

# Time, Clocks, and the Ordering of Events in a Distributed System

by L. Lamport

CS 5204 Operating Systems

Vladimir Glina

Fall 2005

9/14/2005

1

## Overview

- Key Points
- Background
- Partial Ordering
- Extension for Total Ordering
- Further Work
- Key Points Reiteration
- Evaluation
- Discussion

9/14/2005

2

## Key Points

1. The “happens before” relation on the system event set
2. The events partial ordering on the base of the relation
3. The distributed algorithm for logical clock synchronization
4. The algorithm extension to the case of total events ordering
5. The algorithm application for physical clock synchronization

9/14/2005

3

## Background: Distributed System Features

- Spatially separated processes
- Processes communicate through messages
- Message delays are considerable
- Absence of the single timer leads to synchronization problems
  - Example: totally ordered multicast

9/14/2005

4

## Background: Synchronization Approaches

- Physical Clock Adjustment
  - All clocks show the same actual time
  - Problems:
    - **Most important:** backward time flow possible
    - Sophisticated time services (i.e. WWV); or
    - Reliance on a human operator
- Logical Clock Adjustment
  - Consistency is important, not actual time

9/14/2005

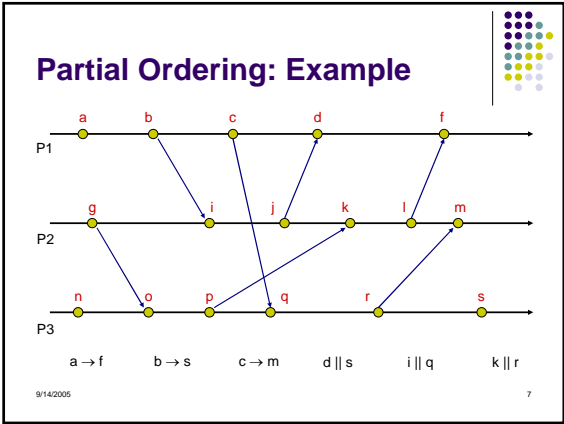
5

## Partial Ordering: Basics

- A system is a set of processes  $P_i$
- A process is a set of events  $a, b, \dots$  with total ordering
- “Happened before” ( $\rightarrow$ ) relation:
  - $(a \in P) \ \&\& \ (b \in P) \ \&\& \ (a \text{ comes before } b) \Rightarrow a \rightarrow b$
  - $(P_1 \text{ sends } a \text{ to } P_2) \ \&\& \ (b \text{ is the receipt of } P_2 \text{ for } a) \Rightarrow a \rightarrow b$
  - $(a \rightarrow b) \ \&\& \ (b \rightarrow c) \Rightarrow a \rightarrow c$
- $\!(a \rightarrow b) \ \&\& \ \!(b \rightarrow a) \Rightarrow a \text{ and } b \text{ are concurrent}$
- $\!(a \rightarrow a) \ \forall a$ , so “happened before” is an irreflexive partial ordering on the set of all the system events

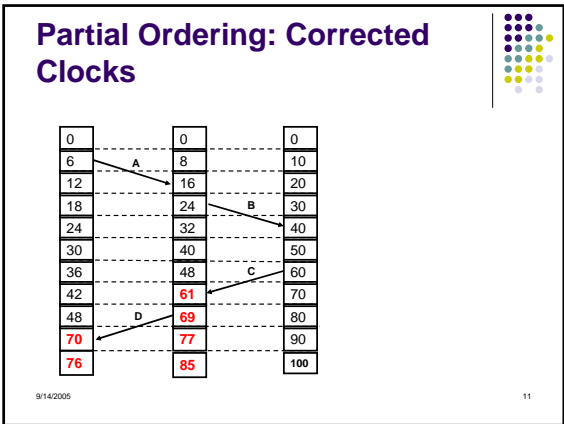
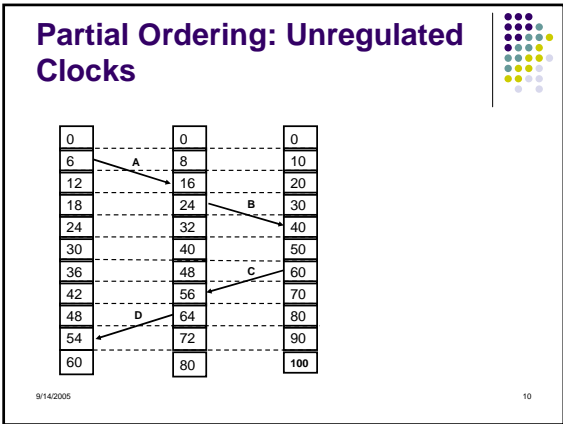
9/14/2005

6



- ### Partial Ordering: Synchronization
- Logical clock:  $C(a) = C_j(a)$  if  $a \in P_j$
  - **Check condition:** for  $\forall a, b$   
 $a \rightarrow b \Rightarrow C(a) < C(b)$  (**not vice versa**)
  - The check condition is satisfied if
    - C1. ( $a, b \in P_j$ ) && ( $a$  comes before  $b$ )  
 $\Rightarrow C_i(a) < C_i(b)$
    - C2. ( $P_i$  sends  $a$  to  $P_j$ ) && ( $b$  is the receipt of  $P_j$  to  $a$ )  
 $\Rightarrow C_i(a) < C_j(b)$
  - **C never decreases!**
- 9/14/2005 8

- ### Partial Ordering: Implementation Rules
- **IR1.** Each  $P_i$  increments  $C_i$  between any two successive events.
  - **IR2.**
    - a) If  $a$  is the sending of a message  $m$  by  $P_i$ , then  $m$  contains a timestamp  $T_m = C_i(a)$ ; **and**
    - b) Upon receiving  $m$ ,  $P_j$  sets  $C_j$  greater than or equal to its present value and greater than  $T_m$
- 9/14/2005 9



- ### Total Ordering: Definition
- $\prec$  is an **arbitrary** total ordering of processes
  - "Happen before" for total ordering ( $\implies$ ):  
 $(a \in P_i) \ \&\& \ (b \in P_j) \Rightarrow a \implies b$  iff
    - $C_i(a) < C_j(b)$ , **or**
    - $P_i \prec P_j$
  - The total ordering depends on  $C_i$  and is not unique
- 9/14/2005 12

## Total Ordering: Synchronization

- $P_i$  broadcasts the message  $T_m;P_i$  (request resource) and puts it on its request queue.
- When  $P_j$  receives  $T_m;P_i$ , it puts the message on its request queue and sends the acknowledgment to  $P_i$ .
- To release the resource,  $P_i$  removes  $T_m;P_i$  from its queue, broadcasts a timestamped release message.
- When  $P_j$  receives the release message, it removes  $T_m;P_i$  from its queue.
- $P_i$  is granted the resource when
  - It has  $T_m;P_i$  in its queue ordered before any other request in the queue by the relation  $\Rightarrow$ ; and
  - $P_i$  has received a message from every other process timestamped later than  $T_m$ .

9/14/2005

13

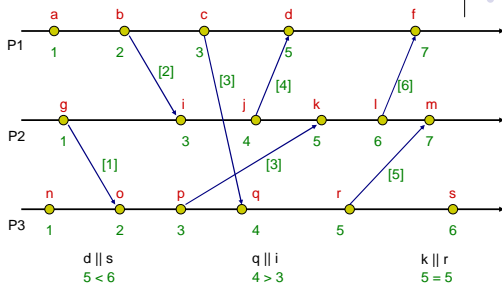
## Further Work: Vector Timestamps

- Lamport clock is:
  - Consistent:  $a \rightarrow b \Rightarrow C(a) < C(b)$
  - Not:**  $C(a) < C(b) \Leftrightarrow a \rightarrow b$  (not strongly consistent)
- Vector timestamps (VT) are strongly consistent
- VT address **potential** causality
  - Allow to say if a **happened** before  $b$ , but not if a **caused**  $b$
- VT say how many events have occurred so far at all processes
- VT solve the totally-ordered multicasting problem

9/14/2005

14

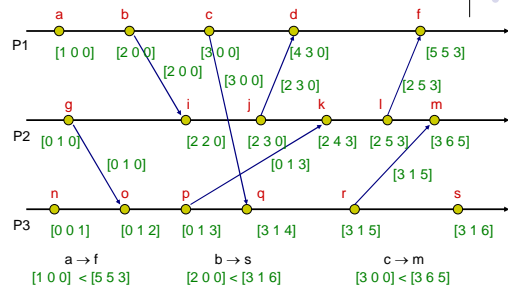
## Lack of Strong Consistency



9/14/2005

15

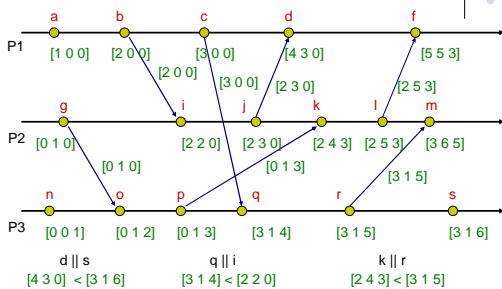
## Vector Clocks (1)



9/14/2005

16

## Vector Clocks (2)



9/14/2005

17

## Key Points Reiteration

- The "happens before" relation on the system event set
- The events partial ordering on the base of the relation
- The distributed algorithm for logical clock synchronization
- The algorithm extension to a case of total events ordering
- The algorithm application for physical clock synchronization

9/14/2005

18

## Evaluation



- The logical clocks idea is very appealing
- Virtually no revision on previous work
- Nice to have more mathematically strict extension on total ordering, if possible

9/14/2005

19

## Discussion



Thank you!

Any questions?

9/14/2005

20