

# Distributed Systems

CS425/ECE428

*Instructor: Radhika Mittal*

*Acknowledgements for some of materials: Indy Gupta and Nikita Borisov*

# Logistics

- Note about exams on CampusWire:
  - Midterm 1 (Mar 4-6), Midterm 2 (April 12-14), Finals (May 7-15).
  - Reservation via PrairieTest.
    - You can reserve a slot for Midterm 1 starting Feb 19.
  - If you need DRES accommodations, please upload your Letter of Accommodations on the CBTF website.
  
- HW1 is due next week Friday.
  - You should be able to solve all questions by now.

# Today's agenda

- **Multicast**
  - Chapter 15.4

# Ordered Multicast

- Three popular flavors implemented by several multicast protocols:
  1. FIFO ordering
  2. Causal ordering
  3. Total ordering

# I. FIFO Order

- Multicasts from each sender are delivered in the order they are sent, at all receivers.
- Don't care about multicasts from different senders.
- More formally
  - *If a correct process issues  $\text{multicast}(g,m)$  and then  $\text{multicast}(g,m')$ , then every correct process that delivers  $m'$  will have already delivered  $m$ .*

## 2. Causal Order

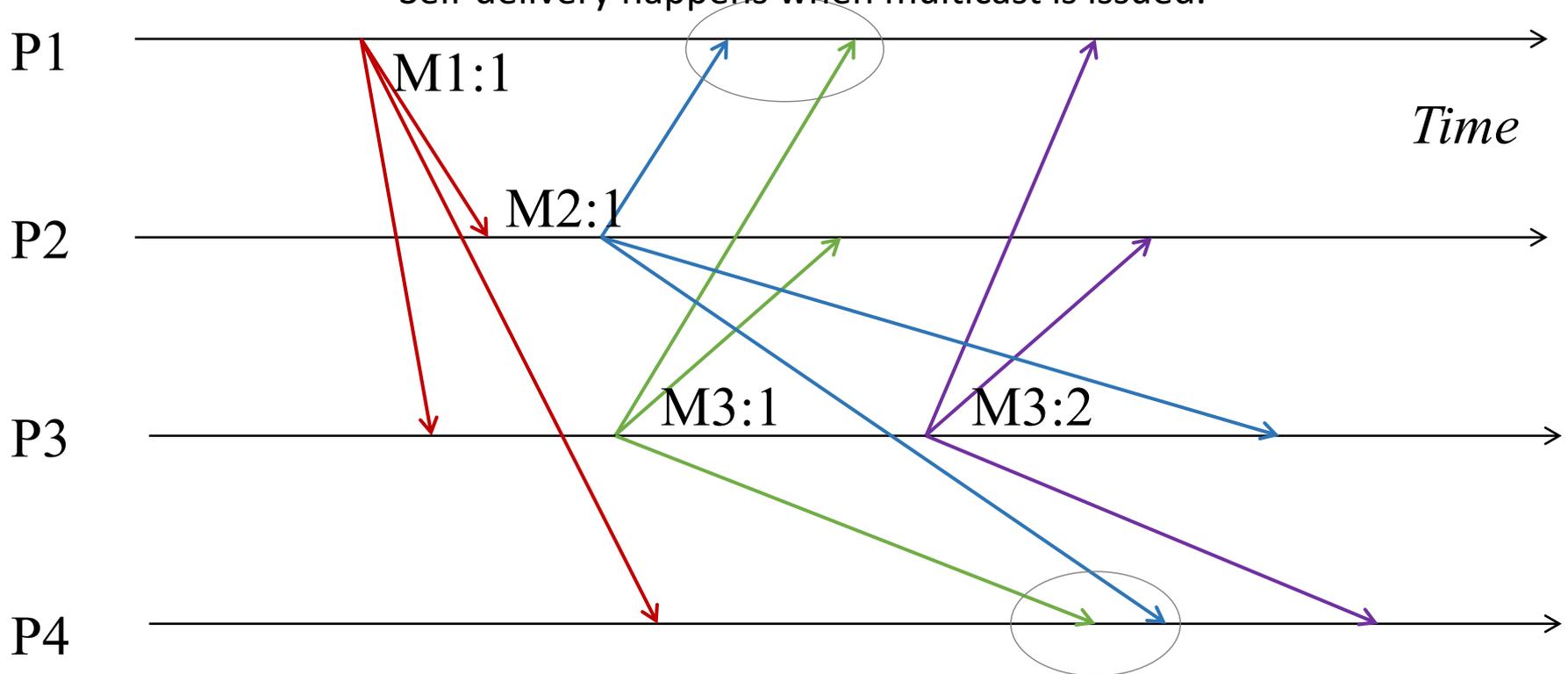
- Multicasts whose send events are causally related, must be delivered in the same causality-obeying order at all receivers.
- More formally
  - *If  $\text{multicast}(g,m) \rightarrow \text{multicast}(g,m')$  then any correct process that delivers  $m'$  will have already delivered  $m$ .*
  - $\rightarrow$  is Lamport's happens-before
  - $\rightarrow$  is induced only by multicast messages in group  $g$ , and when they are **delivered** to the application, rather than all network messages.

# HB Relationship for Causal Ordering

- HB rules in causal ordered multicast:
  - If  $\exists p_i, e \rightarrow_i e'$  then  $e \rightarrow e'$ .
    - If  $\exists p_i, \text{multicast}(g,m) \rightarrow_i \text{multicast}(g,m')$ , then  $\text{multicast}(g,m) \rightarrow \text{multicast}(g,m')$
    - If  $\exists p_i, \text{delivery}(m) \rightarrow_i \text{multicast}(g,m')$ , then  $\text{delivery}(m) \rightarrow \text{multicast}(g,m')$
    - ...
  - ~~For any message  $m$ ,  $\text{send}(m) \rightarrow \text{receive}(m)$~~ 
    - For any *multicast* message  $m$ ,  $\text{multicast}(g,m) \rightarrow \text{delivery}(m)$
  - If  $e \rightarrow e'$  and  $e' \rightarrow e''$  then  $e \rightarrow e''$ 
    - $\text{multicast}(g,m)$  at  $p_i \rightarrow \text{delivery}(m)$  at  $p_j$
    - $\text{delivery}(m)$  at  $p_j \rightarrow \text{multicast}(g,m')$  at  $p_j$
    - $\text{multicast}(g,m)$  at  $p_i \rightarrow \text{multicast}(g,m')$  at  $p_j$
- *Application can only see when messages are “multicast” by the application and “delivered” to the application, and not when they are sent or received by the protocol.*

# Causal Order: Example

Message delivery indicated by arrow endings.  
Self-delivery happens when multicast is issued.



M3:1 → M3:2, M1:1 → M2:1, M1:1 → M3:1 and so should be delivered in that order at each receiver.

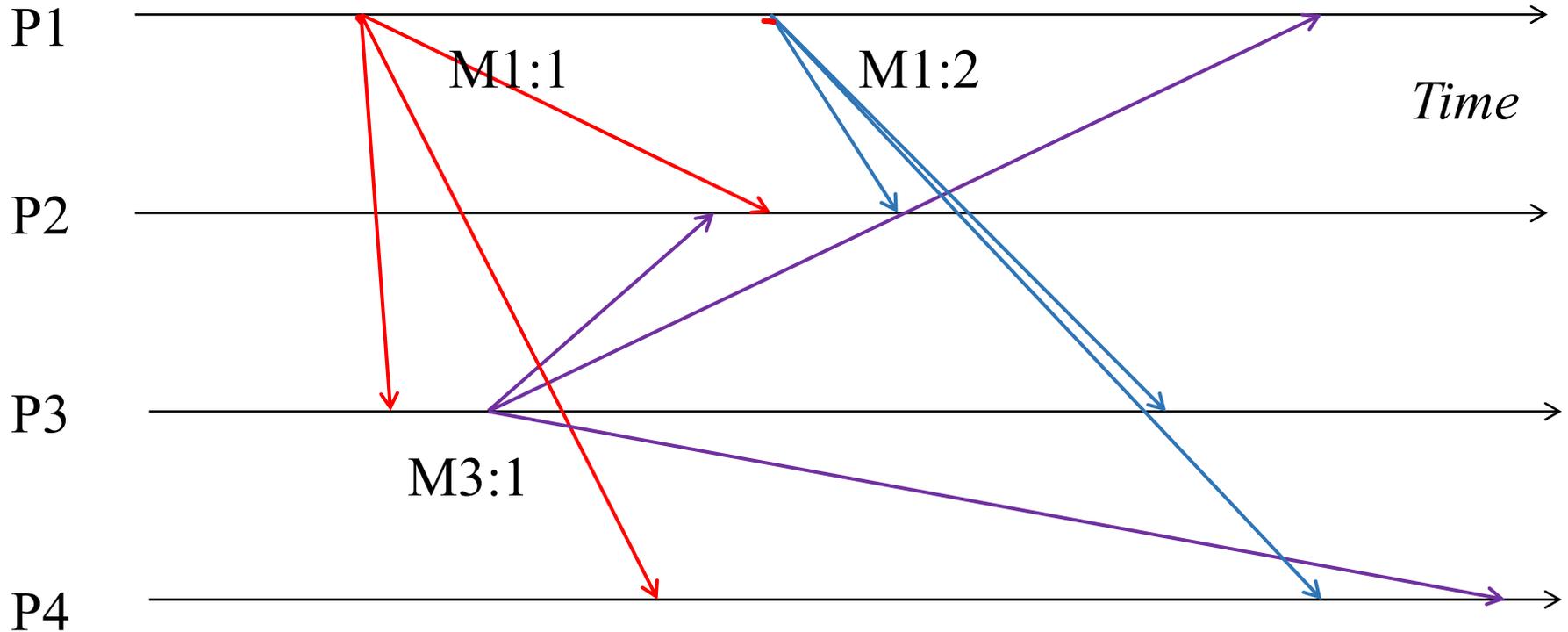
M3:1 and M2:1 are concurrent and thus ok to be delivered in any (and even different) orders at different receivers.

# Causal vs FIFO

- Does Causal Ordering imply FIFO Ordering?
  - Yes
  
- Does FIFO Order imply Causal Order?
  - No

# Example

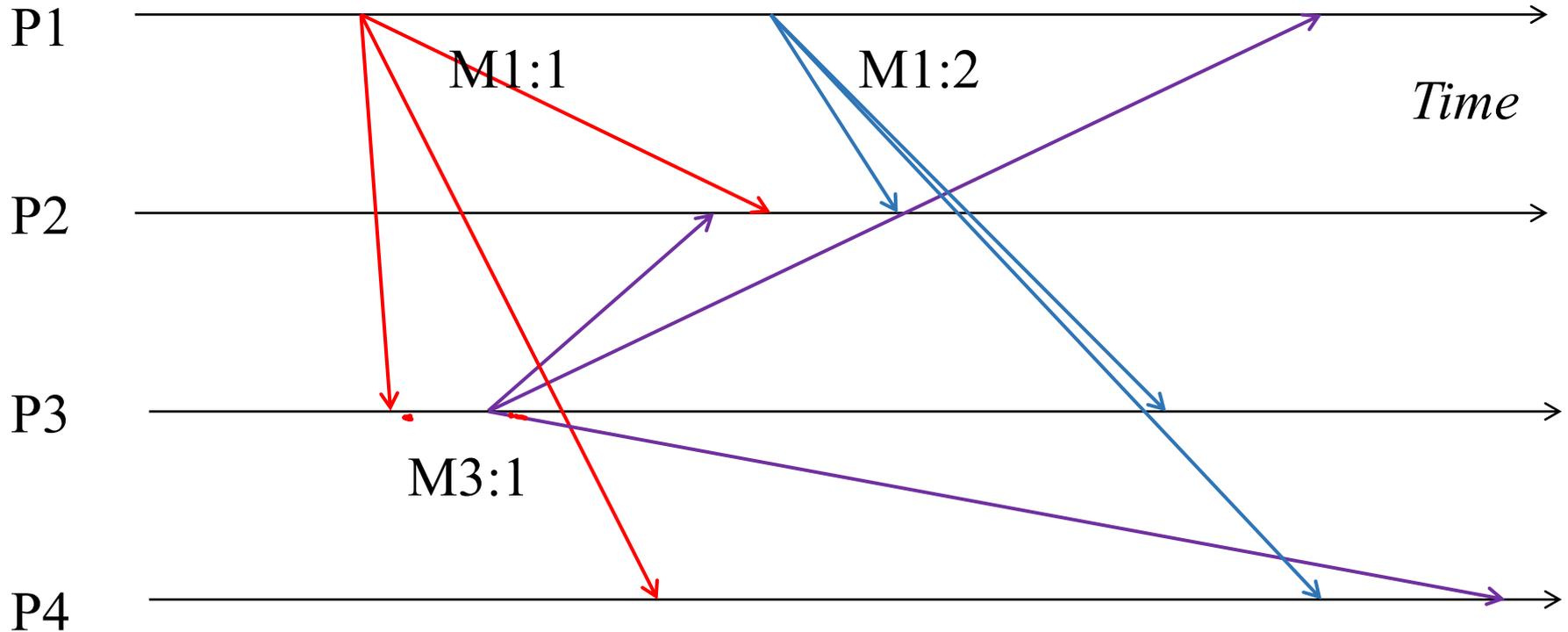
Message delivery indicated by arrow endings.  
Self-delivery happens when multicast is issued.



Does this satisfy FIFO order?

# Example

Message delivery indicated by arrow endings.  
Self-delivery happens when multicast is issued.

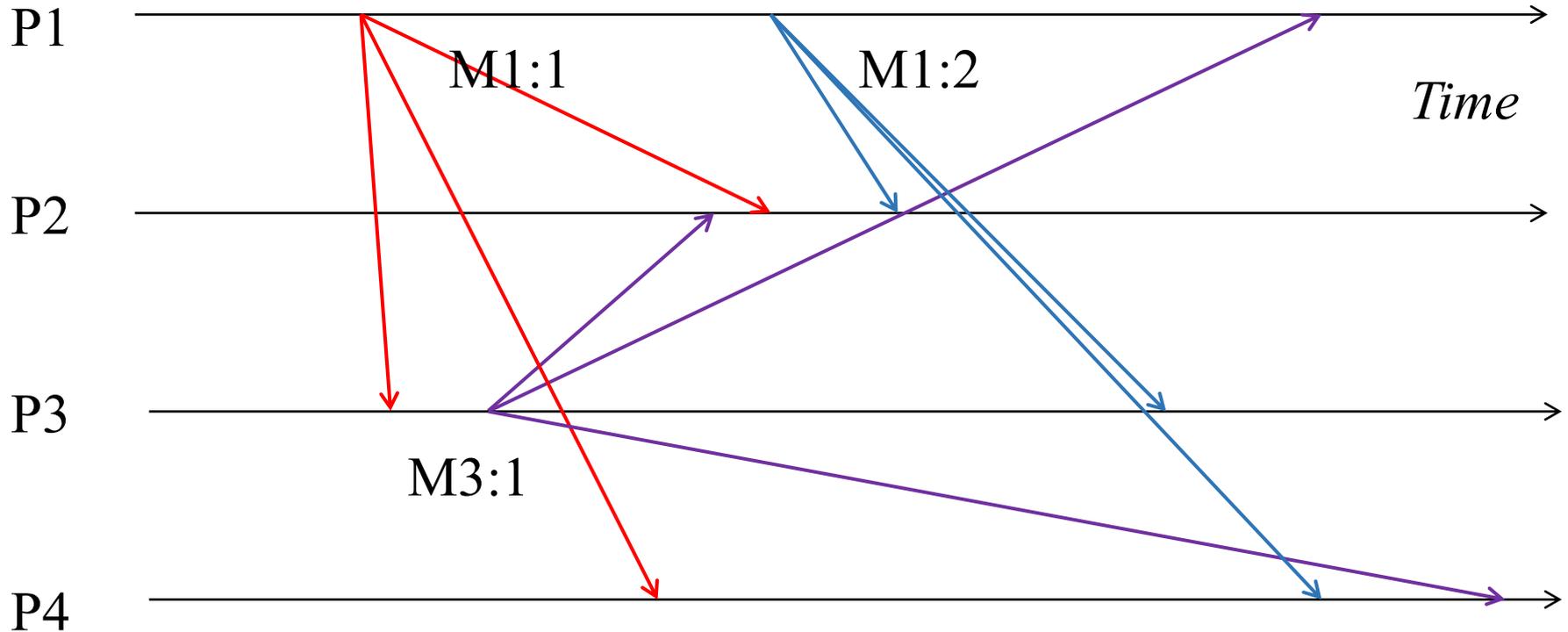


Does this satisfy FIFO order?

Yes

# Example

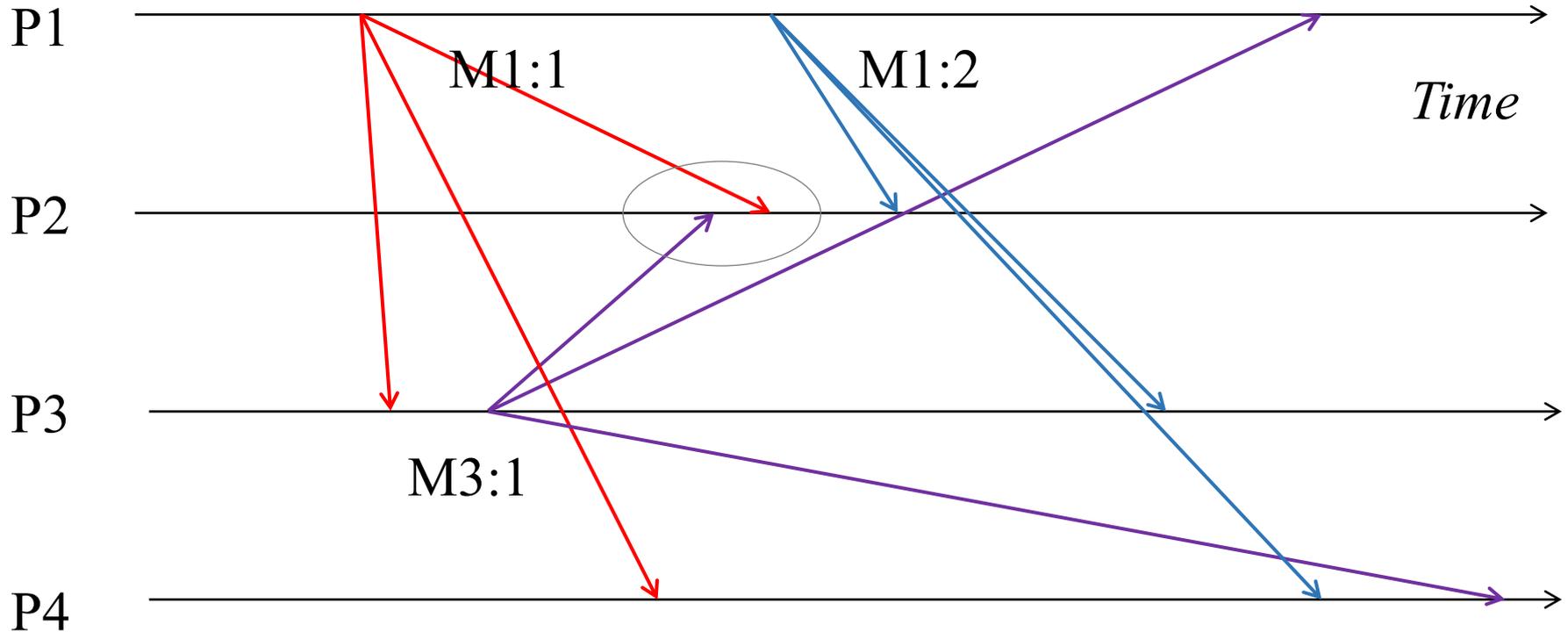
Message delivery indicated by arrow endings.  
Self-delivery happens when multicast is issued.



Does this satisfy causal order?

# Example

Message delivery indicated by arrow endings.  
Self-delivery happens when multicast is issued.

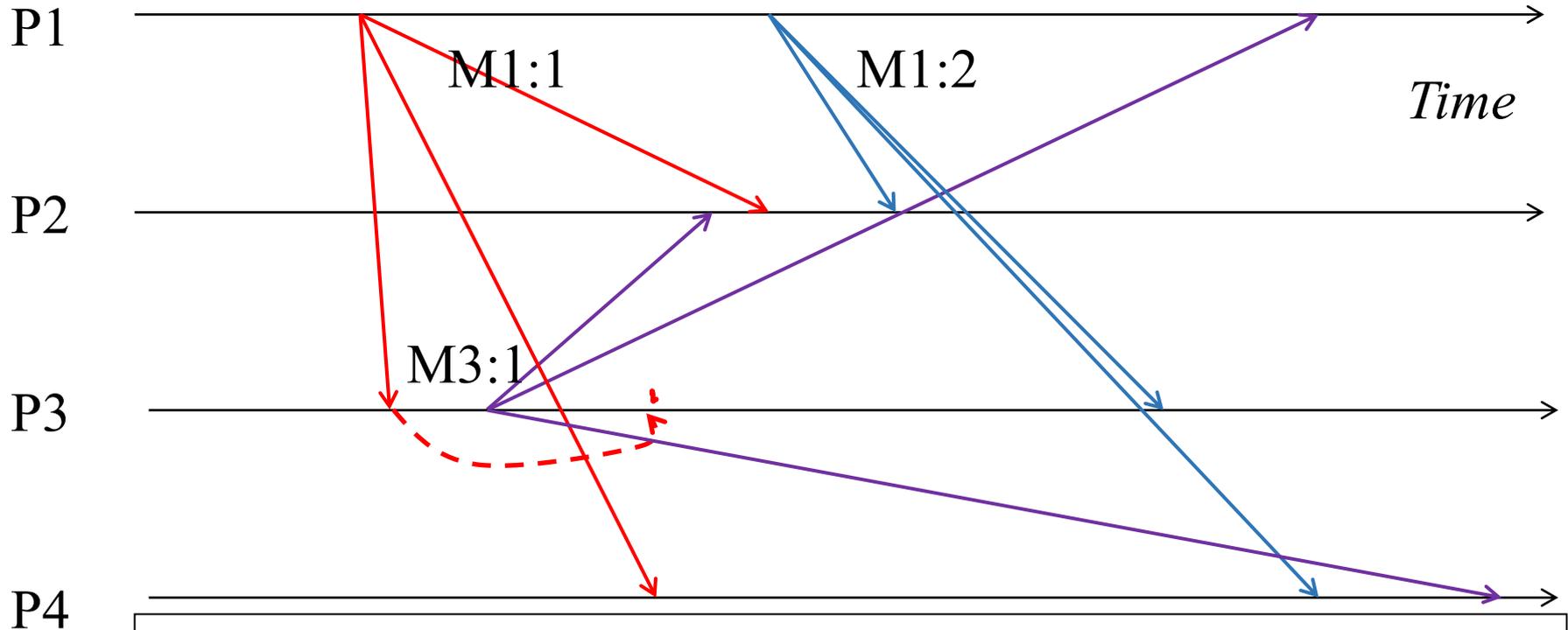


Does this satisfy causal order?

**No**

# Example

Message delivery indicated by (extended) arrow endings.  
Self-delivery happens when multicast is issued.



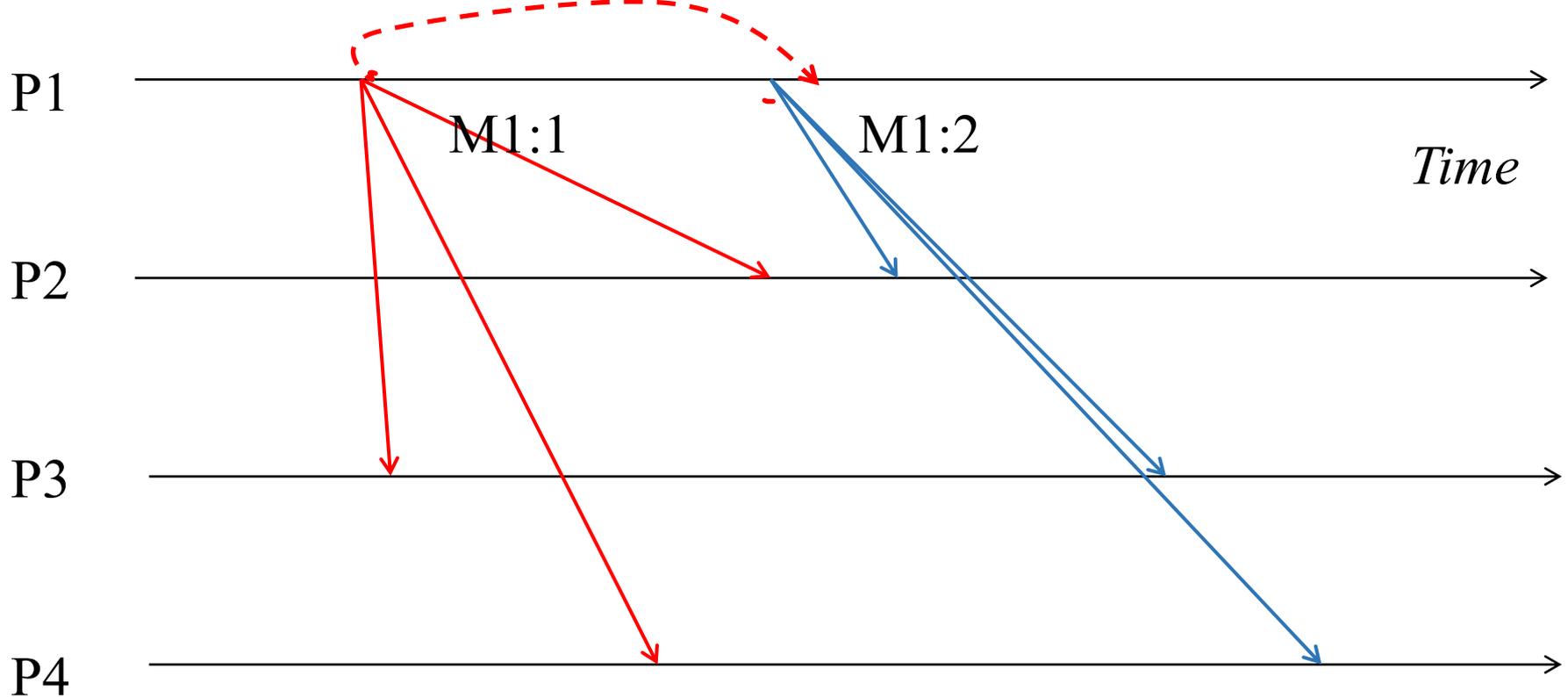
M1:1 is delivered at P3 after M3:1's multicast.

**Does this satisfy causal order?**

**Yes**

# Example

Message delivery indicated by (extended) arrow endings.  
Self-delivery happens when multicast is issued when no extra arrow.

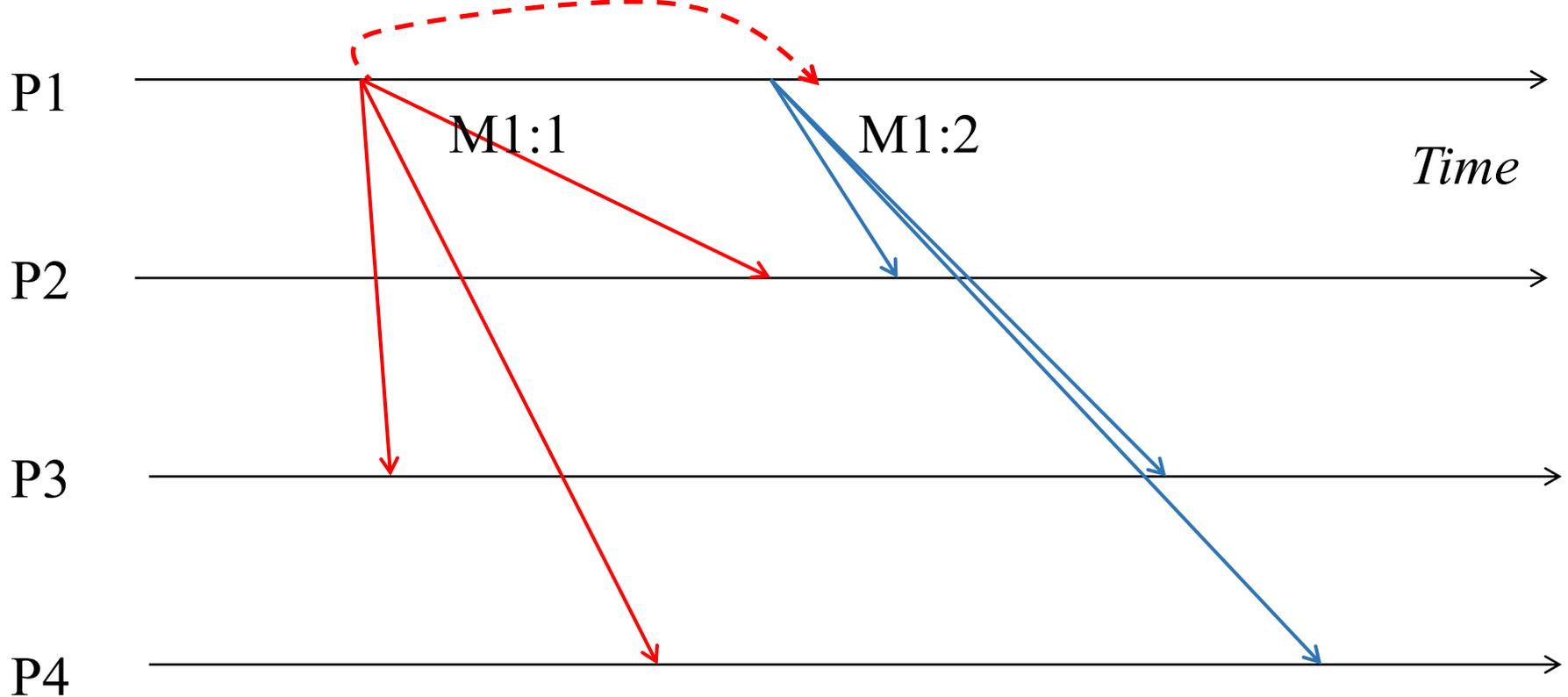


Does this satisfy causal order?

No

# Example

Message delivery indicated by (extended) arrow endings.  
Self-delivery happens when multicast is issued when no extra arrow.



Does this satisfy FIFO order?

No

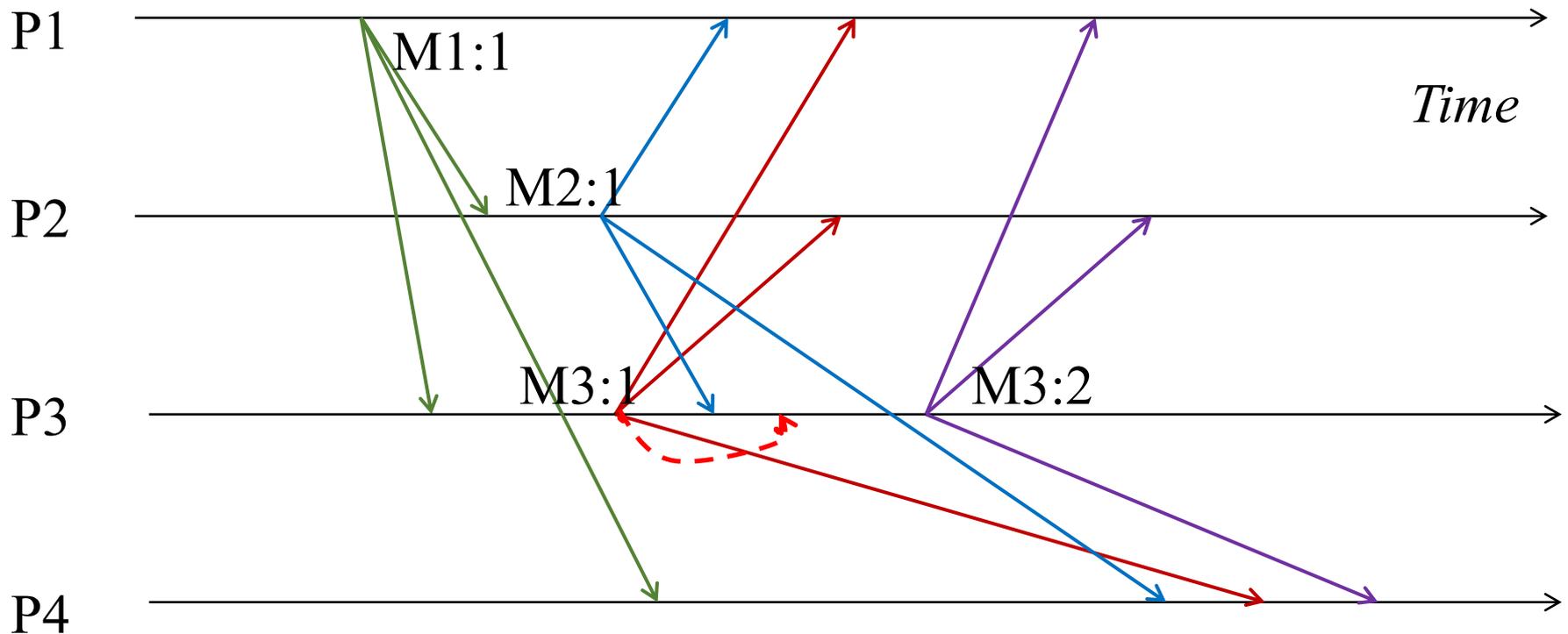
# 3. Total Order

- Ensures all processes deliver all multicasts in the same order.
- Unlike FIFO and causal, this does not pay attention to order of multicast sending.
- Formally
  - If a correct process delivers message  $m$  before  $m'$  (independent of the senders), then any other correct process that delivers  $m'$  will have already delivered  $m$ .

# Total Order: Example

Message delivery indicated by (extended) arrow endings.

Self-delivery happens when multicast is issued (when no extra arrow).



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.

# Causal vs Total

- Total ordering does not imply causal ordering.
- Causal ordering does not imply total ordering.

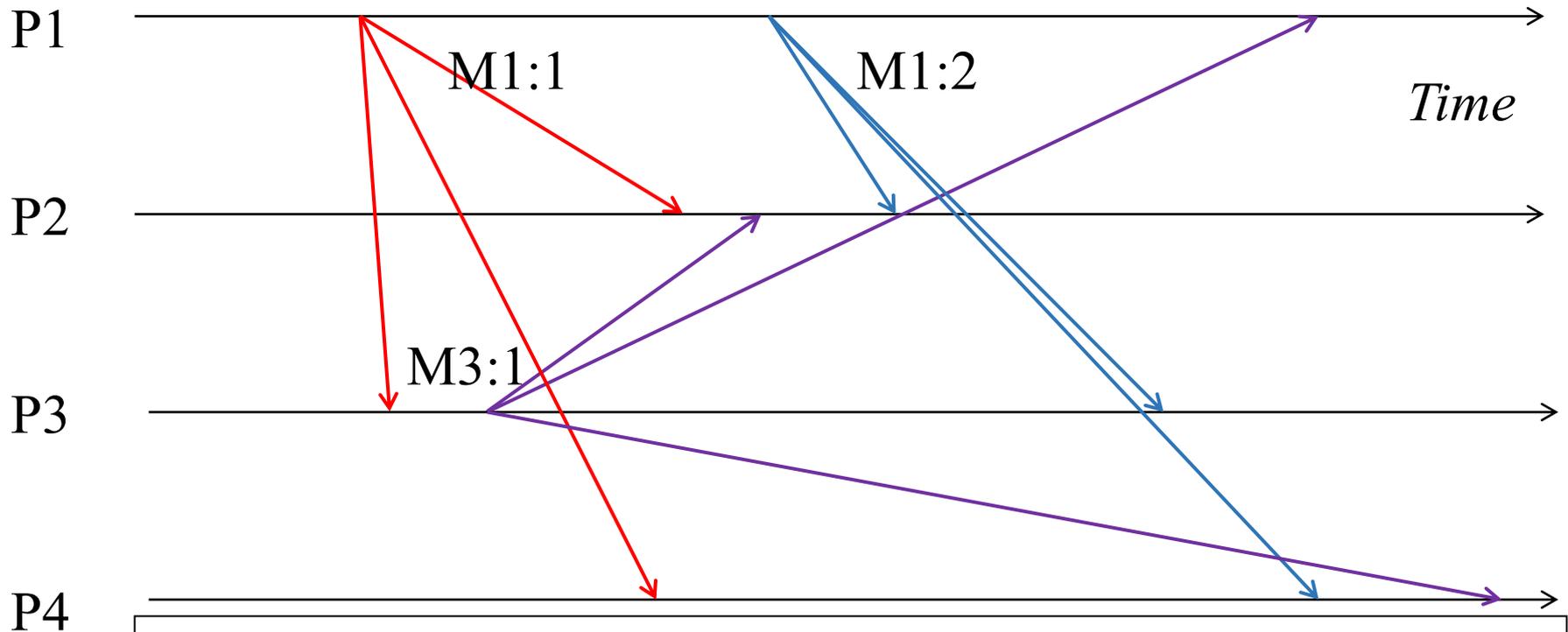
# Hybrid variants

- We can have hybrid ordering protocols:
  - Causal-total hybrid protocol satisfies both Causal and total orders.

# Example

Message delivery indicated by (extended) arrow endings.

Self-delivery happens when multicast is issued (when no extra arrow).



Does this satisfy causal (and FIFO) order?

Yes

# Example

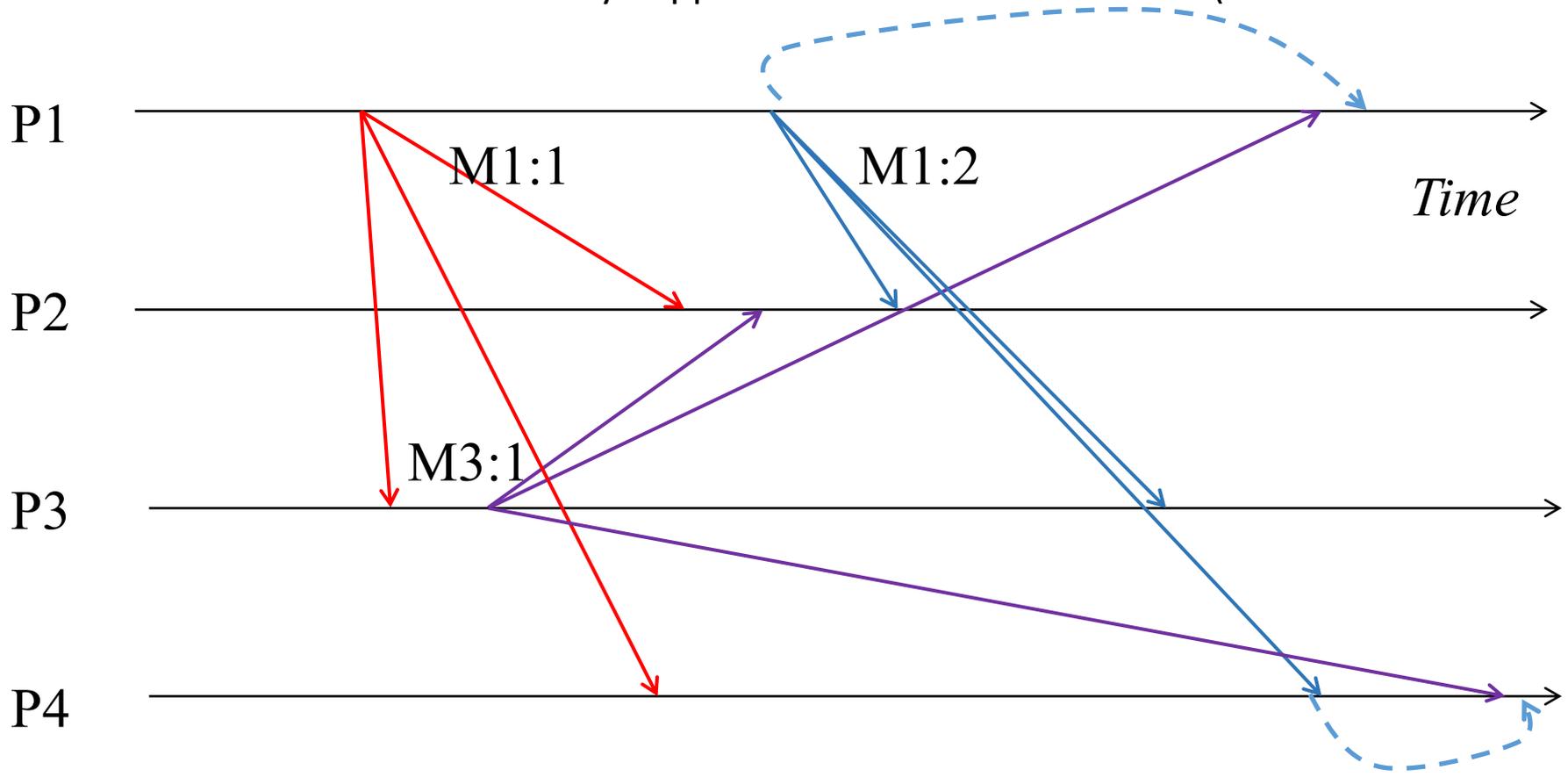
Message delivery indicated by (extended) arrow endings.

Self-delivery happens when multicast is issued (when no extra arrow).



# Example

Message delivery indicated by (extended) arrow endings.  
Self-delivery happens when multicast is issued (when no extra arrow).



Does this satisfy total order?

Yes

# Ordered Multicast

- **FIFO ordering:** If a correct process issues  $\text{multicast}(g,m)$  and then  $\text{multicast}(g,m')$ , then every correct process that delivers  $m'$  will have already delivered  $m$ .
- **Causal ordering:** If  $\text{multicast}(g,m) \rightarrow \text{multicast}(g,m')$  then any correct process that delivers  $m'$  will have already delivered  $m$ .
  - Note that  $\rightarrow$  counts messages **delivered** to the application, rather than all network messages.
- **Total ordering:** If a correct process delivers message  $m$  before  $m'$  (independent of the senders), then any other correct process that delivers  $m'$  will have already delivered  $m$ .

# Next Question

*How do we implement ordered multicast?*

# Ordered Multicast

- **FIFO ordering**

- If a correct process issues  $\text{multicast}(g,m)$  and then  $\text{multicast}(g,m')$ , then every correct process that delivers  $m'$  will have already delivered  $m$ .

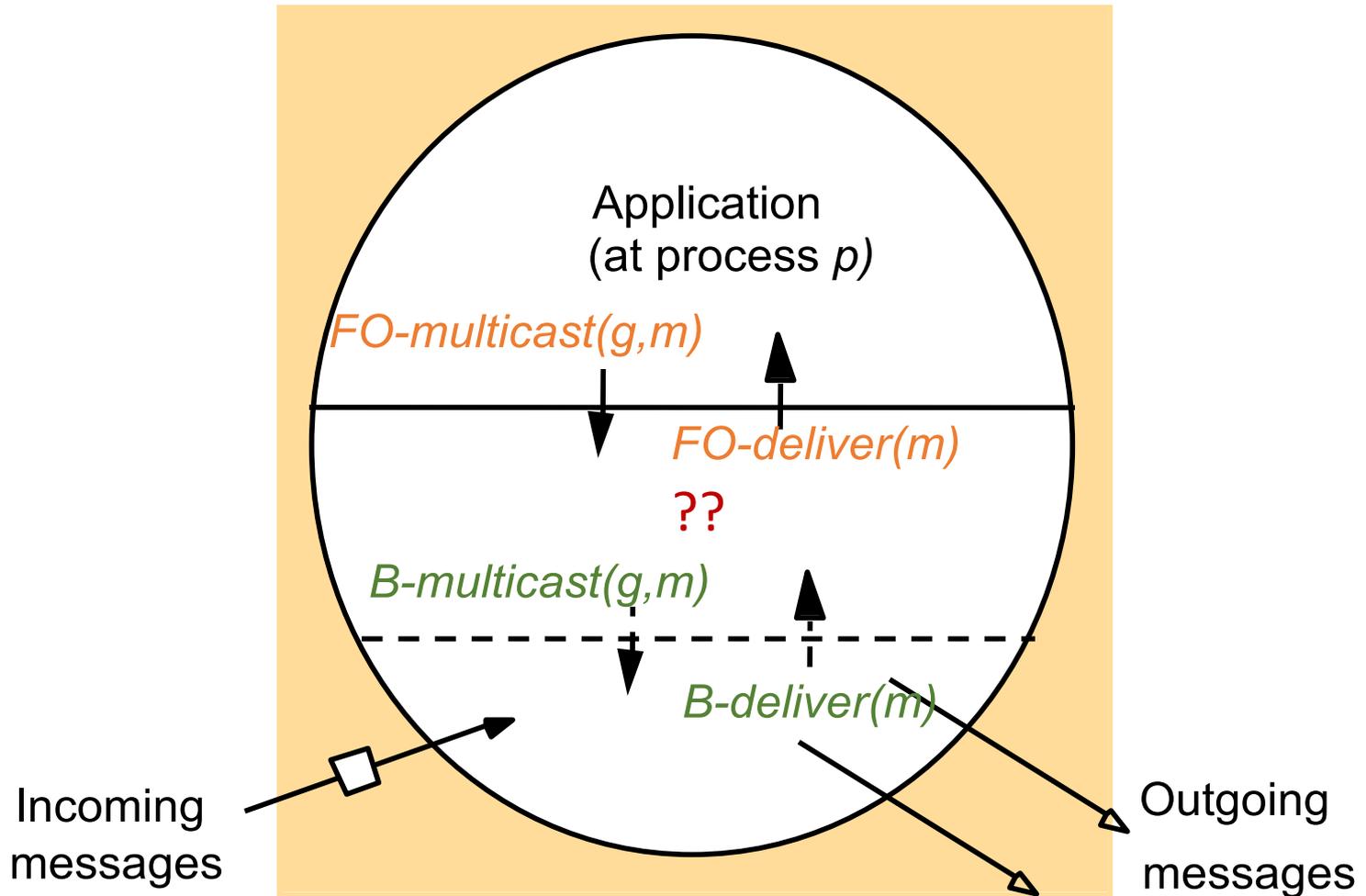
- **Causal ordering**

- If  $\text{multicast}(g,m) \rightarrow \text{multicast}(g,m')$  then any correct process that delivers  $m'$  will have already delivered  $m$ .
- Note that  $\rightarrow$  counts messages **delivered** to the application, rather than all network messages.

- **Total ordering**

- If a correct process delivers message  $m$  before  $m'$  (independent of the senders), then any other correct process that delivers  $m'$  will have already delivered  $m$ .

# Implementing FIFO order multicast



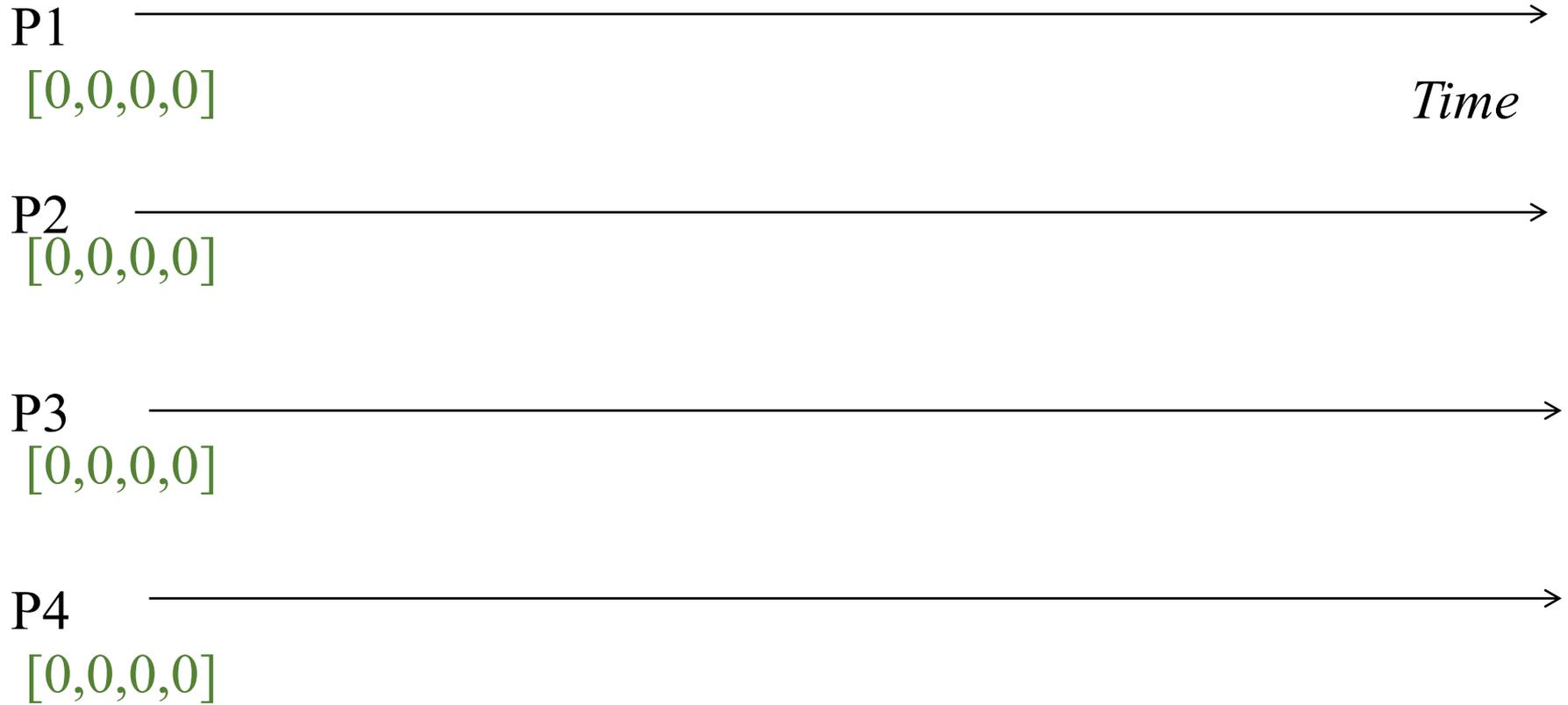
# Implementing FIFO order multicast

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

# Implementing FIFO order multicast

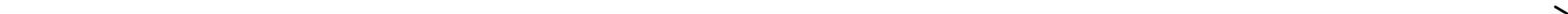
- On FO-multicast( $g, m$ ) at process  $P_j$ :
  - set  $P_j[j] = P_j[j] + 1$
  - piggyback  $P_j[j]$  with  $m$  as its sequence number.
  - B-multicast( $g, \{m, P_j[j]\}$ )
- On B-deliver( $\{m, S\}$ ) at  $P_i$  from  $P_j$ : *If  $P_i$  receives a multicast from  $P_j$  with sequence number  $S$  in message*
  - if ( $S == P_i[j] + 1$ ) then
    - FO-deliver( $m$ ) to application
    - set  $P_i[j] = P_i[j] + 1$
  - else buffer this multicast until above condition is true

# FIFO order multicast execution



# FIFO order multicast execution

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

P2   
[0,0,0,0]

P3   
[0,0,0,0]

P4   
[0,0,0,0]

## Sequence Vector

*Do not confuse with vector timestamps!*

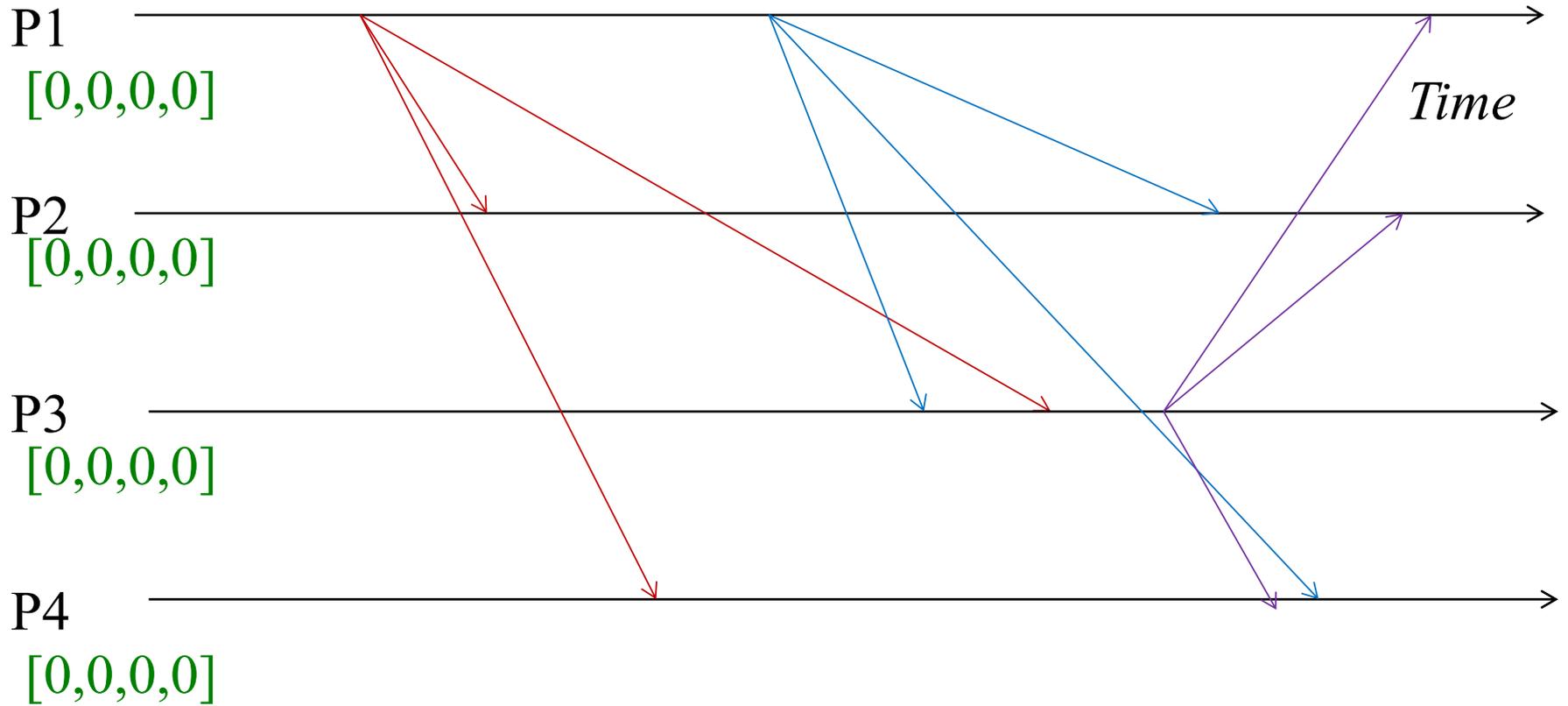
$P_i[i]$ , is the no. of messages  $P_i$  multicast (and delivered to itself).

$P_i[j] \forall j \neq i$  is no. of messages delivered at  $P_i$  from  $P_j$ .

# FIFO order multicast execution

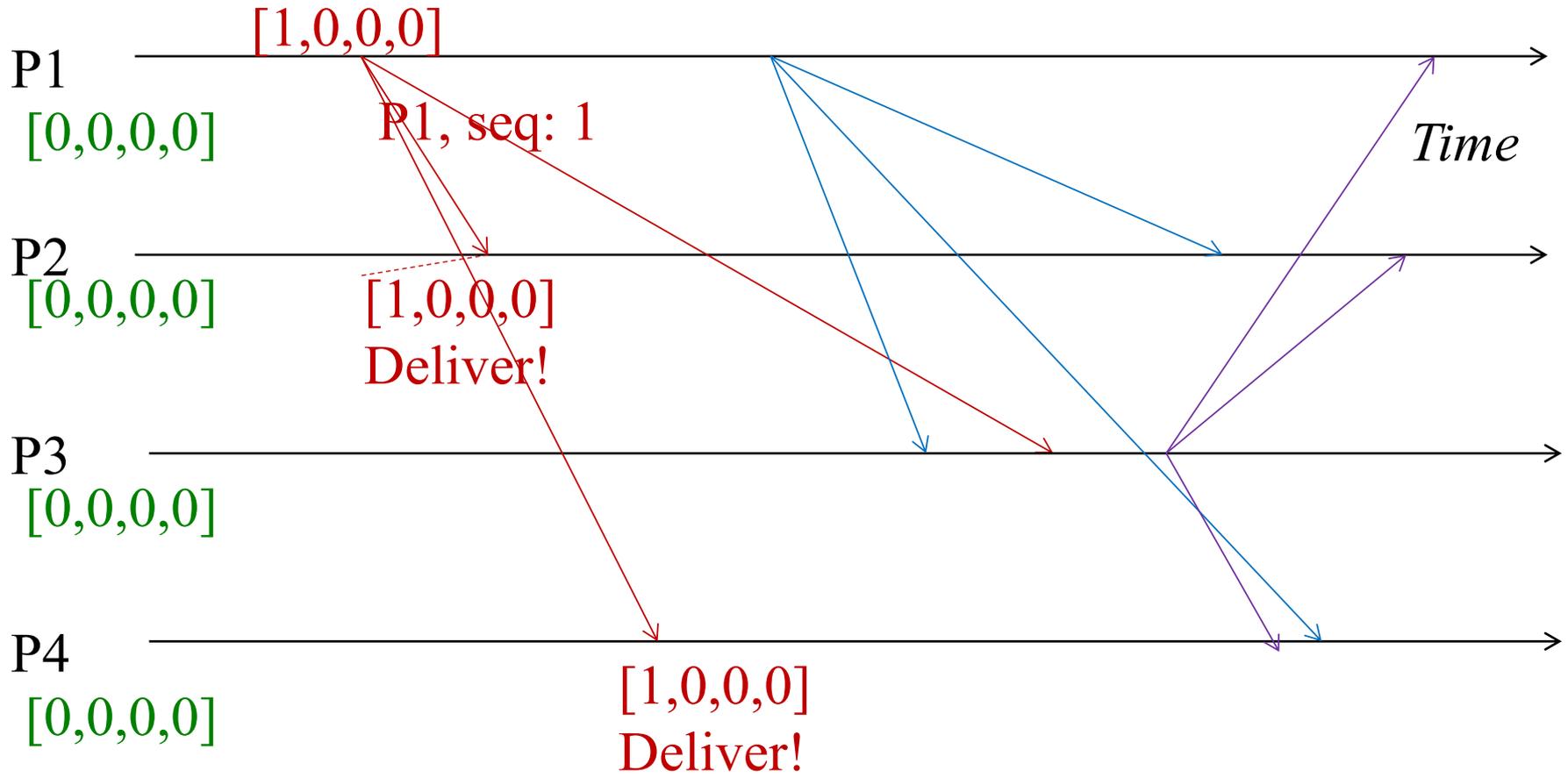


# FIFO order multicast execution

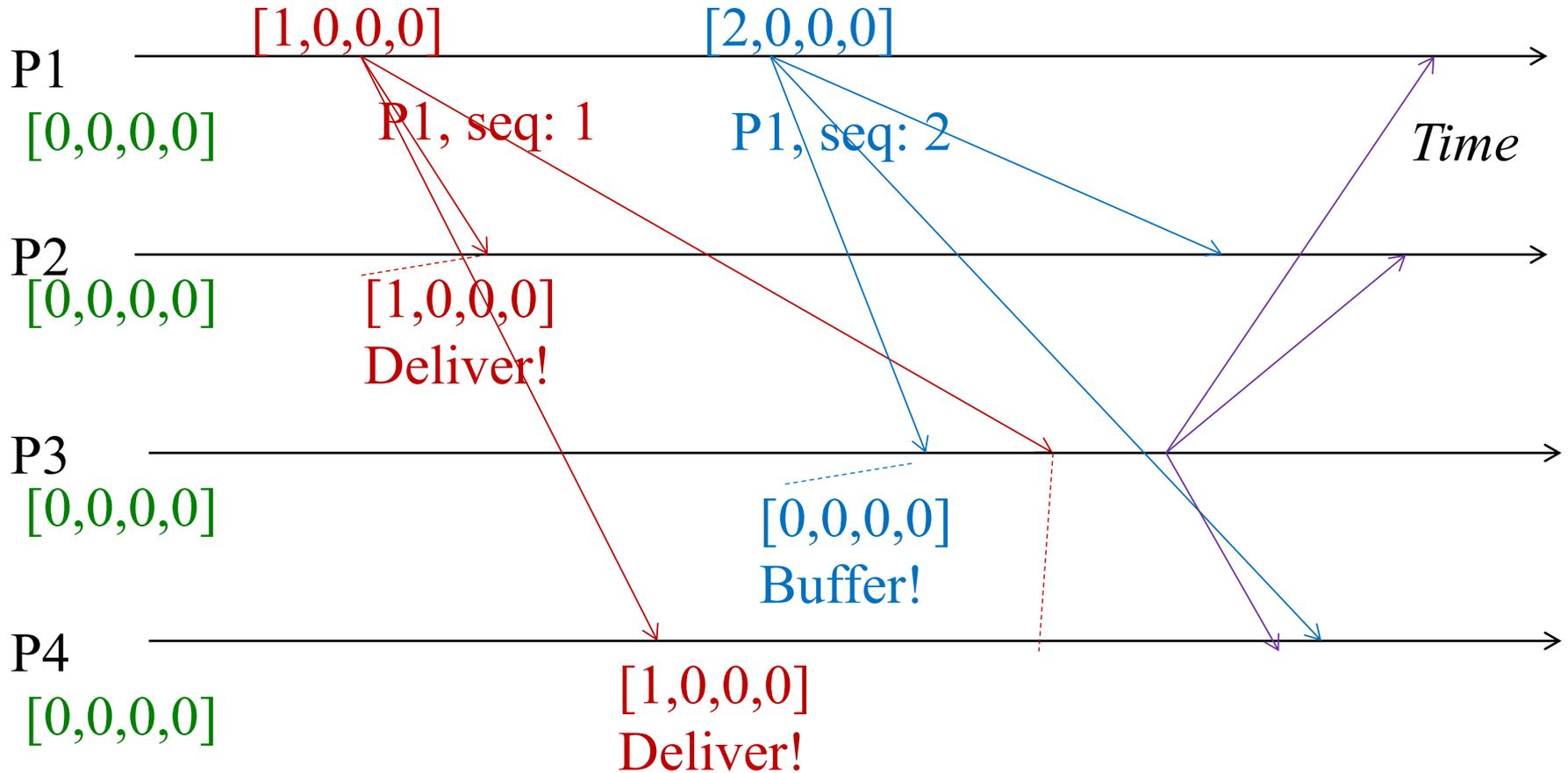


Self-deliveries omitted for simplicity.

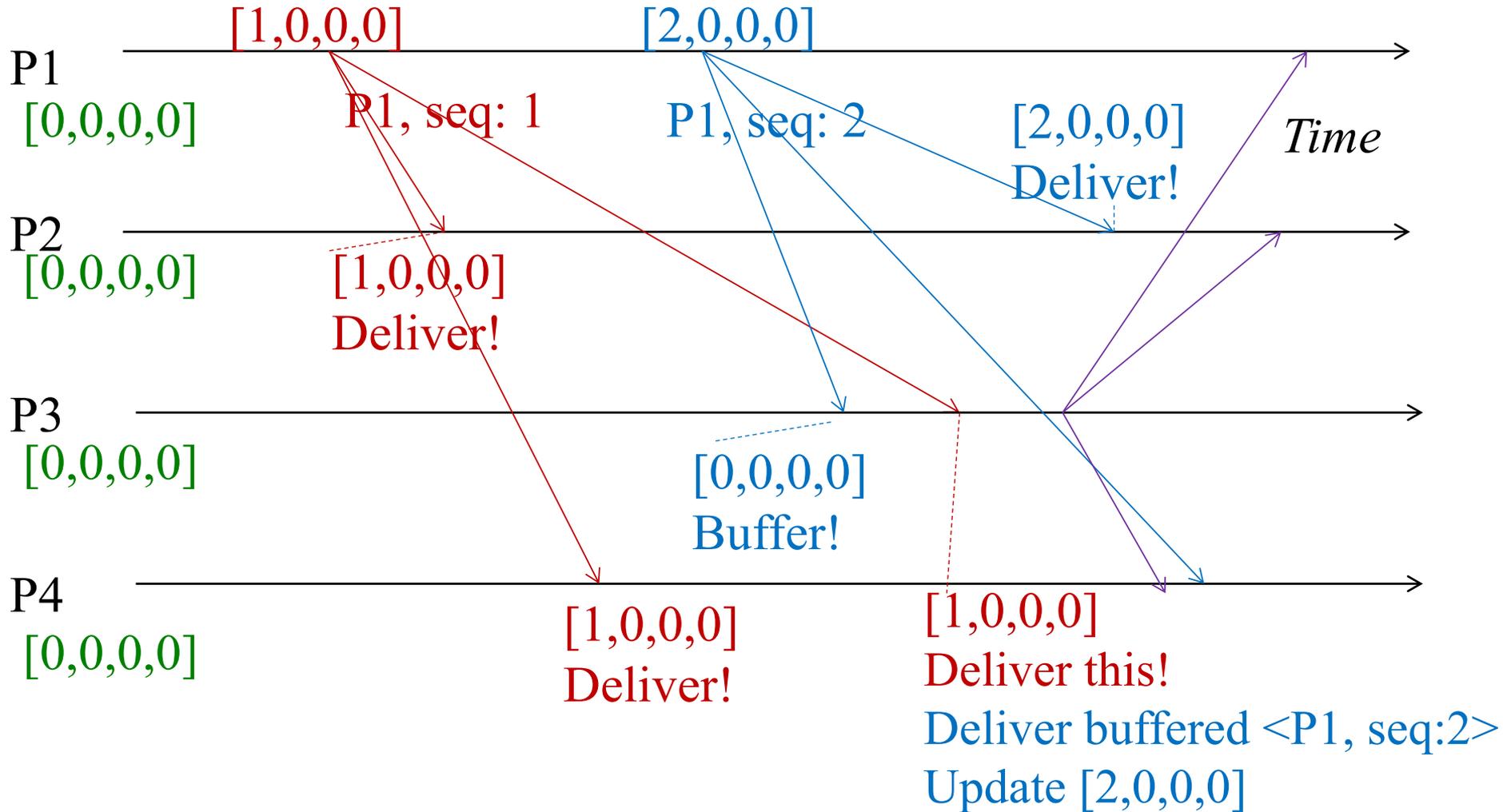
# FIFO order multicast execution



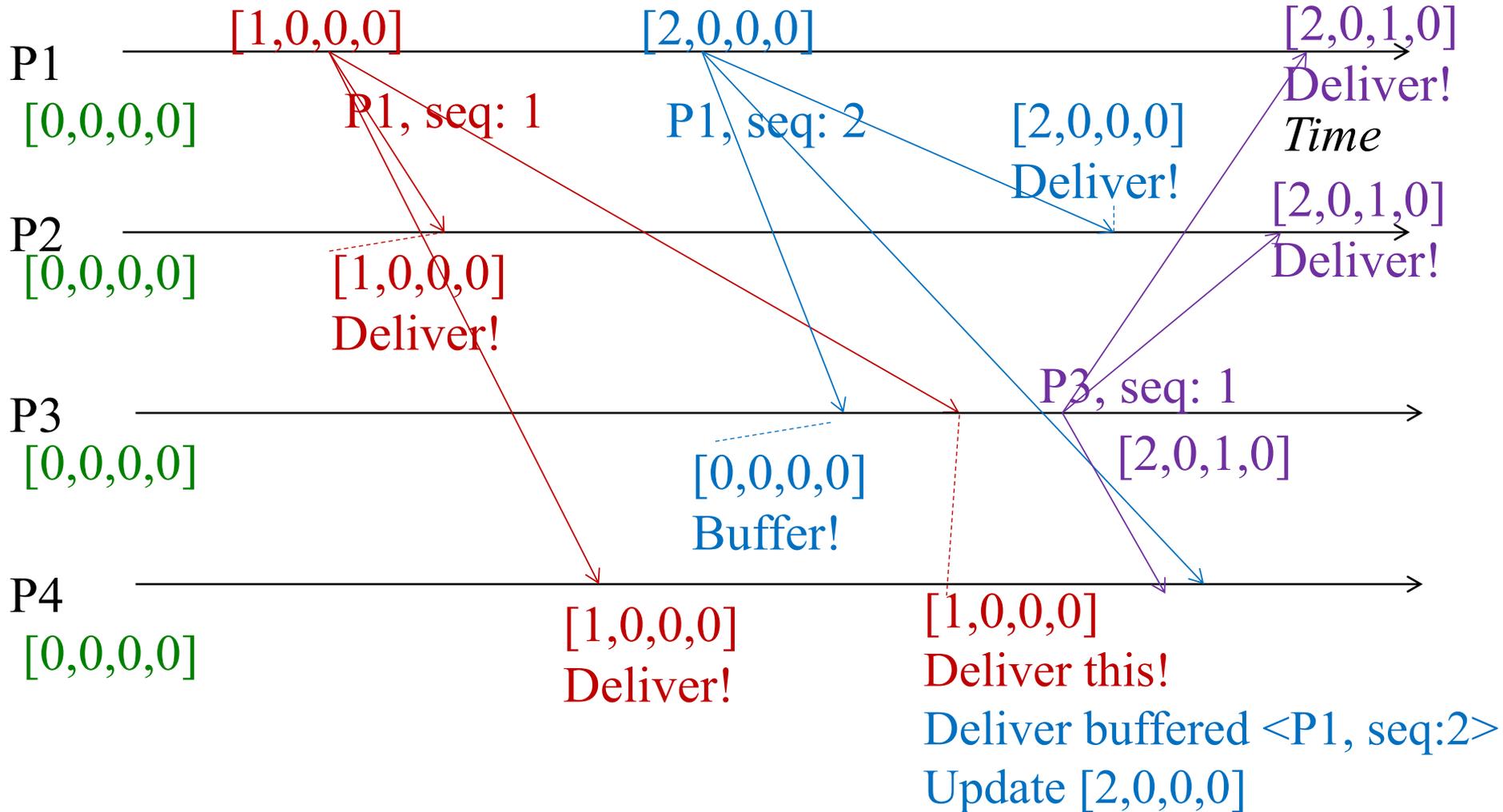
# FIFO order multicast execution



# FIFO order multicast execution



# FIFO order multicast execution





# Implementing FIFO order multicast

- On FO-multicast( $g, m$ ) at process  $P_j$ :
  - set  $P_j[j] = P_j[j] + 1$
  - piggyback  $P_j[j]$  with  $m$  as its sequence number.
  - B-multicast( $g, \{m, P_j[j]\}$ )
- On B-deliver( $\{m, S\}$ ) at  $P_i$  from  $P_j$ : *If  $P_i$  receives a multicast from  $P_j$  with sequence number  $S$  in message*
  - if ( $S == P_i[j] + 1$ ) then
    - FO-deliver( $m$ ) to application
    - set  $P_i[j] = P_i[j] + 1$
  - else buffer this multicast until above condition is true

# Implementing FIFO reliable multicast

- On FO-multicast( $g, m$ ) at process  $P_j$ :
  - set  $P_j[j] = P_j[j] + 1$
  - piggyback  $P_j[j]$  with  $m$  as its sequence number.
  - R-multicast( $g, \{m, P_j[j]\}$ )**
- On **R-deliver( $\{m, S\}$ )** at  $P_i$  from  $P_j$ : *If  $P_i$  receives a multicast from  $P_j$  with sequence number  $S$  in message*
  - if ( $S == P_i[j] + 1$ ) then
    - FO-deliver( $m$ ) to application
    - set  $P_i[j] = P_i[j] + 1$
  - else buffer this multicast until above condition is true

# Ordered Multicast

- **FIFO ordering:** If a correct process issues  $\text{multicast}(g,m)$  and then  $\text{multicast}(g,m')$ , then every correct process that delivers  $m'$  will have already delivered  $m$ .
- **Causal ordering:** If  $\text{multicast}(g,m) \rightarrow \text{multicast}(g,m')$  then any correct process that delivers  $m'$  will have already delivered  $m$ .
  - Note that  $\rightarrow$  counts messages **delivered** to the application, rather than all network messages.
- **Total ordering:** If a correct process delivers message  $m$  before  $m'$  (independent of the senders), then any other correct process that delivers  $m'$  will have already delivered  $m$ .

# Implementing total order multicast

- Basic idea:
  - Same sequence number counter across different processes.
  - Instead of different sequence number counter for each process.
- Two types of approach
  - Using a centralized sequencer
  - A decentralized mechanism (ISIS)

# Implementing total order multicast

- Basic idea:
  - Same sequence number counter across different processes.
  - Instead of different sequence number counter for each process.
- Two types of approach
  - **Using a centralized sequencer**
  - A decentralized mechanism (ISIS)

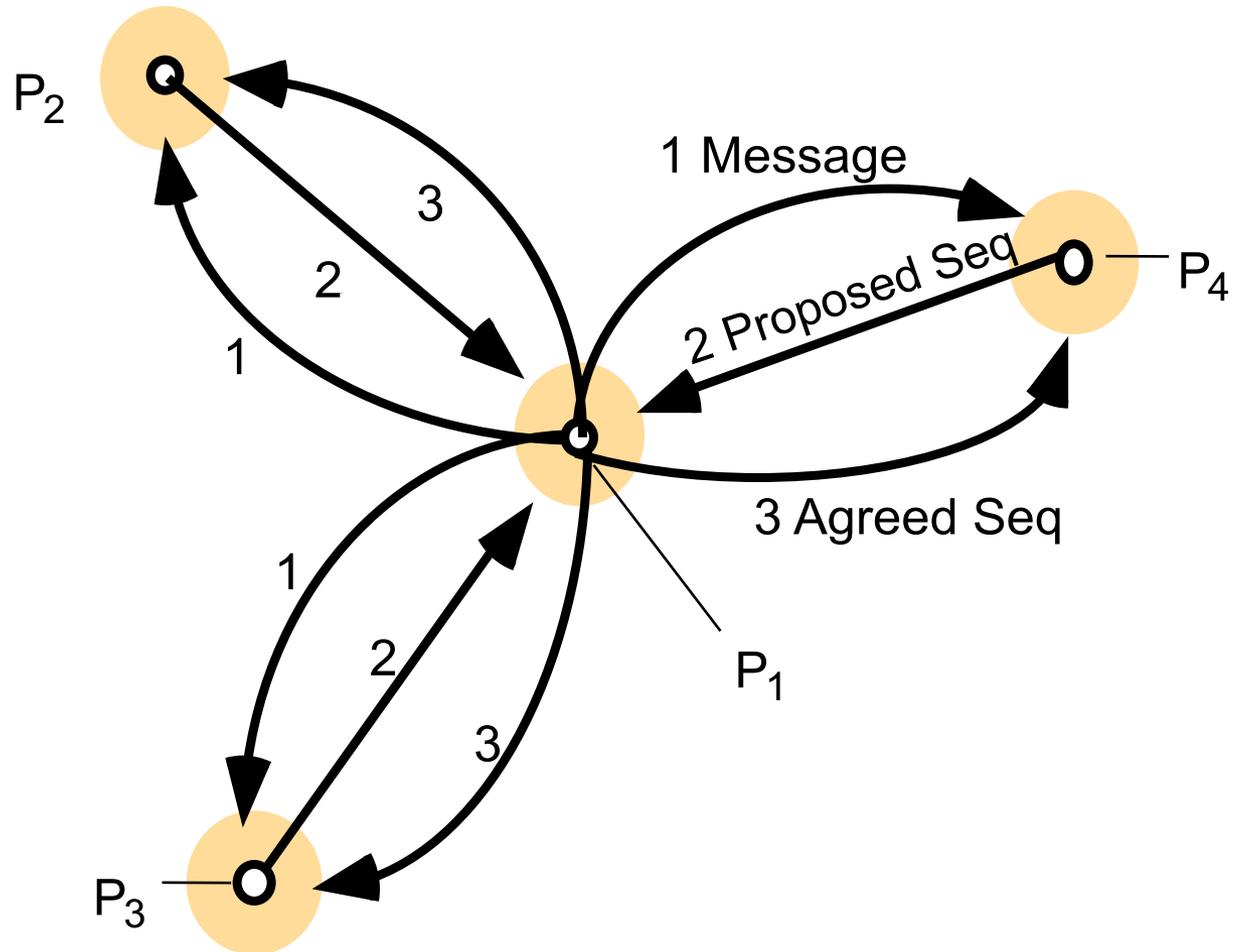
# Sequencer based total ordering

- Special process elected as leader or sequencer.
- TO-multicast( $g, m$ ) at  $P_i$ :
  - Send multicast message  $m$  to group  $g$  and the sequencer
- Sequencer:
  - Maintains a global sequence number  $S$  (initially 0)
  - When a multicast message  $m$  is B-delivered to it:
    - sets  $S = S + 1$ , and B-multicast( $g, \{\text{"order"}, m, S\}$ )
- Receive multicast at process  $P_i$ :
  - $P_i$  maintains a local received global sequence number  $S_i$  (initially 0)
  - On B-deliver( $m$ ) at  $P_i$  from  $P_j$ , it buffers it until both conditions satisfied
    1. B-deliver( $\{\text{"order"}, m, S\}$ ) at  $P_i$  from sequencer, and
    2.  $S_i + 1 = S$
    - Then TO-deliver( $m$ ) to application and set  $S_i = S_i + 1$

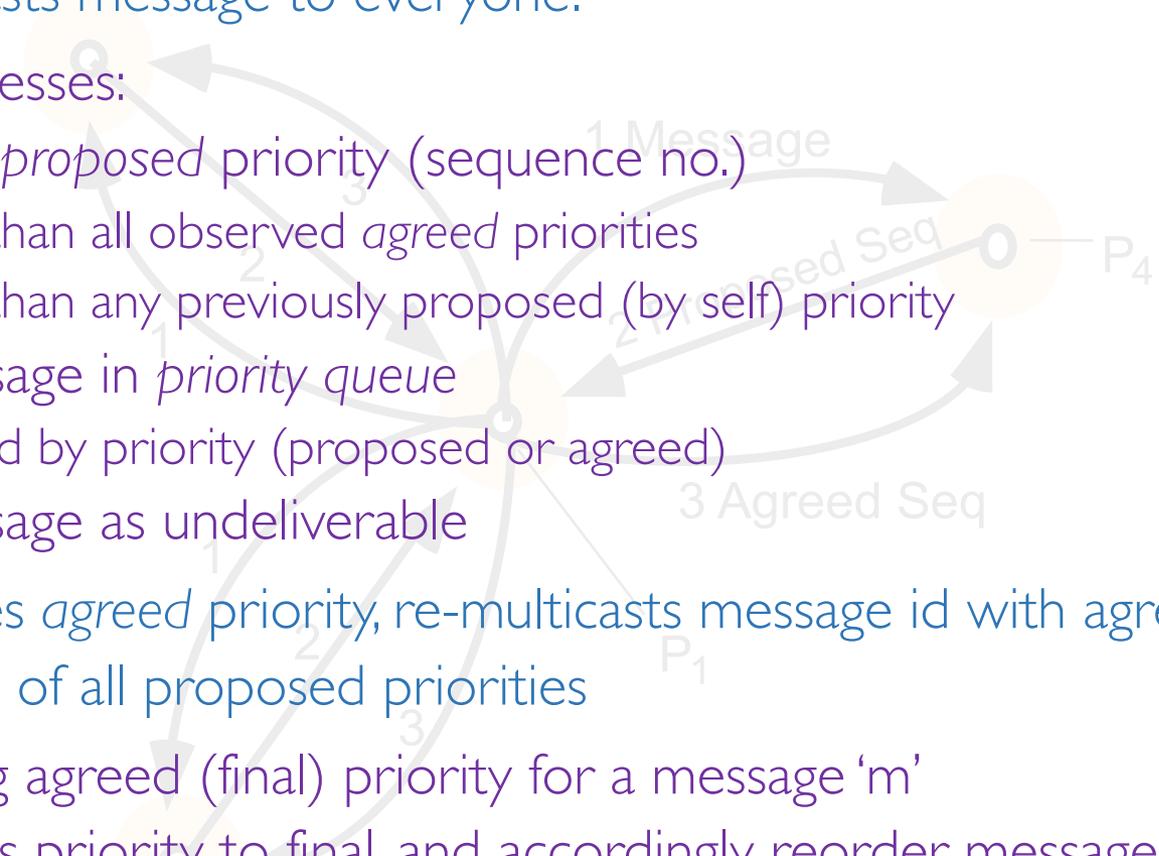
# Implementing total order multicast

- Basic idea:
  - Same sequence number counter across different processes.
  - Instead of different sequence number counter for each process.
- Two types of approach
  - Using a centralized sequencer
  - **A decentralized mechanism (ISIS)**

# ISIS algorithm for total ordering



# ISIS algorithm for total ordering

- Sender multicasts message to everyone.
  - Receiving processes:
    - reply with *proposed* priority (sequence no.)
      - larger than all observed *agreed* priorities
      - larger than any previously proposed (by self) priority
    - store message in *priority queue*
      - ordered by priority (proposed or agreed)
    - mark message as undeliverable
  - Sender chooses *agreed* priority, re-multicasts message id with agreed priority
    - maximum of all proposed priorities
  - Upon receiving agreed (final) priority for a message 'm'
    - Update m's priority to final, and accordingly reorder messages in queue.
    - mark the message m as deliverable.
    - deliver any deliverable messages at front of priority queue.
- 

# ISIS algorithm for total ordering

- We'll discuss an example in the next class.

