i when for each process x,

 $DELIV[x] \ge ST[x, i]$

(2c) \forall x, y, SENT[x, y] = max(SENT[x, y], ST[x, y])

(2d) DELIV[j] = DELIV[j] + 1

Raynal—Schiper—Toueg Algorithm Complexity: space requirement at each process: $O(n^2)$ integers space overhead per message: n^2 integers time complexity at each process for each send and deliver event: $O(n^2)$

Raynal-Schiper-Toueg Algorithm

- P₁ sent 4 messages to P₂
- P₂ sent 3 messages to P₁
- P₁ → DELIV₁[0 2] // 2 messages from P₂ have been delivered to P₁
- P₂→ DELIV₂[4 0] // all 4 messages from P₁ have been delivered to P₂
- · Now if,
 - P₁ sends a message to P₂ → no buffering required
 - P₂ sends a message to P₁ → buffering required

Course ID: SS ZG526, Title: Distributed Computing

DITE Disail Madesaked Com-

Birman-Schiper-Stephenson Protocol



- A message is delivered to a process only if the message preceding it has been delivered
- buffer is needed for pending deliveries
- each message has a vector that contains information for the recipient to determine if another message preceded it
- all messages are broadcast

se ID: SS 2G526, Title: Distributed Computing 17 BITS Pilani, Hyc

Birman-Schiper-Stephenson Protocol



- C_i = vector clock of P_i
- $C_i[j] = j^{th}$ element of C_i
 - contains P_i's latest value for the current time in P_i
- ❖ tm = vector timestamp for message m
- stamped after local clock is incremented

Up to 70000 Table Distributed Computing

Birman-Schiper-Stephenson Protocol

- ☐ P; sends a message m to P;
- ☐ P_i increments C_i[i]
- \square P_i sets the timestamp tm = C_i for message m

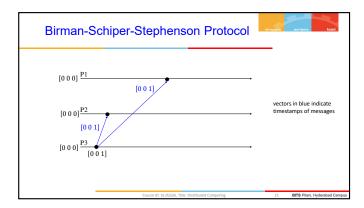
Birman-Schiper-Stephenson Protocol

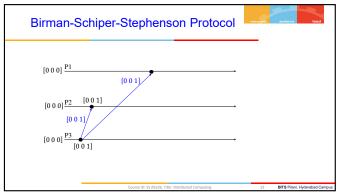


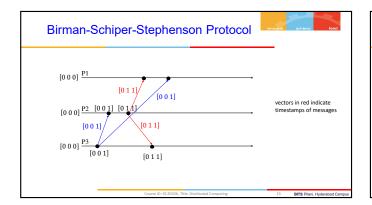
P_i receives a message m from P_i

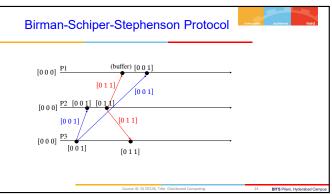
- when P_j (j ≠ i) receives m with timestamp tm, it delays message delivery until:
 - C_i[i] = tm[i] 1 //P_i has received all preceeding messages sent by P_i
 - ∀ k <= n and k ≠ i, C_i[k] ≥ tm[k] //P_j has received all the messages that were received at P_i from other processes before P_i sent m
- when m is delivered to P_j, update P_j's vector clock
 - $\forall i, C_j[i] = \max(C_j[i], tm[i])$
- check buffered messages to see if any can be delivered

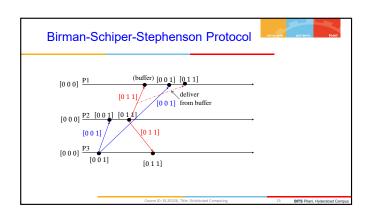
BITS Pilani, Hyderabad Camp

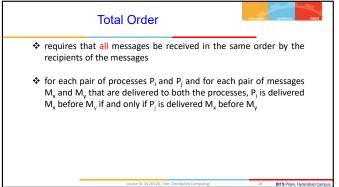












Classification of Application-Level Multicast Algorithms Communication history-based algorithms use a part of the communication history to guarantee ordering requirements eg., RST (Raynal, Schiper. and Toueg), KS (Kshemkalyani and Singhal) algorithms -> causal ordering Lamport's algorithm -> scalar timestamps NewTop protocol - extension of Lamport's algorithm

Classification of Application-Level Multicast Algorithms

- innovate achieve
- Privilege-based algorithms
- A token circulates among sender processes
- The token carries the sequence number for the next message to be multicast
- only the token-holder can multicast
- after a multicast send event, the sequence number is updated
- destination processes deliver messages in the order of increasing sequence numbers
- senders need to know the other senders, hence closed groups are assumed
- can provide total ordering and causal ordering using closed group configuration
- not scalable as do not permit concurrent send events

Course ID- SS 7CF3S. Title: Distributed Committee

BITE Diesi Mulaushad Commi

Classification of Application-Level Multicast Algorithms Senders Privilege rotates Destinations Privilege-based algorithm

Classification of Application-Level Multicast Algorithms



Moving sequencer algorithms

- sender sends the message to all the sequencers to multicast a message
- sequencers circulate a token among themselves
- token carries
- · a sequence number
- a list of all the messages for which a sequence number has already been assigned – such messages have been sent already
- when a sequencer receives the token
 - it assigns a sequence number to all received but un-sequenced messages
 - it sends the newly sequenced messages to the destinations
 - inserts these messages into the token list
 - passes the token to the next sequencer

Classification of Application-Level Multicast Algorithms Moving sequencer algorithms contd... destination processes deliver messages received in the order of increasing sequence number sequencers Token Token

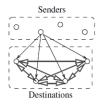
Classification of **Application-Level Multicast Algorithms** Fixed sequencer algorithms simplified version of the Senders previous class of algorithms a single sequencer is present Fixed sequencer centralized algorithms stee., propagation tree approach, ISIS sequencer, Amoeba, Phoenix, Destinations Newtop's asymmetric algorithm

Classification of Application-Level Multicast Algorithms

innovate achieve lead

Destination agreement algorithms

- destinations receive the messages with some limited ordering information
- ☐ destination processes then exchange information among themselves to define an order
- Two sub-classes
 - uses timestamps
 - uses an agreement or consensus protocol among the processes



urse ID: SS ZG526, Title: Distributed Computing 33 BIT

Termination detection using distributed snapshots



- consistent snapshot of a distributed system captures stable properties
- termination of a distributed computation is a stable property
- if a consistent snapshot of a distributed computation is taken after the distributed computation has terminated, the snapshot will capture the termination of the computation
- all processes will have become idle
- there will be no in-transit messages

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Carr

Termination detection by weight throwing



- a process called controlling agent monitors the computation
- communication channel exists between each of the processes and the controlling agent and also between every pair of processes
- initially, all processes are in the idle state
- weight at each process is zero and the weight at the controlling agent is 1
- computation starts when the controlling agent sends a basic message to one of the processes
- process becomes active and the computation starts
- $\ ^{\diamondsuit}$ non-zero weight W (0 <W \leq 1) is assigned to each process in the active state

Course ID: SS ZG526, Title: Distributed Computing

BITS Pilani, Hyderabad Campu

Termination detection by weight throwing contd...



- when a process sends a message, it sends a part of its weight in the message
- when a process receives a message, it add the weight received in the message to its weight
- > sum of weights on all the processes and on all the messages in transit 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
- > controlling agent concludes termination if its weight becomes 1

Course ID: SS ZG526, Title: Distributed Computin

BITS Pilani, Hyderabad Cam

Spanning-tree-based termination detection



- ❖ Assumes N processes P_i , $0 \le i \le N$
- ❖ processes → nodes, channels → edges
- uses a fixed spanning tree of the graph with process P₀ at its root which is responsible for termination detection
- P₀ communicates with other processes to determine their states
- messages used for this purpose are called signals
- all leaf nodes report to their parents, if they have terminated
- an intermediate node will report to its parent when it has completed processing and all of its immediate children have terminated, and so
- root concludes that termination has occurred, if it has terminated and all of <u>its immediate children</u> have also terminated

BITS Pilani, Hyderabad Campus

Recap Quiz

1. In the multicast algorithms, the number of groups may grow to the order of (a) logn (b) log²n (c) n^2 (d) 2^n

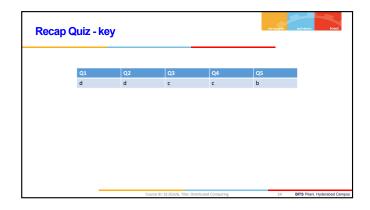
In which of the following message ordering paradigms, the channel acts as a se
 (a) Synchronous (b) causal order (c) FIFO (d) non-FIFO

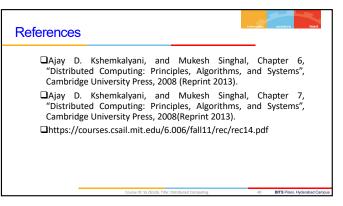
3. In the Raynal–Schiper–Toueg Algorithm, the channel is organized as (a) Synchronous (b) causal order (c) FIFO (d) non-FIFO

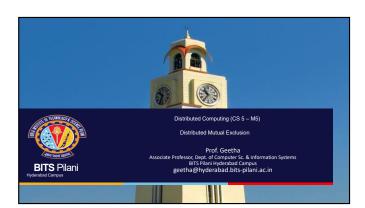
4. Which of the following is not a classification of application-level multicast algorithms?
(a) Fixed sequencer (b) moving sequencer (c) source agreement (d) destination agreement

5. The messages used in the spanning tree based termination detection algorithms are called (a) Markers (b) signals (c) triggers (d) events

BITS Pilani, Hyderabad Camp







What is mutual exclusion?



- ❖Mutual Exclusion also known as Mutex was first identified by Dijkstra
- When a process is accessing a shared variable, it is said to be in a critical section (code segment)
- When no two processes can be in Critical Section at the same time, this state is known as Mutual Exclusion which is a property of concurrency control and is used to prevent race condition (happens when the result of multiple thread execution in critical section differs according to the order in which the threads execute)

Mutual Exclusion Devices:-

Locks, recursive locks, semaphores, monitor, message passing etc.,

Course ID: SS ZG526, Title: Distributed Computing

Mutual exclusion in synchronization



- Mutual exclusion is a property of process synchronization which states that "no two processes can exist in the critical section at any given point of time"
- Any process synchronization technique being used must satisfy the property of mutual exclusion, without which it would not be possible to get rid of a race condition
- During concurrent execution of processes, processes need to enter the critical section (or the section
 of the program shared across processes) at times for execution
- It might so happen that because of the execution of multiple processes at once, the values stored in the critical section become inconsistent; in other words, the values depend on the sequence of execution of instructions – also known as a race condition.
- * The primary task of process synchronization is to get rid of race conditions while executing the

Mutual Exclusion in Distributed Computing



Mutual exclusion is a concurrency control property which is introduced to prevent race conditions; the requirement that a process can not enter its critical section while another concurrent process is currently present or executing in its critical section i.e., only one process is allowed to execute the critical section at any given instance of time

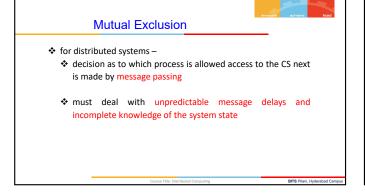
In single computer system, memory and other resources are shared between different processes. The status of shared resources and the status of users is easily available in the shared memory so with the help of shared variable (for example: Semaphores) mutual exclusion problem can be easily solved.

In Distributed systems, we neither have shared memory nor a common physical clock and therefore we can not solve mutual exclusion problem using shared variables. To solve the mutual exclusion problem in distributed systems, message passing is used

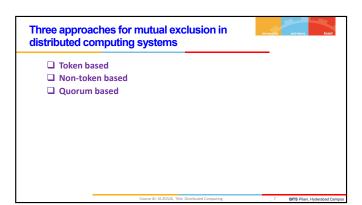
A site in distributed system does not have the complete information of the state of the system due to lack of shared memory and a common physical clock

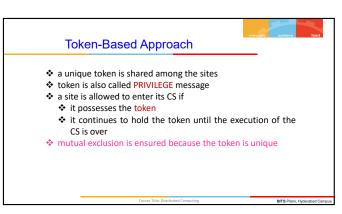
Course ID: SS ZG526, Title: Distributed Computing

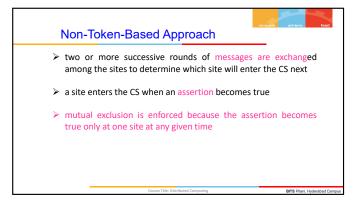
BITS Pilani, Hyderabad Camp

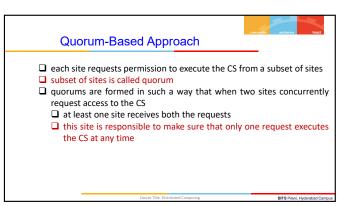


Mutual Exclusion | ensures that concurrent access of processes to a shared resource or data is serialized | executed in a mutually exclusive manner | for distributed system – | only one process is allowed to execute the critical section (CS) at any given time | semaphores or a local kernel cannot be used to implement mutual exclusion











- ✓ system consists of N sites, S₁, S₂,...., S_N
- without loss of generality, assume that a single process is running on each site
- process at site S_i is p_i
- processes communicate asynchronously over an underlying communication network
- ✓ any process wishing to enter the CS
 - ✓ requests all other or a subset of processes by sending REQUEST
- ✓ waits for appropriate replies before entering the CS
- ✓ while waiting the process is not allowed to make further requests to enter the CS

System Model for Mutual Exclusion ☐ site can be in one of the following 3 states: requesting the CS executing the CS neither of the 2 ☐ "requesting the CS" state - site is blocked and cannot make further requests for the CS "idle" state - site is executing outside the CS ☐ for token-based algorithms $\hfill \square$ a site can also be in a state where a site holding the token is executing outside the CS

System Model for mutual exclusion contd.



- 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
- nature of channels (FIFO or not) is algorithm specific
- assume that:
 - channels reliably deliver all messages
 - sites do not crash
 - network does not get partitioned
- * timestamps are used to decide the priority of requests in case of a
- ❖ general rule smaller the timestamp of a request, the higher its priority to execute the CS

System Model contd..

☐ such state is called *idle token state*



Notations:

- ☐ N number of processes or sites involved in invoking the critical section
- ☐ T average message delay
- ☐ E average critical section execution time

Requirements of mutual exclusion algorithms



No Deadlock:
Two or more sites should not endlessly wait for any message that will never arrive

UN o starvation:
Every site which wants to execute critical section should get an opportunity to execute it in finite time.
Any site should not wait indefinitely to execute critical section while other sites are repeatedly executing the critical section

LFairness:
Each site should get a fair chance to execute the critical section. Any request to execute the critical section must be executed in the order they are made i.e Critical section execution requests should be executed in the order of their arrival in the system

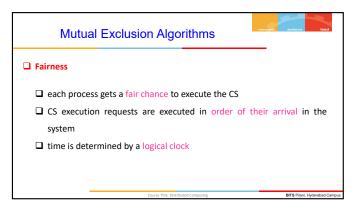
In case of failure, it should be able to recognize it by itself in order to continue functioning without any

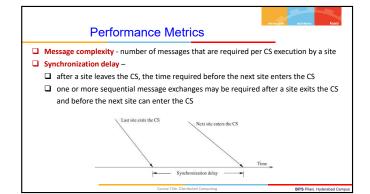
Properties of Mutual Exclusion Algorithms

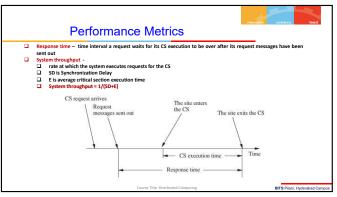


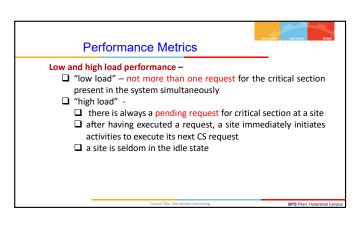
- Safety Property absolutely necessary
- Liveness Property important
- ❖ Fairness important

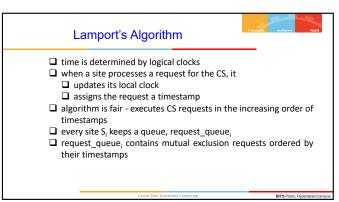
Requirements of Mutual Exclusion Algorithms Safety property at any instant, only one process can execute the critical section absolutely necessary property Liveness property absence of deadlock and starvation a site must not wait indefinitely to execute the CS while other sites are repeatedly executing the CS every requesting site should get an opportunity to execute the CS in finite time











Lamport's Algorithm



- messages in FIFO order

 when a site removes a request from its request queue
 - its own request may come at the top of the queue
 - enables it to enter the CS
- when a site receives a REQUEST, REPLY, or RELEASE message
 - $\ensuremath{ \stackrel{\bullet}{\bullet} }$ it updates its clock using the timestamp in the message

urse Title: Distributed Computing

BITS Pilani, Hyderabad Campu

Lamport's 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, it places site S_i's request on request_queue_j and returns a timestamped REPLY message to S_i

Course Title: Distributed Computing

BITS Pilani, Hyderabad Can

Lamport's Algorithm

☐ Executing the critical section:

- ☐ Site S_i enters the CS when the following two conditions hold:
 - $\hfill\Box$ 11: S_i has received a message with timestamp larger than \hfill (ts, i) from all other sites
 - ☐ L2: S_i's request is at the top of the 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
- \square When a site S_j receives a RELEASE message from site S_i , it removes S_i 's request from its request queue

Course Title: Distributed Computing

BITS Pilani, Hyderabad Cam

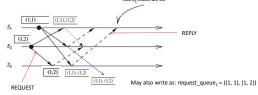
Lamport's Algorithm



Sites S₁ and S₂ make requests for the CS

BITS Pilani, Hyderabad Cam

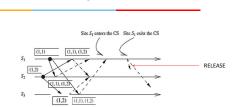
Lamport's Algorithm



Site S₁ enters the CS

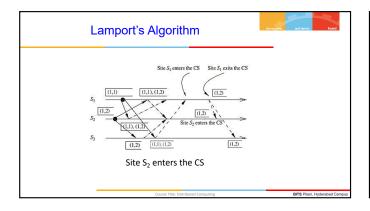
BITS Pilani, Hyderabad Campu

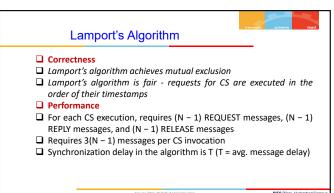
Lamport's Algorithm



Site S₁ exits the CS and sends RELEASE messages

e Title: Distributed Computing





Ricart-Agrawala Algorithm

- communication channels are not required to be FIFO
- uses two types of messages:
 - ❖ REQUEST
 - ❖ 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

Ricart-Agrawala Algorithm

- *timestamps are used to decide the priority of requests in case of conflict -
- if a process p_i that is waiting to execute the CS receives a REQUEST message from process \boldsymbol{p}_{j} , then if the priority of $\boldsymbol{p}_{j}\text{'s}$ request is lower, p_i defers the REPLY to p_i and sends a REPLY message to p_i only after executing the CS for its pending request
- •otherwise, p_i sends a REPLY message to p_i immediately, provided it is currently not executing the CS
- •if several processes are requesting execution of the CS, the highest priority request succeeds in collecting all the needed REPLY messages and gets to execute the CS

Ricart-Agrawala Algorithm

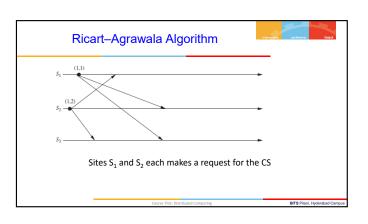
- each process p_i maintains the request-deferred array, RD_i
- size of RD_i = no. of processes in the system
- ❖ initially, $\forall i \forall j$: $RD_i[j] = 0$
- whenever p_i defers the request sent by p_i, it sets RD_i[j] = 1,
- after it has sent a REPLY message to p_i, it sets RD_i[j] = 0
- when 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
 - assigns a timestamp to the request
- execution of the CS requests is always in the order of their timestamps

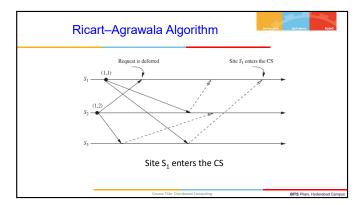
Ricart-Agrawala Algorithm

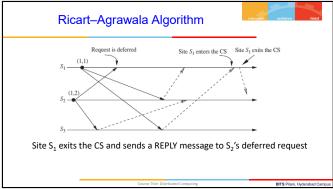
Requesting the critical section:

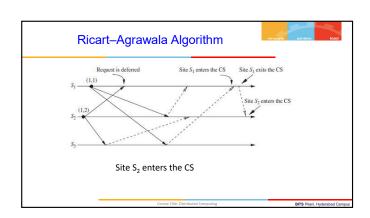
- (a) When a site \boldsymbol{S}_{i} wants to enter the CS, it broadcasts a timestamped REQUEST message to all other sites
- (b) When site S_i receives a REQUEST message from site S_i, it sends a REPLY message to site S_i if site S_i is neither requesting nor executing the CS, or if the site S_j is requesting 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_i sets $RD_i[i] = 1$

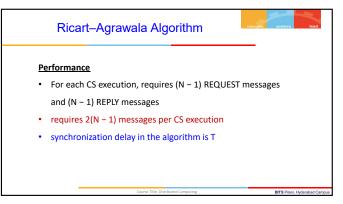
Ricart—Agrawala Algorithm Executing the critical section: (c) Site S_i enters the CS after it has received a REPLY message from every site it sent a REQUEST message to Releasing the critical section: (d) When site S_i exits the CS, it sends all the deferred REPLY messages: ∀j if RD_i [j] = 1, then S_i sends a REPLY message to S_j and sets RD_i [j] = 0 Correctness •Ricart—Agrawala algorithm achieves mutual exclusion











Maekawa's Algorithm

- innovate schieve lead
- quorum-based mutual exclusion algorithm
- request sets for sites (i.e., quorums) are constructed to satisfy the following conditions:
 - M1: $(\forall i \ \forall j : i \neq j, \ 1 \leq i, j \leq N :: R_i \cap R_i \neq \phi)$
 - M2: $(\forall i : 1 \le i \le N :: S_i \in R_i)$
 - M3: $(\forall i : 1 \le i \le N :: |R_i| = K \text{ for some } K)$
 - M4: Any site S_i is contained in K number of R_i 's, $1 \le i, j \le N$
- Maekawa showed that N = K(K 1) + 1
- This relation gives $|R_i| = K = \sqrt{N}$ (square root of N)

Course Title: Distributed Compu

BITS Pilani, Hyderabad Campi

Maekawa's Algorithm

- ☐ there is at least one common site between the request sets of any two sites (condition M1)
- every pair of sites has a common site which mediates conflicts between the pair
- □ a site can have only one outstanding REPLY message at any time; that is,
 □ it grants permission to an incoming request if it has not granted
 - permission to some other site mutual exclusion is guaranteed
- $\hfill \Box$ requires delivery of messages to be in the order they are sent between every pair of sites

Course Title: Distributed Computing

BUT 01 111 1 1 10

Maekawa's Algorithm



- ☐ conditions M1 and M2 are necessary for correctness
- $\hfill \Box$ M3 (vi:1sisN::|R|=K for some K) states that the size of the requests sets of all sites must be equal
 - equal amount of work to invoke mutual exclusion
- ☐ M4 (Any site S_j is contained in K number of R_i's, 1 ≤ i, j ≤ N) enforces that exactly the same number of sites should request permission from any site
 ☐ equal responsibility

Course Title: Distributed Computing

BITS Pilani, Hyderabad Campi

Maekawa's Algorithm



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 $\boldsymbol{S_i}$ executes the CS only after it has received a REPLY message from every site in $\boldsymbol{R_i}.$

Course Title: Distributed Computing

BITS Pilani, Hyderabad Ca

Maekawa's Algorithm



Releasing the critical section:

- (d) After the execution of the CS is over, site S_i sends a RELEASE(i) message to every site in R_i .
- (e) When a site S_j receives a RELEASE(i) message from site S_i , it sends a REPLY message to the next site 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.

Maekawa's Algorithm



Correctness: Maekawa's algorithm achieves mutual exclusion

Performance:

- size of a request set is \sqrt{N}
- an execution of the CS requires √N REQUEST, √N REPLY, and √N RELEASE messages, resulting in 3√N messages per CS execution
- synchronization delay 2T

Course Title: Distributed Computing

BITS Pilani, Hyderabad Cam

Reference

❖Ajay D. Kshemkalyani, and Mukesh Singhal, Chapter 9, "Distributed Computing: Principles, Algorithms, and Systems", Cambridge University Press, 2008 (reprint: 2013).



Suzuki-Kasami's Broadcast Algorithm



- ❖ if a site that wants to enter the CS does not have the token, it broadcasts a REQUEST message for the token to all other sites
- The site that 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 CS execution
- Two design issues exist:
 - $\ensuremath{ \diamondsuit}$ how to distinguish an outdated REQUEST message from a current REQUEST message ?
 - * how to determine which site has an outstanding request for the CS?

Suzuki-Kasami's Broadcast Algorithm



How to distinguish an outdated REQUEST message from a current **REQUEST message?**

- a site may receive a token request message after the corresponding request has been satisfied
- ☐ if a site cannot 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

Suzuki-Kasami's Broadcast Algorithm



- > 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
- > after the corresponding request for the CS has been satisfied at S_i, an issue is how to inform site S_i and all other sites efficiently

Suzuki-Kasami's Broadcast Algorithm



- outdated REQUEST messages are distinguished from current REQUEST messages in the following manner:
 - sn (sn = 1, 2....) is a sequence number that indicates that site S_i is requesting its snth CS execution
 - ❖ a site S_i keeps an array of integers RN_i[1, ... ,n] where RN_i[j] is the largest sequence number received in a REQUEST message so far
 - ❖ when site S_i receives a REQUEST(j, sn) message, it sets RN_i[j] = $max(RN_i[j], sn)$
 - ❖ when a site S_i receives a REQUEST(j,sn) message, the request is outdated if $RN_i[j] > sn$

Suzuki-Kasami's Broadcast Algorithm



- sites with outstanding requests for the CS are determined 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_i 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
 - token array LN[1, ...,n] permits a site to determine if a site has an outstanding request for the CS
 - at site S_i, if RN[i] = LN[i] + 1, then site S_i is currently requesting a token
 - * after executing the CS, a site checks this condition for all the j's to determine all the sites that are requesting the token and places their IDs in queue Q if these IDs are not already present in Q
 - finally, the site sends the token to the site whose ID is at the head of Q

Suzuki-Kasami's Broadcast 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_i receives this message, it sets $RN_i[i]$ to $max(RN_i[i], sn)$. If S_i has the idle token, then it sends the token to S_i if $RN_i[i] = LN[i]+1$.

Executing the critical section:

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

Suzuki-Kasami's Broadcast Algorithm



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

(d) It sets LN[i] element of the token array equal to RN_i[i].

- (e) For every site S_i whose ID is not in the token queue, it appends its ID to the token queue if $RN_i[j] = LN[j] + 1$.
- (f) If the token queue is non-empty 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

Suzuki-Kasami's Broadcast Algorithm



- Correctness
- ✓ mutual exclusion is guaranteed because there is only one token in the system and a site holds the token during the CS execution
- ✓ a requesting site enters the CS in finite time
- ✓ Performance
- √ if a site holds the token, no message is required
- √ if a site does not hold the token, N messages are required
- ✓ synchronization delay is 0 or T

Raymond's Tree-Based Algorithm

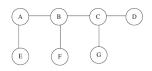


- ☐ uses a spanning tree of the computer network to reduce the number of messages exchanged per critical section execution
- $\hfill \square$ assumes that the underlying network guarantees message delivery
- ☐ time or order of message arrival cannot be predicted ☐ all nodes of the network are completely reliable
- ☐ a spanning tree of a network of N nodes will be a tree that contains all N nodes
- $\hfill \square$ a minimal spanning tree is a spanning tree with minimum cost
- cost function is based on the network link characteristics

Raymond's Tree-Based Algorithm



- > messages between nodes traverse along the undirected edges of the tree
- node needs to hold information about and communicate only to its immediate-neighboring nodes



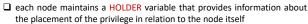
- de C holds information about and communica y to nodes B, D, and G de C does not need to know about the ot des A, E, and F for the operation of the algorith

Raymond's Tree-Based Algorithm



- ❖ only one node can be in possession of the privilege (called the privileged node) at any time
 - except when the privilege is in transit from one node to another in the form of a PRIVILEGE message
- when there are no nodes requesting for the privilege, it remains in possession of the node that last used it

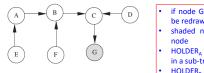
Raymond's Tree-Based Algorithm



- ☐ a node stores in its HOLDER variable the identity of a node that it thinks has the privilege or leads to the node having the privilege
- ☐ HOLDER variables of all nodes maintain directed paths from each node to the node in possession of the privilege
- ☐ for two nodes X and Y, if HOLDER_x = Y, the undirected edge between X and Y can be redrawn as a directed edge from X to Y

Raymond's Tree-Based Algorithm





- if node G holds the privilege, figure can be redrawn
- shaded node represents the privileged
- HOLDER_A = B (as the privilege is located in a sub-tree of A denoted by B)
- $HOLDER_B = C$
- $HOLDER_C = G$ $HOLDER_D = C$
- $HOLDER_E = A$
- HOLDER_F = B HOLDER_G = self

Raymond's Tree-Based Algorithm

G



- HOLDER_c, i.e., G privileged node G, if it no longer needs the privilege, sends the PRIVILEGE message to its neighbor C, which made a request for the
- privilege, and resets HOLDER₆ to C C forwards the PRIVILEGE to B, since it had requested the privilege on behalf of B
- C also resets HOLDER_C to B

Raymond's Tree-Based Algorithm



■ Data Structures

■ HOLDER

- lacktriangle possible values "self" or the identity of one of the immediate neighbors
- - possible values true or false
 - $\hfill \square$ indicates if the current node is executing the critical section

- possible values true or false
- indicates if node has sent a request for the privilege
- prevents the sending of duplicate requests for privilege

Raymond's Tree-Based Algorithm



Data Structures

- REQUEST_Q
 - FIFO queue that could contain "self " or the identities of immediate neighbors as elements
 - REQUEST_Q of a node consists of the identities of those immediate neighbors that have requested for privilege but have not yet been sent the privilege
 - maximum size of REQUEST_Q of a node is the number of immediate neighbors + 1 (for "self")

