

Vector Clocks

Lec8

Problems with Lamport Clocks

- Lamport timestamps do not capture **causality**.
- With Lamport's clocks, one cannot directly compare the timestamps of two events to determine their precedence relationship.
- The main problem is that a simple integer clock cannot order both events within a process and events in different processes.
 - Knowing that $C(a) < C(b)$ is true does not allow us to conclude that $a \rightarrow b$ is true.

Vector Clocks

- Happened-before relation is a partial order
 - But the domain of Lamport clocks (natural numbers) is a total order (wrt $<$)
 - Do not provide complete information about the happened-before relation
- Key idea: assign vector timestamps to events
 - $s.v$ refers to the vector timestamp assigned to state s
- For all s, t : $s \rightarrow t$ IFF $s.v < t.v$
- Want the timestamps also to be partial order
- For two vectors x, y of dimension N :

- $x < y =$

$$\forall k (1 \leq k \leq N): x[k] \leq y[k] \quad \text{and}$$
$$\exists j (1 \leq j \leq N): x[j] < y[j]$$

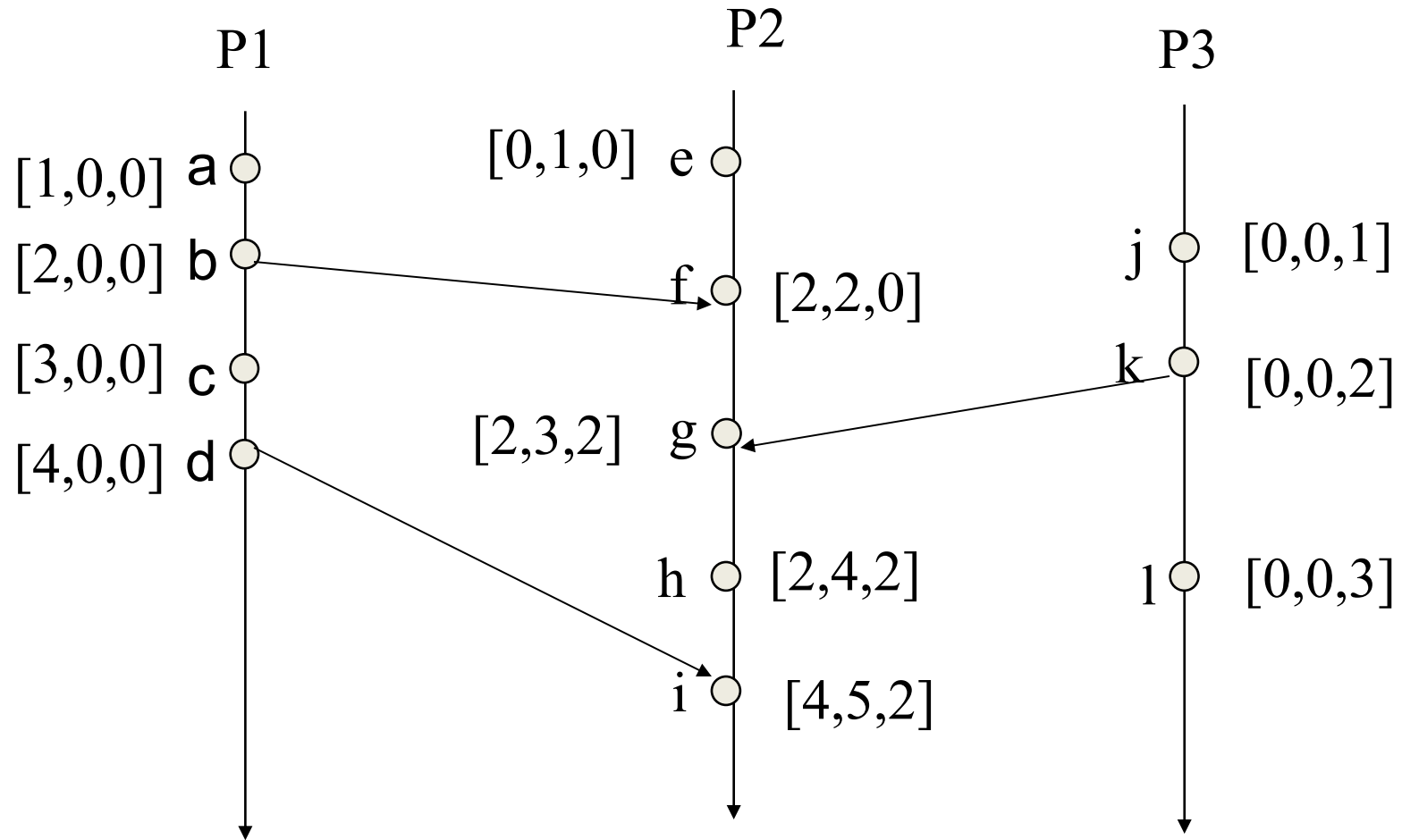
Properties of Vector Timestamps

- $v_i[i]$ is the number of events that have occurred so far at P_i
- If $v_i[j] = k$ then P_i knows that k events have occurred at P_j

Vector Clock Update Rules

1. Initially all clock values are set to the smallest value (e.g., 0).
 2. The local clock value is incremented at least once before each primitive event in a process i.e., $v_i[i] = v_i[i] + 1$
 3. The current value of the entire logical clock vector is delivered to the receiver for every outgoing message.
 4. Values in the timestamp vectors are never decremented.
 5. Upon receiving a message, the receiver sets the value of each entry in its local timestamp vector to the maximum of the two corresponding values in the local vector and in the remote vector received.
- Let v_q be piggybacked on the message sent by process q to process p ; We then have:
 - For $i = 1$ to num-processes {
 $v_p[i] = \max(v_p[i], v_q[i])$
}
 - $v_p[p] = v_p[p] + 1$;

Vector Clock Example



Vector Clock Comparison

Vectors are compared by each corresponding element

$s = t$ if all elements equal

$s \geq t$ if all elements of $s \geq$ all elements of t , and at least one element of $s > t$

$s > t$ if all elements $s >$ all elements of t

Vector Clock Proof

$s \rightarrow t$ IFF $s.v < t.v$

Vector Clock Refinements

Direct Dependency Clocks

- Simplification of vector clocks.
- Each process still maintains vector timestamps
- But, only sends its logical timestamp component with messages (i.e., process- i sends $v[i]$, instead of entire vector v)
- On receive, process- j updates only two VC components:
 - Its own $v[j]$, is simply incremented
 - The sender's component, $v[i]$, is updated using max of earlier value and value sent in the message
- Direct dependency clocks can capture “directly happened before” relationships
 - Message path contains at most one message.

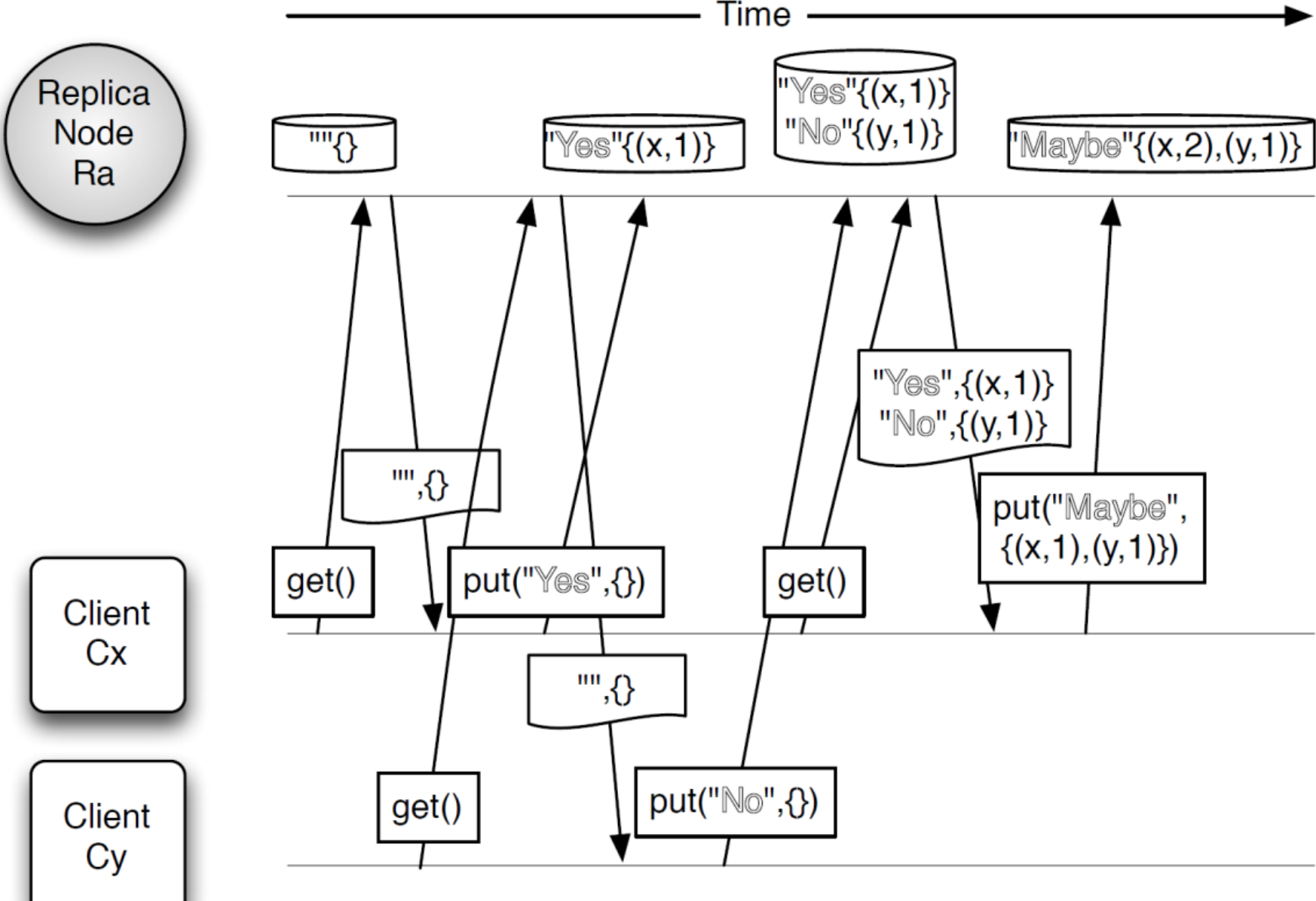
Singhal-Kshemkalyani Optimization

- Goal: reduce the size of VC sent along with messages
- Main idea: Send compressed timestamps (only those components that have changed since the last message to the same recipient)
- BUT: Each process must remember the VC they last sent to other processes.
 - $O(n^2)$ storage overhead
- Singhal-Kshemkalyani technique cuts this overhead down to $O(n)$
 - Maintain Last-sent and Last-updated vectors, which contain only single timestamps last sent/updated to/from each process, instead of entire vectors

Version Vectors

- Conventional VC: One entry for every process, including clients!
 - Can be in the millions
- Version vectors: maintained by each server replica for each object
- N replicas for an object, then $V_j[i]$ is number of updates generated at replica i that is known by replica-j
- Can be expanded/pruned as replicas leave and join
 - Zero-valued entries add no distinguishing information
- Dynamic Version Vector Maintenance. David Ratner, Peter Reiher, and Gerald Popek

Version Vector Example



Causal Order Broadcast

Vector Timestamps and Causality

- We have looked at *total order* of messages where all messages are processed in the same order at each process.
- It is possible to have *any order* where all you care about is that a message reaches all processes, but you don't care about the order of execution.
- *Causal order* is used when a message received by a process can potentially affect any subsequent message sent by that process. Those messages should be received in that order at all processes. Unrelated messages may be delivered in any order.

Display from a Bulletin Board Program

- Users run bulletin board applications which multicast messages
- One multicast group per topic (e.g. *os.interesting*)
- Require reliable multicast - so that all members receive messages
- Ordering:

total (makes
the numbers
the same at
all sites)

Bulletin board: <i>os.interesting</i>		
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		

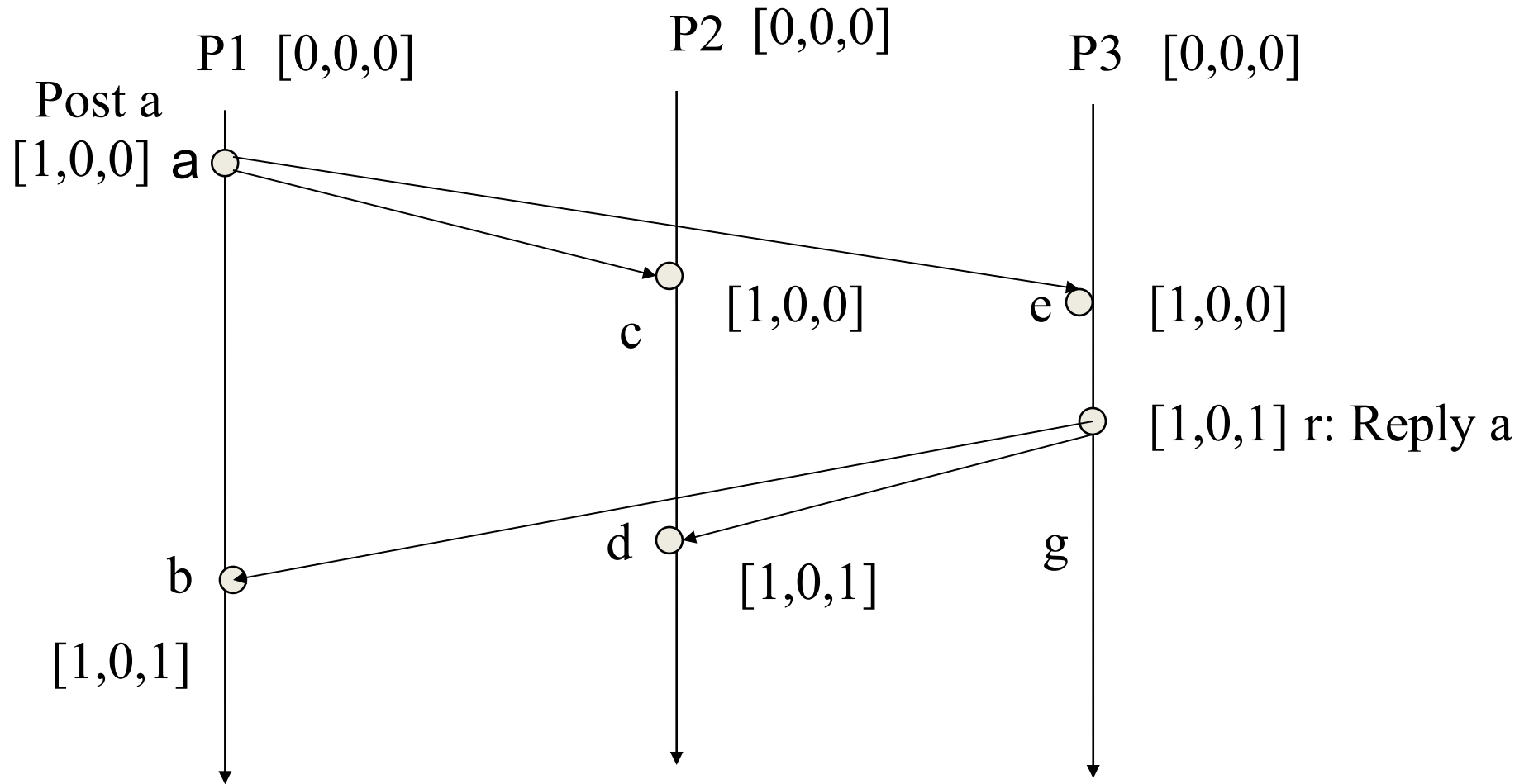
causal (makes
replies come after
original message)

FIFO (gives sender order)

Example Application: Bulletin Board

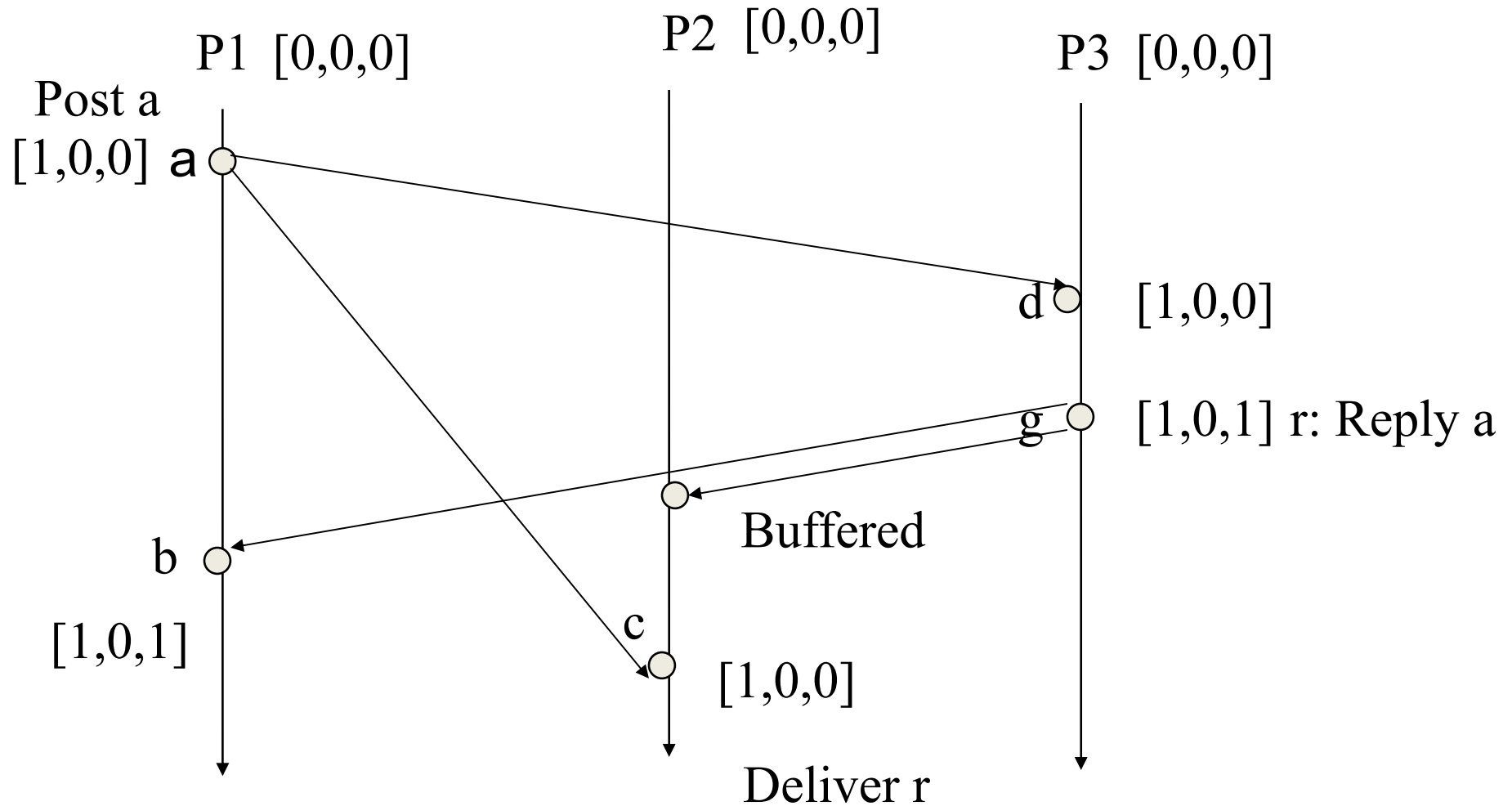
- A totally-ordered multicasting scheme does not imply that if message B is delivered after message A, that B is a reaction to A.
- Totally-ordered multicasting is too strong in this case.
- The receipt of an article causally precedes the posting of a reaction. The receipt of the reaction to an article should always follow the receipt of the article.

Example Application: Bulletin Board



Message *a* arrives at P2 before the reply *r* from P3 does

Example Application: Bulletin Board



The message *a* arrives at P2 after the reply from P3; The reply is not delivered right away.

END