

Unit 6-Coordination and Synchronization

Compiled by Prashant Gautam

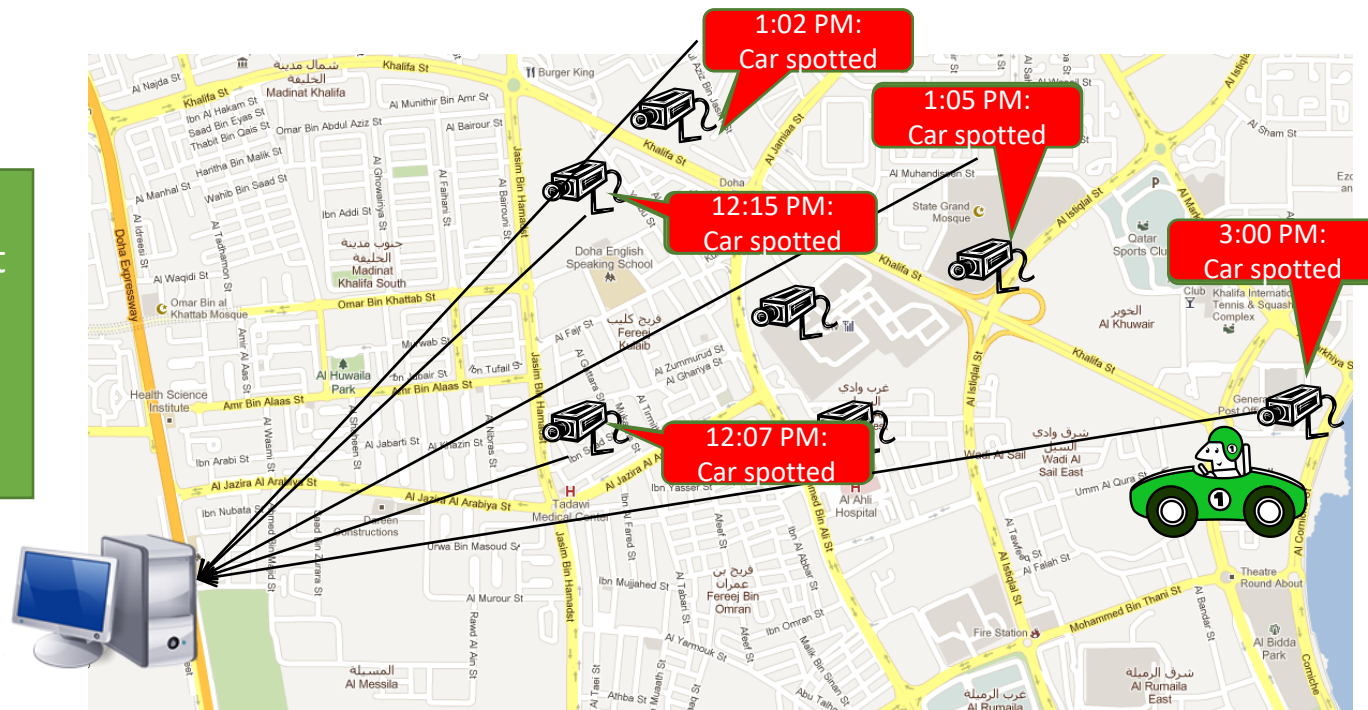
Synchronization

- Until now, we have looked at:
 - how entities are named and identified
 - how entities communicate with each other
- In addition to the above requirements, entities in DS often have to *cooperate* and *synchronize* to solve the problem correctly
 - e.g., In a distributed file system, processes have to synchronize and cooperate such that two processes are not allowed to write to the same part of a file

Need for Synchronization – Example 1

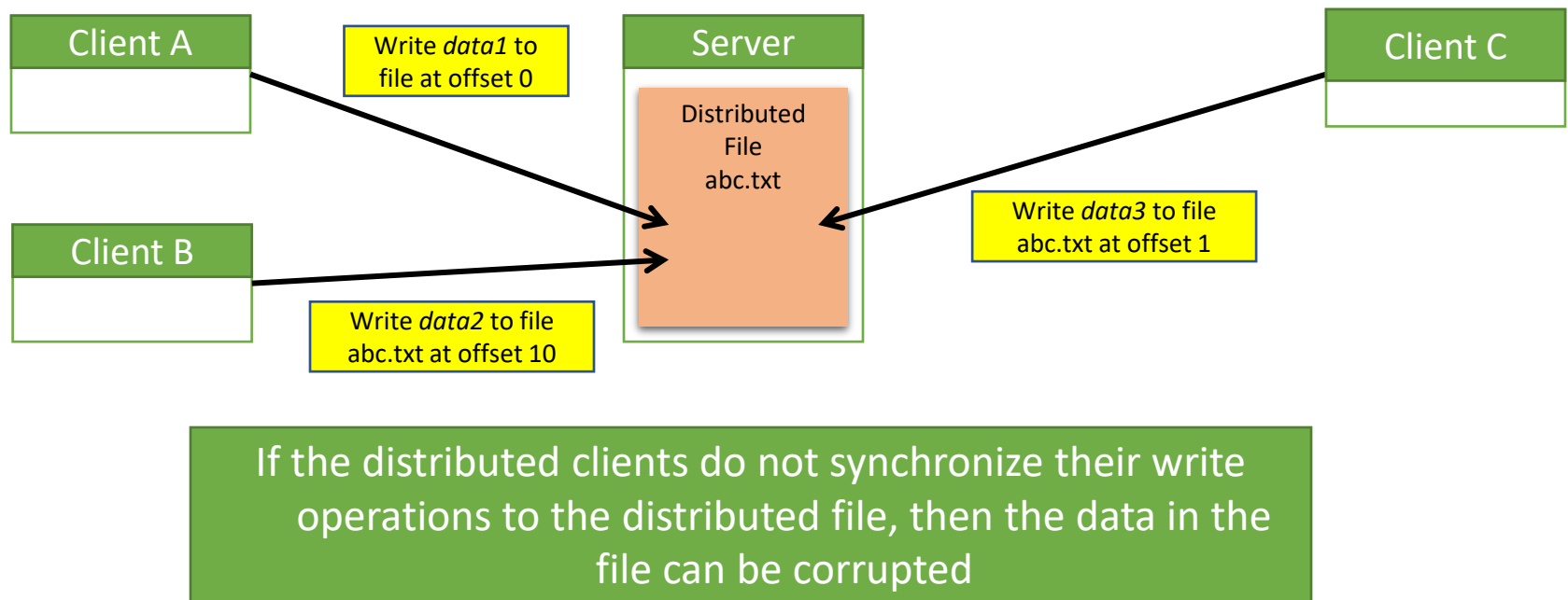
- Vehicle tracking in a City Surveillance System using a Distributed Sensor Network of Cameras
 - *Objective:* To keep track of suspicious vehicles
 - Camera Sensor Nodes are deployed over the city
 - Each Camera Sensor that detects the vehicle reports the time to a central server
 - Server tracks the movement of the suspicious vehicle

If the sensor nodes do not have a consistent version of the time, the vehicle cannot be reliably tracked



Need for Synchronization – Example 2

- Writing a file in a Distributed File System



A Broad Taxonomy of Synchronization

Reason for synchronization and cooperation	Entities have to agree on ordering of events	Entities have to share common resources
Examples	e.g., Vehicle tracking in a Camera Sensor Network, Financial transactions in Distributed eCommerce Systems	e.g., Reading and writing in a Distributed File System
Requirement for entities	Entities should have a common understanding of time across different computers	Entities should coordinate and agree on when and how to access resources
Topics we will study	Time Synchronization	Mutual Exclusion

Overview

- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
- Election Algorithms

Clock Synchronization

- Clock Synchronization is a mechanism to synchronize the time of all the computers in a DS
- We will study
 - Coordinated Universal Time
 - Tracking Time on a Computer
 - Clock Synchronization Algorithms
 - Cristian's Algorithm
 - Berkeley Algorithm
 - Network Time Protocol

Clock Synchronization

- Coordinated Universal Time
- Tracking Time on a Computer
- Clock Synchronization Algorithms
 - Cristian's Algorithm
 - Network Time Protocol
 - Berkeley Algorithm

Coordinated Universal Time (UTC)

- All the computers are generally synchronized to a standard time called Coordinated Universal Time (UTC)
 - UTC is the primary time standard by which the world regulates clocks and time
- UTC is broadcasted via the satellites
 - UTC broadcasting service provides an accuracy of 0.5 msec
- Computer servers and online services with **UTC receivers** can be synchronized by satellite broadcasts
 - Many popular synchronization protocols in distributed systems use UTC as a reference time to synchronize clocks of computers

Clock Synchronization

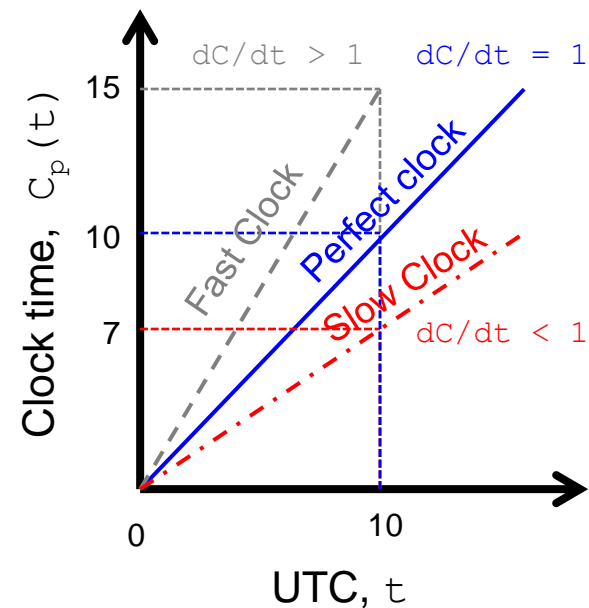
- Coordinated Universal Time
- Tracking Time on a Computer
- Clock Synchronization Algorithms
 - Cristian's Algorithm
 - Berkeley Algorithm
 - Network Time Protocol

Tracking Time on a Computer

- How does a computer keep track of its time?
 - Each computer has a hardware *timer*
 - The timer causes an interrupt 'H' times a second
 - The interrupt handler adds 1 to its Software Clock (C)
- Issues with clocks on a computer
 - In practice, the hardware timer is imprecise
 - It does not interrupt 'H' times a second due to material imperfections of the hardware and temperature variations
 - The computer counts the time slower or faster than actual time
 - Loosely speaking, **Clock Skew** is the skew between:
 - the computer clock and the actual time (e.g., UTC)

Clock Skew

- When the UTC time is t , let the clock on the computer have a time $C(t)$
- Three types of clocks are possible
 - Perfect clock:
 - The timer ticks 'H' interrupts a second
 $dC/dt = 1$
 - Fast clock:
 - The timer ticks more than 'H' interrupts a second
 $dC/dt > 1$
 - Slow clock:
 - The timer ticks less than 'H' interrupts a second
 $dC/dt < 1$



Clock Skew (cont'd)

- **Frequency** of the clock is defined as the ratio of the number of seconds counted by the software clock for every UTC second

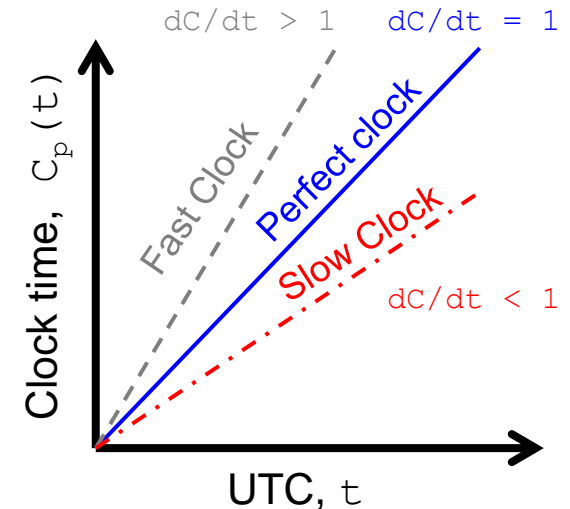
$$\text{Frequency} = dC/dt$$

- **Skew** of the clock is defined as the extent to which the frequency differs from that of a perfect clock

$$\text{Skew} = dC/dt - 1$$

- Hence,

$$\text{Skew} \begin{cases} > 0 & \text{for a fast clock} \\ = 0 & \text{for a perfect clock} \\ < 1 & \text{for a slow clock} \end{cases}$$



Maximum Drift Rate of a Clock

- The manufacturer of the timer specifies the upper and the lower bound that the clock skew may fluctuate. This value is known as *maximum drift rate* (ρ)

$$1 - \rho \leq dC/dt \leq 1 + \rho$$

- How far can two clocks drift apart?
 - If two clocks were synchronized Δt seconds before to UTC, then the two clocks can be as much as $2\rho\Delta t$ seconds apart
- Guaranteeing maximum drift between computers in a DS
 - If maximum drift permissible in a DS is δ seconds, then clocks of every computer has to resynchronize at least $\delta/2\rho$ seconds

Clock Synchronization

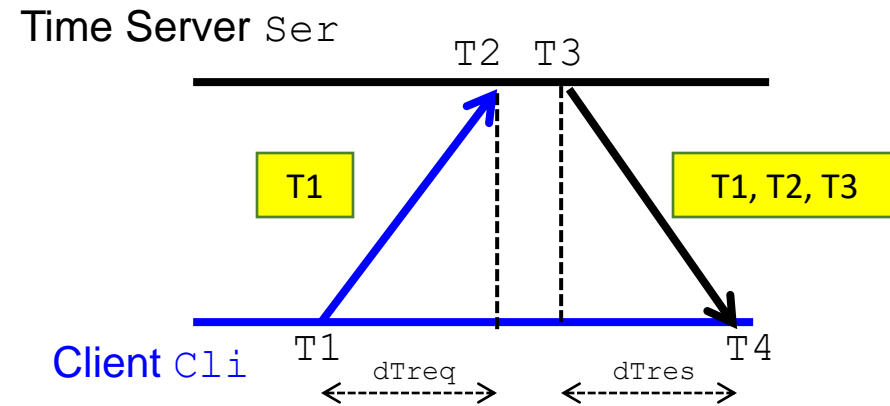
- Coordinated Universal Time
- Tracking Time on a Computer
- Clock Synchronization Algorithms
 - Cristian's Algorithm
 - Berkeley Algorithm
 - Network Time Protocol

Cristian's Algorithm

- Flaviu Cristian (in 1989) provided an algorithm to synchronize networked computers with a time server
- The basic idea:
 - Identify a network time server that has an accurate source for time (e.g., the time server has a UTC receiver)
 - All the clients contact the network time server for synchronization
- However, the network delays incurred while the client contacts the time server will have outdated the reported time
 - The algorithm estimates the network delays and compensates for it

Cristian's Algorithm – Approach

- + Client *Cli* sends a request to Time Server *Ser*, time stamped its local clock time *T1*
- + *S* will record the time of receipt *T2* according to its local clock
 - + *dTreq* is network delay for request transmission



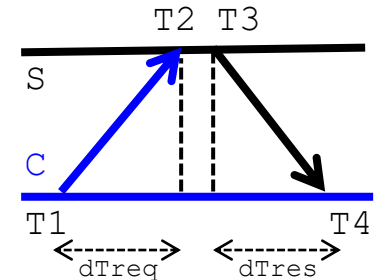
- *Ser* replies to *Cli* at its local time *T3*, piggybacking *T1* and *T2*
- *Cli* receives the reply at its local time *T4*
 - *dTres* is the network delay for response transmission
- Now *Cli* has the information *T1*, *T2*, *T3* and *T4*
- **Assuming that the transmission delay from *Cli* → *Ser* and *Ser* → *Cli* are the same**

$$T2 - T1 \approx T4 - T3$$

Christian's Algorithm – Synchronizing Client Time

- + Client C estimates its offset θ relative to Time Server S

$$\begin{aligned}\theta &= T3 + dT_{res} - T4 \\ &= T3 + ((T2 - T1) + (T4 - T3)) / 2 - T4 \\ &= ((T2 - T1) + (T3 - T4)) / 2\end{aligned}$$



- + If $\theta > 0$ or $\theta < 0$, then the client time should be incremented or decremented by θ seconds

Gradual Time Synchronization at the client

- Instead of changing the time drastically by θ seconds, typically the time is gradually synchronized
 - The software clock is updated at a lesser/greater rate whenever timer interrupts

Note: Setting clock backward (say, if $\theta < 0$) is not allowed in a DS since decrementing a clock at any computer has adverse effects on several applications (e.g., *make* program)

Cristian's Algorithm – Discussion

1. Assumption about packet transmission delays

- Cristian's algorithm assumes that the round-trip times for messages exchanged over the network is reasonably short
- The algorithm assumes that the delay for the request and response are equal

- Will the trend of increasing Internet traffic decrease the accuracy of the algorithm?
- Can the algorithm handle delay asymmetry that is prevalent in the Internet?
- Can the clients be mobile entities with intermittent connectivity?

Cristian's algorithm is intended for synchronizing computers within intranets

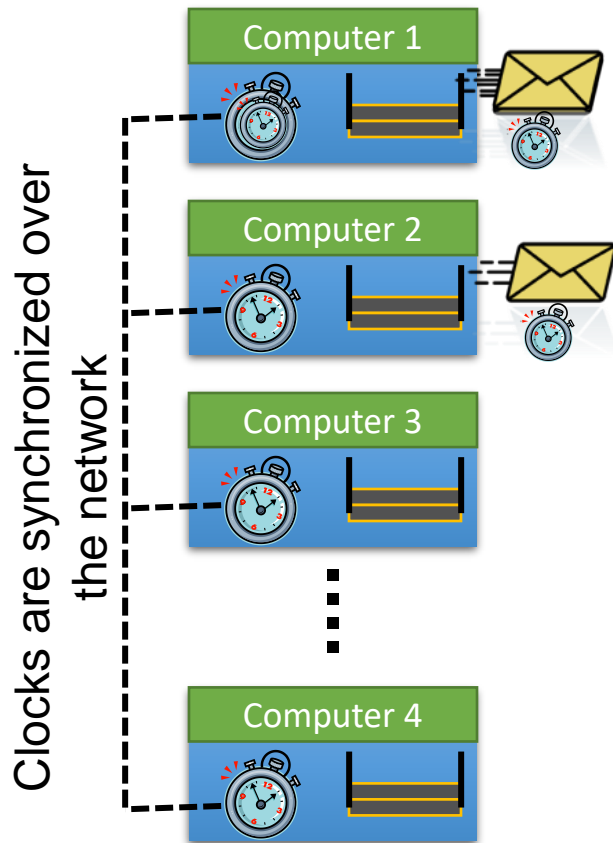
2. A probabilistic approach for calculating delays

- There is no tight bound on the maximum drift between clocks of computers

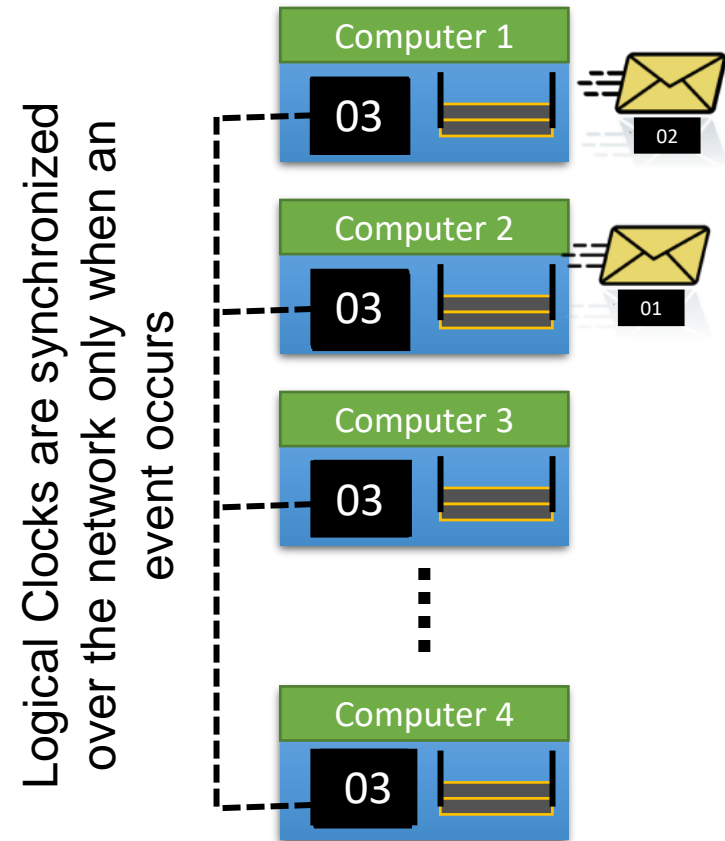
3. Time server failure or faulty server clock

- Faulty clock on the time server leads to inaccurate clocks in the entire DS
- Failure of the time server will render synchronization impossible

Types of Time Synchronization



Clock-based Time Synchronization



Event-based Time Synchronization

Overview

- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
- Election Algorithms

Clock Synchronization

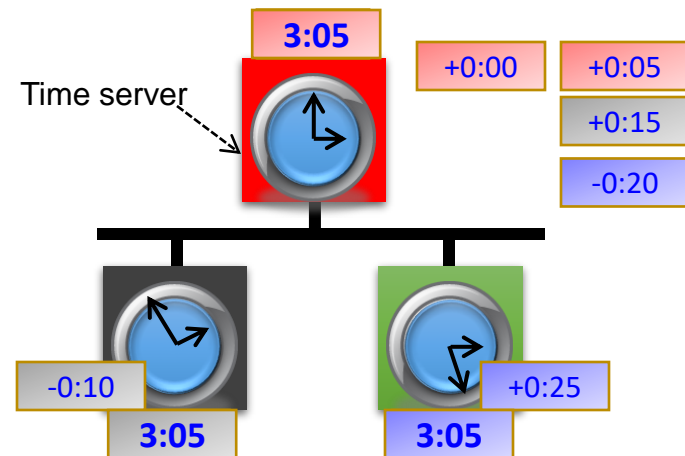
- Coordinated Universal Time
- Tracking Time on a Computer
- Clock Synchronization Algorithms
 - Cristian's Algorithm
 - Berkeley Algorithm
 - Network Time Protocol

Berkeley Algorithm

+ Berkeley Algorithm is a distributed approach for time synchronization

- Approach:

1. A time server periodically (approx. once in 4 minutes) sends its time to all the computers and polls them for the time difference
2. The computers compute the time difference and then reply
3. The server computes an average time difference for each computer
4. The server commands all the computers to update their time (by gradual time synchronization)



Berkeley Algorithm – Discussion

1. Assumption about packet transmission delays

- Berkeley's algorithm predicts network delay (similar to Cristian's algorithm)
- Hence, it is effective in intranets, and not accurate in wide-area networks

2. No UTC Receiver is necessary

- The clocks in the system synchronize by averaging all the computer's times

3. Decreases the effect of faulty clocks

- Fault-tolerant averaging, where outlier clocks are ignored, can be easily performed in Berkeley Algorithm

4. Time server failures can be masked

- If a time server fails, another computer can be elected as a time server

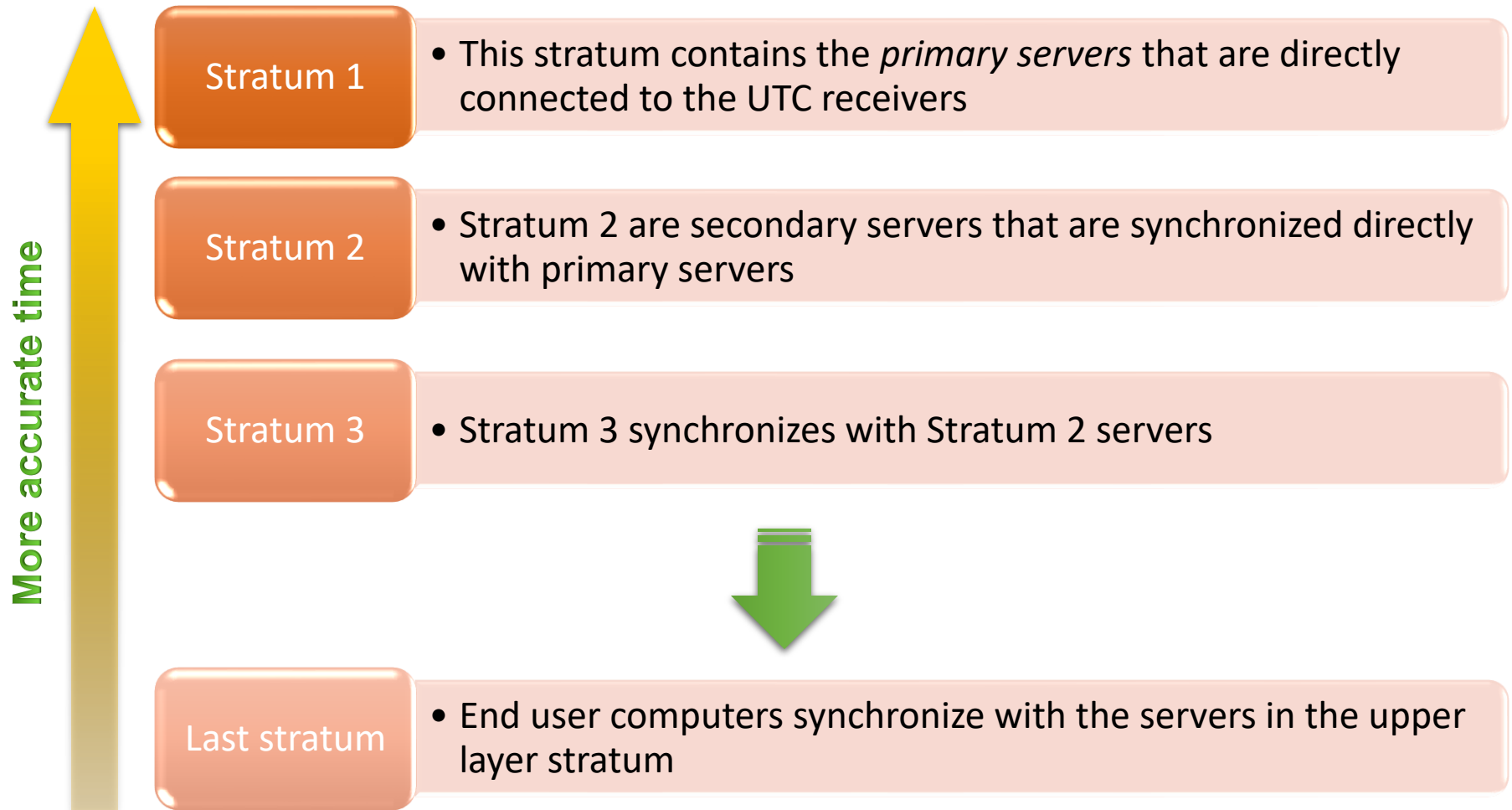
Clock Synchronization

- Coordinated Universal Time
- Tracking Time on a Computer
- Clock Synchronization Algorithms
 - Cristian's Algorithm
 - Berkeley Algorithm
 - Network Time Protocol

Network Time Protocol (NTP)

- NTP defines an architecture for a time service and a protocol to distribute time information over the Internet
- In NTP, servers are connected in a logical hierarchy called *synchronization subnet*
- The levels of synchronization subnet is called *strata*
 - Stratum 1 servers have most accurate time information (connected to a UTC receiver)
 - Servers in each stratum act as time servers to the servers in the lower stratum

Hierarchical organization of NTP Servers



Operation of NTP Protocol

- When a time server **A** contacts time server **B** for synchronization
 - If **stratum(A) ≤ stratum(B)**, then **A** does not synchronize with **B**
 - If **stratum(A) > stratum(B)**, then:
 - Time server **A** synchronizes with **B**
 - An algorithm similar to Cristian's algorithm is used to synchronize. However, larger statistical samples are taken before updating the clock
 - Time server **A** updates its stratum
$$\text{stratum(A)} = \text{stratum(B)} + 1$$

Discussion of NTP Design

Accurate synchronization to UTC time

- NTP enables clients across the Internet to be synchronized accurately to the UTC
- Large and variable message delays are tolerated through statistical filtering of timing data from different servers

Scalability

- NTP servers are hierarchically organized to speed up synchronization, and to scale to a large number of clients and servers

Reliability and Fault-tolerance

- There are redundant time servers, and redundant paths between the time servers
- The architecture provides reliable service that can tolerate lengthy losses of connectivity
- A synchronization subnet can reconfigure as servers become unreachable. For example, if Stratum 1 server fails, then it can become a Stratum 2 secondary server

Security

- NTP protocol uses authentication to check of the timing message originated from the claimed trusted sources

Summary of Clock Synchronization

- Physical clocks on computers are not accurate
- Clock synchronization algorithms provide mechanisms to synchronize clocks on networked computers in a DS
 - Computers on a local network use various algorithms for synchronization
 - Some algorithms (e.g, Cristian's algorithm) synchronize time with by contacting centralized time servers
 - Some algorithms (e.g., Berkeley algorithm) synchronize in a distributed manner by exchanging the time information on various computers
 - NTP provides architecture and protocol for time synchronization over wide-area networks such as Internet

Overview

- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
- Election Algorithms

Why Logical Clocks?

- Lamport (in 1978) showed that:
 - Clock synchronization is not necessary in all scenarios
 - If two processes do not interact, it is not necessary that their clocks are synchronized
 - Many times, it is sufficient if processes agree on the order in which the events has occurred in a DS
 - For example, for a distributed *make* utility, it is sufficient to know if an input file was modified *before* or *after* its object file

Logical Clocks

- Logical clocks are used to define an order of events without measuring the physical time at which the events occurred
- We will study two types of logical clocks
 1. Lamport's Logical Clock (or simply, Lamport's Clock)
 2. Vector Clock

Logical Clocks

- We will study two types of logical clocks
 1. Lamport's Clock
 2. Vector Clock

Lamport's Logical Clock

- Lamport advocated maintaining logical clocks at the processes to keep track of the order of events
- To synchronize logical clocks, Lamport defined a relation called “**happened-before**”
- The expression $\mathbf{a} \rightarrow \mathbf{b}$ (read as “**a** happened before **b**”) means that all entities in a DS agree that event **a** occurred before event **b**

Happened-before Relation

- The **happened-before** relation can be observed directly in two situations:
 1. If **a** and **b** are events in the same process, and **a** occurs before **b**, then **$a \rightarrow b$** is true
 2. If **a** is an event of message **m** being sent by a process, and **b** is the event of the message **m** being received by another process, the **$a \rightarrow b$** is true.
- The **happened-before** relation is transitive
 - If **$a \rightarrow b$** and **$b \rightarrow c$** , then **$a \rightarrow c$**

Time values in Logical Clocks

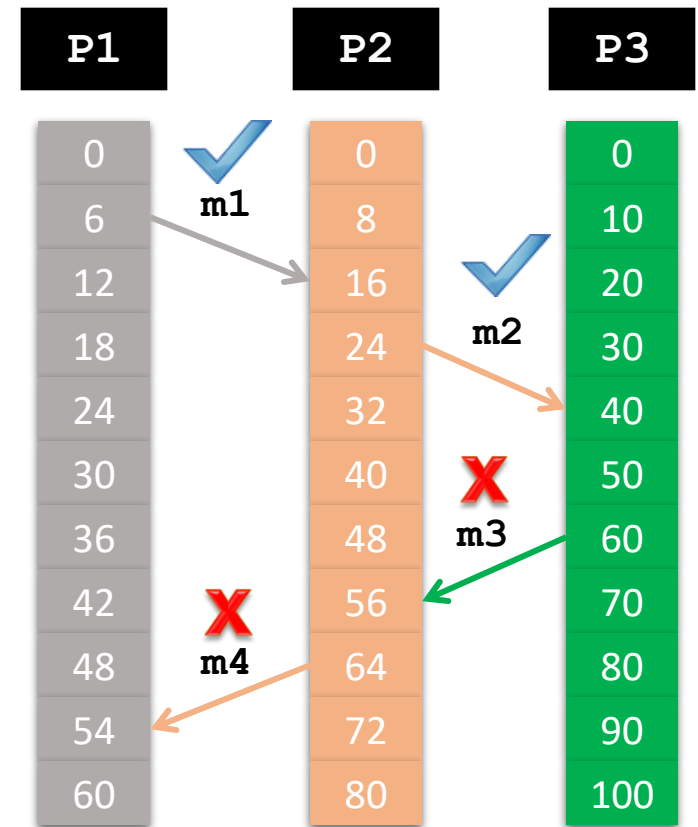
- For every event **a**, assign a logical *time value* **C (a)** on which all processes agree
- Time value for events have the property that
 - If **a** \rightarrow **b**, then **C (a) < C (b)**

Properties of Logical Clock

- From **happened-before** relation, we can infer that:
 - If two events **a** and **b** occur within the same process and $\mathbf{a \rightarrow b}$, then assign $\mathbf{C(a)}$ and $\mathbf{C(b)}$ such that $\mathbf{C(a) < C(b)}$
 - If **a** is the event of sending the message **m** from one process, and **b** is the event of receiving the message **m**, then
 - the time values $\mathbf{C(a)}$ and $\mathbf{C(b)}$ are assigned such that all processes agree that $\mathbf{C(a) < C(b)}$
 - The clock time **C** must always go forward (increasing), and never backward (decreasing)

Synchronizing Logical Clocks

- Three processes **P1**, **P2** and **P3** running at different rates
- If the processes communicate between each other, there might be discrepancies in agreeing on the event ordering
 - Ordering of sending and receiving messages **m1** and **m2** are correct
 - However, **m3** and **m4** violate the happens-before relationship



Lamport's Clock Algorithm

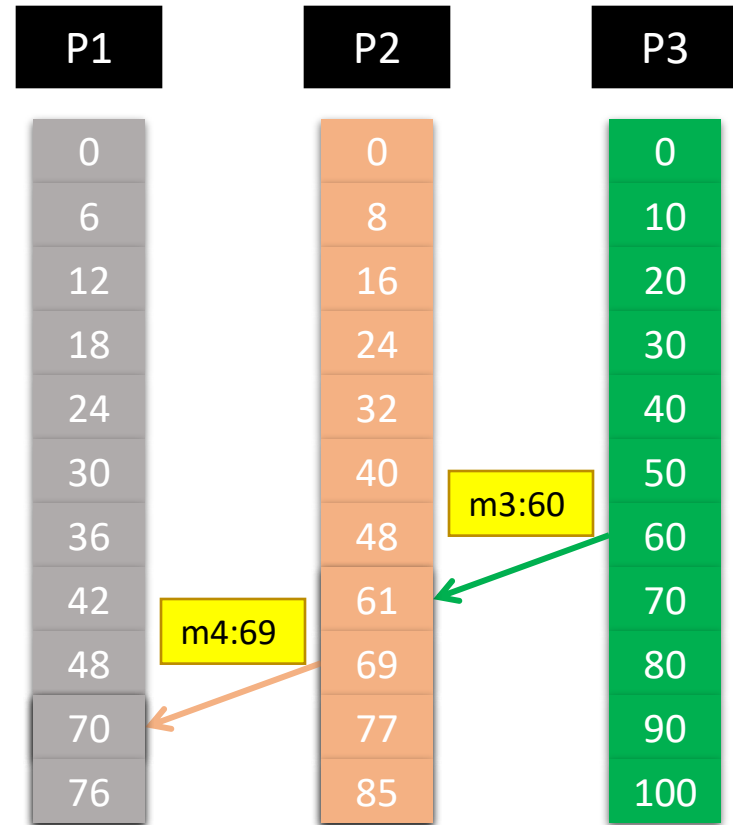
When a message is being sent:

- Each message carries a **timestamp** according to the sender's logical clock

When a message is received:

- If the receiver logical clock is less than message sending time in the packet, then adjust the receiver's clock such that

$$\text{currentTime} = \text{timestamp} + 1$$

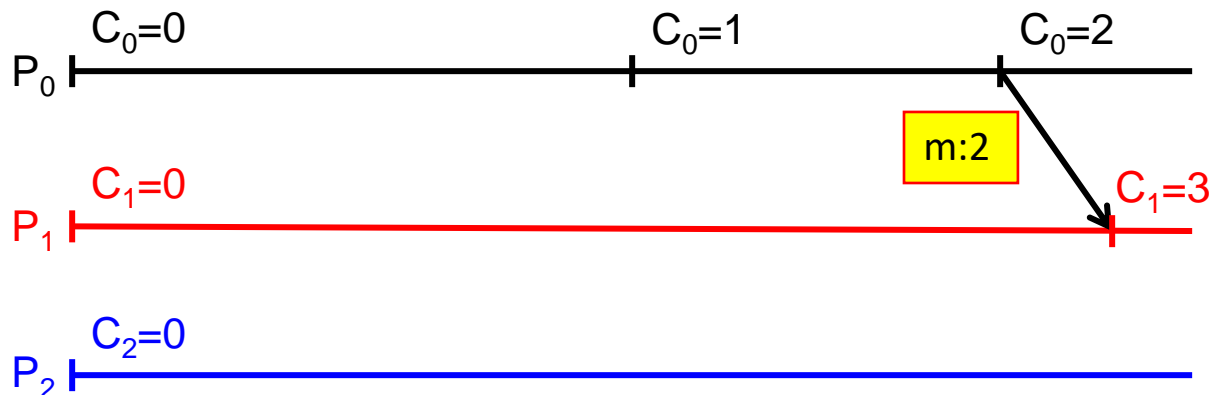


Logical Clock Without a Physical Clock

- Previous examples assumed that there is a physical clock at each computer (probably running at different rates)
- How to attach a time value to an event when there is no global clock?

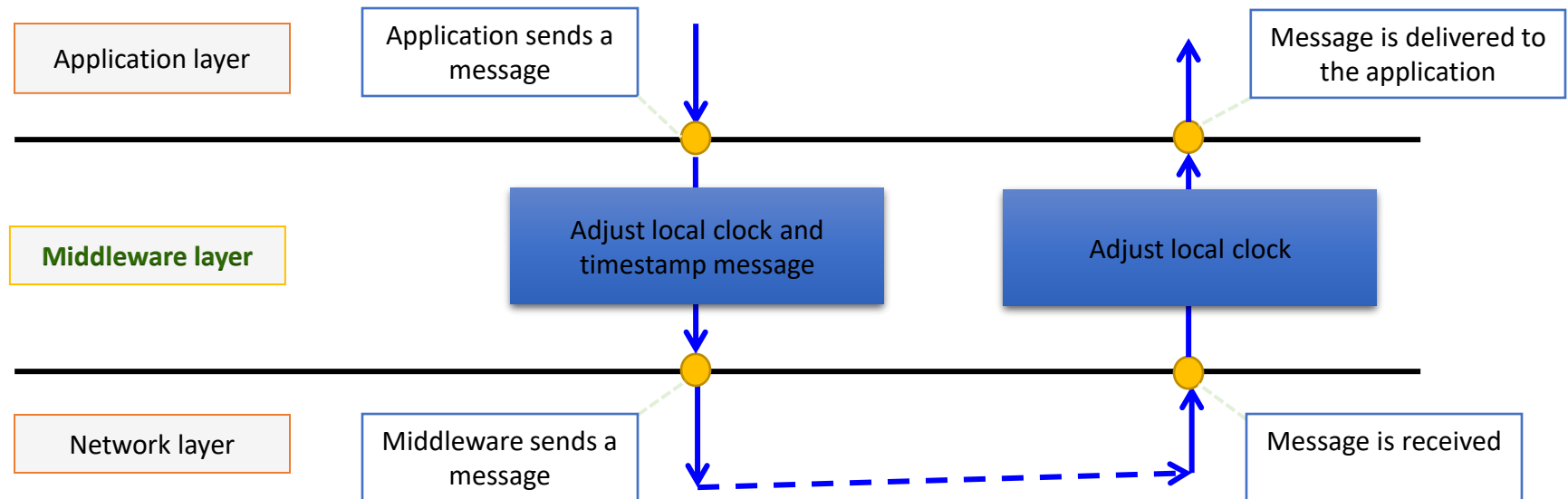
Implementation of Lamport's Clock

- Each process P_i maintains a local counter C_i and adjusts this counter according to the following rules:
 - For any two successive events that take place within P_i , C_i is incremented by 1
 - Each time a message m is sent by process P_i , the message receives a timestamp $ts(m) = C_i$
 - Whenever a message m is received by a process P_j , P_j adjusts its local counter C_j to $\max(C_j, ts(m)) + 1$



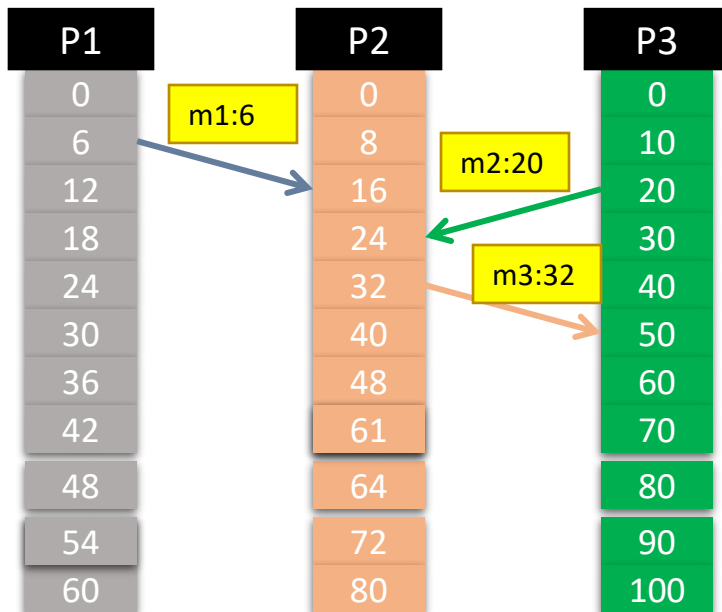
Placement of Logical Clock

- In a computer, several processes use Logical Clocks
 - Similar to how several processes on a computer use one physical clock
- Instead of each process maintaining its own Logical Clock, Logical Clocks can be implemented as a middleware for time service



Limitation of Lamport's Clock

- Lamport's Clock ensures that if $a \rightarrow b$, then $C(a) < C(b)$
- However, it does not say anything about any two events a and b by comparing their time values
 - For any two events a and b , $C(a) < C(b)$ does not mean that $a \rightarrow b$
- Example:



Compare m1 and m3

P2 can infer that $m1 \rightarrow m3$

Compare m1 and m2

P2 **cannot** infer that $m1 \rightarrow m2$ or $m2 \rightarrow m1$

Summary of Lamport's Clock

- Lamport advocated using logical clocks
 - Processes synchronize based on their time values of the logical clock rather than the absolute time on the physical time
- Which applications in DS need logical clocks?
 - Applications with provable ordering of events
 - Perfect physical clock synchronization is hard to achieve in practice. Hence we cannot provably order the events
 - Applications with rare events
 - Events are rarely generated, and physical clock synchronization overhead is not justified
- However, Lamport's clock cannot guarantee perfect ordering of events by just observing the time values of two arbitrary events

Logical Clocks

- We will study two types of logical clocks
 1. Lamport's Clock
 2. Vector Clocks

Vector Clocks

- Vector Clocks was proposed to overcome the limitation of Lamport's clock: the fact that $C(a) < C(b)$ does not mean that $a \rightarrow b$
 - The property of inferring that a occurred before b is called as **causality** property
- A Vector clock for a system of N processes is an array of N integers
- Every process P_i stores its own vector clock VC_i
 - Lamport's time value for events are stored in VC_i
 - $VC_i(a)$ is assigned to an event a
- If $VC_i(a) < VC_i(b)$, then we can infer that $a \rightarrow b$

Updating Vector Clocks

- Vector clocks are constructed by the following two properties:

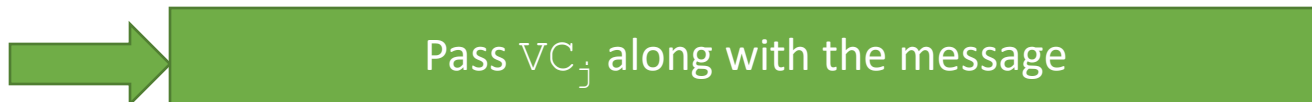
1. $VC_i[i]$ is the number of events that have occurred at process P_i so far

- $VC_i[i]$ is the local logical clock at process P_i



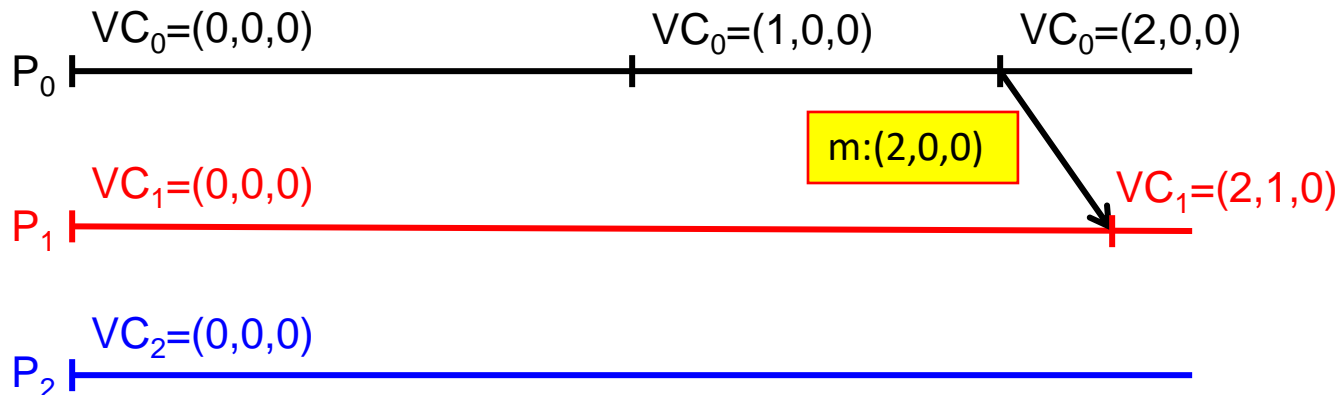
2. If $VC_i[j] = k$, then P_i knows that k events have occurred at P_j

- $VC_i[j]$ is P_i 's knowledge of local time at P_j



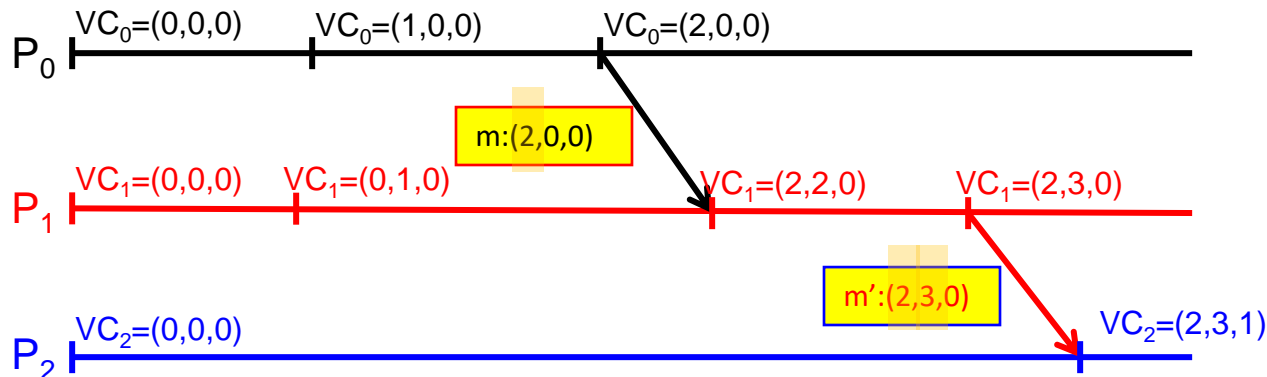
Vector Clock Update Algorithm

- Whenever there is a new event at P_i , increment $VC_i[i]$
- When a process P_i sends a message m to P_j :
 - Increment $VC_i[i]$
 - Set m 's timestamp $ts(m)$ to the vector VC_i
- When message m is received process P_j :
 - $VC_j[k] = \max(VC_j[k], ts(m)[k])$; (for all k)
 - Increment $VC_j[j]$



Inferring Events with Vector Clocks

- Let a process P_i send a message m to P_j with timestamp $ts(m)$, then:
 - P_j knows the number of events at the sender P_i that causally precede m
 - $(ts(m)[i] - 1)$ denotes the number of events at P_i
 - P_j also knows the minimum number of events at other processes P_k that causally precede m
 - $(ts(m)[k] - 1)$ denotes the minimum number of events at P_k

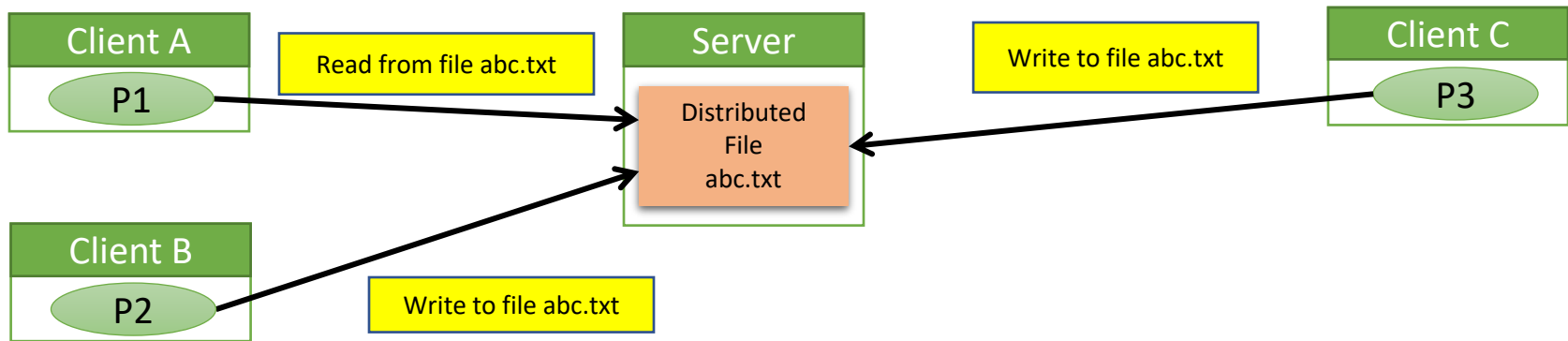


Overview

- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
- Election Algorithms

Need for Mutual Exclusion

- Distributed processes need to coordinate to access shared resources
- Example: Writing a file in a Distributed File System



In uniprocessor systems, mutual exclusion to a shared resource is provided through shared variables or operating system support.

However, such support is insufficient to enable mutual exclusion of distributed entities

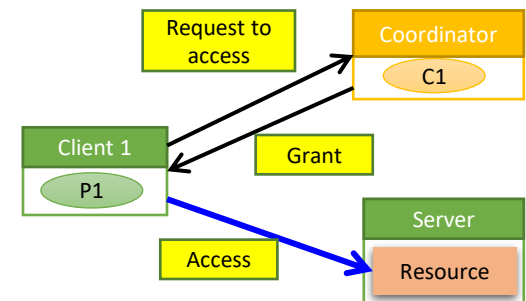
In Distributed System, processes coordinate access to a shared resource by passing messages to enforce *distributed mutual exclusion*

Types of Distributed Mutual Exclusion

- Mutual exclusion algorithms are classified into two categories

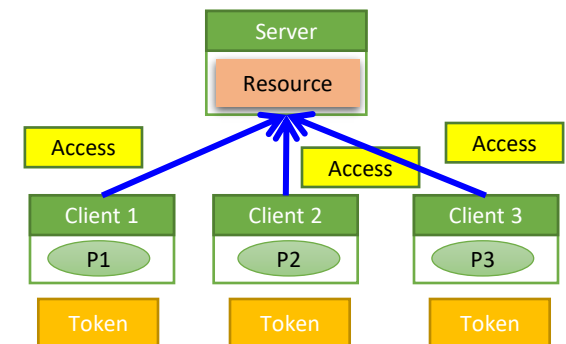
1. Permission-based Approaches

- + A process, which wants to access a shared resource, requests the permission from one or more coordinators



2. Token-based Approaches

- + Each shared resource has a token
- + Token is circulated among all the processes
- + A process can access the resource if it has the token



Overview

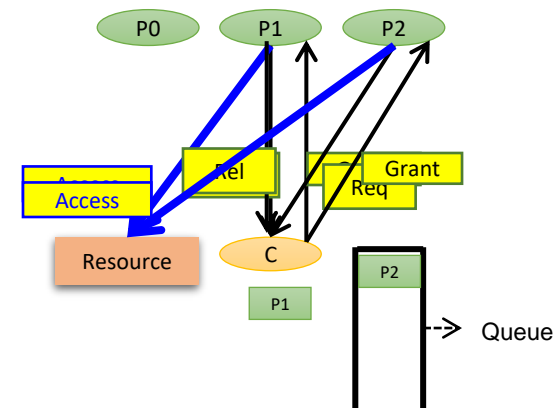
- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
 - Permission-based Approaches
 - Token-based Approaches
- Election Algorithms

Permission-based Approaches

- There are two types of permission-based mutual exclusion algorithms
 - a. Centralized Algorithms
 - b. Decentralized Algorithms
- We will study an example of each type of algorithm

a. A Centralized Algorithm

- One process is elected as a coordinator (**C**) for a shared resource
- Coordinator maintains a **Queue** of access requests
- Whenever a process wants to access the resource, it sends a request message to the coordinator to access the resource
- When the coordinator receives the request:
 - If no other process is currently accessing the resource, it grants the permission to the process by sending a “grant” message
 - If another process is accessing the resource, the coordinator queues the request, and does not reply to the request
- The process releases the exclusive access after accessing the resource
- The coordinator will then send the “grant” message to the next process in the queue



Discussion about Centralized Algorithm

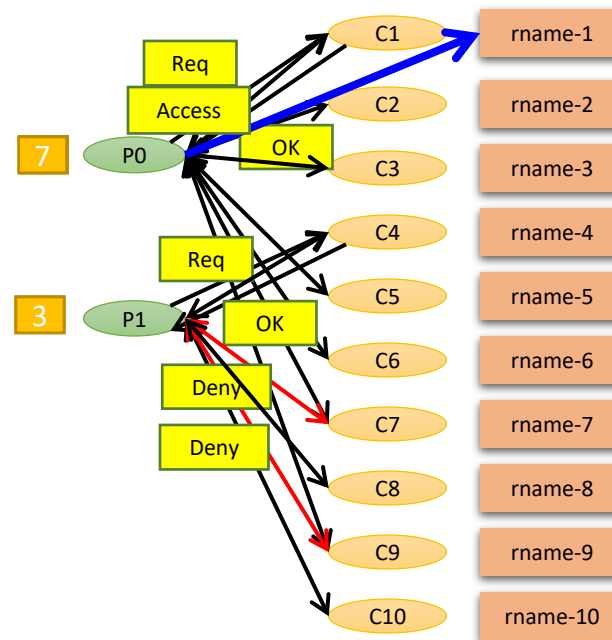
- Blocking vs. non-blocking requests
 - The coordinator can block the requesting process until the resource is free
 - Otherwise, the coordinator can send a “permission-denied” message back to the process
 - The process can poll the coordinator at a later time, or
 - The coordinator queues the request. Once the resource is released, the coordinator will send an explicit “grant” message to the process
- The algorithm guarantees mutual exclusion, and is simple to implement
- Fault-tolerance:
 - Centralized algorithm is vulnerable to a single-point of failure (at coordinator)
 - Processes cannot distinguish between dead coordinator and request blocking
- Performance bottle-neck:
 - In a large system, single coordinator can be overwhelmed with requests

b. A Decentralized Algorithm

- To avoid the drawbacks of the centralized algorithm, Lin *et al.* [1] advocated a decentralized mutual exclusion algorithm
- Assumptions
 - Distributed processes are in a Distributed Hash Table (DHT) based system
 - Each resource is replicated **n** times
 - The **i^{th}** replica of a resource **$rname$** is named as **$rname-i$**
 - Every replica has its own coordinator for controlling access
 - The coordinator for **$rname-i$** is determined by using a hash function
- Approach:
 - Whenever a process wants to access the resource, it will have to get a majority vote from **$m > n/2$** coordinators
 - If a coordinator does not want to vote for a process (because it has already voted for another process), it will send a “permission-denied” message to the process

A Decentralized Algorithm – An Example

- If $n=10$ and $m=7$, then a process needs at-least 7 votes to access the resource



rname-x = xth replica of a resource
Cj = Coordinator j Pi = Process i n = Number of votes gained

Fault-tolerance in Decentralized Algorithm

- The decentralized algorithm assumes that the coordinator recovers quickly from a failure
- However, the coordinator would have reset its state after recovery
 - Coordinator could have forgotten any vote it had given earlier
- Hence, the coordinator may incorrectly grant permission to the processes
 - Mutual exclusion cannot be deterministically guaranteed
 - But, the algorithm *probabilistically* guarantees mutual exclusion

Probabilistic Guarantees in the Decentralized Algorithm

- What is the minimum number of coordinators who should fail for violating mutual exclusion?
 - At least $2m-n$ coordinators should fail
- Let the probability of violating mutual exclusion be P_v
- Derivation of P_v
 - Let T be the lifetime of the coordinator
 - Let $p=\Delta t/T$ be the probability that coordinator crashes during time-interval Δt
 - Let $P[k]$ be the probability that k out of m coordinators crash during the same interval
 - We compute the mutual exclusion violation probability P_v by:

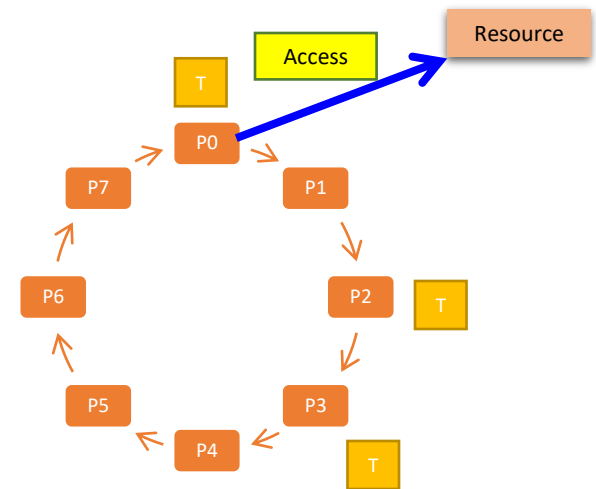
$$P_v = \sum_{k=2m-n}^m P[k] p^k (1-p)^{m-k}$$
- In practice, this probability should be very small
 - For $T=3$ hours, $\Delta t=10$ s, $n=32$, and $m=0.75n$: $P_v=10^{-40}$

Overview

- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
 - Permission-based Approaches
 - Token-based Approaches
- Election Algorithms

Token Ring

- In the Token Ring algorithm, each resource is associated with a *token*
- The token is circulated among the processes
- The process with the token can access the resource
- Circulating the token among processes:
 - + All processes form a logical ring where each process knows its next process
 - + One process is given a *token* to access the resource
 - + The process with the token has the right to access the resource
 - + If the process has finished accessing the resource OR does not want to access the resource:
 - + it passes the token to the next process in the ring



Discussion about Token Ring

- ✓ Token ring approach provides deterministic mutual exclusion
 - There is one token, and the resource cannot be accessed without a token
- ✓ Token ring approach avoids starvation
 - Each process will receive the token
- ✗ Token ring has a high-message overhead
 - When no processes need the resource, the token circulates at a high-speed
- ✗ If the token is lost, it must be regenerated
 - Detecting the loss of token is difficult since the amount of time between successive appearances of the token is unbounded
- ✗ Dead processes must be purged from the ring
 - ACK based token delivery can assist in purging dead processes

Comparison of Mutual Exclusion Algorithms

Algorithm	Delay before a process can access the resource (in message times)	Number of messages required for a process to access and release the shared resource	Problems
Centralized	2	3	<ul style="list-style-type: none"> Coordinator crashes
Decentralized	$2mk$	$2mk + m; k=1,2,\dots$	<ul style="list-style-type: none"> Large number of messages
Token Ring	0 to $(n-1)$	1 to ∞	<ul style="list-style-type: none"> Token may be lost Ring can cease to exist since processes crash

- Assume that:

n = Number of processes in the distributed system

For the Decentralized algorithm:

m = minimum number of coordinators who have to agree for a process to access a resource

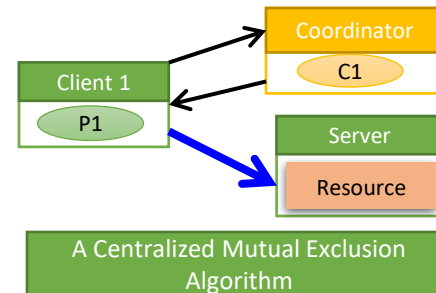
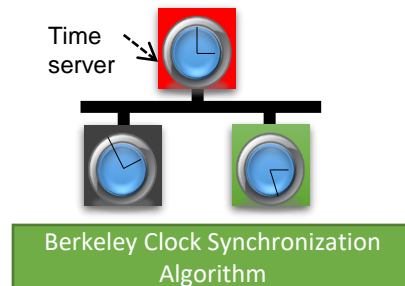
k = average number of requests made by the process to a coordinator to request for a vote

Overview

- Time Synchronization
 - Clock Synchronization
 - Logical Clock Synchronization
- Mutual Exclusion
 - Permission-based Approaches
 - Token-based Approaches
- Election Algorithms

Election in Distributed Systems

- Many distributed algorithms require one process to act as a coordinator
 - Typically, it does not matter which process is elected as the coordinator
- Example algorithms where coordinator election is required



Election Process

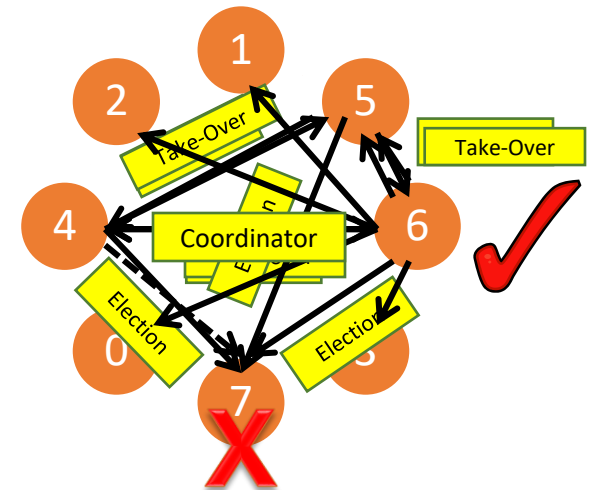
- Any process P_i in the DS can initiate the election algorithm that elects a new coordinator
- At the termination of the election algorithm, the elected coordinator process should be unique
- Every process *may* know the process ID of every other processes, but it does not know which processes have crashed
- Generally, we require that the coordinator is the process with the largest process ID
 - The idea can be extended to elect *best* coordinator
 - Example: Election of a coordinator with least computational load
 - If the computational load of process P_i denoted by \mathbf{load}_i , then coordinator is the process with highest $1/\mathbf{load}_i$. Ties are broken by sorting process ID.

Election Algorithms

- We will study two election algorithms
 1. Bully Algorithm
 2. Ring Algorithm

1. Bully Algorithm

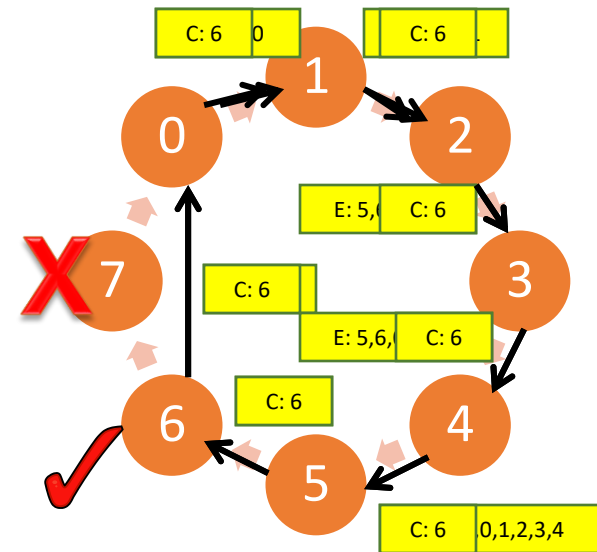
- A process initiates election algorithm when it notices that the existing coordinator is not responding
- Process P_i calls for an election as follows:
 1. P_i sends an “Election” message to all processes with higher process IDs
 2. When process P_j with $j > i$ receives the message, it responds with a “Take-over” message. P_j no more contests in the election
 - i. Process P_j re-initiates another call for election. Steps 1 and 2 continue
 3. If no one responds, P_i wins the election. P_i sends “Coordinator” message to every process



2. Ring Algorithm

- This algorithm is generally used in a ring topology
- When a process P_i detects that the coordinator has crashed, it initiates an election algorithm

1. P_i builds an “Election” message (E), and sends it to its next node. It inserts its ID into the Election message
2. When process P_j receives the message, it appends its ID and forwards the message
 - i. If the next node has crashed, P_j finds the next alive node
3. When the message gets back to the process that started the election:
 - i. it elects process with highest ID as coordinator, and
 - ii. changes the message type to “Coordination” message (C) and circulates it in the ring



Comparison of Election Algorithms

Algorithm	Number of Messages for Electing a Coordinator	Problems
Bully Algorithm	$O(n^2)$	<ul style="list-style-type: none">• Large message overhead
Ring Algorithm	$2n$	<ul style="list-style-type: none">• An overlay ring topology is necessary

- Assume that:
n = Number of processes in the distributed system

Summary of Election Algorithms

- Election algorithms are used for choosing a unique process that will coordinate an activity
- At the end of the election algorithm, all nodes should uniquely identify the coordinator
- We studied two algorithms for election
 - Bully algorithm
 - Processes communicate in a distributed manner to elect a coordinator
 - Ring algorithm
 - Processes in a ring topology circulate election messages to choose a coordinator