

CprE 450/550X  
Distributed Systems and Middleware

# Synchronization

Yong Guan  
3216 Coover  
Tel: (515) 294-8378  
Email: [guan@ee.iastate.edu](mailto:guan@ee.iastate.edu)

March 27, 2003

2

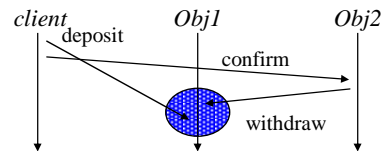
## Readings for Today's Lecture

---

- References
  - Chapter 5 of "Distributed Systems: Principles and Paradigms"

## Synchronization: Introduction

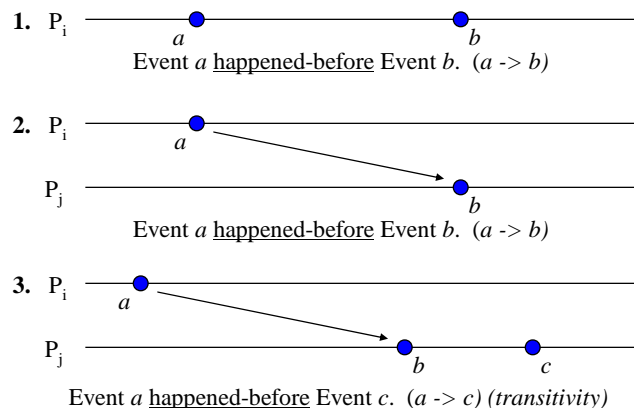
◆ A scary scenario:



- **Synchronization:** temporal ordering of sets of events produced by concurrent processes in time.
  - Synchronization between senders and receivers of messages.
  - Control of joint activity.
  - Serialization of concurrent access to shared objects/resources.
- Why not Semaphores ?!
  - centralized systems: shared memory, central clock
  - distributed system: message passing, no global clock
- Events cannot be totally ordered!

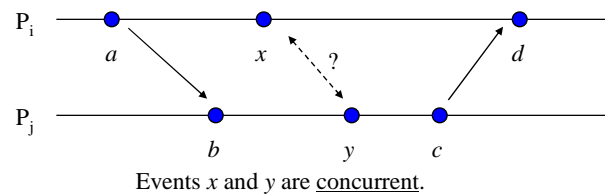
## A Partial Event Ordering for Distributed Systems (Lamport 1978)

- Absence of central time means: no notion of *happened-when* (no total ordering of events)
- But can generate a *happened-before* notion (partial ordering of events)
- *Happened-Before* relation:



## *happened-before* Relation

- What when no *happened-before* relation exists between two events?



- Problem:
  - only approximate knowledge of state of other processes
- Need global time:
  - common clock
  - synchronized clocks

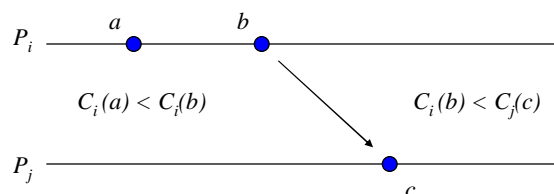
## Logical Clocks

- ◆ Absolute time?
- ◆ Is chronological ordering necessary?
- ◆ Logical clock: assigns a number to each local event.

### Clock Condition

$\forall$  Events  $a, b$  : if  $a \rightarrow b$ , then  $C(a) < C(b)$

- In Other Words:



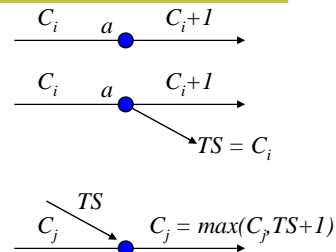
## Total Ordering with Logical Clocks

◆ Rules:

Rule 1: increment  $C_i$  after every local event.

Rule 2: timestamp outgoing messages with current local clock

Rule 3: Upon receiving message with timestamp  $TS$ ,  $P_j$  updates local clock  $C_j$  to be  
 $C_j = \max(C_j, TS+1)$

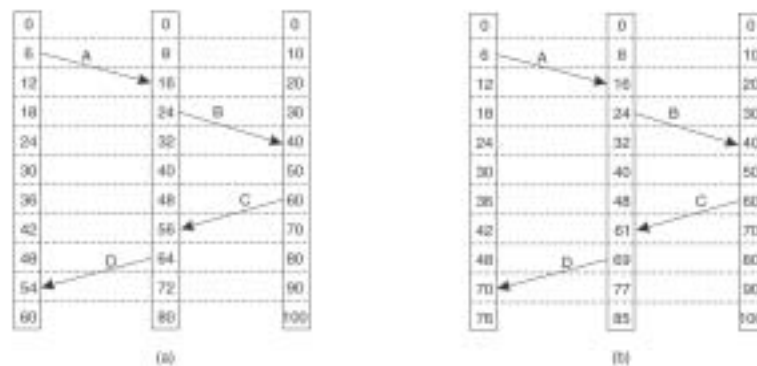


- Total ordering of events: assuming that clocks satisfy Clock Condition, define following relation:

$$a \Rightarrow b \Leftrightarrow \begin{array}{l} C_i(a) < C_j(b) \\ \text{or} \\ C_i(a) = C_j(b) \text{ and } i < j \end{array}$$

for events  $a$  on  $P_i$  and  $b$  on  $P_j$ .

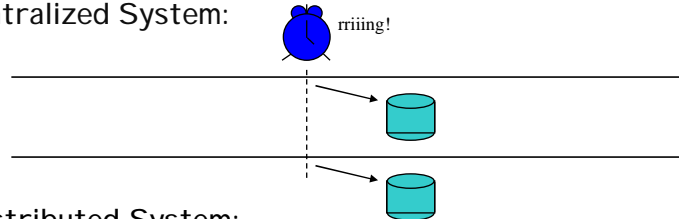
## Lamport Timestamps



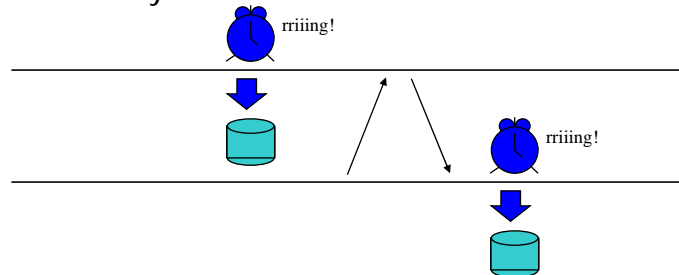
- Three processes, each with its own clock. The clocks run at different rates.
- Lamport's algorithm corrects the clocks.

## Example: Distributed Checkpointing

- ◆ "At 5pm everybody writes its state to stable storage!"
- ◆ Centralized System:

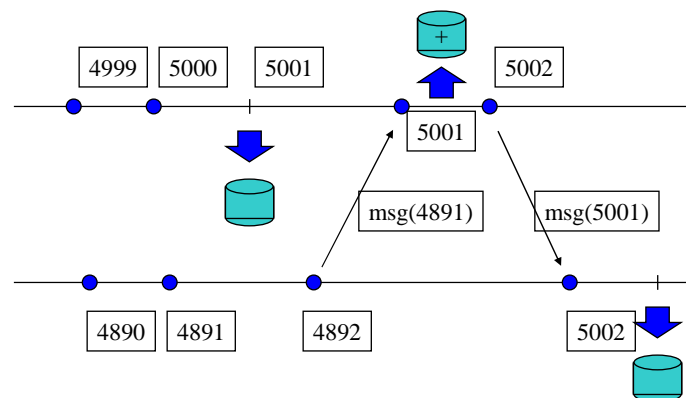


- Distributed System:

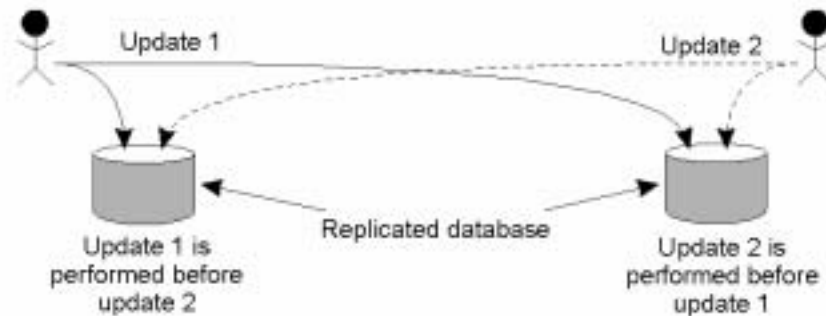


## Distributed Checkpointing and Logical Clocks

"At logical-clock time 5000 write state to stable storage!"



## Another Example: Totally-Ordered Multicasting



Updating a replicated database and leaving it in an inconsistent state.

## Vector Timestamps

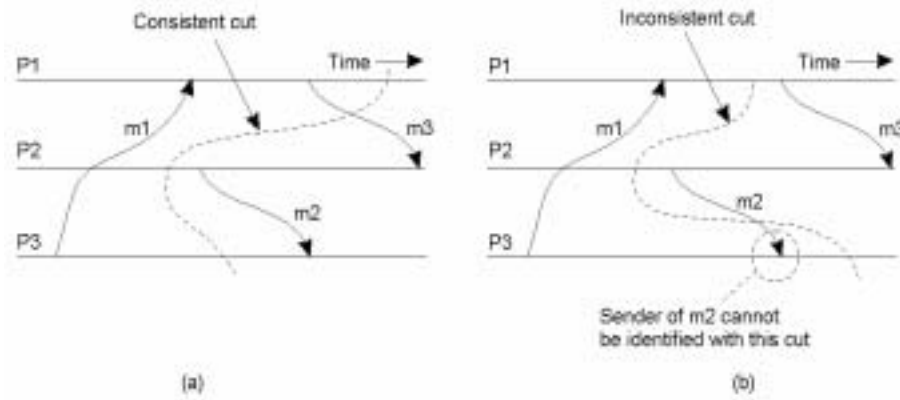
- ◆ Lamport timestamps: Can we say sth if  $C(a) < C(b)$ ?
- ◆ Example: BBS message A and B, if totally ordered multicast is used, no way to say whether A is a reaction to B, or A and B are completely independent.
  - The problem of Lamport timestamps does not capture causality.
- ◆ **Causality can be captured by Vector Timestamps**
  - $VT(a)$ : A vector timestamp assigned to event a.
  - If  $VT(a) < VT(b)$ , then event a is known to causally precede event b.
  - Vector timestamp are constructed by letting each process  $P_i$  maintain a vector  $V_i$  with the following properties:
    1.  $V_i[i]$  is the number of events happened so far at  $P_i$
    2. If  $V_i[j] = k$ , then  $P_i$  knows that k events have occurred at  $P_j$ .

## Global State (1)

---

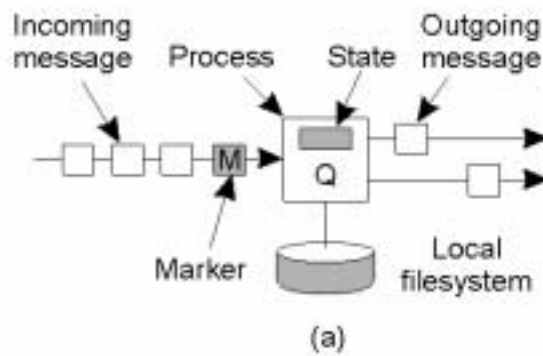
- ◆ Knowing the global state in distributed systems is useful on many occasions.
- ◆ The global state consists of the local state of each process, together with the messages-in-transit.
- ◆ Distributed Snapshot (Chandy and Lamport'85)

## Global State (2)



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

## Global State (3)

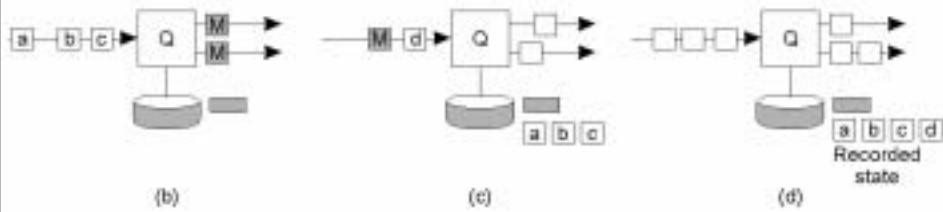


- a) Organization of a process and channels for a distributed snapshot



## Global State (4)

---

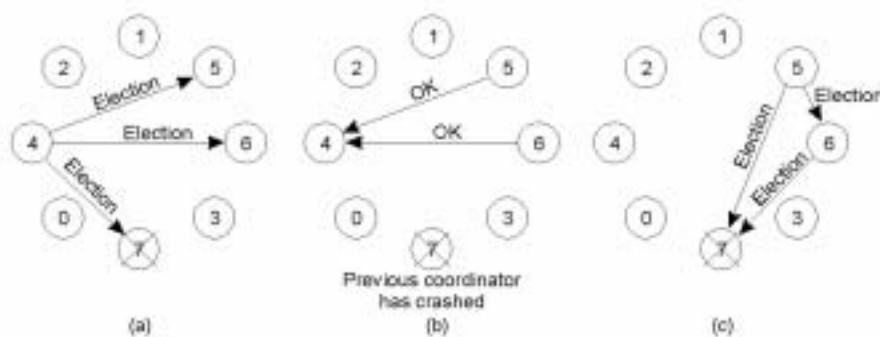


- b) Process Q receives a marker for the first time and records its local state
- c) Q records all incoming message
- d) Q receives a marker for its incoming channel and finishes recording the state of the incoming channel

## Election Algorithms

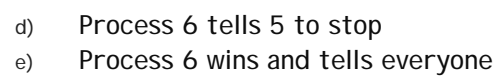
- ◆ Many distributed algorithms requires one process in the system acts as a leader (coordinator, initiator).
- ◆ It does not matter which process it is, but one of them has to do it.
- ◆ The goal of election algorithm is to ensure that when an election starts, it concludes with all processes agreeing on who the new coordinator is to be.

## The Bully Algorithm (1)



The bully election algorithm

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



## Election algorithm using a ring.

---

Next class, we will discuss mutual exclusion and distributed transaction model.