

Distributed Systems

CS425/ECE428

Instructor: Radhika Mittal

Logistics Related

- HWI has been released.
 - You can solve first 5 questions right away
 - You can solve last two questions hopefully by end of this week.
- MP0 due on Wednesday.

Today's agenda

- **Global State**

- Chapter 14.5

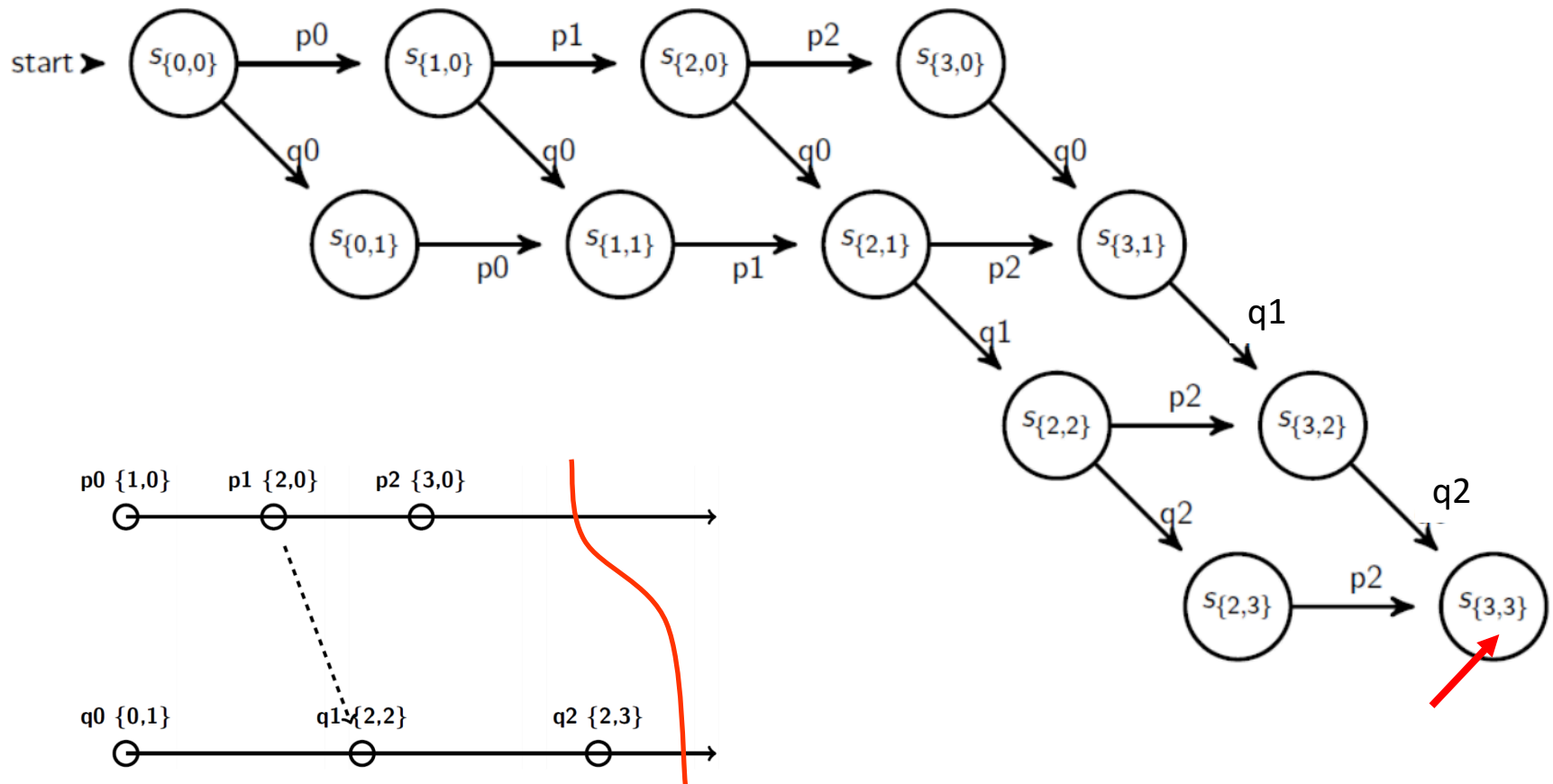
- Goal: reason about how to capture the state across all processes of a distributed system without requiring time synchronization.

- **Multicast**

Recap

- State of each process (and each channel) in the system *at a given instant of time*.
 - Difficult to capture -- requires precisely synchronized time.
- Relax the problem: find a consistent global state.
- Chandy-Lamport algorithm to calculate global state.
 - Obeys causality (creates a consistent cut).
 - Does not interrupt the running distributed application.
 - Can be used to detect global properties.

State Transitions: Example



More notations and definitions

- H = set of all events across all processes.
- A **run** is a total ordering of events in H that is consistent with each h_i 's ordering.
- A **linearization** is a run consistent with happens-before (\rightarrow) relation in H .
- Linearizations pass through consistent global states.
- A global state S_k is reachable from global state S_i , if there is a linearization that passes through S_i and then through S_k .
- The distributed system evolves as a series of transitions between global states S_0, S_1, \dots

Global State Predicates

- A global-state-predicate is a property that is *true* or *false* for a global state.
 - Is there a deadlock?
 - Has the distributed algorithm terminated?
- Two ways of reasoning about predicates (or system properties) as global state gets transformed by events.
 - Liveness
 - Safety

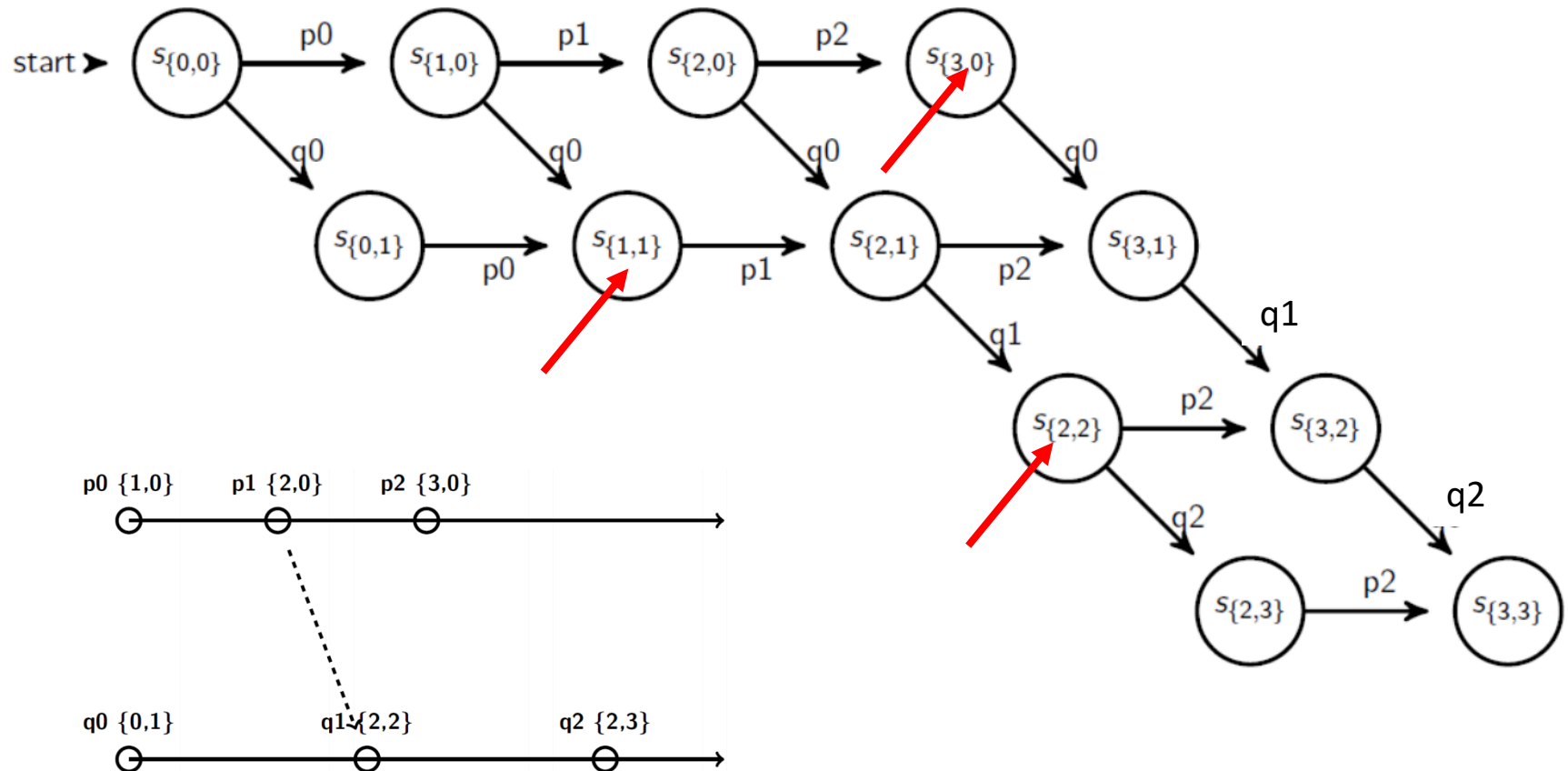
Liveness

- **Liveness** = guarantee that something **good** will happen, **eventually**
- **Examples:**
 - A distributed computation will terminate.
 - “Completeness” in failure detectors: the failure will be detected.
 - All processes will eventually decide on a value.
- A global state S_0 satisfies a **liveness** property P iff:
 - For all linearizations starting from S_0 , P is true for **some** state S_L reachable from S_0 .
 - $\text{liveness}(P(S_0)) \equiv \forall L \in \text{linearizations from } S_0, L \text{ passes through a } S_L \ \& \ P(S_L) = \text{true}$

Liveness Example

If predicate is true only in the marked states, does it satisfy liveness?

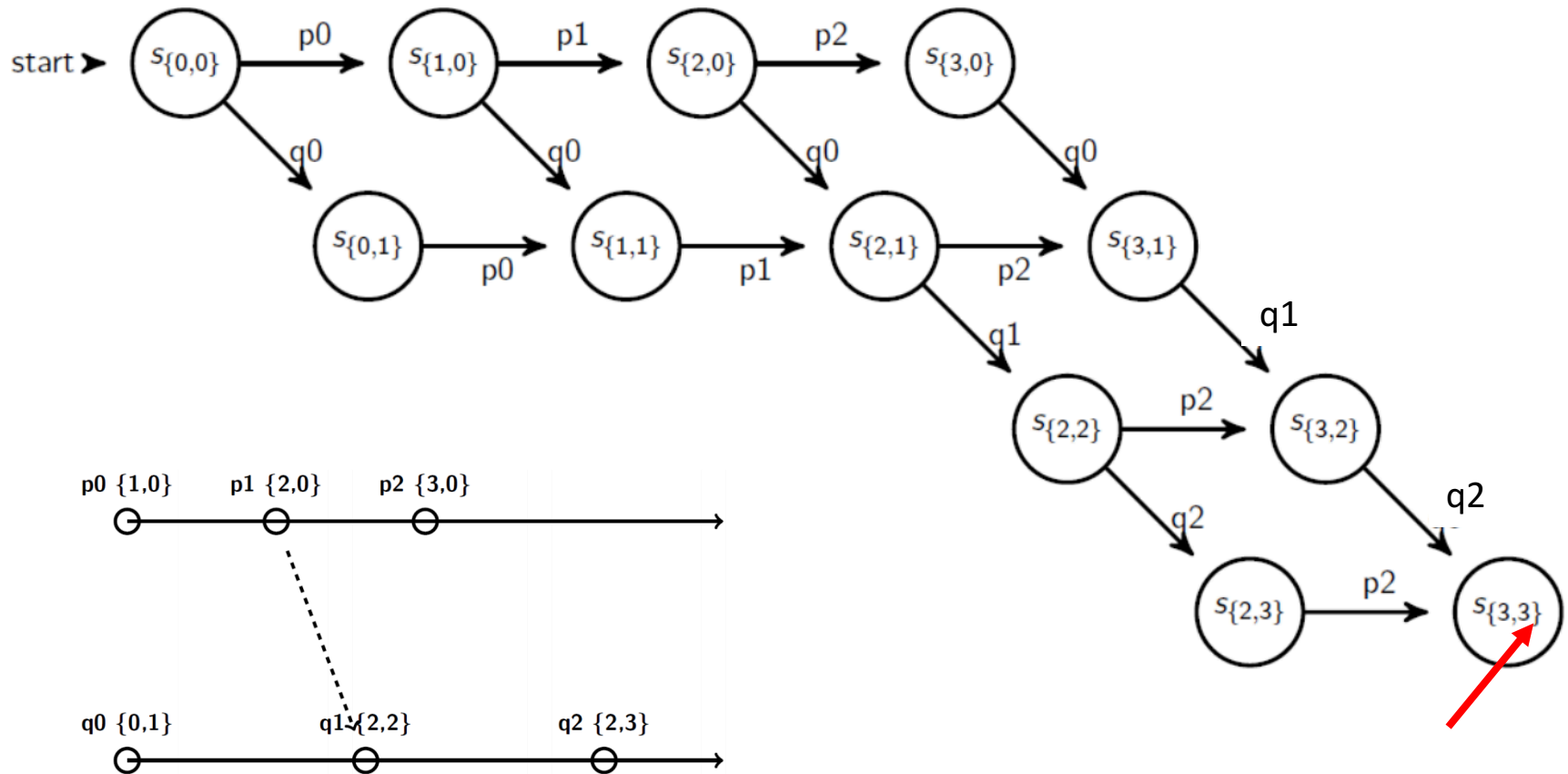
No



Liveness Example

If predicate is true only in the marked states, does it satisfy liveness?

Yes



Liveness

- **Liveness** = guarantee that something **good** will happen, **eventually**
- **Examples:**
 - A distributed computation will terminate.
 - “Completeness” in failure detectors: the failure will be detected.
 - All processes will eventually decide on a value.
- A global state S_0 satisfies a **liveness** property P iff:
 - $\text{liveness}(P(S_0)) \equiv \forall L \in \text{linearizations from } S_0, L \text{ passes through a } S_L \text{ \& } P(S_L) = \text{true}$
 - For any linearization starting from S_0 , P is true for **some** state S_L reachable from S_0 .

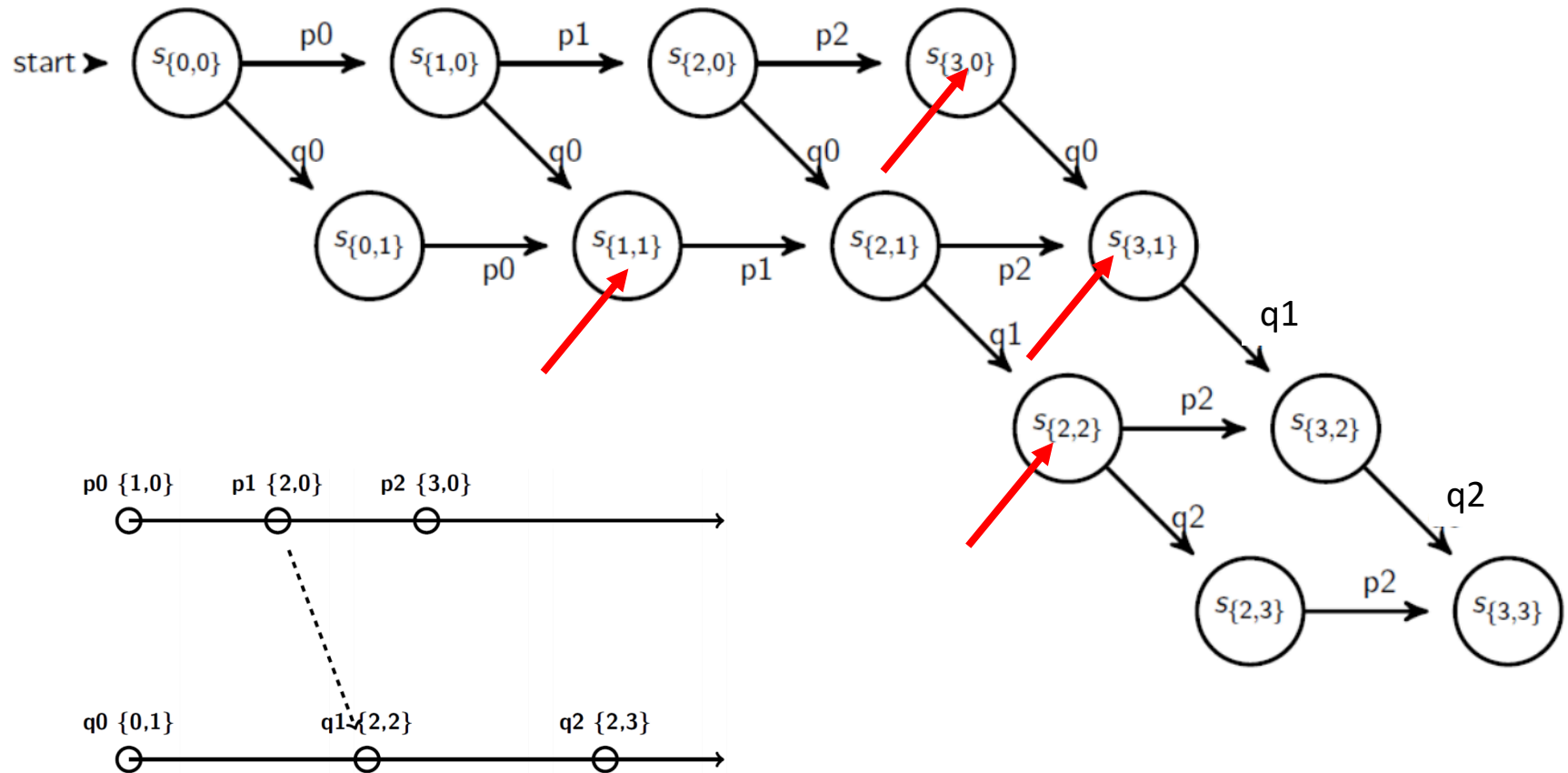
Safety

- **Safety** = guarantee that something **bad** will **never** happen.
- **Examples:**
 - There is no deadlock in a distributed transaction system.
 - “Accuracy” in failure detectors: an alive process is not detected as failed.
 - No two processes decide on different values.
- A global state S_0 satisfies a **safety** property P iff:
 - For **all** states S reachable from S_0 , $P(S)$ is true.
 - $\text{safety}(P(S_0)) \equiv \forall S \text{ reachable from } S_0, P(S) = \text{true}.$

Safety Example

If predicate is true only in the marked states, does it satisfy safety?

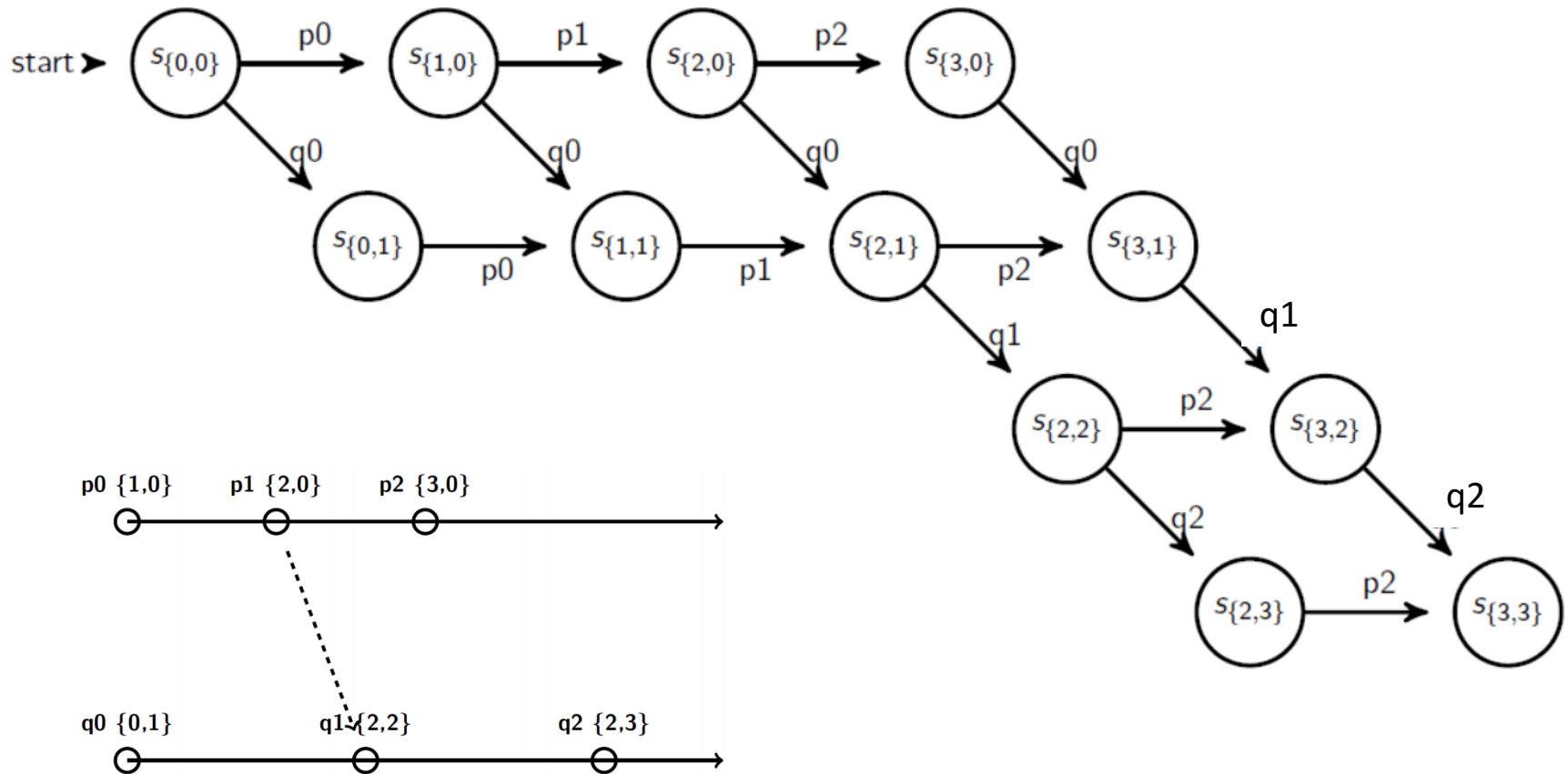
No



Safety Example

If predicate is true only in the **unmarked** states, does it satisfy safety?

Yes

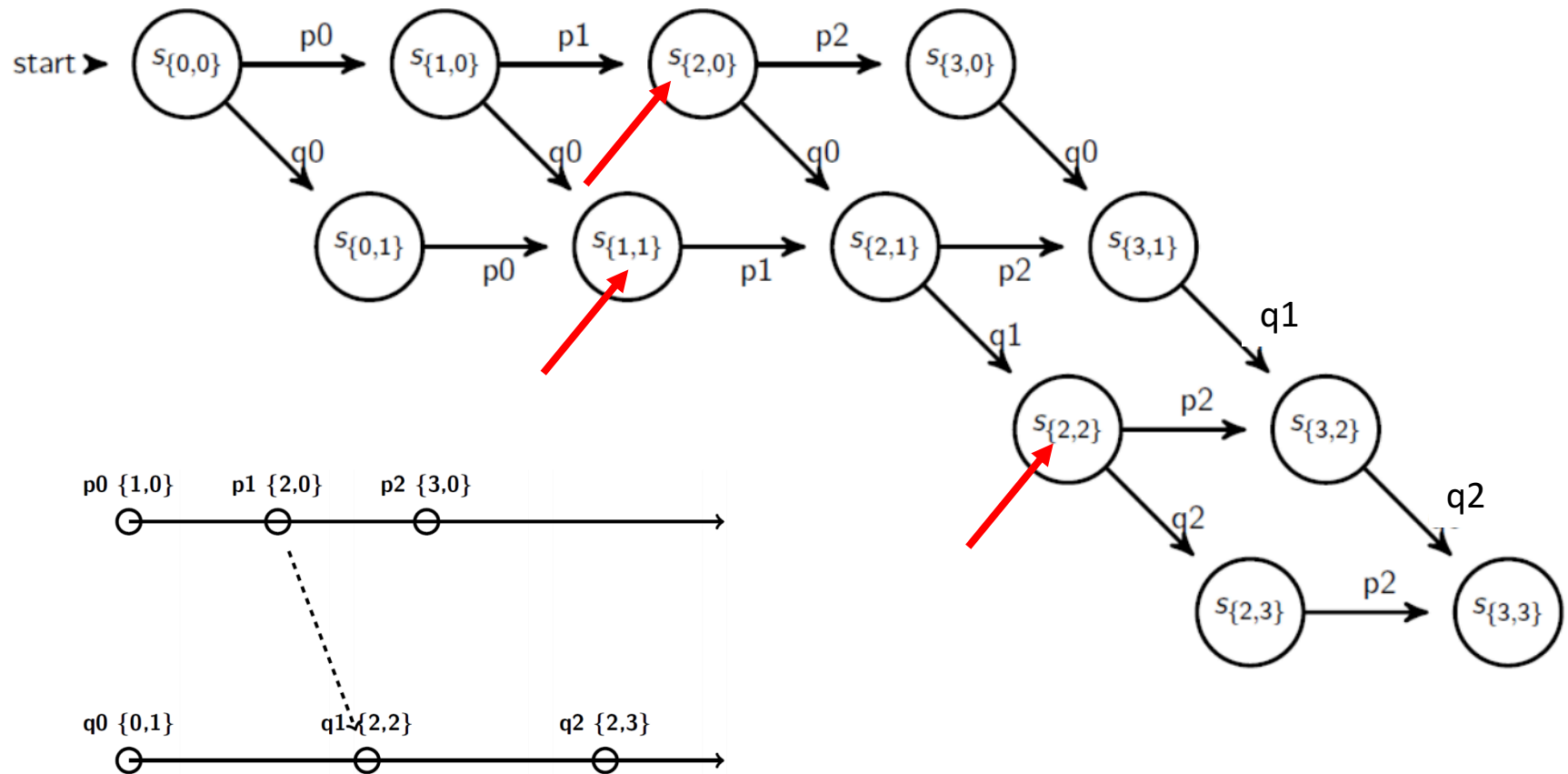


Safety

- **Safety** = guarantee that something **bad** will **never** happen.
- **Examples:**
 - There is no deadlock in a distributed transaction system.
 - “Accuracy” in failure detectors: an alive process is not detected as failed.
 - No two processes decide on different values.
- A global state S_0 satisfies a **safety** property P iff:
 - $\text{safety}(P(S_0)) \equiv \forall S \text{ reachable from } S_0, P(S) = \text{true}.$
 - For **all** states S reachable from S_0 , $P(S)$ is true.

Liveness Example

Technically satisfies liveness, but difficult to capture or reason about.



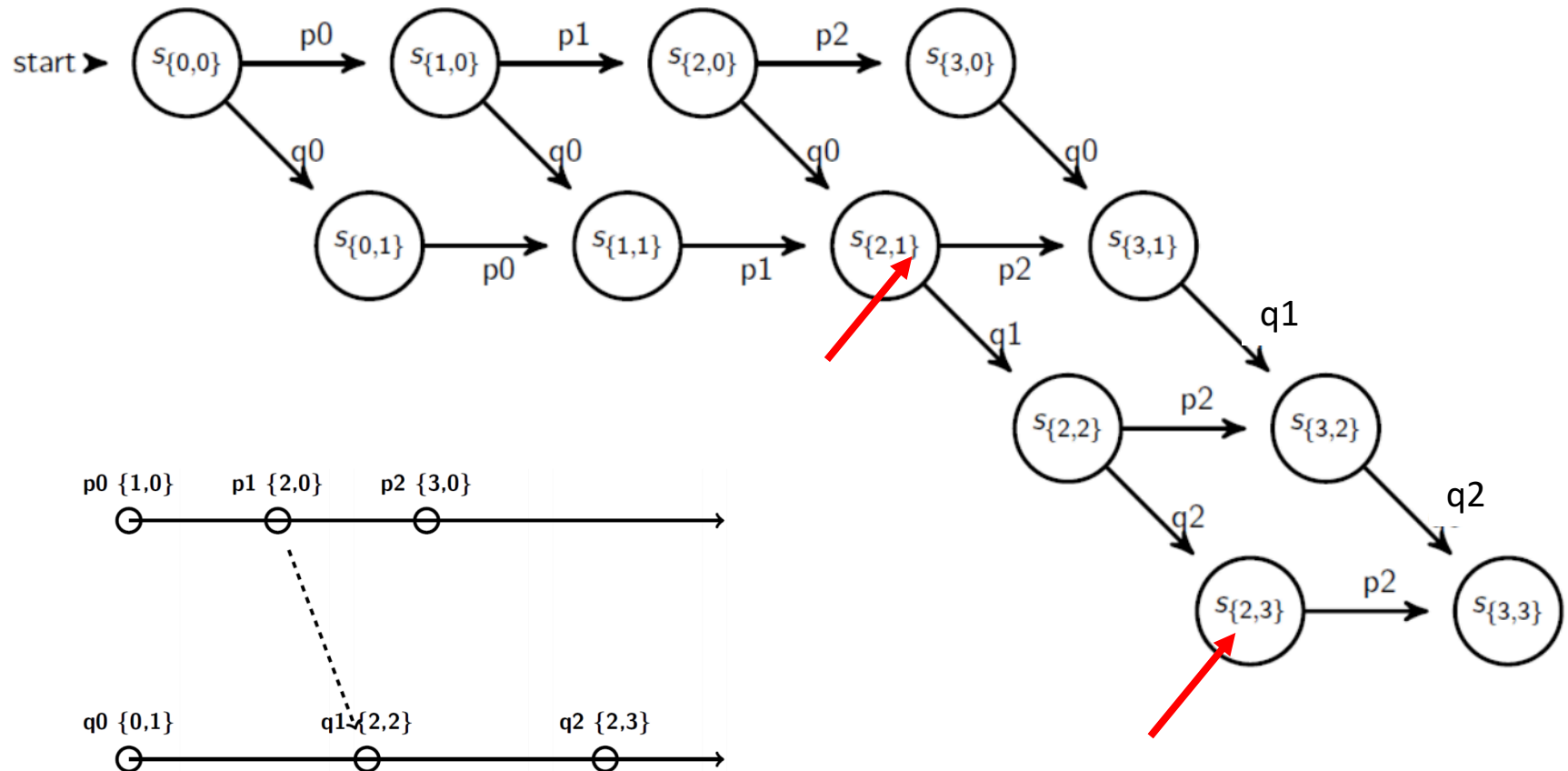
Stable Global Predicates

- once true, stays true forever afterwards (for stable liveness)

Stable Global Predicates

If predicate is true only in the marked states, is it stable?

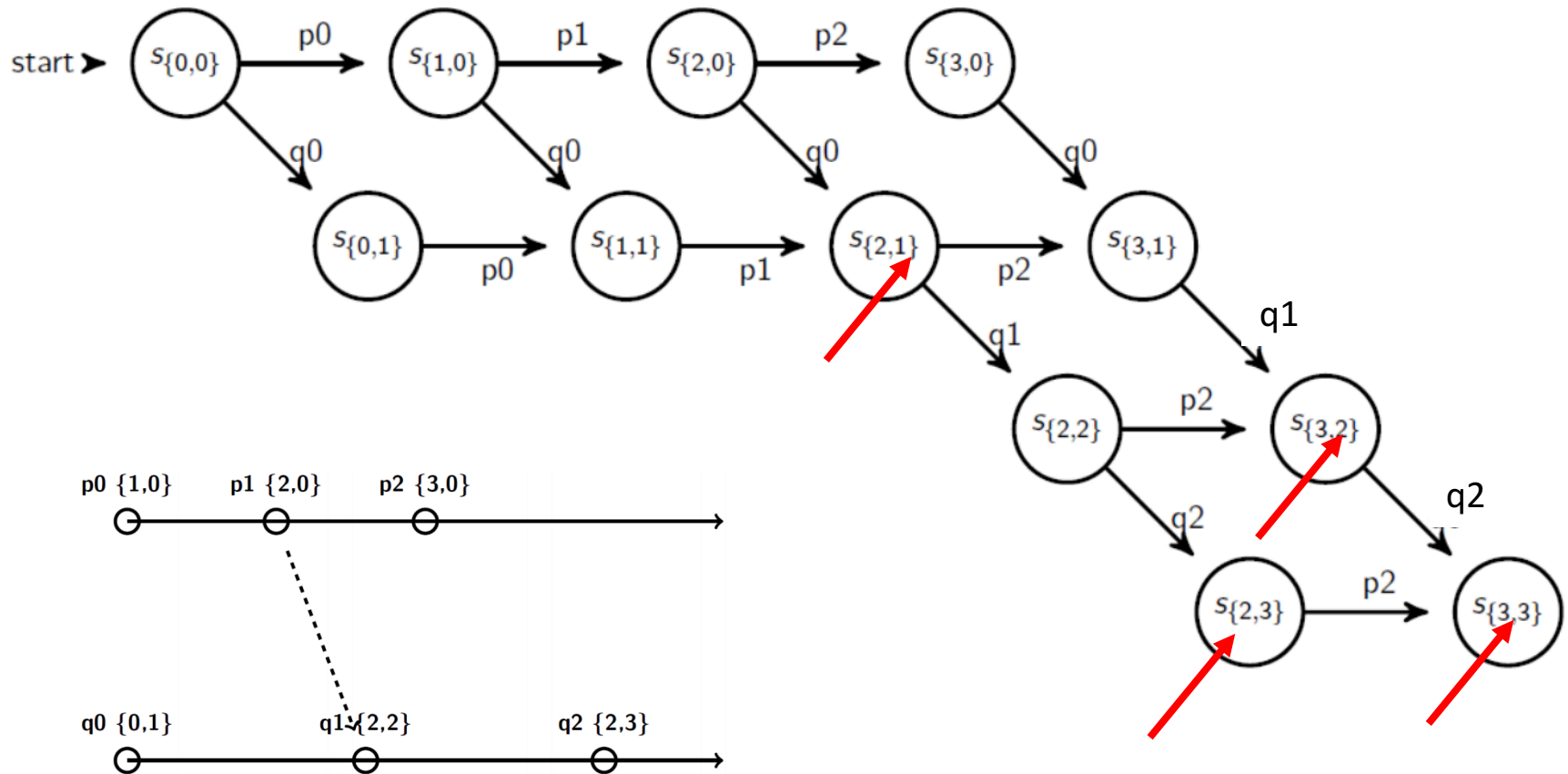
No



Stable Global Predicates

If predicate is true only in the marked states, is it stable?

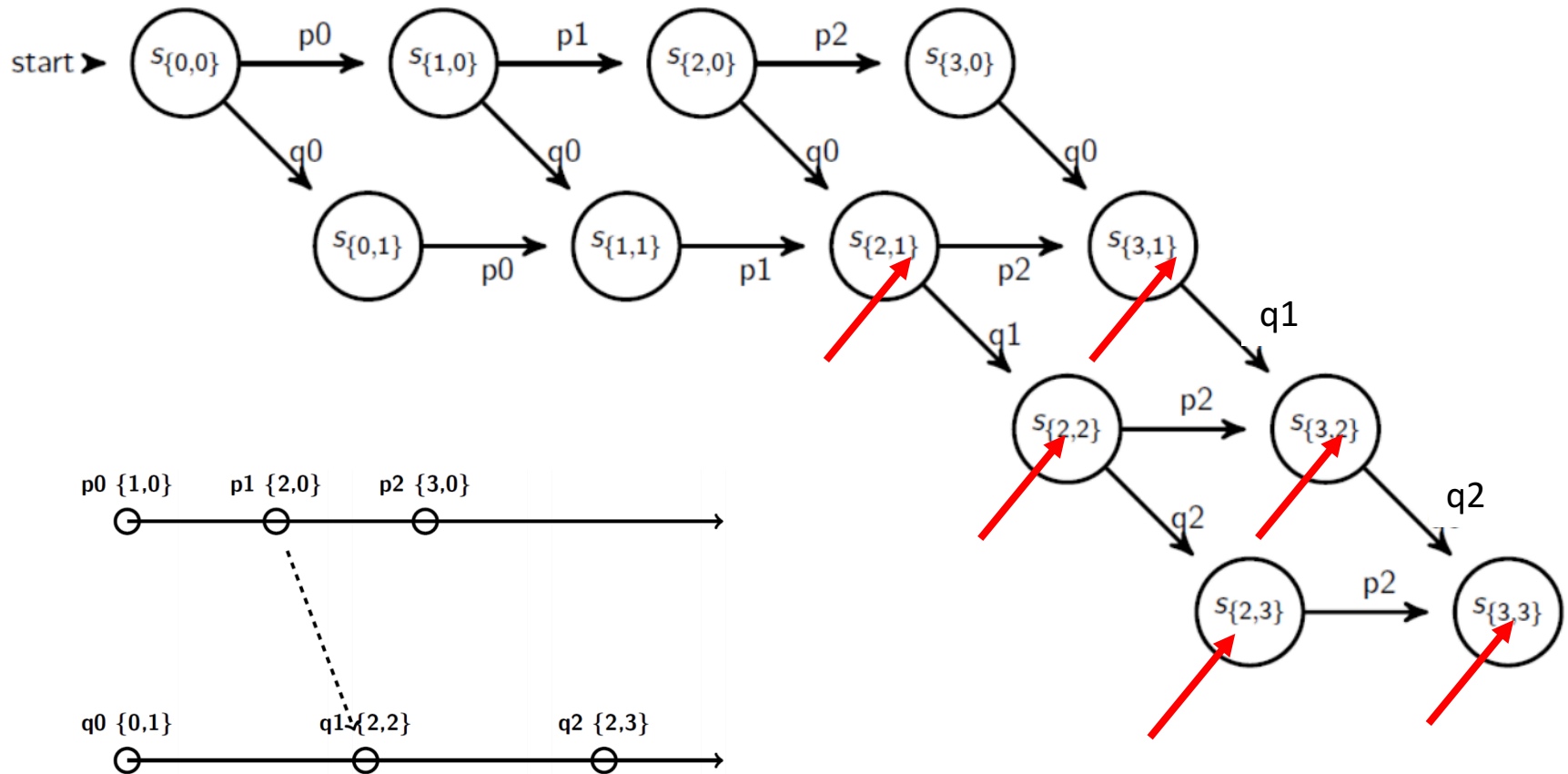
No



Stable Global Predicates

If predicate is true only in the marked states, is it stable?

Yes



Stable Global Predicates

- once true for a state S , stays true for all states reachable from S (for stable liveness)
- once false for a state S , stays false for all states reachable from S (for stable non-safety)
- Stable liveness examples (once true, always true)
 - Computation has terminated.
- Stable non-safety examples (once false, always false)
 - There is no deadlock.
 - An object is not orphaned.
- *All stable global properties can be detected using the Chandy-Lamport algorithm.*

Global Snapshot Summary

- The ability to calculate global snapshots in a distributed system is very important.
- But don't want to interrupt running distributed application.
- Chandy-Lamport algorithm calculates global snapshot.
- Obeys causality (creates a consistent cut).
- Can be used to detect global properties.
- Safety vs. Liveness.

Rest of today's agenda

- **Multicast**
 - Chapter 15.4
- **Goal:** reason about desirable properties for message delivery among a group of processes.

Communication modes

- **Unicast**
 - Messages are sent from exactly one process to one process.
- **Broadcast**
 - Messages are sent from exactly one process to all processes on the network.
- **Multicast**
 - Messages broadcast within a group of processes.
 - A multicast message is sent from any one process to a group of processes on the network.

Where is multicast used?

- Distributed storage
 - Write to an object are multicast across replica 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.
-

Communication modes

- **Unicast**

- Messages are sent from exactly one process to one process.
 - *Best effort*: if a message is delivered it would be intact; no reliability guarantees.
 - *Reliable*: guarantees delivery of messages.
 - *In order*: messages will be delivered in the same order that they are sent.

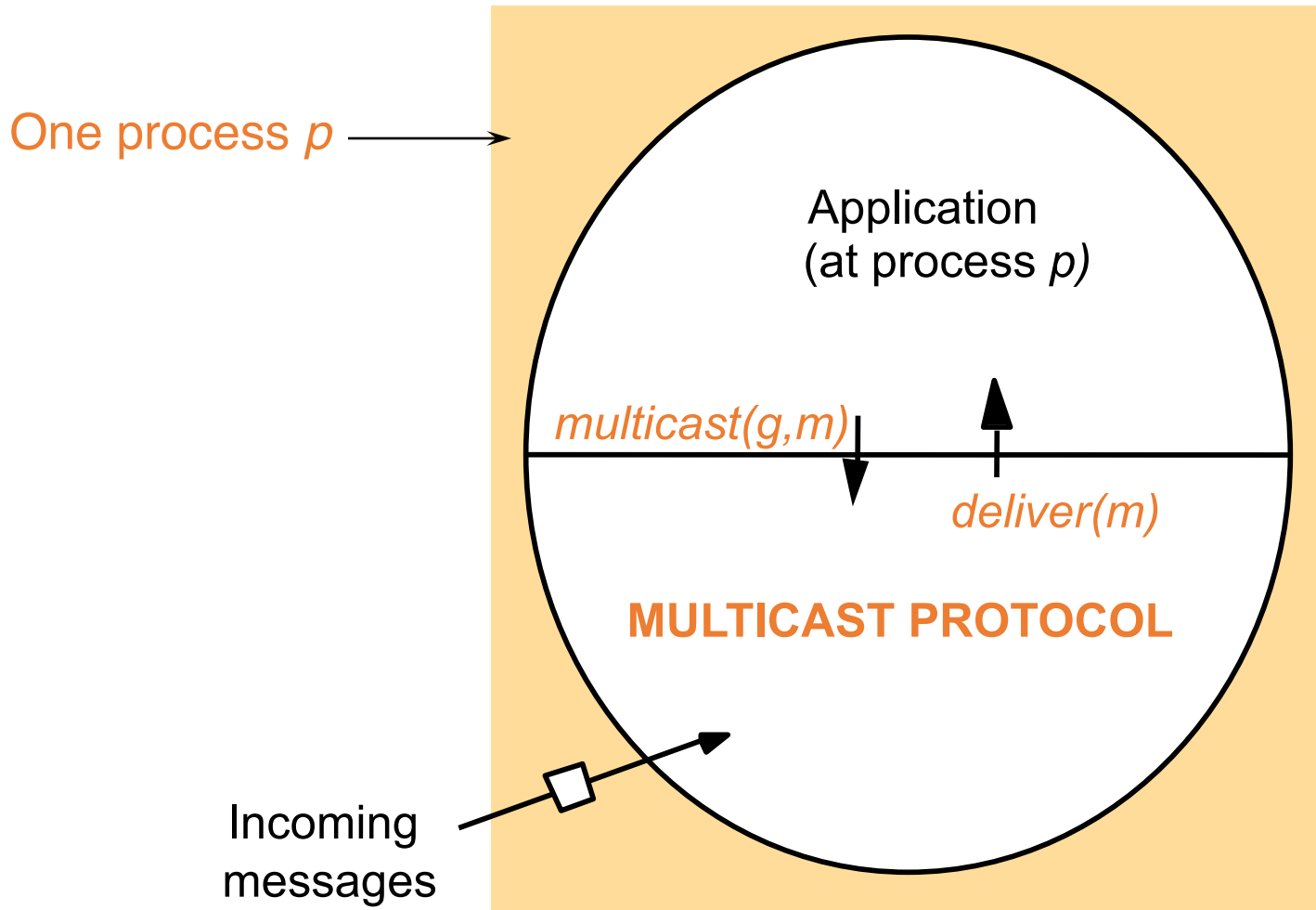
- **Broadcast**

- Messages are sent from exactly one process to all processes on the network.

- **Multicast**

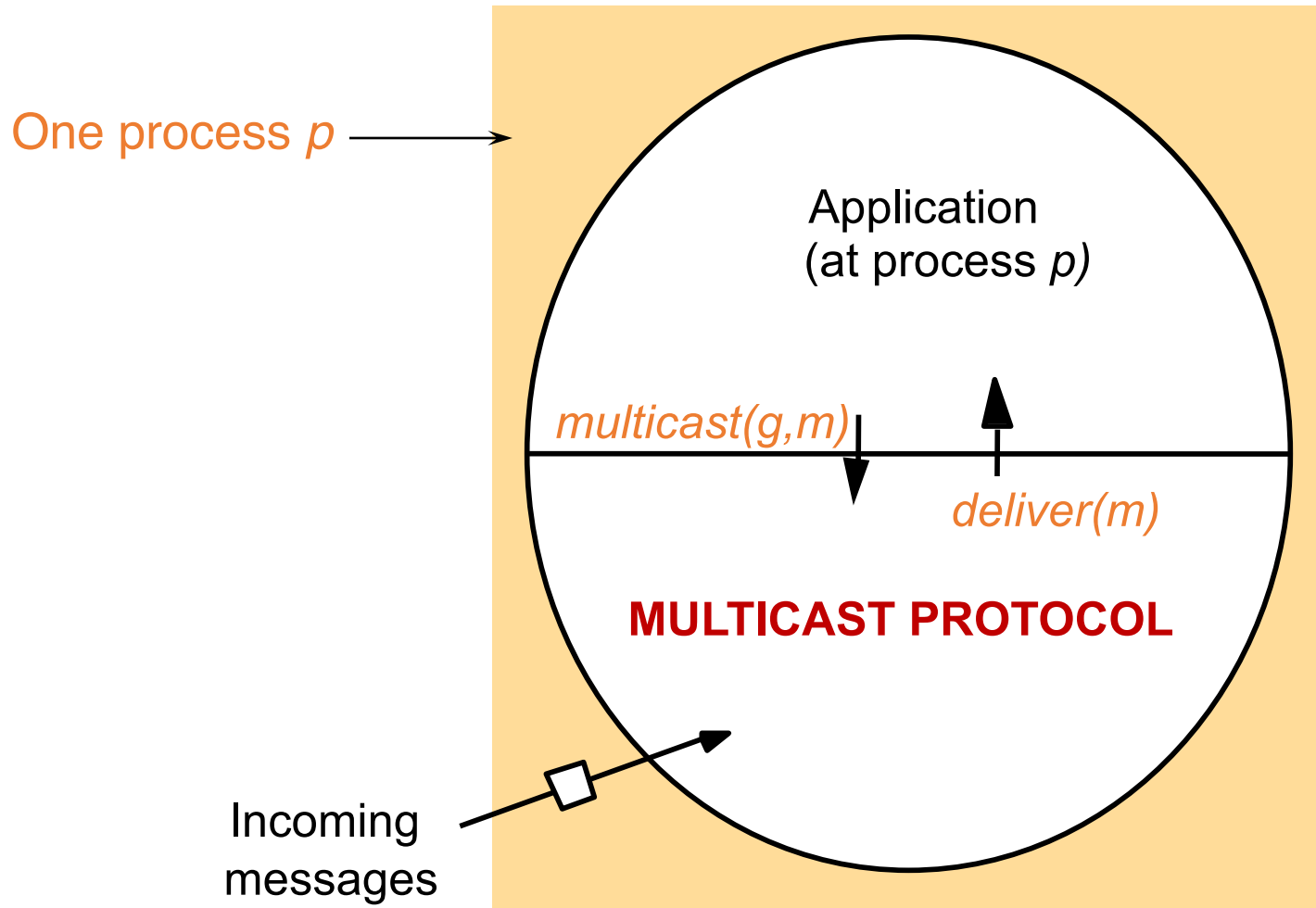
- Messages broadcast within a group of processes.
- A multicast message is sent from any one process to the group of processes on the network.
- *How do we define (and achieve) reliable or ordered multicast?*

What we are designing in this class?



' g ' is a multicast group that also includes the process ' p '.

What we are designing in this class?



' g ' is a multicast group that also includes the process ' p '.

Basic Multicast (B-Multicast)

- Straightforward way to implement B-multicast:
 - use a reliable one-to-one send (unicast) operation:
B-multicast(group g , message m):
for each process p in g , send (p,m).
receive(m): B-deliver(m) at p .
- Guarantees: message is eventually delivered to the group if:
 - Processes are non-faulty.
 - The unicast “send” is reliable.
 - *Sender does not crash.*
- *Can we provide reliable delivery even after sender crashes?*
 - *What does this mean?*

Reliable Multicast (R-Multicast)

- **Integrity:** A *correct* (i.e., non-faulty) process p delivers a message m at most once.
 - *Assumption: no process sends **exactly** the same message twice*
- **Validity:** If a *correct* process multicasts (sends) message m , then it will *eventually* deliver m to itself.
 - *Liveness for the sender.*
- **Agreement:** If a *correct* process delivers message m , then all the other *correct* processes in $\text{group}(m)$ will *eventually* deliver m .
 - *All or nothing.*
- Validity and agreement together ensure overall liveness: if some correct process multicasts a message m , then, all correct processes deliver m too.

Reliable Multicast (R-Multicast)

- **Integrity:** A *correct* (i.e., non-faulty) process p delivers a message m at most once.

- Assur

- **Validity:** If a process multicasts a message m , then it will eventually deliver m .

- Liveness

- **Agreement:** If a process multicasts a message m , then all correct processes will eventually deliver m .

- All or

What happens if a process initiates B-multicasts of a message but fails after unicasting to a subset of processes in the group?

Agreement is violated! R-multicast not satisfied.

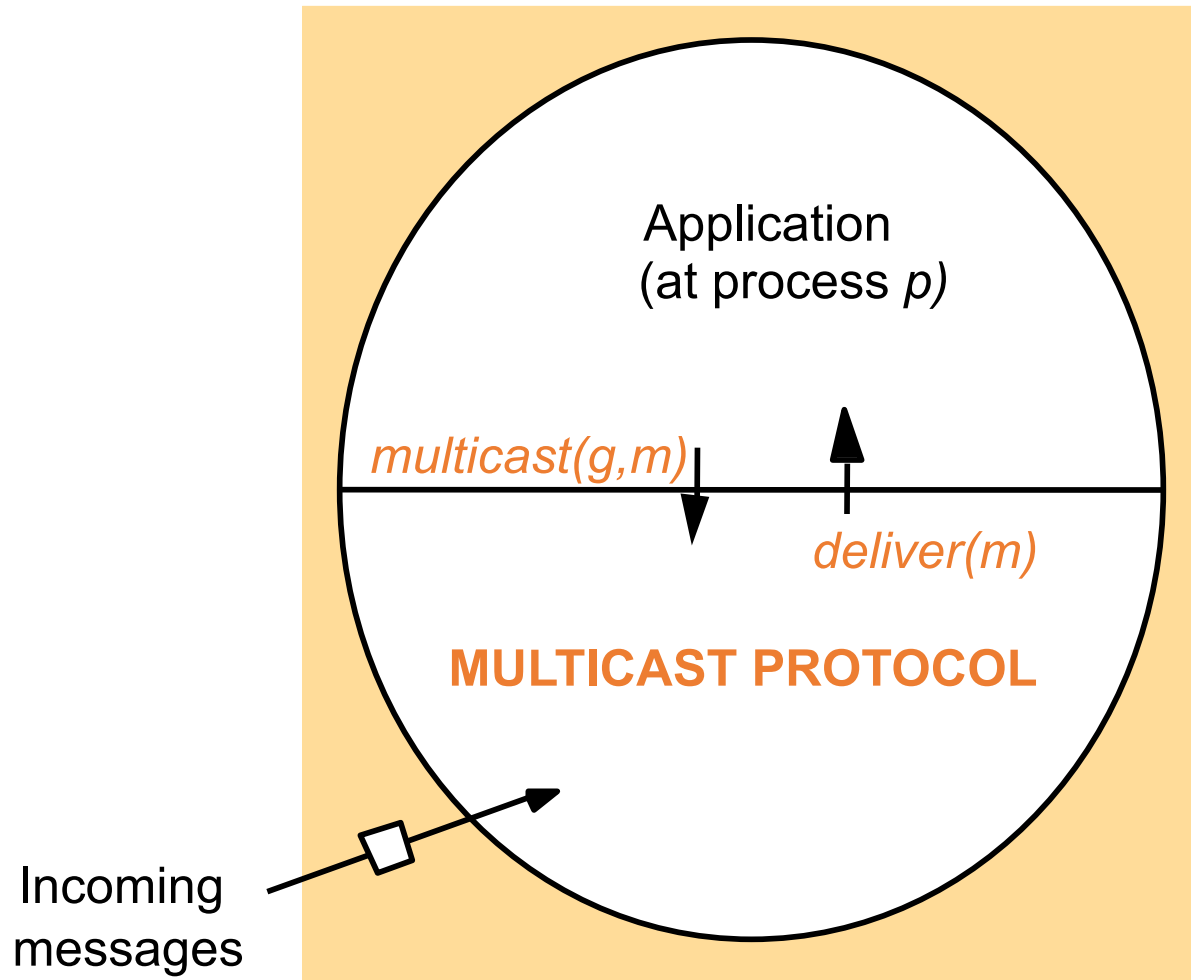
- Validity and agreement together ensure overall liveness: if some correct process multicasts a message m , then, all correct processes deliver m too.

twice

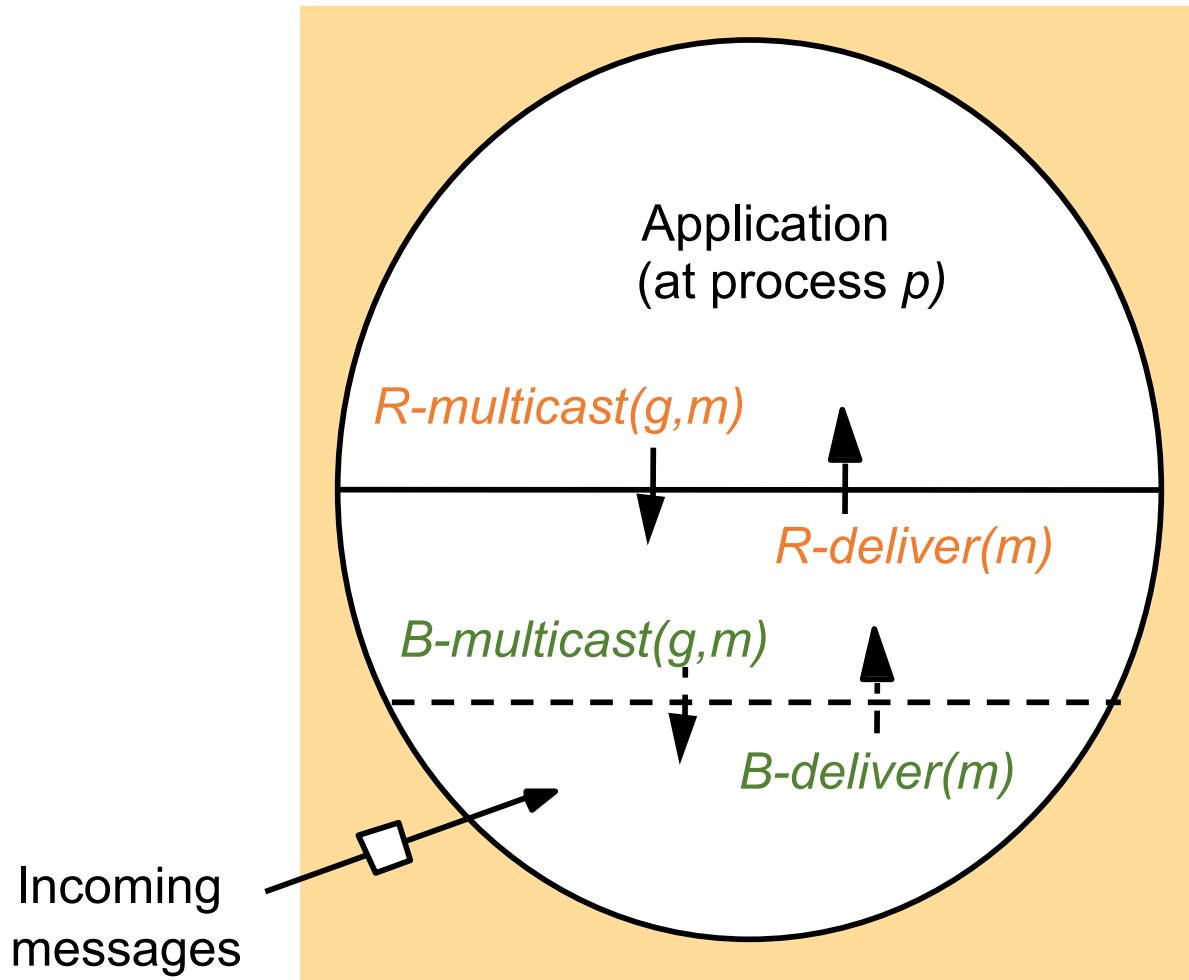
then it will

the other

Implementing R-Multicast



Implementing R-Multicast



Implementing R-Multicast

On initialization

Received := {};

For process p to R-multicast message m to group g

B-multicast(g, m); ($p \in g$ is included as destination)

On B-deliver(m) at process q in $g = \text{group}(m)$

if ($m \notin \text{Received}$):

Received := Received \cup { m };

if ($q \neq p$): B-multicast(g, m);

R-deliver(m)

Reliable Multicast (R-Multicast)

- **Integrity:** A *correct* (i.e., non-faulty) process p delivers a message m at most once.
 - *Assumption: no process sends **exactly** the same message twice*
- **Validity:** If a *correct* process multicasts (sends) message m , then it will *eventually* deliver m to itself.
 - *Liveness for the sender.*
- **Agreement:** If a *correct* process delivers message m , then all the other *correct* processes in $\text{group}(m)$ will *eventually* deliver m .
 - *All or nothing.*
- Validity and agreement together ensure overall liveness: if some correct process multicasts a message m , then, all correct processes deliver m too.

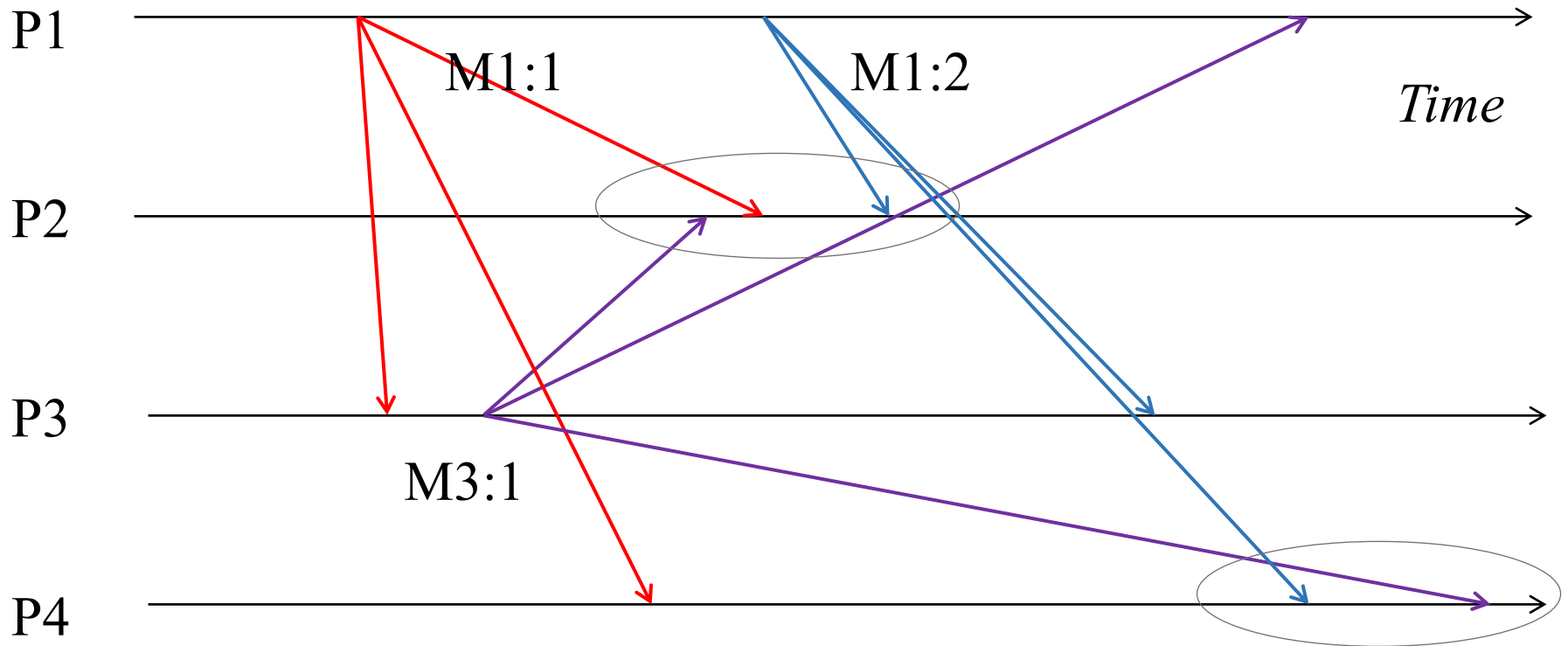
Ordered Multicast

- 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 .*

FIFO Order: Example



M1:1 and M1:2 should be delivered in that order at each receiver.
Order of delivery of M3:1 and M1:2 could be different at different receivers.

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.

Where is causal ordering useful?

- 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.

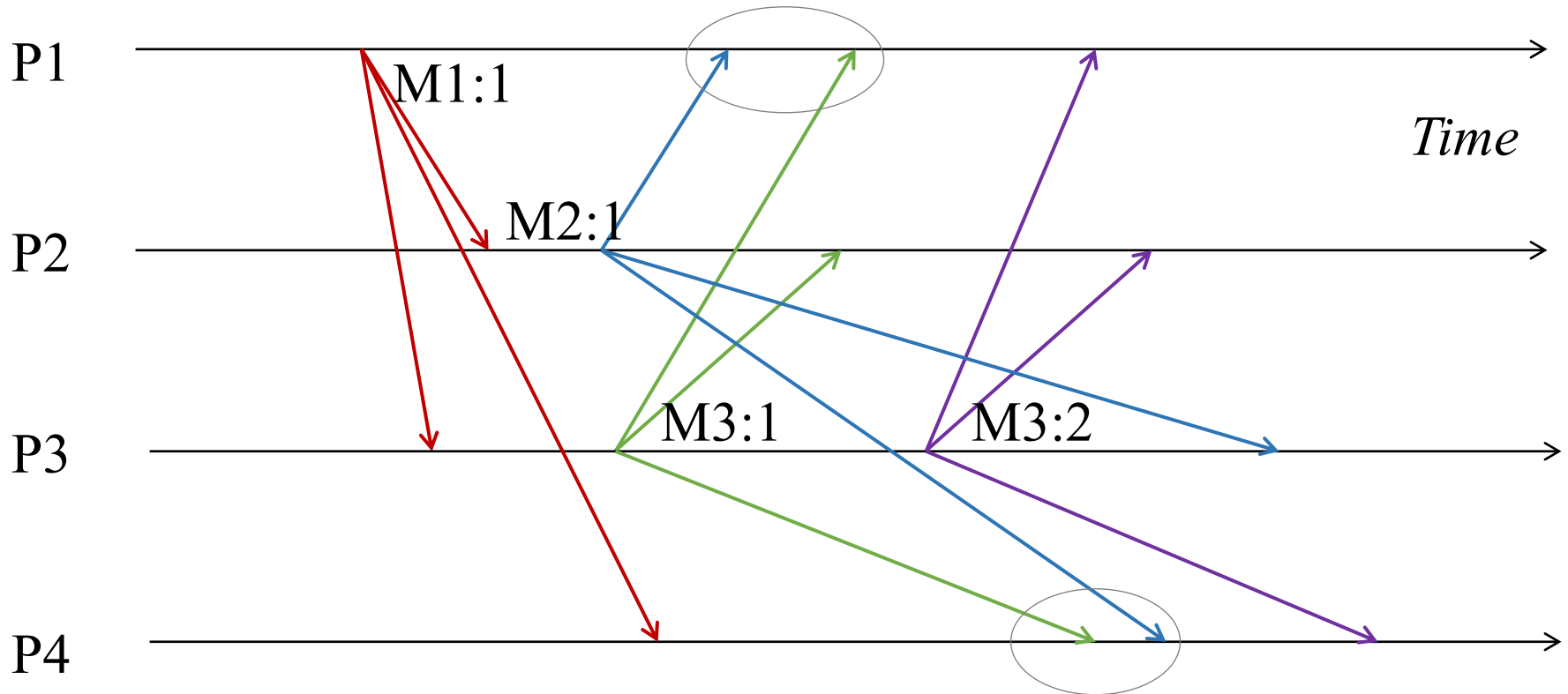
HB Relationship for Causal Ordering

- 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 , **send(m) \rightarrow receive(m)**

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



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

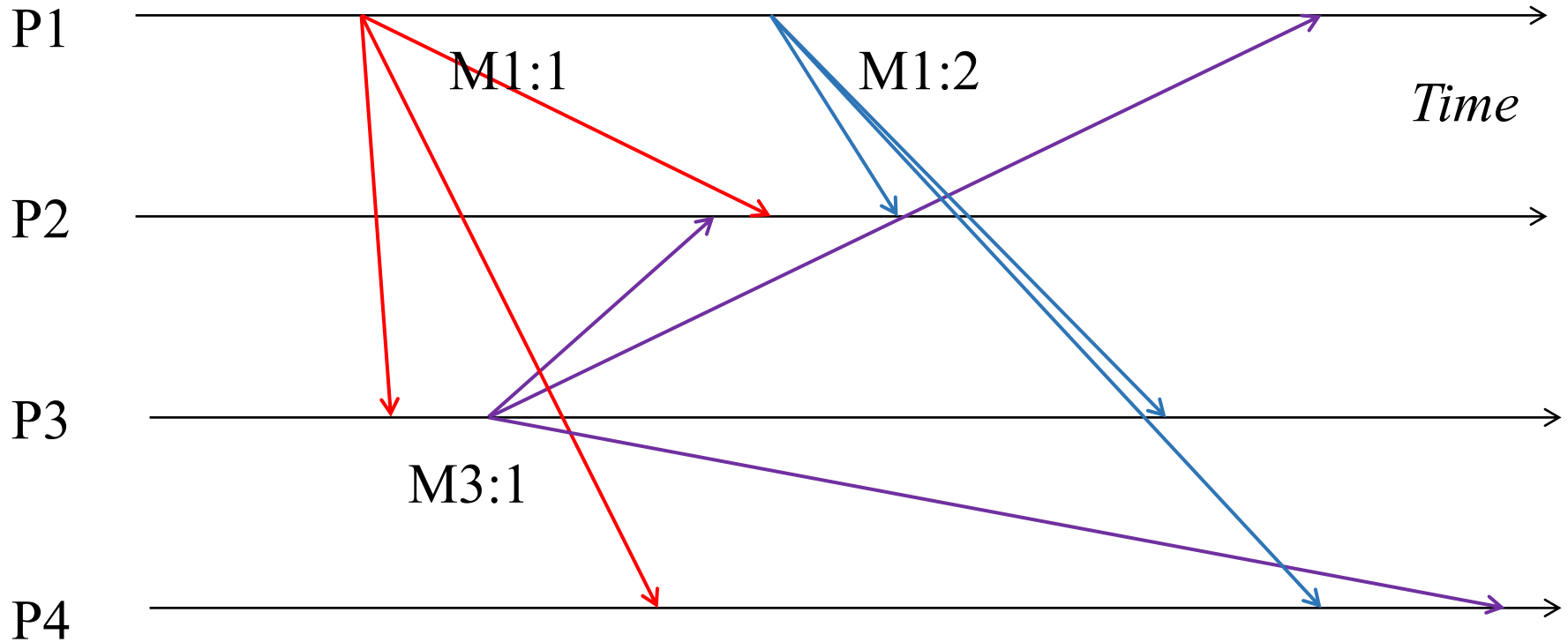
$M3:1$ and $M2:1$ are concurrent and thus ok to be delivered in any (and even different) orders at different receivers.

Causal vs FIFO

- Does Causal Ordering imply FIFO Ordering?
 - Yes

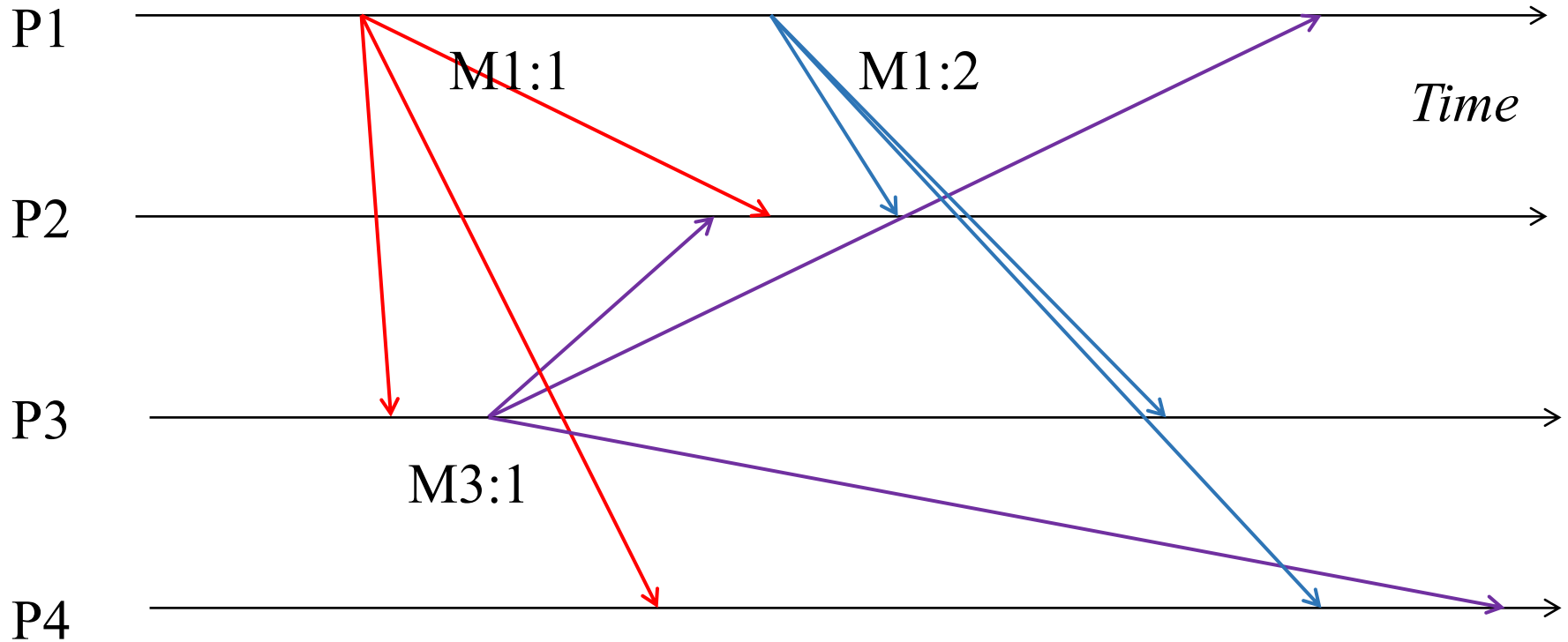
- Does FIFO Order imply Causal Order?
 - No

Example



Does this satisfy FIFO order?

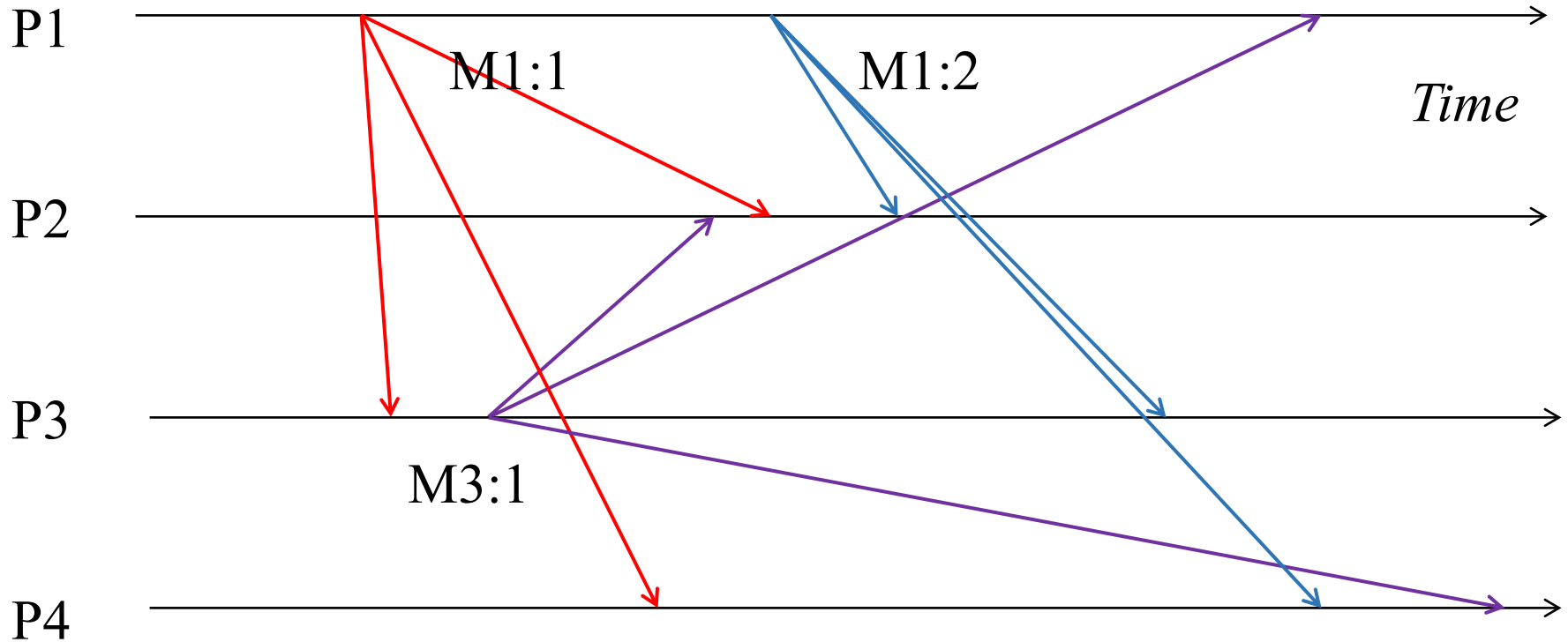
Example



Does this satisfy FIFO order?

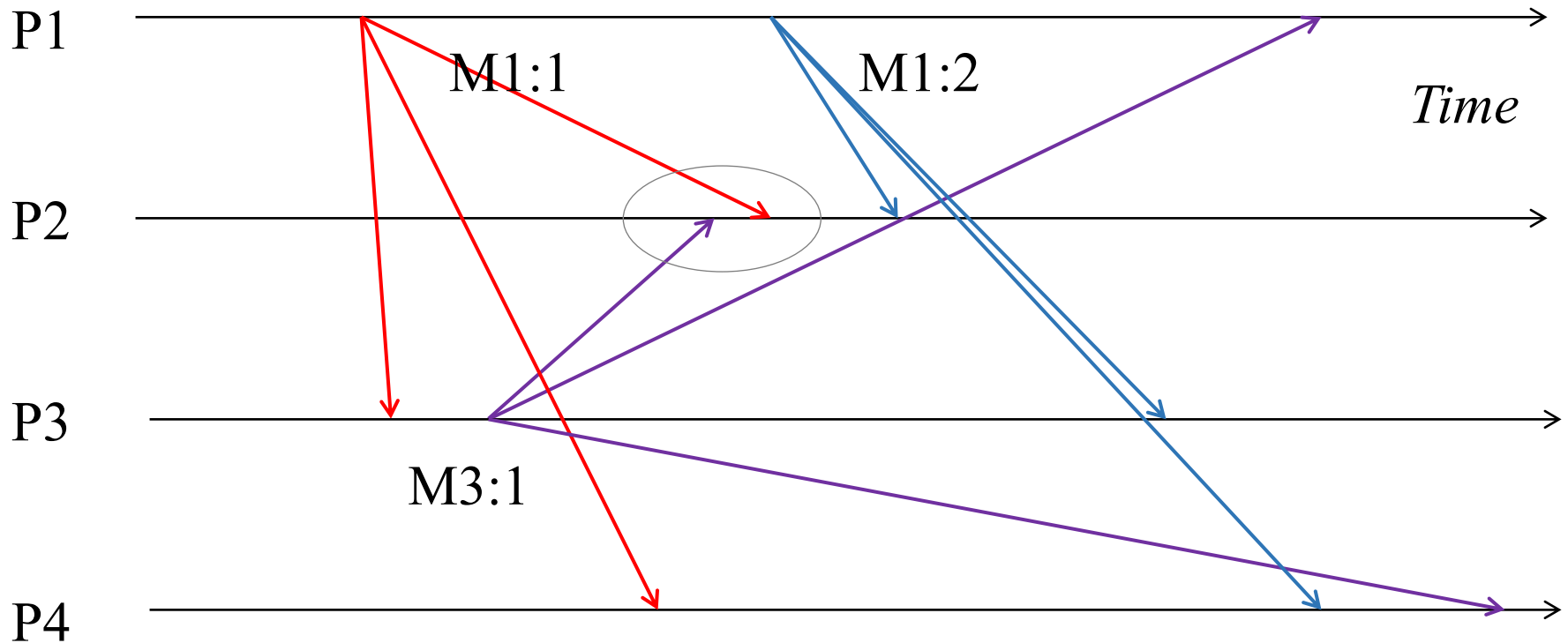
Yes

Example



Does this satisfy causal order?

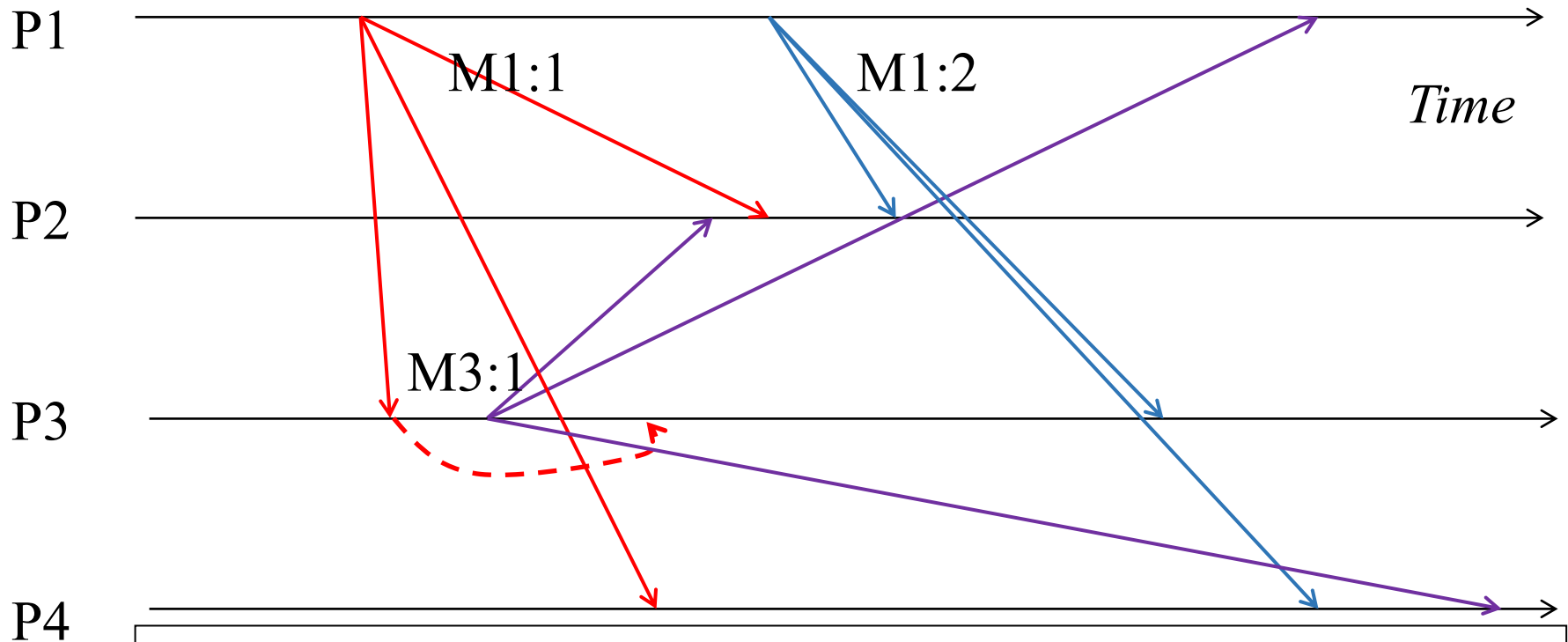
Example



Does this satisfy causal order?

No

Example

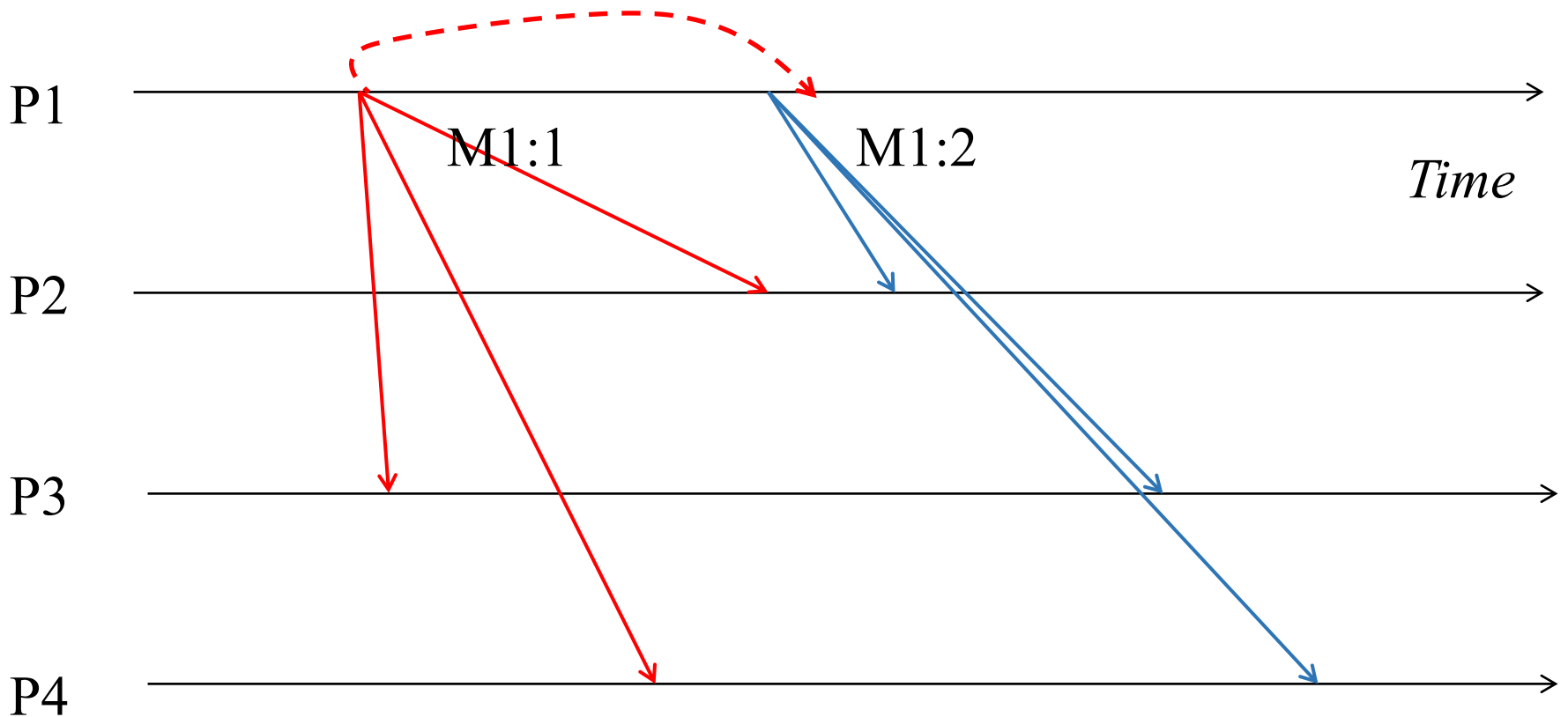


M1:1 is delivered at P3 after M3:1's multicast.

Does this satisfy causal order?

Yes

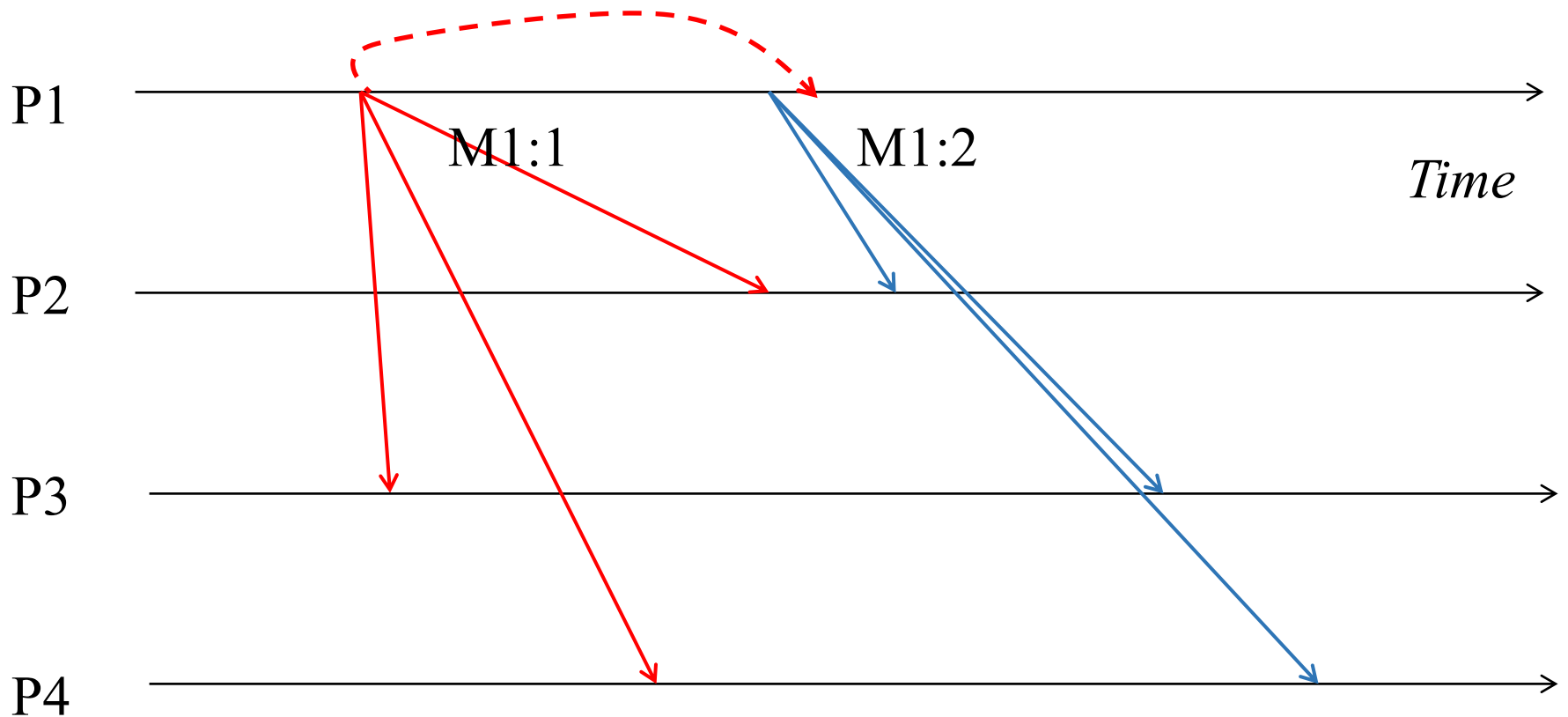
Example



Does this satisfy causal order?

No

Example



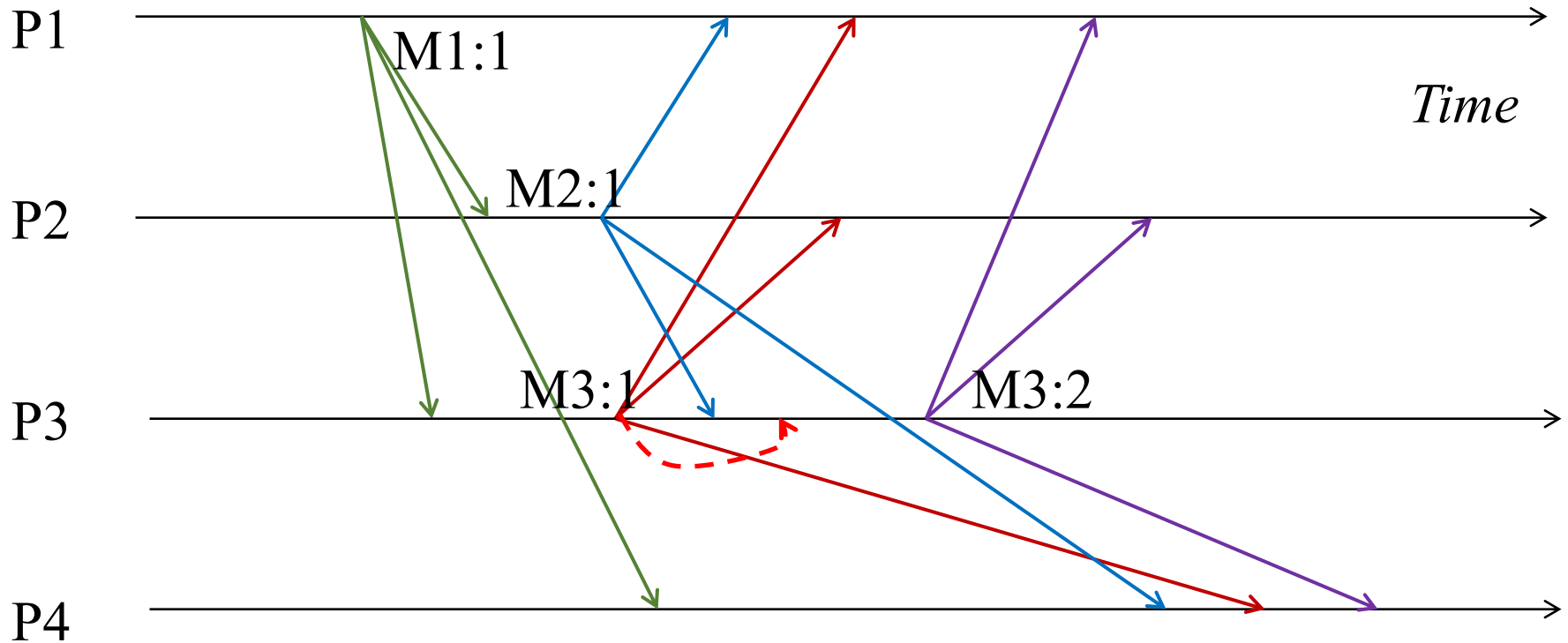
Does this satisfy FIFO order?

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



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.

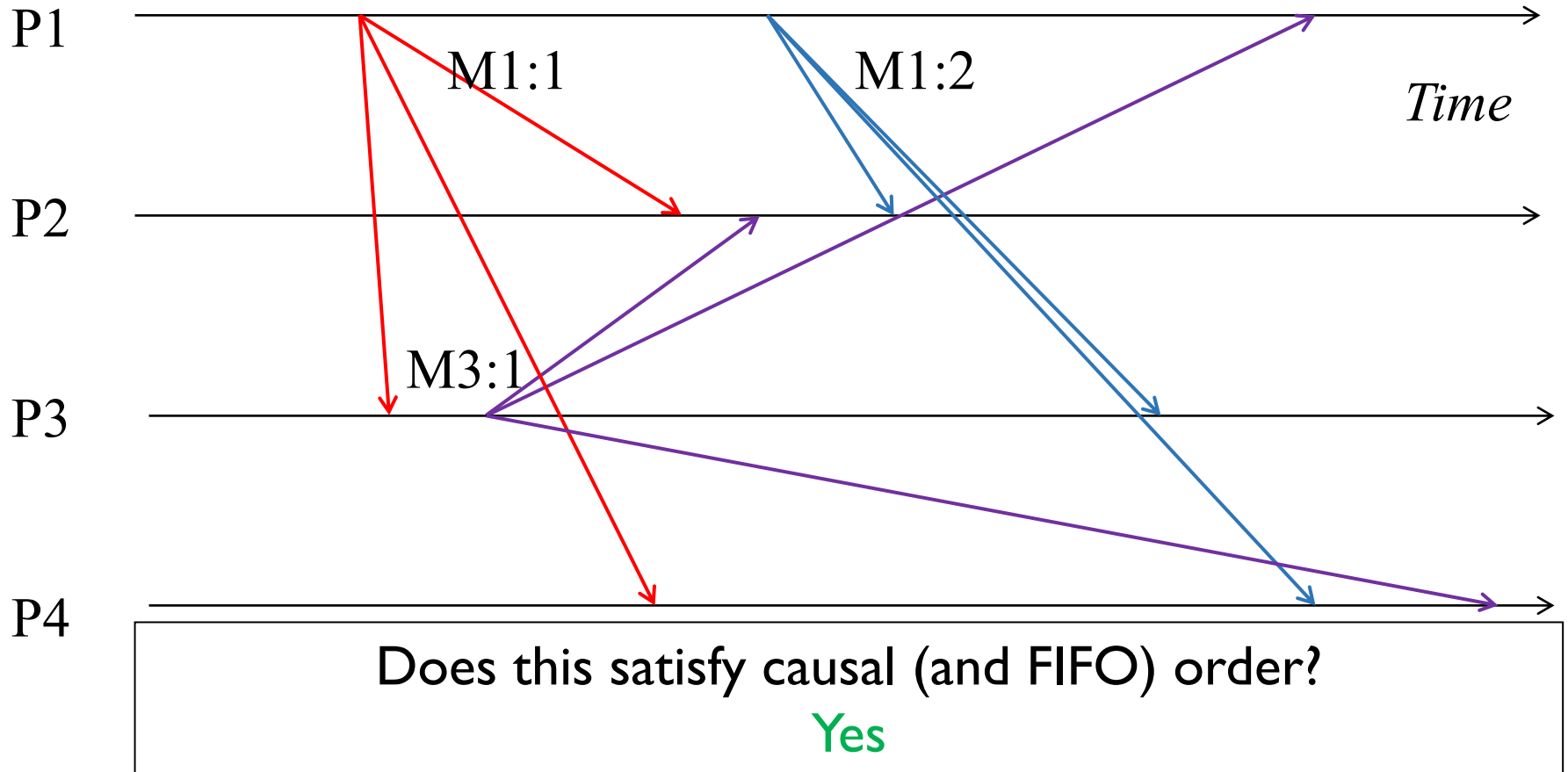
Causal vs Total

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

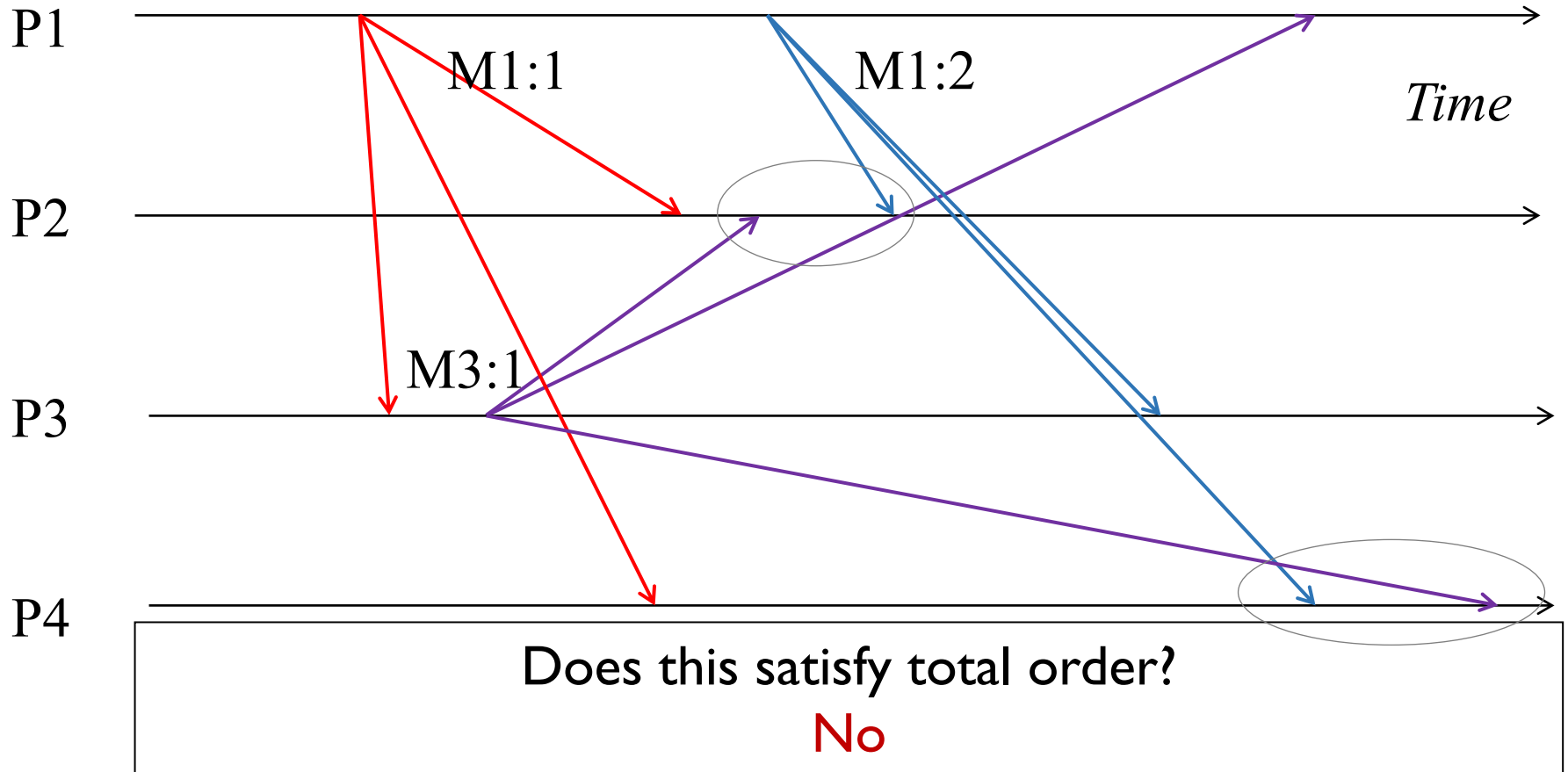
Hybrid variants

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

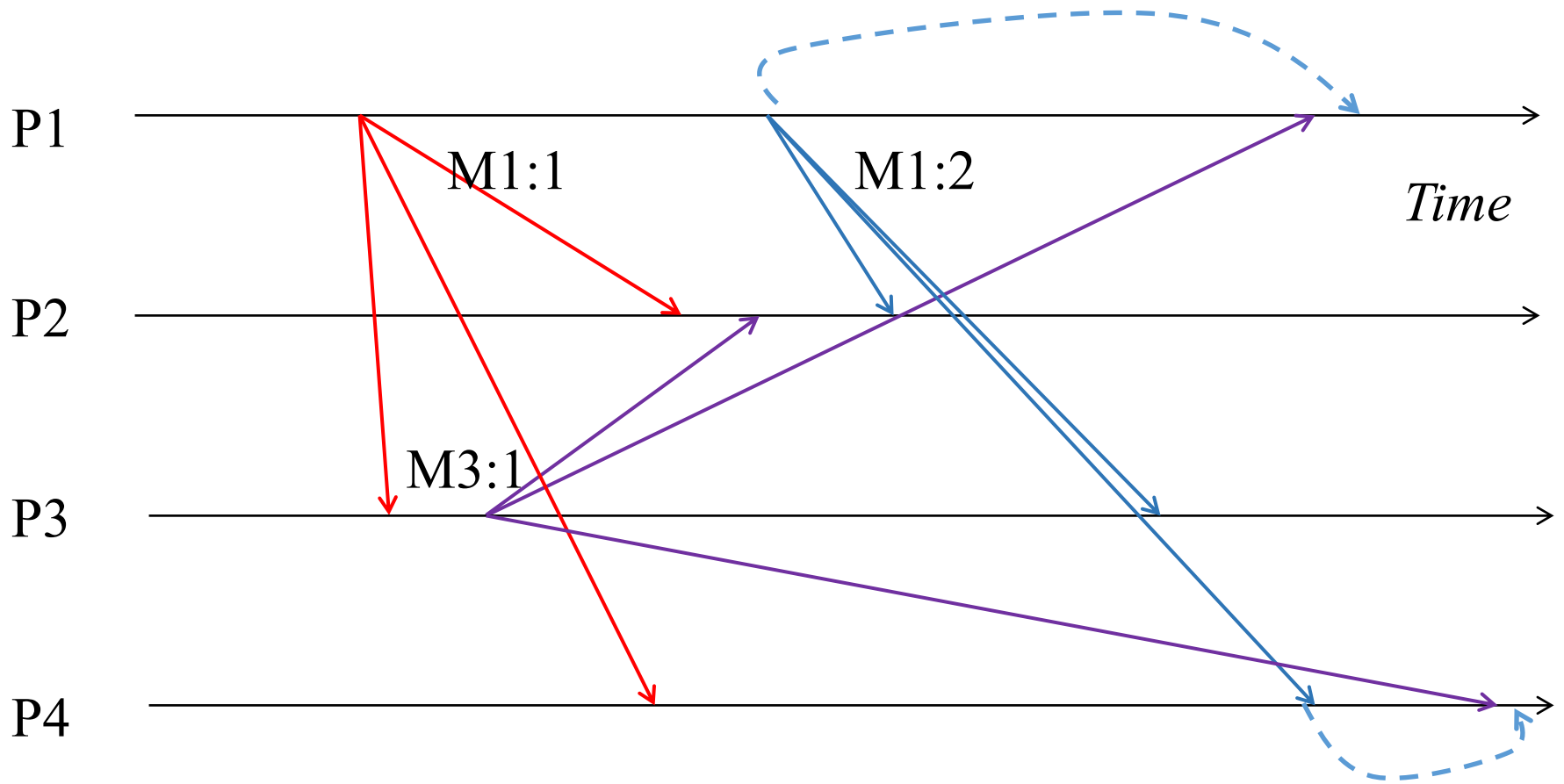
Example



Example



Example



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 .

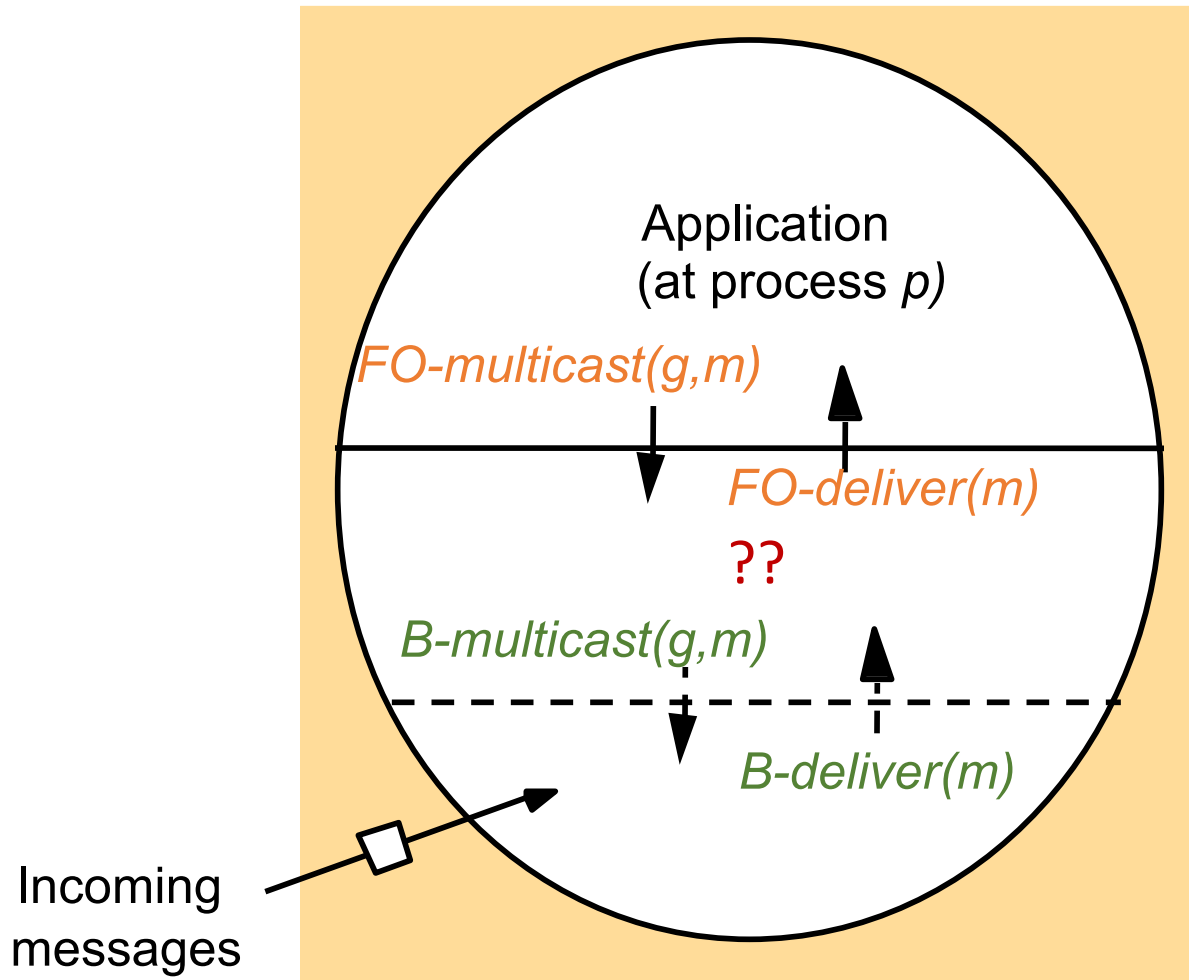
- **Causal ordering**

- 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



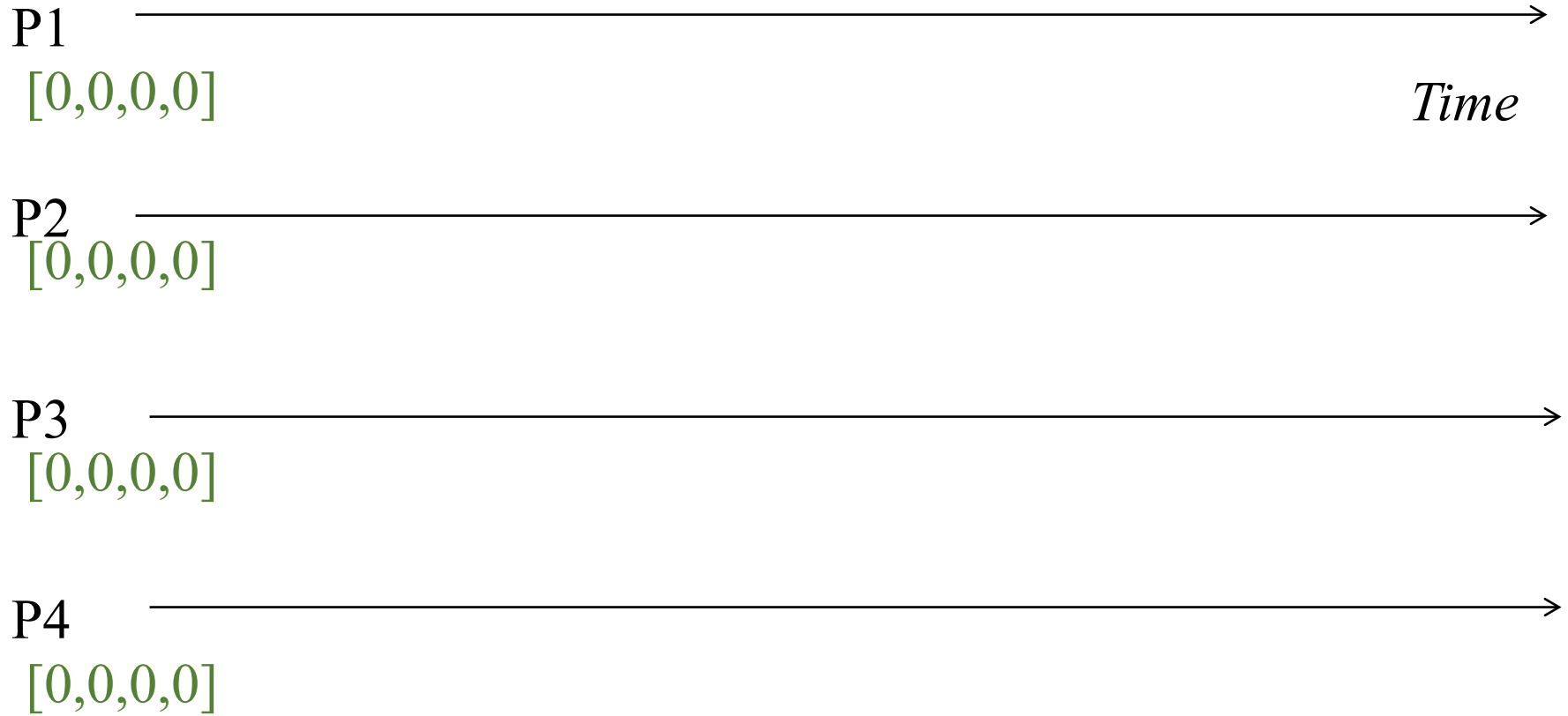
Implementing FIFO order multicast

- 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



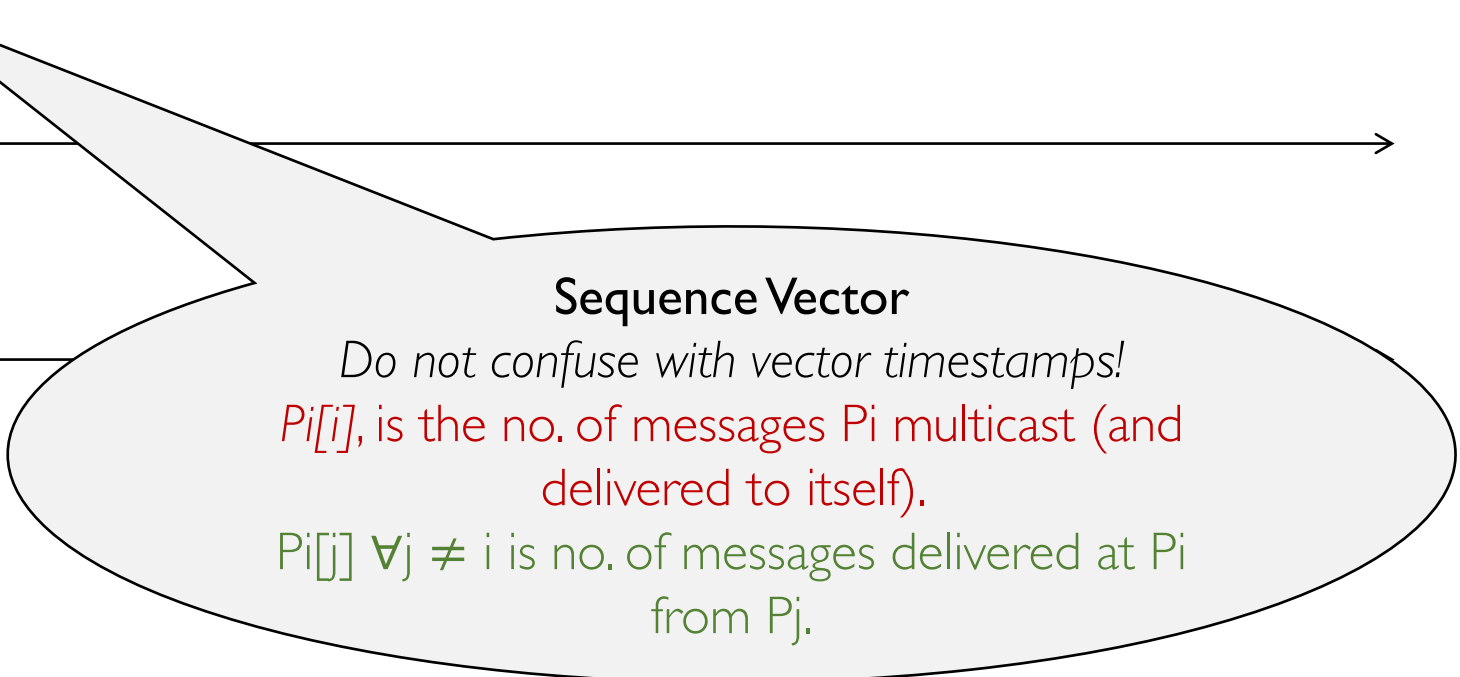
FIFO order multicast execution

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

P2  *Time*
[0,0,0,0]

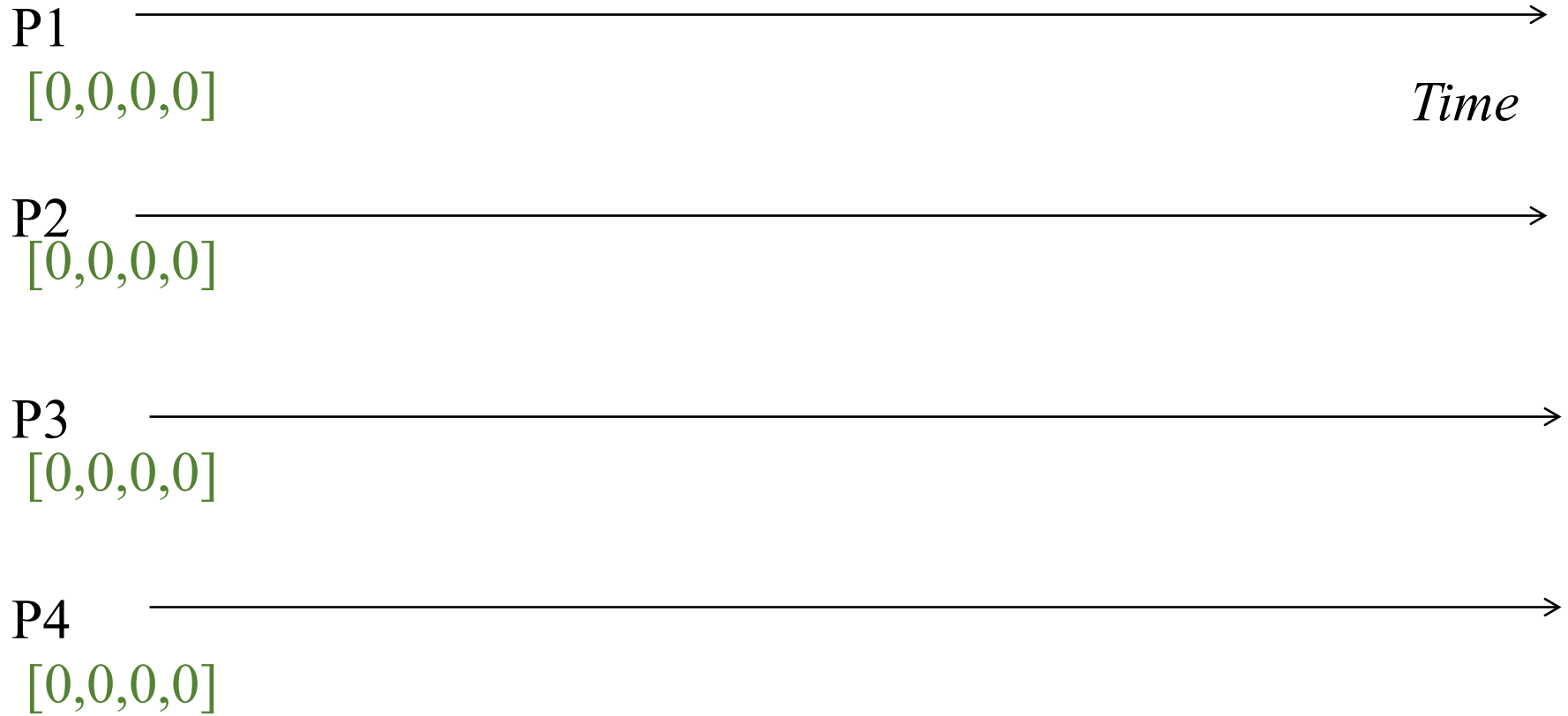
P3  *Time*
[0,0,0,0]

P4  *Time*
[0,0,0,0]

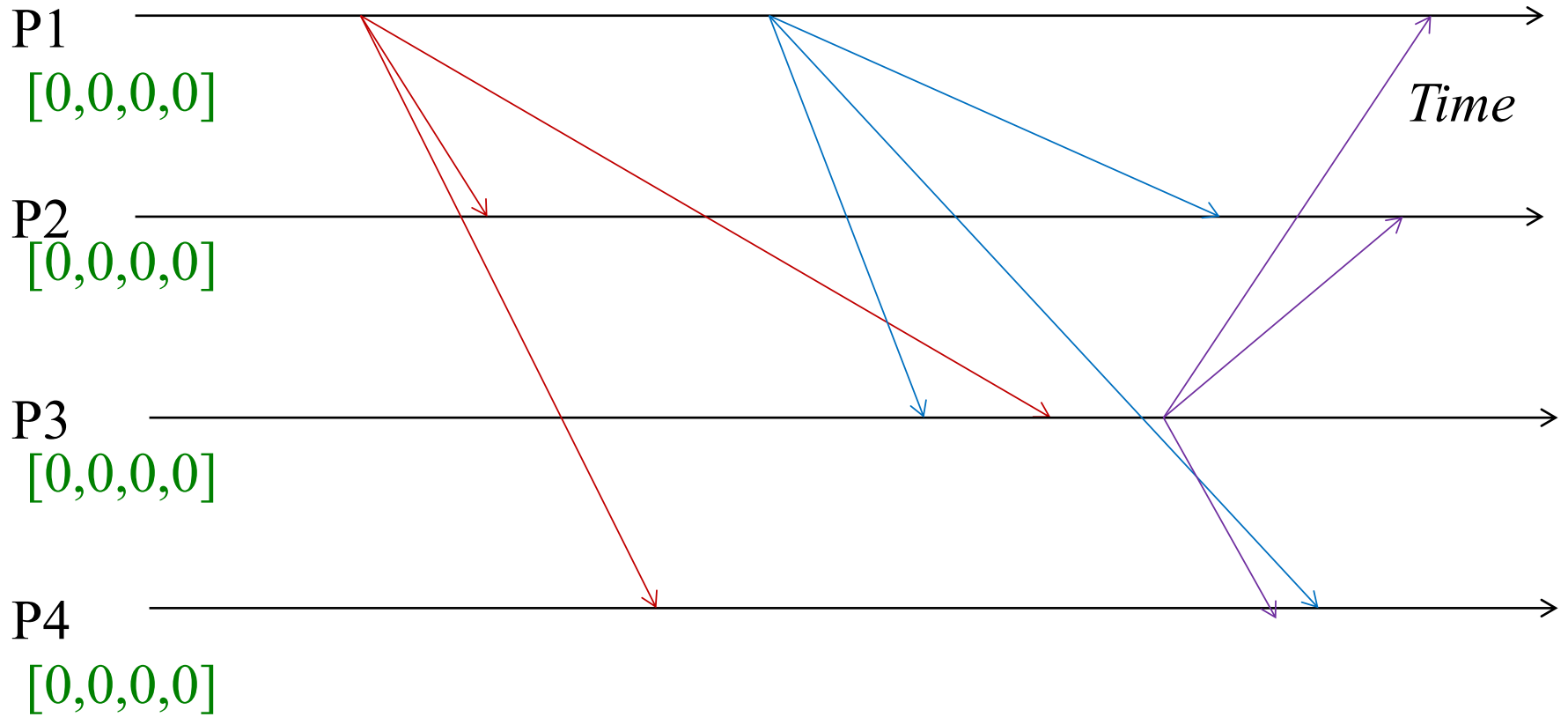


Sequence Vector
Do not confuse with vector timestamps!
 $P_i[i]$, is the no. of messages P_i multicast (and delivered to itself).
 $P_i[j] \forall j \neq i$ is no. of messages delivered at P_i from P_j .

FIFO order multicast execution

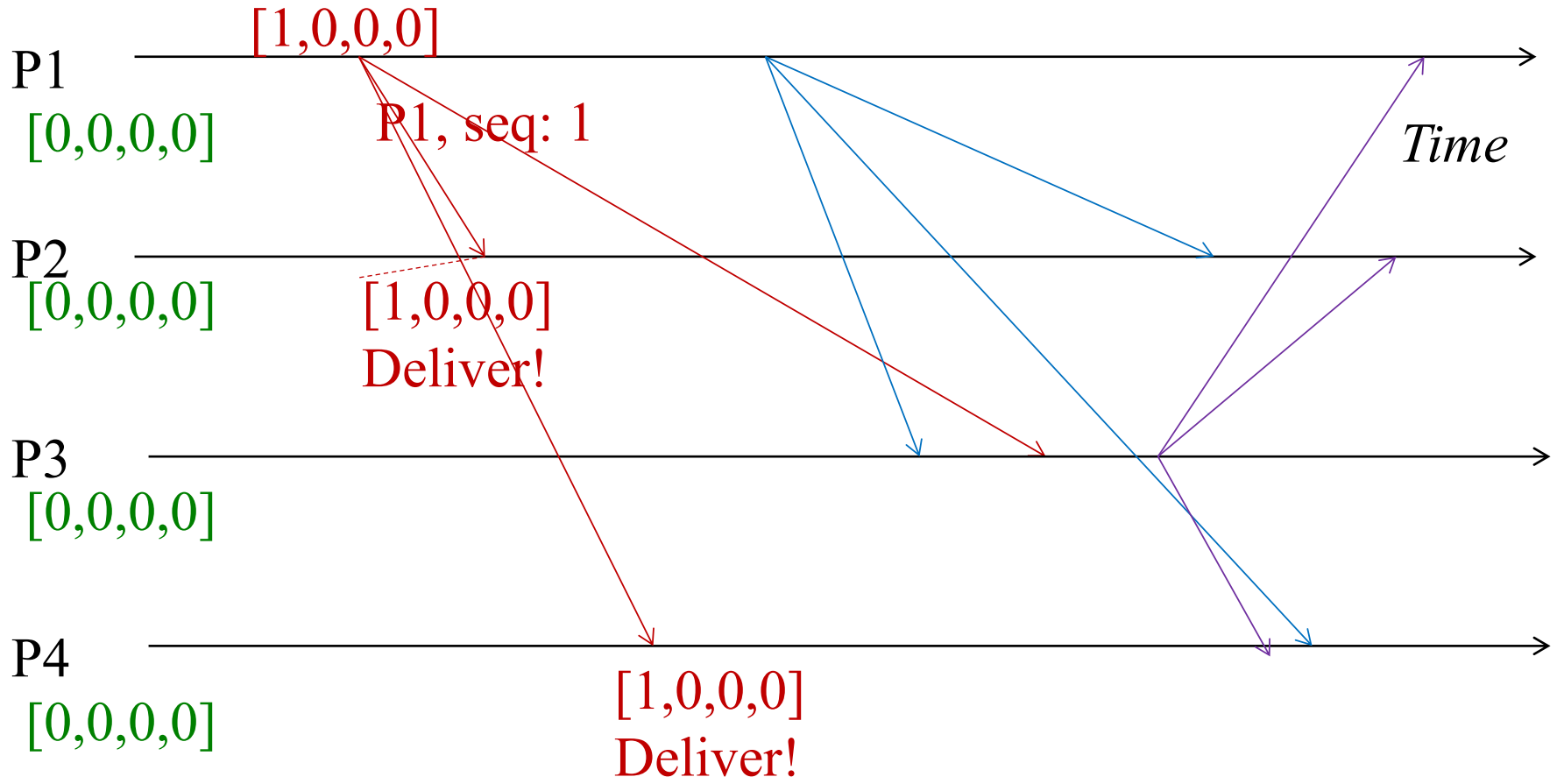


FIFO order multicast execution

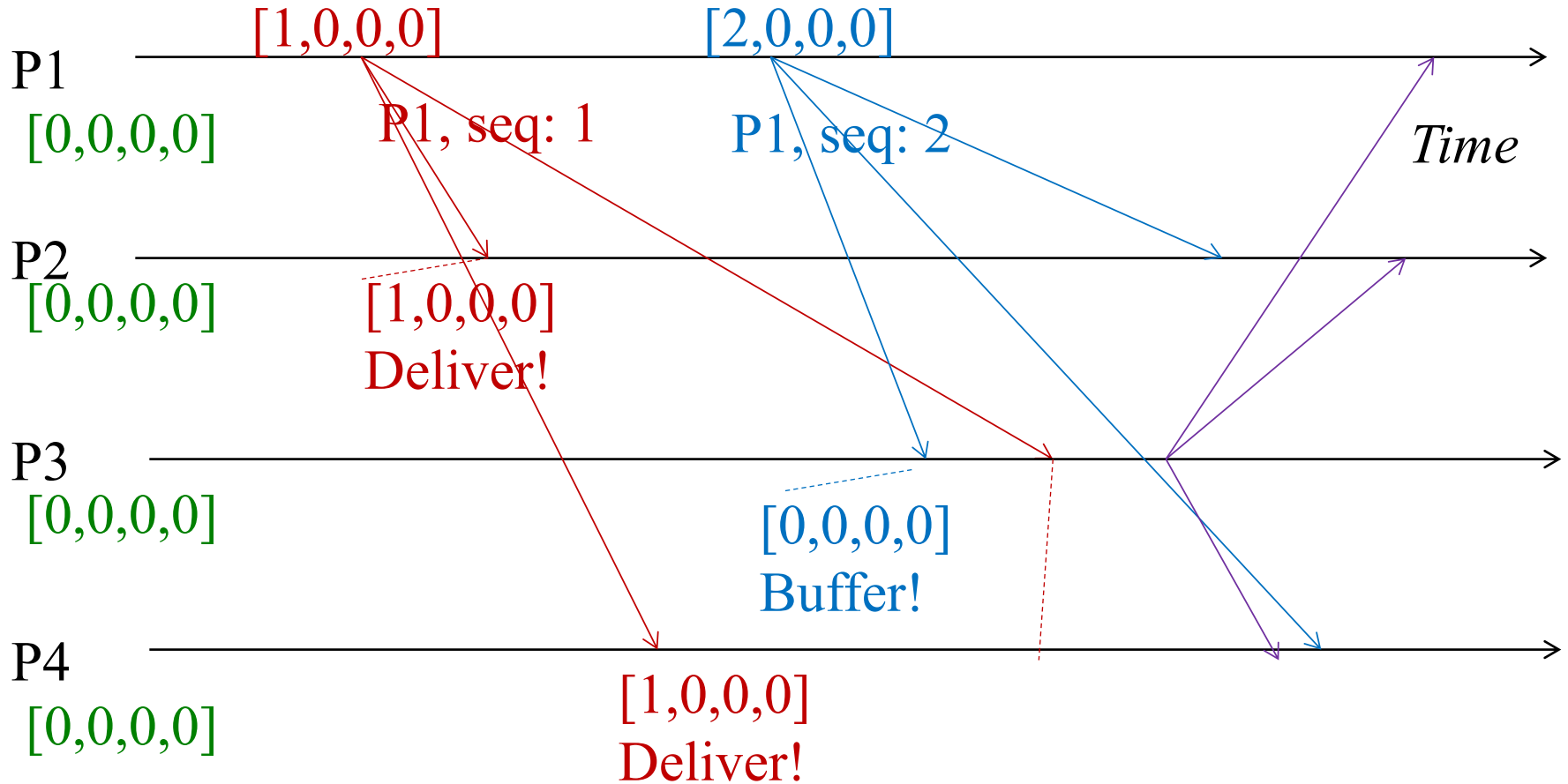


Self-deliveries omitted for simplicity.

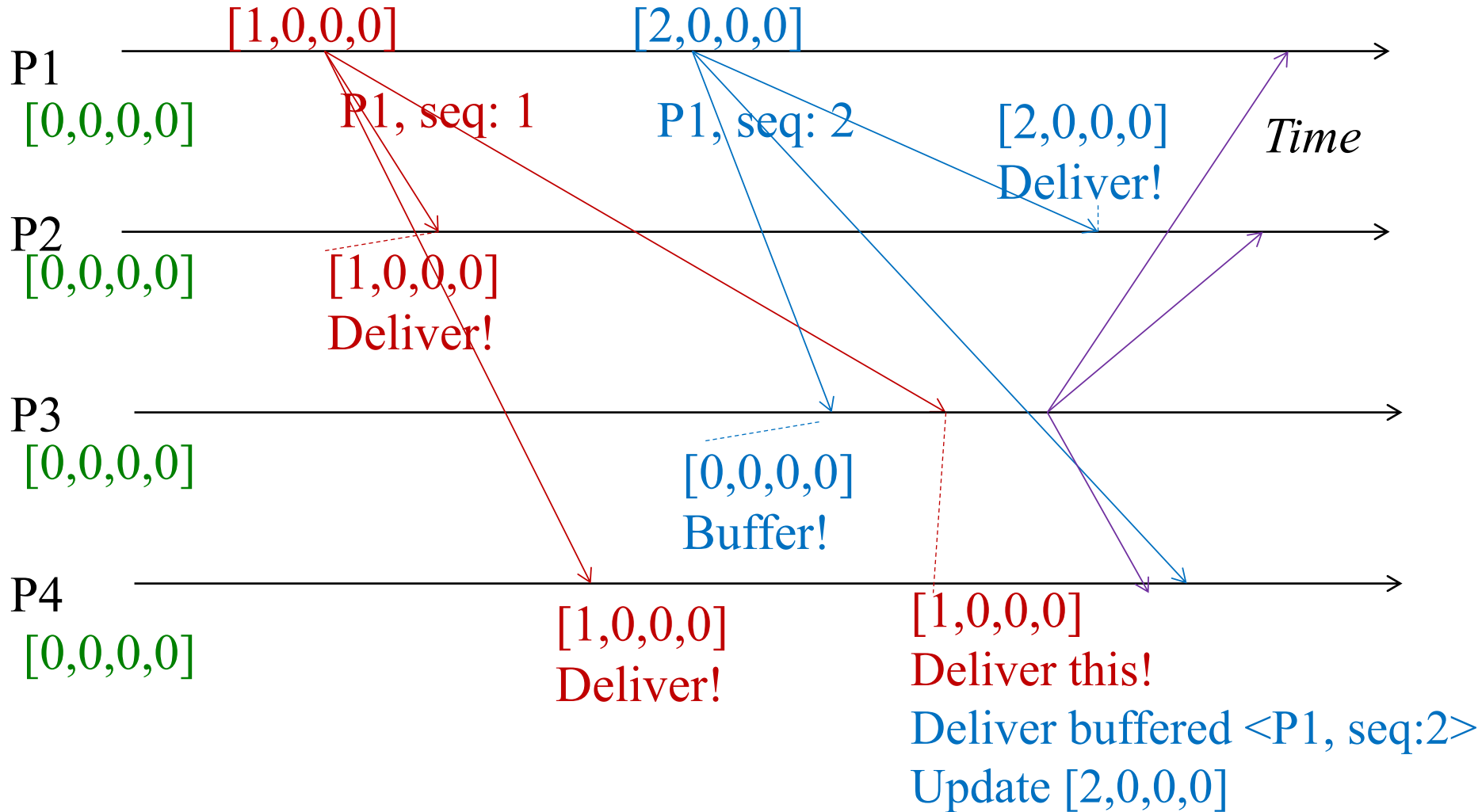
FIFO order multicast execution



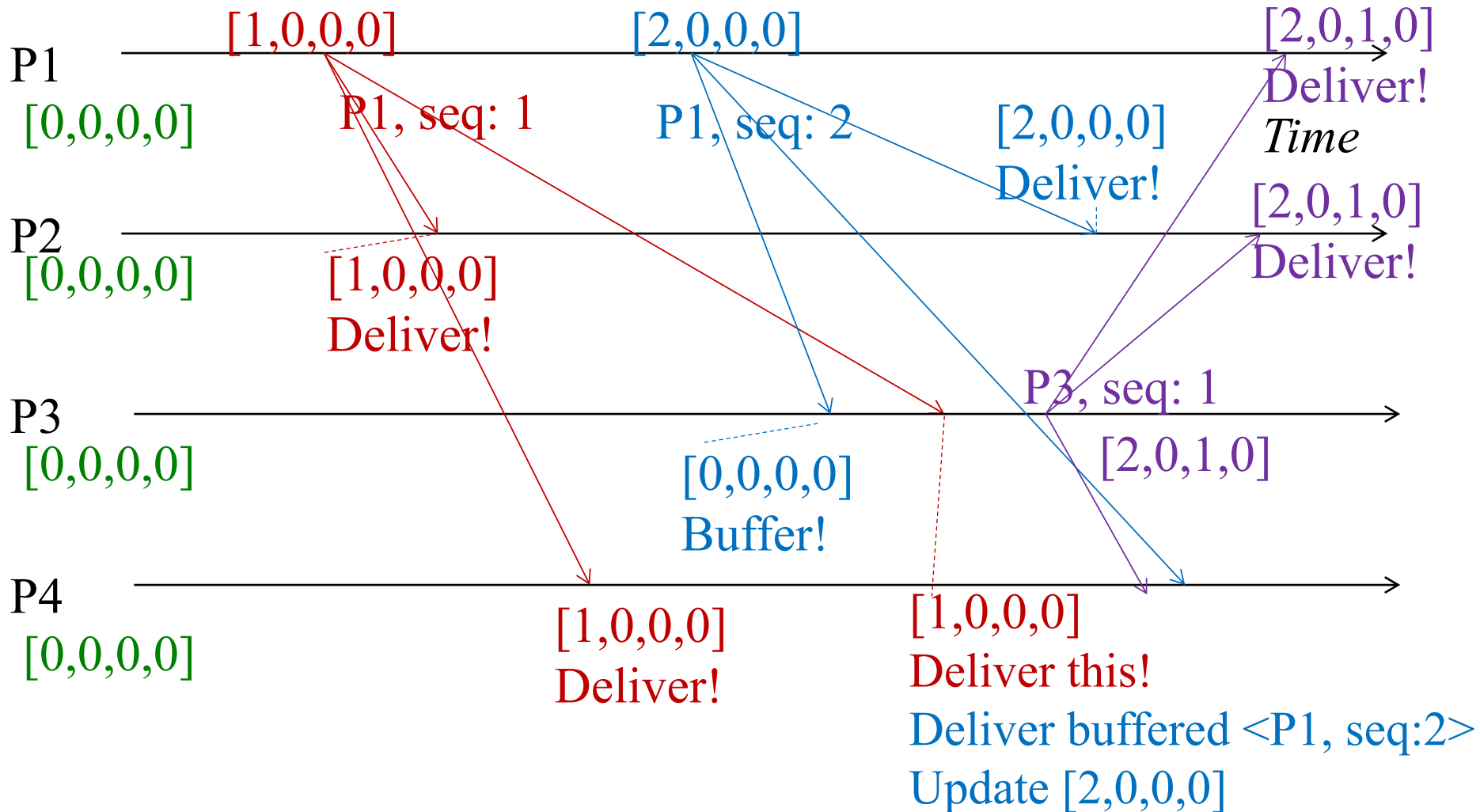
FIFO order multicast execution



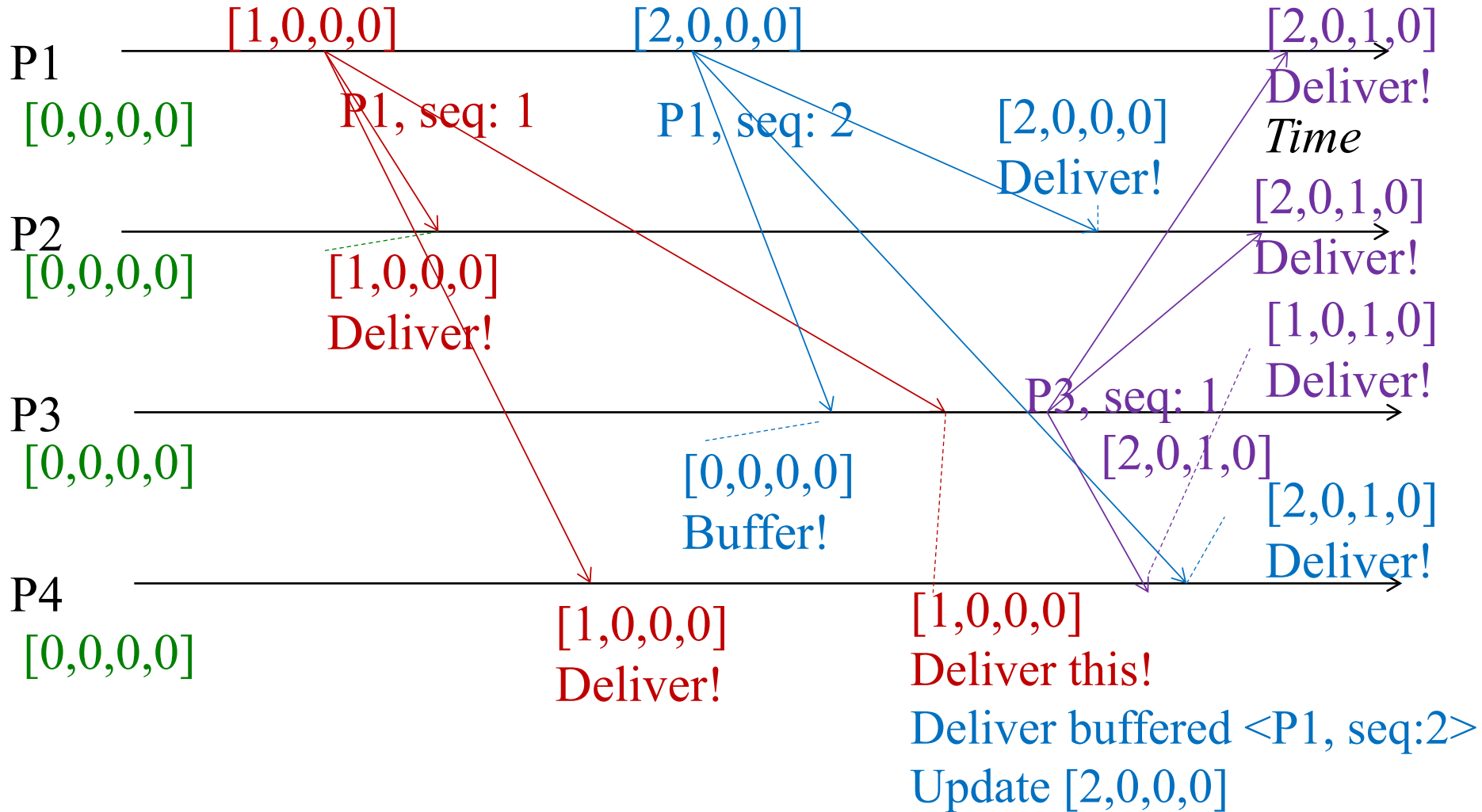
FIFO order multicast execution



FIFO order multicast execution



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 .

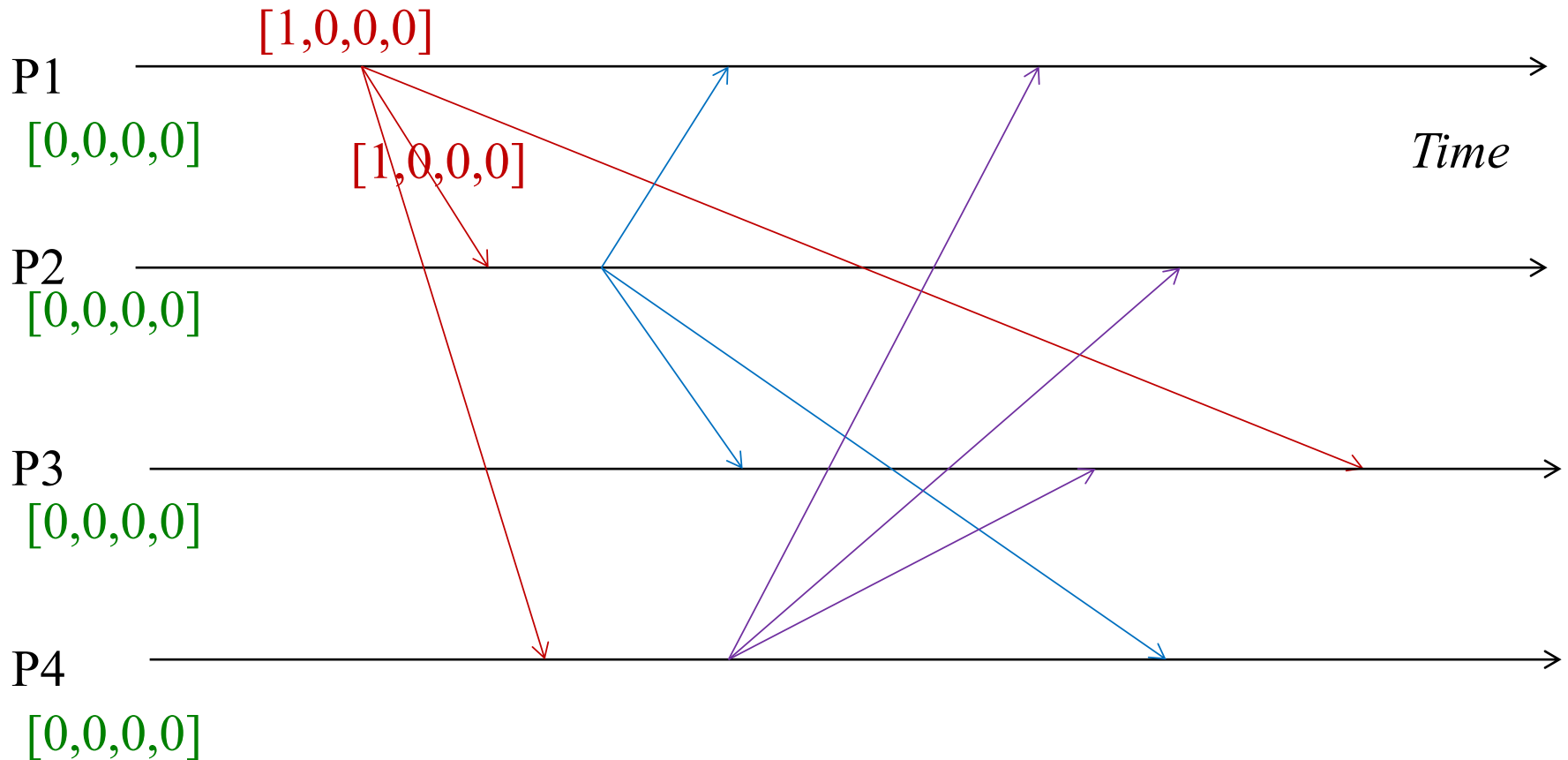
Implementing causal order multicast

- Similar to FIFO Multicast
 - What you send with a message differs.
 - Updating rules differ.
- Each receiver maintains a vector of per-sender sequence numbers (integers)
 - 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 .
- Ignores other network messages. Only looks at multicast messages delivered to the application.

Implementing causal order multicast

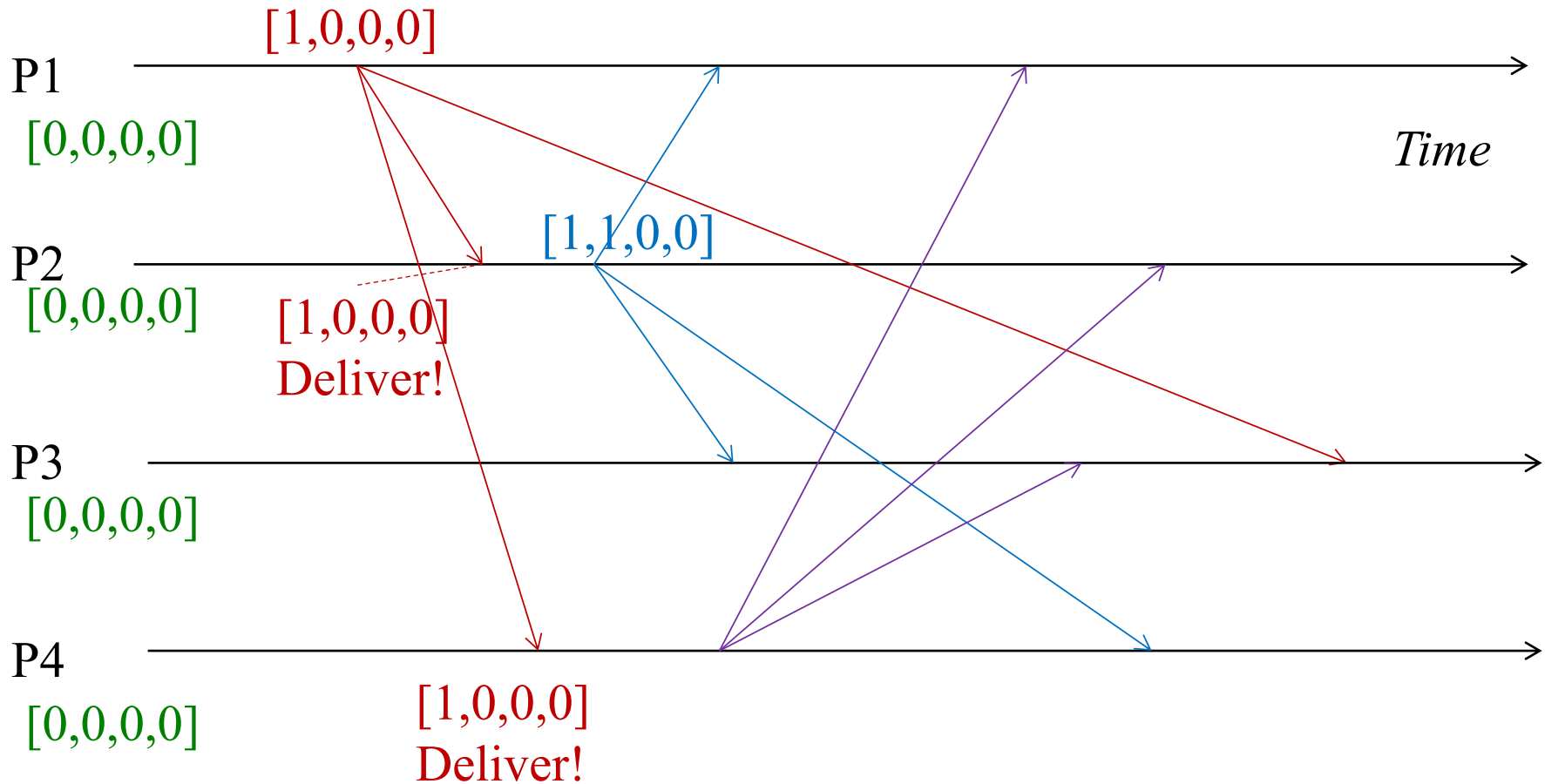
- *CO-multicast*(g, m) at P_j :
 - set $P_j[j] = P_j[j] + 1$
 - piggyback entire vector $P_j[1 \dots N]$ with m as its sequence no.
 - B-multicast($g, \{m, P_j[1 \dots N]\}$)
- On B-deliver($\{m, V[1 \dots N]\}$) at P_i from P_j : If P_i receives a multicast from P_j with sequence vector $V[1 \dots N]$, buffer it until both:
 1. This message is the next one P_i is expecting from P_j , i.e.,
$$V[j] = P_i[j] + 1$$
 2. All multicasts, anywhere in the group, which happened-before m have been received at P_i , i.e.,
$$\text{For all } k \neq j: V[k] \leq P_i[k]$$When above two conditions satisfied,
CO-deliver(m) and set $P_i[j] = V[j]$

Causal order multicast execution

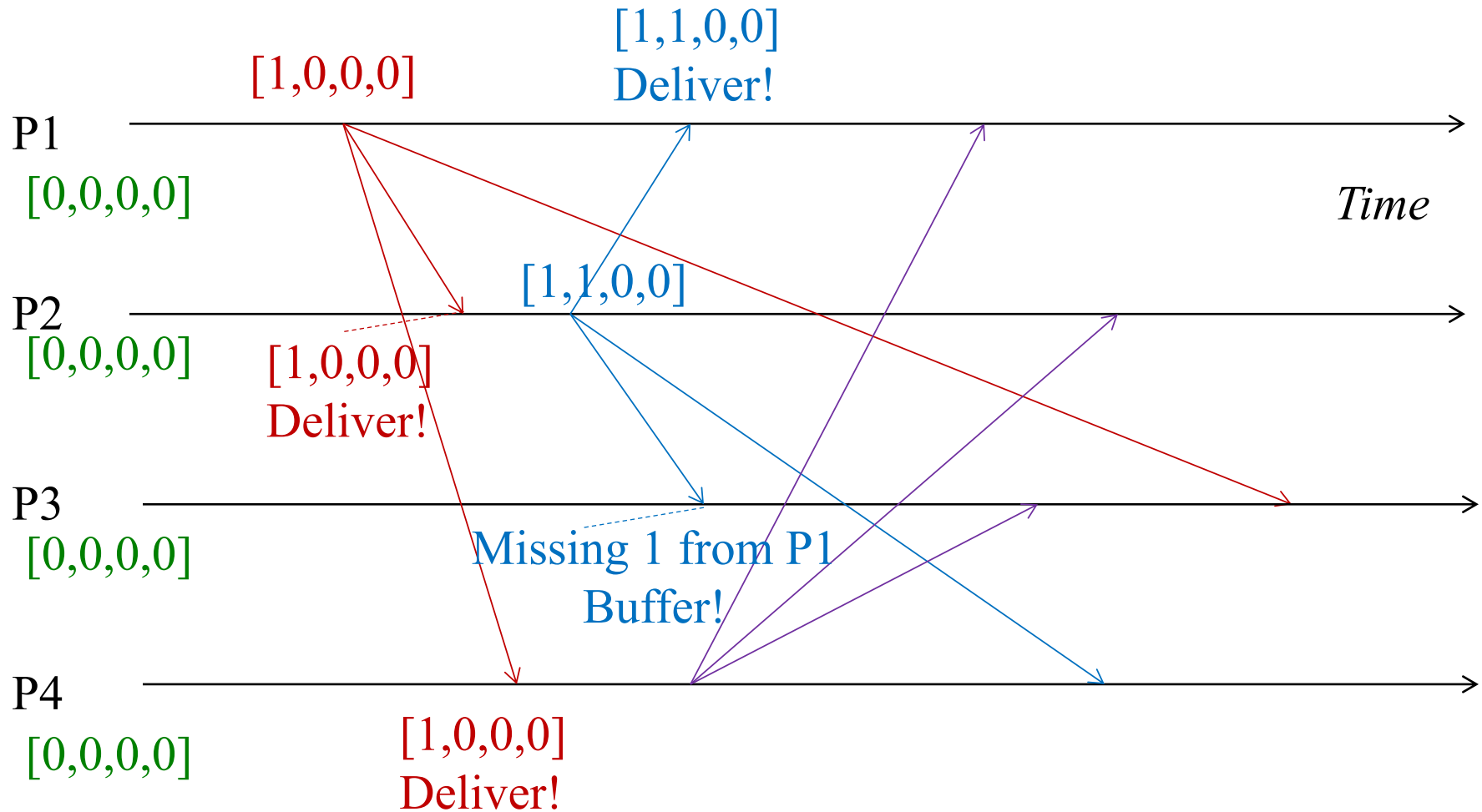


Self-deliveries omitted for simplicity.

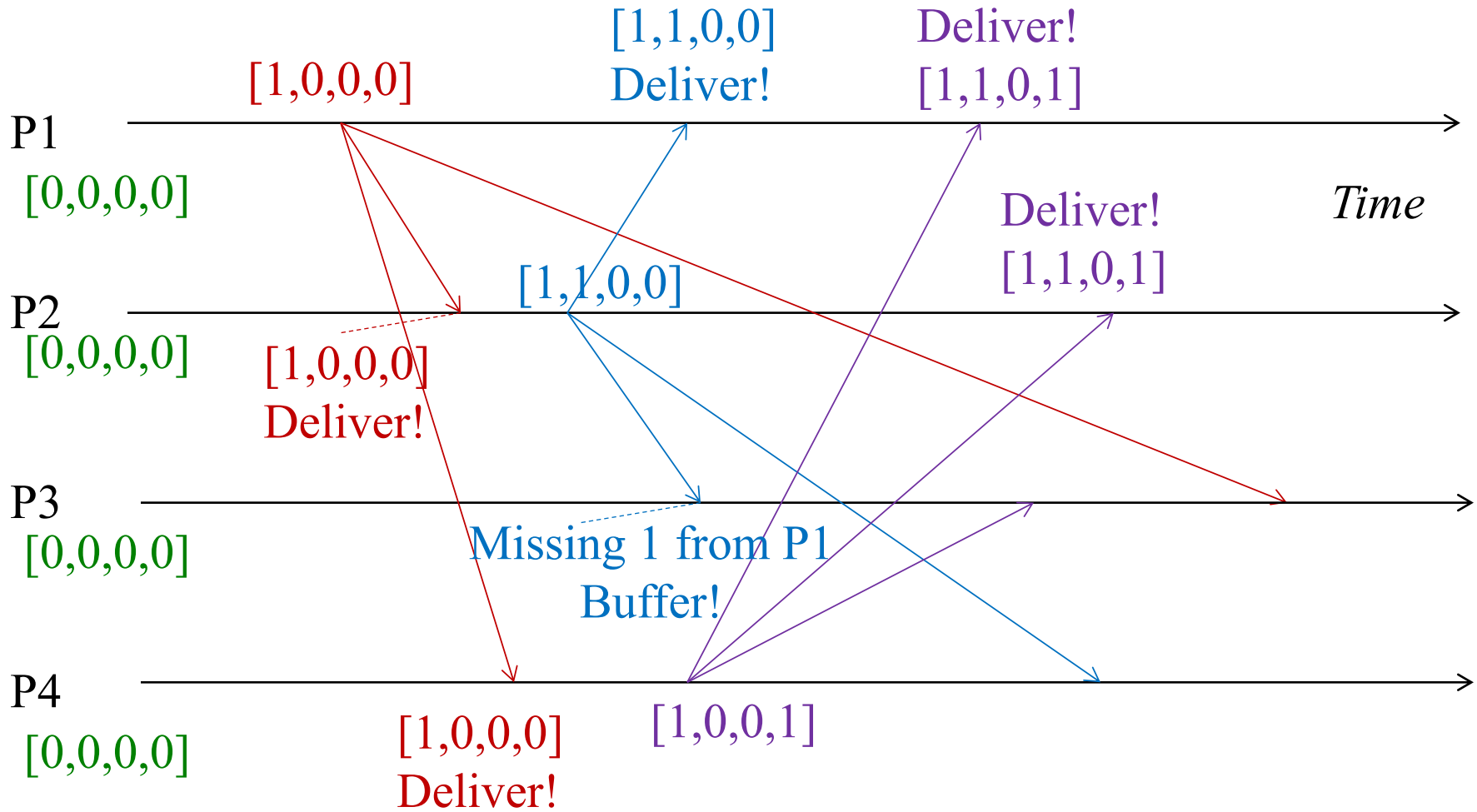
Causal order multicast execution



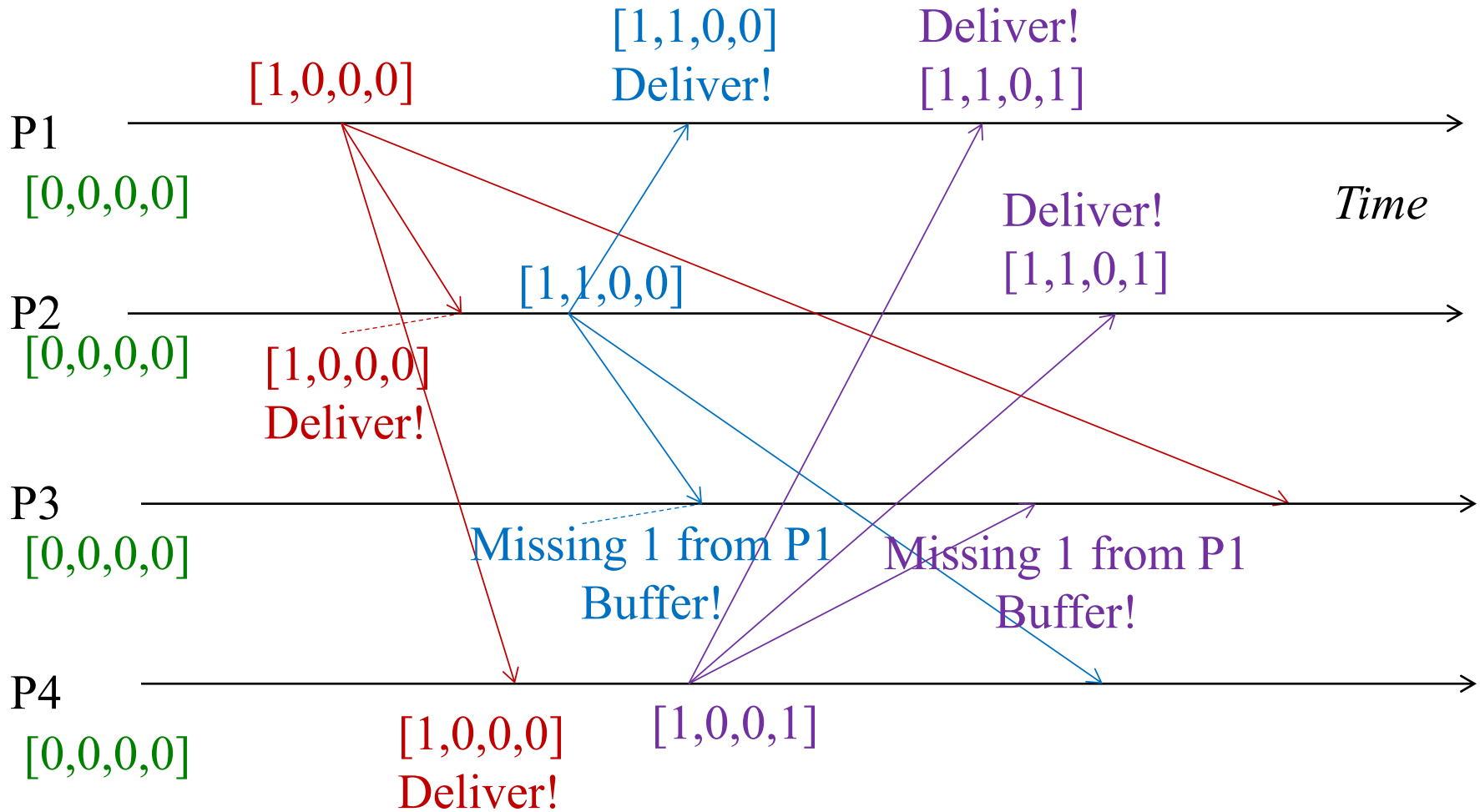
Causal order multicast execution



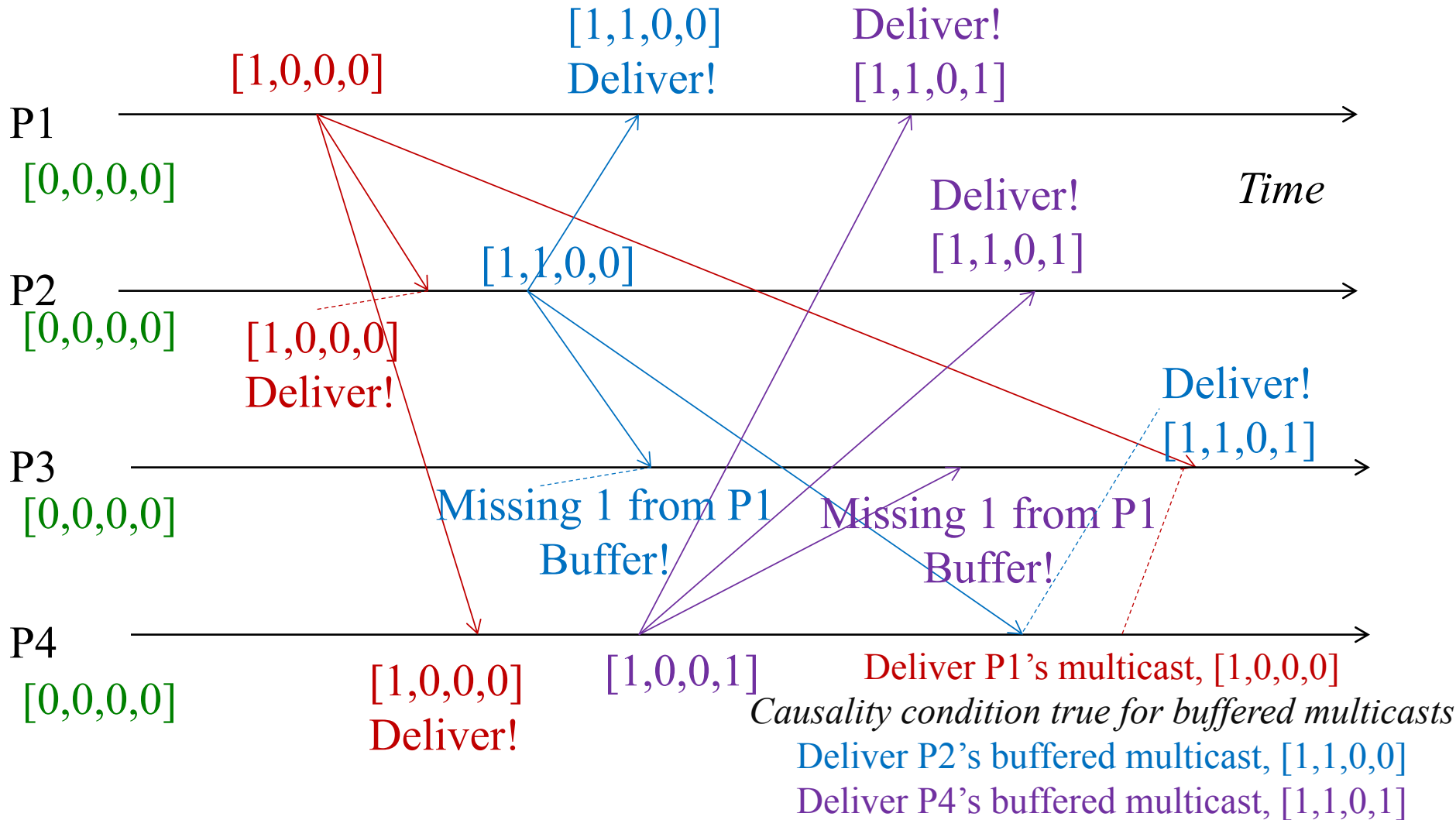
Causal order multicast execution



Causal order multicast execution



Causal order multicast execution



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)

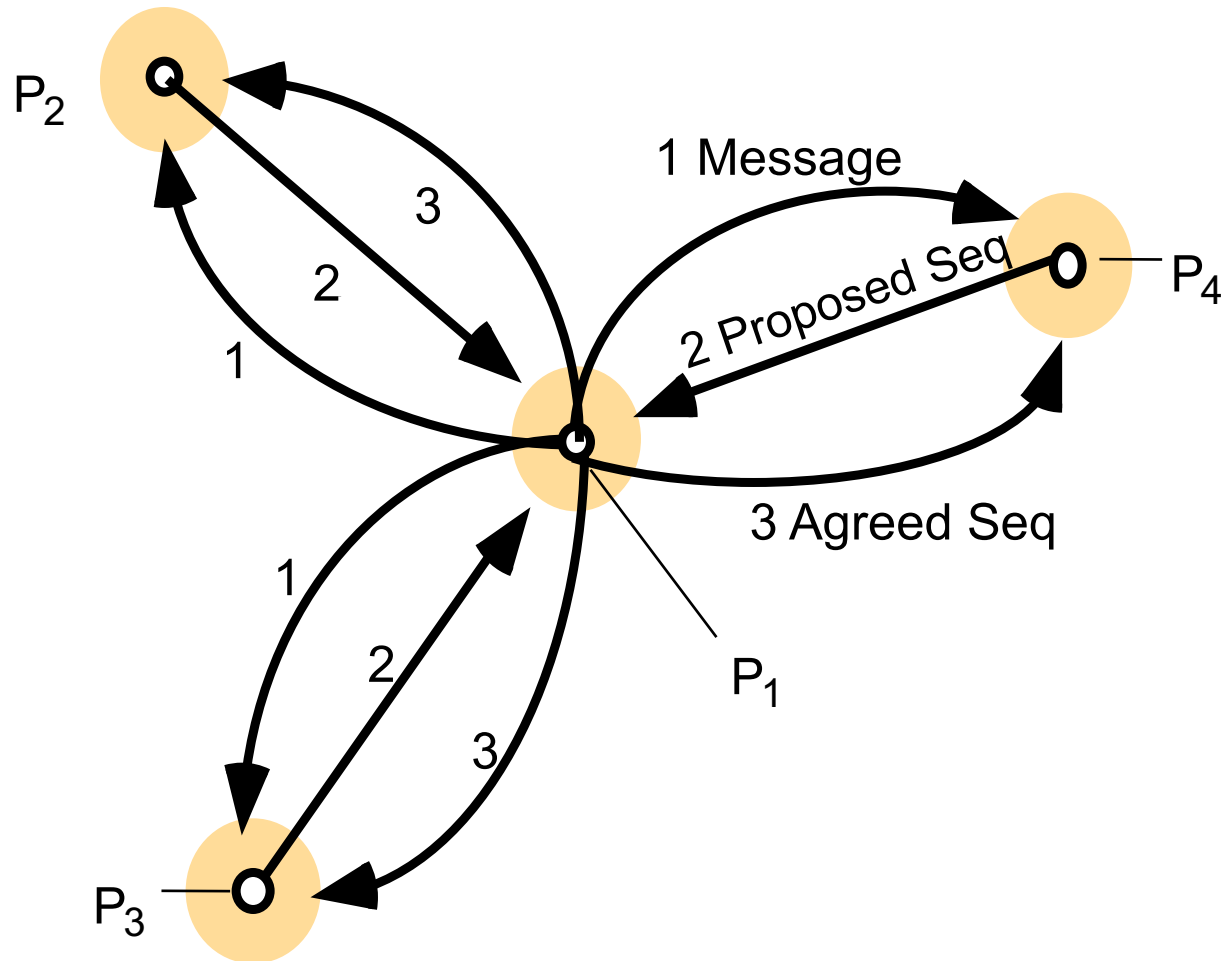
Sequencer based total ordering

- 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 $S_i = S_i + 1$

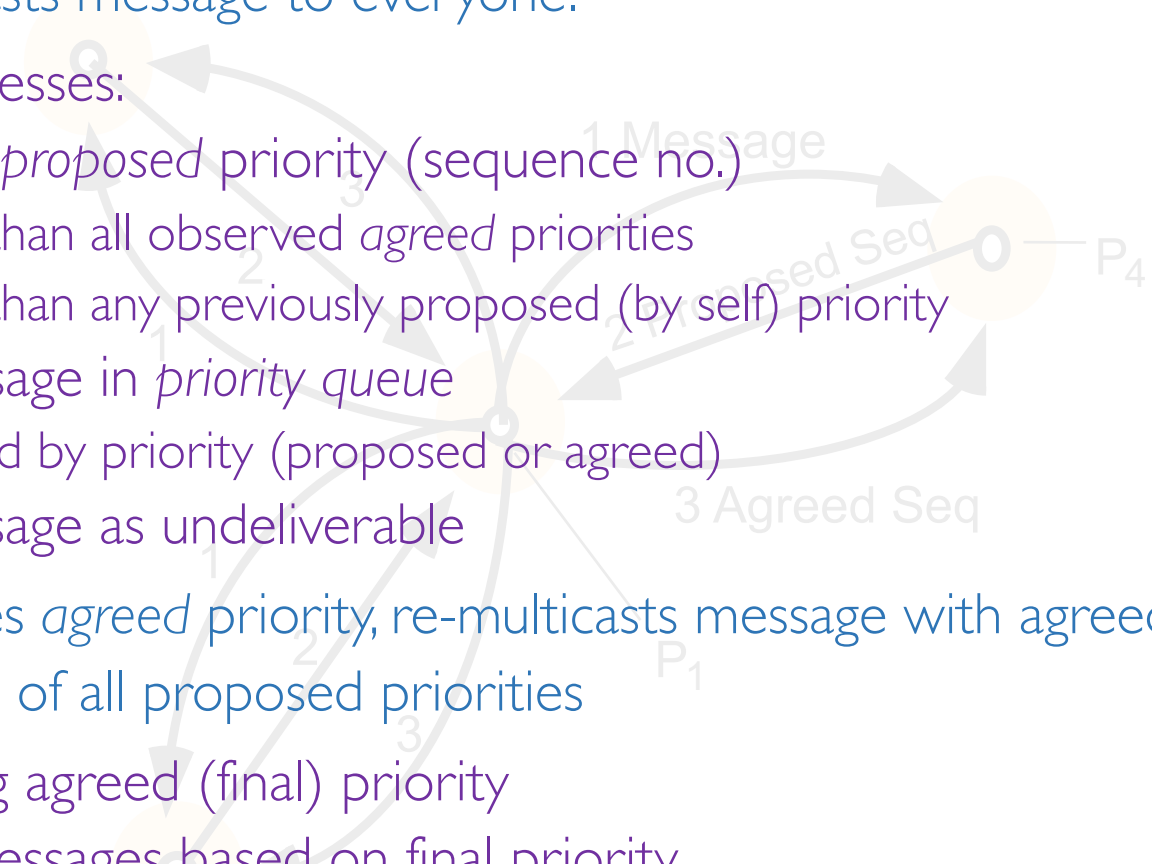
Implementing total order multicast

- 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)**

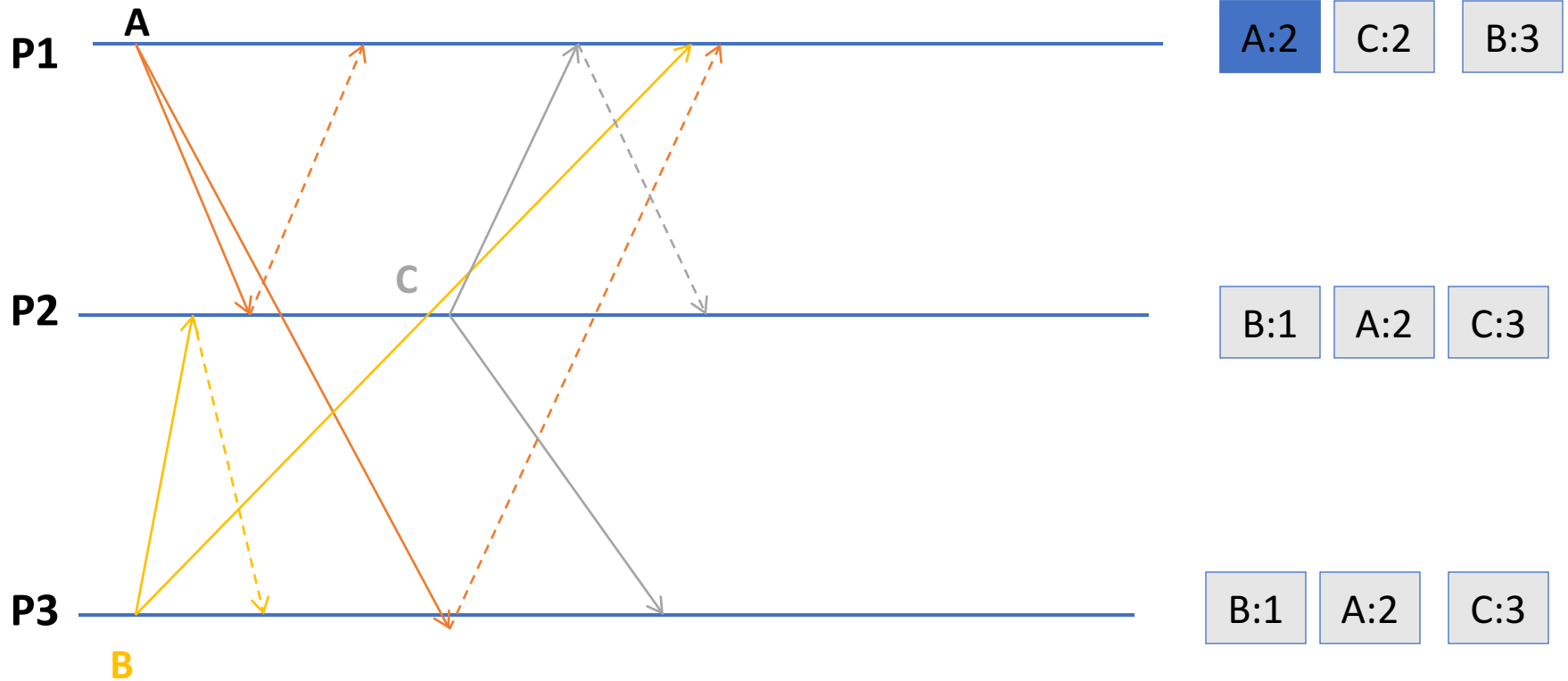
ISIS algorithm for total ordering



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 with agreed priority
 - maximum of all proposed priorities
 - Upon receiving agreed (final) priority
 - reorder messages based on final priority.
 - mark the message as deliverable.
 - deliver any deliverable messages at front of priority queue.
- 

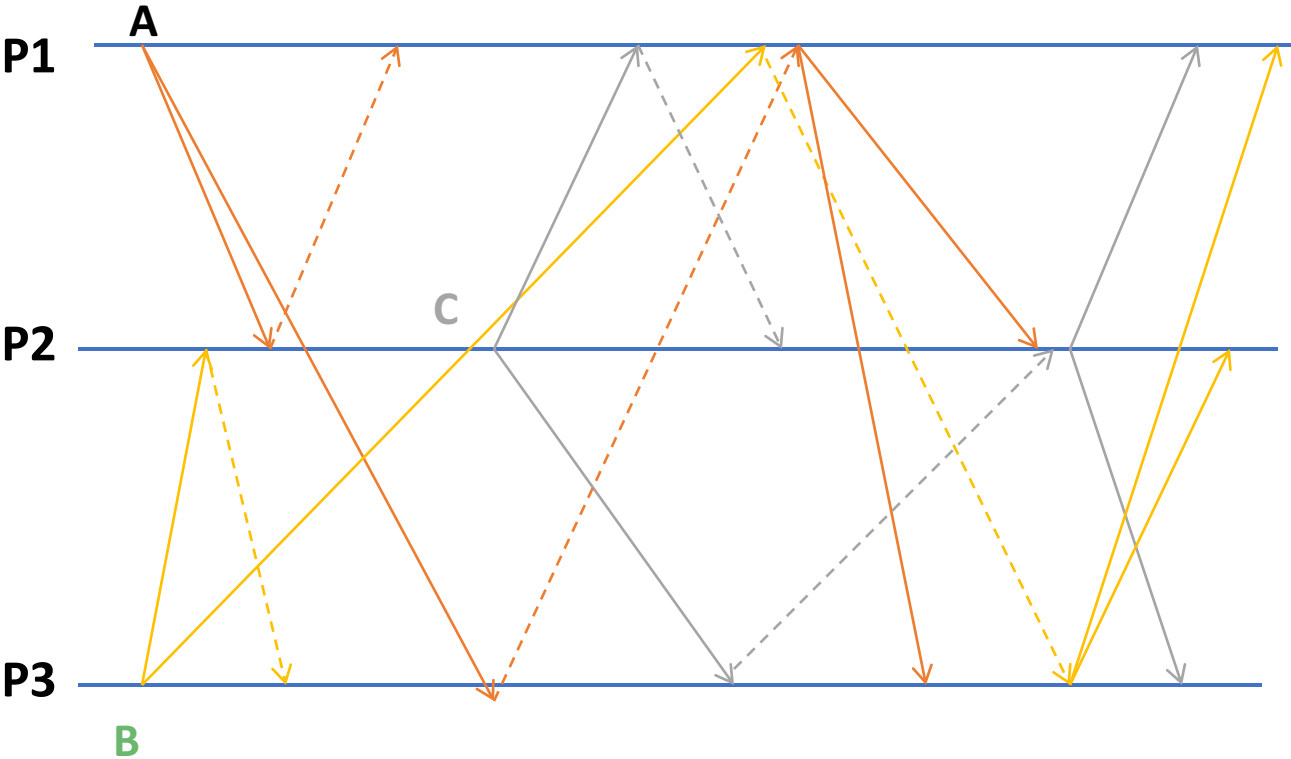
Example: ISIS algorithm



How do we break ties?

- Problem: priority queue requires unique priorities.
- Solution: add process # to suggested priority.
 - *priority.(id of the process that proposed the priority)*
 - i.e., 3.2 == process 2 proposed priority 3
- Compare on priority first, use process # to break ties.
 - $2.1 > 1.3$
 - $3.2 > 3.1$

Example: ISIS algorithm



A:2.3 ✓	C:2.1 ✓	B:3.1 ✓
---------	---------	---------

B:3.1 ✓	A:2.3 ✓	C:3.3 ✓
---------	---------	---------

B:3.1 ✓	A:2.3 ✓	C:3.3 ✓
---------	---------	---------

Proof of total order with ISIS

- Consider two messages, m_1 and m_2 , and two processes, p and p' .
- Suppose that p delivers m_1 before m_2 .
- When p delivers m_1 , it is at the head of the queue. m_2 is either:
 - Already in p 's queue, and deliverable, so
 - $\text{finalpriority}(m_1) < \text{finalpriority}(m_2)$
 - Already in p 's queue, and not deliverable, so
 - $\text{finalpriority}(m_1) < \text{proposedpriority}(m_2) \leq \text{finalpriority}(m_2)$
 - Not yet in p 's queue:
 - same as above, since proposed priority $>$ priority of any delivered message
- Suppose p' delivers m_2 before m_1 , by the same argument:
 - $\text{finalpriority}(m_2) < \text{finalpriority}(m_1)$
 - Contradiction!

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.

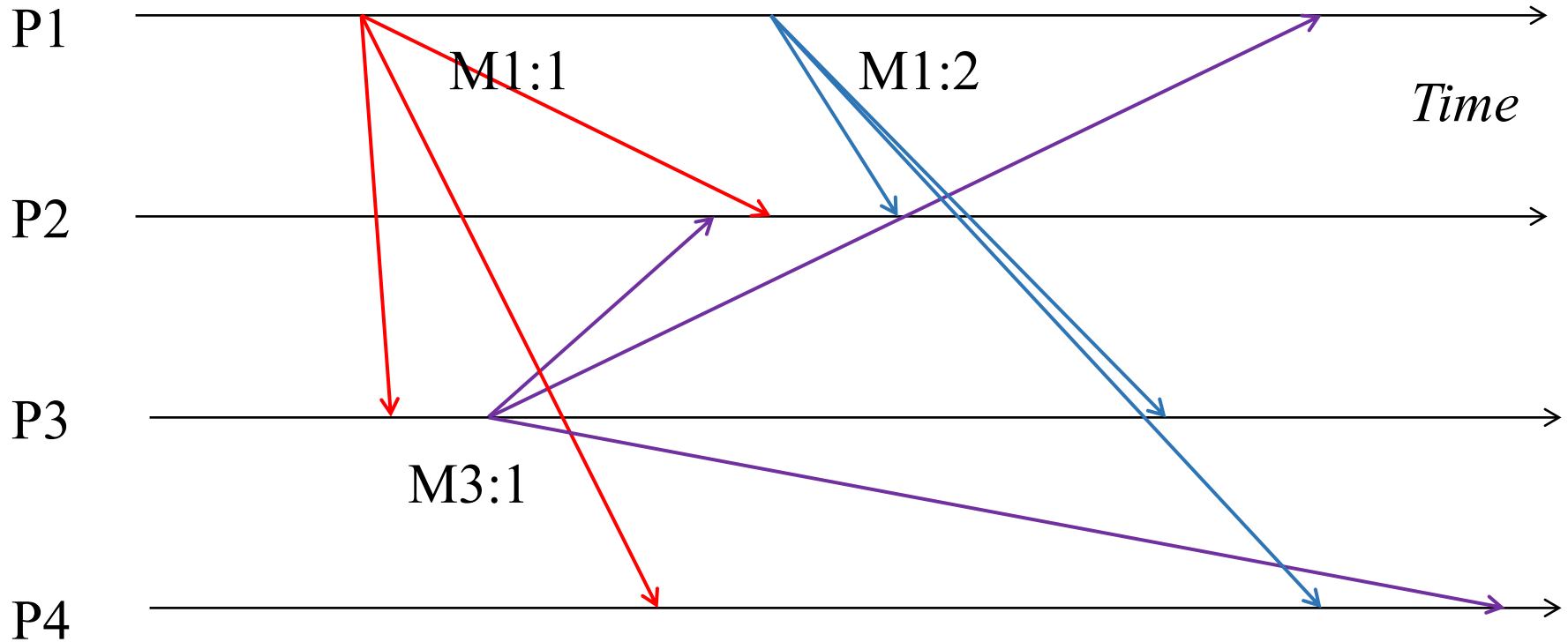
- **Total ordering**

- 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 .

Summary

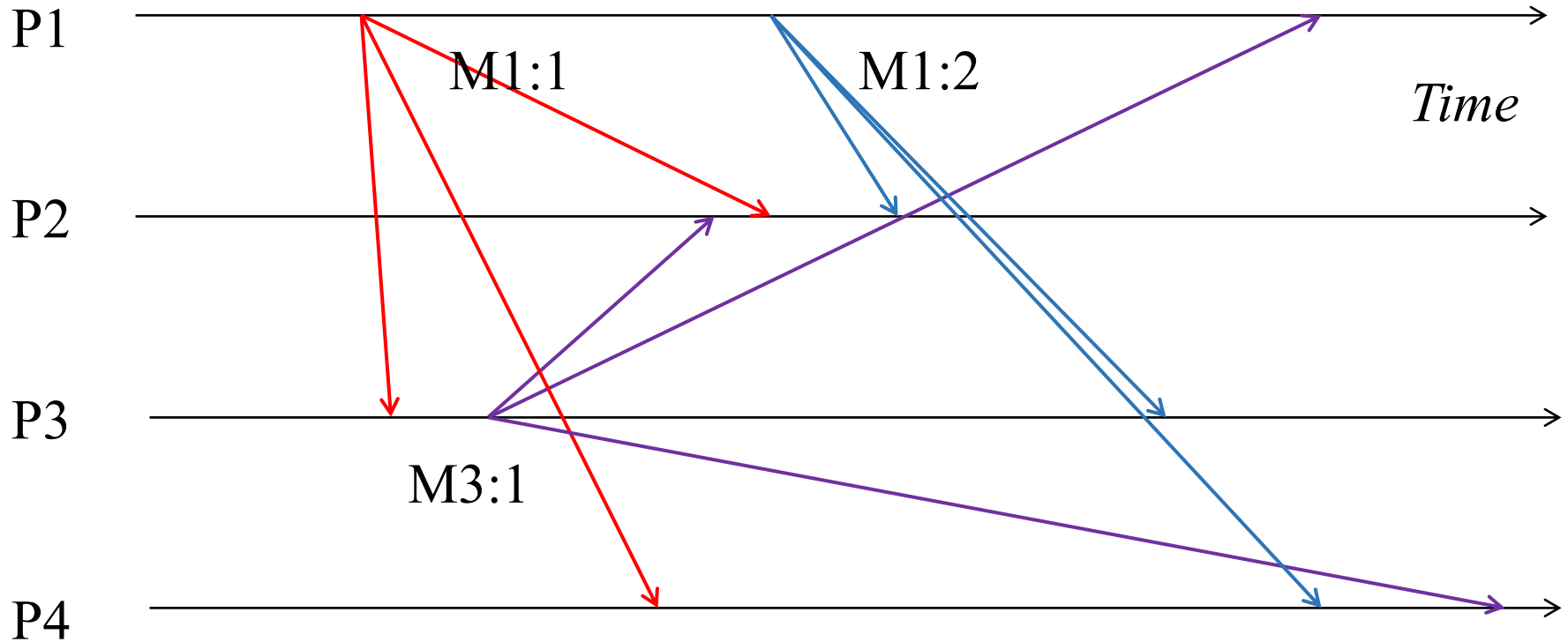
- Multicast is an important communication mode in distributed systems.
- Applications may have different requirements:
 - Reliability
 - Ordering: FIFO, Causal, Total
 - Combinations of the above.

Example



Does this satisfy FIFO order?

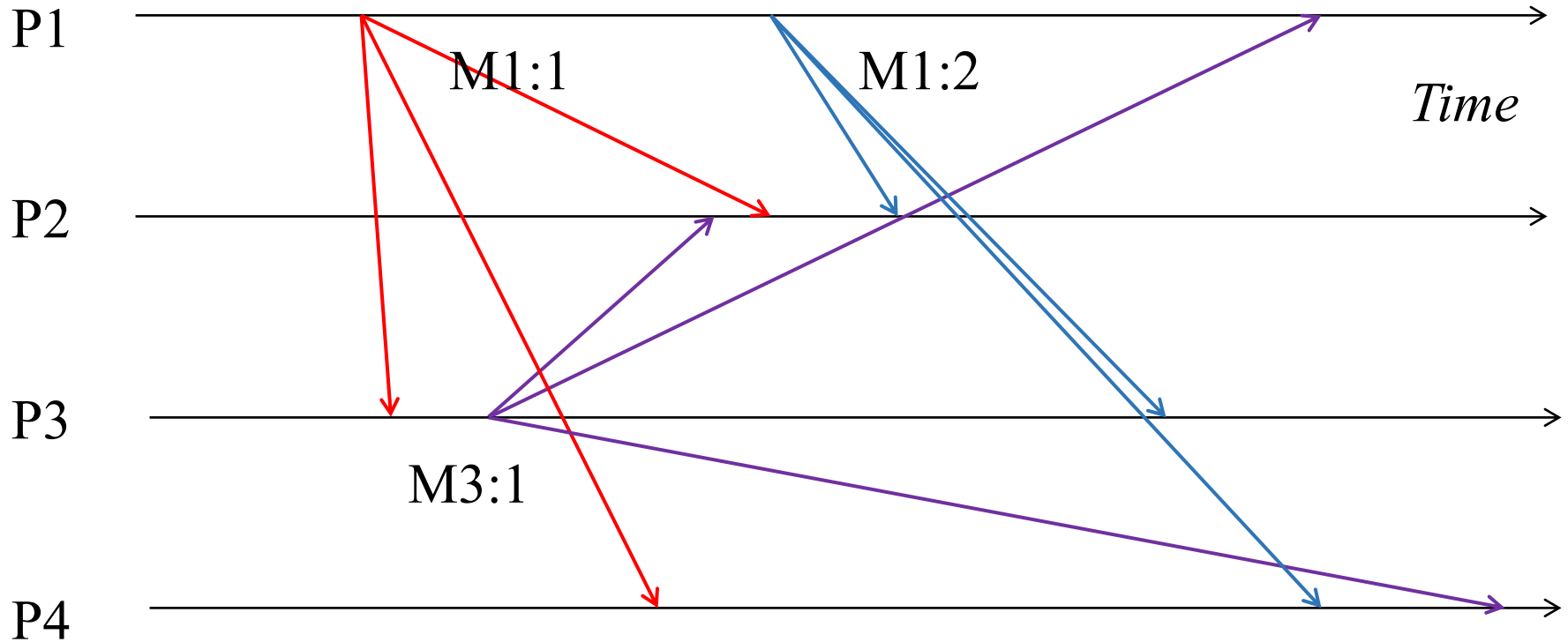
Example



Does this satisfy FIFO order?

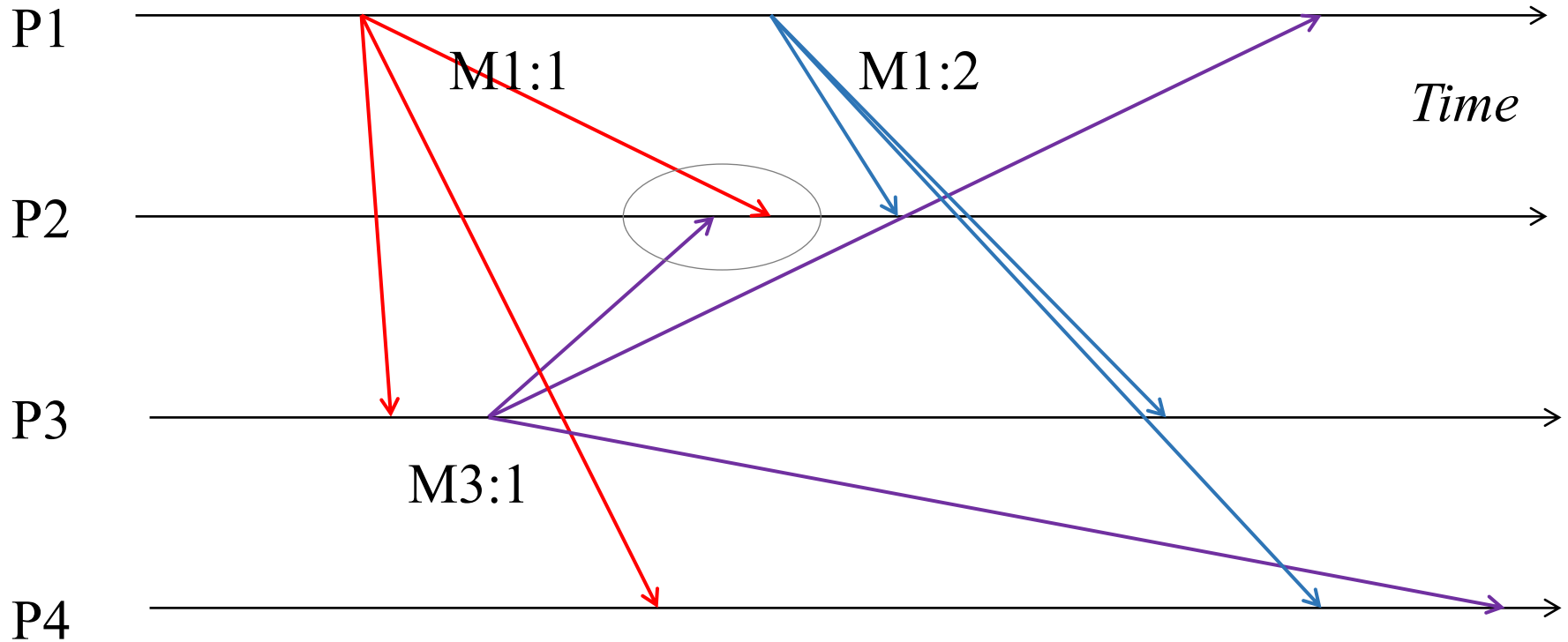
Yes

Example



Does this satisfy causal order?

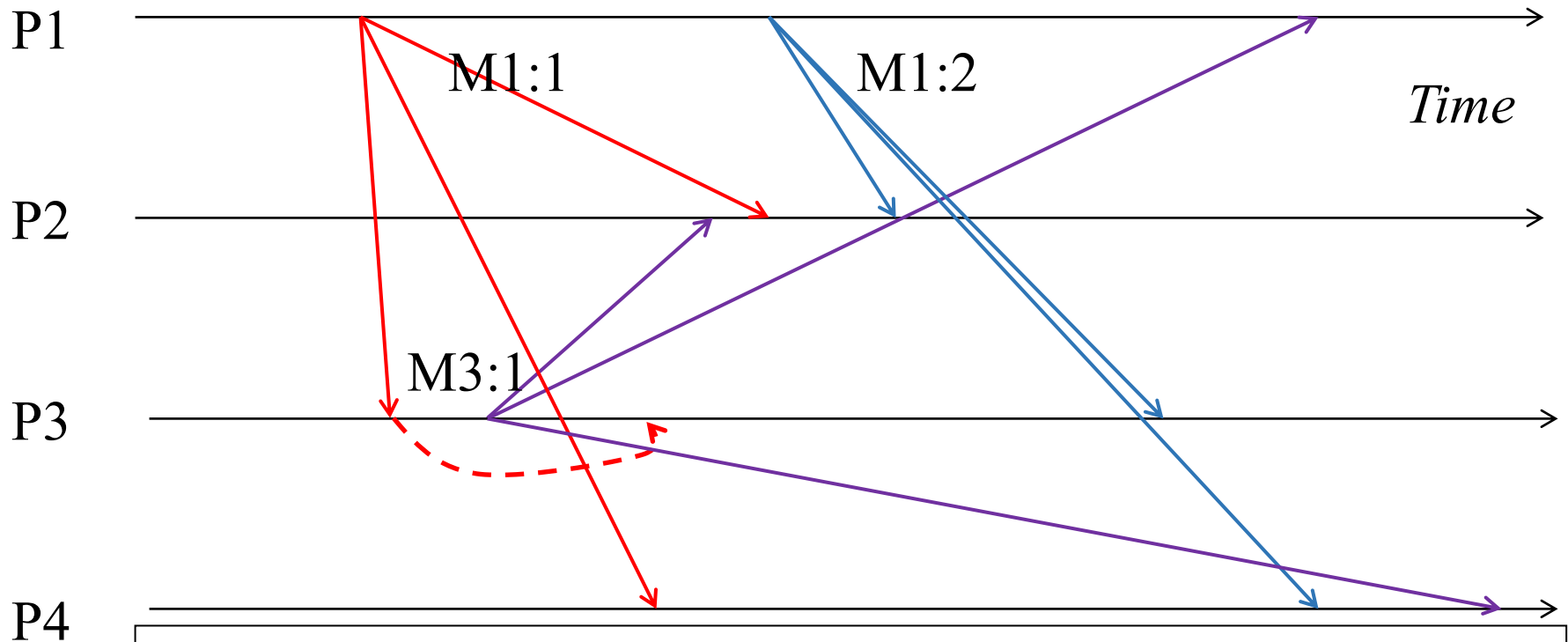
Example



Does this satisfy causal order?

No

Example

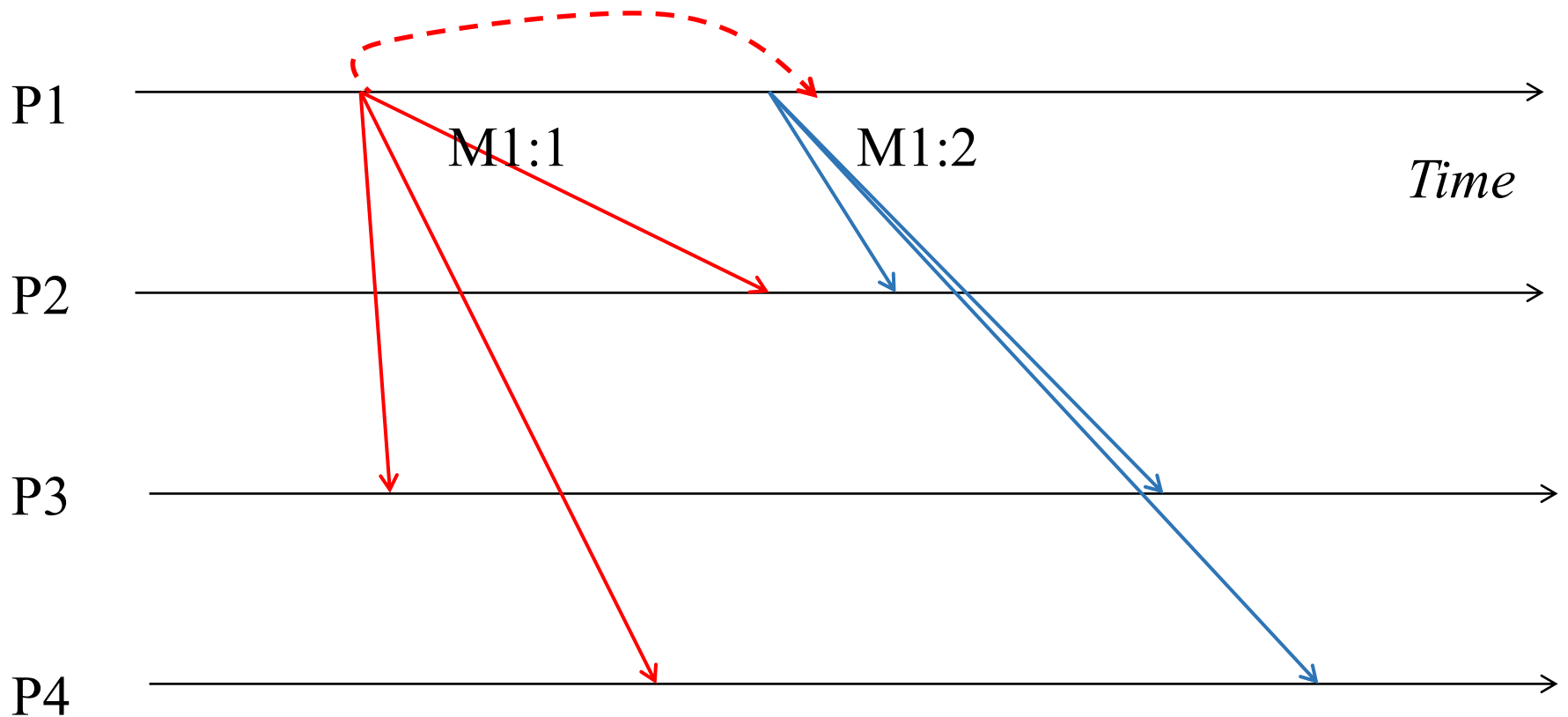


M1:1 is delivered at P3 after M3:1's multicast.

Does this satisfy causal order?

Yes

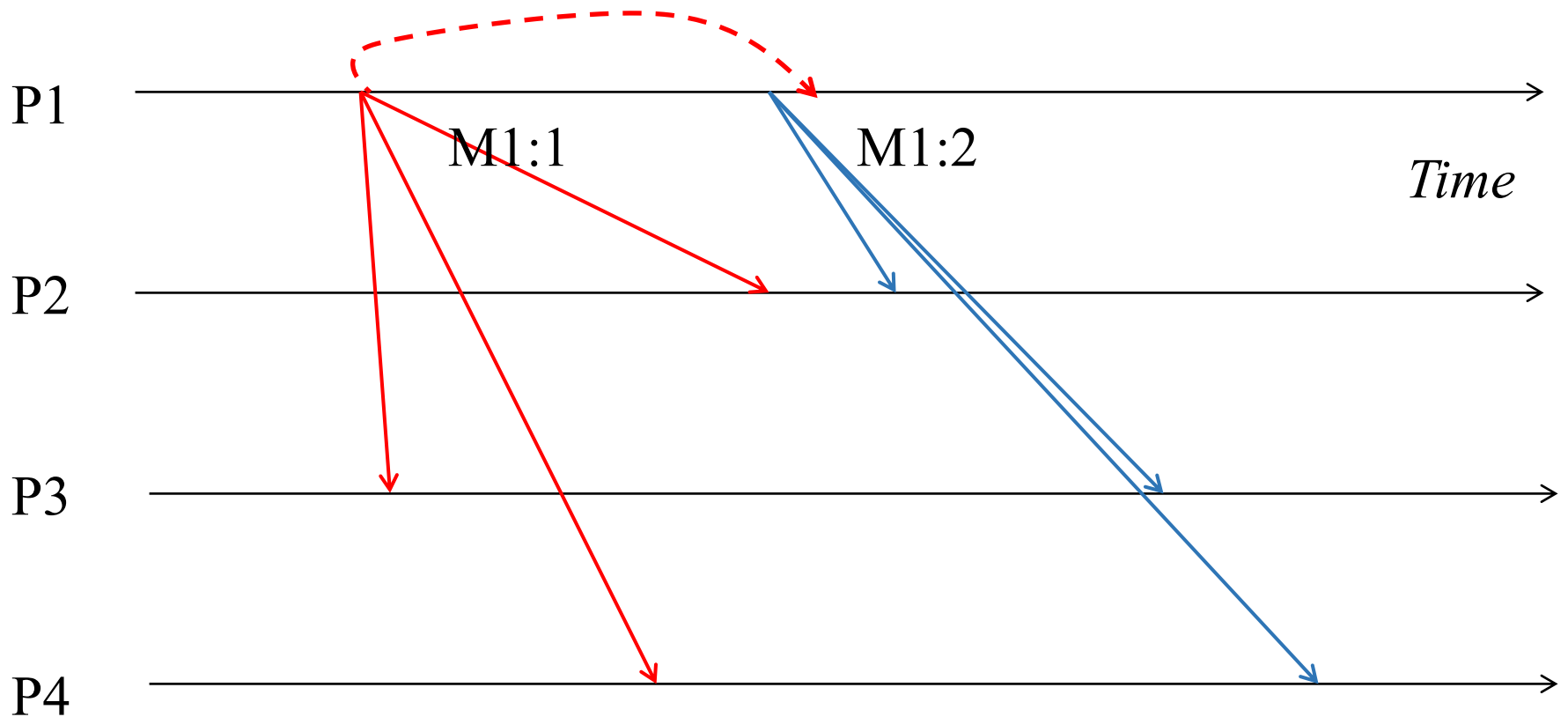
Example



Does this satisfy causal order?

No

Example



Does this satisfy FIFO order?

No

3. Total Order

Next class!

Also in next class,

How do we implement ordered multicast?