

# Advanced Topics in Operating Systems

MSc in Computer Science  
UNYT-UoG

Dr. Marenglen Biba  
18-19-20 December 2009



# Lesson 6

01: Introduction

02: Architectures

03: Processes

04: Communication

05: Naming

**06: Synchronization**

07: Consistency & Replication

08: Fault Tolerance

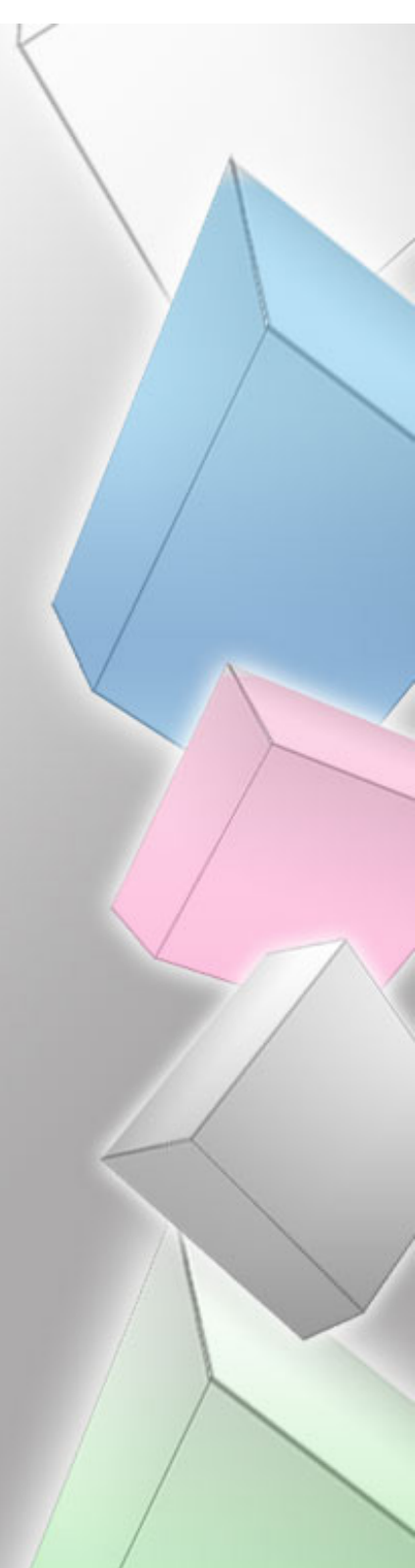
09: Security

10: Distributed Object-Based Systems

11: Distributed File Systems

12: Distributed Web-Based Systems

13: Distributed Coordination-Based Systems

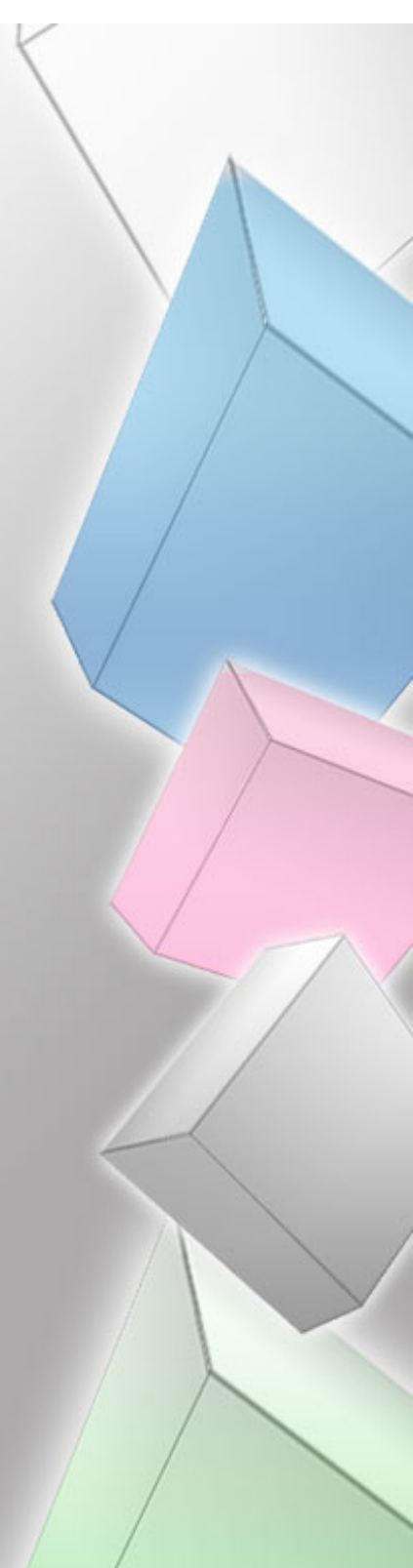


# Synchronization

- It is important that multiple processes **do not simultaneously access** a shared resource, such as printer, but instead cooperate in granting each other temporary exclusive access.
- Moreover multiple processes may sometimes need to agree on the ordering of events, such as whether **message  $m1$**  from process  $P$  was sent before or after message  $m2$  from process  $Q$ .
- While communication is important, it is not everything.
  - **Closely related is how processes cooperate and synchronize with one another.**
- Synchronization in distributed systems is **often much more difficult** compared to synchronization in uniprocessor or multiprocessor systems.

# CLOCK SYNCHRONIZATION

- Physical clocks
- Logical clocks
- Vector clocks



# Clock problem

- In a centralized system, time is **unambiguous**.
  - When a process wants to know the time, it makes a system call and the kernel tells it.
- If process *A* asks for the time, and then a little later process *B* asks for the time, the value that *B* gets will be higher than (or possibly equal to) the value *A* got.
  - It will certainly not be lower.
- In a distributed system, achieving agreement on time is not **trivial**.
- Is it possible to synchronize all the clocks in a distributed system?
  - **The answer is surprisingly complicated.**

# Computer's clock (timer)

- All computers have a circuit for keeping track of time. Despite the widespread use of the word "clock" to refer to these devices, they are not actually clocks in the usual sense.
  - **Timer** is perhaps a better word.
- A computer timer is usually a precisely **machined quartz crystal**.
- When kept under tension, quartz crystals oscillate at a well-defined frequency that depends on the kind of crystal, how it is cut, and the amount of tension.
- Associated with each crystal are two registers, a **counter** and a **holding register**.
- Each oscillation of the crystal decrements the counter by one. When the counter gets to zero, an **interrupt** is generated and the counter is reloaded from the holding register.
- In this way, it is possible to program a timer to generate an interrupt **60 times a second**, or at any other desired frequency.
- Each interrupt is called one **clock tick**.

# Clocks on different machines

- With a single computer and a single clock, it does not matter much if this clock is off by a small amount.
  - Since all processes on the machine use the same. clock, they will still be internally consistent.
- As soon as multiple CPUs are introduced, each with its own clock, the situation **changes radically**.
  - Although the frequency at which a crystal oscillator runs is usually **fairly stable**, it is impossible to guarantee that the crystals in different computers all run at exactly the same frequency.
- In practice, when a system has  $n$  computers, all  $n$  crystals will run at slightly **different rates**, causing the (software) clocks gradually to get out of synch and give different values when read out.
  - This difference in time values is called **clock skew**.

# External Physical Clocks

- In some systems (e.g., real-time systems), the actual clock time is important.
- Under these circumstances, external physical clocks are needed.
- For reasons of efficiency and redundancy, multiple physical clocks are generally considered desirable, which yields two problems:
  - (1) How do we synchronize them with realworld clocks and
  - (2) How do we synchronize the clocks with each other?

# Solar second

- Since the invention of mechanical clocks in the 17th century, time has been measured **astronomically**.
- Every day, the sun appears to rise on the eastern horizon, then climbs to a maximum height in the sky, and finally sinks in the west.
- The event of the sun's reaching its highest apparent point in the sky is called the **transit of the sun**.
- This event occurs at about noon each day.
- The interval between two consecutive transits of the sun is called the **solar day**.
- Since there are 24 hours in a day, each containing 3600 seconds, the **solar second** is defined as exactly  $1/86400$ th of a solar day.

# Solar second

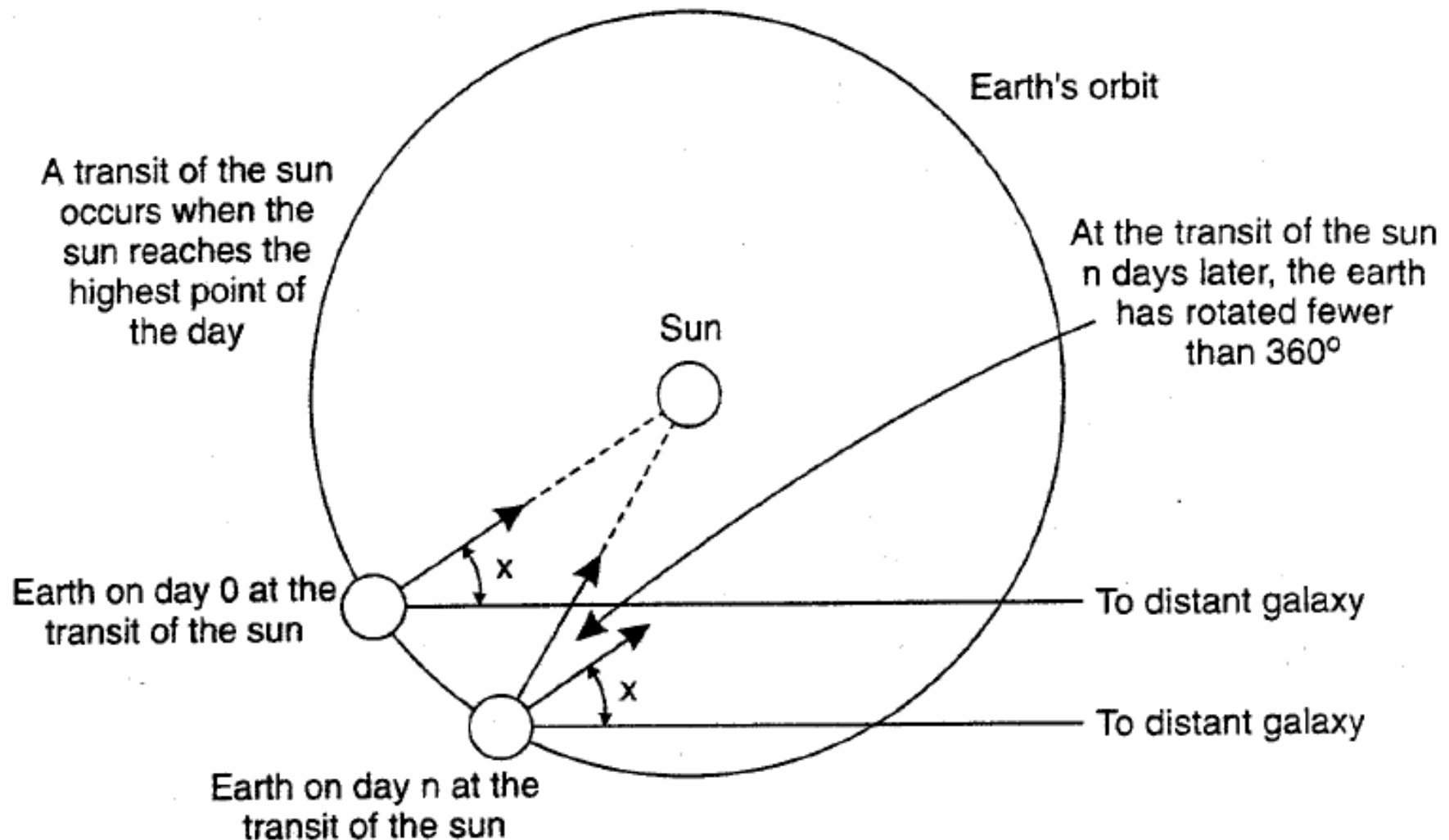


Figure 6-2. Computation of the mean solar day.

# Mean Solar Second

- In the 1940s, it was established that the period of the earth's rotation is **not constant**.
- The earth is slowing down due to tidal friction and atmospheric drag.
- Based on studies of growth patterns in ancient coral, geologists now believe that 300 million years ago there were about 400 days per year.
- The length of the year (the time for one trip around the sun) is not thought to have changed; the day **has simply become longer**.
- In addition to this long-term trend, **short-term variations** in the length of the day also occur, probably caused by turbulence deep in the earth's core of molten iron.
- These revelations led astronomers to compute the length of the day by measuring a large number of days and taking the average before dividing by 86,400.
- The resulting quantity was called the **mean solar second**.

# TAI

- With the invention of the atomic clock in 1948, it became possible to measure time much more accurately, and independent of the wiggling and wobbling of the earth, by counting transitions of the cesium 133 atom.
  - The physicists took over the job of timekeeping from the astronomers and defined the second to be the time it takes the cesium 133 atom to make exactly 9,192,631,770 transitions.
- The choice of 9,192,631,770 was made to make the atomic second equal to the mean solar second in the year of its introduction.
- Currently, several laboratories around the world have cesium 133 clocks.
- Periodically, each laboratory tells the **Bureau International de l'Heure (BIR)** in Paris how many times its clock has ticked. The BIR averages these to produce International Atomic Time, which is abbreviated TAI.
- Thus TAI is just the mean number of ticks of the cesium 133 clocks since midnight on Jan. 1, 1958 (the beginning of time) divided by 9,192,631,770.

# Time really matters

- Although TAI is highly stable and available to anyone who wants to go to the trouble of **buying a cesium clock**, there is a serious problem with it;
  - 86,400 TAI seconds is now about 3 msec **less than a mean solar day** (because the mean solar day is getting longer all the time).
- Using TAI for keeping time would mean that over the course of the years, noon would get **earlier and earlier**, until it would eventually occur in the wee hours of the morning.
- People might notice this and we could have the same kind of situation as occurred in 1582 when Pope Gregory XIII decreed that 10 days be omitted from the calendar.
- This event caused riots in the streets because landlords demanded a full month's rent and bankers a full month's interest, while employers refused to pay workers for the 10 days they did not work, to mention only a few of the conflicts.
  - The Protestant countries, as a matter of principle, refused to have anything to do with papal decrees and did not accept the Gregorian calendar for 170 years. 😊

# TAI Vs Solar Seconds

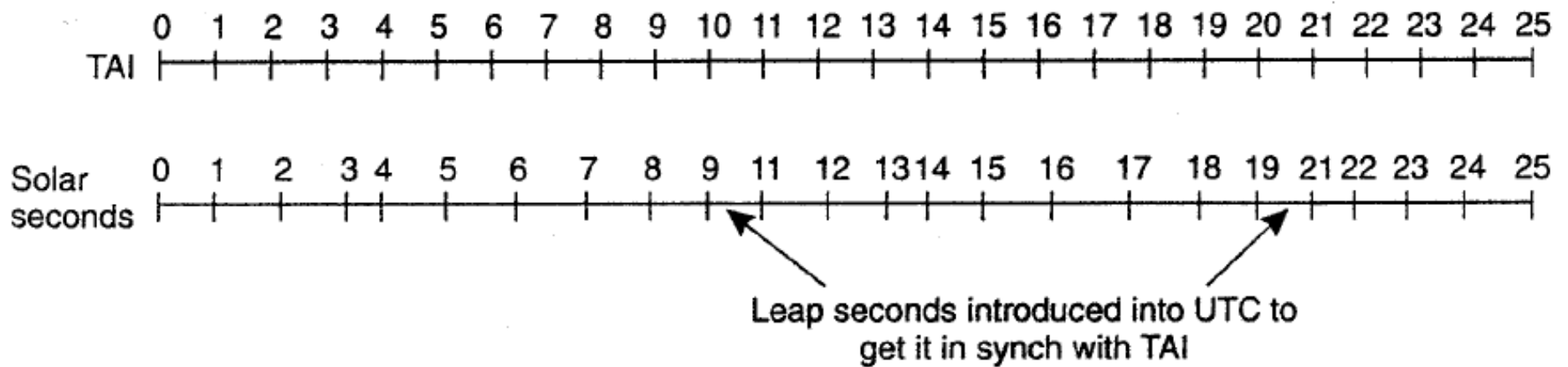


Figure 6-3. TAI seconds are of constant length, unlike solar seconds. Leap seconds are introduced when necessary to keep in phase with the sun.

BIR solves the problem by introducing leap seconds whenever the discrepancy between TAI and solar time grows to 800 msec. The use of leap seconds as a correction gives rise to a time system based on constant TAI seconds but which stays in phase with the apparent motion of the sun.

It is called **Universal Coordinated Time**, but is abbreviated as **UTC**. UTC is the basis of all modern civil timekeeping. It has essentially replaced the old standard, Greenwich Mean Time which is astronomical time.

# Simple hardware correction

- Most electric power companies **synchronize** the timing of their 60-Hz or 50-Hz clocks to UTC, so when BIH announces a leap second, the power companies raise their frequency to 61 Hz or 51 Hz for 60 or 50 sec. to advance all the clocks in their distribution area.
- Since 1 sec is a **noticeable interval** for a computer, an operating system that needs to keep accurate time over a period of years must have special software to account for leap seconds as they are announced (unless they use the power line for time, which is usually too crude).
- The total number of leap seconds introduced into UTC so far is about 30.

# NIST

- To provide UTC to people who need precise time, the **National Institute of Standard Time (NIST)** operates a shortwave radio station with call letters WWV from Fort Collins, Colorado.
- WWV broadcasts a short pulse at the start of each UTC second.
- The accuracy of WWV itself is about  $\pm 1$  msec, but due to random atmospheric fluctuations that can affect the length of the signal path, in practice the accuracy is no better than  $\pm 10$  msec.
- Several earth satellites also offer a UTC service. The Geostationary Environment Operational Satellite can provide UTC accurately to **0.5 msec**, and some other satellites do even better.

# Global Positioning System

- As a step toward actual **clock synchronization** problems, we first consider a related problem, namely determining one's geographical position anywhere on Earth.
- This positioning problem is by itself solved through a highly specific dedicated distributed system, namely GPS, which is an acronym for **global positioning system**.
- GPS is a **satellite-based distributed system** that was launched in 1978.
- Although it has been used mainly for military applications, in recent years it has found its way to many civilian applications, notably for traffic navigation.
- However, many more application domains exist. For example, GPS phones now allow to let callers track each other's position, a feature which may show to be extremely handy when you are lost or in trouble.

# GPS

- GPS uses 29 satellites each circulating in an orbit at a height of approximately 20,000 km.
- Each satellite has up to **four atomic clocks**, which are regularly calibrated from special stations on Earth.
- A satellite **continuously** broadcasts its position, and time stamps each message with its local time.
- This broadcasting allows every receiver on Earth to accurately compute its own position using, in principle, **only three satellites**.



# GPS

- Let's assume that all clocks, including the receiver's, are synchronized.
- In order to compute a position, consider first the two-dimensional case, as shown here, in which two satellites are drawn, along with the circles representing points at the same distance from each respective satellite.
- The *y-axis* represents the height, while the *x-axis* represents a straight line along the Earth's surface at sea level. Ignoring the highest point,
- We see that the intersection of the two circles is a unique point

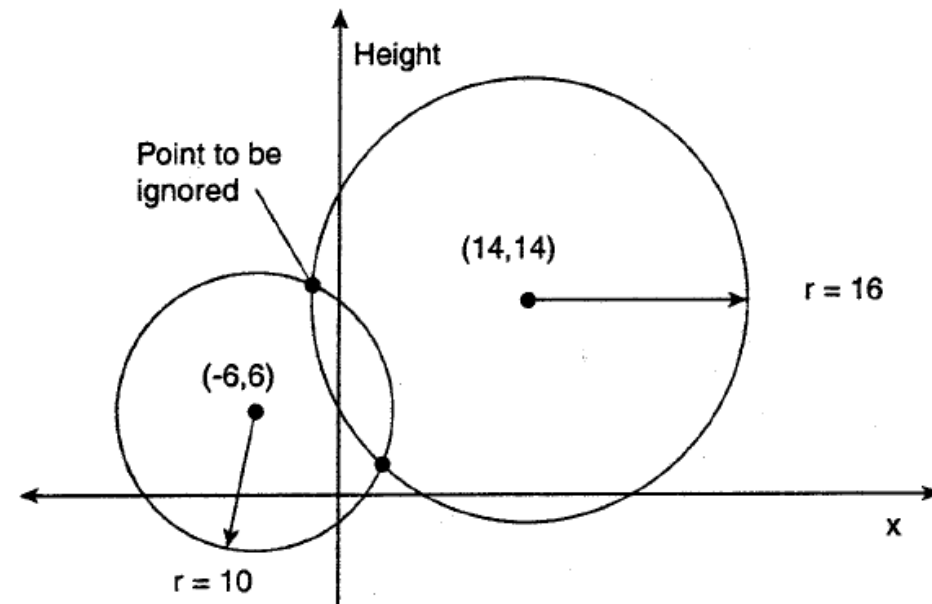


Figure 6-4. Computing a position in a two-dimensional space.

# GPS and synchronization

- The principle of intersecting circles can be expanded to three dimensions, meaning that we need three satellites to determine the longitude, latitude, and altitude of a receiver on Earth.
- This positioning is all fairly straightforward, but matters become complicated when we can **no longer assume that all clocks are perfectly synchronized.**
- There are two important real-world facts that we need to take into account:
  1. **It takes a while before data on a satellite's position reaches the receiver,**
  2. **The receiver's clock is generally not in synch with that of a satellite.**

# GPS: into the unknown

## Principal operation

- $\Delta_r$ : unknown deviation of the receiver's clock.
- $x_r, y_r, z_r$ : unknown coordinates of the receiver.
- $T_i$ : timestamp on a message from satellite  $i$
- $\Delta_j = (T_{now} - T_i) + \Delta_r$ : measured delay of the message sent by satellite  $i$ .
- Measured distance to satellite  $i$ :  $c \times \Delta_j$   
( $c$  is speed of light)
- Real distance is

$$d_i = c\Delta_j - c\Delta_r = \sqrt{(x_i - x_r)^2 + (y_i - y_r)^2 + (z_i - z_r)^2}$$

## Observation

4 satellites  $\Rightarrow$  4 equations in 4 unknowns (with  $\Delta_r$  as one of them)

# Machine Clocks and UTC

- There is a clock in machine  $p$  that ticks on each timer interrupt.
- Denote the value of that clock by  $C_p(t)$ , where  $t$  is UTC time.
- Ideally, we have that for each machine  $p$ ,  $C_p(t) = t$ , or, in other words,  $dC/dt = 1$ .
- The constant  $\rho$  is specified by the manufacturer and is known as the maximum drift rate.

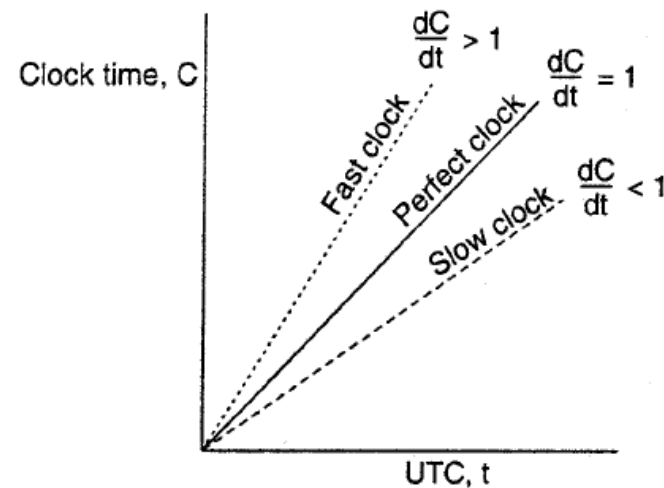


Figure 6-5. The relation between clock time and UTC when clocks tick at different rates.

In practice:  $1 - \rho \leq \frac{dC}{dt} \leq 1 + \rho$ .

## Goal

Never let two clocks in any system differ by more than  $\delta$  time units  $\Rightarrow$  synchronize at least every  $\delta / (2\rho)$  seconds.

# Clock synchronization: Principle I

- **Principle I**
  - Every machine asks a time server for the accurate time at least once every  $\delta/2p$  seconds (Network Time Protocol).
- **Note**
  - Okay, but you need an **accurate measure** of round trip delay, including interrupt handling and processing incoming messages.
- If the operating system designers want to guarantee that no two clocks ever differ by more than  $\delta$ , clocks must be resynchronized (in software) at least every  $\delta/2p$  seconds.
- **The various algorithms differ in precisely how this resynchronization is done.**

# Clock synchronization: Principle II

## Principle II

- Let the time server scan all machines **periodically**, calculate an average, and inform each machine how it should adjust its time relative to its present time.

## Note

- Okay, you'll probably get every machine in sync. You don't even need to propagate UTC time.



# Network Time Protocol

- A common approach in many protocols is to let clients contact a time server.
  - The latter can accurately provide the current time, for example, because it is equipped with a WWV receiver or an accurate clock.
- The problem, of course, is that when contacting the server, **message delays** will have outdated the reported time.
  - **The trick is to find a good estimation for these delays.**

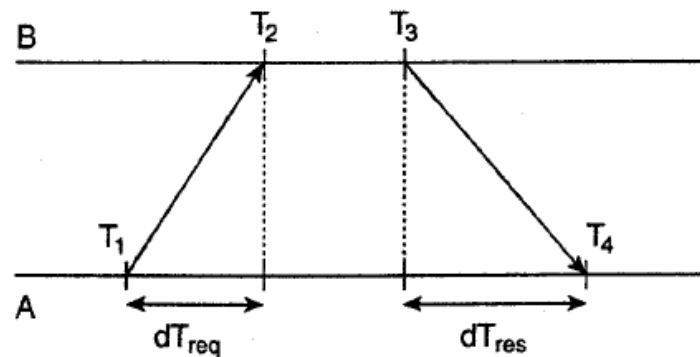


Figure 6-6. Getting the current time from a time server.

## Time Corrections

You'll have to take into account that setting the time back is never allowed => smooth adjustments.

# Berkeley algorithm

- The time server (actually, a time daemon) is **active**, polling every machine from time to time to ask what time it is there.
- Based on the answers, it computes an **average time** and tells all the other machines to advance their clocks to the new time or slow their clocks down until some specified reduction has been achieved.
- This method is suitable for a system in which no machine has a WWV receiver.
- **The time daemon's time must be set manually by the operator periodically.**

# Berkeley algorithm

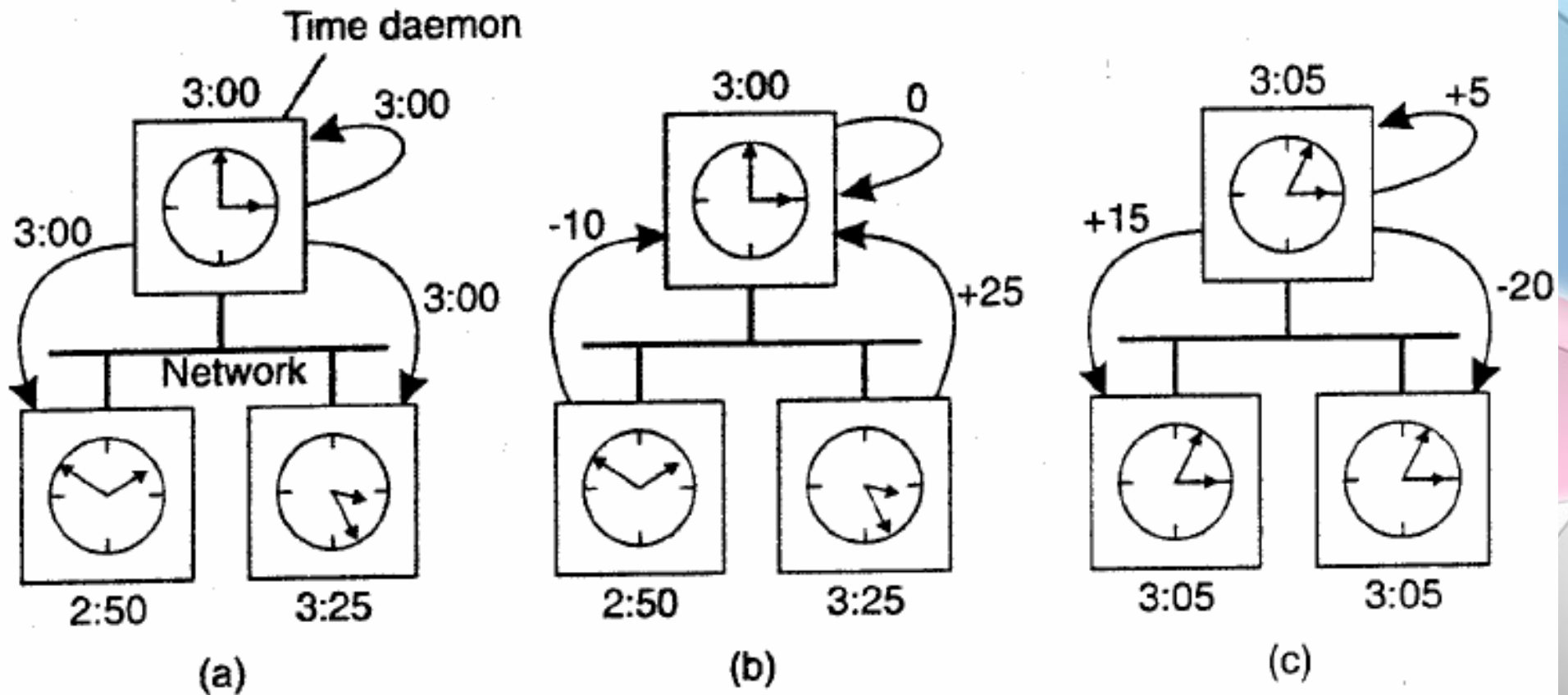


Figure 6-7. (a) The time daemon asks all the other machines for their clock values. (b) The machines answer. (c) The time daemon tells everyone how to adjust their clock.

# Clock Synchronization in Wireless Networks

- In many wireless networks, notably sensor networks, nodes are **resource constrained**, and multihop routing is expensive.
- In addition, it is often important to optimize algorithms for energy consumption.
- These and other observations have led to the design of very different clock synchronization algorithms for wireless networks.
- **Reference broadcast synchronization** (RBS) is a clock synchronization protocol that is quite different from other proposals.
- First, the protocol **does not assume that there is a single node** with an accurate account of the actual time available.
- Instead of aiming to provide all nodes UTC time, it aims at merely internally synchronizing the clocks, just as the Berkeley algorithm does.

# Logical Clocks

- So far, we have assumed that clock synchronization is naturally **related to real time**.
- However, it may be sufficient that every node agrees on a current time, without that time necessarily being the same as the real time.
- For most algorithms, it is conventional to speak of the clocks as **logical clocks**.
- We showed that although clock synchronization is possible, **it need not be absolute**.
- If two processes do not interact, it is not necessary that their clocks be synchronized because the lack of synchronization would not be observable and thus could not cause problems.
- Furthermore, Lamport pointed out that what usually matters is not that all processes agree on exactly what time it is, but rather that they **agree on the order** in which events occur.
- **Lamport Algorithm for synchronizing logical clocks**

# The Happened-before relationship

## Problem

- We first need to introduce a notion of ordering before we can order anything.

## The happened-before relation

- If  $a$  and  $b$  are two events in the same process, and  $a$  comes before  $b$ , then  $a \rightarrow b$ .
- If  $a$  is the sending of a message, and  $b$  is the receipt of that message, then  $a \rightarrow b$
- If  $a \rightarrow b$  and  $b \rightarrow c$ , then  $a \rightarrow c$

## Note

- This introduces a partial ordering of events in a system with concurrently operating processes.

# Logical clocks

## Problem

- How do we maintain a global view on the system's behavior that is consistent with the happened-before relation?

## Solution

- Attach a **timestamp**  $C(e)$  to each event  $e$ , satisfying the following properties:
  - P1** If  $a$  and  $b$  are two events in the same process, and  $a \rightarrow b$ , then we demand that  $C(a) < C(b)$ .
  - P2** If  $a$  corresponds to sending a message  $m$ , and  $b$  to the receipt of that message, then also  $C(a) < C(b)$ .

## Problem

- How to attach a timestamp to an event when there's no global clock  
=> maintain a **consistent** set of logical clocks, one per process.

# Lamport algorithm

## Solution

Each process  $P_i$  maintains a local counter  $C_i$  and adjusts this counter according to the following rules:

- 1:** For any two successive events that take place within  $P_i$ ,  $C_i$  is incremented by 1.
- 2:** Each time a message  $m$  is sent by process  $P_i$ , the message receives a timestamp  $ts(m) = C_i$ .
- 3:** Whenever a message  $m$  is received by a process  $P_j$ ,  $P_j$  adjusts its local counter  $C_j$  to  $\max\{C_j ; ts(m)\}$ ; then executes step 1 before passing  $m$  to the application.

## Notes

- Property P1 is satisfied by (1); Property P2 by (2) and (3).

# Lamport algorithm

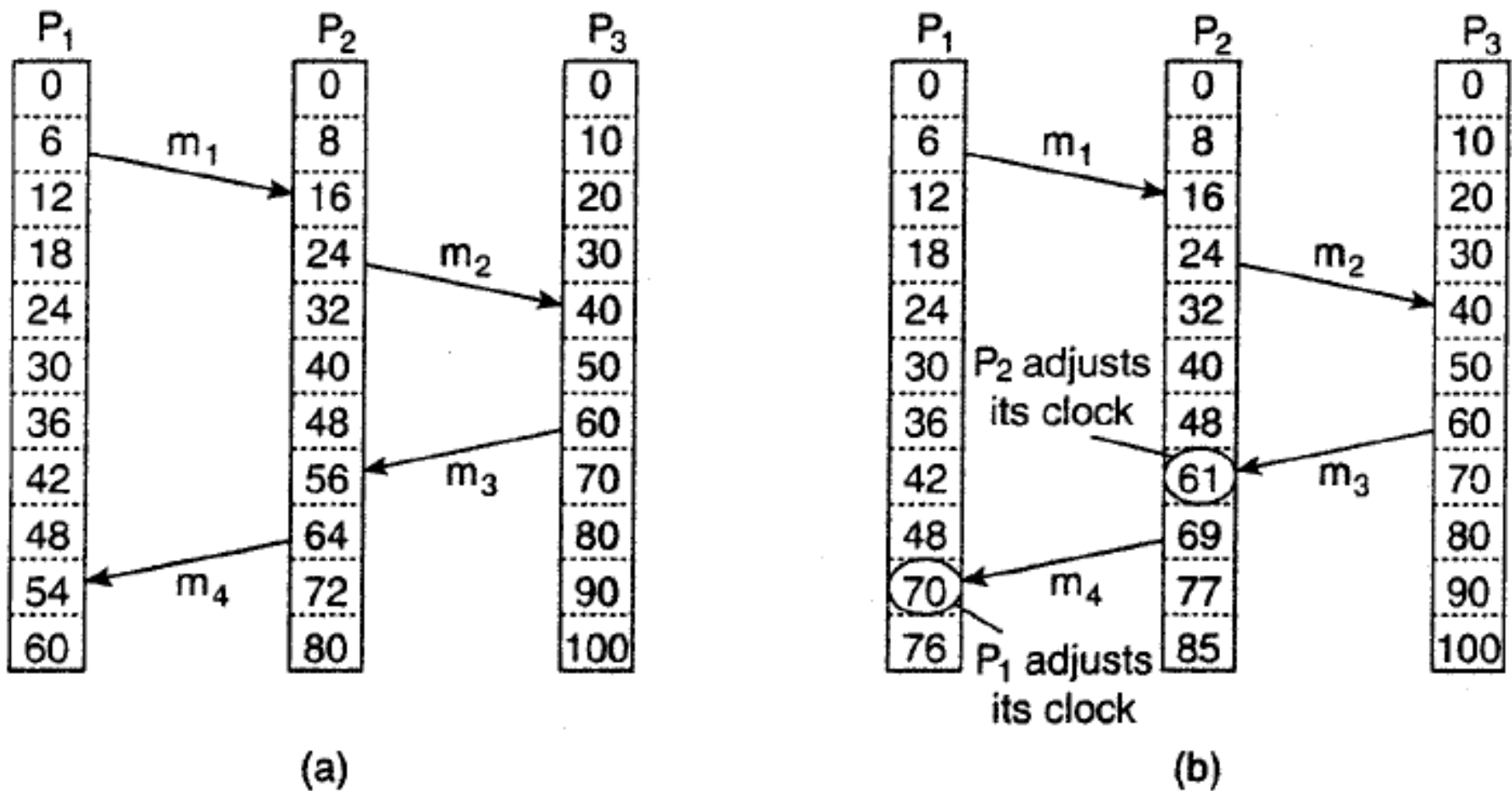


Figure 6-9. (a) Three processes, each with its own clock. The clocks run at different rates. (b) Lamport's algorithm corrects the clocks.

# Logical clocks among layers

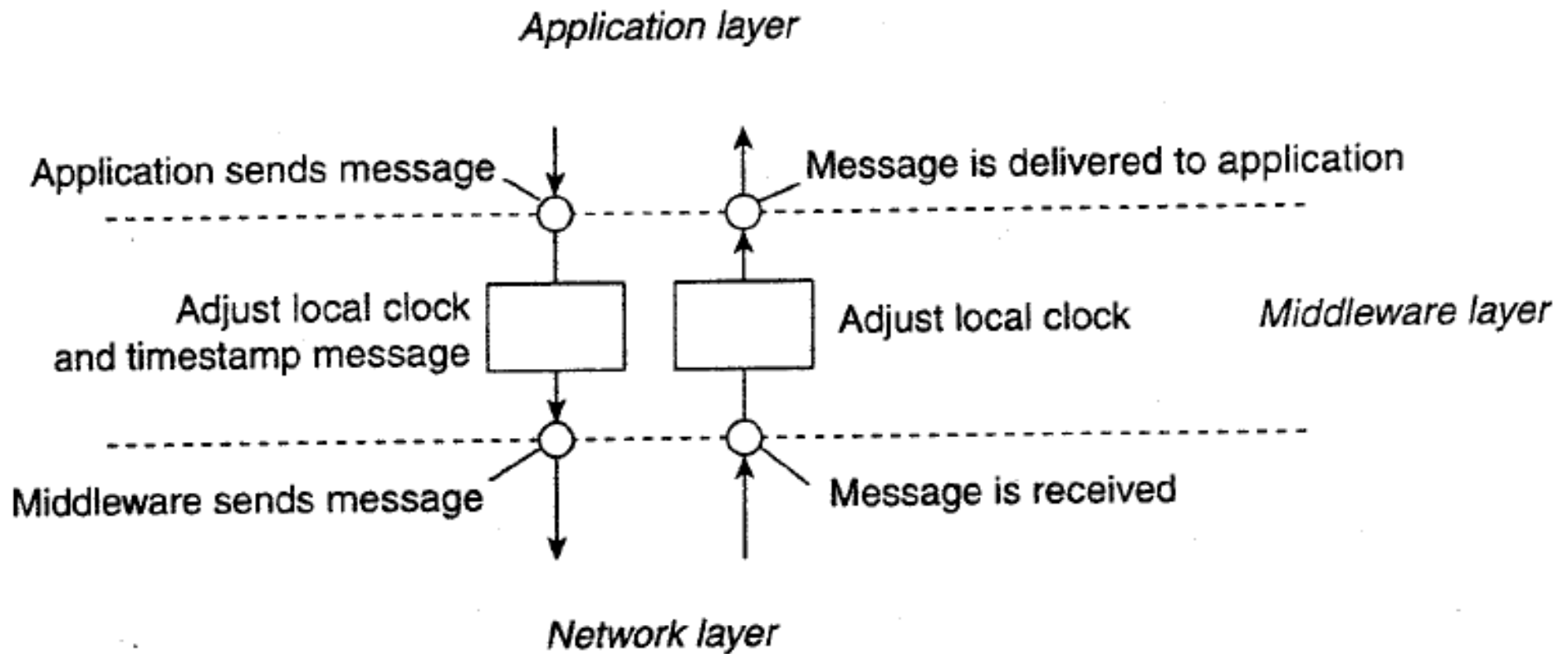


Figure 6-10. The positioning of Lamport's logical clocks in distributed systems.

# Example: Totally ordered multicast

- For example, to improve query performance, a bank may place copies of an account database in two different cities, say New York and San Francisco.
- **A query is always forwarded to the nearest copy.**
  - The price for a fast response to a query is partly paid in higher update costs, because each update operation must be carried out at each replica.
- In fact, there is a more stringent requirement with respect to updates. Assume a customer in San Francisco wants to add \$100 to his account, which currently contains \$1,000. At the same time, a bank employee in New York initiates an update by which the customer's account is to be increased with **1 percent interest**.
  - Both updates should be carried out at **both copies** of the database.
- However, due to **communication delays** in the underlying network, the updates may arrive in the order => next slide:

# Example: Totally ordered multicast

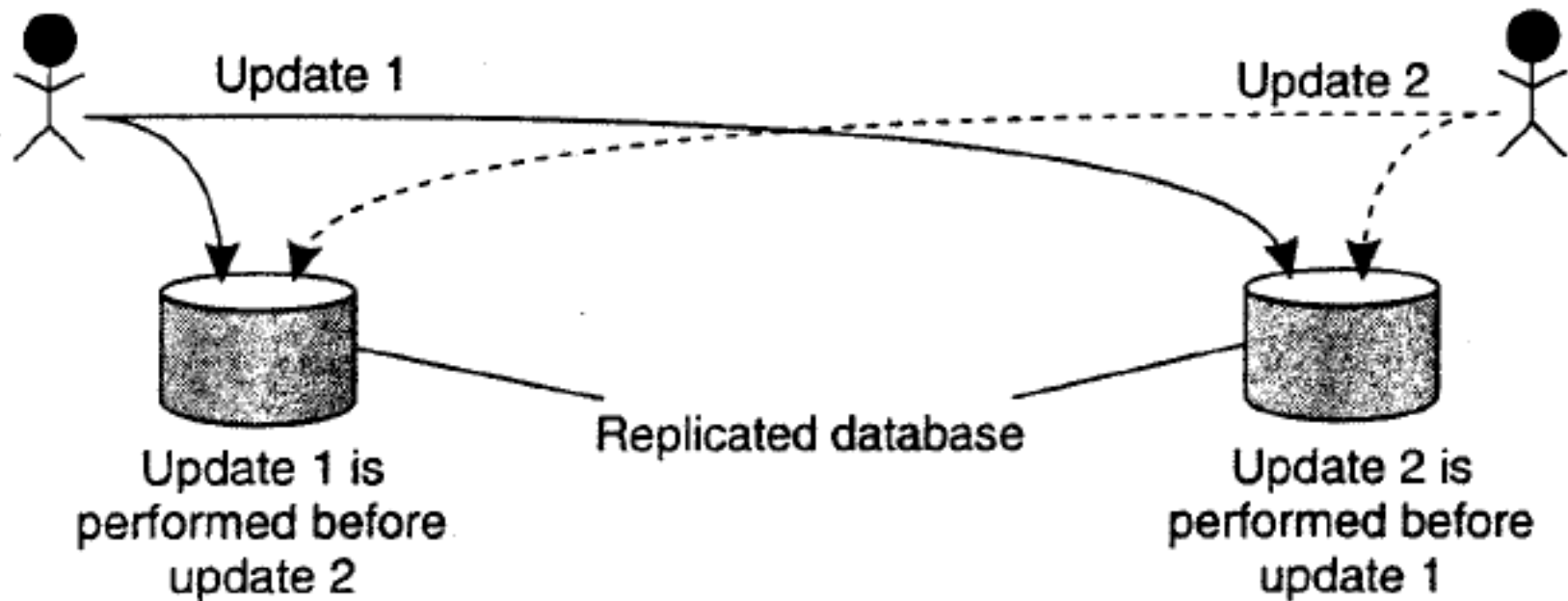


Figure 6-11. Updating a replicated database and leaving it in an inconsistent state.

# Vector Clocks

- Lamport's logical clocks lead to a situation where all events in a distributed system are totally ordered with the property that if event  $a$  happened before event  $b$ , then  $a$  will also be positioned in that ordering before  $b$ , that is,  $C(a) < C(b)$ .
- However, with Lamport clocks, nothing can be said about the relationship between two events  $a$  and  $b$  by merely comparing their time values  $C(a)$  and  $C(b)$ , respectively.
- In other words, if  $C(a) < C(b)$ , then this does not necessarily imply that  $a$  indeed happened before  $b$ .
- Something more is needed for that.

# Vector Clocks

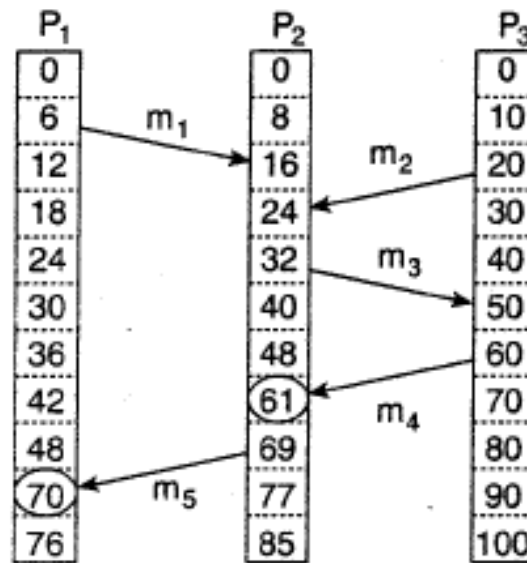


Figure 6-12. Concurrent message transmission using logical clocks.

In the case of  $m_1$  and  $m_3$ , we know that these values correspond to events that took place at process  $P_2$ , meaning that  $m_3$  was indeed sent after the receipt of message  $m_1$ .

This may indicate that the sending of message  $m_3$  depended on what was received through message  $m_1$ .

However, we also know that  $Trcv(m_1) < Tsnd(m_2)$ . However, the sending of  $m_2$  has nothing to do with the receipt of  $m_1$ .

# Vector Clocks

- The problem is that Lamport clocks do not capture **causality**. Causality can be captured by means of **vector** clocks. A vector clock  $VC(a)$  assigned to an event  $a$  has the property that if  $VC(a) < VC(b)$  for some event  $b$ , then event  **$a$  is known to causally precede event  $b$** .
- Vector clocks are constructed by letting each process  $P_i$  maintain a vector  $VC_i$  with the following two properties:
  1.  $VC_i[i]$  is the number of events that have occurred so far at  $P_i$ . In other words,  $VC_i[i]$  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  $P_j$ . It is thus  $P_i$ 's knowledge of the local time at  $P_j$ .

# Vector Clocks

Steps carried out to accomplish property 2 of previous slide:

1. Before executing an event,  $P_i$  executes  $VC_i[i] \leftarrow VC_i[i] + 1$ .
2. When process  $P_i$  sends a message  $m$  to  $P_j$ , it sets  $m$ 's (vector) timestamp  $ts(m)$  equal to  $VC_i$  after having executed the previous step.
3. Upon the receipt of a message  $m$ , process  $P_j$  adjusts its own vector by setting  $VC_j[k] \leftarrow \max\{VC_j[k], ts(m)[k]\}$  for each  $k$ , after which it executes the first step and delivers the message to the application.

# Enforcing Causal Communication

- Using vector clocks, it is now possible to ensure that a message is delivered only if all messages that **causally precede** it have also been received as well.
- To enable such a scheme, we will assume that messages are multicast within a group of processes.
  - Note that this causally-ordered multicasting is weaker than the totally-ordered multicasting we discussed earlier.
- Specifically, if two messages are not in any way related to each other, we do not care in which order they are delivered to applications.
- They may even be delivered in different order at different locations.

# Enforcing Causal Communication

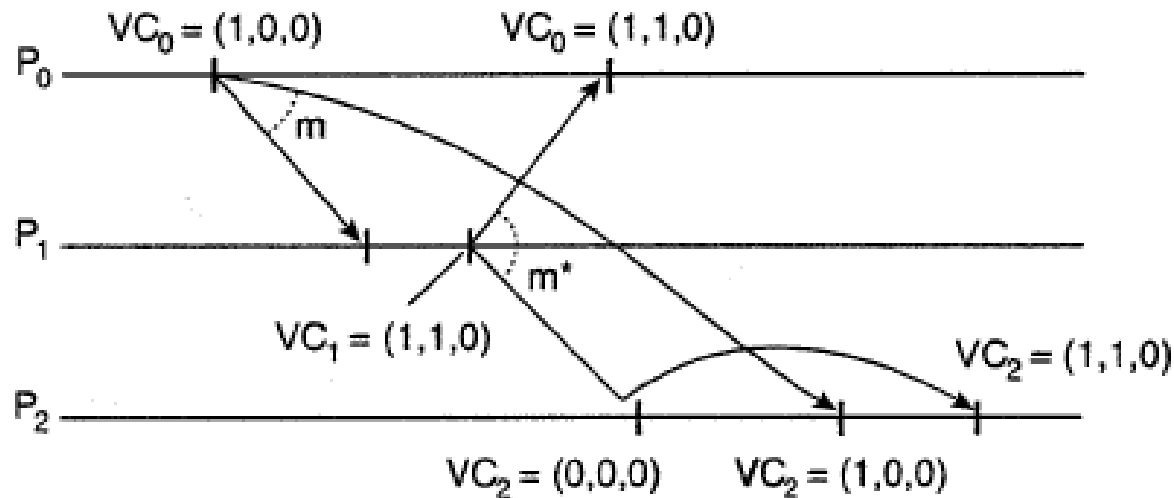
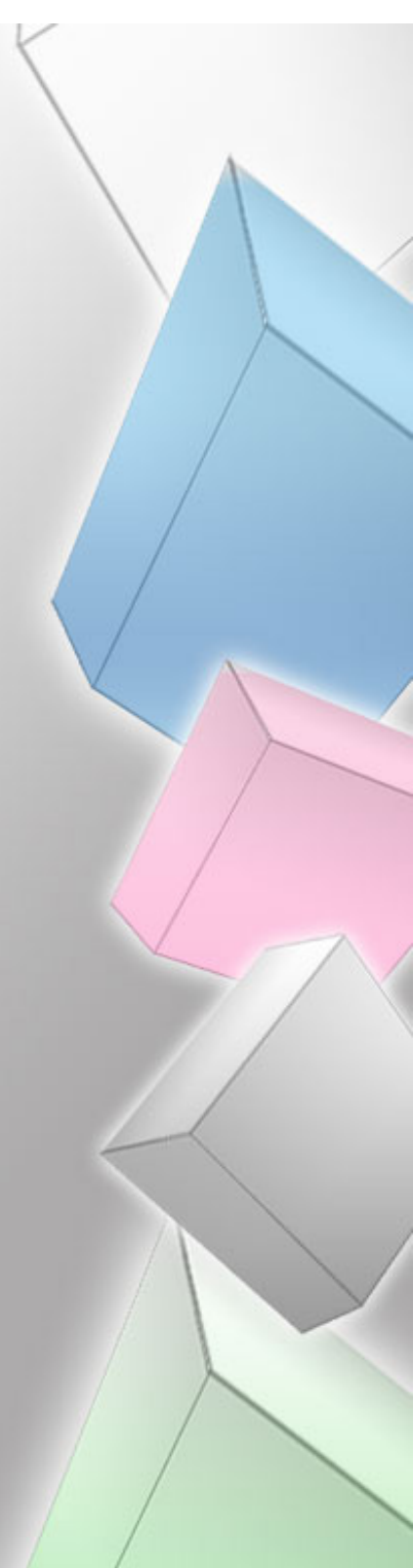


Figure 6-13. Enforcing causal communication.

- At local time  $(1,0,0)$ ,  $P_0$  sends message  $m$  to the other two processes.
- After its receipt by  $P_1$ , the latter decides to send  $m^*$ , which arrives at  $P_2$  **sooner** than  $m$ .
- At that point, **the delivery of  $m^*$  is delayed** by  $P_2$  until  $m$  has been received and delivered to  $P_2$ 's application layer

# Concurrency in Distributed Systems



# Concurrency in Distributed Systems

- Fundamental to distributed systems is the **concurrency** and **collaboration** among multiple processes.
- In many cases, this also means that processes will need to **simultaneously access** the same resources.
- To prevent that such concurrent accesses corrupt the resource, or make it inconsistent, solutions are needed to grant **mutual exclusive access** by processes.
- In this section, we take a look at some of the more important **distributed algorithms** that have been proposed.



# Token-based solutions

- In token-based solutions **mutual exclusion** is achieved by passing a special message between the processes, known as a **token**.
- There is only **one token available** and who ever has that token is allowed to access the shared resource.
- When finished, the token is passed on to a next process.
- If a process having the token is **not interested** in accessing the resource, it simply passes it on.

# Token-based solutions

- Token-based solutions have a few important properties.
- First, depending on how the processes are organized, they can fairly easily ensure that **every process will get a chance** at accessing the resource.
  - **In other words, they avoid starvation.**
- Second, **deadlocks** by which several processes are waiting for each other to proceed, can easily be avoided, contributing to their simplicity.
- Unfortunately, the main drawback of token-based solutions is a rather serious one:
  - when the token is **lost** (e.g., because the process holding it crashed), an intricate distributed procedure needs to be started to ensure that a new token is created, but above all, that it is also the only token.

# Permission-based approach

- As an alternative, many distributed mutual exclusion algorithms follow a permission-based approach.
- In this case a process wanting to access the resource **first requires the permission of other processes.**
- There are many different ways toward granting such a permission:
  - Via a centralized server.
  - Completely decentralized, using a peer-to-peer system.
  - Completely distributed, with no topology imposed.
  - Completely distributed along a (logical) ring.

# Centralized Algorithm

- The most straightforward way to achieve mutual exclusion in a distributed system is to simulate how it is done in a **one-processor system**.
- One process is elected as the **coordinator**. Whenever a process wants to access a shared resource, it sends a request message to the coordinator stating which resource it wants to access and asking for permission.
- If **no other process is currently accessing** that resource, the coordinator sends back a reply granting permission.

# Centralized Algorithm

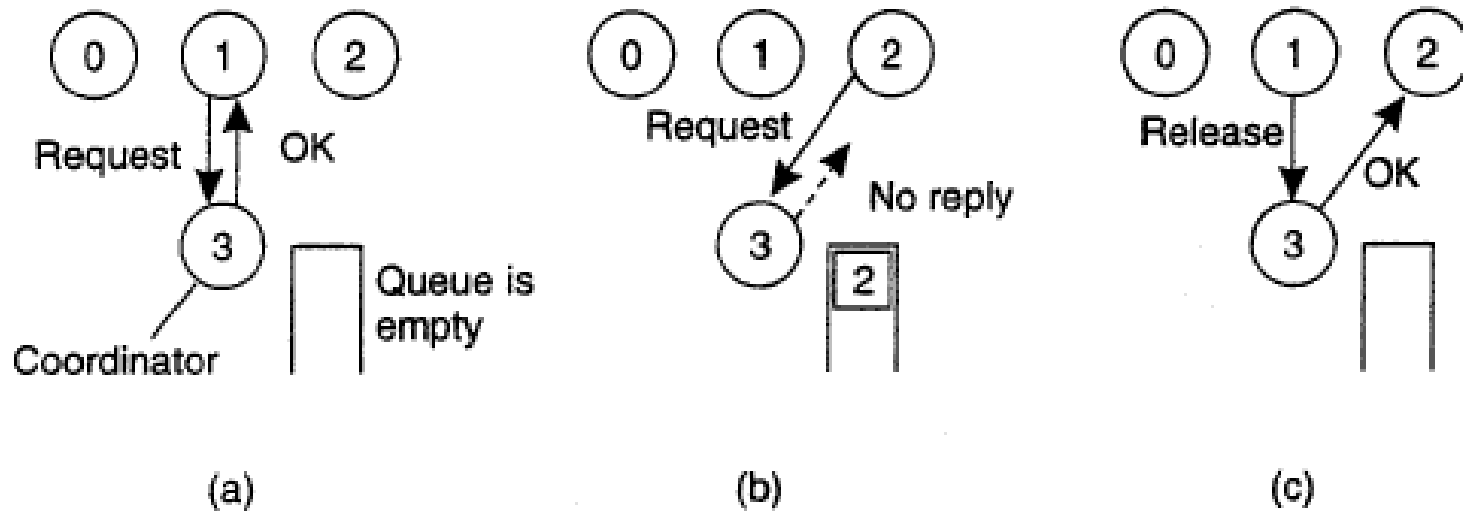


Figure 6-14. (a) Process 1 asks the coordinator for permission to access a shared resource. Permission is granted. (b) Process 2 then asks permission to access the same resource. The coordinator does not reply. (c) When process 1 releases the resource, it tells the coordinator, which then replies to 2.

# Centralized Algorithm

- The centralized approach has shortcomings.
- The coordinator is a **single point of failure**, so if it crashes, the entire system may go down.
- If processes normally block after making a request, they cannot distinguish a **dead coordinator** from "**permission denied**" since in both cases no message comes back.
- In addition, in a large system, a single coordinator can become a **performance bottleneck**.
- Nevertheless, the benefits coming from its simplicity outweigh in many cases the potential drawbacks.
- Moreover, distributed solutions are not necessarily better, as we will see next.

# Decentralized Algorithm

## Principle

- Assume every resource is **replicated  $n$  times**, with each replica having its own coordinator => access requires a majority vote from  $m > n/2$  coordinators.
- A coordinator always responds immediately to a request.

## Assumption

When a coordinator crashes, it will recover quickly, but will have forgotten about permissions it had granted.

- It has been shown that this algorithm is **probabilistically correct**: the probability of violating correctness is much smaller than the availability of any resource.

# Distributed algorithm

- To many, having a probabilistically correct algorithm is just not good enough.
- So researchers have looked for deterministic distributed mutual exclusion algorithms.
  - Lamport's 1978 paper on clock synchronization presented the first one.
  - Ricart and Agrawala (1981) made it more efficient.
- This algorithm requires that there be a total ordering of all events in the system.
- That is, for any pair of events, such as messages, it must be unambiguous which one actually happened first.
- Lamport's algorithm is one way to achieve this ordering and can be used to provide timestamps for distributed mutual exclusion.

# Distributed algorithm

- The algorithm works as follows. When a process wants to access a shared resource, it **builds a message** containing the name of the resource, its process number, and the current (logical) time.
- It then sends the message to all other processes, conceptually including itself. **The sending of messages is assumed to be reliable; that is, no message is lost.**
- When a process receives a request message from another process, the action it takes depends on its own state with respect to the resource named in the message.

## Three different cases:

1. If the receiver is not accessing the resource and does not want to access it, it sends back an OK message to the sender.
2. If the receiver already has access to the resource, it simply does not reply. Instead, it queues the request.
3. If the receiver wants to access the resource as well but has not yet done so, it compares the timestamp of the incoming message with the one contained in the message that it has sent everyone. **The lowest one wins.**

# Distributed algorithm

- After sending out requests asking permission, a process sits back and waits until everyone else has given permission.
- As soon as all the permissions are in, it may go ahead.
- When it is finished, it sends *OK* messages to all processes on its queue and deletes them all from the queue.

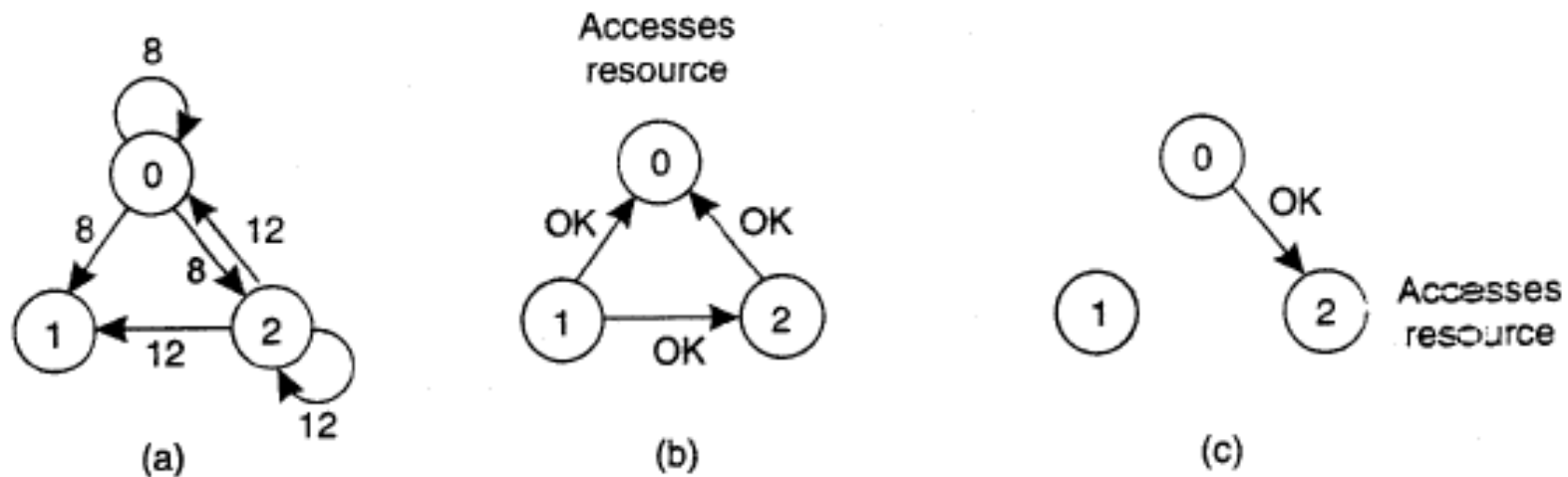


Figure 6-15. (a) Two processes want to access a shared resource at the same moment., (b) Process 0 has the lowest timestamp. so it wins. (c) When process 0 is done, it sends an *OK* also, so 2 can now go ahead.

# A Token Ring Algorithm

- In software, a logical ring is constructed in which each process is assigned a position in the ring.
- The ring positions may be allocated in numerical order of network addresses or some other means.
- It does not matter what the ordering is. All that matters is that each process knows who is next in line after itself.

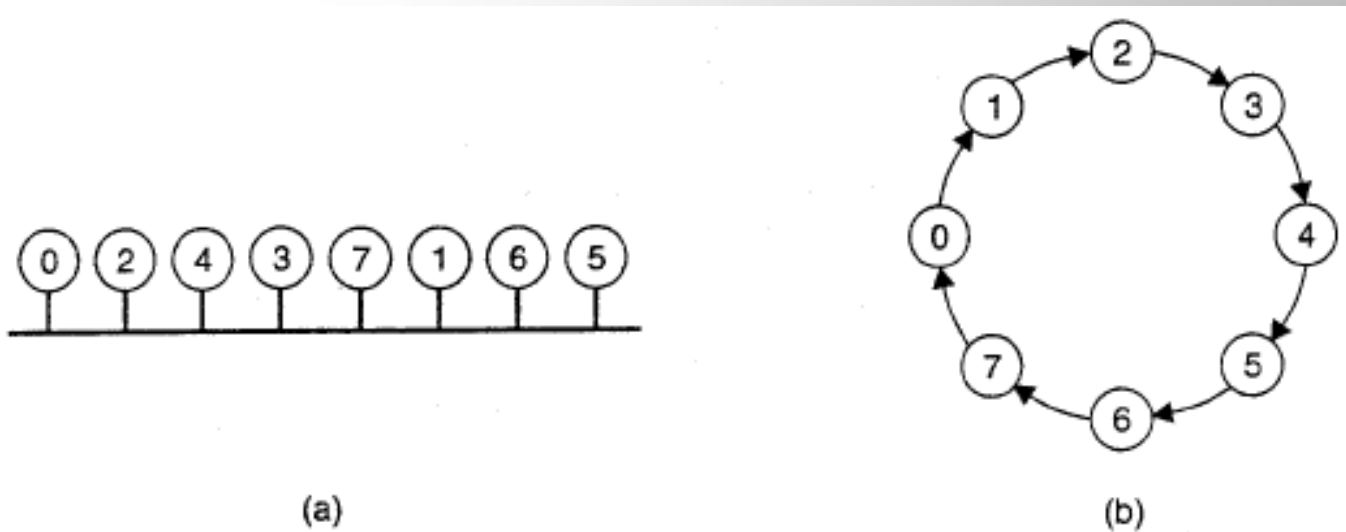


Figure 6-16. (a) An unordered group of processes on a network. (b) A logical ring constructed in software.

# A Token Ring Algorithm

- When the ring is initialized, process 0 is given a token.
- The token circulates around the ring. It is passed from process  $k$  to process  $k + 1$  in point-to-point messages.
- When a process acquires the token from its neighbor, it checks to see if it needs to access the shared resource.
- If so, the process goes ahead, does all the work it needs to, and releases the resources. After it has finished, it passes the token along the ring.
- It is not permitted to immediately enter the resource again using the same token.

# A Token Ring Algorithm

- As usual, this algorithm has problems too. If the token is **ever lost**, it must be regenerated.
  - In fact, detecting that it is lost is difficult, since the amount of time between successive appearances of the token on the network is unbounded. The fact that the token has not been spotted for an hour does not mean that it has been lost; somebody may still be using it.
- The algorithm also runs into trouble if a process crashes, but recovery is easier than in the other cases. If we require a process receiving the token to acknowledge receipt, a **dead process** will be detected when its neighbor tries to give it the token and fails.
- At that point the dead process can be **removed from the group**, and the token holder can throw the token over the head of the dead process to the next member down the line, or the one after that, if necessary.
- Of course, doing so requires that everyone maintain the current ring configuration.

# Algorithms Comparison

<b>Algorithm</b>	<b>Messages per entry/exit</b>	<b>Delay before entry (in message times)</b>	<b>Problems</b>
Centralized	3	2	Coordinator crash
Decentralized	$3mk, k = 1,2,\dots$	$2m$	Starvation, low efficiency
Distributed	$2(n - 1)$	$2(n - 1)$	Crash of any process
Token ring	1 to $\infty$	0 to $n - 1$	Lost token, process crash

# Global positioning of nodes

## Problem

- How can a single node efficiently estimate the latency between any two other nodes in a distributed system?

## • Solution

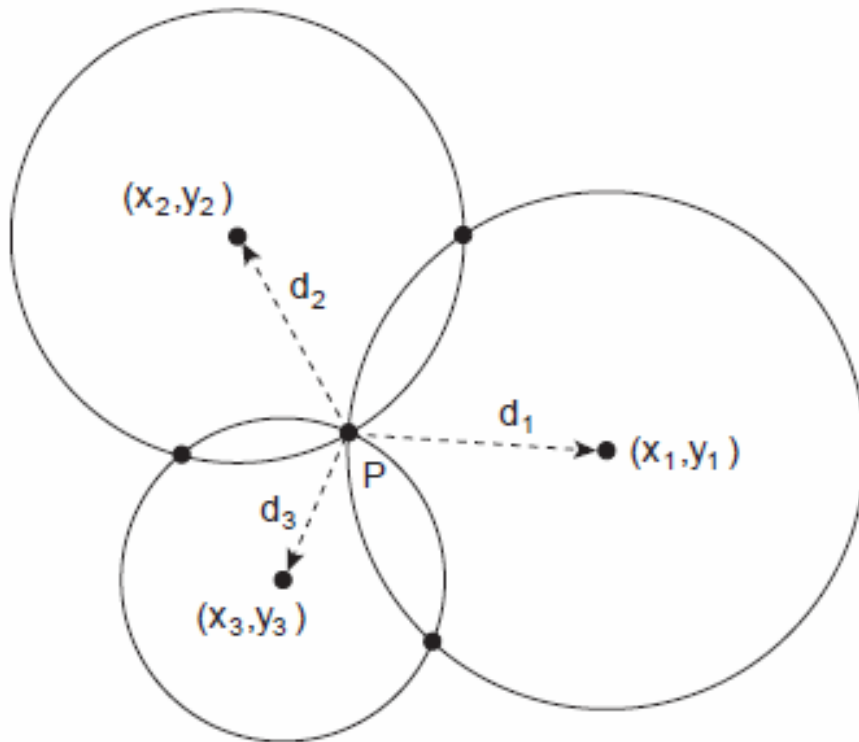
- Construct a **geometric overlay network**, in which the distance  $d(P;Q)$  reflects the actual latency between P and Q.



# Computing position

## Observation

- A node  $P$  needs  $m+1$  landmarks to compute its own position in a  $m$ -dimensional space. Consider two-dimensional case.



## Solution

$P$  needs to solve three equations in two unknowns  $(x_P, y_P)$ :

$$d_i = \sqrt{(x_i - x_P)^2 + (y_i - y_P)^2}$$

# Applications of Geometric Overlay Networks

- There are many applications of **geometric overlay networks**.
- Consider the situation where a Web site at server  $O$  has been replicated to multiple servers  $S_1 \dots S_k$  on the Internet.
- When a client  $C$  requests a page from  $O$ , the latter may decide to redirect that request to the server closest to  $C$ , that is, the one that will give the best response time.
- If the geometric location of  $C$  is known, as well as those of each replica server,  $O$  can then simply pick that server  $S_i$ , for which  $d(C, S_i)$  is minimal.
- Note that such a selection requires only local processing at  $O$ .

# Applications of Geometric Overlay Networks

- Another example is optimal replica placement.
- Consider again a Web site that has gathered the positions of its clients.
- If the site were to replicate its content to  $K$  servers, it can compute the  $K$  best positions where to place replicas such that the average client-to replica response time is minimal.
- Performing such computations is almost trivially feasible if clients and servers have geometric positions that reflect inter node latencies



# Election Algorithms

- Many distributed algorithms require one process to act as coordinator, initiator, or otherwise perform some special role. In general, it does not matter which process takes on this special responsibility, but one of them has to do it.
  - Algorithms for electing a coordinator (using this as a generic name for the special process).
- If all processes are exactly the same, with no distinguishing characteristics, there is no way to select one of them to be special .
  - We will assume that each process has a unique number, for example, its network address (for simplicity, we will assume one process per machine).
- In general, election algorithms attempt to locate the process with the highest process number and designate it as coordinator.
- The algorithms differ in the way they do the location.

# The Bully Algorithm

When any process notices that the coordinator is no longer responding to requests, it initiates an election. A process,  $P$ , holds an election as follows:

1.  $P$  sends an *ELECTION* message to all processes with higher numbers.
2. If no one responds,  $P$  wins the election and becomes coordinator.
3. If one of the higher-ups answers, it takes over.  $P$ 's job is done.

# Election by bullying

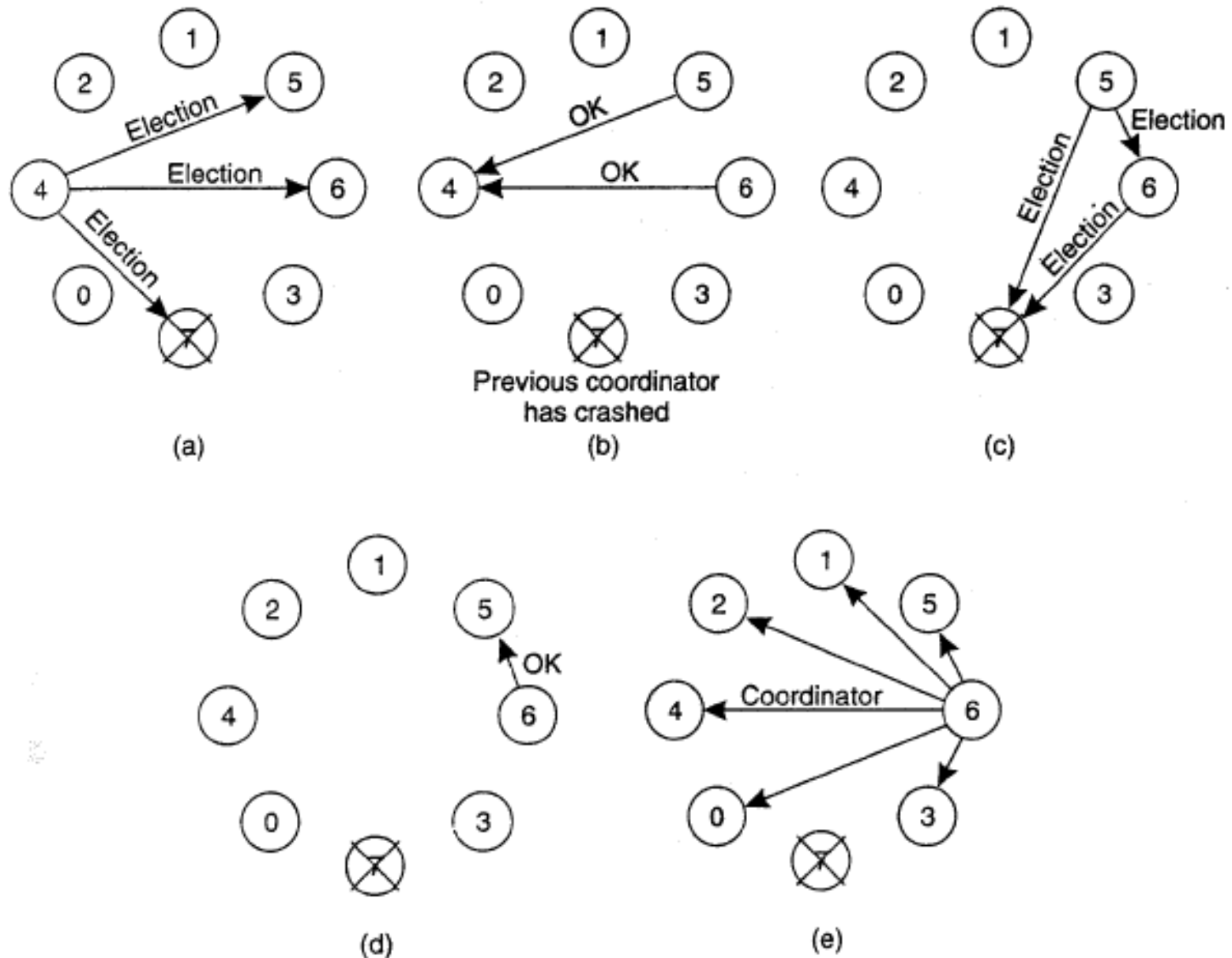


Figure 6-20. The bully election algorithm. (a) Process 4 holds an election. (b) Processes 5 and 6 respond, telling 4 to stop. (c) Now 5 and 6 each hold an election. (d) Process 6 tells 5 to stop. (e) Process 6 wins and tells everyone.

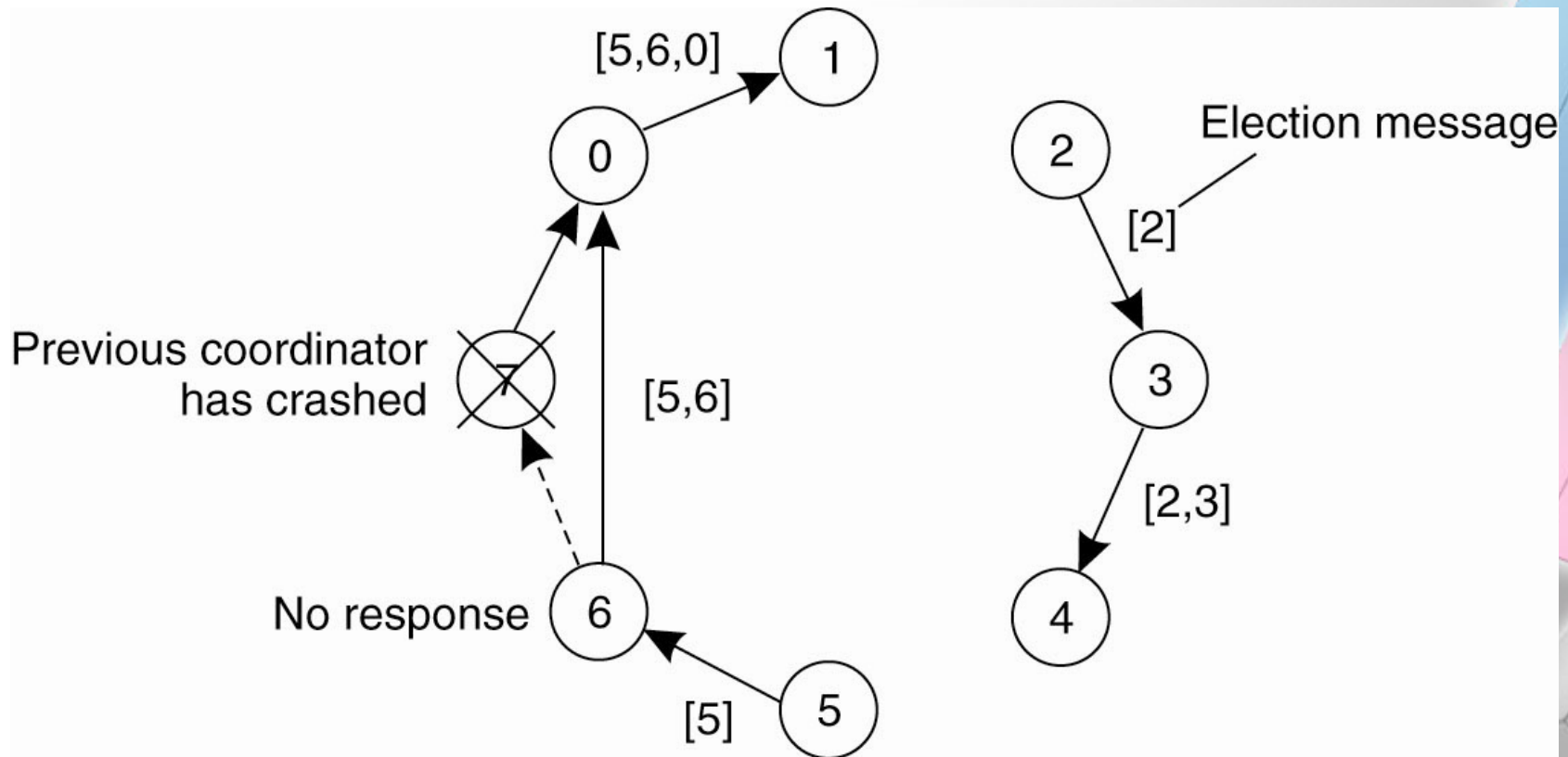
# The Ring algorithm

## Principle

Process priority is obtained by organizing processes into a (logical) ring. Process with the highest priority should be elected as coordinator.

- Any process can start an election by sending an election message to its successor. If a successor is down, the message is passed on to the next successor.
- If a message is passed on, the sender adds itself to the list. When it gets back to the initiator, everyone had a chance to make its presence known.
- The initiator sends a coordinator message around the ring containing a list of all living processes. **The one with the highest priority is elected as coordinator.**

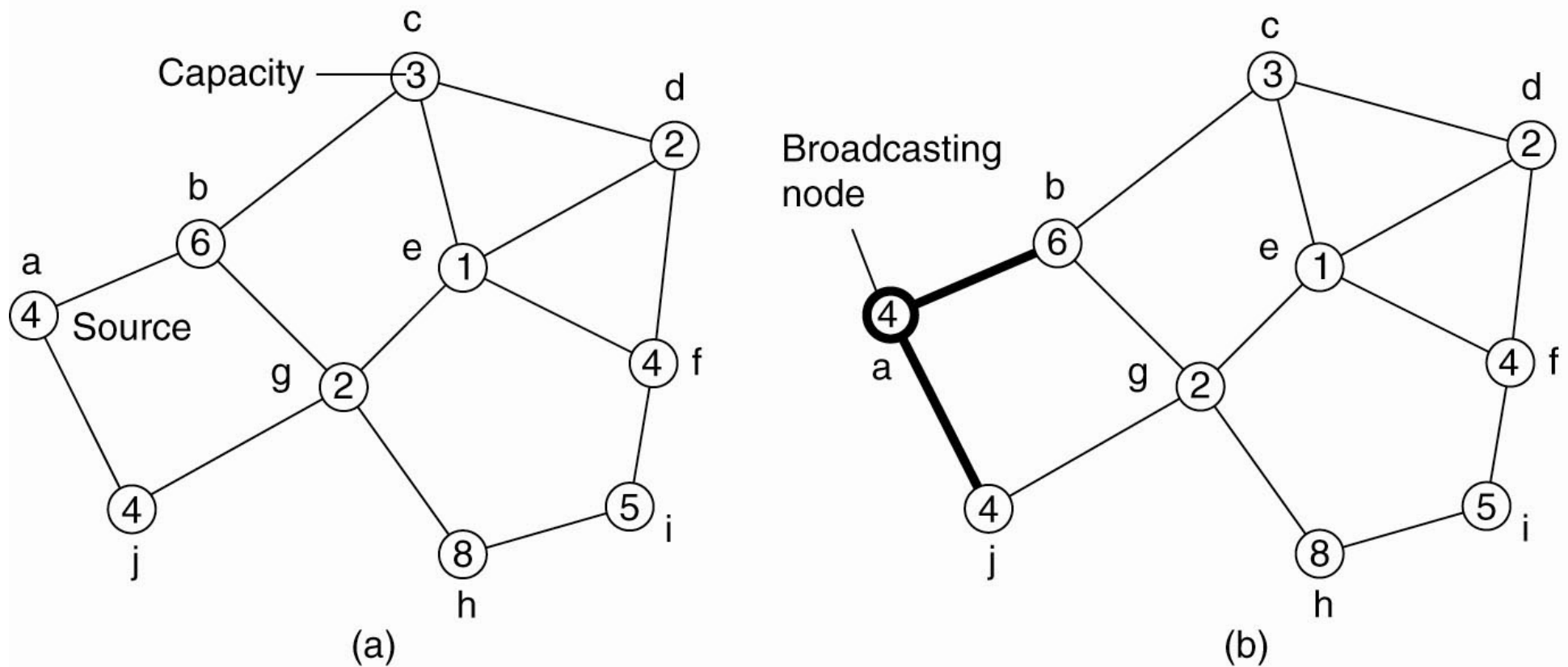
# The Ring algorithm



# Election in Wireless Networks

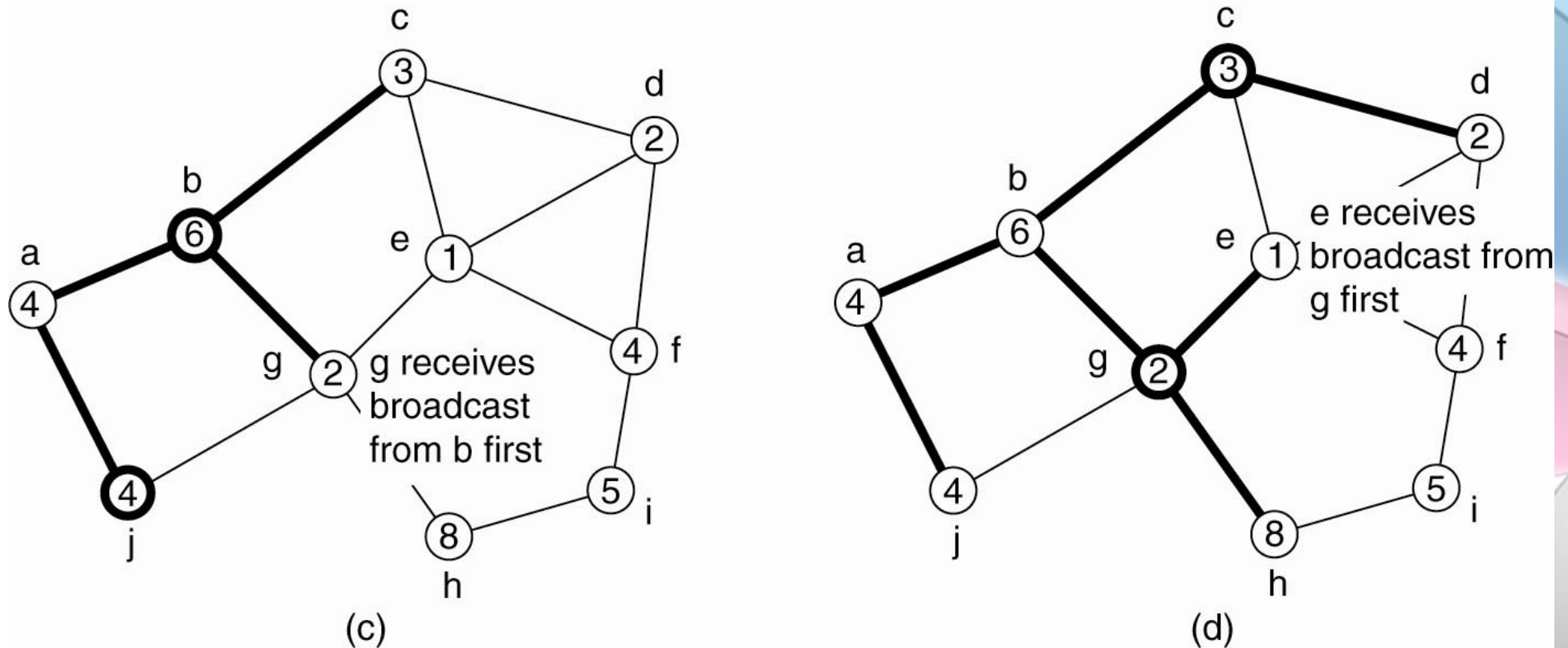
- Traditional election algorithms are generally based on assumptions that are not realistic in wireless environments.
  - For example, they assume that message passing is reliable and that the topology of the network does not change.
- These **assumptions are false** in most wireless environments, especially those for mobile ad hoc networks.
- Only few protocols for elections have been developed that work in ad hoc networks.
- An important property of some solutions is that the **best leader can be elected rather than just a random** as was more or less the case in the previously discussed solutions.

# Elections in Wireless Environments (1)



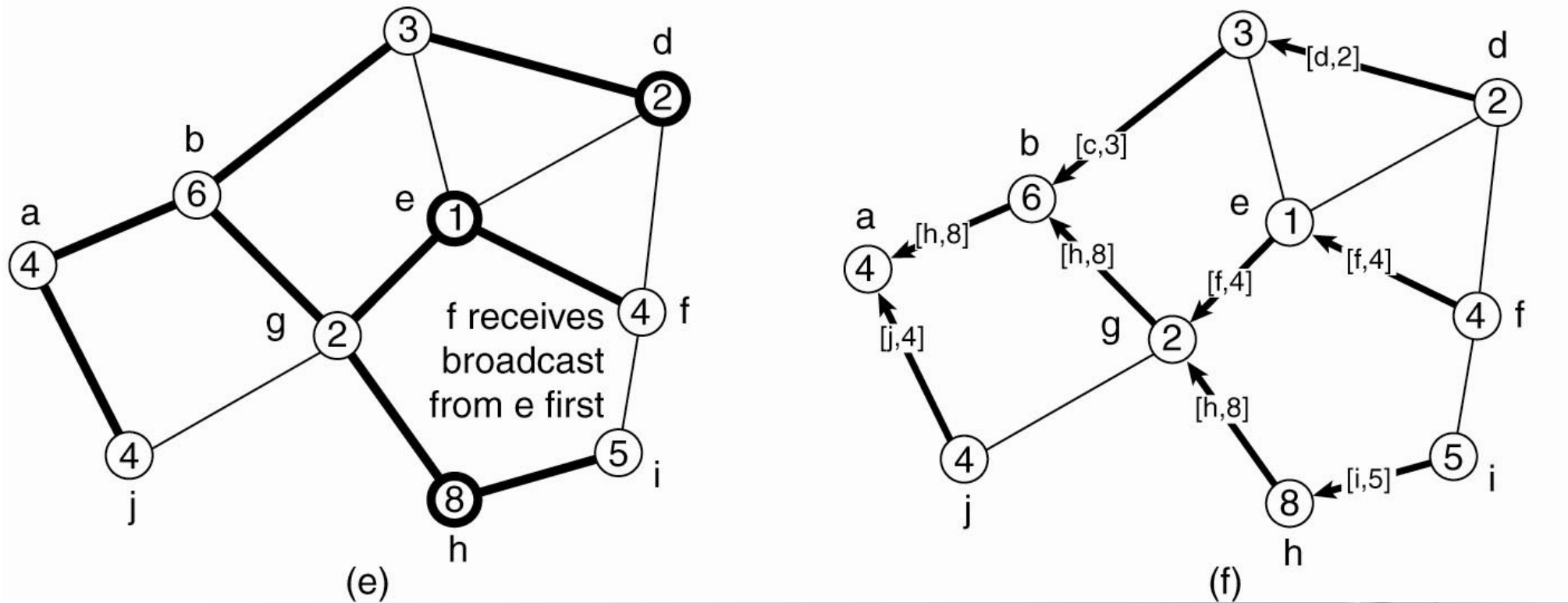
- Figure 6-22. Election algorithm in a wireless network, with node a as the source. (a) Initial network. (b)–(e) The build-tree phase

# Elections in Wireless Environments (2)



- Figure 6-22. Election algorithm in a wireless network, with node a as the source. (a) Initial network. (b)–(e) The build-tree phase

# Elections in Wireless Environments (3)



- Figure 6-22. (e) The build-tree phase.  
(f) Reporting of best node to source.

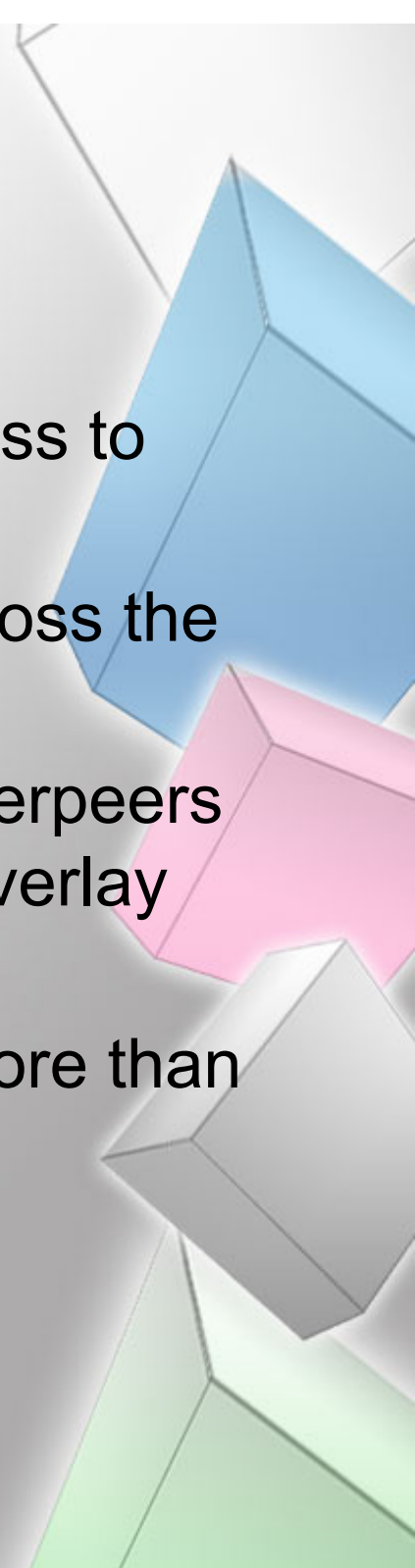
# Large Scale Systems

- The algorithms we have been discussing so far generally apply to **relatively small distributed systems**.
- Moreover, the algorithms concentrate on the selection of **only a single node**.
- There are situations when several nodes should actually be selected, such as in the case of **superpeers** in peer-to-peer networks.
- Problem: how do we select superpeers?

# Selecting Superpeers

## Requirements:

1. Normal nodes should have low-latency access to superpeers.
2. Superpeers should be evenly distributed across the overlay network.
3. There should be a predefined portion of superpeers relative to the total number of nodes in the overlay network.
4. Each superpeer should not need to serve more than a fixed number of normal nodes.



# Superpeer election

## DHTs

- Reserve a fixed part of the ID space for superpeers. Example if  $S$  superpeers are needed for a system that uses  $m$ -bit identifiers, simply reserve the  $k = \lceil \log_2 S \rceil$  leftmost bits for superpeers,
- With  $N$  nodes, we'll have, on average,  $2^{k-m}N$  superpeers.

## Routing to superpeer

- Send message for key  $p$  to node responsible for  $p$  AND  
11 1100 00

# Elections in Large-Scale Systems

- A completely different approach is based on positioning nodes in an  $m$ -dimensional geometric space as we discussed above.
- In this case, assume we need to place  $N$  superpeers *evenly* throughout the overlay.
- **The basic idea is simple: a total of  $N$  tokens are spread across  $N$  randomly-chosen nodes.**
  - No node can hold more than one token.
  - Each token represents a repelling force by which another token is inclined to move away.
  - The net effect is that if all tokens exert the same repulsion force, they will move away from each other and spread themselves evenly in the geometric space.

# End of Lesson 6

- Readings
  - Distributed Systems, Chapter 6.

