# CS 425 / ECE 428 Distributed Systems Fall 2016

Indranil Gupta (Indy)

Sep 29, 2016

*Lecture 12: Time and Ordering*

All slides © IG

# WHY SYNCHRONIZATION?

- **You want to catch a bus at 6.05 pm, but your watch is off by 15 minutes**
  - What if your watch is Late by 15 minutes?
    - You'll miss the bus!
  - What if your watch is Fast by 15 minutes?
    - You'll end up unfairly waiting for a longer time than you intended
- **Time synchronization is required for both**
  - Correctness
  - Fairness

2

# Synchronization In The Cloud

- Cloud airline reservation system
- Server A receives a client request to purchase last ticket on flight ABC 123.
- Server A timestamps purchase using local clock 9h:15m:32.45s, and logs it. Replies ok to client.
- That was the last seat. Server A sends message to Server B saying "flight full."
- B enters "Flight ABC 123 full" + its own local clock value (which reads 9h:10m:10.11s) into its log.
- Server C queries A's and B's logs. Is confused that a client purchased a ticket at A after the flight became full at B.
- This may lead to further incorrect actions by C

# WHY IS IT CHALLENGING?

- **End hosts in Internet-based systems (like clouds)**
  - Each have their own clocks
  - Unlike processors (CPUs) within one server or workstation which share a system clock
- **Processes in Internet-based systems follow an *asynchronous* system model**
  - No bounds on
    - Message delays
    - Processing delays
  - Unlike multi-processor (or parallel) systems which follow a *synchronous* system model

# Some Definitions

- An Asynchronous Distributed System consists of a number of processes.

- Each process has a state (values of variables).

- Each process takes actions to change its state, which may be an instruction or a communication action (send, receive).

- An event is the occurrence of an action.

- Each process has a local clock – events *within* a process can be assigned timestamps, and thus ordered linearly.

- But – in a distributed system, we also need to know the time order of events *across* different processes.

# Clock Skew vs. Clock Drift

- **Each process (running at some end host) has its own clock.**
- **When comparing two clocks at two processes**:
  - Clock Skew = Relative Difference in clock *values* of two processes
    - Like distance between two vehicles on a road
  - Clock Drift = Relative Difference in clock *frequencies (rates)* of two processes
    - Like difference in speeds of two vehicles on the road
- **A non-zero clock skew implies clocks are not synchronized.**
- **A non-zero clock drift causes skew to increase (eventually).**
  - If faster vehicle is ahead, it will drift away
  - If faster vehicle is behind, it will catch up and then drift away

# HOW OFTEN TO SYNCHRONIZE?

- Maximum Drift Rate (MDR) of a clock
- Absolute MDR is defined relative to Coordinated Universal Time (UTC). UTC is the "correct" time at any point of time.
  - MDR of a process depends on the environment.
- Max drift rate between two clocks with similar MDR is 2 * MDR
- Given a maximum acceptable skew M between any pair of clocks, need to synchronize at least once every: M / (2 * MDR) time units
  - Since time = distance/speed

# External vs Internal Synchronization

- **Consider a group of processes**
- **External Synchronization**
    - Each process C(i)'s clock is within a bound D of a well-known clock S external to the group
    - $|C(i) - S| < D$ at all times
    - External clock may be connected to UTC (Universal Coordinated Time) or an atomic clock
    - E.g., Cristian's algorithm, NTP
- **Internal Synchronization**
    - Every pair of processes in group have clocks within bound D
    - $|C(i) - C(j)| < D$ at all times and for all processes i, j
    - E.g., Berkeley algorithm

# External vs Internal Synchronization (2)

- **External Synchronization with D => Internal Synchronization with 2*D**

- **Internal Synchronization does not imply External Synchronization**
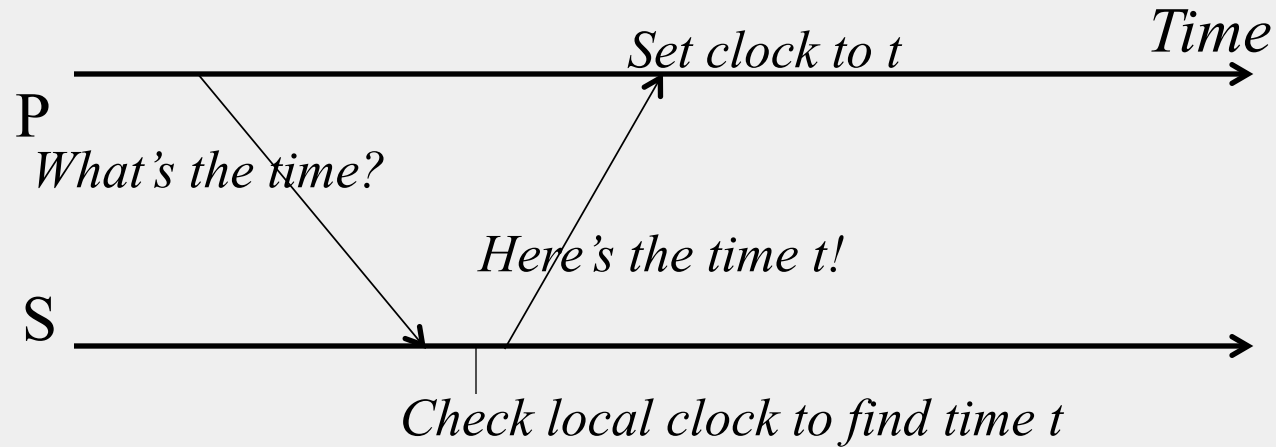  - In fact, the entire system may drift away from the external clock S!

# NEXT

- Algorithms for Clock Synchronization

# CRISTIAN'S ALGORITHM

# BASICS

- **External time synchronization**
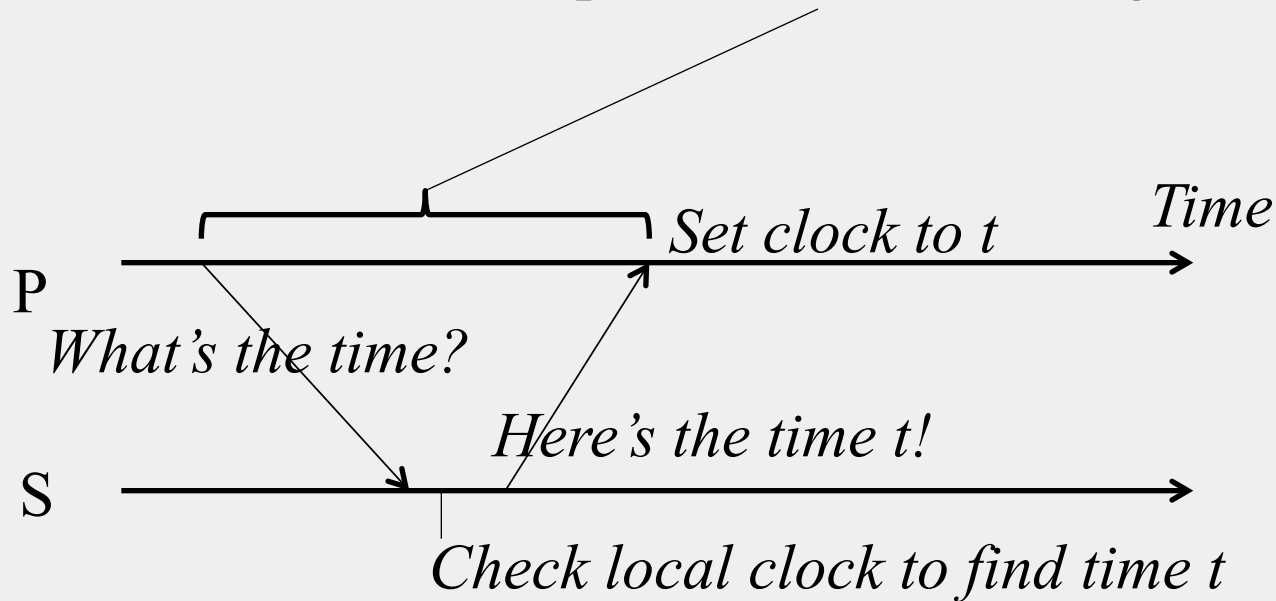- **All processes P synchronize with a time server S**



P

*What's the time?*

*Set clock to t*

*Time*

S

*Here's the time t!*

*Check local clock to find time t*

# What's Wrong

- By the time response message is received at P, time has moved on

- P's time set to $t$ is inaccurate!

- Inaccuracy a function of message latencies

- Since latencies unbounded in an asynchronous system, the inaccuracy cannot be bounded

# Cristian's Algorithm

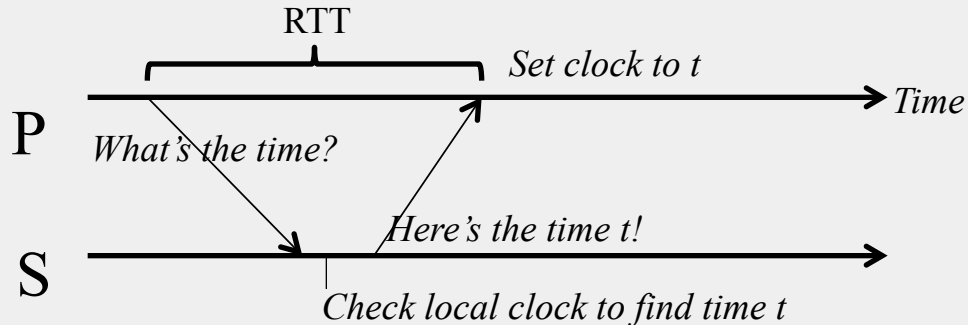- P measures the round-trip-time RTT of message exchange

# CRISTIAN'S ALGORITHM (2)

- **P measures the round-trip-time RTT of message exchange**
- **Suppose we know the minimum P → S latency min1**
- **And the minimum S → P latency min2**
  - min1 and min2 depend on Operating system overhead to buffer messages, TCP time to queue messages, etc.

RTT

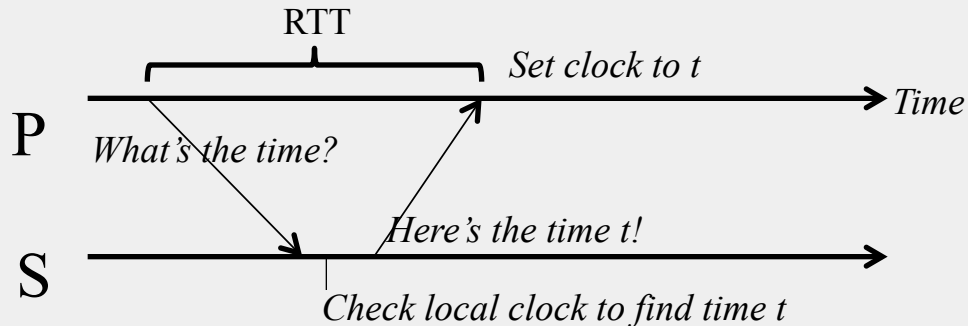P ——————————————————————→ Time

*What's the time?*

*Set clock to t*

S ——————————————————————→

*Here's the time t!*

*Check local clock to find time t*

# CRISTIAN'S ALGORITHM (3)

- P measures the round-trip-time RTT of message exchange
- Suppose we know the minimum P → S latency min1
- And the minimum S → P latency min2
  - min1 and min2 depend on Operating system overhead to buffer messages, TCP time to queue messages, etc.
- The actual time at P when it receives response is between [t+min2, t+RTT-min1]

RTT

Set clock to t

P       Time

*What's the time?*

S

*Here's the time t!*

*Check local clock to find time t*

# CRISTIAN'S ALGORITHM (4)

- The actual time at P when it receives response is between [t+min2, t+RTT-min1]

- P sets its time to halfway through this interval
    - To: t + (RTT+min2-min1)/2

- Error is at most (RTT-min2-min1)/2
    - Bounded!

RTT

*Set clock to t*

P                                        *Time*

*What's the time?*

*Here's the time t!*

S

*Check local clock to find time t*

# GOTCHAS

- **Allowed to increase clock value but should never decrease clock value**
  - May violate ordering of events within the same process

- **Allowed to increase or decrease speed of clock**

- **If error is too high, take multiple readings and average them**

# NTP = Network Time Protocol

- NTP Servers organized in a tree
- Each Client = a leaf of tree
- Each node synchronizes with its tree parent



Primary servers

Secondary servers

Tertiary servers

Client

# NTP Protocol

# What the Child Does

- Child calculates *offset* between its clock and parent's clock

- Uses *ts1, tr1, ts2, tr2*

- Offset is calculated as

$$o = (tr1 - tr2 + ts2 - ts1)/2$$

# WHY o = (tr1 - tr2 + ts2 - ts1)/2?

- **Offset $o = (tr1 - tr2 + ts2 - ts1)/2$**
- **Let's calculate the error**
- **Suppose real offset is *oreal***
  - Child is ahead of parent by *oreal*
  - Parent is ahead of child by *-oreal*
- **Suppose one-way latency of Message 1 is *L1* (*L2* for Message 2)**
- **No one knows *L1* or *L2*!**
- **Then**

  $tr1 = ts1 + L1 + oreal$

  $tr2 = ts2 + L2 - oreal$

# Why o = (tr1 - tr2 + ts2 - ts1)/2? (2)

- **Then**

  $tr1 = ts1 + L1 + oreal$

  $tr2 = ts2 + L2 - oreal$

- **Subtracting second equation from the first**

  $oreal = (tr1 - tr2 + ts2 - ts1)/2 + (L2 - L1)/2$

  $=> oreal = o + (L2 - L1)/2$

  $=> |oreal - o| < |(L2 - L1)/2| < |(L2 + L1)/2|$

  – Thus, the error is bounded by the round-trip-time

23

# AND YET...

- **We still have a non-zero error!**
- **We just can't seem to get rid of error**
  - Can't, as long as message latencies are non-zero
- **Can we avoid synchronizing clocks altogether, and still be able to order events?**

# LAMPORT TIMESTAMPS

# Ordering Events in a Distributed System

- To order events across processes, trying to sync clocks is one approach
- What if we instead assigned timestamps to events that were not *absolute* time?
- As long as these timestamps obey *causality*, that would work

    If an event A causally happens before another event B, then timestamp(A) < timestamp(B)

    Humans use causality all the time

    E.g., I enter a house only after I unlock it

    E.g., You receive a letter only after I send it

# Logical (or Lamport) Ordering

- Proposed by Leslie Lamport in the 1970s
- Used in almost all distributed systems since then
- Almost all cloud computing systems use some form of logical ordering of events

27

# Logical (or Lamport) Ordering(2)

- Define a logical relation *Happens-Before* among pairs of events
- Happens-Before denoted as →
- Three rules
1. On the same process: $a$ → $b$, if *time(a) < time(b)* (using the local clock)
2. If p1 sends *m* to p2: *send(m)* → *receive(m)*
3. (Transitivity) If $a$ → $b$ *and* $b$ → $c$ then $a$ → $c$
- Creates a *partial order* among events
  - Not all events related to each other via →

# EXAMPLE



P1 — A, B, C, D, E

P2 — E, F, G

P3 — H, I, J

*Time*

While P1 and P3 each have an event labeled E, these are different events as they occur at different processes.

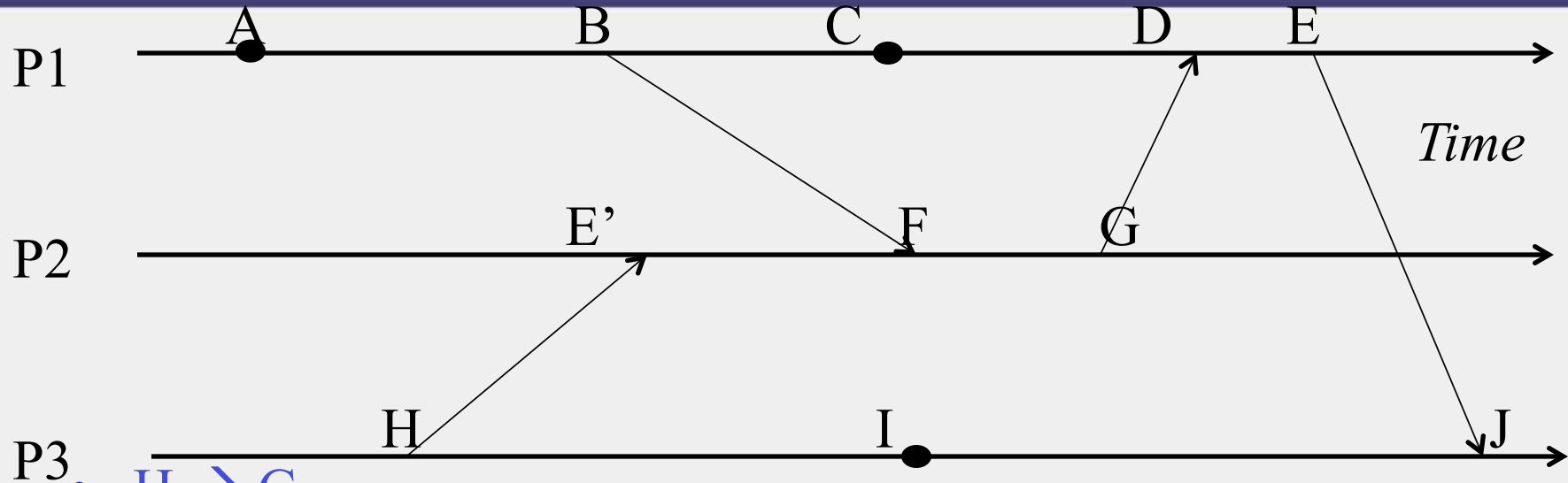| | |
|---|---|
| ● | *Instruction or step* |
| → | *Message* |

29

# Happens-Before



- A ➔ B
- B ➔ F
- A ➔ F

# Happens-Before (2)



P1 — A — B — C — D — E

P2 — E' — F — G

P3 — H — I — J

- H → G
- F → J
- H → J
- C → J

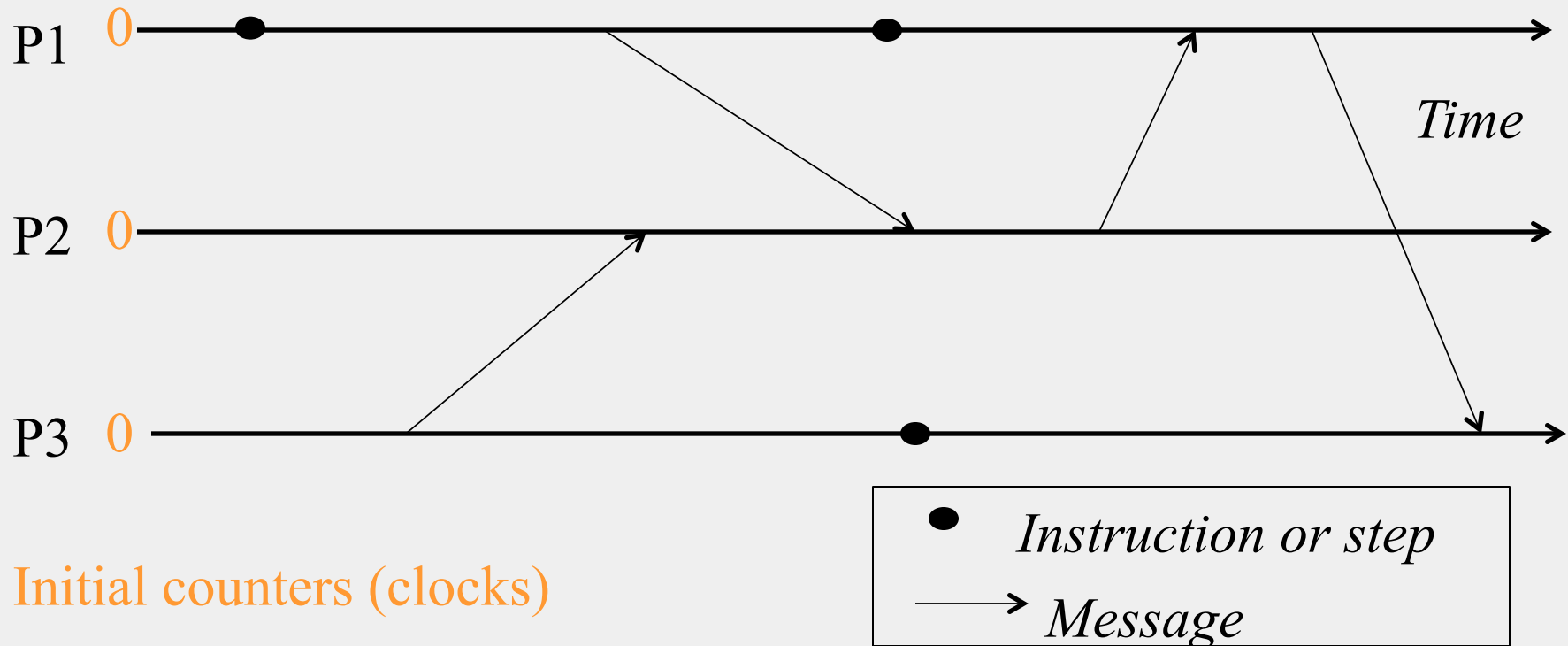Time

Legend:
- Instruction or step
- → Message

# In practice: Lamport timestamps

- **Goal: Assign logical (Lamport) timestamp to each event**
- **Timestamps obey causality**
- **Rules**
  - Each process uses a local counter (clock) which is an integer
    - initial value of counter is zero
  - A process increments its counter when a send or an instruction happens at it. The counter is assigned to the event as its timestamp.
  - A send (message) event carries its timestamp
  - For a receive (message) event the counter is updated by
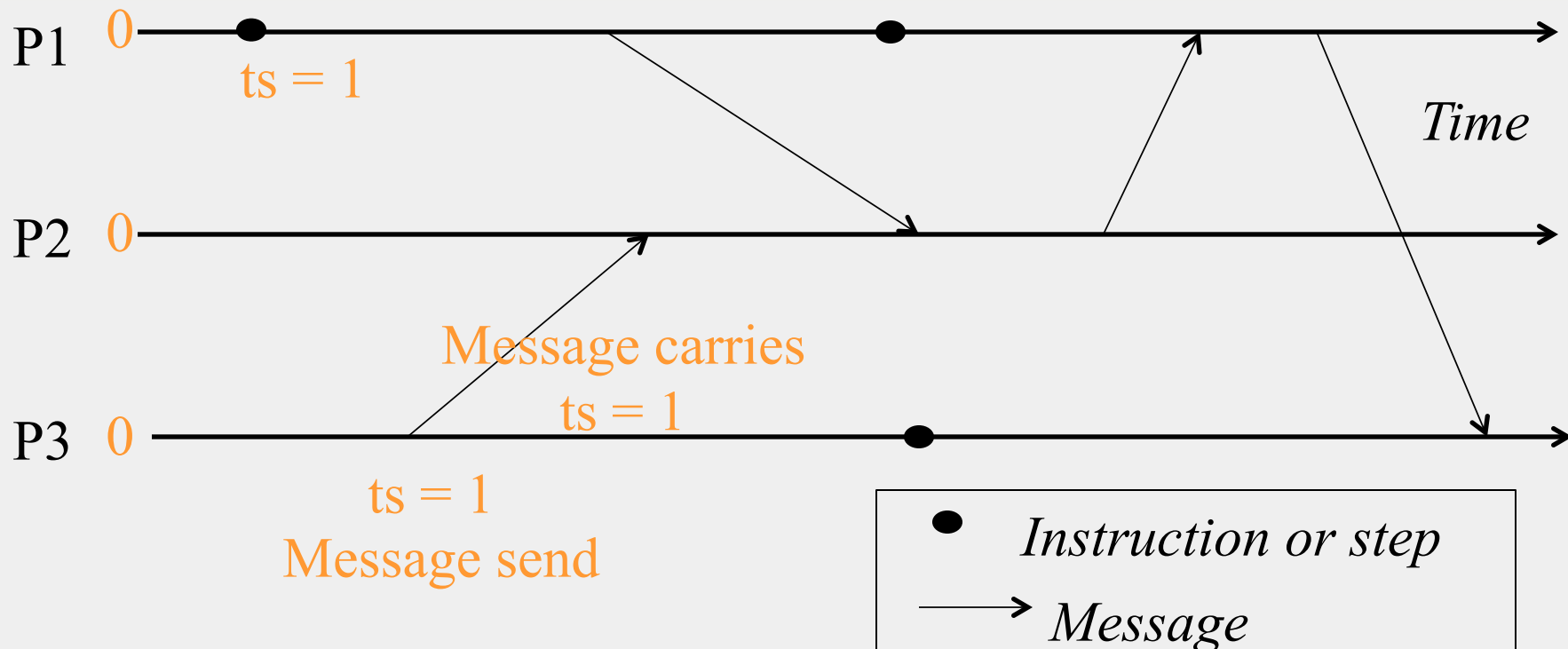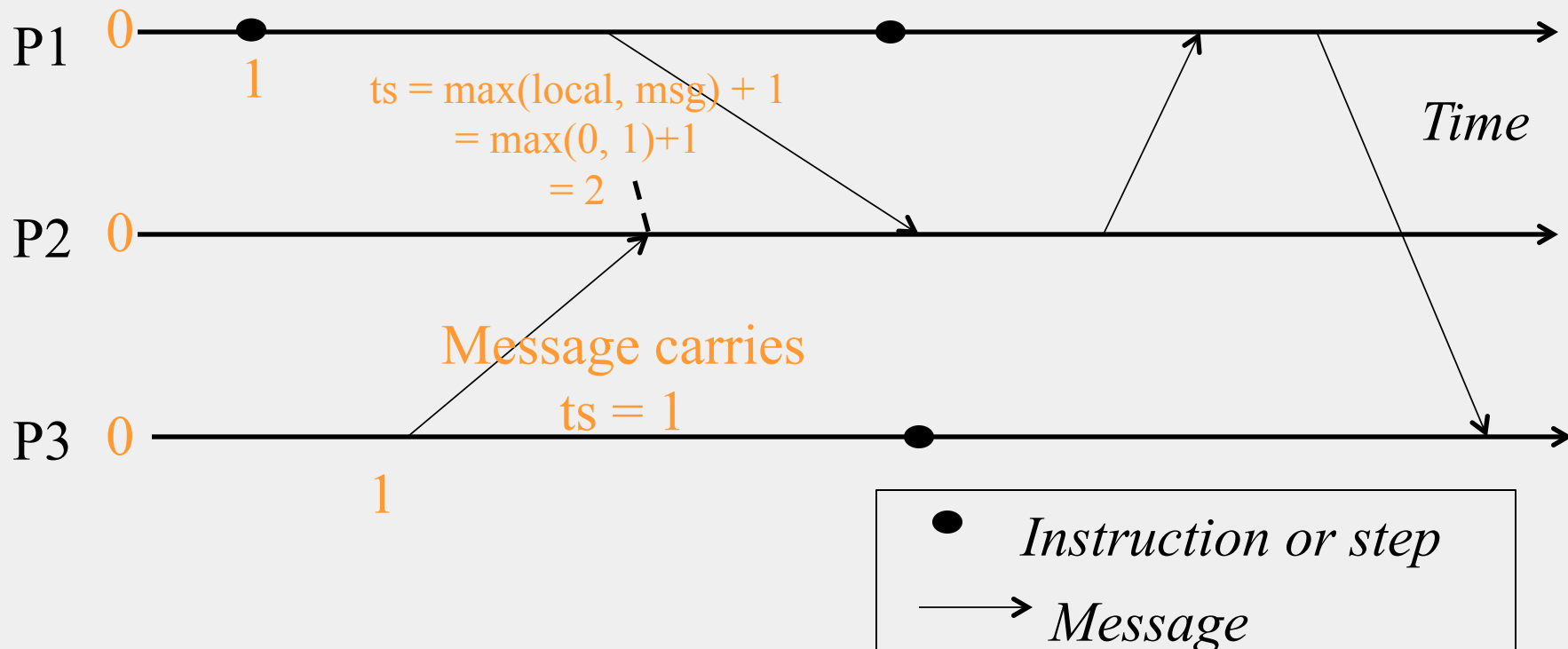
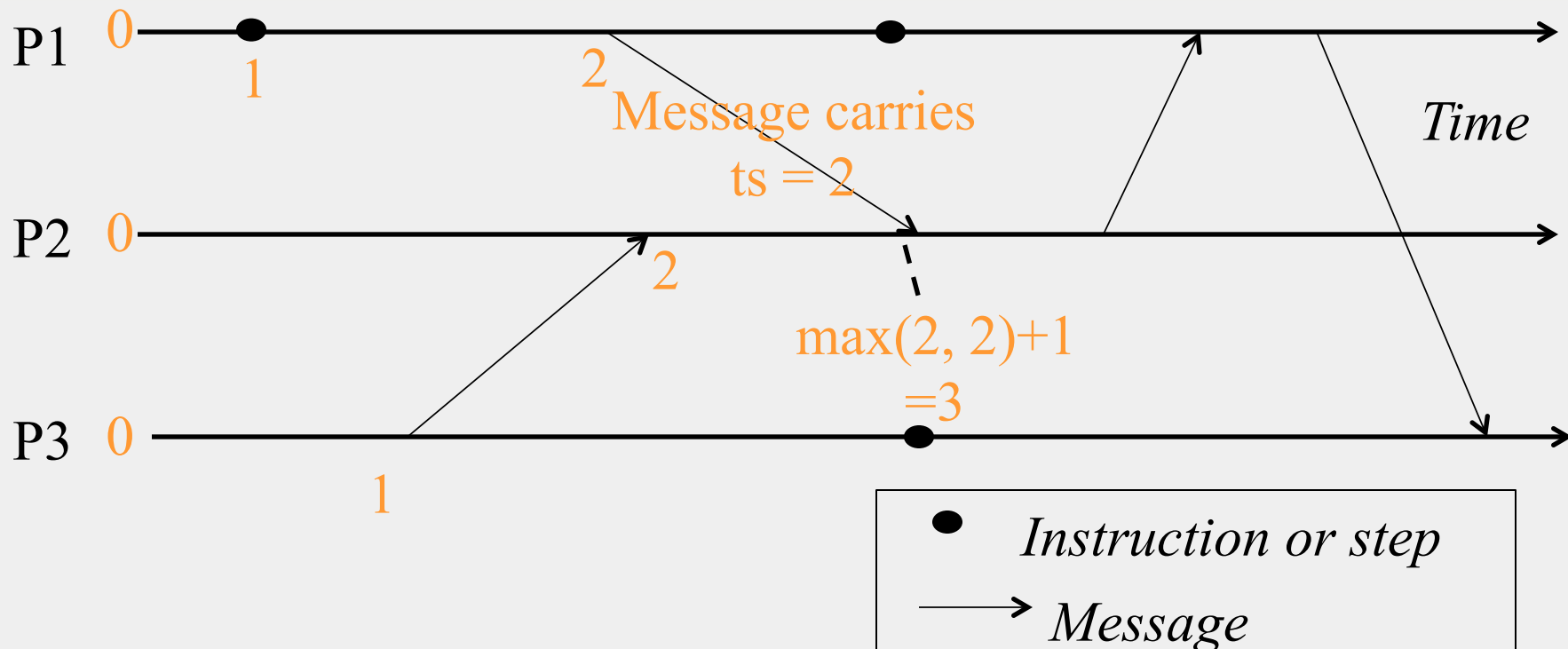$$max(local\ clock,\ message\ timestamp) + 1$$

# EXAMPLE



P1

P2

P3

*Time*

| | |
|---|---|
| ● | *Instruction or step* |
| → | *Message* |

33

# LAMPORT TIMESTAMPS

P1  0

P2  0

P3  0

*Time*

Initial counters (clocks)

● *Instruction or step*

→ *Message*

# Lamport Timestamps



P1  0

ts = 1

Time

P2  0

Message carries
ts = 1

P3  0

ts = 1
Message send

Instruction or step

Message

# Lamport Timestamps



P1  0

1

$ts = max(local, msg) + 1$
$= max(0, 1)+1$
$= 2$

P2  0

*Time*

Message carries
$ts = 1$

P3  0

1

●    *Instruction or step*

→    *Message*

# LAMPORT TIMESTAMPS



P1    0    1    2    Message carries    ts = 2    Time

P2    0    2    max(2, 2)+1 =3

P3    0    1

● Instruction or step
→ Message

# LAMPORT TIMESTAMPS

# Lamport Timestamps

# Obeying Causality



**P1**   0     A     B     C     D     E

1    2    3    5    6

*Time*

**P2**   0     E'     F     G

2    3    4

**P3**   0     H     I     J

1    2    7

- A → B :: 1 < 2
- B → F :: 2 < 3
- A → F :: 1 < 3

*Instruction or step*

→ *Message*

# Obeying Causality (2)



- H ➜ G :: 1 < 4
- F ➜ J  :: 3 < 7
- H ➜ J :: 1 < 7
- C ➜ J :: 3 < 7

# NOT ALWAYS *IMPLYING* CAUSALITY



- ? C → F ? :: 3 = 3
- ? H → C ? :: 1 < 3
- (C, F) and (H, C) are pairs of *concurrent* events

Instruction or step

→ Message

# CONCURRENT EVENTS

- **A pair of concurrent events doesn't have a causal path from one event to another (either way, in the pair)**

- **Lamport timestamps not guaranteed to be ordered or unequal for concurrent events**

- **Ok, since concurrent events are not causality related!**

- **Remember**

E1 → E2 ⟹ timestamp(E1) < timestamp (E2), BUT

timestamp(E1) < timestamp (E2) ⟹

{E1 → E2} OR {E1 and E2 concurrent}

# Next

- Can we have causal or logical timestamps from which we can tell if two events are concurrent or causally related?

# Vector Timestamps

- Used in key-value stores like Riak
- Each process uses a vector of integer clocks
- Suppose there are N processes in the group 1…N
- Each vector has N elements
- Process *i maintains vector* $V_i[1…N]$
- *j*th element of vector clock at process *i,* $V_i[j]$, is *i*'s knowledge of latest events at process *j*
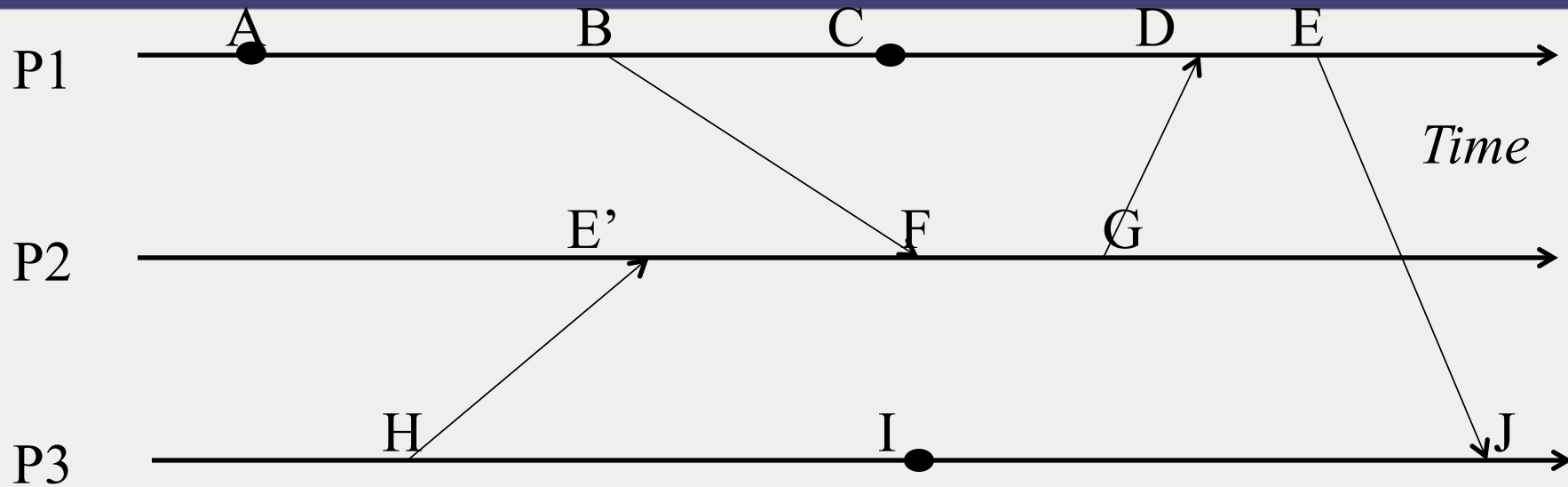
# Assigning Vector Timestamps

- Incrementing vector clocks
1. On an instruction or send event at process $i$, it increments only its $i$th element of its vector clock
2. Each message carries the send-event's vector timestamp $V_{message}[1...N]$
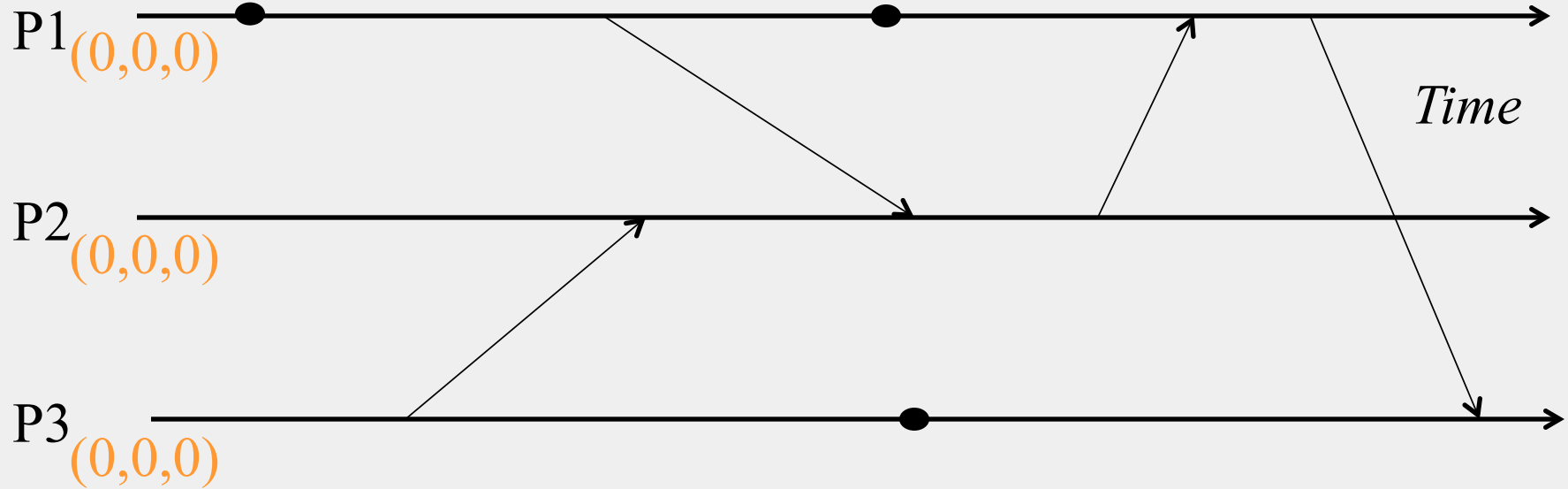3. On receiving a message at process $i$:

   $V_i[i] = V_i[i] + 1$

   $V_i[j] = max(V_{message}[j], V_i[j])$ for $j \neq i$

# EXAMPLE

# Vector Timestamps

P1 (0,0,0)

P2 (0,0,0)

P3 (0,0,0)

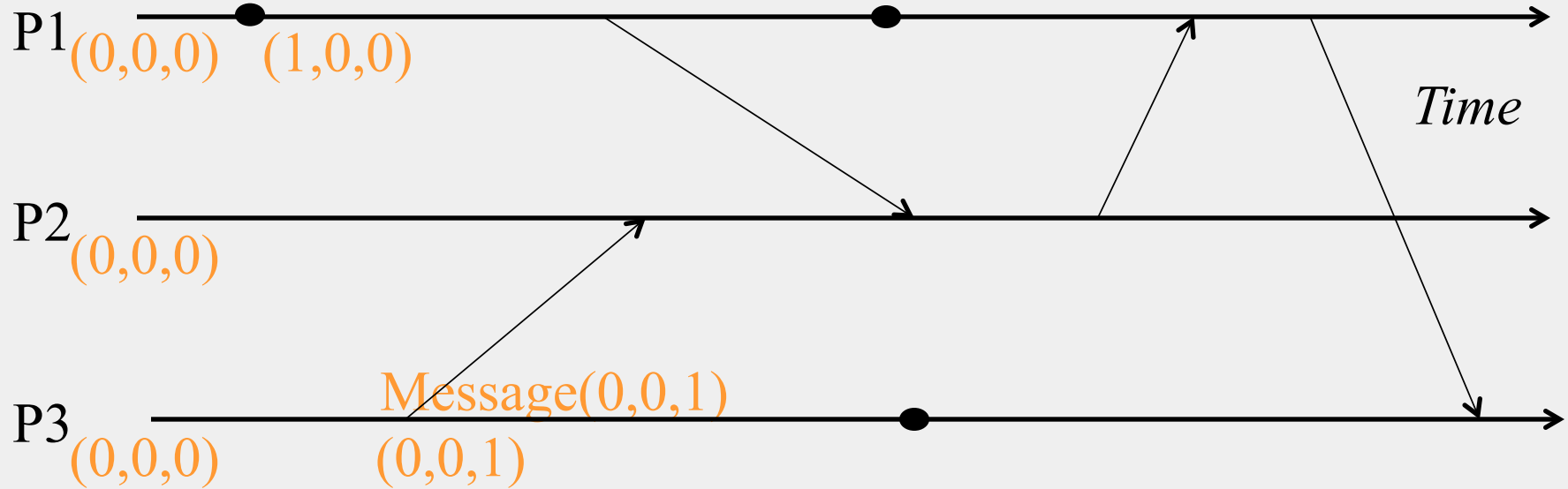*Time*

Initial counters (clocks)

# Vector Timestamps



P1 (0,0,0)  (1,0,0)

*Time*

P2 (0,0,0)

P3 (0,0,0)

Message(0,0,1)

(0,0,1)

# Vector Timestamps



P1 (0,0,0)  (1,0,0)

*Time*

P2 (0,0,0)  (0,1,1)

Message(0,0,1)

P3 (0,0,0)  (0,0,1)

# Vector Timestamps



P1   (0,0,0)   (1,0,0)   (2,0,0)

Message(2,0,0)

P2   (0,0,0)   (0,1,1)   (2,2,1)

P3   (0,0,0)   (0,0,1)

*Time*

# Vector Timestamps

P1 (0,0,0) (1,0,0) (2,0,0) (3,0,0) (4,3,1) (5,3,1)

*Time*

P2 (0,0,0) (0,1,1) (2,2,1) (2,3,1)

P3 (0,0,0) (0,0,1) (0,0,2) (5,3,3)

# Causally-Related ...

- $VT_1 = VT_2,$
  - *iff* (if and only if)
    - $VT_1[i] = VT_2[i]$, for all $i = 1, \ldots, N$
- $VT_1 \leq VT_2,$
  - *iff* $VT_1[i] \leq VT_2[i]$, for all $i = 1, \ldots, N$
- Two events are causally related *iff*
  - $VT_1 < VT_2,$ i.e.,
    - *iff* $VT_1 \leq VT_2$ &
      - there exists $j$ such that
        - $1 \leq j \leq N$ & $VT_1[j] < VT_2[j]$
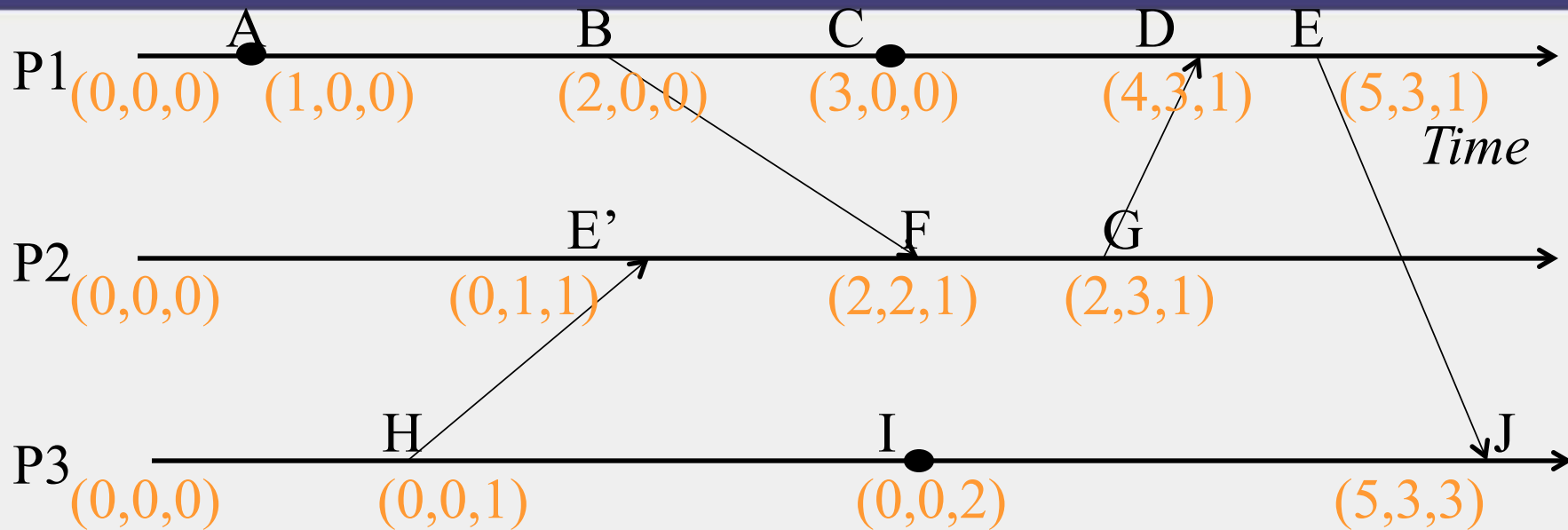
# … OR NOT CAUSALLY-RELATED

- Two events $VT_1$ and $VT_2$ are <span style="color:orange">concurrent</span>

  *iff*

    NOT $(VT_1 \leq VT_2)$  AND NOT $(VT_2 \leq VT_1)$

    We'll denote this as $VT_2 \parallel\!\parallel VT_1$

# Obeying Causality

P1 (0,0,0) — A (1,0,0) — B (2,0,0) — C (3,0,0) — D (4,3,1) — E (5,3,1)

*Time*

P2 (0,0,0) — E' (0,1,1) — F (2,2,1) — G (2,3,1)

P3 (0,0,0) — H (0,0,1) — I (0,0,2) — J (5,3,3)

- A → B :: (1,0,0) < (2,0,0)
- B → F :: (2,0,0) < (2,2,1)
- A → F :: (1,0,0) < (2,2,1)

# Obeying Causality (2)



P1  A (1,0,0)   B (2,0,0)   C (3,0,0)   D (4,3,1)   E (5,3,1)
(0,0,0)

*Time*

P2  E' (0,1,1)   F (2,2,1)   G (2,3,1)
(0,0,0)

P3  H (0,0,1)   I (0,0,2)   J (5,3,3)
(0,0,0)

- H → G :: (0,0,1) < (2,3,1)
- F → J  :: (2,2,1) < (5,3,3)
- H → J :: (0,0,1) < (5,3,3)
- C → J :: (3,0,0) < (5,3,3)

# IDENTIFYING CONCURRENT EVENTS



P1
A (1,0,0)   B (2,0,0)   C (3,0,0)   D (4,3,1)   E (5,3,1)
(0,0,0)

*Time*

P2
E' (0,1,1)   F (2,2,1)   G (2,3,1)
(0,0,0)

P3
H (0,0,1)   I (0,0,2)   J (5,3,3)
(0,0,0)

- C & F :: (3,0,0) ||| (2,2,1)
- H & C :: (0,0,1) ||| (3,0,0)
- (C, F) and (H, C) are pairs of *concurrent* events

# Logical Timestamps: Summary

- **Lamport timestamps**
  - Integer clocks assigned to events
  - Obeys causality
  - Cannot distinguish concurrent events
- **Vector timestamps**
  - Obey causality
  - By using more space, can also identify concurrent events

# Time and Ordering: Summary

- **Clocks are unsynchronized in an asynchronous distributed system**
- **But need to order events, across processes!**
- **Time synchronization**
  - Cristian's algorithm
  - NTP
  - Berkeley algorithm
  - But error a function of round-trip-time

- **Can avoid time sync altogether by instead assigning logical timestamps to events**

# HW1 Statistics

(min, max, median, average, stdev)

- On-campus Undergrads:   0, 80, 71, 65.1, 19.6
- On-campus Grads:   0, 80, 74, 69.7, 16.4
- MCS-online:   0, 80, 74.5, 59.9, 31.9
- MCS-DS:   0, 80, 63, 56.2, 23.1

# REMINDERS

- (4 cr students) MP2 due this Sunday, Demos on Monday

  – Signup sheet soon on Piazza

- (All) HW2 due next Tuesday