# Synchronization Algorithms

© Riccardo Sisto, Politecnico di Torino

Reference for study:

Van Steen, Tanenbaum, "Distributed Systems", chapter 6

# Synchronization

- Processes in a DS are asynchronous

- Sometimes, we need to synchronize them
  - a process must wait until another process completes another operation
  - some ordering of events must be enforced
  - some timing requirements must be satisfied

- Two possible approaches for decentralized systems:
  - synchronize physical clocks
  - logical clocks

# Physical Clock Synchronization

- Different possibilities exist
  - UTC Receivers
    - synchronization can be very accurate (up to 50ns precision) but expensive
  - NTP
    - cheaper solution but not so accurate

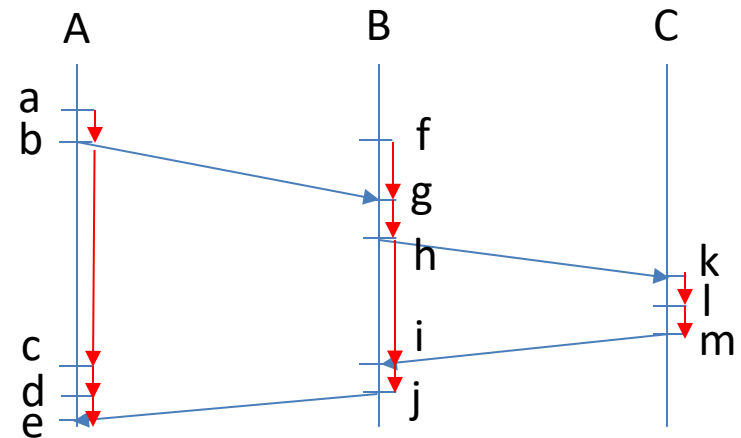| | NTP accuracy |
|---|---|
| LAN | <1ms |
| public internet | 10-50ms |
| public internet under congestion | >100ms |

# Logical Clocks

- For synchronization purposes, agreement on **ordering of events** is often sufficient

- Logical clocks are distributed algorithms that can be used to achieve this sort of agreement
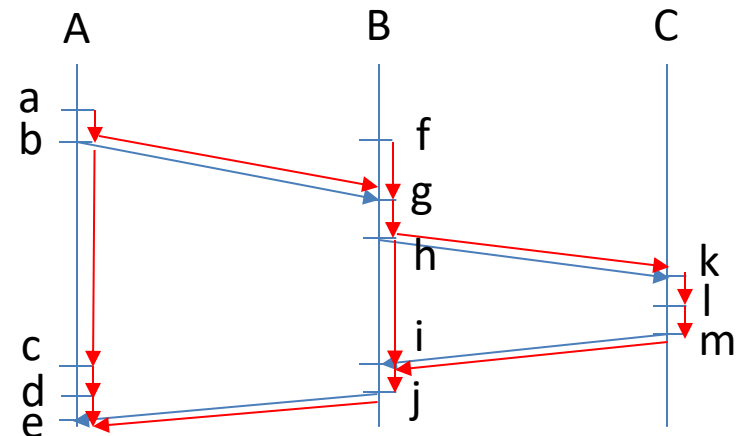
# The Happens-before Relation in a DS

- Definition:

  **x → y**        event **x** happens before event **y**

- Baseline: in a distributed system based on message passing, **x → y** can  be observed in these cases:

  – **x** happens before **y** in the same process

# The Happens-before Relation in a DS

- Definition: the Happens-before relation

    $x \rightarrow y$    event $x$ happens before event $y$

- Baseline: in a distributed system based on message passing, $x \rightarrow y$ can be observed in these cases:

    - $x$ happens before $y$ in the same process

    - $x$ and $y$ are the events of sending and receiving the same message (by different processes)

    - transitive property holds

# Happens-before and Causality

- The happens-before relation also captures the possibility of **causal relation** (causality) between two events

  – if $x \rightarrow y$            then **x** and **y** **may be** <u>causally related</u>

    in cause-effect relation (**x** cause, **y** effect)
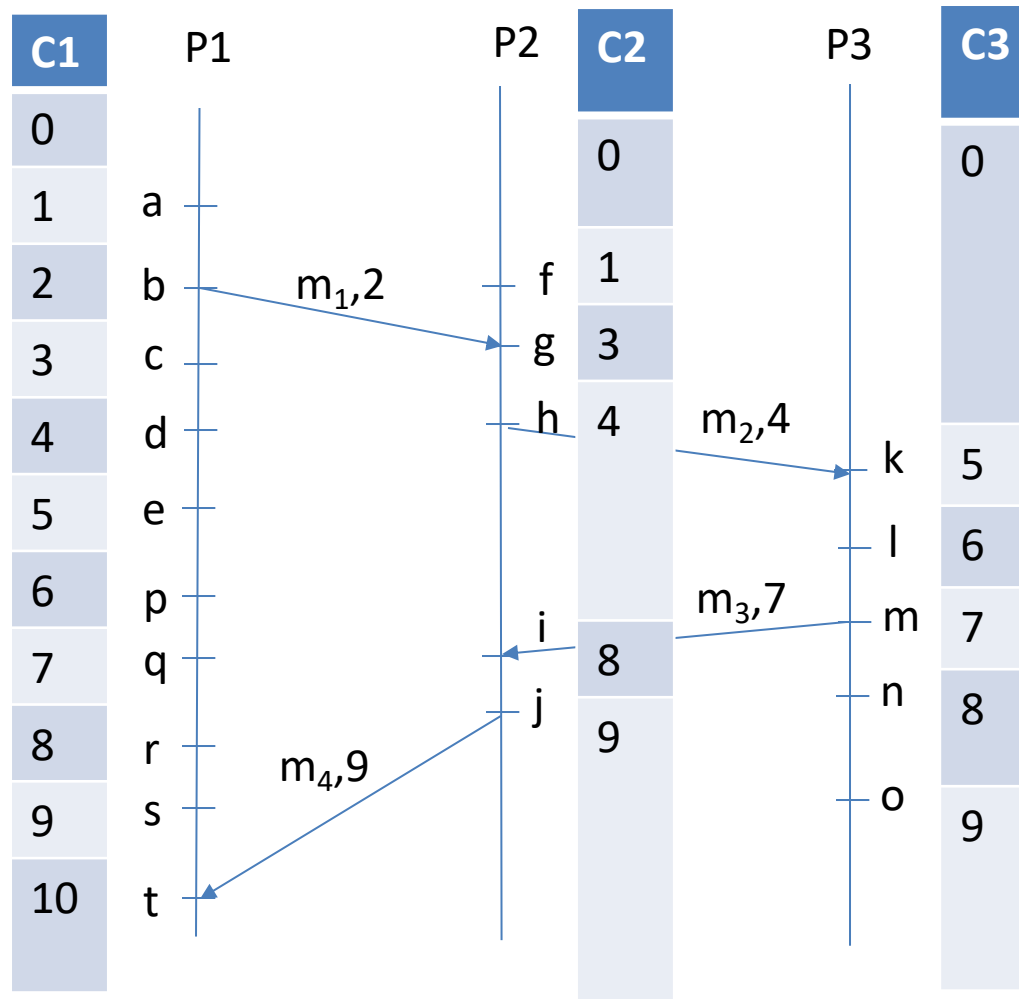    OR  **y**  depends on/is influenced by  **x**

  – if $x \nrightarrow y$
    and $y \nrightarrow x$        then **x** and **y** **are not** causally related

    (we say they are **concurrent**)

# Lamport Clocks

- Goal: assign timestamps (i.e., clock values) to events in such a way that   if   $x \rightarrow y$   then   $C(x) < C(y)$

- Algorithm:
  - each process Pi keeps a local time counter $C_i$
  - Pi timestamps each local event $x$ with the current $C_i$ value $C_i(x)$, after having incremented $C_i$
  - Pi timestamps each message it sends with the timestamp $C_i(s)$ assigned to the send event $s$
  - when Pi receives a message with timestamp $C_r$,
    - if $C_r > C_i$, Pi sets $C_i = C_r$ (clock adjustment) and then it timestamps the receive event as usual
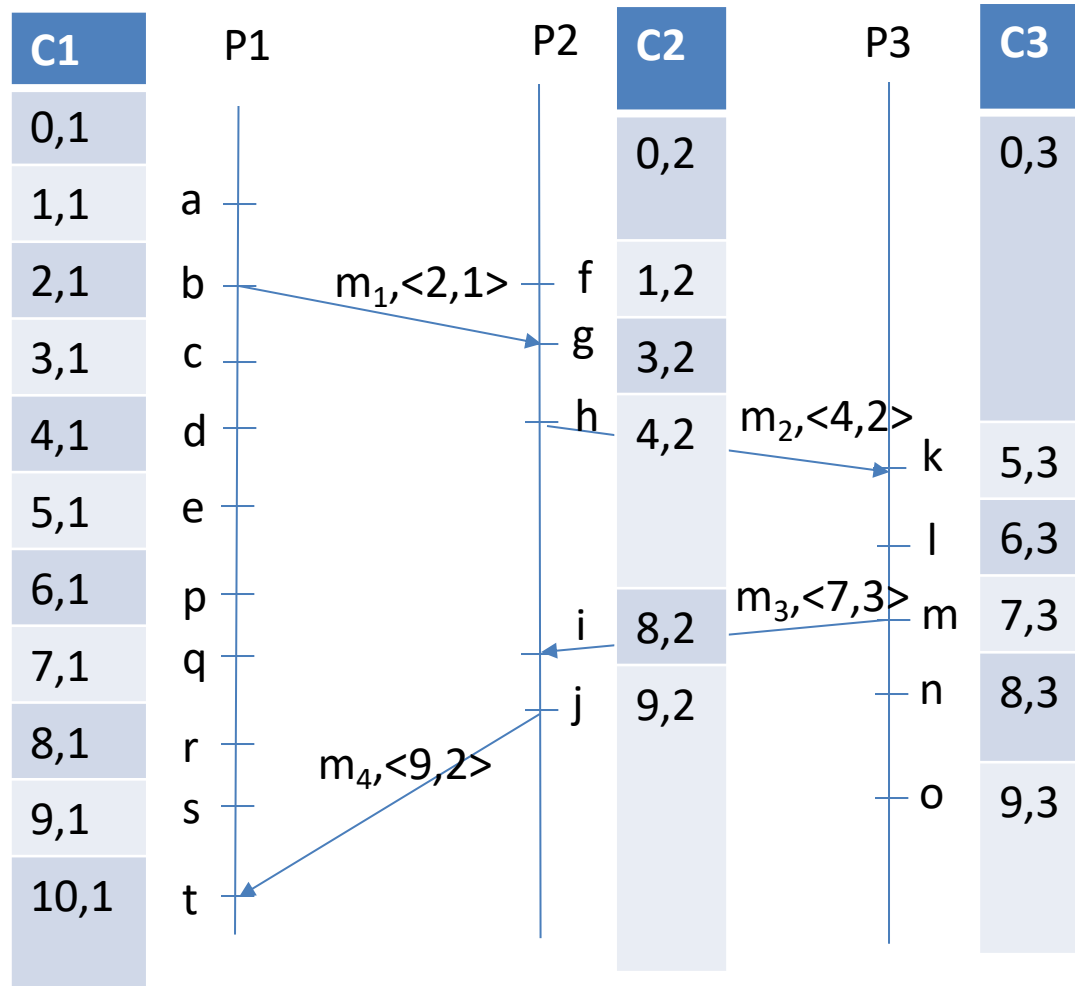
      equivalent to $C_i = \max(C_i, C_r)$

# Example

# Discussion

- Some events (in different processes) are assigned the same timestamp
  - If we want to avoid it, we can include the process id into the timestamp: $C(\mathbf{x})=<C_i(\mathbf{x}),i>$ if event $\mathbf{x}$ occurs in $P_i$
    - assumption: process ids are unique
    - $<C_i(\mathbf{x}),i> < <C_j(\mathbf{y}),j>$ iff $C_i(\mathbf{x})<C_j(\mathbf{y})$ or $C_i(\mathbf{x})=C_j(\mathbf{y})$ and $i<j$

- $\mathbf{x} \rightarrow \mathbf{y}$ => $C(\mathbf{x})<C(\mathbf{y})$ but **not** vice versa

- The algorithm can be implemented by
  - a timestamping function that increments and returns the counter
  - hooks in the send/receive functions that timestamp send/receive events and messages and adjust the counter
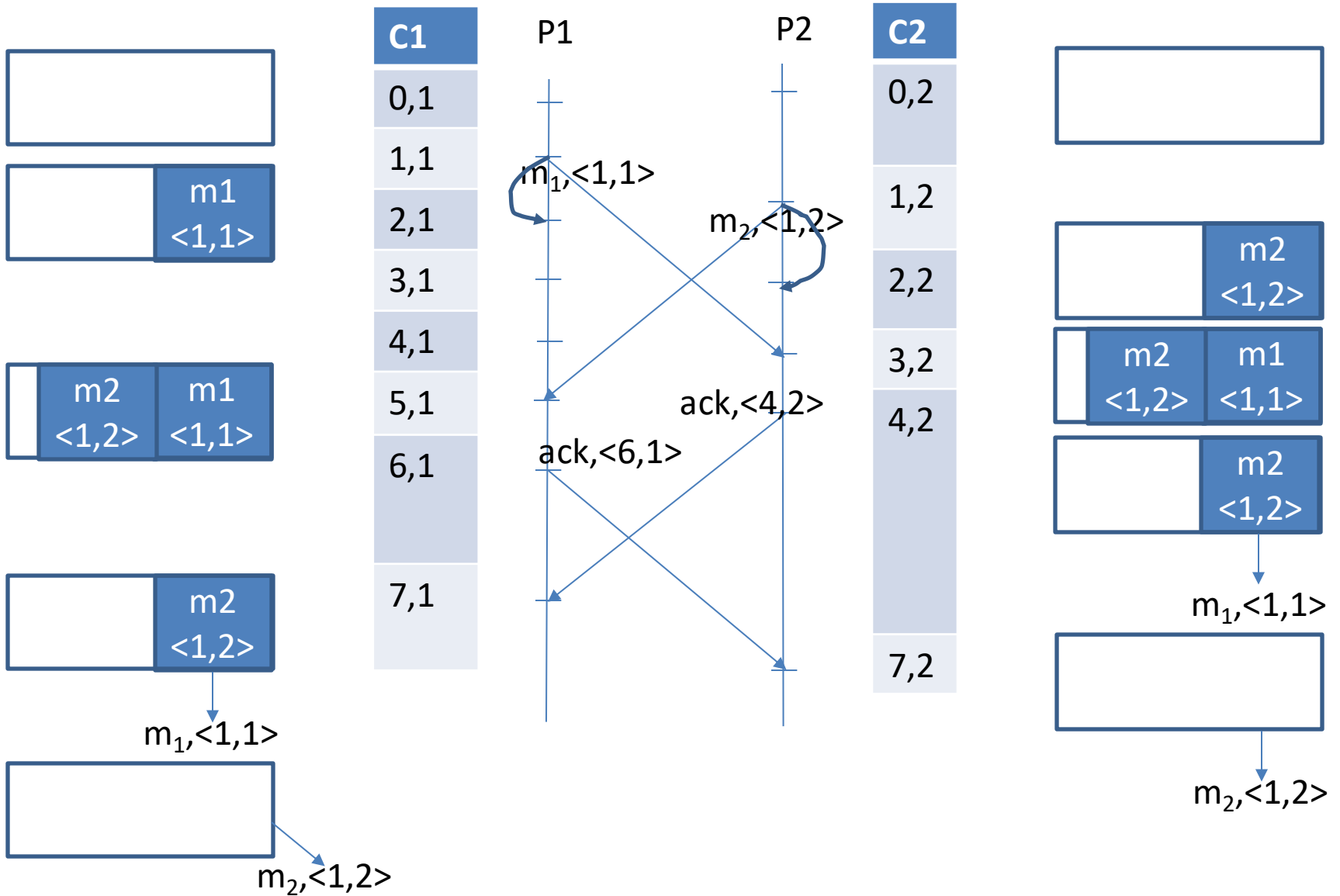
# Example with Total Ordering

# Application: Total-ordered Multicast

- Problem statement:
  - a set of processes that send multicast messages
  - all multicast messages must be delivered in the same order to each receiver (total-ordered multicast)

- Example
  - 2 or more identical replicated servers
    - replication goals: fault-tolerance and latency reduction
  - service requests can be sent to any server in the group
  - the receiving server multicasts the request to all servers
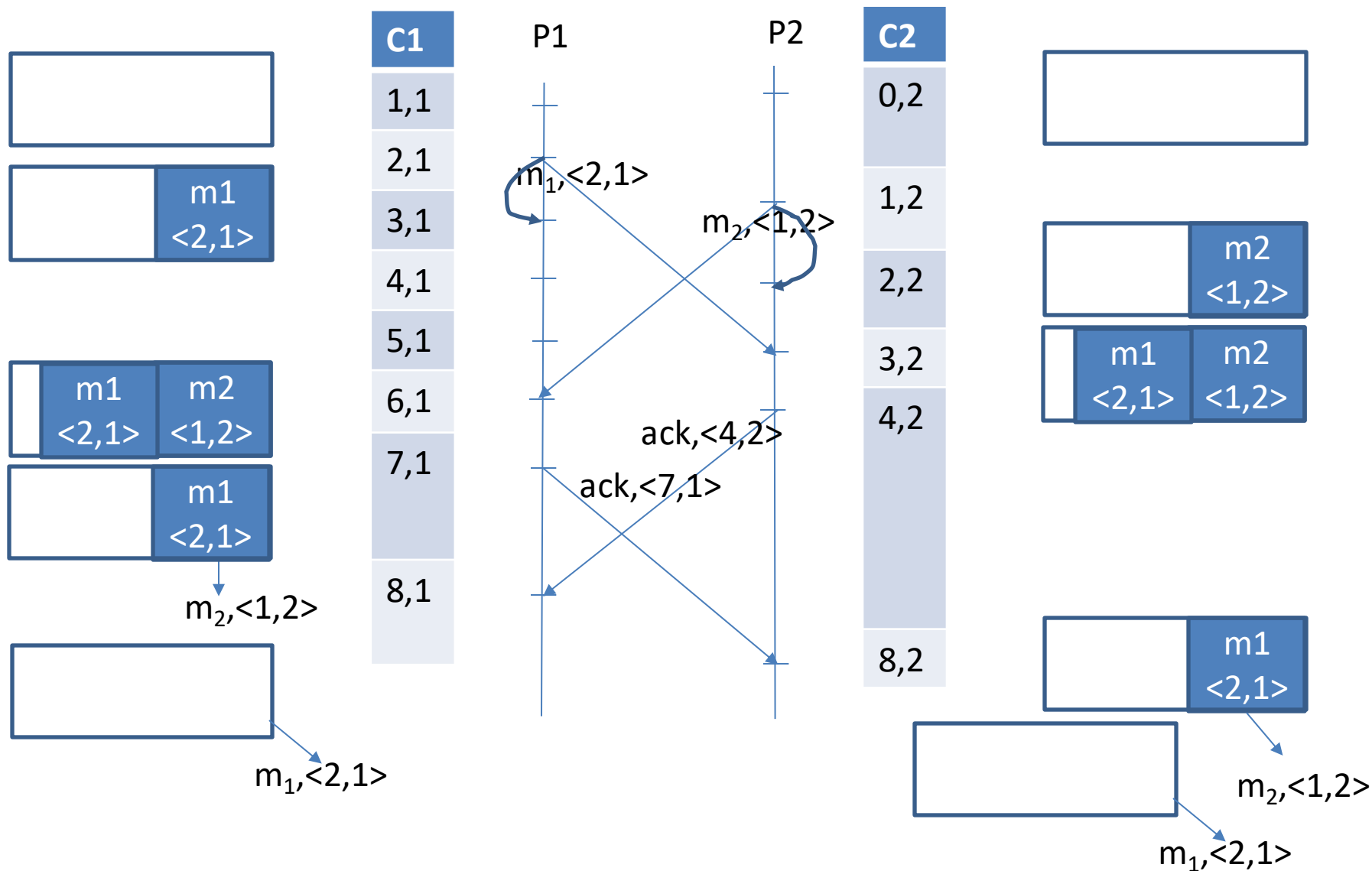  - operations must be executed in each replica in the same order

# Total-ordered Multicast Algorithm

- Lamport clocks are maintained in each process.

  => Each multicast message is timestamped by the sender with the send event timestamp

- Received messages are queued, ordered by timestamp

- Each receiver acknowledges message reception to the other receivers

- When acknowledgments for the message at the head of the queue have been received from all other receivers, that message is dequeued and handed to the application

# Example

# Example: Other Scenario

# Total-ordered Multicast: Discussion

- The algorithm has been shown to deliver messages to all destination processes in the same order under the following assumptions:

  - no messages are lost

  - messages from the same sender are received in the same order they were sent

- In case the messages have the meaning of operations to be executed on replicated data in the same order, the algorithm implements state machine replication
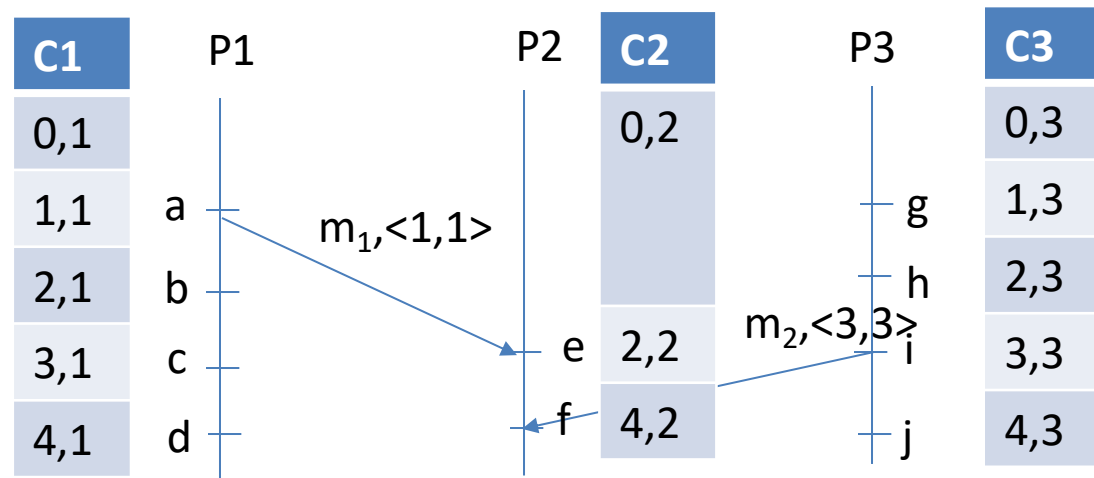
# Vector Clocks

- Lamport clocks provide a **total ordering** of events, but this ordering does not necessarily imply *causality*

    $C(\mathbf{x}) < C(\mathbf{y})$  ?  $\mathbf{x} \rightarrow \mathbf{y}$

- Example:

    $C(\mathbf{e}) < C(\mathbf{i})$

    but  not $\mathbf{e} \rightarrow \mathbf{i}$



- Vector clocks are a way to provide a **partial ordering** of events that captures causality, i.e.

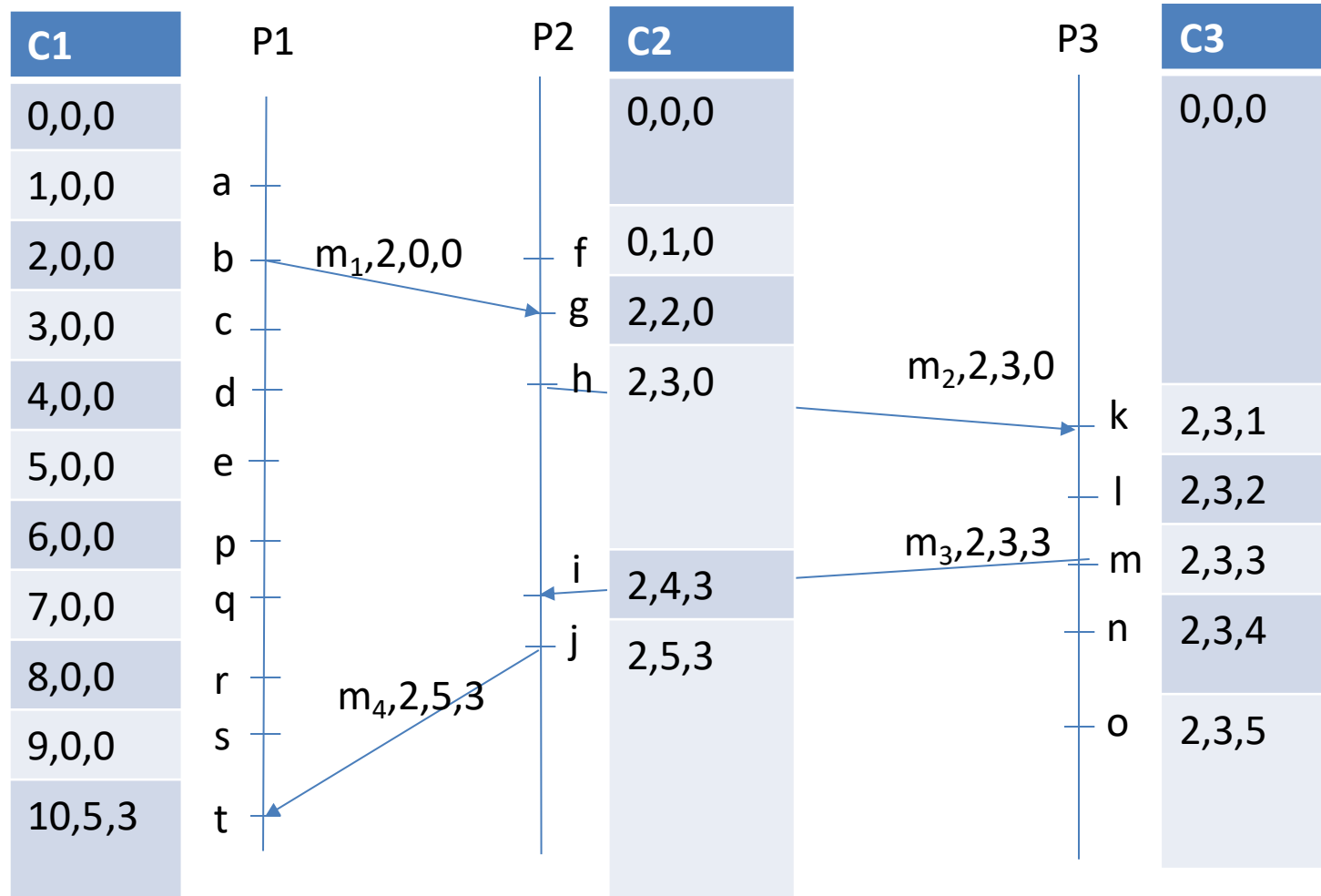    $C(\mathbf{x}) < C(\mathbf{y})$  <=>  $\mathbf{x} \rightarrow \mathbf{y}$

# Vector Clocks Algorithm

- Each process Pi keeps a local **vector** of time counters VCi[]
  - VCi[i] is the local event counter (i.e. local time at Pi, as with Lamport)
  - VCi[j] is Pi's knowledge of the local time at Pj

- Pi timestamps each local event **x** with the current VCi value ts(**x**), after having incremented VCi[i]

- Pi timestamps each message it sends with the timestamp ts(**s**) assigned to the send event **s**

- when Pi receives a message with timestamp tsr,
  - Pi sets VCi[k]=max(VCi[k],tsr[k]) for each k (vector clock adjustment) and then it timestamps the receive event as usual

# Example: How does this scenario change with Vector Clocks?

# Example with Vector Clocks

| C1 | P1 | | P2 | C2 | | P3 | C3 |
|---|---|---|---|---|---|---|---|
| 0,0,0 | | | | 0,0,0 | | | 0,0,0 |
| 1,0,0 | a | | | | | | |
| 2,0,0 | b | $m_1,2,0,0$ | f | 0,1,0 | | | |
| 3,0,0 | c | | g | 2,2,0 | | | |
| 4,0,0 | d | | h | 2,3,0 | $m_2,2,3,0$ | k | 2,3,1 |
| 5,0,0 | e | | | | | l | 2,3,2 |
| 6,0,0 | p | | i | 2,4,3 | $m_3,2,3,3$ | m | 2,3,3 |
| 7,0,0 | q | | | | | n | 2,3,4 |
| 8,0,0 | r | | j | 2,5,3 | | | |
| 9,0,0 | s | $m_4,2,5,3$ | | | | o | 2,3,5 |
| 10,5,3 | t | | | | | | |

# Vector Clock Application Example

- Causal Ordered Multicast
  - each multicast message must be delivered after the delivery of all other causally related messages that were sent before
  - non causally related messages that were sent before may be received after

- Use Example
  - in a replicated chat, we may accept that the ordering of posts is different in each replica provided each reply to a post p does never come before p