

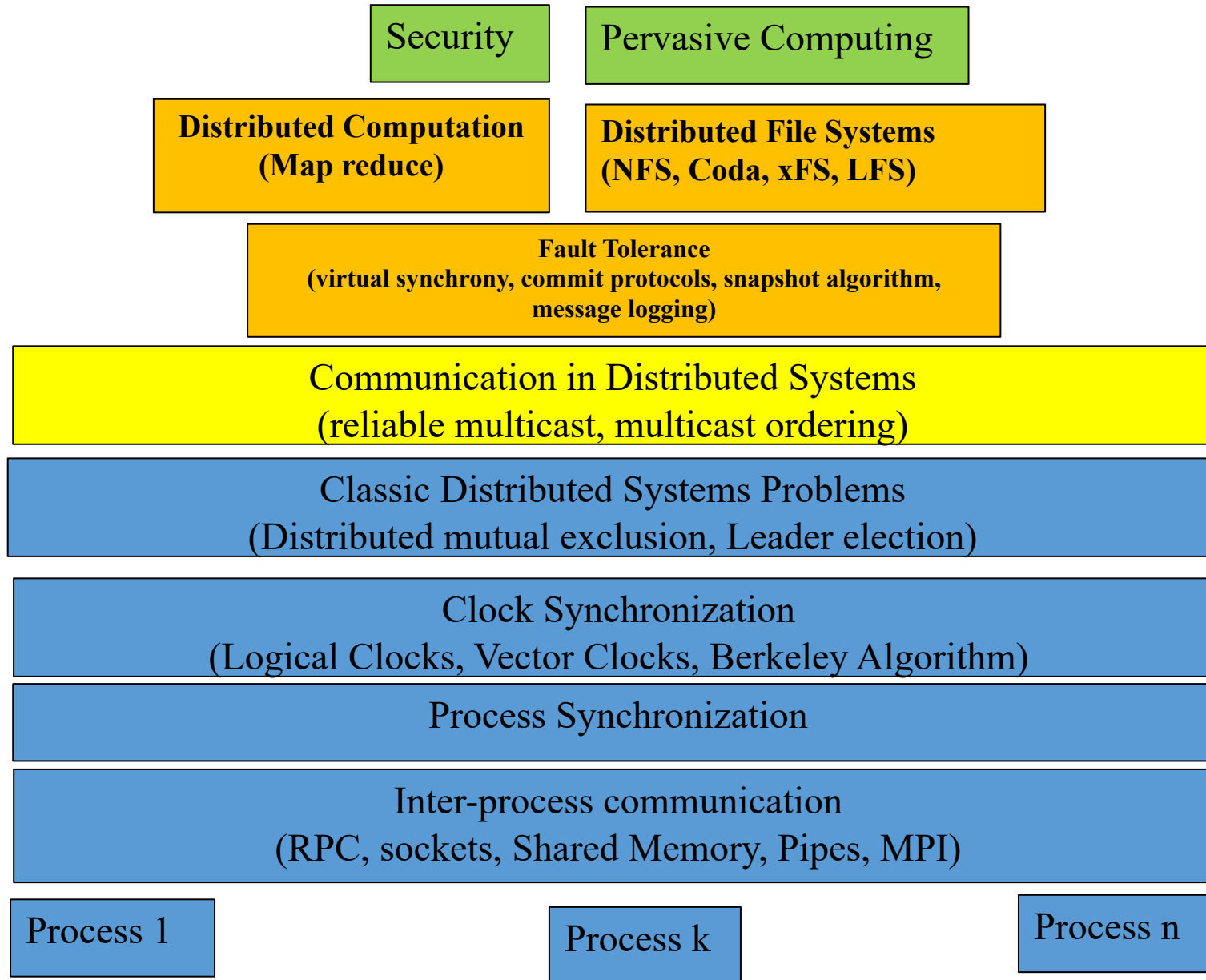
# CMSC621: Advanced Operating Systems

**Nilanjan Banerjee**

*Associate Professor, University of Maryland*  
Baltimore County  
nilanb@umbc.edu

**Advanced Operating Systems**

# Overview of the course



# Other Communication Forms

- **Multicast** → message sent to a group of processes
- **Broadcast** → message sent to all processes (anywhere)
- **Unicast** → message sent from one sender process to one receiver process

# Multicast Problem

Node with a piece of information  
to be communicated to everyone



Distributed Group  
of "Nodes" =  
Processes at  
Internet-based host

# Who Uses Multicast?

- A widely-used abstraction by almost all cloud systems
- Storage systems a database
  - Replica servers for a key: writes/reads to the key are multicast within the replica group
  - All servers: membership information (e.g., heartbeats) is multicast across all servers in cluster
- Online scoreboards (ESPN, French Open, FIFA World Cup)
  - Multicast to group of clients interested in the scores
- Stock Exchanges
  - Group is the set of broker computers
  - Groups of computers for High frequency Trading
- Air traffic control system
  - All controllers need to receive the same updates in the same order

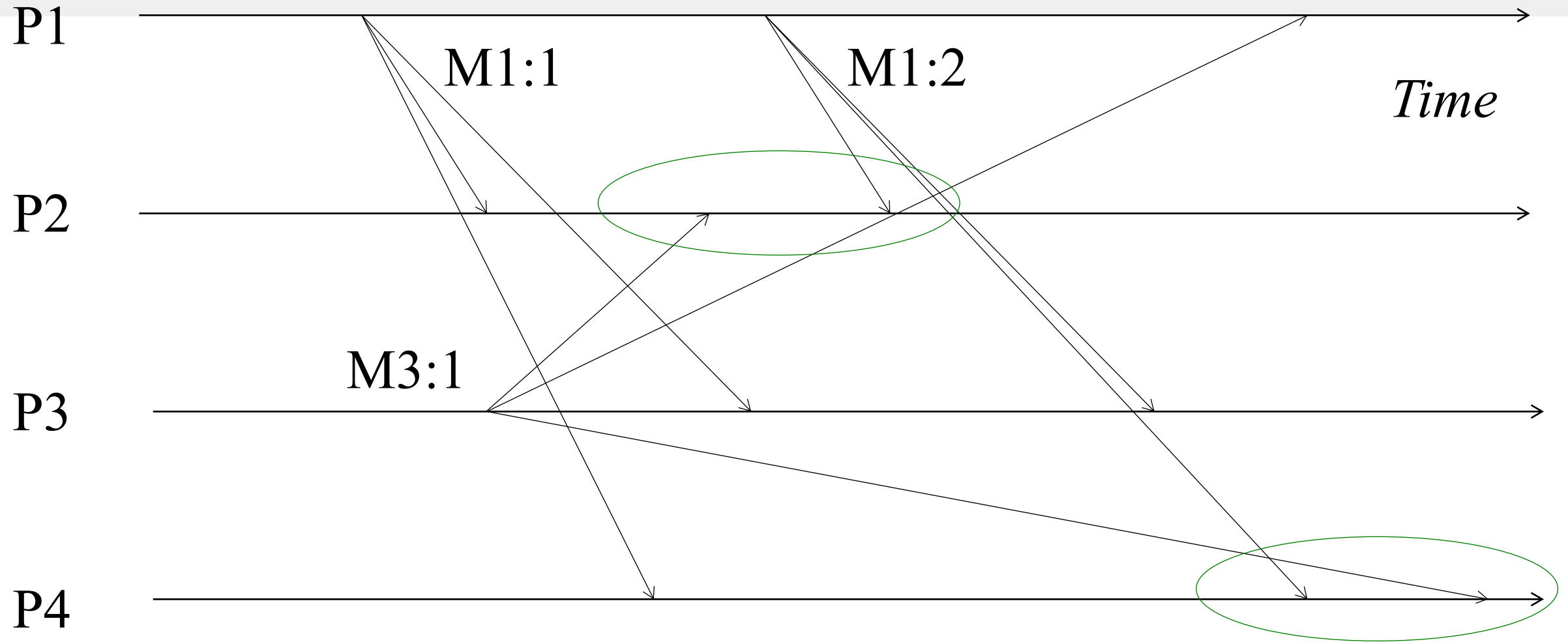
# Multicast Ordering

- Determines the meaning of “same order” of multicast delivery at different processes in the group
- Three popular flavors implemented by several multicast protocols
  1. FIFO ordering
  2. Causal ordering
  3. Total ordering

# FIFO ordering

- Multicasts from each sender are received in the order they are sent, at all receivers
- **Don't worry** about multicasts from different senders
- More formally
  - If a correct process issues (sends) **multicast(g,m)** to group **g** and **then multicast(g,m')**, then **every correct process that receives m' would already have received m**

# FIFO Ordering: Example



M1:1 and M1:2 should be received in that order at each receiver

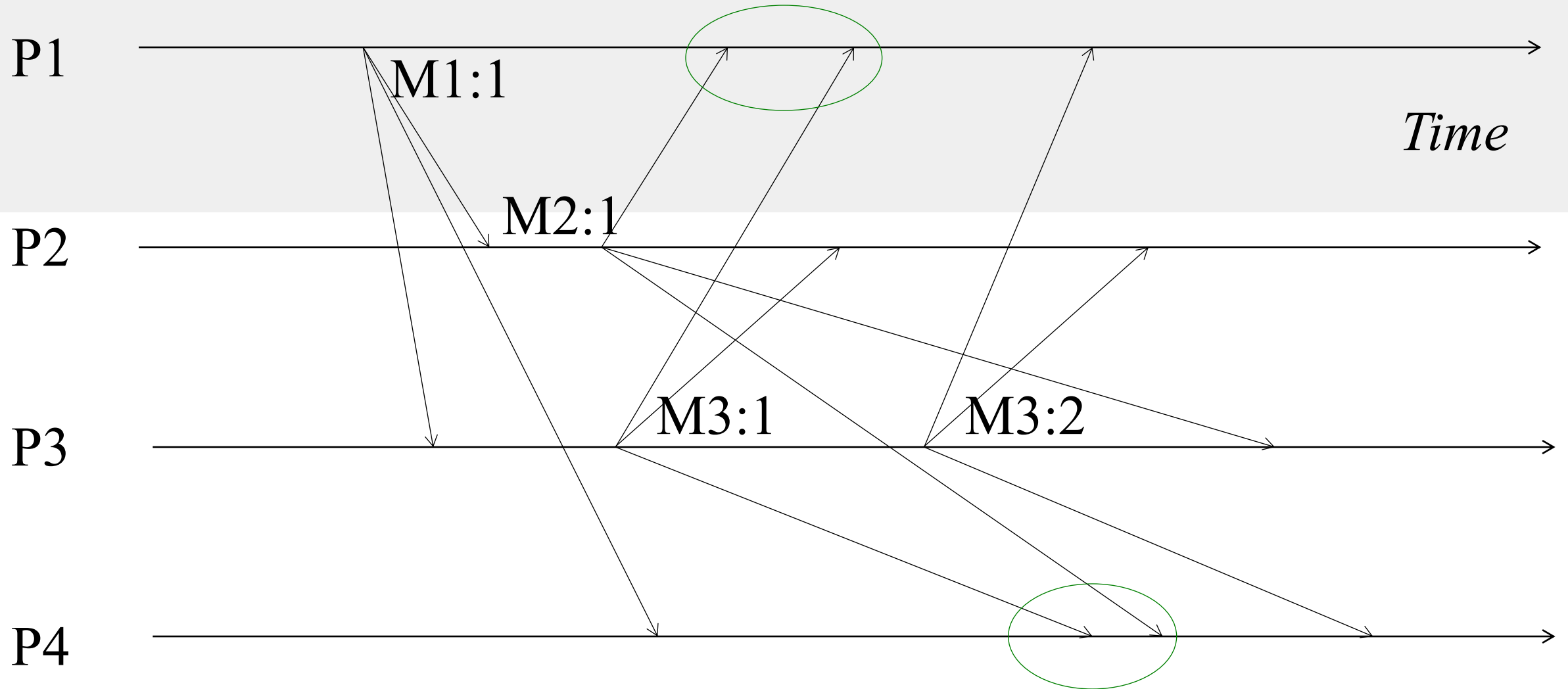
**Order of delivery of M3:1 and M1:2 could be different at different receivers**



# Causal Ordering

- Multicasts whose send events are causally related, must be received in the same causality-obeying order at all receivers
- Formally
  - If **multicast(g,m)  $\rightarrow$  multicast(g,m')** then any correct process that delivers m' would already have delivered m.
  - ( $\rightarrow$  is Lamport's happens-before)

# Causal Ordering: Example



M3:1  $\rightarrow$  M3:2, and so should be received in that order at each receiver

M1:1  $\rightarrow$  M3:1, and so should be received in that order at each receiver

**M3:1 and M2:1 are concurrent and thus ok to be received in different orders at different receivers**

# Causal vs. FIFO

- Causal Ordering  $\Rightarrow$  FIFO Ordering
- Why?
  - If two multicasts  $M$  and  $M'$  are sent by the same process  $P$ , and  $M$  was sent before  $M'$ , then  $M \rightarrow M'$
  - Then a multicast protocol that implements causal ordering will obey FIFO ordering since  $M \rightarrow M'$
- Reverse is not true! FIFO ordering does not imply causal ordering.

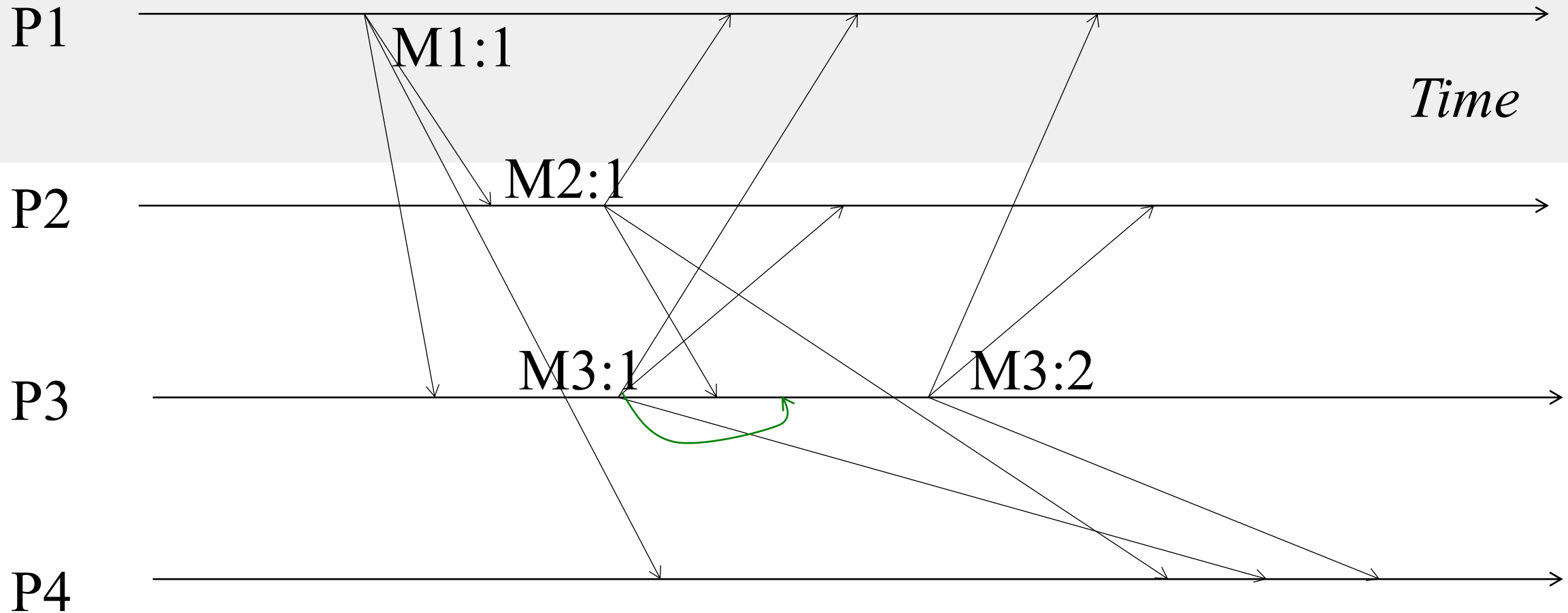
# Why do you need Causal at All?

- Group = set of your friends on a social network
- A friend sees your message  $m$ , and she posts a response (comment)  $m'$  to it
  - If friends receive  $m'$  before  $m$ , it wouldn't make sense
  - But if two friends post messages  $m''$  and  $n''$  concurrently, then they can be seen in any order at receivers
- A variety of systems implement causal ordering: Social networks, bulletin boards, comments on websites, etc.

# Total Ordering

- Also known as “Atomic Broadcast”
- Unlike FIFO and causal, this does not pay attention to order of multicast sending
- Ensures all receivers receive all multicasts in the same order
- Formally
  - If a correct process  $P$  delivers message  $m$  before  $m'$  (independent of the senders), then any other correct process  $P'$  that receives  $m'$  would already have received  $m$ .

# Total Ordering: Example



The order of receipt of multicasts is the same at all processes.

**M1:1, then M2:1, then M3:1, then M3:2**

**May need to delay delivery of some messages**

# Hybrid Variants

- Since FIFO/Causal are orthogonal to Total, can have hybrid ordering protocols too
  - FIFO-total hybrid protocol satisfies both FIFO and total orders
  - Causal-total hybrid protocol satisfies both Causal and total orders

# Implementation?

- That was *what* ordering is.
- But *how* do we implement each of these orderings?




# FIFO Multicast: Data Structures


- Each receiver maintains a per-sender sequence number (integers)
  - Processes  $P_1$  through  $P_N$
  - $P_i$  maintains a vector of sequence numbers  $P_i[1...N]$  (initially all zeroes)
  - $P_i[j]$  is the latest sequence number  $P_i$  has received from  $P_j$

# FIFO Multicast: Updating Rules

- Send multicast at process  $P_j$ :
  - Set  $P_j[j] = P_j[j] + 1$
  - Include new  $P_j[j]$  in multicast message as its sequence number
- Receive multicast: If  $P_i$  receives a multicast from  $P_j$  with sequence number  $S$  in message
  - if ( $S == P_i[j] + 1$ ) then
    - deliver message to application
    - Set  $P_i[j] = P_i[j] + 1$
  - else buffer this multicast until above condition is true

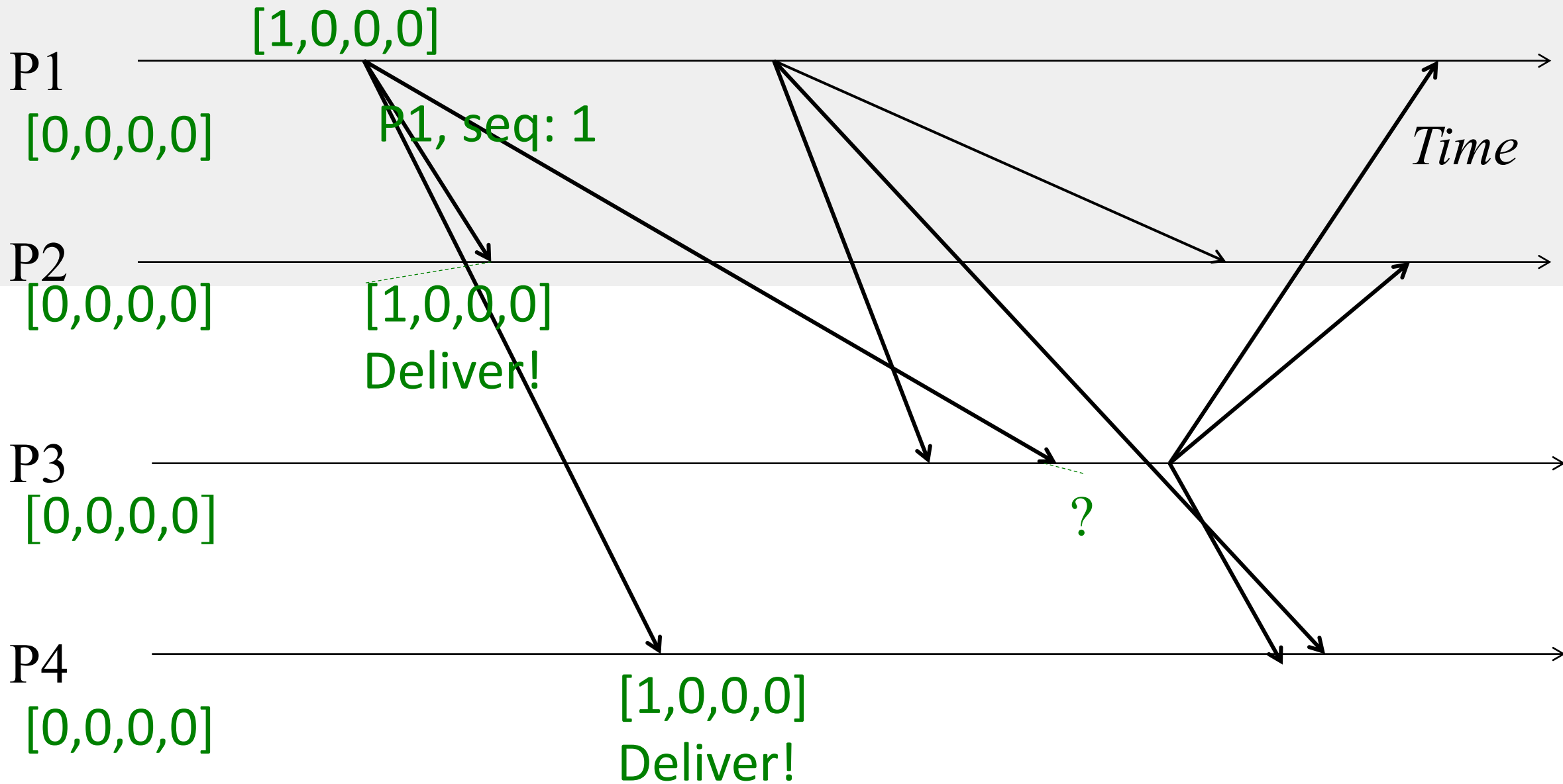
# FIFO Ordering: Example

P1   
[0,0,0,0] *Time*

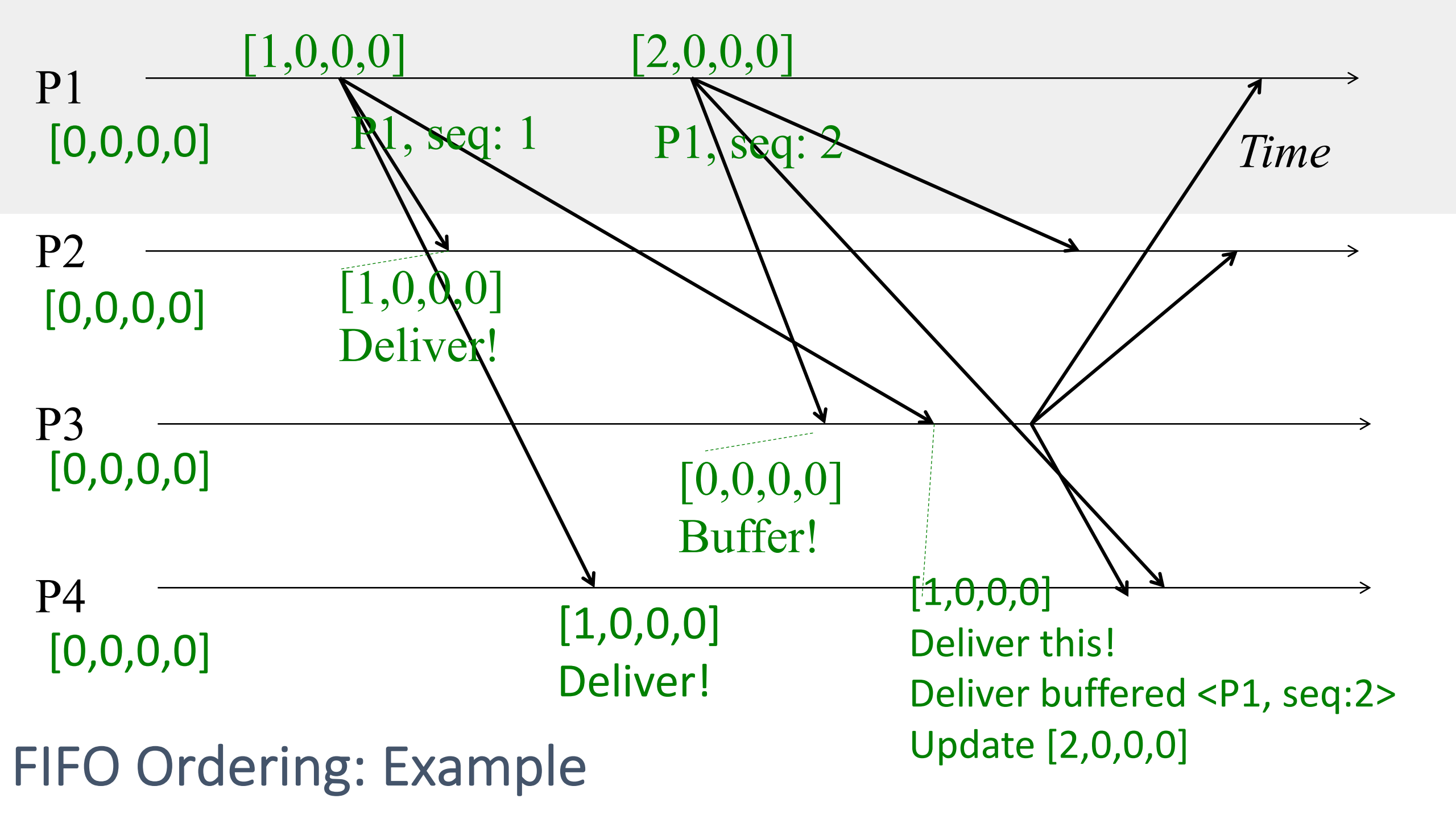
P2   
[0,0,0,0]

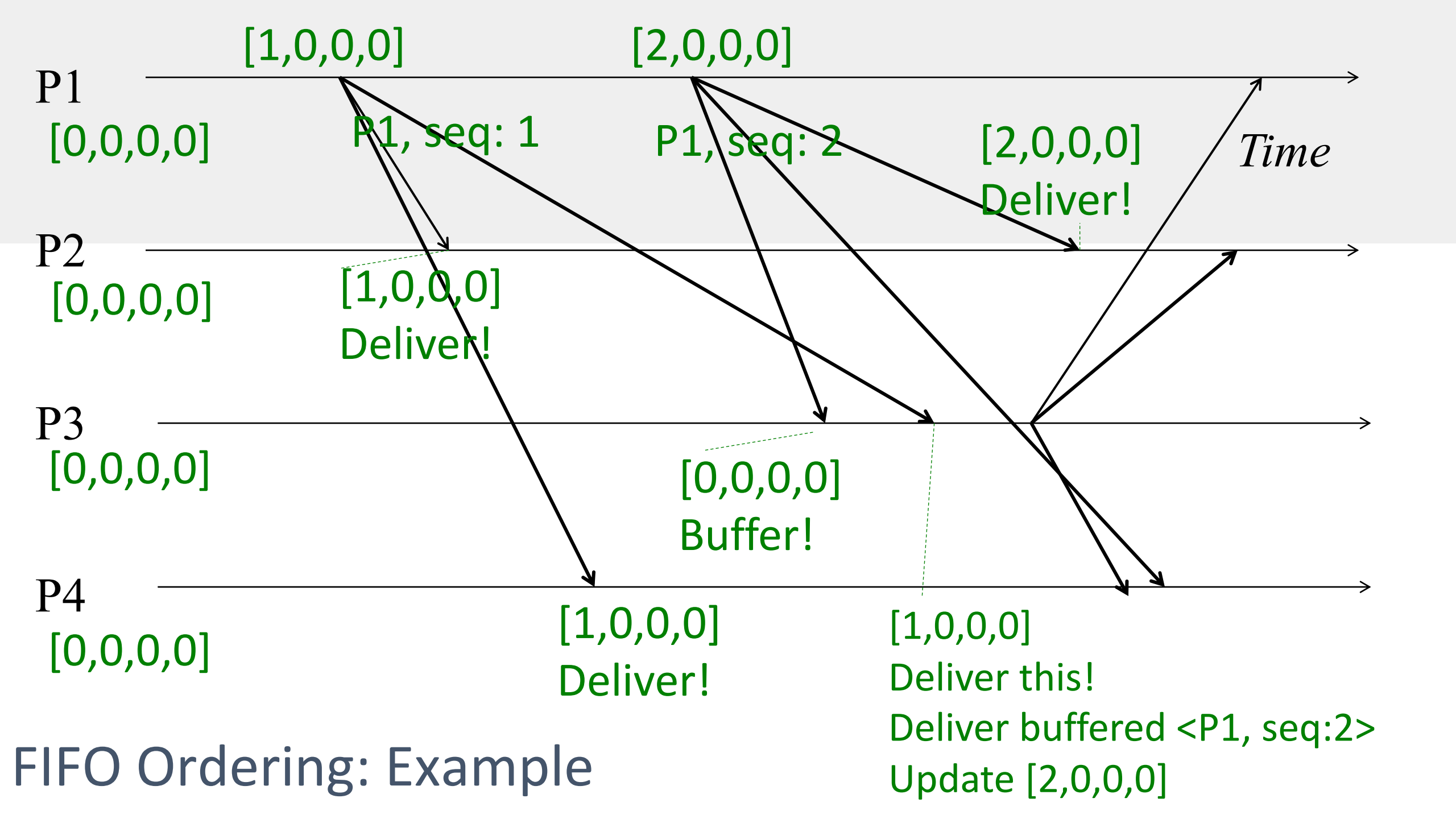
P3   
[0,0,0,0]

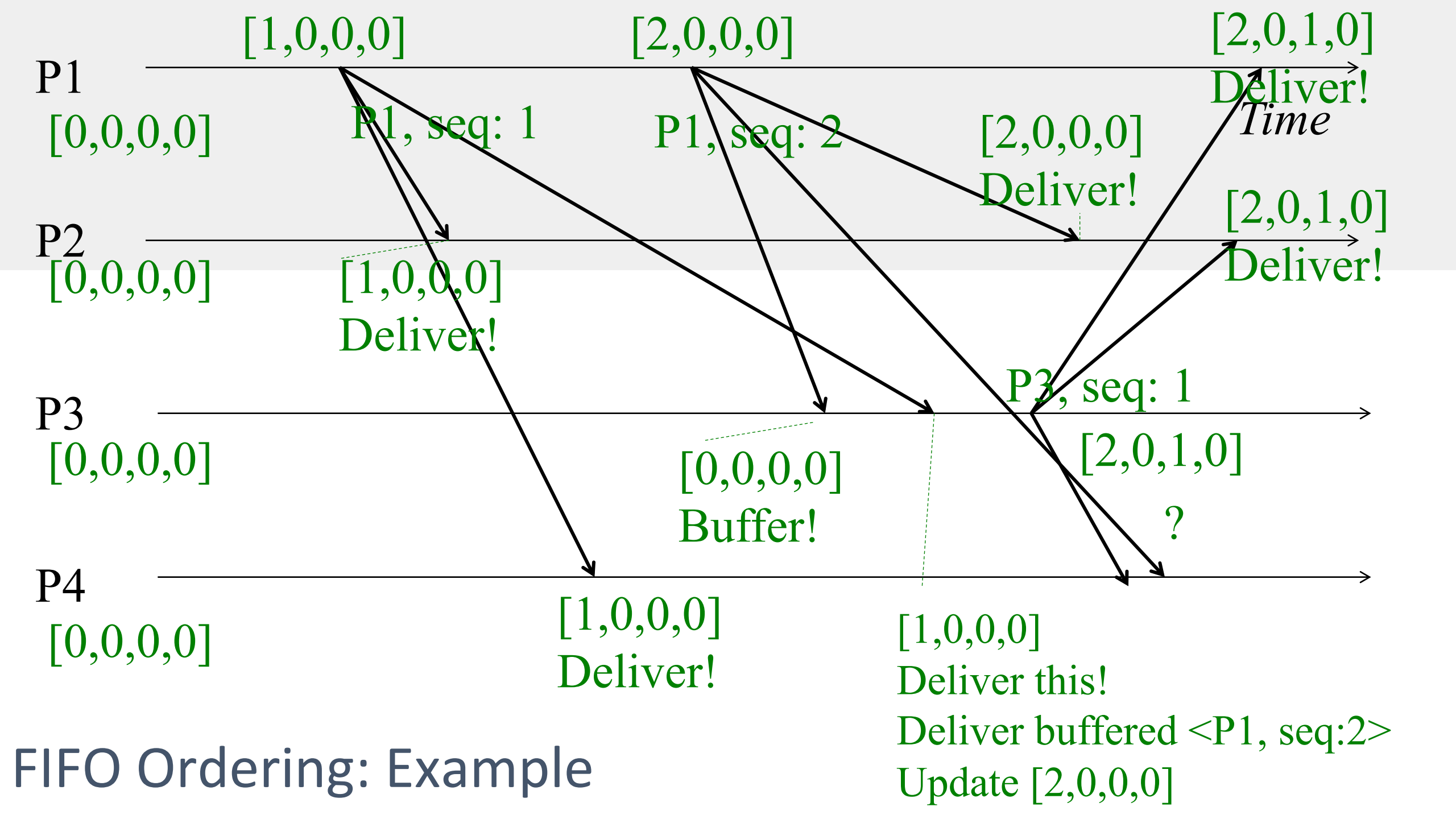
P4   
[0,0,0,0]

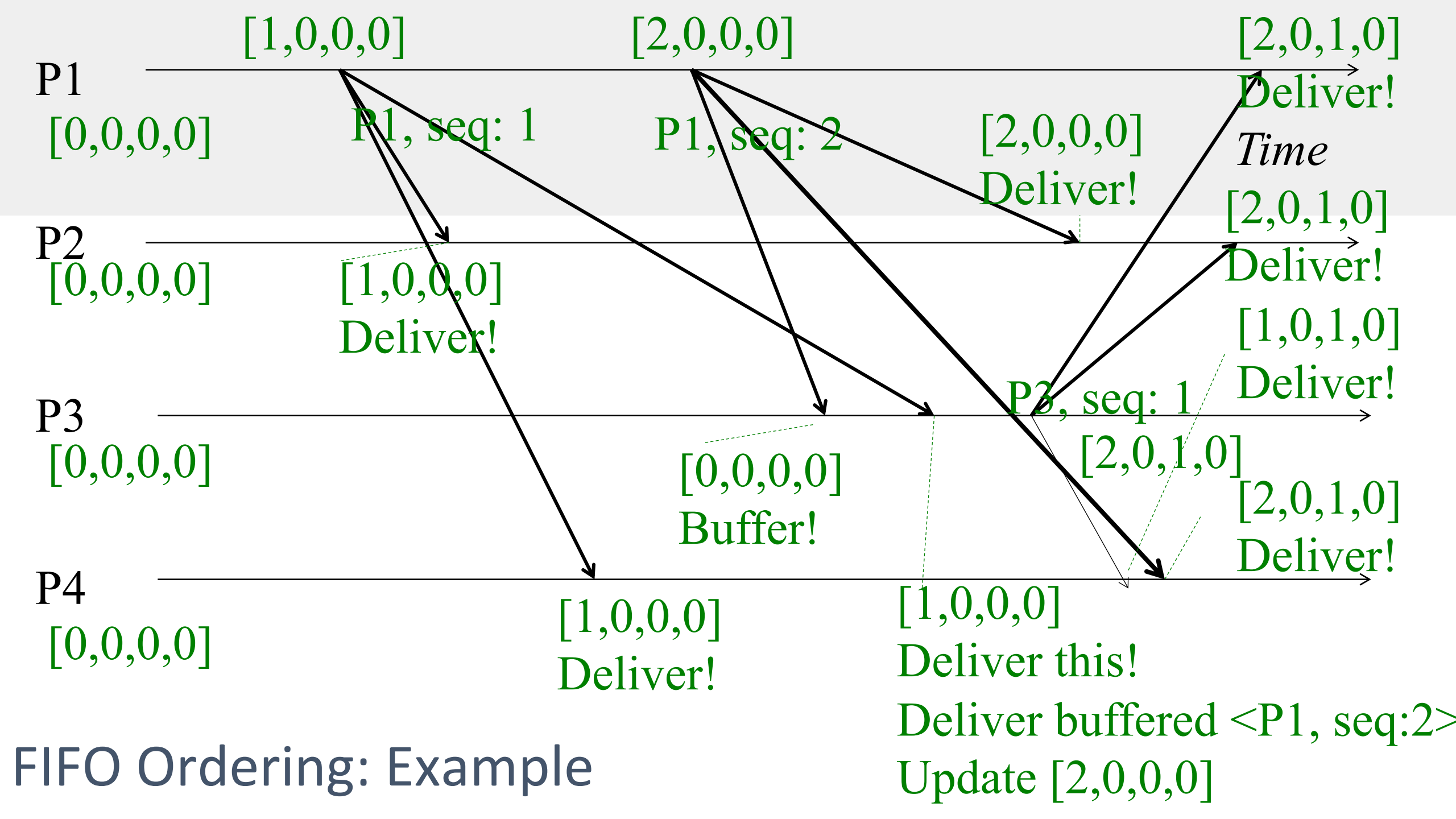


FIFO Ordering: Example











# Causal Ordering

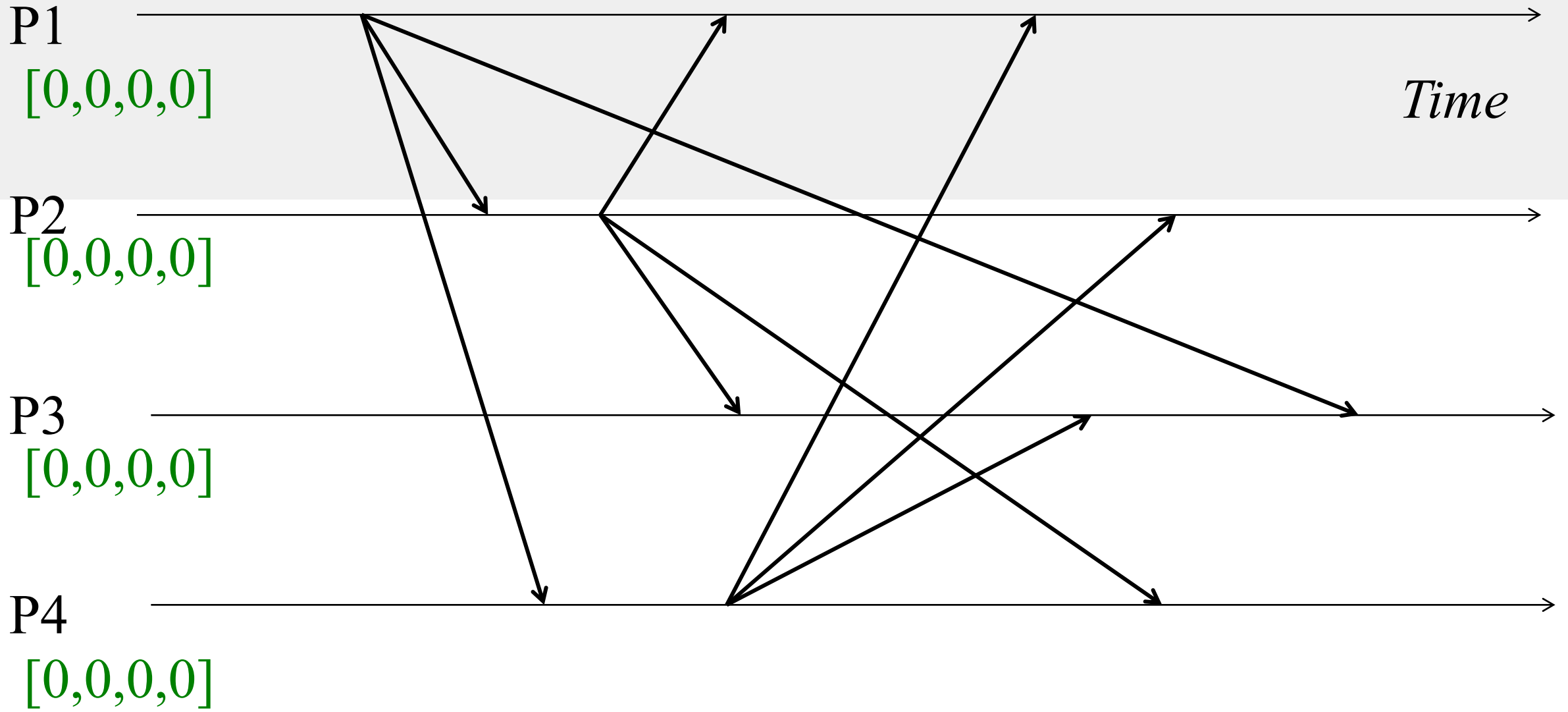
- Multicasts whose send events are causally related, must be received in the same causality-obeying order at all receivers
- Formally
  - If  $\text{multicast}(g,m) \rightarrow \text{multicast}(g,m')$  then any correct process that delivers  $m'$  would already have delivered  $m$ .
  - ( $\rightarrow$  is Lamport's happens-before)

# Causal Multicast: Data structures

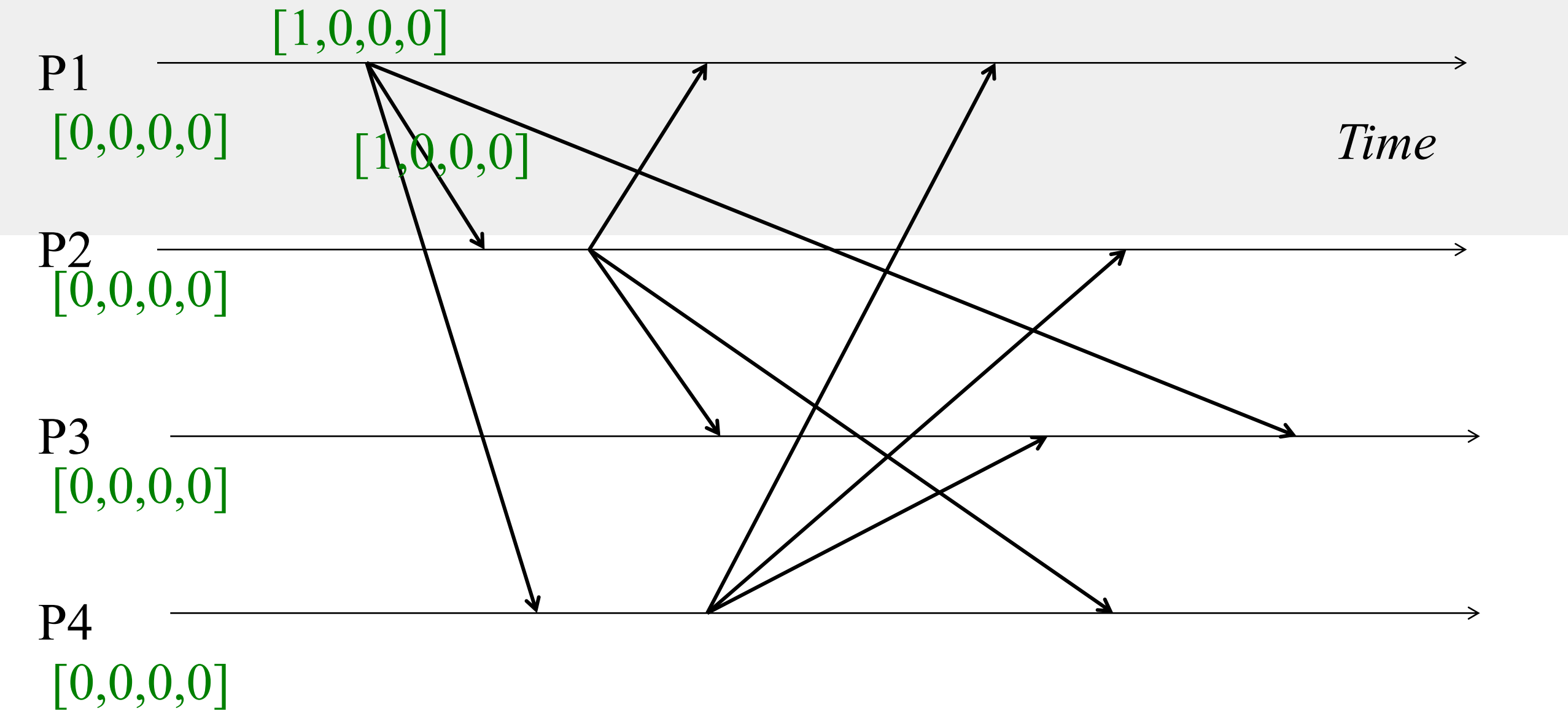
- Each receiver maintains a vector of per-sender sequence numbers (integers)
  - Similar to FIFO Multicast, but updating rules are different
  - Processes  $P_1$  through  $P_N$
  - $P_i$  maintains a vector  $P_i[1...N]$  (initially all zeroes)
  - $P_i[j]$  is the latest sequence number  $P_i$  has received from  $P_j$

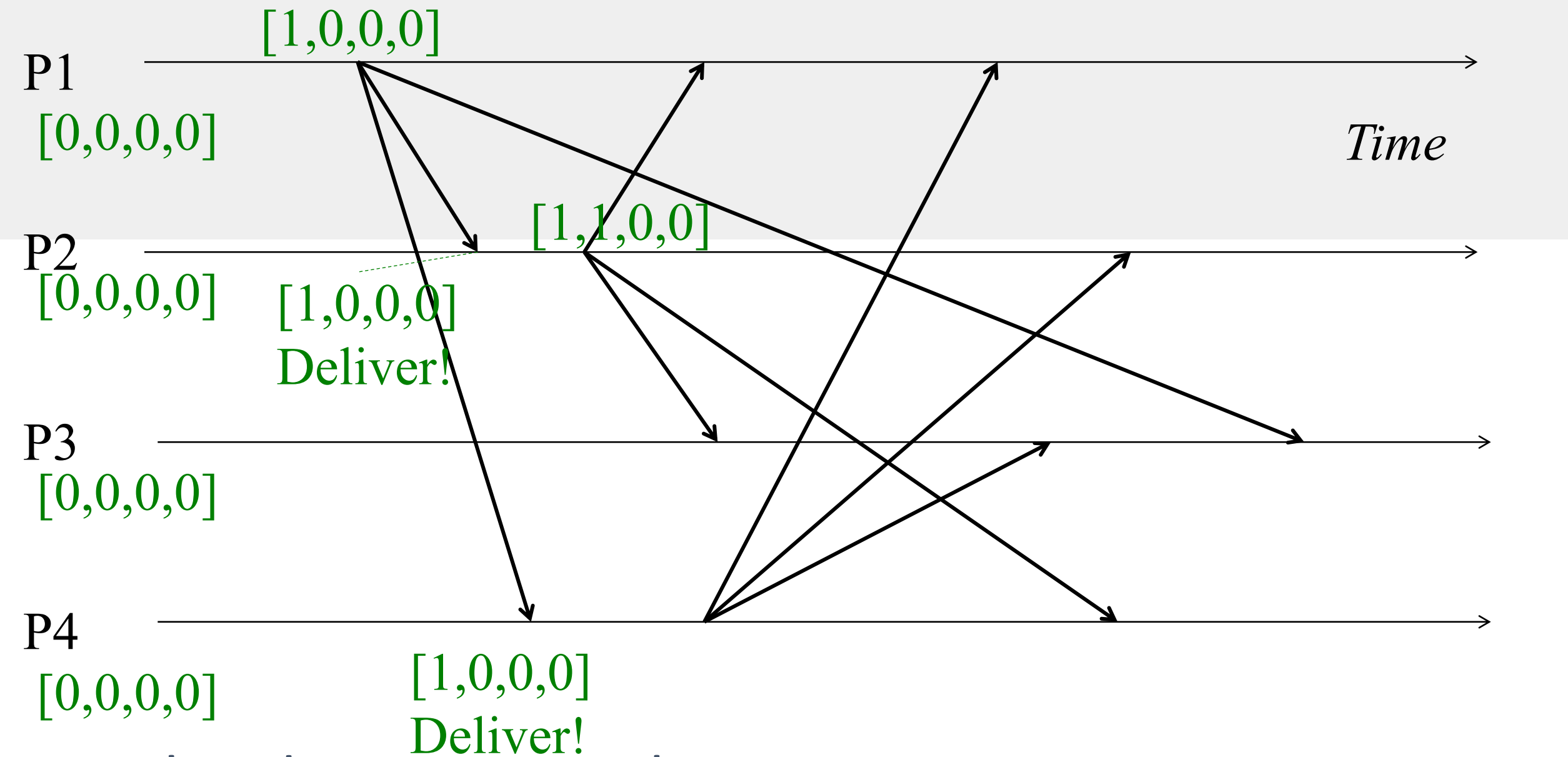
# Causal Multicast: Updating Rules

- Send multicast at process  $P_j$ :
  - Set  $P_j[j] = P_j[j] + 1$
  - Include new entire vector  $P_j[1...N]$  in multicast message as its sequence number
- Receive multicast: If  $P_i$  receives a multicast from  $P_j$  with vector  $M[1...N]$  ( $= P_j[1...N]$ ) in message, buffer it until both:
  - This message is the next one  $P_i$  is expecting from  $P_j$ , i.e.,
$$M[j] = P_i[j] + 1$$
  - All multicasts, anywhere in the group, which happened-before  $M$  have been received at  $P_i$ , i.e.,
$$\text{For all } k \neq j: M[k] \leq P_i[k]$$
i.e., ***Receiver satisfies causality***
  - When above two conditions satisfied, deliver  $M$  to application and set  $P_i[j] = M[j]$

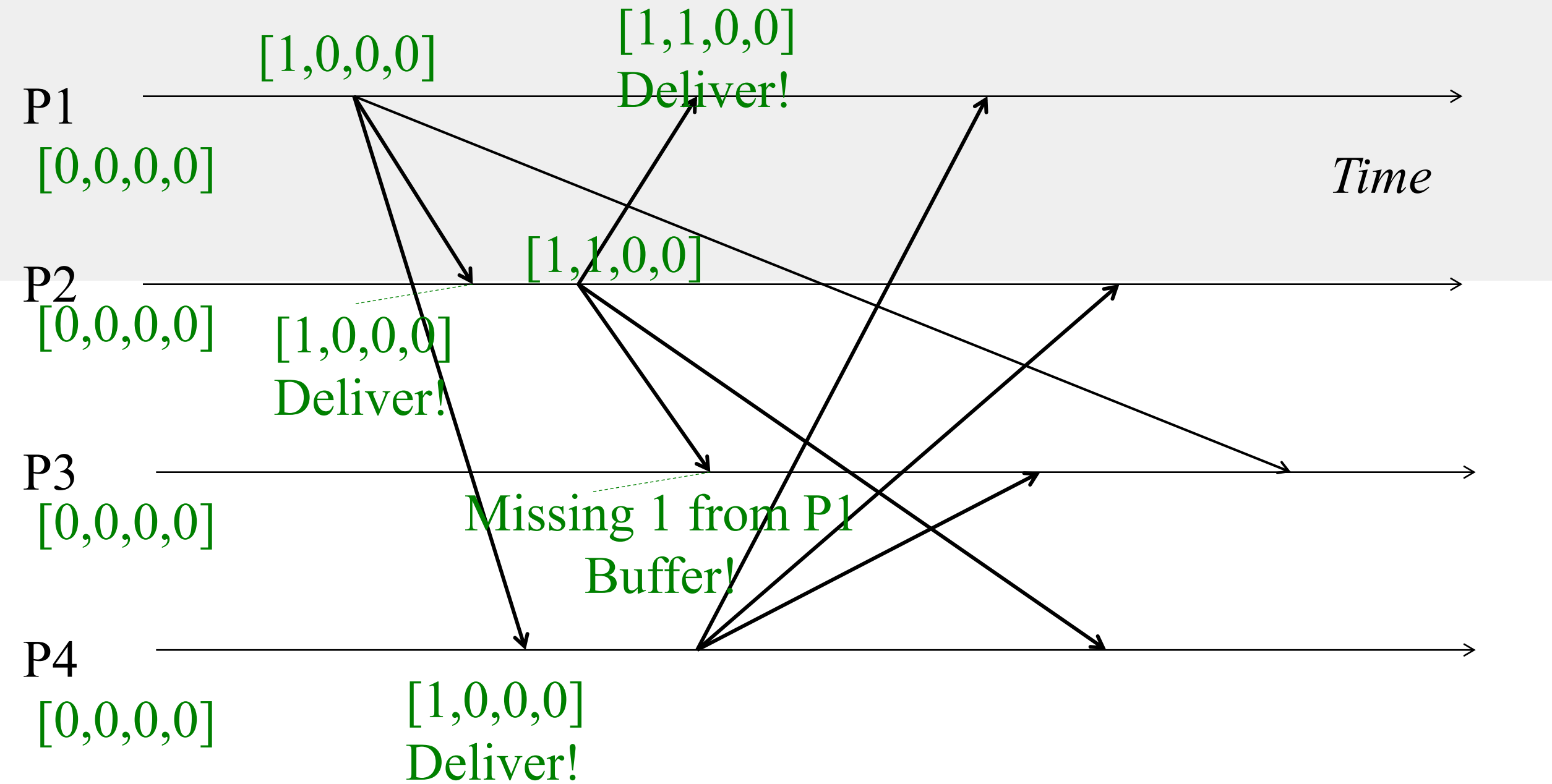


# Causal Ordering: Example

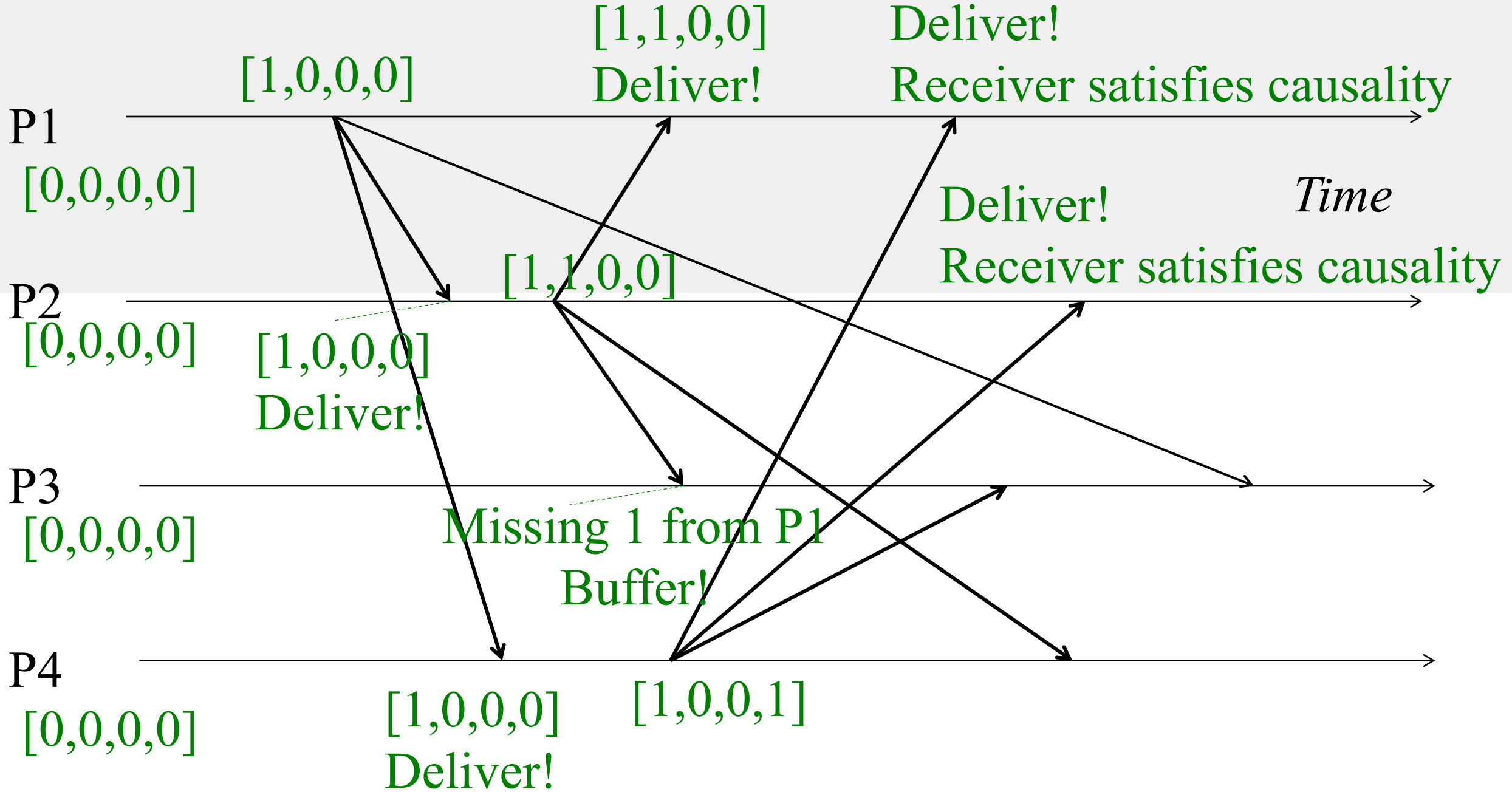




Causal Ordering: Example

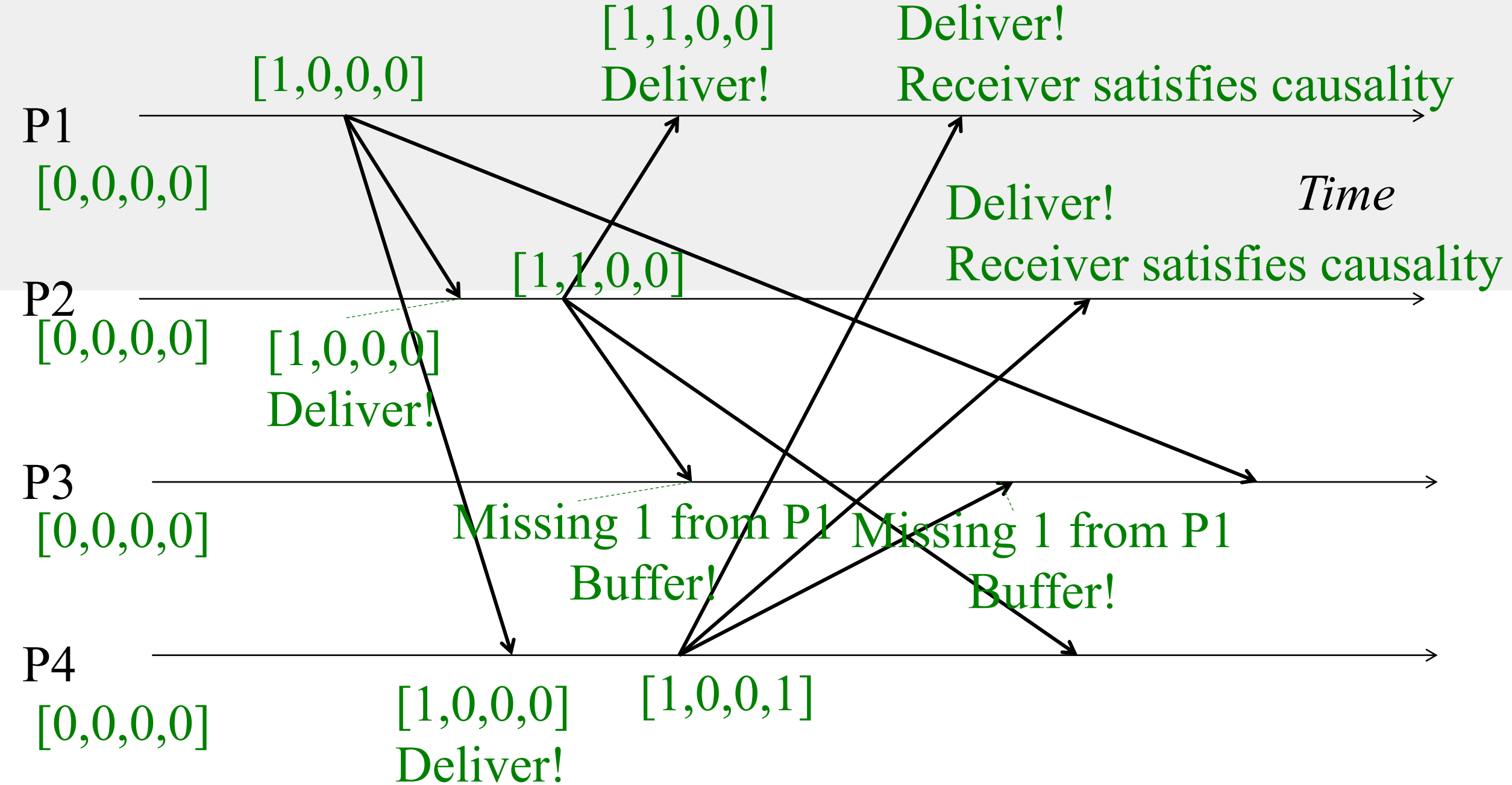


Causal Ordering: Example

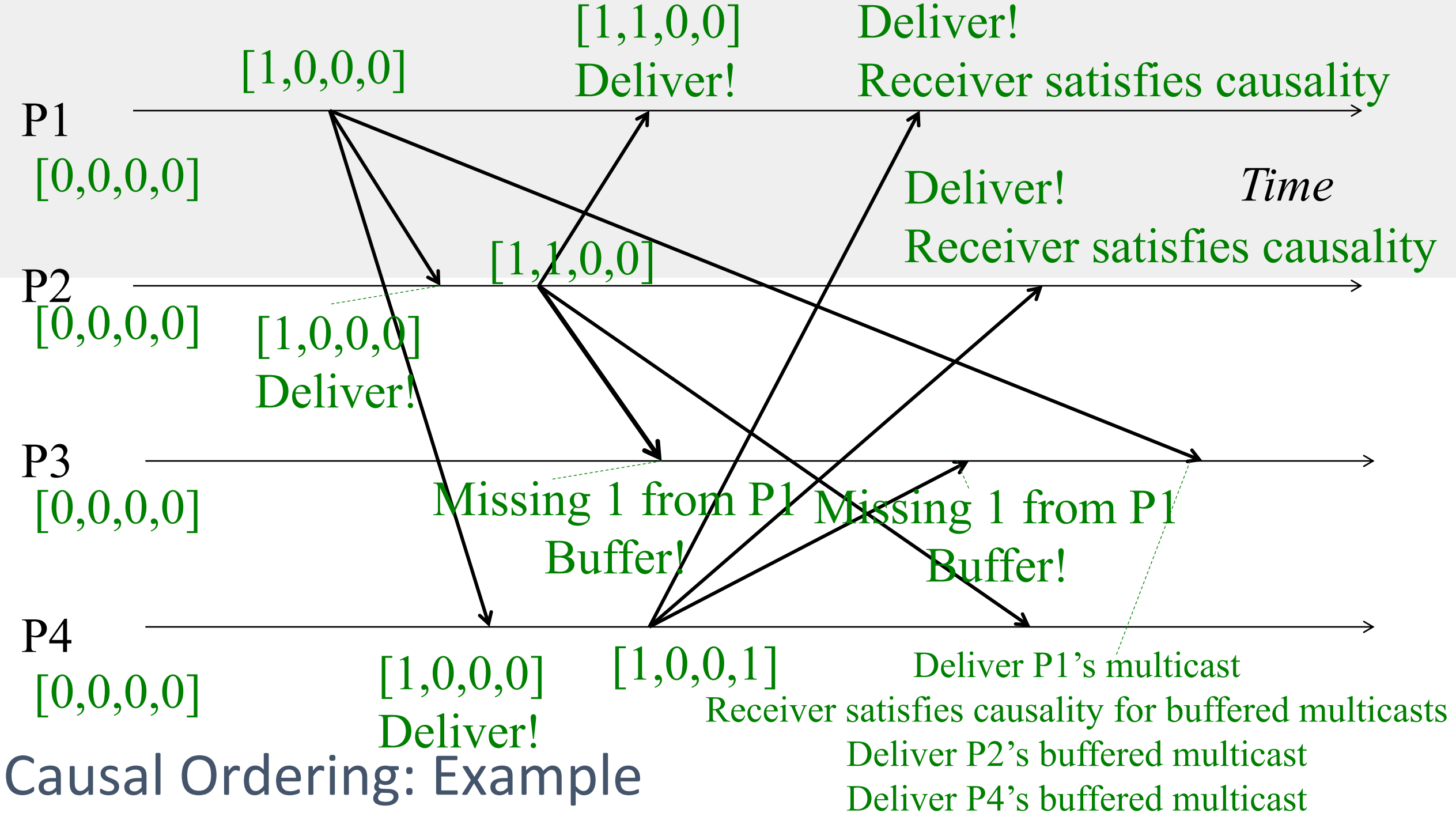


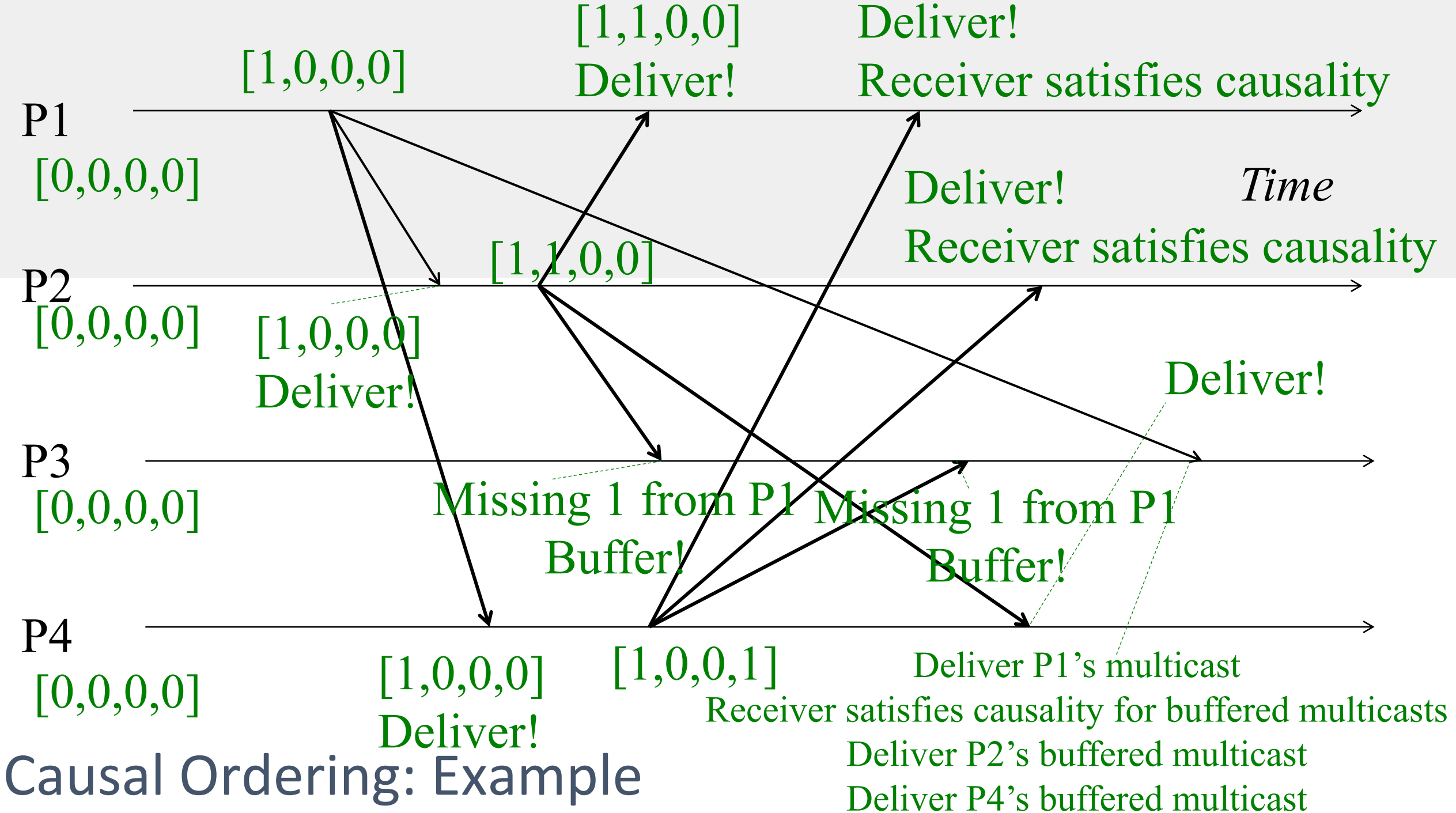
Causal Ordering: Example





Causal Ordering: Example





# Total Ordering

- Ensures all receivers receive all multicasts in the same order
- Formally
  - If a correct process  $P$  delivers message  $m$  before  $m'$  (independent of the senders), then any other correct process  $P'$  that delivers  $m'$  would already have delivered  $m$ .

# Sequencer-based Approach

- Special process elected as leader or sequencer
- Send multicast at process  $P_i$ :
  - Send multicast message  $M$  to group and sequencer
- Sequencer:
  - Maintains a global sequence number  $S$  (initially 0)
  - When it receives a multicast message  $M$ , it sets  $S = S + 1$ , and multicasts  $\langle M, S \rangle$
- Receive multicast at process  $P_i$ :
  - $P_i$  maintains a local received global sequence number  $S_i$  (initially 0)
  - If  $P_i$  receives a multicast  $M$  from  $P_j$ , it buffers it until it both
    1.  $P_i$  receives  $\langle M, S(M) \rangle$  from sequencer, and
    2.  $S_i + 1 = S(M)$
  - Then deliver its message to application and set  $S_i = S_i + 1$

# Summary: Multicast Ordering

- Ordering of multicasts affects correctness of distributed systems using multicasts
- Three popular ways of implementing ordering
  - FIFO, Causal, Total
- And their implementations
- What about reliability of multicasts?
- What about failures?