

Algorithm 3.15 *The Birman-Schiper-Stephenson algorithm for causal message ordering of broadcast messages* [5].

Idea: This algorithm uses vector logical clocks implemented by the variable V in every process, initialized to all zeroes, for the sole purpose of causal message ordering, but not for affixing timestamps to other events in the system. The links do not have to be FIFO. Every process numbers its broadcasts consecutively in its own component of its vector logical clock and sends this vector along with the message to all other processes. If on receipt of a message a process finds that it cannot yet be delivered, it puts it into a buffer B of pending messages (along with their timestamps). We define the following condition for a message m that carries timestamp V_m and is received from process P_j to indicate whether it can be delivered:

$$D_j(m) = (V + e_j \geq V_m).$$

Condition $D_j(m)$ says that message m is the message expected next from P_j , and that the receiving process is at least as up to date with respect to all other processes as P_j was when it sent m .

Implementation:

I. Broadcasting a message

$V \leftarrow V + e_i$
broadcast(m, V)

III. Delivering a message from P_j

deliver(m):
 $\text{deliver}(m)$
 $V \leftarrow V + e_j$
 remove (m, V_m) from B

II. Receiving a message from P_j

upon receipt of (m, V_m) **do**
 if $D_j(m)$ **then**
 $\text{deliver}(m)$
 while ($\{(m, k, V_m) \in B \mid D_k(m)\} \neq \emptyset$) **do**
 $\text{deliver}(m)$ with (m, V_m) $\in B$ such that $D_j(m)$
 else add (m, j, V_m) to B

Correctness: First we prove the safety of the algorithm, that is, we prove that if $m(m_1) \rightarrow m(m_2)$, then $d_i(m_1) \rightarrow d_i(m_2)$ for $i = 1, \dots, n$. Because predicate D_j already checks the order of the broadcasts from the same process, we can assume that m_1 is broadcast by process P_j and m_2 is broadcast by process P_k , with $j \neq k$. Because $m(m_1) \rightarrow m(m_2)$, $V(m_1)[j] \leq V(m_2)[j]$. By code fragment III, the only modification made to the vector logical time in the receiving process when a message from process P_j is delivered is adding e_j to its vector clock. Now m_2 can only be delivered in P_i when $V_i[j] \geq V(m_2)[j]$, but $V_i[j]$ can only attain this value after m_1 has been delivered (V_i is here the vector clock in the receiving process P_i).

The proof of the liveness of the algorithm is left as an exercise for the reader. \square

