

Distributed Systems

CS425/ECE428

02/14/2020

Today's agenda

- **Multicast (contd.)**

- *Chapter 15.4*
- Implementing ordered multicast.

- **Acknowledgement:**

- Materials derived from Prof. Indy Gupta, Prof. Nitin Vaidya, and Prof. Nikita Borisov.

Logistics

- Midterm on March 2nd 7-9pm.
 - Please let us know of any conflicts by Monday.
- HW2 will be released tonight.
 - Due on Feb 27th.
- Still have people who do not have CampusWire access!
 - Please email the instructors and make sure you have access.

Recap: Multicast

- Useful communication mode in distributed systems:
 - Writing an object across replica servers.
 - Group messaging.
 -
- Basic multicast (B-multicast): unicast send to each process in the group.
 - Does not guarantee consistent message delivery if sender fails.
- Reliable multicast (R-multicast):
 - Defined by three properties: *integrity, validity, agreement*.
 - If some correct process multicasts a message **m**, then all other correct processes deliver the **m** (exactly once).
 - When a process receives a message 'm' for the first time, it re-multicasts it again to other processes in the group.

Recap: 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 .

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 23 and 24, does that satisfy FIFO order?

Yes

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 24 and 25, does that satisfy FIFO order?

Yes

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 24 and 25, does that satisfy causal order?

No

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 23 and 24 for one process displaying the bulletin and not for another, does that satisfy FIFO order?

Yes

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 23 and 24 for one process displaying the bulletin and not for another, does that satisfy total order?

No

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 24 and 25 for all processes displaying the bulletin does that satisfy causal order?

No

Multicast Ordering Example

Online bulletin board		
Item	From	Subject
23	A.Hanlon	Mach
24	G.Joseph	Microkernels
25	A.Hanlon	Re: Microkernels
26	T.L'Heureux	RPC performance
27	M.Walker	Re: Mach
end		

If we swap items 24 and 25 for all processes displaying the bulletin does that satisfy total order?

Yes

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 .

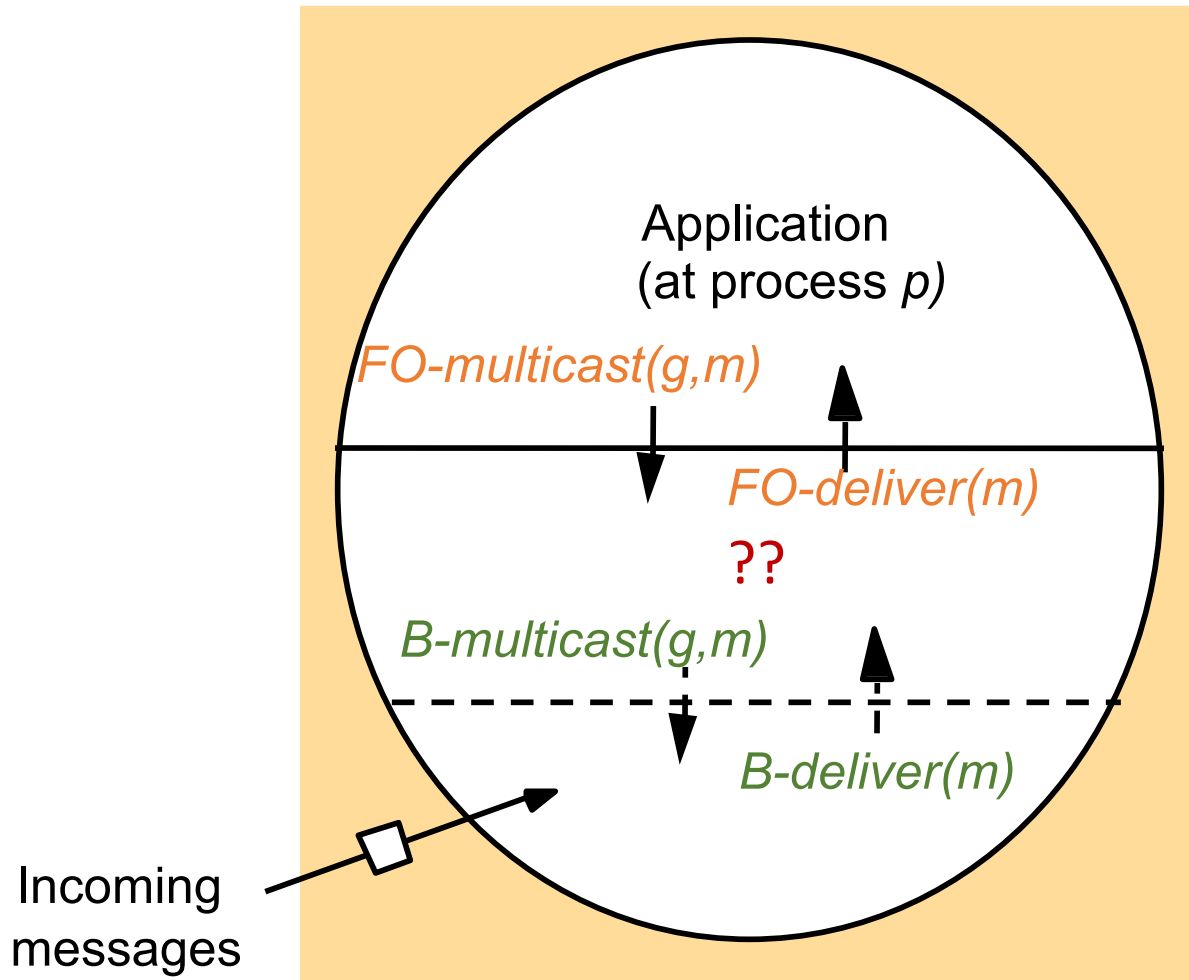
- **Causal ordering**

- 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



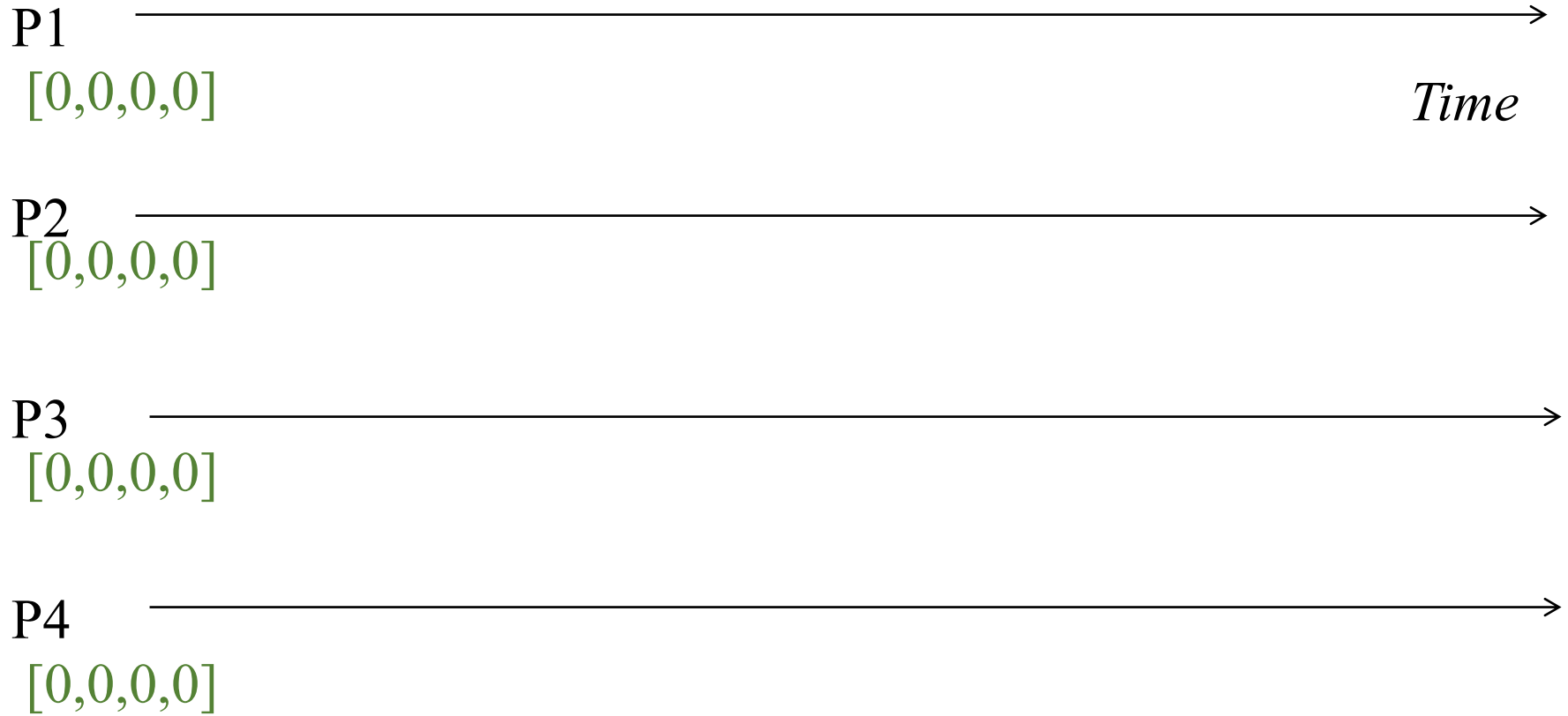
Implementing FIFO order multicast

- 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]

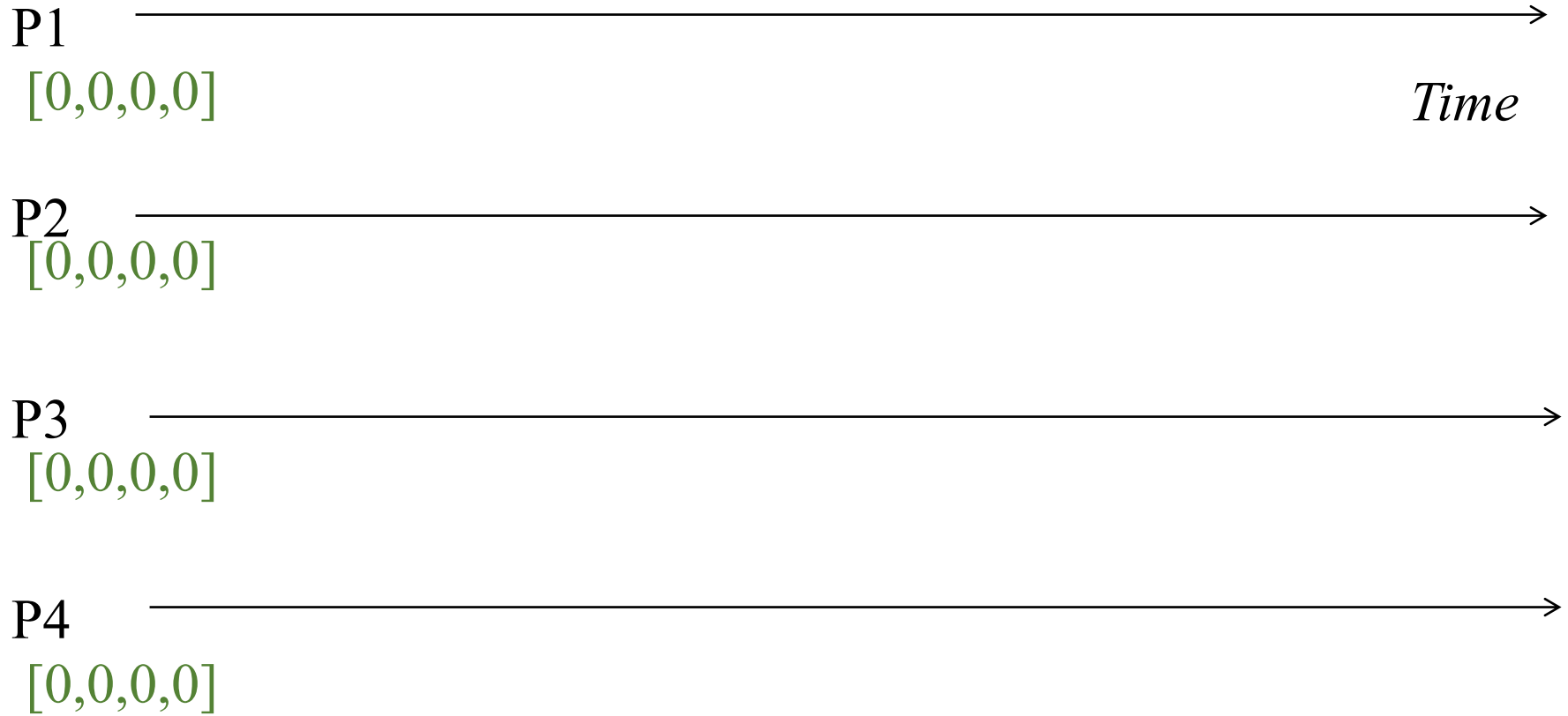
Sequence Vector

Do not confuse with vector timestamps!

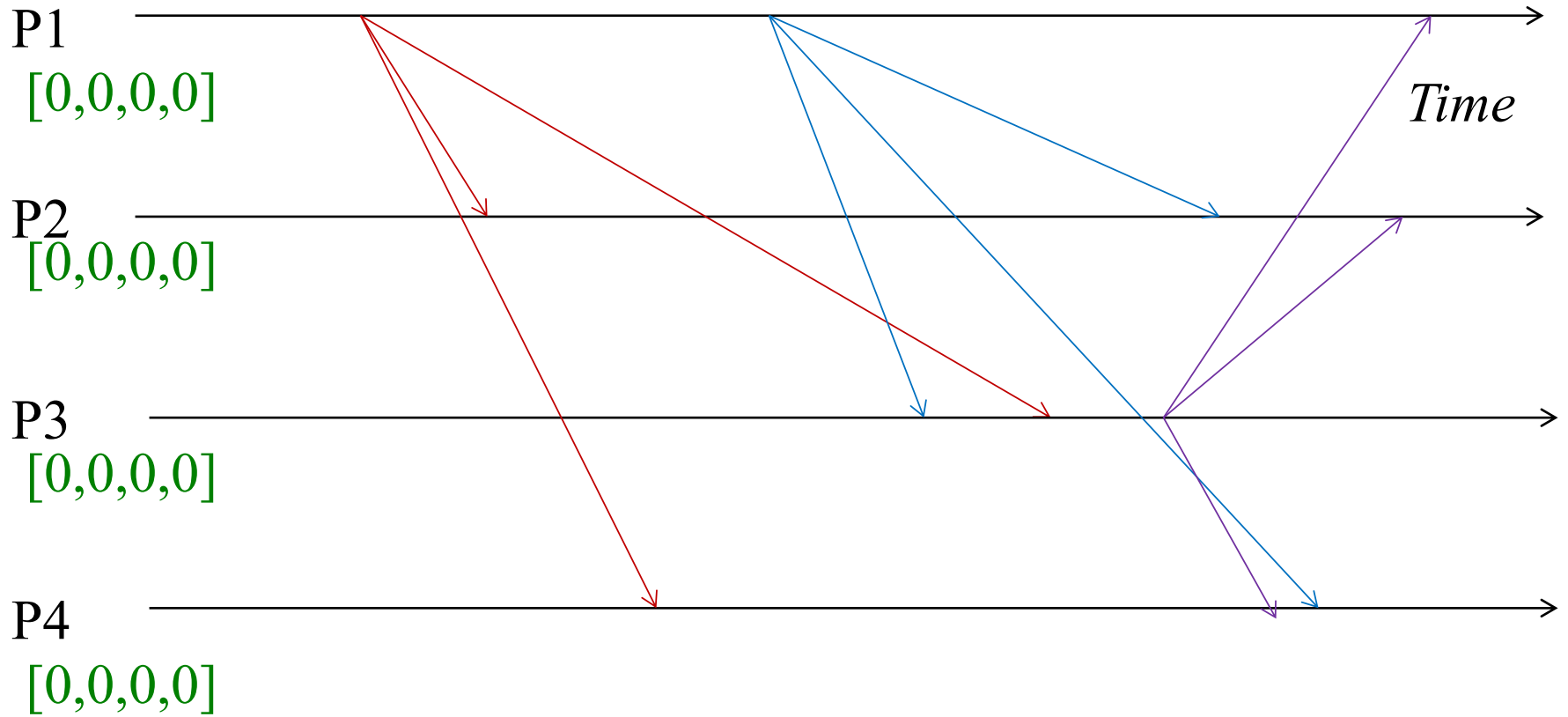
$P_i[i]$, is the no. of messages P_i multicast (and delivered to itself).

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

FIFO order multicast execution

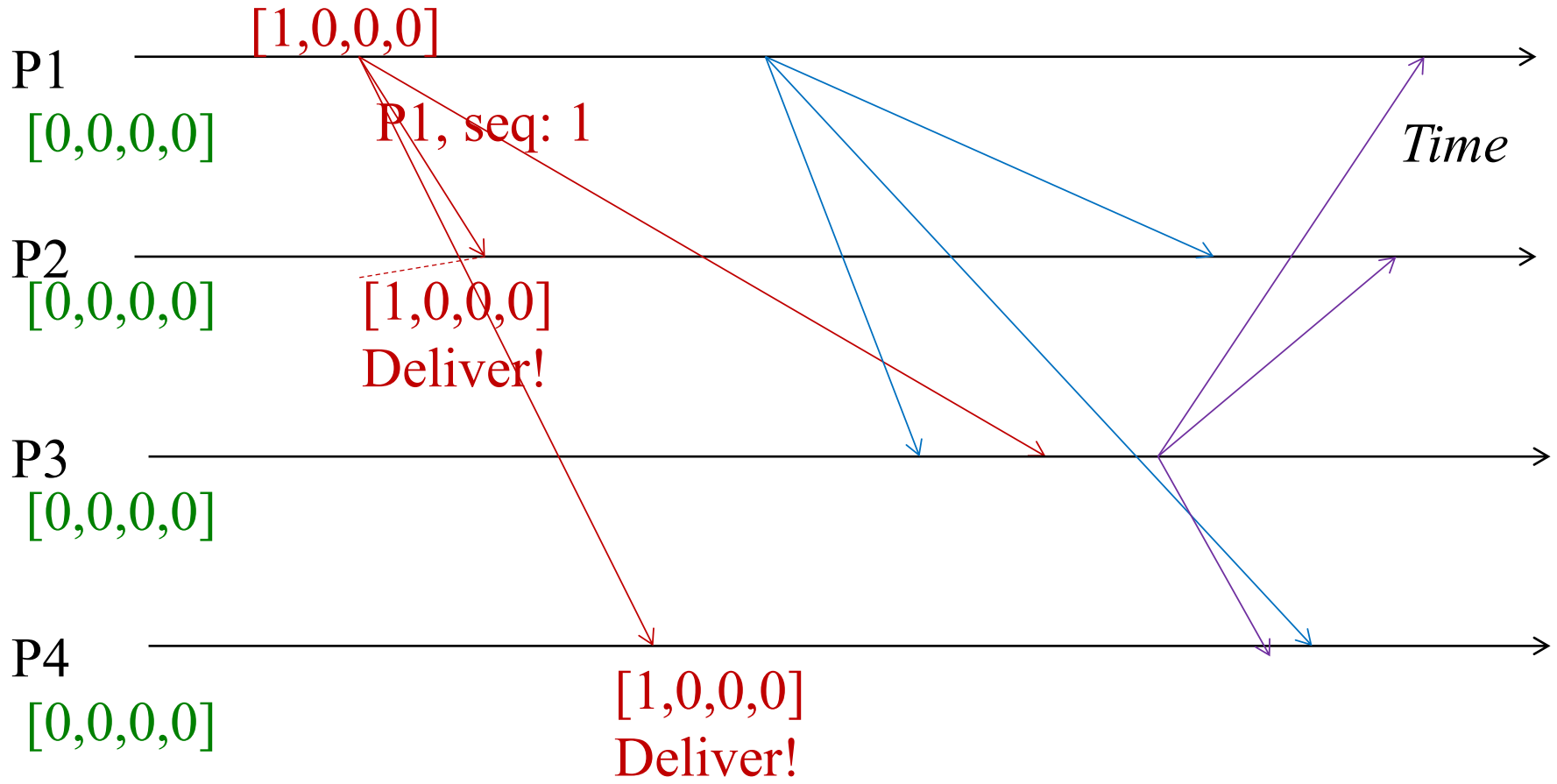


FIFO order multicast execution

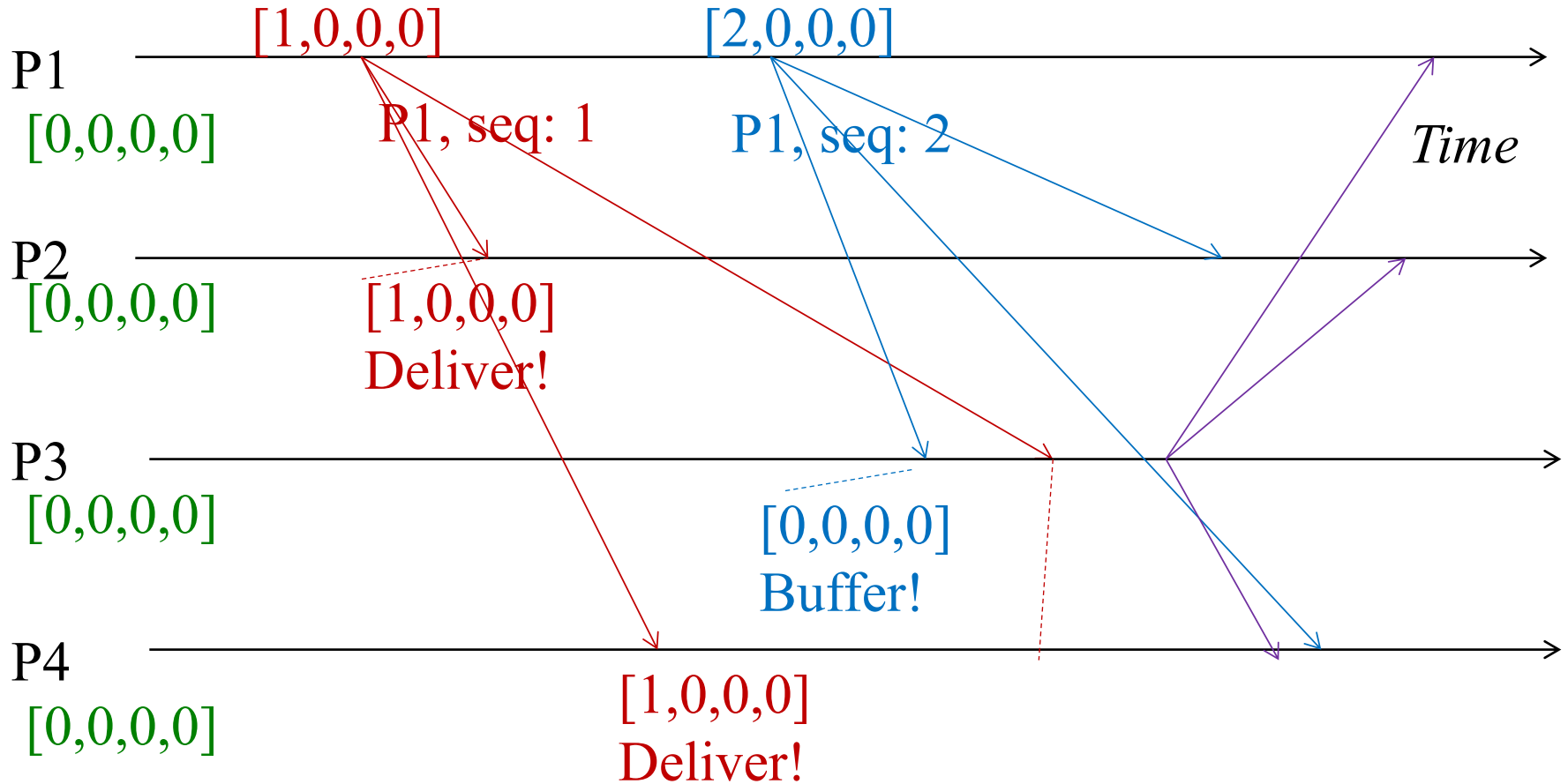


Self-deliveries omitted for simplicity.

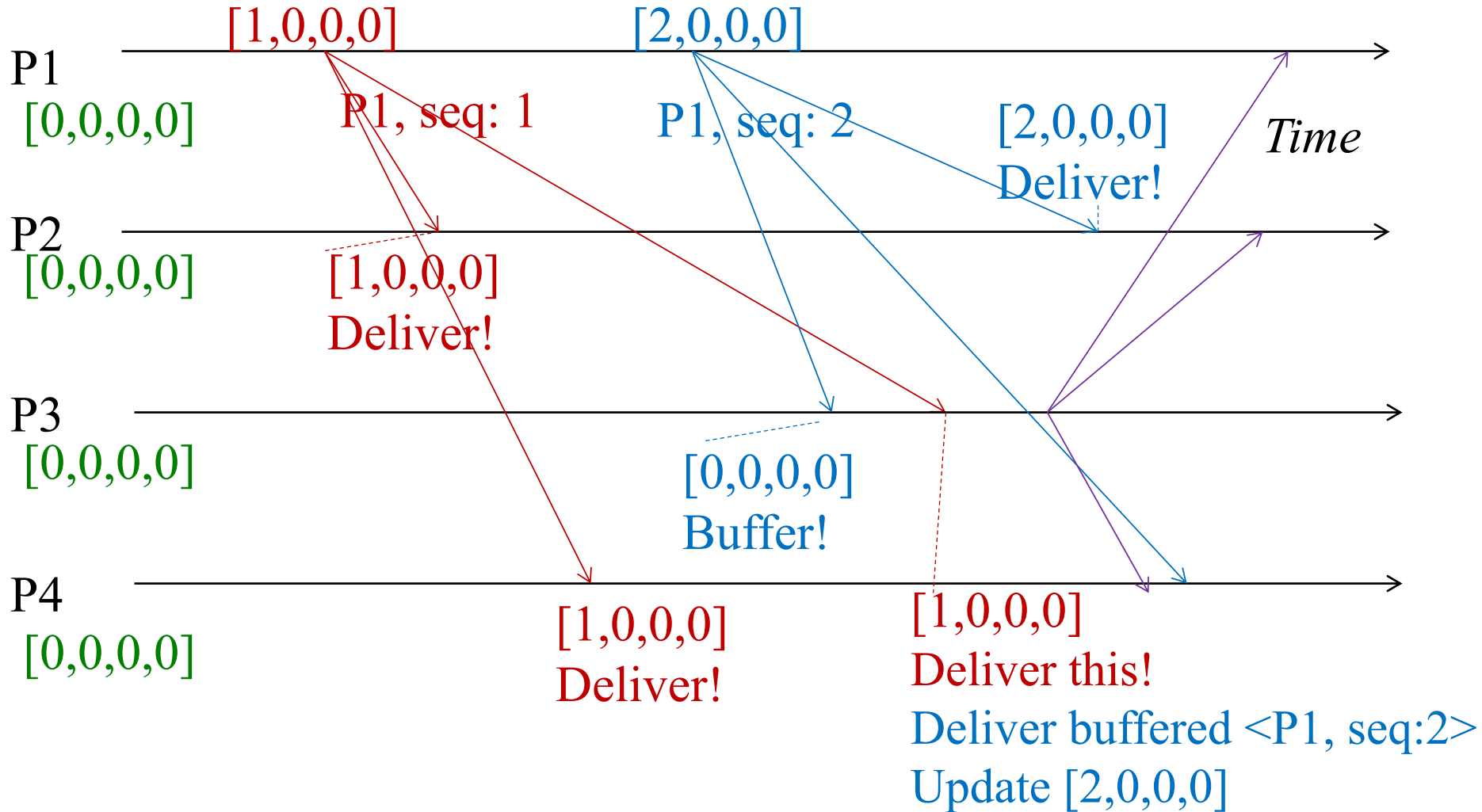
FIFO order multicast execution



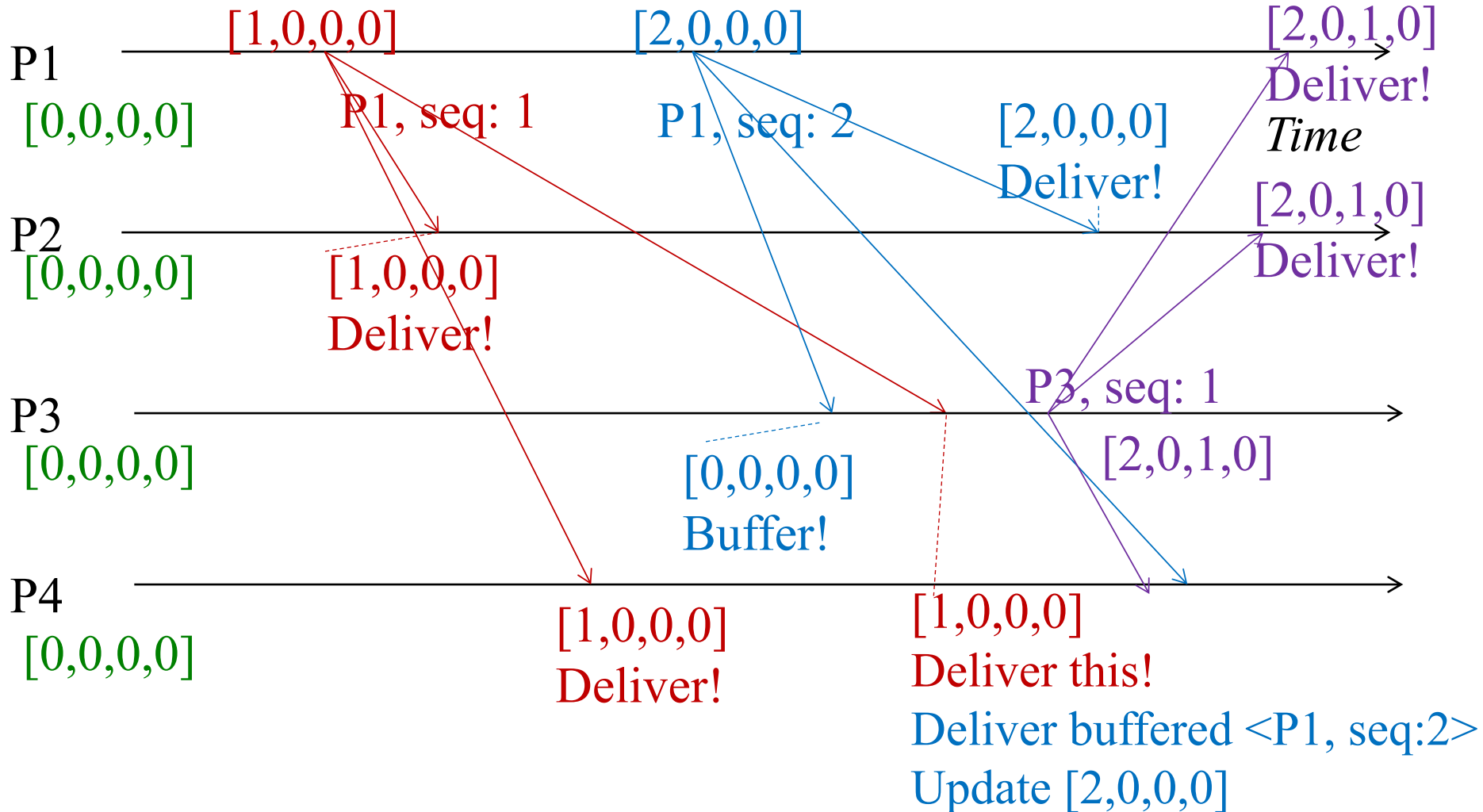
FIFO order multicast execution



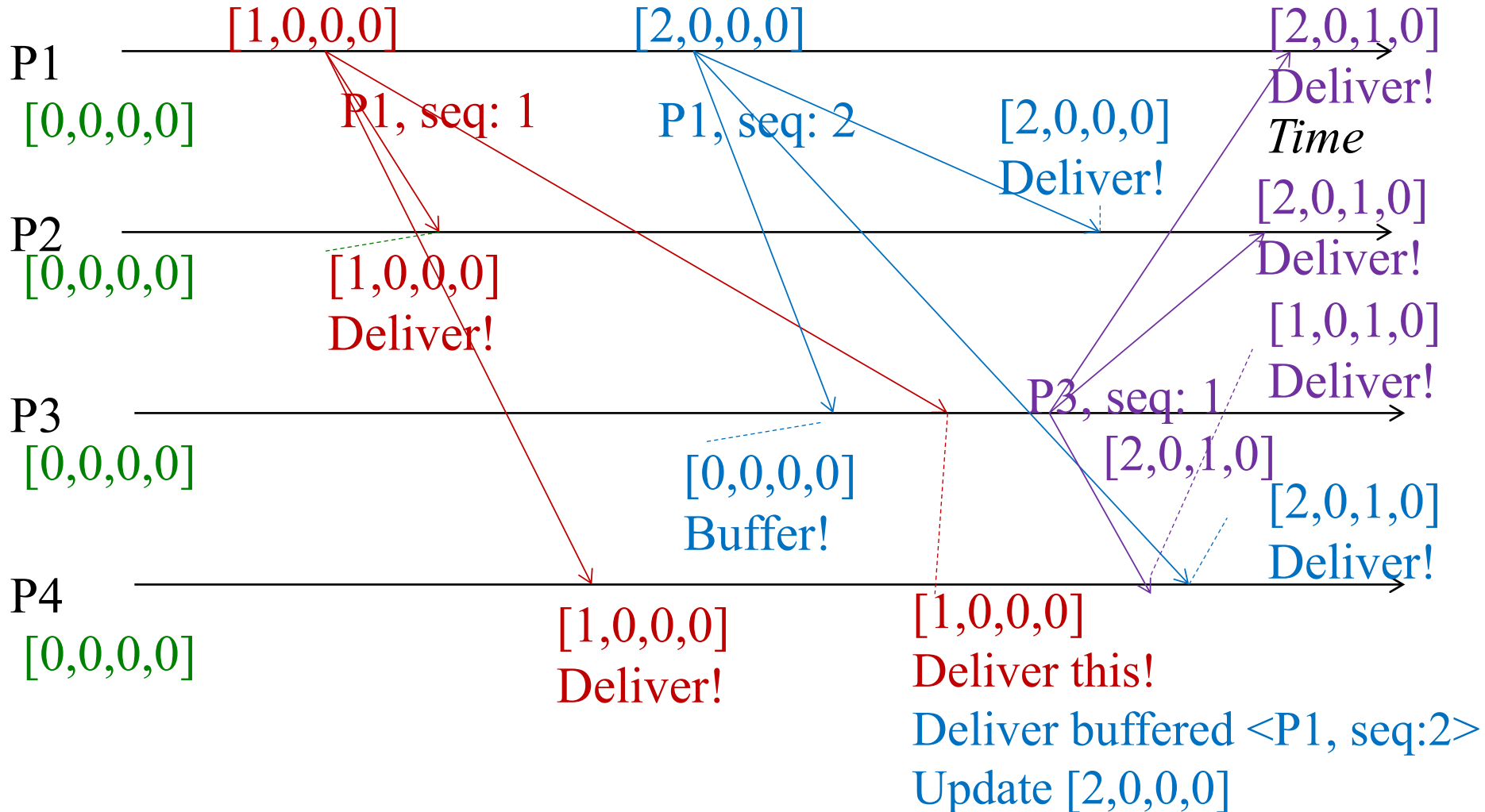
FIFO order multicast execution



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 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)

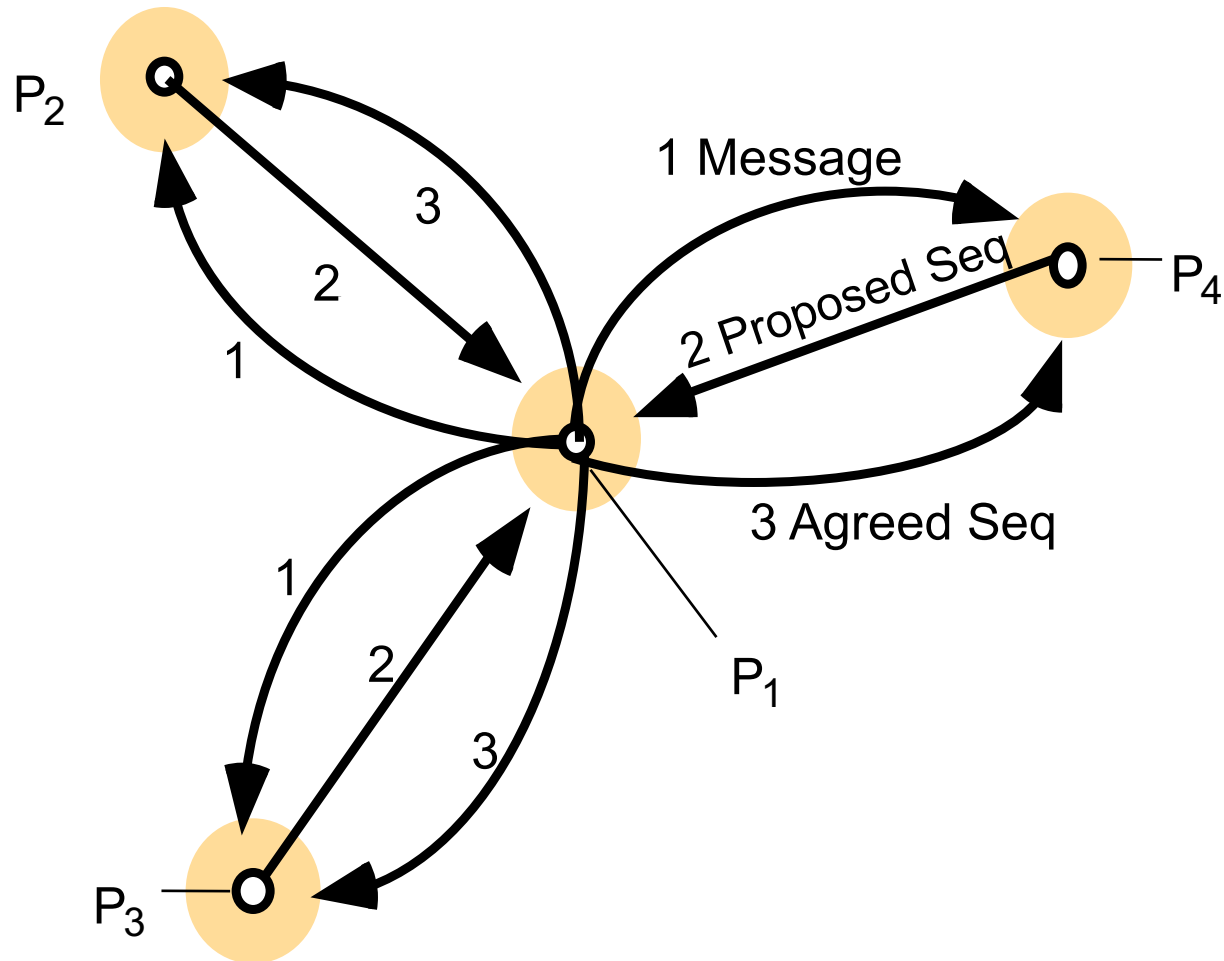
Sequencer based total ordering

- 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 $S_i = S_i + 1$

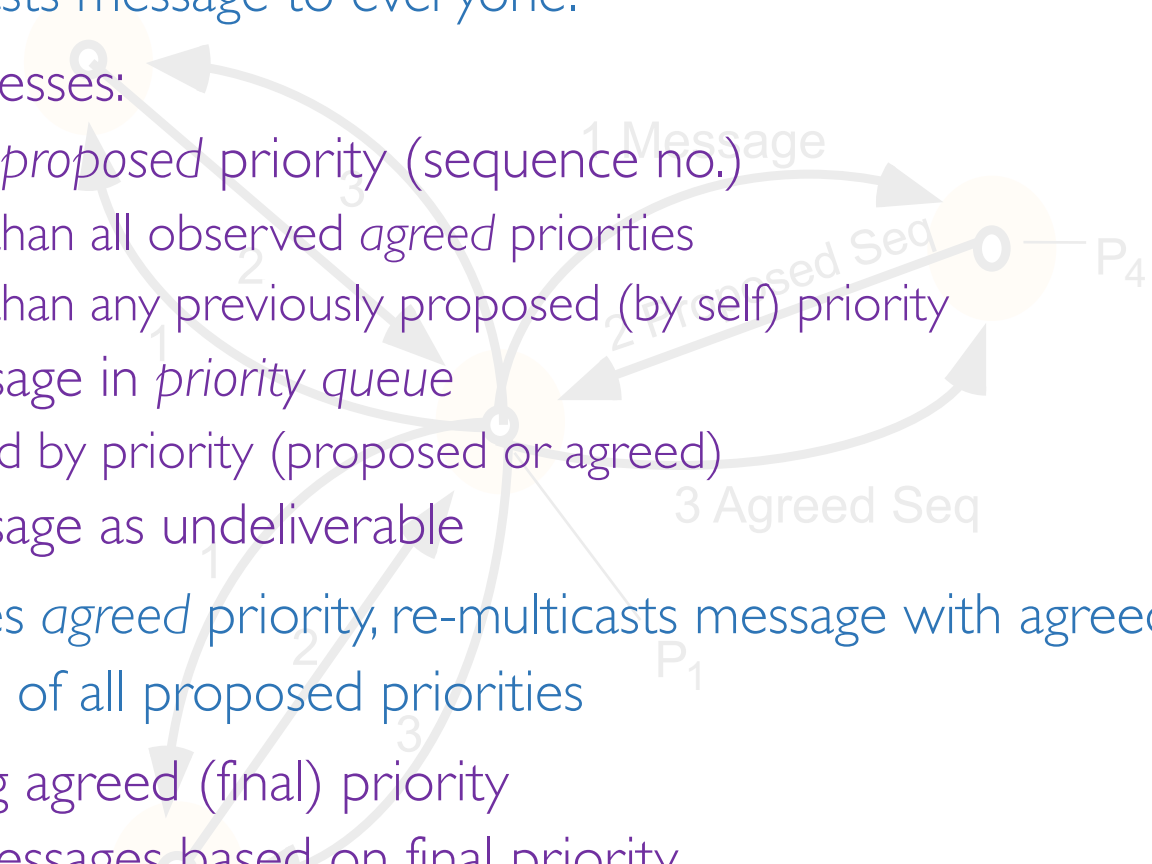
Implementing total order multicast

- 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)**

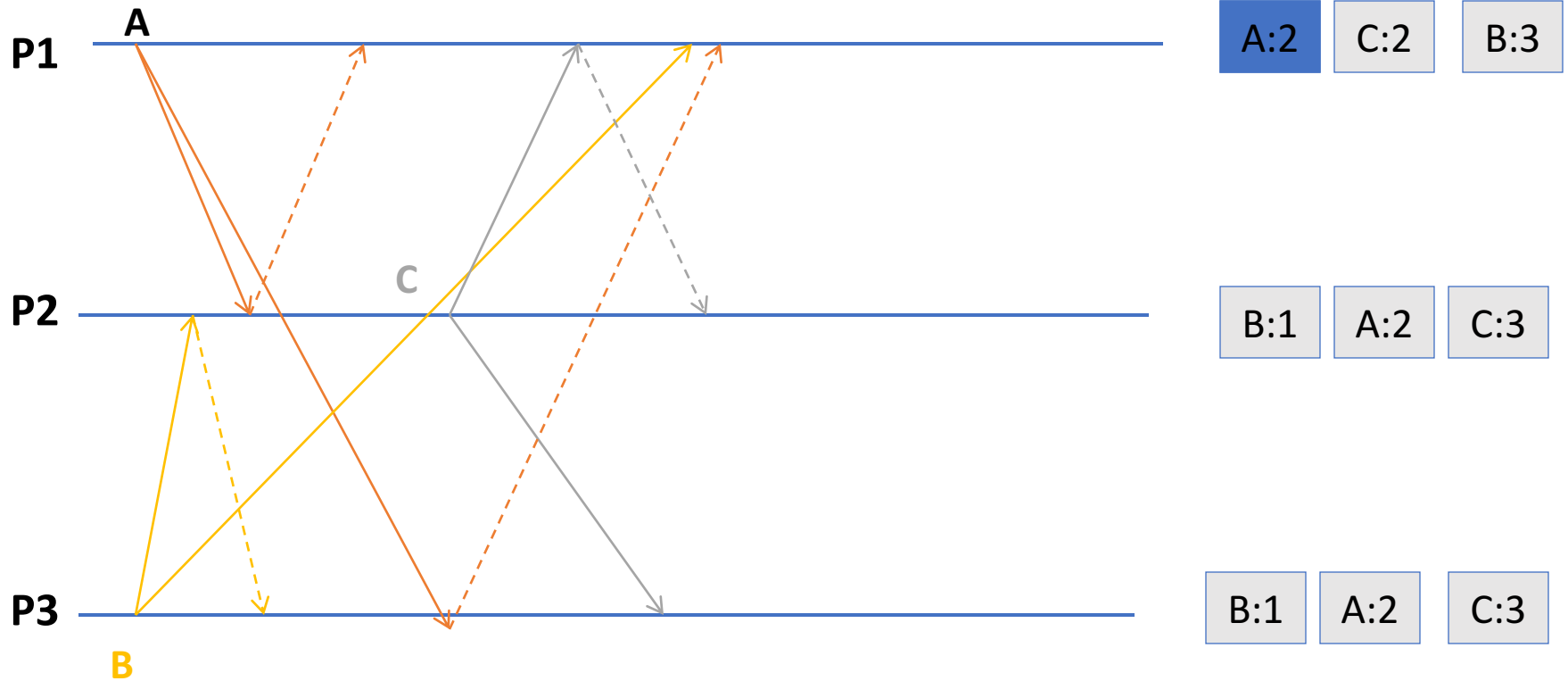
ISIS algorithm for total ordering



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.
- 

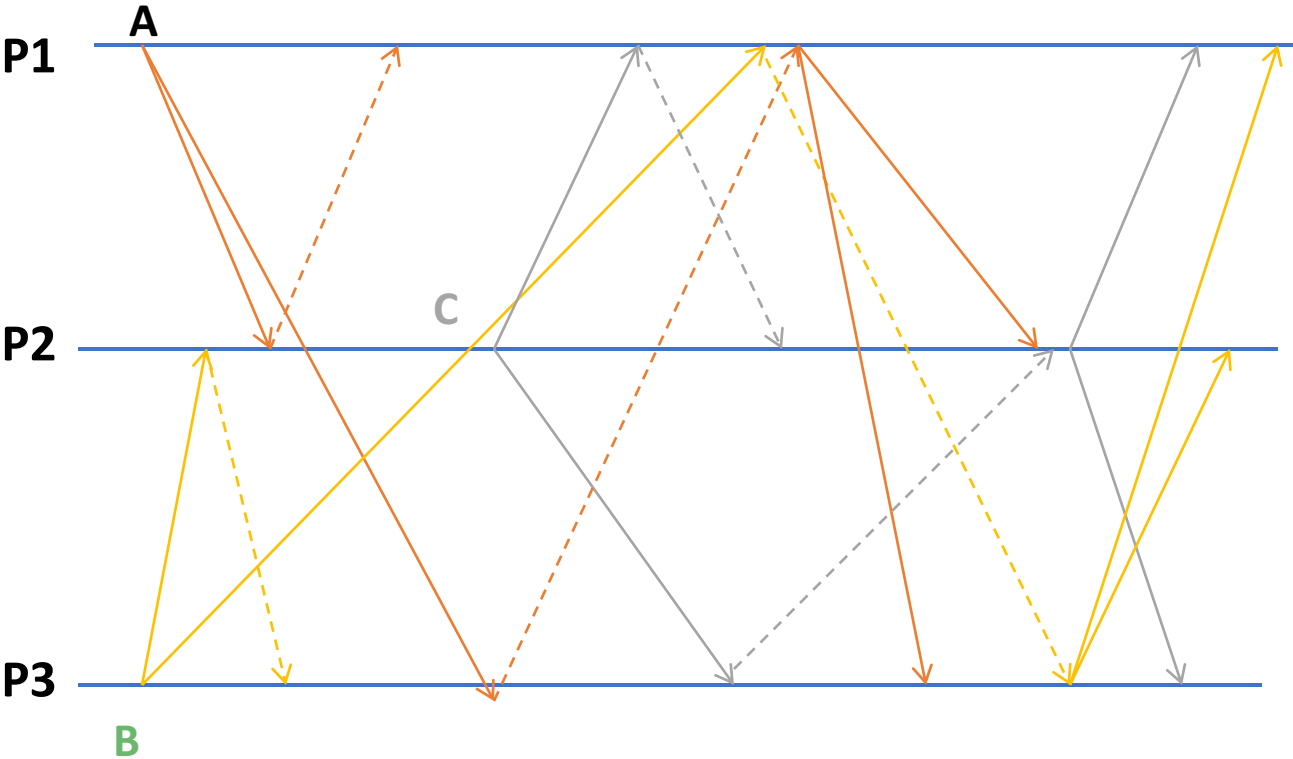
Example: ISIS algorithm



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 .

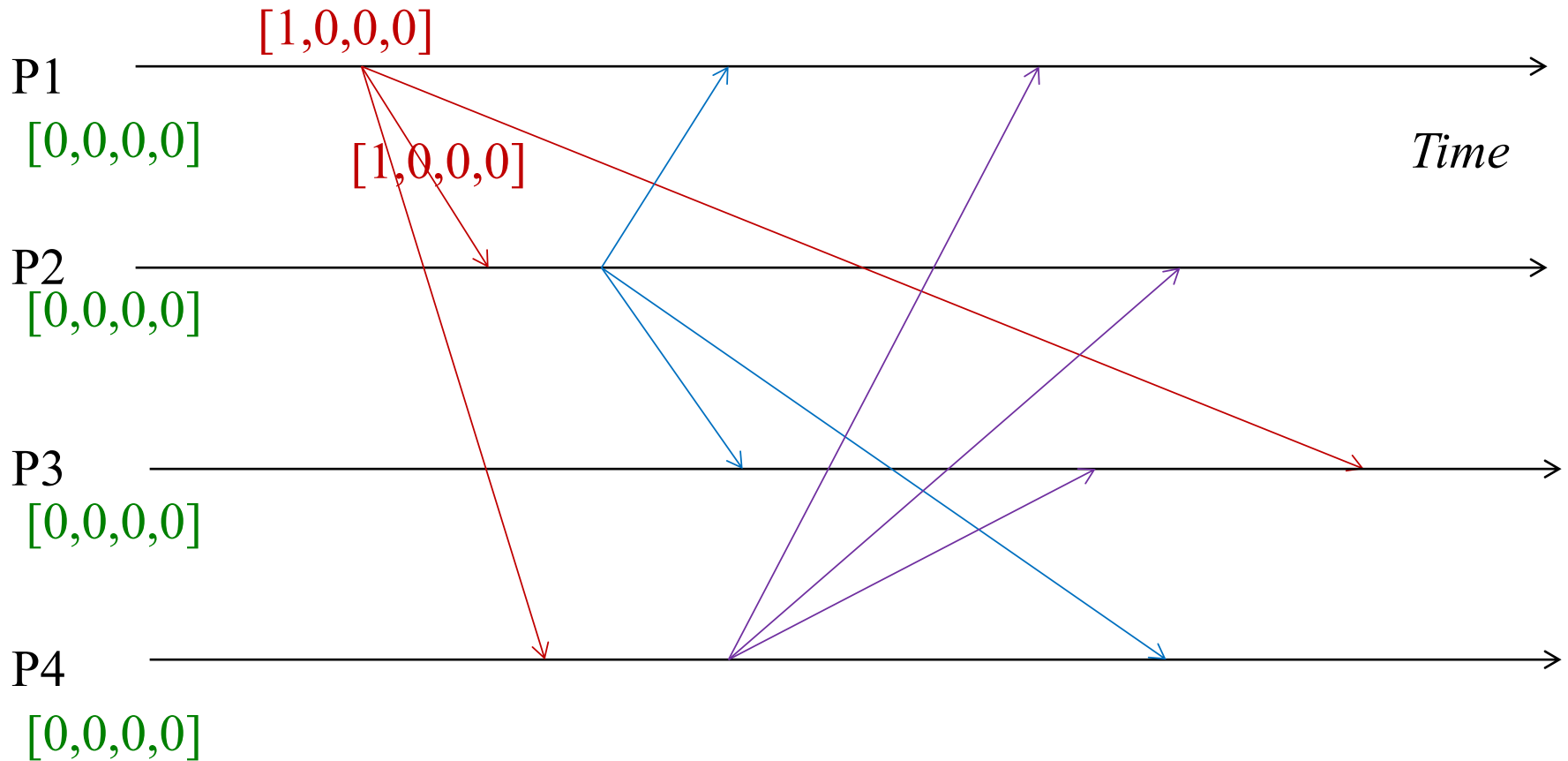
Implementing causal order multicast

- 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 .

Implementing causal order multicast

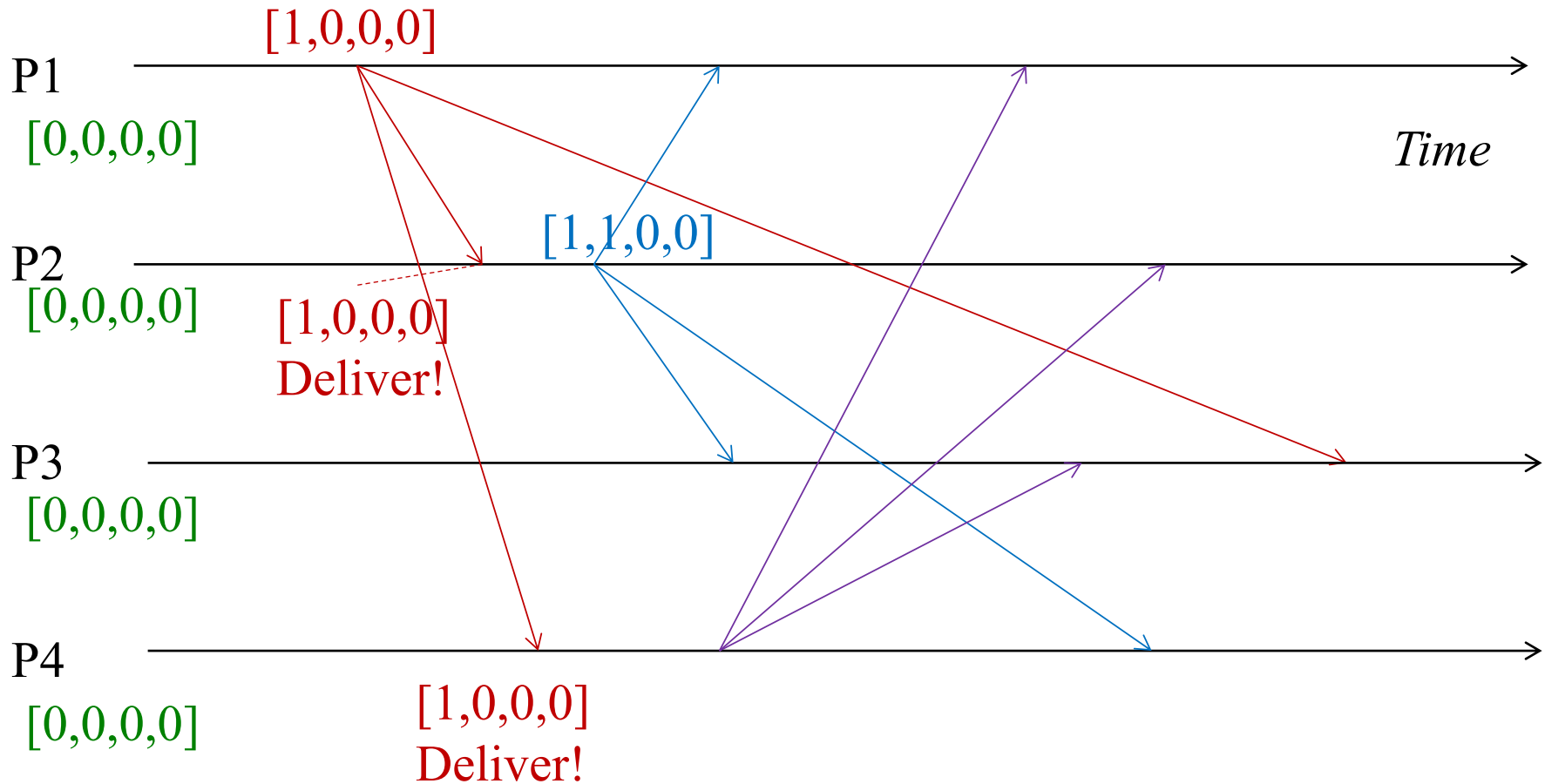
- *CO-multicast*(g, m) at P_j :
 - set $P_j[j] = P_j[j] + 1$
 - piggyback entire vector $P_j[1 \dots N]$ with m as its sequence no.
 - B-multicast($g, \{m, P_j[1 \dots N]\}$)
- On B-deliver($\{m, V[1 \dots N]\}$) at P_i from P_j : If P_i receives a multicast from P_j with sequence vector $V[1 \dots N]$, 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

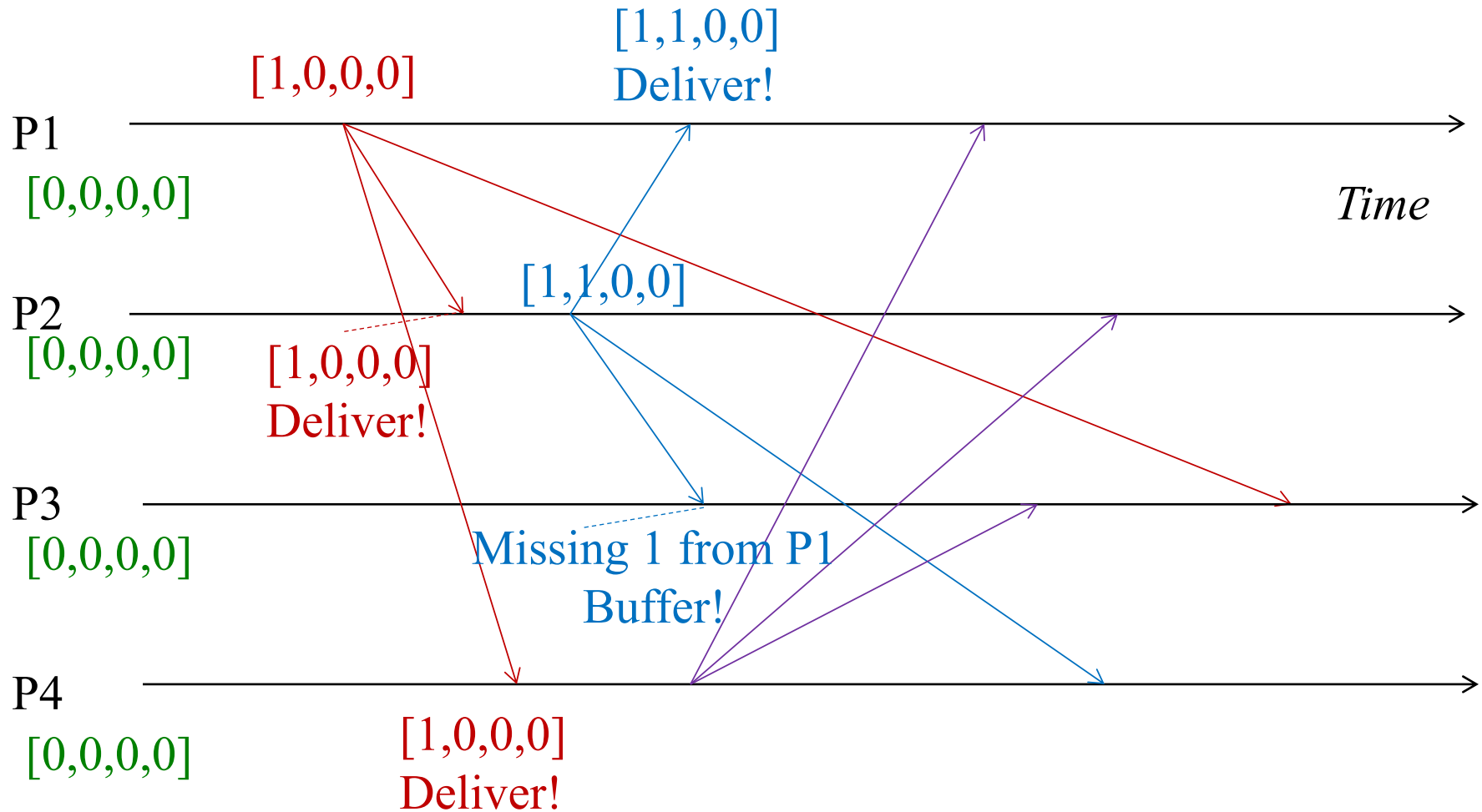


Self-deliveries omitted for simplicity.

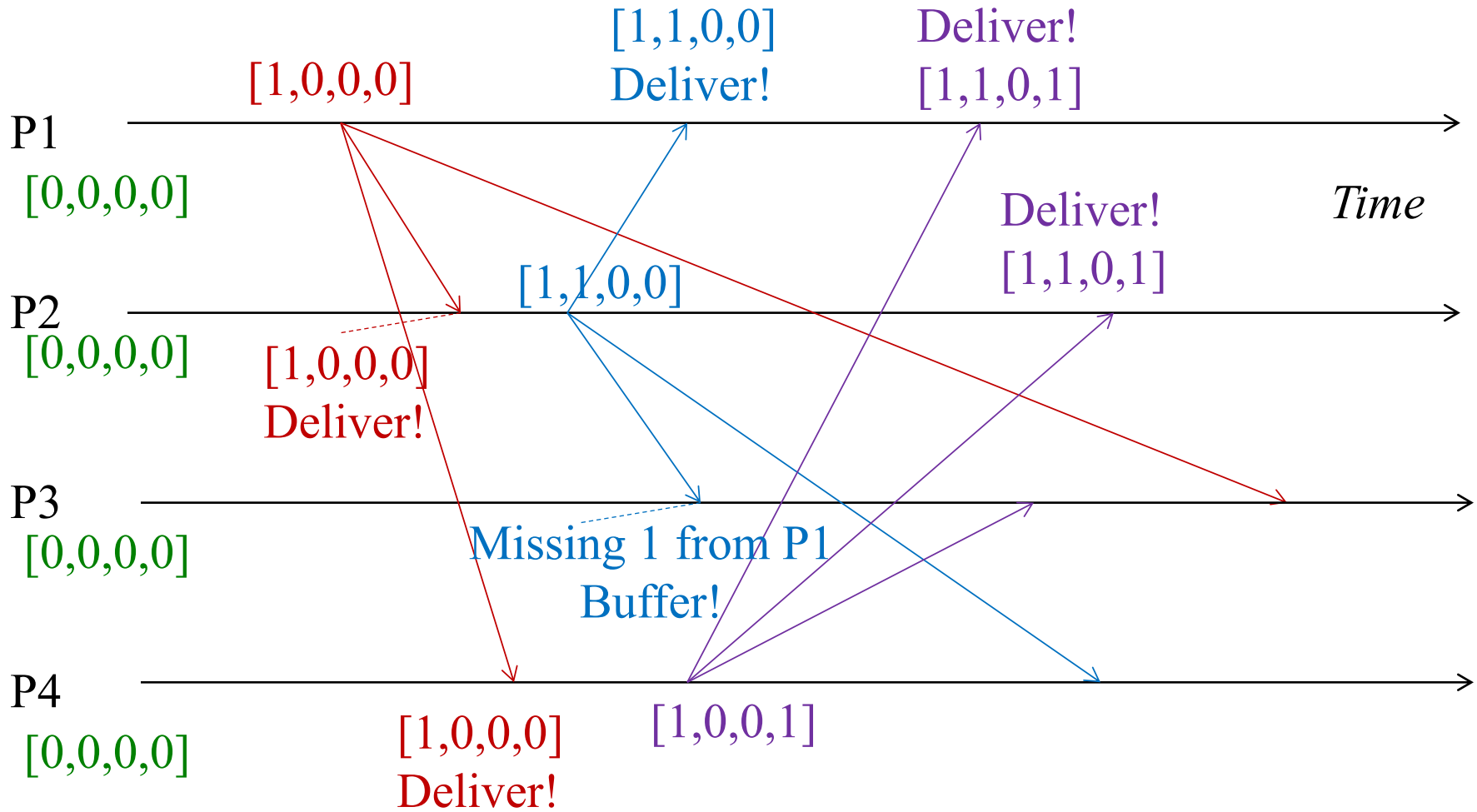
Causal order multicast execution



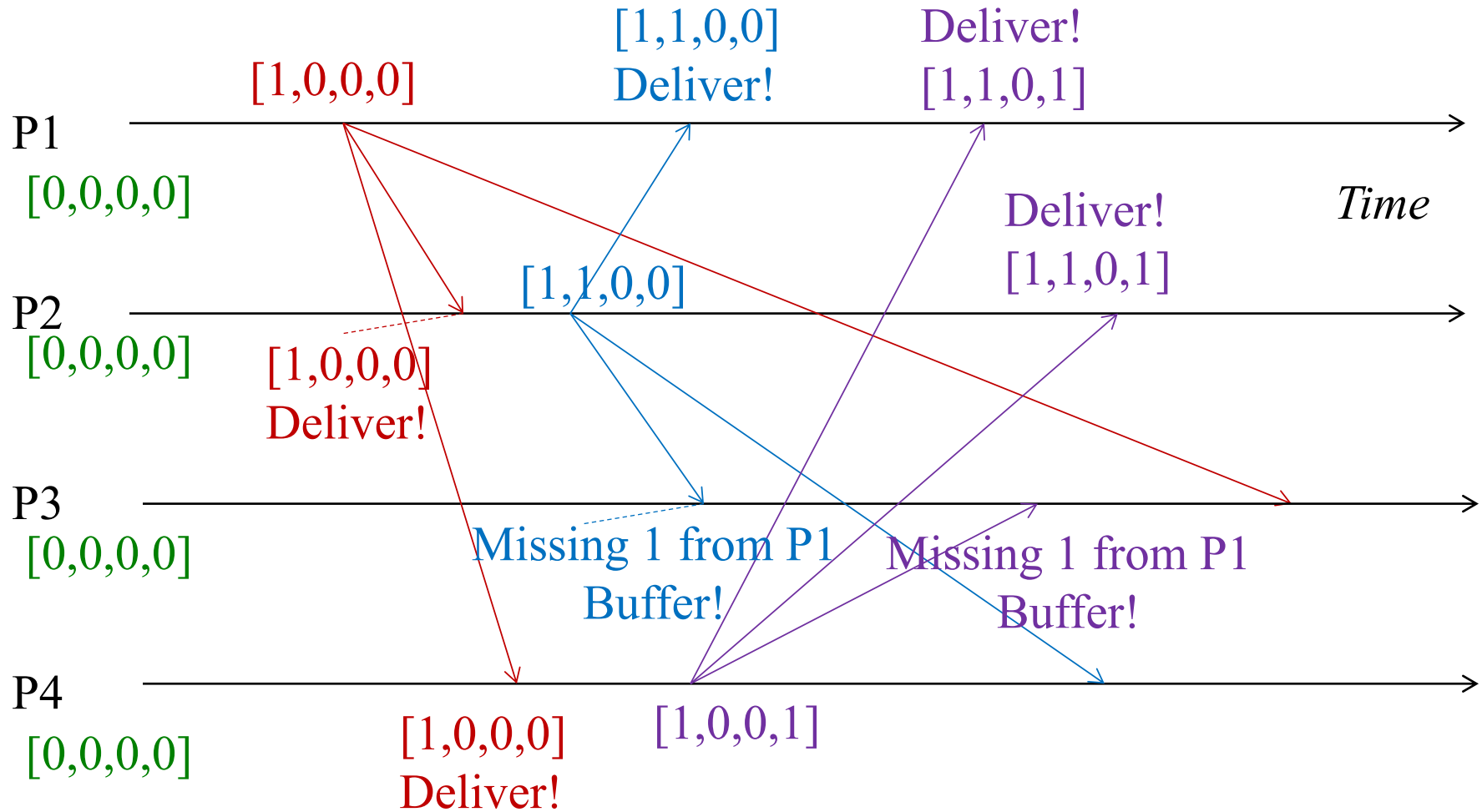
Causal order multicast execution



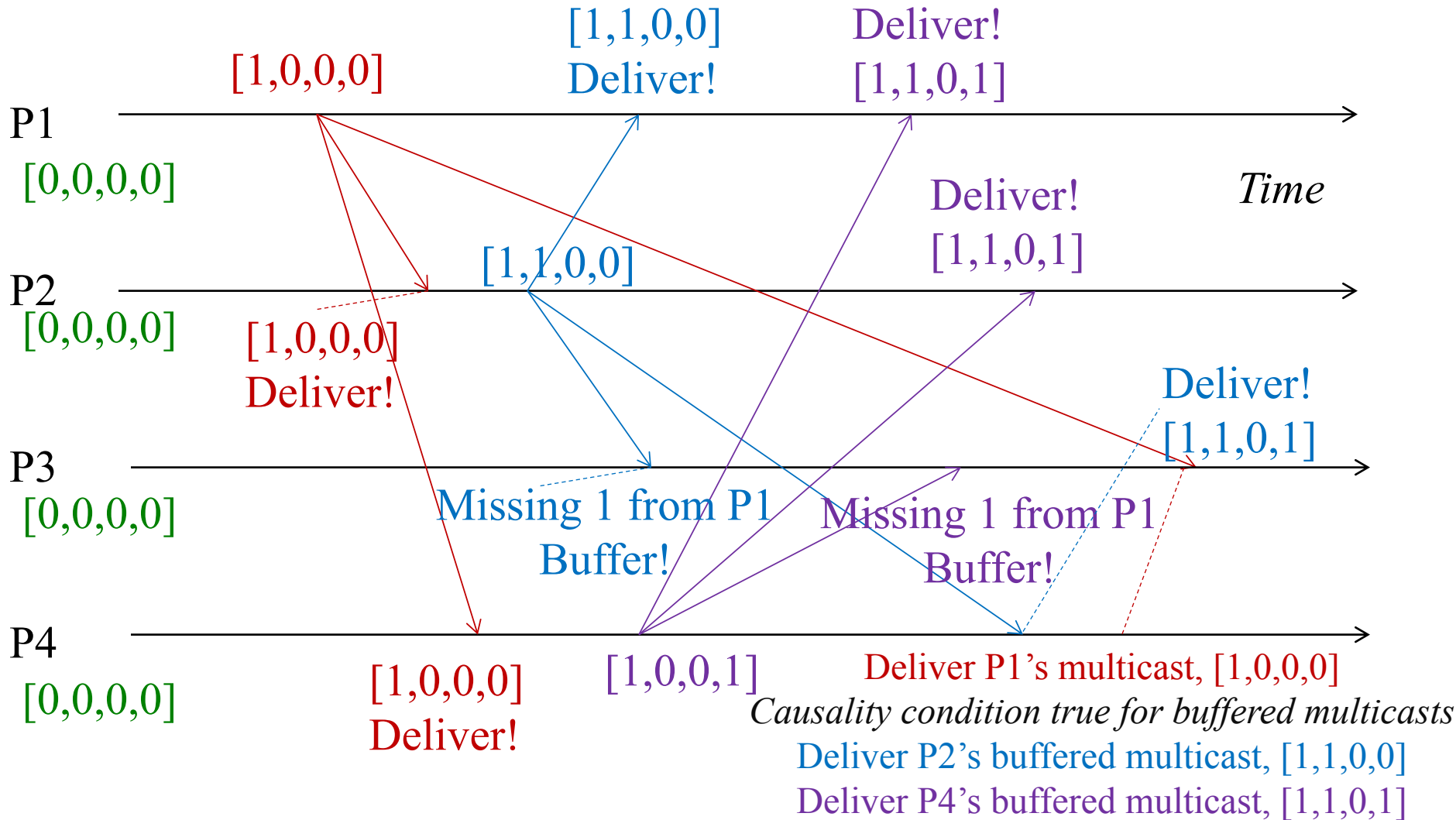
Causal order multicast execution



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 .

Summary

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