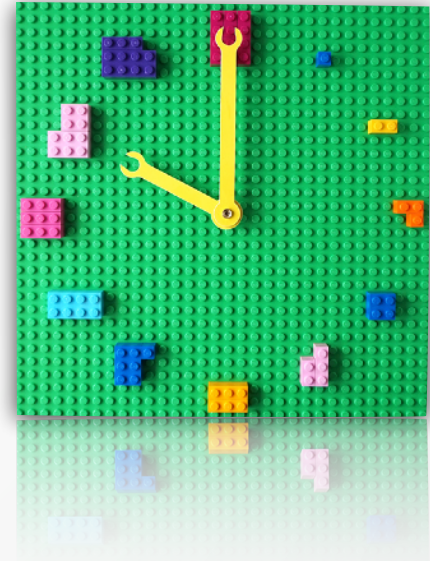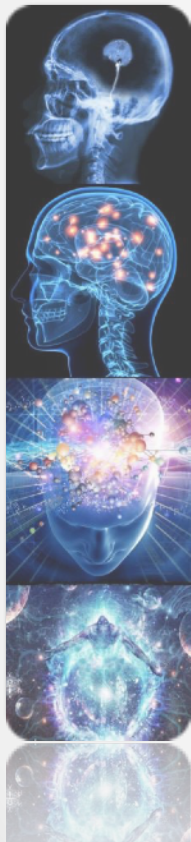ID2203

KTH-2023

**Advanced Course**
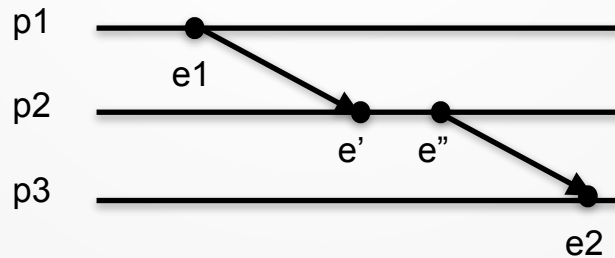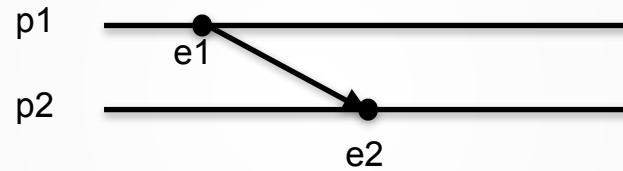# Distributed Systems

## Logical Clocks

Paris Carbone

# COURSE TOPICS

▶ Intro to Distributed Systems

▶ Fundamental Abstractions and Failure Detectors

▶ Reliable and Causal Order Broadcast

▶ Distributed Shared Memory-CRDTs

▶ Consensus (Paxos)

▶ Replicated State Machines (OmniPaxos, Raft, Zab etc.)

▶ Time Abstractions and Interval Clocks (Spanner etc.)

▶ Consistent Snapshotting (Stream Data Management)

▶ Distributed ACID Transactions (Cloud DBs)

ID2203

KTH-2023

# RECAP: CAUSAL ORDER

- Given an execution trace $\beta$,

  two events $a, b \in \beta$ are causally ordered ($a \rightarrow_\beta b$) iff either:

  - a occurs before b on the same process

  - a is a send(m) and b deliver(m) event

  - there exists a sequence of causally ordered events from a to b (transitive) e.g. If $a \rightarrow_\beta c$ and $c \rightarrow_\beta b$

- Two events, a and b, are concurrent if **not** a $\rightarrow_\beta$ b and **not** b $\rightarrow_\beta$ a

  - Concurrent events are denoted as such: $a || b$

ID2203

KTH-2023

# Key Observations

- We <u>cannot</u> order any two events in a D.S. using physical time.
  - <u>Multiple</u> Physical Clocks cannot keep precise physical time.

- We <u>can</u> order any two events in a D.S. using logical time, i.e., causal order.
  - <u>Logical</u> Clocks can capture causality.

Operating Systems     R. Stockton Gaines Editor

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

Leslie Lamport
Massachusetts Computer Associates, Inc.

The concept of one event happening before another in a distributed system is examined, and is shown to define a partial ordering of the events. A distributed algorithm is given for synchronizing a system of logical clocks which can be used to totally order the events. The use of the total ordering is illustrated with a method for solving synchronization problems. The algorithm is then specialized for synchronizing physical clocks, and a bound is derived on how far out of synchrony the clocks can become.

Key Words and Phrases: distributed systems, computer networks, clock synchronization, multiprocess systems

CR Categories: 4.32, 5.29

A distributed system consists of a collection of distinct processes which are spatially separated, and which communicate with one another by exchanging messages. A network of interconnected computers, such as the ARPA net, is a distributed system. A single computer can also be viewed as a distributed system in which the central control unit, the memory units, and the input-output channels are separate processes. A system is distributed if the message transmission delay is not negligible compared to the time between events in a single process.

We will concern ourselves primarily with systems of spatially separated computers. However, many of our remarks will apply more generally. In particular, a multiprocessing system on a single computer involves problems similar to those of a distributed system because of the unpredictable order in which certain events can occur.

In a distributed system, it is sometimes impossible to say that one of two events occurred first. The relation "happened before" is therefore only a partial ordering of the events in the system. We have found that problems often arise because people are not fully aware of this fact and its implications.

In this paper, we discuss the partial ordering defined by the "happened before" relation, and give a distributed algorithm for e_____
of all the even_____
mechanism for_____
illustrate its us_____
chronization p_____
ior can occur if_____
differs from th_____
avoided by intro_____
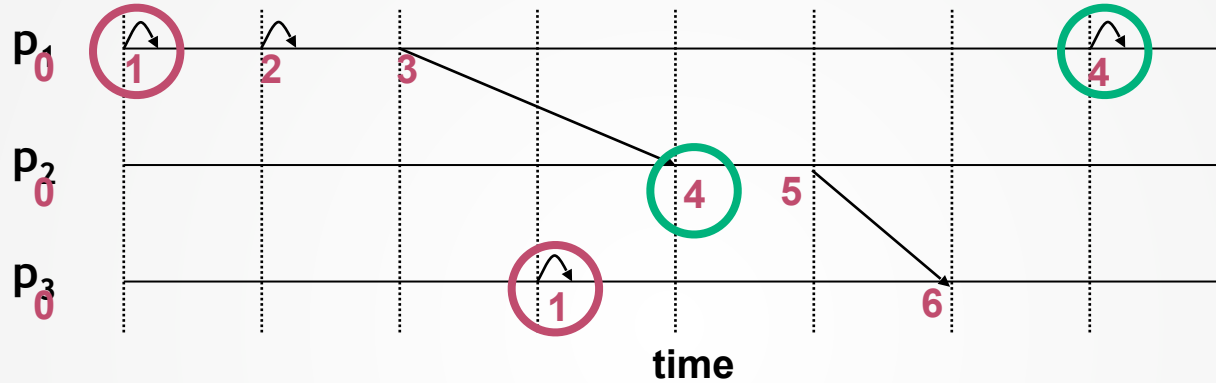
# LOGICAL CLOCK INTUITION

- A logical clock is :
    - an algorithm that assigns a timestamp to each event occurring in a distributed system. $t(a), t(b)$, etc.

- Timestamps can be used to derive a relation between events.


- We are interested in preserving the happen-before relation.

    If a $\rightarrow_\beta$ b then t(a) < t(b)

- **Two types of clocks: Lamport and Vector clocks**

# Lamport Clocks

# LAMPORT CLOCKS

- Each process has a local logical Lamport clock, kept in variable $t_p$,

- initially $t_p = 0$

  - A process p piggybacks $t_p$ on every message sent

- On **internal event** $a$:

  - $t_p := t_p + 1$ ; perform internal event $a$

- **Sending event** of message m:

  - $t_p := t_p + 1$ ; send(m, $t_p$))

- Receiving/Delivering **event** of message m with timestamp $t_q$ from q:

  - $t_p := \max(t_p, t_q) + 1$ ; perform delivery event $a$

# LAMPORT LOGICAL CLOCKS



Lamport logical clocks guarantee that:

If $a \rightarrow_\beta b$, then $\mathbf{t}(a) < \mathbf{t}(b)$,

if $\mathbf{t}(a) \geq \mathbf{t}(b)$, then not $(a \rightarrow_\beta b)$

Can we do better?

# TOTAL ORDER WITH LAMPORT CLOCKS

- We can timestamp with process identifiers.
- The pair (t, p) is unique

Total Order Relation $(\prec)$ : $(t_p, p) \prec (t_q, q)$ iff either

- $t_p < t_q$
- $t_p = t_q \wedge p < q$

i.e. break ties using process identifiers

e.g. $(5, p_5) < (7, p_2)$, $(4, p_2) < (4, p_3)$

ID2203

KTH-2023

# LAMPORT CLOCKS

- Each process has a local logical lamport clock, kept in variable $t_p$,

- initially $t_p = 0$

  - A process p piggybacks $(t_p, p)$ on every message sent

- On **internal event** *a*:

  - $t_p := t_p + 1$ ; perform internal event *a*

- **Sending event** of message m:

  - $t_p := t_p + 1$ ; send(m, $(t_p, p)$)

- Receiving/Delivering **event** of message m with timestamp $(t_q, q)$ from q:

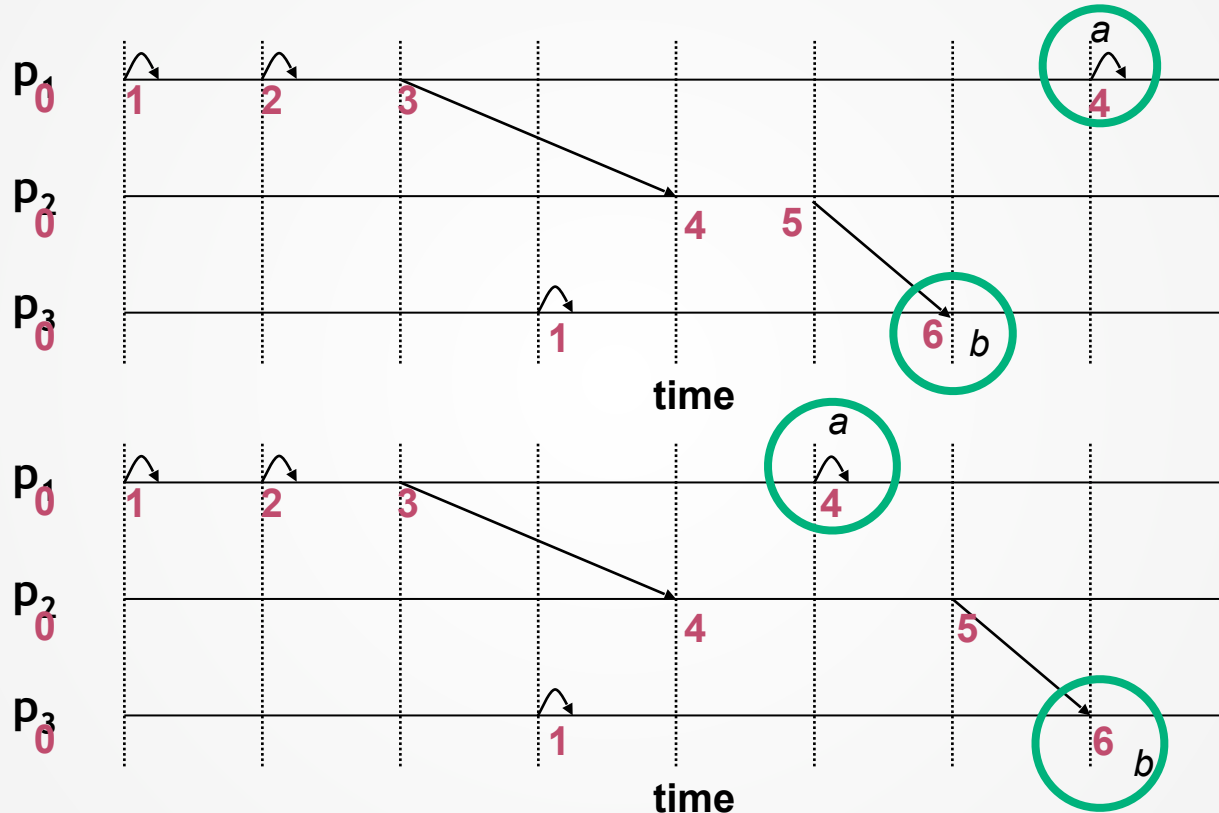  - $t_p := \max(t_p, t_q) + 1$ ; perform delivery event *a*

# REMARKS

- The total order ($\prec$) defined on $(t, p_i)$ pairs is:

    - A method to deterministically <u>derive</u> a total order of events using local and relative process information.

    - Always respecting causal order.

    - A convention rather than the <u>actual</u> order of events (whatever that means).

# Vector Clocks

- Two events $a$ and $b$ are <span style="color:red">concurrent</span>  $(a \parallel_\beta b)$ in an execution E $(\text{trace}(E) = \beta)$ if

    - **not** $a \rightarrow_\beta b$ and **not** $b \rightarrow_\beta a$

- Computation theorem implies that if $(a \parallel_\beta b)$ in $\beta$ then there are <span style="color:red">two executions</span> (with traces $\beta_1$ and $\beta_2$) that are <span style="color:red">similar</span> where $a$ occurs before $b$ in $\beta_1$, $b$ occurs before $a$ in $\beta_2$
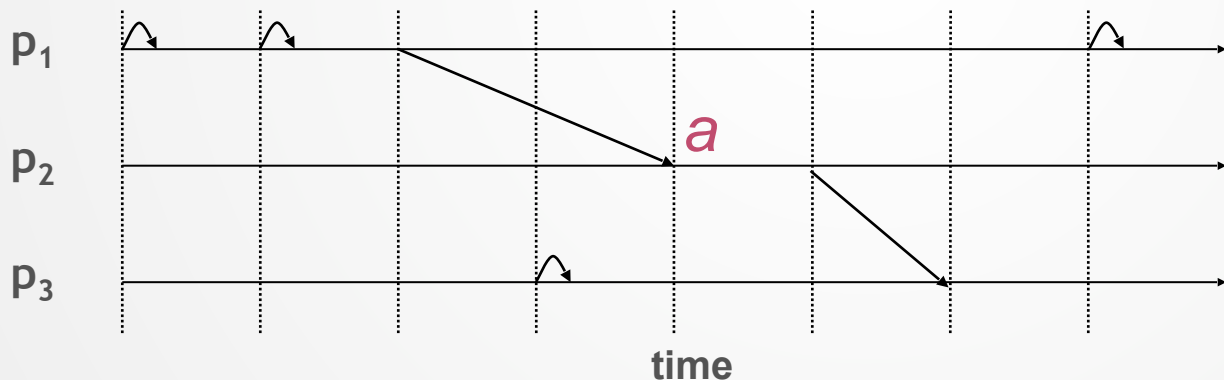
time

time

# VECTOR CLOCKS

- **Lamport Clock Limitation:** We cannot tell by looking to the timestamps of event $a$ and $b$ whether there is a causal relation between the events, or they are concurrent.

- **Vector Clocks** guarantee that:

    - 1. if $\mathbf{v}(a) < \mathbf{v}(b)$ then $a \longrightarrow_\beta b$

    - 2. if $a \longrightarrow_\beta b$ then $\mathbf{v}(a) < \mathbf{v}(b)$ (same as Lamport Clocks)

        - where $\mathbf{v}(a)$ is a vector clock of event $a$

- Vector clock for an event $a$

  - $\mathbf{v}(a) = [x_1, \ldots, x_n]$

  - $x_i$ is <u>the number of events</u> at $p_i$ that happened-before $a$



**p₁**

**3 events at p1**
**1 event at p2**
**0 events at p3**
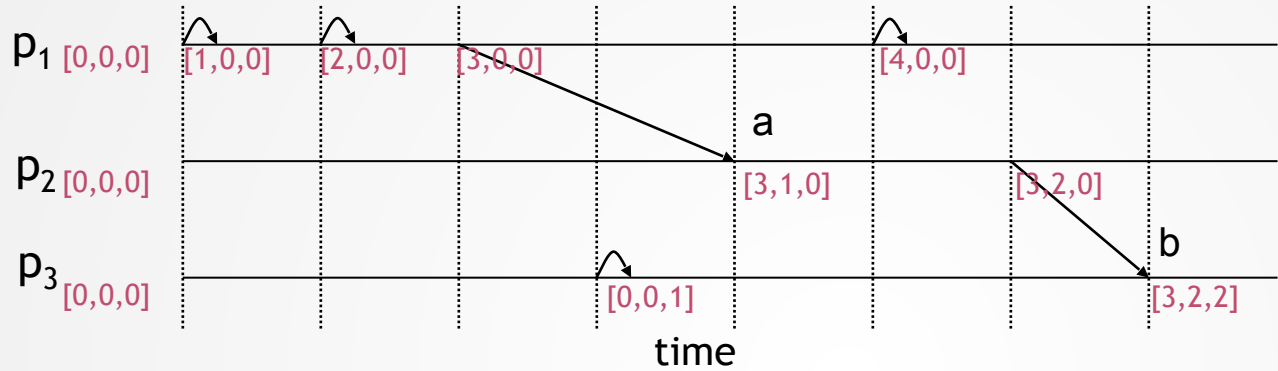
**p₂**

$a$

**p₃**

**[ 3, 1, 0 ]**

**time**

KTH-2023

# VECTOR TIME IMPLEMENTATION

- Processes $p_1, \ldots, p_n \in P$

- Each process $p_i$ has local vector **v** of size **n** (number of processes)

  - **v**[i] = 0 for all i in 1…n

  - Piggyback **v** on every sent message

- For each transition (on each event) update local **v** at $p_i$:

  - **v**[i] := **v**[i] + 1   (internal, send or deliver)

  - **v**[j] := max( **v**[j], **v**$_q$[j] ), for all j ≠ i $\in |P|$ (deliver)

    - where **v**$_q$ is clock in message received from process q

ID2203

KTH-2023

# COMPARING VECTOR CLOCKS

- $v_p \leq v_q$ iff
  - $v_p[i] \leq v_q[i]$ for $\forall i \in |P|$

- $v_p < v_q$ iff
  - $v_p \leq v_q$ and for some i, $v_p[i] < v_q[i]$

- $v_p$ and $v_q$ are concurrent $(v_p \| v_q)$ iff
  - not $v_p < v_q$, and not $v_q < v_p$

- Vector clocks guarantee
  - If $v(a) < v(b)$ then $a \rightarrow b$, and
  - If $a \rightarrow b$, then $v(a) < v(b)$
    - where $v(a)$ is the vector clock of event a

$[3,0,0] \leq [3,1,0]$

$[3,0,0] < [3,1,0]$

$[3,1,0] <> [4,0,0]$

# EXAMPLE OF VECTOR CLOCKS



$p_1$ [0,0,0]  [1,0,0]  [2,0,0]  [3,0,0]  [4,0,0]

a

$p_2$ [0,0,0]  [3,1,0]  [3,2,0]

b

$p_3$ [0,0,0]  [0,0,1]  [3,2,2]

time

v(a) < v(b) implies a → b

a

$p_1$ [0,0,0]  [1,0,0]  [2,0,0]  [3,0,0]  [4,0,0]

$p_2$ [0,0,0]  [3,1,0]  [3,2,0]

b

$p_3$ [0,0,0]  [0,0,1]  [3,2,2]

time

v(a) <> v(b) implies a || b

# LIMITATIONS OF VECTOR CLOCKS?

p$_1$ [0,0,0]  [1,0,0]  [2,0,0]  [3,0,0]  [4,0,0]

p$_2$ [0,0,0]  [3,1,0]  [3,2,0]

p$_3$ [0,0,0]  [0,0,1]  [3,2,2]

time

- Vectors need to be statically defined of size n

- Insufficient for problems that require total  event ordering

# ORDERING - SUMMARY

- the relation $\rightarrow_\beta$ on events in executions

  - Partial: $\rightarrow_\beta$ doesn't order concurrent events

- the relation $\prec$ on Lamport logical clocks

  - Total: any two distinct clock values are ordered (adding pid)

- the relation $\prec$ on vector timestamps

  - Partial: timestamp of concurrent events not ordered

# LOGICAL CLOCK INVARIANTS

**Lamport clock**

$$\text{If } a \xrightarrow{}_\beta b \text{ then } t(a) < t(b) \qquad (1)$$

$$\text{If } t(a) < t(b) \text{ then } \mathbf{not\ b} \xrightarrow{}_\beta \mathbf{a} \quad (2)$$

**Vector clock**

$$\text{If } a \xrightarrow{}_\beta b \text{ then } v(a) < v(b) \qquad (1)$$

$$\text{If } v(a) < v(b) \text{ then } a \xrightarrow{}_\beta b \qquad (2)$$

ID2203

KTH-2023