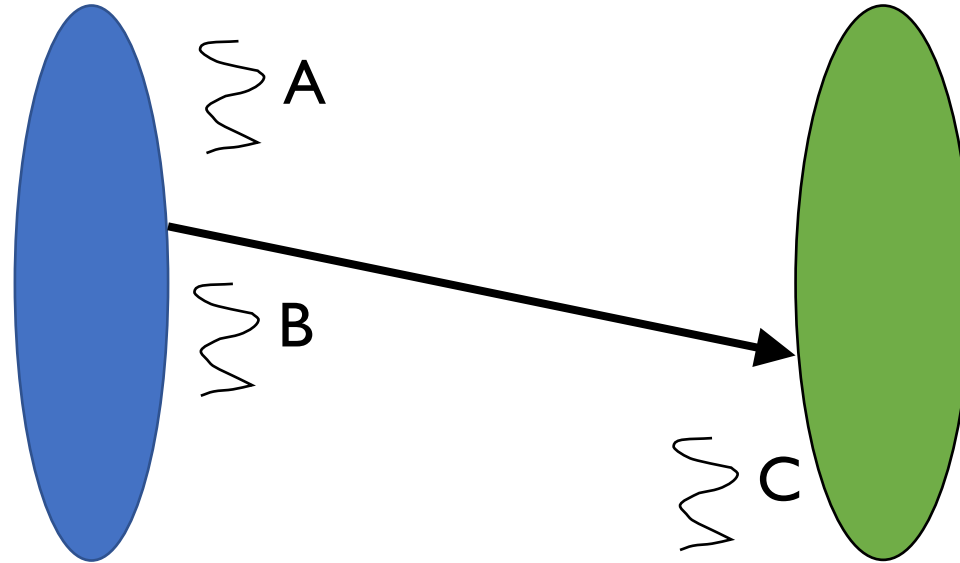# Distributed Systems

## CS425/ECE428

*Instructor: Radhika Mittal*

# While we wait…



Local processing event 〜〜

Message Delivery →

- Can we conclude that event A occurred before event C?
- Can we conclude that event B occurred before event C?

The clocks of blue and green processes cannot be perfectly synchronized. Can we simply compare timestamps of these events?

# Logistics Related

- VM clusters have been assigned!

- Newly registered students:
  - Please make sure you have access to Campuswire and Gradescope
  - If you are in 4 credits, make sure you have been allocated a VM cluster for the MPs.
  - Email Neel (netid: neeld2) to get the required access.

- *Can you please say your name before speaking up in class?*

# Today's agenda

- Logical Clocks and Timestamps
    - Chapter 14.4

- Global State (if time)
    - Chapter 14.5

# Event Ordering

- A usecase of synchronized clocks:
    - Reasoning about order of events.

- Why is it useful?
    - Debugging distributed applications
    - Reconciling updates made to an object in a distributed datastore.
    - Rollback recovery during failures:
        1. *Checkpoint state of the system; 2. Log events (with timestamps); 3. Rollback to checkpoint and replay events in order if system crashes.*

    - ….

- Can we reason about order of events without synchronized clocks?

# Process, state, events

- Consider a system with **n** processes: $<p_1, p_2, p_3, …., p_n>$

- Each process $p_i$ is described by its *state* $s_i$ that gets transformed over time.
  - State includes values of all local variables, affected files, etc.

- $s_i$ gets transformed when an *event* occurs.

- Three types of events:
  - Local computation.
  - Sending a message.
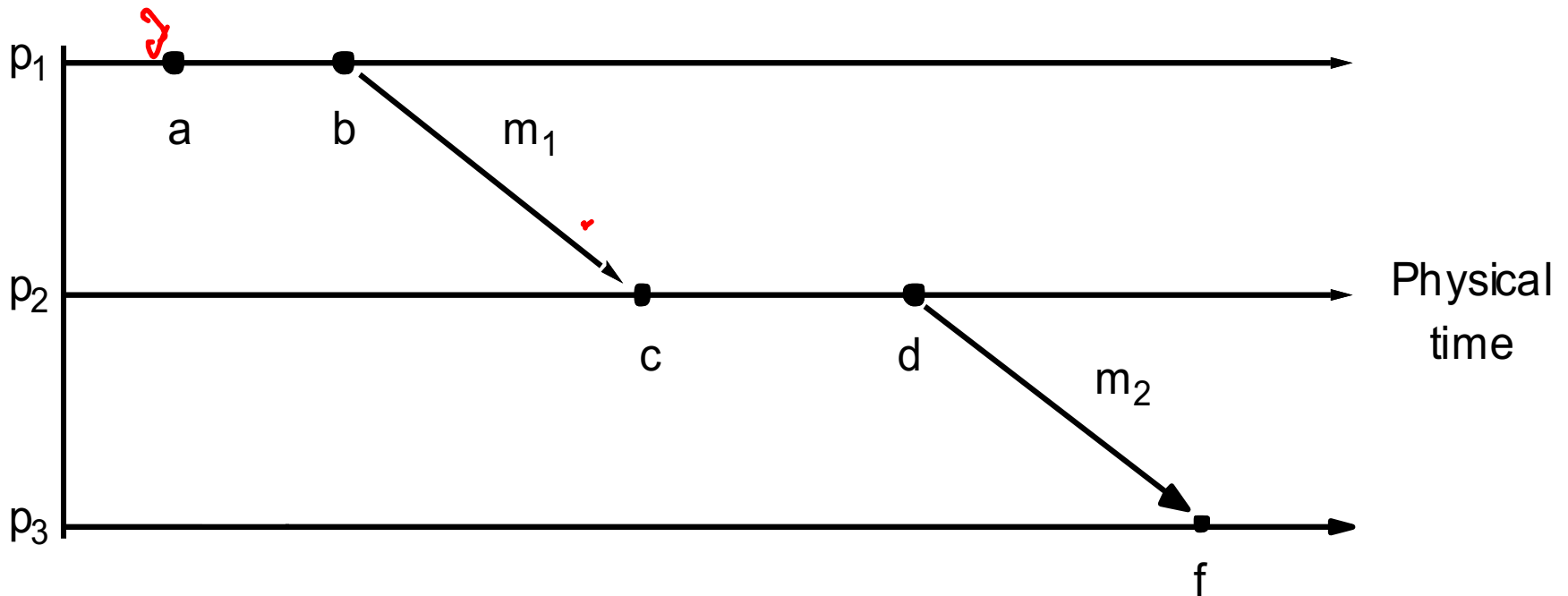  - Receiving a message.

# Event Ordering

- Easy to order events within a single process $p_i$, based on their time of occurrence.

- How do we reason about events across processes?
  - A message must be *sent* before it gets *received* at another process.

- These two notions help define *happened-before* (HB) relationship denoted by →.
  - **e** → **e'** means **e** *happened before* **e'**.

# Happened-Before Relationship

- *Happened-before* (HB) relationship denoted by $\rightarrow$.
  - $e \rightarrow e'$ means $e$ *happened before* $e'$.
  - $e \rightarrow_i e'$ means $e$ *happened before* $e'$, as observed by $p_i$.

- HB rules:
  - If $\exists\ p_i$ , $e \rightarrow_i e'$ then $e \rightarrow e'$.
  - For any message m, **send(m)** $\rightarrow$ **receive(m)**
  - If $e \rightarrow e'$ and $e' \rightarrow e''$ then $e \rightarrow e''$

$$p_i \rightarrow p_j$$

- Also called *"causal"* or *"potentially causal"* ordering.
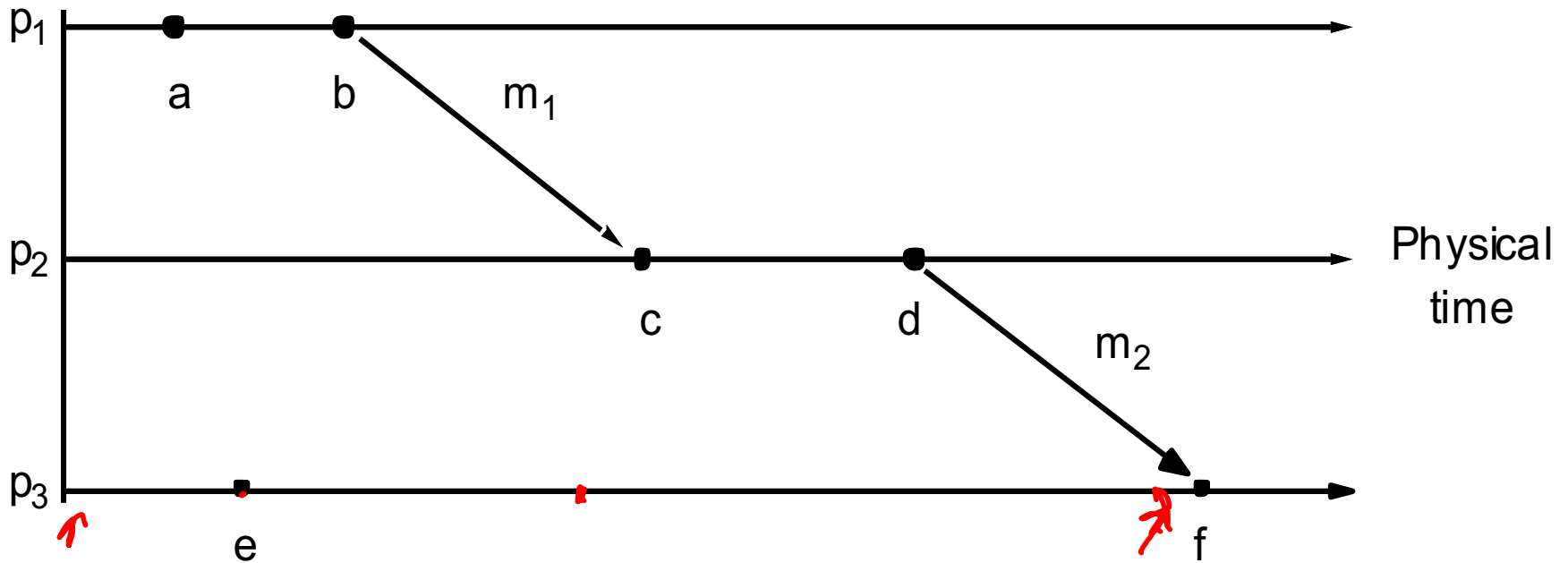
# Event Ordering: Example



Which event happened first?

$a \rightarrow b$ and $b \rightarrow c$ and $c \rightarrow d$ and $d \rightarrow f$

$a \rightarrow b$ and $a \rightarrow c$ and $a \rightarrow d$ and $a \rightarrow f$

# Event Ordering: Example



$p_1$

a     b     $m_1$

$p_2$                                 Physical
time

c     d

$m_2$

$p_3$

e                f

What can we say about **e**?

$e \rightarrow f$

$e \| b, \ e \| c, \ e \| d$

$a \nrightarrow e$ and $e \nrightarrow a$

$a \| e$

**a** and **e** are *concurrent*.

# Event Ordering: Example



What can we say about **e** and **d**?
e || d

# Logical Timestamps: Example



$p_1$

a       b      $m_1$

**h**

$p_2$      Physical time
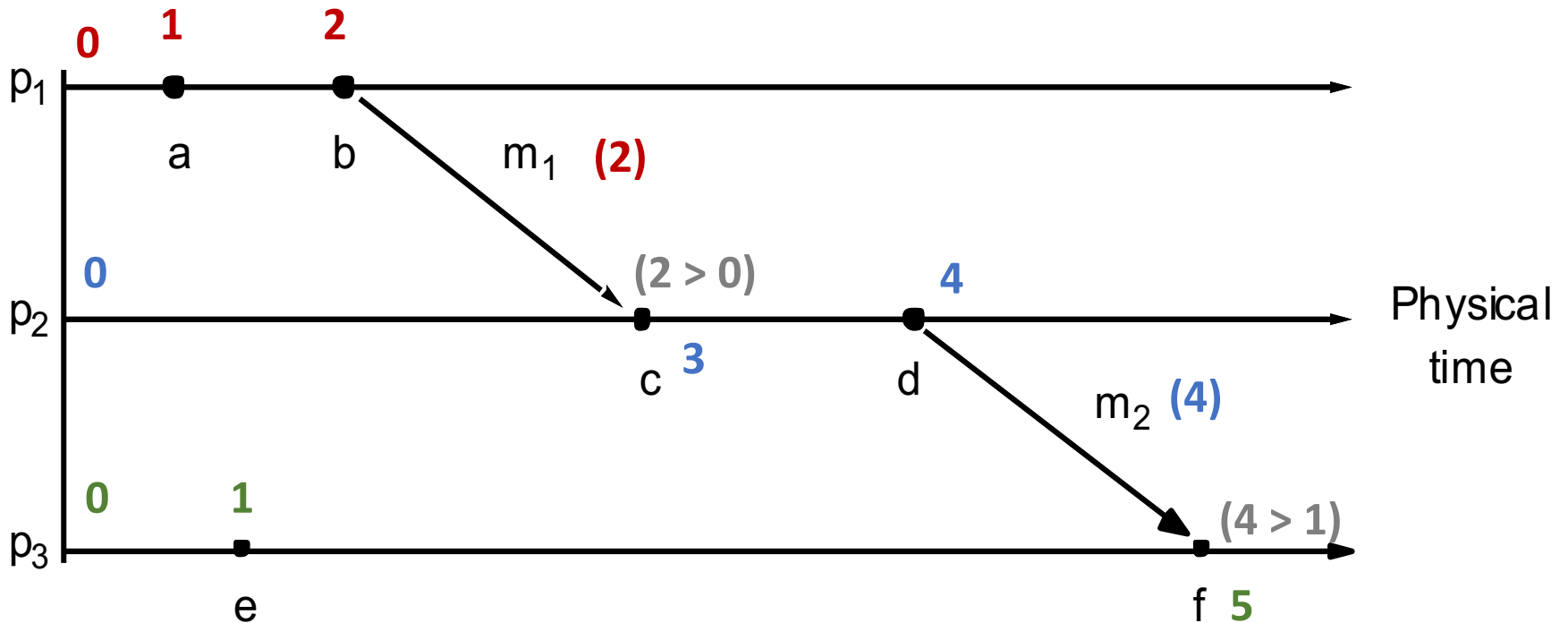
c      d    $m_2$

$p_3$

e      **g**      f

What can we say about **e** and **d**?

e → d

# Lamport's Logical Clock

- Logical timestamp for each event that captures the *happened-before* relationship.

- *Algorithm:* Each process $p_i$
    1. initializes local clock $L_i = 0$.
    2. increments $L_i$ before timestamping each event.
    3. piggybacks $L_i$ when sending a message.
        - (i.e. sends $L_i$ along with the message)
    4. upon receiving a message with clock value $t$
        - sets $L_i = max(t, L_i)$
        - increments $L_i$ before timestamping the receive event (as per step 2).
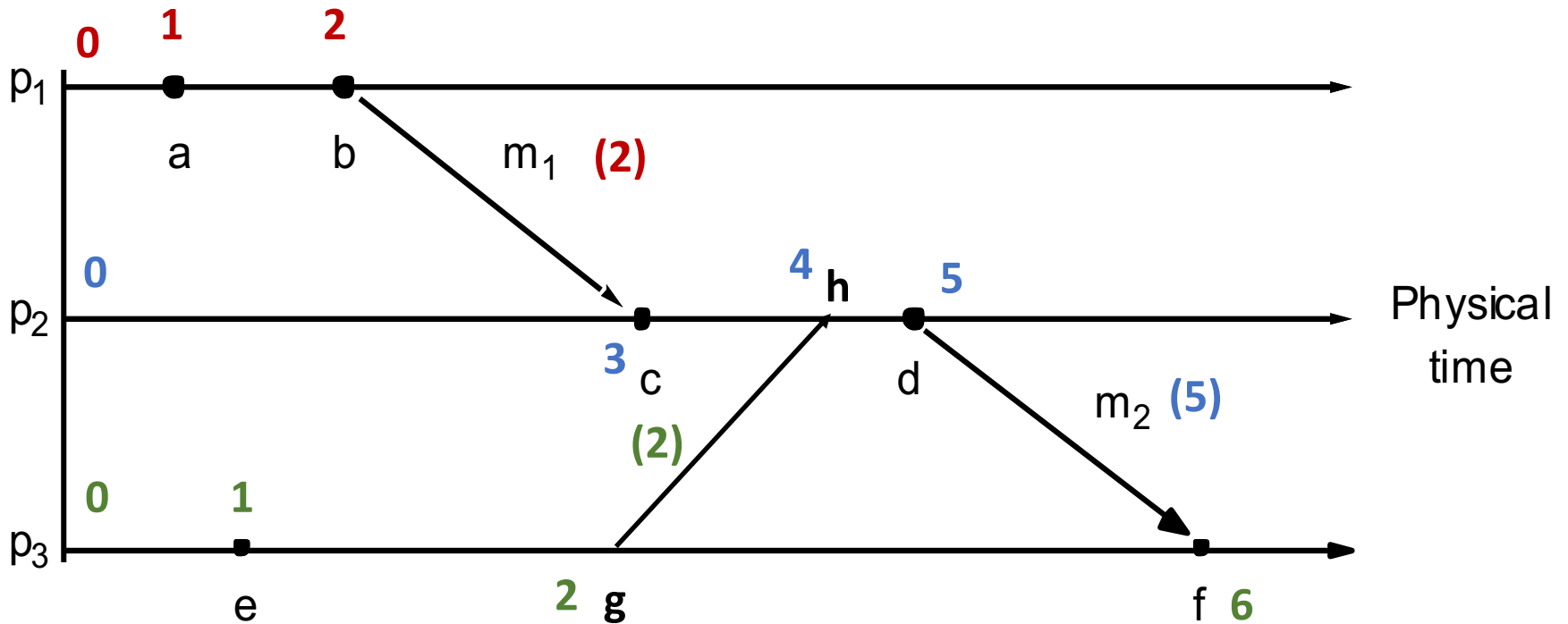
# Logical Timestamps: Example

# Lamport's Logical Clock

- Logical timestamp for each event that captures the *happened-before* relationship.
- *Algorithm:* Each process $p_i$
    1. initializes local clock $L_i = 0$.
    2. increments $L_i$ before timestamping each event.
    3. piggybacks $L_i$ when sending a message.
        - (i.e. sends $L_i$ along with the message)
    4. upon receiving a message with clock value $t$
        - sets $L_i = \max(t, L_i)$
        - increments $L_i$ before timestamping the receive event (as per step 2).
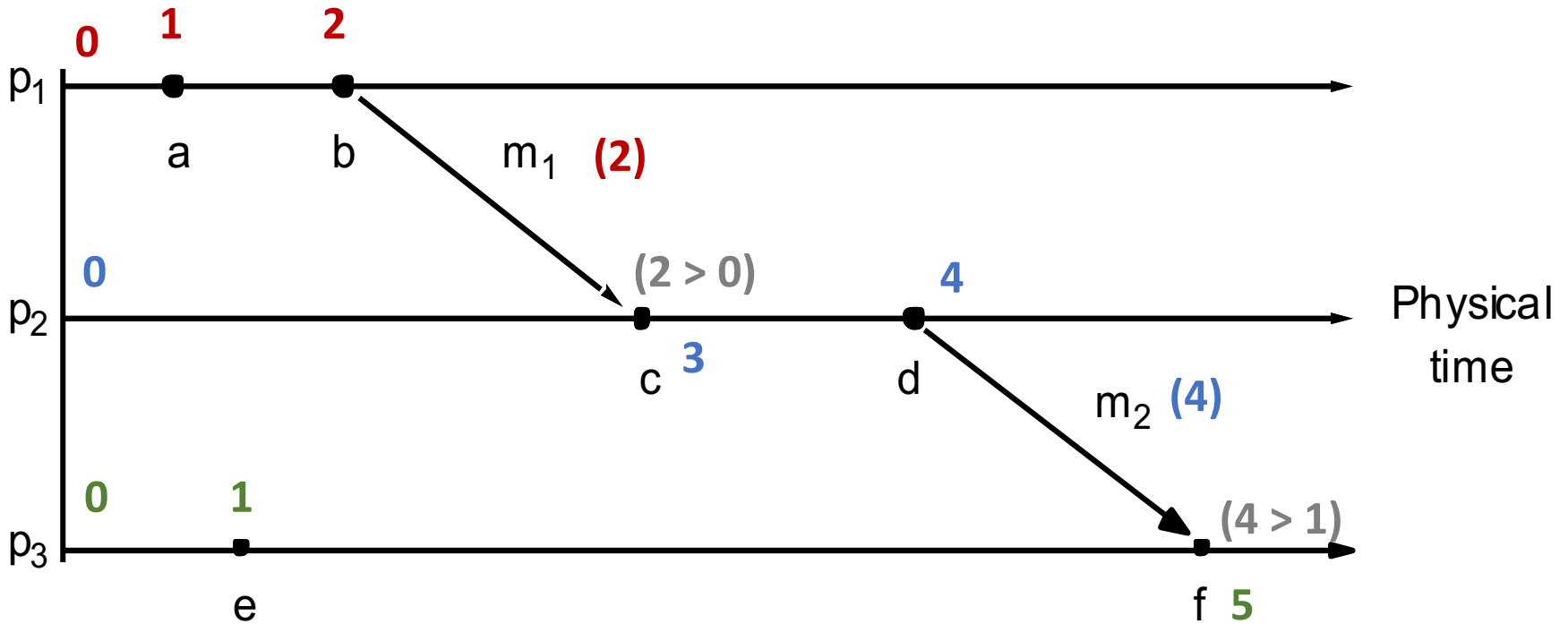
# Logical Timestamps: Example

# Lamport's Logical Clock

- Logical timestamp for each event that captures the *happened-before* relationship.

- If **e → e'** then
  - **L(e) < L(e')**

- What can we conclude if **L(e) < L(e')**?
  - We cannot say that **e → e'**
  - We can say: **e' ↛ e**
  - Either **e → e'** or **e || e'**

# Logical Timestamps: Example

$e \to e'$    $L(e) < L(e')$

$L(e) = L(e')$
$e(e) \& L(e')$
$e \not\to e'$
$e' \not\to e$
$e || e'$



L(e) < L(d), e || d          L(e) < L(f), e → f

# Vector Clocks

- Each event associated with a vector timestamp.

- Each process $p_i$ maintains vector of clocks $V_i$

- The size of this vector is the same as the no. of processes.
    - $V_i[j]$ is the clock for process $p_j$ as maintained by $p_i$

- Algorithm: each process $p_i$:

# Vector Clocks

- Each event associated with a vector timestamp.

- Each process $p_i$ maintains vector of clocks $V_i$

- The size of this vector is the same as the no. of processes.
    - $V_i[j]$ is the clock for process $p_j$ as maintained by $p_i$

- Algorithm: each process $p_i$:
    1. initializes local clock $V_i[j] = 0$

# Vector Clocks

- Each event associated with a vector timestamp.

- Each process $p_i$ maintains vector of clocks $V_i$

- The size of this vector is the same as the no. of processes.
    - $V_i[j]$ is the clock for process $p_j$ as maintained by $p_i$

- Algorithm: each process $p_i$:
    1. initializes local clock $V_i[j] = 0$
    2. increments $V_i[i]$ before timestamping each event.

# Vector Clocks

- Each event associated with a vector timestamp.

- Each process $p_i$ maintains vector of clocks $V_i$

- The size of this vector is the same as the no. of processes.
    - $V_i[j]$ is the clock for process $p_j$ as maintained by $p_i$

- Algorithm: each process $p_i$:
    1. initializes local clock $V_i[j] = 0$
    2. increments $V_i[i]$ before timestamping each event.
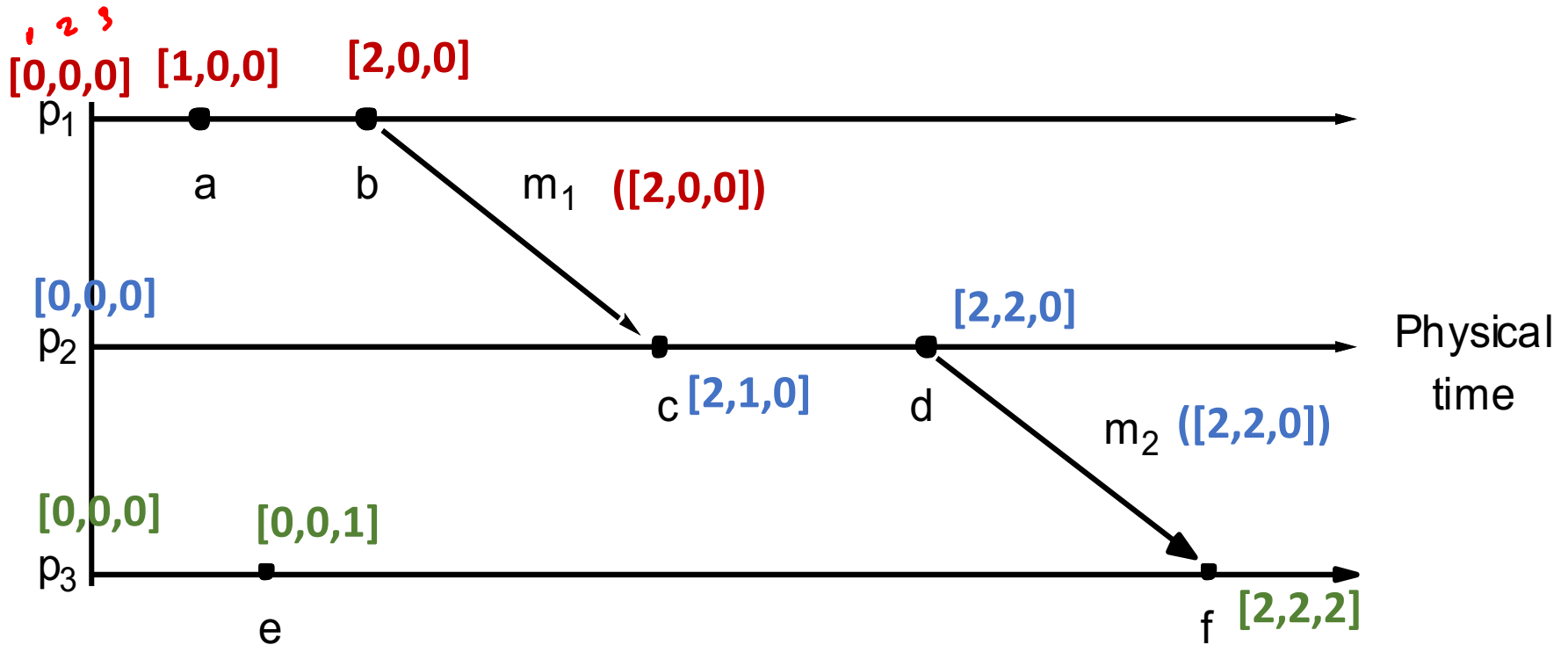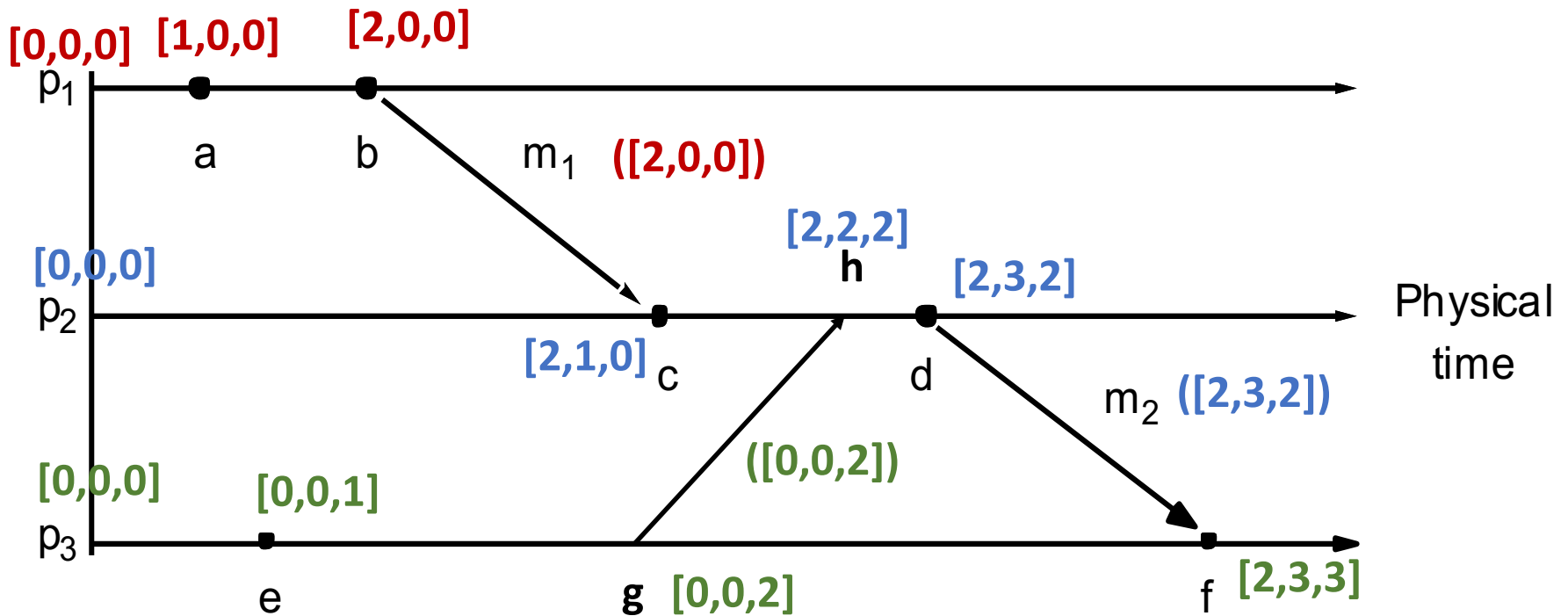    3. piggybacks $V_i$ when sending a message.

# Vector Clocks

- Each event associated with a vector timestamp.

- Each process $p_i$ maintains vector of clocks $V_i$

- The size of this vector is the same as the no. of processes.
    - $V_i[j]$ is the clock for process $p_j$ as maintained by $p_i$

- Algorithm: each process $p_i$:
    1. initializes local clock $V_i[j] = 0$
    2. increments $V_i[i]$ before timestamping each event.
    3. piggybacks $V_i$ when sending a message.
    4. upon receiving a message with vector clock value $v$
        - sets $V_i[j] = \max(V_i[j], v[j])$ for all $j=1\ldots n$.
        - increments $V_i[i]$ before timestamping receive event (as per step 2).

# Vector Timestamps: Example



, 2 3

[0,0,0] [1,0,0] [2,0,0]

$p_1$

a          b

$m_1$  ([2,0,0])

[0,0,0]

[2,2,0]

$p_2$                                          Physical
                                               time

c [2,1,0]          d

$m_2$  ([2,2,0])

[0,0,0]  [0,0,1]

$p_3$

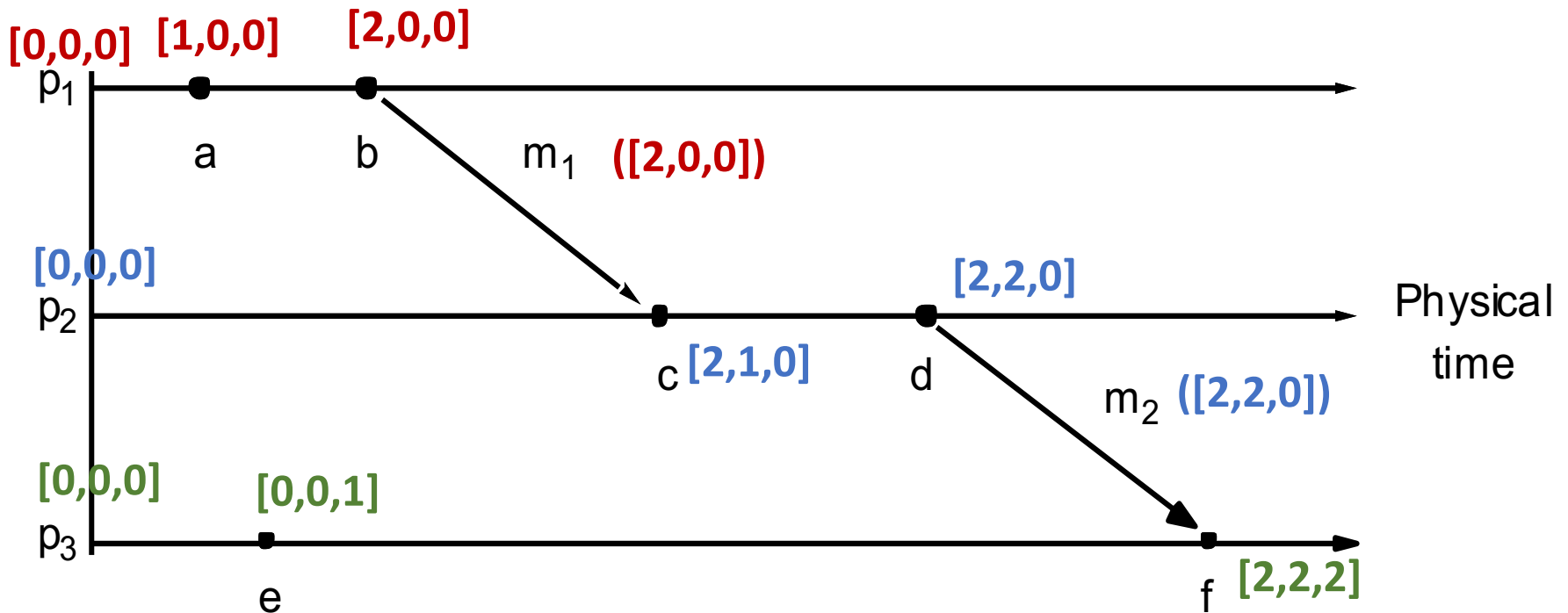e                                          f [2,2,2]

# Vector Timestamps: Example
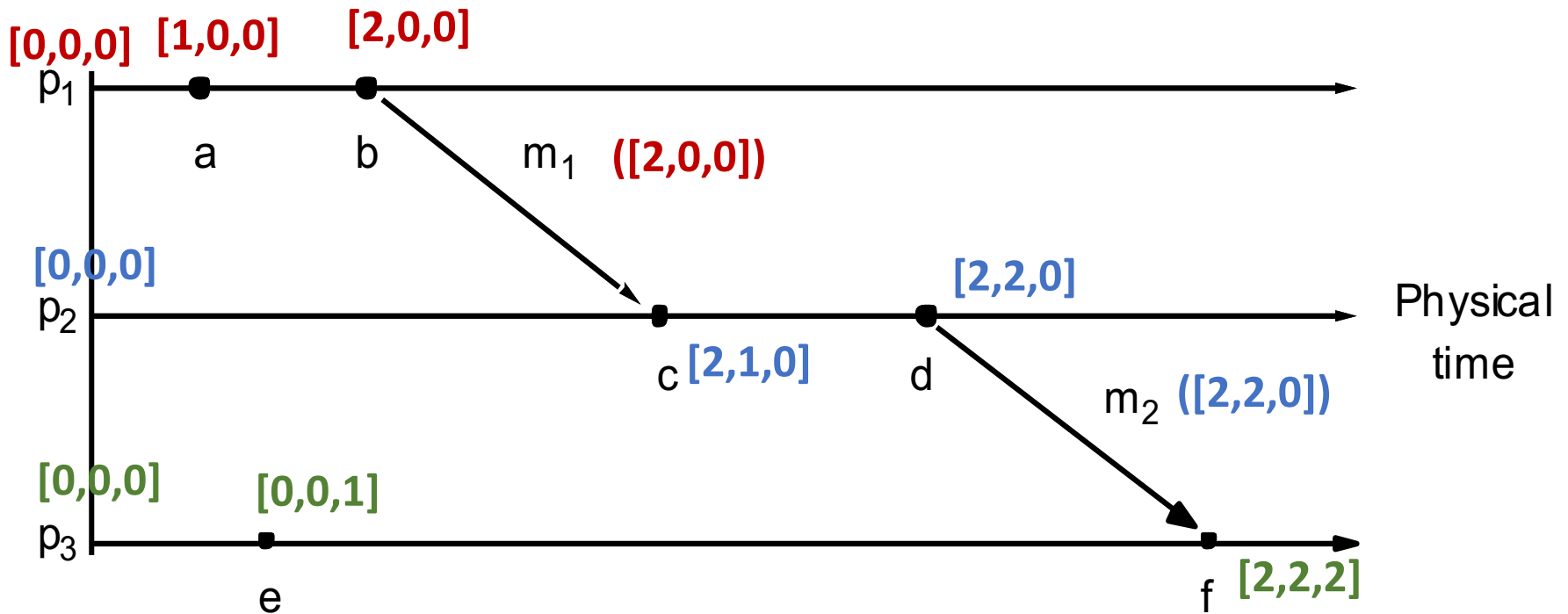
# Comparing Vector Timestamps

- Let  $V(e) = V$  and  $V(e') = V'$

- $V = V'$,  iff  $V[i] = V'[i]$, for all $i = 1, \ldots, n$
- $V \leq V'$,  iff  $V[i] \leq V'[i]$, for all $i = 1, \ldots, n$
- $V < V'$,  iff  $V \leq V'$ & $V \neq V'$

  iff  $V \leq V'$ & $\exists \, j$ such that $(V[j] < V'[j])$

- $e \rightarrow e'$ iff $V < V'$
  - ($V < V'$ implies $e \rightarrow e'$) and ($e \rightarrow e'$ implies $V < V'$)
- $e \, || \, e'$ iff ($V \nless V'$ and $V' \nless V$)

# Vector Timestamps: Example



What can we say about e & f based on their vector timestamps?

# Vector Timestamps: Example



$V(e) < V(f), e \rightarrow f$

# Vector Timestamps: Example



**[0,0,0]** **[1,0,0]** **[2,0,0]**

$p_1$

a     b     $m_1$ **([2,0,0])**

**[0,0,0]**

$p_2$

    c **[2,1,0]**     d     **[2,2,0]**     $m_2$ **([2,2,0])**

**[0,0,0]** **[0,0,1]**

$p_3$

    e     f **[2,2,2]**

Physical time
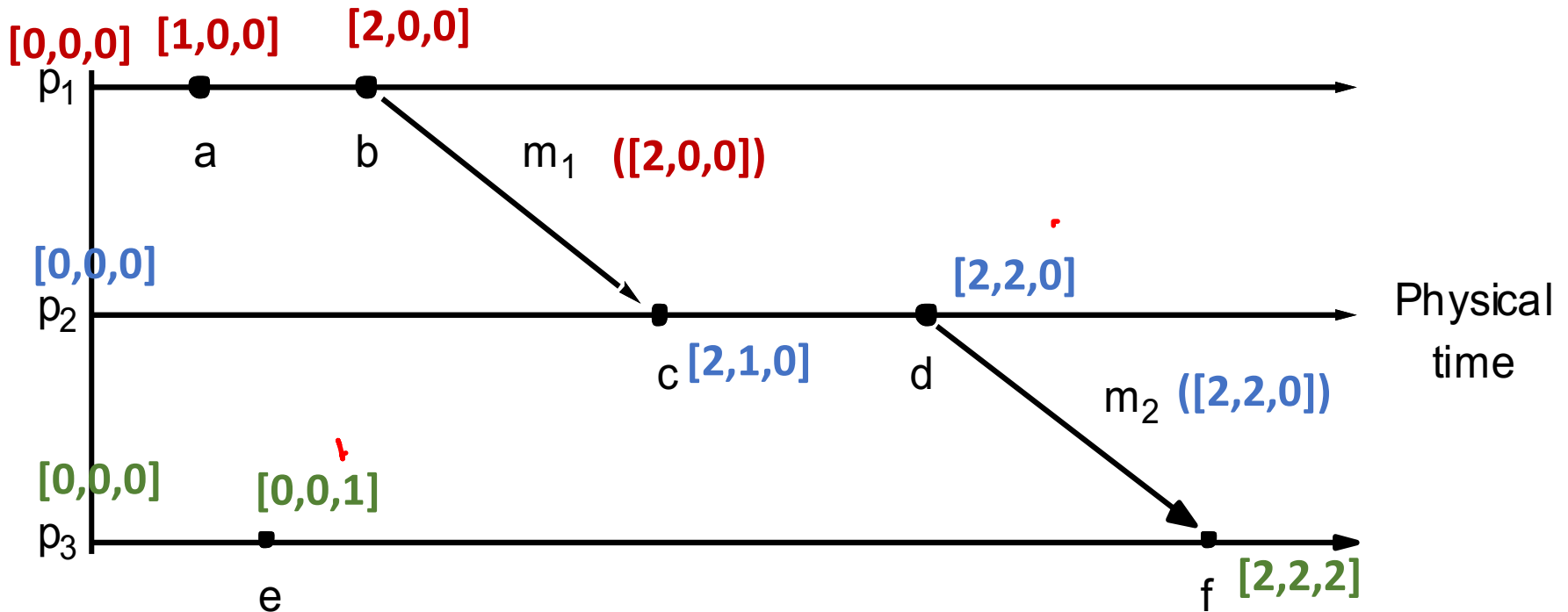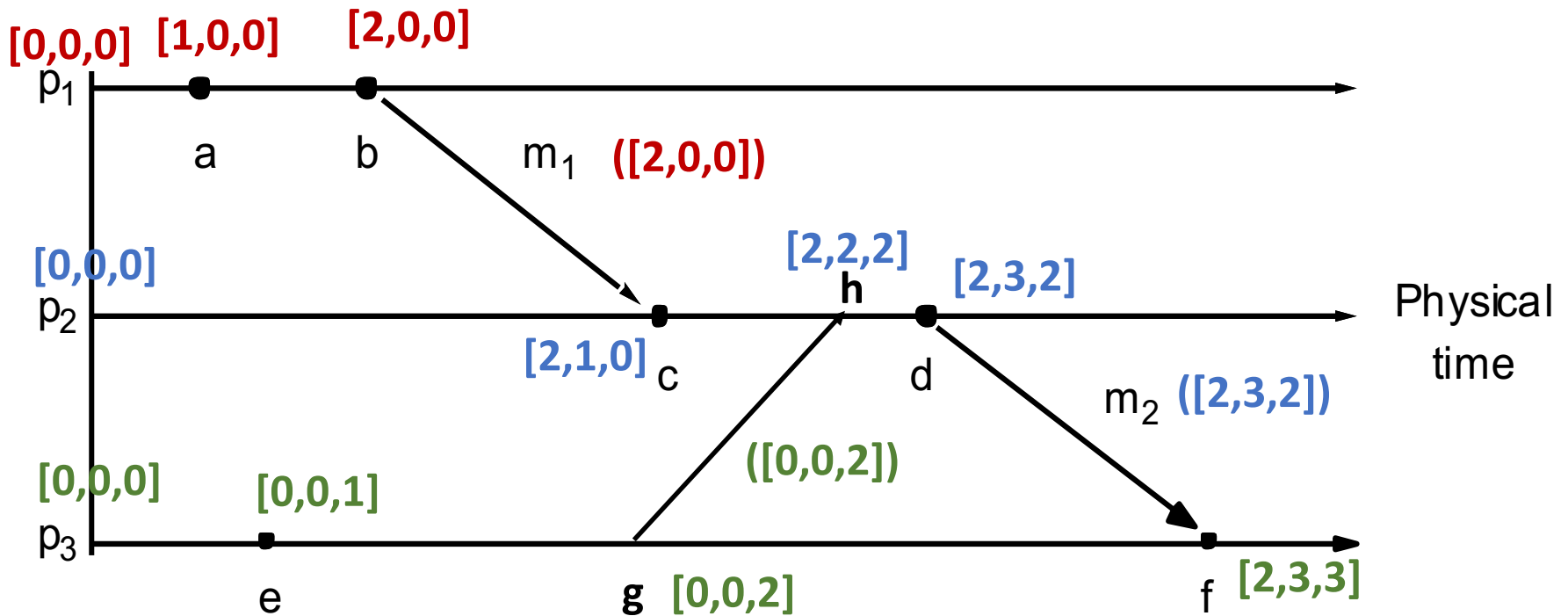
What can we say about e & d based on their vector timestamps?

# Vector Timestamps: Example



$V(e) \not< V(d)$ and $V(d) \not< V(e)$, e || d

# Vector Timestamps: Example



How about now?

# Vector Timestamps: Example



$[0,0,0]$ $[1,0,0]$ $[2,0,0]$

$p_1$

a  b  $m_1$ ([2,0,0])

$[0,0,0]$

$[2,2,2]$ **h** $[2,3,2]$

$p_2$

Physical time

$[2,1,0]$ c  d  $m_2$ ([2,3,2])

([0,0,2])

$[0,0,0]$ $[0,0,1]$

$p_3$

e  **g** $[0,0,2]$  f $[2,3,3]$

$V(e) < V(f), e \rightarrow f$
$V(e) < V(d), e \rightarrow d$

# Timestamps Summary

- Comparing timestamps across events is useful.
  - Reconciling updates made to an object in a distributed datastore.
  - Rollback recovery during failures:

    *1. Checkpoint state of the system; 2. Log events (with timestamps);*
    *3. Rollback to checkpoint and replay events in order if system crashes.*

- How to compare timestamps across different processes?
  - **Physical timestamp:** requires clock synchronization.
    - Google's Spanner Distributed Database uses "TrueTime".
  - **Lamport's timestamps:** cannot fully differentiate between causal and concurrent ordering of events.
    - Oracle uses "System Change Numbers" based on Lamport's clock.
  - **Vector timestamps:** larger message sizes.
    - Amazon's DynamoDB uses vector clocks.

# Timestamps Summary

- Comparing timestamps across events is useful.
  - Reconciling updates made to an object in a distributed datastore.
  - Rollback recovery during failures:
    - *1. Checkpoint state of the system; 2. Log events (with timestamps); 3. Rollback to checkpoint and replay events in order if system crashes.*

- How to compare timestamps across different processes?
  - **Physical timestamp:** requires clock synchronization.
    - Google's Spanner Distributed Database uses "TrueTime".
  - **Lamport's timestamps:** cannot fully differentiate between causal and concurrent ordering of events.
    - Oracle uses "System Change Numbers" based on Lamport's clock.
  - **Vector timestamps:** larger message sizes.
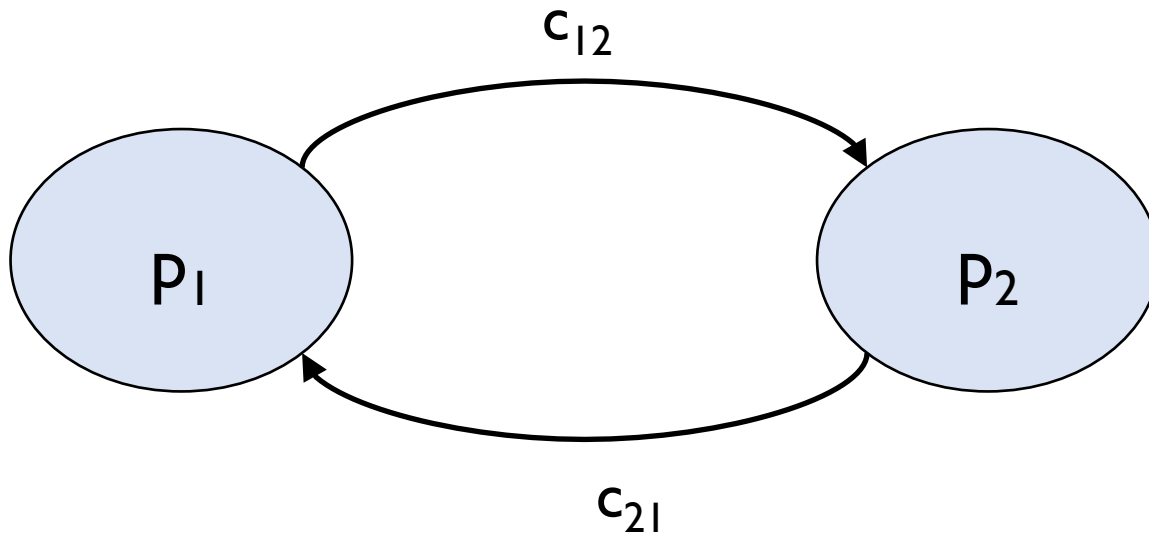    - Amazon's DynamoDB uses vector clocks.

# Today's agenda

- Logical Clocks and Timestamps
  - Chapter 14.4

- **Global State**
  - Chapter 14.5

# Process, state, events

- Consider a system with **n** processes: $<p_1, p_2, p_3, \ldots, p_n>$.

- Each process $p_i$ is associated with *state* $s_i$.
    - State includes values of all local variables, affected files, etc.

- Each channel can also be associated with a state.
    - Which messages are currently *pending* on the channel.
    - Can be computed from process' state:
        - Record when a process sends and receives messages.
        - if $p_i$ sends a message that $p_j$ has not yet received, it is pending on the channel.

- State of a process (or a channel) gets transformed when an *event* occurs. 3 types of events:
    - local computation, sending a message, receiving a message.

# Global State (or Global Snapshot)

- State of each process (and each channel) in the system at a given instant of time.
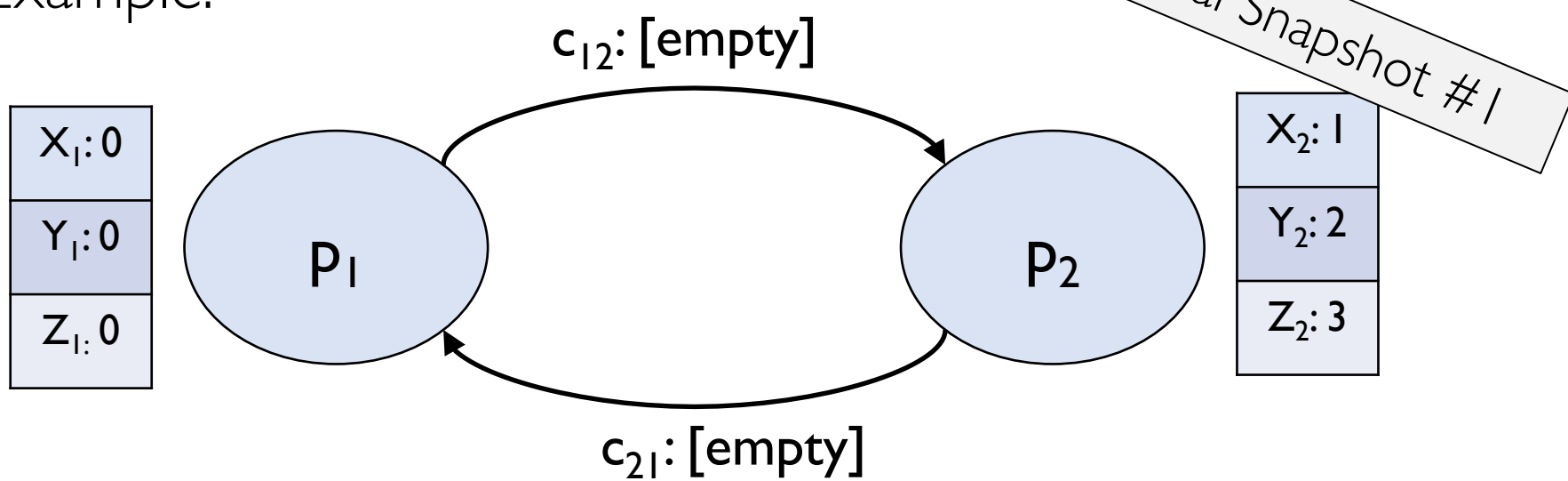
- Example:

$c_{12}$

$P_1$ $P_2$

$c_{21}$

Two processes: $p_1$ and $p_2$.

$c_{12}$: channel from $p_1$ to $p_2$.          $c_{21}$: channel from $p_2$ to $p_1$.

# Global State (or Global Snapshot)

- State of each process (and each channel) in the system at a given instant of time.

- Example:

$c_{12}$: [empty]

Global Snapshot #1

| $X_1$: 0 |
|---|
| $Y_1$: 0 |
| $Z_1$: 0 |

$P_1$

$P_2$

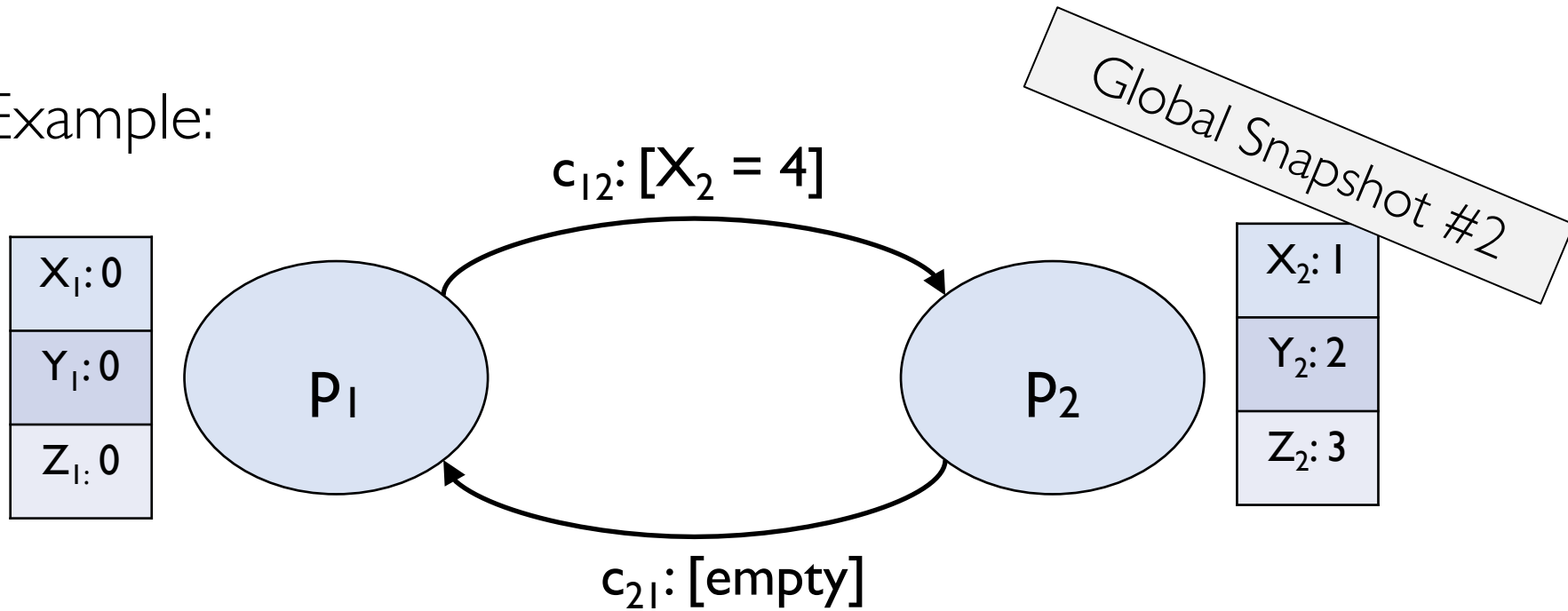| $X_2$: 1 |
|---|
| $Y_2$: 2 |
| $Z_2$: 3 |

$c_{21}$: [empty]

Process state for $p_1$ and $p_2$.
No pending messages on the channels.

# Global State (or Global Snapshot)

- State of each process (and each channel) in the system at a given instant of time.
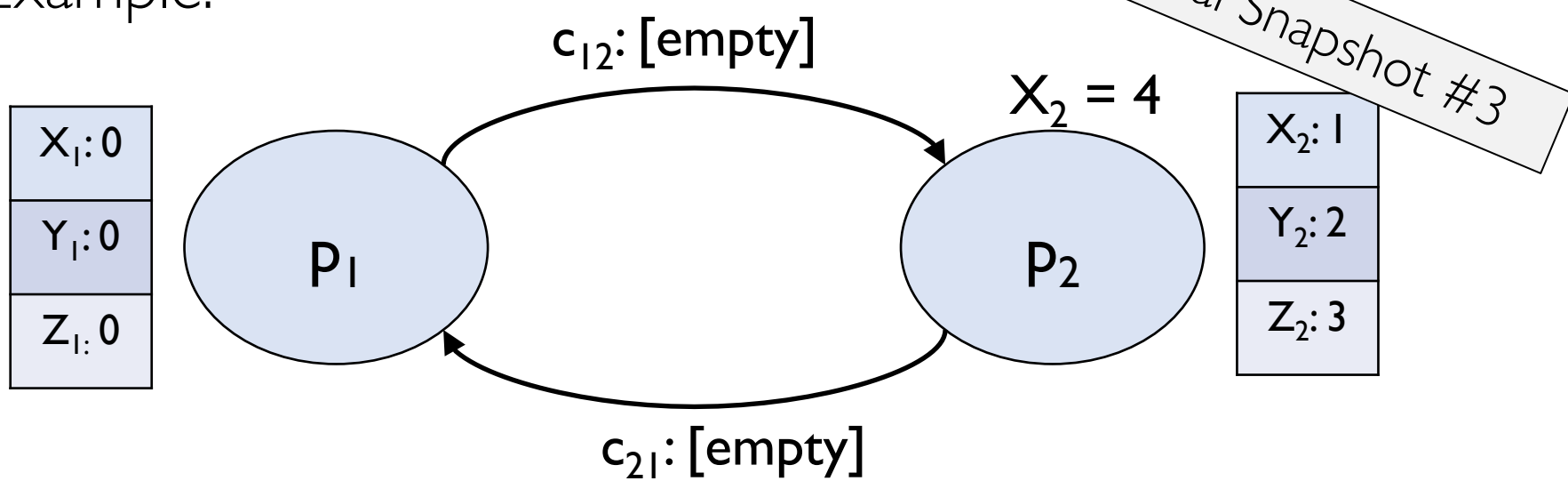
- Example:

$c_{12}$: [$X_2$ = 4]

| |
|---|
| $X_1$: 0 |
| $Y_1$: 0 |
| $Z_1$: 0 |

$P_1$

$P_2$

| |
|---|
| $X_2$: 1 |
| $Y_2$: 2 |
| $Z_2$: 3 |

$c_{21}$: [empty]

event 1: $p_1$ sends a message to $p_2$ asking it to set $X_2$ = 4

# Global State (or Global Snapshot)

- State of each process (and each channel) in the system at a given instant of time.
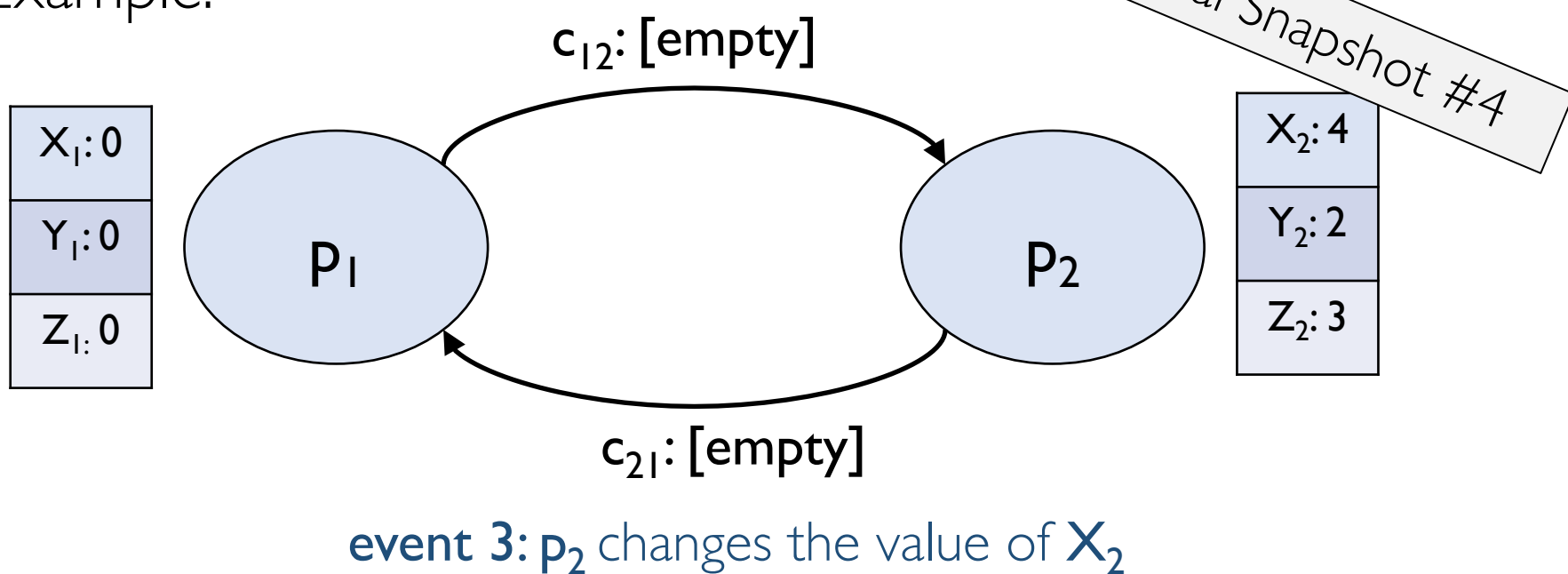
- Example:

$c_{12}$: [empty]

$X_2 = 4$

| $X_1$: 0 |
|---|
| $Y_1$: 0 |
| $Z_1$: 0 |

$P_1$

$P_2$

| $X_2$: 1 |
|---|
| $Y_2$: 2 |
| $Z_2$: 3 |

$c_{21}$: [empty]

**event 2: $p_2$** receives the message.

# Global State (or Global Snapshot)

- State of each process (and each channel) in the system at a given instant of time.

- Example:

$c_{12}$: [empty]

Global Snapshot #4

| $X_1$: 0 |
|---|
| $Y_1$: 0 |
| $Z_1$: 0 |

$P_1$

$P_2$

| $X_2$: 4 |
|---|
| $Y_2$: 2 |
| $Z_2$: 3 |

$c_{21}$: [empty]

event 3: $p_2$ changes the value of $X_2$

# Capturing a global snapshot

- Useful to capture a global snapshot of the system:
    - *Checkpointing* the system state.
    - Reasoning about unreferenced objects (for garbage collection).
    - Deadlock detection.
    - Distributed debugging.

  *To be continued in next class….*