

# DISTRIBUTED SYSTEMS

## Principles and Paradigms

Second Edition

ANDREW S. TANENBAUM

MAARTEN VAN STEEN

# Chapter 6

## Synchronization

# Clock Synchronization

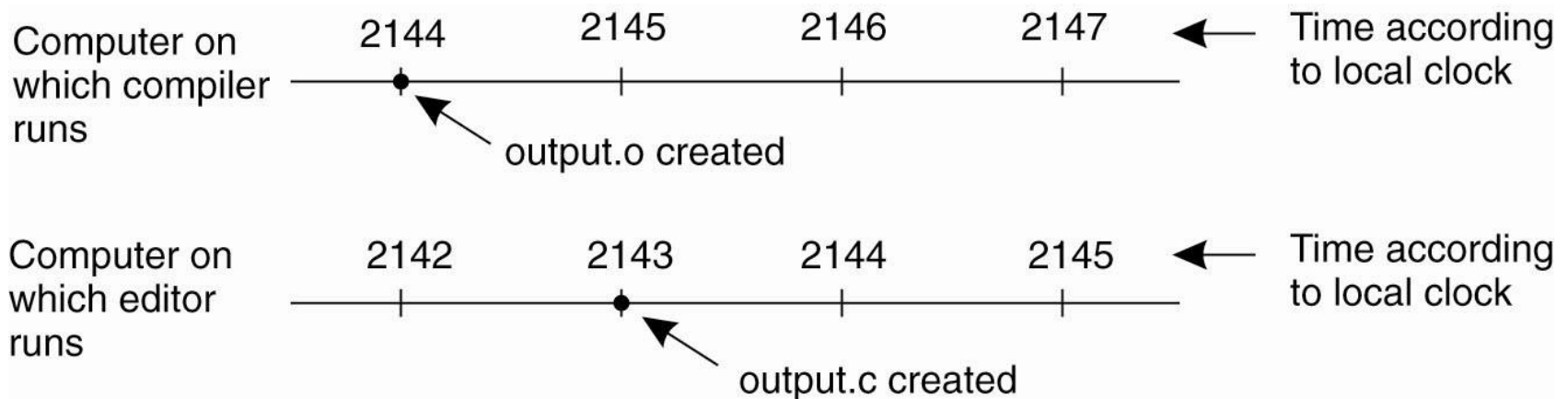


Figure 6-1. When each machine has its own clock, an event that occurred after another event may nevertheless be assigned an earlier time.

# Clock Synchronization Algorithms

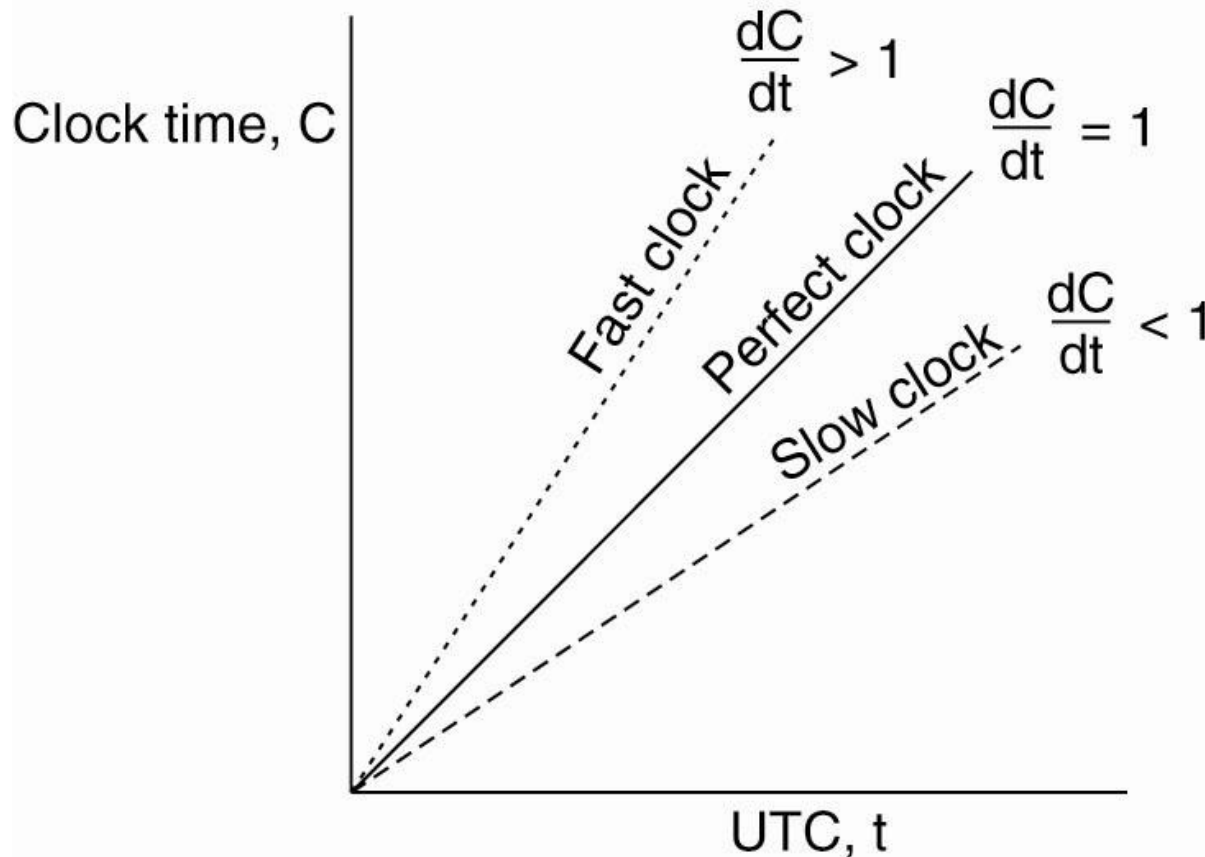


Figure 6-5. The relation between clock time and UTC when clocks tick at different rates.

# Network Time Protocol

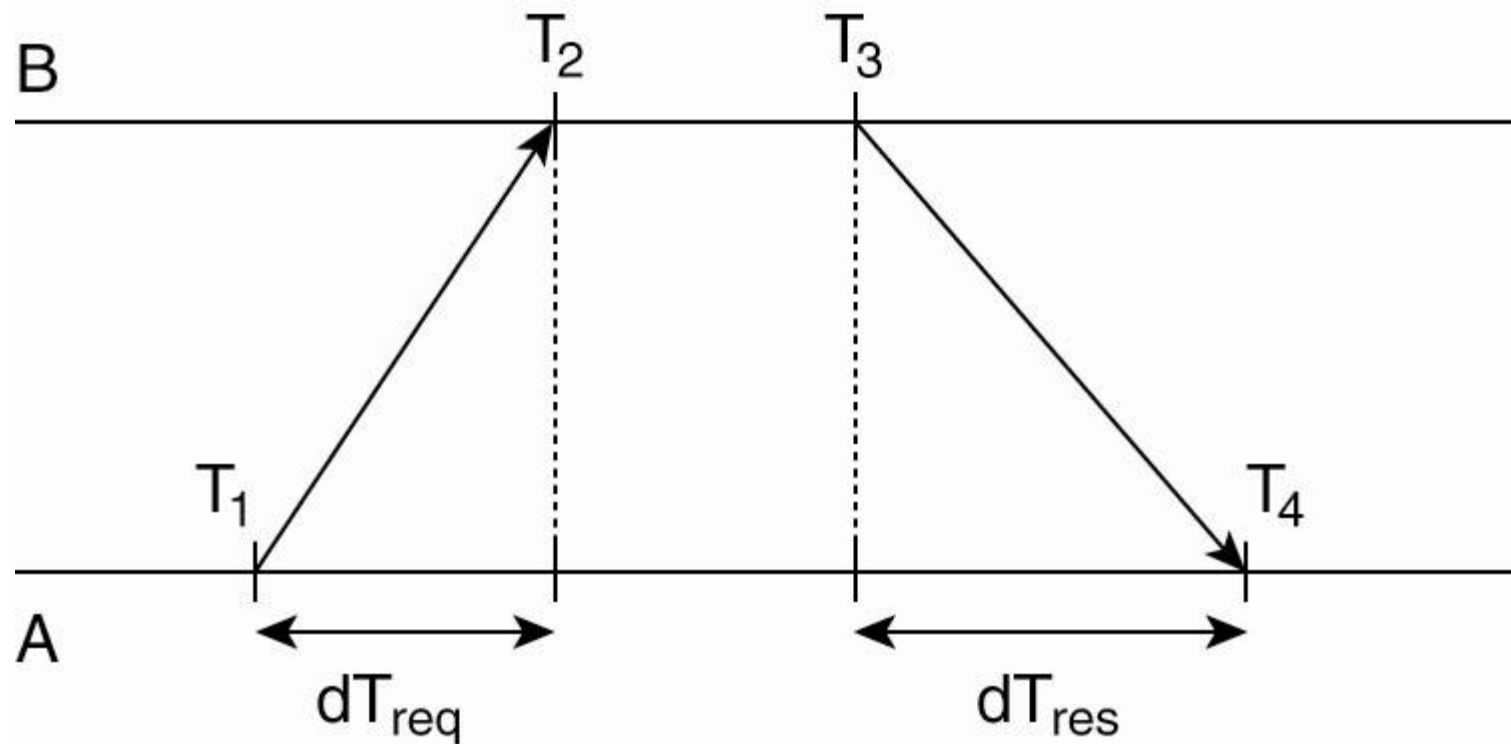
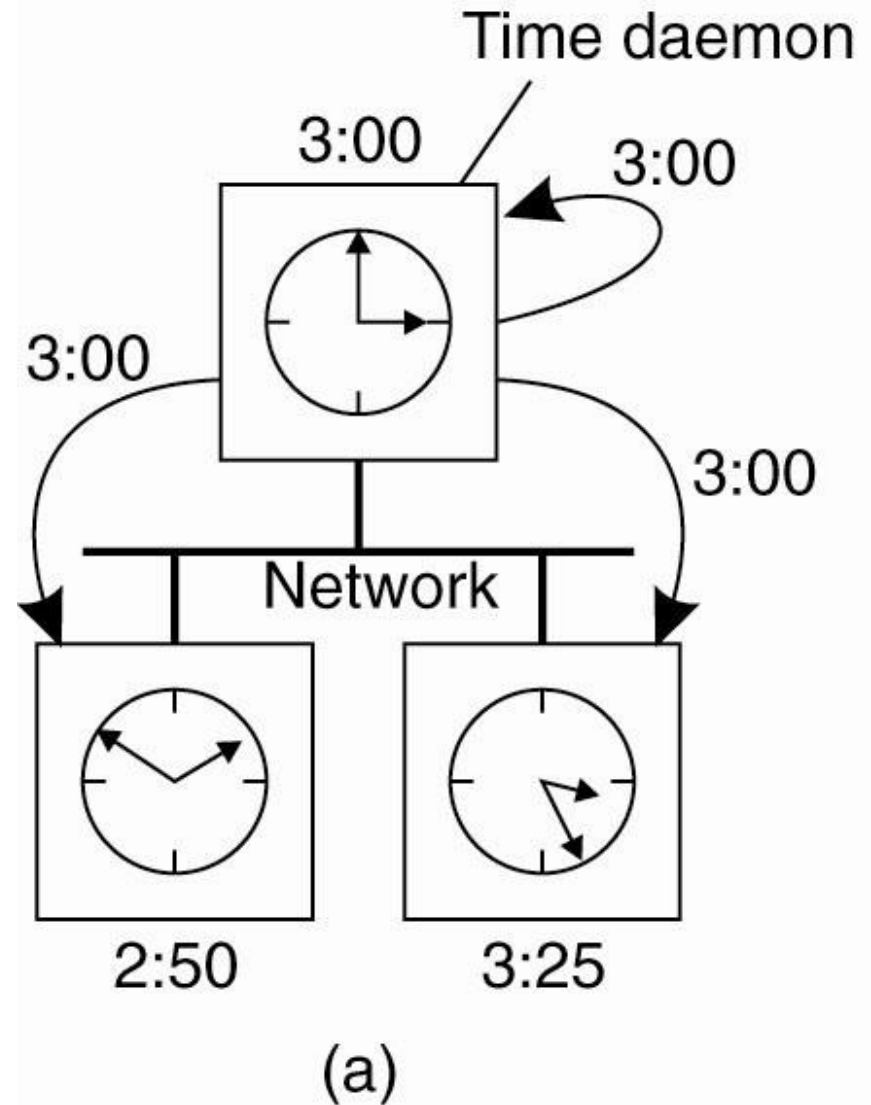


Figure 6-6. Getting the current time from a time server.

# The Berkeley Algorithm (1)

Figure 6-7. (a) The time daemon asks all the other machines for their clock values.



# The Berkeley Algorithm (2)

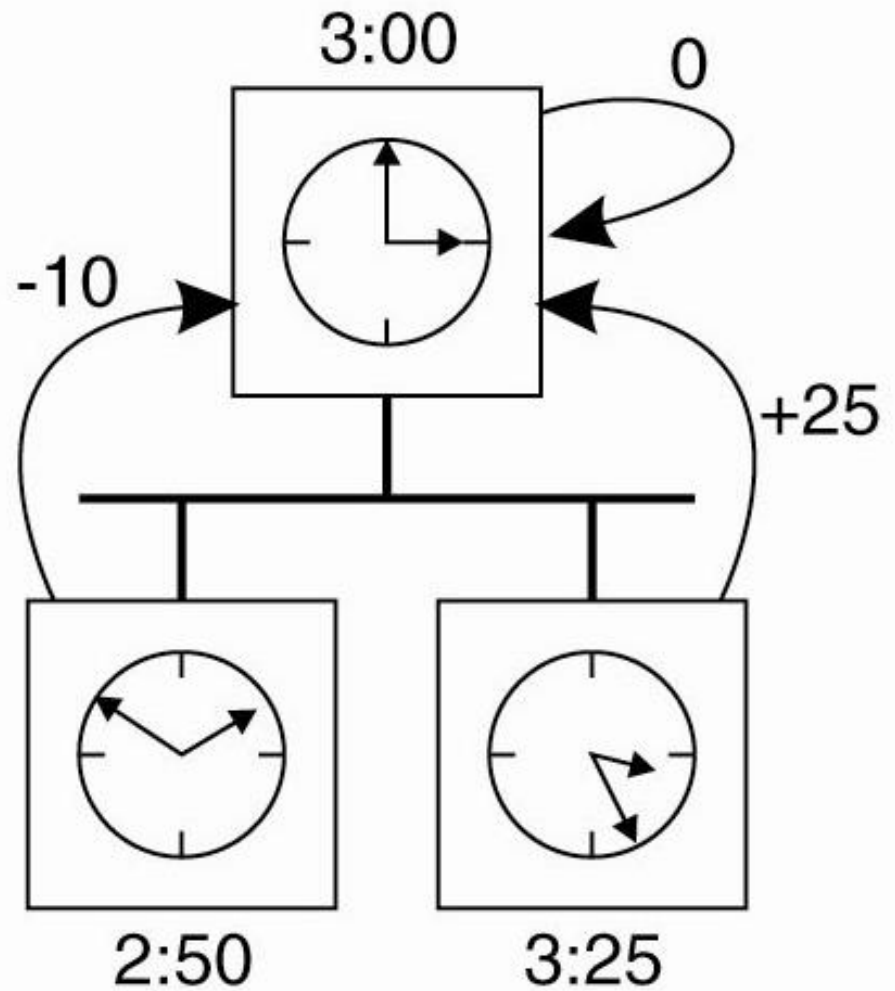
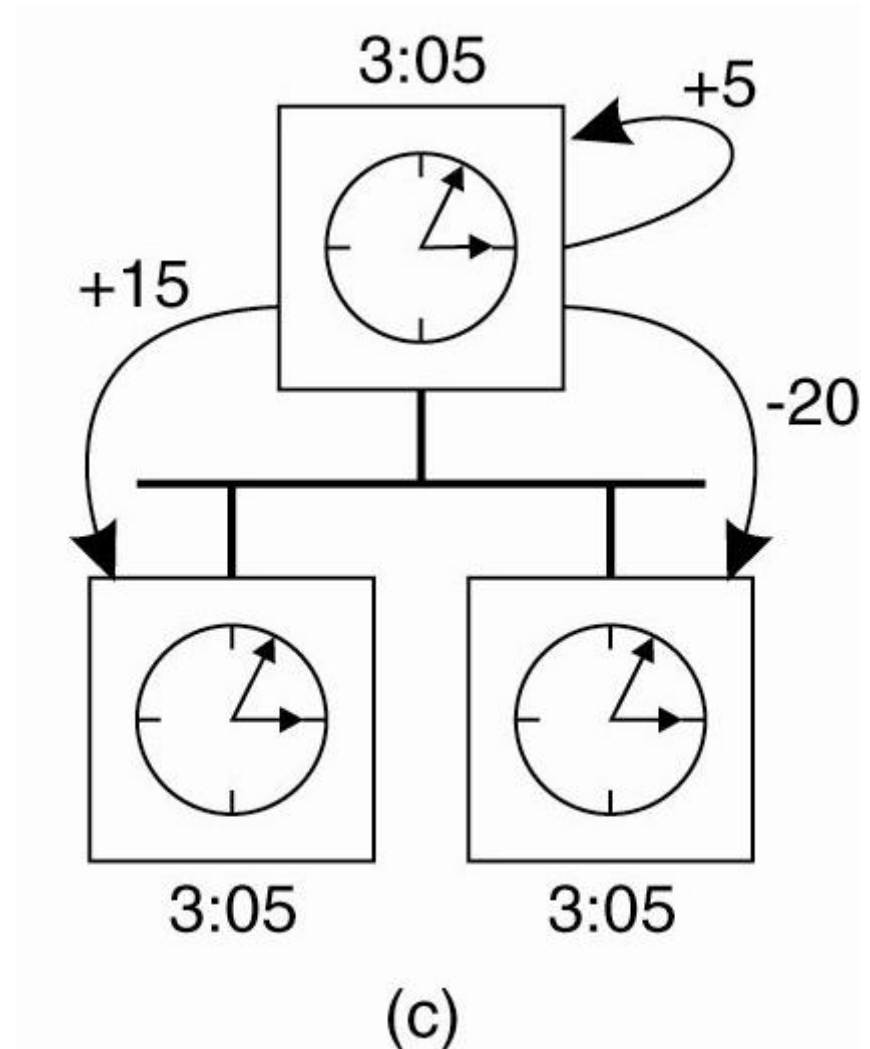


Figure 6-7.  
(b) The machines answer.

(b)

# The Berkeley Algorithm (3)

Figure 6-7. (c) The time daemon tells everyone how to adjust their clock.



# Clock Synchronization in Wireless Networks (1)

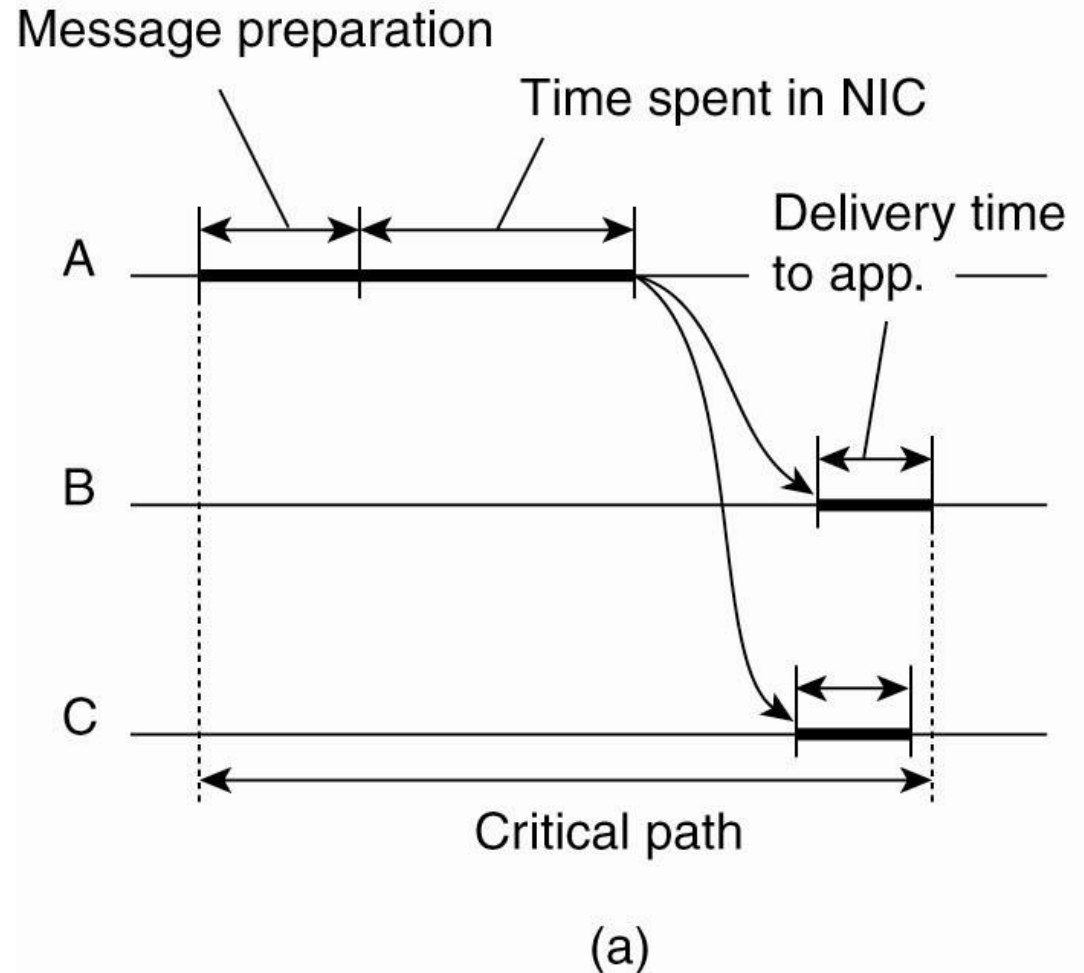


Figure 6-8. (a) The usual critical path in determining network delays.



# Clock Synchronization in Wireless Networks (2)

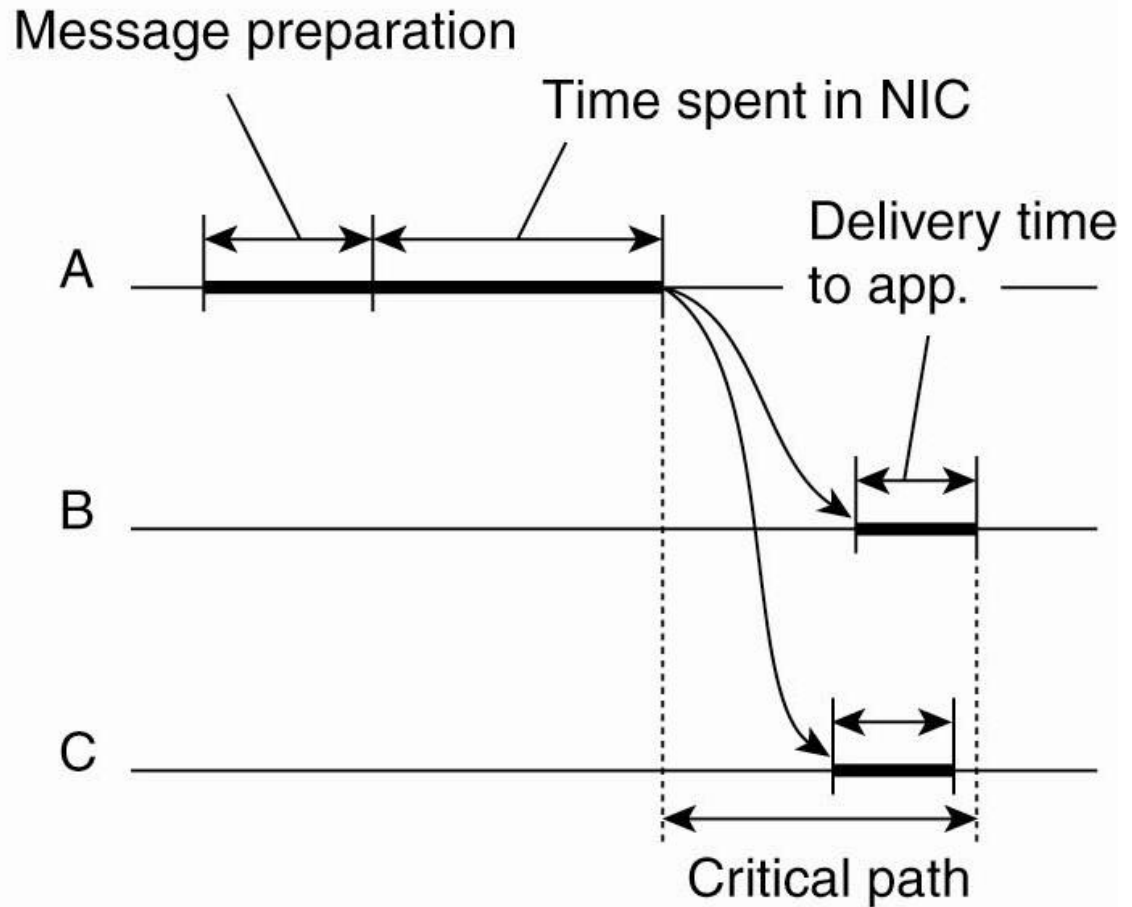


Figure 6-8. (b) The critical path in the case of RBS.

(b)

# Lamport's Logical Clocks (1)

The "happens-before" relation  $\rightarrow$  can be observed directly in two situations:

- If  $a$  and  $b$  are events in the same process, and  $a$  occurs before  $b$ , then  $a \rightarrow b$  is true.
- If  $a$  is the event of a message being sent by one process, and  $b$  is the event of the message being received by another process, then  $a \rightarrow b$

# Lamport's Logical Clocks (2)

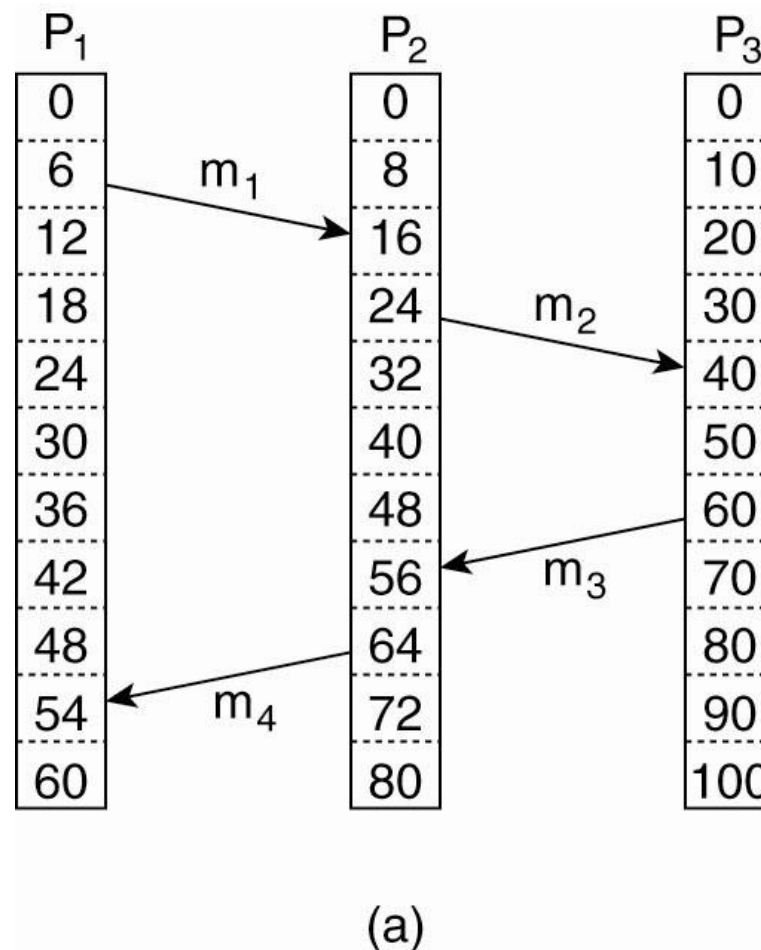


Figure 6-9. (a) Three processes, each with its own clock. The clocks run at different rates.

# Lamport's Logical Clocks (3)

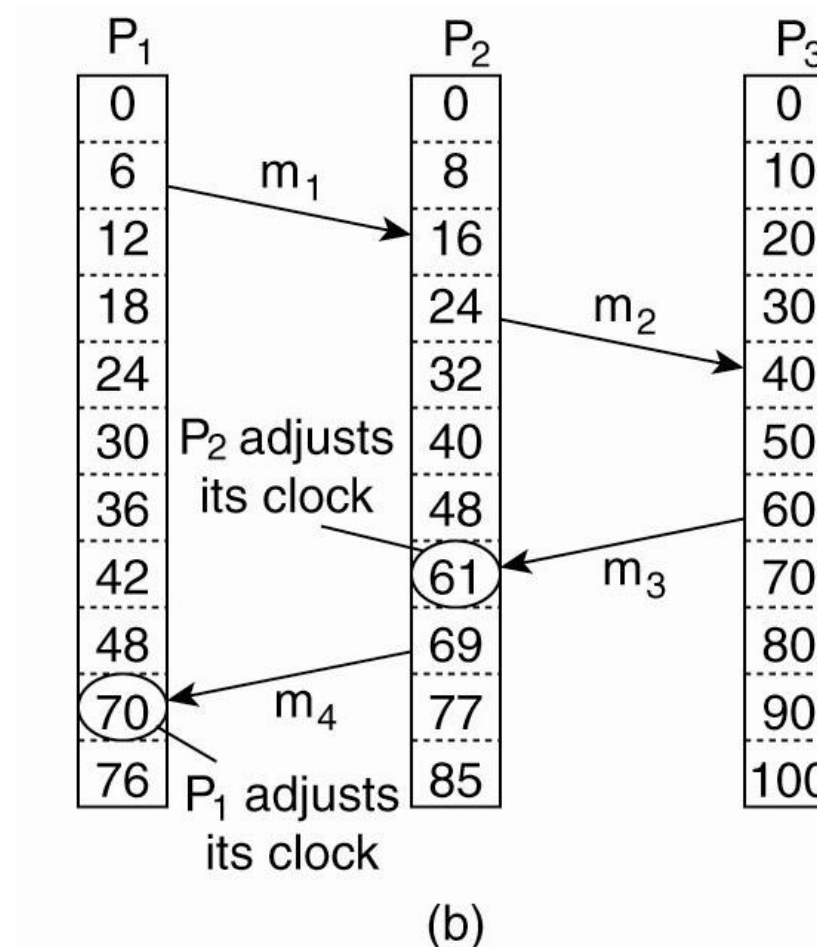


Figure 6-9. (b) Lamport's algorithm corrects the clocks.

# Lamport's Logical Clocks (4)

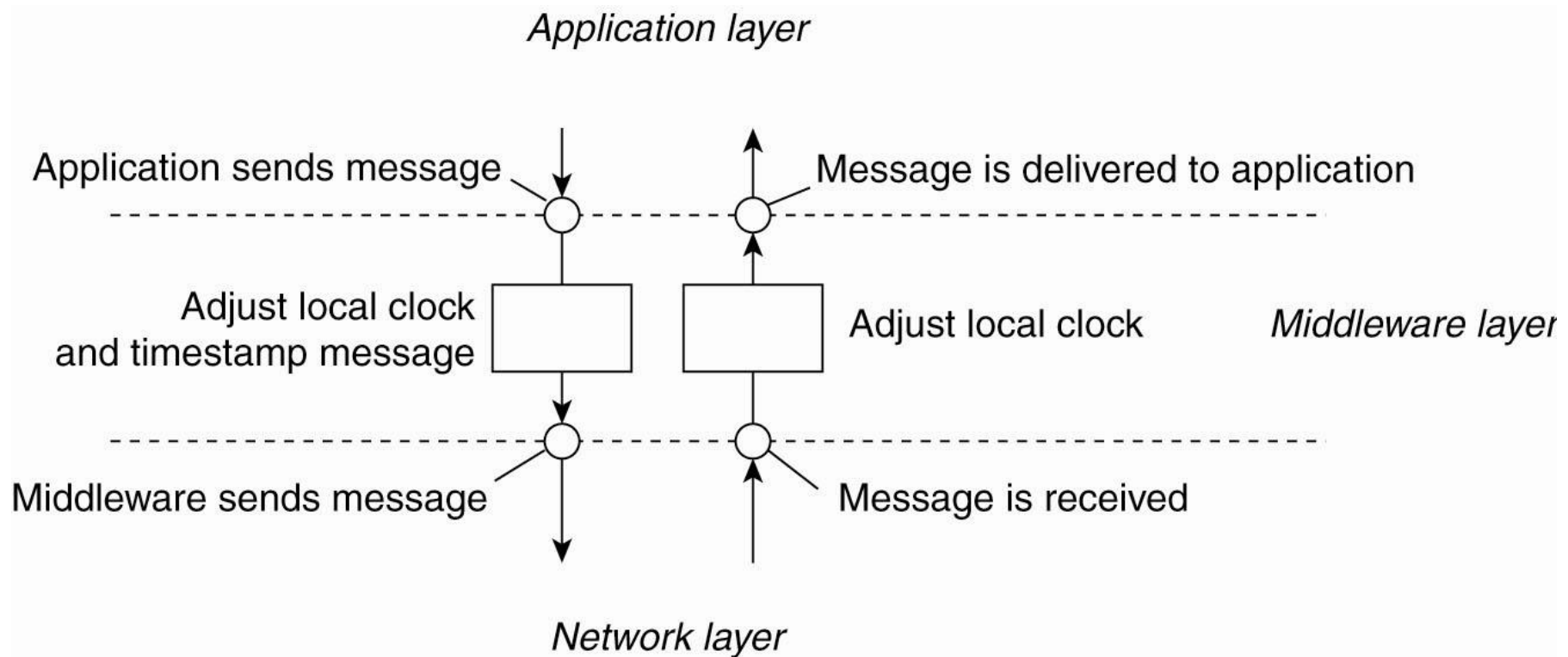


Figure 6-10. The positioning of Lamport's logical clocks in distributed systems.

# Lamport's Logical Clocks (5)

Updating counter  $C_i$  for process  $P_i$

1. Before executing an event  $P_i$  executes  $C_i \leftarrow C_i + 1$ .
2. When process  $P_i$  sends a message  $m$  to  $P_j$ , it sets  $m$ 's timestamp  $ts(m)$  equal to  $C_i$  after having executed the previous step.
3. Upon the receipt of a message  $m$ , process  $P_j$  adjusts its own local counter as  $C_j \leftarrow \max\{C_j, ts(m)\}$ , after which it then executes the first step and delivers the message to the application.

# Example: Totally Ordered Multicasting

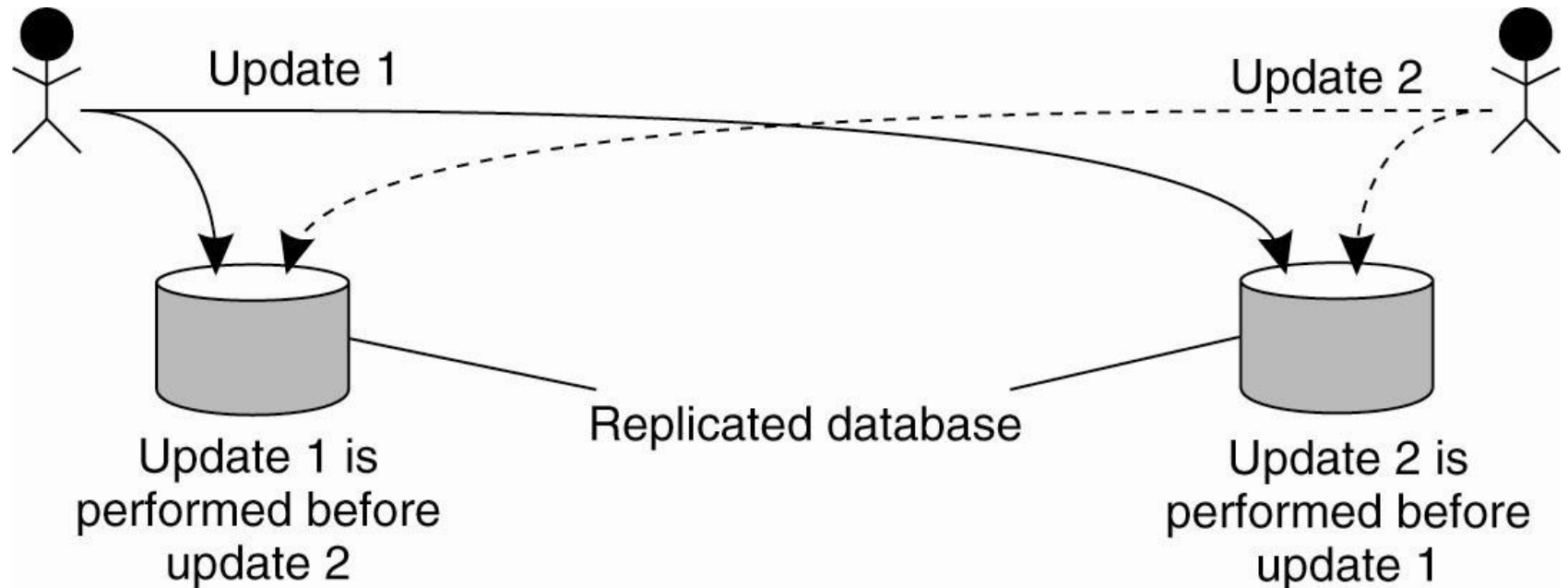


Figure 6-11. Updating a replicated database and leaving it in an inconsistent state.

# Vector Clocks (1)

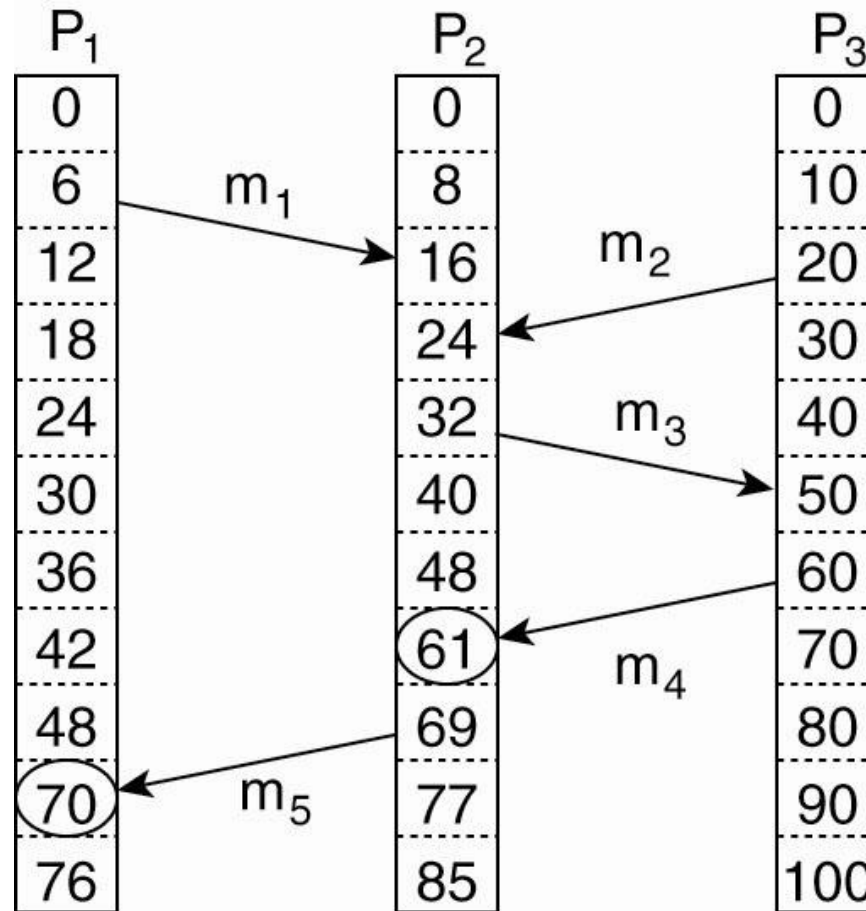


Figure 6-12. Concurrent message transmission using logical clocks.



# Vector Clocks (2)

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 (3)

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

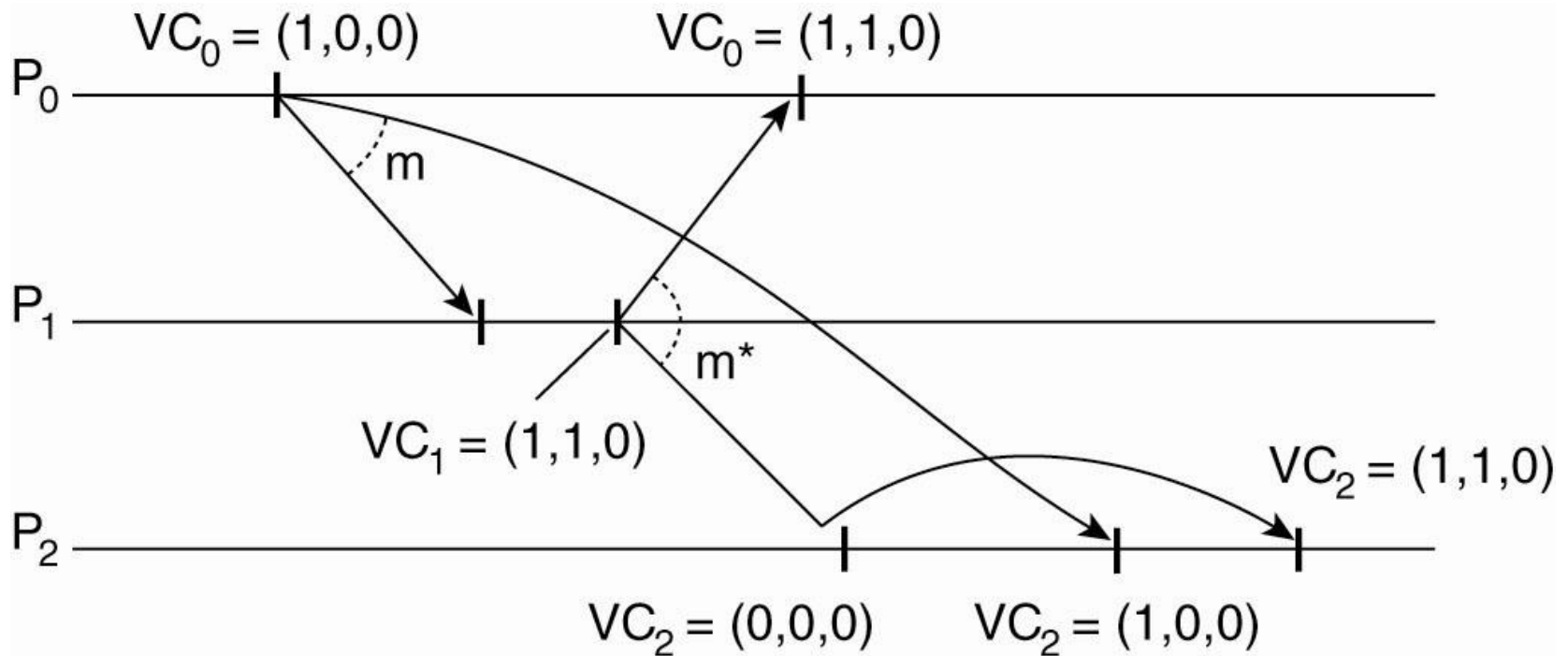


Figure 6-13. Enforcing causal communication.

# Mutual Exclusion

## A Centralized Algorithm (1)

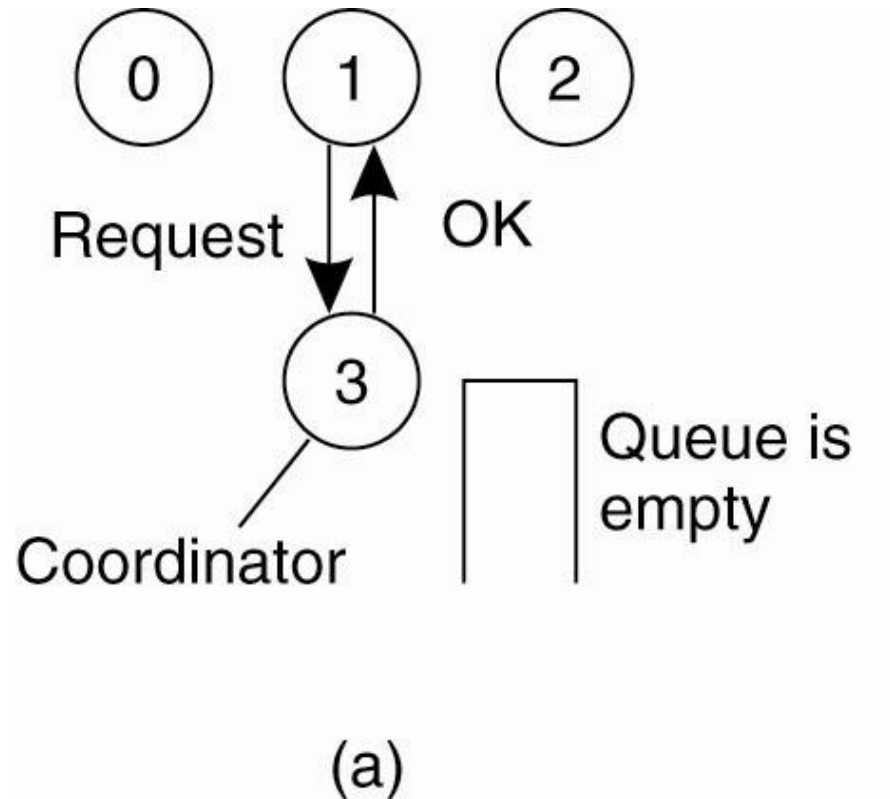


Figure 6-14. (a) Process 1 asks the coordinator for permission to access a shared resource. Permission is granted.

# Mutual Exclusion

## A Centralized Algorithm (2)

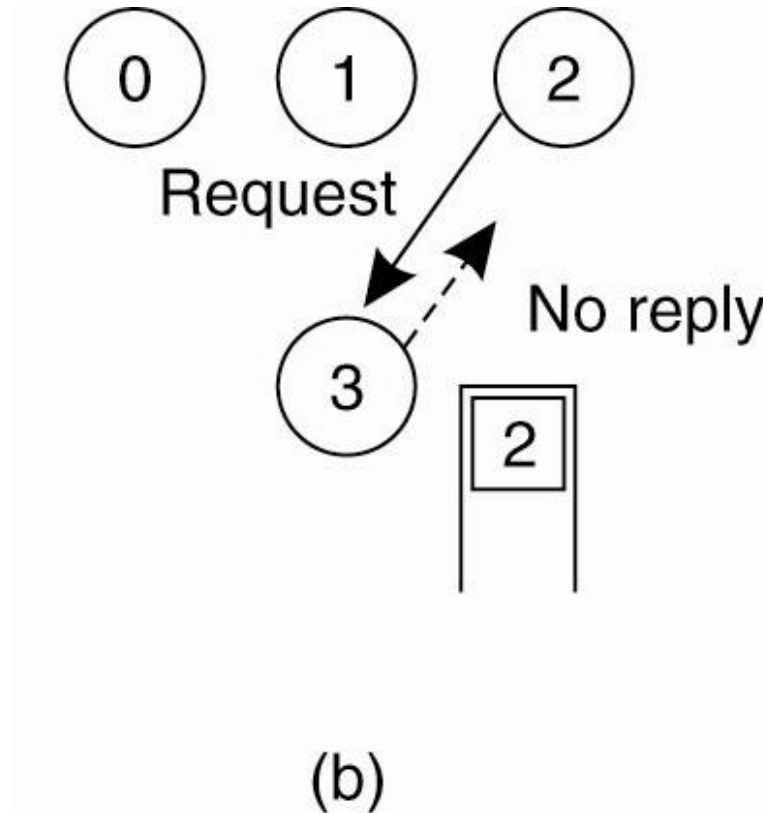
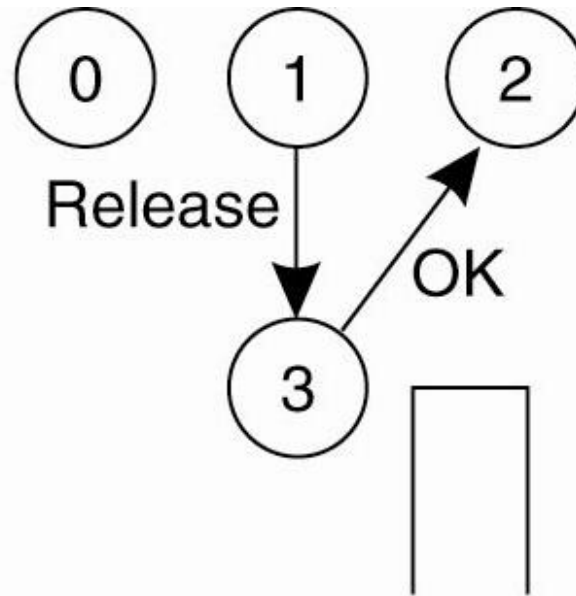


Figure 6-14. (b) Process 2 then asks permission to access the same resource. The coordinator does not reply.

# Mutual Exclusion

## A Centralized Algorithm (3)



(c)

Figure 6-14. (c) When process 1 releases the resource, it tells the coordinator, which then replies to 2.

# A Distributed Algorithm (1)

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.

# A Distributed Algorithm (2)

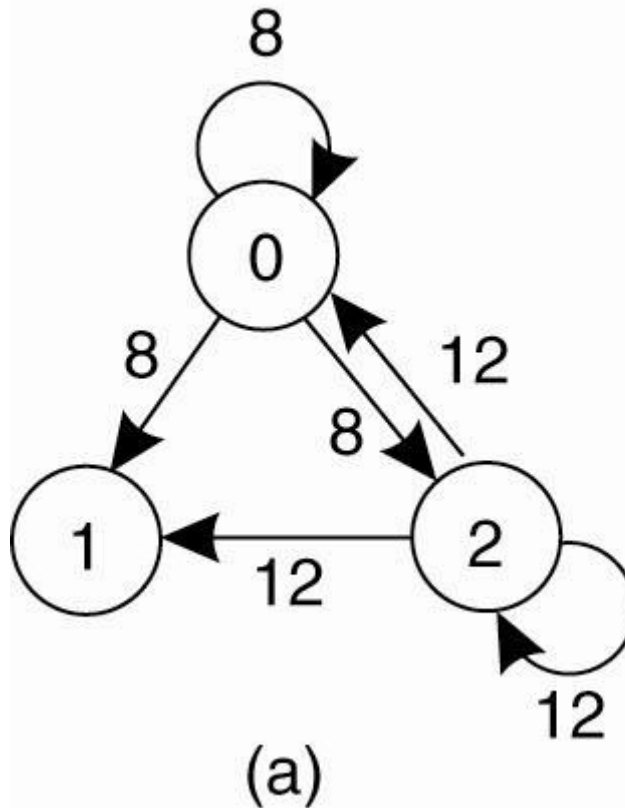


Figure 6-15. (a) Two processes want to access a shared resource at the same moment.



# A Distributed Algorithm (3)

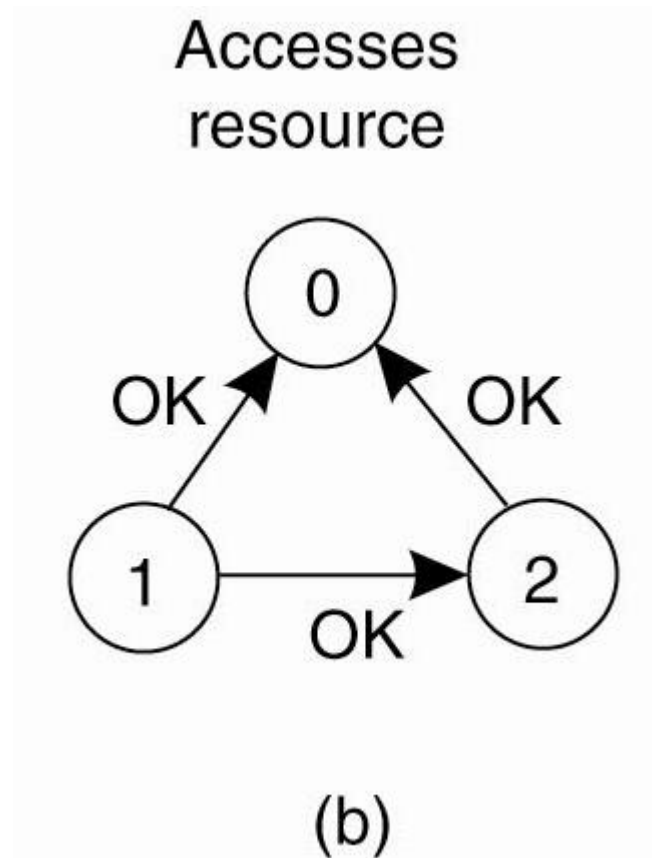


Figure 6-15. (b) Process 0 has the lowest timestamp, so it wins.

# A Distributed Algorithm (4)

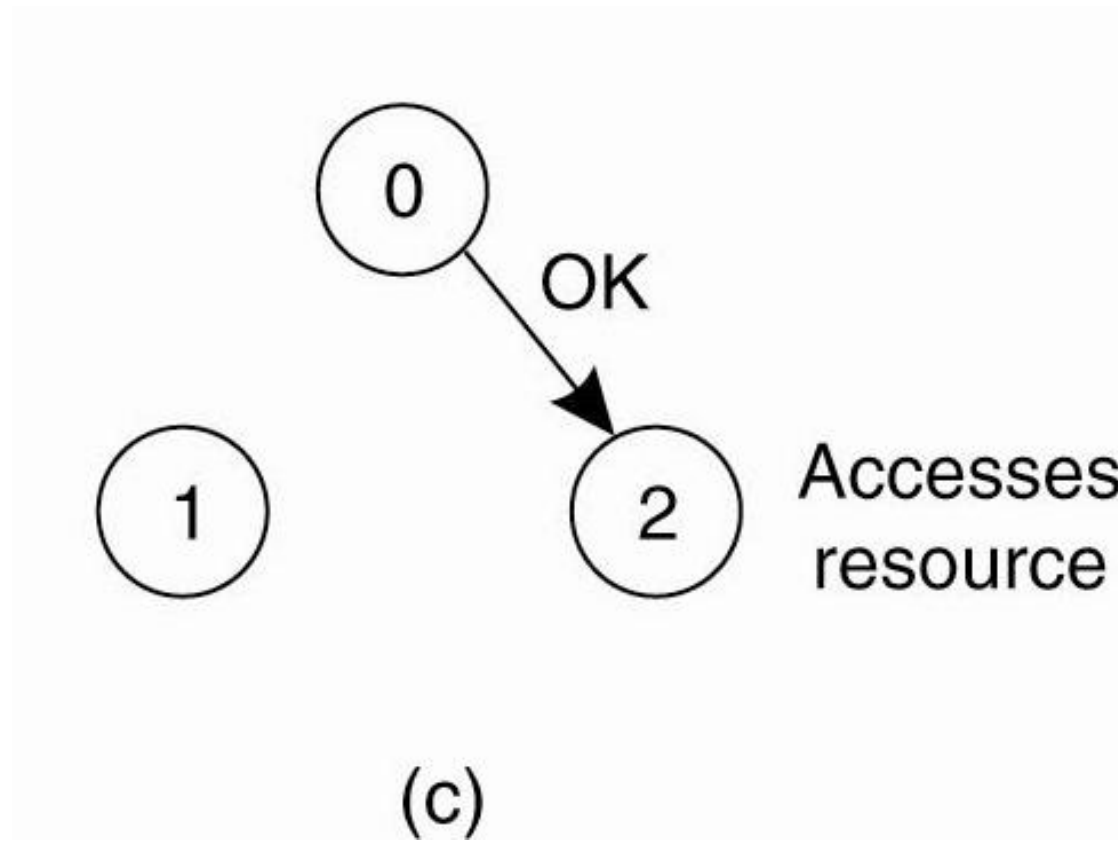


Figure 6-15. (c) When process 0 is done, it sends an OK also, so 2 can now go ahead.

# A Token Ring Algorithm

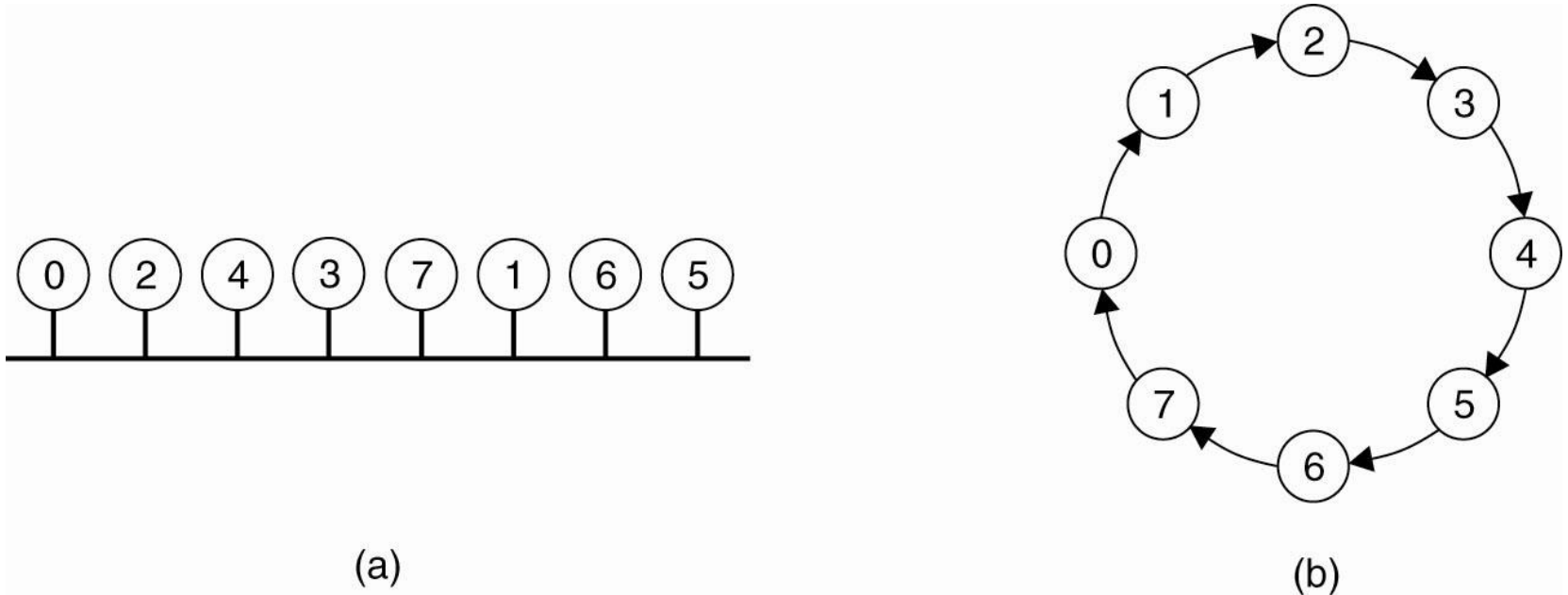


Figure 6-16. (a) An unordered group of processes on a network.  
(b) A logical ring constructed in software.

# A Comparison of the Four Algorithms

Algorithm	Messages per entry/exit	Delay before entry (in message times)	Problems
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

Figure 6-17. A comparison of three mutual exclusion algorithms.

# Election Algorithms

## The Bully Algorithm

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.

# The Bully Algorithm (1)

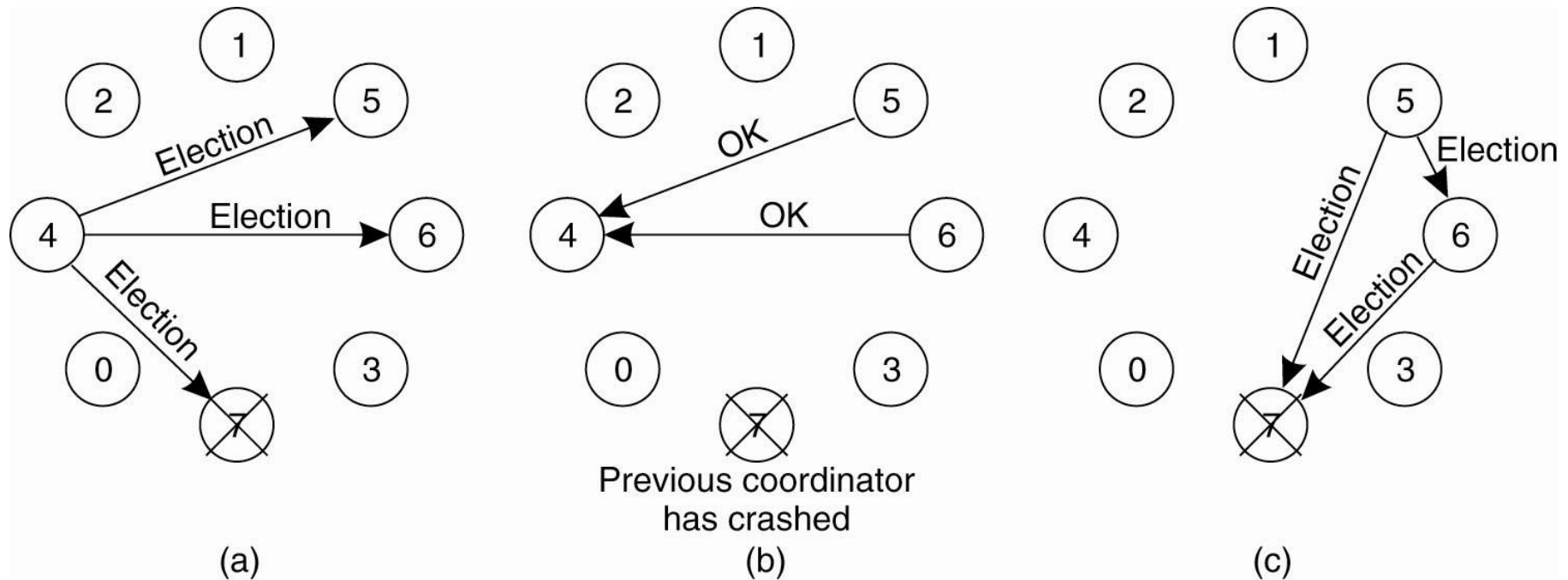


Figure 6-20. The bully election algorithm. (a) Process 4 holds an election. (b) Processes 5 and 6 respond, telling 4 to stop. (c) Now 5 and 6 each hold an election.

# The Bully Algorithm (2)

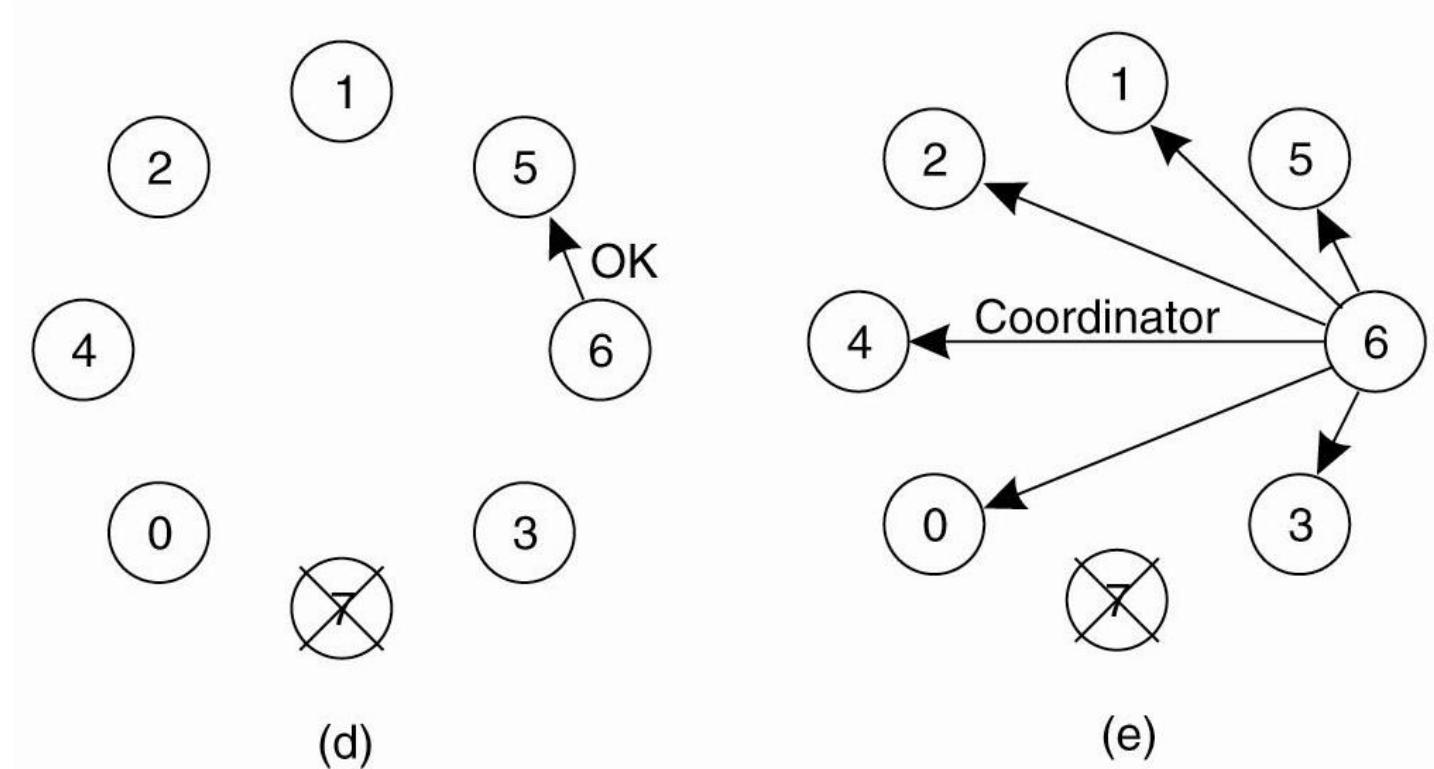


Figure 6-20. The bully election algorithm. (d) Process 6 tells 5 to stop. (e) Process 6 wins and tells everyone.

# A Ring Algorithm

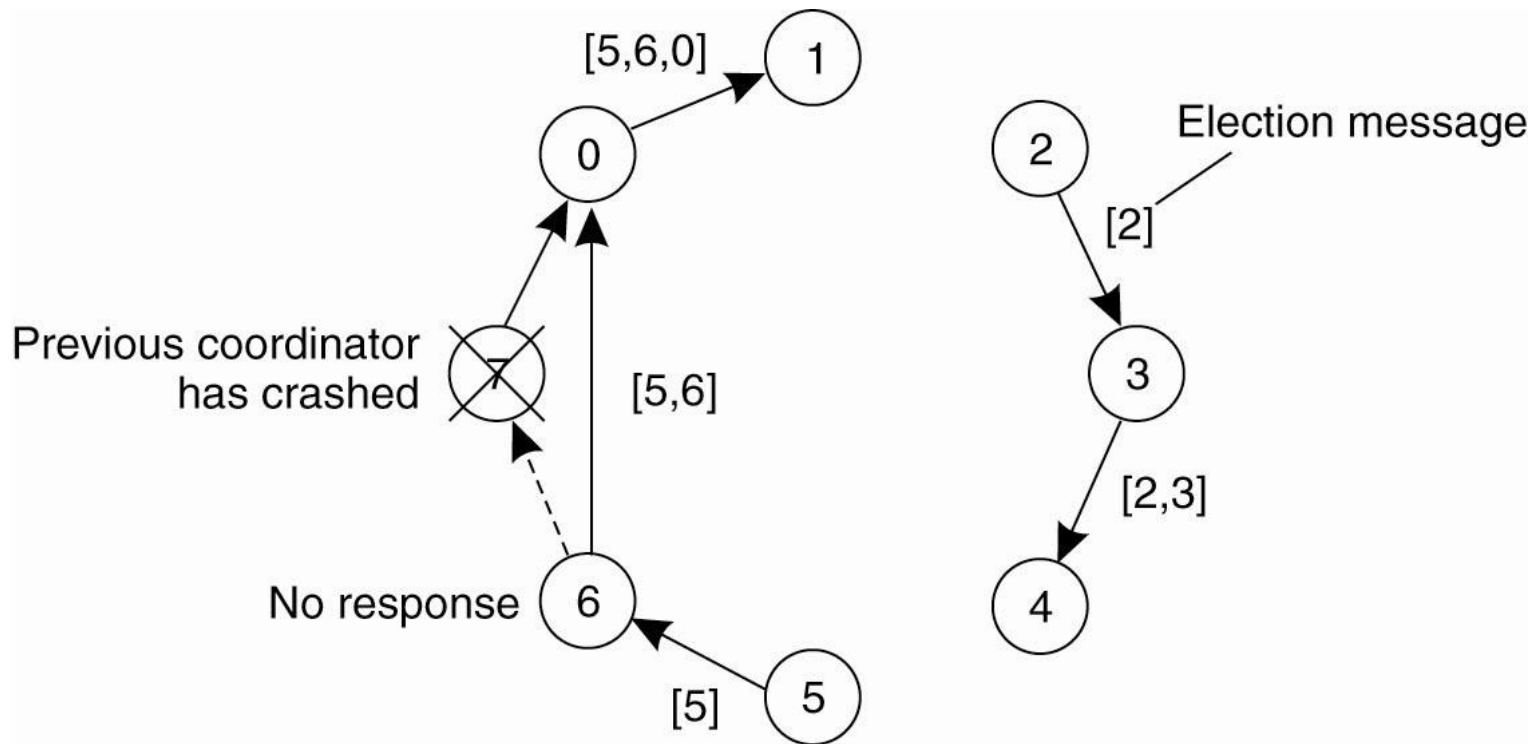


Figure 6-21. Election algorithm using a ring.



# Elections in Wireless Environments (1)

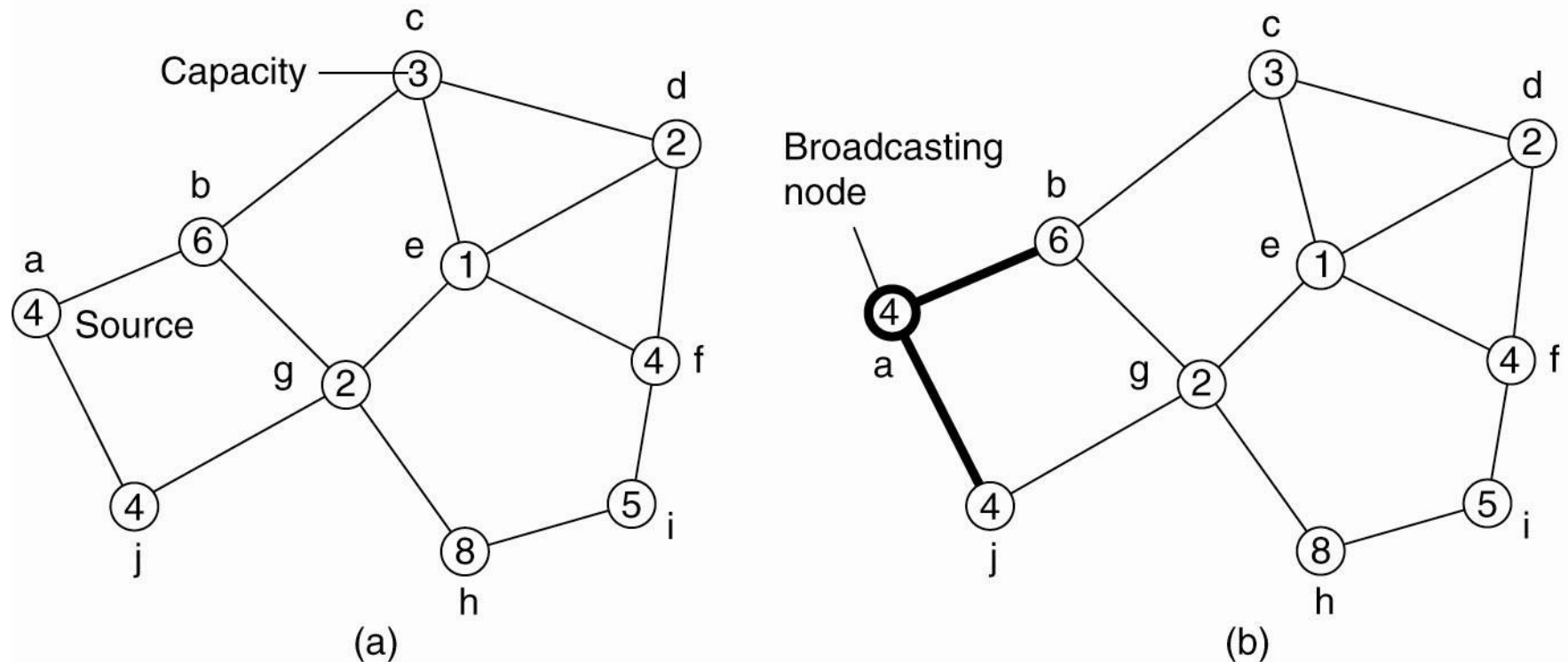


Figure 6-22. Election algorithm in a wireless network, with node a as the source. (a) Initial network. (b)–(e) The build-tree phase

# Elections in Wireless Environments (2)

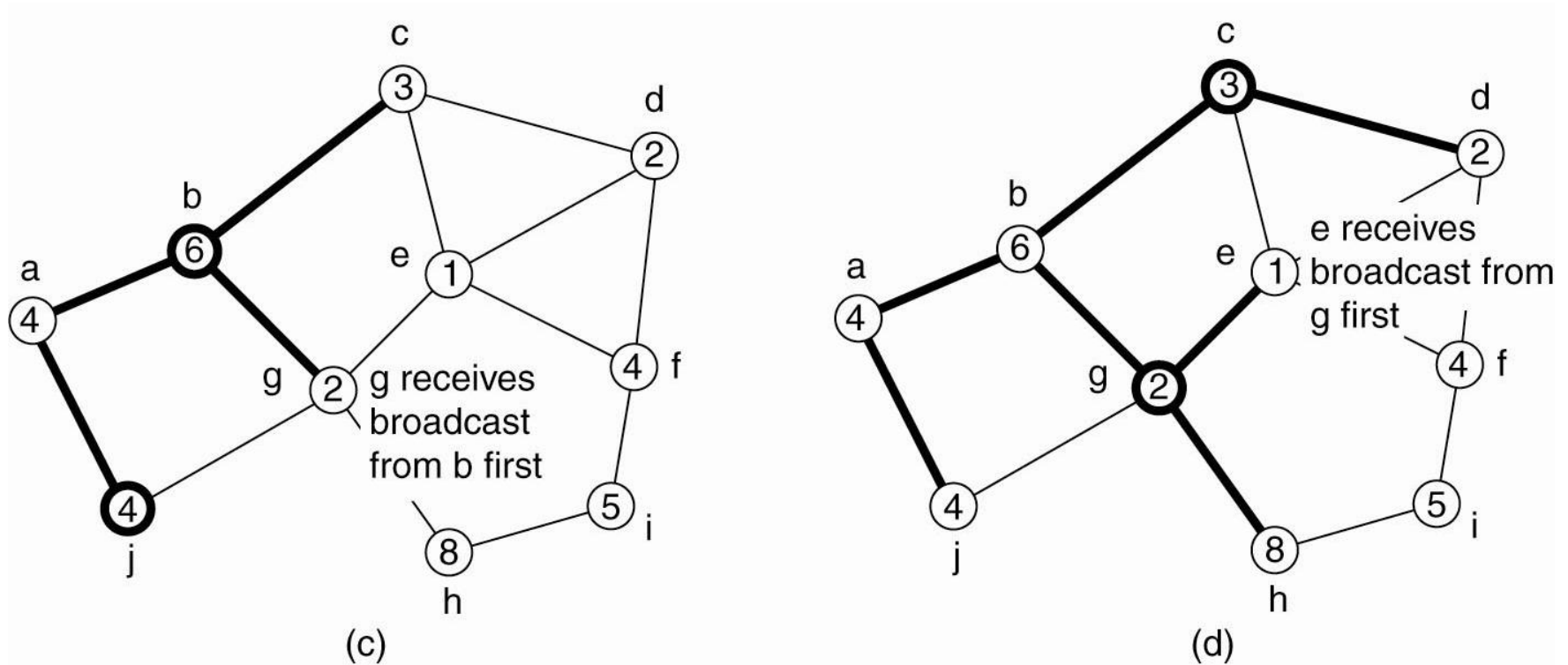


Figure 6-22. Election algorithm in a wireless network, with node a as the source. (a) Initial network. (b)–(e) The build-tree phase

# Elections in Wireless Environments (3)

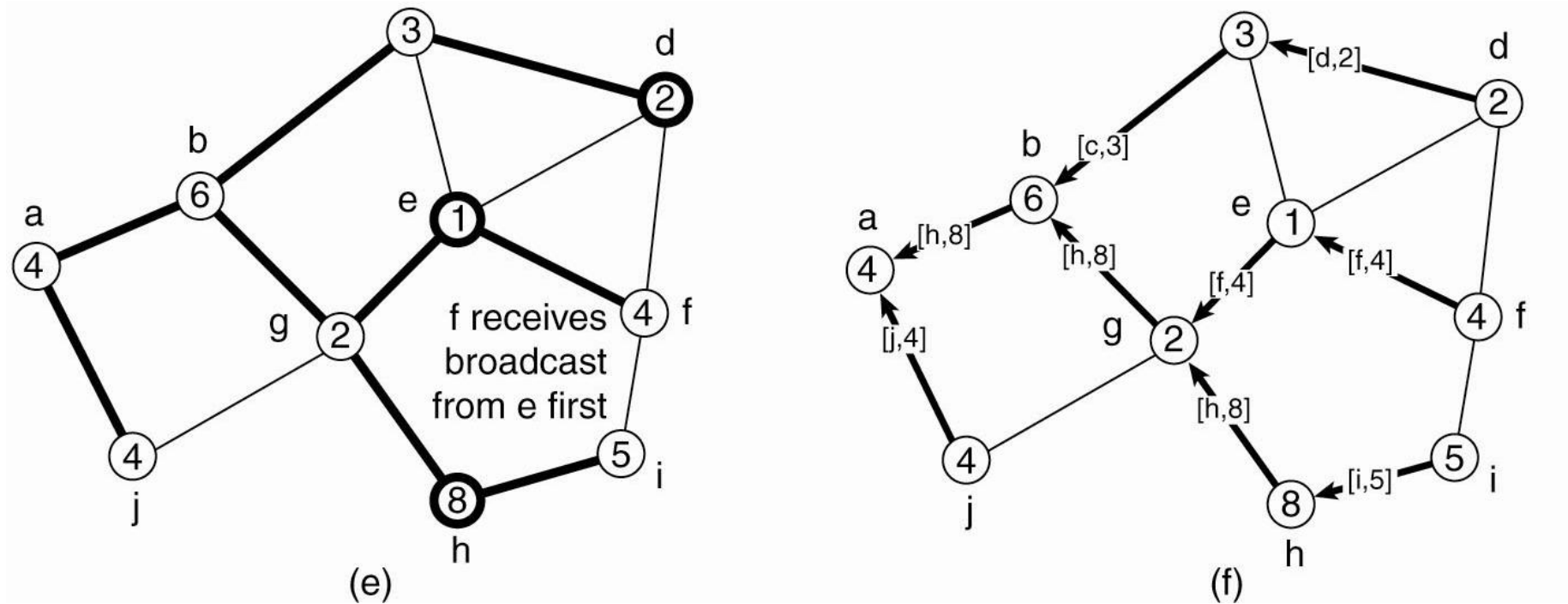


Figure 6-22. (e) The build-tree phase.  
(f) Reporting of best node to source.