

ECE151 – Lecture 7

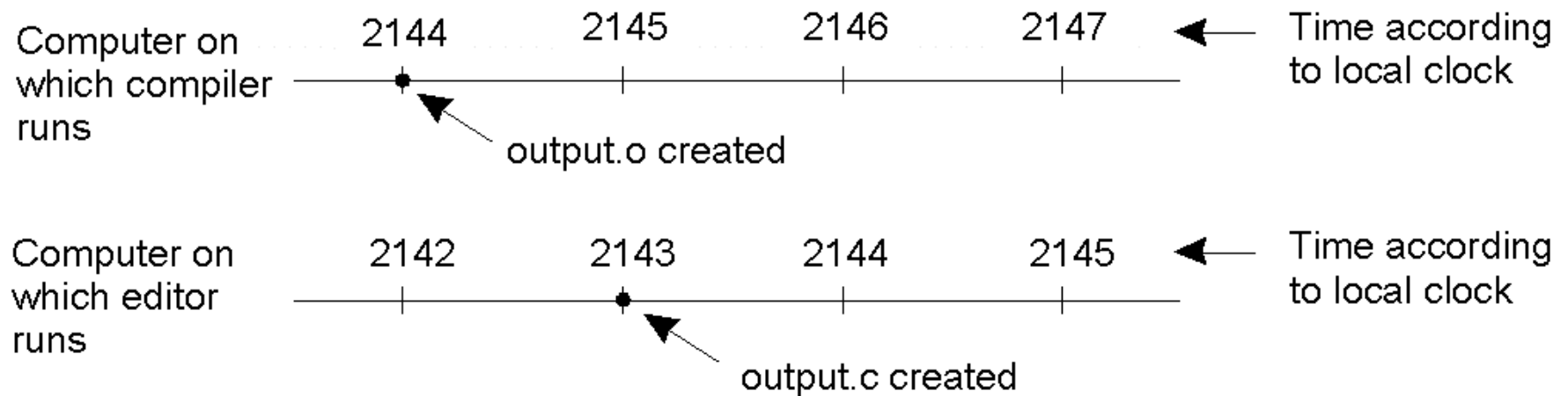
Chapter 5 Synchronization

Clock Synchronization

Physical clocks

Logical clocks

Vector clocks



When each machine has its own clock, an event that occurred after another event may nevertheless be assigned an earlier time.

Physical Clocks

Problem: Sometimes we simply need the exact time, not just an ordering.

Solution: Universal Coordinated Time (UTC):

Based on the number of transitions per second of the cesium 133 atom (pretty accurate).

At present, the real time is taken as the average of some 50 cesium-clocks around the world.

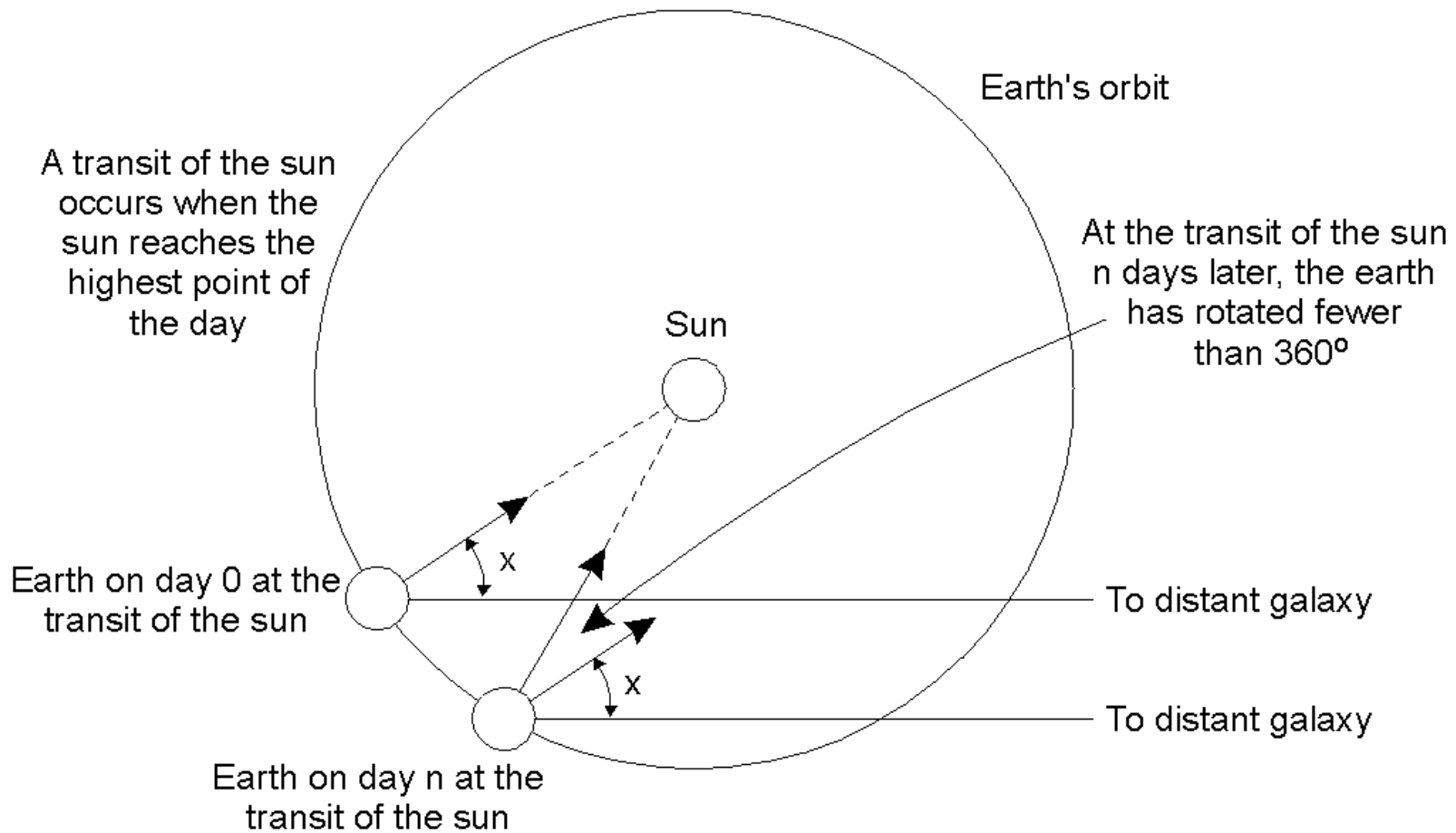
Introduces a leap second from time to time to compensate that days are getting longer.

UTC is **broadcast** through short wave radio and satellite.

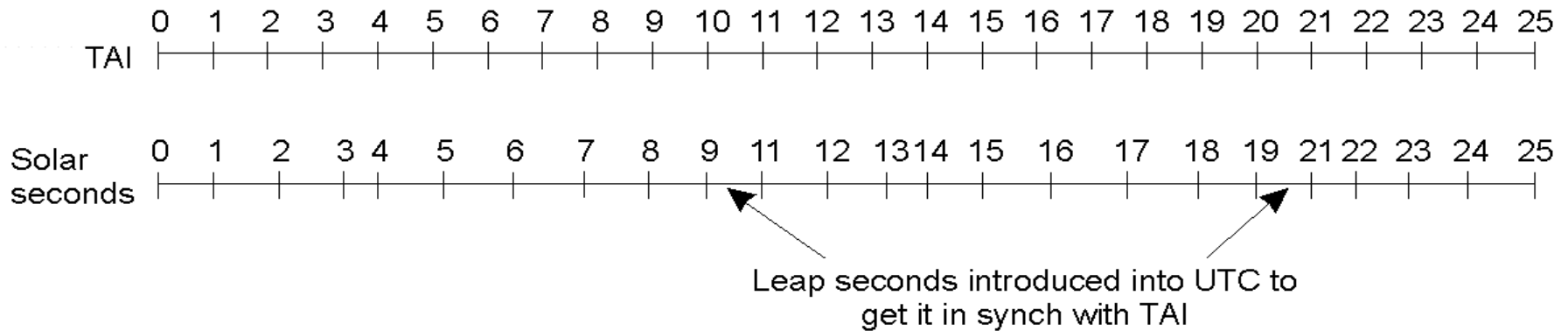
Satellites can give an accuracy of about 0.5 ms.

Question: Does this solve all our problems? Don't we now have some global timing mechanism?

Physical Clocks (1)



Physical Clocks (2)



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

Physical Clocks

Problem: Suppose we have a distributed system with a UTC-receiver somewhere in it we still have to distribute its time to each machine.

Basic principle:

Every machine has a timer that generates an interrupt H times per second.

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, $dC/dt = 1$

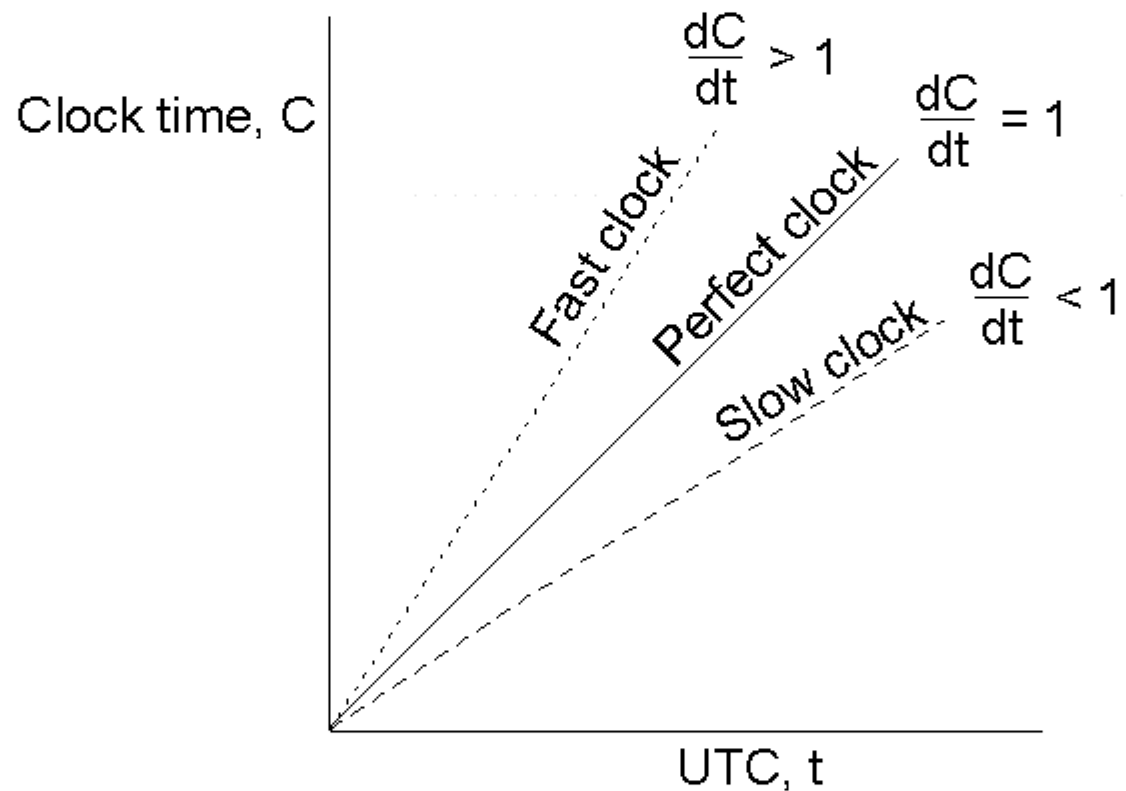
Clock Synchronization Algorithms

The relation between clock time and UTC when clocks tick at different rates.

In practice:

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

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



Clock Synchronization Principles

Principle I: Every machine asks a **time server** for the accurate time at least once every $\delta/(2\rho)$ seconds.

Okay, but you need an accurate measure of round trip delay, including interrupt handling and processing incoming messages.

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.

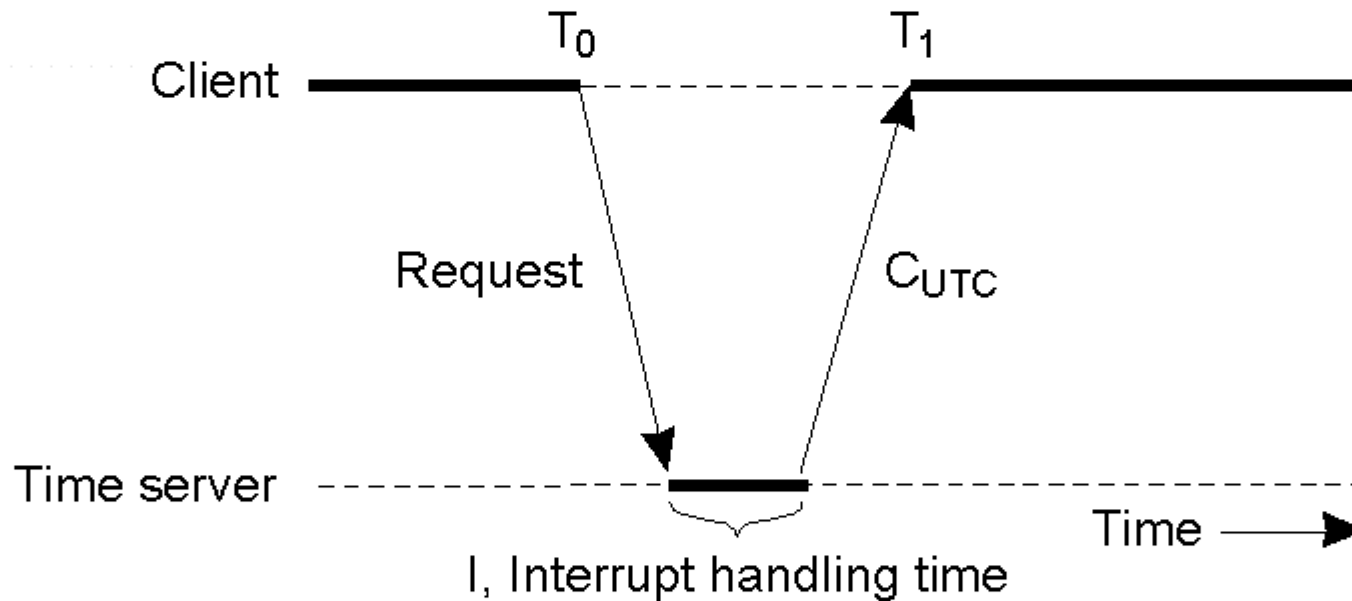
Okay, you'll probably get every machine in sync. You don't even need to propagate UTC time (why not?)

Fundamental problem: You'll have to ensure that setting time back is **never** allowed (smooth adjustments)

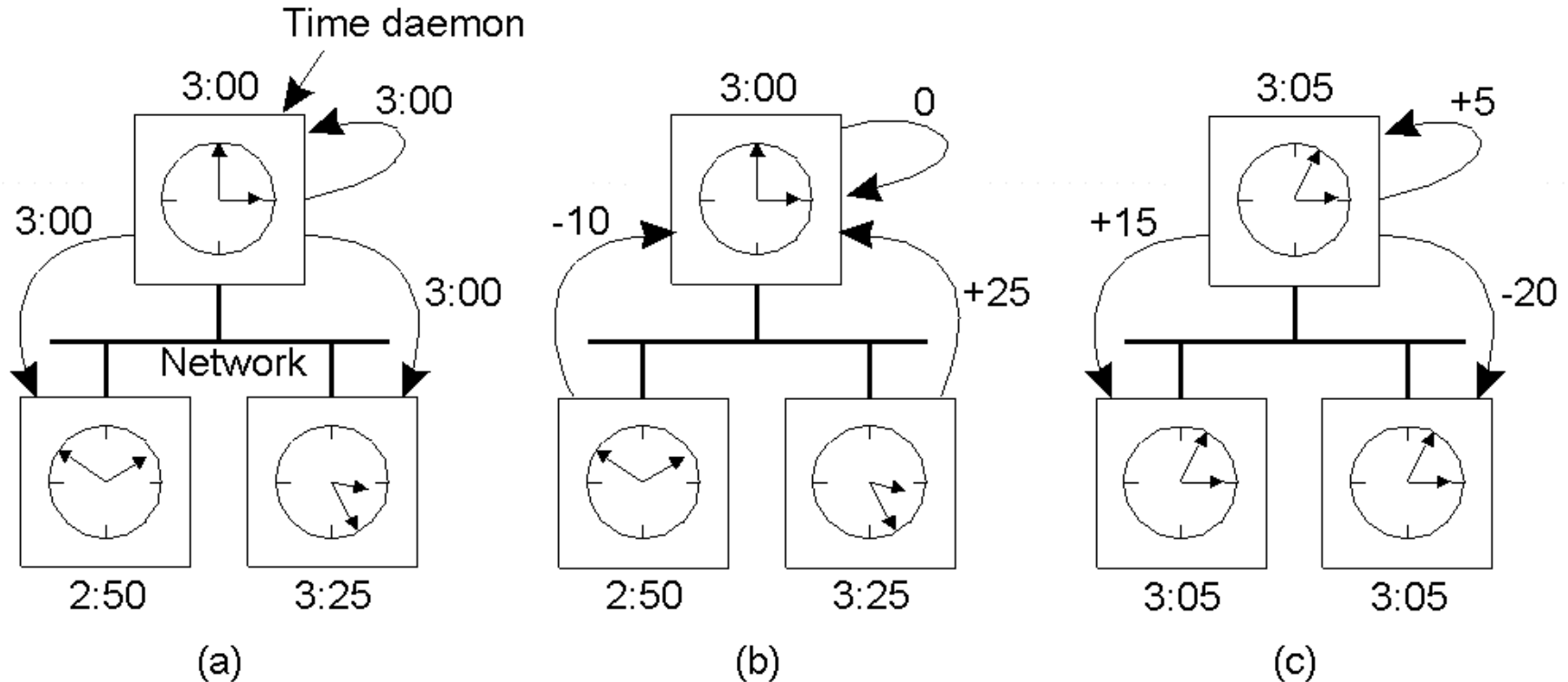
Cristian's Algorithm

Getting the current time from a time server.

Both T_0 and T_1 are measured with the same clock



The Berkeley Algorithm



- 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

The Happened-Before Relationship

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

The **happened-before** relation on the set of events in a distributed system is the smallest relation satisfying:

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 happenedbefore 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.

Logical Clocks

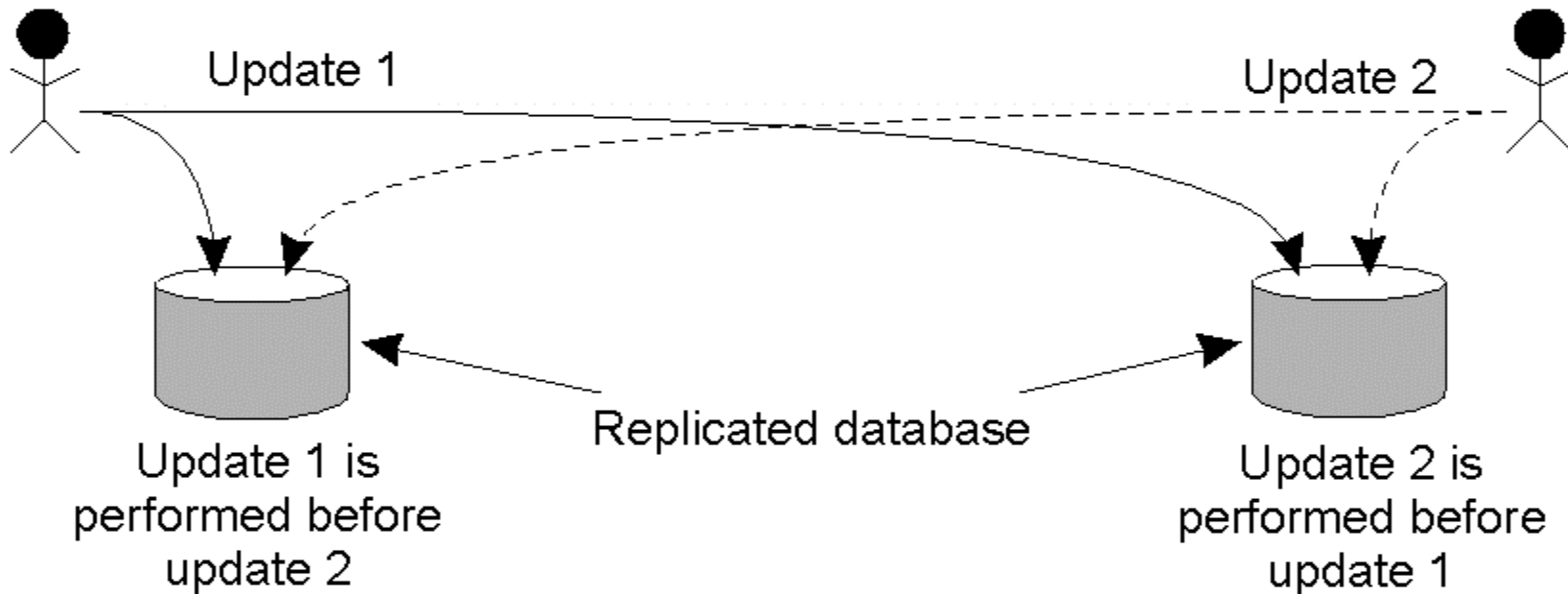
Each process P_i maintains a **local** counter C_i and adjusts this counter according to the following rules:

- 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 $T_m = C_i$.
- 3: Whenever a message m is received by a process P_j , P_j adjusts its local counter C_j :

$$C_j = \max(C_j + 1, T_m + 1)$$

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

Lamport Timestamps



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

Total Ordering with Logical Clocks

Problem: it can still occur that two events happen at the same time. Avoid this by attaching a process number to an event:

P_i timestamps event *e* with $C_i(e)$

Then: $C_i(a)$ before $C_j(b)$ if and only if:

1: $C_i(a) < C_j(a)$ or

2: $C_i(a) < C_j(b)$ and $i = j$

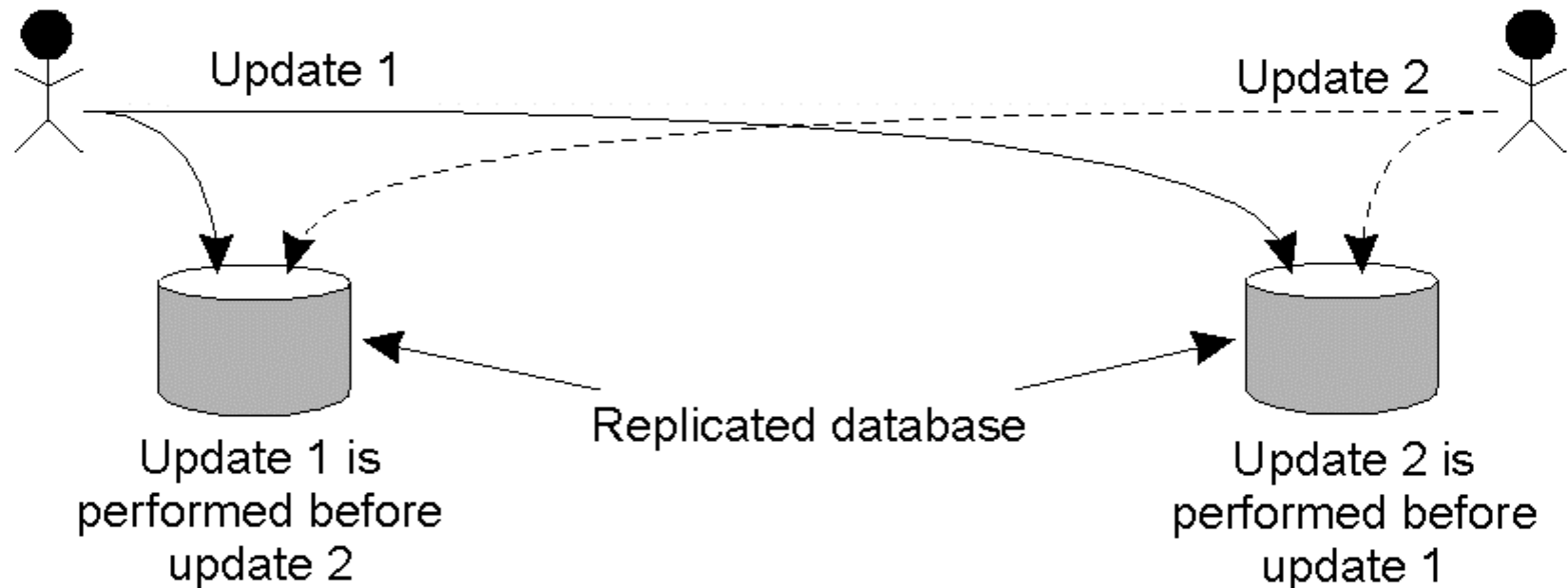
Example: Totally-Ordered Multicasting

Problem: We sometimes need to guarantee that concurrent updates on a replicated database are seen in the same order everywhere:

Process P_1 adds \$100 to an account (initial value: \$1000)

Process P_2 increments account by 1%

There are two replicas



Outcome: in absence of proper synchronization, replica #1 will end up with \$1111, while replica #2 ends up with \$1110.

Example: Totally-Ordered Multicast

Process P_i sends timestamped message msg_i to all others. The message itself is put in a local queue $queue_i$.

Any incoming message at P_j is queued in $queue_j$, according to its timestamp.

P_j passes a message msg_i to its application if:

- (1) msg_i is at the head of $queue_j$
- (2) for each process P_k , there is a message msg_k in $queue_j$ with a larger timestamp.

Note: We are assuming that communication is reliable and FIFO ordered.

Extension to Multicasting: Vector Timestamps

Observation: Lamport timestamps do not guarantee that if $C(a) < C(b)$ that a indeed happened before b .

We need **vector timestamps** for that.

Each process P_i has an array $V_i[1..n]$, where $V_i[j]$ denotes the number of events that process P_i knows have taken place at process P_j .

When P_i sends a message m , it adds 1 to $V_i[I]$, and sends V_i along with m as **vector timestamp** $vt(m)$. Result: upon arrival, each other process knows P_i 's timestamp.

Question: What does $V_i[j] = k$ mean in terms of messages sent and received?

Extension to Multicasting: Vector Timestamps

When a process P_j receives a message m from P_i with vector timestamp $vt(m)$, it

- (1) updates each $V_j[k]$ to $\max(V_j[k], V(m)[k])$, and
- (2) increments V_j j by 1. **NOTE:** Book is wrong!

To support causal delivery of messages, assume you increment your own component only when sending a message. Then, P_j postpones delivery of m until:

- $vt(m)[i] = V_j[i] + 1$.
- $vt(m)[k] \leq V_j[k]$ for $k \neq i$.

Example: Take $V_3 = [0, 2, 2]$, $vt(m) = [1, 3, 0]$.

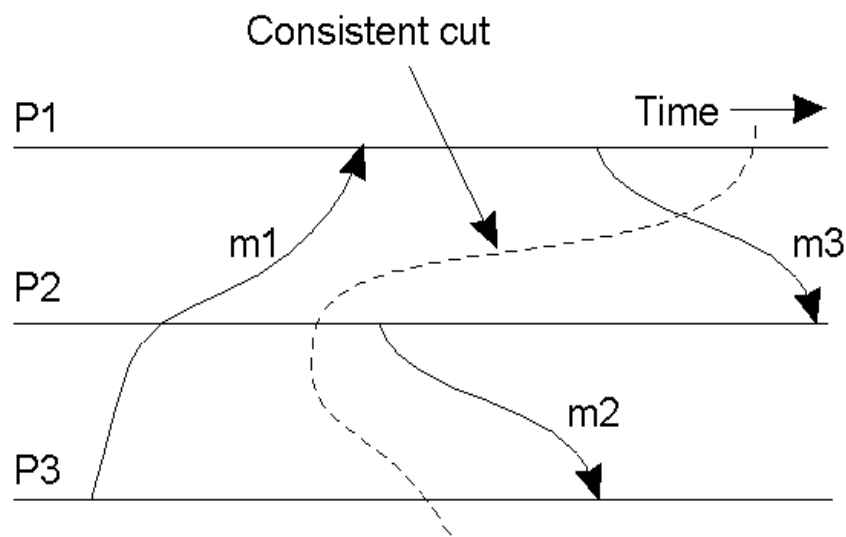
What information does P_3 have, and what will it do when receiving m (from P_1)?

Global State

Basic Idea: Sometimes you want to collect the current state of a distributed computation, called a **distributed snapshot**.

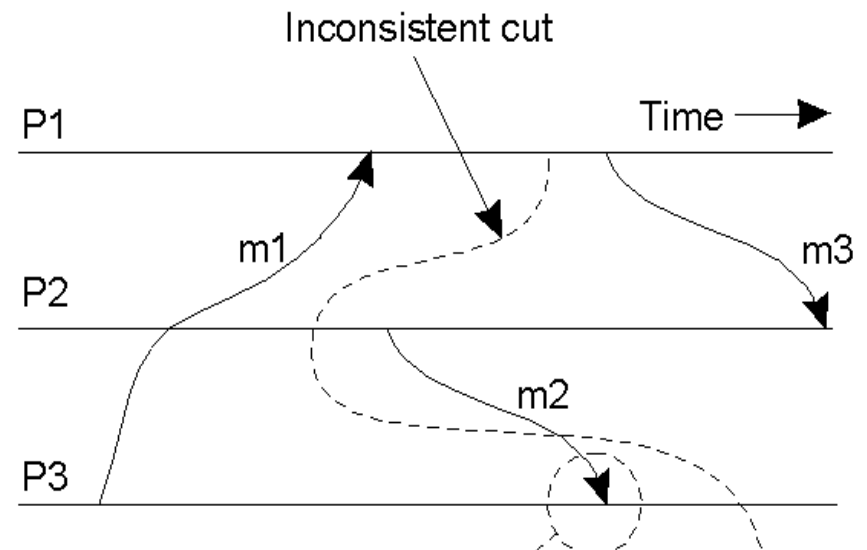
It consists of all local states and messages in transit.

Important: A distributed snapshot should reflect a **consistent** state:



- a) A consistent cut
- b) An inconsistent cut

(a)



Sender of m2 cannot be identified with this cut

(b)

Global State

Any process P can initiate taking a distributed snapshot

P starts by recording its own local state

P sends a marker along each of its outgoing channels

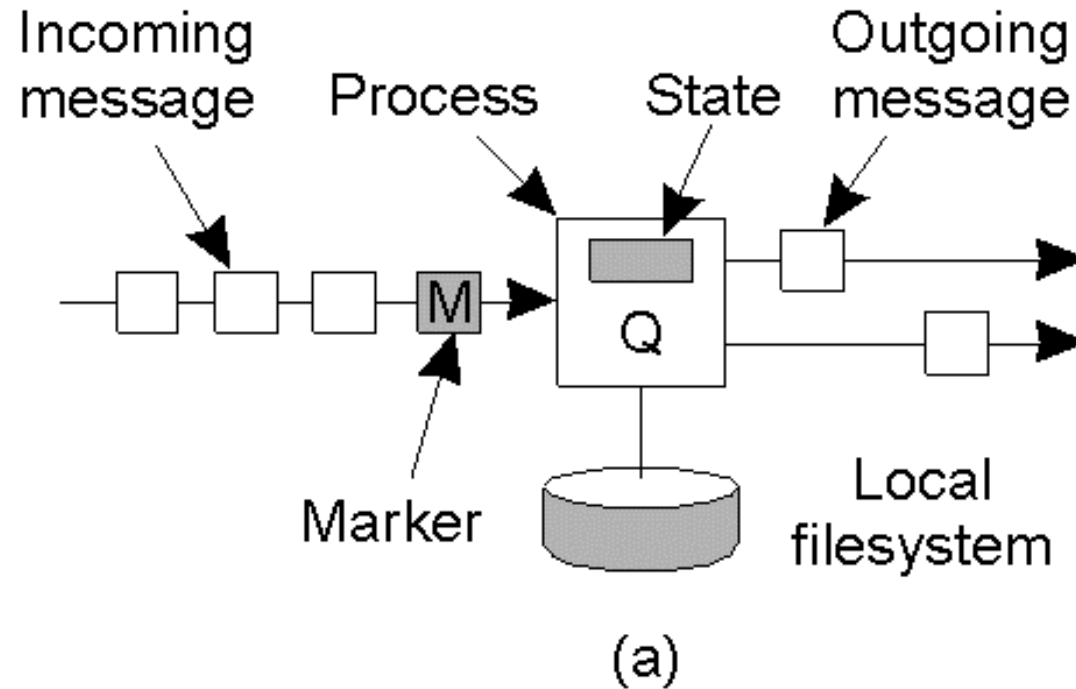
When Q receives a marker through channel C , its action depends on whether it had already recorded its local state:

- Not yet recorded: it records its local state, and sends the marker along each of its outgoing channels
- Already recorded: the marker on C indicates that the channel's state should be recorded:

All messages received before this marker and
after Q recorded its own state.

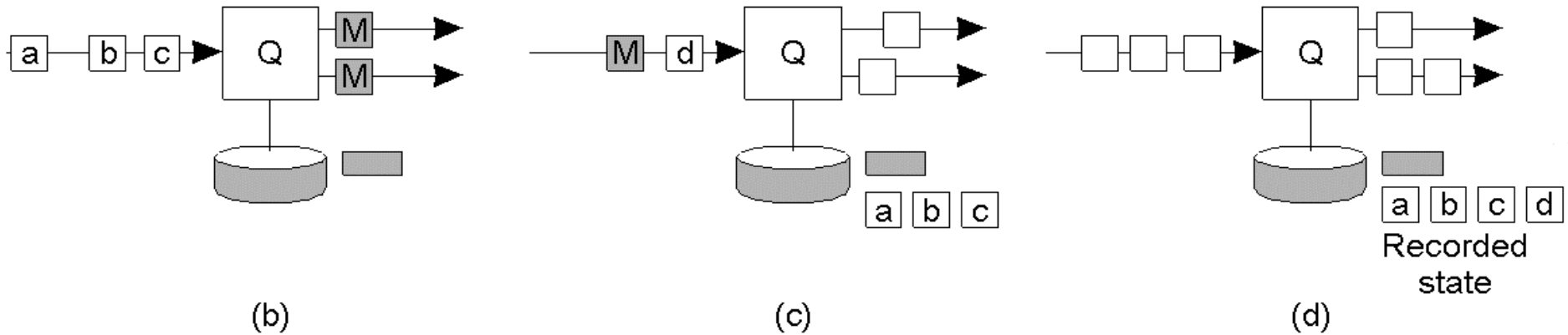
Q is finished when it has received a marker along each of its incoming channels

Global State



- a) Organization of a process and channels for a distributed snapshot

Global State



- b)** Process Q receives a marker for the first time and records its local state
- c)** Q records all incoming messages
- d)** Q receives a marker for its incoming channel and finishes recording the state of the incoming channel

Election Algorithms

Principle: An algorithm requires that some process acts as a coordinator. The question is how to select this special process **dynamically**.

Note: In many systems the coordinator is chosen by hand (e.g. file servers). This leads to centralized solutions with a single point of failure.

Question: If a coordinator is chosen dynamically, to what extent can we speak about a centralized or distributed solution?

Question: Is a fully distributed solution, i.e. one without a coordinator, always more robust than any centralized/coordinated solution?

Election by Bullying

Principle: Each process has an associated priority (weight). The process with the highest priority should always be elected as the coordinator.

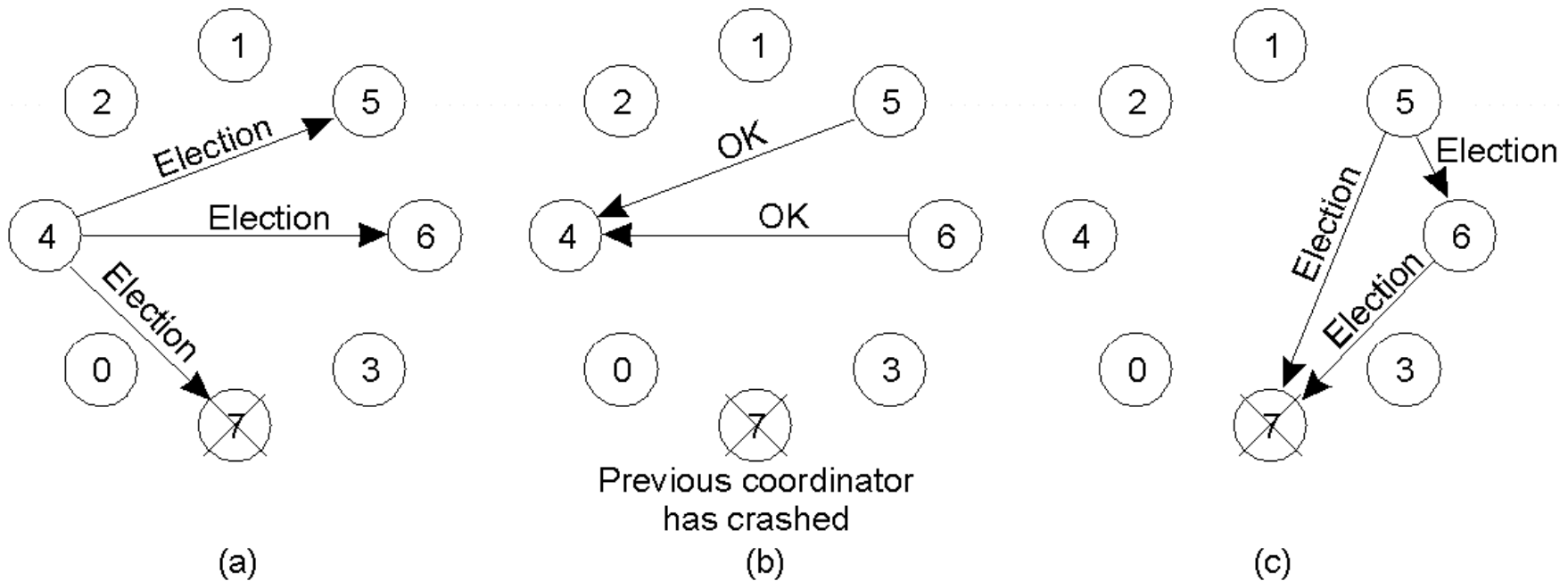
Issue: How do we find the heaviest process?

Any process can just start an election by sending an election message to all other processes (assuming you don't know the weights of the others).

If a process P_{heavy} receives an election message from a lighter process P_{light} , it sends a take-over message to P_{light} . P_{light} is out of the race.

If a process doesn't get a take-over message back, it wins, and sends a victory message to all other processes.

The Bully Algorithm

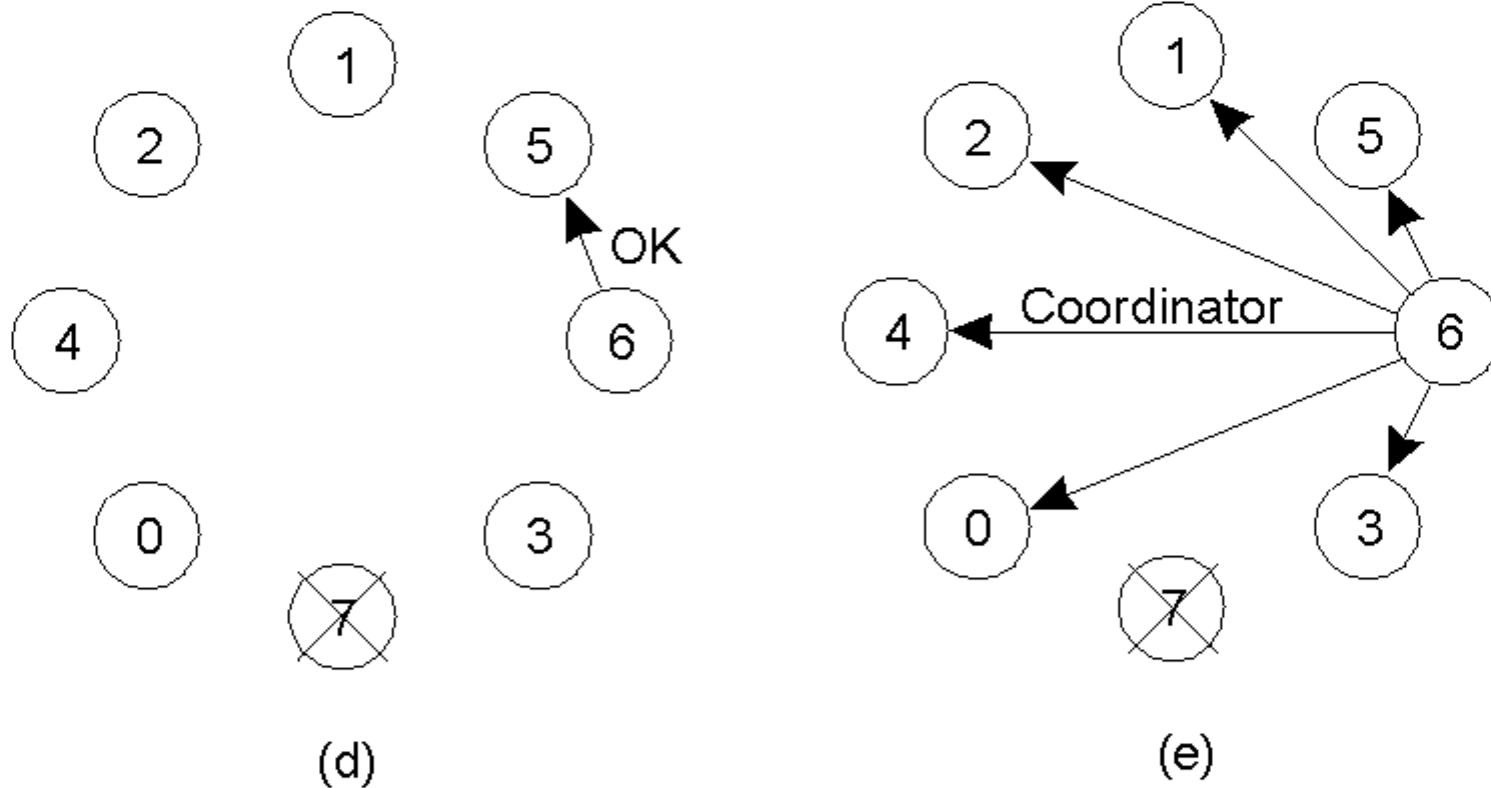


The bully election algorithm

- a) Process 4 holds an election
- b) Process 5 and 6 respond, telling 4 to stop
- c) Now 5 and 6 each hold an election

Global State

- d) Process 6 tells 5 to stop
- e) Process 6 wins and tells everyone



Election in a Ring

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.

Does it matter if two processes initiate an election?

What happens if a process crashes *during* the election?

A Ring Algorithm

